Introduction to Pattern Recognition MTSC 852, Fall 2015

Sokratis Makrogiannis, Ph.D.

Department of Mathematical Sciences Delaware State University smakrogiannis@desu.edu

October 20, 2015

S. Makrogiannis (DSU)

Introduction to Pattern Recognition

October 20, 2015 1 / 30

< □ > < □ > < 三 > < 三 > < 三 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Outline

- 1 Pattern Recognition Basics
- 2 Pattern Recognition Systems
- 3 Pattern Recognition System Design Cycle
- 4 Learning Techniques

S. Makrogiannis (DSU)

Introduction to Pattern Recognition

October 20, 2015 2 / 30

< □ > < □ > < 三 > < 三 > < 三 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Pattern Recognition

Definitions

- The research field that solves problems for recognition and classification of objects represented by sets of measurements
- \circ "The process of giving names ω to observations x" -Schümann
- "A problem of estimating density functions in a high-dimensional space and dividing the space into regions of catergories or classes" -Fukunaga
- "The assignment of a physical object or event to one of pre-specified categories" -Duda and Hart

S. Makrogiannis (DSU)

Introduction to Pattern Recognition

October 20, 2015 3 / 30

イロト イポト イヨト イヨト

Components of a Pattern Recognition System



Pattern Recognition Systems

Main components of Brain Atrophy Recognizer



S. Makrogiannis (DSU)

Introduction to Pattern Recognition

< □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ♪ < ○ へ () n October 20, 2015 5 / 30

Sensing

- Different sensors mainly transducers can be used such as cameras, microphone arrays, MRI or CT scanners, fingerprint devices, etc
- Characteristics and limitations of transducers may pose challenges to pattern recognition
- Usual characteristics are bandwidth, resolution, quantization, electronic noise, distortions, signal-to-noise ratio, etc

S. Makrogiannis (DSU)

Introduction to Pattern Recognition

October 20, 2015 6 / 30

Sac

A = A = A = A = A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

Segmentation and Grouping

- The goal of this step is to separate the different objects
- That is, to delineate visual objects in images or videos, separate characters in manual scripts, separate phonemes for voice recognition, etc
- Segmentation is a complicated problem that depends on the domain



S. Makrogiannis (DSU)

Introduction to Pattern Recognition

October 20, 2015 7 / 30

Feature Extraction/Dimensionality Reduction

- Goal of a feature extractor is to create object representations that are similar for objects of the same class and dissimilar for objects of different classes
- Therefore, we are seeking to compute distinguishing, or discriminant, features





Feature Extraction/Dimensionality Reduction

- Another desired property of features mainly extracted from images is invariance to geometric transformations such as translation, rotation scaling, or nonlinear deformations, and resiliency to occlusions
- In speech recognition we are looking for features that are invariant to time translations and to changes in amplitude
- Feature computation is also domain-specific
- A related research area deals with the selection/computation of the most discriminant features for classification; this is called dimensionality reduction

S. Makrogiannis (DSU)

Introduction to Pattern Recognition

Sac

イロト イポト イヨト イヨト 二日

Feature Extraction Example

• Computation of texture features from an abdominal CT scan



< □ > < 同 >

S. Makrogiannis (DSU)

Introduction to Pattern Recognition

October 20, 2015 10 / 30

Э

Sac

Classification

- The classifier uses as input the feature vector that represents an object and produces a category label for the object
- Classification is a major part of this work



S. Makrogiannis (DSU)

Introduction to Pattern Recognition

October 20, 2015 11 / 30

Sac

イロト イポト イヨト イヨト

Classification

- One challenge in classification is the variability in feature values that reduce the separation between different categories
- Variability may be caused by complexity or noise
- The term noise refers to any randomness in measurements that our model cannot account for



S. Makrogiannis (DSU)

Introduction to Pattern Recognition

October 20, 2015 12 / 30

Classification Decision and Performance Evaluation

- This component uses the classifier's output to make a decision
- One criterion for making a classification decision is the error rate. The error rate is the percentage of new patterns assigned to the wrong category. Therefore we seek to minimize the error rate in our decision
- Often times we associate a risk with each classification decision. This can be used to compute a total cost that will help make a decision
- More recent approaches multiple classifiers and combine their decisions to reduce the classification error rates



S. Makrogiannis (DSU)

Introduction to Pattern Recognition

October 20, 2015 13 / 30

Pattern Recognition Problems

- Machine vision Visual inspection Security
- Character recognition Automated mail sorting Signature verification
- Computer-aided diagnosis Breast cancer diagnosis EEG, ECG signal analysis
- Speech recognition Human Computer Interaction systems



S. Makrogiannis (DSU)

Introduction to Pattern Recognition

October 20, 2015 14 / 30

Sar

Pattern Recognition System Design Cycle

- In this section we outline the procedure that is followed to design a pattern classification system
- It is called a cycle because the design sequence is repeated until we reach the desired outcome



S. Makrogiannis (DSU)

Introduction to Pattern Recognition

October 20, 2015 15 / 30

Sar

Data Collection

- Data collection has a significant impact on the classifier design in terms of feasibility and costs
- We often conduct a feasibility study with a small set of data to design our system, then perform a full scale study to estimate performance
- One crucial question is, how many training and testing samples are needed to design a system that can be used for real applications?

S. Makrogiannis (DSU)

Introduction to Pattern Recognition

October 20, 2015

イロト イ理ト イヨト イヨト

16 / 30

Pattern Recognition System Design Cycle

What Features to Compute

- This is a critical and domain-specific problem
- Prior knowledge may help in this stage. For example, when we design a face recognition system, we can use our knowledge of face structure like topological relations between the nose, lips, and eyes

S. Makrogiannis (DSU)

Introduction to Pattern Recognition

October 20, 2015 17 / 30

・ロト ・ 同ト ・ ヨト ・ ヨト

Which Models to Use

- In this stage we develop criteria for the selection of mathematical models to represent each class
- The main question is how can we measure or predict if a hypothesized model approximates well the underlying true model of our patterns?
- Usually the answer depends on the domain; we can study the complexity, dimensionality, variability, and population of patterns

S. Makrogiannis (DSU)

Introduction to Pattern Recognition

< □ ▶ < @ ▶ < 重 ▶ < 重 ▶ 三 のへの n October 20, 2015 18 / 30

Training Stage

- Training is the process of determining the classifier using data that we call training data
- These are called learning from example techniques

S. Makrogiannis (DSU)

Introduction to Pattern Recognition

October 20, 2015 19 / 30

< ∃ >

I > <
 I >
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I
 I

Performance Evaluation

- This is a necessary component for determining the feature set, and classification model.
- An additional objective is to identify the potential for further improvement and to address overfitting



S. Makrogiannis (DSU)

Introduction to Pattern Recognition

October 20, 2015 20

20 / 30

Computational Complexity

- One consideration is how an algorithm scales as a function of the feature dimensionality, number of objects, or number of categories
- Sometimes the performance of a recognizer is excellent but it is not possible to implement
- We are mostly concerned with the complexity for making a decision; not so much for the complexity of learning stage that is executed off-line

S. Makrogiannis (DSU)

Introduction to Pattern Recognition

October 20, 2015 21 / 30

200

イロト イポト イヨト イヨト 二日

Learning Techniques

Main categories

- Supervised learning
- Unsupervised learning
- Reinforcement learning

S. Makrogiannis (DSU)

Introduction to Pattern Recognition

October 20, 2015 22 / 30

3

코 ト - モ ト

Supervised Learning

 In supervised learning, we use a training data set with patterns of known categories so we can fit our models and/or adapt the classifier parameters



S. Makrogiannis (DSU)

Introduction to Pattern Recognition

October 20, 2015 23 / 30

Sac

イロト イポト イヨト イヨト

Unsupervised Learning

- In unsupervised learning there exists no training set and the system learns forms clusters of the input patterns using optimization algorithms
- Relevant questions:

How do we select the number of clusters? How do we select the cluster models?



S. Makrogiannis (DSU)

Introduction to Pattern Recognition

October 20, 2015 24 / 30

Reinforcement Learning

• In this type of learning, the only given feedback is that the tentative category is right or wrong, without further information

2 SO HAA rs :: Select Region - Alcount Internet Explorer			NERS ro :: Content based res		
the Ed. yes Parates July 2	*	*	fill yes Fportes (sell)		
Q 🗤 + 🖸 + 💽 🗟 🐔	👷 taudes 🥑 😥 - 🍒 🗟 - 🔜 🕮 🖄	0	54 + 🔘 - 🗷 🗟 🐔	👷 Taratas 🤣 🔒 😼 🔛 🖢	E 3
Address 🛃 May (Steeland) op Schematicity	dans/udjaar0.Midling.ph-12006	🗶 🚺 🏟 1992 * Appen	ii 🛃 His (hela li affidenati);	dans) n/aarð Milling (ör i 2008kaljoffin Hayda	-essy-be 🛛 🖉 🕼 Loka
SCHEMAxm		SAN		Content-based results	
home about contact			Former A		
Espin + Narss + Ross + Espin + Espin + Adver +	Select Region:		Adout 10 Adout 10 Durante 10 Adout 10	Colora Inaz	100.010
		MC (10)			*
	× ×	~		10604.34 NEAR DECEMENT NEX-10(a).349 Red Server Red Server	309 WAYN GALER AN MANN
	84(MA) H2(34) V8			Y	~ *
(W. LINCOL) (W. LINCOL)	Inspire an open pro-	Cont Carp. Andres Agency Social State Press Agency		1000-Fair 0471% 06014 and 100-101a-0-0 800-101a- Rad Service Red Service	10.714 00003 pg 00.074 309 00-056-09 Feet Sense
				*	-
				10800-last 04.81% 000-1 Jay 103-12 Jay, 310 Rad Sanjar Rad Sanjar	49 0011 yr 0011 yr 00.414
				shrine (\$2245 Retr	Ontroppe 2 Min
Vip (heda A. p(bhenal () yitem) o	Annoh-Northad John	· iterat			Contract Contract

S. Makrogiannis (DSU)

Introduction to Pattern Recognition

October 20, 2015 25 / 30

イロト イ理ト イヨト イヨト

Sac

Problem Description

- Suppose that we want to build a system for automatic recognition of sea bass and salmon in a fish packing plant
- The fish are placed on a conveyor belt and a fixed camera acquires images of fish on the moving conveyor belt
- We decide to use the physical characteristics such as length, width, brightness, position of the mouth of fish to separate the two species



S. Makrogiannis (DSU)

Introduction to Pattern Recognition

October 20, 2015 26 / 30

- In pattern recognition terms, we are asked to build a classifier of fish into two categories, that is salmon and sea bass
- The physical characteristics that we want to measure are called features
- To describe each fish type we develop mathematical models
- We find the best matching model for each fish to perform classification

S. Makrogiannis (DSU)

Introduction to Pattern Recognition

October 20, 2015 27 / 30

Main stages

- Pre-processing: this stage can be used to improve image quality, reduce image noise, apply segmentation to separate the fish, so we can compute more accurate features
- Feature computation or extraction is the process of measuring the characteristics of fish from the sensor data, for example lightness, length, width, number of fins, etc
- The model refers to the probability density model for sea bass and for salmon and the type of boundary between these classes, that is linear, quadratic, multi-modal, etc
- Performance evaluation is the process of calculating the probability of erroneous classification decisions
- Generalization studies if the selected recognizer will sustain a good performance under real conditions, i.e., a diverse set of unlabeled patterns

Feature Extraction

• Feature computation or extraction is the process of measuring the characteristics of fish from the sensor data, for example lightness, length, width, number of fins, etc



S. Makrogiannis (DSU)

Introduction to Pattern Recognition

October 20, 2015 29 / 30

Classifier Generalization

- The model refers to the probability density model for sea bass and for salmon and the type of boundary between these classes, that is linear, quadratic, multi-modal, etc
- Generalization studies if the selected recognizer will maintain a good performance level under real conditions, i.e., a diverse set of unlabeled patterns



S. Makrogiannis (DSU)

Introduction to Pattern Recognition

October 20, 2015 30 / 30

Probability Theory Review MTSC 852, Fall 2015

Sokratis Makrogiannis, Ph.D.

Department of Mathematical Sciences Delaware State University smakrogiannis@desu.edu

October 20, 2015

S. Makrogiannis (DSU)

Probability Theory Review

October 20, 2015 1 / 12

Outline

- Notation and Preliminaries
- Inner Product 2
- 3 Outer Product
- 4 Derivatives of Matrices
- (5) Determinant, Trace, Rank
- 6 Matrix Inversion
- **Eigenvectors and Eigenvalues** (7)

S. Makrogiannis (DSU)

Probability Theory Review

3 October 20, 2015 2 / 12

Sac

イロト イポト イヨト イヨト

Notation and Preliminaries

• d-dimensional vector:
$$\mathbf{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$$

Transpose of \mathbf{x} : $\mathbf{x}^T = (x_1 x_2 \dots x_d)$
• $n \times d$ matrix: $M = \begin{pmatrix} m_{11} & m_{12} & \dots & m_{1d} \\ m_{21} & m_{22} & \dots & m_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ m_{n1} & m_{n2} & \dots & m_{nd} \end{pmatrix}$
Transpose of M : $M^T = \begin{pmatrix} m_{11} & m_{21} & \dots & m_{n1} \\ m_{12} & m_{22} & \dots & m_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ m_{1d} & m_{2d} & \dots & m_{nd} \end{pmatrix}$
Symmetric matrix, if $m_{ij} = m_{ji}$
Antisymmetric matrix, if $m_{ij} = -m_{ji}$

October 20, 2015 3 / 12

▲ 臣 ▶ 臣 • • ○ < (~

Notation and Preliminaries

• Multiply vector by matrix: $M\mathbf{x} = \mathbf{y}$

$\binom{m_{11}}{m_{11}}$	m_{12}		m_{1d}	$\left(\begin{array}{c} x_1 \end{array} \right)$		(y_1)
m_{21}	m_{22}	•••	m _{2d}	<i>x</i> ₂		<i>y</i> 2
÷	÷		:		=	÷
$\int m_{n1}$	m _{n2}		m _{nd})	$\left(x_n \right)$		\ y _n /

or,

$$y_i = \sum_{j=1}^d m_{ij} x_j$$

S. Makrogiannis (DSU)

Probability Theory Review

October 20, 2015 4 / 12

Inner Product

• The inner product of two vectors **x**, **y** is a scalar:

$$\mathbf{x} \cdot \mathbf{y} = \mathbf{x}^T \mathbf{y} = \mathbf{y}^T \mathbf{x} = \sum_{i=1}^d x_i y_i$$

- Euclidean norm: $\|\mathbf{x}\| = (\mathbf{x}^T \mathbf{x})^{1/2}$
- Angle between **x** and **y**: $\cos \theta = \frac{\mathbf{x}^T \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|}$
- $\mathbf{x}^T \mathbf{y} = \mathbf{0} \rightarrow \text{vectors are orthogonal}$
- $\|\mathbf{x}^{\mathsf{T}}\mathbf{y}\| = \|x\|\|y\| \rightarrow \text{vectors are colinear}$
- Cauchy Schwartz inequality: $\|\mathbf{x}^T \mathbf{y}\| \le \|x\| \|y\|$
- Vectors $[x_1, x_2, ..., x_n]$ are linearly independent, if no vector in the set can be written as a linear combination of any of the others

S. Makrogiannis (DSU)

Probability Theory Review

< □ ト < □ ト < 三 ト < 三 ト < 三 ト う Q () October 20, 2015 5 / 12

Outer Product

• Multiply vector by matrix: $M = \mathbf{x}\mathbf{y}^T$

$$\begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix} (y_1 y_2 \dots y_n) = \begin{pmatrix} x_1 y_1 & x_1 y_2 & \dots & x_1 y_n \\ x_2 y_1 & x_2 y_2 & \dots & x_2 y_n \\ \vdots & \vdots & \ddots & \vdots \\ x_d y_1 & x_d y_2 & \dots & x_d y_n \end{pmatrix}$$

S. Makrogiannis (DSU)

Probability Theory Review

October 20, 2015 6 / 12
Derivatives of Matrices

• Let
$$f(\mathbf{x})$$
 be a scalar-valued function of d variables $\mathbf{x} = (x_1 x_2 \dots x_d)^T$
Gradient of $f: \nabla f(\mathbf{x}) = \begin{pmatrix} \partial f/\partial x_1 \\ \vdots \\ \partial f/\partial x_d \end{pmatrix}$
• Let $\mathbf{f}(\mathbf{x}): \mathbb{R}^d \to \mathbb{R}^n$. Its Jacobian matrix is:
 $\mathbf{J}(\mathbf{x}) = \frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}} = \begin{pmatrix} \partial f_1/\partial x_1 & \dots & \partial f_1/\partial x_d \\ \vdots & \vdots & \ddots \\ \partial f_n/\partial x_1 & \dots & \partial f_n/\partial x_d \end{pmatrix}$
• Let matrix M with elements dependent on θ . Its derivative in θ is :
 $\frac{\partial M}{\partial \theta} = \begin{pmatrix} \partial m_{11}/\partial \theta & \dots & \partial m_{1d}/\partial \theta \\ \vdots & \vdots & \ddots \\ \partial m_{n1}/\partial \theta & \dots & \partial m_{nd}/\partial \theta \end{pmatrix}$
Derivative of inverse: $\frac{\partial}{\partial \theta} M^{-1} = -M^{-1}\frac{\partial M}{\partial \theta} M^{-1}$
Derivative identities:
 $\frac{\partial}{\partial \mathbf{x}}[M\mathbf{x}] = M, \quad \frac{\partial}{\partial \mathbf{x}}[\mathbf{y}^T\mathbf{x}] = \frac{\partial}{\partial \mathbf{x}}[\mathbf{x}^T\mathbf{y}] = \mathbf{y}, \quad \frac{\partial}{\partial \mathbf{x}}[\mathbf{x}^T M\mathbf{x}] = [M + M^T]\mathbf{x}$
• S. Makrogiannis (DSU) Probability Theory Review October 20, 2015 7/12

Taylor Expansions

• Taylor expansion for f(x) about x_0 :

$$f(x) = f(x_0) + \frac{df}{dx} \bigg|_{x=x_0} (x-x_0) + \frac{1}{2!} \frac{d^2 f}{dx^2} \bigg|_{x=x_0} (x-x_0)^2 + O\left((x-x_0)^3\right)$$

• Taylor expansion for vector $f(\mathbf{x})$ about \mathbf{x}_0 :

$$f(\mathbf{x}) = f(\mathbf{x}_0) + \left[\frac{df}{d\mathbf{x}}\right]_{\mathbf{x}=\mathbf{x}_0}^T (\mathbf{x}-\mathbf{x}_0) + \frac{1}{2!} (\mathbf{x}-\mathbf{x}_0)^T \left[\frac{d^2f}{d\mathbf{x}^2}\right]_{\mathbf{x}=\mathbf{x}_0}^T (\mathbf{x}-\mathbf{x}_0) + O\left(\|\mathbf{x}-\mathbf{x}_0\|^3\right)_{\mathbf{x}=\mathbf{x}_0}^T (\mathbf{$$

S. Makrogiannis (DSU)

Probability Theory Review

October 20, 2015 8 / 12

Determinant, Trace, Rank

Determinant of a d×d matrix A: |A| = ∑_{k=1}^d a_{ik}|A_{i|k}|(-1)^{k+i} A_{i|k}: minor formed by removing the *i*th row and *k*th column of A
 Trace is the sum of diagonal elements:

$$tr[A] = \sum_{k=1}^{d} a_{kk}$$

- Rank of a matrix is the number of linearly independent rows or columns
- A square matrix is non-singular ↔ its rank equals the number of rows (or columns)
- A non-singular matrix has a non-zero determinant

S. Makrogiannis (DSU)

Probability Theory Review

October 20, 2015 9 / 12

◆□▶ ◆母▶ ◆ヨ▶ ◆ヨ▶ ヨー つくや

Matrix Inversion

- Suppose square matrix $M_{d \times d}$
- Inverse of M is M^{-1} : $MM^{-1} = M^{-1}M = \mathbb{I}$ M^{-1} exists $\leftrightarrow M$ is non-singular M^{-1} can be also written as $M^{-1} = \frac{Adj[M]}{|M|}$

Adj[M]: matrix whose i, j entry equals to cofactor $C_{i,j} = (-1)^{j+i} M_{j|i}$

• Pseudo-inverse of M, denoted by M^{\dagger} , is used when A^{-1} does not exist, or when M is not square If $M^{T}M$ is non-singular, pseudoinverse M^{\dagger} is $M^{\dagger} = [M^{T}M]^{-1}M^{T}$ Observe that $M^{\dagger}M = \mathbb{I}$ We use pseudo-inverse to solve least squares problem Inverse of product of square matrices M, N: $[MN]^{-1} = N^{-1}M^{-1}$

S. Makrogiannis (DSU)

Probability Theory Review

Eigenvectors and Eigenvalues

- Suppose square matrix $M_{d \times d}$
- Frequently, we need to solve linear equations of the form

$$M\mathbf{x} = \lambda \mathbf{x} \Rightarrow (M - \lambda \mathbb{I})\mathbf{x} = 0 \Rightarrow egin{casel} \mathbf{x} = 0 & ext{trivial case} \ M - \lambda \mathbb{I} = 0 & ext{non-trivial case} \end{bmatrix}$$

$$M - \lambda \mathbb{I} = 0 \Rightarrow |M - \lambda \mathbb{I}| = 0 \Rightarrow \lambda^d + a_1 \lambda^{d-1} + \ldots + a_{d-1} \lambda + a_d = 0$$

For each root we solve the set of linear equations

- Solution vectors e_i are called eigenvectors, while corresponding solution scalars λ_i are called eigenvalues
- M is real and symmetric \rightarrow there exist d solution vectors $\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_d\}$ and d corresponding eigenvalues $\{\lambda_1, \lambda_2, \dots, \lambda_d\}$
- When multiplied by M, the eigenvectors change in magnitude but not orientation: $M\mathbf{e}_i = \lambda_i \mathbf{e}_i$

S. Makrogiannis (DSU)

Probability Theory Review

イロト (母) (ヨ) (ヨ) (ヨ) () () October 20, 2015

11 / 12

Eigenvectors and Eigenvalues Properties

- Suppose that we have found the eigendecomposition: $M\mathbf{e}_i = \lambda_i \mathbf{e}_i$. The following properties apply
 - $tr[M] = \sum_{i=1}^{d} \lambda_i$

•
$$|M| = \prod_{i=1}^d \lambda_i$$

- M is non-singular ightarrow all eigenvalues are non-zero
- M is real and symmetric \rightarrow all eigenvalues are real and eigenvectors are orthogonal
- M is positive definite ightarrow all eigenvalues are positive

S. Makrogiannis (DSU)

Probability Theory Review

October 20, 2015 12 / 12

イロト (母) (ヨ) (ヨ) (ヨ) () ()

Probability Theory Review MTSC 852, Fall 2015

Sokratis Makrogiannis, Ph.D.

Department of Mathematical Sciences Delaware State University smakrogiannis@desu.edu

September 27, 2015

S. Makrogiannis (DSU)

Probability Theory Review

September 27, 2015 1 / 14



Discrete Random Variables

S. Makrogiannis (DSU)

Probability Theory Review

September 27, 2015 2 / 14

▲ロト ▲園 ト ▲ 臣 ト ▲ 臣 ト 二 臣 … の Q ()



2 Pairs of Discrete Random Variables

S. Makrogiannis (DSU)

Probability Theory Review

September 27, 2015 2 / 14

- 1 Discrete Random Variables
- 2 Pairs of Discrete Random Variables
- ③ Vector Random Variables

S. Makrogiannis (DSU)

Probability Theory Review

September 27, 2015 2 / 14

- 1 Discrete Random Variables
- 2 Pairs of Discrete Random Variables
- 3 Vector Random Variables
- 4 Continuous Random Variables

S. Makrogiannis (DSU)

Probability Theory Review

September 27, 2015 2 / 14

イロト (母) (ヨ) (ヨ) (ヨ) () ()

- 1 Discrete Random Variables
- 2 Pairs of Discrete Random Variables
- 3 Vector Random Variables
- 4 Continuous Random Variables
- 5 Normal Distributions

S. Makrogiannis (DSU)

Probability Theory Review

September 27, 2015 2 / 14

Discrete Random Variables

Let x be a random variable with x ∈ X and X = {v₁, v₂,..., v_m}
Then p_i is probability of x = v_i

$$p_i = Pr[x = v_i], i = 1, \dots, m$$

Axioms

$$p_i \ge 0$$

 $\sum_{i=1}^m p_i = 1$

Probability Mass Functions

 $P(x) \ge 0$ $\sum_{x \in \mathscr{X}} P(x) = 1$

S. Makrogiannis (DSU)

Probability Theory Review

September 27, 2015 3 / 14

Expected Values

• Mean of random variable x:

$$\mathbf{E}[x] = \mu = \sum_{x \in \mathscr{X}} x P(x) = \sum_{i=1}^{m} v_i p_i$$

• For a function f(x):

$$\mathbf{E}[f(x)] = \sum_{x \in \mathscr{X}} f(x) P(x)$$

• Variance
$$E\left[(x-\mu)^2\right]$$

 $E\left[(x-\mu)^2\right] = \sum_{x \in \mathscr{X}} (x-\mu)^2 P(x) = Var[x] = \sigma^2$

We can show that

$$Var[x] = E[x^2] - [E[x]]^2$$

S. Makrogiannis (DSU)

Probability Theory Review

September 27, 2015 4 / 14

Pairs of Discrete Random Variables

- Let x, y be random variables such that $x \in \mathscr{X} = \{v_1, v_2, \dots, v_m\}$ and $y \in \mathscr{Y} = \{w_1, w_2, \dots, w_n\}$
- Joint probability:

$$p_{ij} = Pr[x = v_i, y = w_j]$$

Axioms:

$$p_{ij} \geq 0, \quad \sum_{i=1}^m \sum_{j=1}^n p_{i,j} = 1$$

• Joint probability mass function:

$$P(x,y) \ge 0, \quad \sum_{x \in \mathscr{X}} \sum_{y \in \mathscr{Y}} P(x,y) = 1$$

• Marginal distributions:

$$P_x(x) = \sum_{y \in \mathscr{Y}} P(x, y), \quad P_y(y) = \sum_{x \in \mathscr{X}} P(x, y)$$

S. Makrogiannis (DSU)

Probability Theory Review

September 27, 2015 5 / 14

◆□▶ ◆母▶ ◆ヨ▶ ◆ヨ▶ ヨー つくや

Statistical Independence

Definition (Statistically independent variables) Random variables x, y are statistically independent, if and only if

$$P(x,y) = P_x(x) \cdot P_y(y)$$

S. Makrogiannis (DSU)

Probability Theory Review

September 27, 2015 6 / 14

< ロト < 同ト < ヨト < ヨト

Expected Values of Functions of Two Variables

• For a function f(x,y):

$$\mathbb{E}[f(x,y)] = \sum_{x \in \mathscr{X}} \sum_{y \in \mathscr{Y}} f(x,y) P(x,y)$$

Means:

$$E[\mathbf{x}] = \mu_{\mathbf{x}} = \sum_{\mathbf{x} \in \mathscr{X}} \sum_{\mathbf{y} \in \mathscr{Y}} \mathbf{x} P(\mathbf{x}, \mathbf{y}), \quad E[\mathbf{y}] = \mu_{\mathbf{y}} = \sum_{\mathbf{x} \in \mathscr{X}} \sum_{\mathbf{y} \in \mathscr{Y}} \mathbf{y} P(\mathbf{x}, \mathbf{y})$$
$$E[\mathbf{x}] = \mu = \sum_{\mathbf{x} \in \mathscr{X} \times \mathscr{Y}} \mathbf{x} P(\mathbf{x})$$

Variances:

$$Var[x] = \sigma_x^2 = E\left[(x - \mu_x)^2\right], \quad Var[y] = \sigma_y^2 = E\left[(y - \mu_y)^2\right]$$

Covariance:

$$\sigma_{xy}^2 = \mathbb{E}\left[(x - \mu_x)(y - \mu_y)\right] = \sum_{x \in \mathscr{X}} \sum_{y \in \mathscr{Y}} (x - \mu_x)(y - \mu_y) P(x, y)$$

S. Makrogiannis (DSU)

Probability Theory Review

September 27, 2015 7 / 14

Expected Values of Functions of Two Variables

Corollary

x,y stat. independent $ightarrow \sigma_{xy} = 0$

Corollary

 $\sigma_{xy} = 0
ightarrow x, y$ uncorrelated

Result $x, y \text{ stat. independent} \rightarrow$ E[f(x)g(y)] = E[f(x)]E[g(y)]

- Cauchy-Schwartz inequality:
 - $\sigma_{xy}^2 \leq \sigma_x^2 \sigma_y^2, \quad \rho \in [-1, 1]$
- Correlation Coefficient ρ:

$$\sigma = \frac{\sigma_{xy}}{\sigma_x \sigma_y}$$

• Covariance matrix:

$$\Sigma = \mathrm{E}\left[(\mathbf{x} - \mu_{\mathbf{x}})(\mathbf{y} - \mu_{\mathbf{y}})^{T}\right]$$

S. Makrogiannis (DSU)

Probability Theory Review

Conditional Probability

• Conditional probability of x given y:

$$Pr[x = v_i | y = w_j] = \frac{Pr[x = v_i, y = w_j]}{Pr[y = w_j]}$$

By use of probability mass:

$$P(x|y) = \frac{P(x,y)}{P_y(y)}$$

Result

x, y stat. independent $\rightarrow P(x|y) = P(x)$

S. Makrogiannis (DSU)

Probability Theory Review

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Law of Total Probability and Bayes Rule

Law of total probability: P(y) = ∑_{x∈X} P(x,y) = ∑_{x∈X} P(y|x)P(x)
Based on conditional probability definition:

$$P(x,y) = P(y|x)P(x)$$

$$P(x,y) = P(x|y)P(y) \} \Rightarrow P(x|y) = \frac{P(y|x)P(x)}{P(y)}$$

Bayes rule:

$$P(x|y) = rac{P(y|x)P(x)}{P(y)}$$
 posterior $= rac{likelihood imes prior}{evidence}$

S. Makrogiannis (DSU)

Probability Theory Review

September 27, 2015 10 / 14

◆□▶ ◆母▶ ◆ヨ▶ ◆ヨ▶ ヨー つくや

Vector Random Variables

- We use vector notation to generalize to more variables, i.e. $\mathbf{x} = \{x_1, x_2, \dots, x_d\}$
- Joint probability mass function $P(\mathbf{x})$ Axioms

$$P(\mathbf{x}) \ge 0, \quad \sum_{\mathbf{x} \in \mathbb{R}^d} P(\mathbf{x}) = 1$$

• Statistically independent $\mathbf{x}_i: P(\mathbf{x}) = \prod_{i=1}^d P_{\mathbf{x}_i}(\mathbf{x}_i)$

• Bayes rule:
$$P(\mathbf{x_1}|\mathbf{x_2}) = \frac{P(\mathbf{x_2}|\mathbf{x_1})P(\mathbf{x_1})}{P(\mathbf{x_2})} = \frac{P(\mathbf{x_2}|\mathbf{x_1})P(\mathbf{x_1})}{\sum_{\mathbf{x_1} \in \mathscr{X}_1} P(\mathbf{x_2}|\mathbf{x_1})P(\mathbf{x_1})}$$

Example

$$P(x_1, x_4) = \sum_{x_2} \sum_{x_3} \sum_{x_5} P(x_1, x_2, x_3, x_4, x_5)$$
$$P(x_1, x_2 | x_3) = \frac{P(x_1, x_2, x_3)}{P(x_3)}$$

S. Makrogiannis (DSU)

Probability Theory Review

Continuous Random Variables

Single random variable

- Let x be a continuous random variable
- Probability density function (PDF) p(x): $P_r[x \in (a,b)] = \int_a^b p(x) dx$
- PDF must satisfy: $p(x) \ge 0$, $\int_{-\infty}^{\infty} p(x) dx = 1$
- Expected value: $E[f(x)] = \int_{-\infty}^{\infty} f(x)P(x)dx$

• Mean:
$$E[x] = \int_{-\infty}^{\infty} x P(x) dx$$

• Variance:
$$\mathrm{E}\left[(x-\mu)^2\right] = \int_{-\infty}^{\infty} (x-\mu)^2 P(x) dx$$

S. Makrogiannis (DSU)

Probability Theory Review

September 27, 2015 12 / 14

Continuous Random Variables

Multivariate case

- Let **x** be vector of continuous random variables $\mathbf{x} = \{x_1, x_2, \dots, x_d\}$
- Probability density function (PDF) $p(\mathbf{x})$: $P_r[\mathbf{x} \in \mathscr{X}_1 \times \ldots \times \mathscr{X}_d] = \int_{\mathbf{x} \in \mathscr{X}_1 \times \ldots \times \mathscr{X}_d} p(\mathbf{x}) d\mathbf{x}$
- PDF must satisfy: $p(\mathbf{x}) \geq 0$, $\int_{-\infty}^{\infty} p(\mathbf{x}) d\mathbf{x} = 1$
- Expected value: $E[f(\mathbf{x})] = \int_{-\infty}^{\infty} f(\mathbf{x}) P(\mathbf{x}) d\mathbf{x} = \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} f(\mathbf{x}) P(\mathbf{x}) dx_1 \cdots dx_d$
- Mean: $\mathrm{E}[\mathbf{x}] = \int_{-\infty}^{\infty} \mathbf{x} P(\mathbf{x}) d\mathbf{x}$
- Covariance matrix: $\Sigma = E\left[(\mathbf{x} - \mu_{\mathbf{x}})(\mathbf{y} - \mu_{\mathbf{y}})^T \right] = \int_{-\infty}^{\infty} (\mathbf{x} - \mu)(\mathbf{x} - \mu)^T P(\mathbf{x}) d\mathbf{x}$
- Stat. independent variables in $\mathbf{x} \Rightarrow p(\mathbf{x}) = \prod_{i=1}^{d} P(x_i)$, and Σ is diagonal

S. Makrogiannis (DSU)

Probability Theory Review

<□ ト < □ ト < □ ト < 三 ト < 三 ト < 三 ト ○ Q ○ September 27, 2015 13 / 14

Normal Distributions

Normal density (1-D)

- Probability density function: $p(x) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} = N(\mu, \sigma)$
- Expectations: $E[1] = \int_{-\infty}^{\infty} P(x) dx = 1$, $E[x] = \int_{-\infty}^{\infty} x P(x) dx = \mu$ $E\left[(x-\mu)^2\right] = \int_{-\infty}^{\infty} (x-\mu)^2 P(x) dx = \sigma^2$
- Approximations of cumulative probabilities: $Pr[|x \mu| \le \sigma] \simeq 0.68$, $Pr[|x - 2\mu| \le \sigma] \simeq 0.95$, $Pr[|x - 3\mu| \le \sigma] \simeq 0.997$
- Standardized random variable: $u = \frac{x-\mu}{\sigma} \Rightarrow p(u) = \frac{1}{\sqrt{2\pi\sigma}\sigma}e^{-\frac{u^2}{2}} = N(0,1)$

Theorem (Central Limit Theorem)

The distribution of the sum of d independent random variables approaches the normal distribution (under various conditions)

S. Makrogiannis (DSU)

Probability Theory Review

Bayesian Decision Theory MTSC 852, Fall 2015

Sokratis Makrogiannis, Ph.D.

Department of Mathematical Sciences Delaware State University smakrogiannis@desu.edu

September 3, 2015

S. Makrogiannis (DSU)

Bayesian Decision Theory

September 3, 2015 1 / 20



S. Makrogiannis (DSU)

Bayesian Decision Theory

September 3, 2015 2 / 20

▲ロト ▲園 ト ▲ 臣 ト ▲ 臣 ト 二 臣 … の Q ()





Bayesian Decision Theory for Continuous Features

S. Makrogiannis (DSU)

Bayesian Decision Theory

September 3, 2015 2 / 20

▲ロト ▲撮 ト ▲ 臣 ト ▲ 臣 ト 二 臣 … の � @





Bayesian Decision Theory for Continuous Features



Minimum-Error-Rate Classification

S. Makrogiannis (DSU)

Bayesian Decision Theory

September 3, 2015 2 / 20

Sac

イロト イポト イヨト イヨト 一日

Preliminaries

- The Bayesian Decision Theory is a statistical approach to pattern classification
- This approach assigns a class label to a pattern by quantifying the tradeoffs among the different decisions
- Bayesian Decision Theory poses the decision problem in probabilistic terms
- It assumes that all relevant probability information is known

S. Makrogiannis (DSU)

Bayesian Decision Theory

September 3, 2015 3 / 20

4 E 5 4 E 5

Main concept

- The decisions minimize a total expected "risk"
- Straightforward risk is the classification error
- The risk is associated with the cost of making a specific decision

S. Makrogiannis (DSU)

Bayesian Decision Theory

September 3, 2015 4 / 20

イロト イポト イヨト イヨト

Terminology

Consider the problem: an observer watches fish arrive on a conveyor belt. There are two kinds of fish: salmon and sea bass. The objective is to predict the type of fish that will appear next, assuming that the sequence is random

- State of nature ω is a random variable that describes the class, that is ω_1 for salmon and ω_2 for sea bass
- Prior probabilities P(ω_j): probabilities for states of nature (categories)
- Feature vector x: vector consisting of measurements for patterns, e.g. fish lightness



Terminology

Consider the problem: an observer watches fish arrive on a conveyor belt. There are two kinds of fish: salmon and sea bass. The objective is to predict the type of fish that will appear next, assuming that the sequence is random

- Probability density function p(x): How frequently a pattern with value x occurs
- Class-conditional probability density $p(x|\omega_j)$: Frequency of occurrence of a feature x for a specific class label ω_j . Known as likelihood. Probability of a lightness value for a specific fish type
- Posterior probability density p(ω_j|x): Probability of class ω_j, given a measurement x. Probability of occurrence of a specific fish type, given the lightness value

S. Makrogiannis (DSU)

Bayesian Decision Theory

< □ ▶ < □ ▶ < 三 ▶ < 三 ▶ < 三 ▶ ○ Q (* September 3, 2015 6 / 20

Class-conditional Densities

A key process in pattern classification is to estimate the class-conditional density for each category, also referred to as likelihood



S. Makrogiannis (DSU)

Bayesian Decision Theory

Decision Based on Priors Only

Decide ω_1 if $P(\omega_1) > P(\omega_2)$, otherwise decide ω_2

- This rule always makes the same decision
- It favors the most likely class
- $P(error) = min[P(\omega_1), P(\omega_2)]$
- It may be applied if we know only the prior probabilities $P(\omega_j)$, where ω_j is the state of nature

S. Makrogiannis (DSU)

Bayesian Decision Theory

September 3, 2015 8 / 20

A E > A E >

Bayes formula

- Suppose that in addition to $P(\omega_j)$, we have obtained measurements x and we know the likelihoods $p(x|\omega_j)$
- Then the Bayes rule yields

$$P(\omega_j|x) = rac{p(x|\omega_j)P(\omega_j)}{p(x)}, ext{ where } p(x) = \sum_{j=1}^2 p(x|\omega_j)P(\omega_j)$$

• This formula be expressed as $posterior = \frac{likelihood \times prior}{evidence}$

S. Makrogiannis (DSU)

Bayesian Decision Theory

September 3, 2015 9 / 20

A E F A E F

< □ > < 同 >

Bayes Decision Rule

Decide ω_1 if $P(\omega_1|x) > P(\omega_2|x)$; otherwise decide ω_2

Decide ω_1 if $p(x|\omega_1)P(\omega_1) > p(x|\omega_2)P(\omega_2)$; otherwise decide ω_2



Figure: Posterior probabilities corresponding to class-conditionals of previous figure for $P(\omega_1) = 2/3$ and $P(\omega_2) = 1/3$

S. Makrogiannis (DSU)

Bayesian Decision Theory
Probability of Error

 The probability of error when we make a decision using the Bayesian rule is

$$P(error|x) = \begin{cases} P(\omega_1|x), & \text{if we decide } \omega_2 \\ P(\omega_2|x), & \text{if we decide } \omega_1 \end{cases}$$

- Hence, $P(error|x) = min[P(\omega_1|x), P(\omega_2|x)]$
- This rule minimizes the average probability for error because

$$P(error) = \int_{-\infty}^{\infty} P(error, x) dx = \int_{-\infty}^{\infty} P(error|x) p(x) dx$$

S. Makrogiannis (DSU)

Bayesian Decision Theory

September 3, 2015 11 / 20

Bayesian Decision Theory Generalization

- Let $\Omega = \{\omega_1, \omega_2, ..., \omega_c\}$ be the set of *c* states of nature or classes, and $A = \{a_1, a_2, ..., a_c\}$ be the set of actions. Then $\lambda(\alpha_i | \omega_j)$ is the loss caused by taking action α_i for a true class ω_j .
- Let $\mathbf{x} \in \mathbb{R}^D$ be a D-dimensional feature vector corresponding to measurements
- Then the Bayes formula is defined as

$$P(\omega_j | \mathbf{x}) = rac{p(\mathbf{x} | \omega_j) P(\omega_j)}{p(\mathbf{x})}, ext{ where } p(\mathbf{x}) = \sum_{j=1}^c p(\mathbf{x} | \omega_j) P(\omega_j)$$

S. Makrogiannis (DSU)

Bayesian Decision Theory

September 3, 2015 12 / 20

イロト イポト イヨト イヨト 二日

Loss Function

- We assume that after an observation **x** we take action α_i , when the true state of nature is ω_j . Then the incurred loss for this action is $\lambda(\alpha_i | \omega_j)$
- We define as conditional risk $R(\alpha_i | \mathbf{x})$ the expected loss given by

$$R(\alpha_i|\mathbf{x}) = \sum \lambda(\alpha_i|\omega_j) P(\omega_j|\mathbf{x})$$

S. Makrogiannis (DSU)

Bayesian Decision Theory

September 3, 2015 13 / 20

イロト 不得下 イヨト イヨト

Bayes Risk

- General decision rule is a function $lpha({\sf x})$ with $lpha:\mathbb{R}^D o A$
- The overall risk *R* is the expected loss incurred by a specific decision estimated by

$$R = \int R(\alpha(\mathbf{x})|\mathbf{x})p(\mathbf{x})d\mathbf{x}$$

- Hence, to minimize overall risk, we compute the conditional risk $R(\alpha_i|\mathbf{x})$ for $i = 1, ..., \alpha$ and choose the action α_i that minimizes $R(\alpha_i|\mathbf{x})$
- The corresponding overall risk *R** is called Bayes risk and denotes the optimal classification performance

S. Makrogiannis (DSU)

Bayesian Decision Theory

September 3, 2015 14 / 20

< ロト < 同ト < ヨト < ヨト

Two-category Classification

Assumptions

• Classification problem with two classes ω_1 and ω_2 .

•
$$lpha_i$$
: decide ω_i , for $i=1,2$

• $\lambda_{ij} = \lambda(\alpha_i | \omega_j)$

Bayes Decision Rule

•
$$R(\alpha_1|\mathbf{x}) = \lambda_{11}P(\omega_1|\mathbf{x}) + \lambda_{12}P(\omega_2|\mathbf{x})$$

•
$$R(\alpha_2|\mathbf{x}) = \lambda_{21}P(\omega_1|\mathbf{x}) + \lambda_{22}P(\omega_2|\mathbf{x})$$

• Rule: Decide ω_1 if $R(\alpha_1|\mathbf{x}) < R(\alpha_2|\mathbf{x})$; otherwise decide ω_2

ullet Decide ω_1 if

$$\lambda_{11} P(\boldsymbol{\omega}_1 | \mathbf{x}) + \lambda_{12} P(\boldsymbol{\omega}_2 | \mathbf{x}) < \lambda_{21} P(\boldsymbol{\omega}_1 | \mathbf{x}) + \lambda_{22} P(\boldsymbol{\omega}_2 | \mathbf{x})$$

$$(\lambda_{21} - \lambda_{11})P(\omega_1 | \mathbf{x}) > (\lambda_{12} - \lambda_{22})P(\omega_2 | \mathbf{x})$$

Likelihood Ratio Test (LRT) for Two Categories

Bayes Decision Rule

ullet Decide ω_1 if

$$\lambda_{11}P(\omega_1|\mathbf{x}) + \lambda_{12}P(\omega_2|\mathbf{x}) < \lambda_{21}P(\omega_1|\mathbf{x}) + \lambda_{22}P(\omega_2|\mathbf{x}) \Leftrightarrow$$

$$(\lambda_{21} - \lambda_{11}) P(\omega_1 | \mathbf{x}) > (\lambda_{12} - \lambda_{22}) P(\omega_2 | \mathbf{x}) \Leftrightarrow$$

$$(\lambda_{21} - \lambda_{11}) p(\mathbf{x}|\omega_1) P(\omega_1) > (\lambda_{12} - \lambda_{22}) p(\mathbf{x}|\omega_2) P(\omega_2) \Leftrightarrow$$

$$\frac{p(\mathsf{x}|\omega_1)}{p(\mathsf{x}|\omega_2)} > \frac{(\lambda_{12} - \lambda_{22})}{(\lambda_{21} - \lambda_{11})} \frac{P(\omega_2)}{P(\omega_1)}$$

S. Makrogiannis (DSU)

Bayesian Decision Theory

September 3, 2015 16 / 20

3

イロト イポト イヨト イヨト

Likelihood Ratio Example



Figure: Bayesian decision theory components: class-conditionals, posterior probabilities, and likelihood ratio (left to right)

S. Makrogiannis (DSU)

Bayesian Decision Theory

September 3, 2015 17 / 20

Minimum Error-rate Classification

Assumptions

• Classification problem with c classes in $\Omega = \{ \omega_1, \omega_2, ..., \omega_c \}$

•
$$lpha_i$$
: decide ω_i , for $i=1,..,c$

•
$$\lambda_{ij} = \lambda(lpha_i | \omega_j)$$
, where $\lambda(lpha_i | \omega_j) = \begin{cases} 0, \text{ if } i = j \\ 1, \text{ if } i \neq j \end{cases}$ for $i, j = 1, ..., c$

Conditional Risk and Bayes Decision Rule

•
$$R(\alpha_i | \mathbf{x}) = \sum_{j=1}^{c} \lambda_{ij} P(\omega_j | \mathbf{x}) = \sum_{j \neq i}^{c} P(\omega_j | \mathbf{x}) = 1 - P(\omega_i | \mathbf{x})$$

Rule: Decide ω_i if

$$R(\alpha_i | \mathbf{x}) < R(\alpha_j | \mathbf{x}), \forall j \neq i \Leftrightarrow$$

$$1 - P(\omega_i | \mathbf{x}) < 1 - P(\omega_j | \mathbf{x}), \forall j \neq i \Leftrightarrow$$

 $P(\omega_i | \mathbf{x}) > P(\omega_j | \mathbf{x}), \forall j \neq i$

Maximum A Posteriori (MAP) Criterion

 According to the previous result the decision rule for Minimum Error-rate Classification is:

Decide
$$\omega_i, ext{ if } P(\omega_i | \mathbf{x}) > P(\omega_j | \mathbf{x}), orall j
eq i \Leftrightarrow$$

Decide
$$\omega_i$$
, if $\frac{P(\omega_i | \mathbf{x})}{P(\omega_j | \mathbf{x})} > 1, \forall j \neq i$

• This is also known as Maximum A Posteriori (MAP) criterion as it seeks to maximize the posterior probability $P(\omega_i|\mathbf{x})$

S. Makrogiannis (DSU)

Bayesian Decision Theory

September 3, 2015 19 / 20

Maximum Likelihood (ML) Criterion

Γ

• If in addition to a binary loss function we have equal priors, that is, $P(\omega_i) = P(\omega_j), \forall i \neq j$, then the decision rule becomes

Decide
$$\omega_i$$
, if $\frac{P(\omega_i | \mathbf{x})}{P(\omega_j | \mathbf{x})} > 1, \forall j \neq i \Leftrightarrow$

Decide
$$\omega_i$$
, if $rac{p(\mathbf{x}|\omega_i)P(\omega_i)}{p(\mathbf{x}|\omega_j)P(\omega_j)} > 1, \forall j \neq i \Leftrightarrow$
Decide ω_i , if $rac{p(\mathbf{x}|\omega_i)}{p(\mathbf{x}|\omega_j)} > 1, \forall j \neq i$

 This is also known as Maximum Likelihood (ML) criterion as it seeks to maximize the likelihood p(x|wi)

S. Makrogiannis (DSU)

Bayesian Decision Theory

September 3, 2015 20 / 20

Bayesian Decision Theory MTSC 852, Fall 2015

Sokratis Makrogiannis, Ph.D.

Department of Mathematical Sciences Delaware State University smakrogiannis@desu.edu

October 6, 2015

S. Makrogiannis (DSU)

Bayesian Decision Theory

October 6, 2015 1 / 25

▲ロト ▲掃ト ▲ヨト ▲ヨト ニヨー のへで



S. Makrogiannis (DSU)

Bayesian Decision Theory

October 6, 2015 2 / 25

▲ロト ▲園 ト ▲ 臣 ト ▲ 臣 ト 二 臣 … の Q ()

- Classifiers, Discriminant Functions, Decision Surfaces
- 2 The Normal Density

Bayesian Decision Theory

October 6, 2015 2 / 25

▲ロト ▲撮 ト ▲ 臣 ト ▲ 臣 ト 二 臣 … の � @

- Classifiers, Discriminant Functions, Decision Surfaces
- 2 The Normal Density
- ③ Discriminant Functions for the Normal Density

- Classifiers, Discriminant Functions, Decision Surfaces
- 2 The Normal Density
- 3 Discriminant Functions for the Normal Density
- 4 Error Probabilities and Integrals

- Classifiers, Discriminant Functions, Decision Surfaces
- 2 The Normal Density
- 3 Discriminant Functions for the Normal Density
- ④ Error Probabilities and Integrals
- 5 Receiver Operating Characteristics

S. Makrogiannis (DSU)

Bayesian Decision Theory

October 6, 2015 2 / 25

イロト (母) (ヨ) (ヨ) (ヨ) () ()

- Classifiers, Discriminant Functions, Decision Surfaces
- 2 The Normal Density
- 3 Discriminant Functions for the Normal Density
- ④ Error Probabilities and Integrals
- 5 Receiver Operating Characteristics
- 6 Bayesian Decision Theory for Discrete Features

S. Makrogiannis (DSU)

Bayesian Decision Theory

October 6, 2015 2 / 25

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 ◇◇◇

Discriminant Functions

- A classifier can be represented by a set of discriminant functions $g_i(x), i = 1, ..., c$ corresponding to each state of nature or category.
- Then the decision rule becomes:

Assign feature vector **x** to class ω_i , if $g_i(\mathbf{x}) > g_j(\mathbf{x})$ $\forall j \neq i$

• Then the classifier can be considered to be a network or machine that computes c discriminant functions and makes a decision using a maximum operator.



Bayes Classifier with Multiple Categories

- A Bayes classifier can also be described using discriminant functions.
- For the general case that minimizes conditional risks $g_i(\mathbf{x}) = -R(lpha_i|\mathbf{x})$
- For the MAP -or minimum-error-rate- criterion g_i(x) = P(ω_i|x) With some more operations we can produce other equivalent MAP discriminant functions

$$g_{i}(\mathbf{x}) = P(\omega_{i}|\mathbf{x}) = \frac{p(\mathbf{x}|\omega_{i})P(\omega_{i})}{\sum_{j=1}^{c} p(\mathbf{x}|\omega_{j})P(\omega_{j})}$$
$$g_{i}(\mathbf{x}) = p(\mathbf{x}|\omega_{i})P(\omega_{i})$$
$$g_{i}(\mathbf{x}) = \ln p(\mathbf{x}|\omega_{i}) + \ln P(\omega_{i})$$

• For the ML criterion $g_i(\mathbf{x}) = p(\mathbf{x}|\omega_i)$

S. Makrogiannis (DSU)

Bayesian Decision Theory

October 6, 2015 4 / 25

イロト イポト イヨト イ

Two Categories

- Suppose a problem with two categories ω_1 and ω_2 .
- Then we can define a single discriminant function by

$$g(\mathbf{x}) = g_1(\mathbf{x}) - g_2(\mathbf{x})$$

• The decision rule is:

Decide ω_1 if $g(\mathbf{x}) > 0$; otherwise decide ω_2

• For the MAP -or minimum-error-rate- criterion it follows that

$$g(\mathbf{x}) = P(\omega_1 | \mathbf{x}) - P(\omega_2 | \mathbf{x}) \Leftrightarrow$$
$$g(\mathbf{x}) = \ln p(\mathbf{x} | \omega_1) + \ln P(\omega_1) - \ln p(\mathbf{x} | \omega_2) - \ln P(\omega_2) \Leftrightarrow$$
$$g(\mathbf{x}) = \ln \frac{p(\mathbf{x} | \omega_1)}{p(\mathbf{x} | \omega_2)} + \ln \frac{P(\omega_1)}{P(\omega_2)}$$

S. Makrogiannis (DSU)

Bayesian Decision Theory

Decision Regions and Boundaries

• Decision rules divide the feature space in Decision Regions $\mathscr{R}_i, i = 1, 2, ..., c$, where c is the number of categories



S. Makrogiannis (DSU)

Bayesian Decision Theory

October 6, 2015 6 / 25

3

200

Decision Regions and Boundaries

- Decision rules divide the feature space in Decision Regions $\mathscr{R}_i, i = 1, 2, ..., c$, where c is the number of categories
- The Decision Regions $(\mathcal{R}_i, \mathcal{R}_j)$ are separated by Decision Boundaries on which $g_i(\mathbf{x}) = g_j(\mathbf{x})$.



S. Makrogiannis (DSU)

Bayesian Decision Theory

October 6, 2015 6 / 25

Introduction

- As explained in previous sections, we can design a Bayesian classifier if we know the likelihood p(x|ω_i) and priors P(ω_i) for each category i
- Among all density functions the multivariate Gaussian model is very popular
- Normal density models are popular mainly because of the Central Limit Theorem, and the fact that the normal density is analytically tractable

S. Makrogiannis (DSU)

Bayesian Decision Theory

October 6, 2015 7 / 25

イロト イポト イヨト イヨト

Univariate Normal Density

- Density function also symbolized by $N(\mu, \sigma)$: $p(x) = \frac{1}{\sqrt{2\pi\sigma}}e^{-\frac{(x-\mu)^2}{2\sigma^2}}$
- Expected value: $\mu \equiv \mathscr{E}[x] = \int_{-\infty}^{\infty} x p(x) dx$
- Variance: $\sigma \equiv \mathscr{E}[(x-\mu)^2] = \int_{-\infty}^{\infty} (x-\mu)^2 p(x) dx$
- Entropy: $H(p(x)) = -\int_{-\infty}^{\infty} p(x) \ln p(x) dx$

S. Makrogiannis (DSU)

Bayesian Decision Theory

October 6, 2015 8 / 25

イロト 不得下 イヨト イヨト 二子

Multivariate Normal Density

- Density function (also symbolized by $N(\mu, \Sigma)$) : $p(\mathbf{x}) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} e^{\frac{1}{2} (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)}$
- Expected value: $\mu \equiv \mathscr{E}[\mathbf{x}] = \int_{-\infty}^{\infty} \mathbf{x} p(\mathbf{x}) d\mathbf{x}$
- Covariance matrix: $\Sigma \equiv \mathscr{E}[(\mathbf{x} - \mu)(\mathbf{x} - \mu)^T] = \int_{-\infty}^{\infty} (\mathbf{x} - \mu)(\mathbf{x} - \mu)^T p(\mathbf{x}) d\mathbf{x}$

S. Makrogiannis (DSU)

Bayesian Decision Theory

October 6, 2015 9 / 25

Introduction

• According to previous sections the use of MAP criterion yields the following discriminant functions

$$g_i(\mathbf{x}) = \ln p(\mathbf{x}|\omega_i) + \ln P(\omega_i)$$

• For the case of multivariate normal densities for the likelihood, i.e. when $p(\mathbf{x}|\omega_i) = N(\mu, \Sigma)$, it follows that

 $g_i(\mathbf{x}) = -(1/2)(\mathbf{x} - \mu_i)^T \Sigma_i^{-1}(\mathbf{x} - \mu_i) - (d/2) \ln 2\pi - (1/2) \ln |\Sigma_i| + \ln P(\omega_i)$

• Next, we examine three different Covariance Matrix models

Identical Diagonal Covariance Matrices with $\Sigma_i = \sigma^2 \mathbf{I}$

- We assume that features are statistically independent and have equal standard deviations
- The discriminant function can be simplified as

$$g_i(\mathbf{x}) = -\frac{(\mathbf{x} - \mu_i)^T (\mathbf{x} - \mu_i)}{2\sigma^2} + \ln P(\omega_i) \Leftrightarrow$$

$$g_i(\mathbf{x}) = -\frac{1}{2\sigma^2} [\mathbf{x}^T \mathbf{x} - 2\mu_i^T \mathbf{x} + \mu_i^T \mu_i] + \ln P(\omega_i)$$

 Further simplification -by discarding category-invariant terms - yields a linear discriminant function

$$g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + b_{i0}$$

where

$$\mathbf{w}_i = (1/\sigma^2)\mu_i, b_{i0} = -(1/2\sigma^2)\mu_i^{\mathsf{T}}\mu_i + \ln P(\omega_i)$$

S. Makrogiannis (DSU)

Bayesian Decision Theory

Identical Diagonal Covariance Matrices with $\Sigma_i = \sigma^2 \mathbf{I}$

- A classifier defined by linear discriminant functions is called a linear machine.
- The decision surfaces defined by $g_i(\mathbf{x}) g_j(\mathbf{x})$ are hyperplanes in the feature space
- Here the decision surface is described by the equation

$$\mathbf{w}^{\mathsf{T}}(\mathbf{x} - \mathbf{x}_0) = 0$$

with

$$w = \mu_i - \mu_j$$

$$x_0 = (1/2)(\mu_i + \mu_j) - \frac{\sigma^2}{\|\mu_i - \mu_j\|^2} \ln \frac{P(\omega_i)}{P(\omega_j)} (\mu_i - \mu_j)$$

$$\|\mu_i - \mu_j\|^2 = (\mu_i - \mu_j)^T (\mu_i - \mu_j)$$

• Therefore, the decision surface is a hyperplane passing through x_0 and orthogonal to the line linking the category means $(x - x_0)$

Bayesian Decision Theory

Identical Diagonal Covariance Matrices with $\Sigma_i = \sigma^2 \mathbf{I}$

Additional notes

• If all priors $P(\omega_i)$ are equal, then the discriminant function can be simplified to

$$g_i(\mathbf{x}) = \|\mathbf{x} - \mu_i\|$$

- Hence, each pattern will be classified to the category with the closest mean using a Euclidean norm
- This is called a minimum-distance classifier



S. Makrogiannis (DSU)

Bayesian Decision Theory

October 6, 2015 13 / 25

Identical Arbitrary Covariance Matrices with $\Sigma_i = \Sigma$

 The discriminant function can be simplified by discarding category-invariant terms

$$g_i(\mathbf{x}) = -(1/2)(\mathbf{x} - \mu_i)^T \Sigma^{-1}(\mathbf{x} - \mu_i) + \ln P(\omega_i)$$

• Let the prior probabilities $P(\omega_i)$ be equal for all classes. Then

$$g_i(\mathbf{x}) = -(1/2)(\mathbf{x}-\mu_i)^T \Sigma^{-1}(\mathbf{x}-\mu_i)$$

After expansion of Mahalanobis distance and simplification

$$g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + b_{i0}$$

where

$$\mathbf{w}_i = \Sigma^{-1} \mu_i, b_{i0} = -(1/2) \mu_i^T \Sigma^{-1} \mu_i + \ln P(\omega_i)$$

S. Makrogiannis (DSU)

Bayesian Decision Theory

October 6, 2015 14 / 25

Э

イロト イポト イヨト イヨト

Identical Arbitrary Covariance Matrices with $\Sigma_i = \Sigma$

- The decision surfaces defined by $g_i(\mathbf{x}) g_j(\mathbf{x})$ are hyperplanes in the feature space
- Here the decision surface is described by the equation

$$\mathbf{w}^{\mathsf{T}}(\mathbf{x} - \mathbf{x}_0) = \mathbf{0}$$

with

$$\mathbf{w} = \Sigma^{-1}(\mu_i - \mu_j)$$
$$\mathbf{x}_0 = (1/2)(\mu_i + \mu_j) - \frac{\ln \frac{P(\omega_i)}{P(\omega_j)}}{(\mu_i - \mu_j)^T \Sigma^{-1}(\mu_i - \mu_j)}(\mu_i - \mu_j)$$

• Therefore, the decision surface is a hyperplane passing through x_0 but it is not necessarily orthogonal to the line linking the category means $(x - x_0)$

S. Makrogiannis (DSU)

Bayesian Decision Theory

October 6, 2015 15 / 25

< 市計

Discriminant Functions for the Normal Density

Identical Arbitrary Covariance Matrices with $\Sigma_i = \Sigma$



S. Makrogiannis (DSU)

Bayesian Decision Theory

<ロ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Arbitrary Covariance Matrices

- In this case the covariance matrices are different for each category
- The discriminant function takes the form

$$g_i(\mathbf{x}) = \mathbf{x}^T \mathbf{W}_i \mathbf{x} + \mathbf{w}_i^T \mathbf{x} + b_{i0}$$

where

$$\mathbf{W}_{i} = -(1/2)\boldsymbol{\Sigma}_{i}^{-1}$$
$$\mathbf{w}_{i} = \boldsymbol{\Sigma}_{i}^{-1}\boldsymbol{\mu}_{i}$$
$$b_{i0} = -(1/2)\boldsymbol{\mu}_{i}^{T}\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}_{i} - (1/2)\ln|\boldsymbol{\Sigma}_{i}| + \ln P(\boldsymbol{\omega}_{i})$$

- This is a quadratic form
- In the two-category case the decision surfaces are hyperquadrics assuming any of the following forms: hyperplanes, hyperspheres, hyperellipsoids, hyperparaboloids, or hyperhyperboloids

S. Makrogiannis (DSU)

Bayesian Decision Theory

October 6, 2015 17 / 25

Discriminant Functions for the Normal Density

Arbitrary Covariance Matrices



S. Makrogiannis (DSU)

Bayesian Decision Theory

October 6, 2015 18 / 25

Bayes Error Rate



S. Makrogiannis (DSU)

Bayesian Decision Theory

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Problem Definition

- Suppose that we want to detect a signal. We measure a feature x-let's say voltage- that has a mean μ_1 when a signal is present, and μ_2 when there is no signal
- We assume two normal distributions with different means but the same variance, that is $N(\mu_1, \sigma)$ and $N(\mu_2, \sigma)$ for classes ω_1 and ω_2 respectively
- For the detection we use a threshold value x* -that is unknown to usto classify a feature x into signal or no-signal categories
- We seek a measure of *separability*, i.e., how easy it is to separate the two categories and make a prediction
- Suppose we do not know $\mu_1, \mu_2, \sigma, x*$, but we know the sate of nature and the system's decision for each

S. Makrogiannis (DSU)

Bayesian Decision Theory

October 6, 2015 20 / 25

3

< ロト < 同ト < ヨト < ヨト
Probabilities

- $P(x>x*|x\in \omega_2)$: hit, we detect a signal and the signal is present
- $P(x > x * | x \in \omega_1)$: false alarm, we detect a signal but no signal is present
- $P(x < x * | x \in \omega_2)$: miss, we do not detect a signal but the signal is present
- $P(x < x * | x \in \omega_1)$: correct rejection, we do not detect a signal and no signal is present
- If we have enough samples we can compute the probabilities experimentally
- If x* changes, then the above probabilities will change



Receiver Operating Characteristics (ROC)

- An ROC graph is a plot of probability of hit vs. probability of false alarm for different values of x*
- The curve can be used to measure separability, that is the capacity of our system to detect the signal
- A usual measure of separability is the area under the curve, frequently denoted by AUC ROC
- In this context it is important to choose a measure x whose variations will significantly change the probability values

S. Makrogiannis (DSU)

Bayesian Decision Theory

October 6, 2015 22 / 25

イロト イポト イヨト イヨト 二日

ROC Example

ROC curve for varying discriminability



S. Makrogiannis (DSU)

Bayesian Decision Theory

<□ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

ROC Example 2

ROC curves are not necessarily symmetric



S. Makrogiannis (DSU)

Bayesian Decision Theory

<□ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Bayesian Decision Theory for Discrete Features

- Here we assume that the components of feature vector x can assume only m discrete values v₁, v₂,..., v_m
- The Bayes formula involves probabilities only

$$P(\omega_j|\mathbf{x}) = rac{P(\mathbf{x}|\omega_j)P(\omega_j)}{P(\mathbf{x})}$$

with

$$P(x) = \sum_{i=1}^{c} P(\mathbf{x}|\omega_i) P(\omega_i)$$

- Integrals are replaced by sums in definitions of expectation, variance, entropy and other statistical measures
- Fundamental Bayes decision rule is the same To minimize the overall risk, select the action α_i for which the conditional risk R(α_i|x) is minimized

$$\alpha * = \arg\min_{i} R(\alpha_{i} | \mathbf{x})$$

S. Makrogiannis (DSU)

Bayesian Decision Theory

October 6, 2015 25 / 25

MTSC 852 - Pattern Recognition Lab Session Parametric Estimation

Sokratis Makrogiannis, Ph.D.

October 17, 2015

Contents

1	Bay	esian Parameter Estimation for the Gaussian Density	1
	1.1	Univariate Normal Density	1
	1.2	Estimate $p(\mu \mathcal{D})$	2
	1.3	Estimate $p(x \mathcal{D})$	3
2	\mathbf{Fish}	er Linear Discriminant	7
	2.1	Discriminant Analysis	7
	2.2	Problem Definition	0
	2.3	Class Separability	0
	2.4	Criterion Function	0
	2.5	Scatter Matrices	.1
	2.6	Optimizing the Criterion Function	2
	2.7	Classification Rule	.3

1 Bayesian Parameter Estimation for the Gaussian Density

1.1 Univariate Normal Density

Find the class-conditional density $p(x|\mathcal{D})$ using Bayesian estimation assuming that $p(x|\mu) \sim N(\mu, \sigma^2)$, $p(\mu) \sim N(\mu_0, \sigma_0)$ and σ is known.

- Estimate $p(\mu|\mathcal{D})$ using Bayes rule
- Estimate $p(x|\mathcal{D})$ by integration over the parameter space

Find the class-conditional density $p(x|\mathcal{D})$ using Bayesian estimation assuming that $p(x|\mu) \sim N(\mu, \sigma^2)$, $p(\mu) \sim N(\mu_0, \sigma_0)$ and σ is known.

1.2 Estimate $p(\mu | \mathcal{D})$

- According to previous analysis, we seek to estimate $p(\mathbf{x}|\mathcal{D})$ for each class
- This is achieved by integration over the parameter space

$$p(\mathbf{x}|\mathcal{D}) = \int p(\mathbf{x}, \theta|\mathcal{D}) d\theta$$

- From definition of joint probability: $p(\mathbf{x}|\mathcal{D}) = \int p(\mathbf{x}|\theta, \mathcal{D}) p(\theta|\mathcal{D}) d\theta$
- We use Bayes rule to estimate $p(\theta|\mathcal{D})$: $p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)P(\theta)}{\int p(\mathcal{D}|\theta)P(\theta)d\theta}$
- If the samples are independently drawn, then: $p(\mathcal{D}|\theta) = \prod_{i=1}^{n} p(\mathbf{x}_i|\theta)$
- We use Bayes rule to estimate posterior parameter density $p(\mu|\mathcal{D})$:

$$p(\mu|\mathcal{D}) = \frac{p(\mathcal{D}|\mu)p(\mu)}{\int p(\mathcal{D}|\mu)p(\mu)d\mu}$$

- Let samples $\mathcal{D} = \{\mathbf{x_1}, \mathbf{x_2}, ..., \mathbf{x_n}\}$ be independently drawn. Then: $p(\mathcal{D}|\mu) = \prod_{i=1}^n p(x_i|\mu)$
- Then: $p(\mu|\mathcal{D}) = \alpha \cdot \prod_{i=1}^{n} p(x_i|\mu) p(\mu)$
- According to assumptions:

$$p(x_i|\mu) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}}, \quad p(\mu) = \frac{1}{\sqrt{2\pi\sigma_0}} e^{-\frac{(\mu-\mu_0)^2}{2\sigma_0^2}}$$

So: $p(\mu|\mathcal{D}) = \alpha \cdot \prod_{i=1}^n \left[\frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}} \frac{1}{\sqrt{2\pi\sigma_0}} e^{-\frac{(\mu-\mu_0)^2}{2\sigma_0^2}} \right]$

- After some more manipulations we can show that $p(\mu|\mathcal{D})$ is an exponential function of a quadratic function, hence it has the form $N(\mu_n, \sigma_n)$
- Therefore

$$p(\mu|\mathcal{D}) = \frac{1}{\sqrt{2\pi\sigma_n}} e^{-\frac{(\mu-\mu_n)^2}{2\sigma_n^2}}$$

where,

$$\mu_n = \left(\frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2}\right)\hat{\mu}_n + \frac{\sigma^2}{n\sigma_0^2 + \sigma^2}\mu_0, \quad \sigma_n^2 = \frac{\sigma_0^2\sigma^2}{n\sigma_0^2 + \sigma^2}$$
$$\hat{\mu}_n = (1/n)\sum_{i=1}^n x_i$$

- σ_n decreases as $n \to \infty$ with $\lim_{n \to \infty} \sigma_n^2 = \frac{\sigma^2}{n}$
- We observe that as the number of training samples increases, $p(\mu|\mathcal{D})$ becomes sharper around μ_n . This process is called *Bayesian learning*
- If $\sigma_0 \neq 0$, then μ_n approaches the sample mean $\lim_{n\to\infty} \mu_n = \hat{\mu}_n$.



1.3 Estimate $p(x|\mathcal{D})$

• According to Bayesian estimation process

$$p(x|\mathcal{D}) = \int p(x|\mu, \mathcal{D}) p(\mu|\mathcal{D}) d\mu \Leftrightarrow$$
$$p(x|\mathcal{D}) = \int \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \frac{1}{\sqrt{2\pi}\sigma_n} e^{-\frac{(\mu-\mu_n)^2}{2\sigma_n^2}} d\mu \Leftrightarrow$$

$$p(x|\mathcal{D}) = \frac{1}{2\pi\sigma\sigma_n} f(\sigma, \sigma_n) e^{-\frac{(x-\mu_n)^2}{2(\sigma^2+\sigma_n^2)}}$$

where $f(\sigma, \sigma_n)$ has an integral form:

$$f(\sigma, \sigma_n) = \int \exp\left[(-1/2)\frac{\sigma^2 + \sigma_n^2}{\sigma^2 \sigma_n^2} \left(\mu - \frac{\sigma_n^2 x + \sigma^2 \mu_n}{\sigma^2 + \sigma_n^2}\right)^2\right] d\mu$$

- Observe that $p(x|\mathcal{D}) \sim N(\mu_n, \sigma^2 + \sigma_n^2)$
- The above result gives the class-conditional density $p(x|\omega_i, \mathcal{D}_i)$ based on the posterior parameter mean estimate μ_n and the posterior parameter variance estimate increased by the uncertainty in x that we assume to be known

Exercise 1. Consider Bayesian estimation of the mean of a one-dimensional Gaussian. Suppose you are given the prior for the mean is $p(\mu) \sim N(\mu_0, \sigma_0)$.

- 1. Write a program that plots the density $p(x|\mathcal{D})$ given, μ_0, σ_0, σ and training set $\mathcal{D} = \{x_1, x_2, \ldots, x_n\}$.
- 2. Estimate σ for the x_2 component of ω_3 in Table 1 and in file ch3_dhs_samples.dat. Now assume $\mu_0 = -1$ and plot your estimated densities $p(x|\mathcal{D})$ for each of the following values of the dogmatism $\sigma^2/\sigma_0^2 : 0.1, 1, 10, 100$.
- 3. Repeat above process but this time generate a dense sample set with the same mean and standard deviation as in the real dataset.

	ω_1			ω_2			ω_3		
point	x_1	x_2	x_3	x_1	x_2	x_3	x_1	x_2	x_3
1	0.42	-0.087	0.58	-0.4	0.58	0.089	0.83	1.6	-0.014
2	-0.2	-3.3	-3.4	-0.31	0.27	-0.04	1.1	1.6	0.48
3	1.3	-0.32	1.7	0.38	0.055	-0.035	-0.44	-0.41	0.32
4	0.39	0.71	0.23	-0.15	0.53	0.011	0.047	-0.45	1.4
5	-1.6	-5.3	-0.15	-0.35	0.47	0.034	0.28	0.35	3.1
6	-0.029	0.89	-4.7	0.17	0.69	0.1	-0.39	-0.48	0.11
7	-0.23	1.9	2.2	-0.011	0.55	-0.18	0.34	-0.079	0.14
8	0.27	-0.3	-0.87	-0.27	0.61	0.12	-0.3	-0.22	2.2
9	-1.9	0.76	-2.1	-0.065	0.49	0.0012	1.1	1.2	-0.46
10	0.87	-1.0	-2.6	-0.12	0.054	-0.063	0.18	-0.11	-0.49

Table 1: Three-dimensional data sampled from three categories.

```
function x_density_given_d = ch3_bayesian_estimation_1d(X, sigma,
    mu_0, sigma_0)
% Bayesian parameter estimation for a univariate normal
    distribution.
% S. Makrogiannis, Delaware State Univ, 10/2015.
% Initial parameters and calculations.
```

```
X = sort(X);
n = numel(X);
```

 $hat_mu_n = sum(X) / n;$

```
normal_density = @(x, mu, sigma) ( (1/(sqrt(2*pi)*sigma)) * exp(
   (-0.5) * ( (x - mu) / sigma )^2 ) );
% For a range of values of our random variable x:
for i=1:n
   x_density_given_d(i) = 0;
   for mu = mu_0-4*sigma_0:mu_0+4*sigma_0
      % Estimate p(mu|D) using Bayesian technique.
       [mu_density_given_d(i), mu_n, sigma_n] = ...
          bayesian_parameter_density(mu, sigma, mu_0, sigma_0,
             hat_mu_n, n);
      % Compute p(x|mu)
      x_density_given_mu(i) = normal_density(X(i), mu, sigma);
      % Compute p(x|mu) * p(mu|D)
      % Add up to approximate integral.
      x_density_given_d(i) = x_density_given_d(i) +
          (x_density_given_mu(i) * mu_density_given_d(i));
   end
end
figure, plot(X, x_density_given_d); title('p(X|D)')
end
X_____X
function [mu_density, mu_n, sigma_n] = ...
   bayesian_parameter_density(mu, sigma, mu_0, sigma_0, hat_mu_n,
      n)
% Compute mu_n and sigma_n
normal_density = @(x, mu, sigma) ( (1/(sqrt(2*pi)*sigma)) * exp(
   (-0.5) * ( (x - mu) / sigma )^2 ) );
```

```
mu_n = ( n * sigma_0^2 / ( n * sigma_0^2 + sigma^2 ) ) * hat_mu_n
   + ...
   ( sigma^2 / ( n * sigma_0^2 + sigma^2 ) ) * mu_0;
var_n = (sigma_0^2 * sigma^2) / (( n * sigma_0^2 + sigma^2 ));
sigma_n = sqrt(var_n);
mu_density = normal_density(mu, mu_n, sigma_n);
end
% Bayesian estimation for 1-D Gaussian distributions.
% Load data.
A = load('ch3_dhs_samples.dat');
% Initialize parameters and compute sigma.
dogmatism = [0.1, 1, 10, 100];
n_runs = numel(dogmatism);
sigma = std(A(:,8));
mu_0 = -1;
% Perform density estimation.
for i=1:n_runs
sigma_0 = sqrt(sigma^2/dogmatism(i));
x_density_given_d = ch3_bayesian_estimation_1d(A(:,8), sigma,
   mu_0, sigma_0);
```

2 Fisher Linear Discriminant

2.1 Discriminant Analysis

end

• PCA finds optimal data representations in the least square sense, however this does not imply that the transformed features will produce increased class separability



Figure 1: Bayesian parameter estimation example.



Figure 2: Bayesian parameter estimation example over a densely sampled space.

• On the other hand discriminant analysis techniques look for directions that distinguish between classes

2.2 Problem Definition

- Let's consider the problem of projecting data from d dimensions onto a line
- Let $\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n$ be the set of n points in a d dimensional space divided into subsets \mathcal{D}_i belonging to categories ω_i with cardinalities n_i , where i = 1, 2.
- Then the projections on to the direction determined by \boldsymbol{w} with $|\boldsymbol{w}| = 1$ are

$$y = \boldsymbol{w}^T \boldsymbol{x}$$

- The projections produce a set of n samples y_i with i = 1, ..., n divided into subsets \mathcal{Y}_1 and \mathcal{Y}_2
- Our problem is to find the direction of \boldsymbol{w} that will maximize the separation between the projected points in \mathcal{Y}_1 and \mathcal{Y}_2

2.3 Class Separability

2.4 Criterion Function

• Fisher Linear Discriminant seeks maximization of $J(\boldsymbol{w})$ defined as

$$J(\boldsymbol{w}) = \frac{|m_{y1} - m_{y2}|^2}{s_{y1}^2 + s_{y2}^2}$$

where

 m_{yi} is the sample mean of ω_i in the projected space:

$$m_{yi} = (1/n_i) \sum_{y \in \mathcal{Y}_i} y = (1/n_i) \sum_{x \in \mathcal{D}_i} \boldsymbol{w}^T \boldsymbol{x} = \boldsymbol{w}^T (1/n_i) \sum_{x \in \mathcal{D}_i} \boldsymbol{x} = \boldsymbol{w}^T m_{xi}$$

 s_{yi}^2 is the scatter for projected samples of ω_i :

$$s_{yi}^2 = \sum_{y \in \mathcal{Y}_i} \left(y - m_{yi} \right)^2$$



Figure 3: Projection of data onto different directions defined by \boldsymbol{w} . Observe that projection displayed in the right figure produces greater separability than the projection displayed in the left figure

2.5 Scatter Matrices

Further we define

- Scatter matrices: $S_i = \sum_{x \in D_i} (\boldsymbol{x} \boldsymbol{m}_{xi}) (\boldsymbol{x} \boldsymbol{m}_{xi})^T$, i = 1, 2
- Within-class scatter matrix: $S_W = S_1 + S_2$
- Because

$$s_{yi}^{2} = \sum_{y \in \mathcal{Y}_{i}} (y - m_{yi})^{2} = \sum_{y \in \mathcal{Y}_{i}} (\boldsymbol{w}^{T} \boldsymbol{x} - \boldsymbol{w}^{T} \boldsymbol{m}_{xi})^{2}$$
$$= \boldsymbol{w}^{T} \sum_{y \in \mathcal{Y}_{i}} \left[(\boldsymbol{x} - \boldsymbol{m}_{xi}) (\boldsymbol{x} - \boldsymbol{m}_{xi})^{T} \right] \boldsymbol{w} = \boldsymbol{w}^{T} S_{i} \boldsymbol{w},$$
$$s_{y1}^{2} + s_{y2}^{2} = \boldsymbol{w}^{T} S_{W} \boldsymbol{w}$$

• Consider the numerator of $J(\boldsymbol{w})$:

$$|m_{y1} - m_{y2}|^{2} = (m_{y1} - m_{y2})^{2} = (\boldsymbol{w}^{T} \boldsymbol{m_{x1}} - \boldsymbol{w}^{T} \boldsymbol{m_{x2}})^{2}$$
$$= \boldsymbol{w}^{T} (\boldsymbol{m_{x1}} - \boldsymbol{m_{x2}})^{2} = \boldsymbol{w}^{T} (\boldsymbol{m_{x1}} - \boldsymbol{m_{x2}}) (\boldsymbol{m_{x1}} - \boldsymbol{m_{x2}})^{T} \boldsymbol{w}$$
$$= \boldsymbol{w}^{T} S_{B} \boldsymbol{w},$$

where S_B is the between-class scatter matrix

$$S_B = (m_{x1} - m_{x2})(m_{x1} - m_{x2})^T$$

- Proportional to the sample covariance matrix
- Symmetric and positive-semidefinite
- Nonsingular if n > d
- Outer product of two vectors
- Symmetric and positive-semidefinite
- Its rank is at most 1

2.6 Optimizing the Criterion Function

• We use the scatter matrix definitions it follow that the criterion function is:

$$J(\boldsymbol{w}) = \frac{\boldsymbol{w}^T S_B \boldsymbol{w}}{\boldsymbol{w}^T S_W \boldsymbol{w}}$$

- This is a Rayleigh quotient
- The \boldsymbol{w} that maximizes $J(\boldsymbol{w})$ must satisfy $S_B \boldsymbol{w} = \lambda S_W \boldsymbol{w}$ (generalized eigenvalue problem)
- If S_W is nonsingular, we have the conventional eigenvalue problem

$$S_W^{-1}S_B\boldsymbol{w} = \lambda \boldsymbol{w}$$

• We do not need to solve

$$S_W^{-1}S_B\boldsymbol{w} = \lambda \boldsymbol{w}$$

- Recall that $S_B \boldsymbol{w}$ is at the direction of $\boldsymbol{m}_1 \boldsymbol{m}_2$
- Hence the solution is:

$$\boldsymbol{w} = S_W^{-1}(\boldsymbol{m}_1 - \boldsymbol{m}_2)$$

• After the projection, we make a decision in the unidimensional space

2.7 Classification Rule

• Assuming multivariate normal class-conditional densities $p(\boldsymbol{x}|\omega_i)$ with equal covariance matrices Σ , we recall from Ch. 2 that at the decision boundary

 $\boldsymbol{w}^T\boldsymbol{x} + w_0 = 0,$

where

$$\boldsymbol{w} = \Sigma^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$

- When we use the sample means and sample covariance matrix it follows that \boldsymbol{w} is the one that maximizes the Fisher linear discriminant
- In this case, to classify we apply a threshold to Fisher's linear discriminant

Exercise 2. Consider the Fisher linear discriminant method

- 1. Write a general program to calculate the optimal direction \boldsymbol{w} for a Fisher linear discriminant based on three-dimensional data.
- 2. Find the optimal \boldsymbol{w} for categories ω_2 and ω_3 in Table 1.
- 3. Plot a line representing your optimal direction \boldsymbol{w} and mark on it the positions of the projected points.
- 4. In this subspace, fit each distribution with a univariate Gaussian, and find the resulting decision boundary.
- 5. What is the training error (the error on training points themselves) in the optimal subspace you found in part (2)?
- 6. For comparison, repeat parts (4) and (5) using instead the nonoptimal direction $\boldsymbol{w} = (1.0, 2.0, -1.5)^T$. What is the training error in this nonoptimal subspace?

```
function [Y_class, w, X_class] = ch3_fisher_linear_discriminant(X,
   total_classes, class_numbers)
% Compute discriminant and project data to it.
% S. Makrogiannis, Delaware State Univ, 10/2015.
% Get number of classes
c = total_classes;
[n, c_times_d] = size(X);
d = c_times_d / c;
class_numbers_length = numel(class_numbers);
% Compute Sw and its inverse.
Sw = zeros(d, d);
X_{class} = cell(c, 1);
for i=class_numbers(1):class_numbers(class_numbers_length)
   % Compute scatter matrix for each class.
   first_column = d*(i-1)+1;
   last_column = d*i;
   X_class{i} = X(:, first_column:last_column)';
   mean_vector(:, i) = mean(X_class{i}, 2);
   M = repmat(mean_vector(:, i), 1, n);
```

```
S{i} = (X_class{i} - M) * (X_class{i} - M)';
  % Add scatter matrices.
  Sw = Sw + S{i};
end
% Compute vector w
Sw_Inv = inv(Sw);
w = Sw_Inv * ...
  ( mean_vector( :, class_numbers(1) ) - ...
  mean_vector( :, class_numbers(class_numbers_length) ) );
% Project data to w.
for i=class_numbers(1):class_numbers(class_numbers_length)
  Y_class{i} = w' * X_class{i};
end
end
```

```
end
```

% Fisher linear discriminant.

```
% Load data.
A = load('ch3_dhs_samples.dat');
c = 3;
[n, c_times_d] = size(A);
d = c_times_d / c;
% Find optimal w for categories omega1 and omega2.
class_numbers = [2, 3];
[Y_class, w, X_class] = ch3_fisher_linear_discriminant(A, 3,
   class_numbers);
% Plot a line representing w and the positions of the plotted
   points.
figure, plot3(X_class{2}(1,:), X_class{2}(2,:), X_class{2}(3,:),
   'bo', 'linewidth', 4); hold on;
plot3(X_class{3}(1,:), X_class{3}(2,:), X_class{3}(3,:), 'gx',
   'linewidth', 4);
% plot3([0, w(1)], [0, w(2)], [0, w(3)], 'k-', 'linewidth', 4);
   grid on;
```

```
title('Points and discriminant vector', 'fontsize', 18);
saveas(gcf, 'Fisher_Linear_Discriminant_Lab.png')
Projection_Vector = cell(3, 1);
for i=1:n
   Projection_Vector{2}(1:c,i) = Y_class{2}(i) * w(1);
   Projection_Vector{3}(1:c,i) = Y_class{3}(i) * w(1);
end
figure, plot3(Projection_Vector{2}(1,:),
   Projection_Vector{2}(2,:), Projection_Vector{2}(3,:), ...
   'bo', 'linewidth', 4); hold on;
plot3(Projection_Vector{3}(1,:), Projection_Vector{3}(2,:),
   Projection_Vector{3}(3,:), ...
   'gx', 'linewidth', 4);
% plot3([0, w(1)], [0, w(2)], [0, w(3)], 'k-', 'linewidth', 4);
   grid on;
title('Projection onto line', 'fontsize', 18);
saveas(gcf, 'Fisher_Linear_Discriminant_Lab_02.png')
figure, plot(Y_class{2}, ones(n, 1), 'bo'); hold on;
plot(Y_class{3}, ones(n, 1), 'gx', 'linewidth', 4); grid on;
title('Points in 1-d space', 'fontsize', 18);
saveas(gcf, 'Fisher_Linear_Discriminant_Lab_03.png')
% Fit each distribution with a univariate Gaussian.
mu_2 = mean(Y_class{2});
sigma_2 = std(Y_class{2});
mu_3 = mean(Y_class{3});
sigma_3 = std(Y_class{3});
% Find decision boundary.
y_0 = 0.06;
% Calculate training error.
Y_Data = [Y_class{2}, Y_class{3}];
L_Data = [2* ones(1,n), 3* ones(1,n)];
Decision = Y_Data < y_0;</pre>
```

```
Decision = Decision + 2;
classification_rate = 100 * (sum( Decision == L_Data) /
    numel(L_Data));
fprintf('Overall classification rate = %f\n', classification_rate);
% Repeat above process for w = [1, 2, -1.5]' and compute the
    training
% error.
```



Points and discriminant vector

Figure 4: Fisher linear discriminant example.

Parameter Estimation MTSC 852, Fall 2015

Sokratis Makrogiannis, Ph.D.

Department of Mathematical Sciences Delaware State University smakrogiannis@desu.edu

September 22, 2015

S. Makrogiannis (DSU)

Parameter Estimation

September 22, 2015 1 / 14

▲ロト ▲掃ト ▲ヨト ▲ヨト ニヨー のへで

Outline



1 Parameter Estimation



(2) Maximum Likelihood Estimation



S. Makrogiannis (DSU)

Parameter Estimation

September 22, 2015 2 / 14

Challenges

- According to Bayesian Decision Theory we can build a classifier when the prior probability $P(\omega)$ and likelihood $p(\mathbf{x}|\omega_i)$ are known
- In the real-world these probabilities and densities are unknown and must be estimated from data
- In specific, the estimation of likelihood (density) may be challenging especially, when the number of samples is limited and the dimensionality is high
- There are two main approaches to this density estimation: parametric and nonparametric

S. Makrogiannis (DSU)

Parameter Estimation

イロト イポト イヨト イヨト

Solutions for Density Estimation

Parameter Estimation

- We assume a model for the density function, for example, Gaussian, Rician, etc, and we need to estimate the parameters of the function, for example, mean and variance.
- Main techniques for parameter estimation
 - Maximum Likelihood
 - Bayesian Estimation

Nonparametric Estimation

- We make no assumptions about the form of the density
- Main techniques for nonparametric estimation
 - Parzen Kernels
 - Nearest Neighbor Rule

S. Makrogiannis (DSU)

Parameter Estimation

< ロト < 同ト < ヨト < ヨト

The Parameter Estimation Problem

Assumptions

- A set of *n* training samples/patterns $\mathscr{D} = \{\mathbf{x}_1, \mathbf{x}_2, ..., \mathbf{x}_n\}$
- Samples are drawn from a distribution $p(\mathbf{x}|\omega_i)$ with a known parametric form, e.g. $p(\mathbf{x}|\omega_i) \sim N(\mu_i, \Sigma_i)$
- The parameters are collectively represented by $heta,\; heta=(\mu_i,\Sigma_i)$
- Therefore, density can be written as $p(\mathbf{x}|\mathbf{ heta})$

Objective

Given a set of n training samples/patterns \mathscr{D} , estimate θ

S. Makrogiannis (DSU)

Parameter Estimation

Maximum Likelihood vs Bayesian Estimation

Maximum Likelihood (ML) Estimation

- Parameters heta are assumed to be fixed
- The solution is found as set that yields the best fitting model

 $\hat{ heta} = rg \max p(\mathscr{D}| heta)$

Bayesian Estimation

- \bullet Parameters θ are considered to be random variables with known prior distribution
- Given the observations \mathscr{D} we estimate the posterior $p(heta|\mathscr{D})$

S. Makrogiannis (DSU)

Parameter Estimation

September 22, 2015 6 / 14

Maximum Likelihood Estimation

- As explained before, we seek to estimate p(x|ω_i, θ)
- To achieve this we look for the parameters θ̂ that best describe the *n* samples
 D = {x₁, x₂, ..., x_n}
- This is equivalent to finding the value $\hat{\theta}$, such that $\hat{\theta} = \arg \max p(\mathscr{D}|\theta)$



Maximum Likelihood Estimation

Assuming that samples in *D* are drawn independently,

$$p(\mathscr{D}|\theta) = \prod_{k=1}^{n} p(\mathbf{x}_{\mathbf{k}}|\theta)$$

If f(θ) = p(𝒫|θ) is a differentiable function, we can use differential calculus to find the maximizer from

$$abla_{ heta}f(heta) = 0$$

• Let $heta = (heta_1, heta_2, ..., heta_p)$. Then $abla_{ heta} = [\frac{\partial}{\partial heta_1} \quad \frac{\partial}{\partial heta_2} \quad ... \quad \frac{\partial}{\partial heta_p}]^T$

S. Makrogiannis (DSU)

Parameter Estimation

September 22, 2015 8 / 14

< ロト < 同ト < ヨト < ヨト

Maximum Likelihood Estimation

• For analytical tractability reasons let us optimize the logarithm of *f*. Then

$$\hat{ heta} = rg\max \mathsf{ln}\, f(heta) = rg\max \mathsf{ln}\, \prod_{k=1}^n p(\mathbf{x_k}| heta) = rg\max \sum_{k=1}^n \mathsf{ln}\, p(\mathbf{x_k}| heta)$$

 According to previous treatment we obtain solution from a set of p equations

$$abla_{ heta} \sum_{k=1}^{n} \ln p(\mathbf{x_k}|m{ heta}) = 0$$

S. Makrogiannis (DSU)

Parameter Estimation

September 22, 2015 9 / 14

Case 1: Gaussian with unknown μ

Problem statement

- \bullet We assume a multivariate normal population with unknown mean μ and known covariance matrix Σ
- ullet We seek to estimate μ
- Log likelihood for each sample: $\ln p_k(\mathbf{x_k}|\mu) = (-1/2) \ln \left[(2\pi)^d |\Sigma| \right] - (1/2) (\mathbf{x}_k - \mu)^T \Sigma^{-1} (\mathbf{x}_k - \mu)$
- To optimize we solve:

$$\sum_{k=1}^{n} \nabla_{\mu} \ln p_{k}(\mathbf{x}_{k} | \hat{\mu}) = 0 \Leftrightarrow \sum_{k=1}^{n} \Sigma^{-1}(\mathbf{x}_{k} - \hat{\mu}) = 0$$

Hence:

$$\hat{\mu} = (1/n) \sum_{k=1}^{n} \mathbf{x}_k$$

S. Makrogiannis (DSU)

Parameter Estimation

10 / 14

Case 2: Gaussian with unknown μ and unknown σ

Problem statement

- ${\, \circ \,}$ We assume a univariate normal population with unknown mean μ and unknown covariance matrix σ
- We seek to estimate $heta=(\mu,\sigma)$
- Log likelihood for each sample: $\ln p_k(x_k|\theta) = (-1/2)\ln (2\pi\sigma^2) - (1/2\sigma^2)(\mathbf{x}_k - \mu)^2$ • Derivative is $\nabla_{\theta} \ln p_k(x_k|\theta) = \left[(1/\sigma^2)(x_k - \mu), \frac{-1}{2\sigma^2} + \frac{(x_k - \mu)^2}{2\sigma^4}\right]^T$ • We solve: $\sum_{k=1}^n (1/\hat{\sigma}^2)(x_k - \hat{\mu}) = 0$ and $\sum_{k=1}^n \left[\frac{-1}{2\hat{\sigma}^2} + \frac{(x_k - \hat{\mu})^2}{2\hat{\sigma}^4}\right] = 0$ • After some rearranging we get:

$$\hat{\mu} = (1/n) \sum_{k=1}^{n} x_k, \quad \hat{\sigma}^2 = (1/n) \sum_{k=1}^{n} (x_k - \hat{\mu})^2$$

S. Makrogiannis (DSU)

Case 3: Gaussian with unknown μ and unknown Σ

Problem statement

 \bullet We assume a multivariate normal population with unknown mean μ and unknown covariance matrix Σ

• We seek to estimate $heta=(\mu,\Sigma)$

- Log likelihood for each sample: $\ln p_k(\mathbf{x_k}|\mu) = (-1/2) \ln \left[(2\pi)^d |\Sigma| \right] - (1/2) (\mathbf{x}_k - \mu)^T \Sigma^{-1} (\mathbf{x}_k - \mu)$
- Next, we have to solve a system similar to Case 2: $\sum_{k=1}^{n} \nabla_{\theta} \ln p(\mathbf{x}_{k}|\theta) = 0$
- The maximum-likelihood solution is: $\hat{\mu} = (1/n) \sum_{k=1}^{n} \mathbf{x}_{\mathbf{k}}, \quad \hat{\Sigma} = (1/n) \sum_{k=1}^{n} (\mathbf{x}_{\mathbf{k}} - \hat{\mu}) (\mathbf{x}_{\mathbf{k}} - \hat{\mu})^{T}$

S. Makrogiannis (DSU)

Parameter Estimation

< □ ト < □ ト < 三 ト < 三 ト < 三 ト 三 つ Q (* September 22, 2015 12 / 14 Bias

Bias of Maximum-Likelihood Estimation Technique

 Maximum likelihood estimates for a Gaussian with unknown μ and unknown Σ:

$$\hat{\mu} = (1/n)\sum_{k=1}^{n} \mathbf{x}_{\mathbf{k}}, \quad \hat{\Sigma} = (1/n)\sum_{k=1}^{n} (\mathbf{x}_{\mathbf{k}} - \hat{\mu})(\mathbf{x}_{\mathbf{k}} - \hat{\mu})^{T}$$

• Sample mean and sample covariance matrix:

$$\mu = (1/n) \sum_{k=1}^{n} \mathbf{x}_{\mathbf{k}}, \quad C = \frac{1}{n-1} \sum_{k=1}^{n} (\mathbf{x}_{\mathbf{k}} - \hat{\mu}) (\mathbf{x}_{\mathbf{k}} - \hat{\mu})^{\mathsf{T}}$$

• Hence $\hat{\mu} = \mu$, $\hat{\Sigma} = \frac{n-1}{n}C$

Therefore μ̂ is an unbiased estimate of the mean, but Σ̂ is biased
 Σ̂ → C when n→∞, therefore Σ̂ is called asymptotically unbiased

S. Makrogiannis (DSU)

Parameter Estimation

September 22, 2015 13 / 14

I D I I I I I
Further Comments on Maximum-Likelihood Estimation

- Maximum likelihood estimation is usually simpler than other estimation methods
- The estimates become more accurate as the number of samples increases
- If the likelihood model $p(x|\theta)$ is accurate, it yields very good results
- Model selection is a key process that is studied by the field of Algorithm-Independent Machine Learning

S. Makrogiannis (DSU)

Parameter Estimation

September 22, 2015 14 / 14

Parameter Estimation MTSC 852, Fall 2015

Sokratis Makrogiannis, Ph.D.

Department of Mathematical Sciences Delaware State University smakrogiannis@desu.edu

September 29, 2015

S. Makrogiannis (DSU)

Parameter Estimation

September 29, 2015 1 / 23

▲ロト ▲掃ト ▲ヨト ▲ヨト ニヨー のへで



Bayesian Parameter Estimation

2 Bayesian Parameter Estimation for the Gaussian Density

3 Bayesian Parameter Estimation for Arbitrary Density Models

S. Makrogiannis (DSU)

Parameter Estimation

September 29, 2015 2 / 23

Sac

イロト 不得下 不良下 不良下 一度

Bayesian Estimation

- We consider the parameter vector heta to be a random variable with prior p(heta)
- We use training data D to estimate the posterior probability density p(θ|D)
- Finally we integrate over the parameter space to estimate the class-conditional density $p(\mathbf{x}|\boldsymbol{\omega}_i, \mathcal{D})$

イロト イポト イヨト イヨト

Class-conditional Densities

- In Bayesian classification we face the problem of making a decision using posterior probabilities $P(\omega_i | \mathbf{x})$
- We need to know the likelihood p(x|ω_i) and priors P(ω_i) to calculate the posteriors
- We have seen that the estimation of class-conditional density (or likelihood) is non-trivial in real-world applications
- To solve this problem we first assume a parametric functional form for the likelihood and then use training samples to estimate the likelihood for each class

S. Makrogiannis (DSU)

Parameter Estimation

September 29, 2015 4 / 23

< ロト < 同ト < ヨト < ヨト

Class-conditional Densities

• Given a training sample set $\mathscr{D} = \{\mathbf{x}_1, \mathbf{x}_2, ..., \mathbf{x}_n\}$, the Bayesian rule becomes $P(\boldsymbol{\omega}_i | \mathbf{x}, \mathscr{D}) = \frac{p(\mathbf{x} | \boldsymbol{\omega}_i, \mathscr{D}) P(\boldsymbol{\omega}_i | \mathscr{D})}{p(\mathbf{x} | \boldsymbol{\omega}_i, \mathscr{D}) P(\boldsymbol{\omega}_i | \mathscr{D})}$

$$P(\omega_i | \mathbf{x}, \mathscr{D}) = \frac{p(\mathbf{x} | \omega_i, \mathscr{D}) P(\omega_i | \mathscr{D})}{\sum_{j=1}^{c} p(\mathbf{x} | \omega_j, \mathscr{D}) P(\omega_j | \mathscr{D})}$$

• If we assume that the class-conditionals of different classes are statistically independent it follows that

$$P(\omega_i | \mathbf{x}, \mathscr{D}_i) = \frac{p(\mathbf{x} | \omega_i, \mathscr{D}_i) P(\omega_i)}{\sum_{j=1}^{c} p(\mathbf{x} | \omega_j, \mathscr{D}_i) P(\omega_j)}$$

• Now we need to solve *c* likelihood estimation problems

S. Makrogiannis (DSU)

Parameter Estimation

September 29, 2015 5 / 23

Bayesian Estimation Steps

- According to previous analysis, we seek to estimate p(x|D) for each class
- This is achieved by integration over the parameter space

$$p(\mathbf{x}|\mathscr{D}) = \int p(\mathbf{x}, \theta|\mathscr{D}) d\theta$$

• Using definition of joint probability:

$$p(\mathbf{x}|\mathscr{D}) = \int p(\mathbf{x}|\boldsymbol{ heta}, \mathscr{D}) p(\boldsymbol{ heta}|\mathscr{D}) d\boldsymbol{ heta}$$

- Bayes rule to estimate $p(\theta|\mathscr{D})$: $p(\theta|\mathscr{D}) = \frac{p(\mathscr{D}|\theta)P(\theta)}{\int p(\mathscr{D}|\theta)P(\theta)d\theta}$
- If samples are independently drawn, then:

$$p(\mathscr{D}|\theta) = \prod_{i=1}^{n} p(\mathbf{x}_i|\theta)$$

S. Makrogiannis (DSU)

Parameter Estimation

September 29, 2015 6 / 23

◆□▶ ◆母▶ ◆ヨ▶ ◆ヨ▶ ヨー つくや

Univariate Normal Density

Problem

Find the class-conditional density $p(x|\mathscr{D})$ using Bayesian estimation assuming that $p(x|\mu) \sim N(\mu, \sigma^2)$, $p(\mu) \sim N(\mu_0, \sigma_0)$ and σ is known.

Solution steps

- Estimate $p(\mu|\mathscr{D})$ using Bayes rule
- Estimate $p(x|\mathscr{D})$ by integration over the parameter space

S. Makrogiannis (DSU)

Parameter Estimation

September 29, 2015 7 / 23

Univariate Normal Density

Problem

Find the class-conditional density $p(x|\mathscr{D})$ using Bayesian estimation assuming that $p(x|\mu) \sim N(\mu, \sigma^2)$, $p(\mu) \sim N(\mu_0, \sigma_0)$ and σ is known.

Estimate $p(\mu|\mathscr{D})$

- According to previous analysis, we seek to estimate p(x|D) for each class
- This is achieved by integration over the parameter space

$$p(\mathbf{x}|\mathscr{D}) = \int p(\mathbf{x}, \theta|\mathscr{D}) d\theta$$

- From definition of joint probability: $p(\mathbf{x}|\mathscr{D}) = \int p(\mathbf{x}|\theta, \mathscr{D}) p(\theta|\mathscr{D}) d\theta$
- We use Bayes rule to estimate $p(\theta|\mathscr{D})$: $p(\theta|\mathscr{D}) = \frac{p(\mathscr{D}|\theta)P(\theta)}{\int p(\mathscr{D}|\theta)P(\theta)d\theta}$
- If the samples are independently drawn, then: $p(\mathscr{D}|\theta) = \prod_{i=1}^{n} p(\mathbf{x}_i|\theta)$

• We use Bayes rule to estimate posterior parameter density $p(\mu|\mathscr{D})$:

$$p(\mu|\mathscr{D}) = \frac{p(\mathscr{D}|\mu)p(\mu)}{\int p(\mathscr{D}|\mu)p(\mu)d\mu}$$

- Let samples $\mathscr{D} = \{\mathbf{x_1}, \mathbf{x_2}, ..., \mathbf{x_n}\}$ be independently drawn. Then: $p(\mathscr{D}|\mu) = \prod_{i=1}^n p(\mathbf{x_i}|\mu)$
- Then: $p(\mu|\mathscr{D}) = lpha \cdot \prod_{i=1}^n p(x_i|\mu) p(\mu)$
- According to assumptions:

$$p(x_i|\mu) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}}, \quad p(\mu) = \frac{1}{\sqrt{2\pi\sigma_0}} e^{-\frac{(\mu-\mu_0)^2}{2\sigma_0^2}}$$

So: $p(\mu|\mathscr{D}) = \alpha \cdot \prod_{i=1}^n \left[\frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}} \frac{1}{\sqrt{2\pi\sigma_0}} e^{-\frac{(\mu-\mu_0)^2}{2\sigma_0^2}} \right]$
S. Makrogiannis (DSU) Parameter Estimation September 29, 2015 9 / 23

- After some more manipulations we can show that $p(\mu|\mathscr{D})$ is an exponential function of a quadratic function, hence it has the form $N(\mu_n, \sigma_n)$
- Therefore

$$p(\mu|\mathscr{D}) = \frac{1}{\sqrt{2\pi}\sigma_n} e^{-\frac{(\mu-\mu_n)^2}{2\sigma_n^2}}$$

where,

$$\mu_n = \left(\frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2}\right)\hat{\mu}_n + \frac{\sigma^2}{n\sigma_0^2 + \sigma^2}\mu_0, \quad \sigma_n^2 = \frac{\sigma_0^2\sigma^2}{n\sigma_0^2 + \sigma^2}$$
$$\hat{\mu}_n = (1/n)\sum_{i=1}^n x_i$$

S. Makrogiannis (DSU)

Parameter Estimation

September 29, 2015 10 / 23

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 ◇◇◇

Bayesian learning

- σ_n decreases as $n \to \infty$ with $\lim_{n \to \infty} \sigma_n^2 = \frac{\sigma^2}{n}$
- We observe that as the number of training samples increases, $p(\mu|\mathscr{D})$ becomes sharper around μ_n . This process is called *Bayesian learning*
- If $\sigma_0
 eq 0$, then μ_n approaches the sample mean $\lim_{n o \infty} \mu_n = \hat{\mu}_n$.



Estimate $p(x|\mathscr{D})$

• According to Bayesian estimation process

$$p(x|\mathscr{D}) = \int p(x|\mu,\mathscr{D})p(\mu|\mathscr{D})d\mu \Leftrightarrow$$

$$p(x|\mathscr{D}) = \int \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \frac{1}{\sqrt{2\pi\sigma_n}} e^{-\frac{(\mu-\mu_n)^2}{2\sigma_n^2}} d\mu \Leftrightarrow$$

$$p(x|\mathscr{D}) = \frac{1}{2\pi\sigma\sigma_n} f(\sigma,\sigma_n) e^{-\frac{(x-\mu_n)^2}{2(\sigma^2+\sigma_n^2)}}$$

where $f(\sigma, \sigma_n)$ has an integral form:

$$f(\sigma,\sigma_n) = \int \exp\left[(-1/2)\frac{\sigma^2 + \sigma_n^2}{\sigma^2 \sigma_n^2} \left(\mu - \frac{\sigma_n^2 x + \sigma^2 \mu_n}{\sigma^2 + \sigma_n^2}\right)^2\right] d\mu$$

S. Makrogiannis (DSU)

Parameter Estimation

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Estimate $p(x|\mathscr{D})$

- Observe that $p(x|\mathscr{D}) \sim \mathcal{N}(\mu_n, \sigma^2 + \sigma_n^2)$
- The above result gives the class-conditional density $p(x|\omega_i, \mathcal{D}_i)$ based on the posterior parameter mean estimate μ_n and the posterior parameter variance estimate increased by the uncertainty in x that we assume to be known

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 ◇◇◇

Multivariate Normal Density

Problem

Find the class-conditional density $p(\mathbf{x}|\mathscr{D})$ using Bayesian estimation assuming that $p(\mathbf{x}|\mu) \sim N(\mu, \Sigma)$, $p(\mu) \sim N(\mu_0, \Sigma_0)$, and that Σ , μ_0 , Σ_0 are known.

Solution steps

- Estimate $p(\mu|\mathscr{D})$ using Bayes rule
- Estimate $p(\mathbf{x}|\mathscr{D})$ by integration over the parameter space

S. Makrogiannis (DSU)

Parameter Estimation

September 29, 2015 14 / 23

 Similarly to the unidimensional case we use the Bayes rule to estimate the posterior parametric density, and the assumption for idependent samples x_i ∈ 𝔅 with |𝔅| = n

$$p(\mu|\mathscr{D}) = lpha \prod_{i=1}^n p(\mathbf{x}_i|\mathscr{D}) p(\mu) \Leftrightarrow$$

$$p(\mu|\mathscr{D}) = \alpha'' e^{(-1/2)(\mu-\mu_n)^T \sum_n^{-1} (\mu-\mu_n)}$$

- Hence: $p(\mu | \mathscr{D}) \sim N(\mu_n, \Sigma_n)$
- By equating coefficients we get:

$$\Sigma_n^{-1} = n\Sigma^{-1} + \Sigma_0^{-1}, \quad \Sigma_n^{-1}\mu_n = n\Sigma^{-1}\hat{\mu}_n + \Sigma_0^{-1}\mu_0$$

where $\hat{\mu}_n = (1/n) \sum_{k=1}^n \mathsf{x}_k$

S. Makrogiannis (DSU)

September 29, 2015 15 / 23

イロト (母) (ヨ) (ヨ) (ヨ) () ()

• Finally after some manipulations it follows that

$$\mu_n = \Sigma_0 [\Sigma_0 + (1/n)\Sigma]^{-1} \hat{\mu}_n + (1/n)\Sigma [\Sigma_0 + (1/n)\Sigma]^{-1} \mu_0$$

$$\Sigma_n = \Sigma_0 [\Sigma_0 + (1/n)\Sigma]^{-1} (1/n)\Sigma$$

$$\hat{\mu}_n = (1/n) \sum_{i=1}^n \mathbf{x}_i$$

S. Makrogiannis (DSU)

Parameter Estimation

September 29, 2015 16 / 23

Estimate $p(\mathbf{x}|\mathscr{D})$

• To estimate the class-conditional density we must compute the integral:

$$p(\mathbf{x}|\mathscr{D}) = \int p(\mathbf{x}|\mu) p(\mu|\mathscr{D}) d\mu$$

Finally, we can show that p(x|𝒴) ~ N(μ_n, Σ + Σ_n), either
 (a) by performing integration, or
 (b) by using the central limit theorem and the observation that x can be considered to be the sum of random variables μ with p(μ|𝒴) ~ N(μ_n, Σ_n) and y with p(y) ~ N(0, Σ)

S. Makrogiannis (DSU)

Parameter Estimation

September 29, 2015 17 / 23

◆□▶ ◆母▶ ◆ヨ▶ ◆ヨ▶ ヨー つくや

Bayesian Parameter Estimation for Arbitrary Density Models

We can generalize the Bayesian technique to arbitrary density models $p(\mathbf{x}|\mathscr{D})$ with parameters θ

Main assumptions

- The form of $p(\mathbf{x}| heta)$ is known but the parameter vector heta is unknown
- We know prior density $p(\theta)$
- We can learn more about θ from a set \mathscr{D} of statistically independent n samples \mathbf{x}_k , for k = 1, ..., n, following the unknown $p(\mathbf{x})$

S. Makrogiannis (DSU)

Parameter Estimation

September 29, 2015 18 / 23

イロト イポト イヨト イヨト 二日

Bayesian Parameter Estimation for Arbitrary Density Models

How to compute the class-conditional

- Using definition of joint probability: $p(\mathbf{x}|\mathscr{D}) = \int p(\mathbf{x}|\theta, \mathscr{D}) p(\theta|\mathscr{D}) d\theta$
- Bayes rule to estimate $p(\theta|\mathscr{D})$: $p(\theta|\mathscr{D}) = \frac{p(\mathscr{D}|\theta)P(\theta)}{\int p(\mathscr{D}|\theta)P(\theta)d\theta}$
- If samples are independently drawn, then: $p(\mathscr{D}|\theta) = \prod_{i=1}^n p(\mathsf{x}_i|\theta)$

Comments

- If $p(\mathscr{D}|\theta)$ reaches a sharp peak at $\theta = \hat{\theta}$ and $p(\theta)$ is not zero at $\theta = \hat{\theta}$, then $p(\theta|\mathscr{D})$ also peaks at $\theta = \hat{\theta}$
- Then $p(\mathbf{x}|\mathscr{D}) \simeq p(\mathbf{x}|\hat{\theta})$, hence the ML and Bayesian estimates will be close
- Still, Bayesian estimation uses more information for the class-conditional density estimate than the ML technique

S. Makrogiannis (DSU)

Parameter Estimation

<□ ト < □ ト < □ ト < 三 ト < 三 ト < 三 ト ○ Q ○ September 29, 2015 19 / 23

Recursive Bayes Incremental Learning

- Suppose that we are given a set of training samples $\mathscr{D}^n = {x_1, x_2, ..., x_n}$, where n > 1
- Then $p(\mathscr{D}^n|\theta) = p(\mathbf{x}_n|\theta)p(\mathscr{D}^{n-1}|\theta)$
- Using the Bayes rule: $p(\theta|\mathscr{D}^n) = \frac{p(\mathbf{x}_n|\theta)p(\mathscr{D}^{n-1}|\theta)P(\theta)}{\int p(\mathbf{x}_n|\theta)p(\mathscr{D}^{n-1}|\theta)P(\theta)d\theta} = \frac{p(\mathbf{x}_n|\theta)p(\theta|\mathscr{D}^{n-1})}{\int p(\mathbf{x}_n|\theta)p(\theta|\mathscr{D}^{n-1})d\theta}$
- Starting from $p(\theta | \mathscr{D}^0) = p(\theta)$ we can recursively compute $p(\theta | \mathbf{x}_1), p(\theta | \mathbf{x}_1, \mathbf{x}_2), ..., p(\theta | \mathbf{x}_1, ..., \mathbf{x}_n)$
- This is an on-line learning method that updates the model as more training data are collected and is called *recursive Bayes approach to parameter estimation*
- When the computation converges to a Dirac delta function centered at the true parameter value, we have achieved *Bayesian Learning*

S. Makrogiannis (DSU)

Parameter Estimation

< □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ ▶ < □ > ○ Q ()
September 29, 2015 20 / 23

Identifiability in Recursive Bayes Incremental Learning

- For most typical density models $p(\mathbf{x}|\theta)$, the sequence of posterior densities converges to a delta function, and a true value for θ can be found. Then $p(\mathbf{x}|\theta)$ is called *identifiable*
- However, in some cases there are multiple values of θ that explain the data. Fortunately, the integration operation for estimating $p(\mathbf{x}|\mathcal{D}^n)$ will still converge to $p(\mathbf{x})$ because all optimizing values of θ will yield the same $p(\mathbf{x}|\theta)$
- Non-identifiability may become an issue in unsupervised learning techniques

S. Makrogiannis (DSU)

Parameter Estimation

イロト (母) (ヨ) (ヨ) (ヨ) () ()

Comparing Maximum-Likelihood and Bayesian Parameter Estimation

- ML methods are computationally simpler than Bayesian estimation techniques because they employ differential calculus or gradient search techniques to find $\hat{\theta}$. On the other hand, Bayesian estimation involves complex multidimensional integration
- ML solutions correspond to single optimal models, hence they are easier to interpret and understand than Bayesian solutions that are produced by weighted sums of models
- Bayesian estimation techniques use more information and can produce better results for non-uniform and non-symmetric $p(\theta|\mathscr{D})$ distributions
- Bayesian techniques show the problem of bias and variance that depends on the number of training samples
- Overall Bayesian techniques are more theoretically sound, but ML techniques are simpler to implement and are nearly as accurate

S. Makrogiannis (DSU)

Parameter Estimation

September 29, 2015 22 / 23

Sac

Sources of Classification Error

When developing a classifier, we first estimate the posterior densities for each category, then classify a test sample using a maximum decision rule. The sources of error in such a system are:

- *Bayes or Indistinguishability Error* happens because of overlapping class-conditional densities
- *Model Error* caused by selection of the wrong model. This is not affected by the parameter estimation technique
- *Estimation Error* due to the finite number of training samples. It decreases with increasing cardinality of the data set

S. Makrogiannis (DSU)

Parameter Estimation

September 29, 2015 23 / 23

イロト (母) (ヨ) (ヨ) (ヨ) () ()

Parameter Estimation MTSC 852, Fall 2015

Sokratis Makrogiannis, Ph.D.

Department of Mathematical Sciences Delaware State University smakrogiannis@desu.edu

October 12, 2015

S. Makrogiannis (DSU)

Parameter Estimation

October 12, 2015 1 / 14

イロト (母) (ヨ) (ヨ) (ヨ) () ()

$\mathsf{Outline}$

1 Difficulties caused by Dimensionality

S. Makrogiannis (DSU)

Parameter Estimation

October 12, 2015 2 / 14

Difficulties Caused by Dimensionality

- Many classification problems involve samples in very high dimensional spaces, i.e. with hundreds, or thousands of features
- In this section we discuss the following considerations:
 - relation between dimensionality, number of training samples and classification accuracy
 - 2 computational complexity of classifier
 - ③ overfitting

S. Makrogiannis (DSU)

Parameter Estimation

October 12, 2015 3 / 14

Classification Accuracy vs. Number of Features

- We are usually looking for statistically independent features based on theoretical grounds
- Suppose that we have 2-classes of multivariate normal densities with equal covariance matrices, that is $p(x|\omega_j) \sim N(\boldsymbol{\mu}_j, \boldsymbol{\Sigma})$ with all $P(\omega_j)$ equal for j = 1, 2
- Then the Bayes error is: $P(e) = \frac{1}{\sqrt{2\pi}} \int_{r/2}^{\infty} e^{-u^2/2} du$, where $r^2 = (\boldsymbol{\mu}_1 \boldsymbol{\mu}_2)^T \Sigma^{-1} (\boldsymbol{\mu}_1 \boldsymbol{\mu}_2)$
- For conditionally independent features: $\Sigma = diag(\sigma_1^2, ..., \sigma_D^2)$, therefore $r^2 = \sum_{i=1}^{D} \left(\frac{\mu_{i1} \mu_{i2}}{\sigma_i}\right)^2$
- Observe that the addition of features is expected to increase r^2 , therefore reducing the Bayes error

S. Makrogiannis (DSU)

Parameter Estimation

< □ ト < □ ト < □ ト < 三 ト < 三 ト 三 の Q C October 12, 2015 4 / 14

Classification Accuracy vs. Number of Features

- Based on the previous result, the addition of features with unequal class-conditional means will increase the separability of the data
- So it reasonable to add new features if the classification performance with a given feature set is not good enough



Figure: Projection of data to spaces of reduced dimensionality increases the Bayes error.

S. Makrogiannis (DSU)

Parameter Estimation

October 12, 2015 5 / 14

Classification Accuracy vs. Number of Features and Number of Training Samples

- However, in practice, the addition of features may reduce the classification performance
- This may happen because (1) the learning model is wrong, or (2) the number of samples is insufficient for the estimation of the class-conditional densities
- This is a significant consideration in classifier design that we will revisit in Ch. 9

S. Makrogiannis (DSU)

Parameter Estimation

October 12, 2015 6 / 14

Computational Complexity

The computational load is a factor to be considered in classifier design
To express the computational load we use the order of function f(x)

Definition

Let f and h be two functions of x. We say that f(x) is of the order of h(x), denoted by f(x) = O(h(x)) and read as "big oh of h(x)", if there exist constants c and x_0 such that $|f(x)| \le c|h(x)|$, $\forall x > x_0$

Example

We assume that $f(x) = a_2x^2 + a_1x + a_0$. Then $f(x) = O(x^2)$ because for sufficiently large x, we can choose c and x_0 such that $|f(x)| \le c |x^2|$, $\forall x > x_0$

S. Makrogiannis (DSU)

Parameter Estimation

October 12, 2015 7 / 14

3

イロト イポト イヨト イヨト

Difficulties caused by Dimensionality

Computational Complexity

- When estimating the computational complexity of an algorithm we are interested in the number of additions, multiplications and divisions, or in the time and memory requirements
- Let's see the example of computational complexity for a Bayes classifier

S. Makrogiannis (DSU)

Parameter Estimation

October 12, 2015 8 / 14

Computational Complexity

ML estimation (learning stage)

•
$$\hat{\mu} = (1/n) \sum_{k=1}^{n} \mathbf{x}_{k} = O(dn)$$

• $\hat{\Sigma} = (1/n) \sum_{k=1}^{n} (\mathbf{x}_{k} - \hat{\mu}) (\mathbf{x}_{k} - \hat{\mu})^{T} = O(d^{2}n)$
• $g_{i}(\mathbf{x}) = -(1/2)(\mathbf{x} - \mu_{i})^{T} \sum_{i}^{-1} (\mathbf{x} - \mu_{i}) - (d/2) \ln 2\pi - (1/2) \ln |\Sigma_{i}| + \ln P(\omega_{i})$
• $-(1/2)(\mathbf{x} - \mu_{i})^{T} \sum_{i}^{-1} (\mathbf{x} - \mu_{i}) = O(dn) + O(d^{3})$
• $(d/2) \ln 2\pi = O(1)$
• $(1/2) \ln |\Sigma_{i}| = O(d^{2}n)$
• $\ln P(\omega_{i}) = O(n)$

 $\, \bullet \,$ Computations are repeated c times, hence the learning complexity is

$$O(cd^2n)\simeq O(d^2n),$$

assuming that $c \ll d$ or n and n > d

S. Makrogiannis (DSU)

Parameter Estimation

October 12, 2015 9 / 14

< A

Classification Complexity

Decision Rule

•
$$g_i(\mathbf{x}) = -(1/2)(\mathbf{x} - \mu_i)^T \Sigma_i^{-1}(\mathbf{x} - \mu_i) - (d/2) \ln 2\pi - (1/2) \ln |\Sigma_i| + \ln P(\omega_i)$$

• $g_i(\mathbf{x}) = -(1/2)(\mathbf{x} - \mu_i)^T \Sigma_i^{-1}(\mathbf{x} - \mu_i) : O(d^2)$
• $-(d/2) \ln 2\pi - (1/2) \ln |\Sigma_i| + \ln P(\omega_i) : O(1)$

• Computations are repeated *c* times, hence the classification complexity is

$$O(cd^2)\simeq O(d^2),$$

assuming that $c \ll d$

 Hence, classification stage is computationally simpler than learning stage

S. Makrogiannis (DSU)

October 12, 2015 10 / 14

- - E - b

< □ > < A >

Overfitting

Problem

• When the number of samples is too small for the dimensionality, the class models and decision surfaces may not be optimal

Possible Solutions

- We can reduce the dimensionality by (1) selecting a subset of our original features, or (2) computing features from the existing set
- 2 We can assume that all samples come from the same covariance matrix and pool the training data
- 3 We can find a better estimate of Σ (1) by using prior estimate Σ_0 to get a Bayesian type of estimate $\lambda \Sigma_0 + (1 \lambda)\hat{\Sigma}$ or (2) by applying a threshold to covariances or even assuming statistical independence (heuristic solution)

S. Makrogiannis (DSU)

3

イロト イポト イヨト イヨト

Overfitting Example

- We have 10 data points obtained by adding zero-mean, Gaussian noise to a parabola
- The two curves correspond to a parabola and a 10-th degree polynomial
- While the 10-th degree polynomial produces the best fit for the training data, the parabola is the closest model



S. Makrogiannis (DSU)

Parameter Estimation

October 12, 2015 12 / 14
Assume the Same Covariance Matrix for All Classes

Problem

 We are asked to design a classifier with insufficient data for distributions N(μ₁,Σ₁) and N(μ₁,Σ₂)

Solution

- We make the simplification that Σ₁ = Σ₂ = Σ, where Σ is estimated over both classes
- ${\scriptstyle \bullet}$ We normalize the data before estimating ${\scriptstyle \Sigma}$

S. Makrogiannis (DSU)

Parameter Estimation

October 12, 2015 13 / 14

The Shrinkage Technique

 This technique uses a weighted sum of individual covariances that "shrink" toward a common estimate

$$\Sigma_i(\alpha) = \frac{(1-\alpha)n_i\Sigma_i + \alpha n\Sigma}{(1-\alpha)n_i + \alpha n}$$

where $0 < \alpha < 1$ is the regularizing parameter, n_i is the number of samples for each class, n is the total number of samples, i = 1, ..., c and c is the number of classes

 We can also shrink the common covariance matrix to an identity matrix

$$\Sigma(\beta) = (1-\beta)\Sigma + \beta I,$$

where $0 < \beta < 1$

S. Makrogiannis (DSU)

Parameter Estimation

October 12, 2015 14 / 14

Parameter Estimation MTSC 852, Fall 2015

Sokratis Makrogiannis, Ph.D.

Department of Mathematical Sciences Delaware State University smakrogiannis@desu.edu

October 22, 2015

S. Makrogiannis (DSU)

Parameter Estimation

October 22, 2015 1 / 28

Outline



Component Analysis and Discriminants

- Principal Component Analysis
- Fisher Linear Discriminant



Expectation-Maximization (EM)

S. Makrogiannis (DSU)

Parameter Estimation

October 22, 2015 2 / 28

3

Sac

< ロト < 同ト < ヨト < ヨト

Component Analysis for Dimensionality Reduction

- As seen before, a large number of features can improve class separability, but also complicate density estimation, increase computational complexity, and the system may be more susceptible to overfitting
- One solution is to reduce the original dimensionality by linear combinations of features
- Linear techniques are very useful because they are analytically tractable and computationally efficient

S. Makrogiannis (DSU)

Parameter Estimation

October 22, 2015 3 / 28

イモトイモト

Component Analysis for Dimensionality Reduction

- Two popular techniques are (1) Principal Component Analysis, and (2) Fisher Linear Discriminant along with its multidimensional generalization called Multiple Discriminant Analysis
- Both techniques project data to lower dimensional spaces
- Principal Component Analysis seeks the projections that produce the more accurate representation of the data in a least-squares sense
- Fisher Linear Discriminant seeks the projections that best separate the data into the different classes in a least-squares sense

S. Makrogiannis (DSU)

Parameter Estimation

October 22, 2015 4 / 28

イロト イポト イヨト イヨト

Zero-dimensional Representation

- First stage: representing Data Set by a Single Vector x_0
- Suppose a set of n d-dimensional samples x_1, \ldots, x_n
- We want to find a vector x_0 to represent the data
- We require that x_0 minimizes the squared error criterion function

$$J_0(\mathbf{x}_0) = \sum_{k=1}^n \|\mathbf{x}_0 - \mathbf{x}_k\|^2$$

S. Makrogiannis (DSU)

October 22, 2015 5 / 28

< 口 > < 同 >

Zero-dimensional Representation

• We express $J_0(x_0)$ in terms of sample mean $m = (1/n) \sum_{k=1}^n x_k$ and it follows that

$$J_0(x_0) = \sum_{k=1}^n \|(x_0 - m) - (x_k - m)\|^2 = \dots$$

$$= \sum_{k=1}^{n} \|(\mathbf{x}_{0} - \mathbf{m})\|^{2} - \sum_{k=1}^{n} \|(\mathbf{x}_{k} - \mathbf{m})\|^{2}$$

• Hence $J_0(x_0)$ is minimized when $x_0 = m$

S. Makrogiannis (DSU)

Parameter Estimation

October 22, 2015 6 / 28

イロト イポト イヨト イヨト 二日

One-dimensional Representation

- To obtain a representation of variability of data we project data onto a line running through the sample mean
- Assuming a unit vector ${m e}$, line equation is: ${m x}={m m}+\alpha{m e}$, where lpha is scalar distance from sample mean
- Next, use $\mathbf{x}_k = \mathbf{m} + \alpha_k \mathbf{e}$, k = 1, ..., n, and find the optimal set of α_k by minimizing function

$$J_1(\alpha_1,\ldots,\alpha_n,\boldsymbol{e}) = \sum_{k=1}^n \|(\boldsymbol{m}+\alpha_k\boldsymbol{e})-\boldsymbol{x}_k\|^2$$

• We set $\frac{\partial J_1}{\partial \alpha_k} = 0$, and obtain

$$\alpha_k = \boldsymbol{e}^T (\boldsymbol{x}_k - \boldsymbol{m})$$

This is the projection of x_k - m on to the line that passes through the mean with the direction of e

S. Makrogiannis (DSU)

Parameter Estimation

October 22, 2015 7 / 28

Finding the Best Line Direction

- Next step is to find the direction *e* that produces the best representation
- To do this we first substitute the $lpha_k$ in J_1

$$J_1(\boldsymbol{e}) = \sum_{k=1}^n \alpha_k^2 - 2 \sum_{k=1}^n \alpha_k^2 + \sum_{k=1}^n \|(\boldsymbol{x}_k - \boldsymbol{m})\|^2,$$

and use the definition for scatter matrix $S = \sum_{k=1}^{n} (\mathbf{x}_k - \mathbf{m}) (\mathbf{x}_k - \mathbf{m})^T$ to obtain

$$J_1(\boldsymbol{e}) = -\boldsymbol{e}^T S \boldsymbol{e} + \sum_{k=1}^n \|\boldsymbol{x}_k - \boldsymbol{m}\|^2$$

S. Makrogiannis (DSU)

Parameter Estimation

October 22, 2015 8 / 28

イロト イポト イヨト イヨト

Finding the Best Line Direction

• Previously we obtained:

$$J_1(\boldsymbol{e}) = -\boldsymbol{e}^T S \boldsymbol{e} + \sum_{k=1}^n \|\boldsymbol{x}_k - \boldsymbol{m}\|^2$$

- Hence, we look to maximize $e^T S e$, subject to the constraint $\|e\| = 1$
- By the Langrangian multiplier optimization technique it follows that we should optimize

$$u = \boldsymbol{e}^{\mathsf{T}} \boldsymbol{S} \boldsymbol{e} - \lambda (\boldsymbol{e}^{\mathsf{T}} \boldsymbol{e} - 1)$$

Then solve

$$\frac{\partial u}{\partial \boldsymbol{e}} = \boldsymbol{0} \Rightarrow 2S\boldsymbol{e} - 2\lambda\boldsymbol{e} = \boldsymbol{0} \Rightarrow S\boldsymbol{e} = \lambda\boldsymbol{e}$$

S. Makrogiannis (DSU)

Parameter Estimation

<ロト < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Finding the Best Line Direction

- From previous result: $S \boldsymbol{e} = \lambda \boldsymbol{e}$
- We observe that: $\boldsymbol{e}^T \boldsymbol{S} \boldsymbol{e} = \lambda \boldsymbol{e}^T \boldsymbol{e} = \lambda$
- Hence, we are looking for the eigenvector e_{max} with maximum eigenvalue λ_{max} in order to find the best vector direction
- To perform the analysis we project the data onto a line through sample mean with the orientation defined by e_{max}

S. Makrogiannis (DSU)

Parameter Estimation

October 22, 2015 10 / 28

< ロト < 同ト < ヨト < ヨト

Multi-dimensional Case

• For a projection to a d'-dimensional space the hyperplane is defined as:

$$m{x} = m{m} + \sum_{i=1}^{d'} lpha_i m{e}_i, ext{ where } d' < d'$$

• The criterion function is:

$$J_{d'} = \sum_{k=1}^{n} \left\| (\boldsymbol{m} + \sum_{i=1}^{d'} \alpha_{ki} \boldsymbol{e}_i) - \boldsymbol{x}_k \right\|^2$$

- $J_{d'}$ is minimized when e_i for the eigenvectors of scatter matrix $S_{d'}$ with the d' largest eigenvalues
- Because S_d is real and symmetric the eigenvectors are orthogonal
- Coefficients α_i are called principal components

S. Makrogiannis (DSU)

Parameter Estimation

October 22, 2015 11 / 28

I > <
I >
I

Discriminant Analysis

- PCA finds optimal data representations in the least square sense, however this does not imply that the transformed features will produce increased class separability
- On the other hand discriminant analysis techniques look for directions that distinguish between classes

S. Makrogiannis (DSU)

Parameter Estimation

October 22, 2015 12 / 28

Problem Definition

- Let's consider the problem of projecting data from *d* dimensions onto a line
- Let x_1, \ldots, x_n be the set of *n* points in a *d* dimensional space divided into subsets \mathcal{D}_i belonging to categories ω_i with cardinalities n_i , where i = 1, 2.
- Then the projections on to the direction determined by $oldsymbol{w}$ with $|oldsymbol{w}|=1$ are

$$y = w^T x$$

- The projections produce a set of n samples y_i with i = 1,..., n divided into subsets *Y*₁ and *Y*₂
- Our problem is to find the direction of \boldsymbol{w} that will maximize the separation between the projected points in \mathscr{Y}_1 and \mathscr{Y}_2

S. Makrogiannis (DSU)

3

< ロト < 同ト < ヨト < ヨト

Class Separability



Figure: Projection of data onto different directions defined by w. Observe that projection displayed in the right figure produces greater separability than the projection displayed in the left figure

S. Makrogiannis (DSU)

Parameter Estimation

October 22, 2015 14 / 28

Criterion Function

• Fisher Linear Discriminant seeks maximization of J(w) defined as

$$J(\boldsymbol{w}) = \frac{|m_{y1} - m_{y2}|^2}{s_{y1}^2 + s_{y2}^2}$$

where

 $m_{\mathrm{y}i}$ is the sample mean of ω_i in the projected space:

$$m_{yi} = (1/n_i) \sum_{y \in \mathscr{Y}_i} y = (1/n_i) \sum_{x \in \mathscr{D}_i} \boldsymbol{w}^T \boldsymbol{x} = \boldsymbol{w}^T (1/n_i) \sum_{x \in \mathscr{D}_i} \boldsymbol{x} = \boldsymbol{w}^T m_{xi}$$

 s_{vi}^2 is the scatter for projected samples of ω_i :

$$s_{yi}^2 = \sum_{y \in \mathscr{Y}_i} (y - m_{yi})^2$$

S. Makrogiannis (DSU)

Parameter Estimation

▶ < ≣ > < ≣ > = つへの October 22, 2015 15 / 28

< A

Scatter Matrices

Further we define

- Scatter matrices: $S_i = \sum_{x \in \mathscr{D}_i} (x m_{xi}) (x m_{xi})^T$, i = 1, 2
- Within-class scatter matrix: $S_W = S_1 + S_2$

Because

$$s_{yi}^{2} = \sum_{y \in \mathscr{Y}_{i}} (y - m_{yi})^{2} = \sum_{y \in \mathscr{Y}_{i}} (\boldsymbol{w}^{T} \boldsymbol{x} - \boldsymbol{w}^{T} \boldsymbol{m}_{xi})^{2}$$
$$= \boldsymbol{w}^{T} \sum_{y \in \mathscr{Y}_{i}} \left[(\boldsymbol{x} - \boldsymbol{m}_{xi}) (\boldsymbol{x} - \boldsymbol{m}_{xi})^{T} \right] \boldsymbol{w} = \boldsymbol{w}^{T} S_{i} \boldsymbol{w},$$
$$s_{y1}^{2} + s_{y2}^{2} = \boldsymbol{w}^{T} S_{W} \boldsymbol{w}$$

S. Makrogiannis (DSU)

Parameter Estimation

October 22, 2015 16 / 28

イロト (母) (ヨ) (ヨ) (ヨ) () ()

Between-class Scatter Matrix

• Consider the numerator of J(w):

$$|m_{y1} - m_{y2}|^{2} = (m_{y1} - m_{y2})^{2} = (w^{T} m_{x1} - w^{T} m_{x2})^{2}$$
$$= w^{T} (m_{x1} - m_{x2})^{2} = w^{T} (m_{x1} - m_{x2}) (m_{x1} - m_{x2})^{T} w$$
$$= w^{T} S_{B} w,$$

where S_B is the between-class scatter matrix

$$S_B = (\boldsymbol{m}_{x1} - \boldsymbol{m}_{x2})(\boldsymbol{m}_{x1} - \boldsymbol{m}_{x2})^T$$

S. Makrogiannis (DSU)

Parameter Estimation

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Scatter Matrices Notes

Within-class scatter matrix S_W

- Proportional to the sample covariance matrix
- Symmetric and positive-semidefinite
- Nonsingular if n > d

Between-class scatter matrix S_B

- Outer product of two vectors
- Symmetric and positive-semidefinite
- Its rank is at most 1

S. Makrogiannis (DSU)

Parameter Estimation

October 22, 2015 18 / 28

Optimizing the Criterion Function

 We use the scatter matrix definitions it follow that the criterion function is:

$$J(\boldsymbol{w}) = \frac{\boldsymbol{w}^T S_B \boldsymbol{w}}{\boldsymbol{w}^T S_W \boldsymbol{w}}$$

- This is a Rayleigh quotient
- The \boldsymbol{w} that maximizes $J(\boldsymbol{w})$ must satisfy $S_B \boldsymbol{w} = \lambda S_W \boldsymbol{w}$ (generalized eigenvalue problem)
- If S_W is nonsingular, we have the conventional eigenvalue problem

$$S_W^{-1}S_B$$
 w $=\lambda$ w

S. Makrogiannis (DSU)

Parameter Estimation

October 22, 2015 19 / 28

Optimizing the Criterion Function

We do not need to solve

$$S_W^{-1}S_B$$
 w $=\lambda$ w

- Recall that $S_B oldsymbol{w}$ is at the direction of $oldsymbol{m}_1 oldsymbol{m}_2$
- Hence the solution is:

$$\boldsymbol{w} = S_W^{-1}(\boldsymbol{m}_{x1} - \boldsymbol{m}_{x2})$$

• After the projection, we make a decision in the unidimensional space

S. Makrogiannis (DSU)

Parameter Estimation

October 22, 2015 20 / 28

Classification Rule

• Assuming multivariate normal class-conditional densities $p(\mathbf{x}|\omega_i)$ with equal covariance matrices Σ , we recall from Ch. 2 that at the decision boundary

$$\boldsymbol{w}^{T}\boldsymbol{x}+\boldsymbol{w}_{0}=0,$$

where

$$oldsymbol{w} = \Sigma^{-1} (oldsymbol{\mu}_1 - oldsymbol{\mu}_2)$$

- When we use the sample means and sample covariance matrix it follows that w is the one that maximizes the Fisher linear discriminant
- In this case, to classify we apply a threshold to Fisher's linear discriminant

S. Makrogiannis (DSU)

Parameter Estimation

October 22, 2015 21 / 28

イロト 不得下 イヨト イヨト

Main Concept

- Expectation-Maximization can be used to learn distribution parameters • $\boldsymbol{\theta}$ when some features are missing
- EM iteratively estimates the likelihood given the data that we have •
- Suppose a training set $\mathscr{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$
- Let $\mathbf{x}_k = \{\mathbf{x}_{kg}, \mathbf{x}_{kb}\}$ be a sample consisting of good and bad features
- Let \mathscr{D}_{σ} be the set of good features and \mathscr{D}_{b} be the set of bad features •
- Likelihood function $I(\cdot)$: $I(\mathscr{D}|\boldsymbol{\theta}) = \ln p(\mathscr{D}_{\boldsymbol{g}}, \mathscr{D}_{\boldsymbol{b}}|\boldsymbol{\theta})$ •

• We define the function $Q(\boldsymbol{\theta}; \boldsymbol{\theta}^{i})$:

$$Q(\boldsymbol{\theta};\boldsymbol{\theta}^{i}) = \mathscr{E}_{\mathscr{D}_{b}}\left[I(\mathscr{D}|\boldsymbol{\theta})|\mathscr{D}_{g};\boldsymbol{\theta}^{i}\right]$$

$$=\mathscr{E}_{\mathscr{D}_{b}}\left[\ln\left[p(\mathscr{D}_{g},\mathscr{D}_{b}|\boldsymbol{\theta})\right]p(\mathscr{D}_{b}|\mathscr{D}_{g},\boldsymbol{\theta}^{i})\right]$$

S. Makrogiannis (DSU)

Parameter Estimation

Э

Main Concept

• In $Q(\boldsymbol{\theta}; \boldsymbol{\theta}^i)$:

$$Q(\boldsymbol{\theta}; \boldsymbol{\theta}^{i}) = \mathscr{E}_{\mathscr{D}_{b}}\left[p(\mathscr{D}_{b}|\mathscr{D}_{g}, \boldsymbol{\theta}^{i}) \ln p(\mathscr{D}_{g}, \mathscr{D}_{b}|\boldsymbol{\theta})\right]$$

 $\boldsymbol{\theta}^i$: current best estimate for parameter vector $\boldsymbol{\theta}$: candidate for improved estimate

- In EM we calculate the likelihood of data for different $\boldsymbol{\theta}$ s and select the candidate that maximizes $Q(\boldsymbol{\theta}; \boldsymbol{\theta}^i)$, then set the best candidate to $\boldsymbol{\theta}^{i+1}$
- The procedure is continued till convergence

S. Makrogiannis (DSU)

Parameter Estimation

October 22, 2015 23 / 28

Algorithm

Algorithm 1 (Expectation-Maximization)

<u>1</u>	$\underline{\text{pegin initialize}} \theta^0, T, i = 0$
2	$\underline{do} i \leftarrow i + 1$
3	E step : compute $Q(\theta; \theta^i)$
5	M step : $\theta^{i+1} \leftarrow \arg \max_{\theta} Q(\theta; \theta^i)$
6	<u>until</u> $Q(\theta^{i+1}; \theta^i) - Q(\theta^i; \theta^{i-1}) \le T$
7	$\underline{\operatorname{return}} \ \hat{\theta} \leftarrow \theta^{i+1}$
8	end

Figure: Main steps of EM technique



Figure: Starting from an initial estimate $\boldsymbol{\theta}^{0}$, EM find optimal $\boldsymbol{\theta}^{1}$ in M step. Then we hold $\boldsymbol{\theta}^{1}$ constant and find value $\boldsymbol{\theta}^{2}$ that optimizes $Q(\boldsymbol{\theta}; \boldsymbol{\theta}^{1})$. The procedure continues till convergence. It is different from gradient search

イロト イポト イヨト イヨト

Parameter Estimation

October 22, 2015 24 / 28

EM for Gaussian Mixtures

- In several problems we may need to build a richer class of density models than a single Gaussian
- Then we can use a Gaussian mixture model. This model is a linear superposition of Gaussian components: $p(\mathbf{x}) = \sum_{k=1}^{K} \pi_k N(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$, where $0 \le p_k \le 1$ and $\sum_{k=1}^{K} \pi_k = 1$
- We define a K-dimensional random variable $m{z}$, such that $p(z_k=1)=\pi_k$
- Another important quantity is $p(z_k = 1 | \mathbf{x})$, which we call responsibility and denote by $\gamma(z_k)$
- A generalized version of EM is frequently used to estimate parameters of a Gaussian mixture model

S. Makrogiannis (DSU)

Parameter Estimation

October 22, 2015 25 / 28

▲ロト ▲母 ▶ ▲ヨ ▶ ▲ヨ ▶ ヨ ■ つのの

EM for Gaussian Mixtures - Algorithm

- Initialize means µ_k, covariances Σ_k and mixing coefficients π_k, and evaluate the initial value of log-likelihood
- ② E step Evaluate the responsibilities using current parameter values

$$\gamma(z_{nk}) = \frac{\pi_k N(\boldsymbol{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j N(\boldsymbol{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}$$

③ M step Re-estimate the parameters using current responsibilities

$$\boldsymbol{\mu}_{k}^{new} = (1/N_k) \sum_{n=1}^{N} \gamma(z_{nk}) \boldsymbol{x}_n$$

$$\begin{split} \boldsymbol{\Sigma}_{k}^{new} = (1/N_k) \sum_{n=1}^{N} \boldsymbol{\gamma}(z_{nk}) \left(\mathbf{x}_n - \boldsymbol{\mu}_{k}^{new} \right) \left(\mathbf{x}_n - \boldsymbol{\mu}_{k}^{new} \right)^T \\ \pi_{k}^{new} = N_k/N \text{ where, } N_k = \sum_{n=1}^{N} \boldsymbol{\gamma}(z_{nk}) \end{split}$$

④ Evaluate log likelihood

$$\ln p(\boldsymbol{X}|\boldsymbol{\mu},\boldsymbol{\Sigma},\boldsymbol{\pi}) = \sum_{n=1}^{N} \left\{ \sum_{k=1}^{K} \pi_k N(\boldsymbol{x}_n|\boldsymbol{\mu}_k,\boldsymbol{\Sigma}_k) \right\}$$

S Finish if convergence was reached, otherwise go to step 2.

S. Makrogiannis (DSU)

Parameter Estimation

n=1

October 22, 2015 26 / 28

Expectation-Maximization (EM)

Example: EM for Gaussian Mixtures



Figure: EM algorithm convergence for a mixture of 3 Gaussian distributions.

S. Makrogiannis (DSU)

Parameter Estimation

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Notes on EM

- EM is very useful when optimization of $Q(\cdot; \cdot)$ is simpler than optimization of ML function $I(\cdot)$
- In EM it is guaranteed that the log likelihood of data will increase monotonically

S. Makrogiannis (DSU)

Parameter Estimation

October 22, 2015 28 / 28

イロト 不得下 イヨト イヨト

Nonparametric Techniques MTSC 852, Fall 2015

Sokratis Makrogiannis, Ph.D.

Department of Mathematical Sciences Delaware State University smakrogiannis@desu.edu

October 30, 2015

S. Makrogiannis (DSU)

Nonparametric Techniques

October 30, 2015 1 / 33

イロト (母) (ヨ) (ヨ) (ヨ) () ()

Outline



Parzen Windows 2



3 Probabilistic Neural Networks (PNN)

S. Makrogiannis (DSU)

Nonparametric Techniques

October 30, 2015 2 / 33

990

(日) (四) (코) (코) (코)

Introduction

- In the previous chapter we estimated the class-conditional densities assuming parametric forms
- However unimodal parametric forms or products of functions may not approximate closely the underlying density
- It may be advantageous to use nonparametric techniques that make no assumptions about the forms of the underlying densities
- Here we discuss two nonparametric techniques:
 - (1) estimating class-conditional densities $p(x|\omega_j)$ from sample patterns, for example using Parzen kernels
 - 2 estimating directly the posterior probabilities $P(\omega_j | \mathbf{x})$. One example is nearest neighbor classification

S. Makrogiannis (DSU)

Nonparametric Techniques

October 30, 2015 3 / 33

< ロト < 同ト < ヨト < ヨト

Density Estimation

• Probability P for vector x in region \mathscr{R} :

$$P = \int_{\mathscr{R}} p(\mathbf{x}') d\mathbf{x}'$$

- Let x_1, x_2, \ldots, x_n , be *n* i.i.d. samples following p(x).
- The probability for k of them to be inside \mathscr{R} is:

$$P_k = \binom{n}{k} P^k (1-P)^{n-k}$$

• The expected value and variance are:

$$\mathscr{E}[k] = nP, \quad Var[k] = nP(1-P)$$

 $\mathscr{E}[k/n] = P, \quad Var[k/n] = \frac{P(1-P)}{n}$

Density Estimation

- Binomial distribution peaks sharply about the mean. In that case $P \simeq k/n$, especially for very large n.
- Assuming that x does not change much in \mathcal{R} , it follows that:

$$P=\int_{\mathscr{R}}p(\mathbf{x'})d\mathbf{x'}\simeq p(\mathbf{x})V,$$

where

- V: volume of \mathscr{R}
- \boldsymbol{x} : point in \mathcal{R}
- We combine previous results to arrive at:

$$p(\mathbf{x}) \simeq \frac{k/n}{V}$$

S. Makrogiannis (DSU)

Nonparametric Techniques

October 30, 2015 5 / 33

< 17 ▶

Density Estimation



Figure: Binomial distribution convergence. The relative probability P_k peaks more sharply at the true probability as n increases. For $n \rightarrow \infty$ the estimate will yield the true probability P = 0.7

S. Makrogiannis (DSU)

Nonparametric Techniques

< □ ト < □ ト < □ ト < 三 ト < 三 ト 三 の Q () October 30, 2015 6 / 33
Density Estimation Considerations

 If we fix V and increase n, then P estimate becomes more accurate, but p(x) estimate becomes an average over V:

$$\frac{P}{V} = \frac{\int_{\mathscr{R}} p(\mathbf{x'}) d\mathbf{x'}}{\int_{\mathscr{R}} d\mathbf{x'}}$$

- If we decrease V for fixed n to estimate p(x), then number of samples k→0, hence our estimate becomes useless
- So, we cannot reduce V very much in practice

S. Makrogiannis (DSU)

Nonparametric Techniques

October 30, 2015 7 / 33

Density Estimation Considerations

- In theory, assume that we have an unlimited number of samples
- We define \mathscr{R}_k regions, where k = 1, ... with *n* increasing with *k*. Let V_n , k_n , $p_n(\mathbf{x})$, be the volume of \mathscr{R}_n , the number of samples in \mathscr{R}_n , and the density estimate in \mathscr{R}_n respectively. It then follows that:

$$p_n(\mathbf{x}) = \frac{k_n/n}{V_n}$$

• To reach an estimate of p(x) the following must be satisfied:

$$\lim_{n\to\infty} V_n = 0, \quad \lim_{n\to\infty} k_n = \infty, \quad \lim_{n\to\infty} k_n/n = 0$$

• These conditions ensure that we can estimate p(x) accurately, that the frequency ratio will converge to P, and that $p_n(x)$ will converge in general

S. Makrogiannis (DSU)

Density Estimation Approaches

We can attempt to satisfy these conditions by the following approaches:

- **1** Shrink \mathscr{R}_n by defining V_n as a function of n
 - Determine k from the data
 - We must ensure that $p_n(x)$ converges to p(x)
 - Parzen windows use this principle
- 2 Define k_n as a function of n
 - Then V_n will have to grow till it encloses k_n neighbors of x
 - This is the k_n nearest neighbor estimation

S. Makrogiannis (DSU)

Nonparametric Techniques

October 30, 2015 9 / 33

Density Estimation Approaches



Figure: Two density estimation approaches: reduce V_n to estimate $p_n(x)$ more accurately (top), or increase V_n to increase k_n to estimate P_k more accurately (bottom)

S. Makrogiannis (DSU)

Nonparametric Techniques

14 I.S. October 30, 2015 10 / 33

Sar

< □ > < 同 >

- The Parzen window method defines a window that may be a function of the number of data points
- More specifically, \mathscr{R}_n is a *d*-dimensional hypercube
- The volume of the hypecube is:

$$V_n = h_n^d$$

where h_n : edge length of cube

• To yield the number of points in \mathscr{R}_n denoted by k_n we use a window function:

$$\phi(oldsymbol{u}) = egin{cases} 1 & |u_j| \leq 1/2 & j=1,\ldots,d \ 0 & ext{otherwise} \end{cases}$$

S. Makrogiannis (DSU)

Nonparametric Techniques

October 30, 2015 11 / 33

• We defined the window function as:

$$\phi(oldsymbol{u}) = egin{cases} 1 & |u_j| \leq 1/2 & j=1,\ldots,d \ 0 & ext{otherwise} \end{cases}$$

• Then, the number of points inside the hypercube centered at \boldsymbol{x} is given by:

$$k_n = \sum_{i=1}^n \phi\left(\frac{\mathbf{x} - \mathbf{x}_i}{h_n}\right)$$

• From density estimation we have that:

$$p_n(\boldsymbol{x}) = \frac{k_n/n}{V_n}$$

• By substitution it follows that: $p_n(x) = (1/n) \sum_{i=1}^n \frac{1}{V_n} \phi\left(\frac{x-x_i}{h_n}\right)$

S. Makrogiannis (DSU)

Nonparametric Techniques

October 30, 2015 12 / 33

• Before, we arrived at:

$$p_n(\mathbf{x}) = (1/n) \sum_{i=1}^n \frac{1}{V_n} \phi\left(\frac{\mathbf{x} - \mathbf{x}_i}{h_n}\right)$$

- This is a general form for estimating density
- This operation can also be considered as interpolation
- We can choose different function types for $\phi(\cdot)$ to handle discontinuities that may be caused by the original window function

S. Makrogiannis (DSU)

Nonparametric Techniques

October 30, 2015 13 / 33

< □ > < 同 >

- Next, we require that $p_n(x)$ be a density function
- Hence, we need $\phi(x)$ to be a density function with conditions

 $\phi(\mathbf{x}) \geq 0$

and

$$\int \phi(\boldsymbol{u}) d\boldsymbol{u} = 1$$

• If in addition $V_n = h_n^d$, then $p_n(x)$ is a density function too

S. Makrogiannis (DSU)

Nonparametric Techniques

October 30, 2015 14 / 33

▲ロト ▲掃ト ▲ヨト ▲ヨト ニヨー のへで

- From previous analysis we observed that the size of $V_n = h_n^d$ can affect density estimation
- Now we focus our interest on h_n
- First, we define $\delta_n(x)$ as: $\delta_n(x) = (1/n) rac{1}{V_n} \phi\left(rac{x}{h_n}
 ight)$
- $p_n(\mathbf{x})$ becomes: $p_n(\mathbf{x}) = (1/n) \sum_{i=1}^n \delta_n(\mathbf{x} \mathbf{x}_i)$

S. Makrogiannis (DSU)

Nonparametric Techniques

▶ < ≧ ▶ < ≧ ▶ Ξ ∽ Q (? October 30, 2015 15 / 33

- If h_n is too large, then
 - $\delta_n(\mathbf{x} \mathbf{x}_i)$ becomes a slowly varying function with a low peak
 - $p_n(x)$ becomes a smooth and "out-of-focus" estimate of p(x)
- If h_n is too small, then
 - $\delta_n(\mathbf{x}-\mathbf{x}_i)$ has a sharper peak and changes faster with distance
 - $p_n(x)$ becomes a sum of sharp kernels centered at the training samples affected by noisy samples
- We also note that:

$$\int \delta_n(\mathbf{x} - \mathbf{x}_i) d\mathbf{x} = \int (1/n) \frac{1}{V_n} \phi\left(\frac{\mathbf{x}}{h_n}\right) d\mathbf{x} = \int \phi(\mathbf{u}) d\mathbf{u} = 1$$

• Hence, the distribution is normalized

S. Makrogiannis (DSU)

Nonparametric Techniques

October 30, 2015 16 / 33

イロト イポト イヨト イヨト 二日

- If h_n is too large, then the density estimate will be less accurate
- If *h_n* is too small, then he density estimate will be sensitive to statistical variability
- In real world problems, we must choose the kernel size as a trade-off between the two weaknesses
- In an ideal world with unlimited training samples we could reduce the kernel size and still get an accurate density $p_n(x) = p(x)$

S. Makrogiannis (DSU)

Nonparametric Techniques

October 30, 2015 17 / 33

イロト 不得下 イヨト イヨト

Convergence Definition

- We consider p_n(x) to be a random variable because it is a function of n random samples x_i, i = 1,...,n
- We denote the mean and variance of $p_n(x)$ by $\bar{p}_n(x)$ and $\sigma_n^2(x)$
- Then $p_n(x)$ will converge to p(x), if:

$$\lim_{n\to\infty}\bar{p}_n(\boldsymbol{x})=p(\boldsymbol{x})$$

and

$$\lim_{n\to\infty}\sigma_n^2(\boldsymbol{x})=0$$

S. Makrogiannis (DSU)

Nonparametric Techniques

October 30, 2015 18 / 33

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Convergence Conditions

- To prove convergence, we place conditions on $p({m x}), \ \phi({m u}),$ and h_n
- Previous conditions are that $p(\cdot)$ is continuous at \pmb{x} , and that $\phi(\cdot)$ is a density function
- We will show that the following conditions have to be met as well:

$$\sup_{\boldsymbol{u}} \phi(\boldsymbol{u}) < \infty, \quad \lim_{\|\boldsymbol{u}\| \to \infty} \phi(\boldsymbol{u}) \prod_{i=1}^{d} u_i = 0$$

$$\lim_{n\to\infty}V_n=0,\quad \lim_{n\to\infty}nV_n=\infty$$

• These conditions ensure that $\phi(\boldsymbol{u})$ is bounded and that V_n must approach zero but at a slower rate that 1/n

S. Makrogiannis (DSU)

Nonparametric Techniques

< □ ト < □ ト < 三 ト < 三 ト < 三 ト 三 のへで October 30, 2015 19 / 33

Convergence of the Mean

• First, we compute $\bar{p}_n(x)$:

$$\bar{p}_n(\mathbf{x}) = \mathscr{E}[p_n(\mathbf{x})] = (1/n) \sum_{i=1}^n \mathscr{E}\left[\frac{1}{V_n} \phi\left(\frac{\mathbf{x} - \mathbf{x}_i}{h_n}\right)\right]$$
$$= \int \frac{1}{V_n} \phi\left(\frac{\mathbf{x} - \mathbf{v}}{h_n}\right) p(\mathbf{v}) d\mathbf{v} = \int \delta_n(\mathbf{x} - \mathbf{v}) p(\mathbf{v}) d\mathbf{v}$$

- We observe that $\bar{p}_n(x)$ is the result of convolving the window function with the density function, therefore it is a blurred version of the unknown density
- However, when $V_n
 ightarrow 0$, $\delta_n({m x}-{m v})$ becomes a delta function at ${m x}$.
- Hence, if $m{p}$ is continuous at $m{x}$, then $\lim_{n o \infty} ar{p}_n(m{x}) = p(m{x})$

S. Makrogiannis (DSU)

イロト イポト イヨト イヨト 二日

Convergence of the Variance

• To avoid a noisy estimate, we need to ensure the convergence of variance:

$$\sigma_n^2(\mathbf{x}) = \sum_{i=1}^n \mathscr{E}\left[\left(\frac{1}{nV_n}\phi\left(\frac{\mathbf{x}-\mathbf{x}_i}{h_n}\right) - (1/n)\bar{p}_n(\mathbf{x})\right)^2\right]$$
$$= n\mathscr{E}\left[\frac{1}{n^2V_n^2}\phi^2\left(\frac{\mathbf{x}-\mathbf{x}_i}{h_n}\right)\right] - (1/n)\bar{p}_n^2(\mathbf{x})$$
$$= \frac{1}{nV_n}\int \frac{1}{V_n}\phi^2\left(\frac{\mathbf{x}-\mathbf{v}}{h_n}\right)p(\mathbf{v})d\mathbf{v} - (1/n)\bar{p}_n^2(\mathbf{x})$$

• If we drop the second term, bound ϕ , and use above expression for $\bar{p}_n(x)$ it follows that:

$$\sigma_n^2(\mathbf{x}) \leq \frac{\sup\left(\phi(\cdot)\right)\bar{p}_n(\mathbf{x})}{nV_n}$$

S. Makrogiannis (DSU)

Nonparametric Techniques

<ロト < @ ト < 注 ト く 注 ト 注 の Q C October 30, 2015 21 / 33

Convergence of the Variance

• Previously we arrived at:

$$\sigma_n^2(\mathbf{x}) \leq \frac{\sup\left(\phi(\cdot)\right)\bar{p}_n(\mathbf{x})}{nV_n}$$

- We need to start with a large V_n to reach a small variance
- Still, because the numerator is finite w.r.t *n*, we can still let $V_n \rightarrow 0$ as long as $nV_n \rightarrow \infty$ to obtain zero variance
- This means that V_n can be V_1/\sqrt{n} , for example

S. Makrogiannis (DSU)

Nonparametric Techniques

October 30, 2015 22 / 33

イロト イポト イヨト イヨト 二日

Gaussian Kernel Example

- Suppose that the true density p(x) is univariate normal, with zero mean, and unit variance
- Suppose we use a Gaussian kernel for density estimation given by:

$$\phi(u)=\frac{1}{\sqrt{2\pi}}e^{-u^2/2}$$

• The density estimate at x is:

$$p_n(x) = (1/n) \sum_{i=1}^n \frac{1}{h_n} \phi\left(\frac{x-x_i}{h_n}\right)$$

- where $h_n = h_1/\sqrt{n}$
- S. Makrogiannis (DSU)

Figure: Parzen kernel density estimation for a univariate normal distribution versus the number of samples and window width. The contribution of each point to the density is more visible for smaller window widths. Larger *n* improves density estimation

Nonparametric Techniques

October 30, 2015 23 / 33

Gaussian Kernel Example



Figure: Parzen kernel density estimation for a bivariate normal distribution versus the number of samples and window width. Smaller window width produces "noisier" estimates for fixed n

S. Makrogiannis (DSU)

Nonparametric Techniques

October 30, 2015 24 / 33

Gaussian Kernel Example



Figure: Parzen kernel density estimation for a mixture of a uniform and a triangular distribution. Observe that more samples improve estimation

S. Makrogiannis (DSU)

Nonparametric Techniques

◆□ → < □ → < 三 → < 三 → < 三 → へ ○ October 30, 2015 25 / 33

Parzen Kernel-based Classification

- In Parzen kernel-based classification we estimate the class-conditional density at each test point and make a decision using Maximum-a-Posteriori rule
- In this classifier we can reduce the training error as much as we wish, but we may cause overfitting
- Gaussian windows are reasonable choices but it may take some experimentation to find the window size

S. Makrogiannis (DSU)

Nonparametric Techniques

October 30, 2015 26 / 33

Parzen Kernel-based Classification

Strengths:

 * no assumption about the density functions – we use the same procedure

Weaknesses:

* requirement for a large number of samples

* computationally demanding because all training samples are used to estimate densities each time

* when the dimensionality grows the demand for a large number of samples grows exponentially (curse of dimensionality)

S. Makrogiannis (DSU)

Nonparametric Techniques

October 30, 2015 27 / 33

Parzen Kernel-based Classification



Figure: Decision boundaries versus window width *h*. Smaller *h* produces more complicated decision boundaries (left column) than larger *h* (right column)

S. Makrogiannis (DSU)

Nonparametric Techniques

October 30, 2015 28 / 33

Sar

Probabilistic Neural Network (PNN) Example

- We will show how Parzen density-based classifier can be implemented as a neural network using PNN
- Suppose we have *n* patterns, in a *d*-dimensional space divided into *c* classes
- PNN will have:
 - an input layer with *d* input units
 - an intermediate layer with *n* pattern units. Each input unit is connected to all *n* pattern units
 - a category layer with *c* category units. Each pattern unit is connected to only one category unit
- Connections from input to pattern units carry weights w, which need to be learned

S. Makrogiannis (DSU)

Nonparametric Techniques

October 30, 2015 29 / 33

イロト 不得下 イヨト イヨト

PNN Topology



Figure: PNN will have: an input layer with d input units; an intermediate layer with n pattern units. Each input unit is connected to all n pattern units; a category layer with c category units. Each pattern unit is connected to only one category unit

S. Makrogiannis (DSU)

Nonparametric Techniques

<□ ト < □ ト < □ ト < 三 ト < 三 ト < 三 つ Q () October 30, 2015 30 / 33

PNN Training

- For each train pattern x, we normalize the weight s.t. $\sum_{i=1}^{d} x_i^2 = 1$
- We link the input unit and the corresponding pattern unit with a link that has a weight $w_j = x_j, j = 1, ..., n$
- We denote the components of the jth pattern by $x_{jk}, k = 1, \dots, d$

$$\frac{1}{2} \frac{\text{begin initialize } j = 0, n = \# \text{patterns}}{\frac{1}{2} \frac{1}{2} \frac{1$$

Figure: PNN training algorithm

S. Makrogiannis (DSU)

Nonparametric Techniques

October 30, 2015 31 / 33

イロト イポト イヨト イヨト 二日

PNN Classification

- We normalize the test pattern x
- We compute the net activation by inner product: $net_k = \boldsymbol{w}_k^T \boldsymbol{x}$
- Each output unit sums the contributions from all pattern units using the nonlinear function $e^{\frac{net_k-1}{\sigma^2}}$, σ : user-defined parameter
- We use the window function to define activation function:

$$\phi(\frac{\mathbf{x}-\mathbf{w}_k}{h_n}) \sim e^{-(\mathbf{x}-\mathbf{w}_k)^T(\mathbf{x}-\mathbf{w}_k)/2\sigma^2} = e^{\frac{net_k-1}{\sigma^2}}$$

- Each pattern unit contributes to its associated category unit a signal equal to type of probability of association with a training point
- The sum of estimates gives the discriminant function $g_i(x)$. This is a Parzen window estimate for the distribution
- We apply a maximum likelihood rule to assign the pattern to a class

S. Makrogiannis (DSU)

イロト 不良下 イヨト イヨト

PNN Testing

Algorithm 2 (PNN classification)

$$\begin{array}{c|c} i & \underline{\text{begin initialize}} \ k = 0, \mathbf{x} = \text{test pattern} \\ z & \underline{\text{do}} \ k \leftarrow k + 1 \\ z & z_k \leftarrow w_k^t \mathbf{x} \\ \downarrow & \underline{\text{if}} \ a_{kc} = 1 \underline{\text{then}} \ g_c \leftarrow g_c + \exp[(z_k - 1)/\sigma^2] \\ z & \underline{\text{until}} \ k = n \\ z & \underline{\text{ntil}} \ k = n \\ z & \underline{\text{return class}} \leftarrow \arg\max_i g_i(\mathbf{x}) \\ \gamma & \underline{\text{end}} \end{array}$$

Figure: PNN Testing algorithm

S. Makrogiannis (DSU)

Nonparametric Techniques

October 30, 2015 33 / 33

▲ロト ▲園 ト ▲ 臣 ト ▲ 臣 ト 二 臣 … の Q ()

Nonparametric Techniques MTSC 852, Fall 2015

Sokratis Makrogiannis, Ph.D.

Department of Mathematical Sciences Delaware State University smakrogiannis@desu.edu

October 29, 2015

S. Makrogiannis (DSU)

Nonparametric Techniques

October 29, 2015 1 / 25

▲ロト ▲掃ト ▲ヨト ▲ヨト ニヨー のへで







2 Probabilistic Neural Networks (PNN)

S. Makrogiannis (DSU)

Nonparametric Techniques

October 29, 2015 2 / 25

996

< ロト < 団ト < 注ト < 注ト = 注:</p>

- The Parzen window method defines a window that may be a function of the number of data points
- More specifically, \mathscr{R}_n is a *d*-dimensional hypercube
- The volume of the hypecube is:

$$V_n = h_n^d$$

where h_n : edge length of cube

• To yield the number of points in \mathscr{R}_n denoted by k_n we use a window function:

$$\phi(oldsymbol{u}) = egin{cases} 1 & |u_j| \leq 1/2 & j=1,\ldots,d \ 0 & ext{otherwise} \end{cases}$$

S. Makrogiannis (DSU)

Nonparametric Techniques

October 29, 2015 3 / 25

• We defined the window function as:

$$\phi(oldsymbol{u}) = egin{cases} 1 & |u_j| \leq 1/2 & j=1,\ldots,d \ 0 & ext{otherwise} \end{cases}$$

• Then, the number of points inside the hypercube centered at \boldsymbol{x} is given by:

$$k_n = \sum_{i=1}^n \phi\left(\frac{\mathbf{x} - \mathbf{x}_i}{h_n}\right)$$

• From density estimation we have that:

$$p_n(\boldsymbol{x}) = \frac{k_n/n}{V_n}$$

• By substitution it follows that: $p_n(x) = (1/n) \sum_{i=1}^n \frac{1}{V_n} \phi\left(\frac{x-x_i}{h_n}\right)$

S. Makrogiannis (DSU)

Nonparametric Techniques

October 29, 2015 4 / 25

• Before, we arrived at:

$$p_n(\mathbf{x}) = (1/n) \sum_{i=1}^n \frac{1}{V_n} \phi\left(\frac{\mathbf{x} - \mathbf{x}_i}{h_n}\right)$$

- This is a general form for estimating density
- This operation can also be considered as interpolation
- We can choose different function types for $\phi(\cdot)$ to handle discontinuities that may be caused by the original window function

S. Makrogiannis (DSU)

Nonparametric Techniques

October 29, 2015 5 / 25

- Next, we require that $p_n(x)$ be a density function
- Hence, we need $\phi(\mathbf{x})$ to be a density function with conditions •

 $\phi(\mathbf{x}) > 0$

and

$$\int \phi(\boldsymbol{u}) d\boldsymbol{u} = 1$$

• • If in addition $V_n = h_n^d$, then $p_n(\mathbf{x})$ is a density function too

S. Makrogiannis (DSU)

Nonparametric Techniques

3 October 29, 2015 6 / 25

< ロト < 同ト < ヨト < ヨト

- From previous analysis we observed that the size of $V_n = h_n^d$ can affect density estimation
- Now we focus our interest on h_n
- First, we define $\delta_n(x)$ as: $\delta_n(x) = (1/n) rac{1}{V_n} \phi\left(rac{x}{h_n}
 ight)$
- $p_n(\mathbf{x})$ becomes: $p_n(\mathbf{x}) = (1/n) \sum_{i=1}^n \delta_n(\mathbf{x} \mathbf{x}_i)$

S. Makrogiannis (DSU)

Nonparametric Techniques

October 29, 2015 7 / 25

- If h_n is too large, then
 - $\delta_n(\mathbf{x} \mathbf{x}_i)$ becomes a slowly varying function with a low peak
 - $p_n(x)$ becomes a smooth and "out-of-focus" estimate of p(x)
- If h_n is too small, then
 - $\delta_n(\mathbf{x}-\mathbf{x}_i)$ has a sharper peak and changes faster with distance
 - $p_n(x)$ becomes a sum of sharp kernels centered at the training samples affected by noisy samples
- We also note that:

$$\int \delta_n(\mathbf{x} - \mathbf{x}_i) d\mathbf{x} = \int (1/n) \frac{1}{V_n} \phi\left(\frac{\mathbf{x}}{h_n}\right) d\mathbf{x} = \int \phi(\mathbf{u}) d\mathbf{u} = 1$$

• Hence, the distribution is normalized

S. Makrogiannis (DSU)

Nonparametric Techniques

October 29, 2015 8 / 25

< ロト < 同ト < ヨト < ヨト

- If h_n is too large, then the density estimate will be less accurate
- If *h_n* is too small, then he density estimate will be sensitive to statistical variability
- In real world problems, we must choose the kernel size as a trade-off between the two weaknesses
- In an ideal world with unlimited training samples we could reduce the kernel size and still get an accurate density $p_n(x) = p(x)$

S. Makrogiannis (DSU)

Nonparametric Techniques

October 29, 2015 9 / 25

< ロト < 同ト < ヨト < ヨト
Convergence Definition

- We consider p_n(x) to be a random variable because it is a function of n random samples x_i, i = 1,...,n
- We denote the mean and variance of $p_n(x)$ by $\bar{p}_n(x)$ and $\sigma_n^2(x)$
- Then $p_n(x)$ will converge to p(x), if:

$$\lim_{n\to\infty}\bar{p}_n(\boldsymbol{x})=p(\boldsymbol{x})$$

and

$$\lim_{n\to\infty}\sigma_n^2(\boldsymbol{x})=0$$

S. Makrogiannis (DSU)

Nonparametric Techniques

October 29, 2015 10 / 25

イロト 不同下 イヨト イヨト ヨー ろくつ

Convergence Conditions

- To prove convergence, we place conditions on $p({m x}), \ \phi({m u}),$ and h_n
- Previous conditions are that $p(\cdot)$ is continuous at \pmb{x} , and that $\phi(\cdot)$ is a density function
- We will show that the following conditions have to be met as well:

$$\sup_{\boldsymbol{u}} \phi(\boldsymbol{u}) < \infty, \quad \lim_{\|\boldsymbol{u}\| \to \infty} \phi(\boldsymbol{u}) \prod_{i=1}^{d} u_i = 0$$

$$\lim_{n\to\infty}V_n=0,\quad \lim_{n\to\infty}nV_n=\infty$$

• These conditions ensure that $\phi(\boldsymbol{u})$ is bounded and that V_n must approach zero but at a slower rate that 1/n

S. Makrogiannis (DSU)

Nonparametric Techniques

October 29, 2015 11 / 25

Convergence of the Mean

• First, we compute $\bar{p}_n(x)$:

$$\bar{p}_n(\mathbf{x}) = \mathscr{E}[p_n(\mathbf{x})] = (1/n) \sum_{i=1}^n \mathscr{E}\left[\frac{1}{V_n} \phi\left(\frac{\mathbf{x} - \mathbf{x}_i}{h_n}\right)\right]$$
$$= \int \frac{1}{V_n} \phi\left(\frac{\mathbf{x} - \mathbf{v}}{h_n}\right) p(\mathbf{v}) d\mathbf{v} = \int \delta_n(\mathbf{x} - \mathbf{v}) p(\mathbf{v}) d\mathbf{v}$$

- We observe that $\bar{p}_n(x)$ is the result of convolving the window function with the density function, therefore it is a blurred version of the unknown density
- However, when $V_n
 ightarrow 0$, $\delta_n({m x}-{m v})$ becomes a delta function at ${m x}$.
- Hence, if $m{p}$ is continuous at $m{x}$, then $\lim_{n o \infty} ar{p}_n(m{x}) = p(m{x})$

S. Makrogiannis (DSU)

イロト イポト イヨト イヨト 二日

Convergence of the Variance

• To avoid a noisy estimate, we need to ensure the convergence of variance:

$$\sigma_n^2(\mathbf{x}) = \sum_{i=1}^n \mathscr{E}\left[\left(\frac{1}{nV_n}\phi\left(\frac{\mathbf{x}-\mathbf{x}_i}{h_n}\right) - (1/n)\bar{p}_n(\mathbf{x})\right)^2\right]$$
$$= n\mathscr{E}\left[\frac{1}{n^2V_n^2}\phi^2\left(\frac{\mathbf{x}-\mathbf{x}_i}{h_n}\right)\right] - (1/n)\bar{p}_n^2(\mathbf{x})$$
$$= \frac{1}{nV_n}\int \frac{1}{V_n}\phi^2\left(\frac{\mathbf{x}-\mathbf{v}}{h_n}\right)p(\mathbf{v})d\mathbf{v} - (1/n)\bar{p}_n^2(\mathbf{x})$$

• If we drop the second term, bound ϕ , and use above expression for $\bar{p}_n(x)$ it follows that:

$$\sigma_n^2(\mathbf{x}) \leq \frac{\sup\left(\phi(\cdot)\right)\bar{p}_n(\mathbf{x})}{nV_n}$$

S. Makrogiannis (DSU)

Nonparametric Techniques

◆□ ト < □ ト < □ ト < 三 ト < 三 ト < 三 ・ つ Q (C) October 29, 2015 13 / 25

Convergence of the Variance

• Previously we arrived at:

$$\sigma_n^2(\mathbf{x}) \leq \frac{\sup\left(\phi(\cdot)\right)\bar{p}_n(\mathbf{x})}{nV_n}$$

- We need to start with a large V_n to reach a small variance
- Still, because the numerator is finite w.r.t *n*, we can still let $V_n \rightarrow 0$ as long as $nV_n \rightarrow \infty$ to obtain zero variance
- This means that V_n can be V_1/\sqrt{n} , for example

S. Makrogiannis (DSU)

Nonparametric Techniques

October 29, 2015 14 / 25

Gaussian Kernel Example

- Suppose that the true density p(x) is univariate normal, with zero mean, and unit variance
- Suppose we use a Gaussian kernel for density estimation given by:

$$\phi(u)=\frac{1}{\sqrt{2\pi}}e^{-u^2/2}$$

• The density estimate at x is:

$$p_n(x) = (1/n) \sum_{i=1}^n \frac{1}{h_n} \phi\left(\frac{x-x_i}{h_n}\right)$$

- where $h_n = h_1/\sqrt{n}$
- S. Makrogiannis (DSU)

Figure: Parzen kernel density estimation for a univariate normal distribution versus the number of samples and window width. The contribution of each point to the density is more visible for smaller window widths. Larger *n* improves density estimation

Nonparametric Techniques

October 29, 2015 15 / 25

Gaussian Kernel Example



Figure: Parzen kernel density estimation for a bivariate normal distribution versus the number of samples and window width. Smaller window width produces "noisier" estimates for fixed n

S. Makrogiannis (DSU)

Nonparametric Techniques

October 29, 2015 16 / 25

Gaussian Kernel Example



Figure: Parzen kernel density estimation for a mixture of a uniform and a triangular distribution. Observe that more samples improve estimation

S. Makrogiannis (DSU)

Nonparametric Techniques

Parzen Kernel-based Classification

- In Parzen kernel-based classification we estimate the class-conditional density at each test point and make a decision using Maximum-a-Posteriori rule
- In this classifier we can reduce the training error as much as we wish, but we may cause overfitting
- Gaussian windows are reasonable choices but it may take some experimentation to find the window size

S. Makrogiannis (DSU)

Nonparametric Techniques

October 29, 2015 18 / 25

Parzen Kernel-based Classification

Strengths:

 * no assumption about the density functions – we use the same procedure

Weaknesses:

* requirement for a large number of samples

* computationally demanding because all training samples are used to estimate densities each time

* when the dimensionality grows the demand for a large number of samples grows exponentially (curse of dimensionality)

S. Makrogiannis (DSU)

Nonparametric Techniques

October 29, 2015 19 / 25

Parzen Kernel-based Classification



Figure: Decision boundaries versus window width *h*. Smaller *h* produces more complicated decision boundaries (left column) than larger *h* (right column)

S. Makrogiannis (DSU)

Nonparametric Techniques

October 29, 2015 20 / 25

Sar

Probabilistic Neural Network (PNN) Example

- We will show how Parzen density-based classifier can be implemented as a neural network using PNN
- Suppose we have *n* patterns, in a *d*-dimensional space divided into *c* classes
- PNN will have:
 - an input layer with d input units
 - an intermediate layer with *n* pattern units. Each input unit is connected to all *n* pattern units
 - a category layer with *c* category units. Each pattern unit is connected to only one category unit
- Connections from input to pattern units carry weights w, which need to be learned

S. Makrogiannis (DSU)

Nonparametric Techniques

October 29, 2015 21 / 25

イロト 不得下 イヨト イヨト

PNN Topology



Figure: PNN will have: an input layer with d input units; an intermediate layer with n pattern units. Each input unit is connected to all n pattern units; a category layer with c category units. Each pattern unit is connected to only one category unit

S. Makrogiannis (DSU)

Nonparametric Techniques

◆□ ト < □ ト < 三 ト < 三 ト < 三 つ へ () October 29, 2015 22 / 25

PNN Training

- For each train pattern x, we normalize the weight s.t. $\sum_{i=1}^{d} x_i^2 = 1$
- We link the input unit and the corresponding pattern unit with a link that has a weight $w_j = x_j, j = 1, ..., n$
- We denote the components of the jth pattern by $x_{jk}, \, k=1,\ldots,d$

$$\frac{1}{2} \frac{\text{begin initialize } j = 0, n = \# \text{patterns}}{\frac{1}{2} \frac{1}{2} \frac{1$$

Figure: PNN training algorithm

S. Makrogiannis (DSU)

Nonparametric Techniques

October 29, 2015 23 / 25

イロト 不得下 イヨト イヨト

PNN Classification

- We normalize the test pattern x
- We compute the net activation by inner product: $net_k = \boldsymbol{w}_k^T \boldsymbol{x}$
- Each output unit sums the contributions from all pattern units using the nonlinear function $e^{\frac{net_k-1}{\sigma^2}}$, σ : user-defined parameter
- We use the window function to define activation function:

$$\phi(\frac{\mathbf{x}-\mathbf{w}_k}{h_n}) \sim e^{-(\mathbf{x}-\mathbf{w}_k)^T(\mathbf{x}-\mathbf{w}_k)/2\sigma^2} = e^{\frac{net_k-1}{\sigma^2}}$$

- Each pattern unit contributes to its associated category unit a signal equal to type of probability of association with a training point
- The sum of estimates gives the discriminant function $g_i(x)$. This is a Parzen window estimate for the distribution
- We apply a maximum likelihood rule to assign the pattern to a class

S. Makrogiannis (DSU)

PNN Testing

Algorithm 2 (PNN classification)

$$\begin{array}{c|c} i & \underline{\text{begin initialize}} \ k = 0, \mathbf{x} = \text{test pattern} \\ z & \underline{\text{do}} \ k \leftarrow k + 1 \\ z & z_k \leftarrow w_k^t \mathbf{x} \\ \downarrow & \underline{\text{if}} \ a_{kc} = 1 \underline{\text{then}} \ g_c \leftarrow g_c + \exp[(z_k - 1)/\sigma^2] \\ z & \underline{\text{until}} \ k = n \\ z & \underline{\text{ntil}} \ k = n \\ z & \underline{\text{return class}} \leftarrow \arg\max_i g_i(\mathbf{x}) \\ \gamma & \underline{\text{end}} \end{array}$$

Figure: PNN Testing algorithm

S. Makrogiannis (DSU)

Nonparametric Techniques

October 29, 2015 25 / 25

▲ロト ▲園 ト ▲ 臣 ト ▲ 臣 ト 二 臣 … の Q ()

Nonparametric Techniques MTSC 852, Fall 2015

Sokratis Makrogiannis, Ph.D.

Department of Mathematical Sciences Delaware State University smakrogiannis@desu.edu

November 6, 2015

S. Makrogiannis (DSU)

Nonparametric Techniques

November 6, 2015 1 / 26

イロト (母) (ヨ) (ヨ) (ヨ) () ()



Nearest Neighbor Estimation





Metrics and Nearest Neighbor Classification

S. Makrogiannis (DSU)

Nonparametric Techniques

November 6, 2015 2 / 26

200

(日) (四) (코) (코) (코)

k_n Nearest Neighbor Estimation

- In this method, to estimate p(x) from *n* training samples we grow the region \mathcal{R}_n with volume V_n around x such that it encloses k_n samples
- The samples enclosed by V_n are the k_n nearest-neighbors of x
- We estimate density by

$$p_n(\boldsymbol{x}) = \frac{k_n/n}{V_n}$$

- We can show that the conditions $\lim_{n\to\infty} k_n = \infty$ and $\lim_{n\to\infty} k_n/n = 0$ are necessary and sufficient for $p_n(x)$ to converge to p(x) at points where p(x) is continuous
- Assume that $k_n = \sqrt{n}$. Then for a very large *n* we have that $V_n \simeq V = 1/(\sqrt{n}p(\mathbf{x}))$, following the form V_1/\sqrt{n} that we discussed before
- While p_n(x) is continuous, the gradient is not. Still, the points of discontinuity are more frequently not close to the training points

S. Makrogiannis (DSU)

Nearest Neighbor Estimation

The k-Nearest Neighbor (k-NN) Rule Error Rate



Figure: k_n density estimation for uni- and multi-variate cases

S. Makrogiannis (DSU)

Nonparametric Techniques

November 6, 2015 4 / 26

イロト イロト イヨト

k_n Nearest Neighbor (NN) and Parzen-window Estimation

- We will compare the estimates by Parzen and k_n NN
- For n=1 and $k_n=\sqrt{n}=1$ the k_n NN estimate becomes

$$p_n(x) = \frac{1}{2|x-x_1|}$$

- This is a poor estimate with integral diverging to infinity
- For larger *n*, the estimate becomes more accurate but the integral remains infinite
- In Parzen window approach we can change k_n according to $k_n = k_1 \sqrt{n}$ by varying k_1

S. Makrogiannis (DSU)

Nonparametric Techniques

November 6, 2015 5 / 26

イロト イポト イヨト イヨト

Nearest Neighbor Estimation

k_n Nearest Neighbor (NN) and Parzen-window Estimation



Figure: k_n density estimation for variable n, k_n . Larger k_n improves density estimates

S. Makrogiannis (DSU)

Nonparametric Techniques

November 6, 2015 6 / 26

Э

A Posteriori Probability Estimation

- Suppose that we have *n* training points. n_i of them that belong to category ω_i and that we want to estimate the posterior probability $P(\omega_i|\mathbf{x})$ at a point \mathbf{x}
- We can use the Bayesian decision rule for this purpose. We place a volume V around x. Let k be the total number of samples inside volume V, and k_i be the number of samples inside V that belong to ω_i
- Then the likelihood is: $p(\boldsymbol{x}|\boldsymbol{\omega}_i) = rac{k_i/n_i}{V}$

• Unconditional density:
$$p(x) = rac{k/n}{V}$$

- Priors: $p(\omega_i) = n_i/n$
- By use of Bayes rule:

$$P_n(\omega_i|\mathbf{x}) = \frac{p(\mathbf{x}|\omega_i)p(\omega_i)}{p(\mathbf{x})} = \frac{\frac{k_i/n_i}{V}\frac{n_i}{n}}{\frac{k/n}{V}} = \frac{k_i}{k}$$

S. Makrogiannis (DSU)

The Nearest Neighbor Rule

- Let Dⁿ be a set of training points, or prototypes, Dⁿ = {x₁,...,x_n} and let x be a test point that is closest to the training point x' ∈ Dⁿ
- The nearest neighbor rule will classify $m{x}$ to the class of $m{x}'$
- This is a suboptimal procedure; it yields an error rate that is greater than the Bayes rate

S. Makrogiannis (DSU)

Nonparametric Techniques

November 6, 2015 8 / 26

イロト イポト イヨト イヨト

The Nearest Neighbor Rule

- We consider the prototype labels to be random variables with probabilities equal to posteriors $P(\omega_i|x')$
- Assuming that **x** and **x'** are sufficiently close, it follows that $P(\omega_i | \mathbf{x}) \simeq P(\omega_i | \mathbf{x}')$
- Then the category ω_m of test point $m{x}$ is found by:

$$\omega_m = \arg\max_i P(\omega_i | \mathbf{x})$$

- This rule will partition the feature space into regions defined by a neighbor similarity measure
- This is called Voronoi tesselation



Figure: Voronoi Tesselation using NN rule

S. Makrogiannis (DSU)

Nonparametric Techniques

November 6, 2015 9 / 26

The k-Nearest Neighbor (k-NN) Rule

- This rule classifies a point *x* by examining the labels of the *k* nearest neighbors and applying majority voting
- We consider the 2-class problem with k being odd
- Labels ω_i in each of the *k* neighbors are random variables with probabilities $P(\omega_i | \mathbf{x}), i = 1, 2$
- k-NN rule selects class ω_m , if a majority of k nearest neighbors are labeled ω_m . This event has a probability of

$$\sum_{i=(k+1)/2}^{k} \binom{k}{i} P(\omega_m | \mathbf{x})^{i} [1 - P(\omega_m | \mathbf{x})]^{k-i}$$

when k increases, the probability for selecting ω_m increases

S. Makrogiannis (DSU)

Nonparametric Techniques

November 6, 2015 10 / 26

The Nearest Neighbor Rule

The k-Nearest Neighbor (k-NN) Rule



Figure: k NN rule example for k = 5

S. Makrogiannis (DSU)

Nonparametric Techniques

The Nearest Neighbor Rule

The k-Nearest Neighbor (k-NN) Rule Error Rate



Figure: k NN rule error rate for variable k

S. Makrogiannis (DSU)

Nonparametric Techniques

 $\exists \rightarrow$ November 6, 2015 12 / 26

Considerations for the k-Nearest Neighbor (k-NN) Rule

- This rule can be considered as an attempt to estimate posterior probabilities P(\overline{\overlin}\overline{\overline{\overline{\overline{\ove
- A large value of k reduces the error rate
- However, we still need to keep the neighbors \mathbf{x}' around \mathbf{x} to be small enough so that $P(\boldsymbol{\omega}_i | \mathbf{x}) \simeq P(\boldsymbol{\omega}_i | \mathbf{x}')$
- In practice, this means that k should be a small fraction of n

S. Makrogiannis (DSU)

Nonparametric Techniques

November 6, 2015 13 / 26

Computational Complexity of the k-Nearest Neighbor Rule (k-NN)

- Suppose n samples with d dimensions. We are looking for the nearest neighbor to a point x
- The simplest approach for this is:
 - (1) Calculate Euclidean distances from all training points to x'
 - 2 Find the point with minimum distance
- Computational complexity is $O(dn^2)$

S. Makrogiannis (DSU)

Nonparametric Techniques

November 6, 2015 14 / 26

The Nearest Neighbor Rule

Reducing Computational Complexity of the k-Nearest Neighbor Rule (k-NN)

Basic approaches

- Partial distance
- Prestructuring with a search tree
- 3 Editing prototypes by eliminating non-informative prototypes during training

S. Makrogiannis (DSU)

Nonparametric Techniques

→ 4 差 ト 4 差 ト 差 の Q C November 6, 2015 15 / 26

Reducing Computational Complexity of the k-Nearest Neighbor Rule (k-NN)

Partial distance

(1) We calculate distances using a subset of the full d-dimensional space

$$D_r(\boldsymbol{a}, \boldsymbol{b}) = \left(\sum_{k=1}^r (a_k - b_k)^2\right)^{1/2}, r \leq d$$

2 If the distance exceeds a limit, then we do not compute further

S. Makrogiannis (DSU)

Nonparametric Techniques

November 6, 2015 16 / 26

Reducing Computational Complexity of the k-Nearest Neighbor Rule (k-NN)

Prestructuring with a search tree

- I Here we create a search tree with linked prototypes
- In classification stage, we compute distances from x to linked "entry" prototypes
- 3 We find the closest prototype and find its linked prototypes
- We compute distances from x to these prototypes
- 5 Repeat until we reach no other tree elements

S. Makrogiannis (DSU)

Nonparametric Techniques

November 6, 2015 17 / 26

Reducing Computational Complexity of the k-Nearest Neighbor Rule (k-NN)

Editing prototypes

We eliminate non-informative prototypes during training Algorithm 3 (Nearest-neighbor editing)

1	<u>begin initialize</u> $j = 0, \mathcal{D} = \text{data set}, n = \# \text{prototypes}$
2	construct the full Voronoi diagram of \mathcal{D}
3	<u>do</u> $j \leftarrow j + 1$; for each prototype \mathbf{x}'_j
4	Find the Voronoi neighbors of \mathbf{x}'_i
5	<u>if</u> any neighbor is not from the same class as \mathbf{x}'_i then mark \mathbf{x}'_i
6	$\underline{\text{until }} j = n$
7	Discard all points that are not marked
8	Construct the Voronoi diagram of the remaining (marked) prototypes
g end	

S. Makrogiannis (DSU)

Nonparametric Techniques

Metrics and Nearest Neighbor Classification

- The central component of a nearest neighbor classifier is the distance function $D(\cdot, \cdot)$ between patterns
- $D(\cdot, \cdot)$ is usually a metric

S. Makrogiannis (DSU)

Nonparametric Techniques

November 6, 2015 19 / 26

Sac

Properties of Metrics

Let a, b, c be three vector data points in a vector space \mathbb{R}^d with dimensionality d. Then a metric $D : \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$ must have the following properties

- Nonnegativity: $D(\boldsymbol{a}, \boldsymbol{b}) \geq 0$
- Reflexivity: $D(\boldsymbol{a}, \boldsymbol{b}) = 0 \Leftrightarrow \boldsymbol{a} = \boldsymbol{b}$
- Symmetry: $D(\boldsymbol{a}, \boldsymbol{b}) = D(\boldsymbol{b}, \boldsymbol{a})$
- Triangle inequality: $D(\boldsymbol{a},\boldsymbol{b}) + D(\boldsymbol{b},\boldsymbol{c}) \geq D(\boldsymbol{a},\boldsymbol{c})$

S. Makrogiannis (DSU)

Nonparametric Techniques

November 6, 2015 20 / 26
Metrics

Minkowski Metric

A general class of metrics for *d*-dimensional patterns is the Minkowski metric given by

$$L_k(\boldsymbol{a}, \boldsymbol{b}) = \left(\sum_{i=1}^d |a_i - b_i|^k\right)^{1/k}$$

This is also called the L_k norm

- L₁ norm: Manhattan or city block distance. This is the shortest path between *a* and *b*. In this path each segment is parallel to the coordinate axes
- L₂ norm: Euclidean distance
- L_∞ norm: corresponds to the maximum of the distances between the projections of *a* and *b* onto each of the *d* coordinate axes

S. Makrogiannis (DSU)

Nonparametric Techniques

 Metrics and Nearest Neighbor Classification

Minkowski Metrics Example



Figure: Isosurfaces for L_1 (white), L_2 (light gray), L_4 (dark gray), and L_{∞} (pink)

S. Makrogiannis (DSU)

Nonparametric Techniques

November 6, 2015 22 / 26

イロト イ戸ト イヨト

Sac

Metrics

Tanimoto Metric

This is mostly used in taxonomy to find a distance between two sets S_1 and S_2 :

$$D_{Tanimoto}(S_1, S_2) = rac{n_1 + n_2 - 2n_{12}}{n_1 + n_2 - n_{12}}$$

where n_1, n_2 are the cardinalities of sets S_1 and S_2 , and n_{12} is the cardinality of $S_1 \cap S_2$

S. Makrogiannis (DSU)

Nonparametric Techniques

November 6, 2015 23 / 26

Tangent Distance

- On several occasions the computation of a specific metric in an NN classifier may be sensitive to the problem of invariance
- For example, suppose we need to classify digits using 10 × 10 pixel grayscale images
- Then a small translation of a digit by a few pixels may change the metric computation significantly
- One can reduce this problem by use of the tangent distance



Figure: In this example the digit 8 is closer to the prototype of digit 5 than the shifted digit 5

Tangent Distance

- The general approach for addressing variability is to construct tangent vectors for all transformations
- Let x' be a prototype, α_i be a transformation parameter, and F_i(x'; α_i) be the transformed prototype
- Then for each transformation we can construct the tangent vector

 $TV_i = F_i(\mathbf{x}'; \alpha_i) - \mathbf{x}'$

 This process produces an r×d matrix
 T for each vector x' that corresponds to the tangent vectors at x'



Figure: Generation of tangent vectors using rotation and line thinning

イロト イポト イヨト イヨト

S. Makrogiannis (DSU)

Nonparametric Techniques

November 6, 2015 25 / 26

Sar

Tangent Distance

- To reduce the computation load we can sample the tangent space to create a linear subspace
- Then we compute the tangent distance between x' and x

$$D_{tan}(\mathbf{x}',\mathbf{x}) = \min_{\mathbf{a}} \| (\mathbf{x}' + \mathbf{T} \mathbf{\alpha}) - \mathbf{x} \|$$

- So we are looking for the smallest Euclidean distance from *x'* to *x*
- In classification we need to find optimizing value of α_i
- In that case we find the optimal *a* using gradient descent or matrix methods





Figure: Euclidean space produced by tangent vectors TV_1 and TV_2

November 6, 2015

26 / 26

イロト イポト イヨト イ

Linear Discriminant Functions MTSC 852, Fall 2015

Sokratis Makrogiannis, Ph.D.

Department of Mathematical Sciences Delaware State University smakrogiannis@desu.edu

November 12, 2015

S. Makrogiannis (DSU)

Linear Discriminant Functions

November 12, 2015 1 / 17

▲ロト ▲掃ト ▲ヨト ▲ヨト ニヨー のへで

Outline



- Linear Discriminant Functions and Decision Surfaces
 - Two-category Case
 - Multicategory Case



Generalized Linear Discriminant Functions

S. Makrogiannis (DSU)

Linear Discriminant Functions

November 12, 2015 2 / 17

3

Sac

イロト イポト イヨト イヨト

Introduction

- Pattern classification approaches can be categorized into *generative* and *discriminant*
- Generative approaches first estimate models of class-conditional probability density that generate the training patterns, then define the discriminant function. These methods usually assume some parametric form of density
- *Discriminant approaches* find the discriminant function directly without making any assumptions about the distribution of the training patterns
- This chapter deals with discriminant functions that are either linear in components of pattern vector **x**, or linear in a function of **x**

S. Makrogiannis (DSU)

Linear Discriminant Functions

November 12, 2015 3 / 17

Introduction

- Linear discriminant functions are analytically tractable and computationally efficient
- They can be used for trial classifiers
- We can find a linear discriminant function by minimizing a criterion function
- The usual criterion function is related to the training error, that is the average loss for classifying the training patterns
- A major consideration is to ensure that the discriminant function will classify at a good rate unknown, new patterns

S. Makrogiannis (DSU)

Linear Discriminant Functions

November 12, 2015 4 / 17

Linear Discriminant Function

• A linear discriminant function g(x) is typically of the form:

$$g(\boldsymbol{x}) = \boldsymbol{w}^T \boldsymbol{x} + w_0$$

where \boldsymbol{x} is the pattern vector, \boldsymbol{w} is the weight vector, and w_0 is the bias or threshold weight

 In general, when we have c classes, we need to find c discriminant functions

S. Makrogiannis (DSU)

Linear Discriminant Functions

November 12, 2015 5 / 17

Two categories

- ullet Suppose a problem with two categories ω_1 and ω_2
- Then we can define a single discriminant function by

$$g(\boldsymbol{x}) = g_1(\boldsymbol{x}) - g_2(\boldsymbol{x})$$

• The decision rule is:

Decide
$$\omega_1$$
 if $g(\mathbf{x}) > 0$; decide ω_2 if $g(\mathbf{x}) < 0$

- The decision surface is represented by g(x) = 0. For a linear discriminant function the decision surface is a hyperplane H
- Let points x_1, x_2 on the decision surface. Then:

$$\boldsymbol{w}^{T}\boldsymbol{x}_{1} = \boldsymbol{w}^{T}\boldsymbol{x}_{2} \Leftrightarrow \boldsymbol{w}^{T}(\boldsymbol{x}_{1} - \boldsymbol{x}_{2}) = 0$$

We observe that w is normal to the hyperplane

Two categories

- g(x) also gives a measure of the distance of the point x to the hyperplane
- We can write x as:

$$\boldsymbol{x} = \boldsymbol{x}_p + r \frac{\boldsymbol{w}}{\|\boldsymbol{w}\|}$$

where x_p is the normal projection of x to H, r is the algebraic distance • Given that $g(x_p) = 0$, it follows that

$$g(\mathbf{x}) = \mathbf{w}^{T}\mathbf{x} + w_{0} = \mathbf{w}^{T}(\mathbf{x}_{p} + r\frac{\mathbf{w}}{\|\mathbf{w}\|}) + w_{0} = \mathbf{w}^{T}\mathbf{x}_{p} + r\frac{\mathbf{w}^{T}\mathbf{w}}{\|\mathbf{w}\|} + w_{0}$$
$$\Leftrightarrow g(\mathbf{x}) = \mathbf{w}^{T}\mathbf{x}_{p} + r\|\mathbf{w}\| + w_{0} \Leftrightarrow g(\mathbf{x}) = r\|\mathbf{w}\| \Leftrightarrow r = \frac{g(\mathbf{x})}{\|\mathbf{w}\|}$$

S. Makrogiannis (DSU)

Linear Discriminant Functions

November 12, 2015 7 / 17

Two categories

- A linear discriminant function produces a hyperplane decision surface
- The weight vector w determines the orientation and the bias w₀ determines the location of the hyperplane
- g(x) is proportional to the distance of point x from the hyperplane
- *w* points toward the positive side of the hyperplane where g(x) > 0



Figure: Linear decision boundary for two categories

- To address classification into multiple c categories one could
 i) make c decisions for ω_i versus not ω_i, or ii) define c(c-1)/2 discriminants for all possible class pairs
- However these strategies can lead to regions of undefined classification



Figure: One-against-all and taking all pairwise discriminant may lead to Linear Discriminant Functions November 12, 2015 9 / 17

S. Makrogiannis (DSU)

• To avoid these shortcomings we can use a set of discriminant functions $g_i(x), i = 1, ..., c$ corresponding to each category, with a decision rule:

Assign feature vector **x** to class ω_i , if $g_i(\mathbf{x}) > g_j(\mathbf{x})$ $\forall j \neq i$

 This classifier can be considered as a network or linear machine that computes c discriminant functions and makes a decision using a maximum operator



S. Makrogiannis (DSU)

Linear Discriminant Functions

November 12, 2015 10 / 17

- $\bullet~$ Let $\mathscr{R}_i, \mathscr{R}_j$ be the decision regions corresponding to categories $\pmb{\omega}_i, \pmb{\omega}_j$
- If $\mathscr{R}_i, \mathscr{R}_j$ are contiguous, then at the boundary H_{ij} between $\mathscr{R}_i, \mathscr{R}_j$ we have that

$$g_i(\mathbf{x}) = g_j(\mathbf{x}) \Leftrightarrow (\mathbf{w}_i - \mathbf{w}_j)^T \mathbf{x} + (w_{i0} - w_{j0}) = 0$$

- We observe that w_i w_j is normal to H_{ij} and similarly to the two-category case the distance from x to H_{ij} is (g_i(x) g_j(x)/||w_i w_j||
- In a linear machine the decision regions are convex, and every decision region is singly connected, therefore this classifier is best suited for unimodal class-conditional densities $p(\mathbf{x}|\omega_i)$

S. Makrogiannis (DSU)

Linear Discriminant Functions

◆□ ト ◆□ ト ◆ 三 ト ◆ 三 ト ◆ 三 ・ つ へ ○ November 12, 2015 11 / 17



Decision Boundaries for linear machines Figure:

S. Makrogiannis (DSU)

Linear Discriminant Functions

∃ ⊳ November 12, 2015 12 / 17

990

3

Higher Order Discriminant Functions

Linear discriminant:

$$g(\mathbf{x}) = w_0 + \sum_{i=1}^d w_i x_i$$

where w_i : components of x

• Quadratic discriminant includes second order terms x_ix_j:

$$g(\mathbf{x}) = w_0 + \sum_{i=1}^d w_i x_i + \sum_{i=1}^d \sum_{j=1}^d w_{ij} x_i x_j$$

• By adding higher order terms we generate a polynomial discriminant

S. Makrogiannis (DSU)

Linear Discriminant Functions

November 12, 2015 13 / 17

Generalized Linear Discriminant Functions

• The generalized discriminant function is defined as:

$$g(\mathbf{x}) = \sum_{i=1}^{\hat{d}} \alpha_i y_i(\mathbf{x}) = \mathbf{\alpha}^T \mathbf{y}$$

where $\pmb{\alpha}$: \hat{d} -dimensional vector

- $y_i(x)$: arbitrary functions of x, also called ϕ functions
- g(x) may be nonlinear in x but is linear in y
- Functions $y_i(x)$ map points from a *d*-dimensional space to a \hat{d} -dimensional space
- These mappings may simplify the original problem to that of finding a linear discriminant function

S. Makrogiannis (DSU)

Linear Discriminant Functions

November 12, 2015 14 / 17

Generalized Linear Discriminant Functions - Example

• Let g(x) be:

$$g(\mathbf{x}) = \alpha_1 + \alpha_2 x + \alpha_3 x^2$$

• Then **y** is :

$$\boldsymbol{y} = \left(\begin{array}{c} 1\\ x\\ x^2 \end{array}\right)$$

- The inherent dimensionality of the data is 1
- The decision regions in y-space are convex, but they are non-convex in x-space



Figure: The ϕ functions transform a line to a parabola and creates a non-simply connected decision region in the x-space

S. Makrogiannis (DSU)

Linear Discriminant Functions

November 12, 2015 15 / 17

Generalized Linear Discriminant Functions

- The mapping to a higher dimensional space while interesting in concept, it implies a requirement for a lot more training data because of the curse of dimensionality
- We can address this drawback by imposing a constraint of large margins between the training samples. This approach is used in Support Vector Machines (SVM)

S. Makrogiannis (DSU)

Linear Discriminant Functions

November 12, 2015 16 / 17

Linear Discriminant Functions and Augmented Vectors

• For the linear discriminant we have that:

$$g(\mathbf{x}) = w_0 + \sum_{i=1}^d w_i x_i = \sum_{i=0}^d w_i x_i$$

where
$$x_0 = 1$$

Then $g(\mathbf{x}) = \boldsymbol{\alpha}^T \mathbf{y}$

$$\boldsymbol{y} = \begin{pmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{pmatrix} = \begin{pmatrix} 1 \\ \boldsymbol{x} \end{pmatrix} \text{ and } \boldsymbol{\alpha} = \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{pmatrix} = \begin{pmatrix} w_0 \\ \boldsymbol{w} \end{pmatrix}$$

- This is a transformation from a d-dimensional to a d+1-dimensional space creates a hyperplane passing from the origin
- We reduce the problem of finding a weight vector w and a threshold weight w₀ to that of finding just a weight vector d
- S. Makrogiannis (DSU)

Linear Discriminant Functions

November 12, 2015 17 / 17

Linear Discriminant Functions MTSC 852, Fall 2015

Sokratis Makrogiannis, Ph.D.

Department of Mathematical Sciences Delaware State University smakrogiannis@desu.edu

November 17, 2015

S. Makrogiannis (DSU)

Linear Discriminant Functions

November 17, 2015 1 / 12



2 Perceptron Criterion Function

S. Makrogiannis (DSU)

Linear Discriminant Functions

Linearly separable samples

- Suppose a set of samples $\mathscr{D} = \{y_1, y_2, \dots, y_n\}$
- We are looking for a weight vector \boldsymbol{a} in a linear discriminant function $g(\boldsymbol{y}) = \boldsymbol{a}^T \boldsymbol{y}$
- If there exists a weight vector that classifies all samples correctly, then the samples are said to be linearly separable
- By the decision rule, a sample y_i is classified correctly to ω_1 , if $a^T y_i > 0$, and is classified correctly to ω_2 , if $a^T y_i < 0$
- If we negate the samples in ω_2 -a process called normalization- then we are looking for a weight vector \boldsymbol{a} , such that $\boldsymbol{a}^T \boldsymbol{y}_i > 0, \forall \boldsymbol{y}_i \in \mathscr{D}$
- Such a vector is called a separating vector, or solution vector

S. Makrogiannis (DSU)

- The weight vector can be considered as a point in the weight space
- Each sample places a constraint on the location of the solution vector
- Assuming normalization, the solution vector must be located on the positive side of the separating hyperplane
- The solution vector lies on the intersection of *n* half-spaces. This region is called the solution region





Margins

- Based on previous discussion, we can tell that the solution vector is not unique
- We can impose additional constraints: one is to look for the unit vector that maximizes the minimum distance from the samples to the separating hyperplane another is to look for the minimum length weight vector **a** such that $\mathbf{a}^T \mathbf{y} \ge b$, where b is a positive constant called margin
- The objective is to find a solution vector that is closer to the middle of the solution region, therefore the classifier will perform well for unlabeled samples



Figure:

S. Makrogiannis (DSU)

Linear Discriminant Functions

Gradient Descent

- To find a solution of $\mathbf{a}^T \mathbf{y}_i > 0$ we can define a criterion function $J(\mathbf{a})$ and seek the minimizing vector \mathbf{a}
- Function minimization can be achieved by the gradient descent technique
- The gradient descent begins from an initial value of *a* and moves along the negative of the gradient:

$$\boldsymbol{a}(k+1) = \boldsymbol{a}(k) - \boldsymbol{n}(k)\boldsymbol{\nabla}J(\boldsymbol{a}(k))$$

where n is a positive scale factor or learning rate (step size)

• This method is dependent on *n*(*k*): if it's too small, convergence is slow, whereas if it's too large, the algorithm may overshoot or diverge

Algorithm 1 (Basic gradient descent)

```
\begin{array}{c} 1 \ \underline{\text{begin initialize}}\\ 2 \ \underline{do} \ k \leftarrow k+1\\ 3 \ \underline{a} \leftarrow a - \eta(k) \nabla J(\mathbf{a})\\ 4 \ \underline{\text{until}}\\ \eta(k) \nabla J(\mathbf{a}) < \theta\\ 5 \ \underline{\text{return }}\\ 6 \ \underline{\text{end}}\\ \end{array}
```

Setting the Learning Rate n(k)

• The second order Taylor expansion of the criterion function J(a) at a is:

$$J(\boldsymbol{a}) \simeq J(\boldsymbol{a}(k)) + \boldsymbol{\nabla} J^{\mathsf{T}}(\boldsymbol{a} - \boldsymbol{a}(k)) + (1/2)(\boldsymbol{a} - \boldsymbol{a}(k))^{\mathsf{T}} \boldsymbol{H}(\boldsymbol{a} - \boldsymbol{a}(k))$$

• By using the gradient descent equation we get:

$$J(\boldsymbol{a}(k+1)) \simeq J(\boldsymbol{a}(k)) + n(k) \|\boldsymbol{\nabla}J\|^2 + (1/2)n^2(k)\boldsymbol{\nabla}J^T \boldsymbol{H}\boldsymbol{\nabla}J$$

where **H** is the Hessian matrix of second order derivatives $\frac{\partial^2 J}{\partial a_i \partial a_j}$ • We can show that $J(\mathbf{a}(k+1))$ is minimized by:

$$n(k) = \frac{\|\boldsymbol{\nabla}J\|^2}{\boldsymbol{\nabla}J^T \boldsymbol{H} \boldsymbol{\nabla}J}$$

• If the criterion function is quadratic in the region of interest, then H is a constant and n is a constant as well

S. Makrogiannis (DSU)

Newton's Algorithm

• In this method we choose a(k+1) to minimize the second order expansion • This yields the equation:

$$oldsymbol{a}(k\!+\!1) = oldsymbol{a}(k) - oldsymbol{H}^{-1}oldsymbol{
abla} J$$

- Newton's algorithm achieves faster convergence
- However it requires that the Hessian is non-singular and implies increased computational cost $O(d^3)$

Algorithm 2 (Newton descent)

$$\begin{array}{c|c} t & \underline{\text{begin initialize}} & \text{a, criterion } \theta \\ \hline z & \underline{\text{do}} \\ s & \mathbf{a} \leftarrow \mathbf{a} - \mathbf{H}^{-1} \nabla J(\mathbf{a}) \\ 4 & \underline{\text{until}} & \mathbf{H}^{-1} \nabla J(\mathbf{a}) < \theta \\ \hline s & \underline{\text{return }} \mathbf{a} \\ \epsilon & \underline{\text{end}} \end{array}$$

Figure:

S. Makrogiannis (DSU)

Gradient Descent vs. Newton's Algorithm



Figure: Newton's algorithm (black line) converges in fewer steps than gradient descent (red line), however it is more computational intensive as it requires inversion of the Hessian

S. Makrogiannis (DSU)

Linear Discriminant Functions

The Perceptron Criterion Function

- The only missing piece is now the criterion function J(a)
- Since we need to minimize this function, one obvious option is the number of misclassified samples. This function is discontinuous, therefore gradient calculation may be problematic
- We can use the linear discriminant function to define the Perceptron criterion function after normalizing the patterns:

$$J_p(\boldsymbol{a}) = \sum_{\boldsymbol{y} \in \mathscr{Y}} (-\boldsymbol{a}^T \boldsymbol{y})$$

where $\mathscr{Y}(a)$ is the set of samples misclassified by a

• This criterion function is always nonnegative because $a^T y \leq 0$ for any misclassified sample

S. Makrogiannis (DSU)

Linear Discriminant Functions

November 17, 2015 10 / 12

The Perceptron Criterion Function

• We have that:

$$abla J_p = \sum_{oldsymbol{y} \in \mathscr{Y}} (-oldsymbol{y})$$

• The update rule in gradient descent becomes:

$$oldsymbol{a}(k+1) = oldsymbol{a}(k) + n(k) \sum_{oldsymbol{y} \in \mathscr{Y}_k} oldsymbol{y}$$

where \mathscr{Y}_k is the set of samples misclassified by $\boldsymbol{a}(k)$

In the Batch Perceptron algorithm the next vector is computed by adding the a multiple of the sum of misclassified samples to the weight vector from previous iteration Algorithm 3 (Batch Perceptron)

$$\begin{array}{c} i \ \underline{\text{begin initialize}} \\ \underline{a} \ \underline{begin (initialize)} \\ \underline{begin (initialize)} \\$$

S. Makrogiannis (DSU)

Linear Discriminant Functions

Perceptron Learning

If classes are linearly separable then the perceptron rule will converge to a valid solution

- A version of the perceptron rule uses variable learning rate
- However, this technique will converge only under specific conditions

On the other hand, if the two classes are not linearly separable, the perceptron rule will not converge

- Because there will always be at least one misclassified sample, the corrections in perceptron rule will continue with no end
- One approach to solve this problem is to use variable rates n(k) approaching zero as $k \to \infty$

Linear Discriminant Functions MTSC 852, Fall 2015

Sokratis Makrogiannis, Ph.D.

Department of Mathematical Sciences Delaware State University smakrogiannis@desu.edu

November 19, 2015

S. Makrogiannis (DSU)

Linear Discriminant Functions
Outline



1 Minimum Squared Error Techniques



Minimum Squared Error Techniques

Minimum Squared Error Techniques

- The previous criterion functions use the misclassified samples
- Here we introduce techniques that use all samples for estimation of the weight vector
- Now we are looking to solve $\mathbf{a}^T \mathbf{y}_i = b_i$, where b_i is a set of arbitrary constants

Minimum Squared Error (MSE) Techniques

• So we are looking for a solution of a system of linear equations

• Using matrix notation, we are looking to find a weight vector that satisfies

$$\begin{pmatrix} y_{10} & y_{12} & \dots & y_{1d} \\ y_{20} & y_{22} & \dots & y_{2d} \\ \vdots & \vdots & & \vdots \\ \vdots & \vdots & & \vdots \\ y_{n0} & y_{n2} & \dots & y_{nd} \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \\ \vdots \\ a_d \end{pmatrix} = \begin{pmatrix} b_0 \\ b_1 \\ \vdots \\ \vdots \\ b_n \end{pmatrix}, \quad \mathbf{Ya} = \mathbf{b}$$

• If **Y** were non-singular, we could solve for $\mathbf{a} = \mathbf{Y}^{-1}b$

• But \boldsymbol{Y} is rectangular with more rows than columns, therefore our system is overdetermined and usually there is no exact solution

Minimum Squared Error (MSE) Techniques

• Still, we can find an approximate solution of $\mathbf{Y}\mathbf{a} = \mathbf{b}$ by minimizing

$$e = Ya - b$$

• We can define the equivalent criterion function $J_s(a)$:

$$J_s(\boldsymbol{a}) = \|\boldsymbol{Y}\boldsymbol{a} - \boldsymbol{b}\|^2 = \sum_{i=1}^n (\boldsymbol{a}^T \boldsymbol{y}_i - \boldsymbol{b}_i)^2$$

We can solve this problem using a gradient search technique
Otherwise, we form the gradient \(\nabla J_s\) and set it to zero:

$$\nabla J_s = \sum_{i=1}^n 2(\boldsymbol{a}^T \boldsymbol{y}_i - \boldsymbol{b}_i) \boldsymbol{y}_i = 2 \boldsymbol{Y}^T (\boldsymbol{Y} \boldsymbol{a} - \boldsymbol{b})$$
$$\boldsymbol{Y}^T (\boldsymbol{Y} \boldsymbol{a} - \boldsymbol{b}) = 0 \Leftrightarrow \boldsymbol{Y}^T \boldsymbol{Y} \boldsymbol{a} = \boldsymbol{Y}^T \boldsymbol{b}$$

S. Makrogiannis (DSU)

Linear Discriminant Functions

Minimum Squared Error (MSE) Techniques

- In the equation Y^TYa = Y^Tb, we observe that Y^TY is a d×d square matrix that is often nonsingular
- Then the solution is:

$$\boldsymbol{a} = \left(\boldsymbol{Y}^{T} \boldsymbol{Y}\right)^{-1} \boldsymbol{Y}^{T} \boldsymbol{b}$$
$$\Leftrightarrow \boldsymbol{a} = \boldsymbol{Y}^{\dagger} \boldsymbol{b}$$

where the $d \times n$ matrix $\mathbf{Y}^{\dagger} \equiv (\mathbf{Y}^{T}\mathbf{Y})^{-1}\mathbf{Y}^{T}$ is called the pseudoinverse of \mathbf{Y} with $\mathbf{Y}^{\dagger}\mathbf{Y} = \mathbf{I}$, but in general $\mathbf{Y}\mathbf{Y}^{\dagger} \neq \mathbf{I}$

Y[†] is defined as:

$$oldsymbol{Y}^{\dagger}\equiv\lim_{arepsilon
ightarrow0}(oldsymbol{Y}^{\,T}oldsymbol{Y}+arepsilonoldsymbol{I})^{-1}oldsymbol{Y}^{\,T}$$

- We can show that this limit always exists and that $a = Y^{\dagger}b$ is a solution to the MSE problem
- We note that the MSE solution depends on **b**. If **b** is fixed, the MSE solution does not necessarily yield the separating vector in a linearly separable case
- But MSE may yield a useful discriminant function in both separable and nonseparable cases

The LMS Procedure

The MSE criterion function $J_s(\mathbf{a}) = \|\mathbf{Y}\mathbf{a} - \mathbf{b}\|^2$ can be minimized by gradient descent

- It addresses problems of singularity of $\mathbf{Y}^T \mathbf{Y}$
- We don't have to deal with large matrices in data of high dimensionality

The LMS Procedure

From previous result, the gradient of the MSE criterion function is:

$$\boldsymbol{\nabla} J_{s} = 2 \, \boldsymbol{Y}^{T} (\, \boldsymbol{Y} \boldsymbol{a} - \boldsymbol{b})$$

Then the update rule becomes

$$\boldsymbol{a}(k+1) = \boldsymbol{a}(k) - \boldsymbol{n}(k)\boldsymbol{\nabla}J(\boldsymbol{a}(k)) = \boldsymbol{a}(k) + \boldsymbol{n}(k)\boldsymbol{Y}^{\mathsf{T}}(\boldsymbol{b} - \boldsymbol{Y}\boldsymbol{a}(k))$$

• We can show that if n(k) = n(1)/k, where n(1) is a positive constant, then this rule converges to a solution of $\mathbf{Y}^T(\mathbf{Ya}(k) - \mathbf{b}) = 0$

• Widrow-Hoff or LMS rule We can consider each sample sequentially to derive:

$$\boldsymbol{a}(k+1) = \boldsymbol{a}(k) + \boldsymbol{n}(k)(\boldsymbol{b}(k) - \boldsymbol{a}^{T}(k)\boldsymbol{y}^{k})\boldsymbol{y}^{k}$$

Algorithm 10 (LMS)

 $\begin{array}{c|c} 1 & \underline{\text{begin initialize}} & \mathbf{a}, \mathbf{b}, \text{criterion } \theta, \eta(\cdot), k = 0 \\ 2 & \underline{\text{do}} & k \leftarrow k + 1 \\ 3 & \mathbf{a} \leftarrow \mathbf{a} + \eta(k)(b_k - \mathbf{a}^t \mathbf{y}^k)\mathbf{y}^k \\ 4 & \underline{\text{until}} & \eta(k)(b_k - \mathbf{a}^t \mathbf{y}^k)\mathbf{y}^k < \theta \\ 5 & \underline{\text{return }} \mathbf{a} \\ 6 & \underline{\text{end}} \end{array}$

S. Makrogiannis (DSU)

Linear Discriminant Functions

Minimum Squared Error Techniques

Comparing Perceptron with MSE

- The perceptron rule always finds a solution if classes are linearly separable but does not converge if classes are nonseparable
- The MSE approach will converge to a solution but it may not be the separating hyperplane if classes are linearly separable



Figure: The LMS algorithm does not necessarily converge to the separating hyperplane even if the classes are linearly separable

Ho-Kashyap Techniques

- We noted that the MSE solution does not necessarily converge to the separating hyperplance for linearly separable classes
- This is mainly due to the selection of a fixed arbitrary margin vector ${m b}$ in MSE
- However by definition of linear separability, if the classes are linearly separable, then $\exists \hat{a}, \hat{b} : Y \hat{a} = \hat{b} > 0$
- One approach for finding the separating hyperplane using MSE would be to find both the separating vector **a** and margin vector **b**
- This can be achieved in the following steps
 - Use gradient descent to find b
 - ② Use MSE approach to find a
 - Iterate until convergence

Ho-Kashyap Descent Procedure

• The gradients of MSE criterion function $J_s(a) = \|Ya - b\|^2$ in a and b are

$$\nabla_{\boldsymbol{a}} J_{\boldsymbol{s}} = \sum_{i=1}^{n} 2(\boldsymbol{a}^{T} \boldsymbol{y}_{i} - \boldsymbol{b}_{i}) \boldsymbol{y}_{i} = 2 \boldsymbol{Y}^{T} (\boldsymbol{Y} \boldsymbol{a} - \boldsymbol{b})$$

$$\nabla_{\boldsymbol{b}} J_{\boldsymbol{s}} = \sum_{i=1}^{n} (-2) (\boldsymbol{a}^{T} \boldsymbol{y}_{i} - \boldsymbol{b}_{i}) = -2 (\boldsymbol{Y} \boldsymbol{a} - \boldsymbol{b})$$

- For a specific **b**, the MSE solution for **a** is $\mathbf{a} = \mathbf{Y}^{\dagger} \mathbf{b}$, where $\mathbf{Y}^{\dagger} \equiv (\mathbf{Y}^{T} \mathbf{Y})^{-1} \mathbf{Y}^{T}$
- We use gradient descent to find \boldsymbol{b} ; however we must keep the constraint $\boldsymbol{b} > 0$

Ho-Kashyap Descent Procedure

• We initialize **b** to a positive value and set to zero all positive components of $\nabla_{\boldsymbol{b}} J_s$ by use of $\nabla_{\boldsymbol{b}} J_s^- = 1/2 [\nabla_{\boldsymbol{b}} J_s - |\nabla_{\boldsymbol{b}} J_s|]$ in the update rule:

$$\boldsymbol{b}(k+1) = \boldsymbol{b}(k) - (n/2) [\boldsymbol{\nabla}_{\boldsymbol{b}} J_{\boldsymbol{s}} - |\boldsymbol{\nabla}_{\boldsymbol{b}} J_{\boldsymbol{s}}|]$$

We compute the weight vector
 a(k+1) = Y[†]b(k+1)

Algorithm 11 (Ho-Kashyap)

```
 \begin{array}{c} i \ \underline{\text{begin}} \ \underline{\text{initialize}} \ \mathbf{a}, \mathbf{b}, \eta(\cdot) < 1, \text{criteria} \ b_{min}, k_{max} \\ \hline x & \underline{\mathbf{d}} \ k \leftarrow k + 1 \\ \hline y & \mathbf{e} \leftarrow \mathbf{Y} \mathbf{a} - \mathbf{b} \\ \downarrow & \mathbf{e}^+ \leftarrow 1/2(\mathbf{e} + Abs[\mathbf{e}]) \\ \hline s & \mathbf{b} \leftarrow \mathbf{a} + 2\eta(k)\mathbf{e}^+ \\ \hline \epsilon & \mathbf{a} \leftarrow \mathbf{Y} \mathbf{b} \\ \hline \epsilon & \mathbf{a} \leftarrow \mathbf{Y} \mathbf{b} \\ \hline \tau & \underline{\text{if}} \ Abs[\mathbf{e}] \leq b_{min} \ \underline{\text{then}} \ \underline{\text{return}} \ \mathbf{a}, \mathbf{b} \ and \ \underline{\text{exit}} \\ \hline s & \underline{\text{nutil}} \ k = k_{max} \\ g & \text{Print NO SOLUTION FOUND} \\ \hline \rho \ \text{end} \end{array}
```