

# Introduction to Machine Learning

## Introduction to Bayesian Classification

Andres Mendez-Vazquez

June 3, 2020

# Outline

1

## Introduction

- Supervised Learning
- Handling Noise in Classification
- Models of Classification
- Naive Bayes
  - Examples
  - The Naive Bayes Model
  - The Multi-Class Case

2

## Discriminant Functions and Decision Surfaces

- Introduction
- Gaussian Distribution
- Influence of the Covariance  $\Sigma$
- Example
- Maximum Likelihood Principle
- Maximum Likelihood on a Gaussian
- Some Remarks

3

## Introduction

- A first solution for the Maximum A Posteriori (MAP)
- Maximum Likelihood Vs Maximum A Posteriori
- Properties of the MAP

4

## Exercises

- Some Stuff you can try

# Outline

1

## Introduction

### ● Supervised Learning

- Handling Noise in Classification
- Models of Classification
- Naive Bayes
  - Examples
  - The Naive Bayes Model
  - The Multi-Class Case

2

## Discriminant Functions and Decision Surfaces

- Introduction
- Gaussian Distribution
- Influence of the Covariance  $\Sigma$
- Example
- Maximum Likelihood Principle
- Maximum Likelihood on a Gaussian
- Some Remarks

3

## Introduction

- A first solution for the Maximum A Posteriori (MAP)
- Maximum Likelihood Vs Maximum A Posteriori
- Properties of the MAP

4

## Exercises

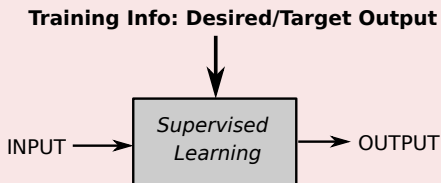
- Some Stuff you can try

# Classification Problem

## Goal

Given  $x_{new}$ , provide  $f(x_{new})$

## The Machinery in General looks...



# Outline

1

## Introduction

- Supervised Learning
- **Handling Noise in Classification**
- Models of Classification
- Naive Bayes
  - Examples
  - The Naive Bayes Model
  - The Multi-Class Case

2

## Discriminant Functions and Decision Surfaces

- Introduction
- Gaussian Distribution
- Influence of the Covariance  $\Sigma$
- Example
- Maximum Likelihood Principle
- Maximum Likelihood on a Gaussian
- Some Remarks

3

## Introduction

- A first solution for the Maximum A Posteriori (MAP)
- Maximum Likelihood Vs Maximum A Posteriori
- Properties of the MAP

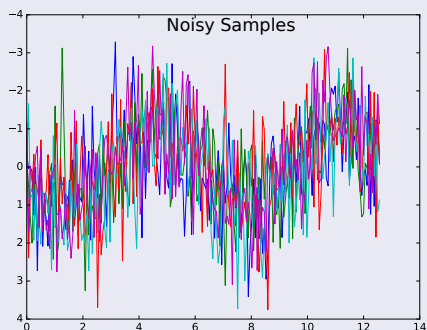
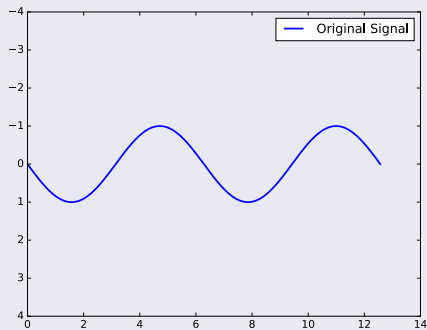
4

## Exercises

- Some Stuff you can try

# How do we handle Noise?

Imagine the following signal from  $\sin(\theta)$



## What if we know the noise?

Given a series of observed samples  $\{\hat{x}_1, \hat{x}_2, \dots, \hat{x}_N\}$  with noise  $\epsilon \sim N(0, 1)$

We could use our knowledge on the noise, for example additive:

$$\hat{x}_i = x_i + \epsilon$$

We can use our knowledge of probability to remove such noise

$$E[\hat{x}_i] = E[x_i + \epsilon] = E[x_i] + E[\epsilon]$$

Then, because  $E[\epsilon] = 0$

$$E[x_i] = E[\hat{x}_i] \approx \frac{1}{N} \sum_{i=1}^N \hat{x}_i$$

## What if we know the noise?

Given a series of observed samples  $\{\hat{\mathbf{x}}_1, \hat{\mathbf{x}}_2, \dots, \hat{\mathbf{x}}_N\}$  with noise  $\epsilon \sim N(0, 1)$

We could use our knowledge on the noise, for example additive:

$$\hat{\mathbf{x}}_i = \mathbf{x}_i + \epsilon$$

We can use our knowledge of probability to remove such noise

$$E[\hat{\mathbf{x}}_i] = E[\mathbf{x}_i + \epsilon] = E[\mathbf{x}_i] + E[\epsilon]$$

Then, because  $E[\epsilon] = 0$

$$E[\mathbf{x}_i] = E[\hat{\mathbf{x}}_i] \approx \frac{1}{N} \sum_{i=1}^N \hat{\mathbf{x}}_i$$



## What if we know the noise?

Given a series of observed samples  $\{\hat{\mathbf{x}}_1, \hat{\mathbf{x}}_2, \dots, \hat{\mathbf{x}}_N\}$  with noise  $\epsilon \sim N(0, 1)$

We could use our knowledge on the noise, for example additive:

$$\hat{\mathbf{x}}_i = \mathbf{x}_i + \epsilon$$

We can use our knowledge of probability to remove such noise

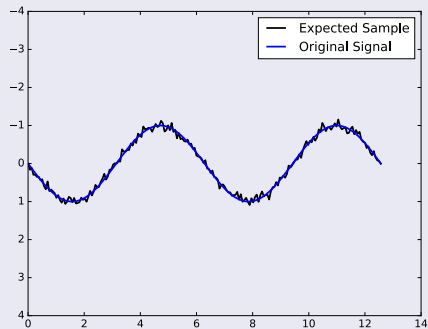
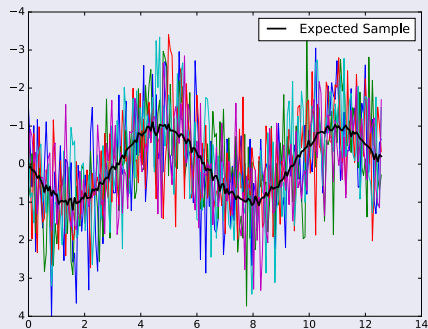
$$E[\hat{\mathbf{x}}_i] = E[\mathbf{x}_i + \epsilon] = E[\mathbf{x}_i] + E[\epsilon]$$

Then, because  $E[\epsilon] = 0$

$$E[\mathbf{x}_i] = E[\hat{\mathbf{x}}_i] \approx \frac{1}{N} \sum_{i=1}^N \hat{\mathbf{x}}_i$$

# In our example

We have a nice result



Therefore, we have

## The Bayesian Models

- They allow to deal with noise from the samples

Quite different from the deterministic models so far

- Unless Samples are Preprocessed to Reduce the Noise

Something that people in areas as Control tend to do

- The importance of Filters as Kalman Filters

Therefore, we have

## The Bayesian Models

- They allow to deal with noise from the samples

## Quite different from the deterministic models so far

- Unless Samples are Preprocessed to Reduce the Noise

something that people in areas as Control tend to do

- The importance of Filters as Kalman Filters

Therefore, we have

## The Bayesian Models

- They allow to deal with noise from the samples

## Quite different from the deterministic models so far

- Unless Samples are Preprocessed to Reduce the Noise

## Something that people in area as Control tend to do

- The importance of Filters as Kalman Filters

# Outline

1

## Introduction

- Supervised Learning
- Handling Noise in Classification
- **Models of Classification**
- Naive Bayes
  - Examples
  - The Naive Bayes Model
  - The Multi-Class Case

2

## Discriminant Functions and Decision Surfaces

- Introduction
- Gaussian Distribution
- Influence of the Covariance  $\Sigma$
- Example
- Maximum Likelihood Principle
- Maximum Likelihood on a Gaussian
- Some Remarks

3

## Introduction

- A first solution for the Maximum A Posteriori (MAP)
- Maximum Likelihood Vs Maximum A Posteriori
- Properties of the MAP

4

## Exercises

- Some Stuff you can try

# Example

## Given a Spoken Language

The task is to determine the language that someone is speaking

# Example

## Given a Spoken Language

The task is to determine the language that someone is speaking

## Generative Models

- They try to learn each language.
- Therefore, they try to determine the spoken language based in such learning.



# Example

## Given a Spoken Language

The task is to determine the language that someone is speaking

## Generative Models

- They try to learn each language.
- Therefore, they try to determine the spoken language based in such learning.

## Discriminative Models

- They try to determine the linguistic differences without learning any language!!!
- Quite easier!!!

# Example

## Given a Spoken Language

The task is to determine the language that someone is speaking

## Generative Models

- They try to learn each language.
- Therefore, they try to determine the spoken language based in such learning.

## Discriminative Models

- They try to determine the linguistic differences without learning any language!!!
- Quite easier!!!

# Example

## Given a Spoken Language

The task is to determine the language that someone is speaking

## Generative Models

- They try to learn each language.
- Therefore, they try to determine the spoken language based in such learning.

## Discriminative Models

- They try to determine the linguistic differences without learning any language!!!
- Quite easier!!!

# Therefore

## Generative Methods

- 1 Model class-conditional pdfs and prior probabilities.

- 2 "Generative" since sampling can generate synthetic data points.

# Therefore

## Generative Methods

- 1 Model class-conditional pdfs and prior probabilities.
- 2 “Generative” since sampling can generate synthetic data points.

## Examples

- Gaussians, Naïve Bayes, Mixtures of Multinomials.
- Mixtures of Gaussians, Mixtures of Experts, Hidden Markov Models (HMM).
- Sigmoidal Belief Networks, Bayesian Networks, Markov Random Fields.

# Therefore

## Generative Methods

- 1 Model class-conditional pdfs and prior probabilities.
- 2 “Generative” since sampling can generate synthetic data points.

## Examples

- Gaussians, Naïve Bayes, Mixtures of Multinomials.
- Mixtures of Gaussians, Mixtures of Experts, Hidden Markov Models (HMM).
- Sigmoidal Belief Networks, Bayesian Networks, Markov Random Fields

# Therefore

## Generative Methods

- 1 Model class-conditional pdfs and prior probabilities.
- 2 “Generative” since sampling can generate synthetic data points.

## Examples

- Gaussians, Naïve Bayes, Mixtures of Multinomials.
- Mixtures of Gaussians, Mixtures of Experts, Hidden Markov Models (HMM).
- Sigmoidal Belief Networks, Bayesian Networks, Markov Random Fields

# Therefore

## Generative Methods

- 1 Model class-conditional pdfs and prior probabilities.
- 2 “Generative” since sampling can generate synthetic data points.

## Examples

- Gaussians, Naïve Bayes, Mixtures of Multinomials.
- Mixtures of Gaussians, Mixtures of Experts, Hidden Markov Models (HMM).
- Sigmoidal Belief Networks, Bayesian Networks, Markov Random Fields.



# Furthermore

## Discriminative Methods

- 1 Directly estimate posterior probabilities.
- 2 No attempt to model underlying probability distributions.
- 3 Focus computational resources on given task for better performance.

# Furthermore

## Discriminative Methods

- 1 Directly estimate posterior probabilities.
  - 2 No attempt to model underlying probability distributions.
- Focus computational resources on given task for better performance.

## Popular models

- Logistic regression, SVMs.
- Traditional neural networks, Nearest neighbor.
- Conditional Random Fields (CRF).

# Furthermore

## Discriminative Methods

- 1 Directly estimate posterior probabilities.
- 2 No attempt to model underlying probability distributions.
- 3 Focus computational resources on given task for better performance.

## Popular models

- Logistic regression, SVMs.
- Traditional neural networks, Nearest neighbor.
- Conditional Random Fields (CRF).

# Furthermore

## Discriminative Methods

- 1 Directly estimate posterior probabilities.
- 2 No attempt to model underlying probability distributions.
- 3 Focus computational resources on given task for better performance.

## Popular models

- Logistic regression, SVMs.
- Traditional neural networks, Nearest neighbor.
- Conditional Random Fields (CRF).

# Furthermore

## Discriminative Methods

- 1 Directly estimate posterior probabilities.
- 2 No attempt to model underlying probability distributions.
- 3 Focus computational resources on given task for better performance.

## Popular models

- Logistic regression, SVMs.
- Traditional neural networks, Nearest neighbor.
- Conditional Random Fields (CRF).

# Furthermore

## Discriminative Methods

- 1 Directly estimate posterior probabilities.
- 2 No attempt to model underlying probability distributions.
- 3 Focus computational resources on given task for better performance.

## Popular models

- Logistic regression, SVMs.
- Traditional neural networks, Nearest neighbor.
- Conditional Random Fields (CRF).

# Outline

1

## Introduction

- Supervised Learning
- Handling Noise in Classification
- Models of Classification
- **Naive Bayes**
  - Examples
  - The Naive Bayes Model
  - The Multi-Class Case

2

## Discriminant Functions and Decision Surfaces

- Introduction
- Gaussian Distribution
- Influence of the Covariance  $\Sigma$
- Example
- Maximum Likelihood Principle
- Maximum Likelihood on a Gaussian
- Some Remarks

3

## Introduction

- A first solution for the Maximum A Posteriori (MAP)
- Maximum Likelihood Vs Maximum A Posteriori
- Properties of the MAP

4

## Exercises

- Some Stuff you can try

# Naive Bayes Model

## Task for two classes

Let  $\omega_1, \omega_2$  be the two classes in which our samples belong.



# Naive Bayes Model

## Task for two classes

Let  $\omega_1, \omega_2$  be the two classes in which our samples belong.

There is a prior probability of belonging to that class

- $P(\omega_1)$  for Class 1.
- $P(\omega_2)$  for Class 2.

# Naive Bayes Model

## Task for two classes

Let  $\omega_1, \omega_2$  be the two classes in which our samples belong.

There is a prior probability of belonging to that class

- $P(\omega_1)$  for Class 1.
- $P(\omega_2)$  for Class 2.

The Rule for classification is the following one:

$$P(\omega_i | \mathbf{x}) = \frac{P(\mathbf{x} | \omega_i) P(\omega_i)}{P(\mathbf{x})} \quad (1)$$

Remark: Bayes to the next level.

# Naive Bayes Model

## Task for two classes

Let  $\omega_1, \omega_2$  be the two classes in which our samples belong.

There is a prior probability of belonging to that class

- $P(\omega_1)$  for Class 1.
- $P(\omega_2)$  for Class 2.

The Rule for classification is the following one

$$P(\omega_i | \mathbf{x}) = \frac{P(\mathbf{x} | \omega_i) P(\omega_i)}{P(\mathbf{x})} \quad (1)$$

**Remark:** Bayes to the next level.

## In Informal English

We have that

$$\textit{posterior} = \frac{\textit{likelihood} \times \textit{prior-information}}{\textit{evidence}} \quad (2)$$

## In Informal English

We have that

$$\textit{posterior} = \frac{\textit{likelihood} \times \textit{prior-information}}{\textit{evidence}} \quad (2)$$

Basically

One: If we can observe  $x$ .

Two: we can convert the prior information into the posterior information.

## In Informal English

We have that

$$posterior = \frac{likelihood \times prior-information}{evidence} \quad (2)$$

Basically

One: If we can observe  $x$ .

Two: we can convert the prior-information into the posterior information.

## We have the following terms...

### Likelihood

We call  $p(\mathbf{x}|\omega_i)$  the likelihood of  $\omega_i$  given  $\mathbf{x}$ :

- This indicates that given a category  $\omega_i$ : If  $p(\mathbf{x}|\omega_i)$  is “large”, then  $\omega_i$  is the “likely” class of  $\mathbf{x}$ .

## We have the following terms...

### Likelihood

We call  $p(\mathbf{x}|\omega_i)$  the likelihood of  $\omega_i$  given  $\mathbf{x}$ :

- This indicates that given a category  $\omega_i$ : If  $p(\mathbf{x}|\omega_i)$  is “large”, then  $\omega_i$  is the “likely” class of  $\mathbf{x}$ .

### Prior Probability

It is the known probability of a given class.

Remark: Because, we lack information about this class, we tend to use the uniform distribution.

However: We can use other tricks for it.



## We have the following terms...

### Likelihood

We call  $p(\mathbf{x}|\omega_i)$  the likelihood of  $\omega_i$  given  $\mathbf{x}$ :

- This indicates that given a category  $\omega_i$ : If  $p(\mathbf{x}|\omega_i)$  is “large”, then  $\omega_i$  is the “likely” class of  $\mathbf{x}$ .

### Prior Probability

It is the known probability of a given class.

Remark: Because, we lack information about this class, we tend to use the uniform distribution.

However: We can use other tricks for it.

### Evidence

The evidence factor can be seen as a scale factor that guarantees that the posterior probability sum to one.

## We have the following terms...

### Likelihood

We call  $p(\mathbf{x}|\omega_i)$  the likelihood of  $\omega_i$  given  $\mathbf{x}$ :

- This indicates that given a category  $\omega_i$ : If  $p(\mathbf{x}|\omega_i)$  is “large”, then  $\omega_i$  is the “likely” class of  $\mathbf{x}$ .

### Prior Probability

It is the known probability of a given class.

**Remark:** Because, we lack information about this class, we tend to use the uniform distribution.

However, We can use other tricks for it.

### Evidence

The evidence factor can be seen as a scale factor that guarantees that the posterior probability sum to one.

## We have the following terms...

### Likelihood

We call  $p(\mathbf{x}|\omega_i)$  the likelihood of  $\omega_i$  given  $\mathbf{x}$ :

- This indicates that given a category  $\omega_i$ : If  $p(\mathbf{x}|\omega_i)$  is “large”, then  $\omega_i$  is the “likely” class of  $\mathbf{x}$ .

### Prior Probability

It is the known probability of a given class.

**Remark:** Because, we lack information about this class, we tend to use the uniform distribution.

**However:** We can use other tricks for it.

### Evidence

The evidence factor can be seen as a scale factor that guarantees that the posterior probability sum to one.

## We have the following terms...

### Likelihood

We call  $p(\mathbf{x}|\omega_i)$  the likelihood of  $\omega_i$  given  $\mathbf{x}$ :

- This indicates that given a category  $\omega_i$ : If  $p(\mathbf{x}|\omega_i)$  is “large”, then  $\omega_i$  is the “likely” class of  $\mathbf{x}$ .

### Prior Probability

It is the known probability of a given class.

**Remark:** Because, we lack information about this class, we tend to use the uniform distribution.

**However:** We can use other tricks for it.

### Evidence

The evidence factor can be seen as a scale factor that guarantees that the posterior probability sum to one.

The most important term in all this

The factor

*likelihood  $\times$  prior-information*

(3)

# Outline

1

## Introduction

- Supervised Learning
- Handling Noise in Classification
- Models of Classification
- **Naive Bayes**
  - **Examples**
    - The Naive Bayes Model
    - The Multi-Class Case

2

## Discriminant Functions and Decision Surfaces

- Introduction
- Gaussian Distribution
- Influence of the Covariance  $\Sigma$
- Example
- Maximum Likelihood Principle
- Maximum Likelihood on a Gaussian
- Some Remarks

3

## Introduction

- A first solution for the Maximum A Posteriori (MAP)
- Maximum Likelihood Vs Maximum A Posteriori
- Properties of the MAP

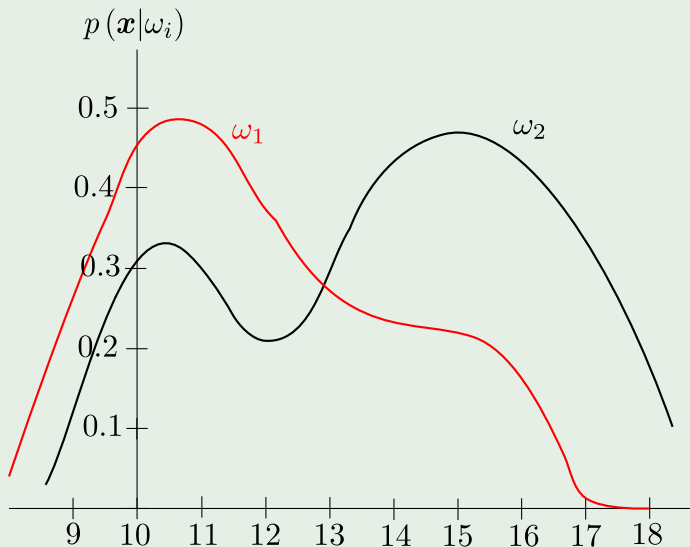
4

## Exercises

- Some Stuff you can try

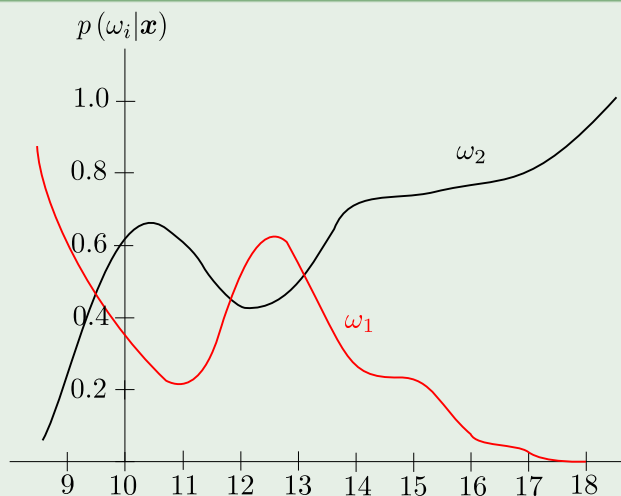
## Example

We have the likelihood of two classes



## Example

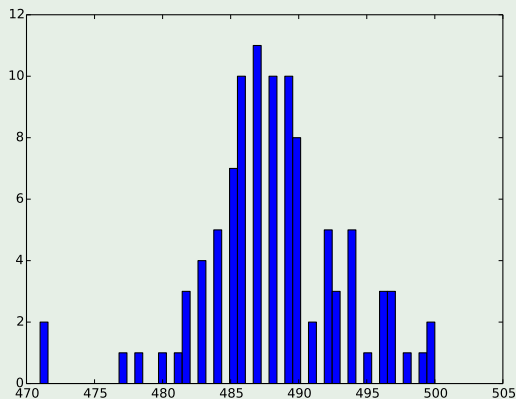
We have the posterior of two classes when  $P(\omega_1) = \frac{2}{3}$  and  $P(\omega_2) = \frac{1}{3}$





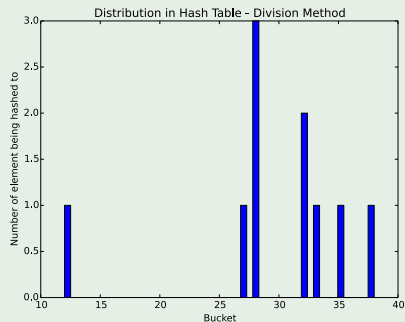
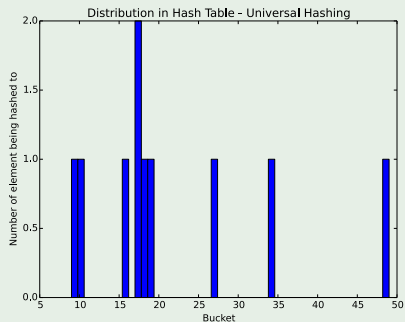
## Example of key distribution

Example, mean = 488.5 and dispersion = 5



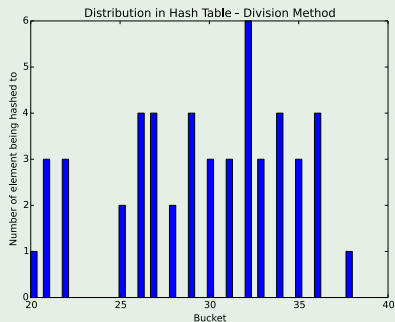
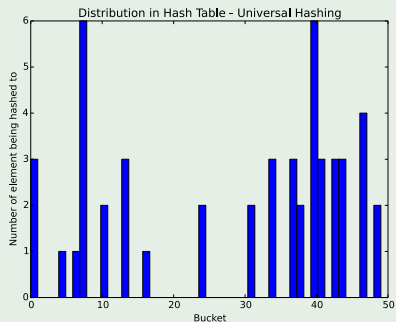
# Example with 10 keys

## Universal Hashing Vs Division Method



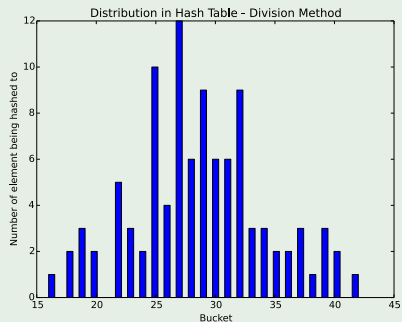
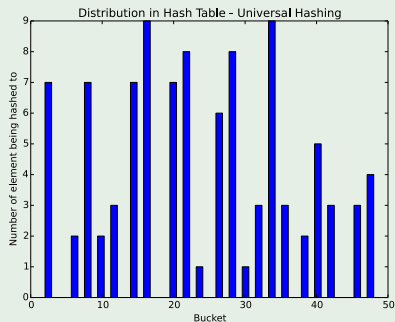
# Example with 50 keys

## Universal Hashing Vs Division Method



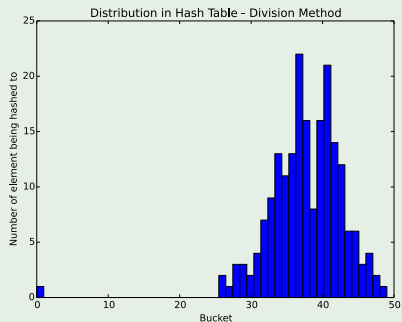
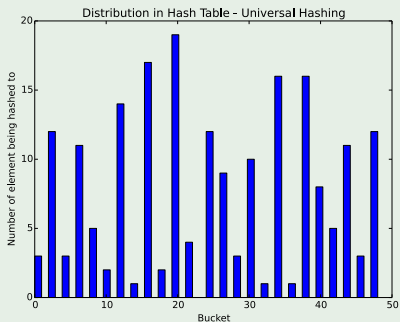
# Example with 100 keys

## Universal Hashing Vs Division Method



# Example with 200 keys

## Universal Hashing Vs Division Method



# Outline

1

## Introduction

- Supervised Learning
- Handling Noise in Classification
- Models of Classification
- **Naive Bayes**
  - Examples
  - **The Naive Bayes Model**
  - The Multi-Class Case

2

## Discriminant Functions and Decision Surfaces

- Introduction
- Gaussian Distribution
- Influence of the Covariance  $\Sigma$
- Example
- Maximum Likelihood Principle
- Maximum Likelihood on a Gaussian
- Some Remarks

3

## Introduction

- A first solution for the Maximum A Posteriori (MAP)
- Maximum Likelihood Vs Maximum A Posteriori
- Properties of the MAP

4

## Exercises

- Some Stuff you can try

# Naive Bayes Model

In the case of two classes, we can use demarginalization

$$P(\mathbf{x}) = \sum_{i=1}^2 p(\mathbf{x}, \omega_i) = \sum_{i=1}^2 p(\mathbf{x}|\omega_i) P(\omega_i) \quad (4)$$

## Error in this rule

We have that

$$P(\text{error}|\mathbf{x}) = \begin{cases} P(\omega_1|\mathbf{x}) & \text{if we decide } \omega_2 \\ P(\omega_2|\mathbf{x}) & \text{if we decide } \omega_1 \end{cases} \quad (5)$$

Thus, we have that

$$P(\text{error}) = \int_{-\infty}^{\infty} P(\text{error}, \mathbf{x}) d\mathbf{x} = \int_{-\infty}^{\infty} P(\text{error}|\mathbf{x}) p(\mathbf{x}) d\mathbf{x} \quad (6)$$



## Error in this rule

We have that

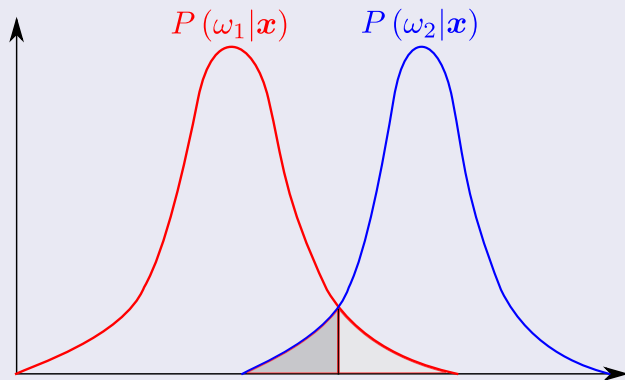
$$P(\text{error}|\mathbf{x}) = \begin{cases} P(\omega_1|\mathbf{x}) & \text{if we decide } \omega_2 \\ P(\omega_2|\mathbf{x}) & \text{if we decide } \omega_1 \end{cases} \quad (5)$$

Thus, we have that

$$P(\text{error}) = \int_{-\infty}^{\infty} P(\text{error}, \mathbf{x}) d\mathbf{x} = \int_{-\infty}^{\infty} P(\text{error}|\mathbf{x}) p(\mathbf{x}) d\mathbf{x} \quad (6)$$

# Graphically

We have



$$P(\text{error}) = \int_{-\infty}^{\infty} P(\text{error}, \mathbf{x}) d\mathbf{x}$$

# Classification Rule

Thus, we have the Bayes Classification Rule

1 If  $P(\omega_1|\mathbf{x}) > P(\omega_2|\mathbf{x})$   $\mathbf{x}$  is classified to  $\omega_1$

2 If  $P(\omega_1|\mathbf{x}) < P(\omega_2|\mathbf{x})$   $\mathbf{x}$  is classified to  $\omega_2$

# Classification Rule

Thus, we have the Bayes Classification Rule

- 1 If  $P(\omega_1|\mathbf{x}) > P(\omega_2|\mathbf{x})$   $\mathbf{x}$  is classified to  $\omega_1$
- 2 If  $P(\omega_1|\mathbf{x}) < P(\omega_2|\mathbf{x})$   $\mathbf{x}$  is classified to  $\omega_2$

## What if we remove the normalization factor?

Remember

$$P(\omega_1|\mathbf{x}) + P(\omega_2|\mathbf{x}) = 1 \quad (7)$$

## What if we remove the normalization factor?

### Remember

$$P(\omega_1|\mathbf{x}) + P(\omega_2|\mathbf{x}) = 1 \quad (7)$$

### We are able to obtain the new Bayes Classification Rule

1 If  $P(\mathbf{x}|\omega_1)p(\omega_1) > P(\mathbf{x}|\omega_2)P(\omega_2)$   $\mathbf{x}$  is classified to  $\omega_1$

2 If  $P(\mathbf{x}|\omega_1)p(\omega_1) < P(\mathbf{x}|\omega_2)P(\omega_2)$   $\mathbf{x}$  is classified to  $\omega_2$

## What if we remove the normalization factor?

### Remember

$$P(\omega_1|\mathbf{x}) + P(\omega_2|\mathbf{x}) = 1 \quad (7)$$

### We are able to obtain the new Bayes Classification Rule

- 1 If  $P(\mathbf{x}|\omega_1)p(\omega_1) > P(\mathbf{x}|\omega_2)P(\omega_2)$   $\mathbf{x}$  is classified to  $\omega_1$
- 2 If  $P(\mathbf{x}|\omega_1)p(\omega_1) < P(\mathbf{x}|\omega_2)P(\omega_2)$   $\mathbf{x}$  is classified to  $\omega_2$

## We have several cases

If for some  $\mathbf{x}$  we have  $P(\mathbf{x}|\omega_1) = P(\mathbf{x}|\omega_2)$

The final decision relies completely from the prior probability.

On the Other hand if  $P(\omega_1) = P(\omega_2)$ , the 'state' is equally probable

In this case the decision is based entirely on the likelihoods  $P(\mathbf{x}|\omega_i)$ .



## We have several cases

If for some  $\mathbf{x}$  we have  $P(\mathbf{x}|\omega_1) = P(\mathbf{x}|\omega_2)$

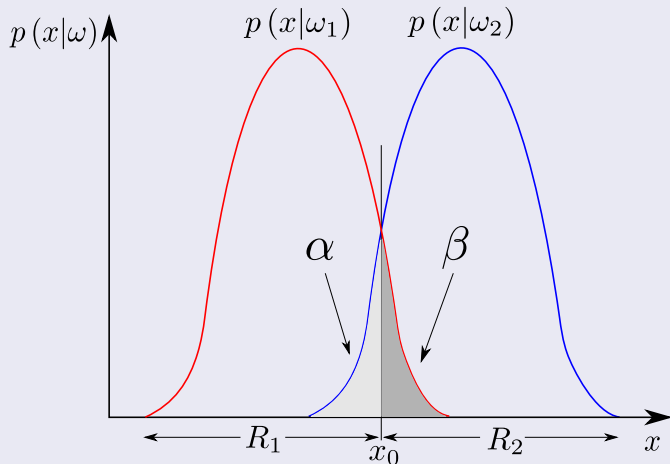
The final decision relies completely from the prior probability.

On the Other hand if  $P(\omega_1) = P(\omega_2)$ , the “state” is equally probable

In this case the decision is based entirely on the likelihoods  $P(\mathbf{x}|\omega_i)$ .

## How the Rule looks like

If  $P(\omega_1) = P(\omega_2)$  the Rule depends on the term  $p(x|\omega_i)$



## Error in Naive Bayes

Error in equiprobable classes  $p(\omega_1) = p(\omega_2) = \frac{1}{2}$

$$\begin{aligned} P_e &= \int_{-\infty}^{\infty} P(\mathbf{x}, \text{error}) d\mathbf{x} \\ &= \int_{-\infty}^{x_0} p(x, \omega_2) dx + \int_{x_0}^{\infty} p(x, \omega_1) dx \\ &= \int_{-\infty}^{x_0} p(x|\omega_2) P(\omega_2) dx + \int_{x_0}^{\infty} p(x|\omega_1) P(\omega_1) dx \\ &= P(\omega_2) \int_{-\infty}^{x_0} p(x|\omega_2) dx + P(\omega_1) \int_{x_0}^{\infty} p(x|\omega_1) dx \\ &= \frac{1}{2} \int_{-\infty}^{x_0} p(x|\omega_2) dx + \frac{1}{2} \int_{x_0}^{\infty} p(x|\omega_1) dx \end{aligned}$$

## Error in Naive Bayes

Error in equiprobable classes  $p(\omega_1) = p(\omega_2) = \frac{1}{2}$

$$\begin{aligned} P_e &= \int_{-\infty}^{\infty} P(\mathbf{x}, \text{error}) d\mathbf{x} \\ &= \int_{-\infty}^{x_0} p(\mathbf{x}, \omega_2) d\mathbf{x} + \int_{x_0}^{\infty} p(\mathbf{x}, \omega_1) d\mathbf{x} \\ &= \int_{-\infty}^{x_0} p(\mathbf{x}|\omega_2) P(\omega_2) d\mathbf{x} + \int_{x_0}^{\infty} p(\mathbf{x}|\omega_1) P(\omega_1) d\mathbf{x} \\ &= P(\omega_2) \int_{-\infty}^{x_0} p(\mathbf{x}|\omega_2) d\mathbf{x} + P(\omega_1) \int_{x_0}^{\infty} p(\mathbf{x}|\omega_1) d\mathbf{x} \\ &= \frac{1}{2} \int_{-\infty}^{x_0} p(\mathbf{x}|\omega_2) d\mathbf{x} + \frac{1}{2} \int_{x_0}^{\infty} p(\mathbf{x}|\omega_1) d\mathbf{x} \end{aligned}$$

## Error in Naive Bayes

Error in equiprobable classes  $p(\omega_1) = p(\omega_2) = \frac{1}{2}$

$$\begin{aligned}P_e &= \int_{-\infty}^{\infty} P(\mathbf{x}, \text{error}) d\mathbf{x} \\&= \int_{-\infty}^{x_0} p(x, \omega_2) dx + \int_{x_0}^{\infty} p(x, \omega_1) dx \\&= \int_{-\infty}^{x_0} p(x|\omega_2) P(\omega_2) dx + \int_{x_0}^{\infty} p(x|\omega_1) P(\omega_1) dx \\&= P(\omega_2) \int_{-\infty}^{x_0} p(x|\omega_2) dx + P(\omega_1) \int_{x_0}^{\infty} p(x|\omega_1) dx \\&= \frac{1}{2} \int_{-\infty}^{x_0} p(x|\omega_2) dx + \frac{1}{2} \int_{x_0}^{\infty} p(x|\omega_1) dx\end{aligned}$$

## Error in Naive Bayes

Error in equiprobable classes  $p(\omega_1) = p(\omega_2) = \frac{1}{2}$

$$\begin{aligned}P_e &= \int_{-\infty}^{\infty} P(\mathbf{x}, \text{error}) d\mathbf{x} \\&= \int_{-\infty}^{x_0} p(x, \omega_2) dx + \int_{x_0}^{\infty} p(x, \omega_1) dx \\&= \int_{-\infty}^{x_0} p(x|\omega_2) P(\omega_2) dx + \int_{x_0}^{\infty} p(x|\omega_1) P(\omega_1) dx \\&= P(\omega_2) \int_{-\infty}^{x_0} p(x|\omega_2) dx + P(\omega_1) \int_{x_0}^{\infty} p(x|\omega_1) dx \\&= \frac{1}{2} \int_{-\infty}^{x_0} p(x|\omega_2) dx + \frac{1}{2} \int_{x_0}^{\infty} p(x|\omega_1) dx\end{aligned}$$

## Error in Naive Bayes

Error in equiprobable classes  $p(\omega_1) = p(\omega_2) = \frac{1}{2}$

$$\begin{aligned} P_e &= \int_{-\infty}^{\infty} P(\mathbf{x}, \text{error}) d\mathbf{x} \\ &= \int_{-\infty}^{x_0} p(x, \omega_2) dx + \int_{x_0}^{\infty} p(x, \omega_1) dx \\ &= \int_{-\infty}^{x_0} p(x|\omega_2) P(\omega_2) dx + \int_{x_0}^{\infty} p(x|\omega_1) P(\omega_1) dx \\ &= P(\omega_2) \int_{-\infty}^{x_0} p(x|\omega_2) dx + P(\omega_1) \int_{x_0}^{\infty} p(x|\omega_1) dx \\ &= \frac{1}{2} \int_{-\infty}^{x_0} p(x|\omega_2) dx + \frac{1}{2} \int_{x_0}^{\infty} p(x|\omega_1) dx \end{aligned}$$

## Error in Naive Bayes

### Something Notable

**Bayesian classifier is optimal with respect to minimizing the classification error probability.**



# Proof

## Step 1

- $R_1$  be the region of the feature space in which we decide in favor of  $\omega_1$
- $R_2$  be the region of the feature space in which we decide in favor of  $\omega_2$

# Proof

## Step 1

- $R_1$  be the region of the feature space in which we decide in favor of  $\omega_1$
- $R_2$  be the region of the feature space in which we decide in favor of  $\omega_2$

## Step 2

$$P_e = P(x \in R_2, \omega_1) + P(x \in R_1, \omega_2) \quad (8)$$

# Proof

## Step 1

- $R_1$  be the region of the feature space in which we decide in favor of  $\omega_1$
- $R_2$  be the region of the feature space in which we decide in favor of  $\omega_2$

## Step 2

$$P_e = P(x \in R_2, \omega_1) + P(x \in R_1, \omega_2) \quad (8)$$

This

$$\begin{aligned} P_e &= P(x \in R_2 | \omega_1) P(\omega_1) + P(x \in R_1 | \omega_2) P(\omega_2) \\ &= P(\omega_1) \int_{R_2} p(x | \omega_1) dx + P(\omega_2) \int_{R_1} p(x | \omega_2) dx \end{aligned}$$

# Proof

## Step 1

- $R_1$  be the region of the feature space in which we decide in favor of  $\omega_1$
- $R_2$  be the region of the feature space in which we decide in favor of  $\omega_2$

## Step 2

$$P_e = P(x \in R_2, \omega_1) + P(x \in R_1, \omega_2) \quad (8)$$

## Thus

$$\begin{aligned} P_e &= P(x \in R_2 | \omega_1) P(\omega_1) + P(x \in R_1 | \omega_2) P(\omega_2) \\ &= P(\omega_1) \int_{R_2} p(x | \omega_1) dx + P(\omega_2) \int_{R_1} p(x | \omega_2) dx \end{aligned}$$

# Proof

## Step 1

- $R_1$  be the region of the feature space in which we decide in favor of  $\omega_1$
- $R_2$  be the region of the feature space in which we decide in favor of  $\omega_2$

## Step 2

$$P_e = P(x \in R_2, \omega_1) + P(x \in R_1, \omega_2) \quad (8)$$

## Thus

$$\begin{aligned} P_e &= P(x \in R_2 | \omega_1) P(\omega_1) + P(x \in R_1 | \omega_2) P(\omega_2) \\ &= P(\omega_1) \int_{R_2} p(x | \omega_1) dx + P(\omega_2) \int_{R_1} p(x | \omega_2) dx \end{aligned}$$

# Proof

It is more

$$P_e = P(\omega_1) \int_{R_2} \frac{p(\omega_1, x)}{P(\omega_1)} dx + P(\omega_2) \int_{R_1} \frac{p(\omega_2, x)}{P(\omega_2)} dx \quad (9)$$

Finally

$$P_e = \int_{R_2} p(\omega_1|x) p(x) dx + \int_{R_1} p(\omega_2|x) p(x) dx \quad (10)$$

Now, we choose the Bayes Classification Rule

$$R_1 : P(\omega_1|x) > P(\omega_2|x)$$

$$R_2 : P(\omega_2|x) > P(\omega_1|x)$$

# Proof

It is more

$$P_e = P(\omega_1) \int_{R_2} \frac{p(\omega_1, x)}{P(\omega_1)} dx + P(\omega_2) \int_{R_1} \frac{p(\omega_2, x)}{P(\omega_2)} dx \quad (9)$$

Finally

$$P_e = \int_{R_2} p(\omega_1|x) p(x) dx + \int_{R_1} p(\omega_2|x) p(x) dx \quad (10)$$

Now, we choose the Bayes Classification Rule:

$$R_1 : P(\omega_1|x) > P(\omega_2|x)$$

$$R_2 : P(\omega_2|x) > P(\omega_1|x)$$

## Proof

It is more

$$P_e = P(\omega_1) \int_{R_2} \frac{p(\omega_1, x)}{P(\omega_1)} dx + P(\omega_2) \int_{R_1} \frac{p(\omega_2, x)}{P(\omega_2)} dx \quad (9)$$

Finally

$$P_e = \int_{R_2} p(\omega_1|x) p(x) dx + \int_{R_1} p(\omega_2|x) p(x) dx \quad (10)$$

Now, we choose the Bayes Classification Rule

$$R_1 : P(\omega_1|x) > P(\omega_2|x)$$

$$R_2 : P(\omega_2|x) > P(\omega_1|x)$$



# Proof

Thus

$$P(\omega_1) = \int_{R_1} p(\omega_1|x) p(x) dx + \int_{R_2} p(\omega_1|x) p(x) dx \quad (11)$$

Now, we have

$$P(\omega_1) - \int_{R_1} p(\omega_1|x) p(x) dx = \int_{R_2} p(\omega_1|x) p(x) dx \quad (12)$$

Then

$$P_e = P(\omega_1) - \int_{R_1} p(\omega_1|x) p(x) dx + \int_{R_1} p(\omega_2|x) p(x) dx \quad (13)$$

# Proof

Thus

$$P(\omega_1) = \int_{R_1} p(\omega_1|x) p(x) dx + \int_{R_2} p(\omega_1|x) p(x) dx \quad (11)$$

Now, we have...

$$P(\omega_1) - \int_{R_1} p(\omega_1|x) p(x) dx = \int_{R_2} p(\omega_1|x) p(x) dx \quad (12)$$

Then

$$P_e = P(\omega_1) - \int_{R_1} p(\omega_1|x) p(x) dx + \int_{R_1} p(\omega_2|x) p(x) dx \quad (13)$$

# Proof

Thus

$$P(\omega_1) = \int_{R_1} p(\omega_1|x) p(x) dx + \int_{R_2} p(\omega_1|x) p(x) dx \quad (11)$$

Now, we have...

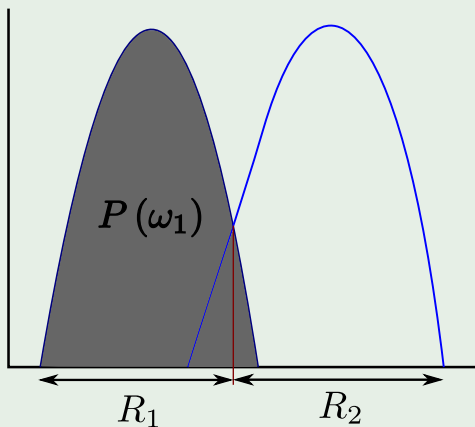
$$P(\omega_1) - \int_{R_1} p(\omega_1|x) p(x) dx = \int_{R_2} p(\omega_1|x) p(x) dx \quad (12)$$

Then

$$P_e = P(\omega_1) - \int_{R_1} p(\omega_1|x) p(x) dx + \int_{R_1} p(\omega_2|x) p(x) dx \quad (13)$$

# Graphically $P(\omega_1)$ : Thanks Edith 2013 Class!!!

In Gray

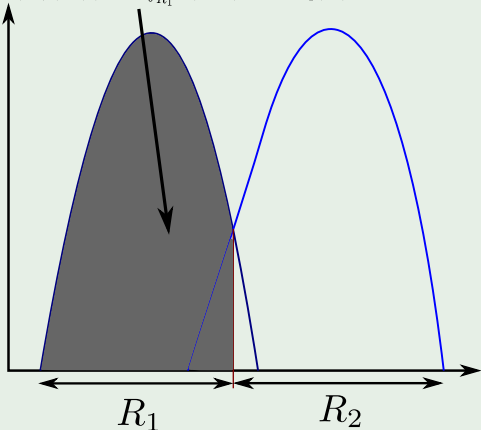


Thus we have

$$\int_{R_1} p(\omega_1|x) p(x) dx = \int_{R_1} p(\omega_1, x) dx = P_{R_1}(\omega_1)$$

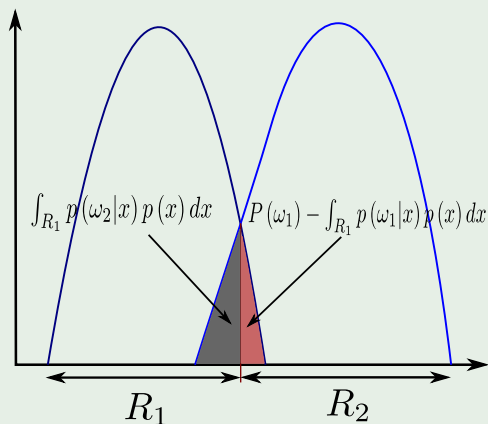
Thus

$$\int_{R_1} p(\omega_1|x) p(x) dx = \int_{R_1} p(\omega_1, x) dx = P_{R_1}(\omega_1)$$



Finally  $P_e$

A great idea Edith!!!



Thus

Finally

$$P_e = P(\omega_1) - \int_{R_1} [p(\omega_1|x) - p(\omega_2|x)] p(x) dx \quad (14)$$

Thus

The probability of error is minimized at the region of space in which  $R_1 : P(\omega_1|x) > P(\omega_2|x)$ .

Thus

Finally

$$P_e = P(\omega_1) - \int_{R_1} [p(\omega_1|x) - p(\omega_2|x)] p(x) dx \quad (14)$$

Thus

The probability of error is minimized at the region of space in which  $R_1 : P(\omega_1|x) > P(\omega_2|x)$ .



# Finally

## Similarly

$$P_e = P(\omega_2) - \int_{R_2} [p(\omega_2|x) - p(\omega_1|x)] p(x) dx \quad (15)$$

Thus

The probability of error is minimized at the region of space in which  $R_2 : P(\omega_2|x) > P(\omega_1|x)$ .

Thus

The Naive Bayes Rule minimizes the error.

## Finally

### Similarly

$$P_e = P(\omega_2) - \int_{R_2} [p(\omega_2|x) - p(\omega_1|x)] p(x) dx \quad (15)$$

### Thus

The probability of error is minimized at the region of space in which  $R_2 : P(\omega_2|x) > P(\omega_1|x)$ .

### Thus

The Naive Bayes Rule minimizes the error.

## Finally

### Similarly

$$P_e = P(\omega_2) - \int_{R_2} [p(\omega_2|x) - p(\omega_1|x)] p(x) dx \quad (15)$$

### Thus

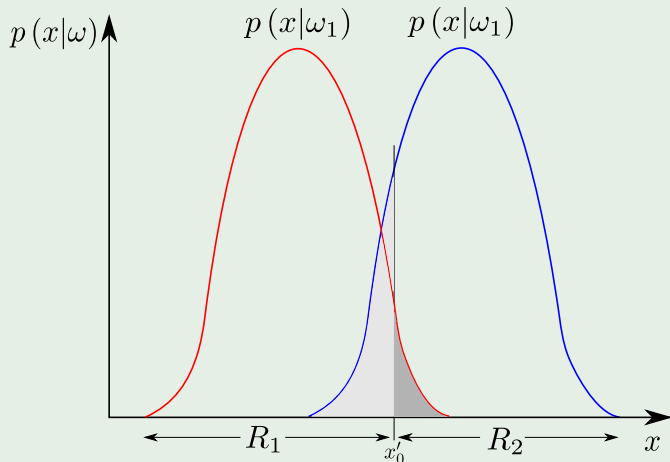
The probability of error is minimized at the region of space in which  $R_2 : P(\omega_2|x) > P(\omega_1|x)$ .

### Thus

The Naive Bayes Rule minimizes the error.

After all!!!

If you choose any other  $x'_0$



# Outline

1

## Introduction

- Supervised Learning
- Handling Noise in Classification
- Models of Classification
- **Naive Bayes**
  - Examples
  - The Naive Bayes Model
  - **The Multi-Class Case**

2

## Discriminant Functions and Decision Surfaces

- Introduction
- Gaussian Distribution
- Influence of the Covariance  $\Sigma$
- Example
- Maximum Likelihood Principle
- Maximum Likelihood on a Gaussian
- Some Remarks

3

## Introduction

- A first solution for the Maximum A Posteriori (MAP)
- Maximum Likelihood Vs Maximum A Posteriori
- Properties of the MAP

4

## Exercises

- Some Stuff you can try

For  $M$  classes  $\omega_1, \omega_2, \dots, \omega_M$

We have that vector  $\mathbf{x}$  is in  $\omega_i$

$$P(\omega_i|\mathbf{x}) > P(\omega_j|\mathbf{x}) \quad \forall j \neq i \quad (16)$$

Something Notable

It turns out that such a choice also minimizes the classification error probability.

For  $M$  classes  $\omega_1, \omega_2, \dots, \omega_M$

We have that vector  $\mathbf{x}$  is in  $\omega_i$

$$P(\omega_i|\mathbf{x}) > P(\omega_j|\mathbf{x}) \quad \forall j \neq i \quad (16)$$

### Something Notable

It turns out that such a choice also minimizes the classification error probability.

# Outline

1

## Introduction

- Supervised Learning
- Handling Noise in Classification
- Models of Classification
- Naive Bayes
  - Examples
  - The Naive Bayes Model
  - The Multi-Class Case

2

## Discriminant Functions and Decision Surfaces

- **Introduction**
- Gaussian Distribution
- Influence of the Covariance  $\Sigma$
- Example
- Maximum Likelihood Principle
- Maximum Likelihood on a Gaussian
- Some Remarks

3

## Introduction

- A first solution for the Maximum A Posteriori (MAP)
- Maximum Likelihood Vs Maximum A Posteriori
- Properties of the MAP

4

## Exercises

- Some Stuff you can try



## Decision Surface

Because the  $R_1$  and  $R_2$  are contiguous

The separating surface between both of them is described by

$$P(\omega_1|x) - P(\omega_2|x) = 0 \quad (17)$$

Thus, we define the decision function as

$$g_{12}(x) = P(\omega_1|x) - P(\omega_2|x) = 0 \quad (18)$$

## Decision Surface

Because the  $R_1$  and  $R_2$  are contiguous

The separating surface between both of them is described by

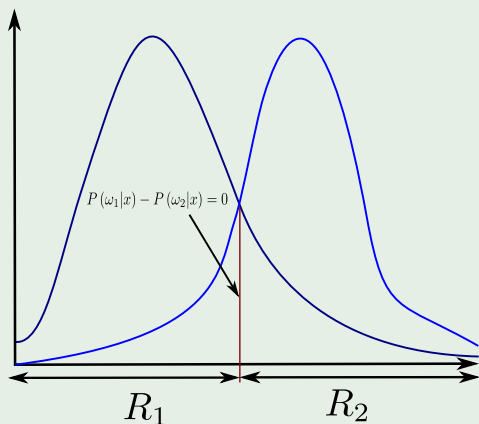
$$P(\omega_1|x) - P(\omega_2|x) = 0 \quad (17)$$

Thus, we define the decision function as

$$g_{12}(x) = P(\omega_1|x) - P(\omega_2|x) = 0 \quad (18)$$

# Which decision function for the Naive Bayes

A single number in this case



## In general

### First

Instead of working with probabilities, we work with an equivalent function of them  $g_i(\mathbf{x}) = f(P(\omega_i|\mathbf{x}))$ .

- Classic Example the Monotonically increasing  $f(P(\omega_i|\mathbf{x})) = \ln P(\omega_i|\mathbf{x})$ .

## In general

### First

Instead of working with probabilities, we work with an equivalent function of them  $g_i(\mathbf{x}) = f(P(\omega_i|\mathbf{x}))$ .

- Classic Example the Monotonically increasing  $f(P(\omega_i|\mathbf{x})) = \ln P(\omega_i|\mathbf{x})$ .

The decision is now:

classify  $\mathbf{x}$  in  $\omega_i$  if  $g_i(\mathbf{x}) > g_j(\mathbf{x}) \forall j \neq i$ .

## In general

### First

Instead of working with probabilities, we work with an equivalent function of them  $g_i(\mathbf{x}) = f(P(\omega_i|\mathbf{x}))$ .

- Classic Example the Monotonically increasing  $f(P(\omega_i|\mathbf{x})) = \ln P(\omega_i|\mathbf{x})$ .

### The decision test is now

classify  $\mathbf{x}$  in  $\omega_i$  if  $g_i(\mathbf{x}) > g_j(\mathbf{x}) \forall j \neq i$ .

$$g_{ij}(\mathbf{x}) = g_i(\mathbf{x}) - g_j(\mathbf{x}) \quad i, j = 1, 2, \dots, M \quad i \neq j$$

## In general

### First

Instead of working with probabilities, we work with an equivalent function of them  $g_i(\mathbf{x}) = f(P(\omega_i|\mathbf{x}))$ .

- Classic Example the Monotonically increasing  $f(P(\omega_i|\mathbf{x})) = \ln P(\omega_i|\mathbf{x})$ .

The decision test is now

classify  $\mathbf{x}$  in  $\omega_i$  if  $g_i(\mathbf{x}) > g_j(\mathbf{x}) \forall j \neq i$ .

The decision surfaces, separating contiguous regions, are described by

$$g_{ij}(\mathbf{x}) = g_i(\mathbf{x}) - g_j(\mathbf{x}) \quad i, j = 1, 2, \dots, M \quad i \neq j$$

# Outline

1

## Introduction

- Supervised Learning
- Handling Noise in Classification
- Models of Classification
- Naive Bayes
  - Examples
  - The Naive Bayes Model
  - The Multi-Class Case

2

## Discriminant Functions and Decision Surfaces

- Introduction
- **Gaussian Distribution**
- Influence of the Covariance  $\Sigma$
- Example
- Maximum Likelihood Principle
- Maximum Likelihood on a Gaussian
- Some Remarks

3

## Introduction

- A first solution for the Maximum A Posteriori (MAP)
- Maximum Likelihood Vs Maximum A Posteriori
- Properties of the MAP

4

## Exercises

- Some Stuff you can try



# Gaussian Distribution

We can use the Gaussian distribution

$$p(\mathbf{x}|\omega_i) = \frac{1}{(2\pi)^{l/2} |\Sigma_i|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^T \Sigma_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) \right\} \quad (19)$$

Example

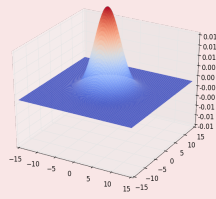
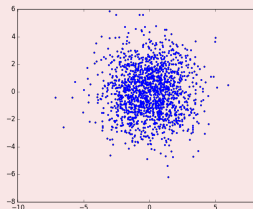
# Gaussian Distribution

We can use the Gaussian distribution

$$p(\mathbf{x}|\omega_i) = \frac{1}{(2\pi)^{l/2} |\Sigma_i|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^T \Sigma_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) \right\} \quad (19)$$

Example

$$\Sigma = \begin{bmatrix} 3 & 0 \\ 0 & 3 \end{bmatrix}$$



# Some Properties

## About $\Sigma$

It is the covariance matrix between variables.

### Thus

- It is positive semi-definite.
- Symmetric.
- The inverse exists.

# Some Properties

## About $\Sigma$

It is the covariance matrix between variables.

## Thus

- It is positive semi-definite.
- Symmetric.
- The inverse exists.

# Outline

1

## Introduction

- Supervised Learning
- Handling Noise in Classification
- Models of Classification
- Naive Bayes
  - Examples
  - The Naive Bayes Model
  - The Multi-Class Case

2

## Discriminant Functions and Decision Surfaces

- Introduction
- Gaussian Distribution
- **Influence of the Covariance  $\Sigma$**
- Example
- Maximum Likelihood Principle
- Maximum Likelihood on a Gaussian
- Some Remarks

3

## Introduction

- A first solution for the Maximum A Posteriori (MAP)
- Maximum Likelihood Vs Maximum A Posteriori
- Properties of the MAP

4

## Exercises

- Some Stuff you can try

## Influence of the Covariance $\Sigma$

Look at the following Covariance

$$\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

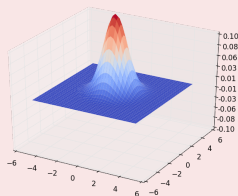
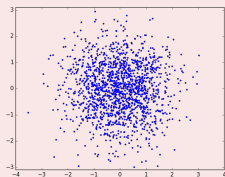
It's simple the unit Gaussian with mean  $\mu$

# Influence of the Covariance $\Sigma$

Look at the following Covariance

$$\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

It's simple the unit Gaussian with mean  $\mu$



## The Covariance $\Sigma$ as a Rotation

Look at the following Covariance

$$\Sigma = \begin{bmatrix} 16 & 0 \\ 0 & 1 \end{bmatrix}$$

Actually, it flattens the circle through the  $x$ -axis

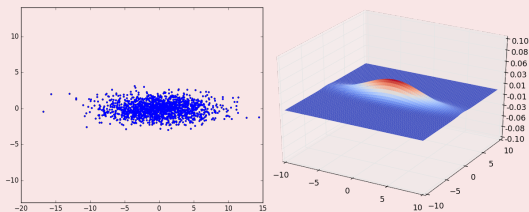


# The Covariance $\Sigma$ as a Rotation

Look at the following Covariance

$$\Sigma = \begin{bmatrix} 16 & 0 \\ 0 & 1 \end{bmatrix}$$

Actually, it flattens the circle through the  $x$  -  $axis$



## Influence of the Covariance $\Sigma$

Look at the following Covariance

$$\Sigma_a = R\Sigma_b R^T \text{ with } R = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}$$

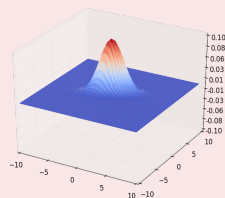
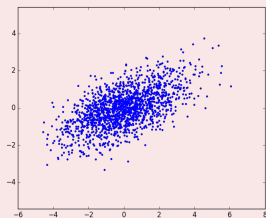
It allows to rotate the axes

# Influence of the Covariance $\Sigma$

Look at the following Covariance

$$\Sigma_a = R\Sigma_b R^T \text{ with } R = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}$$

It allows to rotate the axes



## Now For Two Classes

Then, we use the following trick for two Classes  $i = 1, 2$

We know that the pdf of correct classification is

$$p(x, \omega_i) = p(x|\omega_i) P(\omega_i)!!!$$

Wait

It is possible to generate the following decision function:

$$g_i(x) = \ln [p(x|\omega_i) P(\omega_i)] = \ln p(x|\omega_i) + \ln P(\omega_i) \quad (20)$$

Thus

$$g_i(x) = -\frac{1}{2} (x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i) + \ln P(\omega_i) + c_i \quad (21)$$

## Now For Two Classes

Then, we use the following trick for two Classes  $i = 1, 2$

We know that the pdf of correct classification is

$$p(x, \omega_i) = p(x|\omega_i) P(\omega_i)!!!$$

Thus

It is possible to generate the following decision function:

$$g_i(\mathbf{x}) = \ln [p(x|\omega_i) P(\omega_i)] = \ln p(x|\omega_i) + \ln P(\omega_i) \quad (20)$$

Thus

$$g_i(\mathbf{x}) = -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) + \ln P(\omega_i) + c_i \quad (21)$$

## Now For Two Classes

Then, we use the following trick for two Classes  $i = 1, 2$

We know that the pdf of correct classification is

$$p(x, \omega_1) = p(x|\omega_i) P(\omega_i)!!!$$

Thus

It is possible to generate the following decision function:

$$g_i(\mathbf{x}) = \ln [p(x|\omega_i) P(\omega_i)] = \ln p(x|\omega_i) + \ln P(\omega_i) \quad (20)$$

Thus

$$g_i(\mathbf{x}) = -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) + \ln P(\omega_i) + c_i \quad (21)$$

# Outline

1

## Introduction

- Supervised Learning
- Handling Noise in Classification
- Models of Classification
- Naive Bayes
  - Examples
  - The Naive Bayes Model
  - The Multi-Class Case

2

## Discriminant Functions and Decision Surfaces

- Introduction
- Gaussian Distribution
- Influence of the Covariance  $\Sigma$
- **Example**
- Maximum Likelihood Principle
- Maximum Likelihood on a Gaussian
- Some Remarks

3

## Introduction

- A first solution for the Maximum A Posteriori (MAP)
- Maximum Likelihood Vs Maximum A Posteriori
- Properties of the MAP

4

## Exercises

- Some Stuff you can try

## We can work one of the possible decision surfaces

Assume first that  $\Sigma_i = \sigma^2 I$

- The features are statistically independent
- Each feature has the same variance



## We can work one of the possible decision surfaces

Assume first that  $\Sigma_i = \sigma^2 I$

- The features are statistically independent
- Each feature has the same variance

Therefore:

- The samples fall in equal size spherical clusters!!!
- Each Cluster centered at mean vector  $\mu_i$ .

## We can work one of the possible decision surfaces

Assume first that  $\Sigma_i = \sigma^2 I$

- The features are statistically independent
- Each feature has the same variance

Therefore

- The samples fall in equal size spherical clusters!!!
- Each Cluster centered at mean vector  $\mu_i$ .

## We can work one of the possible decision surfaces

Assume first that  $\Sigma_i = \sigma^2 I$

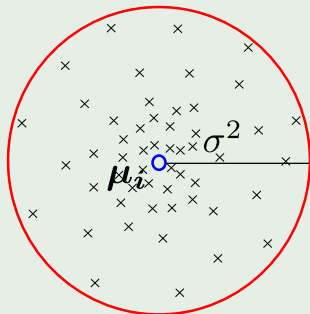
- The features are statistically independent
- Each feature has the same variance

Therefore

- The samples fall in equal size spherical clusters!!!
- Each Cluster centered at mean vector  $\mu_i$ .

## For Example

We have



## Now

We have that

$$|\Sigma_i| = \sigma^{2d} \text{ and } \Sigma_i^{-1} = \left(\frac{1}{\sigma^2}\right) I$$

Something Notable

- Gaussian Multivariate function after the log

$$g_i(\mathbf{x}) = -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^T \Sigma_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) + \ln P(\omega_i) - \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_i|$$

The term  $-\frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_i|$

It is unimportant therefore it can be ignored!!!

## Now

We have that

$$|\Sigma_i| = \sigma^{2d} \text{ and } \Sigma_i^{-1} = \left(\frac{1}{\sigma^2}\right) I$$

Something Notable

- Gaussian Multivariate function after the log

$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \Sigma_i^{-1}(\mathbf{x} - \boldsymbol{\mu}_i) + \ln P(\omega_i) - \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_i|$$

This term

is unimportant therefore it can be ignored!!!

## Now

We have that

$$|\Sigma_i| = \sigma^{2d} \text{ and } \Sigma_i^{-1} = \left(\frac{1}{\sigma^2}\right) I$$

Something Notable

- Gaussian Multivariate function after the log

$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \Sigma_i^{-1}(\mathbf{x} - \boldsymbol{\mu}_i) + \ln P(\omega_i) - \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_i|$$

The term  $-\frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_i|$

It is unimportant therefore it can be ignored!!!

Then

We have the following discriminant functions

$$g_i(\mathbf{x}) = -\frac{\underbrace{\|\mathbf{x} - \boldsymbol{\mu}_i\|^2}_{(\mathbf{x} - \boldsymbol{\mu}_i)^T (\mathbf{x} - \boldsymbol{\mu}_i)}}{2\sigma^2} + \ln P(\omega_i) \quad (22)$$

Then, we have that

$$g_i(\mathbf{x}) = -\frac{1}{2\sigma^2} \left[ \mathbf{x}^T \mathbf{x} - 2\boldsymbol{\mu}_i^T \mathbf{x} + \boldsymbol{\mu}_i^T \boldsymbol{\mu}_i \right] + \ln P(\omega_i)$$



Then

We have the following discriminant functions

$$g_i(\mathbf{x}) = -\frac{\underbrace{\|\mathbf{x} - \boldsymbol{\mu}_i\|^2}_{(\mathbf{x} - \boldsymbol{\mu}_i)^T (\mathbf{x} - \boldsymbol{\mu}_i)}}{2\sigma^2} + \ln P(\omega_i) \quad (22)$$

Then, we have that

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

We can then...

Do you notice that  $\mathbf{x}^T \mathbf{x}$  is actually the same for all  $g_i$ ?

Then, we can ignore that term thus, we get

$$g_i(\mathbf{x}) = \frac{1}{\sigma^2} \underbrace{\boldsymbol{\mu}_i^T}_{\mathbf{w}_i^T} \mathbf{x} - \frac{1}{2\sigma^2} \underbrace{\boldsymbol{\mu}_i^T \boldsymbol{\mu}_i}_{w_{i0}} + \ln P(\omega_i)$$

Or if you want

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

We can then...

Do you notice that  $\mathbf{x}^T \mathbf{x}$  is actually the same for all  $g_i$ ?

Then, we can ignore that term thus, we get

$$g_i(\mathbf{x}) = \frac{1}{\sigma^2} \underbrace{\boldsymbol{\mu}_i^T \mathbf{x}}_{\mathbf{w}_i^T} - \frac{1}{2\sigma^2} \underbrace{\boldsymbol{\mu}_i^T \boldsymbol{\mu}_i}_{w_{i0}} + \ln P(\omega_i)$$

Or if you want

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

# Outline

1

## Introduction

- Supervised Learning
- Handling Noise in Classification
- Models of Classification
- Naive Bayes
  - Examples
  - The Naive Bayes Model
  - The Multi-Class Case

2

## Discriminant Functions and Decision Surfaces

- Introduction
- Gaussian Distribution
- Influence of the Covariance  $\Sigma$
- Example
- **Maximum Likelihood Principle**
- Maximum Likelihood on a Gaussian
- Some Remarks

3

## Introduction

- A first solution for the Maximum A Posteriori (MAP)
- Maximum Likelihood Vs Maximum A Posteriori
- Properties of the MAP

4

## Exercises

- Some Stuff you can try

Given a series of classes  $\omega_1, \omega_2, \dots, \omega_M$

We assume for each class  $\omega_j$

The samples are drawn independently according to the probability law  $p(\mathbf{x}|\omega_j)$

We call these samples as

i.i.d. — independent identically distributed random variables.

We assume in addition

$p(\mathbf{x}|\omega_j)$  has a known parametric form with vector  $\theta_j$  of parameters.

Given a series of classes  $\omega_1, \omega_2, \dots, \omega_M$

We assume for each class  $\omega_j$

The samples are drawn independently according to the probability law  $p(\mathbf{x}|\omega_j)$

We call those samples as

i.i.d. — independent identically distributed random variables.

We assume in addition

$p(\mathbf{x}|\omega_j)$  has a known parametric form with vector  $\theta_j$  of parameters.

Given a series of classes  $\omega_1, \omega_2, \dots, \omega_M$

We assume for each class  $\omega_j$

The samples are drawn independently according to the probability law  $p(\mathbf{x}|\omega_j)$

We call those samples as

i.i.d. — independent identically distributed random variables.

We assume in addition

$p(\mathbf{x}|\omega_j)$  has a known parametric form with vector  $\theta_j$  of parameters.

Given a series of classes  $\omega_1, \omega_2, \dots, \omega_M$

For example

$$p(\mathbf{x}|\omega_j) \sim N(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) \quad (23)$$

In our case

We will assume that there is no dependence between classes!!!



Given a series of classes  $\omega_1, \omega_2, \dots, \omega_M$

For example

$$p(\mathbf{x}|\omega_j) \sim N(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) \quad (23)$$

In our case

We will assume that there is no dependence between classes!!!

## Now

Suppose that  $\omega_j$  contains  $n$  samples  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$

$$p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n | \theta_j) = \prod_{j=1}^n p(\mathbf{x}_j | \theta_j) \quad (24)$$

We can see then the function  $p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n | \theta_j)$  as a function of

$$L(\theta_j) = \prod_{j=1}^n p(\mathbf{x}_j | \theta_j) \quad (25)$$

Now

Suppose that  $\omega_j$  contains  $n$  samples  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$

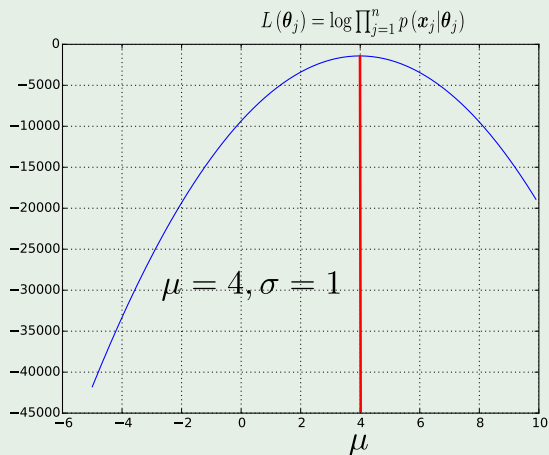
$$p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n | \boldsymbol{\theta}_j) = \prod_{j=1}^n p(\mathbf{x}_j | \boldsymbol{\theta}_j) \quad (24)$$

We can see then the function  $p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n | \boldsymbol{\theta}_j)$  as a function of

$$L(\boldsymbol{\theta}_j) = \prod_{j=1}^n p(\mathbf{x}_j | \boldsymbol{\theta}_j) \quad (25)$$

# Example

$$L(\theta_j) = \log \prod_{j=1}^n p(\mathbf{x}_j | \theta_j)$$



# Outline

1

## Introduction

- Supervised Learning
- Handling Noise in Classification
- Models of Classification
- Naive Bayes
  - Examples
  - The Naive Bayes Model
  - The Multi-Class Case

2

## Discriminant Functions and Decision Surfaces

- Introduction
- Gaussian Distribution
- Influence of the Covariance  $\Sigma$
- Example
- Maximum Likelihood Principle
- **Maximum Likelihood on a Gaussian**
- Some Remarks

3

## Introduction

- A first solution for the Maximum A Posteriori (MAP)
- Maximum Likelihood Vs Maximum A Posteriori
- Properties of the MAP

4

## Exercises

- Some Stuff you can try

## Maximum Likelihood on a Gaussian

Then, using the log!!!

$$\ln L(\omega_i) = -\frac{n}{2} \ln |\Sigma_i| - \frac{1}{2} \left[ \sum_{j=1}^n (\mathbf{x}_j - \boldsymbol{\mu}_i)^T \Sigma_i^{-1} (\mathbf{x}_j - \boldsymbol{\mu}_i) \right] + c_2 \quad (26)$$

We know that

$$\frac{d\mathbf{x}^T A \mathbf{x}}{d\mathbf{x}} = A\mathbf{x} + A^T \mathbf{x}, \quad \frac{dA\mathbf{x}}{d\mathbf{x}} = A \quad (27)$$

Thus, we expand equation 26

$$-\frac{n}{2} \ln |\Sigma_i| - \frac{1}{2} \sum_{j=1}^n \left[ \mathbf{x}_j^T \Sigma_i^{-1} \mathbf{x}_j - 2\mathbf{x}_j^T \Sigma_i^{-1} \boldsymbol{\mu}_i + \boldsymbol{\mu}_i^T \Sigma_i^{-1} \boldsymbol{\mu}_i \right] + c_2 \quad (28)$$

## Maximum Likelihood on a Gaussian

Then, using the log!!!

$$\ln L(\omega_i) = -\frac{n}{2} \ln |\Sigma_i| - \frac{1}{2} \left[ \sum_{j=1}^n (\mathbf{x}_j - \boldsymbol{\mu}_i)^T \Sigma_i^{-1} (\mathbf{x}_j - \boldsymbol{\mu}_i) \right] + c_2 \quad (26)$$

We know that

$$\frac{d\mathbf{x}^T A \mathbf{x}}{d\mathbf{x}} = A\mathbf{x} + A^T \mathbf{x}, \quad \frac{dA\mathbf{x}}{d\mathbf{x}} = A \quad (27)$$

Thus, we expand equation 26

$$-\frac{n}{2} \ln |\Sigma_i| - \frac{1}{2} \sum_{j=1}^n \left[ \mathbf{x}_j^T \Sigma_i^{-1} \mathbf{x}_j - 2\mathbf{x}_j^T \Sigma_i^{-1} \boldsymbol{\mu}_i + \boldsymbol{\mu}_i^T \Sigma_i^{-1} \boldsymbol{\mu}_i \right] + c_2 \quad (28)$$

## Maximum Likelihood on a Gaussian

Then, using the log!!!

$$\ln L(\omega_i) = -\frac{n}{2} \ln |\Sigma_i| - \frac{1}{2} \left[ \sum_{j=1}^n (\mathbf{x}_j - \boldsymbol{\mu}_i)^T \Sigma_i^{-1} (\mathbf{x}_j - \boldsymbol{\mu}_i) \right] + c_2 \quad (26)$$

We know that

$$\frac{d\mathbf{x}^T A \mathbf{x}}{d\mathbf{x}} = A\mathbf{x} + A^T \mathbf{x}, \quad \frac{dA\mathbf{x}}{d\mathbf{x}} = A \quad (27)$$

Thus, we expand equation 26

$$-\frac{n}{2} \ln |\Sigma_i| - \frac{1}{2} \sum_{j=1}^n \left[ \mathbf{x}_j^T \Sigma_i^{-1} \mathbf{x}_j - 2\mathbf{x}_j^T \Sigma_i^{-1} \boldsymbol{\mu}_i + \boldsymbol{\mu}_i^T \Sigma_i^{-1} \boldsymbol{\mu}_i \right] + c_2 \quad (28)$$



# Maximum Likelihood

Then

$$\frac{\partial \ln L(\omega_i)}{\partial \mu_i} = \sum_{j=1}^n \Sigma_i^{-1} (\mathbf{x}_j - \mu_i) = 0$$

$$n \Sigma_i^{-1} \left[ -\mu_i + \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j \right] = 0$$

$$\hat{\mu}_i = \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j$$

# Maximum Likelihood

Then

$$\frac{\partial \ln L(\omega_i)}{\partial \boldsymbol{\mu}_i} = \sum_{j=1}^n \boldsymbol{\Sigma}_i^{-1} (\mathbf{x}_j - \boldsymbol{\mu}_i) = 0$$
$$n \boldsymbol{\Sigma}_i^{-1} \left[ -\boldsymbol{\mu}_i + \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j \right] = 0$$

$$\hat{\boldsymbol{\mu}}_i = \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j$$

# Maximum Likelihood

Then

$$\frac{\partial \ln L(\omega_i)}{\partial \boldsymbol{\mu}_i} = \sum_{j=1}^n \boldsymbol{\Sigma}_i^{-1} (\mathbf{x}_j - \boldsymbol{\mu}_i) = 0$$
$$n \boldsymbol{\Sigma}_i^{-1} \left[ -\boldsymbol{\mu}_i + \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j \right] = 0$$
$$\hat{\boldsymbol{\mu}}_i = \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j$$

# Maximum Likelihood

Then, we derive with respect to  $\Sigma_i$

For this we use the following tricks:

- 1  $\frac{\partial \log|\Sigma|}{\partial \Sigma^{-1}} = -\frac{1}{|\Sigma|} \cdot |\Sigma| (\Sigma)^T = -\Sigma$
- 2  $\frac{\partial \text{Tr}[AB]}{\partial A} = \frac{\partial \text{Tr}[BA]}{\partial A} = B^T$
- 3 Trace(of a number)=the number
- 4  $\text{Tr}(A^T B) = \text{Tr}(B A^T)$

Thus

$$f(\Sigma_i) = -\frac{n}{2} \ln |\Sigma_i| - \frac{1}{2} \sum_{j=1}^n [(\mathbf{x}_j - \boldsymbol{\mu}_i)^T \Sigma_i^{-1} (\mathbf{x}_j - \boldsymbol{\mu}_i)] + c_1 \quad (29)$$

# Maximum Likelihood

Thus

$$f(\Sigma_i) = -\frac{n}{2} \ln |\Sigma_i| - \frac{1}{2} \sum_{j=1}^n \left[ \text{Trace} \left\{ (\mathbf{x}_j - \boldsymbol{\mu}_i)^T \Sigma_i^{-1} (\mathbf{x}_j - \boldsymbol{\mu}_i) \right\} \right] + c_1 \quad (30)$$

Tricks!!!

$$f(\Sigma_i) = -\frac{n}{2} \ln |\Sigma_i| - \frac{1}{2} \sum_{j=1}^n \left[ \text{Trace} \left\{ \Sigma_i^{-1} (\mathbf{x}_j - \boldsymbol{\mu}_i) (\mathbf{x}_j - \boldsymbol{\mu}_i)^T \right\} \right] + c_1 \quad (31)$$

# Maximum Likelihood

Thus

$$f(\Sigma_i) = -\frac{n}{2} \ln |\Sigma_i| - \frac{1}{2} \sum_{j=1}^n \left[ \text{Trace} \left\{ (\mathbf{x}_j - \boldsymbol{\mu}_i)^T \Sigma_i^{-1} (\mathbf{x}_j - \boldsymbol{\mu}_i) \right\} \right] + c_1 \quad (30)$$

Tricks!!!

$$f(\Sigma_i) = -\frac{n}{2} \ln |\Sigma_i| - \frac{1}{2} \sum_{j=1}^n \left[ \text{Trace} \left\{ \Sigma_i^{-1} (\mathbf{x}_j - \boldsymbol{\mu}_i) (\mathbf{x}_j - \boldsymbol{\mu}_i)^T \right\} \right] + c_1 \quad (31)$$

# Maximum Likelihood

Derivative with respect to  $\Sigma$

$$\frac{\partial f(\Sigma_i)}{\partial \Sigma_i} = \frac{n}{2} \Sigma_i - \frac{1}{2} \sum_{j=1}^n [(\mathbf{x}_j - \boldsymbol{\mu}_i)(\mathbf{x}_j - \boldsymbol{\mu}_i)^T]^T \quad (32)$$

Thus, when making it equal to zero

$$\hat{\Sigma}_i = \frac{1}{n} \sum_{j=1}^n (\mathbf{x}_j - \boldsymbol{\mu}_i)(\mathbf{x}_j - \boldsymbol{\mu}_i)^T \quad (33)$$

# Maximum Likelihood

Derivative with respect to  $\Sigma$

$$\frac{\partial f(\Sigma_i)}{\partial \Sigma_i} = \frac{n}{2} \Sigma_i - \frac{1}{2} \sum_{j=1}^n [(\mathbf{x}_j - \boldsymbol{\mu}_i)(\mathbf{x}_j - \boldsymbol{\mu}_i)^T]^T \quad (32)$$

Thus, when making it equal to zero

$$\hat{\Sigma}_i = \frac{1}{n} \sum_{j=1}^n (\mathbf{x}_j - \boldsymbol{\mu}_i)(\mathbf{x}_j - \boldsymbol{\mu}_i)^T \quad (33)$$



# Therefore

## Step 1 - Assume a Gaussian Distribution over each class

- The So Called Model Selection

## Step 2

- Adjust the Gaussian Distribution, for each class, using the previous Maximum Likelihood

## Step 3

$$R_1 : P(\omega_1|x) > P(\omega_2|x)$$

$$R_2 : P(\omega_2|x) > P(\omega_1|x)$$

# Therefore

## Step 1 - Assume a Gaussian Distribution over each class

- The So Called Model Selection

## Step 2

- Adjust the Gaussian Distribution, for each class, using the previous Maximum Likelihood

Step 3

$$R_1 : P(\omega_1|x) > P(\omega_2|x)$$

$$R_2 : P(\omega_2|x) > P(\omega_1|x)$$

# Therefore

## Step 1 - Assume a Gaussian Distribution over each class

- The So Called Model Selection

## Step 2

- Adjust the Gaussian Distribution, for each class, using the previous Maximum Likelihood

## Step 3

$$R_1 : P(\omega_1|x) > P(\omega_2|x)$$

$$R_2 : P(\omega_2|x) > P(\omega_1|x)$$

# Outline

1

## Introduction

- Supervised Learning
- Handling Noise in Classification
- Models of Classification
- Naive Bayes
  - Examples
  - The Naive Bayes Model
  - The Multi-Class Case

2

## Discriminant Functions and Decision Surfaces

- Introduction
- Gaussian Distribution
- Influence of the Covariance  $\Sigma$
- Example
- Maximum Likelihood Principle
- Maximum Likelihood on a Gaussian
- **Some Remarks**

3

## Introduction

- A first solution for the Maximum A Posteriori (MAP)
- Maximum Likelihood Vs Maximum A Posteriori
- Properties of the MAP

4

## Exercises

- Some Stuff you can try

## In the case of Bayesian Model

We have

$$P(Y_n = i | \mathbf{x}_n) = \frac{P(\mathbf{x}_n | Y_n = i) P(Y_n = i)}{P(\mathbf{x}_n)}$$

In the Generative Model

- We model two distribution  $P(\mathbf{x}_n | Y_n = 1)$  and  $P(Y_n = i)$

In the Discriminative Model

- We model a single distribution  $P(Y_n = i)$

## In the case of Bayesian Model

We have

$$P(Y_n = i | \mathbf{x}_n) = \frac{P(\mathbf{x}_n | Y_n = i) P(Y_n = i)}{P(\mathbf{x}_n)}$$

In the Generative Model

- We model two distribution  $P(\mathbf{x}_n | Y_n = 1)$  and  $P(Y_n = i)$

In the Discriminative Model

- We model a single distribution  $P(Y_n = i)$

## In the case of Bayesian Model

We have

$$P(Y_n = i | \mathbf{x}_n) = \frac{P(\mathbf{x}_n | Y_n = i) P(Y_n = i)}{P(\mathbf{x}_n)}$$

In the Generative Model

- We model two distribution  $P(\mathbf{x}_n | Y_n = 1)$  and  $P(Y_n = i)$

In the Discriminative Model

- We model a single distribution  $P(Y_n = i)$

# Therefore

## We have

- In the Generative Model, we discover the distribution from  $X$  and  $Y$

## Therefore

Although discriminative models tend to be faster and less complex, they cannot model the joint  $P(X, Y)$ .

## Thus

- We have a decision problem
  - ▶ Do we want to know the joint distribution?



## Therefore

### We have

- In the Generative Model, we discover the distribution from  $X$  and  $Y$

### Therefore

Although discriminative models tend to be faster and less complex, they cannot model the joint  $P(X, Y)$ .

### Issues

- We have a decision problem
  - ▶ Do we want to know the joint distribution?

# Therefore

## We have

- In the Generative Model, we discover the distribution from  $X$  and  $Y$

## Therefore

Although discriminative models tend to be faster and less complex, they cannot model the joint  $P(X, Y)$ .

## Thus

- We have a decision problem
  - ▶ Do we want to know the joint distribution?

# Outline

1

## Introduction

- Supervised Learning
- Handling Noise in Classification
- Models of Classification
- Naive Bayes
  - Examples
  - The Naive Bayes Model
  - The Multi-Class Case

2

## Discriminant Functions and Decision Surfaces

- Introduction
- Gaussian Distribution
- Influence of the Covariance  $\Sigma$
- Example
- Maximum Likelihood Principle
- Maximum Likelihood on a Gaussian
- Some Remarks

3

## Introduction

- **A first solution for the Maximum A Posteriori (MAP)**
- Maximum Likelihood Vs Maximum A Posteriori
- Properties of the MAP

4

## Exercises

- Some Stuff you can try

# Introduction

We go back to the Bayesian Rule

$$p(\Theta|\mathcal{X}) = \frac{p(\mathcal{X}|\Theta)p(\Theta)}{p(\mathcal{X})} \quad (34)$$

We now seek that value for  $\Theta$ , called  $\Theta_{\text{MAP}}$ .

It allows to maximize the posterior  $p(\Theta|\mathcal{X})$ .

# Introduction

We go back to the Bayesian Rule

$$p(\Theta|\mathcal{X}) = \frac{p(\mathcal{X}|\Theta)p(\Theta)}{p(\mathcal{X})} \quad (34)$$

We now seek that value for  $\Theta$ , called  $\hat{\Theta}_{MAP}$

It allows to maximize the posterior  $p(\Theta|\mathcal{X})$

## Development of the solution

We look to maximize  $\hat{\Theta}_{MAP}$

$$\begin{aligned}\hat{\Theta}_{MAP} &= \underset{\Theta}{\operatorname{argmax}} p(\Theta|\mathcal{X}) \\ &= \underset{\Theta}{\operatorname{argmax}} \frac{p(\mathcal{X}|\Theta) p(\Theta)}{P(\mathcal{X})} \\ &\approx \underset{\Theta}{\operatorname{argmax}} p(\mathcal{X}|\Theta) p(\Theta) \\ &= \underset{\Theta}{\operatorname{argmax}} \prod_{x_i \in \mathcal{X}} p(x_i|\Theta) p(\Theta)\end{aligned}$$

$P(\mathcal{X})$  can be removed because it has no functional relation with  $\Theta$ .

## Development of the solution

We look to maximize  $\hat{\Theta}_{MAP}$

$$\begin{aligned}\hat{\Theta}_{MAP} &= \operatorname{argmax}_{\Theta} p(\Theta|\mathcal{X}) \\ &= \operatorname{argmax}_{\Theta} \frac{p(\mathcal{X}|\Theta)p(\Theta)}{P(\mathcal{X})} \\ &\approx \operatorname{argmax}_{\Theta} p(\mathcal{X}|\Theta)p(\Theta) \\ &= \operatorname{argmax}_{\Theta} \prod_{x_i \in \mathcal{X}} p(x_i|\Theta)p(\Theta)\end{aligned}$$

$P(\mathcal{X})$  can be removed because it has no functional relation with  $\Theta$ .

## Development of the solution

We look to maximize  $\hat{\Theta}_{MAP}$

$$\begin{aligned}\hat{\Theta}_{MAP} &= \underset{\Theta}{\operatorname{argmax}} p(\Theta|\mathcal{X}) \\ &= \underset{\Theta}{\operatorname{argmax}} \frac{p(\mathcal{X}|\Theta) p(\Theta)}{P(\mathcal{X})} \\ &\approx \underset{\Theta}{\operatorname{argmax}} p(\mathcal{X}|\Theta) p(\Theta) \\ &= \underset{\Theta}{\operatorname{argmax}} \prod_{x_i \in \mathcal{X}} p(x_i|\Theta) p(\Theta)\end{aligned}$$

$P(\mathcal{X})$  can be removed because it has no functional relation with  $\Theta$ .



## Development of the solution

We look to maximize  $\hat{\Theta}_{MAP}$

$$\begin{aligned}\hat{\Theta}_{MAP} &= \operatorname{argmax}_{\Theta} p(\Theta|\mathcal{X}) \\ &= \operatorname{argmax}_{\Theta} \frac{p(\mathcal{X}|\Theta) p(\Theta)}{P(\mathcal{X})} \\ &\approx \operatorname{argmax}_{\Theta} p(\mathcal{X}|\Theta) p(\Theta) \\ &= \operatorname{argmax}_{\Theta} \prod_{x_i \in \mathcal{X}} p(x_i|\Theta) p(\Theta)\end{aligned}$$

$P(\mathcal{X})$  can be removed because it has no functional relation with  $\Theta$ .

## Development of the solution

We look to maximize  $\hat{\Theta}_{MAP}$

$$\begin{aligned}\hat{\Theta}_{MAP} &= \underset{\Theta}{\operatorname{argmax}} p(\Theta|\mathcal{X}) \\ &= \underset{\Theta}{\operatorname{argmax}} \frac{p(\mathcal{X}|\Theta)p(\Theta)}{P(\mathcal{X})} \\ &\approx \underset{\Theta}{\operatorname{argmax}} p(\mathcal{X}|\Theta)p(\Theta) \\ &= \underset{\Theta}{\operatorname{argmax}} \prod_{x_i \in \mathcal{X}} p(x_i|\Theta)p(\Theta)\end{aligned}$$

$P(\mathcal{X})$  can be removed because it has no functional relation with  $\Theta$ .

We can make this easier

Use logarithms

$$\hat{\Theta}_{MAP} = \operatorname{argmax}_{\Theta} \left[ \sum_{x_i \in \mathcal{X}} \log p(x_i | \Theta) + \log p(\Theta) \right] \quad (35)$$

# Outline

1

## Introduction

- Supervised Learning
- Handling Noise in Classification
- Models of Classification
- Naive Bayes
  - Examples
  - The Naive Bayes Model
  - The Multi-Class Case

2

## Discriminant Functions and Decision Surfaces

- Introduction
- Gaussian Distribution
- Influence of the Covariance  $\Sigma$
- Example
- Maximum Likelihood Principle
- Maximum Likelihood on a Gaussian
- Some Remarks

3

## Introduction

- A first solution for the Maximum A Posteriori (MAP)
- **Maximum Likelihood Vs Maximum A Posteriori**
- Properties of the MAP

4

## Exercises

- Some Stuff you can try

# What can we do?

We can specify a distribution

Then, learn the parameters

Remember the Bayesian Rule

$$p(\theta|\mathcal{X}) = \frac{p(\mathcal{X}|\theta)p(\theta)}{p(\mathcal{X})} \quad (36)$$

We seek that value for  $\theta$ , called  $\theta_{MAP}$

It allows to maximize the posterior  $p(\theta|\mathcal{X})$

# What can we do?

We can specify a distribution

Then, learn the parameters

Remember the Bayesian Rule

$$p(\Theta|\mathcal{X}) = \frac{p(\mathcal{X}|\Theta)p(\Theta)}{p(\mathcal{X})} \quad (36)$$

We seek that value for  $\Theta$ , called  $\Theta_{ML}$

It allows to maximize the posterior  $p(\Theta|\mathcal{X})$

# What can we do?

We can specify a distribution

Then, learn the parameters

Remember the Bayesian Rule

$$p(\Theta|\mathcal{X}) = \frac{p(\mathcal{X}|\Theta)p(\Theta)}{p(\mathcal{X})} \quad (36)$$

We seek that value for  $\Theta$ , called  $\hat{\Theta}_{MAP}$

It allows to maximize the posterior  $p(\Theta|\mathcal{X})$

Therefore

We can use this idea of maximizing the posterior

To obtain the distribution through the Maximum a Posteriori



## Development of the solution

We look to maximize  $\hat{\Theta}_{MAP}$

$$\begin{aligned}\hat{\Theta}_{MAP} &= \underset{\Theta}{\operatorname{argmax}} p(\Theta|\mathcal{X}) \\ &= \underset{\Theta}{\operatorname{argmax}} \frac{p(\mathcal{X}|\Theta) p(\Theta)}{P(\mathcal{X})} \\ &\approx \underset{\Theta}{\operatorname{argmax}} p(\mathcal{X}|\Theta) p(\Theta) \\ &= \underset{\Theta}{\operatorname{argmax}} \prod_{x_i \in \mathcal{X}} p(x_i|\Theta) p(\Theta)\end{aligned}$$

$P(\mathcal{X})$  can be removed because it has no functional relation with  $\Theta$ .

## Development of the solution

We look to maximize  $\hat{\Theta}_{MAP}$

$$\begin{aligned}\hat{\Theta}_{MAP} &= \operatorname{argmax}_{\Theta} p(\Theta|\mathcal{X}) \\ &= \operatorname{argmax}_{\Theta} \frac{p(\mathcal{X}|\Theta)p(\Theta)}{P(\mathcal{X})} \\ &\approx \operatorname{argmax}_{\Theta} p(\mathcal{X}|\Theta)p(\Theta) \\ &= \operatorname{argmax}_{\Theta} \prod_{x_i \in \mathcal{X}} p(x_i|\Theta)p(\Theta)\end{aligned}$$

$P(\mathcal{X})$  can be removed because it has no functional relation with  $\Theta$ .

## Development of the solution

We look to maximize  $\hat{\Theta}_{MAP}$

$$\begin{aligned}\hat{\Theta}_{MAP} &= \underset{\Theta}{\operatorname{argmax}} p(\Theta|\mathcal{X}) \\ &= \underset{\Theta}{\operatorname{argmax}} \frac{p(\mathcal{X}|\Theta) p(\Theta)}{P(\mathcal{X})} \\ &\approx \underset{\Theta}{\operatorname{argmax}} p(\mathcal{X}|\Theta) p(\Theta) \\ &= \underset{\Theta}{\operatorname{argmax}} \prod_{x_i \in \mathcal{X}} p(x_i|\Theta) p(\Theta)\end{aligned}$$

$P(\mathcal{X})$  can be removed because it has no functional relation with  $\Theta$ .

## Development of the solution

We look to maximize  $\hat{\Theta}_{MAP}$

$$\begin{aligned}\hat{\Theta}_{MAP} &= \underset{\Theta}{\operatorname{argmax}} p(\Theta|\mathcal{X}) \\ &= \underset{\Theta}{\operatorname{argmax}} \frac{p(\mathcal{X}|\Theta) p(\Theta)}{P(\mathcal{X})} \\ &\approx \underset{\Theta}{\operatorname{argmax}} p(\mathcal{X}|\Theta) p(\Theta) \\ &= \underset{\Theta}{\operatorname{argmax}} \prod_{x_i \in \mathcal{X}} p(x_i|\Theta) p(\Theta)\end{aligned}$$

$P(\mathcal{X})$  can be removed because it has no functional relation with  $\Theta$ .

## Development of the solution

We look to maximize  $\hat{\Theta}_{MAP}$

$$\begin{aligned}\hat{\Theta}_{MAP} &= \underset{\Theta}{\operatorname{argmax}} p(\Theta|\mathcal{X}) \\ &= \underset{\Theta}{\operatorname{argmax}} \frac{p(\mathcal{X}|\Theta) p(\Theta)}{P(\mathcal{X})} \\ &\approx \underset{\Theta}{\operatorname{argmax}} p(\mathcal{X}|\Theta) p(\Theta) \\ &= \underset{\Theta}{\operatorname{argmax}} \prod_{x_i \in \mathcal{X}} p(x_i|\Theta) p(\Theta)\end{aligned}$$

$P(\mathcal{X})$  can be removed because it has no functional relation with  $\Theta$ .

We can make this easier

Use logarithms

$$\hat{\Theta}_{MAP} = \operatorname{argmax}_{\Theta} \left[ \sum_{x_i \in \mathcal{X}} \log p(x_i | \Theta) + \log p(\Theta) \right] \quad (37)$$

# What Does the MAP Estimate Get?

## Something Notable

The MAP estimate allows us to inject into the estimation calculation our prior beliefs regarding the parameters values in  $\Theta$ .

# What Does the MAP Estimate Get?

## Something Notable

The MAP estimate allows us to inject into the estimation calculation our prior beliefs regarding the parameters values in  $\Theta$ .

## For example

Let's conduct  $N$  independent trials of the following Bernoulli experiment with  $q$  parameter:

- We will ask each individual we run into in the hallway whether they will vote PRI or PAN in the next presidential election.



# What Does the MAP Estimate Get?

## Something Notable

The MAP estimate allows us to inject into the estimation calculation our prior beliefs regarding the parameters values in  $\Theta$ .

## For example

Let's conduct  $N$  independent trials of the following Bernoulli experiment with  $q$  parameter:

- We will ask each individual we run into in the hallway whether they will vote PRI or PAN in the next presidential election.

Where the values of  $x_i$  is either PRI or PAN.

# What Does the MAP Estimate Get?

## Something Notable

The MAP estimate allows us to inject into the estimation calculation our prior beliefs regarding the parameters values in  $\Theta$ .

## For example

Let's conduct  $N$  independent trials of the following Bernoulli experiment with  $q$  parameter:

- We will ask each individual we run into in the hallway whether they will vote PRI or PAN in the next presidential election.

With probability  $q$  to vote PRI

Where the values of  $x_i$  is either PRI or PAN.

## First the Maximum Likelihood Estimate

### Samples

$$\mathcal{X} = \left\{ x_i = \begin{cases} PAN \\ PRI \end{cases} \quad i = 1, \dots, N \right\} \quad (38)$$

The log likelihood function

## First the Maximum Likelihood Estimate

### Samples

$$\mathcal{X} = \left\{ x_i = \begin{cases} PAN \\ PRI \end{cases} \quad i = 1, \dots, N \right\} \quad (38)$$

### The log likelihood function

$$\begin{aligned} \log p(\mathcal{X}|q) &= \sum_{i=1}^N \log p(x_i|q) \\ &= \sum_i \log p(x_i = PRI|q) + \dots \\ &\quad \sum_i \log p(x_i = PAN|1-q) \\ &= n_{PRI} \log(q) + (N - n_{PRI}) \log(1-q) \end{aligned}$$

Where  $n_{PRI}$  are the numbers of individuals who are planning to vote PRI this fall

## First the Maximum Likelihood Estimate

### Samples

$$\mathcal{X} = \left\{ x_i = \begin{cases} PAN \\ PRI \end{cases} \quad i = 1, \dots, N \right\} \quad (38)$$

### The log likelihood function

$$\begin{aligned} \log p(\mathcal{X}|q) &= \sum_{i=1}^N \log p(x_i|q) \\ &= \sum_i \log p(x_i = PRI|q) + \dots \\ &\quad \sum_i \log p(x_i = PAN|1-q) \\ &= n_{PRI} \log(q) + (N - n_{PRI}) \log(1-q) \end{aligned}$$

Where  $n_{PRI}$  are the numbers of individuals who are planning to vote PRI this fall

## First the Maximum Likelihood Estimate

### Samples

$$\mathcal{X} = \left\{ x_i = \begin{cases} PAN \\ PRI \end{cases} \quad i = 1, \dots, N \right\} \quad (38)$$

### The log likelihood function

$$\begin{aligned} \log p(\mathcal{X}|q) &= \sum_{i=1}^N \log p(x_i|q) \\ &= \sum_i \log p(x_i = PRI|q) + \dots \\ &\quad \sum_i \log p(x_i = PAN|1-q) \\ &= n_{PRI} \log(q) + (N - n_{PRI}) \log(1-q) \end{aligned}$$

Where  $n_{PRI}$  are the numbers of individuals who are planning to vote PRI this fall

## First the Maximum Likelihood Estimate

### Samples

$$\mathcal{X} = \left\{ x_i = \begin{cases} PAN \\ PRI \end{cases} \quad i = 1, \dots, N \right\} \quad (38)$$

### The log likelihood function

$$\begin{aligned} \log p(\mathcal{X}|q) &= \sum_{i=1}^N \log p(x_i|q) \\ &= \sum_i \log p(x_i = PRI|q) + \dots \\ &\quad \sum_i \log p(x_i = PAN|1 - q) \\ &= n_{PRI} \log(q) + (N - n_{PRI}) \log(1 - q) \end{aligned}$$

Where  $n_{PRI}$  are the numbers of individuals who are planning to vote PRI this fall

## We use our classic tricks

By setting

$$\mathcal{L} = \log p(\mathcal{X}|q) \quad (39)$$

We have that

$$\frac{\partial \mathcal{L}}{\partial q} = 0 \quad (40)$$

Thus

$$\frac{n_{PRI}}{q} - \frac{(N - n_{PRI})}{(1 - q)} = 0 \quad (41)$$



## We use our classic tricks

By setting

$$\mathcal{L} = \log p(\mathcal{X}|q) \quad (39)$$

We have that

$$\frac{\partial \mathcal{L}}{\partial q} = 0 \quad (40)$$

Thus

$$\frac{n_{PRI}}{q} - \frac{(N - n_{PRI})}{(1 - q)} = 0 \quad (41)$$

## We use our classic tricks

By setting

$$\mathcal{L} = \log p(\mathcal{X}|q) \quad (39)$$

We have that

$$\frac{\partial \mathcal{L}}{\partial q} = 0 \quad (40)$$

Thus

$$\frac{n_{PRI}}{q} - \frac{(N - n_{PRI})}{(1 - q)} = 0 \quad (41)$$

## Final Solution of ML

We get

$$\hat{q}_{PRI} = \frac{n_{PRI}}{N} \quad (42)$$

Thus

If we say that  $N = 20$  and if 12 are going to vote PRI, we get  $\hat{q}_{PRI} = 0.6$ .

## Final Solution of ML

We get

$$\hat{q}_{PRI} = \frac{n_{PRI}}{N} \quad (42)$$

Thus

If we say that  $N = 20$  and if 12 are going to vote PRI, we get  $\hat{q}_{PRI} = 0.6$ .

## Building the MAP estimate

Obviously we need a prior belief distribution

We have the following constraints:

- The prior for  $q$  must be zero outside the  $[0, 1]$  interval.
- Within the  $[0, 1]$  interval, we are free to specify our beliefs in any way we wish.
- In most cases, we would want to choose a distribution for the prior beliefs that peaks somewhere in the  $[0, 1]$  interval.

# Building the MAP estimate

## Obviously we need a prior belief distribution

We have the following constraints:

- The prior for  $q$  must be zero outside the  $[0, 1]$  interval.
- Within the  $[0, 1]$  interval, we are free to specify our beliefs in any way we wish.
- In most cases, we would want to choose a distribution for the prior beliefs that peaks somewhere in the  $[0, 1]$  interval.

We assume the following

- The state of Colima has traditionally voted PRI in presidential elections.
- However, on account of the prevailing economic conditions, the voters are more likely to vote PAN in the election in question.

## Building the MAP estimate

### Obviously we need a prior belief distribution

We have the following constraints:

- The prior for  $q$  must be zero outside the  $[0, 1]$  interval.
- Within the  $[0, 1]$  interval, we are free to specify our beliefs in any way we wish.
- In most cases, we would want to choose a distribution for the prior beliefs that peaks somewhere in the  $[0, 1]$  interval.

We assume the following

- The state of Colima has traditionally voted PRI in presidential elections.
- However, on account of the prevailing economic conditions, the voters are more likely to vote PAN in the election in question.

## Building the MAP estimate

### Obviously we need a prior belief distribution

We have the following constraints:

- The prior for  $q$  must be zero outside the  $[0, 1]$  interval.
- Within the  $[0, 1]$  interval, we are free to specify our beliefs in any way we wish.
- In most cases, we would want to choose a distribution for the prior beliefs that peaks somewhere in the  $[0, 1]$  interval.

We assume the following:

- The state of Colima has traditionally voted PRI in presidential elections.
- However, on account of the prevailing economic conditions, the voters are more likely to vote PAN in the election in question.



## Building the MAP estimate

### Obviously we need a prior belief distribution

We have the following constraints:

- The prior for  $q$  must be zero outside the  $[0, 1]$  interval.
- Within the  $[0, 1]$  interval, we are free to specify our beliefs in any way we wish.
- In most cases, we would want to choose a distribution for the prior beliefs that peaks somewhere in the  $[0, 1]$  interval.

### We assume the following

- The state of Colima has traditionally voted PRI in presidential elections.
- However, on account of the prevailing economic conditions, the voters are more likely to vote PAN in the election in question.

## What prior distribution can we use?

We could use a Beta distribution being parametrized by two values  $\alpha$  and  $\beta$

$$p(q) = \frac{1}{B(\alpha, \beta)} q^{\alpha-1} (1-q)^{\beta-1}. \quad (43)$$

### Where

We have  $B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$  is the beta function where  $\Gamma$  is the generalization of the notion of factorial in the case of the real numbers.

### Properties

When both the  $\alpha, \beta > 0$  then the beta distribution has its mode (Maximum value) at

$$\frac{\alpha-1}{\alpha+\beta-2}. \quad (44)$$

## What prior distribution can we use?

We could use a Beta distribution being parametrized by two values  $\alpha$  and  $\beta$

$$p(q) = \frac{1}{B(\alpha, \beta)} q^{\alpha-1} (1-q)^{\beta-1}. \quad (43)$$

### Where

We have  $B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$  is the beta function where  $\Gamma$  is the generalization of the notion of factorial in the case of the real numbers.

### Properties

When both the  $\alpha, \beta > 0$  then the beta distribution has its mode (Maximum value) at

$$\frac{\alpha-1}{\alpha+\beta-2}. \quad (44)$$

## What prior distribution can we use?

We could use a Beta distribution being parametrized by two values  $\alpha$  and  $\beta$

$$p(q) = \frac{1}{B(\alpha, \beta)} q^{\alpha-1} (1-q)^{\beta-1}. \quad (43)$$

### Where

We have  $B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$  is the beta function where  $\Gamma$  is the generalization of the notion of factorial in the case of the real numbers.

### Properties

When both the  $\alpha, \beta > 0$  then the beta distribution has its mode (Maximum value) at

$$\frac{\alpha - 1}{\alpha + \beta - 2}. \quad (44)$$

We then do the following

We do the following

We can choose  $\alpha = \beta$  so the beta prior peaks at 0.5.

As a further expression of our belief

We make the following choice  $\alpha = \beta = 5$ .

Why? Look at the variance of the beta distribution

$$\frac{\alpha\beta}{(\alpha + \beta)^2 (\alpha + \beta + 1)} \quad (45)$$

We then do the following

We do the following

We can choose  $\alpha = \beta$  so the beta prior peaks at 0.5.

As a further expression of our belief

We make the following choice  $\alpha = \beta = 5$ .

Why? Look at the variance of the beta distribution

$$\frac{\alpha\beta}{(\alpha + \beta)^2 (\alpha + \beta + 1)} \quad (45)$$

We then do the following

We do the following

We can choose  $\alpha = \beta$  so the beta prior peaks at 0.5.

As a further expression of our belief

We make the following choice  $\alpha = \beta = 5$ .

Why? Look at the variance of the beta distribution

$$\frac{\alpha\beta}{(\alpha + \beta)^2 (\alpha + \beta + 1)}. \quad (45)$$

Thus, we have the following nice properties

We have a variance with  $\alpha = \beta = 5$

$$\text{Var}(q) \approx 0.025$$

Thus, the standard deviation

$sd \approx 0.16$  which is a nice dispersion at the peak point!!!



Thus, we have the following nice properties

We have a variance with  $\alpha = \beta = 5$

$$\text{Var}(q) \approx 0.025$$

Thus, the standard deviation

$sd \approx 0.16$  which is a nice dispersion at the peak point!!!

Now, our MAP estimate for  $\hat{p}_{MAP}$ ...

We have then

$$\hat{p}_{MAP} = \underset{\Theta}{\operatorname{argmax}} \left[ \sum_{x_i \in \mathcal{X}} \log p(x_i|q) + \log p(q) \right] \quad (46)$$

Plugging back the ML

$$\hat{p}_{MAP} = \underset{\Theta}{\operatorname{argmax}} [n_{PRI} \log q + (N - n_{PRI}) \log(1 - q) + \log p(q)] \quad (47)$$

Where

$$\log p(q) = \log \left( \frac{1}{B(\alpha, \beta)} q^{\alpha-1} (1 - q)^{\beta-1} \right) \quad (48)$$

Now, our MAP estimate for  $\hat{p}_{MAP}$ ...

We have then

$$\hat{p}_{MAP} = \underset{\Theta}{\operatorname{argmax}} \left[ \sum_{x_i \in \mathcal{X}} \log p(x_i | q) + \log p(q) \right] \quad (46)$$

Plugging back the ML

$$\hat{p}_{MAP} = \underset{\Theta}{\operatorname{argmax}} [n_{PRI} \log q + (N - n_{PRI}) \log(1 - q) + \log p(q)] \quad (47)$$

Where

$$\log p(q) = \log \left( \frac{1}{B(\alpha, \beta)} q^{\alpha-1} (1-q)^{\beta-1} \right) \quad (48)$$

Now, our MAP estimate for  $\hat{p}_{MAP}$ ...

We have then

$$\hat{p}_{MAP} = \underset{\Theta}{\operatorname{argmax}} \left[ \sum_{x_i \in \mathcal{X}} \log p(x_i|q) + \log p(q) \right] \quad (46)$$

Plugging back the ML

$$\hat{p}_{MAP} = \underset{\Theta}{\operatorname{argmax}} [n_{PRI} \log q + (N - n_{PRI}) \log(1 - q) + \log p(q)] \quad (47)$$

Where

$$\log p(q) = \log \left( \frac{1}{B(\alpha, \beta)} q^{\alpha-1} (1 - q)^{\beta-1} \right) \quad (48)$$

## The log of $p(q)$

We have that

$$\log p(q) = (\alpha - 1) \log q + (\beta - 1) \log(1 - q) - \log B(\alpha, \beta) \quad (49)$$

Now taking the derivative with respect to  $q$ , we get

$$\frac{n_{PRI}}{q} - \frac{(N - n_{PRI})}{(1 - q)} - \frac{\beta - 1}{1 - q} + \frac{\alpha - 1}{q} = 0 \quad (50)$$

Thus

$$\hat{q}_{MAP} = \frac{n_{PRI} + \alpha - 1}{N + \alpha + \beta - 2} \quad (51)$$

## The log of $p(q)$

We have that

$$\log p(q) = (\alpha - 1) \log q + (\beta - 1) \log(1 - q) - \log B(\alpha, \beta) \quad (49)$$

Now taking the derivative with respect to  $p$ , we get

$$\frac{n_{PRI}}{q} - \frac{(N - n_{PRI})}{(1 - q)} - \frac{\beta - 1}{1 - q} + \frac{\alpha - 1}{q} = 0 \quad (50)$$

Thus

$$\hat{q}_{MAP} = \frac{n_{PRI} + \alpha - 1}{N + \alpha + \beta - 2} \quad (51)$$

## The log of $p(q)$

We have that

$$\log p(q) = (\alpha - 1) \log q + (\beta - 1) \log(1 - q) - \log B(\alpha, \beta) \quad (49)$$

Now taking the derivative with respect to  $p$ , we get

$$\frac{n_{PRI}}{q} - \frac{(N - n_{PRI})}{(1 - q)} - \frac{\beta - 1}{1 - q} + \frac{\alpha - 1}{q} = 0 \quad (50)$$

Thus

$$\hat{q}_{MAP} = \frac{n_{PRI} + \alpha - 1}{N + \alpha + \beta - 2} \quad (51)$$

Now

With  $N = 20$  with  $n_{PRI} = 12$  and  $\alpha = \beta = 5$

$$\hat{q}_{MAP} = 0.571$$



# Outline

1

## Introduction

- Supervised Learning
- Handling Noise in Classification
- Models of Classification
- Naive Bayes
  - Examples
  - The Naive Bayes Model
  - The Multi-Class Case

2

## Discriminant Functions and Decision Surfaces

- Introduction
- Gaussian Distribution
- Influence of the Covariance  $\Sigma$
- Example
- Maximum Likelihood Principle
- Maximum Likelihood on a Gaussian
- Some Remarks

3

## Introduction

- A first solution for the Maximum A Posteriori (MAP)
- Maximum Likelihood Vs Maximum A Posteriori
- **Properties of the MAP**

4

## Exercises

- Some Stuff you can try

# Properties

## First

**MAP** estimation “pulls” the estimate toward the prior.

## Second

The more focused our prior belief, the larger the pull toward the prior.

## Example

If  $\alpha = \beta$  =equal to large value

- It will make the MAP estimate to move closer to the prior.

# Properties

## First

**MAP** estimation “pulls” the estimate toward the prior.

## Second

The more focused our prior belief, the larger the pull toward the prior.

## Example

If  $\alpha = \beta$  =equal to large value

- It will make the MAP estimate to move closer to the prior.

# Properties

## First

**MAP** estimation “pulls” the estimate toward the prior.

## Second

The more focused our prior belief, the larger the pull toward the prior.

## Example

If  $\alpha = \beta =$ equal to large value

- It will make the MAP estimate to move closer to the prior.

# Properties

## Third

In the expression we derived for  $\hat{q}_{MAP}$ , the parameters  $\alpha$  and  $\beta$  play a “smoothing” role vis-a-vis the measurement  $n_{PRI}$ .

## Fourth

Since we referred to  $q$  as the parameter to be estimated, we can refer to  $\alpha$  and  $\beta$  as the hyper-parameters in the estimation calculations.

# Properties

## Third

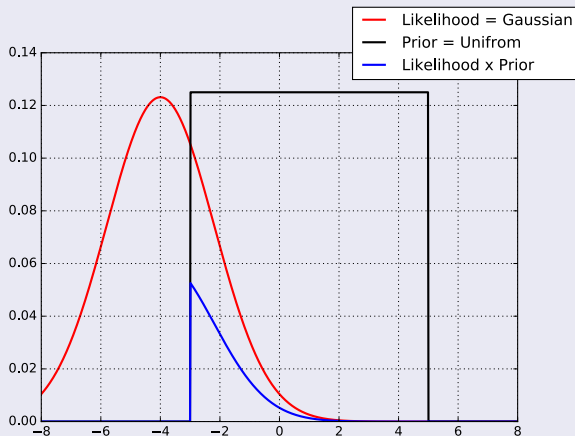
In the expression we derived for  $\hat{q}_{MAP}$ , the parameters  $\alpha$  and  $\beta$  play a “smoothing” role vis-a-vis the measurement  $n_{PRI}$ .

## Fourth

Since we referred to  $q$  as the parameter to be estimated, we can refer to  $\alpha$  and  $\beta$  as the hyper-parameters in the estimation calculations.

# Basically the MAP

It is using the power of Likelihood  $\times$  Prior to obtain more information from the data



## Beyond simple derivation

### In the previous technique

We took an logarithm of the **likelihood**  $\times$  **the prior** to obtain a function that can be derived in order to obtain each of the parameters to be estimated.

What if we cannot derive the likelihood  $\times$  the prior?

For example when we have something like  $\{\theta_i\}$ .

We can try the following:

EM + MAP to be able to estimate the sought parameters.



## Beyond simple derivation

### In the previous technique

We took an logarithm of the **likelihood**  $\times$  **the prior** to obtain a function that can be derived in order to obtain each of the parameters to be estimated.

### What if we cannot derive the **likelihood** $\times$ **the prior**?

For example when we have something like  $|\theta_i|$ .

We can try the following:

EM + MAP to be able to estimate the sought parameters.

## Beyond simple derivation

### In the previous technique

We took an logarithm of the **likelihood**  $\times$  **the prior** to obtain a function that can be derived in order to obtain each of the parameters to be estimated.

### What if we cannot derive the **likelihood** $\times$ **the prior**?

For example when we have something like  $|\theta_i|$ .

### We can try the following

EM + MAP to be able to estimate the sought parameters.

# Outline

1

## Introduction

- Supervised Learning
- Handling Noise in Classification
- Models of Classification
- Naive Bayes
  - Examples
  - The Naive Bayes Model
  - The Multi-Class Case

2

## Discriminant Functions and Decision Surfaces

- Introduction
- Gaussian Distribution
- Influence of the Covariance  $\Sigma$
- Example
- Maximum Likelihood Principle
- Maximum Likelihood on a Gaussian
- Some Remarks

3

## Introduction

- A first solution for the Maximum A Posteriori (MAP)
- Maximum Likelihood Vs Maximum A Posteriori
- Properties of the MAP

4

## Exercises

- Some Stuff you can try

# Exercises

## Duda and Hart

### Chapter 3

- 3.1, 3.2, 3.3, 3.13

## Theodoulidis

### Chapter 2

- 2.5, 2.7, 2.10, 2.12, 2.14, 2.17

# Exercises

## Duda and Hart

### Chapter 3

- 3.1, 3.2, 3.3, 3.13

## Theodoridis

### Chapter 2

- 2.5, 2.7, 2.10, 2.12, 2.14, 2.17