

Classification

Cédric Archambeau

cedrica@amazon.com



Peyresq Summer School
France, July 2016

Overview

- 1 Classification (2.5 hours)
- 2 Clustering (1.5 hours)
- 3 Practical sessions (1 hour)

Overview

- 1 Classification (2.5 hours)
- 2 Clustering (1.5 hours)
- 3 Practical sessions (1 hour)

LEARNING GOALS

- Understand what is a classification problem and when it can be applied.
- Being able to reason about new models and derive learning algorithms.

Overview

- 1 Classification (2.5 hours)
- 2 Clustering (1.5 hours)
- 3 Practical sessions (1 hour)

LEARNING GOALS

- Understand what is a classification problem and when it can be applied.
- Being able to reason about new models and derive learning algorithms.
- **Being able to learn more by yourself!**

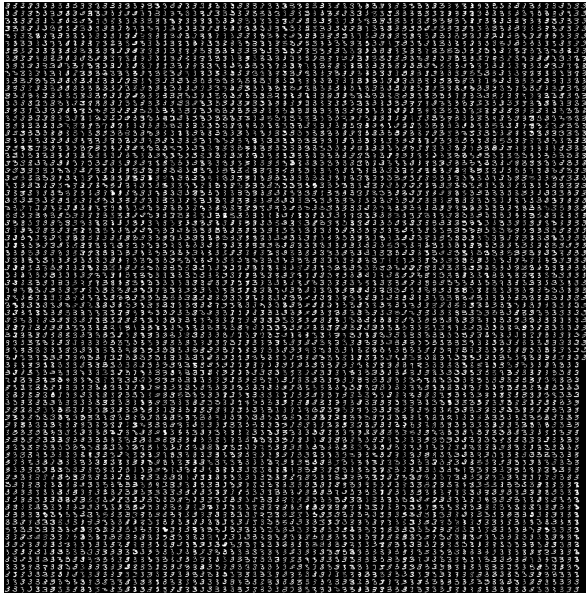
Outline

- 1 What is classification?
- 2 Decision theory
- 3 Generative classifiers
- 4 Discriminative classifiers
- 5 Summary
- 6 Exercises

Outline

- 1 What is classification?
- 2 Decision theory
- 3 Generative classifiers
- 4 Discriminative classifiers
- 5 Summary
- 6 Exercises

A first example: digit classification



MNIST handwritten digits

A first example: digit classification



MNIST handwritten digit sample

A first example: digit classification



MNIST handwritten digit sample

Given an image, can we predict which digit it is (i.e., which label it has)?

A first example: digit classification



MNIST handwritten digit sample

Given an image, can we predict which digit it is (i.e., which label it has)?

Pre-processed data set of handwritten digits: $\mathcal{D} = \{(\mathbf{x}_i, t_i) | i = 1, \dots, n\}$.

- <http://yann.lecun.com/exdb/mnist>

A first example: digit classification



MNIST handwritten digit sample

Given an image, can we predict which digit it is (i.e., which label it has)?

Pre-processed data set of handwritten digits: $\mathcal{D} = \{(\mathbf{x}_i, t_i) | i = 1, \dots, n\}$.

- <http://yann.lecun.com/exdb/mnist>
- Instance or *data point* i consists in a 28×28 bitmap image \mathbf{x}_i and a label $t_i \in \{0, \dots, 9\}$.

A first example: digit classification



MNIST handwritten digit sample

Given an image, can we predict which digit it is (i.e., which label it has)?

Pre-processed data set of handwritten digits: $\mathcal{D} = \{(\mathbf{x}_i, t_i) | i = 1, \dots, n\}$.

- <http://yann.lecun.com/exdb/mnist>
- Instance or *data point* i consists in a 28×28 bitmap image \mathbf{x}_i and a label $t_i \in \{0, \dots, 9\}$.
- Each image is represented as a 784-dimensional vector of pixels, quantized to $\{0, \dots, 255\}$.

A first example: digit classification



MNIST handwritten digit sample

Can we learn $f : \mathbf{x} \mapsto f(\mathbf{x}) = t$?

Pre-processed data set of handwritten digits: $\mathcal{D} = \{(\mathbf{x}_i, t_i) | i = 1, \dots, n\}$.

- <http://yann.lecun.com/exdb/mnist>
- Instance or *data point* i consists in a 28×28 bitmap image \mathbf{x}_i and a label $t_i \in \{0, \dots, 9\}$.
- Each image is represented as a 784-dimensional vector of pixels, quantized to $\{0, \dots, 255\}$.

A first example: digit classification



We would like to distinguish digit 8 from digit 9:

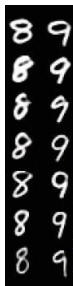
A first example: digit classification



We would like to distinguish digit 8 from digit 9:

$$f : \mathbb{R}^d \rightarrow \{-1, +1\}$$

A first example: digit classification



We would like to distinguish digit 8 from digit 9:

$$f : \mathbb{R}^d \rightarrow \{-1, +1\} \text{ or } f : \mathbb{R}^d \rightarrow \{0, 1\}.$$

A first example: digit classification

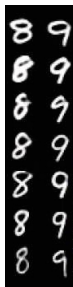


We would like to distinguish digit 8 from digit 9:

$$f : \mathbb{R}^d \rightarrow \{-1, +1\} \text{ or } f : \mathbb{R}^d \rightarrow \{0, 1\}.$$

This is a binary classification problem!

A first example: digit classification



We would like to distinguish digit 8 from digit 9:

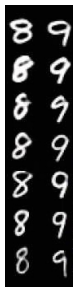
$$f : \mathbb{R}^d \rightarrow \{-1, +1\} \text{ or } f : \mathbb{R}^d \rightarrow \{0, 1\}.$$

This is a binary classification problem!

- Why not a look-up table?

$$f_{\text{LU}}(\mathbf{x}) = \begin{cases} t_i & \text{if } \mathbf{x} = \mathbf{x}_i, i \in \{1, \dots, n\}, \\ \text{Don't know} & \text{if } \mathbf{x} \neq \mathbf{x}_i, i \in \{1, \dots, n\}. \end{cases}$$

A first example: digit classification



We would like to distinguish digit 8 from digit 9:

$$f : \mathbb{R}^d \rightarrow \{-1, +1\} \text{ or } f : \mathbb{R}^d \rightarrow \{0, 1\}.$$

This is a binary classification problem!

- Why not a look-up table?

$$f_{\text{LU}}(\mathbf{x}) = \begin{cases} t_i & \text{if } \mathbf{x} = \mathbf{x}_i, i \in \{1, \dots, n\}, \\ \text{Don't know} & \text{if } \mathbf{x} \neq \mathbf{x}_i, i \in \{1, \dots, n\}. \end{cases}$$

- Why not nearest neighbours?

$$f_{\text{NN}}(\mathbf{x}) = t_i \iff \|\mathbf{x} - \mathbf{x}_i\| \leq \|\mathbf{x} - \mathbf{x}_j\|, j \in \{1, \dots, n\}.$$

A first example: digit classification



We would like to distinguish digit 8 from digit 9:

$$f(\mathbf{x}) = \text{sign}(y(\mathbf{x})) = \begin{cases} +1 & \text{if } y(\mathbf{x}) > 0, \\ -1 & \text{if } y(\mathbf{x}) < 0. \end{cases}$$

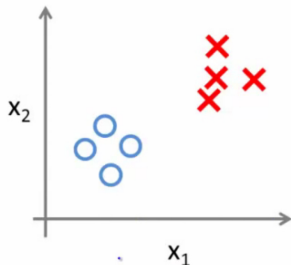
- Why not a look-up table?

$$f_{\text{LU}}(\mathbf{x}) = \begin{cases} t_i & \text{if } \mathbf{x} = \mathbf{x}_i, i \in \{1, \dots, n\}, \\ \text{Don't know} & \text{if } \mathbf{x} \neq \mathbf{x}_i, i \in \{1, \dots, n\}. \end{cases}$$

- Why not nearest neighbours?

$$f_{\text{NN}}(\mathbf{x}) = t_i \iff \|\mathbf{x} - \mathbf{x}_i\| \leq \|\mathbf{x} - \mathbf{x}_j\|, j \in \{1, \dots, n\}.$$

Linear discriminant function (aka linear classifier)

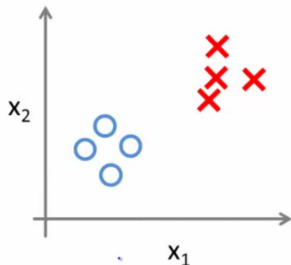


We assume the instances can be separated by a linear subspace (or hyperplane):

$$y(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b, \quad f_{\text{LIN}}(\mathbf{x}) = \text{sign}(y(\mathbf{x})),$$

where $\mathbf{w} \in \mathbb{R}^d \setminus \{\mathbf{0}\}$, $b \in \mathbb{R}$.

Linear discriminant function (aka linear classifier)



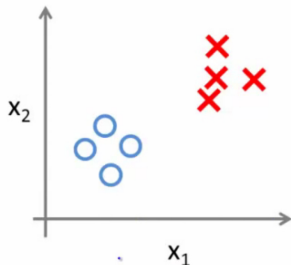
We assume the instances can be separated by a linear subspace (or hyperplane):

$$y(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b, \quad f_{\text{LIN}}(\mathbf{x}) = \text{sign}(y(\mathbf{x})),$$

where $\mathbf{w} \in \mathbb{R}^d \setminus \{\mathbf{0}\}$, $b \in \mathbb{R}$.

- The decision boundary is the set $\{\mathbf{x} : y(\mathbf{x}) = 0\}$.

Linear discriminant function (aka linear classifier)



We assume the instances can be separated by a linear subspace (or hyperplane):

$$y(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b, \quad f_{\text{LIN}}(\mathbf{x}) = \text{sign}(y(\mathbf{x})),$$

where $\mathbf{w} \in \mathbb{R}^d \setminus \{\mathbf{0}\}$, $b \in \mathbb{R}$.

- The decision boundary is the set $\{\mathbf{x} : y(\mathbf{x}) = 0\}$.
- **Learning** is to find \mathbf{w} and b such that $\forall i : f_{\text{LIN}}(\mathbf{x}_i) \approx t_i$.

Relation to nearest neighbour classification?

$$\begin{aligned}\|\mathbf{x}_* - \mathbf{w}_{+1}\| &< \|\mathbf{x}_* - \mathbf{w}_{-1}\| &&\Leftrightarrow \|\mathbf{x}_* - \mathbf{w}_{+1}\|^2 < \|\mathbf{x}_* - \mathbf{w}_{-1}\|^2 \\ \Leftrightarrow \|\mathbf{x}_*\|^2 - 2\mathbf{w}_{+1}^T \mathbf{x}_* + \|\mathbf{w}_{+1}\|^2 &< \|\mathbf{x}_*\|^2 - 2\mathbf{w}_{-1}^T \mathbf{x}_* + \|\mathbf{w}_{-1}\|^2 \\ \Leftrightarrow \mathbf{w}_{+1}^T \mathbf{x}_* - \|\mathbf{w}_{+1}\|^2/2 &> \mathbf{w}_{-1}^T \mathbf{x}_* - \|\mathbf{w}_{-1}\|^2/2 \\ \Leftrightarrow (\mathbf{w}_{+1} - \mathbf{w}_{-1})^T \mathbf{x}_* + \frac{1}{2} (\|\mathbf{w}_{-1}\|^2 - \|\mathbf{w}_{+1}\|^2) &> 0.\end{aligned}$$

How can we picture a linear discriminant function?

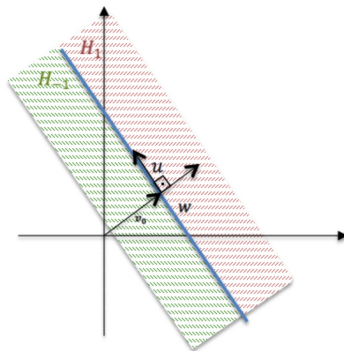


Figure 2.6: Separating hyperplane in \mathbb{R}^2 . The decision boundary (blue) is defined by the normal vector \mathbf{w} and an offset $b \in \mathbb{R}$. \mathbf{v}_0 is a point on the hyperplane, obtained by orthogonal projection of the origin. The plane separates \mathbb{R}^2 into two halfspaces \mathcal{H}_{+1} ($\mathbf{w}^T \mathbf{x} + b > 0$) and \mathcal{H}_{-1} ($\mathbf{w}^T \mathbf{x} + b < 0$), the decision regions of the corresponding linear discriminant.

How can we picture a linear discriminant function?

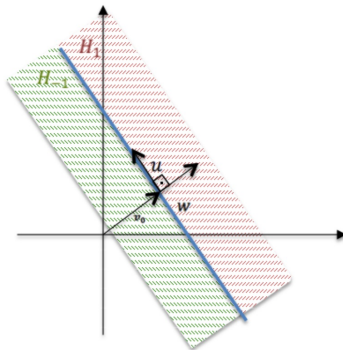


Figure 2.6: Separating hyperplane in \mathbb{R}^2 . The decision boundary (blue) is defined by the normal vector \mathbf{w} and an offset $b \in \mathbb{R}$. \mathbf{v}_0 is a point on the hyperplane, obtained by orthogonal projection of the origin. The plane separates \mathbb{R}^2 into two halfspaces \mathcal{H}_{+1} ($\mathbf{w}^T \mathbf{x} + b > 0$) and \mathcal{H}_{-1} ($\mathbf{w}^T \mathbf{x} + b < 0$), the decision regions of the corresponding linear discriminant.

- Vector \mathbf{w} is orthogonal to any vector \mathbf{u} in the hyperplane: $\mathbf{w}^\top \mathbf{u} = 0$.

How can we picture a linear discriminant function?

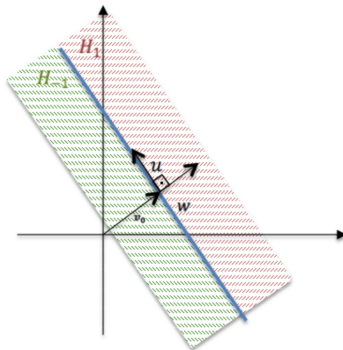


Figure 2.6: Separating hyperplane in \mathbb{R}^2 . The decision boundary (blue) is defined by the normal vector \mathbf{w} and an offset $b \in \mathbb{R}$. \mathbf{v}_0 is a point on the hyperplane, obtained by orthogonal projection of the origin. The plane separates \mathbb{R}^2 into two halfspaces \mathcal{H}_{+1} ($\mathbf{w}^T \mathbf{x} + b > 0$) and \mathcal{H}_{-1} ($\mathbf{w}^T \mathbf{x} + b < 0$), the decision regions of the corresponding linear discriminant.

- Vector \mathbf{w} is orthogonal to any vector \mathbf{u} in the hyperplane: $\mathbf{w}^\top \mathbf{u} = 0$.
- Offset vector $\mathbf{v}_0 = -(b/\|\mathbf{w}\|^2)\mathbf{w}$ is the projection of the origin.

How can we picture a linear discriminant function?

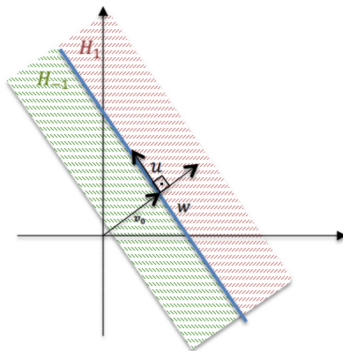


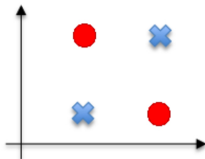
Figure 2.6: Separating hyperplane in \mathbb{R}^2 . The decision boundary (blue) is defined by the normal vector \mathbf{w} and an offset $b \in \mathbb{R}$. \mathbf{v}_0 is a point on the hyperplane, obtained by orthogonal projection of the origin. The plane separates \mathbb{R}^2 into two halfspaces \mathcal{H}_{+1} ($\mathbf{w}^T \mathbf{x} + b > 0$) and \mathcal{H}_{-1} ($\mathbf{w}^T \mathbf{x} + b < 0$), the decision regions of the corresponding linear discriminant.

- Vector \mathbf{w} is orthogonal to any vector \mathbf{u} in the hyperplane: $\mathbf{w}^\top \mathbf{u} = 0$.
- Offset vector $\mathbf{v}_0 = -(b/\|\mathbf{w}\|^2)\mathbf{w}$ is the projection of the origin.
- We restrict weight vectors to be unit norm: $\{\mathbf{w} : \|\mathbf{w}\| = 1\}$.

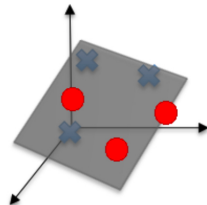
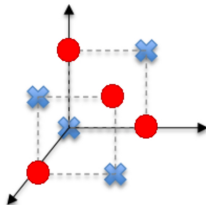
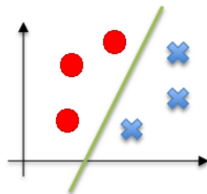
Can instances always be separated?

Can instances always be separated?

Not separable



Separable



Feature maps

- Representing digits as vectors of pixels is arbitrary. Perhaps a transformation would be beneficial?

Feature maps

- Representing digits as vectors of pixels is arbitrary. Perhaps a transformation would be beneficial?
- Let $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^p : \mathbf{x} \mapsto \phi(\mathbf{x})$. The feature function $\phi(\mathbf{x})$ defines a mapping of input space into a feature space.

Feature maps

- Representing digits as vectors of pixels is arbitrary. Perhaps a transformation would be beneficial?
- Let $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^p : \mathbf{x} \mapsto \phi(\mathbf{x})$. The feature function $\phi(\mathbf{x})$ defines a mapping of input space into a feature space.
- We generalise linear classifiers by learning them in the feature space:

$$y(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) + b, \quad f_{\text{LIN}}(\mathbf{x}) = \text{sign}(y(\mathbf{x})).$$

Feature maps

- Representing digits as vectors of pixels is arbitrary. Perhaps a transformation would be beneficial?
- Let $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^p : \mathbf{x} \mapsto \phi(\mathbf{x})$. The feature function $\phi(\mathbf{x})$ defines a mapping of input space into a feature space.
- We generalise linear classifiers by learning them in the feature space:

$$y(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) + b, \quad f_{\text{LIN}}(\mathbf{x}) = \text{sign}(y(\mathbf{x})).$$

Example: linear classifier with quadratic features

$$\phi(\mathbf{x}) = \begin{bmatrix} x_1 \\ \vdots \\ x_d \\ x_1 x_1 \\ \vdots \\ x_1 x_d \\ x_2 x_2 \\ \vdots \\ x_2 x_d \\ \vdots \\ x_d x_d \end{bmatrix} \in \mathbb{R}^{d(d+3)/2}. \quad y(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) + b = \sum_j w_j x_j + \sum_j \sum_{k \leq j} w_{jk} x_j x_k + b$$

How do we estimate \mathbf{w} ?

Perceptron (Rosenblatt, '62):

- Instance i is correctly classified if $\mathbf{w}^\top \phi(\mathbf{x}_i) t_i > 0$.

How do we estimate \mathbf{w} ?

Perceptron (Rosenblatt, '62):

- Instance i is correctly classified if $\mathbf{w}^\top \phi(\mathbf{x}_i) t_i > 0$.
- Perceptron criterion:

$$E_P(\mathbf{w}) = - \sum_{i \in \mathcal{M}} \mathbf{w}^\top \phi(\mathbf{x}_i) t_i,$$

where $\mathcal{M} = \{i : \mathbf{w}^\top \phi(\mathbf{x}_i) t_i < 0\}$.

How do we estimate \mathbf{w} ?

Perceptron (Rosenblatt, '62):

- Instance i is correctly classified if $\mathbf{w}^\top \phi(\mathbf{x}_i)t_i > 0$.
- Perceptron criterion:

$$E_P(\mathbf{w}) = - \sum_{i \in \mathcal{M}} \mathbf{w}^\top \phi(\mathbf{x}_i)t_i,$$

where $\mathcal{M} = \{i : \mathbf{w}^\top \phi(\mathbf{x}_i)t_i < 0\}$.

- The error $E_P(\mathbf{w})$ can be minimised by applying to following rule:

$$\mathbf{w} \leftarrow \mathbf{w} + \phi(\mathbf{x}_i)t_i.$$

How do we estimate \mathbf{w} ?

Perceptron (Rosenblatt, '62):

- Instance i is correctly classified if $\mathbf{w}^\top \phi(\mathbf{x}_i)t_i > 0$.
- Perceptron criterion:

$$E_P(\mathbf{w}) = - \sum_{i \in \mathcal{M}} \mathbf{w}^\top \phi(\mathbf{x}_i)t_i,$$

where $\mathcal{M} = \{i : \mathbf{w}^\top \phi(\mathbf{x}_i)t_i < 0\}$.

- The error $E_P(\mathbf{w})$ can be minimised by applying to following rule:

$$\mathbf{w} \leftarrow \mathbf{w} + \phi(\mathbf{x}_i)t_i.$$

- Why does this algorithm work?

How do we estimate \mathbf{w} ?

Perceptron (Rosenblatt, '62):

- Instance i is correctly classified if $\mathbf{w}^\top \phi(\mathbf{x}_i) t_i > 0$.
- Perceptron criterion:

$$E_P(\mathbf{w}) = - \sum_{i \in \mathcal{M}} \mathbf{w}^\top \phi(\mathbf{x}_i) t_i,$$

where $\mathcal{M} = \{i : \mathbf{w}^\top \phi(\mathbf{x}_i) t_i < 0\}$.

- The error $E_P(\mathbf{w})$ can be minimised by applying to following rule:

$$\mathbf{w} \leftarrow \mathbf{w} + \phi(\mathbf{x}_i) t_i.$$

- Why does this algorithm work?

$$t_i \phi(\mathbf{x}_i)^\top (\mathbf{w} + t_i \phi(\mathbf{x}_i)) = t_i \phi(\mathbf{x}_i)^\top \mathbf{w} + \|\phi(\mathbf{x}_i)\|^2 > t_i \phi(\mathbf{x}_i)^\top \mathbf{w}.$$

How do we estimate \mathbf{w} ?

Perceptron (Rosenblatt, '62):

- Instance i is correctly classified if $\mathbf{w}^\top \phi(\mathbf{x}_i) t_i > 0$.
- Perceptron criterion:

$$E_P(\mathbf{w}) = - \sum_{i \in \mathcal{M}} \mathbf{w}^\top \phi(\mathbf{x}_i) t_i,$$

where $\mathcal{M} = \{i : \mathbf{w}^\top \phi(\mathbf{x}_i) t_i < 0\}$.

- The error $E_P(\mathbf{w})$ can be minimised by applying to following rule:

$$\mathbf{w} \leftarrow \mathbf{w} + \phi(\mathbf{x}_i) t_i.$$

- Why does this algorithm work?

$$t_i \phi(\mathbf{x}_i)^\top (\mathbf{w} + t_i \phi(\mathbf{x}_i)) = t_i \phi(\mathbf{x}_i)^\top \mathbf{w} + \|\phi(\mathbf{x}_i)\|^2 > t_i \phi(\mathbf{x}_i)^\top \mathbf{w}.$$

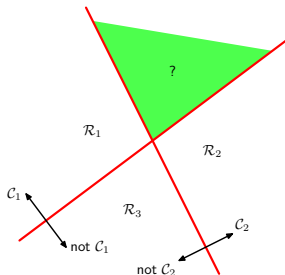
- If the data is not linearly separable, then the perceptron will not converge :-)

How should we handle multi-class problems?

- By combining binary classifiers?

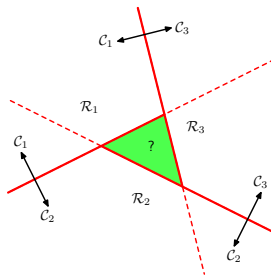
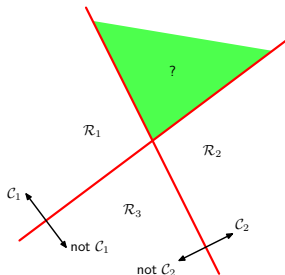
How should we handle multi-class problems?

- By combining binary classifiers?



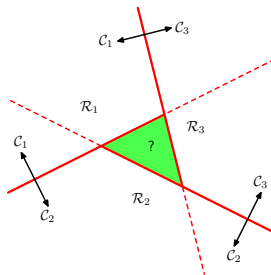
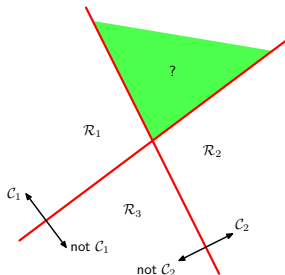
How should we handle multi-class problems?

- By combining binary classifiers?



How should we handle multi-class problems?

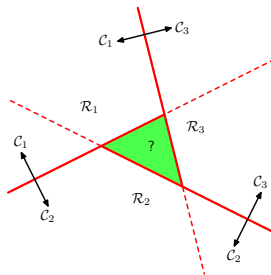
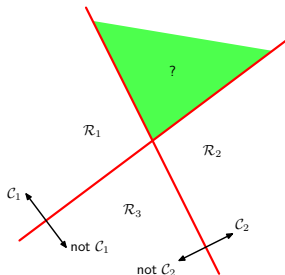
- By combining binary classifiers?



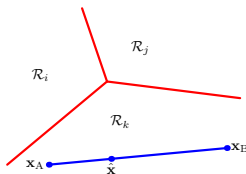
- Are there other ways?

How should we handle multi-class problems?

- By combining binary classifiers?



- Are there other ways?



$$y_k(\mathbf{x}) = \mathbf{w}_k^\top \phi(\mathbf{x}) + b_k,$$

$$f(\mathbf{x}) = \arg \max_k \{y_1(\mathbf{x}), \dots, y_m(\mathbf{x})\}.$$

Examples of classification problems

Examples of classification problems

- Spam detection

Examples of classification problems

- Spam detection
- Fraud detection

Examples of classification problems

- Spam detection
- Fraud detection
- Document categorisation

Examples of classification problems

- Spam detection
- Fraud detection
- Document categorisation
- Sentiment analysis

Examples of classification problems

- Spam detection
- Fraud detection
- Document categorisation
- Sentiment analysis
- Face recognition

Examples of classification problems

- Spam detection
- Fraud detection
- Document categorisation
- Sentiment analysis
- Face recognition
- Object categorisation

Examples of classification problems

- Spam detection
- Fraud detection
- Document categorisation
- Sentiment analysis
- Face recognition
- Object categorisation
- ...

Outline

- 1 What is classification?
- 2 Decision theory
- 3 Generative classifiers
- 4 Discriminative classifiers
- 5 Summary
- 6 Exercises

Bayes' rule

$$P(t|\mathbf{x}) = \frac{p(\mathbf{x}|t)P(t)}{p(\mathbf{x})},$$

$$p(\mathbf{x}) = \sum_t p(\mathbf{x}|t)P(t).$$

Bayes' rule

$$P(t|\mathbf{x}) = \frac{p(\mathbf{x}|t)P(t)}{p(\mathbf{x})},$$

$$p(\mathbf{x}) = \sum_t p(\mathbf{x}|t)P(t).$$

- $P(t)$ is the class prior

Bayes' rule

$$P(t|\mathbf{x}) = \frac{p(\mathbf{x}|t)P(t)}{p(\mathbf{x})},$$

$$p(\mathbf{x}) = \sum_t p(\mathbf{x}|t)P(t).$$

- $P(t)$ is the class prior
- $P(t|\mathbf{x})$ is the class posterior

Bayes' rule

$$P(t|\mathbf{x}) = \frac{p(\mathbf{x}|t)P(t)}{p(\mathbf{x})},$$

$$p(\mathbf{x}) = \sum_t p(\mathbf{x}|t)P(t).$$

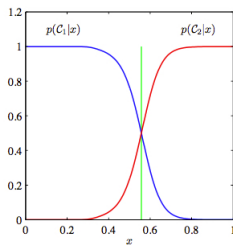
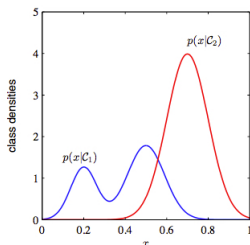
- $P(t)$ is the class prior
- $P(t|\mathbf{x})$ is the class posterior
- $p(\mathbf{x}|t)$ is the class-conditional density (or likelihood)

Bayes' rule

$$P(t|\mathbf{x}) = \frac{p(\mathbf{x}|t)P(t)}{p(\mathbf{x})},$$

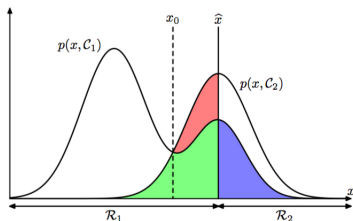
$$p(\mathbf{x}) = \sum_t p(\mathbf{x}|t)P(t).$$

- $P(t)$ is the class prior
- $P(t|\mathbf{x})$ is the class posterior
- $p(\mathbf{x}|t)$ is the class-conditional density (or likelihood)



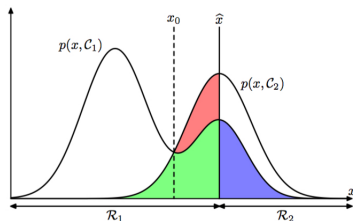
$$\begin{aligned} P(C_1) &= P(t = -1), \\ P(C_1|x) &= P(t = -1|x), \\ p(x, C_1) &= p(x, t = -1). \end{aligned}$$

Minimising the classification error



$$\begin{aligned} p(x, C_1) &= p(x, t = -1) \\ &= p(x|t = -1)P(t = -1). \end{aligned}$$

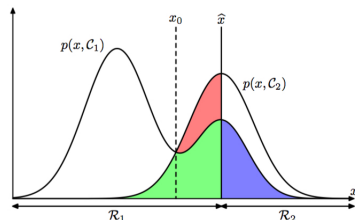
Minimising the classification error



$$\begin{aligned} p(x, C_1) &= p(x, t = -1) \\ &= p(x|t = -1)P(t = -1). \end{aligned}$$

- Let \hat{x} be the **decision threshold**: $\mathcal{R}_1 = \{x : x < \hat{x}\}$ and $\mathcal{R}_2 = \{x : x > \hat{x}\}$.

Minimising the classification error



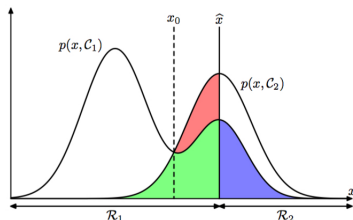
$$\begin{aligned} p(x, C_1) &= p(x, t = -1) \\ &= p(x|t = -1)P(t = -1). \end{aligned}$$

- Let \hat{x} be the **decision threshold**: $\mathcal{R}_1 = \{x : x < \hat{x}\}$ and $\mathcal{R}_2 = \{x : x > \hat{x}\}$.
- The **misclassification rate** is the combined coloured areas:

$$P(\text{error}) = P(x \in \mathcal{R}_1, C_2) + P(x \in \mathcal{R}_2, C_1),$$

$$\text{where } P(x \in \mathcal{R}_1, C_2) = \int_{x \in \mathcal{R}_1} p(x, C_2) dx.$$

Minimising the classification error



$$\begin{aligned} p(x, C_1) &= p(x, t = -1) \\ &= p(x|t = -1)P(t = -1). \end{aligned}$$

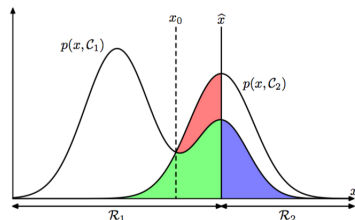
- Let \hat{x} be the **decision threshold**: $\mathcal{R}_1 = \{x : x < \hat{x}\}$ and $\mathcal{R}_2 = \{x : x > \hat{x}\}$.
- The **misclassification rate** is the combined coloured areas:

$$P(\text{error}) = P(x \in \mathcal{R}_1, C_2) + P(x \in \mathcal{R}_2, C_1),$$

where $P(x \in \mathcal{R}_1, C_2) = \int_{x \in \mathcal{R}_1} p(x, C_2) dx$.

- False positives: blue area. False negatives: green+red area.

Minimising the classification error



$$\begin{aligned} p(x, C_1) &= p(x, t = -1) \\ &= p(x|t = -1)P(t = -1). \end{aligned}$$

- Let \hat{x} be the **decision threshold**: $\mathcal{R}_1 = \{x : x < \hat{x}\}$ and $\mathcal{R}_2 = \{x : x > \hat{x}\}$.
- The **misclassification rate** is the combined coloured areas:

$$P(\text{error}) = P(x \in \mathcal{R}_1, C_2) + P(x \in \mathcal{R}_2, C_1),$$

where $P(x \in \mathcal{R}_1, C_2) = \int_{x \in \mathcal{R}_1} p(x, C_2) dx$.

- False positives: blue area. False negatives: green+red area.
- Bayes error** at $x = x_0$: blue+green area.

Example of a Bayes optimal classifier

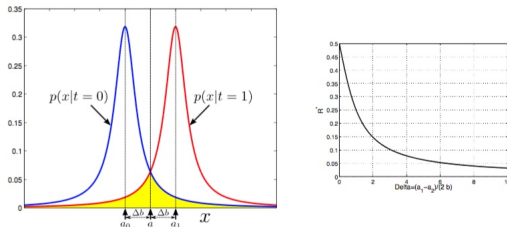
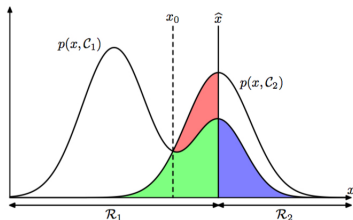


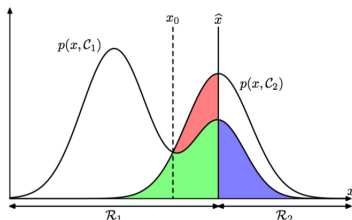
Figure 5.5: Bayes-optimal classifier and Bayes error for two class-conditional Cauchy distributions, centered at a_0 and a_1 . The optimal rule thresholds at the midpoint $a = (a_0 + a_1)/2$. Since the class prior is $P(t = 0) = P(t = 1) = 1/2$, the Bayes error R^* is twice the yellow area. Right plot show R^* as function of separation parameter Δ . The slow decay of R^* is due to the very heavy tails of the Cauchy distributions.

Precision and recall



	$t = 1$	$t = -1$
$x > x_0$	TP	FP
$x < x_0$	FN	TN

Precision and recall



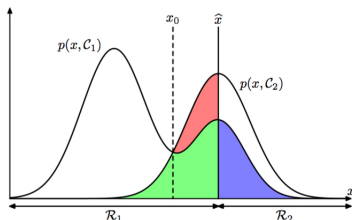
	$t = 1$	$t = -1$
$x > x_0$	TP	FP
$x < x_0$	FN	TN

- The **precision** is the proportion of positives in the instances classified as being positive:

$$P(x) = \frac{TP(x)}{TP(x) + FP(x)},$$

where TP are the true positives and FP the false positives.

Precision and recall



	$t = 1$	$t = -1$
$x > x_0$	TP	FP
$x < x_0$	FN	TN

- The **precision** is the proportion of positives in the instances classified as being positive:

$$P(x) = \frac{TP(x)}{TP(x) + FP(x)},$$

where TP are the true positives and FP the false positives.

- The **recall** is the proportion of correctly classified positives:

$$R(x) = \frac{TP(x)}{TP(x) + FN(x)},$$

where FN are the false negatives.

Should we always minimise the misclassification rate?

Should we always minimise the misclassification rate?

In a hospital, a tissue sample is taken from a patient, giving rise to an input vector \mathbf{x} . A classifier $f(\mathbf{x})$ is to predict whether the patient has cancer ($t = 1$) or not ($t = -1$). Is the cost of predicting that the patient has cancer while he/she has not the same, as predicting that the patient has not contracted cancer while he/she has the disease?

Should we always minimise the misclassification rate?

In a hospital, a tissue sample is taken from a patient, giving rise to an input vector x . A classifier $f(x)$ is to predict whether the patient has cancer ($t = 1$) or not ($t = -1$). Is the cost of predicting that the patient has cancer while he/she has not the same, as predicting that the patient has not contracted cancer while he/she has the disease?

Let us define a **loss function**, which assigns a unique loss to every decision we could take:

	cancer	normal
cancer	0	1000
normal	1	0

Should we always minimise the misclassification rate?

In a hospital, a tissue sample is taken from a patient, giving rise to an input vector x . A classifier $f(x)$ is to predict whether the patient has cancer ($t = 1$) or not ($t = -1$). Is the cost of predicting that the patient has cancer while he/she has not the same, as predicting that the patient has not contracted cancer while he/she has the disease?

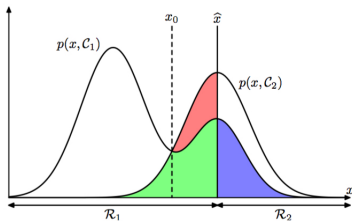
Let us define a **loss function**, which assigns a unique loss to every decision we could take:

	cancer	normal
cancer	0	1000
normal	1	0

The **expected loss** is given by

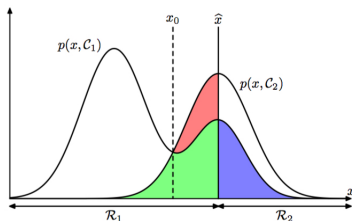
$$\mathbb{E}(L) = \sum_k \sum_l \int_{x \in \mathcal{R}_l} L_{kl} p(x, C_k) dx.$$

How can we make a trade-off?



	$t = 1$	$t = -1$
$x > \hat{x}$	TP	FP
$x < \hat{x}$	FN	TN
	P	N

How can we make a trade-off?

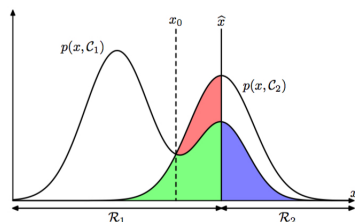


	$t = 1$	$t = -1$
$x > \hat{x}$	TP	FP
$x < \hat{x}$	FN	TN
	P	N

- We defined the decision threshold of a linear classifier as follows:

$$x_0 = \{x : f_{\text{LIN}}(x) = 0\}$$

How can we make a trade-off?



	$t = 1$	$t = -1$
$x > \hat{x}$	TP	FP
$x < \hat{x}$	FN	TN
	P	N

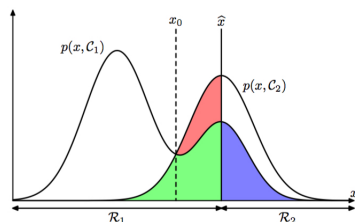
- We defined the decision threshold of a linear classifier as follows:

$$x_0 = \{x : f_{\text{LIN}}(x) = 0\}$$

- We can decide to threshold at any another score α :

$$x_\alpha = \{x : f(x) = \alpha\}$$

How can we make a trade-off?



	$t = 1$	$t = -1$
$x > \hat{x}$	<i>TP</i>	<i>FP</i>
$x < \hat{x}$	<i>FN</i>	<i>TN</i>
	<i>P</i>	<i>N</i>

- We defined the decision threshold of a linear classifier as follows:

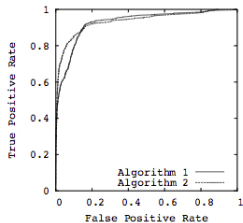
$$x_0 = \{x : f_{\text{LIN}}(x) = 0\}$$

- We can decide to threshold at any another score α :

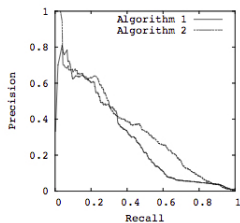
$$x_\alpha = \{x : f(x) = \alpha\}$$

- What is the effect of choosing \hat{x} ?

Area under the curve (AUC)



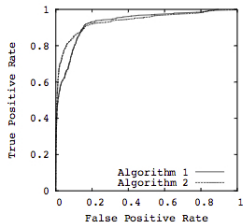
(a) Comparison in ROC space



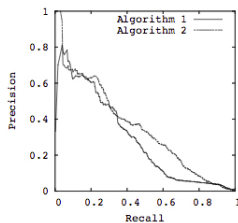
(b) Comparison in PR space

- AUC enables us to compare classifiers irrespective of the decision threshold.

Area under the curve (AUC)



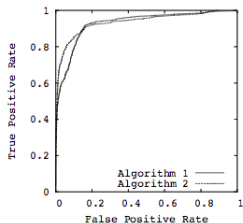
(a) Comparison in ROC space



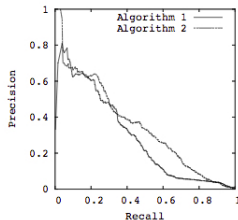
(b) Comparison in PR space

- AUC enables us to compare classifiers irrespective of the decision threshold.
- Receiver-operating characteristic (ROC):
 - ▶ Monotonic
 - ▶ $AUC \approx$ probability of scoring a positive higher than a negative.

Area under the curve (AUC)



(a) Comparison in ROC space



(b) Comparison in PR space

- AUC enables us to compare classifiers irrespective of the decision threshold.
- Receiver-operating characteristic (ROC):
 - ▶ Monotonic
 - ▶ $AUC \approx$ probability of scoring a positive higher than a negative.
- Precision-recall:
 - ▶ Non-monotonic
 - ▶ One to one mapping with ROC

How do we measure the performance in multi-class classification?



- Confusion matrix:

	$t = 1$	$t = 2$	$t = 3$
$f(\mathbf{x}) = 1$	c_{11}	c_{12}	c_{13}
$f(\mathbf{x}) = 2$	c_{21}	c_{22}	c_{23}
$f(\mathbf{x}) = 3$	c_{31}	c_{32}	c_{33}
	P_1	P_2	P_3

How do we measure the performance in multi-class classification?



- Confusion matrix:

	$t = 1$	$t = 2$	$t = 3$
$f(\mathbf{x}) = 1$	c_{11}	c_{12}	c_{13}
$f(\mathbf{x}) = 2$	c_{21}	c_{22}	c_{23}
$f(\mathbf{x}) = 3$	c_{31}	c_{32}	c_{33}
	P_1	P_2	P_3

- What are the expressions of precision and recall?

Outline

- 1 What is classification?
- 2 Decision theory
- 3 Generative classifiers**
- 4 Discriminative classifiers
- 5 Summary
- 6 Exercises

Generative classifiers

Consider the data set $\mathcal{D} = \{(\mathbf{x}_i, t_i) | i = 1, \dots, n\}$. We are interested in the **posterior class probability**:

$$P(t = k | \mathbf{x}) = \frac{\overbrace{p(\mathbf{x} | t = k)}^{\text{class-conditional density}} \overbrace{P(t = k)}^{\text{class prior}}}{p(\mathbf{x})}, \quad k = \{1, \dots, m\}.$$

Generative classifiers

Consider the data set $\mathcal{D} = \{(\mathbf{x}_i, t_i) | i = 1, \dots, n\}$. We are interested in the **posterior class probability**:

$$P(t = k | \mathbf{x}) = \frac{\overbrace{p(\mathbf{x} | t = k)}^{\text{class-conditional density}} \overbrace{P(t = k)}^{\text{class prior}}}{p(\mathbf{x})}, \quad k = \{1, \dots, m\}.$$

- Prior: $P(t = k) = \pi_k$.
- Continuous features: $p(\mathbf{x} | t = k) = \text{Gaussian}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$.

Generative classifiers

Consider the data set $\mathcal{D} = \{(\mathbf{x}_i, t_i) | i = 1, \dots, n\}$. We are interested in the **posterior class probability**:

$$P(t = k | \mathbf{x}) = \frac{\overbrace{p(\mathbf{x} | t = k)}^{\text{class-conditional density}} \overbrace{P(t = k)}^{\text{class prior}}}{p(\mathbf{x})}, \quad k = \{1, \dots, m\}.$$

- Prior: $P(t = k) = \pi_k$.
- Continuous features: $p(\mathbf{x} | t = k) = \text{Gaussian}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$.

How can we learn the parameters $\boldsymbol{\theta} = \{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^m$?

Generative classifiers

Consider the data set $\mathcal{D} = \{(\mathbf{x}_i, t_i) | i = 1, \dots, n\}$. We are interested in the **posterior class probability**:

$$P(t = k | \mathbf{x}) = \frac{\overbrace{p(\mathbf{x} | t = k)}^{\text{class-conditional density}} \overbrace{P(t = k)}^{\text{class prior}}}{p(\mathbf{x})}, \quad k = \{1, \dots, m\}.$$

- Prior: $P(t = k) = \pi_k$.
- Continuous features: $p(\mathbf{x} | t = k) = \text{Gaussian}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$.

How can we learn the parameters $\boldsymbol{\theta} = \{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^m$?

$$\arg \max_{\boldsymbol{\theta}} \ln \prod_i p(t_i | \boldsymbol{\theta}) \quad (\text{maximum likelihood estimation})$$

where $p(t_i | \boldsymbol{\theta}) = \text{Categorical}(\pi_1 \text{Gaussian}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1), \dots, \pi_m \text{Gaussian}(\boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m))$.

Generative classifiers

Consider the data set $\mathcal{D} = \{(\mathbf{x}_i, t_i) | i = 1, \dots, n\}$. We are interested in the **posterior class probability**:

$$P(t = k | \mathbf{x}) = \frac{\overbrace{p(\mathbf{x} | t = k)}^{\text{class-conditional density}} \overbrace{P(t = k)}^{\text{class prior}}}{p(\mathbf{x})}, \quad k = \{1, \dots, m\}.$$

- Prior: $P(t = k) = \pi_k$.
- Continuous features: $p(\mathbf{x} | t = k) = \text{Gaussian}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$.
- Discrete features: $p(\mathbf{x} | t = k) = \text{Multinomial}(\boldsymbol{\mu}_k)$.

How can we learn the parameters $\boldsymbol{\theta} = \{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^m$?

$$\arg \max_{\boldsymbol{\theta}} \ln \prod_i p(t_i | \boldsymbol{\theta}) \quad (\text{maximum likelihood estimation})$$

where $p(t_i | \boldsymbol{\theta}) = \text{Categorical}(\pi_1 \text{Gaussian}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1), \dots, \pi_m \text{Gaussian}(\boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m))$.

where $p(t_i | \boldsymbol{\theta}) = \text{Categorical}(\pi_1 \text{Multinomial}(\boldsymbol{\mu}_1), \dots, \pi_m \text{Multinomial}(\boldsymbol{\mu}_m))$.

Definitions

Multivariate Gaussian probability density:

$$\mathbf{x} \sim \text{Gaussian}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2} |\boldsymbol{\Sigma}|^{1/2}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})}.$$

Multinomial probability distribution:

$$\mathbf{x} \sim \text{Multinomial}(\boldsymbol{\mu}) = \frac{(\sum_j x_j)!}{\prod_j x_j!} \prod_{j=1}^d \mu_j^{x_j}.$$

Categorical probability distribution:

$$t \sim \text{Categorical}(\mathbf{p}) = \prod_{k=1}^m p_k^{\delta_k(t)},$$

where $\delta_z(\cdot)$ is the kronecker delta centred at z .

Binary classification

$$P(t = 1|\mathbf{x}) = \frac{p(\mathbf{x}|t = 1)P(t = 1)}{\sum_k p(\mathbf{x}|t = k)P(t = k)}$$

Binary classification

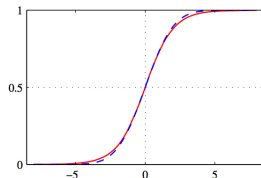
$$\begin{aligned} P(t = 1|x) &= \frac{p(x|t = 1)P(t = 1)}{\sum_k p(x|t = k)P(t = k)} \\ &= \frac{1}{1 + \frac{p(x|t=-1)P(t=-1)}{p(x|t=1)P(t=1)}} \end{aligned}$$

Binary classification

$$\begin{aligned}P(t = 1|x) &= \frac{p(x|t = 1)P(t = 1)}{\sum_k p(x|t = k)P(t = k)} \\&= \frac{1}{1 + \frac{p(x|t=-1)P(t=-1)}{p(x|t=1)P(t=1)}} \\&= \frac{1}{1 + e^{-\ln \frac{p(x|t=1)P(t=1)}{p(x|t=-1)P(t=-1)}}}\end{aligned}$$

Binary classification

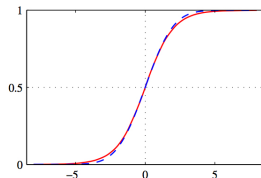
$$\begin{aligned}P(t = 1|x) &= \frac{p(x|t = 1)P(t = 1)}{\sum_k p(x|t = k)P(t = k)} \\&= \frac{1}{1 + \frac{p(x|t=-1)P(t=-1)}{p(x|t=1)P(t=1)}} \\&= \frac{1}{1 + e^{-\ln \frac{p(x|t=1)P(t=1)}{p(x|t=-1)P(t=-1)}}}\end{aligned}$$



Sigmoid: $\sigma(z) = \frac{1}{1+\exp(-z)}$.

Binary classification

$$\begin{aligned}P(t = 1|\mathbf{x}) &= \frac{p(\mathbf{x}|t = 1)P(t = 1)}{\sum_k p(\mathbf{x}|t = k)P(t = k)} \\&= \frac{1}{1 + \frac{p(\mathbf{x}|t=-1)P(t=-1)}{p(\mathbf{x}|t=1)P(t=1)}} \\&= \frac{1}{1 + e^{-\ln \frac{p(\mathbf{x}|t=1)P(t=1)}{p(\mathbf{x}|t=-1)P(t=-1)}}}\end{aligned}$$



Sigmoid: $\sigma(z) = \frac{1}{1+\exp(-z)}$.

Let $\alpha \in [0, 1]$. The classifier is defined as follows:

$$f(\mathbf{x}) = \begin{cases} +1 & \text{if } P(t = 1|\mathbf{x}) > \alpha, \\ -1 & \text{if } P(t = 1|\mathbf{x}) < \alpha. \end{cases}$$

Binary classification with continuous features

- The classifier is defined by the following priors and class-conditionals:

$$\begin{aligned}P(t = 1) &= \pi, & p(\mathbf{x}|t = 1) &= \text{Gaussian}(\boldsymbol{\mu}_{+1}, \boldsymbol{\Sigma}), \\P(t = -1) &= 1 - \pi, & p(\mathbf{x}|t = -1) &= \text{Gaussian}(\boldsymbol{\mu}_{-1}, \boldsymbol{\Sigma}).\end{aligned}$$

where the classes are assumed to share the **same** covariance.

Binary classification with continuous features

- The classifier is defined by the following priors and class-conditionals:

$$\begin{aligned}P(t = 1) &= \pi, & p(\mathbf{x}|t = 1) &= \text{Gaussian}(\boldsymbol{\mu}_{+1}, \boldsymbol{\Sigma}), \\P(t = -1) &= 1 - \pi, & p(\mathbf{x}|t = -1) &= \text{Gaussian}(\boldsymbol{\mu}_{-1}, \boldsymbol{\Sigma}).\end{aligned}$$

where the classes are assumed to share the **same** covariance.

- The **log-likelihood** is given by

$$\begin{aligned}\ln p(\mathbf{t}|\boldsymbol{\theta}) &= \sum_i \delta_{+1}(t_i) (\ln \pi + \ln \text{Gaussian}(\boldsymbol{\mu}_{+1}, \boldsymbol{\Sigma})) \\&\quad + \sum_i \delta_{-1}(t_i) (\ln(1 - \pi) + \ln \text{Gaussian}(\boldsymbol{\mu}_{-1}, \boldsymbol{\Sigma})).\end{aligned}$$

Binary classification with continuous features

- The classifier is defined by the following priors and class-conditionals:

$$\begin{aligned}P(t = 1) &= \pi, & p(\mathbf{x}|t = 1) &= \text{Gaussian}(\boldsymbol{\mu}_{+1}, \boldsymbol{\Sigma}), \\P(t = -1) &= 1 - \pi, & p(\mathbf{x}|t = -1) &= \text{Gaussian}(\boldsymbol{\mu}_{-1}, \boldsymbol{\Sigma}).\end{aligned}$$

where the classes are assumed to share the **same** covariance.

- The **log-likelihood** is given by

$$\begin{aligned}\ln p(\mathbf{t}|\boldsymbol{\theta}) &= \sum_i \delta_{+1}(t_i) (\ln \pi + \ln \text{Gaussian}(\boldsymbol{\mu}_{+1}, \boldsymbol{\Sigma})) \\&\quad + \sum_i \delta_{-1}(t_i) (\ln(1 - \pi) + \ln \text{Gaussian}(\boldsymbol{\mu}_{-1}, \boldsymbol{\Sigma})).\end{aligned}$$

- **Maximum likelihood solution:**

$$\pi = \frac{\sum_i \delta_{+1}(t_i)}{n},$$

Binary classification with continuous features

- The classifier is defined by the following priors and class-conditionals:

$$\begin{aligned}P(t = 1) &= \pi, & p(\mathbf{x}|t = 1) &= \text{Gaussian}(\boldsymbol{\mu}_{+1}, \boldsymbol{\Sigma}), \\P(t = -1) &= 1 - \pi, & p(\mathbf{x}|t = -1) &= \text{Gaussian}(\boldsymbol{\mu}_{-1}, \boldsymbol{\Sigma}).\end{aligned}$$

where the classes are assumed to share the **same** covariance.

- The **log-likelihood** is given by

$$\begin{aligned}\ln p(\mathbf{t}|\boldsymbol{\theta}) &= \sum_i \delta_{+1}(t_i) (\ln \pi + \ln \text{Gaussian}(\boldsymbol{\mu}_{+1}, \boldsymbol{\Sigma})) \\&\quad + \sum_i \delta_{-1}(t_i) (\ln(1 - \pi) + \ln \text{Gaussian}(\boldsymbol{\mu}_{-1}, \boldsymbol{\Sigma})).\end{aligned}$$

- **Maximum likelihood solution:**

$$\pi = \frac{\sum_i \delta_{+1}(t_i)}{n}, \quad \boldsymbol{\mu}_{+1} = \frac{1}{n_{+1}} \sum_{i=1}^N \delta_{+1}(t_i) \mathbf{x}_i,$$

Binary classification with continuous features

- The classifier is defined by the following priors and class-conditionals:

$$\begin{aligned}P(t = 1) &= \pi, & p(\mathbf{x}|t = 1) &= \text{Gaussian}(\boldsymbol{\mu}_{+1}, \boldsymbol{\Sigma}), \\P(t = -1) &= 1 - \pi, & p(\mathbf{x}|t = -1) &= \text{Gaussian}(\boldsymbol{\mu}_{-1}, \boldsymbol{\Sigma}).\end{aligned}$$

where the classes are assumed to share the **same** covariance.

- The **log-likelihood** is given by

$$\begin{aligned}\ln p(\mathbf{t}|\boldsymbol{\theta}) &= \sum_i \delta_{+1}(t_i) (\ln \pi + \ln \text{Gaussian}(\boldsymbol{\mu}_{+1}, \boldsymbol{\Sigma})) \\&\quad + \sum_i \delta_{-1}(t_i) (\ln(1 - \pi) + \ln \text{Gaussian}(\boldsymbol{\mu}_{-1}, \boldsymbol{\Sigma})).\end{aligned}$$

- Maximum likelihood solution:**

$$\pi = \frac{\sum_i \delta_{+1}(t_i)}{n}, \quad \boldsymbol{\mu}_{+1} = \frac{1}{n_{+1}} \sum_{i=1}^N \delta_{+1}(t_i) \mathbf{x}_i, \quad \boldsymbol{\mu}_{-1} = \frac{1}{n_{-1}} \sum_{i=1}^N \delta_{-1}(t_i) \mathbf{x}_i,$$

Binary classification with continuous features

- The classifier is defined by the following priors and class-conditionals:

$$\begin{aligned}P(t = 1) &= \pi, & p(\mathbf{x}|t = 1) &= \text{Gaussian}(\boldsymbol{\mu}_{+1}, \boldsymbol{\Sigma}), \\P(t = -1) &= 1 - \pi, & p(\mathbf{x}|t = -1) &= \text{Gaussian}(\boldsymbol{\mu}_{-1}, \boldsymbol{\Sigma}).\end{aligned}$$

where the classes are assumed to share the **same** covariance.

- The **log-likelihood** is given by

$$\begin{aligned}\ln p(\mathbf{t}|\boldsymbol{\theta}) &= \sum_i \delta_{+1}(t_i) (\ln \pi + \ln \text{Gaussian}(\boldsymbol{\mu}_{+1}, \boldsymbol{\Sigma})) \\&\quad + \sum_i \delta_{-1}(t_i) (\ln(1 - \pi) + \ln \text{Gaussian}(\boldsymbol{\mu}_{-1}, \boldsymbol{\Sigma})).\end{aligned}$$

- Maximum likelihood solution:**

$$\pi = \frac{\sum_i \delta_{+1}(t_i)}{n}, \quad \boldsymbol{\mu}_{+1} = \frac{1}{n_{+1}} \sum_{i=1}^N \delta_{+1}(t_i) \mathbf{x}_i, \quad \boldsymbol{\mu}_{-1} = \frac{1}{n_{-1}} \sum_{i=1}^N \delta_{-1}(t_i) \mathbf{x}_i,$$

$$\boldsymbol{\Sigma} = \frac{n_{+1}}{n} \mathbf{S}_{+1} + \frac{n_{-1}}{n} \mathbf{S}_{-1},$$

Binary classification with continuous features

- The classifier is defined by the following priors and class-conditionals:

$$\begin{aligned}P(t = 1) &= \pi, & p(\mathbf{x}|t = 1) &= \text{Gaussian}(\boldsymbol{\mu}_{+1}, \boldsymbol{\Sigma}), \\P(t = -1) &= 1 - \pi, & p(\mathbf{x}|t = -1) &= \text{Gaussian}(\boldsymbol{\mu}_{-1}, \boldsymbol{\Sigma}).\end{aligned}$$

where the classes are assumed to share the **same** covariance.

- The **log-likelihood** is given by

$$\begin{aligned}\ln p(\mathbf{t}|\boldsymbol{\theta}) &= \sum_i \delta_{+1}(t_i) (\ln \pi + \ln \text{Gaussian}(\boldsymbol{\mu}_{+1}, \boldsymbol{\Sigma})) \\&\quad + \sum_i \delta_{-1}(t_i) (\ln(1 - \pi) + \ln \text{Gaussian}(\boldsymbol{\mu}_{-1}, \boldsymbol{\Sigma})).\end{aligned}$$

- Maximum likelihood solution:**

$$\begin{aligned}\pi &= \frac{\sum_i \delta_{+1}(t_i)}{n}, \quad \boldsymbol{\mu}_{+1} = \frac{1}{n_{+1}} \sum_{i=1}^N \delta_{+1}(t_i) \mathbf{x}_i, \quad \boldsymbol{\mu}_{-1} = \frac{1}{n_{-1}} \sum_{i=1}^N \delta_{-1}(t_i) \mathbf{x}_i, \\ \boldsymbol{\Sigma} &= \frac{n_{+1}}{n} \mathbf{S}_{+1} + \frac{n_{-1}}{n} \mathbf{S}_{-1}, \quad \mathbf{S}_{\pm 1} = \frac{1}{n_{\pm 1}} \sum_{i=1}^N \delta_{\pm 1}(t_i) (\mathbf{x}_i - \boldsymbol{\mu}_{\pm 1})(\mathbf{x}_i - \boldsymbol{\mu}_{\pm 1})^\top.\end{aligned}$$

What is the form of the decision boundary?

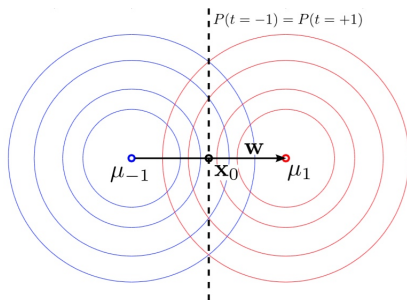
Plugging the priors and the class-conditionals in the posterior leads to

$$P(t = 1|\mathbf{x}) = \sigma \left((\boldsymbol{\mu}_{+1} - \boldsymbol{\mu}_{-1})^\top \boldsymbol{\Sigma}^{-1} \mathbf{x} + -\frac{1}{2} \boldsymbol{\mu}_{+1}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_{+1} + \frac{1}{2} \boldsymbol{\mu}_{-1}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_{-1} + \ln \frac{\pi}{1 - \pi} \right)$$

What is the form of the decision boundary?

Plugging the priors and the class-conditionals in the posterior leads to

$$\begin{aligned} P(t = 1|\mathbf{x}) &= \sigma \left((\boldsymbol{\mu}_{+1} - \boldsymbol{\mu}_{-1})^\top \boldsymbol{\Sigma}^{-1} \mathbf{x} + -\frac{1}{2} \boldsymbol{\mu}_{+1}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_{+1} + \frac{1}{2} \boldsymbol{\mu}_{-1}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_{-1} + \ln \frac{\pi}{1 - \pi} \right) \\ &= \sigma \left(\mathbf{w}^\top \mathbf{x} + b \right). \end{aligned}$$



Binary classification with discrete features

- The classifier is defined by the following priors and class-conditionals:

$$\begin{aligned}P(t = 1) &= \pi, & p(\mathbf{x}|t = 1) &= \text{Multinomial}(\boldsymbol{\mu}_{+1}), \\P(t = -1) &= 1 - \pi, & p(\mathbf{x}|t = -1) &= \text{Multinomial}(\boldsymbol{\mu}_{-1}).\end{aligned}$$

where the features are assumed to be **independent** given the class.

Binary classification with discrete features

- The classifier is defined by the following priors and class-conditionals:

$$\begin{aligned}P(t = 1) &= \pi, & p(\mathbf{x}|t = 1) &= \text{Multinomial}(\boldsymbol{\mu}_{+1}), \\P(t = -1) &= 1 - \pi, & p(\mathbf{x}|t = -1) &= \text{Multinomial}(\boldsymbol{\mu}_{-1}).\end{aligned}$$

where the features are assumed to be **independent** given the class.

- The **log-likelihood** is given by

$$\begin{aligned}\ln p(\mathbf{t}|\boldsymbol{\theta}) &= \sum_i \delta_{+1}(t_i) (\ln \pi + \ln \text{Multinomial}(\boldsymbol{\mu}_{+1})) \\&\quad + \sum_i \delta_{-1}(t_i) (\ln(1 - \pi) + \ln \text{Multinomial}(\boldsymbol{\mu}_{-1})).\end{aligned}$$

Binary classification with discrete features

- The classifier is defined by the following priors and class-conditionals:

$$\begin{aligned} P(t = 1) &= \pi, & p(\mathbf{x}|t = 1) &= \text{Multinomial}(\boldsymbol{\mu}_{+1}), \\ P(t = -1) &= 1 - \pi, & p(\mathbf{x}|t = -1) &= \text{Multinomial}(\boldsymbol{\mu}_{-1}). \end{aligned}$$

where the features are assumed to be **independent** given the class.

- The **log-likelihood** is given by

$$\begin{aligned} \ln p(\mathbf{t}|\boldsymbol{\theta}) &= \sum_i \delta_{+1}(t_i) (\ln \pi + \ln \text{Multinomial}(\boldsymbol{\mu}_{+1})) \\ &\quad + \sum_i \delta_{-1}(t_i) (\ln(1 - \pi) + \ln \text{Multinomial}(\boldsymbol{\mu}_{-1})). \end{aligned}$$

- **Maximum likelihood solution:**

$$\pi = \frac{\sum_i \sum_j \delta_{+1}(t_i) x_{ij}}{\sum_i \sum_j x_{ij}},$$

Binary classification with discrete features

- The classifier is defined by the following priors and class-conditionals:

$$\begin{aligned}P(t = 1) &= \pi, & p(\mathbf{x}|t = 1) &= \text{Multinomial}(\boldsymbol{\mu}_{+1}), \\P(t = -1) &= 1 - \pi, & p(\mathbf{x}|t = -1) &= \text{Multinomial}(\boldsymbol{\mu}_{-1}).\end{aligned}$$

where the features are assumed to be **independent** given the class.

- The **log-likelihood** is given by

$$\begin{aligned}\ln p(\mathbf{t}|\boldsymbol{\theta}) &= \sum_i \delta_{+1}(t_i) (\ln \pi + \ln \text{Multinomial}(\boldsymbol{\mu}_{+1})) \\&\quad + \sum_i \delta_{-1}(t_i) (\ln(1 - \pi) + \ln \text{Multinomial}(\boldsymbol{\mu}_{-1})).\end{aligned}$$

- Maximum likelihood solution:**

$$\pi = \frac{\sum_i \sum_j \delta_{+1}(t_i) x_{ij}}{\sum_i \sum_j x_{ij}},$$

$$\boldsymbol{\mu}_{+1} = \frac{\sum_i \delta_{+1}(t_i) \mathbf{x}_i}{\sum_i \sum_j \delta_{+1}(t_i) x_{ij}},$$

Binary classification with discrete features

- The classifier is defined by the following priors and class-conditionals:

$$\begin{aligned}P(t = 1) &= \pi, & p(\mathbf{x}|t = 1) &= \text{Multinomial}(\boldsymbol{\mu}_{+1}), \\P(t = -1) &= 1 - \pi, & p(\mathbf{x}|t = -1) &= \text{Multinomial}(\boldsymbol{\mu}_{-1}).\end{aligned}$$

where the features are assumed to be **independent** given the class.

- The **log-likelihood** is given by

$$\begin{aligned}\ln p(\mathbf{t}|\boldsymbol{\theta}) &= \sum_i \delta_{+1}(t_i) (\ln \pi + \ln \text{Multinomial}(\boldsymbol{\mu}_{+1})) \\&\quad + \sum_i \delta_{-1}(t_i) (\ln(1 - \pi) + \ln \text{Multinomial}(\boldsymbol{\mu}_{-1})).\end{aligned}$$

- Maximum likelihood solution:**

$$\pi = \frac{\sum_i \sum_j \delta_{+1}(t_i) x_{ij}}{\sum_i \sum_j x_{ij}},$$

$$\boldsymbol{\mu}_{+1} = \frac{\sum_i \delta_{+1}(t_i) \mathbf{x}_i}{\sum_i \sum_j \delta_{+1}(t_i) x_{ij}}, \quad \boldsymbol{\mu}_{-1} = \frac{\sum_i \delta_{-1}(t_i) \mathbf{x}_i}{\sum_i \sum_j \delta_{-1}(t_i) x_{ij}}.$$

What is the form of the decision boundary?

Plugging the priors and the class-conditionals in the posterior leads to

$$P(t = 1|\mathbf{x}) = \sigma \left((\ln \boldsymbol{\mu}_{+1} - \ln \boldsymbol{\mu}_{-1})^\top \mathbf{x} + \ln \frac{\pi}{1 - \pi} \right)$$

What is the form of the decision boundary?

Plugging the priors and the class-conditionals in the posterior leads to

$$\begin{aligned} P(t = 1|\mathbf{x}) &= \sigma \left((\ln \mu_{+1} - \ln \mu_{-1})^\top \mathbf{x} + \ln \frac{\pi}{1 - \pi} \right) \\ &= \sigma \left(\mathbf{w}^\top \mathbf{x} + b \right). \end{aligned}$$

What is the form of the decision boundary?

Plugging the priors and the class-conditionals in the posterior leads to

$$\begin{aligned} P(t = 1|\mathbf{x}) &= \sigma \left((\ln \mu_{+1} - \ln \mu_{-1})^\top \mathbf{x} + \ln \frac{\pi}{1 - \pi} \right) \\ &= \sigma \left(\mathbf{w}^\top \mathbf{x} + b \right). \end{aligned}$$

What if an entry of $\mu_{\pm 1}$ is zero? When can this occur?

What is the form of the decision boundary?

Plugging the priors and the class-conditionals in the posterior leads to

$$\begin{aligned}P(t = 1|\mathbf{x}) &= \sigma \left((\ln \boldsymbol{\mu}_{+1} - \ln \boldsymbol{\mu}_{-1})^\top \mathbf{x} + \ln \frac{\pi}{1 - \pi} \right) \\&= \sigma \left(\mathbf{w}^\top \mathbf{x} + b \right).\end{aligned}$$

What if an entry of $\boldsymbol{\mu}_{\pm 1}$ is zero? When can this occur?

$$\arg \max_{\boldsymbol{\theta}} \ln \prod_i p(t_i|\boldsymbol{\theta}) + \ln p(\boldsymbol{\theta}) \quad (\text{maximum a posteriori estimation})$$

where $p(\boldsymbol{\theta}) = \text{Dirichlet}(\alpha \mathbf{1})$:

What is the form of the decision boundary?

Plugging the priors and the class-conditionals in the posterior leads to

$$\begin{aligned} P(t = 1|\mathbf{x}) &= \sigma \left((\ln \mu_{+1} - \ln \mu_{-1})^\top \mathbf{x} + \ln \frac{\pi}{1 - \pi} \right) \\ &= \sigma \left(\mathbf{w}^\top \mathbf{x} + b \right). \end{aligned}$$

What if an entry of $\mu_{\pm 1}$ is zero? When can this occur?

$$\arg \max_{\boldsymbol{\theta}} \ln \prod_i p(t_i|\boldsymbol{\theta}) + \ln p(\boldsymbol{\theta}) \quad (\text{maximum a posteriori estimation})$$

where $p(\boldsymbol{\theta}) = \text{Dirichlet}(\alpha \mathbf{1})$:

- Parameter α can be interpreted as a pseudo-count. (\star)
- Adding a prior is equivalent to regularisation.

Naive Bayes classifier

- Generative classifier making the simplifying assumption that features are **independent** given the class:

$$P(t = k|\mathbf{x}) = \frac{p(\mathbf{x}|t = k)P(t = k)}{p(\mathbf{x})} \approx \frac{\prod_{j=1}^d p(x_j|t = k)P(t = k)}{p(\mathbf{x})}.$$

Naive Bayes classifier

- Generative classifier making the simplifying assumption that features are **independent** given the class:

$$P(t = k|\mathbf{x}) = \frac{p(\mathbf{x}|t = k)P(t = k)}{p(\mathbf{x})} \approx \frac{\prod_{j=1}^d p(x_j|t = k)P(t = k)}{p(\mathbf{x})}.$$

- Number of parameters scales linearly with the number of features!

Naive Bayes classifier

- Generative classifier making the simplifying assumption that features are **independent** given the class:

$$P(t = k | \mathbf{x}) = \frac{p(\mathbf{x} | t = k)P(t = k)}{p(\mathbf{x})} \approx \frac{\prod_{j=1}^d p(x_j | t = k)P(t = k)}{p(\mathbf{x})}.$$

- Number of parameters scales linearly with the number of features!



Figure 6.8: The Reuters RCV1 collection is a set of 800,000 documents (news articles), with about 200 words per document on average. After standard pre-processing (stop word removal), its dictionary (set of distinct words) is roughly of size 400,000. A common machine learning problem associated with this data is to classify documents into groups (for example: politics, business, sports, science, movies), which are often organized in a hierarchical fashion.

Naive Bayes classifier

- Generative classifier making the simplifying assumption that features are **independent** given the class:

$$P(t = k | \mathbf{x}) = \frac{p(\mathbf{x} | t = k) P(t = k)}{p(\mathbf{x})} \approx \frac{\prod_{j=1}^d p(x_j | t = k) P(t = k)}{p(\mathbf{x})}.$$

- Number of parameters scales linearly with the number of features!



Document categorisation:

$$P(t = k) = \pi_k,$$
$$p(\mathbf{x} | t = k) = \text{Multinomial}(\boldsymbol{\mu}_k).$$

Figure 6.8: The Reuters RCV1 collection is a set of 800,000 documents (news articles), with about 200 words per document on average. After standard pre-processing (stop word removal), its dictionary (set of distinct words) is roughly of size 400,000. A common machine learning problem associated with this data is to classify documents into groups (for example: politics, business, sports, science, movies), which are often organized in a hierarchical fashion.

Naive Bayes classifier

- Generative classifier making the simplifying assumption that features are **independent** given the class:

$$P(t = k | \mathbf{x}) = \frac{p(\mathbf{x} | t = k) P(t = k)}{p(\mathbf{x})} \approx \frac{\prod_{j=1}^d p(x_j | t = k) P(t = k)}{p(\mathbf{x})}.$$

- Number of parameters scales linearly with the number of features!



Figure 6.8: The Reuters RCV1 collection is a set of 800,000 documents (news articles), with about 200 words per document on average. After standard pre-processing (stop word removal), its dictionary (set of distinct words) is roughly of size 400,000. A common machine learning problem associated with this data is to classify documents into groups (for example: politics, business, sports, science, movies), which are often organized in a hierarchical fashion.

Document categorisation:

$$P(t = k) = \pi_k,$$

$$p(\mathbf{x} | t = k) = \text{Multinomial}(\boldsymbol{\mu}_k).$$

- Category/theme/topic k is modelled by a discrete distribution $\boldsymbol{\mu}_k$ over the vocabulary of size d .

Naive Bayes classifier

- Generative classifier making the simplifying assumption that features are **independent** given the class:

$$P(t = k | \mathbf{x}) = \frac{p(\mathbf{x} | t = k) P(t = k)}{p(\mathbf{x})} \approx \frac{\prod_{j=1}^d p(x_j | t = k) P(t = k)}{p(\mathbf{x})}.$$

- Number of parameters scales linearly with the number of features!



Figure 6.8: The Reuters RCV1 collection is a set of 800,000 documents (news articles), with about 200 words per document on average. After standard pre-processing (stop word removal), its dictionary (set of distinct words) is roughly of size 400,000. A common machine learning problem associated with this data is to classify documents into groups (for example: politics, business, sports, science, movies), which are often organized in a hierarchical fashion.

Document categorisation:

$$P(t = k) = \pi_k,$$

$$p(\mathbf{x} | t = k) = \text{Multinomial}(\boldsymbol{\mu}_k).$$

- Category/theme/topic k is modeled by a discrete distribution $\boldsymbol{\mu}_k$ over the vocabulary of size d .
- \mathbf{x}_i represents document i ; it contains the word counts.

Maximum likelihood (ML) estimation: summary

The likelihood is the joint probability of observing i.i.d. data:

$$\ell(\boldsymbol{\theta}; \mathbf{t}) = \ln p(\mathbf{t}; \boldsymbol{\theta}) = \ln \prod_{i=1}^n p(t_i; \boldsymbol{\theta}).$$

Maximum likelihood (ML) estimation: summary

The likelihood is the joint probability of observing i.i.d. data:

$$\ell(\boldsymbol{\theta}; \mathbf{t}) = \ln p(\mathbf{t}; \boldsymbol{\theta}) = \ln \prod_{i=1}^n p(t_i; \boldsymbol{\theta}).$$

The goal is to find the parameters that maximise the log-likelihood function:

$$\boldsymbol{\theta}^* = \arg \max_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; \mathbf{t}).$$

- ML leads to a point estimate of $\boldsymbol{\theta}$ and is asymptotically consistent.
- The likelihood is unbounded, so ML estimator can overfit!

Maximum a posteriori (MAP) estimation: summary

Penalise unreasonable values (\sim regularisation) by imposing a prior distribution on the parameters:

$$p(\boldsymbol{\theta}|\mathbf{X}) \propto p(\mathbf{X}|\boldsymbol{\theta})p(\boldsymbol{\theta}).$$

Maximum a posteriori (MAP) estimation: summary

Penalise unreasonable values (\sim regularisation) by imposing a prior distribution on the parameters:

$$p(\boldsymbol{\theta}|\mathbf{X}) \propto p(\mathbf{X}|\boldsymbol{\theta})p(\boldsymbol{\theta}).$$

The goal is to maximise the penalised log-likelihood :

$$\ell_{\text{MAP}}(\boldsymbol{\theta}; \mathbf{t}) = \ell(\boldsymbol{\theta}; \mathbf{t}) + \ln p(\boldsymbol{\theta}).$$

- MAP leads to a point estimate of $\boldsymbol{\theta}$ and asymptotically agrees with ML estimate.
- MAP is not invariant under reparametrisation!

Outline

- 1 What is classification?
- 2 Decision theory
- 3 Generative classifiers
- 4 Discriminative classifiers**
- 5 Summary
- 6 Exercises

Generative versus discriminative classification

$$f(\mathbf{x}) = \begin{cases} +1 & \text{if } P(t = 1|\mathbf{x}) > \alpha, \\ -1 & \text{if } P(t = 1|\mathbf{x}) < \alpha. \end{cases}$$

Generative versus discriminative classification

$$f(\mathbf{x}) = \begin{cases} +1 & \text{if } P(t = 1|\mathbf{x}) > \alpha, \\ -1 & \text{if } P(t = 1|\mathbf{x}) < \alpha. \end{cases}$$

- Generative classifiers:

$$P(t = k|\mathbf{x}) \propto p(\mathbf{x}|t = k)P(t = k)$$

- ▶ Require explicit class-conditionals
- ▶ Take a linear form in specific cases

Generative versus discriminative classification

$$f(\mathbf{x}) = \begin{cases} +1 & \text{if } P(t = 1|\mathbf{x}) > \alpha, \\ -1 & \text{if } P(t = 1|\mathbf{x}) < \alpha. \end{cases}$$

- Generative classifiers:

$$P(t = k|\mathbf{x}) \propto p(\mathbf{x}|t = k)P(t = k)$$

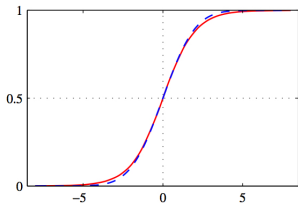
- ▶ Require explicit class-conditionals
- ▶ Take a linear form in specific cases

- Discriminative classifier:

$$P(t = k|\mathbf{x}) = \sigma(y(\mathbf{x}; \boldsymbol{\theta}))$$

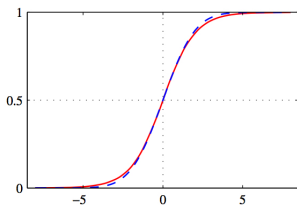
- ▶ Does not rely on class-conditionals
- ▶ Less parameters to learn (or optimise)
- ▶ Easy to change the feature map $\phi(\mathbf{x})$

(Binary) logistic regression



- Linear discriminant: $y(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) + b$.

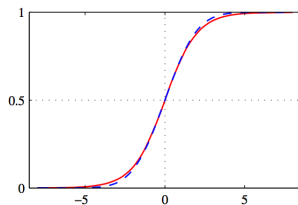
(Binary) logistic regression



- Linear discriminant: $y(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) + b$.
- Logistic link:

$$P(t = +1|\mathbf{x}) = \sigma(y(\mathbf{x})) = \frac{1}{1 + \exp(-y(\mathbf{x}))}.$$

(Binary) logistic regression



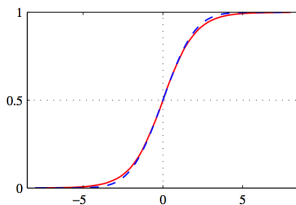
- Linear discriminant: $y(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) + b$.
- Logistic link:

$$P(t = +1|\mathbf{x}) = \sigma(y(\mathbf{x})) = \frac{1}{1 + \exp(-y(\mathbf{x}))}.$$

- Conditional likelihood:

$$t|\mathbf{x} \sim \text{Bernoulli}(\sigma(y(\mathbf{x}))) = \sigma(y(\mathbf{x}))^{\delta_{+1}(t)} (1 - \sigma(y(\mathbf{x})))^{\delta_{-1}(t)}.$$

(Binary) logistic regression



- Linear discriminant: $y(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x}) + b$.
- Logistic link:

$$P(t = +1|\mathbf{x}) = \sigma(y(\mathbf{x})) = \frac{1}{1 + \exp(-y(\mathbf{x}))}.$$

- Conditional likelihood:

$$t|\mathbf{x} \sim \text{Bernoulli}(\sigma(y(\mathbf{x}))) = \sigma(y(\mathbf{x}))^{\delta_{+1}(t)} (1 - \sigma(y(\mathbf{x})))^{\delta_{-1}(t)}.$$

- Alternative formulation: $P(t|\mathbf{x}) = \sigma(ty(\mathbf{x})). \quad (\star)$

How do we learn \mathbf{w} ?

$$\ln p(\mathbf{t}|\mathbf{x}; \mathbf{w}) = \sum_i \ln \text{Bernoulli}(\sigma(y(\mathbf{x}_i))) .$$

How do we learn \mathbf{w} ?

$$\ln p(\mathbf{t}|\mathbf{x}; \mathbf{w}) = \sum_i \ln \text{Bernoulli}(\sigma(y(\mathbf{x}_i))) .$$

Iterative reweighted least squares (IRLS):

$$\mathbf{w} \leftarrow \mathbf{w} + \underbrace{(\Phi^\top \mathbf{R} \Phi)^{-1}}_{=\mathbf{H}^{-1}} \underbrace{\Phi^\top (\boldsymbol{\sigma} - \mathbf{t})}_{=\nabla_{\mathbf{w}} \ln p(\mathbf{t}|\mathbf{w})}$$

where $\Phi = (\phi(\mathbf{x}_1)^\top, \dots, \phi(\mathbf{x}_n)^\top)^\top$, $R_{ii} = \sigma(y(\mathbf{x}_i))(1 - \sigma(y(\mathbf{x}_i)))$, $\boldsymbol{\sigma} = (\sigma(y(\mathbf{x}_1)), \dots, \sigma(y(\mathbf{x}_n)))^\top$ and $\mathbf{t} = (t_1, \dots, t_n)^\top$.

How do we learn \mathbf{w} ?

$$\ln p(\mathbf{t}|\mathbf{x}; \mathbf{w}) = \sum_i \ln \text{Bernoulli}(\sigma(y(\mathbf{x}_i))) .$$

Iterative reweighted least squares (IRLS):

$$\mathbf{w} \leftarrow \mathbf{w} + \underbrace{(\Phi^\top \mathbf{R} \Phi)^{-1}}_{=\mathbf{H}^{-1}} \underbrace{\Phi^\top (\boldsymbol{\sigma} - \mathbf{t})}_{=\nabla_{\mathbf{w}} \ln p(\mathbf{t}|\mathbf{w})}$$

where $\Phi = (\phi(\mathbf{x}_1)^\top, \dots, \phi(\mathbf{x}_n)^\top)^\top$, $R_{ii} = \sigma(y(\mathbf{x}_i))(1 - \sigma(y(\mathbf{x}_i)))$, $\boldsymbol{\sigma} = (\sigma(y(\mathbf{x}_1)), \dots, \sigma(y(\mathbf{x}_n)))^\top$ and $\mathbf{t} = (t_1, \dots, t_n)^\top$.

- ▶ Instantiation of Newton-Raphson
- ▶ Objective is convex!

How do we learn \mathbf{w} ?

$$\ln p(\mathbf{t}|\mathbf{x}; \mathbf{w}) = \sum_i \ln \text{Bernoulli}(\sigma(y(\mathbf{x}_i))) .$$

- 1 Iterative reweighted least squares (IRLS):

$$\mathbf{w} \leftarrow \mathbf{w} + \underbrace{(\Phi^\top \mathbf{R} \Phi)^{-1}}_{=\mathbf{H}^{-1}} \underbrace{\Phi^\top (\boldsymbol{\sigma} - \mathbf{t})}_{=\nabla_{\mathbf{w}} \ln p(\mathbf{t}|\mathbf{w})}$$

where $\Phi = (\phi(\mathbf{x}_1)^\top, \dots, \phi(\mathbf{x}_n)^\top)^\top$, $R_{ii} = \sigma(y(\mathbf{x}_i))(1 - \sigma(y(\mathbf{x}_i)))$, $\boldsymbol{\sigma} = (\sigma(y(\mathbf{x}_1)), \dots, \sigma(y(\mathbf{x}_n)))^\top$ and $\mathbf{t} = (t_1, \dots, t_n)^\top$.

- ▶ Instantiation of Newton-Raphson
- ▶ Objective is convex!

- 2 Alternatives include gradient descent

How do we learn \mathbf{w} ?

$$\ln p(\mathbf{t}|\mathbf{x}; \mathbf{w}) = \sum_i \ln \text{Bernoulli}(\sigma(y(\mathbf{x}_i))) .$$

- 1 Iterative reweighted least squares (IRLS):

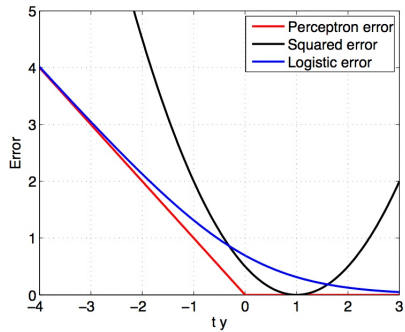
$$\mathbf{w} \leftarrow \mathbf{w} + \underbrace{(\Phi^\top \mathbf{R} \Phi)^{-1}}_{=\mathbf{H}^{-1}} \underbrace{\Phi^\top (\boldsymbol{\sigma} - \mathbf{t})}_{=\nabla_{\mathbf{w}} \ln p(\mathbf{t}|\mathbf{w})}$$

where $\Phi = (\phi(\mathbf{x}_1)^\top, \dots, \phi(\mathbf{x}_n)^\top)^\top$, $R_{ii} = \sigma(y(\mathbf{x}_i))(1 - \sigma(y(\mathbf{x}_i)))$, $\boldsymbol{\sigma} = (\sigma(y(\mathbf{x}_1)), \dots, \sigma(y(\mathbf{x}_n)))^\top$ and $\mathbf{t} = (t_1, \dots, t_n)^\top$.

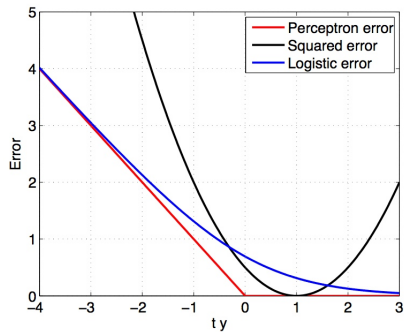
- ▶ Instantiation of Newton-Raphson
- ▶ Objective is convex!

- 2 Alternatives include gradient descent and stochastic gradient descent (\star)

Classification losses



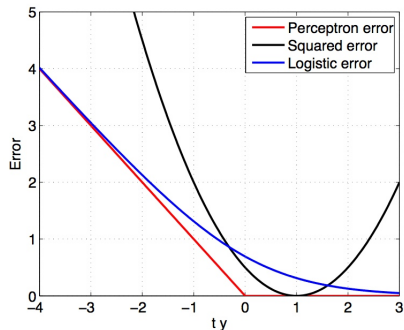
Classification losses



- Perceptron loss:

$$E(x_i) = \begin{cases} 0 & \text{if } t_i y(x_i) > 0, \\ t_i y(x_i) & \text{if } t_i y(x_i) < 0. \end{cases}$$

Classification losses



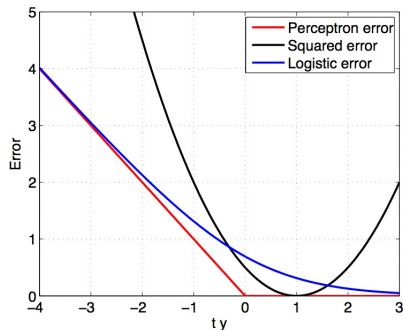
- Perceptron loss:

$$E(x_i) = \begin{cases} 0 & \text{if } t_i y(x_i) > 0, \\ t_i y(x_i) & \text{if } t_i y(x_i) < 0. \end{cases}$$

- Logistic loss:

$$E(x_i) = \ln \left(1 + e^{-t_i y(x_i)} \right).$$

Classification losses



- Perceptron loss:

$$E(x_i) = \begin{cases} 0 & \text{if } t_i y(x_i) > 0, \\ t_i y(x_i) & \text{if } t_i y(x_i) < 0. \end{cases}$$

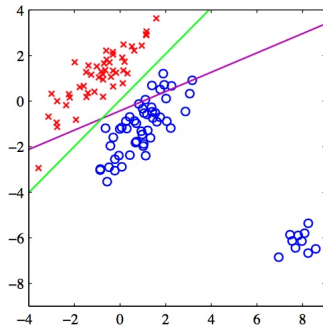
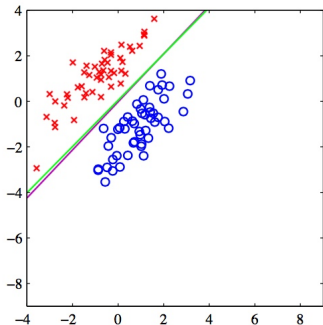
- Logistic loss:

$$E(x_i) = \ln \left(1 + e^{-t_i y(x_i)} \right).$$

- Squared error:

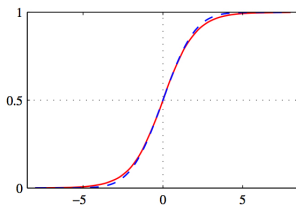
$$E(x_i) = \frac{1}{2} (t_i y(x_i) - 1)^2.$$

Is the squared error suitable for classification?



(Green: perceptron. Magenta: squared error.)

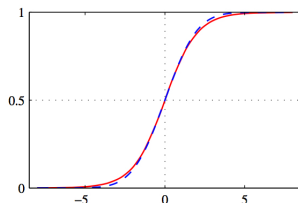
Other link functions?



- Probit regression:

$$\Phi(y(\mathbf{x})) = \int_{-\infty}^{y(\mathbf{x})} \text{Gaussian}(0, 1) dz, \quad y(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b.$$

Other link functions?



- Probit regression:

$$\Phi(y(\mathbf{x})) = \int_{-\infty}^{y(\mathbf{x})} \text{Gaussian}(0, 1) dz, \quad y(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b.$$

- Latent variable view:

$$t|z \sim I(tz > 0), \quad z \sim \text{Gaussian}(\mathbf{y}(\mathbf{x}), 1).$$

Multinomial logistic regression

- Linear discriminant:

$$y_k(\mathbf{x}) = \mathbf{w}_k^\top \phi(\mathbf{x}) + b_k, \quad k \in \{1, \dots, m\}.$$



Multinomial logistic regression



- Linear discriminant:

$$y_k(\mathbf{x}) = \mathbf{w}_k^\top \phi(\mathbf{x}) + b_k, \quad k \in \{1, \dots, m\}.$$

- Softmax:

$$P(t = k|x) = \frac{\exp(y(\mathbf{x}_k))}{\sum_l \exp(y(\mathbf{x}_l))}.$$

Multinomial logistic regression



- Linear discriminant:

$$y_k(\mathbf{x}) = \mathbf{w}_k^\top \phi(\mathbf{x}) + b_k, \quad k \in \{1, \dots, m\}.$$

- Softmax:

$$P(t = k|x) = \frac{\exp(y(\mathbf{x}_k))}{\sum_l \exp(y(\mathbf{x}_l))}.$$

- Conditional likelihood:

$$t|\mathbf{x} \sim \text{Categorical}(\boldsymbol{\mu}),$$

where $\mu_k = P(t = k|x)$.

Outline

- 1 What is classification?
- 2 Decision theory
- 3 Generative classifiers
- 4 Discriminative classifiers
- 5 Summary
- 6 Exercises

Summary



- Linear classifiers:
 - ▶ Perceptron
 - ▶ Naive Bayes
 - ▶ (Multi-nomial) logistic regression

Summary



- Linear classifiers:
 - ▶ Perceptron
 - ▶ Naive Bayes
 - ▶ (Multi-nomial) logistic regression
- Linear classifier can be non-linear!

Summary



- Linear classifiers:
 - ▶ Perceptron
 - ▶ Naive Bayes
 - ▶ (Multi-nomial) logistic regression
- Linear classifier can be non-linear!
- Techniques to learn the parameters

Summary



- Linear classifiers:
 - ▶ Perceptron
 - ▶ Naive Bayes
 - ▶ (Multi-nomial) logistic regression
- Linear classifier can be non-linear!
- Techniques to learn the parameters
- Trade-offs when making decisions

Outline

- 1 What is classification?
- 2 Decision theory
- 3 Generative classifiers
- 4 Discriminative classifiers
- 5 Summary
- 6 Exercises**

Exercise 1

Can you propose a Naive Bayes classifier with continuous features? Derive the maximum likelihood estimates of the parameters.

Exercise 2

Derive the update equations of a generative classifier with discrete binary features.

Exercise 3

What is the form of the decision boundary for a binary classifier with Gaussian features with different covariance matrices?

Exercise 3^{*}

What are the expressions of the precision and the recall in the multi-class case?

References

- J. H. Albert, and S. Chib (1993): *Bayesian Analysis of Binary and Polychotomous Response Data*. Journal of the American Statistical Association 88 (422): 669-679.
- C. Bishop: *Pattern Recognition and Machine Learning*. Springer, 2006.
- J. Davis and M. Goadrich: *The Relationship Between Precision-Recall and ROC Curves*. ICML 2006.
- Y. Ng and M. Jordan: *On Discriminative vs. Generative classifiers: A comparison of Logistic Regression and Naive Bayes*. NIPS 2001.
- M. Seeger: *Pattern Classification and Machine Learning*. Lecture notes EPFL, 2012.

