# Lecture 7
## Probabilistic Models

M. Poel
G. Englebienne

University of Twente

Example (Iris Data)

## Example (Iris Data)

Example (Iris Data)



What is the probability that this is a △

We want to know

$$p(\mathcal{C}|\mathbf{x})$$

Two possibilities:

1. Discrimintative Models: learn a $p(\mathcal{C}|\mathbf{x})$ directly

$$p(\mathcal{C}|\mathbf{x}) = f(\mathbf{x}, \boldsymbol{\theta})$$

For example, logistic regression: $p(\mathcal{C}|\mathbf{x}) = \sigma(\mathbf{x}, \mathbf{w})$

2. Generative model: learn the class-wise distribution of your data, and use Bayes' rule

$$p(\mathcal{C}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C})\, p(\mathcal{C})}{p(\mathbf{x})}$$

$$p(\mathcal{C}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C})\, p(\mathcal{C})}{p(\mathbf{x})}$$

Posterior

Likelihood

Prior

Evidence

Questions:

1. How do we compute $p(\mathbf{x}|\mathcal{C})$?

2. How do we use $p(\mathcal{C}|\mathbf{x})$ to make decisions?

$$p(\mathcal{C}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C})\,p(\mathcal{C})}{p(\mathbf{x})}$$

Posterior — Likelihood — Prior — Evidence

Questions:

1. How do we compute $p(\mathbf{x}|\mathcal{C})$?
2. How do we use $p(\mathcal{C}|\mathbf{x})$ to make decisions?

### Example (Biased Coins)

Consider two otherwise undistinguishable coins; coin A is fair, coin B is biased and has a probability of 0.7 of landing on its head. You pick a coin at random and throw the following sequence:

Tail Head Tail Head Tail Head Head

What is the probability that this is coin A?

### Example (Biased Coins)

First, we notice that the throws do not affect each other:

$$p(t, h, t, h, t, h, h) = p(t)\, p(h)\, p(t)\, p(h)\, p(t)\, p(h)\, p(h)$$
$$= p(t)^3 p(h)^4$$

The throws depend on which coin we used, and we picked one of the two coins with equal probability, so we can use Bayes' theorem to compute the posterior probability that coin B has been used:

$$p(C = a|X = \{t, h, t, h, t, h, h\}) = \frac{p(t, h, t, h, t, h, h|a)\, p(a)}{p(t, h, t, h, t, h, h)}$$
$$= \frac{0.5^3\, 0.5^4\, 0.5}{0.5^3\, 0.5^4\, 0.5 + 0.3^3\, 0.7^4\, 0.5}$$
$$\approx 0.55 > 0.5$$

# The Bernoulli Distribution

The Bernoulli distribution is the probability mass function of a binary event $x \in \{0, 1\}$:

$$p(x|\mu) = \mu^x (1 - \mu)^{1-x}$$

This leads to the following log-likelihood of a set of i.i.d. measurements:

$$\mathcal{L}(x_{1\ldots N}; \mu) = \log \prod_{n=1}^{N} \mu^{x_n} (1 - \mu)^{1-x_n}$$
$$= \sum_{n=1}^{N} x_n \log \mu + (1 - x_n) \log(1 - \mu)$$

The Maximum-Likelihood estimator for the parameter $\mu$ is found by maximising the log-likelihood

We find the maximum of the function by setting its first derivative to zero:

$$\frac{\partial}{\partial \mu} \mathcal{L}(x_{1\dots N}; \mu) = \frac{\partial}{\partial \mu} \sum_{n=1}^{N} x_n \log \mu + (1 - x_n) \log(1 - \mu)$$

$$0 = \sum_{n=1}^{N} \frac{x_n}{\mu} - \frac{1 - x_n}{1 - \mu} = \sum_{n=1}^{N} \frac{x_n (1 - \mu) - (1 - x_n) \mu}{\mu (1 - \mu)}$$

$$0 = \sum_{n=1}^{N} x_n - \mu \mathbf{x}_n - \mu + \mu x_n = \sum_{n=1}^{N} x_n - \mu$$

$$\mu = \frac{\sum_{n=1}^{N} x_n}{N}$$

The Bernouli distribution can be extended to multiple classes

- Using one-of-K encoding ($x = k \in \{1, \ldots, K\} \Rightarrow \mathbf{x} : x_{j \neq k} = 0, x_k = 1$)

$$p(\mathbf{x}) = \prod_{i=1}^{K} p_i^{x_i} \quad \text{where} \quad \sum_{i=1}^{K} p_i = 1$$

- The MLE estimator for $p_i$ is given by:

$$p_i = \frac{\sum_{n=1}^{N} x_i^{(n)}}{N}$$

Example (Waiting for the bus)

What is the probability that you will have to wait *exactly* 5 minutes for your bus?

### Example (Waiting for the bus)

What is the probability that you will have to wait *exactly* 5 minutes for your bus?

zero.

We cannot assign probabilities to values. Instead, we need to assign probabilities to value *ranges*

### Example (A better question would be:)

What is the probability that you will have to wait between 5 and 6 minutes for your bus?

Probability Density Functions (PDFs):

- Are non-negative everywhere, integrate to 1
- Not probabilities, can be interpreted as *relative* probabilities
  - Can be larger than 1

Probability Density Functions (PDFs):

▶ Are non-negative everywhere, integrate to 1
▶ Not probabilities, can be interpreted as *relative* probabilities
  ▶ Can be larger than 1

# Probability Density Functions

Probability Density Functions (PDFs):

- ▶ Are non-negative everywhere, integrate to 1
- ▶ Not probabilities, can be interpreted as *relative* probabilities
  - ▶ Can be larger than 1

Probability Density Functions (PDFs):

- ▶ Are non-negative everywhere, integrate to 1
- ▶ Not probabilities, can be interpreted as *relative* probabilities
  - ▶ Can be larger than 1

Probability Density Functions (PDFs):

- Are non-negative everywhere, integrate to 1
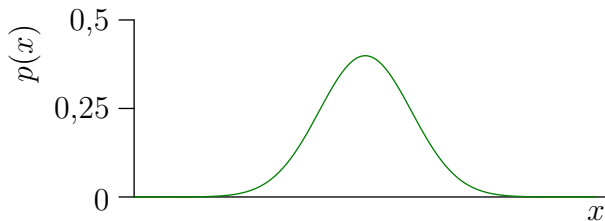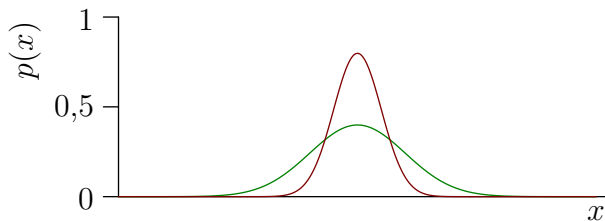- Not probabilities, can be interpreted as *relative* probabilities
  - Can be larger than 1

$$\mathcal{N}(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp{-\frac{(x - \mu)^2}{2\sigma^2}} \qquad (1)$$

# The Gaussian or Normal distribution

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{d/1}|\boldsymbol{\Sigma}|^{1/2}} \exp{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^{\top}\boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})} \qquad (2)$$

The central limit theorem can informally be stated as follows:

### The Central Limit Theorem
The sum of a sufficiently large number of independent, identically distributed variables with finite variance will have an approximately Gaussian distribution.

Notice that no assumption is made about the distribution of these variables

The central limit theorem can informally be stated as follows:

Notice that

The central limit theorem can informally be stated as follows:

### The Central Limit
The sum of a sufficiently large number of independent, identically distributed random variables with finite variance will have an approximately normal distribution.

Notice that no assumption is made about those variables

The central limit theorem can informally be stated as follows:

### The Central Limit Theorem

The sum of a sufficiently large number of independent, identically distributed random variables with finite variance will have an approximately normal distribution.

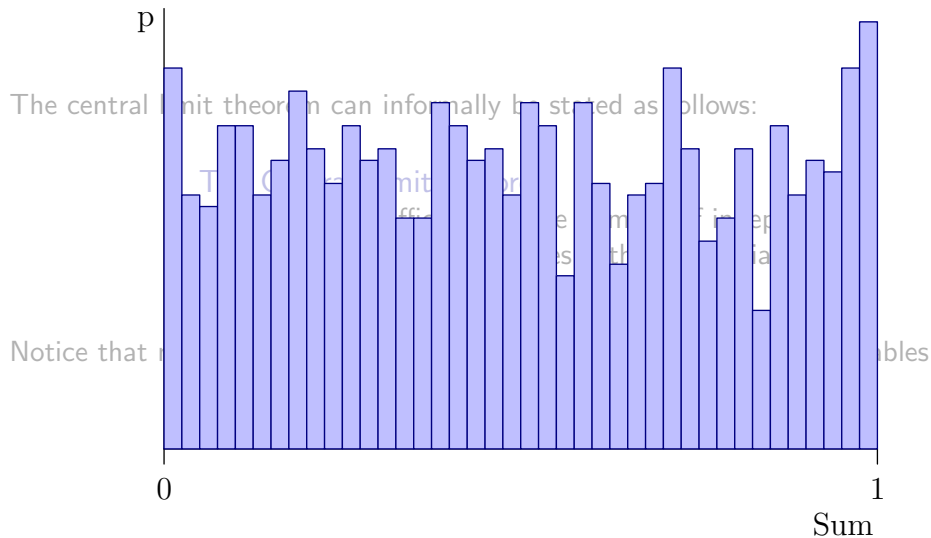Notice that no assumption is made about the distribution of these variables

The central limit theorem can informally be stated as follows:

### The Central Limit Theorem
The sum of a sufficiently large number of independent, identically distributed random variables with finite variance will have an approximately normal distribution.

Notice that no assumption is made about the distribution of these variables

p

The central limit theorem can inform. be stated as follows:

The Central Limit T
The sum of a sufficie number of independent,
identically distribute ith finite variance will
have an approxima distribution.

Notice that no assumption is m tribution of these variables

0                                                                          10
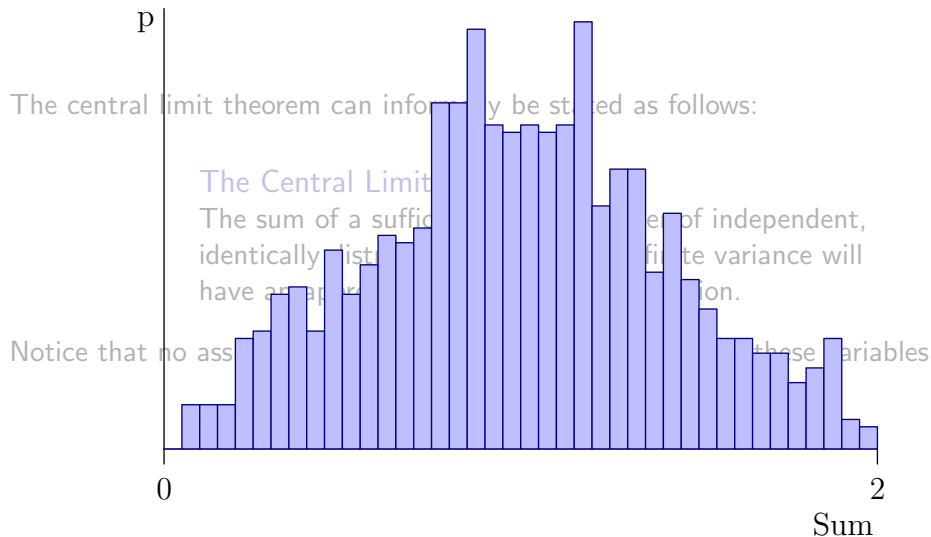
Sum

p

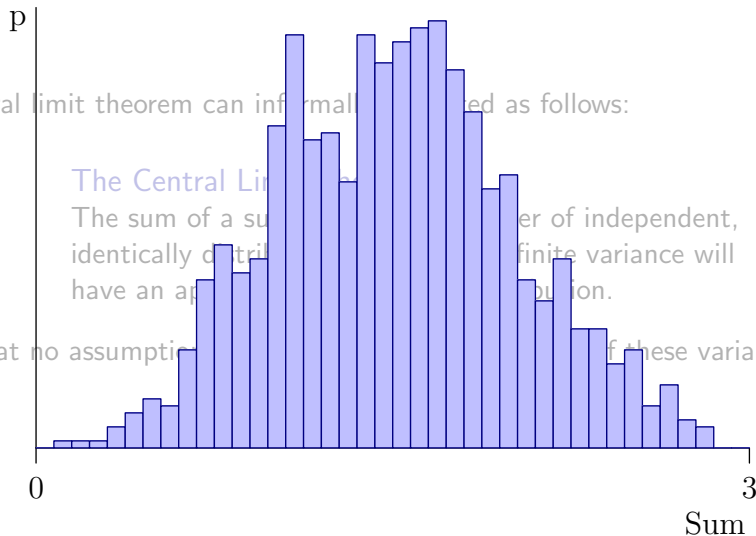The central limit theorem can informally be stated as follows:
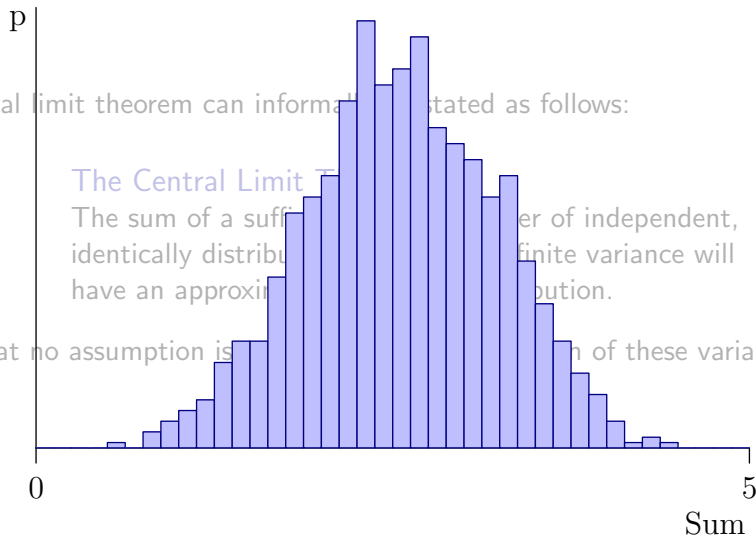
### The Central Limit Theorem
The sum of a sufficiently large number of independent, identically distributed variables with finite variance will have an approximately Gaussian distribution.
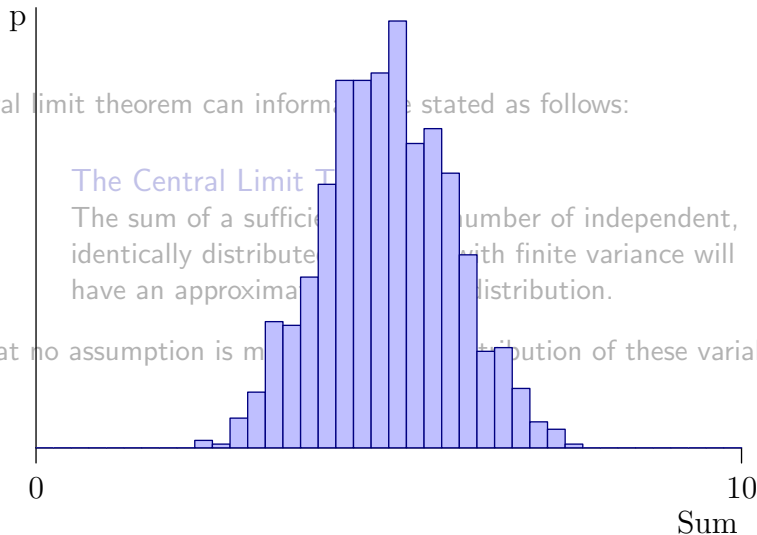
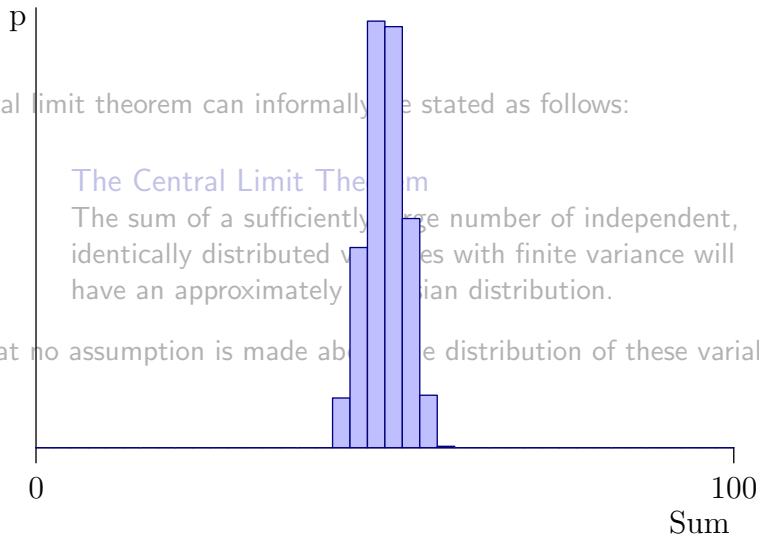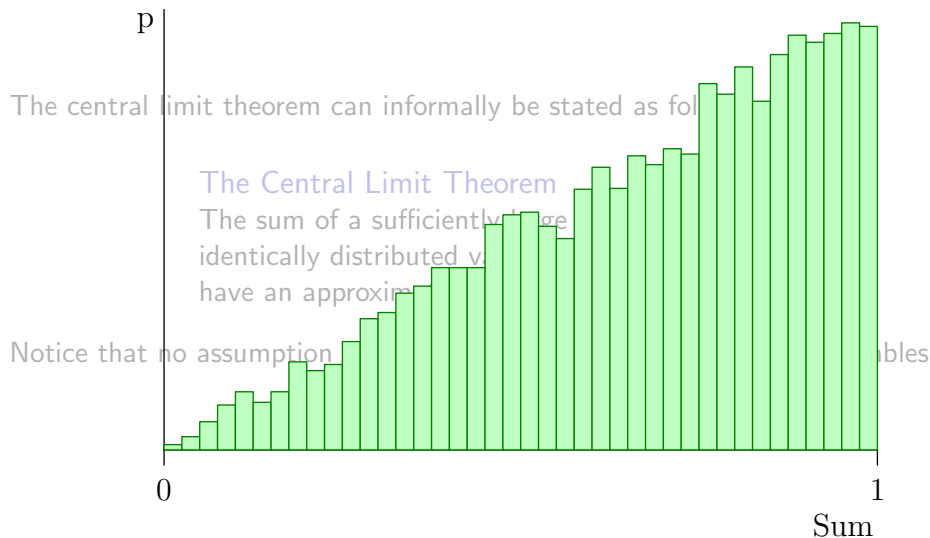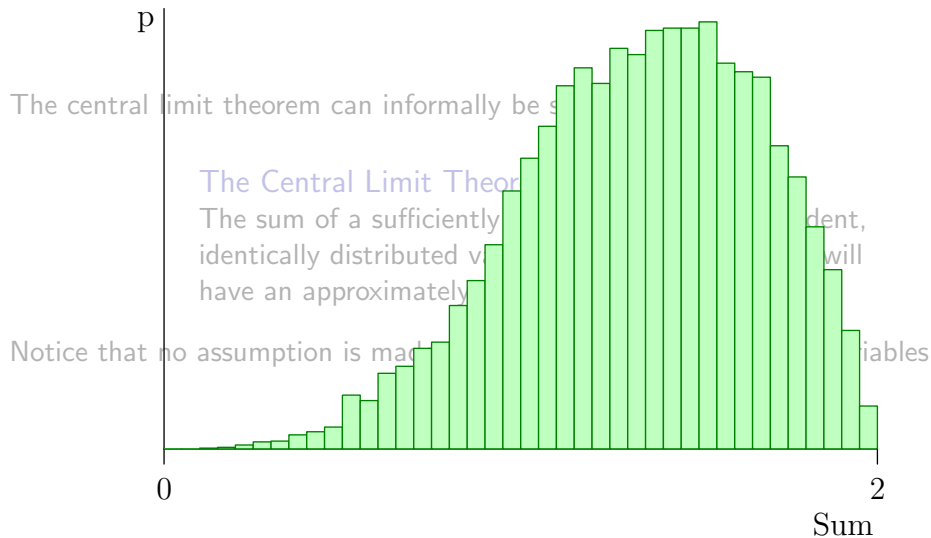Notice that no assumption is made about the distribution of these variables

0                                                                      100

Sum

The central limit theorem can informally be stated as fol

### The Central Limit Theorem
The sum of a sufficiently large
identically distributed va
have an approxim

Notice that no assumption                                                                                     bles

p

The central limit theorem can informally be s

### The Central Limit Theor
The sum of a sufficiently                           dent,
identically distributed v                           will
have an approximately

Notice that no assumption is mad                                    iables

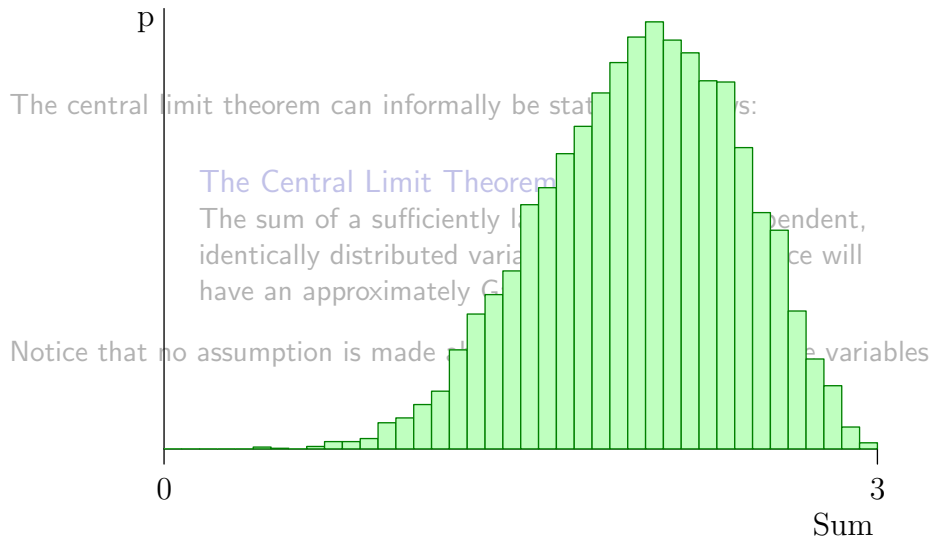0                                                              2

Sum

The central limit theorem can informally be stated as follows:

### The Central Limit Theorem
The sum of a sufficiently large number of independent, identically distributed variables with finite variance will have an approximately Gaussian distribution.

Notice that no assumption is made about the variables

The central limit theorem can informally be stated as follows:

### The Central Limit Theorem
The sum of a sufficiently large number of independent, identically distributed variables with finite variance will have an approximately Gaussian distribution.

Notice that no assumption is made about the distribution of these variables

The central limit theorem can informally be stated as follows:

### The Central Limit Theorem
The sum of a sufficiently large number of independent, identically distributed variables with finite variance will have an approximately Gaussian distribution.

Notice that no assumption is made about the nature of these variables

p

The central limit theorem can informally be stated as follows:

### The Central Limit Theorem
The sum of a sufficiently large number of independent, identically distributed variables with finite variance will have an approximately Gaussian distribution.

Notice that no assumption is made about the distribution of these variables

0        100

Sum

The expectation of a variable (or function of that variable)

$$\mathbb{E}_{p(x)}[f(\mathbf{x})] = \int f(\mathbf{x})\, p(\mathbf{x})\, d\mathbf{x}$$

is the value that the long-term average of the function converges to.

- From our perspective, it's our "best bet" value
- When making decision, we want to minimise our expected error

# Decision threshold

Classification: obtain a feature vector $\mathbf{x}$ and predict the corresponding class $\mathcal{C}$

### Example

Given an X-ray image, predict the health state of the person

## Bayesian decision rule

- Assign class labels so as to minimise the expected error
- Assign an observation $\mathbf{x}$ to class $\mathcal{C}_i$ if

$$p(\mathcal{C}_i|\mathbf{x}) > p(\mathcal{C}_j|\mathbf{x}) \qquad \forall j \neq i \qquad (3)$$

From Bayes' rule, we have that

$$p(\mathcal{C}_k|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)}{\sum_j p(\mathbf{x}|\mathcal{C}_j)p(\mathcal{C}_j)} \quad (4)$$

We want to minimise the probability of a mistake, that is:

$$p(\text{mistake}) = p(\mathbf{x} \in \mathcal{R}_1, \mathcal{C}_2) + p(\mathbf{x} \in \mathcal{R}_2, \mathcal{C}_1) \quad (5)$$
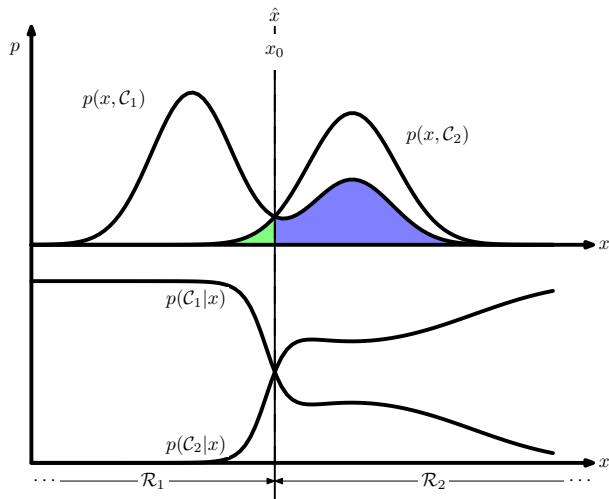
$$= \int_{\mathcal{R}_1} p(\mathbf{x}, \mathcal{C}_2)d\mathbf{x} + \int_{\mathcal{R}_2} p(\mathbf{x}, \mathcal{C}_1) \quad (6)$$

Since $p(\mathbf{x}, \mathcal{C}_k) = p(\mathcal{C}_k|\mathbf{x})p(\mathbf{x})$ and $p(\mathbf{x})$ is the same in both terms, $p(\text{mistake})$ is minimal if each point $\mathbf{x}$ is assigned to the class for which $p(\mathcal{C}_k|\mathbf{x})$ is largest.
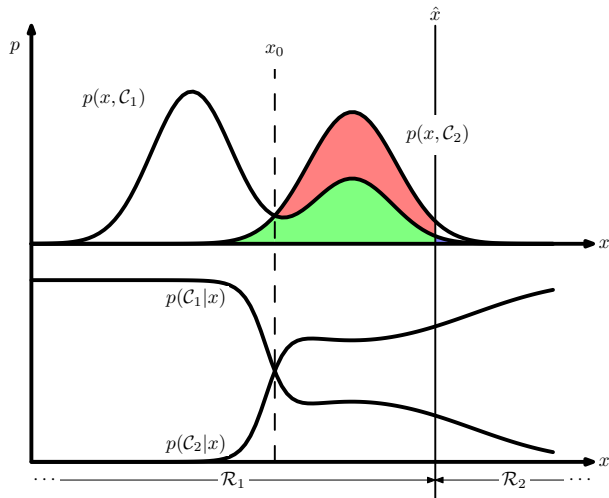
In some cases, the posterior probability $p(\mathcal{C}_k|\mathbf{x})$ of the most likely class may be far less than one.

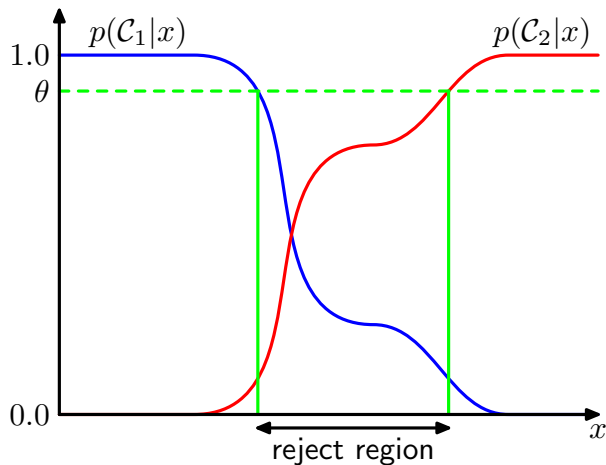- The regions where this is the case lead to most misclassifications

In some cases it is better to avoid making a decision when that is the case, in order to improve the performance on the examples for which a decision is made.

## Example

In medical image classification, it may be suitable to automatically classify images for which we are very confident and leave the difficult cases for a human to evaluate.

# The reject option

Achieved by choosing a threshold, $\theta$, and rejecting datapoints for which the largest $p(\mathcal{C}_k|\mathbf{x}) \leqslant \theta$.

In the case of unbalanced misclassification costs: loss matrix

Cancer classification example

$$L = \begin{pmatrix} 0 & 1000 \\ 1 & 0 \end{pmatrix} \tag{7}$$

The expected loss is then given by

$$\mathbb{E}[L] = \sum_k \sum_j \int_{\mathcal{R}_j} L_{kj} p(\mathbf{x}, \mathcal{C}_k) d\mathbf{x} \tag{8}$$

which is minimised by assigning each datapoint $\mathbf{x}$ to the class $j$ for which

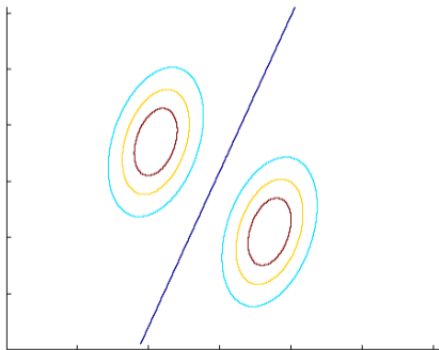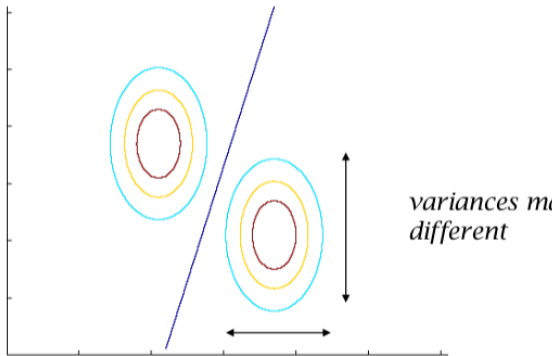$$\sum_k L_{kj} p(\mathcal{C}_k | \mathbf{x}) \tag{9}$$

We often lack data to train a model well by MLE

- ▶ To address this, we can simplify the model
- ▶ Reduce the number of parameters without losing too much discriminating power
- ▶ Examples (Gaussian PDF):
  - ▶ Assume classes share the same covariance
  - ▶ Assume the covariances are diagonal
  - ▶ Assume that covariance matrices are spherical



*likelihoods*

*posterior for $C_1$*

*discriminant:*
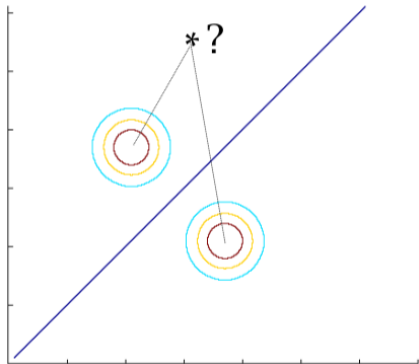$P(C_1|\boldsymbol{x}) = 0.5$

# Simplifying Assumptions

We often lack data to train a model well by MLE

- To address this, we can simplify the model
- Reduce the number of parameters without losing too much discriminating power
- Examples (Gaussian PDF):
  - Assume classes share the same covariance
  - Assume the covariances are diagonal
  - Assume that covariance matrices are spherical

We often lack data to train a model well by MLE

- To address this, we can simplify the model
- Reduce the number of parameters without losing too much discriminating power
- Examples (Gaussian PDF):
  - Assume classes share the same covariance
  - Assume the covariances are diagonal
  - Assume that covariance matrices are spherical



variances m...
different

# Simplifying Assumptions

We often lack data to train a model well by MLE

▶ To address this, we can simplify the model
▶ Reduce the number of parameters without losing too much discriminating power
▶ Examples (Gaussian PDF):
  ▶ Assume classes share the same covariance
  ▶ Assume the covariances are diagonal
  ▶ Assume that covariance matrices are spherical

Naive Bayes: Assume all data dimensions are independent given the class

$$p(\mathbf{x}|\mathcal{C}) = p(x_1|\mathcal{C}) \cdots p(x_N|\mathcal{C})$$
$$= \prod_i p(x_i|\mathcal{C})$$

Features:

- ▶ Scales linearly in the number of features
- ▶ Overly confident if features are not independent
- ▶ Performs surprisingly well in practice
- ▶ Beware: Naive Bayes is not "Bayesian Learning"
- ▶ Notice: conditional independence $\neq$ marginal independence

$$p(x_1, \ldots, x_n) = \sum_{\mathcal{C}} p(x_1|\mathcal{C}) \cdots p(x_N|\mathcal{C}) p(\mathcal{C}) \neq p(x_1) \cdots p(x_N)$$

# Naive Bayes

Naive Bayes: Assume all data dimensions are independent given the class

$$p(\mathbf{x}|\mathcal{C}) = p(x_1|\mathcal{C}) \cdots p(x_N|\mathcal{C})$$
$$= \prod_i p(x_i|\mathcal{C})$$

Features:

- ▶ Scales linearly in the number of features
- ▶ Overly confident if features are not independent
- ▶ Performs surprisingly well in practice
- ▶ Beware: Naive Bayes is not "Bayesian Learning"
- ▶ Notice: conditional independence $\neq$ marginal independence

$$p(x_1, \ldots, x_n) = \sum_{\mathcal{C}} p(x_1|\mathcal{C}) \cdots p(x_N|\mathcal{C})p(\mathcal{C}) \neq p(x_1) \cdots p(x_N)$$
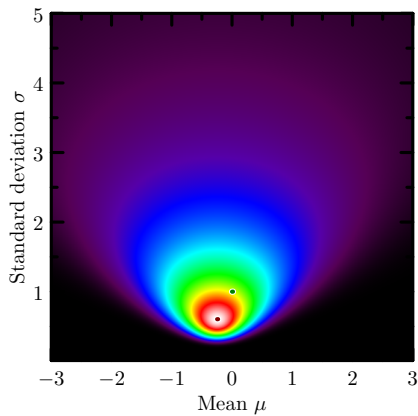
Instead of learning the parameters that maximise the likelihood, why not learn the most likely parameters? Using Bayes' rule, we have:
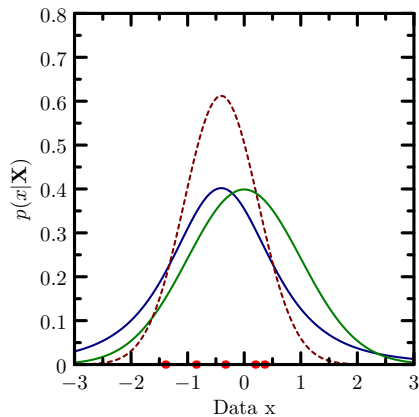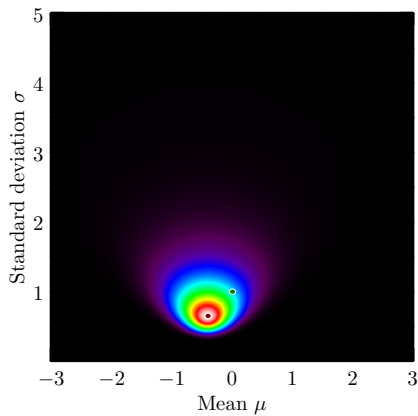
$$p(\boldsymbol{\theta}|\mathbf{x}) = \frac{p(\mathbf{x}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{\int p(\mathbf{x})d\boldsymbol{\theta}} \tag{10}$$
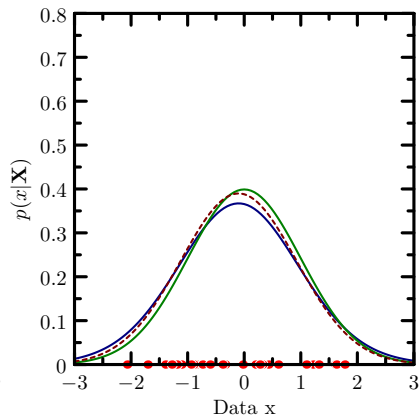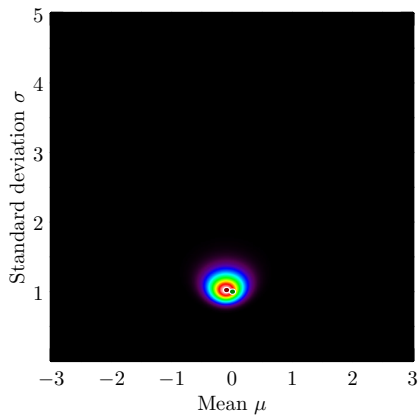
This requires us to place a prior over the parameter values
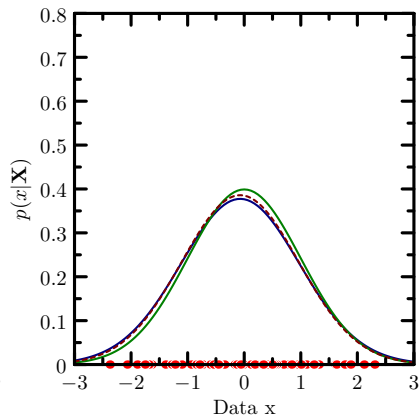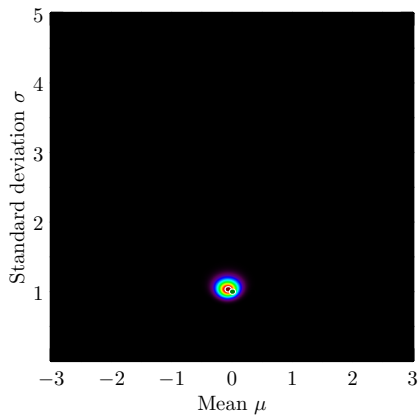
▶ Any prior is possible, choose prior to reflect prior knowledge

▶ If we use a Gaussian distribution with zero mean, this is equivalent to ML learning with parameter shrinkage

▶ The denominator is often intractable to compute but is constant, so that

$$\underset{\boldsymbol{\theta}}{\arg\max} \frac{p(\mathbf{x}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{\int p(\mathbf{x})d\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\arg\max}\, p(\mathbf{x}|\boldsymbol{\theta})p(\boldsymbol{\theta}) \tag{11}$$

In fact, we're not really interested in knowing the original distribution that "generated" the data

- ▶ We'll never know that anyway

What we really want to do, is to use the knowledge that we have in an optimal way. That is, we want

$$p(t|\mathbf{x}, \mathbf{X}, \mathbf{t}) = \int p(t|\mathbf{x}, \boldsymbol{\theta}) p(\boldsymbol{\theta}|\mathbf{X}, \mathbf{t}) \mathrm{d}\boldsymbol{\theta} \tag{12}$$

In effect, we consider all the models (of the form that we have chosen beforehand) that could have generated the data, and weigh their prediction according to how probable they are.

We've talked about Parametric Modelling

- ▶ Generative vs. Discriminative Models
- ▶ Probability Mass and Density Functions
- ▶ Learning the parameters of said functions
- ▶ Effect of Independence assumptions