

## Three lecture block

### Lecture 5: Probabilistic Reasoning

DD2431

Giampiero Salvi

Autumn, 2016

- probabilistic reasoning
- learning as inference
- learning with latent variables

Giampiero Salvi

Lecture 5: Probabilistic Reasoning

Constructive Alignment HT2017

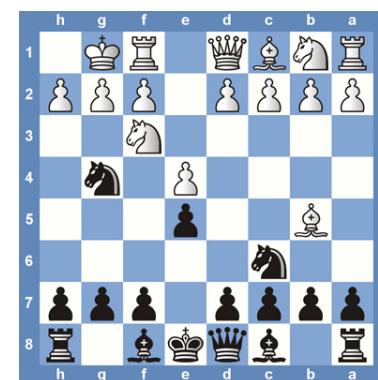
Giampiero Salvi

Lecture 5: Probabilistic Reasoning

Why Machine Learning?

- assignments: part of Lab3 is relevant (Naïve Bayes)
- two exam problems are on probabilistic methods
- the topic is fundamental for the advanced course (DD2434)

AI in the 1970s

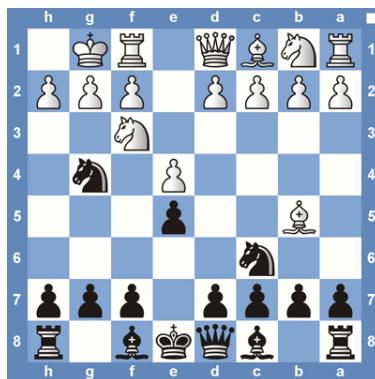


AI today



## Why Machine Learning?

AI in the 1970s



We need to deal with uncertainty!

AI today



## Two views of the universe

### Determinism

If all forces and positions of objects are known, and sufficient computing resources, then no uncertainty on the future (P-S. Laplace, A. Einstein, ...)

Quantum mechanics poses serious doubts to this view

- Heisenberg's uncertainty principle
- Schrödinger's cat

(N. Bohr, W. Heisenberg, S. Hawking, ...)

## Two views of the universe

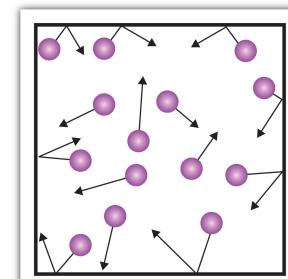
Regardless of the view, all agree we need to deal with uncertainty, because:

- measurements not accurate enough
- not enough computing power
- need to simplify the problems

## Two views of the universe

Regardless of the view, all agree we need to deal with uncertainty, because:

- measurements not accurate enough
- not enough computing power
- need to simplify the problems

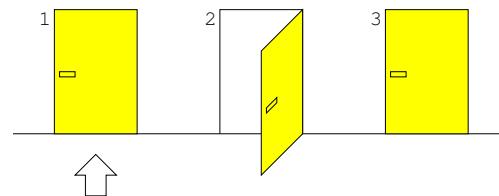


Example: pressure in gases vs particle impacts

## Subjective Uncertainty

- not only a description of randomness
- but rather degree of belief

Example: Monty Hall Problem



Giampiero Salvi

Lecture 5: Probabilistic Reasoning

## Heuristics

### Heuristic

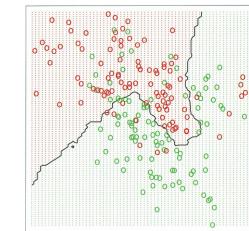
experience-based techniques for problem solving, learning, and discovery that give a solution which is not guaranteed to be optimal (Wikipedia)

Typical examples:

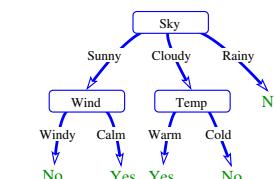
- Artificial Neural Networks
- Decision Trees
- Evolutionary methods
- $k$ -nearest neighbor

## Examples of ML Methods seen so far

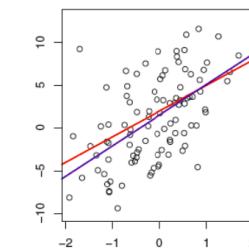
### $k$ Nearest Neighbour



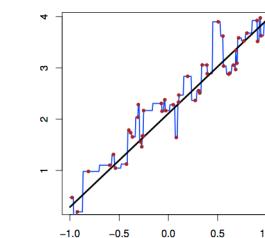
### Decision Trees



### Least squares Regression



### $k$ -NN Regression



Giampiero Salvi

Lecture 5: Probabilistic Reasoning

## Heuristics

### Heuristic

experience-based techniques for problem solving, learning, and discovery that give a solution which is not guaranteed to be optimal (Wikipedia)

Typical examples:

- Artificial Neural Networks
- Decision Trees
- Evolutionary methods
- $k$ -nearest neighbor

we need a more unified theory for ML

## Probability Theory in ML

incorporate probabilistic thinking at all levels

- knowledge is uncertain
- use observations to reduce uncertainty
- uncertainty propagation
- probability distributions as carriers of information<sup>1</sup>

<sup>1</sup>E T Jaynes. *Probability theory: The logic of science*. Ed. by G Larry Bretthorst. Cambridge university press, June 2003.

## Engineering vs Science

Engineering:

- ML as collection of methods
- fine tune aspects to boost the results

Science:

- define unified theory
- give the deepest possible interpretation to the results

reality not 100% clear cut

## Engineering vs Science

Engineering:

- ML as collection of methods
- fine tune aspects to boost the results

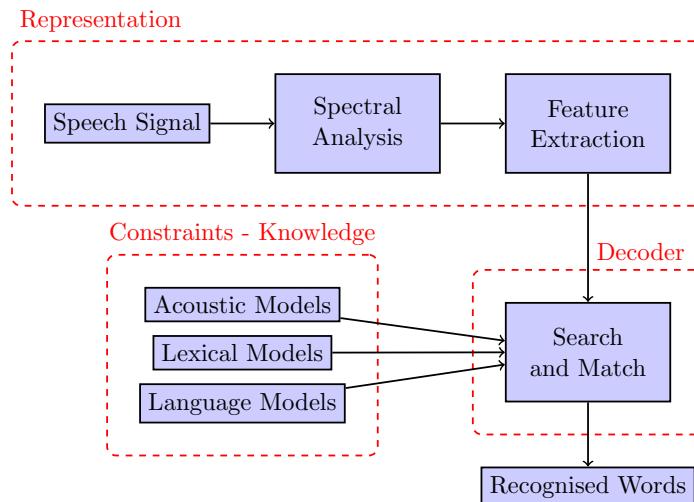
Science:

- define unified theory
- give the deepest possible interpretation to the results

## Advantages of Probability Based Methods

- **Results are interpretable.** More transparent and mathematically rigorous than methods such as *ANN*, *Evolutionary methods*.
- **Tool for interpreting other methods.** And make the assumptions explicit — *concept learning, least squares*.
- **Work with sparse training data.** More powerful than deterministic methods when training data is sparse (framework for including prior knowledge).
- **Belief Propagation:** Easy to merge different parts of a complex system and to update current knowledge with new observations.

## Example: Automatic Speech Recognition



## Advantages of Probability Based Methods, ctnd.

- **Shape a way of thinking.** All aspects of learning, modelling and inference can be cast under the same theory.

## Disadvantages of Probability Based Methods

- **Often hard to derive closed solutions.** Need to resort to heuristic approximations.
- **Inefficient for large data sets.** But many argue that the need for large data set is a flaw in the methods.

## General Formulation

Two (random) variables,  $\mathbf{X}$ ,  $\mathbf{Y}$

- Input data  $\mathbf{x}_i \in \mathbb{R}^q$
- Output data  $\mathbf{y}_i \in \mathbb{R}^D$

Relationship:  $f : \mathbf{X} \rightarrow \mathbf{Y}$

## Example: Linear regression

Model:

$$\mathbf{y}_i = \mathbf{W}\mathbf{x}_i + \epsilon$$

Learning: estimate

$$p(\mathbf{W}|\mathbf{Y}, \mathbf{X})$$

Regression: estimate

$$p(\mathbf{Y}|\mathbf{W}, \mathbf{X})$$

conditional distribution *after* we have observed the data.

## Outline

- ① Probability Theory Basics
- ② Common Distributions
- ③ Bayesian Decision Theory: Classification

**Reading:** Prince, S.J.D., Part I (Chapters 2, 3, 5)

**More about probabilistic Learning:**

Bishop, C. M. Pattern Recognition and Machine Learning,  
Springer.

## Outline

## Different views on probabilities

① Probability Theory Basics

Axiomatic defines axioms and derives properties

② Common Distributions

Classical number of ways something can happen over total  
number of things that can happen (e.g. dice)

Logical same, but weight the different ways

Frequency frequency of success in repeated experiments

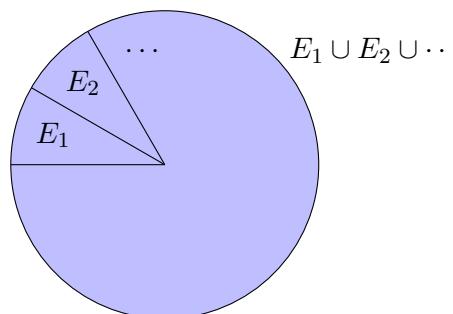
Subjective degree of belief (basis for Bayesian statistics)

③ Bayesian Decision Theory: Classification

## Axiomatic definition of probabilities (Kolmogorov)

Given an event  $E$  in a event space  $F$

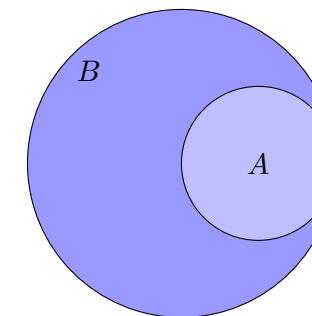
- ①  $P(E) \geq 0$  for all  $E \in F$
- ② sure event  $\Omega$ :  $P(\Omega) = 1$
- ③  $E_1, E_2, \dots$  countable sequence of pairwise disjoint events, then



$$P(E_1 \cup E_2 \cup \dots) = \sum_{i=1}^{\infty} P(E_i)$$

## Consequences

- ① Monotonicity:  $P(A) \leq P(B)$  if  $A \subseteq B$



Example:  $A = \{3\}$ ,  $B = \{\text{odd}\}$

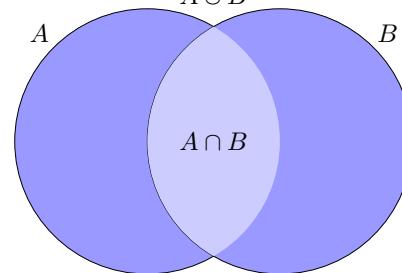
- ② Empty set  $\emptyset$ :  $P(\emptyset) = 0$

Example:  $P(A \cap B)$  where  $A = \{\text{odd}\}, B = \{\text{even}\}$

- ③ Bounds:  $0 \leq P(E) \leq 1$  for all  $E \in F$

## More Consequences: Addition

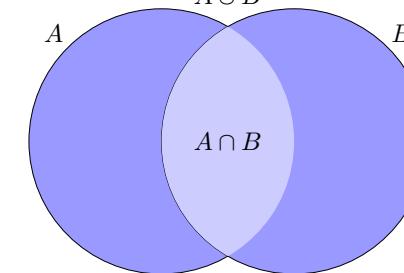
$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$



$$\begin{array}{lll} A & = & \{1, 3, 5\}, & P(A) & = & \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = \frac{1}{2} \\ \text{Example: } B & = & \{5, 6\}, & P(B) & = & \frac{1}{6} + \frac{1}{6} = \frac{1}{3} \end{array}$$

## More Consequences: Addition

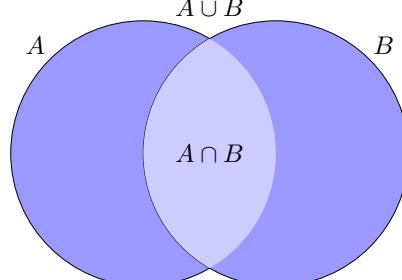
$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$



$$\begin{array}{lll} A & = & \{1, 3, 5\}, & P(A) & = & \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = \frac{1}{2} \\ \text{Example: } B & = & \{5, 6\}, & P(B) & = & \frac{1}{6} + \frac{1}{6} = \frac{1}{3} \\ A \cap B & = & \{5\} & P(A \cap B) & = & \frac{1}{6} \end{array}$$

## More Consequences: Addition

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$



**Example:**

$A = \{1, 3, 5\}$ ,	$P(A) = \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = \frac{1}{2}$
$B = \{5, 6\}$ ,	$P(B) = \frac{1}{6} + \frac{1}{6} = \frac{1}{3}$
$A \cap B = \{5\}$	$P(A \cap B) = \frac{1}{6}$
$A \cup B = \{1, 3, 5, 6\}$	$P(A \cup B) = \frac{1}{6} + \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = \frac{2}{3}$

## Random (Stochastic) Variables

A random variable is a **function** that assigns a number  $x$  to the outcome of an experiment

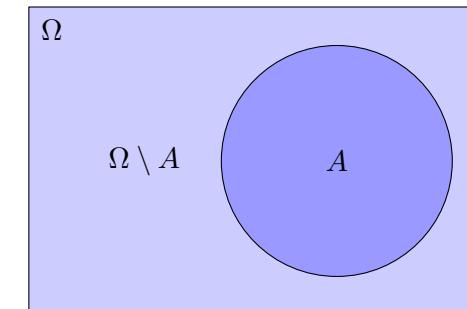
- the result of flipping a coin,
- the result of measuring the temperature

The *probability distribution*  $P(x)$  of a random variable (r.v.) captures the fact that

- the r.v. will have different values when observed **and**
- some values occur more than others.

## More Consequences: Negation

$$P(\bar{A}) = P(\Omega \setminus A) = 1 - P(A)$$



**Example:**

$A = \{1, 2\}$ ,	$P(A) = \frac{1}{6} + \frac{1}{6} = \frac{1}{3}$
$\bar{A} = \{3, 4, 5, 6\}$ ,	$P(\bar{A}) = \frac{1}{6} + \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = 1 - \frac{1}{3}$

## Formal definition of RVs

$$RV = \{f : \mathcal{S}_a \rightarrow \mathcal{S}_b, P(x)\}$$

where:

$\mathcal{S}_a$  = set of possible outcomes of the experiment

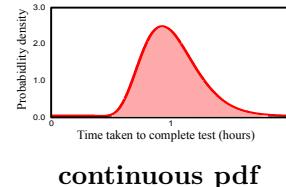
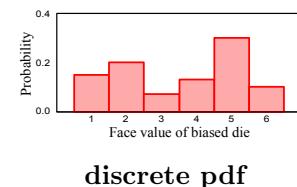
$\mathcal{S}_b$  = domain of the variable

$f : \mathcal{S}_a \rightarrow \mathcal{S}_b$  = function mapping outcomes to values  $x$

$P(x)$  = probability distribution function

## Types of Random Variables

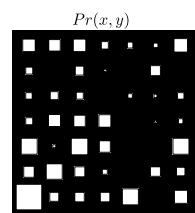
- A **discrete random variable** takes values from a predefined set.
- For a **Boolean discrete random variable** this predefined set has two members -  $\{0, 1\}$ , {yes, no} etc.
- A **continuous random variable** takes values that are real numbers.



Figures taken from **Computer Vision: models, learning and inference** by Simon Prince.

## Joint Probabilities

- Consider two random variables  $x$  and  $y$ .
- Observe multiple paired instances of  $x$  and  $y$ . Some paired outcomes will occur more frequently.
- This information is encoded in the joint probability distribution  $P(x, y)$ .
- $P(\mathbf{x})$  denotes the joint probability of  $\mathbf{x} = (x_1, \dots, x_K)$ .



← discrete joint pdf

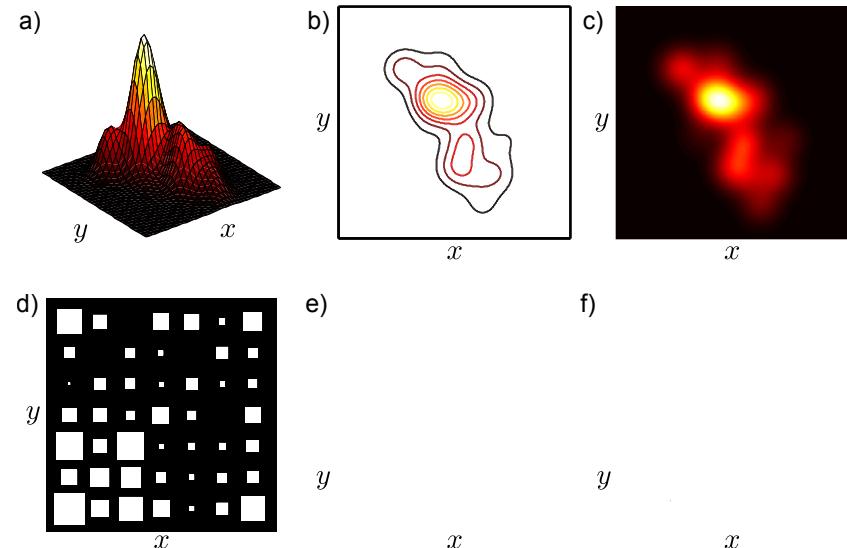
Figure from **Computer Vision: models, learning and inference** by Simon Prince.

## Examples of Random Variables



- Any real number (theoretically infinite)
- Discrete events: either 1, 2, 3, 4, 5, or 6.
- Discrete probability distribution  $p(x) = P(d = x)$
- $P(d = 1) = 1/6$  (fair dice)
- $P(t = 36.6) = 0$
- $P(36.6 < t < 36.7) = 0.1$

## Joint Probabilities (cont.)



## Marginalization

The probability distribution of any single variable can be recovered from a joint distribution by summing for the discrete case

$$P(x) = \sum_y P(x, y)$$

and integrating for the continuous case

$$P(x) = \int_y P(x, y) dy$$

## Marginalization (cont.)

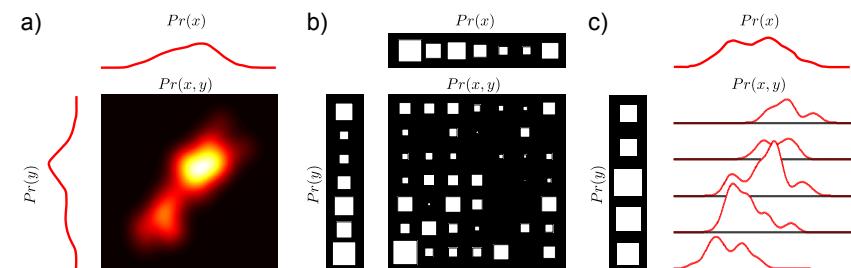


Figure from **Computer Vision: models, learning and inference** by Simon Prince.

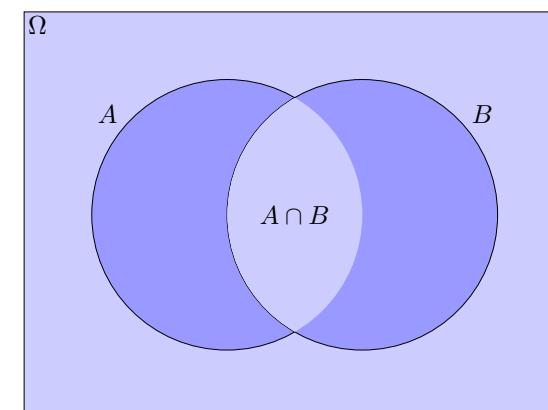
## Conditional Probabilities

$$P(A|B)$$

The probability of event  $A$  when we *know* that event  $B$  has happened

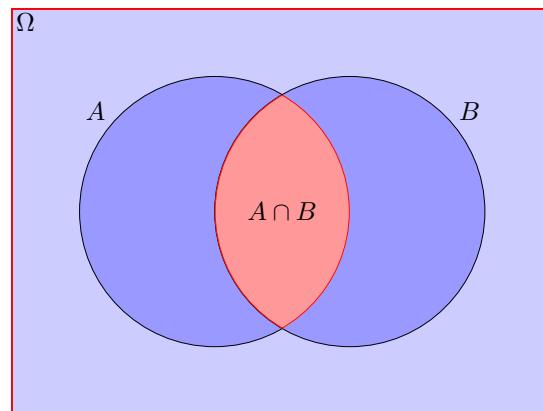
Note: different from the probability that event  $A$  *and* event  $B$  will happen

$$P(A|B) \neq P(A \cap B)$$



## Conditional Probabilities

$$P(A|B) \neq P(A \cap B)$$

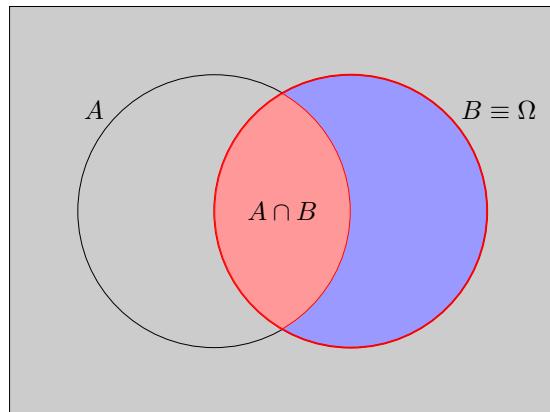


Giampiero Salvi

Lecture 5: Probabilistic Reasoning

## Conditional Probabilities

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

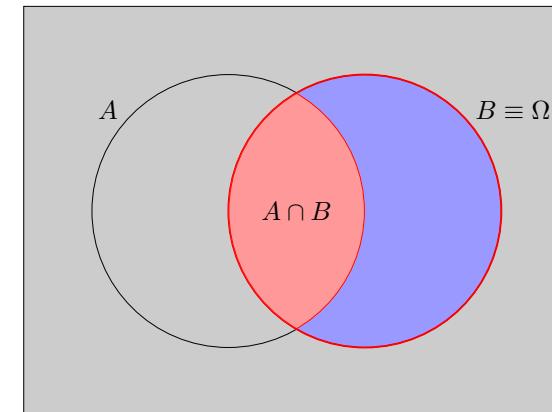


Giampiero Salvi

Lecture 5: Probabilistic Reasoning

## Conditional Probabilities

$$P(A|B) \neq P(A \cap B)$$



Giampiero Salvi

Lecture 5: Probabilistic Reasoning

## Conditional Probability (Random Variables)

- The conditional probability of  $x$  given that  $y$  takes value  $y^*$  indicates the different values of r.v.  $x$  which we'll observe given that  $y$  is fixed to value  $y^*$ .
- The conditional probability can be recovered from the joint distribution  $P(x, y)$ :

$$P(x | y = y^*) = \frac{P(x, y = y^*)}{P(y = y^*)} = \frac{P(x, y = y^*)}{\int_x P(x, y = y^*) dx}$$

- Extract an appropriate slice, and then normalize it.

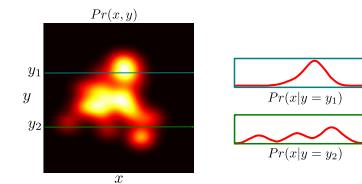


Figure from Computer Vision: models, learning and inference by Simon Prince.

Giampiero Salvi

Lecture 5: Probabilistic Reasoning

## Independence

- two events are independent if the joint distribution can be factorized:  $P(A \cap B) = P(A)P(B)$
- this means that:

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(A)P(B)}{P(B)} = P(A)$$

## Independence

- two events are independent if the joint distribution can be factorized:  $P(A \cap B) = P(A)P(B)$
- this means that:

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(A)P(B)}{P(B)} = P(A)$$

knowing that  $B$  has happened does not tell us anything about  $A$

## Bayes' Rule

if

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

then

$$P(A \cap B) = P(A|B)P(B)$$

## Bayes' Rule

if

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

then

$$P(A \cap B) = P(A|B)P(B) = P(B|A)P(A)$$

## Bayes' Rule

if

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

then

$$P(A \cap B) = P(A|B)P(B) = P(B|A)P(A)$$

and

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

## Bayes' Rule (random variables)

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

Each term in Bayes' rule has a name:

- $P(y|x) \leftarrow \text{Posterior}$  (what we know about  $y$  given  $x$ .)
- $P(y) \leftarrow \text{Prior}$  (what we know about  $y$  before we consider  $x$ .)
- $P(x|y) \leftarrow \text{Likelihood}$  (propensity for observing a certain value of  $x$  given a certain value of  $y$ )
- $P(x) \leftarrow \text{Evidence}$  (a constant to ensure that the l.h.s. is a valid distribution)

## General ML problem (supervised learning)

Data:

$$\{(\mathbf{x}^1, y^1), (\mathbf{x}^2, y^2), \dots, (\mathbf{x}^n, y^n)\}$$

Where  $\mathbf{x}$  are features, and  $y$  is the answer

- if  $y$  is discrete: classification
- if  $y$  is continuous: regression

## General ML problem (supervised learning)

Data:

$$\{(\mathbf{x}^1, y^1), (\mathbf{x}^2, y^2), \dots, (\mathbf{x}^n, y^n)\}$$

Where  $\mathbf{x}$  are features, and  $y$  is the answer

- if  $y$  is discrete: classification
- if  $y$  is continuous: regression
- Learning: we observe several examples of  $\mathbf{x}$  and we know  $y$
- Inference: we want to know  $y$  given a new  $\mathbf{x}$

## Machine Learning with Probabilities

Learning: we observe several examples of  $\mathbf{x}$  and we know  $y$

- we can estimate  $P(y)$  and  $P(\mathbf{x}|y)$

Inference: we want to know  $y$  given a new  $\mathbf{x}$

- we want to estimate  $P(y|\mathbf{x})$

## Machine Learning with Probabilities

Learning: we observe several examples of  $\mathbf{x}$  and we know  $y$

- we can estimate  $P(y)$  and  $P(\mathbf{x}|y)$

Inference: we want to know  $y$  given a new  $\mathbf{x}$

- we want to estimate  $P(y|\mathbf{x})$
- $P(\mathbf{x}|y) \leftarrow \text{Likelihood}$  represents the probability of observing data  $\mathbf{x}$  given the hypothesis  $y$ .
- $P(y) \leftarrow \text{Prior of } y$  represents the background knowledge of hypothesis  $y$  being correct.
- $P(y|\mathbf{x}) \leftarrow \text{Posterior}$  represents the probability that hypothesis  $y$  is true after data  $\mathbf{x}$  has been observed.

## Bayes' Rule

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

With

- $P(\mathbf{x}|y) \leftarrow \text{Likelihood}$  represents the probability of observing data  $\mathbf{x}$  given the hypothesis  $y$ .
- $P(y) \leftarrow \text{Prior of } y$  represents the background knowledge of hypothesis  $y$  being correct.
- $P(y|\mathbf{x}) \leftarrow \text{Posterior}$  represents the probability that hypothesis  $y$  is true after data  $\mathbf{x}$  has been observed.

## Learning and Inference

- **Probabilistic Learning:** The process of learning the likelihood distribution  $P(\mathbf{x}|y)$  and prior probability distribution  $P(y)$  from a set of training points

$$\{(\mathbf{x}^1, y^1), (\mathbf{x}^2, y^2), \dots, (\mathbf{x}^n, y^n)\}$$

- **Probabilistic Inference:** The process of calculating the posterior probability distribution  $P(y|\mathbf{x})$  for certain data  $\mathbf{x}$ .

## Outline

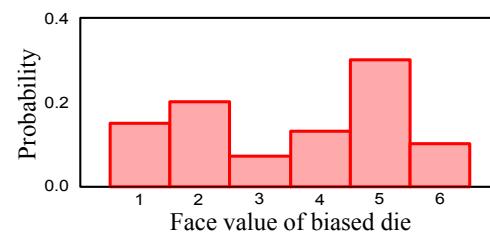
1 Probability Theory Basics

2 Common Distributions

3 Bayesian Decision Theory: Classification

## Categorical

- Domain: discrete variables ( $x \in \{x_1, \dots, x_K\}$ )
- Parameters:  $\lambda = [\lambda_1, \dots, \lambda_K]$
- with  $\lambda_k \in [0, 1]$  and  $\sum_{k=1}^K \lambda_k = 1$

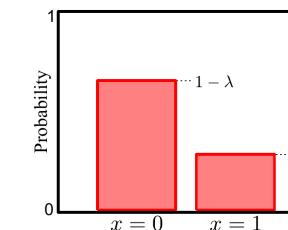


## Bernoulli

- Domain: binary variables ( $x \in \{0, 1\}$ )
- Parameters:  $\lambda = Pr(x = 1)$ ,  $\lambda \in [0, 1]$

Then  $Pr(x = 0) = 1 - \lambda$ , and

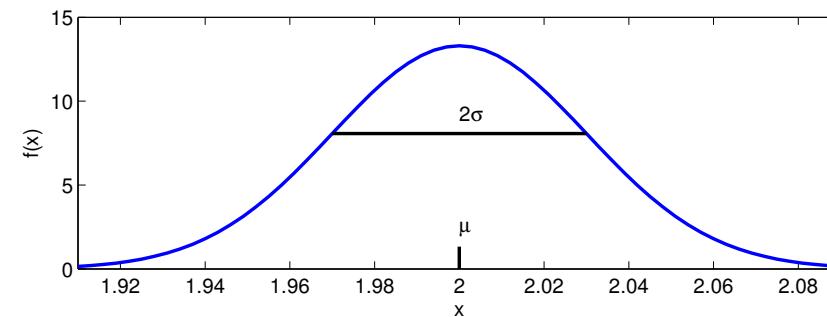
$$Pr(x) = \lambda^x(1 - \lambda)^{1-x} = \begin{cases} \lambda, & \text{if } x = 1, \\ 1 - \lambda, & \text{if } x = 0 \end{cases}$$



## Gaussian distributions: One-dimensional

- aka univariate normal distribution
- Domain: real numbers ( $x \in \mathbb{R}$ )

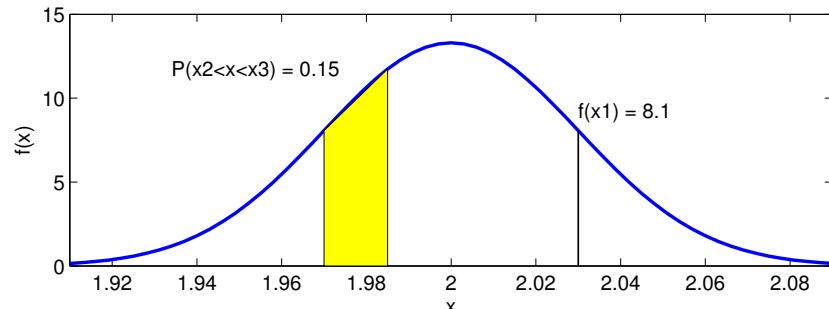
$$f(x|\mu, \sigma^2) = N(\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right]$$



## Gaussian distributions: One-dimensional

- aka univariate normal distribution
- Domain: real numbers ( $x \in \mathbb{R}$ )

$$f(x|\mu, \sigma^2) = N(\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right]$$

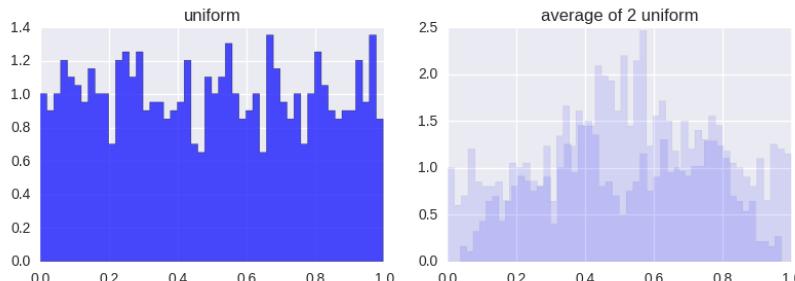


Giampiero Salvi

Lecture 5: Probabilistic Reasoning

## Why Gaussian: Central Limit Theorem

The distribution of the sum (or average) of a large number of independent, identically distributed variables will be approximately normal, regardless of the underlying distribution.<sup>2</sup>



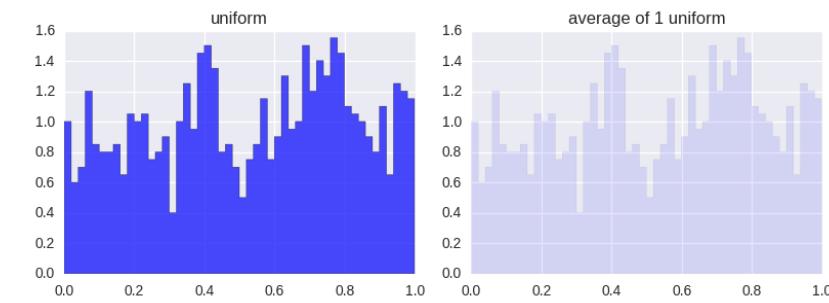
<sup>2</sup>Christopher M Bishop. *Pattern recognition and machine learning*. 2006, p. 78, script inspired by C.E. Ek.

Giampiero Salvi

Lecture 5: Probabilistic Reasoning

## Why Gaussian: Central Limit Theorem

The distribution of the sum (or average) of a large number of independent, identically distributed variables will be approximately normal, regardless of the underlying distribution.<sup>2</sup>.



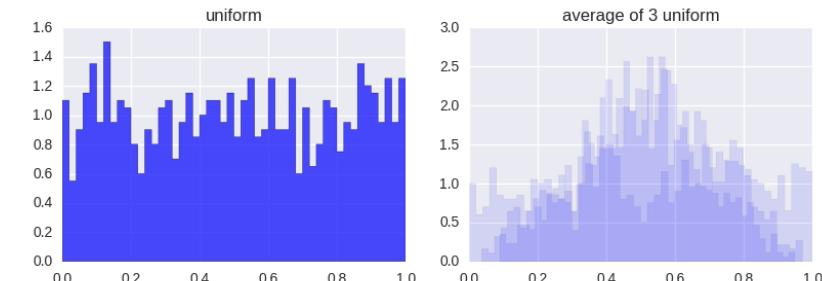
<sup>2</sup>Christopher M Bishop. *Pattern recognition and machine learning*. 2006, p. 78, script inspired by C.E. Ek.

Giampiero Salvi

Lecture 5: Probabilistic Reasoning

## Why Gaussian: Central Limit Theorem

The distribution of the sum (or average) of a large number of independent, identically distributed variables will be approximately normal, regardless of the underlying distribution.<sup>2</sup>.



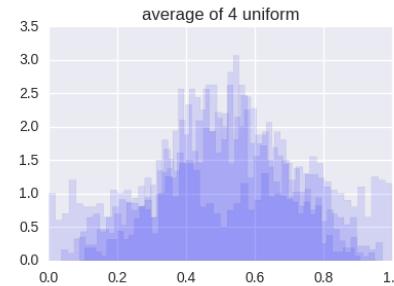
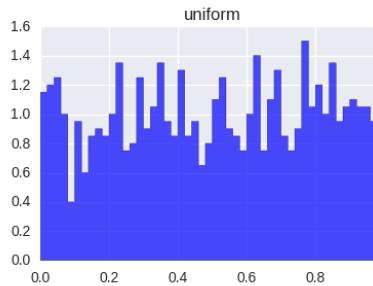
<sup>2</sup>Christopher M Bishop. *Pattern recognition and machine learning*. 2006, p. 78, script inspired by C.E. Ek.

Giampiero Salvi

Lecture 5: Probabilistic Reasoning

## Why Gaussian: Central Limit Theorem

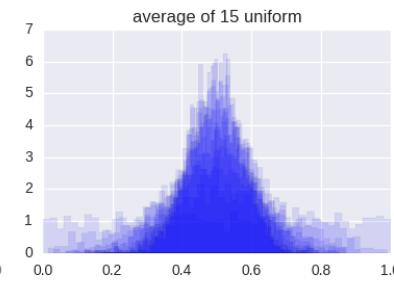
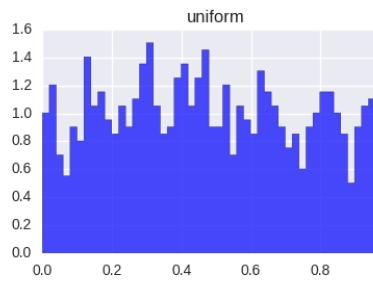
The distribution of the sum (or average) of a large number of independent, identically distributed variables will be approximately normal, regardless of the underlying distribution.<sup>2</sup>.



<sup>2</sup>Christopher M Bishop. *Pattern recognition and machine learning*. 2006, p. 78, script inspired by C.E. Ek.

## Why Gaussian: Central Limit Theorem

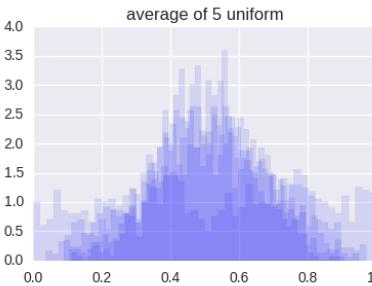
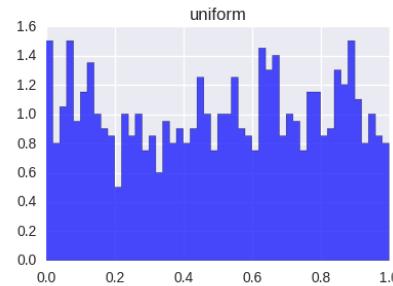
The distribution of the sum (or average) of a large number of independent, identically distributed variables will be approximately normal, regardless of the underlying distribution.<sup>2</sup>.



<sup>2</sup>Christopher M Bishop. *Pattern recognition and machine learning*. 2006, p. 78, script inspired by C.E. Ek.

## Why Gaussian: Central Limit Theorem

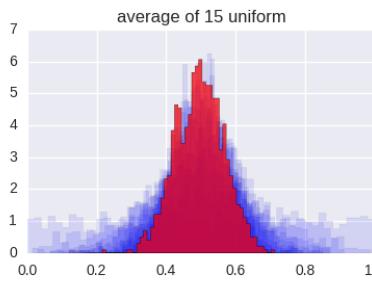
The distribution of the sum (or average) of a large number of independent, identically distributed variables will be approximately normal, regardless of the underlying distribution.<sup>2</sup>.



<sup>2</sup>Christopher M Bishop. *Pattern recognition and machine learning*. 2006, p. 78, script inspired by C.E. Ek.

## Why Gaussian: Central Limit Theorem

The distribution of the sum (or average) of a large number of independent, identically distributed variables will be approximately normal, regardless of the underlying distribution.<sup>2</sup>.



<sup>2</sup>Christopher M Bishop. *Pattern recognition and machine learning*. 2006, p. 78, script inspired by C.E. Ek.

## Gaussian distributions: $D$ Dimensions

- aka multivariate normal distribution
- Domain: real numbers ( $\mathbf{x} \in \mathbb{R}^D$ )

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_D \end{bmatrix} \quad \mu = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \dots \\ \mu_D \end{bmatrix} \quad \Sigma = \begin{bmatrix} \sigma_{11}^2 & \sigma_{12} & \dots & \sigma_{1D} \\ \sigma_{21} & \dots & & \\ \dots & & & \\ \sigma_{D1} & \dots & & \sigma_{DD}^2 \end{bmatrix}$$

$$f(\mathbf{x}|\mu, \Sigma) = \frac{\exp\left[-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)\right]}{(2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}}}$$

## Gaussian distributions

$$f(\mathbf{x}|\mu, \Sigma) = \frac{\exp\left[-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)\right]}{(2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}}}$$

Eigenvalue decomposition of the covariance matrix:

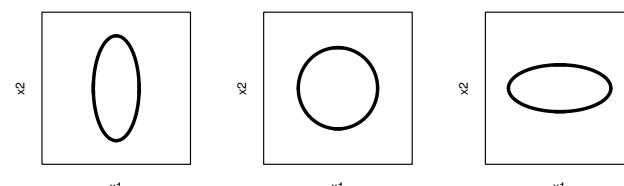
$$\Sigma = \lambda R \Sigma_{\text{diag}} R^T$$

## Gaussian distributions

$$f(\mathbf{x}|\mu, \Sigma) = \frac{\exp\left[-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)\right]}{(2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}}}$$

Eigenvalue decomposition of the covariance matrix:

$$\Sigma = \lambda R \Sigma_{\text{diag}} R^T$$

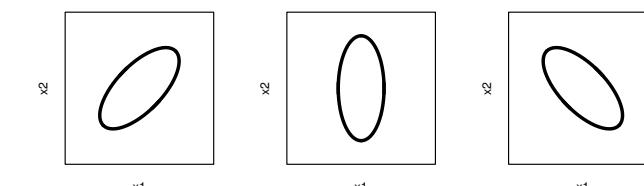


## Gaussian distributions

$$f(\mathbf{x}|\mu, \Sigma) = \frac{\exp\left[-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)\right]}{(2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}}}$$

Eigenvalue decomposition of the covariance matrix:

$$\Sigma = \lambda R \Sigma_{\text{diag}} R^T$$

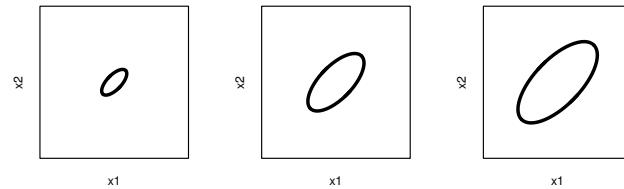


## Gaussian distributions

$$f(\mathbf{x}|\mu, \Sigma) = \frac{\exp\left[-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)\right]}{(2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}}}$$

Eigenvalue decomposition of the covariance matrix:

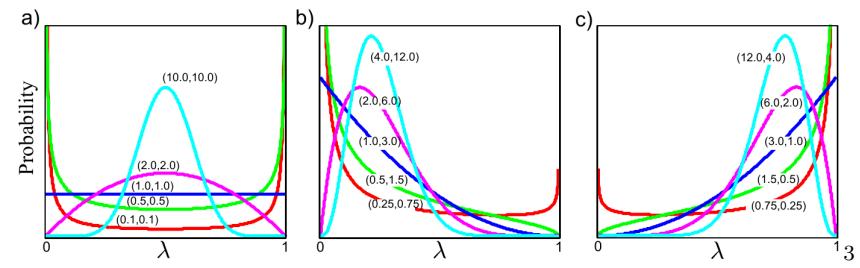
$$\Sigma = \lambda R \Sigma_{\text{diag}} R^T$$



## Beta and Dirichlet (PDF over Probabilities)

Beta

- Domain: real numbers, bounded ( $\lambda \in [0, 1]$ )
- Parameters:  $\alpha, \beta \in \mathbb{R}_+$
- describes probability of parameter  $\lambda$  in Bernoulli



<sup>3</sup>Figure from Computer Vision: models, learning and inference by Simon Prince.

## Beta and Dirichlet (PDF over Probabilities)

Beta

- Domain: real numbers, bounded ( $\lambda \in [0, 1]$ )
- Parameters:  $\alpha, \beta \in \mathbb{R}_+$
- describes probability of parameter  $\lambda$  in Bernoulli

Dirichlet

- Domain:  $K$  real numbers, bounded ( $\lambda_1, \dots, \lambda_K \in [0, 1]$ )
- Parameters:  $\alpha_1, \dots, \alpha_K \in \mathbb{R}_+$
- describes probability of parameters  $\lambda_k$  in Categorical

## Expected value

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

- Shows the “center of gravity” of a distribution
- Sampled expected value (mean)

$$\bar{\mathbf{x}} = \frac{1}{N} \sum_i \mathbf{x}_i$$

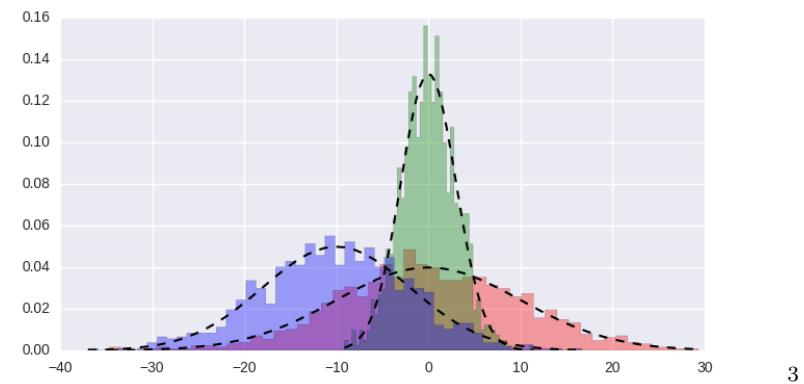
## Variance

$$\sigma^2(\mathbf{x}) = \text{var}[\mathbf{x}] = \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])^2]$$

- Shows the “spread” of a distribution
- Sample variance

$$\overline{\sigma^2(\mathbf{x})} = \frac{1}{N-1} \sum_i^N (\mathbf{x}_i - \mu(\mathbf{x}_i))^2$$

## Examples



<sup>3</sup>Script by C.E. Ek

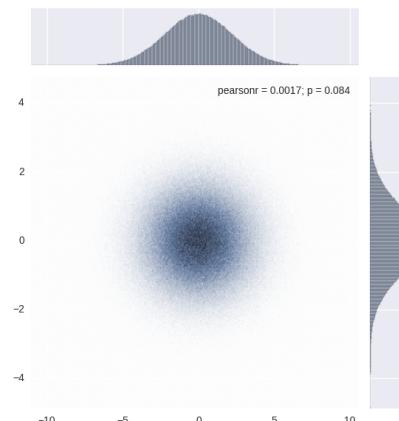
## Covariance

$$\sigma(\mathbf{x}, \mathbf{y}) = \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{y} - \mathbb{E}[\mathbf{y}])]$$

- Shows how the “spread” of how two variables vary *together*
- Sample co-variance

$$\overline{\sigma(\mathbf{x}, \mathbf{y})} = \frac{1}{N-1} \sum_i^N (\mathbf{x}_i - \mu(\mathbf{x}_i))(\mathbf{y}_i - \mu(\mathbf{y}))$$

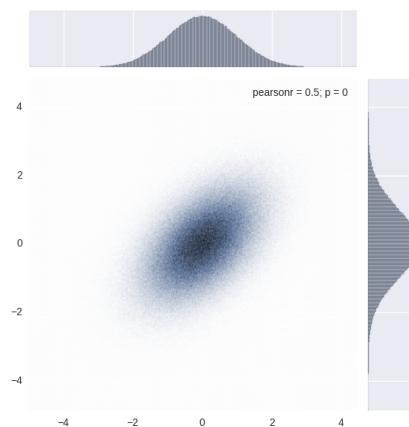
## Examples



<sup>4</sup>

Script by C.E. Ek

## Examples



4

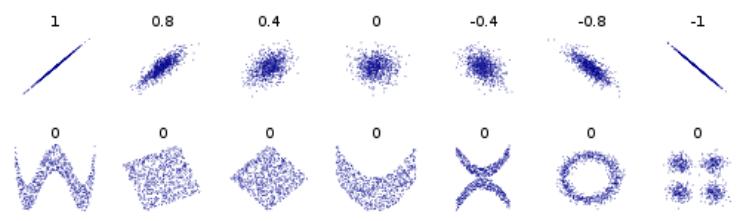
<sup>4</sup>Script by C.E. Ek

Giampiero Salvi

Lecture 5: Probabilistic Reasoning

## Covariance and Independence

- covariance is “linear” dependency
- dependent variables may have zero covariance
- in some distributions zero covariance is equivalent to independence

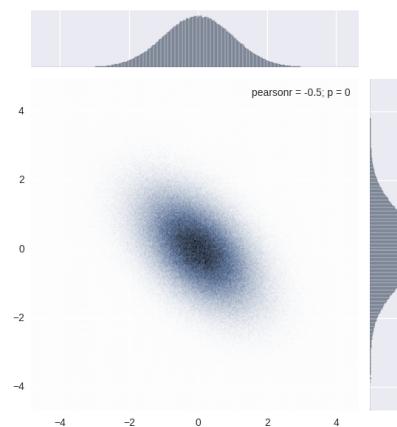


<sup>5</sup>Figure adapted from Wikipedia

Giampiero Salvi

Lecture 5: Probabilistic Reasoning

## Examples



4

<sup>4</sup>Script by C.E. Ek

Giampiero Salvi

Lecture 5: Probabilistic Reasoning

## Covariance and Independence (Gaussian)

- covariance is “linear” dependency
- dependent variables may have zero covariance
- in Gaussian (and few other distribution) zero covariance is equivalent to independence

$$f(\mathbf{x}|\mu, \Sigma) = \frac{\exp\left[-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)\right]}{(2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}}}$$

## Outline

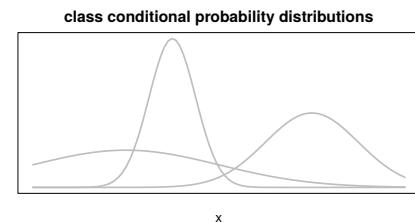
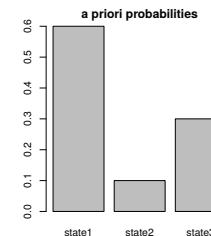
### 1 Probability Theory Basics

### 2 Common Distributions

### 3 Bayesian Decision Theory: Classification

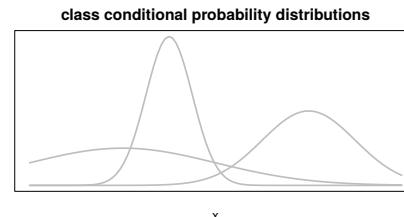
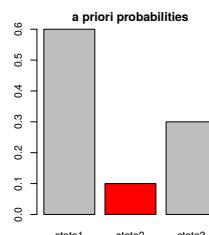
## The Probabilistic Model of Classification

- one of  $c$  states  $y_j$  is selected with *a priori* probability  $P(y_j)$
- When in state  $y_j$ , some observations  $\hat{\mathbf{x}}$  are generated with distribution  $p(\mathbf{x}|y_j)$



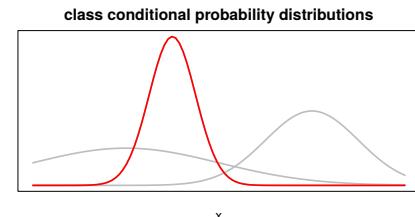
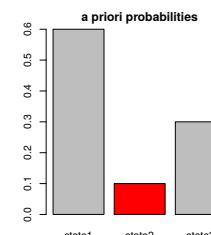
## The Probabilistic Model of Classification

- one of  $c$  states  $y_j$  is selected with *a priori* probability  $P(y_j)$
- When in state  $y_j$ , some observations  $\hat{\mathbf{x}}$  are generated with distribution  $p(\mathbf{x}|y_j)$



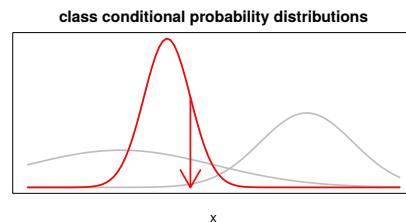
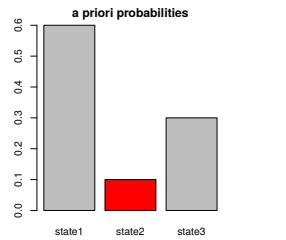
## The Probabilistic Model of Classification

- one of  $c$  states  $y_j$  is selected with *a priori* probability  $P(y_j)$
- When in state  $y_j$ , some observations  $\hat{\mathbf{x}}$  are generated with distribution  $p(\mathbf{x}|y_j)$



## The Probabilistic Model of Classification

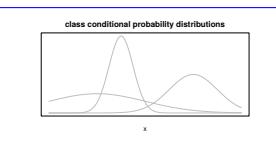
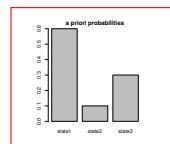
- one of  $c$  states  $y_j$  is selected with *a priori* probability  $P(y_j)$
- When in state  $y_j$ , some observations  $\hat{\mathbf{x}}$  are generated with distribution  $p(\mathbf{x}|y_j)$



## Problem

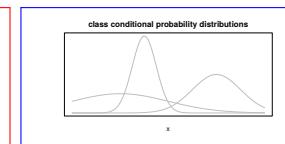
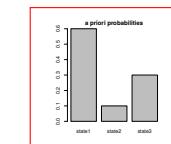
- If I observe  $\hat{\mathbf{x}}$  and I know  $P(y_j)$  and  $p(\mathbf{x}|y_j)$  for each  $j$
- what can I say about the state of the problem  $y_j$ ?

## Bayes decision theory



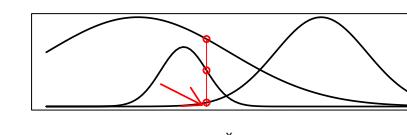
$$P(y_j|\mathbf{x}) = \frac{p(\mathbf{x}|y_j) P(y_j)}{p(\mathbf{x})}$$

## Bayes decision theory

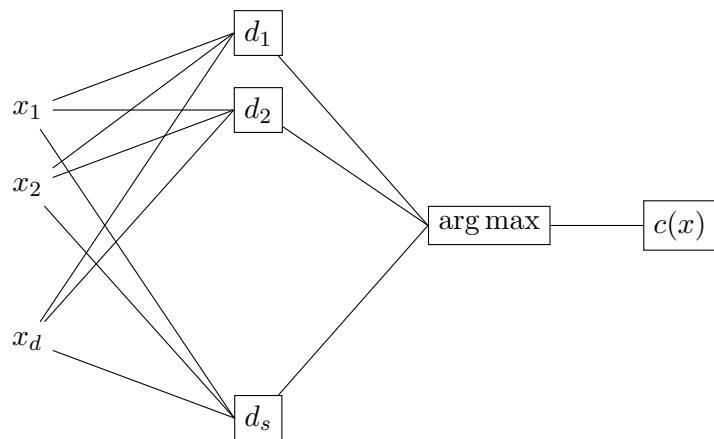


$$P(y_j|\mathbf{x}) = \frac{p(\mathbf{x}|y_j) P(y_j)}{p(\mathbf{x})}$$

posterior probabilities



## Classifiers: Discriminant Functions



$$d_i(\mathbf{x}) = p(\mathbf{x}|y_i) P(y_i)$$

## Example: Which Gender?

**Task:** Determine the gender of a person given their measured hair length.

### Notation:

- Let  $g \in \{\text{'f'}, \text{'m'}\}$  be a r.v. denoting the gender of a person.
- Let  $x$  be the measured length of the hair.

## Example: Which Gender?

**Task:** Determine the gender of a person given their measured hair length.

## Example: Which Gender?

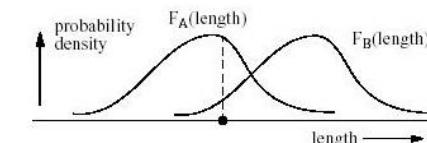
**Task:** Determine the gender of a person given their measured hair length.

### Notation:

- Let  $g \in \{\text{'f'}, \text{'m'}\}$  be a r.v. denoting the gender of a person.
- Let  $x$  be the measured length of the hair.

### Information given:

- The hair length observation was made at a boy's school thus  $P(g = \text{'m'}) = .95, P(g = \text{'f'}) = .05$
- Knowledge of the likelihood distributions  $P(x | g = \text{'f'})$  and  $P(x | g = \text{'m'})$



## Example: Which Gender?

**Task:** Determine the gender of a person given their measured hair length  $\Rightarrow$  calculate  $P(g | x)$ .

**Solution:**

Apply Bayes' Rule to get

$$\begin{aligned} P(g = 'm' | x) &= \frac{P(x | g = 'm') P(g = 'm')}{P(x)} \\ &= \frac{P(x | g = 'm') P(g = 'm')}{P(x | g = 'f') P(g = 'f') + P(x | g = 'm') P(g = 'm')} \end{aligned}$$

Can calculate  $P(g = 'f' | x) = 1 - P(g = 'm' | x)$

## Selecting the most probably hypothesis

- **Maximum A Posteriori (MAP) Estimate:**

Hypothesis with highest probability given observed data

$$\begin{aligned} y_{\text{MAP}} &= \arg \max_{y \in \mathcal{Y}} P(y | \mathbf{x}) \\ &= \arg \max_{y \in \mathcal{Y}} \frac{P(\mathbf{x} | y) P(y)}{P(\mathbf{x})} \\ &= \arg \max_{y \in \mathcal{Y}} P(\mathbf{x} | y) P(y) \end{aligned}$$

Giampiero Salvi      Lecture 5: Probabilistic Reasoning

Probability Theory Basics  
Common Distributions  
Bayesian Decision Theory: Classification

## Selecting the most probably hypothesis

Giampiero Salvi      Lecture 5: Probabilistic Reasoning

Probability Theory Basics  
Common Distributions  
Bayesian Decision Theory: Classification

## Example: Cancer or Not?

- **Maximum A Posteriori (MAP) Estimate:**

Hypothesis with highest probability given observed data

$$\begin{aligned} y_{\text{MAP}} &= \arg \max_{y \in \mathcal{Y}} P(y | \mathbf{x}) \\ &= \arg \max_{y \in \mathcal{Y}} \frac{P(\mathbf{x} | y) P(y)}{P(\mathbf{x})} \\ &= \arg \max_{y \in \mathcal{Y}} P(\mathbf{x} | y) P(y) \end{aligned}$$

- **Maximum Likelihood Estimate (MLE):**

Hypothesis with highest likelihood of generating observed data.

$$y_{\text{MLE}} = \arg \max_{y \in \mathcal{Y}} P(\mathbf{x} | y)$$

Useful if we do not know prior distribution or if it is uniform.

**Scenario:**

A patient takes a lab test and the result comes back positive. The test returns a correct positive result in only 98% of the cases in which the disease is actually present, and a correct negative result in only 97% of the cases in which the disease is not present. Furthermore, 0.8% of the entire population have cancer.

## Example: Cancer or Not?

### Scenario:

A patient takes a lab test and the result comes back positive. The test returns a correct positive result in only 98% of the cases in which the disease is actually present, and a correct negative result in only 97% of the cases in which the disease is not present. Furthermore, 0.8% of the entire population have cancer.

### Scenario in probabilities:

- Priors:

$$P(\text{disease}) = .008 \quad P(\text{not disease}) = .992$$

- Likelihoods:

$$\begin{aligned} P(+ | \text{disease}) &= .98 & P(+ | \text{not disease}) &= .03 \\ P(- | \text{disease}) &= .02 & P(- | \text{not disease}) &= .97 \end{aligned}$$

## Example: Cancer or Not?

### Find MAP estimate:

When test returned a positive result,

$$\begin{aligned} y_{\text{MAP}} &= \arg \max_{y \in \{\text{disease, not disease}\}} P(y | +) \\ &= \arg \max_{y \in \{\text{disease, not disease}\}} P(+ | y) P(y) \end{aligned}$$

Substituting in the correct values get

$$\begin{aligned} P(+ | \text{disease}) P(\text{disease}) &= .98 \times .008 = .0078 \\ P(+ | \text{not disease}) P(\text{not disease}) &= .03 \times .992 = .0298 \end{aligned}$$

Therefore  $y_{\text{MAP}} = \text{"not disease"}$ .

## Example: Cancer or Not?

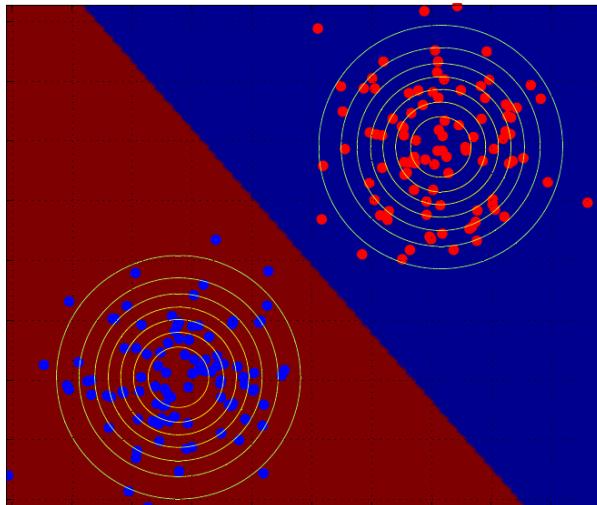
### Find MAP estimate:

When test returned a positive result,

$$\begin{aligned} y_{\text{MAP}} &= \arg \max_{y \in \{\text{disease, not disease}\}} P(y | +) \\ &= \arg \max_{y \in \{\text{disease, not disease}\}} P(+ | y) P(y) \end{aligned}$$

## Demo: Gaussian Decision Boundaries

<https://github.com/giampierosalvi/GaussianDecisionBoundaries>

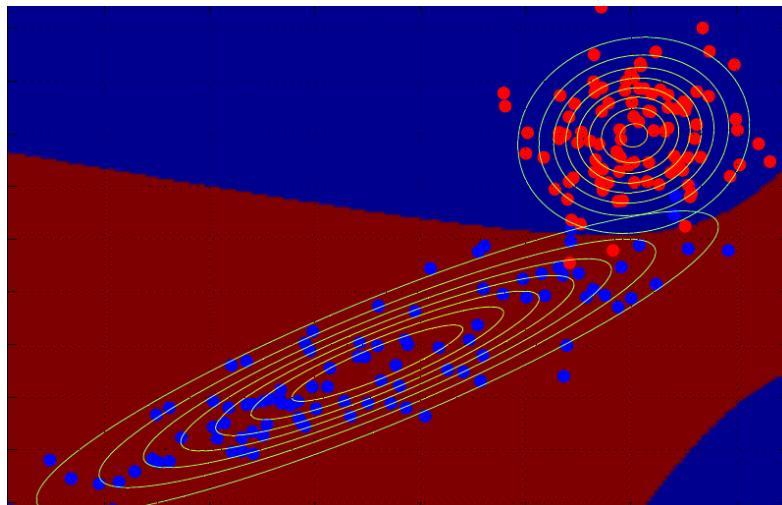


Giampiero Salvi

Lecture 5: Probabilistic Reasoning

## Demo: Gaussian Decision Boundaries

<https://github.com/giampierosalvi/GaussianDecisionBoundaries>

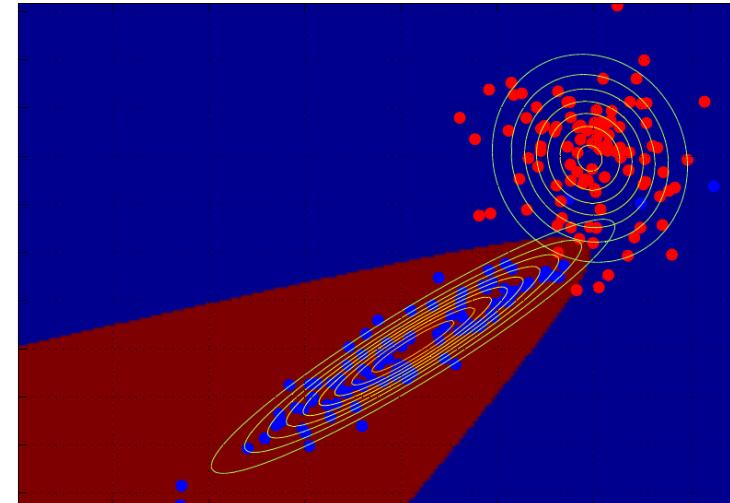


Giampiero Salvi

Lecture 5: Probabilistic Reasoning

## Demo: Gaussian Decision Boundaries

<https://github.com/giampierosalvi/GaussianDecisionBoundaries>

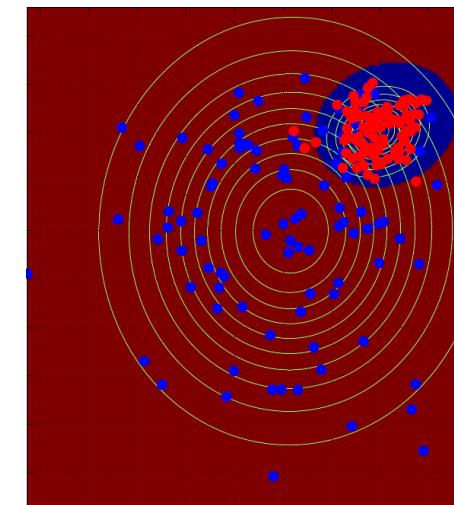


Giampiero Salvi

Lecture 5: Probabilistic Reasoning

## Demo: Gaussian Decision Boundaries

<https://github.com/giampierosalvi/GaussianDecisionBoundaries>



Giampiero Salvi

Lecture 5: Probabilistic Reasoning

## Discriminative vs Generative Models

## Summary

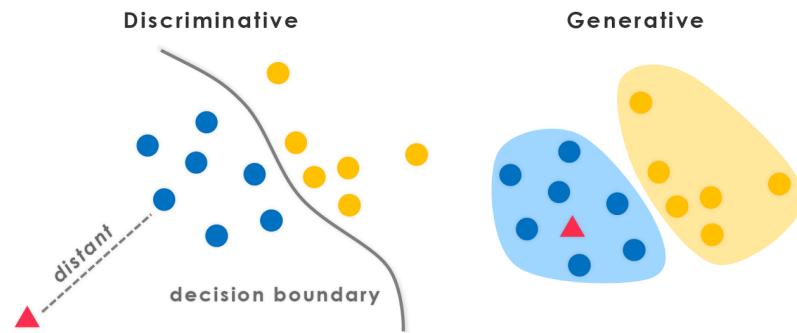


Figure from Nguyen *et al.* 2015. <http://www.evolvingai.org/fooling>

Today:

- 1 Probability Theory Basics
- 2 Common Distributions
- 3 Bayesian Decision Theory: Classification

Next time: How to fit probability models to data