

# **Artificial Intelligence and Machine Learning**

**Barbara Caputo**

## Useful info (1)

Teacher: Barbara Caputo,  
[https://scholar.google.it/citations  
?hl=en&user=mHbdIAwAAAAJ&view\\_op=list\\_works](https://scholar.google.it/citations?hl=en&user=mHbdIAwAAAAJ&view_op=list_works)

Assistant: Paolo Russo

Where/how to find us:

Email: [barbara.caputo@iit.it](mailto:barbara.caputo@iit.it)  
paolo.russo@iit.it

Office: 4E-45 (BC), Lab 7 –R3BC31 (PR)

Q/A time: after the lectures (BC)/by appointment  
You can come any time at your own risk!

## Useful info (1)

Textbooks/Useful books:

K. P. Murphy. Machine Learning: a probabilistic perspective. MIT Press.

B. Schoelkopf, A. Smola. Learning with Kernels. MIT Press.

W. Feller. An introduction to probability theory and its applications. Vol 1. Wiley Ed.

I. Goodfellow, Y. Bengio, A. Courville. Deep Learning. MIT Press.

## Useful info (2)

Exam modality:

Homeworks (lab experiences) +written exam +oral exam

## Useful info (2)

Exam modality:

- Homeworks during the lectures: 3, on three key topics (PCA, svm, cnn???e-m???). Reports to be submitted by end of the lectures (19/01/2019, 23:59 CET). Every experience graded (2 points), if overall points <4 no admission to written exam.

## Useful info (2)

Exam modality:

- Homeworks during the lectures: 3, on three key topics (PCA, svm, cnn???e-m???). Reports to be submitted by end of the lectures (19/01/2019, 23:59 CET). Every experience graded (2 points), if overall points  $<4$  no admission to written exam.
- there will be a cut-off date for submitting the reports for every exam session

## Useful info (2)

### Exam modality:

- Homeworks during the lectures: 3, on three key topics (PCA, svm, cnn???e-m???). Reports to be submitted by end of the lectures (19/01/2019, 23:59 CET). Every experience graded (2 points), if overall points <4 no admission to written exam.
- Written exam: exercises as done during lectures. If grade <18, no admission to oral exam.

## Useful info (2)

### Exam modality:

- Homeworks during the lectures: 3, on three key topics (PCA, svm, cnn??). Reports to be submitted by end of the lectures (19/01/2019, 23:59 CET). Every experience graded (2 points), if overall points  $<4$  no admission to written exam.
- Written exam: exercises as done during lectures. If grade  $<18$ , no admission to oral exam.
- Oral exam: pen and pencil, three questions. You decide the first, the second is decided the day before the exam, the third I decide.



## Useful info (2)

### Exam modality:

- Homeworks during the lectures: 3, on three key topics (PCA, svm, cnn??). Reports to be submitted by end of the lectures (19/01/2019, 23:59 CET). Every experience graded (2 points), if overall points  $<4$  no admission to written exam.
- Written exam: exercises as done during lectures. If grade  $<18$ , no admission to oral exam.
- Oral exam: pen and pencil, three questions. You decide the first, the second is decided the day before the exam, the third I decide.

## Useful info (2)

### Exam modality:

- Written exam: exercises as done during lectures. If grade  $< 18$ , no admission to oral exam.
- Oral exam: pen and pencil, three questions. You decide the first, the second is decided the day before the exam, the third I decide. Final grade average between written and oral exam.

# Artificial Intelligence and Machine Learning

AI: theory and algorithms that enable computers to mimic human intelligence

ML: a subset of AI that includes statistical techniques enabling machines to improve at tasks with experience.

**This course focuses on ML only, and specifically on a specific family of ML algorithms (discriminative methods)**

# Outline

## Theory:

- Probabilities:
  - Probability measures, events, random variables, conditional probabilities, dependence, expectations, etc
- Bayes rule
- Parameter estimation:
  - Maximum Likelihood Estimation (MLE)
  - Maximum a Posteriori (MAP)

# Sample space

**Def:** A **sample space**  $\Omega$  is the set of all possible outcomes of a (conceptual or physical) random experiment. ( $\Omega$  can be finite or infinite.)

## Examples:

–  $\Omega$  may be the set of all possible outcomes of a dice roll (1,2,3,4,5,6)

-Pages of a book opened randomly. (1-157)

-Real numbers for temperature, location, time, etc



# Events

We will ask the question:

**What is the probability of a particular event?**

**Def:** *Event A is a subset of the sample space  $\Omega$*

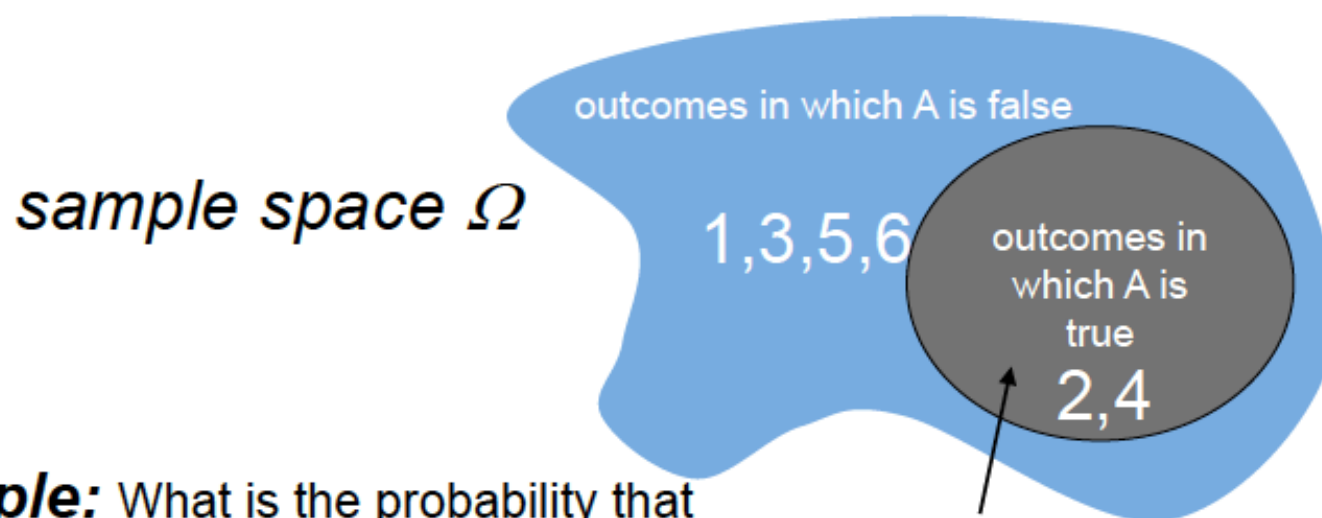
***Examples:***

What is the probability of

- *the book is open at an odd number*
- *rolling a dice the number  $< 4$*
- *a random person's height  $X : a < X < b$*

# Probability

**Def:** Probability  $P(A)$ , the probability that event (subset)  $A$  happens, is a function that maps the event  $A$  onto the interval  $[0, 1]$ .  $P(A)$  is also called the **probability measure** of  $A$ .



**Example:** What is the probability that the number on the dice is 2 or 4?

$P(A)$  is the volume of the area.

# Kolmogorov Axioms

(i) Nonnegativity:  $P(A) \geq 0$  for each  $A$  event.

(ii)  $P(\Omega) = 1$ .

(iii)  $\sigma$ -additivity: For disjoint sets (events)  $A_i$ , we have

$$P(\cup_{i=1}^{\infty} A_i) = \sum_{i=1}^{\infty} P(A_i)$$

## Consequences:

$$P(\emptyset) = 0.$$

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

$$P(A^c) = 1 - P(A).$$



# Random Variables

**Def:** Real valued **random variable** is a function of the outcome of a randomized experiment

$$X : \Omega \rightarrow \mathbb{R}$$

$$P(a < X < b) \doteq P(\omega : a < X(\omega) < b)$$

$$P(X = a) \doteq P(\omega : X(\omega) = a)$$

## Examples:

- **Discrete random variable examples ( $\Omega$  is discrete):**
- $X(\omega) = \text{True}$  if a randomly drawn person ( $\omega$ ) from our class ( $\Omega$ ) is female
- $X(\omega) = \text{The hometown } X(\omega) \text{ of a randomly drawn person } (\omega) \text{ from our class } (\Omega)$

# Discrete Distributions

- Bernoulli distribution:  $\text{Ber}(p)$

$\Omega = \{\text{head}, \text{tail}\}$   $X(\text{head}) = 1$ ,  $X(\text{tail}) = 0$ .

$$P(X = a) = P(\omega : X(\omega) = a) = \begin{cases} p, & \text{for } a = 1 \\ 1 - p, & \text{for } a = 0 \end{cases}$$



- Binomial distribution:  $\text{Bin}(n, p)$

Suppose a coin with head prob.  $p$  is tossed  $n$  times. What is the probability of getting  $k$  heads and  $n-k$  tails?

$\Omega = \{\text{possible } n \text{ long head/tail series}\}$ ,  $|\Omega| = 2^n$

$K(\omega) = \text{number of heads in } \omega = (\omega_1, \dots, \omega_n) \in \{\text{head}, \text{tail}\}^n = \Omega$

$$P(K = k) = P(\omega : K(\omega) = k) = \sum_{\omega: K(\omega)=k} p^k (1-p)^{n-k} = \binom{n}{k} p^k (1-p)^{n-k}$$

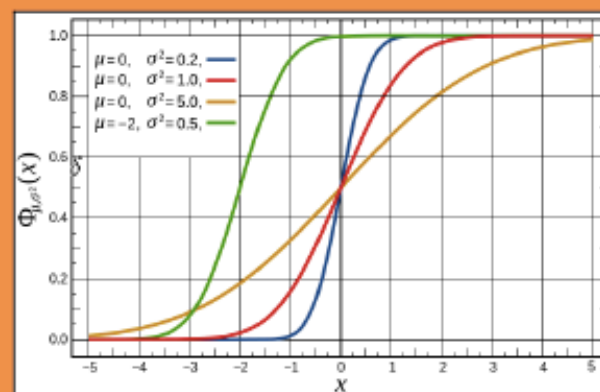
# Continuous Distribution

**Def:** continuous probability distribution: its cumulative distribution function is absolutely continuous.

**Def:** cumulative distribution function

USA:  $F_X(z) = P(X \leq z)$

Hungary:  $F_X(z) = P(X < z)$



**Def:** Let  $F(-\infty) = 0$ .  $F : (-\infty, \infty) \rightarrow \mathbb{R}$  is **absolutely continuous**

$$F(x) = \int_{-\infty}^x f(t) dt \text{ for some function } f.$$

**Def:**  $f$  is called the density of the distribution.

**Properties :**  $\frac{d}{dx}F(x) = f(x)$   $F(x) = \int_{-\infty}^x f(t) dt$

# Probability Density Function (pdf)

## Pdf properties:

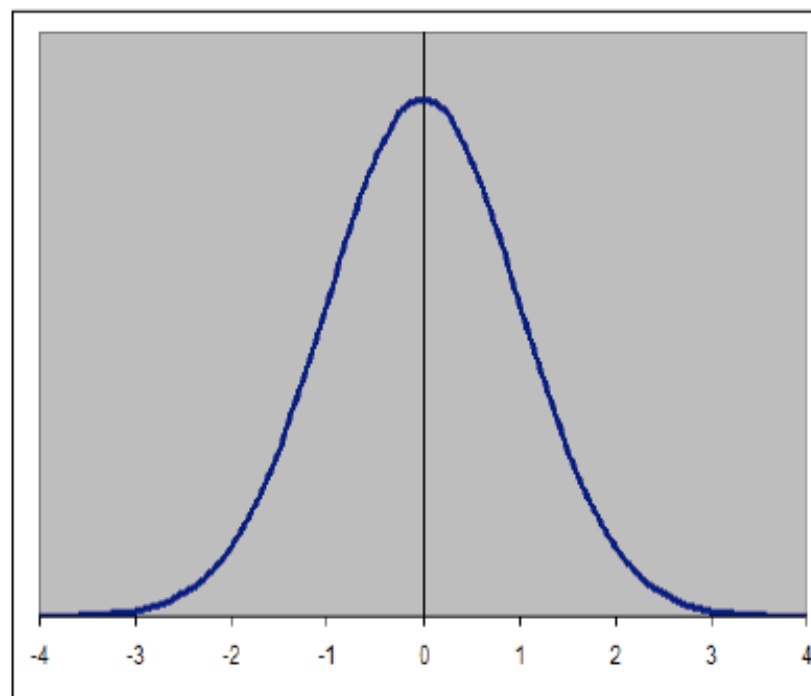
$$f(x) \geq 0$$

$$\int_{-\infty}^{\infty} f(x) dx = 1$$

$$f(x) = \frac{d}{dx} F(x)$$

$$F(x) = \int_{-\infty}^x f(t) dt$$

$$P(a \leq X \leq b) = \int_a^b f(x) dx$$



Intuitively, one can think of  $f(x)dx$  as being the probability of  $X$  falling within the infinitesimal interval  $[x, x + dx]$ .  $P(x < X < x + dx) = f(x)dx$

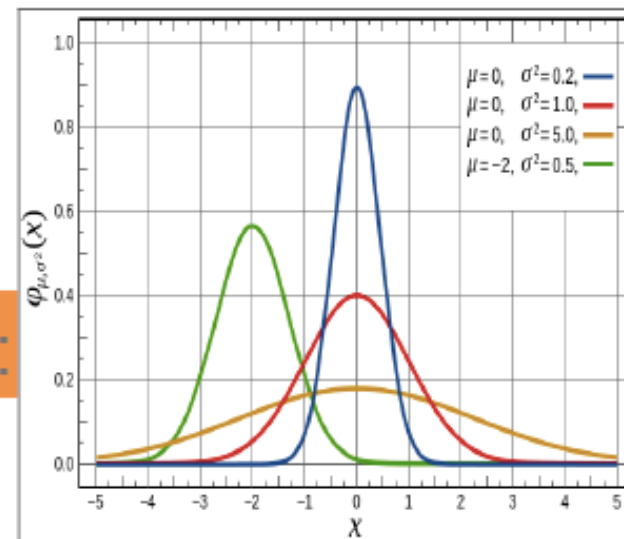
# Moments

**Expectation: average value, mean, 1<sup>st</sup> moment:**

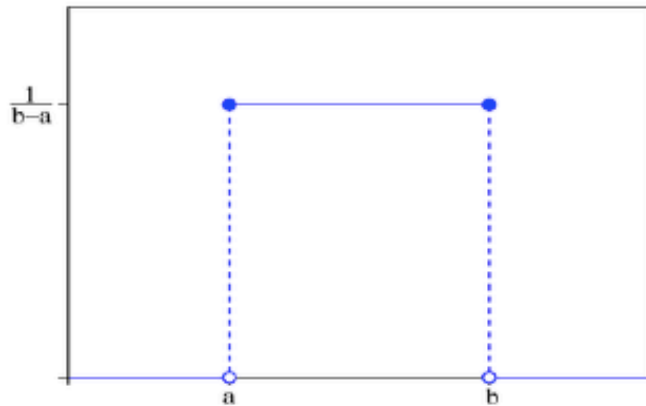
$$E(X) = \begin{cases} \sum_{i \in \Omega} x_i p(x_i) & \text{discrete} \\ \int_{-\infty}^{\infty} x p(x) dx & \text{continuous} \end{cases}$$

**Variance: the spread, 2<sup>nd</sup> moment:**

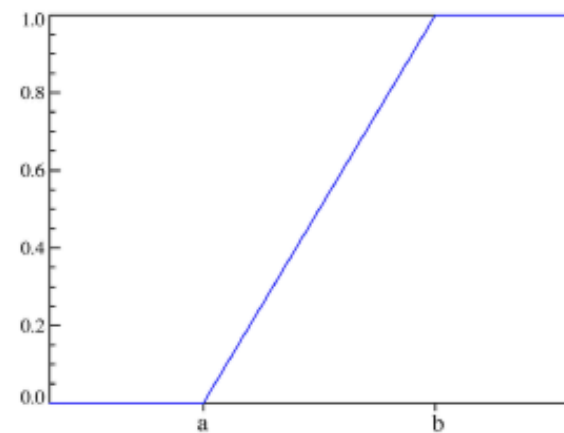
$$E(X) = \begin{cases} \sum_{i \in \Omega} [x_i - E(X)]^2 p(x_i) & \text{discrete} \\ \int_{-\infty}^{\infty} (x - E(x))^2 p(x) dx & \text{continuous} \end{cases}$$



# Uniform Distribution



PDF

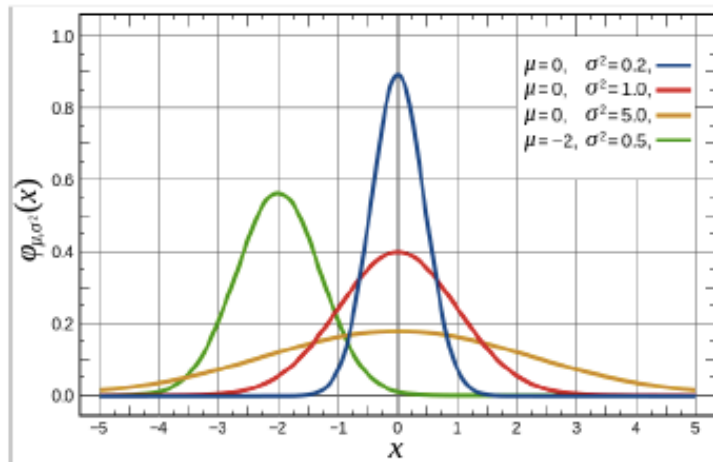


CDF

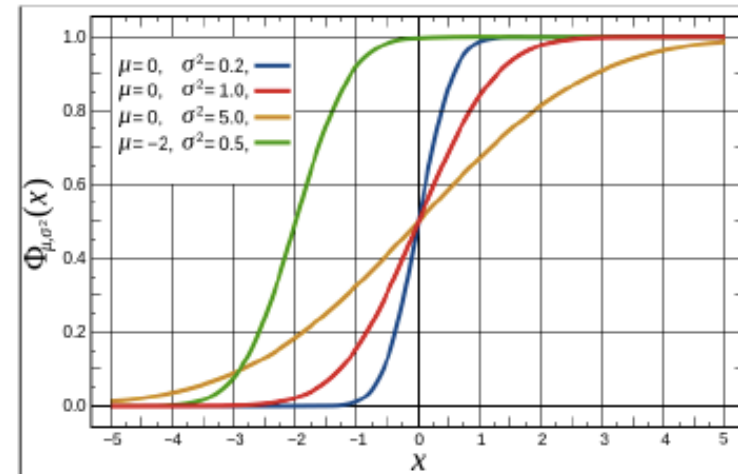
$$f(x) = \begin{cases} \frac{1}{b-a} & a \leq x \leq b \\ 0 & \text{Otherwise} \end{cases}$$

$$F(x) = \begin{cases} 0 & x \leq a \\ \frac{x-a}{b-a} & a < x \leq b \\ 1 & b < x \end{cases}$$

# Normal (Gaussian) Distribution



PDF



CDF

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

$$F(x) = \frac{1}{2} \left[ 1 + \operatorname{erf}\left(\frac{x-\mu}{\sqrt{2\sigma^2}}\right) \right]$$

# Multivariate (Joint) Distribution



# Multivariate (Joint) Distribution

We can generalize the above ideas from 1-dimension to any finite dimensions.

$$P(a \leq X \leq b, c \leq Y \leq d) = ?$$

$$P(a_1 \leq X_1 \leq b_1, \dots, a_d \leq X_d \leq b_d) = ?$$

---

# Multivariate (Joint) Distribution

We can generalize the above ideas from 1-dimension to any finite dimensions.

$$P(a \leq X \leq b, c \leq Y \leq d) = ?$$

$$P(a_1 \leq X_1 \leq b_1, \dots, a_d \leq X_d \leq b_d) = ?$$

---

# Multivariate (Joint) Distribution

We can generalize the above ideas from 1-dimension to any finite dimensions.

$$P(a \leq X \leq b, c \leq Y \leq d) = ?$$

$$P(a_1 \leq X_1 \leq b_1, \dots, a_d \leq X_d \leq b_d) = ?$$

Discrete distribution:

$$P(\text{headache} \wedge \text{no flu}) = 7/80$$

$$P(\text{headache}) = 7/80 + 1/80$$

$$P(X = \text{headache}, Y = \text{flu}) = 1/80$$

	Flu	No Flu
Headache	1/80	7/80
No Headache	1/80	71/80

# Multivariate Gaussian distribution

# Multivariate Gaussian distribution

For  $A \subset \mathbb{R}^d$ ,  $P([X_1, \dots, X_d] \in A) = \int_A f(x_1, \dots, x_d) dx_1 \cdots dx_d$

$$F_X(z_1, \dots, z_d) = \int_{-\infty}^{z_1} \cdots \int_{-\infty}^{z_d} f(x_1, \dots, x_d) dx_1 \cdots dx_d \quad \text{Multivariate CDF}$$

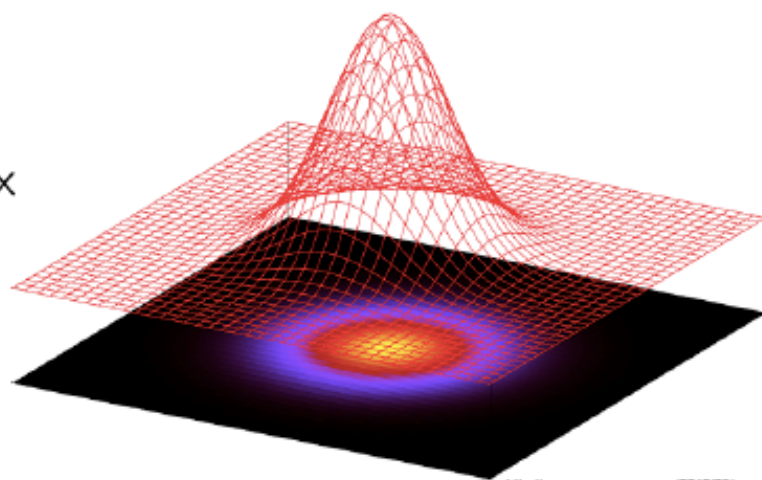
# Multivariate Gaussian distribution

For  $A \subset \mathbb{R}^d$ ,  $P([X_1, \dots, X_d] \in A) = \int_A f(x_1, \dots, x_d) dx_1 \cdots dx_d$

$F_X(z_1, \dots, z_d) = \int_{-\infty}^{z_1} \cdots \int_{-\infty}^{z_d} f(x_1, \dots, x_d) dx_1 \cdots dx_d$       Multivariate CDF

$\mu \in \mathbb{R}^d$  : mean vector

$\Sigma \in \mathbb{R}^{d \times d}$  : covariance matrix



<http://www.mosenware.com/2010/03/computing-your-skill.htm>

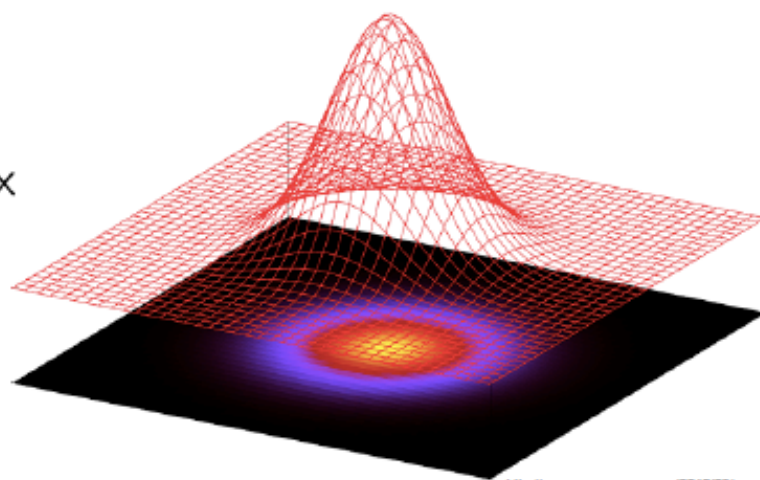
# Multivariate Gaussian distribution

For  $A \subset \mathbb{R}^d$ ,  $P([X_1, \dots, X_d] \in A) = \int_A f(x_1, \dots, x_d) dx_1 \cdots dx_d$

$F_X(z_1, \dots, z_d) = \int_{-\infty}^{z_1} \cdots \int_{-\infty}^{z_d} f(x_1, \dots, x_d) dx_1 \cdots dx_d$       Multivariate CDF

$\mu \in \mathbb{R}^d$  : mean vector

$\Sigma \in \mathbb{R}^{d \times d}$  : covariance matrix



$$f_X(x_1, \dots, x_d) = \frac{1}{\sqrt{|2\pi\Sigma|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu)\right)$$

# Conditional Probability

$P(X|Y)$  = Fraction of worlds in which X event is true given Y event is true.

$$P(X|Y) = \frac{P(X, Y)}{P(Y)}$$



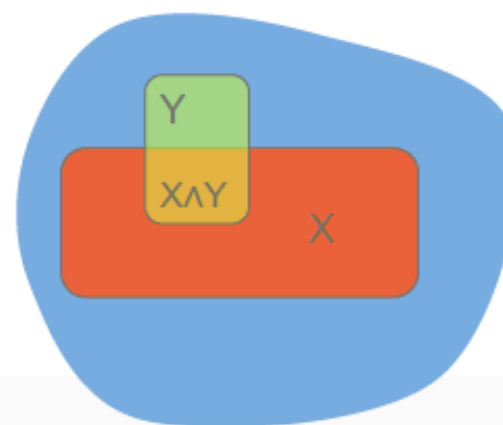
# Conditional Probability

$P(X|Y)$  = Fraction of worlds in which X event is true given Y event is true.

$$P(X|Y) = \frac{P(X, Y)}{P(Y)}$$

$$P(\text{flu}|\text{headache}) = \frac{P(\text{flu, headache})}{P(\text{headache})} = \frac{1/80}{1/80 + 7/80}$$

	Flu	No Flu
Headache	1/80	7/80
No Headache	1/80	71/80



# Independence

# Independence

**Independent random variables:**

$$P(X, Y) = P(X)P(Y)$$

$$P(X|Y) = P(X)$$

# Independence

**Independent random variables:**

$$P(X, Y) = P(X)P(Y)$$

$$P(X|Y) = P(X)$$

Y and X don't contain information about each other.

Observing Y doesn't help predicting X.

Observing X doesn't help predicting Y.

## **Examples:**

Independent: Winning on roulette this week and next week.

Dependent: Russian roulette

# Conditionally Independent

# Conditionally Independent

**Conditionally independent:**

$$P(X, Y|Z) = P(X|Z)P(Y|Z)$$

Knowing  $Z$  makes  $X$  and  $Y$  independent

# Conditionally Independent

**Conditionally independent:**

$$P(X, Y|Z) = P(X|Z)P(Y|Z)$$

Knowing Z makes X and Y independent

**Examples:**

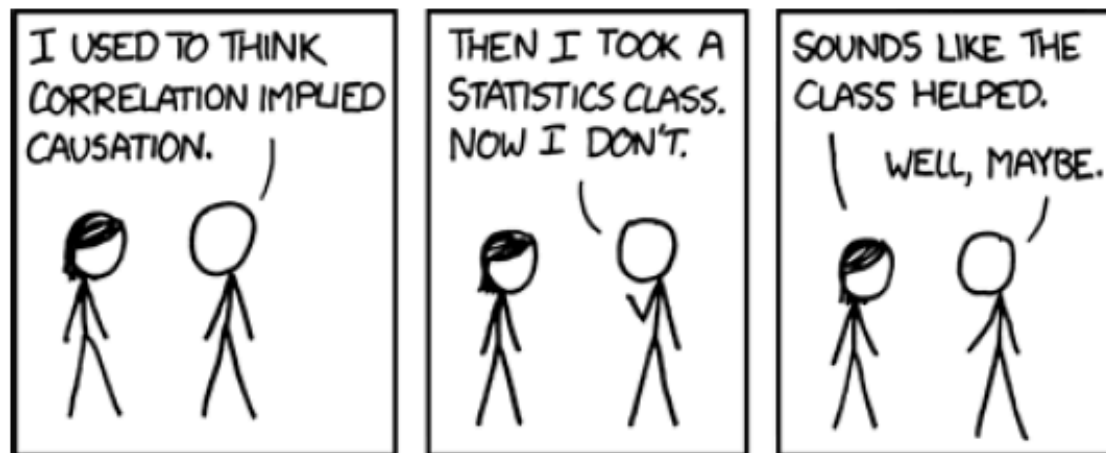
Dependent: show size and reading skills

Conditionally independent: show size and reading skills given age

# Conditionally Independent

**London taxi drivers:** A survey has pointed out a positive and significant correlation between the number of accidents and wearing coats. They concluded that coats could hinder movements of drivers and be the cause of accidents. A new law was prepared to prohibit drivers from wearing coats when driving.

Finally another study pointed out that people wear coats when it rains...





# Conditional Independence

Formally:  $X$  is **conditionally independent** of  $Y$  given  $Z$ :

# Conditional Independence

Formally: X is **conditionally independent** of Y given Z:

$$P(X, Y|Z) = P(X|Z)P(Y|Z)$$

# Conditional Independence

Formally: X is **conditionally independent** of Y given Z:

$$P(X, Y|Z) = P(X|Z)P(Y|Z)$$

$$P(\text{Accidents, Coats} | \text{Rain}) = P(\text{Accidents} | \text{Rain})P(\text{Coats} | \text{Rain})$$

# Conditional Independence

Formally: X is **conditionally independent** of Y given Z:

$$P(X, Y|Z) = P(X|Z)P(Y|Z)$$

$$P(\text{Accidents, Coats} | \text{Rain}) = P(\text{Accidents} | \text{Rain})P(\text{Coats} | \text{Rain})$$

Equivalent to:

$$(\forall x, y, z) P(X = x | Y = y, Z = z) = P(X = x | Z = z)$$

**Let's take a break!**

# Bayes Rule

# Chain Rule & Bayes Rule

---

# Chain Rule & Bayes Rule

Chain rule:

$$P(X, Y) = P(X|Y)P(Y) = P(Y|X)P(X)$$

---



# Chain Rule & Bayes Rule

Chain rule:

$$P(X, Y) = P(X|Y)P(Y) = P(Y|X)P(X)$$

Bayes rule:

$$P(X|Y) = \frac{P(Y|X)P(X)}{P(Y)}$$

---

# Chain Rule & Bayes Rule

Chain rule:

$$P(X, Y) = P(X|Y)P(Y) = P(Y|X)P(X)$$

Bayes rule:

$$P(X|Y) = \frac{P(Y|X)P(X)}{P(Y)}$$

Bayes rule is important for reverse conditioning.

# AIDS test (Bayes rule)

## Data

- ☐ Approximately 0.1% are infected
  - ☐ Test detects all infections
  - ☐ Test reports positive for 1% healthy people
-





# AIDS test (Bayes rule)

## Data

- ❑ Approximately 0.1% are infected
- ❑ Test detects all infections
- ❑ Test reports positive for 1% healthy people

Probability of having AIDS if test is positive:

$$\begin{aligned}P(a = 1|t = 1) &= \frac{P(t = 1|a = 1)P(a = 1)}{P(t = 1)} \\&= \frac{P(t = 1|a = 1)P(a = 1)}{P(t = 1|a = 1)P(a = 1) + P(t = 1|a = 0)P(a = 0)} \\&= \frac{1 \cdot 0.001}{1 \cdot 0.001 + 0.01 \cdot 0.999} = 0.091\end{aligned}$$

Only 9%!...

# Improving the diagnosis

## **Use a follow-up test!**

- Test 2 reports positive for 90% infections
- Test 2 reports positive for 5% healthy people







# Improving the diagnosis

## Use a follow-up test!

- Test 2 reports positive for 90% infections
- Test 2 reports positive for 5% healthy people

$$\begin{aligned} P(a = 0 | t_1 = 1, t_2 = 1) &= \frac{P(t_1 = 1, t_2 = 1 | a = 0)P(a = 0)}{P(t_1 = 1, t_2 = 1 | a = 1)P(a = 1) + P(t_1 = 1, t_2 = 1 | a = 0)P(a = 0)} \\ &= \frac{0.01 \cdot 0.05 \cdot 0.999}{1 \cdot 0.9 \cdot 0.001 + 0.01 \cdot 0.05 \cdot 0.999} = 0.357 \end{aligned}$$

$$P(a = 1 | t_1 = 1, t_2 = 1) = 0.643$$

## Why can't we use Test 1 twice?

Outcomes are **not** independent but tests 1 and 2

are **conditionally independent**  $p(t_1, t_2 | a) = p(t_1 | a) \cdot p(t_2 | a)$

# Naïve Bayes Assumption

---

# Naïve Bayes Assumption

**Naïve Bayes assumption:** Features  $X_1$  and  $X_2$  are conditionally independent given the class label  $Y$ :

$$P(X_1, X_2|Y) = P(X_1|Y)P(X_2|Y)$$

---

# Naïve Bayes Assumption

**Naïve Bayes assumption:** Features  $X_1$  and  $X_2$  are conditionally independent given the class label  $Y$ :

$$P(X_1, X_2|Y) = P(X_1|Y)P(X_2|Y)$$

More generally: 
$$P(X_1 \dots X_d|Y) = \prod_{i=1}^d P(X_i|Y)$$

---

# Naïve Bayes Assumption

**Naïve Bayes assumption:** Features  $X_1$  and  $X_2$  are conditionally independent given the class label  $Y$ :

$$P(X_1, X_2|Y) = P(X_1|Y)P(X_2|Y)$$

More generally: 
$$P(X_1 \dots X_d|Y) = \prod_{i=1}^d P(X_i|Y)$$

How many parameters to estimate?

( $X$  is composed of  $d$  binary features, e.g. presence of “earn”  
 $Y$  has  $K$  possible class labels)

**$(2^d - 1)K$  vs  $(2 - 1)dK$**

# Naïve Bayes Classifier

## Given:

- Class prior  $P(Y)$
- $d$  conditionally independent features  $X_1, \dots, X_d$  given the class label  $Y$
- For each  $X_i$ , we have the conditional likelihood  $P(X_i | Y)$

# Naïve Bayes Classifier

## Given:

- Class prior  $P(Y)$
- $d$  conditionally independent features  $X_1, \dots, X_d$  given the class label  $Y$
- For each  $X_i$ , we have the conditional likelihood  $P(X_i|Y)$

## Decision rule:

$$\begin{aligned} f_{NB}(\mathbf{x}) &= \arg \max_y P(x_1, \dots, x_d | y) P(y) \\ &= \arg \max_y \prod_{i=1}^d P(x_i | y) P(y) \end{aligned}$$



# Naïve Bayes Algorithm for discrete features


# Naïve Bayes Algorithm for discrete features

Training Data:  $\{(X^{(j)}, Y^{(j)})\}_{j=1}^n$      $X^{(j)} = (X_1^{(j)}, \dots, X_d^{(j)})$

$n$   $d$  dimensional features + class labels

$$f_{NB}(\mathbf{x}) = \arg \max_y \prod_{i=1}^d P(x_i|y)P(y)$$

**We need to estimate these probabilities!**




# Naïve Bayes Algorithm for discrete features

Training Data:  $\{(X^{(j)}, Y^{(j)})\}_{j=1}^n$      $X^{(j)} = (X_1^{(j)}, \dots, X_d^{(j)})$

$n$   $d$  dimensional features + class labels

$$f_{NB}(x) = \arg \max_y \prod_{i=1}^d P(x_i|y)P(y)$$

**We need to estimate these probabilities!**



Estimate them with Relative Frequencies!

For Class Prior

$$\hat{P}(y) = \frac{\#\{j : Y^{(j)} = y\}}{n}$$

For Likelihood

$$\frac{\hat{P}(x_i, y)}{\hat{P}(y)} = \frac{\#\{j : X_i^{(j)} = x_i, Y^{(j)} = y\}/n}{\#\{j : Y^{(j)} = y\}/n}$$

# Naïve Bayes Algorithm for discrete features

Training Data:  $\{(X^{(j)}, Y^{(j)})\}_{j=1}^n$      $X^{(j)} = (X_1^{(j)}, \dots, X_d^{(j)})$

$n$   $d$  dimensional features + class labels

$$f_{NB}(x) = \arg \max_y \prod_{i=1}^d P(x_i|y)P(y)$$

**We need to estimate these probabilities!**

Estimate them with Relative Frequencies!

For Class Prior     $\hat{P}(y) = \frac{\#\{j : Y^{(j)} = y\}}{n}$

For Likelihood     $\frac{\hat{P}(x_i, y)}{\hat{P}(y)} = \frac{\#\{j : X_i^{(j)} = x_i, Y^{(j)} = y\}/n}{\#\{j : Y^{(j)} = y\}/n}$

NB Prediction for test data:

$$X = (x_1, \dots, x_d)$$

$$Y = \arg \max_y \hat{P}(y) \prod_{i=1}^d \frac{\hat{P}(x_i, y)}{\hat{P}(y)}$$

# Subtlety: Insufficient training data

What if you never see a training instance where  $X_1 = a$  when  $Y = b$ ?

For example,

there is no  $X_1 = \text{'Earn'}$  when  $Y = \text{'SpamEmail'}$  in our dataset.

$$\Rightarrow P(X_1 = a, Y = b) = 0 \Rightarrow P(X_1 = a | Y = b) = 0$$

---

# Subtlety: Insufficient training data

What if you never see a training instance where  $X_1 = a$  when  $Y = b$ ?

---

# Subtlety: Insufficient training data

What if you never see a training instance where  $X_1 = a$  when  $Y = b$ ?

For example,

there is no  $X_1 = \text{'Earn'}$  when  $Y = \text{'SpamEmail'}$  in our dataset.

$$\Rightarrow P(X_1 = a, Y = b) = 0 \Rightarrow P(X_1 = a | Y = b) = 0$$

# Subtlety: Insufficient training data

What if you never see a training instance where  $X_1 = a$  when  $Y = b$ ?

For example,

there is no  $X_1 = \text{'Earn'}$  when  $Y = \text{'SpamEmail'}$  in our dataset.

$$\Rightarrow P(X_1 = a, Y = b) = 0 \Rightarrow P(X_1 = a | Y = b) = 0$$

$$\Rightarrow P(X_1 = a, X_2 \dots X_n | Y) = P(X_1 = a | Y) \prod_{i=2}^d P(X_i | Y) = 0$$



# Subtlety: Insufficient training data

What if you never see a training instance where  $X_1 = a$  when  $Y = b$ ?

For example,

there is no  $X_1 = \text{'Earn'}$  when  $Y = \text{'SpamEmail'}$  in our dataset.

$$\Rightarrow P(X_1 = a, Y = b) = 0 \Rightarrow P(X_1 = a | Y = b) = 0$$

$$\Rightarrow P(X_1 = a, X_2 \dots X_n | Y) = P(X_1 = a | Y) \prod_{i=2}^d P(X_i | Y) = 0$$

Thus, no matter what the values  $X_2, \dots, X_d$  take:

$$P(Y = b | X_1 = a, X_2, \dots, X_d) = 0$$

What now???

**See you next week!**