#### Probability Theory Refresher

#### Mário A. T. Figueiredo

Instituto Superior Técnico & Instituto de Telecomunicações

Lisboa, Portugal

#### LxMLS 2016: Lisbon Machine Learning School

#### July 21, 2016





• The study of probability has roots in games of chance





• The study of probability has roots in games of chance



 Great names of science: Cardano, Fermat, Pascal, Laplace, Kolmogorov, Bernoulli, Poisson, Cauchy, Boltzman, de Finetti, ...



• The study of probability has roots in games of chance



- Great names of science: Cardano, Fermat, Pascal, Laplace, Kolmogorov, Bernoulli, Poisson, Cauchy, Boltzman, de Finetti, ...
- Natural tool to model uncertainty, information, knowledge, belief, observations, ...



• The study of probability has roots in games of chance



- Great names of science: Cardano, Fermat, Pascal, Laplace, Kolmogorov, Bernoulli, Poisson, Cauchy, Boltzman, de Finetti, ...
- Natural tool to model uncertainty, information, knowledge, belief, observations, ...
- ...thus also learning, decision making, inference, science,...

#### What is probability?

• Classical definition: 
$$\mathbb{P}(A) = \frac{N_A}{N}$$

...with N mutually exclusive equally likely outcomes,  $N_A$  of which result in the occurrence of A. Laplace, 1814

Example:  $\mathbb{P}(\text{randomly drawn card is }) = 13/52.$ 

**Example**:  $\mathbb{P}(\text{getting 1 in throwing a fair die}) = 1/6.$ 

### What is probability?

• Classical definition: 
$$\mathbb{P}(A) = \frac{N_A}{N}$$

...with N mutually exclusive equally likely outcomes,  $N_A$  of which result in the occurrence of A. Laplace, 1814

Example:  $\mathbb{P}(\text{randomly drawn card is }) = 13/52.$ 

Example:  $\mathbb{P}(\text{getting 1 in throwing a fair die}) = 1/6.$ 

• Frequentist definition:  $\mathbb{P}(A) = \lim_{N \to \infty} \frac{N_A}{N}$ 

 $\dots$ relative frequency of occurrence of A in infinite number of trials.

What is probability?

• Classical definition: 
$$\mathbb{P}(A) = \frac{N_A}{N}$$

...with N mutually exclusive equally likely outcomes,  $N_A$  of which result in the occurrence of A. Laplace, 1814

Example:  $\mathbb{P}(\text{randomly drawn card is }) = 13/52.$ 

Example:  $\mathbb{P}(\text{getting 1 in throwing a fair die}) = 1/6.$ 

• Frequentist definition:  $\mathbb{P}(A) = \lim_{N \to \infty} \frac{N_A}{N}$ 

 $\dots$ relative frequency of occurrence of A in infinite number of trials.

• Subjective probability:  $\mathbb{P}(A)$  is a degree of belief. *de Finetti, 1930s* 

...gives meaning to  $\mathbb{P}($  "it will rain tomorrow" ).

- Sample space  $\mathcal{X} =$  set of possible outcomes of a random experiment. Examples:
  - Tossing two coins:  $\mathcal{X} = \{HH, TH, HT, TT\}$
  - Roulette:  $\mathcal{X} = \{1, 2, ..., 36\}$
  - Draw a card from a shuffled deck:  $\mathcal{X} = \{A, 2, ..., Q \diamondsuit, K \diamondsuit\}$ .

- Sample space  $\mathcal{X} =$  set of possible outcomes of a random experiment. Examples:
  - Tossing two coins:  $\mathcal{X} = \{HH, TH, HT, TT\}$
  - Roulette:  $\mathcal{X} = \{1, 2, ..., 36\}$
  - Draw a card from a shuffled deck:  $\mathcal{X} = \{A\clubsuit, 2\clubsuit, ..., Q\diamondsuit, K\diamondsuit\}.$
- An event A is a subset of  $\mathcal{X}$ :  $A \subseteq \mathcal{X}$  (also written  $A \in 2^{\mathcal{X}}$ ).

Examples:

- "exactly one H in 2-coin toss":  $A = \{TH, HT\} \subset \{HH, TH, HT, TT\}.$
- "odd number in the roulette":  $B = \{1, 3, ..., 35\} \subset \{1, 2, ..., 36\}$ .
- "drawn a  $\heartsuit$  card":  $C = \{A\heartsuit, 2\heartsuit, ..., K\heartsuit\} \subset \{A\clubsuit, ..., K\diamondsuit\}$

- Sample space X = set of possible outcomes of a random experiment.
   (More delicate) examples:
  - Distance travelled by tossed die:  $\mathcal{X} = \mathbb{R}_+$
  - Location of the next rain drop on a given square tile:  $\mathcal{X} = \mathbb{R}^2$

- Sample space X = set of possible outcomes of a random experiment. (More delicate) examples:
  - Distance travelled by tossed die:  $\mathcal{X} = \mathbb{R}_+$
  - Location of the next rain drop on a given square tile:  $\mathcal{X} = \mathbb{R}^2$
- Properly handling the continuous case requires deeper concepts:

- Sample space X = set of possible outcomes of a random experiment.
   (More delicate) examples:
  - Distance travelled by tossed die:  $\mathcal{X} = \mathbb{R}_+$
  - Location of the next rain drop on a given square tile:  $\mathcal{X} = \mathbb{R}^2$
- Properly handling the continuous case requires deeper concepts:
  - Let  $\Sigma$  be collection of subsets of  $\mathcal{X}$ :  $\Sigma \subseteq 2^{\mathcal{X}}$

- Sample space X = set of possible outcomes of a random experiment.
   (More delicate) examples:
  - Distance travelled by tossed die:  $\mathcal{X} = \mathbb{R}_+$
  - Location of the next rain drop on a given square tile:  $\mathcal{X} = \mathbb{R}^2$
- Properly handling the continuous case requires deeper concepts:
  - Let  $\Sigma$  be collection of subsets of  $\mathcal{X}$ :  $\Sigma \subseteq 2^{\mathcal{X}}$
  - $\Sigma$  is a  $\sigma$ -algebra if

$$\begin{array}{l} \star \ A\in\Sigma\Rightarrow A^{c}\in\Sigma\\ \star \ A_{1},A_{2},\ldots\in\Sigma\Rightarrow\bigcup_{i=1}^{\infty}A_{i}\in\Sigma\end{array}$$

- Sample space X = set of possible outcomes of a random experiment.
   (More delicate) examples:
  - Distance travelled by tossed die:  $\mathcal{X} = \mathbb{R}_+$
  - Location of the next rain drop on a given square tile:  $\mathcal{X} = \mathbb{R}^2$
- Properly handling the continuous case requires deeper concepts:
  - Let  $\Sigma$  be collection of subsets of  $\mathcal{X}$ :  $\Sigma \subseteq 2^{\mathcal{X}}$
  - $\Sigma$  is a  $\sigma$ -algebra if

★ 
$$A \in \Sigma \Rightarrow A^c \in \Sigma$$
  
★  $A_1, A_2, ... \in \Sigma \Rightarrow \bigcup_{i=1}^{\infty} A_i \in \Sigma$ 

• Corollary: if  $\Sigma \subset 2^{\mathcal{X}}$  is a  $\sigma$ -algebra,  $\emptyset \in \Sigma$  and  $\mathcal{X} \in \Sigma$ 

- Sample space X = set of possible outcomes of a random experiment. (More delicate) examples:
  - Distance travelled by tossed die:  $\mathcal{X} = \mathbb{R}_+$
  - Location of the next rain drop on a given square tile:  $\mathcal{X} = \mathbb{R}^2$
- Properly handling the continuous case requires deeper concepts:
  - Let  $\Sigma$  be collection of subsets of  $\mathcal{X}$ :  $\Sigma \subseteq 2^{\mathcal{X}}$
  - $\Sigma$  is a  $\sigma$ -algebra if

$$\star A \in \Sigma \Rightarrow A^c \in \Sigma$$
$$\star A_1, A_2, \dots \in \Sigma \Rightarrow \bigcup_{i=1}^{\infty} A_i \in \Sigma$$

- Corollary: if  $\Sigma \subset 2^{\mathcal{X}}$  is a  $\sigma$ -algebra,  $\emptyset \in \Sigma$  and  $\mathcal{X} \in \Sigma$
- Example in  $\mathbb{R}^n$ : collection of Lebesgue-measurable sets is a  $\sigma$ -algebra.

Probability is a function that maps events A into the interval [0, 1].
 Kolmogorov's axioms (1933) for probability P : Σ → [0, 1]

• Probability is a function that maps events A into the interval [0, 1]. Kolmogorov's axioms (1933) for probability  $\mathbb{P}: \Sigma \to [0, 1]$ 

• For any A,  $\mathbb{P}(A) \ge 0$ 

Probability is a function that maps events A into the interval [0,1].
 Kolmogorov's axioms (1933) for probability P : Σ → [0, 1]

• For any 
$$A$$
,  $\mathbb{P}(A) \ge 0$ 

$$\blacktriangleright \mathbb{P}(\mathcal{X}) = 1$$

Probability is a function that maps events A into the interval [0, 1].
 Kolmogorov's axioms (1933) for probability P : Σ → [0, 1]

• For any 
$$A$$
,  $\mathbb{P}(A) \ge 0$ 

$$\blacktriangleright \mathbb{P}(\mathcal{X}) = 1$$

• If  $A_1, A_2 ... \subseteq \mathcal{X}$  are disjoint events, then  $\mathbb{P}\left(\bigcup_i A_i\right) = \sum_i \mathbb{P}(A_i)$ 

- Probability is a function that maps events A into the interval [0, 1].
   Kolmogorov's axioms (1933) for probability P : Σ → [0, 1]
  - For any A,  $\mathbb{P}(A) \ge 0$
  - $\blacktriangleright \mathbb{P}(\mathcal{X}) = 1$
  - If  $A_1, A_2 ... \subseteq \mathcal{X}$  are disjoint events, then  $\mathbb{P}\left(\bigcup_i A_i\right) = \sum_i \mathbb{P}(A_i)$
- From these axioms, many results can be derived. Examples:
- $\blacktriangleright \ \mathbb{P}(\emptyset) = 0$
- $\blacktriangleright \ C \subset D \ \Rightarrow \ \mathbb{P}(C) \le \mathbb{P}(D)$
- $\blacktriangleright \mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B) \mathbb{P}(A \cap B)$
- $\mathbb{P}(A \cup B) \le \mathbb{P}(A) + \mathbb{P}(B)$  (union bound)



# Conditional Probability and Independence • If $\mathbb{P}(B) > 0$ , $\mathbb{P}(A|B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)}$ (conditional prob. of A, given B)

# Conditional Probability and Independence • If $\mathbb{P}(B) > 0$ , $\mathbb{P}(A|B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)}$ (conditional prob. of A, given B)

• ...satisfies all of Kolmogorov's axioms:

- For any  $A \subseteq \mathcal{X}$ ,  $\mathbb{P}(A|B) \ge 0$
- $\blacktriangleright \ \mathbb{P}(\mathcal{X}|B) = 1$

• If 
$$A_1, A_2, \dots \subseteq \mathcal{X}$$
 are disjoint, then  $\mathbb{P}\left(\bigcup_i A_i \middle| B\right) = \sum_i \mathbb{P}(A_i \middle| B)$ 



# Conditional Probability and Independence • If $\mathbb{P}(B) > 0$ , $\mathbb{P}(A|B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)}$ (conditional prob. of A, given B)

• ...satisfies all of Kolmogorov's axioms:

- For any  $A \subseteq \mathcal{X}$ ,  $\mathbb{P}(A|B) \ge 0$
- $\blacktriangleright \ \mathbb{P}(\mathcal{X}|B) = 1$

• If 
$$A_1, A_2, \dots \subseteq \mathcal{X}$$
 are disjoint, then  $\mathbb{P}\left(\bigcup_i A_i \middle| B\right) = \sum_i \mathbb{P}(A_i | B)$ 



• Independence: A, B are independent (denoted  $A \perp B$ ) if and only if

$$\mathbb{P}(A \cap B) = \mathbb{P}(A) \mathbb{P}(B).$$

 $\label{eq:conditional Probability and Independence} \ensuremath{\mathsf{o}}\ \ensuremath{\mathsf{If}}\ \ensuremath{\mathbb{P}}(B) > 0, \qquad \ensuremath{\mathbb{P}}(A|B) = \frac{\ensuremath{\mathbb{P}}(A\cap B)}{\ensuremath{\mathbb{P}}(B)}$ 

• If 
$$\mathbb{P}(B) > 0$$
,  $\mathbb{P}(A|B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)}$ 

• Events A, B are independent  $(A \perp B) \Leftrightarrow \mathbb{P}(A \cap B) = \mathbb{P}(A) \mathbb{P}(B)$ .

• If 
$$\mathbb{P}(B) > 0$$
,  $\mathbb{P}(A|B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)}$ 

- Events A, B are independent  $(A \perp B) \Leftrightarrow \mathbb{P}(A \cap B) = \mathbb{P}(A) \mathbb{P}(B)$ .
- Relationship with conditional probabilities:

$$A \perp\!\!\!\perp B \ \Leftrightarrow \ \mathbb{P}(A|B) = \mathbb{P}(A)$$

• If 
$$\mathbb{P}(B) > 0$$
,  $\mathbb{P}(A|B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)}$ 

- Events A, B are independent  $(A \perp B) \Leftrightarrow \mathbb{P}(A \cap B) = \mathbb{P}(A) \mathbb{P}(B)$ .
- Relationship with conditional probabilities:

$$A \perp\!\!\!\perp B \iff \mathbb{P}(A|B) = \mathbb{P}(A)$$

• Example:  $\mathcal{X} = \text{``52 cards''}$ ,  $A = \{3\heartsuit, 3\clubsuit, 3\diamondsuit, 3\clubsuit\}$ , and  $B = \{A\heartsuit, 2\heartsuit, ..., K\heartsuit\}$ ; then,  $\mathbb{P}(A) = 1/13$ ,  $\mathbb{P}(B) = 1/4$  $\mathbb{P}(A \cap B) = \mathbb{P}(\{3\heartsuit\}) = \frac{1}{52}$ 

• If 
$$\mathbb{P}(B) > 0$$
,  $\mathbb{P}(A|B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)}$ 

- Events A, B are independent  $(A \perp B) \Leftrightarrow \mathbb{P}(A \cap B) = \mathbb{P}(A) \mathbb{P}(B)$ .
- Relationship with conditional probabilities:

$$A \perp\!\!\!\perp B \iff \mathbb{P}(A|B) = \mathbb{P}(A)$$

• Example: 
$$\mathcal{X} =$$
 "52 cards",  $A = \{3\heartsuit, 3\clubsuit, 3\diamondsuit, 3\clubsuit\}$ , and  
 $B = \{A\heartsuit, 2\heartsuit, ..., K\heartsuit\}$ ; then,  $\mathbb{P}(A) = 1/13$ ,  $\mathbb{P}(B) = 1/4$   
 $\mathbb{P}(A \cap B) = \mathbb{P}(\{3\heartsuit\}) = \frac{1}{52}$   
 $\mathbb{P}(A) \mathbb{P}(B) = \frac{1}{13} \frac{1}{4} = \frac{1}{52}$ 

• If 
$$\mathbb{P}(B) > 0$$
,  $\mathbb{P}(A|B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)}$ 

- Events A, B are independent  $(A \perp B) \Leftrightarrow \mathbb{P}(A \cap B) = \mathbb{P}(A) \mathbb{P}(B)$ .
- Relationship with conditional probabilities:

$$A \perp\!\!\!\perp B \iff \mathbb{P}(A|B) = \mathbb{P}(A)$$

• Example: 
$$\mathcal{X} = \text{``52 cards''}, A = \{3\heartsuit, 3\clubsuit, 3\diamondsuit, 3\diamondsuit, 3\clubsuit\}, \text{ and}$$
  
 $B = \{A\heartsuit, 2\heartsuit, ..., K\heartsuit\}; \text{ then, } \mathbb{P}(A) = 1/13, \mathbb{P}(B) = 1/4$   
 $\mathbb{P}(A \cap B) = \mathbb{P}(\{3\heartsuit\}) = \frac{1}{52}$   
 $\mathbb{P}(A) \mathbb{P}(B) = \frac{1}{13} \frac{1}{4} = \frac{1}{52}$   
 $\mathbb{P}(A|B) = \mathbb{P}(\text{``3''}|\text{``\S''}) = \frac{1}{13} = \mathbb{P}(A)$ 

#### **Bayes Theorem**

• Law of total probability: if  $A_1, ..., A_n$  are a partition of  $\mathcal X$ 

$$\mathbb{P}(B) = \sum_{i} \mathbb{P}(B|A_i)\mathbb{P}(A_i)$$
$$= \sum_{i} \mathbb{P}(B \cap A_i)$$



#### **Bayes Theorem**

• Law of total probability: if  $A_1, ..., A_n$  are a partition of  $\mathcal{X}$ 

$$\mathbb{P}(B) = \sum_{i} \mathbb{P}(B|A_{i})\mathbb{P}(A_{i})$$

$$= \sum_{i} \mathbb{P}(B \cap A_{i})$$

$$X$$

$$A_{1}$$

$$A_{2}$$

$$A_{2}$$

$$A_{1}$$

$$A_{2}$$

$$A_{1}$$

$$A_{2}$$

$$A_{3}$$

$$A_{2}$$

$$A_{3}$$

$$A_{2}$$

$$A_{3}$$

$$A_{4}$$

$$A_{5}$$

$$A_{6}$$

$$A_{$$

• Bayes' theorem: if  $\{A_1, ..., A_n\}$  is a partition of  $\mathcal{X}$ 

$$\mathbb{P}(A_i|B) = \frac{\mathbb{P}(B \cap A_i)}{\mathbb{P}(B)} = \frac{\mathbb{P}(B|A_i) \mathbb{P}(A_i)}{\sum_i \mathbb{P}(B|A_i) \mathbb{P}(A_i)}$$

#### **Random Variables**

• A (real) random variable (RV) is a function:  $X : \mathcal{X} \to \mathbb{R}$ 



#### **Random Variables**

• A (real) random variable (RV) is a function:  $X : \mathcal{X} \to \mathbb{R}$ 



• Discrete RV: range of X is countable (e.g.,  $\mathbb{N}$  or  $\{0,1\}$ )

#### **Random Variables**

• A (real) random variable (RV) is a function:  $X : \mathcal{X} \to \mathbb{R}$ 



- Discrete RV: range of X is countable (e.g.,  $\mathbb{N}$  or  $\{0,1\}$ )
- Continuous RV: range of X is uncountable (e.g.,  $\mathbb{R}$  or [0, 1])
### **Random Variables**

• A (real) random variable (RV) is a function:  $X : \mathcal{X} \to \mathbb{R}$ 



- Discrete RV: range of X is countable (e.g.,  $\mathbb{N}$  or  $\{0,1\}$ )
- Continuous RV: range of X is uncountable (e.g.,  $\mathbb{R}$  or [0, 1])
- Example: number of head in tossing two coins,  $\mathcal{X} = \{HH, HT, TH, TT\},$  X(HH) = 2, X(HT) = X(TH) = 1, X(TT) = 0.Range of  $X = \{0, 1, 2\}.$

## Random Variables

• A (real) random variable (RV) is a function:  $X : \mathcal{X} \to \mathbb{R}$ 



- Discrete RV: range of X is countable (e.g.,  $\mathbb{N}$  or  $\{0, 1\}$ )
- Continuous RV: range of X is uncountable (e.g.,  $\mathbb{R}$  or [0, 1])
- Example: number of head in tossing two coins,  $\mathcal{X} = \{HH, HT, TH, TT\},\$ X(HH) = 2, X(HT) = X(TH) = 1, X(TT) = 0.Range of  $X = \{0, 1, 2\}.$
- Example: distance traveled by a tossed coin; range of  $X = \mathbb{R}_+$ .

• Distribution function:  $F_X(x) = \mathbb{P}(\{\omega \in \mathcal{X} : X(\omega) \le x\})$ 



• Distribution function:  $F_X(x) = \mathbb{P}(\{\omega \in \mathcal{X} : X(\omega) \le x\})$ 



• Example: number of heads in tossing 2 coins;  $range(X) = \{0, 1, 2\}$ .



• Distribution function:  $F_X(x) = \mathbb{P}(\{\omega \in \mathcal{X} : X(\omega) \le x\})$ 



• Example: number of heads in tossing 2 coins;  $range(X) = \{0, 1, 2\}$ .



 $F_X: \mathbb{R} \to [0,1]$  is the distribution function of some r.v. X iff:

 $F_X: \mathbb{R} \to [0,1]$  is the distribution function of some r.v. X iff:

• it is non-decreasing:  $x_1 < x_2 \Rightarrow F_X(x_1) \le F_X(x_2)$ ;

 $F_X: \mathbb{R} \to [0,1]$  is the distribution function of some r.v. X iff:

• it is non-decreasing:  $x_1 < x_2 \Rightarrow F_X(x_1) \le F_X(x_2)$ ;

• 
$$\lim_{x \to -\infty} F_X(x) = 0;$$

 $F_X: \mathbb{R} \to [0,1]$  is the distribution function of some r.v. X iff:

- it is non-decreasing:  $x_1 < x_2 \Rightarrow F_X(x_1) \le F_X(x_2)$ ;
- $\lim_{x \to -\infty} F_X(x) = 0;$
- $\lim_{x \to +\infty} F_X(x) = 1;$

 $F_X: \mathbb{R} \to [0,1]$  is the distribution function of some r.v. X iff:

- it is non-decreasing:  $x_1 < x_2 \Rightarrow F_X(x_1) \le F_X(x_2)$ ;
- $\lim_{x \to -\infty} F_X(x) = 0;$
- $\lim_{x \to +\infty} F_X(x) = 1;$
- it is right semi-continuous:  $\lim_{x \to z^+} F_X(x) = F_X(z)$

 $F_X: \mathbb{R} \to [0,1]$  is the distribution function of some r.v. X iff:

- it is non-decreasing:  $x_1 < x_2 \Rightarrow F_X(x_1) \le F_X(x_2)$ ;
- $\lim_{x \to -\infty} F_X(x) = 0;$
- $\lim_{x \to +\infty} F_X(x) = 1;$
- it is right semi-continuous:  $\lim_{x \to z^+} F_X(x) = F_X(z)$

Further properties:

 $F_X: \mathbb{R} \to [0,1]$  is the distribution function of some r.v. X iff:

- it is non-decreasing:  $x_1 < x_2 \Rightarrow F_X(x_1) \le F_X(x_2)$ ;
- $\lim_{x \to -\infty} F_X(x) = 0;$
- $\lim_{x \to +\infty} F_X(x) = 1;$
- it is right semi-continuous:  $\lim_{x \to z^+} F_X(x) = F_X(z)$

Further properties:

• 
$$\mathbb{P}(X = x) = f_X(x) = F_X(x) - \lim_{z \to x^-} F_X(z);$$

 $F_X: \mathbb{R} \to [0,1]$  is the distribution function of some r.v. X iff:

- it is non-decreasing:  $x_1 < x_2 \Rightarrow F_X(x_1) \le F_X(x_2)$ ;
- $\lim_{x \to -\infty} F_X(x) = 0;$
- $\lim_{x \to +\infty} F_X(x) = 1;$
- it is right semi-continuous:  $\lim_{x \to z^+} F_X(x) = F_X(z)$

Further properties:

• 
$$\mathbb{P}(X=x) = f_X(x) = F_X(x) - \lim_{z \to x^-} F_X(z);$$

• 
$$\mathbb{P}(z < X \leq y) = F_X(y) - F_X(z);$$

 $F_X: \mathbb{R} \to [0,1]$  is the distribution function of some r.v. X iff:

- it is non-decreasing:  $x_1 < x_2 \Rightarrow F_X(x_1) \le F_X(x_2)$ ;
- $\lim_{x \to -\infty} F_X(x) = 0;$
- $\lim_{x \to +\infty} F_X(x) = 1;$
- it is right semi-continuous:  $\lim_{x \to z^+} F_X(x) = F_X(z)$

Further properties:

• 
$$\mathbb{P}(X = x) = f_X(x) = F_X(x) - \lim_{z \to x^-} F_X(z);$$

• 
$$\mathbb{P}(z < X \leq y) = F_X(y) - F_X(z);$$

•  $\mathbb{P}(X > x) = 1 - F_X(x).$ 

• Uniform:  $X \in \{x_1, ..., x_K\}$ , pmf  $f_X(x_i) = 1/K$ .

- Uniform:  $X \in \{x_1, ..., x_K\}$ , pmf  $f_X(x_i) = 1/K$ .
- Bernoulli RV:  $X \in \{0,1\}$ , pmf  $f_X(x) = \begin{cases} p \iff x = 1\\ 1-p \iff x = 0 \end{cases}$

Can be written compactly as  $f_X(x) = p^x(1-p)^{1-x}$ .

• Uniform:  $X \in \{x_1, ..., x_K\}$ , pmf  $f_X(x_i) = 1/K$ .

• Bernoulli RV: 
$$X \in \{0,1\}$$
, pmf  $f_X(x) = \begin{cases} p \iff x = 1\\ 1-p \iff x = 0 \end{cases}$ 

Can be written compactly as  $f_X(x) = p^x(1-p)^{1-x}$ .

• Binomial RV:  $X \in \{0, 1, ..., n\}$  (sum of n Bernoulli RVs)

$$f_X(x) = \mathsf{Binomial}(x; n, p) = \binom{n}{x} p^x (1-p)^{(n-x)}$$

• Uniform:  $X \in \{x_1, ..., x_K\}$ , pmf  $f_X(x_i) = 1/K$ .

• Bernoulli RV: 
$$X \in \{0,1\}$$
, pmf  $f_X(x) = \begin{cases} p & \Leftarrow x = 1\\ 1-p & \Leftarrow x = 0 \end{cases}$ 

Can be written compactly as  $f_X(x) = p^x(1-p)^{1-x}$ .

• Binomial RV:  $X \in \{0, 1, ..., n\}$  (sum of n Bernoulli RVs)

$$f_X(x) = \mathsf{Binomial}(x; n, p) = \binom{n}{x} p^x (1-p)^{(n-x)}$$



• Geometric(p):  $X \in \mathbb{N}$ , pmf  $f_X(x) = p(1-p)^{x-1}$ . (*e.g.*, number of trials until the first success).

- Geometric(p):  $X \in \mathbb{N}$ , pmf  $f_X(x) = p(1-p)^{x-1}$ . (e.g., number of trials until the first success).
- Poisson( $\lambda$ ):  $X \in \mathbb{N} \cup \{0\}$ , pmf  $f_X(x) = \frac{e^{-\lambda}\lambda^x}{x!}$

Notice that  $\sum_{x=0}^{\infty} \frac{\lambda^x}{x!} = e^{\lambda}$ , thus  $\sum_{x=0}^{\infty} f_X(x) = 1$ .

"...probability of the number of independent occurrences in a fixed (time/space) interval if these occurrences have known average rate"



• Distribution function:  $F_X(x) = \mathbb{P}(\{\omega \in \mathcal{X} : X(\omega) \le x\})$ 



• Distribution function:  $F_X(x) = \mathbb{P}(\{\omega \in \mathcal{X} : X(\omega) \le x\})$ 



• Example: continuous RV with uniform distribution on [a, b].



• Distribution function:  $F_X(x) = \mathbb{P}(\{\omega \in \mathcal{X} : X(\omega) \le x\})$ 



• Example: continuous RV with uniform distribution on [a, b].



• Probability density function (pdf, continuous RV):  $f_X(x)$ 

• Distribution function:  $F_X(x) = \mathbb{P}(\{\omega \in \mathcal{X} : X(\omega) \le x\})$ 



• Example: continuous RV with uniform distribution on [a, b].



• Probability density function (pdf, continuous RV):  $f_X(x)$ 

$$F_X(x) = \int_{-\infty}^x f_X(u) \, du$$

• Distribution function:  $F_X(x) = \mathbb{P}(\{\omega \in \mathcal{X} : X(\omega) \le x\})$ 



• Example: continuous RV with uniform distribution on [a, b].



• Probability density function (pdf, continuous RV):  $f_X(x)$  $F_X(x) = \int_{-\infty}^x f_X(u) \, du, \quad \mathbb{P}(X \in [c, d]) = \int_c^d f_X(x) \, dx,$ 

• Distribution function:  $F_X(x) = \mathbb{P}(\{\omega \in \mathcal{X} : X(\omega) \le x\})$ 



• Example: continuous RV with uniform distribution on [a, b].



• Probability density function (pdf, continuous RV):  $f_X(x)$ 

$$F_X(x) = \int_{-\infty} f_X(u) \, du, \quad \mathbb{P}(X \in [c, d]) = \int_c f_X(x) \, dx, \quad \mathbb{P}(X = x) = 0$$

Important Continuous Random Variables

• Uniform:  $f_X(x) = \text{Uniform}(x; a, b) = \begin{cases} \frac{1}{b-a} & \Leftarrow & x \in [a, b] \\ 0 & \Leftarrow & x \notin [a, b] \end{cases}$  (previous slide).

Important Continuous Random Variables

• Uniform:  $f_X(x) = \text{Uniform}(x; a, b) = \begin{cases} \frac{1}{b-a} & \Leftarrow & x \in [a, b] \\ 0 & \Leftarrow & x \notin [a, b] \end{cases}$  (previous slide).

• Gaussian: 
$$f_X(x) = \mathcal{N}(x;\mu,\sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}}e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



Important Continuous Random Variables

• Uniform:  $f_X(x) = \text{Uniform}(x; a, b) = \begin{cases} \frac{1}{b-a} & \Leftarrow & x \in [a, b] \\ 0 & \Leftarrow & x \notin [a, b] \end{cases}$  (previous slide).

• Gaussian: 
$$f_X(x) = \mathcal{N}(x;\mu,\sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



• Exponential:  $f_X(x) = \mathsf{Exp}(x; \lambda) = \begin{cases} \lambda e^{-\lambda x} & \Leftarrow x \ge 0 \\ 0 & \Leftarrow x < 0 \end{cases}$ 

• Expectation: 
$$\mathbb{E}(X) = \begin{cases} \sum_{i} x_i f_X(x_i) & X \in \{x_1, \dots, x_K\} \subset \mathbb{R} \\ \int_{-\infty}^{\infty} x f_X(x) dx & X \text{ continuous} \end{cases}$$

• Expectation: 
$$\mathbb{E}(X) = \begin{cases} \sum_{i} x_i f_X(x_i) & X \in \{x_1, \dots, x_K\} \subset \mathbb{R} \\ \int_{-\infty}^{\infty} x f_X(x) dx & X \text{ continuous} \end{cases}$$

• Example: Bernoulli,  $f_X(x) = p^x (1-p)^{1-x}$ , for  $x \in \{0, 1\}$ .  $\mathbb{E}(X) = 0 (1-p) + 1 p = p.$ 

• Expectation: 
$$\mathbb{E}(X) = \begin{cases} \sum_{i} x_i f_X(x_i) & X \in \{x_1, \dots, x_K\} \subset \mathbb{R} \\ \int_{-\infty}^{\infty} x f_X(x) dx & X \text{ continuous} \end{cases}$$

• Example: Bernoulli, 
$$f_X(x) = p^x (1-p)^{1-x}$$
, for  $x \in \{0, 1\}$ .  
 $\mathbb{E}(X) = 0 (1-p) + 1 p = p.$ 

• Example: Binomial, 
$$f_X(x) = \binom{n}{x} p^x (1-p)^{n-x}$$
, for  $x \in \{0, ..., n\}$ .  
 $\mathbb{E}(X) = n p.$ 

• Expectation: 
$$\mathbb{E}(X) = \begin{cases} \sum_{i} x_i f_X(x_i) & X \in \{x_1, \dots, x_K\} \subset \mathbb{R} \\ \int_{-\infty}^{\infty} x f_X(x) dx & X \text{ continuous} \end{cases}$$

• Example: Bernoulli, 
$$f_X(x) = p^x (1-p)^{1-x}$$
, for  $x \in \{0, 1\}$ .  
 $\mathbb{E}(X) = 0 (1-p) + 1 p = p.$ 

• Example: Binomial, 
$$f_X(x) = \binom{n}{x} p^x (1-p)^{n-x}$$
, for  $x \in \{0, ..., n\}$ .  
 $\mathbb{E}(X) = n p.$ 

• Example: Gaussian,  $f_X(x) = \mathcal{N}(x; \mu, \sigma^2)$ .  $\mathbb{E}(X) = \mu$ .

• Expectation: 
$$\mathbb{E}(X) = \begin{cases} \sum_{i} x_i f_X(x_i) & X \in \{x_1, \dots, x_K\} \subset \mathbb{R} \\ \int_{-\infty}^{\infty} x f_X(x) dx & X \text{ continuous} \end{cases}$$

• Example: Bernoulli, 
$$f_X(x) = p^x (1-p)^{1-x}$$
, for  $x \in \{0, 1\}$ .  
 $\mathbb{E}(X) = 0 (1-p) + 1 p = p.$ 

• Example: Binomial, 
$$f_X(x) = \binom{n}{x} p^x (1-p)^{n-x}$$
, for  $x \in \{0, ..., n\}$ .  
 $\mathbb{E}(X) = n p.$ 

• Example: Gaussian,  $f_X(x) = \mathcal{N}(x; \mu, \sigma^2)$ .  $\mathbb{E}(X) = \mu$ .

• Linearity of expectation:  $\mathbb{E}(X+Y) = \mathbb{E}(X) + \mathbb{E}(Y); \quad \mathbb{E}(\alpha X) = \alpha \mathbb{E}(X), \ \alpha \in \mathbb{R}$ 

### Expectation of Functions of Random Variables

• 
$$\mathbb{E}(g(X)) = \begin{cases} \sum_{i} g(x_i) f_X(x_i) & X \text{ discrete, } g(x_i) \in \mathbb{R} \\ \int_{-\infty}^{\infty} g(x) f_X(x) \, dx & X \text{ continuous} \end{cases}$$

#### Expectation of Functions of Random Variables

• 
$$\mathbb{E}(g(X)) = \begin{cases} \sum_{i} g(x_i) f_X(x_i) & X \text{ discrete, } g(x_i) \in \mathbb{R} \\ \int_{-\infty}^{\infty} g(x) f_X(x) \, dx & X \text{ continuous} \end{cases}$$

• Example: variance,  $\operatorname{var}(X) = \mathbb{E}((X - \mathbb{E}(X))^2)$
• 
$$\mathbb{E}(g(X)) = \begin{cases} \sum_{i} g(x_i) f_X(x_i) & X \text{ discrete, } g(x_i) \in \mathbb{R} \\ \int_{-\infty}^{\infty} g(x) f_X(x) \, dx & X \text{ continuous} \end{cases}$$

• Example: variance,  $\operatorname{var}(X) = \mathbb{E}((X - \mathbb{E}(X))^2) = \mathbb{E}(X^2) - \mathbb{E}(X)^2$ 

• Example: Bernoulli variance,  $\mathbb{E}(X^2) = \mathbb{E}(X) = p$ 

• 
$$\mathbb{E}(g(X)) = \begin{cases} \sum_{i} g(x_i) f_X(x_i) & X \text{ discrete, } g(x_i) \in \mathbb{R} \\ \int_{-\infty}^{\infty} g(x) f_X(x) dx & X \text{ continuous} \end{cases}$$

• Example: variance,  $\operatorname{var}(X) = \mathbb{E}((X - \mathbb{E}(X))^2) = \mathbb{E}(X^2) - \mathbb{E}(X)^2$ 

• Example: Bernoulli variance,  $\mathbb{E}(X^2) = \mathbb{E}(X) = p$ , thus var(X) = p(1-p).

• 
$$\mathbb{E}(g(X)) = \begin{cases} \sum_{i} g(x_i) f_X(x_i) & X \text{ discrete, } g(x_i) \in \mathbb{R} \\ \int_{-\infty}^{\infty} g(x) f_X(x) dx & X \text{ continuous} \end{cases}$$

• Example: variance,  $\operatorname{var}(X) = \mathbb{E}((X - \mathbb{E}(X))^2) = \mathbb{E}(X^2) - \mathbb{E}(X)^2$ 

- Example: Bernoulli variance,  $\mathbb{E}(X^2) = \mathbb{E}(X) = p$ , thus  $\operatorname{var}(X) = p(1-p)$ .
- Example: Gaussian variance,  $\mathbb{E}((X \mu)^2) = \sigma^2$ .

• 
$$\mathbb{E}(g(X)) = \begin{cases} \sum_{i} g(x_i) f_X(x_i) & X \text{ discrete, } g(x_i) \in \mathbb{R} \\ \int_{-\infty}^{\infty} g(x) f_X(x) dx & X \text{ continuous} \end{cases}$$

• Example: variance,  $\operatorname{var}(X) = \mathbb{E}((X - \mathbb{E}(X))^2) = \mathbb{E}(X^2) - \mathbb{E}(X)^2$ 

- Example: Bernoulli variance,  $\mathbb{E}(X^2) = \mathbb{E}(X) = p$ , thus  $\operatorname{var}(X) = p(1-p)$ .
- Example: Gaussian variance,  $\mathbb{E}((X \mu)^2) = \sigma^2$ .
- Probability as expectation of indicator,  $\mathbf{1}_A(x) = \begin{cases} 1 & \Leftarrow x \in A \\ 0 & \Leftarrow x \notin A \end{cases}$

$$\mathbb{P}(X \in A) = \int_A f_X(x) \, dx = \int \mathbf{1}_A(x) \, f_X(x) \, dx = \mathbb{E}(\mathbf{1}_A(X))$$

• Joint pmf of two discrete RVs:  $f_{X,Y}(x,y) = \mathbb{P}(X = x \land Y = y).$ 

Extends trivially to more than two RVs.

- Joint pmf of two discrete RVs:  $f_{X,Y}(x,y) = \mathbb{P}(X = x \land Y = y)$ . Extends trivially to more than two RVs.
- Joint pdf of two continuous RVs:  $f_{X,Y}(x,y)$ , such that

$$\mathbb{P}((X,Y) \in A) = \iint_A f_{X,Y}(x,y) \, dx \, dy, \qquad A \in \sigma(\mathbb{R}^2)$$

Extends trivially to more than two RVs.

- Joint pmf of two discrete RVs:  $f_{X,Y}(x,y) = \mathbb{P}(X = x \land Y = y)$ . Extends trivially to more than two RVs.
- Joint pdf of two continuous RVs:  $f_{X,Y}(x,y)$ , such that

$$\mathbb{P}((X,Y) \in A) = \iint_A f_{X,Y}(x,y) \, dx \, dy, \qquad A \in \sigma(\mathbb{R}^2)$$

Extends trivially to more than two RVs.

• Marginalization: 
$$f_Y(y) = \begin{cases} \sum_x f_{X,Y}(x,y), & \text{if } X \text{ is discrete} \\ \int_{-\infty}^{\infty} f_{X,Y}(x,y) \, dx, & \text{if } X \text{ continuous} \end{cases}$$

- Joint pmf of two discrete RVs:  $f_{X,Y}(x,y) = \mathbb{P}(X = x \land Y = y)$ . Extends trivially to more than two RVs.
- Joint pdf of two continuous RVs:  $f_{X,Y}(x,y)$ , such that

$$\mathbb{P}((X,Y) \in A) = \iint_A f_{X,Y}(x,y) \, dx \, dy, \qquad A \in \sigma(\mathbb{R}^2)$$

Extends trivially to more than two RVs.

• Marginalization: 
$$f_Y(y) = \begin{cases} \sum_x f_{X,Y}(x,y), & \text{if } X \text{ is discrete} \\ \int_{-\infty}^{\infty} f_{X,Y}(x,y) \, dx, & \text{if } X \text{ continuous} \end{cases}$$

• Independence:

$$X \perp \!\!\!\perp Y \iff f_{X,Y}(x,y) = f_X(x) f_Y(y)$$

.

- Joint pmf of two discrete RVs:  $f_{X,Y}(x,y) = \mathbb{P}(X = x \land Y = y)$ . Extends trivially to more than two RVs.
- Joint pdf of two continuous RVs:  $f_{X,Y}(x,y)$ , such that

$$\mathbb{P}((X,Y) \in A) = \iint_A f_{X,Y}(x,y) \, dx \, dy, \qquad A \in \sigma(\mathbb{R}^2)$$

Extends trivially to more than two RVs.

• Marginalization: 
$$f_Y(y) = \begin{cases} \sum_x f_{X,Y}(x,y), & \text{if } X \text{ is discrete} \\ \int_{-\infty}^{\infty} f_{X,Y}(x,y) \, dx, & \text{if } X \text{ continuous} \end{cases}$$

• Independence:

$$X \perp Y \Leftrightarrow f_{X,Y}(x,y) = f_X(x) f_Y(y) \stackrel{\Rightarrow}{\not\leftarrow} \mathbb{E}(X Y) = \mathbb{E}(X) \mathbb{E}(Y).$$

• Conditional pmf (discrete RVs):  

$$f_{X|Y}(x|y) = \mathbb{P}(X = x|Y = y) = \frac{\mathbb{P}(X = x \land Y = y)}{\mathbb{P}(Y = y)} = \frac{f_{X,Y}(x,y)}{f_Y(y)}.$$

- Conditional pmf (discrete RVs):  $f_{X|Y}(x|y) = \mathbb{P}(X = x|Y = y) = \frac{\mathbb{P}(X = x \land Y = y)}{\mathbb{P}(Y = y)} = \frac{f_{X,Y}(x,y)}{f_Y(y)}.$
- Conditional pdf (continuous RVs):  $f_{X|Y}(x|y) = \frac{f_{X,Y}(x,y)}{f_Y(y)}$ ...the meaning is technically delicate.

- Conditional pmf (discrete RVs):  $f_{X|Y}(x|y) = \mathbb{P}(X = x|Y = y) = \frac{\mathbb{P}(X = x \land Y = y)}{\mathbb{P}(Y = y)} = \frac{f_{X,Y}(x,y)}{f_Y(y)}.$
- Conditional pdf (continuous RVs):  $f_{X|Y}(x|y) = \frac{f_{X,Y}(x,y)}{f_Y(y)}$ ...the meaning is technically delicate.

• Bayes' theorem: 
$$f_{X|Y}(x|y) = \frac{f_{Y|X}(y|x) f_X(x)}{f_Y(y)}$$
 (pdf or pmf).

- Conditional pmf (discrete RVs):  $f_{X|Y}(x|y) = \mathbb{P}(X = x|Y = y) = \frac{\mathbb{P}(X = x \land Y = y)}{\mathbb{P}(Y = y)} = \frac{f_{X,Y}(x,y)}{f_Y(y)}.$
- Conditional pdf (continuous RVs):  $f_{X|Y}(x|y) = \frac{f_{X,Y}(x,y)}{f_Y(y)}$ ...the meaning is technically delicate.

• Bayes' theorem: 
$$f_{X|Y}(x|y) = \frac{f_{Y|X}(y|x) f_X(x)}{f_Y(y)}$$
 (pdf or pmf).

• Also valid in the mixed case (e.g., X continuous, Y discrete).

# Joint, Marginal, and Conditional Probabilities: An Example

• A pair of binary variables  $X, Y \in \{0, 1\}$ , with joint pmf:

$f_{X,Y}(x,y)$	Y = 0	Y = I	
X = 0	1/5	2/5	
X = 1	1/10	3/10	

# Joint, Marginal, and Conditional Probabilities: An Example

• A pair of binary variables  $X, Y \in \{0, 1\}$ , with joint pmf:

$f_{X,Y}(x,y)$	Y = 0	Y = I	
X = 0	1/5	2/5	
X = I	1/10	3/10	

• Marginals: 
$$f_X(0) = \frac{1}{5} + \frac{2}{5} = \frac{3}{5}, \qquad f_X(1) = \frac{1}{10} + \frac{3}{10} = \frac{4}{10},$$
  
 $f_Y(0) = \frac{1}{5} + \frac{1}{10} = \frac{3}{10}, \qquad f_Y(1) = \frac{2}{5} + \frac{3}{10} = \frac{7}{10}.$ 

# Joint, Marginal, and Conditional Probabilities: An Example

• A pair of binary variables  $X, Y \in \{0, 1\}$ , with joint pmf:

$f_{X,Y}(x,y)$	Y = 0	Y = I	
X = 0	1/5	2/5	
X = I	1/10	3/10	

• Marginals: 
$$f_X(0) = \frac{1}{5} + \frac{2}{5} = \frac{3}{5}, \qquad f_X(1) = \frac{1}{10} + \frac{3}{10} = \frac{4}{10},$$
  
 $f_Y(0) = \frac{1}{5} + \frac{1}{10} = \frac{3}{10}, \qquad f_Y(1) = \frac{2}{5} + \frac{3}{10} = \frac{7}{10}.$ 

• Conditional probabilities:

$f_{X Y}(x y)$	Y = 0	Y = I	$f_{Y X}(y x)$	Y = 0	Y = 1
X = 0	2/3	4/7	X = 0	1/3	2/3
X = I	1/3	3/7	X = 1	1/4	3/4

## An Important Multivariate RV: Multinomial

• Multinomial:  $X = (X_1, ..., X_K)$ ,  $X_i \in \{0, ..., n\}$ , such that  $\sum_i X_i = n$ ,

$$f_X(x_1, ..., x_K) = \begin{cases} \binom{n}{x_1 \ x_2 \ \cdots \ x_K} p_1^{x_1} \ p_2^{x_2} \ \cdots \ p_k^{x_K} & \Leftarrow \ \sum_i x_i = n \\ 0 & \Leftarrow \ \sum_i x_i \neq n \end{cases}$$

$$\binom{n}{x_1 \ x_2 \ \cdots \ x_K} = \frac{n!}{x_1! \ x_2! \ \cdots \ x_K!}$$

Parameters:  $p_1, ..., p_K \ge 0$ , such that  $\sum_i p_i = 1$ .

## An Important Multivariate RV: Multinomial

• Multinomial:  $X = (X_1, ..., X_K)$ ,  $X_i \in \{0, ..., n\}$ , such that  $\sum_i X_i = n$ ,

$$f_X(x_1, ..., x_K) = \begin{cases} \binom{n}{x_1 \ x_2 \ \cdots \ x_K} p_1^{x_1} \ p_2^{x_2} \ \cdots \ p_k^{x_K} & \Leftarrow \ \sum_i x_i = n \\ 0 & \Leftarrow \ \sum_i x_i \neq n \end{cases}$$

$$\binom{n}{x_1 \ x_2 \ \cdots \ x_K} = \frac{n!}{x_1! \ x_2! \ \cdots \ x_K!}$$

Parameters:  $p_1, ..., p_K \ge 0$ , such that  $\sum_i p_i = 1$ .

• Generalizes the binomial from binary to K-classes.

# An Important Multivariate RV: Multinomial

• Multinomial:  $X = (X_1, ..., X_K)$ ,  $X_i \in \{0, ..., n\}$ , such that  $\sum_i X_i = n$ ,

$$f_X(x_1, ..., x_K) = \begin{cases} \binom{n}{x_1 \ x_2 \ \cdots \ x_K} p_1^{x_1} \ p_2^{x_2} \ \cdots \ p_k^{x_K} & \Leftarrow \ \sum_i x_i = n \\ 0 & \Leftarrow \ \sum_i x_i \neq n \end{cases}$$

$$\binom{n}{x_1 \ x_2 \ \cdots \ x_K} = \frac{n!}{x_1! \ x_2! \ \cdots \ x_K!}$$

Parameters:  $p_1, ..., p_K \ge 0$ , such that  $\sum_i p_i = 1$ .

- Generalizes the binomial from binary to K-classes.
- Example: tossing n independent fair dice,  $p_1 = \cdots = p_6 = 1/6$ .  $x_i =$  number of outcomes with i dots. Of course,  $\sum_i x_i = n$ .

# An Important Multivariate RV: Gaussian

• Multivariate Gaussian:  $X \in \mathbb{R}^n$ ,

$$f_X(x) = \mathcal{N}(x;\mu,C) = \frac{1}{\sqrt{\det(2\pi C)}} \exp\left(-\frac{1}{2}(x-\mu)^T C^{-1}(x-\mu)\right)$$

# An Important Multivariate RV: Gaussian

• Multivariate Gaussian:  $X \in \mathbb{R}^n$ ,

$$f_X(x) = \mathcal{N}(x;\mu,C) = \frac{1}{\sqrt{\det(2\pi C)}} \exp\left(-\frac{1}{2}(x-\mu)^T C^{-1}(x-\mu)\right)$$

• Parameters: vector  $\mu \in \mathbb{R}^n$  and matrix  $C \in \mathbb{R}^{n \times n}$ . Expected value:  $\mathbb{E}(X) = \mu$ . Meaning of C: next slide.

## An Important Multivariate RV: Gaussian

• Multivariate Gaussian:  $X \in \mathbb{R}^n$ ,

$$f_X(x) = \mathcal{N}(x;\mu,C) = \frac{1}{\sqrt{\det(2\pi C)}} \exp\left(-\frac{1}{2}(x-\mu)^T C^{-1}(x-\mu)\right)$$

• Parameters: vector  $\mu \in \mathbb{R}^n$  and matrix  $C \in \mathbb{R}^{n \times n}$ . Expected value:  $\mathbb{E}(X) = \mu$ . Meaning of C: next slide.



LxMLS 2016: Probability Theory

$$\operatorname{cov}(X,Y) = \mathbb{E}\Big[ \left( X - \mathbb{E}(X) \right) \left( Y - \mathbb{E}(Y) \right) \Big] \; = \; \mathbb{E}(X \, Y) - \mathbb{E}(X) \, \mathbb{E}(Y)$$

• Covariance between two RVs:

$$\operatorname{cov}(X,Y) = \mathbb{E}\Big[ \left( X - \mathbb{E}(X) \right) \left( Y - \mathbb{E}(Y) \right) \Big] \; = \; \mathbb{E}(X \, Y) - \mathbb{E}(X) \, \mathbb{E}(Y)$$

• Relationship with variance: var(X) = cov(X, X).

$$\operatorname{cov}(X,Y) = \mathbb{E}\Big[ \left( X - \mathbb{E}(X) \right) \left( Y - \mathbb{E}(Y) \right) \Big] \; = \; \mathbb{E}(X \, Y) - \mathbb{E}(X) \, \mathbb{E}(Y)$$

- Relationship with variance: var(X) = cov(X, X).
- Correlation:  $\operatorname{corr}(X,Y) = \rho(X,Y) = \frac{\operatorname{cov}(X,Y)}{\sqrt{\operatorname{var}(X)}\sqrt{\operatorname{var}(Y)}} \in [-1,\,1]$

$$\operatorname{cov}(X,Y) = \mathbb{E}\Big[ \left( X - \mathbb{E}(X) \right) \left( Y - \mathbb{E}(Y) \right) \Big] \ = \ \mathbb{E}(X \, Y) - \mathbb{E}(X) \, \mathbb{E}(Y)$$

- Relationship with variance: var(X) = cov(X, X).
- Correlation:  $\operatorname{corr}(X,Y) = \rho(X,Y) = \frac{\operatorname{cov}(X,Y)}{\sqrt{\operatorname{var}(X)}\sqrt{\operatorname{var}(Y)}} \in [-1,\,1]$

• 
$$X \perp Y \Leftrightarrow f_{X,Y}(x,y) = f_X(x) f_Y(y)$$

$$\operatorname{cov}(X,Y) = \mathbb{E}\Big[ \left( X - \mathbb{E}(X) \right) \left( Y - \mathbb{E}(Y) \right) \Big] \ = \ \mathbb{E}(X \, Y) - \mathbb{E}(X) \, \mathbb{E}(Y)$$

- Relationship with variance: var(X) = cov(X, X).
- Correlation:  $\operatorname{corr}(X, Y) = \rho(X, Y) = \frac{\operatorname{cov}(X, Y)}{\sqrt{\operatorname{var}(X)}\sqrt{\operatorname{var}(Y)}} \in [-1, 1]$ •  $X \perp Y \Leftrightarrow f_{X,Y}(x, y) = f_X(x) f_Y(y) \stackrel{\Rightarrow}{\not=} \operatorname{cov}(X, Y) = 0.$

• Covariance between two RVs:

$$\operatorname{cov}(X,Y) = \mathbb{E}\Big[ \left( X - \mathbb{E}(X) \right) \left( Y - \mathbb{E}(Y) \right) \Big] \ = \ \mathbb{E}(X \, Y) - \mathbb{E}(X) \, \mathbb{E}(Y)$$

• Relationship with variance:  $\operatorname{var}(X) = \operatorname{cov}(X, X)$ .

• Correlation: 
$$\operatorname{corr}(X, Y) = \rho(X, Y) = \frac{\operatorname{cov}(X, Y)}{\sqrt{\operatorname{var}(X)}\sqrt{\operatorname{var}(Y)}} \in [-1, 1]$$
  
•  $X \perp Y \Leftrightarrow f_{X,Y}(x, y) = f_X(x) f_Y(y) \stackrel{\Rightarrow}{\not\models} \operatorname{cov}(X, Y) = 0.$ 

• Covariance matrix of multivariate RV,  $X \in \mathbb{R}^n$ :

$$\operatorname{cov}(X) = \mathbb{E}\Big[ \big( X - \mathbb{E}(X) \big) \big( X - \mathbb{E}(X) \big)^T \Big] = \mathbb{E}(X X^T) - \mathbb{E}(X) \mathbb{E}(X)^T$$

• Covariance between two RVs:

$$\operatorname{cov}(X,Y) = \mathbb{E}\Big[ \left( X - \mathbb{E}(X) \right) \left( Y - \mathbb{E}(Y) \right) \Big] \; = \; \mathbb{E}(X \, Y) - \mathbb{E}(X) \, \mathbb{E}(Y)$$

• Relationship with variance:  $\operatorname{var}(X) = \operatorname{cov}(X, X)$ .

• Correlation: 
$$\operatorname{corr}(X, Y) = \rho(X, Y) = \frac{\operatorname{cov}(X, Y)}{\sqrt{\operatorname{var}(X)}\sqrt{\operatorname{var}(Y)}} \in [-1, 1]$$
  
•  $X \perp Y \Leftrightarrow f_{X,Y}(x, y) = f_X(x) f_Y(y) \stackrel{\Rightarrow}{\not\models} \operatorname{cov}(X, Y) = 0.$ 

• Covariance matrix of multivariate RV,  $X \in \mathbb{R}^n$ :

$$\operatorname{cov}(X) = \mathbb{E}\Big[ \big( X - \mathbb{E}(X) \big) \big( X - \mathbb{E}(X) \big)^T \Big] = \mathbb{E}(X \ X^T) - \mathbb{E}(X) \mathbb{E}(X)^T$$

• Covariance of Gaussian RV,  $f_X(x) = \mathcal{N}(x; \mu, C) \Rightarrow \operatorname{cov}(X) = C$ 

Let  $A \in \mathbb{R}^{n \times n}$  be a matrix and  $a \in \mathbb{R}^n$  a vector.

• If  $\mathbb{E}(X) = \mu$  and Y = AX, then  $\mathbb{E}(Y) = A\mu$ ;

• If 
$$\mathbb{E}(X) = \mu$$
 and  $Y = AX$ , then  $\mathbb{E}(Y) = A\mu$ ;

• If 
$$\mathbb{E}(X) = \mu$$
 and  $Y = X - \mu$ , then  $\mathbb{E}(Y) = 0$ ;

• If 
$$\mathbb{E}(X) = \mu$$
 and  $Y = AX$ , then  $\mathbb{E}(Y) = A\mu$ ;

• If 
$$\mathbb{E}(X) = \mu$$
 and  $Y = X - \mu$ , then  $\mathbb{E}(Y) = 0$ ;

• If 
$$cov(X) = C$$
 and  $Y = AX$ , then  $cov(Y) = ACA^T$ ;

• If 
$$\mathbb{E}(X) = \mu$$
 and  $Y = AX$ , then  $\mathbb{E}(Y) = A\mu$ ;

• If 
$$\mathbb{E}(X) = \mu$$
 and  $Y = X - \mu$ , then  $\mathbb{E}(Y) = 0$ ;

• If 
$$\operatorname{cov}(X) = C$$
 and  $Y = AX$ , then  $\operatorname{cov}(Y) = ACA^T$ ;

• If 
$$\operatorname{cov}(X) = C$$
 and  $Y = a^T X \in \mathbb{R}$ , then  $\operatorname{var}(Y) = a^T C a \ge 0$ ;

• If 
$$\mathbb{E}(X) = \mu$$
 and  $Y = AX$ , then  $\mathbb{E}(Y) = A\mu$ ;

• If 
$$\mathbb{E}(X) = \mu$$
 and  $Y = X - \mu$ , then  $\mathbb{E}(Y) = 0$ ;

• If 
$$\operatorname{cov}(X) = C$$
 and  $Y = AX$ , then  $\operatorname{cov}(Y) = ACA^T$ ;

• If 
$$\operatorname{cov}(X) = C$$
 and  $Y = a^T X \in \mathbb{R}$ , then  $\operatorname{var}(Y) = a^T C a \ge 0$ ;

• If 
$$\operatorname{cov}(X) = C$$
 and  $Y = C^{-1/2}X$ , then  $\operatorname{cov}(Y) = I$ ;

Let  $A \in \mathbb{R}^{n \times n}$  be a matrix and  $a \in \mathbb{R}^n$  a vector.

• If 
$$\mathbb{E}(X) = \mu$$
 and  $Y = AX$ , then  $\mathbb{E}(Y) = A\mu$ ;

• If 
$$\mathbb{E}(X) = \mu$$
 and  $Y = X - \mu$ , then  $\mathbb{E}(Y) = 0$ ;

• If 
$$\operatorname{cov}(X) = C$$
 and  $Y = AX$ , then  $\operatorname{cov}(Y) = ACA^T$ ;

• If 
$$\operatorname{cov}(X) = C$$
 and  $Y = a^T X \in \mathbb{R}$ , then  $\operatorname{var}(Y) = a^T C a \ge 0$ ;

• If 
$$\operatorname{cov}(X) = C$$
 and  $Y = C^{-1/2}X$ , then  $\operatorname{cov}(Y) = I$ ;

#### Combining the 2-nd and the 4-th facts is called standardization
Scenario: observed RV Y, depends on unknown variable(s) X.
 Goal: given an observation Y = y, infer X.

- Scenario: observed RV Y, depends on unknown variable(s) X.
  Goal: given an observation Y = y, infer X.
- Two main philosophies:
  Frequentist: X = x is fixed, but unknown;
  Bayesian: X is a RV with pdf/pmf f<sub>X</sub>(x) (the prior) prior ⇔ knowledge about X

- Scenario: observed RV Y, depends on unknown variable(s) X.
  Goal: given an observation Y = y, infer X.
- Two main philosophies:
  Frequentist: X = x is fixed, but unknown;
  Bayesian: X is a RV with pdf/pmf f<sub>X</sub>(x) (the prior)
  prior ⇔ knowledge about X
- In both philosophies, a central object is  $f_{Y|X}(y|x)$ several names: likelihood function, observation model,...

- Scenario: observed RV Y, depends on unknown variable(s) X.
  Goal: given an observation Y = y, infer X.
- Two main philosophies:
  Frequentist: X = x is fixed, but unknown;
  Bayesian: X is a RV with pdf/pmf f<sub>X</sub>(x) (the prior)
  prior ⇔ knowledge about X
- In both philosophies, a central object is  $f_{Y|X}(y|x)$ several names: likelihood function, observation model,...
- This in **not** machine learning!  $f_{Y,X}(y,x)$  is assumed known.

- Scenario: observed RV Y, depends on unknown variable(s) X.
  Goal: given an observation Y = y, infer X.
- Two main philosophies:
  Frequentist: X = x is fixed, but unknown;
  Bayesian: X is a RV with pdf/pmf f<sub>X</sub>(x) (the prior)
  prior ⇔ knowledge about X
- In both philosophies, a central object is  $f_{Y|X}(y|x)$  several names: likelihood function, observation model,...
- This in **not** machine learning!  $f_{Y,X}(y,x)$  is assumed known.
- $\bullet\,$  In the Bayesian philosophy, all the knowledge about X is carried by

$$f_{X|Y}(x|y) = \frac{f_{Y|X}(y|x) f_X(x)}{f_Y(y)} = \frac{f_{Y,X}(y,x)}{f_Y(y)}$$

...the posterior (or a posteriori) pdf/pmf.

• The posterior pdf/pmf  $f_{X|Y}(x|y)$  has all the information/knowledge about X, given Y = y (conditionality principle).

- The posterior pdf/pmf  $f_{X|Y}(x|y)$  has all the information/knowledge about X, given Y = y (conditionality principle).
- How to make an optimal "guess"  $\hat{x}$  about X, given this information?

- The posterior pdf/pmf  $f_{X|Y}(x|y)$  has all the information/knowledge about X, given Y = y (conditionality principle).
- How to make an optimal "guess"  $\widehat{x}$  about X, given this information?
- Need to define "optimal": loss function:  $L(\hat{x}, x) \ge 0$  measures "loss" / "cost" incurred by "guessing"  $\hat{x}$  if truth is x.

- The posterior pdf/pmf  $f_{X|Y}(x|y)$  has all the information/knowledge about X, given Y = y (conditionality principle).
- How to make an optimal "guess"  $\hat{x}$  about X, given this information?
- Need to define "optimal": loss function:  $L(\hat{x}, x) \ge 0$  measures "loss" / "cost" incurred by "guessing"  $\hat{x}$  if truth is x.
- The optimal Bayesian decision minimizes the expected loss:

$$\widehat{x}_{\mathsf{Bayes}} = \arg\min_{\widehat{x}} \mathbb{E}[L(\widehat{x}, X) | Y = y]$$

where

$$\mathbb{E}[L(\widehat{x}, X)|Y = y] = \begin{cases} \int L(\widehat{x}, x) f_{X|Y}(x|y) dx, & \text{continuous (estimation)} \\ \sum_{x} L(\widehat{x}, x) f_{X|Y}(x|y), & \text{discrete (classification)} \end{cases}$$

• Consider that  $X \in \{1, ..., K\}$  (discrete/classification problem).

- Consider that  $X \in \{1, ..., K\}$  (discrete/classification problem).
- Adopt the 0/1 loss:  $L(\widehat{x}, x) = 0$ , if  $\widehat{x} = x$ , and 1 otherwise.

- Consider that  $X \in \{1, ..., K\}$  (discrete/classification problem).
- Adopt the 0/1 loss:  $L(\widehat{x}, x) = 0$ , if  $\widehat{x} = x$ , and 1 otherwise.
- Optimal Bayesian decision:

$$\widehat{x}_{\mathsf{Bayes}} = \arg\min_{\widehat{x}} \sum_{x=1}^{K} L(\widehat{x}, x) f_{X|Y}(x|y)$$
$$= \arg\min_{\widehat{x}} \left(1 - f_{X|Y}(\widehat{x}|y)\right)$$
$$= \arg\max_{\widehat{x}} f_{X|Y}(\widehat{x}|y) \equiv \widehat{x}_{\mathsf{MAP}}$$

MAP = maximum a posteriori criterion.

- Consider that  $X \in \{1, ..., K\}$  (discrete/classification problem).
- Adopt the 0/1 loss:  $L(\widehat{x}, x) = 0$ , if  $\widehat{x} = x$ , and 1 otherwise.
- Optimal Bayesian decision:

$$\begin{aligned} \widehat{x}_{\mathsf{Bayes}} &= \arg\min_{\widehat{x}} \sum_{x=1}^{K} L(\widehat{x}, x) f_{X|Y}(x|y) \\ &= \arg\min_{\widehat{x}} \left( 1 - f_{X|Y}(\widehat{x}|y) \right) \\ &= \arg\max_{\widehat{x}} f_{X|Y}(\widehat{x}|y) \equiv \widehat{x}_{\mathsf{MAP}} \end{aligned}$$

MAP = maximum a posteriori criterion.

• Same criterion can be derived for continuous  $\boldsymbol{X}$ 

LxMLS 2016: Probability Theory

• Consider the MAP criterion  $\widehat{x}_{MAP} = \arg \max_x f_{X|Y}(x|y)$ 

- Consider the MAP criterion  $\widehat{x}_{MAP} = \arg \max_{x} f_{X|Y}(x|y)$
- From Bayes law:

$$\widehat{x}_{\mathsf{MAP}} = \arg\max_{x} \frac{f_{Y|X}(y|x) f_X(x)}{f_Y(y)} = \arg\max_{x} f_{Y|X}(y|x) f_X(x)$$

...only need to know posterior  $f_{X|Y}(x|y)$  up to a normalization factor.

- Consider the MAP criterion  $\widehat{x}_{MAP} = \arg \max_{x} f_{X|Y}(x|y)$
- From Bayes law:

$$\widehat{x}_{\mathsf{MAP}} = \arg\max_{x} \frac{f_{Y|X}(y|x) f_X(x)}{f_Y(y)} = \arg\max_{x} f_{Y|X}(y|x) f_X(x)$$

...only need to know posterior  $f_{X|Y}(x|y)$  up to a normalization factor.

• Also common to write:  $\hat{x}_{MAP} = \arg \max_{x} \left( \log f_{Y|X}(y|x) + \log f_{X}(x) \right)$ 

- Consider the MAP criterion  $\widehat{x}_{MAP} = \arg \max_x f_{X|Y}(x|y)$
- From Bayes law:

$$\widehat{x}_{\mathsf{MAP}} = \arg\max_{x} \frac{f_{Y|X}(y|x) \ f_X(x)}{f_Y(y)} = \arg\max_{x} f_{Y|X}(y|x) \ f_X(x)$$

...only need to know posterior  $f_{X|Y}(x|y)$  up to a normalization factor.

- Also common to write:  $\hat{x}_{MAP} = \arg \max_{x} \left( \log f_{Y|X}(y|x) + \log f_{X}(x) \right)$
- If the prior if flat,  $f_X(x) = C$ , then,

$$\widehat{x}_{\mathsf{MAP}} = \arg\max_{x} f_{Y|X}(y|x) \equiv \widehat{x}_{\mathsf{ML}}$$

ML = maximum likelihood criterion.

• Observed n i.i.d. (independent identically distributed) Bernoulli RVs:  $Y = (Y_1, ..., Y_n)$ , with  $Y_i \in \{0, 1\}$ .

Common pmf  $f_{Y_i|X}(y|x) = x^y(1-x)^{1-y}$ , where  $x \in [0, 1]$ .

- Observed n i.i.d. (independent identically distributed) Bernoulli RVs:  $Y = (Y_1, ..., Y_n)$ , with  $Y_i \in \{0, 1\}$ . Common pmf  $f_{Y_i|X}(y|x) = x^y(1-x)^{1-y}$ , where  $x \in [0, 1]$ .
- Likelihood function:  $f_{Y|X}(y_1, ..., y_n|x) = \prod_{i=1}^{n} x^{y_i} (1-x)^{1-y_i}$ Log-likelihood function:

$$\log f_{Y|X}(y_1, ..., y_n | x) = n \log(1 - x) + \log \frac{x}{1 - x} \sum_{i=1}^n y_i$$

- Observed n i.i.d. (independent identically distributed) Bernoulli RVs:  $Y = (Y_1, ..., Y_n)$ , with  $Y_i \in \{0, 1\}$ . Common pmf  $f_{Y_i|X}(y|x) = x^y(1-x)^{1-y}$ , where  $x \in [0, 1]$ .
- Likelihood function:  $f_{Y|X}(y_1, ..., y_n|x) = \prod_{i=1}^n x^{y_i}(1-x)^{1-y_i}$ Log-likelihood function:

$$\log f_{Y|X}(y_1, ..., y_n | x) = n \log(1 - x) + \log \frac{x}{1 - x} \sum_{i=1}^n y_i$$

• Maximum likelihood:  $\widehat{x}_{ML} = \arg \max_{x} f_{Y|X}(y|x) = \frac{1}{n} \sum_{i=1}^{n} y_i$ 

- Observed *n* i.i.d. (independent identically distributed) Bernoulli RVs:  $Y = (Y_1, ..., Y_n)$ , with  $Y_i \in \{0, 1\}$ . Common pmf  $f_{Y_i|X}(y|x) = x^y(1-x)^{1-y}$ , where  $x \in [0, 1]$ .
- Likelihood function:  $f_{Y|X}(y_1, ..., y_n|x) = \prod_{i=1}^n x^{y_i}(1-x)^{1-y_i}$ Log-likelihood function:

$$\log f_{Y|X}(y_1, ..., y_n | x) = n \log(1 - x) + \log \frac{x}{1 - x} \sum_{i=1}^n y_i$$

- Maximum likelihood:  $\widehat{x}_{ML} = \arg \max_{x} f_{Y|X}(y|x) = \frac{1}{n} \sum_{i=1}^{n} y_i$
- Example: n = 10, observed y = (1, 1, 1, 0, 1, 0, 0, 1, 1, 1),  $\hat{x}_{ML} = 7/10$ .

• Observed n i.i.d. (independent identically distributed) Bernoulli RVs.

- Observed n i.i.d. (independent identically distributed) Bernoulli RVs.
- Likelihood:

$$f_{Y|X}(y_1, ..., y_n|x) = \prod_{i=1}^n x^{y_i} (1-x)^{1-y_i} = x^{\sum_i y_i} (1-x)^{n-\sum_i y_i}$$

- Observed *n* i.i.d. (independent identically distributed) Bernoulli RVs.
- Likelihood:

$$f_{Y|X}(y_1, ..., y_n|x) = \prod_{i=1}^n x^{y_i} (1-x)^{1-y_i} = x^{\sum_i y_i} (1-x)^{n-\sum_i y_i}$$

• How to express knowledge that (e.g.) X is around 1/2? Convenient choice: conjugate prior. Form of the posterior = form of the prior.



- Observed *n* i.i.d. (independent identically distributed) Bernoulli RVs.
- Likelihood:

$$f_{Y|X}(y_1, ..., y_n|x) = \prod_{i=1}^n x^{y_i} (1-x)^{1-y_i} = x^{\sum_i y_i} (1-x)^{n-\sum_i y_i}$$

- How to express knowledge that (e.g.) X is around 1/2? Convenient choice: conjugate prior. Form of the posterior = form of the prior.
- ► In our case, the Beta pdf  $f_X(x) \propto x^{\alpha-1}(1-x)^{\beta-1}, \ \alpha, \beta > 0$



- Observed *n* i.i.d. (independent identically distributed) Bernoulli RVs.
- Likelihood:

$$f_{Y|X}(y_1, ..., y_n|x) = \prod_{i=1}^n x^{y_i} (1-x)^{1-y_i} = x^{\sum_i y_i} (1-x)^{n-\sum_i y_i}$$

- How to express knowledge that (e.g.) X is around 1/2? Convenient choice: conjugate prior. Form of the posterior = form of the prior.
- ▶ In our case, the Beta pdf  $f_X(x) \propto x^{\alpha-1}(1-x)^{\beta-1}, \ \alpha, \beta > 0$
- Posterior:  $f_{X|Y}(x|y) = x^{\alpha-1+\sum_i y_i} (1-x)^{\beta-1+n-\sum_i y_i}$



- Observed *n* i.i.d. (independent identically distributed) Bernoulli RVs.
- Likelihood:

$$f_{Y|X}(y_1, ..., y_n|x) = \prod_{i=1}^n x^{y_i} (1-x)^{1-y_i} = x^{\sum_i y_i} (1-x)^{n-\sum_i y_i}$$

- How to express knowledge that (e.g.) X is around 1/2? Convenient choice: conjugate prior. Form of the posterior = form of the prior.
- ▶ In our case, the Beta pdf  $f_X(x) \propto x^{\alpha-1}(1-x)^{\beta-1}, \ \alpha, \beta > 0$
- Posterior:  $f_{X|Y}(x|y) = x^{\alpha-1+\sum_i y_i} (1-x)^{\beta-1+n-\sum_i y_i}$
- MAP:  $\widehat{x}_{MAP} = \frac{\alpha + \sum_i y_i 1}{\alpha + \beta + n 2}$



- Observed *n* i.i.d. (independent identically distributed) Bernoulli RVs.
- Likelihood:

$$f_{Y|X}(y_1, ..., y_n|x) = \prod_{i=1}^n x^{y_i} (1-x)^{1-y_i} = x^{\sum_i y_i} (1-x)^{n-\sum_i y_i}$$

- How to express knowledge that (e.g.) X is around 1/2? Convenient choice: conjugate prior. Form of the posterior = form of the prior.
- ▶ In our case, the Beta pdf  $f_X(x) \propto x^{\alpha-1}(1-x)^{\beta-1}, \ \alpha, \beta > 0$

• Posterior: 
$$f_{X|Y}(x|y) = x^{\alpha-1+\sum_i y_i} (1-x)^{\beta-1+n-\sum_i y_i}$$

- MAP:  $\widehat{x}_{MAP} = \frac{\alpha + \sum_{i} y_i 1}{\alpha + \beta + n 2}$
- Example:  $\alpha = 4$ ,  $\beta = 4$ , n = 10, y = (1, 1, 1, 0, 1, 0, 0, 1, 1, 1),

 $\widehat{x}_{\mathsf{MAP}} = 0.625 \; (\mathsf{recall}\; \widehat{x}_{\mathsf{ML}} = 0.7)$ 



• Consider that  $X \in \mathbb{R}$  (continuous/estimation problem).

- Consider that  $X \in \mathbb{R}$  (continuous/estimation problem).
- Adopt the squared error loss:  $L(\widehat{x},x)=(\widehat{x}-x)^2$

- Consider that  $X \in \mathbb{R}$  (continuous/estimation problem).
- Adopt the squared error loss:  $L(\widehat{x},x)=(\widehat{x}-x)^2$
- Optimal Bayesian decision:

$$\widehat{x}_{\mathsf{Bayes}} = \arg\min_{\widehat{x}} \mathbb{E}[(\widehat{x} - X)^2 | Y = y]$$
$$= \arg\min_{\widehat{x}} \ \widehat{x}^2 - 2\,\widehat{x}\,\mathbb{E}[X|Y = y]$$
$$= \mathbb{E}[X|Y = y] \ \equiv \ \widehat{x}_{\mathsf{MMSE}}$$

MMSE = minimum mean squared error criterion.

- Consider that  $X \in \mathbb{R}$  (continuous/estimation problem).
- Adopt the squared error loss:  $L(\widehat{x},x)=(\widehat{x}-x)^2$
- Optimal Bayesian decision:

$$\begin{aligned} \widehat{x}_{\mathsf{Bayes}} &= \arg\min_{\widehat{x}} \mathbb{E}[(\widehat{x} - X)^2 | Y = y] \\ &= \arg\min_{\widehat{x}} \ \widehat{x}^2 - 2 \, \widehat{x} \, \mathbb{E}[X | Y = y] \\ &= \mathbb{E}[X | Y = y] \ \equiv \ \widehat{x}_{\mathsf{MMSE}} \end{aligned}$$

MMSE = minimum mean squared error criterion.

• Does not apply to classification problems.

• Observed n i.i.d. (independent identically distributed) Bernoulli RVs.

- Observed n i.i.d. (independent identically distributed) Bernoulli RVs.
- Likelihood:

$$f_{Y|X}(y_1, ..., y_n|x) = \prod_{i=1}^n x^{y_i} (1-x)^{1-y_i} = x^{\sum_i y_i} (1-x)^{n-\sum_i y_i}$$

- Observed n i.i.d. (independent identically distributed) Bernoulli RVs.
- Likelihood:

$$f_{Y|X}(y_1, ..., y_n|x) = \prod_{i=1}^n x^{y_i} (1-x)^{1-y_i} = x^{\sum_i y_i} (1-x)^{n-\sum_i y_i}$$

► In our case, the Beta pdf  $f_X(x) \propto x^{\alpha-1}(1-x)^{\beta-1}, \ \alpha, \beta > 0$ 



- Observed *n* i.i.d. (independent identically distributed) Bernoulli RVs.
- Likelihood:

$$f_{Y|X}(y_1, ..., y_n|x) = \prod_{i=1}^n x^{y_i} (1-x)^{1-y_i} = x^{\sum_i y_i} (1-x)^{n-\sum_i y_i}$$

- ► In our case, the Beta pdf  $f_X(x) \propto x^{\alpha-1}(1-x)^{\beta-1}, \ \alpha, \beta > 0$
- ► Posterior:  $f_{X|Y}(x|y) = x^{\alpha-1+\sum_i y_i} (1-x)^{\beta-1+n-\sum_i y_i}$


#### Back to the Bernoulli Example

- Observed n i.i.d. (independent identically distributed) Bernoulli RVs.
- Likelihood:

$$f_{Y|X}(y_1, ..., y_n|x) = \prod_{i=1}^n x^{y_i} (1-x)^{1-y_i} = x^{\sum_i y_i} (1-x)^{n-\sum_i y_i}$$

- In our case, the Beta pdf  $f_X(x) \propto x^{\alpha-1}(1-x)^{\beta-1}, \ \alpha, \beta > 0$
- Posterior:  $f_{X|Y}(x|y) = x^{\alpha-1+\sum_i y_i} (1-x)^{\beta-1+n-\sum_i y_i}$

• MMSE: 
$$\hat{x}_{\text{MMSE}} = \frac{\alpha + \sum_{i} y_i}{\alpha + \beta + n}$$



### Back to the Bernoulli Example

- Observed n i.i.d. (independent identically distributed) Bernoulli RVs.
- Likelihood:

$$f_{Y|X}(y_1, ..., y_n|x) = \prod_{i=1}^n x^{y_i} (1-x)^{1-y_i} = x^{\sum_i y_i} (1-x)^{n-\sum_i y_i}$$

- ► In our case, the Beta pdf  $f_X(x) \propto x^{\alpha-1}(1-x)^{\beta-1}, \ \alpha, \beta > 0$
- Posterior:  $f_{X|Y}(x|y) = x^{\alpha-1+\sum_i y_i} (1-x)^{\beta-1+n-\sum_i y_i}$

• MMSE: 
$$\widehat{x}_{MMSE} = \frac{\alpha + \sum_{i} y_{i}}{\alpha + \beta + n}$$

• Example:  $\alpha = 4$ ,  $\beta = 4$ , n = 10, y = (1, 1, 1, 0, 1, 0, 0, 1, 1, 1),



 $\widehat{x}_{\text{MMSE}} \simeq 0.611 \text{ (recall that } \widehat{x}_{\text{MAP}} = 0.625, \ \widehat{x}_{\text{ML}} = 0.7 \text{)}$ 

# Back to the Bernoulli Example

- Observed n i.i.d. (independent identically distributed) Bernoulli RVs.
- Likelihood:

$$f_{Y|X}(y_1, ..., y_n|x) = \prod_{i=1}^n x^{y_i} (1-x)^{1-y_i} = x^{\sum_i y_i} (1-x)^{n-\sum_i y_i}$$

- ▶ In our case, the Beta pdf  $f_X(x) \propto x^{\alpha-1}(1-x)^{\beta-1}, \ \alpha, \beta > 0$
- Posterior:  $f_{X|Y}(x|y) = x^{\alpha-1+\sum_i y_i} (1-x)^{\beta-1+n-\sum_i y_i}$

• MMSE: 
$$\widehat{x}_{MMSE} = \frac{\alpha + \sum_{i} y_{i}}{\alpha + \beta + n}$$

• Example:  $\alpha = 4$ ,  $\beta = 4$ , n = 10, y = (1, 1, 1, 0, 1, 0, 0, 1, 1, 1),



 $\widehat{x}_{\text{MMSE}} \simeq 0.611 \text{ (recall that } \widehat{x}_{\text{MAP}} = 0.625, \ \widehat{x}_{\text{ML}} = 0.7 \text{)}$ 

• Conjugate prior equivalent to "virtual" counts; often called *smoothing* in NLP and ML.

Mário A. T. Figueiredo (IST & IT)

LxMLS 2016: Probability Theory

## The Bernstein-Von Mises Theorem

• In the previous example, we had n = 10, y = (1, 1, 1, 0, 1, 0, 0, 1, 1, 1), thus  $\sum_i y_i = 7.$ With a Beta prior with  $\alpha = 4$  and  $\beta = 4$ , we had

$$\widehat{x}_{\mathsf{ML}} = 0.7, \quad \widehat{x}_{\mathsf{MAP}} = \frac{3 + \sum_{i} y_{i}}{6 + n} = 0.625, \quad \widehat{x}_{\mathsf{MMSE}} = \frac{4 + \sum_{i} y_{i}}{8 + n} \simeq 0.612$$

### The Bernstein-Von Mises Theorem

• In the previous example, we had n = 10, y = (1, 1, 1, 0, 1, 0, 0, 1, 1, 1), thus  $\sum_i y_i = 7$ . With a Beta prior with  $\alpha = 4$  and  $\beta = 4$ , we had

$$\widehat{x}_{\mathsf{ML}} = 0.7, \quad \widehat{x}_{\mathsf{MAP}} = \frac{3 + \sum_{i} y_{i}}{6 + n} = 0.625, \quad \widehat{x}_{\mathsf{MMSE}} = \frac{4 + \sum_{i} y_{i}}{8 + n} \simeq 0.61$$

• Consider n=100, and  $\sum_i y_i=70$ , with the same Beta(4,4) prior

$$\widehat{x}_{\mathsf{ML}} = 0.7, \quad \widehat{x}_{\mathsf{MAP}} = \frac{73}{106} \simeq 0.689, \quad \widehat{x}_{\mathsf{MMSE}} = \frac{74}{108} \simeq 0.685$$

... both Bayesian estimates are much closer to the ML.

# The Bernstein-Von Mises Theorem

• In the previous example, we had n = 10, y = (1, 1, 1, 0, 1, 0, 0, 1, 1, 1), thus  $\sum_i y_i = 7$ . With a Beta prior with  $\alpha = 4$  and  $\beta = 4$ , we had

$$\widehat{x}_{\mathsf{ML}} = 0.7, \quad \widehat{x}_{\mathsf{MAP}} = \frac{3 + \sum_{i} y_{i}}{6 + n} = 0.625, \quad \widehat{x}_{\mathsf{MMSE}} = \frac{4 + \sum_{i} y_{i}}{8 + n} \simeq 0.61$$

• Consider n=100, and  $\sum_i y_i=70$ , with the same Beta(4,4) prior

$$\widehat{x}_{\mathsf{ML}} = 0.7, \quad \widehat{x}_{\mathsf{MAP}} = \frac{73}{106} \simeq 0.689, \quad \widehat{x}_{\mathsf{MMSE}} = \frac{74}{108} \simeq 0.685$$

... both Bayesian estimates are much closer to the ML.

• This illustrates an important result in Bayesian inference: the Bernstein-Von Mises theorem; under (mild) conditions,

$$\lim_{n \to \infty} \widehat{x}_{\mathsf{MAP}} = \lim_{n \to \infty} \widehat{x}_{\mathsf{MMSE}} = \widehat{x}_{\mathsf{ML}}$$

message: if you have a lot of data, priors don't matter much.

• Cauchy-Schwartz's inequality for RVs:

 $\mathbb{E}(|XY|) \leq \sqrt{\mathbb{E}(X^2)\mathbb{E}(Y^2)}$ 

• Cauchy-Schwartz's inequality for RVs:

$$\mathbb{E}(|XY|) \le \sqrt{\mathbb{E}(X^2)\mathbb{E}(Y^2)}$$

• Recall that a real function g is convex if, for any x, y, and  $\alpha \in [0, 1]$ 

$$g(\alpha x + (1 - \alpha)y) \le \alpha g(x) + (1 - \alpha)g(y)$$



• Cauchy-Schwartz's inequality for RVs:

$$\mathbb{E}(|XY|) \le \sqrt{\mathbb{E}(X^2)\mathbb{E}(Y^2)}$$

• Recall that a real function g is convex if, for any x, y, and  $\alpha \in [0, 1]$ 



Jensen's inequality: if g is a real convex function, then  $\mathbb{E}(g(X)) \geq g(\mathbb{E}(X))$ 

• Cauchy-Schwartz's inequality for RVs:

$$\mathbb{E}(|XY|) \le \sqrt{\mathbb{E}(X^2)\mathbb{E}(Y^2)}$$

• Recall that a real function g is convex if, for any x,y, and  $\alpha\in[0,1]$ 



Jensen's inequality: if g is a real convex function, then  $\mathbb{E}(g(X)) \geq g(\mathbb{E}(X))$ 

 $\begin{array}{ll} \mbox{Examples: } \mathbb{E}(X)^2 \leq \mathbb{E}(X^2) \ \Rightarrow \ \mbox{var}(X) = \mathbb{E}(X^2) - \mathbb{E}(X)^2 \geq 0. \\ \mathbb{E}(\log X) \leq \log \mathbb{E}(X), \ \ \mbox{for } X \ \mbox{a positive RV}. \end{array}$ 

Mário A. T. Figueiredo (IST & IT)

Entropy of a discrete RV 
$$X \in \{1, ..., K\}$$
:  $H(X) = -\sum_{x=1}^{K} f_X(x) \log f_X(x)$ 

Entropy of a discrete RV  $X \in \{1, ..., K\}$ :  $H(X) = -\sum_{x=1}^{K} f_X(x) \log f_X(x)$ 

• Positivity:  $H(X) \ge 0$ ;  $H(X) = 0 \iff f_X(i) = 1$ , for exactly one  $i \in \{1, ..., K\}$ .

Entropy of a discrete RV  $X \in \{1, ..., K\}$ :  $H(X) = -\sum_{x=1}^{K} f_X(x) \log f_X(x)$ 

• Positivity:  $H(X) \ge 0$ ;  $H(X) = 0 \iff f_X(i) = 1$ , for exactly one  $i \in \{1, ..., K\}$ .

#### • Upper bound: $H(X) \le \log K$ ; $H(X) = \log K \Leftrightarrow f_X(x) = 1/k$ , for all $x \in \{1, ..., K\}$

Entropy of a discrete RV  $X \in \{1, ..., K\}$ :  $H(X) = -\sum_{x=1}^{K} f_X(x) \log f_X(x)$ 

• Positivity:  $H(X) \ge 0$ ;  $H(X) = 0 \iff f_X(i) = 1$ , for exactly one  $i \in \{1, ..., K\}$ .

#### • Upper bound: $H(X) \le \log K$ ; $H(X) = \log K \Leftrightarrow f_X(x) = 1/k$ , for all $x \in \{1, ..., K\}$

• Measure of uncertainty/randomness of X

Entropy of a discrete RV  $X \in \{1, ..., K\}$ :  $H(X) = -\sum_{x=1}^{K} f_X(x) \log f_X(x)$ 

• Positivity:  $H(X) \ge 0$ ;  $H(X) = 0 \iff f_X(i) = 1$ , for exactly one  $i \in \{1, ..., K\}$ .

#### • Upper bound: $H(X) \le \log K$ ; $H(X) = \log K \Leftrightarrow f_X(x) = 1/k$ , for all $x \in \{1, ..., K\}$

• Measure of uncertainty/randomness of X

Continuous RV X, differential entropy:  $h(X) = -\int f_X(x) \log f_X(x) dx$ 

Entropy of a discrete RV  $X \in \{1, ..., K\}$ :  $H(X) = -\sum_{x=1}^{K} f_X(x) \log f_X(x)$ 

• Positivity:  $H(X) \ge 0$ ;  $H(X) = 0 \Leftrightarrow f_X(i) = 1$ , for exactly one  $i \in \{1, ..., K\}$ .

#### • Upper bound: $H(X) \le \log K$ ; $H(X) = \log K \Leftrightarrow f_X(x) = 1/k$ , for all $x \in \{1, ..., K\}$

• Measure of uncertainty/randomness of X

Continuous RV X, differential entropy:  $h(X) = -\int f_X(x) \log f_X(x) dx$ 

• h(X) can be positive or negative. Example, if  $f_X(x) = \text{Uniform}(x; a, b), \ h(X) = \log(b - a).$ 

Entropy of a discrete RV  $X \in \{1, ..., K\}$ :  $H(X) = -\sum_{x=1}^{K} f_X(x) \log f_X(x)$ 

• Positivity:  $H(X) \ge 0$ ;  $H(X) = 0 \Leftrightarrow f_X(i) = 1$ , for exactly one  $i \in \{1, ..., K\}$ .

#### • Upper bound: $H(X) \le \log K$ ; $H(X) = \log K \Leftrightarrow f_X(x) = 1/k$ , for all $x \in \{1, ..., K\}$

• Measure of uncertainty/randomness of X

Continuous RV X, differential entropy:  $h(X) = -\int f_X(x) \log f_X(x) dx$ 

• h(X) can be positive or negative. Example, if  $f_X(x) = \text{Uniform}(x; a, b), \ h(X) = \log(b - a).$ 

• If  $f_X(x) = \mathcal{N}(x; \mu, \sigma^2)$ , then  $h(X) = \frac{1}{2}\log(2\pi e\sigma^2)$ .

Entropy of a discrete RV  $X \in \{1, ..., K\}$ :  $H(X) = -\sum_{x=1}^{K} f_X(x) \log f_X(x)$ 

• Positivity:  $H(X) \ge 0$ ;  $H(X) = 0 \iff f_X(i) = 1$ , for exactly one  $i \in \{1, ..., K\}$ .

#### • Upper bound: $H(X) \le \log K$ ; $H(X) = \log K \Leftrightarrow f_X(x) = 1/k$ , for all $x \in \{1, ..., K\}$

• Measure of uncertainty/randomness of X

Continuous RV X, differential entropy:  $h(X) = -\int f_X(x) \log f_X(x) dx$ 

• h(X) can be positive or negative. Example, if  $f_X(x) = \text{Uniform}(x; a, b), \ h(X) = \log(b - a).$ 

• If  $f_X(x) = \mathcal{N}(x; \mu, \sigma^2)$ , then  $h(X) = \frac{1}{2}\log(2\pi e\sigma^2)$ .

• If 
$$\operatorname{var}(Y) = \sigma^2$$
, then  $h(Y) \leq \frac{1}{2}\log(2\pi e\sigma^2)$ 

Kullback-Leibler divergence (KLD) between two pmf:

$$D(f_X || g_X) = \sum_{x=1}^{K} f_X(x) \log \frac{f_X(x)}{g_X(x)}$$

Kullback-Leibler divergence (KLD) between two pmf:

$$D(f_X || g_X) = \sum_{x=1}^{K} f_X(x) \log \frac{f_X(x)}{g_X(x)}$$

Positivity:  $D(f_X || g_X) \ge 0$  $D(f_X || g_X) = 0 \iff f_X(x) = g_X(x), \text{ for } x \in \{1, ..., K\}$ 

Kullback-Leibler divergence (KLD) between two pmf:

$$D(f_X || g_X) = \sum_{x=1}^{K} f_X(x) \log \frac{f_X(x)}{g_X(x)}$$

Positivity:  $D(f_X || g_X) \ge 0$  $D(f_X || g_X) = 0 \iff f_X(x) = g_X(x), \text{ for } x \in \{1, ..., K\}$ 

KLD between two pdf:

$$D(f_X || g_X) = \int f_X(x) \log \frac{f_X(x)}{g_X(x)} dx$$

Kullback-Leibler divergence (KLD) between two pmf:

$$D(f_X || g_X) = \sum_{x=1}^{K} f_X(x) \log \frac{f_X(x)}{g_X(x)}$$

Positivity:  $D(f_X || g_X) \ge 0$  $D(f_X || g_X) = 0 \iff f_X(x) = g_X(x), \text{ for } x \in \{1, ..., K\}$ 

KLD between two pdf:

$$D(f_X || g_X) = \int f_X(x) \log \frac{f_X(x)}{g_X(x)} dx$$

Positivity: 
$$D(f_X || g_X) \ge 0$$
  
 $D(f_X || g_X) = 0 \iff f_X(x) = g_X(x)$ , almost everywhere

## Mutual information

#### Mutual information (MI) between two random variables:

$$I(X;Y) = D(f_{X,Y}||f_X f_Y)$$

#### Mutual information (MI) between two random variables:

$$I(X;Y) = D(f_{X,Y} || f_X f_Y)$$

Positivity: 
$$I(X;Y) \ge 0$$
  
 $I(X;Y) = 0 \iff X, Y$  are independent.

#### Mutual information (MI) between two random variables:

$$I(X;Y) = D(f_{X,Y}||f_X f_Y)$$

Positivity: 
$$I(X;Y) \ge 0$$
  
 $I(X;Y) = 0 \iff X, Y$  are independent.

MI is a measure of dependency between two random variables

# Other Stuff

Note covered, but also very important for machine learning:

- Exponential families,
- Basic inequalities (Markov, Chebyshev, etc...)
- Stochastic processes (Markov chains, hidden Markov models,...)

# Recommended Reading (Probability and Statistics)

- K. Murphy, "Machine Learning: A Probabilistic Perspective", MIT Press, 2012 (Chapter 2).
- L. Wasserman, "All of Statistics: A Concise Course in Statistical Inference", Springer, 2004.