

## Three Concepts: Information Lecture 2: Mathematical Preliminaries

#### Teemu Roos

#### Complex Systems Computation Group Department of Computer Science, University of Helsinki

#### Fall 2007



・ロト ・部 ト ・ヨト ・ヨト

1

SQR



#### Lecture 2: Mathematical Preliminaries



"I think you should be more explicit here in step two."

◆ロ > ◆母 > ◆臣 > ◆臣 >





- Functions
- Limits and Convergence
- Convexity



3

590

<ロト < 団ト < 巨ト < 巨ト</p>





- Functions
- Limits and Convergence
- Convexity
- Probability
  - Probability Space and Random Variables
  - Joint and Conditional Distributions
  - Expectation
  - Law of Large Numbers



< □ > < 同 >

э





- Functions
- Limits and Convergence
- Convexity
- 2 Probability
  - Probability Space and Random Variables
  - Joint and Conditional Distributions
  - Expectation
  - Law of Large Numbers
- Inequalities
  - Jensen's Inequality
  - Gibbs's Inequality



< A

Functions Limits and Convergence Convexity

# Calculus





G.W. Leibniz, 1646-1716

Isaac Newton, 1643-1727

◆ロ > ◆母 > ◆臣 > ◆臣 >



### Functions

Functions associate with each possible input value x a unique output value y. The set of possible inputs is called the **domain** (*"alphabet"*). The set of possible outputs is called the **codomain**, and the set of actual outcomes is called the **range**. (Usually we just use the term 'range' for both.)



Calculus Probability Inequalities Functions Limits and Convergence Convexity

#### Examples: Exponent



Exponent function exp :  $\mathbb{R} \to \mathbb{R}^+$ , exp  $k = e^k = \overbrace{e \times e \times \dots \times e}^k$ : multiplicative growth (nuclear reaction, "interest on interest", ...)

< ロ > < 同 > < 三 > < 三 >

SQ (P

Functions Limits and Convergence Convexity

#### Examples: Exponent



Exponent function exp :  $\mathbb{R} \to \mathbb{R}^+$ , exp  $k = e^k = e^k \times e \times \dots \times e^k$ : multiplicative growth (nuclear reaction, "interest on interest", ...)

$$\exp x \cdot \exp y = \exp(x+y)$$

< ロ > < 同 > < 三 > < 三 >

SQ (P

Functions Limits and Convergence Convexity

#### Examples: Exponent



Teemu Roos Three Conc

Calculus Probability nequalities Functions Limits and Convergence Convexity

#### Examples: Logarithm



Natural logarithm In :  $\mathbb{R}^+ \to \mathbb{R}$ , ln exp x = x: time to grow to x, number of digits (log<sub>10</sub>).

イロト イポト イヨト イヨト

1

SQR

Functions Limits and Convergence Convexity

## Examples: Logarithm



Natural logarithm In :  $\mathbb{R}^+ \to \mathbb{R}$ , ln exp x = x: time to grow to x, number of digits (log<sub>10</sub>).

General (base *a*) logarithm,  $\log_a a^x = x$ :  $\log_a x = \frac{1}{\ln a} \ln x$ 

・ロッ ・ 日 ・ ・ 日 ・ ・ 日 ・

Calculus Probability Inequalities Functions Limits and Convergence Convexity

## Examples: Logarithm



 $\ln xy = \ln x + \ln y$ 

◆ロ > ◆母 > ◆臣 > ◆臣 >

3

Calculus Probability Inequalities Functions Limits and Convergence Convexity

#### Examples: Logarithm



 $\ln xy = \ln x + \ln y \qquad \ln x^r = r \ln x$ 

◆ロ > ◆母 > ◆臣 > ◆臣 >

3

Calculus Probability nequalities Functions Limits and Convergence Convexity

#### Examples: Logarithm



 $\ln xy = \ln x + \ln y$   $\ln x^{r} = r \ln x$   $\ln \frac{1}{x} = -\ln x$ 

・ロト ・回ト ・モト ・モト

3

Calculus Calculus Probability Inequalities Functions Limits and Convergence Convexity

#### Examples: Logarithm



 $\ln xy = \ln x + \ln y$   $\ln x^{r} = r \ln x$   $\ln \frac{1}{x} = -\ln x$   $\ln \frac{x}{y} = \ln x - \ln y$ 

イロト イポト イヨト イヨト

nar

Calculus Probability nequalities Functions Limits and Convergence Convexity

### Examples: Logarithm



 $\ln xy = \ln x + \ln y \quad \ln x^{r} = r \ln x \quad \ln \frac{1}{x} = -\ln x \quad \ln \frac{x}{y} = \ln x - \ln y$  $\ln x \le x - 1 \text{ with equality if and only if } x = 1$ (NB: doesn't work with  $\log_{a} x$  if  $a \ne e$ )

イロト イボト イヨト イヨト

Calculus Calculus Probability Inequalities Functions Limits and Convergence Convexity

#### Examples: Logarithm



 $\ln xy = \ln x + \ln y \quad \ln x^{r} = r \ln x \quad \ln \frac{1}{x} = -\ln x \quad \ln \frac{x}{y} = \ln x - \ln y$  $\ln x \le x - 1 \text{ with equality if and only if } x = 1 \qquad \frac{d \ln x}{dx} = \frac{1}{x}$ (NB: doesn't work with  $\log_{a} x$  if  $a \ne e$ )  $\frac{d \ln x}{dx} = \frac{1}{x}$ 

イロト イボト イヨト イヨト

Functions Limits and Convergence Convexity

#### Limits and Convergence

• A sequence of values  $(x_i : i \in \mathbb{N})$  converges to limit L,  $\lim_{i\to\infty} x_i = L$ , iff for any  $\epsilon > 0$  there exists a number  $N \in \mathbb{N}$  such that

 $|x_i - L| < \epsilon$  for all  $i \ge N$  .

イロト イポト イヨト イヨト

SQR

Functions Limits and Convergence Convexity

#### Limits and Convergence

• A sequence of values  $(x_i : i \in \mathbb{N})$  converges to limit L,  $\lim_{i\to\infty} x_i = L$ , iff for any  $\epsilon > 0$  there exists a number  $N \in \mathbb{N}$  such that

$$|x_i - L| < \epsilon$$
 for all  $i \ge N$  .

• f(x) has a *limit* L as x approaches c,  $\lim_{x\to c} f(x) = L$ , (from above  $c^+$ /below  $c^-$ ) iff for any  $\epsilon > 0$  there exists a number  $\delta > 0$  such that

$$|f(x) - L| < \epsilon$$
 for all  $\begin{cases} c < x < c + \delta & \text{`above'} \\ c - \delta < x < c & \text{`below'} \\ 0 < |x - c| < \delta & -- \end{cases}$ 

イロト イポト イヨト イヨト

Calculus Probability Inequalities Functions Limits and Convergence Convexity

#### Example: Logarithm Again



Even though  $x \ln x$  is undefined at x = 0, we have (by l'Hôpital's rule):

$$\lim_{x\to 0^+} x \ln x = 0 \;\; .$$

イロト イポト イヨト イヨト

nar



#### Convexity

Function  $f : \mathcal{X} \to \mathbb{R}$  is said to be **convex** iff for any  $x, y \in \mathcal{X}$  and any  $t \in [0, 1]$ , we have



イロト イポト イヨト イヨト

SQR



#### Convexity

Function  $f : \mathcal{X} \to \mathbb{R}$  is said to be **convex** iff for any  $x, y \in \mathcal{X}$  and any  $t \in [0, 1]$ , we have



Function f is **strictly convex** iff the above inequality holds strictly ('<' instead of ' $\leq$ ').

イロト イポト イヨト イヨト



## Convexity

Function  $f : \mathcal{X} \to \mathbb{R}$  is said to be **convex** iff for any  $x, y \in \mathcal{X}$  and any  $t \in [0, 1]$ , we have



Function f is **strictly convex** iff the above inequality holds strictly ('<' instead of ' $\leq$ ').

Function f is (strictly) **concave** iff the above holds for -f.

- 「「「」」(「」)(「」)(「」)

Functions Limits and Convergence Convexity

### Convexity and Derivatives

#### Theorem

If function f has a second derivative f'', and f'' is non-negative  $(\geq 0)$  for all x, then f is convex. If f'' is positive (> 0) for all x, then f is *strictly* convex.

< ロ > < 同 > < 三 > < 三 >

Functions Limits and Convergence Convexity

### Convexity and Derivatives

#### Theorem

If function f has a second derivative f'', and f'' is non-negative  $(\geq 0)$  for all x, then f is convex. If f'' is positive (> 0) for all x, then f is *strictly* convex.



Example:  $f'(x) = \frac{d \exp x}{dx} = \exp x \Rightarrow f''(x) = \exp x > 0$ . Hence exp is strictly convex.

Functions Limits and Convergence Convexity

### Convexity and Derivatives

#### Theorem

If function f has a second derivative f'', and f'' is non-negative  $(\geq 0)$  for all x, then f is convex. If f'' is positive (> 0) for all x, then f is *strictly* convex.



Example:  $f'(x) = \frac{d \exp x}{dx} = \exp x \Rightarrow f''(x) = \exp x > 0$ . Hence exp is strictly convex.

Probability Space and Random Variables Joint and Conditional Distributions Expectation Law of Large Numbers

# **Probability**



#### A.N. Kolmogorov, 1903-1987

Teemu Roos Three Concepts: Information

・ロット 4 日マ 4 田マ 4

э

## **Probability Space**

A probability space  $(\Omega, \mathcal{F}, P)$  is defined by

◆ロ > ◆母 > ◆臣 > ◆臣 >

SQC

## **Probability Space**

A probability space  $(\Omega, \mathcal{F}, P)$  is defined by

• the sample space  $\Omega$  whose elements are called outcomes  $\omega$ ,

イロト イポト イヨト イヨト

SQR

## **Probability Space**

A probability space  $(\Omega, \mathcal{F}, P)$  is defined by

- the sample space  $\Omega$  whose elements are called outcomes  $\omega$ ,
- a sigma algebra *F* of subsets of Ω, whose elements are called events *E*, and

< ロ > < 同 > < 三 > < 三 >

## **Probability Space**

A probability space  $(\Omega, \mathcal{F}, P)$  is defined by

- the sample space  $\Omega$  whose elements are called outcomes  $\omega,$
- a sigma algebra  $\mathcal{F}$  of subsets of  $\Omega$ , whose elements are called **events** E, and
- a measure *P* which determines the **probabilities of events**,  $P : \mathcal{F} \rightarrow [0, 1].$

イロト イポト イヨト イヨト

## **Probability Space**

A probability space  $(\Omega, \mathcal{F}, P)$  is defined by

- the sample space  $\Omega$  whose elements are called outcomes  $\omega,$
- a sigma algebra  $\mathcal{F}$  of subsets of  $\Omega$ , whose elements are called **events** E, and
- a measure *P* which determines the **probabilities of events**, *P* :  $\mathcal{F} \rightarrow [0, 1]$ .

Measure *P* has to satisfy the **probability axioms**:  $P(E) \ge 0$  for all  $E \in \mathcal{F}$ ,  $P(\Omega) = 1$ , and  $P(E_1 \cup E_2 \cup ...) = \sum_i P(E_i)$  if  $(E_i)$  is a countable sequence of *disjoint* events.

イロト イポト イヨト イヨト

## **Probability Space**

A probability space  $(\Omega, \mathcal{F}, P)$  is defined by

- the sample space  $\Omega$  whose elements are called outcomes  $\omega,$
- a sigma algebra  $\mathcal{F}$  of subsets of  $\Omega$ , whose elements are called **events** E, and
- a measure *P* which determines the **probabilities of events**,  $P : \mathcal{F} \rightarrow [0, 1].$

Measure *P* has to satisfy the **probability axioms**:  $P(E) \ge 0$  for all  $E \in \mathcal{F}$ ,  $P(\Omega) = 1$ , and  $P(E_1 \cup E_2 \cup ...) = \sum_i P(E_i)$  if  $(E_i)$  is a countable sequence of *disjoint* events.

These axioms imply the usual rules of **probability calculus**, e.g.,  $P(A \cup B) = P(A) + P(B) - P(A \cap B)$ ,  $P(\Omega \setminus E) = 1 - P(E)$ , etc.

<ロト < 同ト < ヨト < ヨト -

Outline Probability Space and Random Variables Calculus Joint and Conditional Distributions Probability Expectation Inequalities Law of Large Numbers

## Venn Diagrams



Teemu Roos Three Concepts: Information

#### **Probability Calculus**

The conditional probability of event B given that event A occurs is defined as

$$P(B \mid A) = rac{P(A \cap B)}{P(A)}$$

for A such that P(A) > 0.

(日) (同) (三) (三)

SQR
# **Probability Calculus**

The conditional probability of event B given that event A occurs is defined as

$$P(B \mid A) = rac{P(A \cap B)}{P(A)}$$
 for A such that  $P(A) > 0$ .

$$P(A \cap B) = P(A) \times P(B \mid A) = P(B) \times P(A \mid B) .$$

(日) (同) (三) (三)

1

Outline Probability Space and Random Variables Calculus Joint and Conditional Distributions Probability Expectation Law of Large Numbers

# **Probability Calculus**

1

• The conditional probability of event B given that event A occurs is defined as

$$P(B \mid A) = \frac{P(A \cap B)}{P(A)} \quad \text{for } A \text{ such that } P(A) > 0.$$

$$P(A \cap B) = P(A) \times P(B \mid A) = P(B) \times P(A \mid B) \quad .$$

$$\text{Bayes' rule: } P(B \mid A) = \frac{P(A \mid B) \times P(B)}{P(A)} \quad .$$

イロト イポト イヨト イヨト

# Probability Calculus

- The conditional probability of event B given that event A occurs is defined as  $P(B \mid A) = rac{P(A \cap B)}{P(A)}$  for A such that P(A) > 0.  $P(A \cap B) = P(A) \times P(B \mid A) = P(B) \times P(A \mid B) .$ 3 Bayes' rule:  $P(B \mid A) = \frac{P(A \mid B) \times P(B)}{P(A)}$ . Chain rule: Ν  $P(\bigcap_{i=1}^{N} E_i) = \prod P(E_i \mid \bigcap_{i=1}^{i-1} E_j)$ i=1
  - $= P(E_1) \times P(E_2 \mid E_1) \times P(E_3 \mid E_1 \cap E_2) \times \dots \times P(E_N \mid E_1 \cap \dots \cap E_{N-1}) .$

イロト 不得 トイヨト イヨト 二日

SQ (P

#### **Random Variables**

Technically, a random variable is a (measurable) function  $X : \Omega \to \mathbb{R}$  from the sample space to the reals.

イロト イポト イヨト イヨト

-

#### **Random Variables**

 Technically, a random variable is a (measurable) function
 X :  $\Omega \to \mathbb{R}$  from the sample space to the reals.

The probability measure P on  $\Omega$  determines the distribution of X:

$$P_X(A) = \Pr[X \in A] = P(\{\omega : X(\omega) \in A\}) ,$$

where  $A \subseteq \mathbb{R}$ .

(日) (同) (三) (三)

## **Random Variables**

 Technically, a random variable is a (measurable) function
 X :  $\Omega \to \mathbb{R}$  from the sample space to the reals.

The probability measure P on  $\Omega$  determines the distribution of X:

$$P_X(A) = \Pr[X \in A] = P(\{\omega : X(\omega) \in A\}) ,$$

where  $A \subseteq \mathbb{R}$ .

In practice, we often forget about the underlying probability space  $\Omega$ , and just speak of random variable X and its distribution  $P_X$ .

(日) (同) (三) (三)

#### **Random Variables**

The distribution of a random variable can *always* be represented as a *cumulative distribution function* (cdf)  $F_X(x) = \Pr[X \le x]$ .

イロト イポト イヨト イヨト

# **Random Variables**

The distribution of a random variable can *always* be represented as a *cumulative distribution function* (cdf)  $F_X(x) = \Pr[X \le x]$ .

In addition:

A discrete random variable X with countable alphabet X has a probability mass function (pmf) p<sub>X</sub> such that Pr[X = x] = p<sub>X</sub>(x).

<ロト <同ト < 国ト < 国ト

# **Random Variables**

The distribution of a random variable can *always* be represented as a *cumulative distribution function* (cdf)  $F_X(x) = \Pr[X \le x]$ .

In addition:

- A discrete random variable X with countable alphabet  $\mathcal{X}$  has a probability mass function (pmf)  $p_X$  such that  $\Pr[X = x] = p_X(x)$ .
- A continuous random variable Y has a probability density function (pdf)  $f_Y$  such that  $Pr[Y \in A] = \int_A f_Y(x) dy$ .

# **Random Variables**

The distribution of a random variable can *always* be represented as a *cumulative distribution function* (cdf)  $F_X(x) = \Pr[X \le x]$ .

In addition:

- A discrete random variable X with countable alphabet  $\mathcal{X}$  has a probability mass function (pmf)  $p_X$  such that  $\Pr[X = x] = p_X(x)$ .
- A continuous random variable Y has a probability density function (pdf)  $f_Y$  such that  $Pr[Y \in A] = \int_A f_Y(x) dy$ .

There are also *mixed* random variables that are neither discrete nor continuous. They don't have a pmf or pdf, but they do have a cdf.

# **Random Variables**

The distribution of a random variable can *always* be represented as a *cumulative distribution function* (cdf)  $F_X(x) = \Pr[X \le x]$ .

In addition:

- A discrete random variable X with countable alphabet X has a probability mass function (pmf) p<sub>X</sub> such that Pr[X = x] = p<sub>X</sub>(x).
- A continuous random variable Y has a probability density function (pdf)  $f_Y$  such that  $Pr[Y \in A] = \int_A f_Y(x) dy$ .

There are also *mixed* random variables that are neither discrete nor continuous. They don't have a pmf or pdf, but they do have a cdf.

We often omit the subscripts  $X, Y, \ldots$  and write p(x), f(y), etc.

・ロト ・ 同ト ・ ヨト ・ ヨト

-

#### **Random Variables**

Since random variables are functions, we can define more random variables as functions of random variables: if f is a function, and X and Y are r.v.'s, then  $f(X) : \Omega \to \mathbb{R}$  is a r.v., X + Y is a r.v., etc.

< ロ > < 同 > < 三 > < 三 >

## **Random Variables**

Since random variables are functions, we can define more random variables as functions of random variables: if f is a function, and X and Y are r.v.'s, then  $f(X) : \Omega \to \mathbb{R}$  is a r.v., X + Y is a r.v., etc.

Example: Let r.v. X be the outcome of a die. The pmf of X is given by  $p_X(x) = 1/6$  for all  $x \in \{1, 2, 3, 4, 5, 6\}$ .

<ロト <同ト < ヨト < ヨト

## **Random Variables**

Since random variables are functions, we can define more random variables as functions of random variables: if f is a function, and X and Y are r.v.'s, then  $f(X) : \Omega \to \mathbb{R}$  is a r.v., X + Y is a r.v., etc.



Example: Let r.v. X be the outcome of a die.

- The pmf of X is given by  $p_X(x) = 1/6$  for all  $x \in \{1, 2, 3, 4, 5, 6\}.$
- The pmf of r.v.  $X^2$  is given by  $p_{X^2}(x) = 1/6$  for all  $x \in \{1, 4, 9, 16, 25, 36\}.$

# **Random Variables**

Since random variables are functions, we can define more random variables as functions of random variables: if f is a function, and X and Y are r.v.'s, then  $f(X) : \Omega \to \mathbb{R}$  is a r.v., X + Y is a r.v., etc.



- The pmf of X is given by  $p_X(x) = 1/6$  for all  $x \in \{1, 2, 3, 4, 5, 6\}$ .
- The pmf of r.v.  $X^2$  is given by  $p_{X^2}(x) = 1/6$  for all  $x \in \{1, 4, 9, 16, 25, 36\}.$

In particular, a pmf  $p_X$  is a function, and hence,  $p_X(X)$  is also a random variable. Further,  $p_X^2(X)$ ,  $\ln p_X(X)$ , etc. are random variables.

Outline Probability Space and Random Variables Calculus Joint and Conditional Distributions Probability Expectation Inequalities Law of Large Numbers

#### Multivariate Distributions

The probabilistic behavior of two or more random variables is described by multivariate distributions.

The **joint distribution** of r.v.'s X and Y is

$$egin{aligned} & P_{X,Y}(A,B) = \Pr[X \in A \land Y \in B] \ &= P(\{\omega \ : \ X(\omega) \in A, Y(\omega) \in B\}) \end{aligned}$$

イロト イポト イヨト イヨト

# Multivariate Distributions

The probabilistic behavior of two or more random variables is described by multivariate distributions.

The **joint distribution** of r.v.'s X and Y is

$$egin{aligned} & P_{X,Y}(A,B) = \Pr[X \in A \ \land \ Y \in B] \ &= P(\{\omega \ : \ X(\omega) \in A, Y(\omega) \in B\}) \end{aligned}$$

For each multivariate distribution  $P_{X,Y}$ , there are unique marginal distributions  $P_X$  and  $P_Y$  such that

$$P_X(A) = P_{X,Y}(A,\mathbb{R}), \qquad P_Y(B) = P_{X,Y}(\mathbb{R},B)$$

< ロ > < 同 > < 回 > < 回 > :

# Multivariate Distributions

The probabilistic behavior of two or more random variables is described by multivariate distributions.

The **joint distribution** of r.v.'s X and Y is

$$egin{aligned} & P_{X,Y}(A,B) = \Pr[X \in A \ \land \ Y \in B] \ &= P(\{\omega \ : \ X(\omega) \in A, Y(\omega) \in B\}) \end{aligned}$$

For each multivariate distribution  $P_{X,Y}$ , there are unique **marginal** distributions  $P_X$  and  $P_Y$  such that

$$P_X(A) = P_{X,Y}(A,\mathbb{R}), \qquad P_Y(B) = P_{X,Y}(\mathbb{R},B)$$

pmf: 
$$p_Y(y) = \sum_{x \in \mathcal{X}} p_{X,Y}(x,y)$$
 pdf:  $f_Y(y) = \int_{\mathbb{R}} f_{X,Y}(x,y) dx$ .

## Multivariate Distributions

The **conditional distribution** is defined similar to *conditional probability*:

$$P_{Y|X}(B \mid A) = rac{P_{X,Y}(A,B)}{P_X(A)}$$
 for A such that  $P_X(A) > 0$ .

3

# Multivariate Distributions

The **conditional distribution** is defined similar to *conditional probability*:

$$P_{Y\mid X}(B\mid A) = rac{P_{X,Y}(A,B)}{P_X(A)}$$
 for  $A$  such that  $P_X(A) > 0.$ 

For discrete/continuous variables we have:

• discrete r.v.'s:

$$p_{Y|X}(y \mid x) = rac{p_{X,Y}(x,y)}{p_X(x)} , \quad p_X(x) > 0 ,$$

(日) (同) (三) (三)

-

# Multivariate Distributions

The **conditional distribution** is defined similar to *conditional probability*:

$$P_{Y\mid X}(B\mid A) = rac{P_{X,Y}(A,B)}{P_X(A)}$$
 for  $A$  such that  $P_X(A) > 0.$ 

For discrete/continuous variables we have:

• *discrete* r.v.'s:

$$p_{Y|X}(y \mid x) = rac{p_{X,Y}(x,y)}{p_X(x)} \ , \quad p_X(x) > 0 \ ,$$

• continuous r.v.'s:

< ロ > < 同 > < 三 > < 三 >

#### Independence

Variable X is said to be **independent** of variable  $Y(X \perp Y)$  iff

 $P_{X,Y}(A,B) = P_X(A) \times P_Y(B)$  for all  $A,B \subseteq \mathbb{R}$ .

1

#### Independence

Variable X is said to be **independent** of variable Y  $(X \perp Y)$  iff

$$\mathsf{P}_{X,Y}(A,B)=\mathsf{P}_X(A) imes\mathsf{P}_Y(B) \quad ext{for all } A,B\subseteq \mathbb{R}.$$

This is equivalent to

 $P_{X|Y}(A \mid B) = P_X(A)$  for all B such that P(B) > 0,

(日) (同) (三) (三)

#### Independence

and

Variable X is said to be **independent** of variable Y  $(X \perp Y)$  iff

$$P_{X,Y}(A,B)=P_X(A) imes P_Y(B) \quad ext{for all } A,B\subseteq \mathbb{R}.$$

This is equivalent to

 $P_{X|Y}(A \mid B) = P_X(A)$  for all B such that P(B) > 0,

 $P_{Y|X}(B \mid A) = P_Y(B)$  for all A such that P(A) > 0.

In words, knowledge about one variable tells nothing about the other. Note that independence is symmetric,  $X \perp Y \Leftrightarrow Y \perp X$ .

# Expectation

The **expectation** (or expected value, or mean) of a discrete random variable is given by

$$E[X] = \sum_{x \in \mathcal{X}} p(x) x$$
.

イロト イポト イヨト イヨト

# Expectation

The **expectation** (or expected value, or mean) of a discrete random variable is given by

$$E[X] = \sum_{x \in \mathcal{X}} p(x) x \; \; .$$

The expectation of a continuous random variable is given by

$$E[X] = \int_{\mathcal{X}} f(x) \, x \, dx \; \; .$$

# Expectation

The **expectation** (or expected value, or mean) of a discrete random variable is given by

$$E[X] = \sum_{x \in \mathcal{X}} p(x) x \; \; .$$

The expectation of a continuous random variable is given by

$$E[X] = \int_{\mathcal{X}} f(x) \, x \, dx \; \; .$$

In both cases, it is possible that  $E[X] = \pm \infty$ .

# Expectation

The **expectation** (or expected value, or mean) of a discrete random variable is given by

$$E[X] = \sum_{x \in \mathcal{X}} p(x) x \; \; .$$

The expectation of a continuous random variable is given by

$$E[X] = \int_{\mathcal{X}} f(x) \, x \, dx \; \; .$$

In both cases, it is possible that  $E[X] = \pm \infty$ .

 $E[kX] = kE[X] \qquad E[X+Y] = E[X] + E[Y]$ 

E[XY] = E[X]E[Y] if  $X \perp Y$ 

#### Law of Large Numbers

Let  $X_1, X_2, \ldots$  be a sequence of independent outcomes of a die, so that  $p_{X_i}(x) = 1/6$  for all  $i \in \mathbb{N}, x \in \{1, 2, 3, 4, 5, 6\}$ .



Image: A math a math

#### Law of Large Numbers

Let  $X_1, X_2, \ldots$  be a sequence of independent outcomes of a die, so that  $p_{X_i}(x) = 1/6$  for all  $i \in \mathbb{N}, x \in \{1, 2, 3, 4, 5, 6\}$ .



#### Law of Large Numbers

Let  $S_n = \sum_{i=1}^n X_n$  be the sum of the first *n* outcomes.

3

#### Law of Large Numbers

Let  $S_n = \sum_{i=1}^n X_n$  be the sum of the first *n* outcomes.

The distribution of  $S_n$  is given by

$$P_{S_n}(x) = \frac{\# \text{ of ways to get sum } x \text{ with } n \text{ dice}}{6^n}$$

(日) (同) (三) (三)

1

## Law of Large Numbers

Let  $S_n = \sum_{i=1}^n X_n$  be the sum of the first *n* outcomes.

The distribution of  $S_n$  is given by

$$P_{S_n}(x) = \frac{\# \text{ of ways to get sum } x \text{ with } n \text{ dice}}{6^n}$$



・ロト ・ 同ト ・ ヨト ・

э

# Law of Large Numbers

Let  $S_n = \sum_{i=1}^n X_n$  be the sum of the first *n* outcomes.

The distribution of  $S_n$  is given by

$$P_{S_n}(x) = \frac{\# \text{ of ways to get sum } x \text{ with } n \text{ dice}}{6^n}$$



1

# Law of Large Numbers

Let  $S_n = \sum_{i=1}^n X_n$  be the sum of the first *n* outcomes.

The distribution of  $S_n$  is given by

$$P_{S_n}(x) = \frac{\# \text{ of ways to get sum } x \text{ with } n \text{ dice}}{6^n}$$



# Law of Large Numbers

Let  $S_n = \sum_{i=1}^n X_n$  be the sum of the first *n* outcomes.

The distribution of  $S_n$  is given by

$$P_{S_n}(x) = \frac{\# \text{ of ways to get sum } x \text{ with } n \text{ dice}}{6^n}$$


## Law of Large Numbers

Let  $S_n = \sum_{i=1}^n X_n$  be the sum of the first *n* outcomes.

The distribution of  $S_n$  is given by

$$P_{S_n}(x) = \frac{\# \text{ of ways to get sum } x \text{ with } n \text{ dice}}{6^n}$$



1

## Law of Large Numbers

Let  $S_n = \sum_{i=1}^n X_n$  be the sum of the first *n* outcomes.

The distribution of  $S_n$  is given by

$$P_{S_n}(x) = \frac{\# \text{ of ways to get sum } x \text{ with } n \text{ dice}}{6^n}$$



1

## Law of Large Numbers

Let  $S_n = \sum_{i=1}^n X_n$  be the sum of the first *n* outcomes.

The distribution of  $S_n$  is given by

$$P_{S_n}(x) = \frac{\# \text{ of ways to get sum } x \text{ with } n \text{ dice}}{6^n}$$



э

## Law of Large Numbers

Let  $S_n = \sum_{i=1}^n X_n$  be the sum of the first *n* outcomes.

The distribution of  $S_n$  is given by

$$P_{S_n}(x) = \frac{\# \text{ of ways to get sum } x \text{ with } n \text{ dice}}{6^n}$$



Probability Space and Random Variables Joint and Conditional Distributions Expectation Law of Large Numbers

### Law of Large Numbers

### LAW OF LARGE NUMBERS IN AVERAGE OF DIE ROLLS

Outline

Calculus

Probability



3

## Law of Large Numbers

### Weak Law of Large Numbers

For a sequence of independent and identically distributed (i.i.d.) random variables with finite mean  $\mu$ , the average  $\frac{1}{n}S_n$  converges in probability to  $\mu$ :

$$\lim_{n\to\infty} \Pr\left[\left|\frac{S_n}{n}-\mu\right|<\epsilon\right] = 1 \quad \text{for all } \epsilon > 0.$$

We will use the LLN to prove a result known as the Asymptotic Equipartition Property (AEP), which is a central result in information theory (see next lecture).

<ロト <同ト < ヨト < ヨト

Jensen's Inequality Gibbs's Inequality

# Jensen's inequality



### J.L.W.V. Jensen, 1859-1925

Teemu Roos Three Concepts: Information

イロト イポト イヨト イヨト

SQR

Jensen's Inequality Gibbs's Inequality

### Inqualities: Jensen

Jensen's inequality

If f is a convex function and X is a random variable, then

 $E[f(X)] \ge f(E[X])$  .

Moreover, if f is strictly convex, the inequality holds as an equality if and only if X = E[X] with probability 1.

< ロ > < 同 > < 三 > < 三 >

Jensen's Inequality Gibbs's Inequality

### Inqualities: Jensen

Jensen's inequality

If f is a convex function and X is a random variable, then

 $E[f(X)] \ge f(E[X])$  .

Moreover, if f is strictly convex, the inequality holds as an equality if and only if X = E[X] with probability 1.

We give a proof for the first part of the theorem in the special case where X has a finite domain.

イロト イポト イヨト イヨト

Jensen's Inequality Gibbs's Inequality

### Inqualities: Jensen

Jensen's inequality

If f is a convex function and X is a random variable, then

 $E[f(X)] \ge f(E[X])$  .

Moreover, if f is strictly convex, the inequality holds as an equality if and only if X = E[X] with probability 1.

We give a proof for the first part of the theorem in the special case where X has a finite domain.

For two mass points, we have  $p(x_2) = 1 - p(x_1)$ , and the claim holds by definition of convexity:

$$p(x_1) f(x_1) + p(x_2) f(x_2) \ge f(p(x_1) x_1 + p(x_2) x_2)$$
.

イロト イポト イヨト イヨト

Jensen's Inequality Gibbs's Inequality

### Inequalities: Jensen

Induction: Assume that (\*) the theorem holds for N-1 mass points.

$$\sum_{i=1}^{N} p(x_i) f(x_i) = p(x_N) f(x_N) + (1 - p(x_N)) \sum_{i=1}^{N-1} p'(x_i) f(x_i)$$

$$\geq p(x_N) f(x_N) + (1 - p(x_N)) f\left(\sum_{i=1}^{N-1} p'(x_i) x_i\right) (*)$$

$$\geq f\left(p(x_N) x_N + (1 - p(x_N)) \sum_{i=1}^{N-1} p'(x_i) x_i\right) \text{ (convexity)}$$

$$= f\left(\sum_{i=1}^{N} p(x_i) x_i\right) ,$$
where  $p'(x_i) = \frac{p(x_i)}{1 - p(x_N)}.$ 

◆ロ > ◆母 > ◆臣 > ◆臣 >

MQ (P

Jensen's Inequality Gibbs's Inequality

### Inequalities: Jensen

Induction: Assume that (\*) the theorem holds for N-1 mass points.

$$\sum_{i=1}^{N} p(x_i) f(x_i) = p(x_N) f(x_N) + (1 - p(x_N)) \sum_{i=1}^{N-1} p'(x_i) f(x_i)$$
  

$$\geq p(x_N) f(x_N) + (1 - p(x_N)) f\left(\sum_{i=1}^{N-1} p'(x_i) x_i\right) (*)$$
  

$$\geq f\left(p(x_N) x_N + (1 - p(x_N)) \sum_{i=1}^{N-1} p'(x_i) x_i\right) \text{ (convexity)}$$
  

$$= f\left(\sum_{i=1}^{N} p(x_i) x_i\right) ,$$
  
where  $p'(x_i) = \frac{p(x_i)}{1 - p(x_N)}.$ 

・ロト ・回ト ・ヨト ・ヨト

SQC

1

Jensen's Inequality Gibbs's Inequality

## Inequalities: Jensen

Induction: Assume that (\*) the theorem holds for N-1 mass points.

$$\sum_{i=1}^{N} p(x_i) f(x_i) = p(x_N) f(x_N) + (1 - p(x_N)) \sum_{i=1}^{N-1} p'(x_i) f(x_i)$$

$$\geq p(x_N) f(x_N) + (1 - p(x_N)) f\left(\sum_{i=1}^{N-1} p'(x_i) x_i\right) (*)$$

$$\geq f\left(p(x_N) x_N + (1 - p(x_N)) \sum_{i=1}^{N-1} p'(x_i) x_i\right) \text{ (convexity)}$$

$$= f\left(\sum_{i=1}^{N} p(x_i) x_i\right) ,$$
where  $p'(x_i) = \frac{p(x_i)}{1 - p(x_N)}.$ 

・ロト ・部 ト ・ヨト ・ヨト

3

SQC

#### Jensen's Inequality Gibbs's Inequality

### Inequalities: Jensen

Induction: Assume that (\*) the theorem holds for N-1 mass points.

$$\sum_{i=1}^{N} p(x_i) f(x_i) = p(x_N) f(x_N) + (1 - p(x_N)) \sum_{i=1}^{N-1} p'(x_i) f(x_i)$$

$$\geq p(x_N) f(x_N) + (1 - p(x_N)) f\left(\sum_{i=1}^{N-1} p'(x_i) x_i\right) (*)$$

$$\geq f\left(p(x_N) x_N + (1 - p(x_N)) \sum_{i=1}^{N-1} p'(x_i) x_i\right) \text{ (convexity)}$$

$$= f\left(\sum_{i=1}^{N} p(x_i) x_i\right) ,$$
where  $p'(x_i) = \frac{p(x_i)}{1 - p(x_N)}.$ 

・ロト ・回ト ・ヨト ・ヨト

3

SQR

Jensen's Inequality Gibbs's Inequality

### Inequalities: Jensen

Induction: Assume that (\*) the theorem holds for N-1 mass points.

$$\sum_{i=1}^{N} p(x_i) f(x_i) = p(x_N) f(x_N) + (1 - p(x_N)) \sum_{i=1}^{N-1} p'(x_i) f(x_i)$$
  

$$\geq p(x_N) f(x_N) + (1 - p(x_N)) f\left(\sum_{i=1}^{N-1} p'(x_i) x_i\right) (*)$$
  

$$\geq f\left(p(x_N) x_N + (1 - p(x_N)) \sum_{i=1}^{N-1} p'(x_i) x_i\right) \text{ (convexity)}$$
  

$$= f\left(\sum_{i=1}^{N} p(x_i) x_i\right) ,$$

where 
$$p'(x_i) = \frac{p(x_i)}{1 - p(x_N)}$$
.

◆ロ > ◆母 > ◆臣 > ◆臣 >

MQ (P

Jensen's Inequality Gibbs's Inequality

# Gibbs' inequality



### W. Gibbs, 1839-1903

Teemu Roos Three Concepts: Information

イロト イポト イヨト イヨト

SQR

Jensen's Inequality Gibbs's Inequality

### Inqualities: Gibbs

### Gibbs' inequality

For any two discrete probability distributions p and q, we have

$$\sum_{x \in \mathcal{X}} p(x) \log_2 p(x) \ge \sum_{x \in \mathcal{X}} p(x) \log_2 q(x)$$

with equality if and only if p(x) = q(x) for all  $x \in \mathcal{X}$ .

*Proof.* Since  $\log_2 x = \frac{1}{\ln 2} \ln x$ , dividing both sides by  $\ln 2$  changes  $\log_2$  to  $\ln$ .

・ロト ・ 同ト ・ ヨト ・ ヨト

Jensen's Inequality Gibbs's Inequality

### Inqualities: Gibbs

### Gibbs' inequality

For any two discrete probability distributions p and q, we have

$$\sum_{x \in \mathcal{X}} p(x) \ln p(x) \ge \sum_{x \in \mathcal{X}} p(x) \ln q(x)$$

with equality if and only if p(x) = q(x) for all  $x \in \mathcal{X}$ .

*Proof.* Since  $\log_2 x = \frac{1}{\ln 2} \ln x$ , dividing both sides by  $\ln 2$  changes  $\log_2$  to  $\ln$ .

・ロト ・ 同ト ・ ヨト ・ ヨト

Jensen's Inequality Gibbs's Inequality

# Inequalities: Gibbs

$$\sum_{x \in \mathcal{X}} p(x) \ln p(x) \ge \sum_{x \in \mathcal{X}} p(x) \ln q(x)$$

$$\sum_{x \in \mathcal{X}} p(x) \ln q(x) - \sum_{x \in \mathcal{X}} p(x) \ln p(x) = \sum_{x \in \mathcal{X}} p(x) (\ln q(x) - \ln p(x))$$
$$= \sum_{x \in \mathcal{X}} p(x) \ln \frac{q(x)}{p(x)} \qquad \boxed{\ln x - \ln y = \ln \frac{x}{y}}$$
$$\leq \sum_{x \in \mathcal{X}} p(x) \left(\frac{q(x)}{p(x)} - 1\right) \qquad \boxed{\ln x \leq x - 1}$$
$$= \sum_{x \in \mathcal{X}} q(x) - \sum_{x \in \mathcal{X}} p(x) = 1 - 1 = 0 \quad \Box$$

Jensen's Inequality Gibbs's Inequality

# Inequalities: Gibbs

$$\sum_{x \in \mathcal{X}} p(x) \ln p(x) \ge \sum_{x \in \mathcal{X}} p(x) \ln q(x)$$

$$\sum_{x \in \mathcal{X}} p(x) \ln q(x) - \sum_{x \in \mathcal{X}} p(x) \ln p(x) = \sum_{x \in \mathcal{X}} p(x) (\ln q(x) - \ln p(x))$$
$$= \sum_{x \in \mathcal{X}} p(x) \ln \frac{q(x)}{p(x)} \qquad \boxed{\ln x - \ln y = \ln \frac{x}{y}}$$
$$\leq \sum_{x \in \mathcal{X}} p(x) \left(\frac{q(x)}{p(x)} - 1\right) \qquad \boxed{\ln x \le x - 1}$$
$$= \sum_{x \in \mathcal{X}} q(x) - \sum_{x \in \mathcal{X}} p(x) = 1 - 1 = 0 \quad \Box$$

Jensen's Inequality Gibbs's Inequality

# Inequalities: Gibbs

$$\sum_{x \in \mathcal{X}} p(x) \ln p(x) \ge \sum_{x \in \mathcal{X}} p(x) \ln q(x)$$

$$\sum_{x \in \mathcal{X}} p(x) \ln q(x) - \sum_{x \in \mathcal{X}} p(x) \ln p(x) = \sum_{x \in \mathcal{X}} p(x) (\ln q(x) - \ln p(x))$$
$$= \sum_{x \in \mathcal{X}} p(x) \ln \frac{q(x)}{p(x)} \qquad \boxed{\ln x - \ln y = \ln \frac{x}{y}}$$
$$\leq \sum_{x \in \mathcal{X}} p(x) \left(\frac{q(x)}{p(x)} - 1\right) \qquad \boxed{\ln x \le x - 1}$$
$$= \sum_{x \in \mathcal{X}} q(x) - \sum_{x \in \mathcal{X}} p(x) = 1 - 1 = 0 \quad \Box$$

Jensen's Inequality Gibbs's Inequality

# Inequalities: Gibbs

$$\sum_{x \in \mathcal{X}} p(x) \ln p(x) \ge \sum_{x \in \mathcal{X}} p(x) \ln q(x)$$

$$\sum_{x \in \mathcal{X}} p(x) \ln q(x) - \sum_{x \in \mathcal{X}} p(x) \ln p(x) = \sum_{x \in \mathcal{X}} p(x) (\ln q(x) - \ln p(x))$$
$$= \sum_{x \in \mathcal{X}} p(x) \ln \frac{q(x)}{p(x)} \qquad \boxed{\ln x - \ln y = \ln \frac{x}{y}}$$
$$\leq \sum_{x \in \mathcal{X}} p(x) \left(\frac{q(x)}{p(x)} - 1\right) \qquad \boxed{\ln x \leq x - 1}$$
$$= \sum_{x \in \mathcal{X}} q(x) - \sum_{x \in \mathcal{X}} p(x) = 1 - 1 = 0 \quad \Box$$

Outline Calculus Probability Inequalities	Jensen's Inequality Gibbs's Inequality	
--	---	--

# For next week, read Chapter 2 of Cover & Thomas and do **home assignment** (see course web page).

SQC

-