

Section 10. Probability Review

1. Probability functions
2. Random Variables
3. Probability density functions
- 4. Expected values and variance
5. Classical distributions

Instructor: He Wang

Department of Mathematics

Northeastern University

➤ Terminologies:

- **Experiment:** A repeatable procedure with a set of possible results.
- **Sample Space** $S = \{\text{all possible outcomes of an experiment}\}$ Set
- **Event:** A subset of S. $A \subseteq S$

Example 1. Experiment: Flipping a Coin once.

$$S = \{\text{Face}, \text{Tail}\}$$



Example 2. Experiment: Rolling a 6-sided die once.

$$S = \{1, 2, 3, 4, 5, 6\}$$



➤ The Probability Function

$$\begin{aligned} \mathbb{R} &\xrightarrow{f} \mathbb{R} \\ a &\rightarrow f(a) \end{aligned}$$

Definition (1930s, Kolmogorov)

Let S be a finite sample space. A **probability function** P is a function

$$P: \underbrace{S}_{\substack{\text{all subsets} \\ \text{of } S}} \rightarrow [0,1]$$

satisfying the following two axioms:

- 1.) $P(S) = 1$
- 2.) If $A \cap B = \emptyset$, then $P(A \cup B) = P(A) + P(B)$

If S is an infinite set, we need one more axiom:

$$3.) \quad P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i) \quad \text{if } A_i \cap A_j = \emptyset \text{ for any } i \neq j$$

Classical definition of probability (a special case for probability function)

Suppose the outcomes of an experiment are all equally likely, and the total number of all possible outcomes is finite. $\#(S)$

Probability of an event A := $\frac{\text{Number of ways it can happen } \#(A)}{\text{Total number of all possible outcomes } \#(S)}$

That is, $P(A) := \frac{|A|}{|S|}$

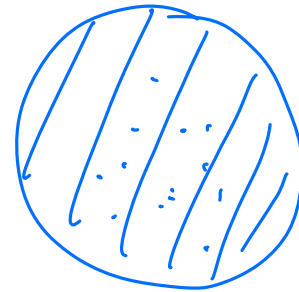
Example 1. Flipping a fair coin 2 times

$S = \{FF, FT, TF, TT\}$ $\rightarrow [0, 1]$



"fair"

Example. Rolling a 6-sided die once.



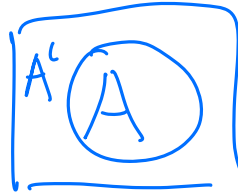
Example. flip an unfair coin once.

$$S = \left\{ \begin{array}{l} F \\ T \end{array} \right\} \xrightarrow{P} P \in [0, 1]$$
$$\qquad \qquad \qquad \xrightarrow{1-P} 1 - P$$

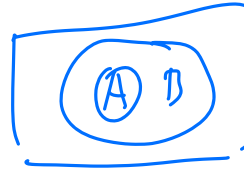


Some properties

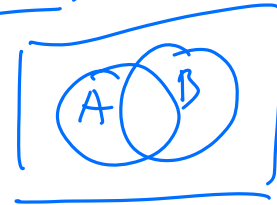
1. $P(A^C) = 1 - P(A)$.



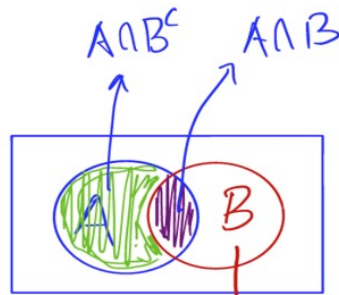
2. If $A \subset B$ then $P(A) \leq P(B)$.



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



4. $P(A) = P(A \cap B^C) + P(A \cap B)$.



➤ Conditional Probability

Definition. Probability that event A occurs given that event B already occurs, denoted by $\mathbf{P(A|B)}$ is a conditional probability, defined by

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

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

Example. Rolling a fair 6-sided die once.

$$S = \{1, 2, 3, 4, 5, 6\}$$

$\{2, 4, 6\} \Rightarrow B$
Given that the result is an even number, what is the probability that the result is 6?

$$A = \{6\}$$



$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{1/6}{3/6} = \frac{1}{3}$$

➤ **Impendence**

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

Definition: The events A and B are called **independent** if

$$\underline{P(A \cap B)} = \underline{P(A)}\underline{P(B)}$$

If A and B are not empty set, A and B are **independent** if and only if

$$P(A|B) = P(A) \text{ if and only if } P(B|A) = P(B)$$

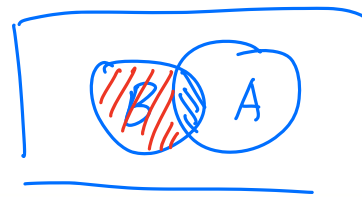
Example. Rolling a 6-sided die twice. $\{(a,b) \mid a \in \{1, \dots, 6\}, b \in \{1, \dots, 6\}\}$

A: the first face is even number
B: the second face is 6.



$$\underline{P(A \cap B)} = P(A)P(B) = \frac{3}{6} \cdot \frac{1}{6}$$

Theorem 1. Law of Total Probability



$$S = A \cup A^c$$

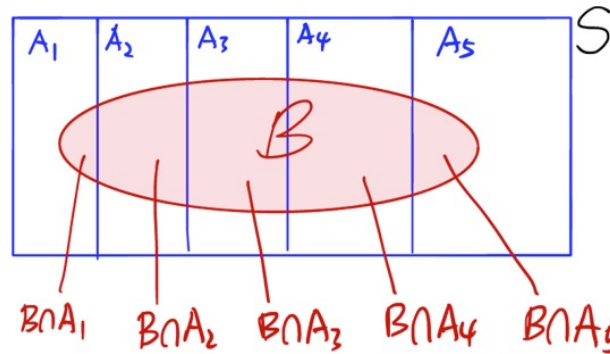
$$P(B) = P(B \cap A) + P(B \cap A^c)$$



Theorem. Law of Total Probability

Let A_1, A_2, \dots, A_n be a sequence of events such that $S = \bigcup_{i=1}^n A_i$ and $A_i \cap A_j = \emptyset$. Then, for any event B ,

$$P(B) = \sum_{i=1}^n P(B|A_i)P(A_i).$$



$$P(B) = P(B \cap A_1) + P(B \cap A_2) + P(B \cap A_3) + P(B \cap A_4) + P(B \cap A_5)$$

$$= P(B|A_1)P(A_1) + P(B|A_2)P(A_2) + P(B|A_3)P(A_3) + P(B|A_4)P(A_4) + P(B|A_5)P(A_5)$$

$P(B|A_i)$

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

Theorem 2. Bayes' Theorem

Theorem. Bayes' Theorem

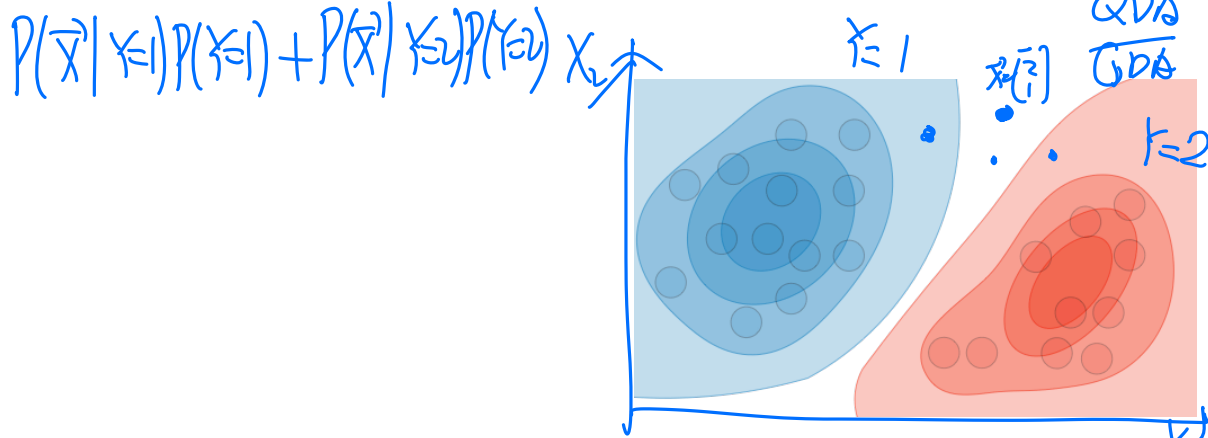
Let A_1, A_2, \dots, A_n be a sequence of events such that $S = \bigcup_{i=1}^n A_i$ and $A_i \cap A_j = \emptyset$. Then, for any event B ,

$$P(A_j | B) = \frac{P(B | A_j) P(A_j)}{\sum_{i=1}^n P(B | A_i) P(A_i)}$$

for any $j = 1, \dots, n$.

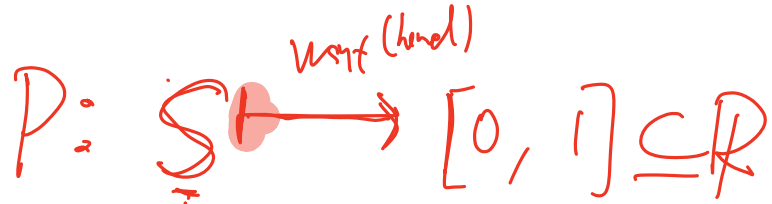
Example: in classification problem.

$$P(Y=1 | \vec{x}) = \frac{P(\vec{x} | Y=1) P(Y=1)}{P(\vec{x} | Y=1) P(Y=1) + P(\vec{x} | Y=2) P(Y=2)}$$



➤ Random Variables

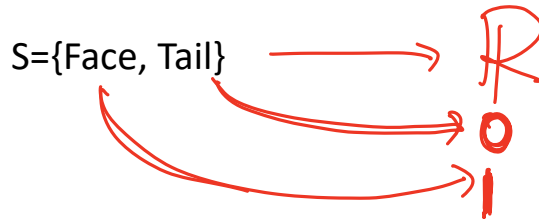
Let S be a sample space.



A random variable is a function



Example 1. Flipping an unfair coin once.



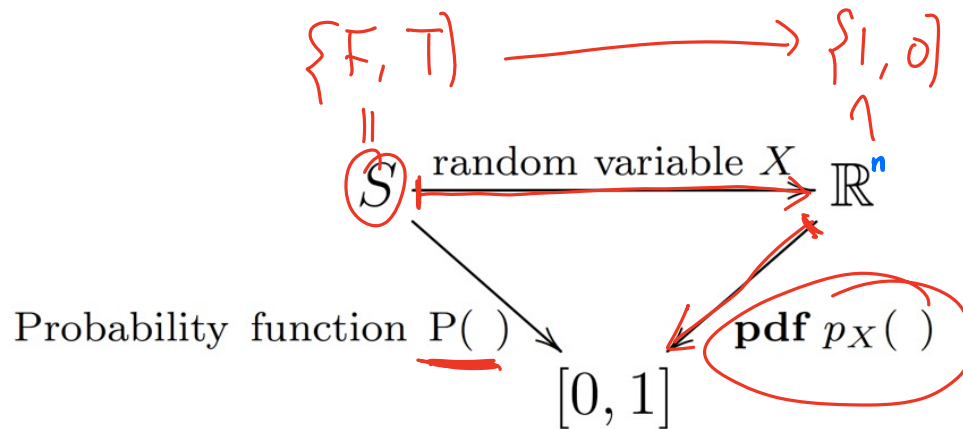
Probability density(mass) function

Definition.

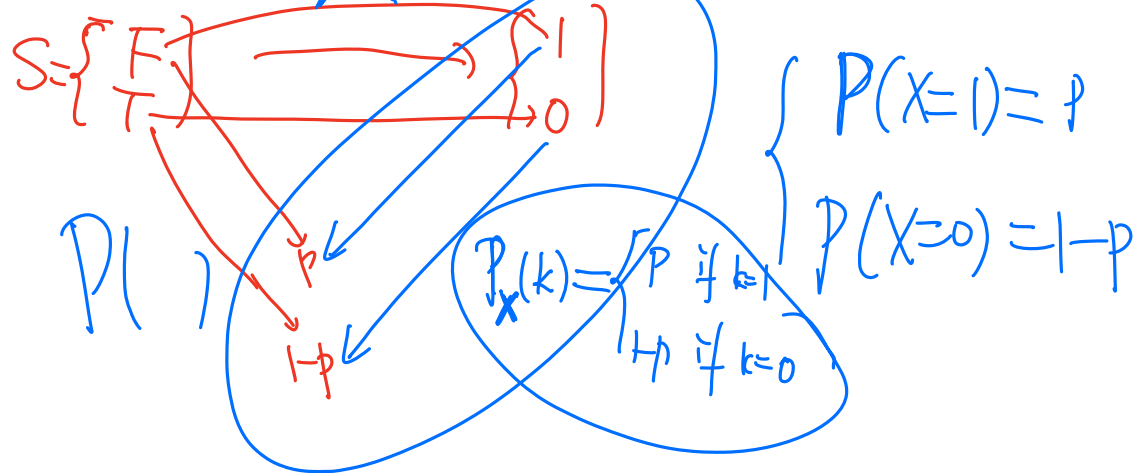
For every discrete random variable X , we define a probability density function (pdf) by

$$p_X(k) = P(X = k) := P(\{s \in S | X(s) = k\})$$

(If $k \notin X(S)$, then $p_X(k) = 0$.)



Example 1. Flipping an unfair coin once.

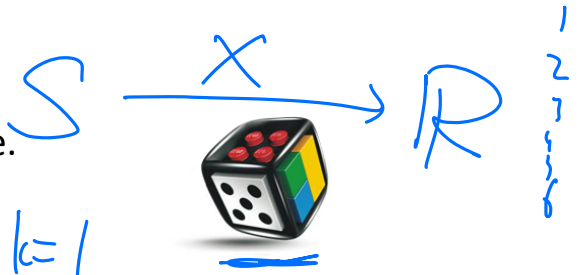


Random Variable $X \sim \text{Bernouli}(\phi)$

$$= p^k (1-p)^{1-k} \quad k \in \{0, 1\}$$

Pdf function $p_X(k) = \phi^k (1 - \phi)^{1-k} = \begin{cases} \phi & \text{if } k = 1 \\ 1 - \phi & \text{if } k = 0 \end{cases}$

Example. Rolling an unfair 6-sided die once.



$$P_X(k) = \begin{cases} \phi_1 & \text{if } k=1 \\ \phi_2 & \vdots \\ \vdots & \vdots \\ \phi_6 & \text{if } k=6 \end{cases}$$

$$\phi_1 + \dots + \phi_6 = 1$$

Assume $Y \sim \text{Categorical}(\phi_1, \dots, \phi_K)$ such that $\phi_1 + \dots + \phi_K = 1$

indicator function

$$p_Y(y) = \phi_1^{\mathbb{I}(y=1)} \phi_2^{\mathbb{I}(y=2)} \dots \phi_K^{\mathbb{I}(y=K)} = \begin{cases} \phi_1 & \text{if } y = 1 \\ \phi_2 & \text{if } y = 2 \\ \vdots & \vdots \\ \phi_K & \text{if } y = K \end{cases}$$

$$\mathbb{I}(\text{true}) = 1$$

$$\mathbb{I}(\text{false}) = 0$$

Continuous Random variables:

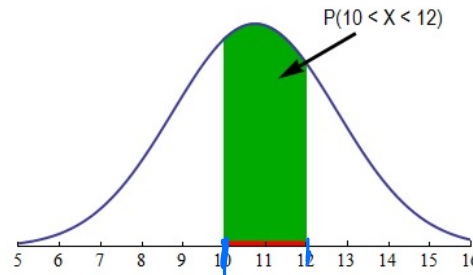
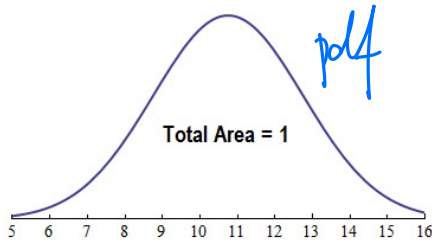


Definition.

The **probability density function (pdf)** of a **continuous** random variable X is a piecewise continuous function $f_X(x)$ satisfying

1. $f_X(x) \geq 0$

2. $\int_{-\infty}^{\infty} f_X(x) dx = 1.$



Definition.

The probability that X is in an interval $[a, b]$ is

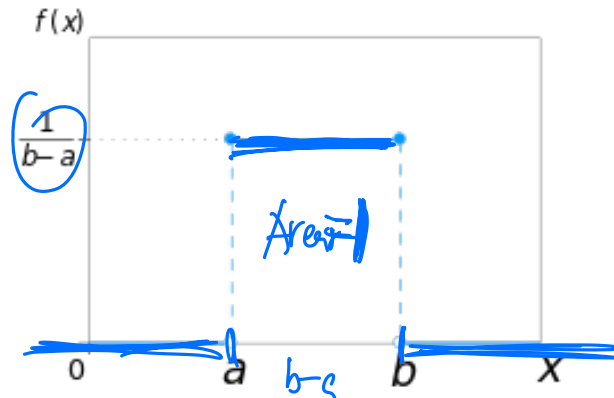
$$P(a \leq X \leq b) = \int_a^b f_X(x) dx$$

Example: Uniform Distribution.

The uniform distribution describes an experiment that choose a number randomly from the interval $[a, b]$.

The *probability density function* of the uniform distribution is

$$f(x) = \begin{cases} \frac{1}{b-a} & \text{for } a \leq x \leq b, \\ 0 & \text{for } x < a \text{ or } x > b \end{cases}$$



➤ Expected Value

Expected value is a generalization of the concept “average”.

Definition: Expected Value of **discrete** random variable

If X is a discrete random variable with probability function $p_X(k)$, then the **expected value** (or Mean) of X is

$$E(X) = \sum_{\text{all } k} k \cdot p_X(k).$$

Definition: (Expected Value of **continuous** random variable)

If X is a continuous random variable with probability function $p_X(x)$, then the **expected value** (or Mean) of X is

$$E(X) = \int_{-\infty}^{\infty} x \cdot p_X(x) dx.$$

Property: $E(aX + b) = aE(X) + b$

➤ Variance and Standard deviation

Definition. (Variance)

The **variance** of a random variable X is

$$\text{Var}(X) := E((X - \mu)^2)$$

Here $\mu = E(X)$ is the mean of X .

The **standard deviation** is $\sigma := \sqrt{\text{Var}(X)}$

Variance is expected squared distance from the mean.

It measures the spread of the data.

Calculation formula: $\text{Var}(X) = E(X^2) - (E(X))^2$

Property: $\text{Var}(aX + b) = a^2 \text{Var}(X)$

Standard Normal Distribution

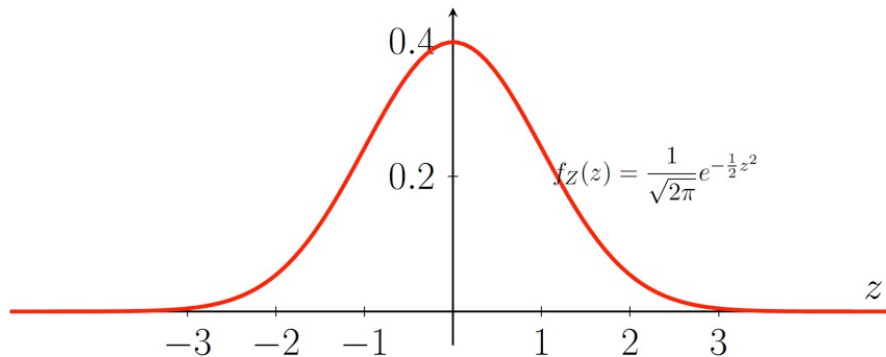
Definition.

The **standard normal distribution** is a continuous **pdf** defined by

$$f_Z(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z^2}$$

for $-\infty < z < \infty$.

The graph is **Gaussian** curve (bell curve).

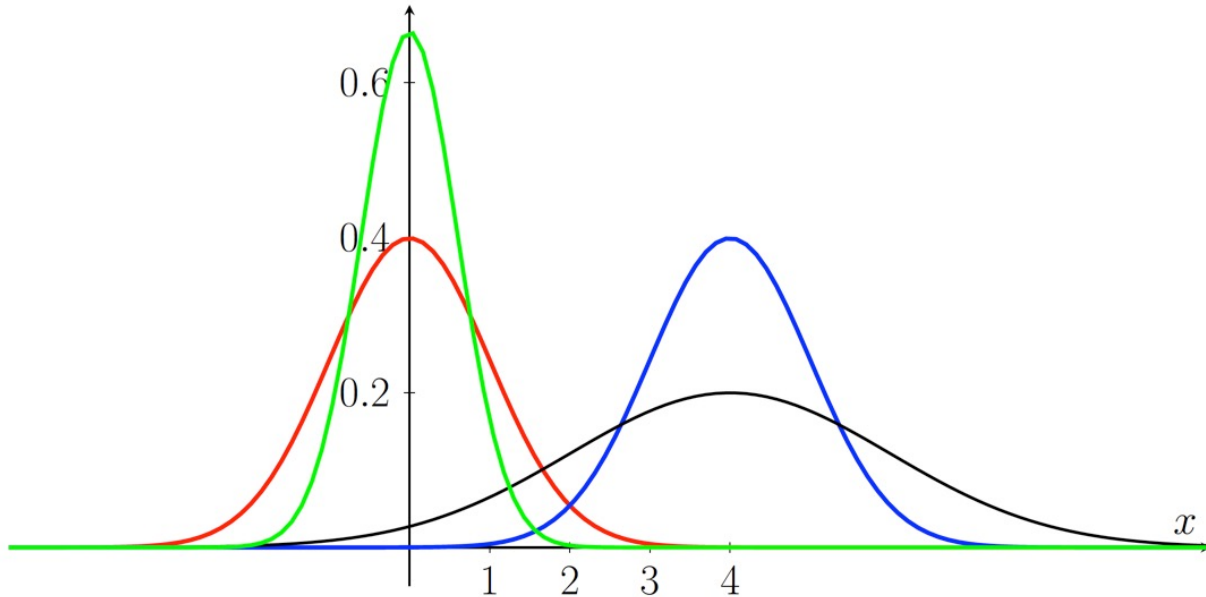


Normal distributions $X = \sigma Z + \mu$:

Definition. (Normal Distribution)

The **Normal Distribution** is a continuous **pdf** function defined as

$$f_X(x) = \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} \quad \text{for } -\infty < x < \infty.$$



Red: $\mu = 0, \sigma = 1$. **Green:** $\mu = 0, \sigma = 0.6$. **Blue:** $\mu = 4, \sigma = 1$. **Black:** $\mu = 4, \sigma = 2$.

➤ Covariance and independence

Suppose X and Y are **any** random variables on the same sample space.

- $E(aX + bY) = aE(X) + bE(Y)$
- $Var(aX + bY) = a^2 Var(X) + b^2 Var(Y) + 2ab Cov(X, Y)$

$Cov(X, Y)$ is the **covariance** of X and Y defined as

$$Cov(X, Y) := E(XY) - E(X)E(Y)$$

If X and Y are independent, then

- $Var(aX + bY) = a^2 Var(X) + b^2 Var(Y)$
- $Cov(X, Y) = 0$
- $E(XY) = E(X)E(Y)$

The converse is not true.

Joint distribution (multi random variables):

Definition. Discrete Joint Density

Let S is a **discrete** sample space. Let X and Y be two random variables on S . The **joint probability density function (joint pdf)** of X and Y is denoted by $p_{X,Y}(x, y)$ defined as

$$p_{X,Y}(x, y) := P(X = x, Y = y).$$

Here, $P(X = x, Y = y)$ is the probability when $X = x$ **and** $Y = y$.

Definition.

If X and Y are **continuous** random variables. the **joint pdf** $f_{X,Y}(x, y)$ of X and Y is a piecewise continuous multi-variable function satisfying

$$(1.) f_{X,Y}(x, y) \geq 0.$$

$$(2.) \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{X,Y}(x, y) dx dy = 1.$$

➤ **Multivariate normal distribution.**

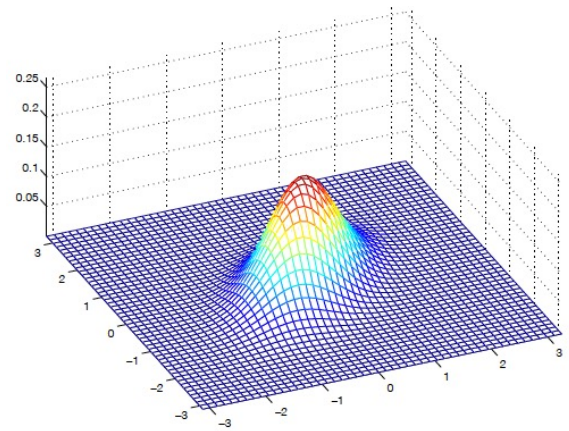
Vector random variable $\vec{X} = \begin{bmatrix} X_1 \\ \vdots \\ X_d \end{bmatrix} \sim \text{Normal}(\vec{\mu}, \Sigma)$

Here $\vec{\mu} \in \mathbb{R}^d$ and Σ is an $d \times d$ symmetric, positive definite matrix.

- The **joint** probability density function (**pdf**) for \vec{X} is

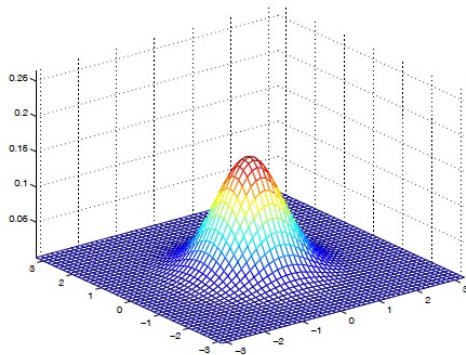
$$f_{\vec{X}}(\vec{x}) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left(-\frac{1}{2}(\vec{x} - \vec{\mu})^T \Sigma^{-1}(\vec{x} - \vec{\mu})\right)$$

$$\int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} f_{\vec{X}}(\vec{x}) \, d\vec{x} = 1$$

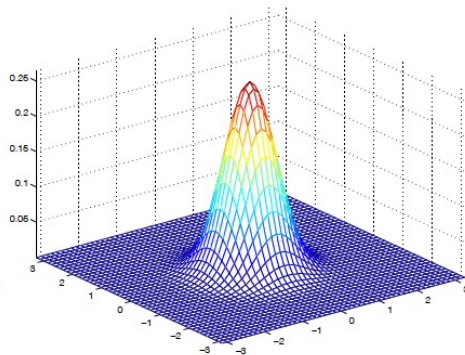


- The **mean vector** of \vec{X} is $E(\vec{X}) = \vec{\mu}$
- The **(co)variance matrix** is $\text{Cov}(\vec{X}) = \Sigma$

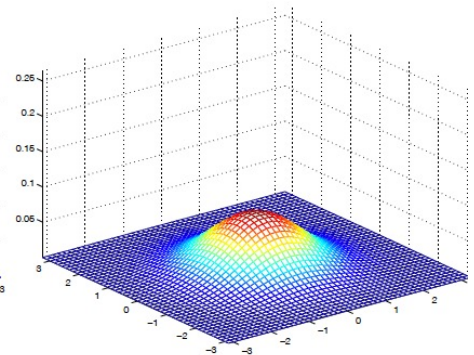
Standard normal

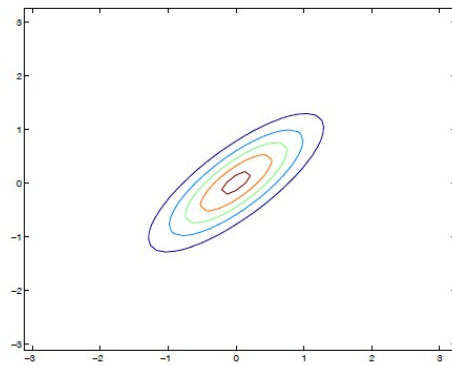
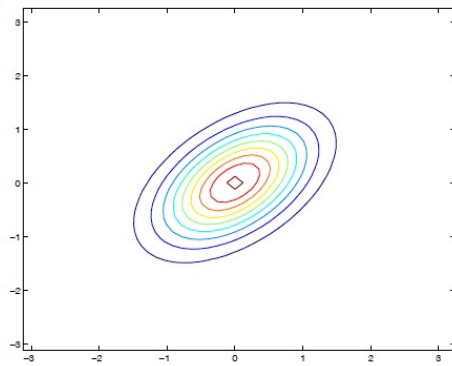
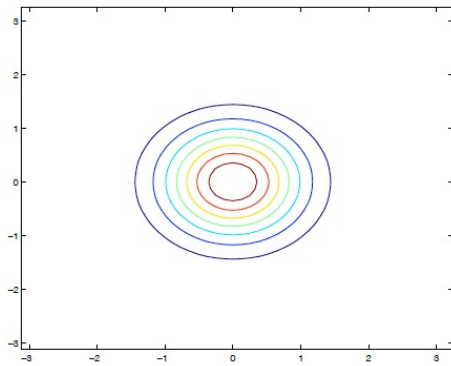
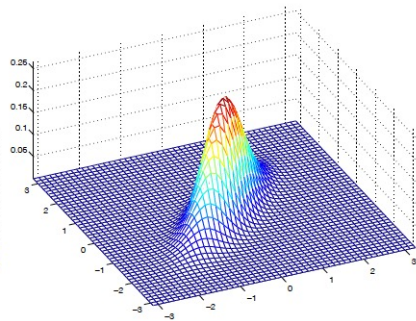
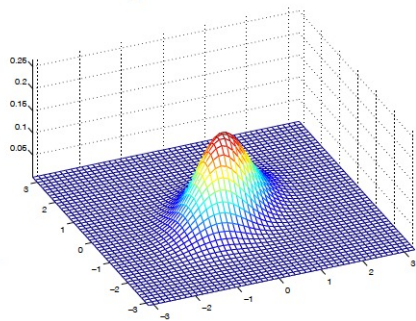
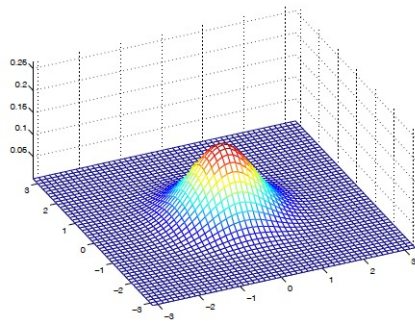


Compressed



Spread-out





➤ **More examples of distributions:**

1. Binomial distribution is a generalization of Bernoulli distribution.

Given a series of n independent trials with two outcomes (T or F) with constant probability p and $1 - p$.

Let X be the number of T appears in the n trials. Then $X \sim \text{Binomial}(n, p)$

$$P(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}$$

For example, flip a coin n times.

2. Multinomial is a generalization of Categorical distribution.

Given a series of n independent trials with m outcomes (O_1, \dots, O_m) with constant probability (ϕ_1, \dots, ϕ_m).

Let \vec{X} be the number of O_i appears in the n trials.

Then $\vec{X} \sim \text{Multinomial}(n, \phi_1, \dots, \phi_m)$

$$P(X_i = n_i) = \frac{n!}{n_1! \dots n_m!} \phi_1^{n_1} \dots \phi_m^{n_m}$$

for each $i = 1, \dots, m$, and each $n_1 + \dots + n_m = n$

For example, Toss a K -side die n times.

3. Exponential random variable is a continuous random variable with pdf given by

$$f_X(x) = \lambda e^{-\lambda x} \text{ for } x \geq 0$$

where λ is a fixed positive number.

- The **mean** and **variance** are given by

$$E(X) = \frac{1}{\lambda}$$

$$\text{Var}(X) = \frac{1}{\lambda^2}$$

- Exponential distribution model the time between occurrences in a time interval.

4. Poisson Distribution

Definition. (Poisson Distribution)

The **Poisson Distribution** $\text{Poisson}(\lambda)$ is a discrete **pdf** function defined as

$$p_X(k) = P(X = k) := \frac{\lambda^k e^{-\lambda}}{k!}$$

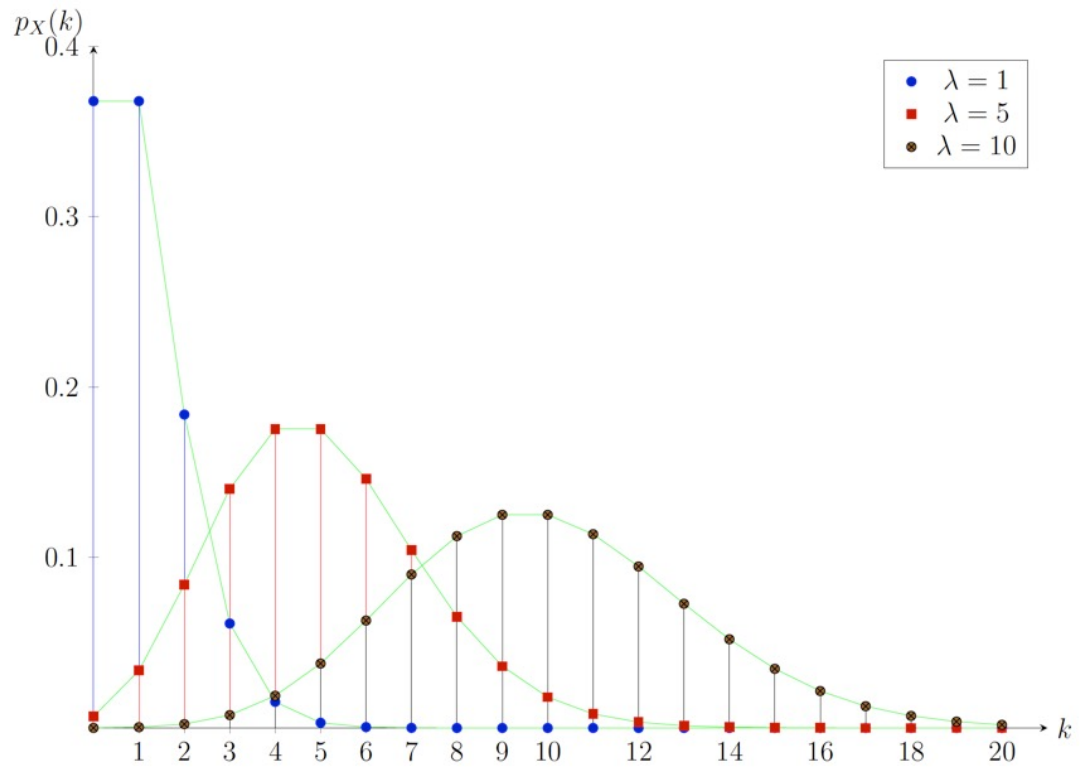
for $k = 0, 1, 2, 3, \dots$. Here, λ is a positive constant.

Theorem.

- (1) It is a well defined **pdf**, i.e., $\sum_k p_X(k) = 1$
- (2) The mean is $E(X) = \lambda$.
- (3) The variance is $\text{Var}(X) = \lambda$.

Applications:

1. Poisson approximation for binomial distribution
2. Poisson Model. The number of occurrences in a time interval with a given rate.



More references:

1. My lecture notes for Math3081:

<https://web.northeastern.edu/he.wang/Teaching/Teaching3081/Math3081.html>

2. Deep Learning by Goodfellow, Bengio, Courville (Chapter 2)

<https://www.deeplearningbook.org/>

3. Pattern Recognition and Machine Learning, by Chris Bishop.
(Chapters 1.2 and 2.1-2.3)

<https://www.microsoft.com/en-us/research/uploads/prod/2006/01/Bishop-Pattern-Recognition-and-Machine-Learning-2006.pdf>

4. Review of Probability Theory

<https://cs229.stanford.edu/section/cs229-prob.pdf>