

Introduction to Machine Learning

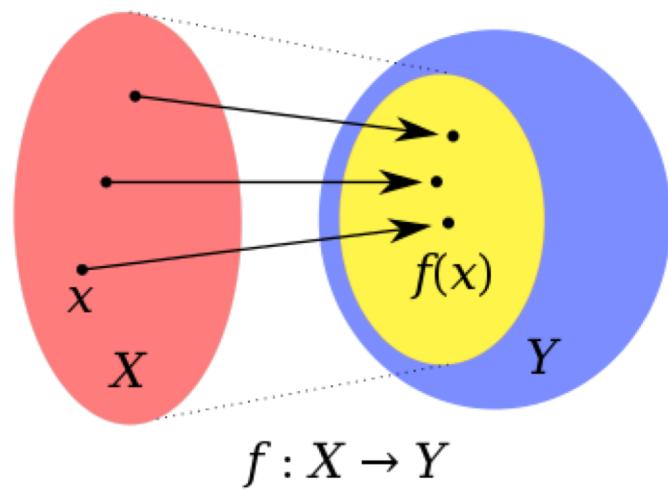
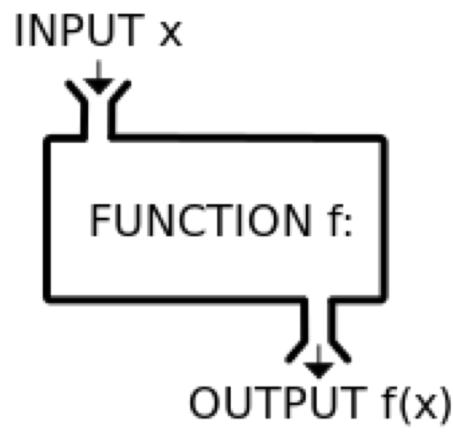
Part 1: Mathematical Foundation of Machine Learning

Zengchang Qin

Function and Data Generalization

Functions

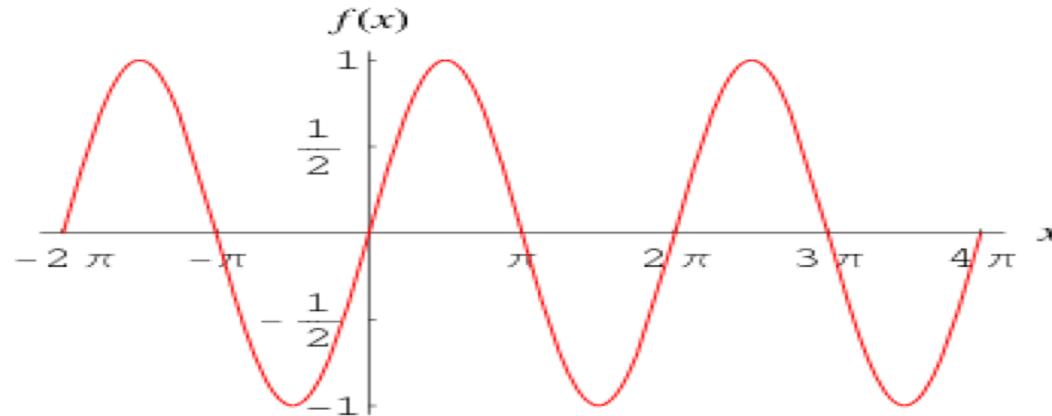
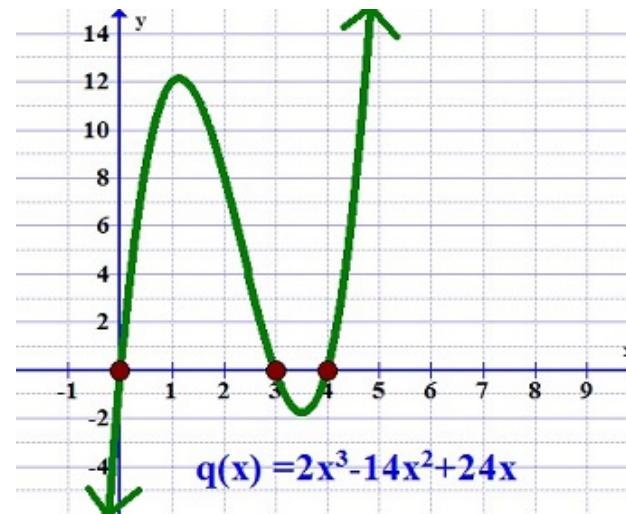
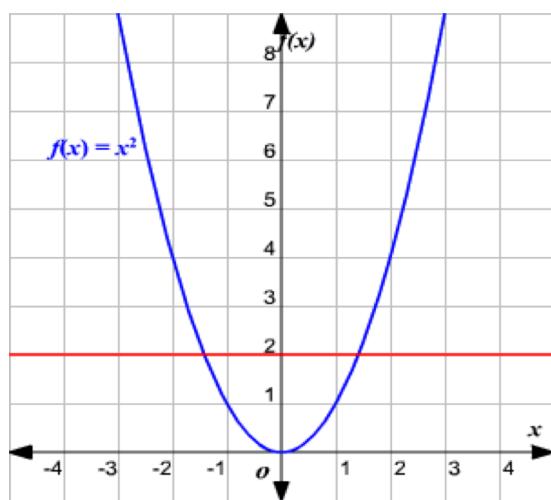
In mathematics, a **function** is a relation between a set of inputs and a set of permissible outputs with the property that each input is related to exactly one output.



A sample function: $f(x) = 2x+3$

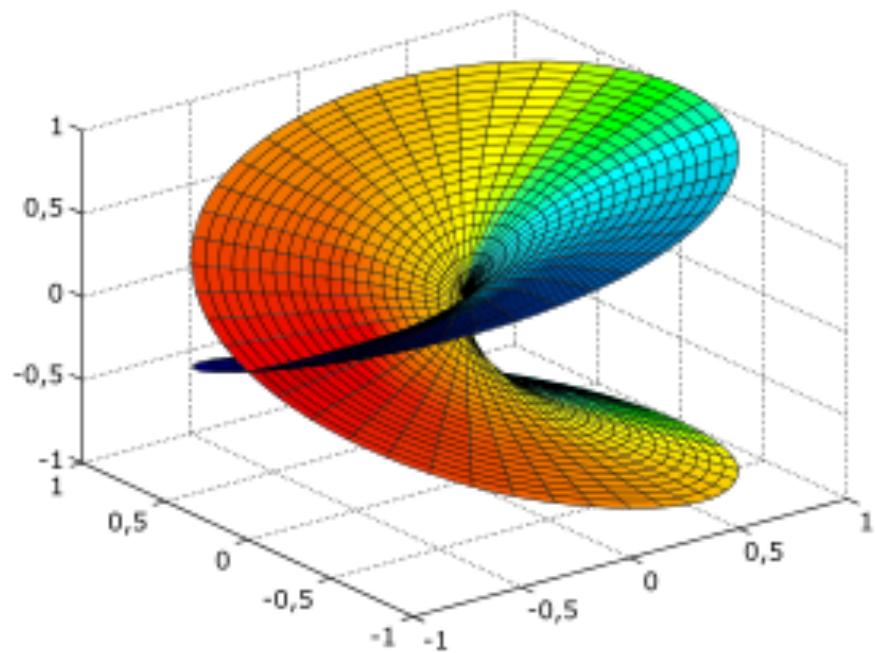
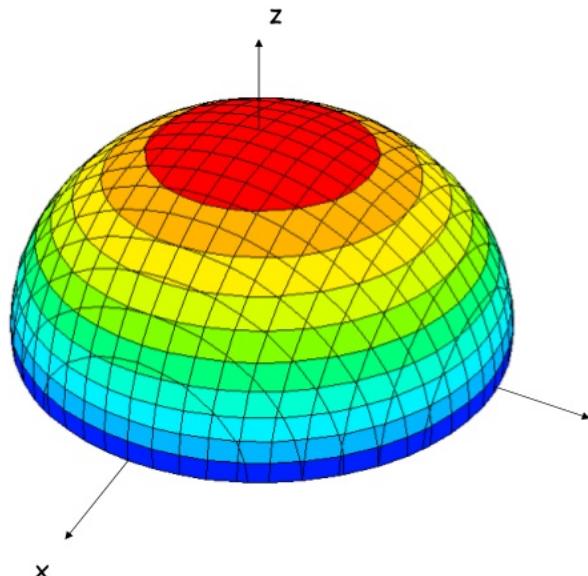
Functions

We have learned different **types** of functions.



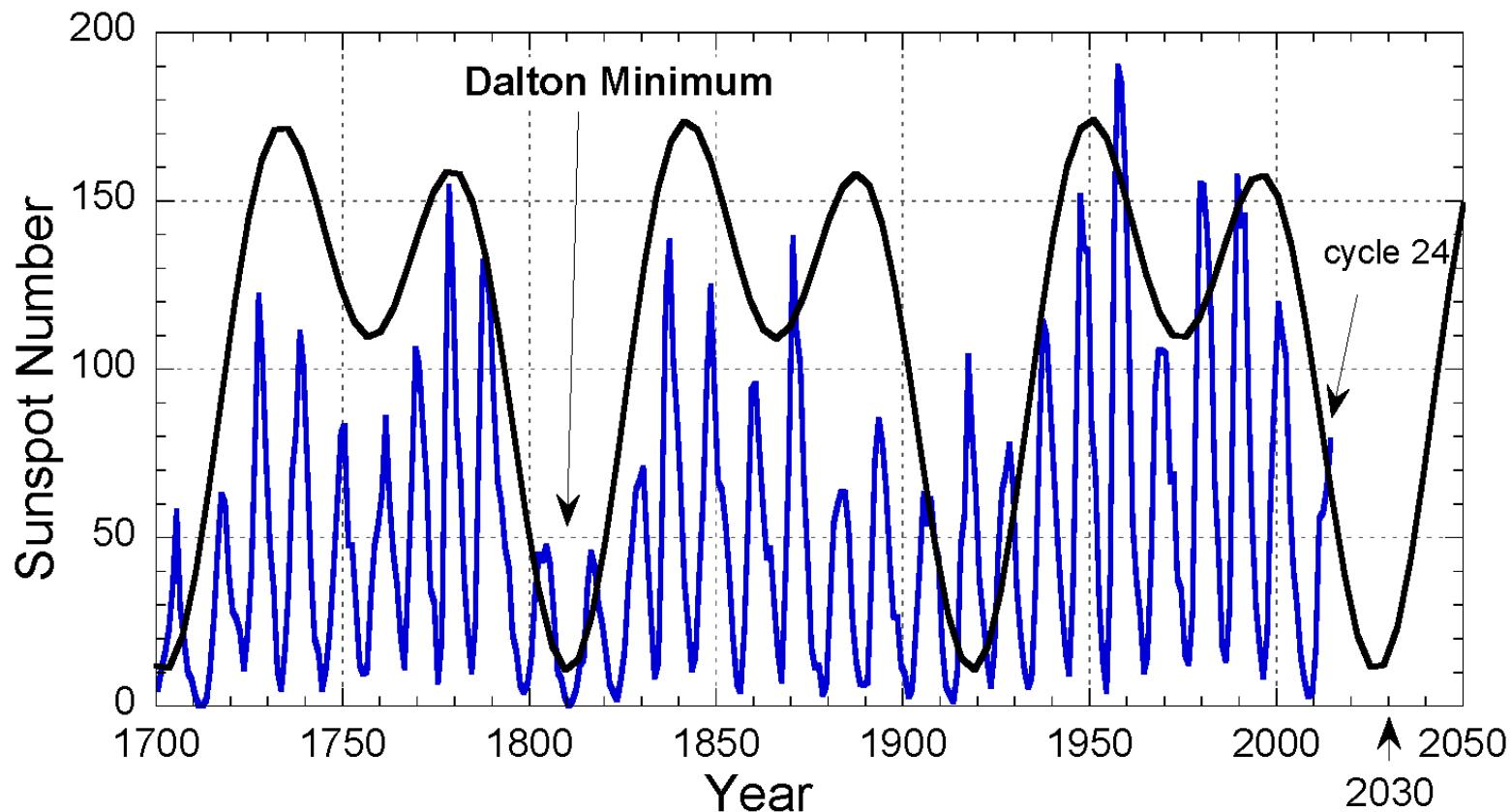
Functions

We have learned different **types** of functions.

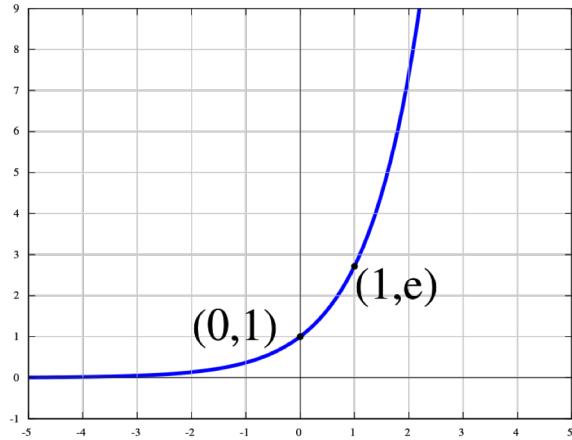


The Real-World Data

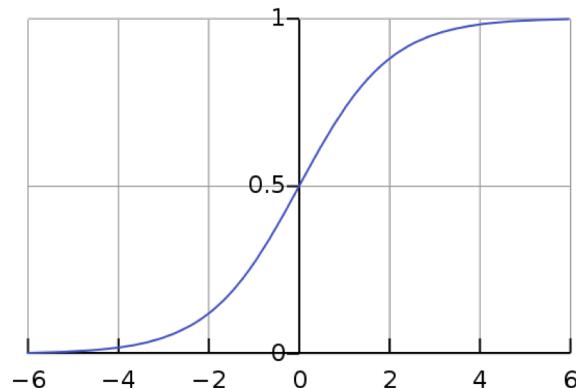
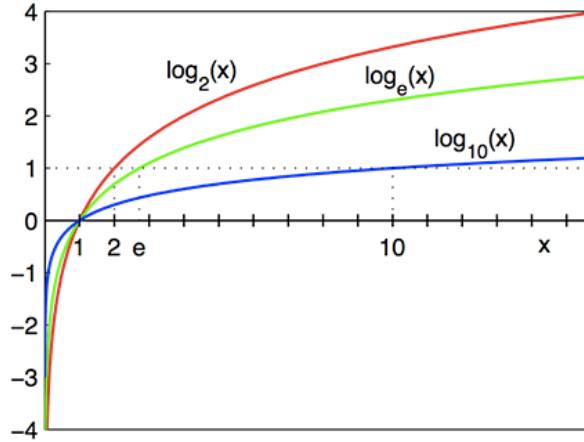
In the real-world, when we are investigating relations, we may find the following:



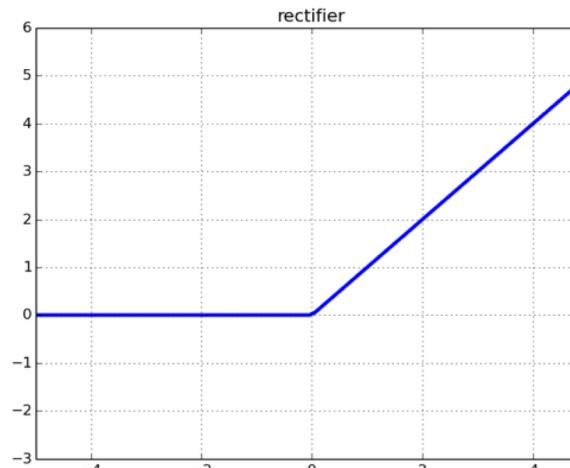
Some Functions



$$y = e^x \quad \text{http://setosa.io/ev/exponentiation/}$$



$$S(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{e^x + 1}$$



$$f(x) = x^+ = \max(0, x)$$

Function Decomposition

"Function Composition" is applying one function to the results of another:
The result of $f()$ is sent through $g()$

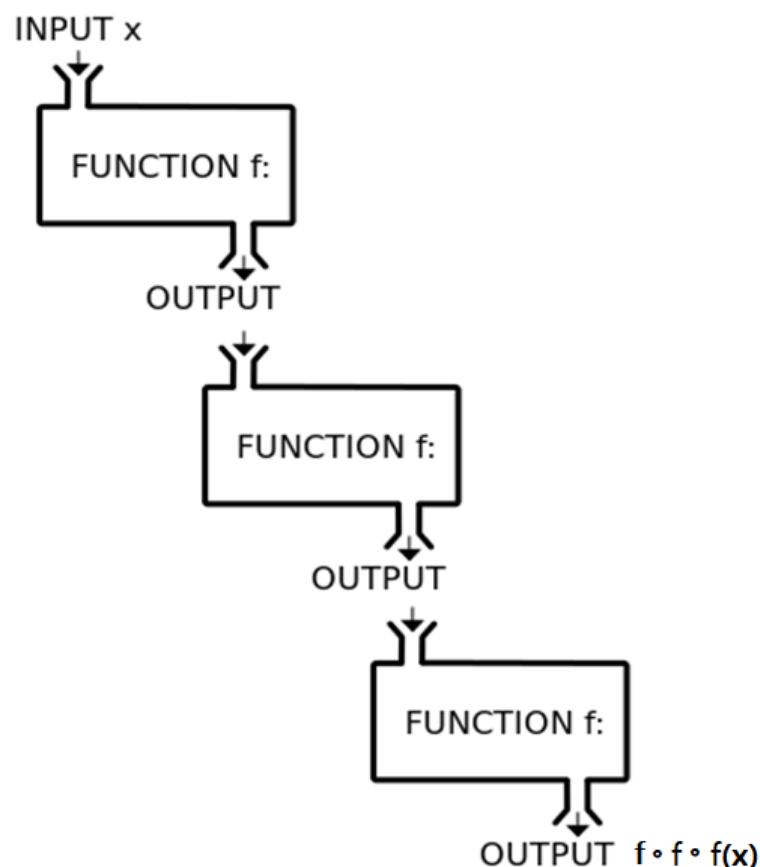
It is written: $(g \circ f)(x)$
Which means: $g(f(x))$

$$f(x) = 2x + 3$$

$$f \circ f(x) = ?$$

$$f \circ f \circ f(x) = ?$$

$$f \circ f \circ f(2) =$$

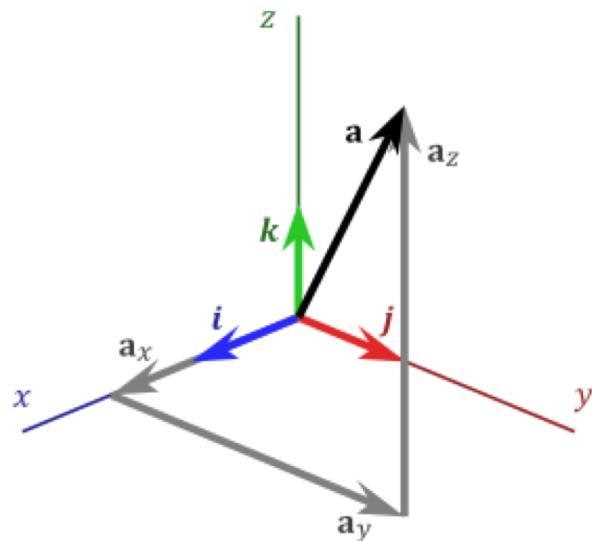


Linear Algebra

Vector

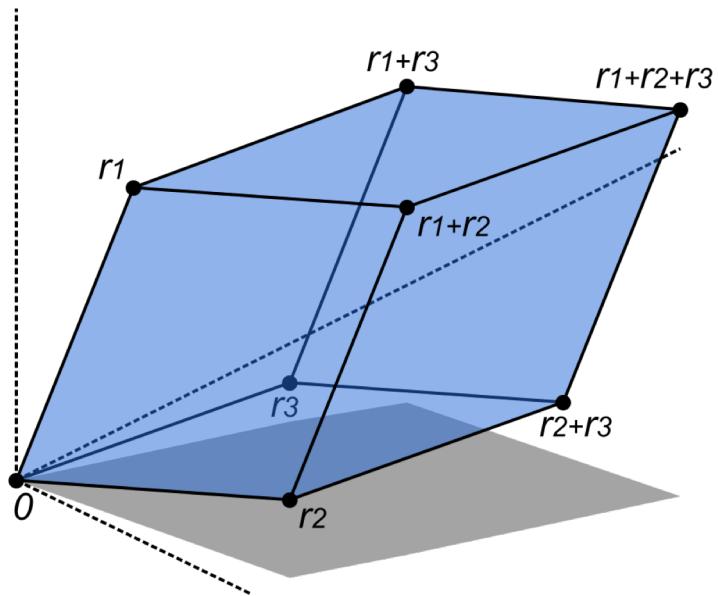
A **vector space** V is a set (the elements of which are called vectors) on which two operations are defined: vectors can be added together, and vectors can be multiplied by real numbers called **scalars**.

Can be written in column form or row form – **Column form is conventional!**


$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad \mathbf{x} + \mathbf{y} = \begin{bmatrix} x_1 + y_1 \\ \vdots \\ x_n + y_n \end{bmatrix}, \quad \alpha\mathbf{x} = \begin{bmatrix} \alpha x_1 \\ \vdots \\ \alpha x_n \end{bmatrix}$$

Vector Space

- Euclidean space is used to mathematically represent physical space, with notions such as distance, length, and angles.
- Although it becomes hard to visualize for $n > 3$, these concepts generalize mathematically in obvious ways.
- Linear relations hold in high dimensional space.



Norm of Vectors

A **norm** on a real vector space V is a function $\|\cdot\| : V \rightarrow \mathbb{R}$ that satisfies

- (i) $\|\mathbf{x}\| \geq 0$, with equality if and only if $\mathbf{x} = \mathbf{0}$
- (ii) $\|\alpha\mathbf{x}\| = |\alpha|\|\mathbf{x}\|$
- (iii) $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$ (the **triangle inequality** again)

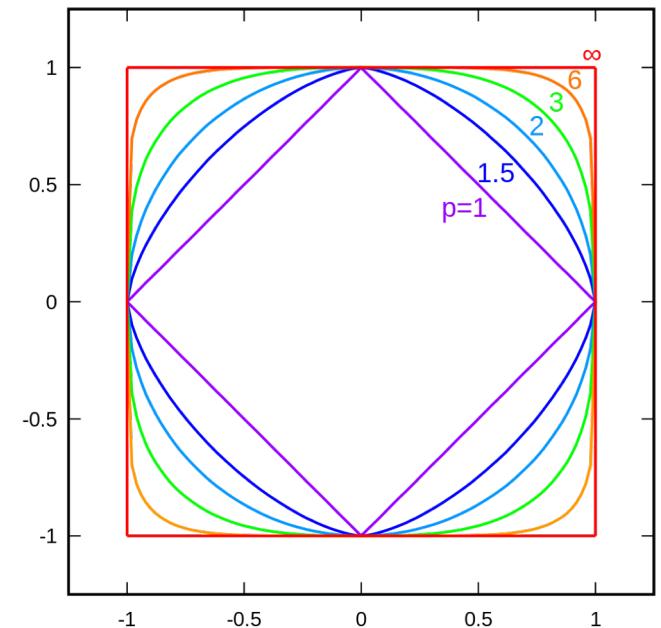
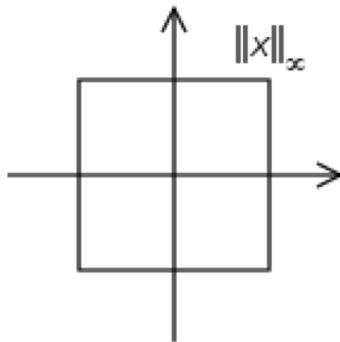
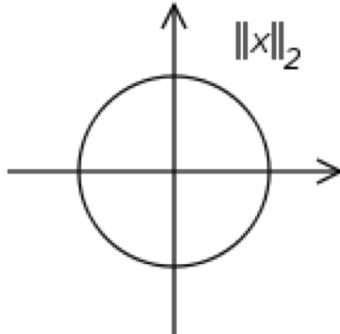
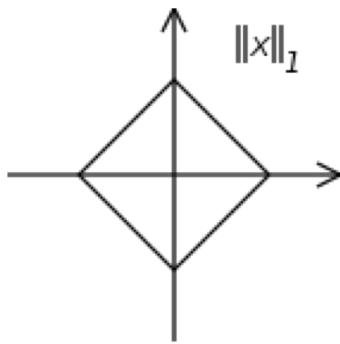
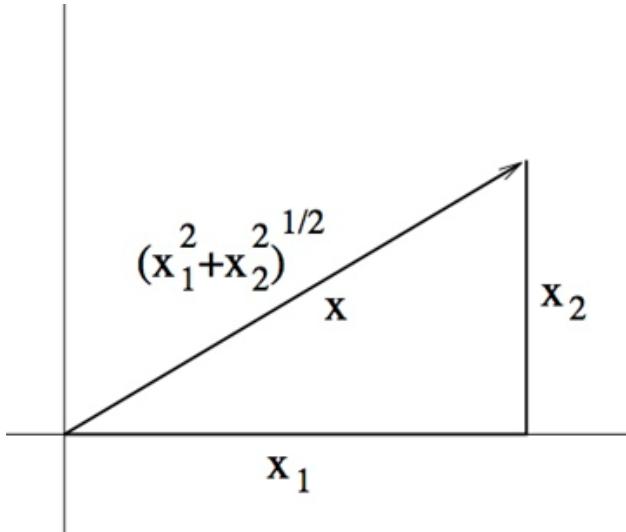
We will typically only be concerned with a few specific norms on \mathbb{R}^n :

$$\begin{aligned}\|\mathbf{x}\|_1 &= \sum_{i=1}^n |x_i| & \|\mathbf{x}\|_p &= \left(\sum_{i=1}^n |x_i|^p \right)^{\frac{1}{p}} & (p \geq 1) \\ \|\mathbf{x}\|_2 &= \sqrt{\sum_{i=1}^n x_i^2} & \|\mathbf{x}\|_\infty &= \max_{1 \leq i \leq n} |x_i|\end{aligned}$$

L-0 to L-infinity Norms

a **norm** is a function that assigns a strictly *positive length* to a vector.

A simple example is two dimensional Euclidean space R2 equipped with the "Euclidean norm"



$$\|\mathbf{x}\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{\frac{1}{p}}$$

Matrix

A vector can be regarded as **special case** of a matrix, where one of matrix dimensions = 1.

$$\mathbf{A} = \begin{pmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \cdots & a_{mn} \end{pmatrix}$$

Matrix transpose (denoted T)

$$\mathbf{A} = \begin{pmatrix} 2 & 7 & -1 & 0 & 3 \\ 4 & 6 & -3 & 1 & 8 \end{pmatrix} \quad \mathbf{A}^T = \begin{pmatrix} 2 & 4 \\ 7 & 6 \\ -1 & -3 \\ 0 & 1 \\ 3 & 8 \end{pmatrix}$$

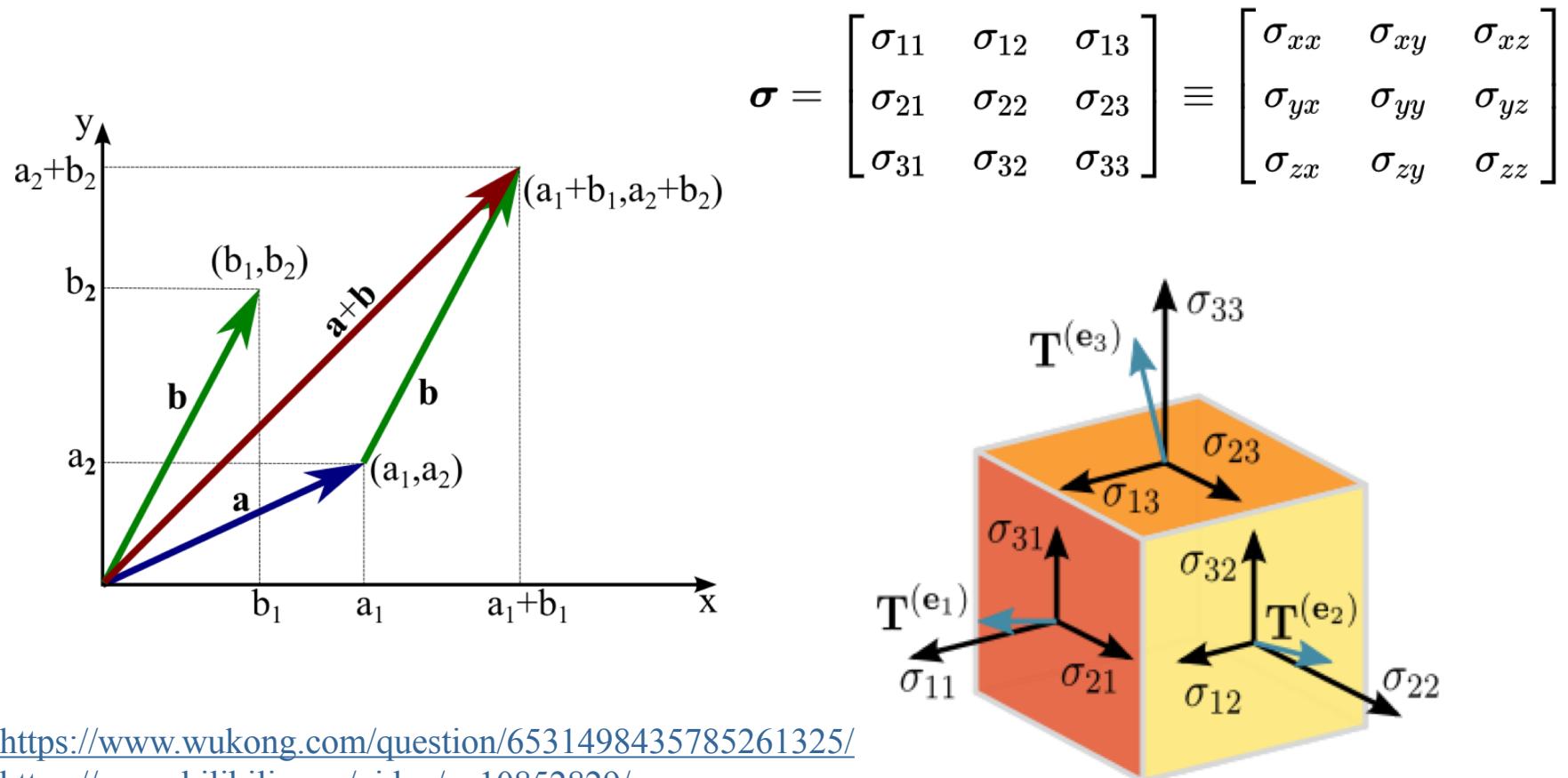
$$C = AB \quad \Leftrightarrow \quad c_{ij} = \sum_{k=1}^n a_{ik} b_{kj},$$

$$\begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \\ a_{31} & a_{32} \end{bmatrix} \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{bmatrix} = \begin{bmatrix} a_{11}b_{11} + a_{12}b_{21} & a_{11}b_{12} + a_{12}b_{22} \\ a_{21}b_{11} + a_{22}b_{21} & a_{21}b_{12} + a_{22}b_{22} \\ a_{31}b_{11} + a_{32}b_{21} & a_{31}b_{12} + a_{32}b_{22} \end{bmatrix}$$

Vector to Tensor

<https://www.quora.com/What-is-a-tensor>

Columns are the stresses (forces per unit area) acting on the \mathbf{e}_1 , \mathbf{e}_2 , and \mathbf{e}_3 faces of the cube.



<https://www.wukong.com/question/6531498435785261325/>

<https://www.bilibili.com/video/av10852829/>

Tensor

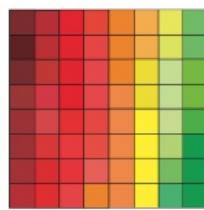
tensor = multidimensional array

vector



$$\mathbf{v} \in \mathbb{R}^{64}$$

matrix

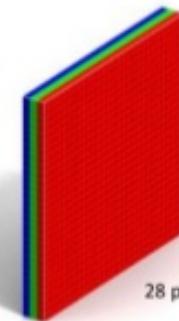
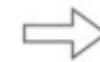


$$X \in \mathbb{R}^{8 \times 8}$$

tensor



$$\mathcal{X} \in \mathbb{R}^{4 \times 4 \times 4}$$



Color Image
(RGB)

3 channels
(RGB)

4	6	1	3
0	9	7	2
2	26	35	19
1	15	22	25

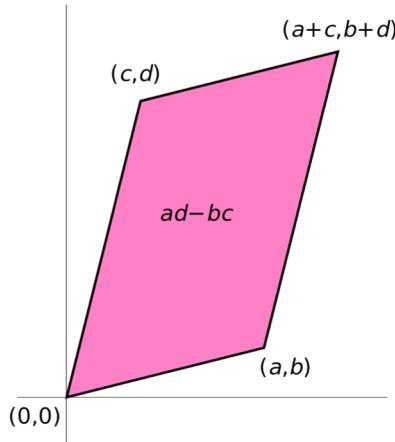
3 Colour Channels

Height: 4 Units
(Pixels)

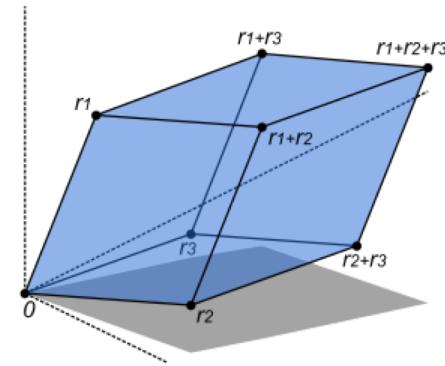
Width: 4 Units
(Pixels)

Determinant

In linear algebra, the determinant is a useful value that can be computed from the elements of a square matrix. The determinant of a matrix A is denoted $\det(A)$, $\det A$, or $|A|$. It can be viewed as the scaling factor of the transformation described by the matrix.



$$|A| = \begin{vmatrix} a & b \\ c & d \end{vmatrix} = ad - bc.$$



$$\begin{aligned} |A| &= \begin{vmatrix} a & b & c \\ d & e & f \\ g & h & i \end{vmatrix} = a \begin{vmatrix} e & f \\ h & i \end{vmatrix} - b \begin{vmatrix} d & f \\ g & i \end{vmatrix} + c \begin{vmatrix} d & e \\ g & h \end{vmatrix} \\ &= aei + bfg + cdh - ceg - bdi - afh. \end{aligned}$$

Eigenvector and Eigenvalue

For a square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$, there may be vectors which, when \mathbf{A} is applied to them, are simply scaled by some constant. We say that a nonzero vector $\mathbf{x} \in \mathbb{R}^n$ is an **eigenvector** of \mathbf{A} corresponding to **eigenvalue** λ if

$$\mathbf{Ax} = \lambda\mathbf{x}$$

The zero vector is excluded from this definition because $\mathbf{A}\mathbf{0} = \mathbf{0} = \lambda\mathbf{0}$ for every λ .

We now give some useful results about how eigenvalues change after various manipulations.

The **trace** of a square matrix is the sum of its diagonal entries:

$$\text{tr}(\mathbf{A}) = \sum_{i=1}^n A_{ii}$$

<http://setosa.io/ev/eigenvectors-and-eigenvalues/>

Singular Value Decomposition

Singular Value Decomposition:

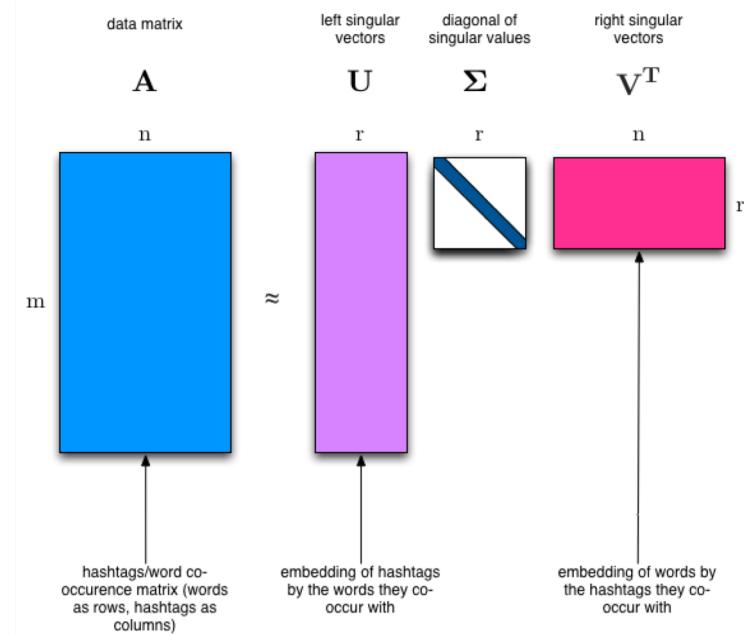
Formally, the SVD of a real $m \times n$ matrix A is a factorization of the form $A = U \Sigma V^T$, where U is an $m \times m$ orthogonal matrix of left singular vectors, Σ is an $m \times n$ diagonal matrix of singular values, and V^T is an $n \times n$ orthogonal matrix of right singular vectors.

$$M = \begin{bmatrix} 1 & 0 & 0 & 0 & 2 \\ 0 & 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 \end{bmatrix}$$

$$U = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & -1 \\ 1 & 0 & 0 & 0 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} 2 & 0 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 & 0 \\ 0 & 0 & \sqrt{5} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$V^* = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ \sqrt{0.2} & 0 & 0 & 0 & \sqrt{0.8} \\ 0 & 0 & 0 & 1 & 0 \\ -\sqrt{0.8} & 0 & 0 & 0 & \sqrt{0.2} \end{bmatrix}$$



Jacobian and Hessian Matrices

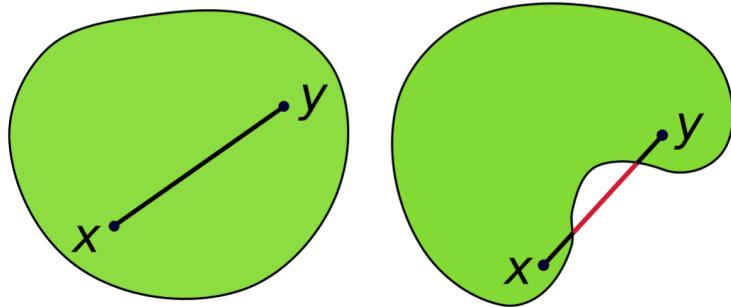
The **Jacobian** of $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a matrix of first-order partial derivatives:

$$\mathbf{J}_f = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} \quad \text{i.e.} \quad [\mathbf{J}_f]_{ij} = \frac{\partial f_i}{\partial x_j} \quad \text{Note the special case } m = 1, \text{ where } \nabla f = \mathbf{J}_f^\top.$$

The **Hessian** matrix of $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is a matrix of second-order partial derivatives:

$$\nabla^2 f = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix} \quad \text{i.e.} \quad [\nabla^2 f]_{ij} = \frac{\partial^2 f}{\partial x_i \partial x_j}$$

Convex Set and Function



A function f is **convex** if

$$f(t\mathbf{x} + (1 - t)\mathbf{y}) \leq tf(\mathbf{x}) + (1 - t)f(\mathbf{y})$$

for all $\mathbf{x}, \mathbf{y} \in \text{dom } f$ and all $t \in [0, 1]$.

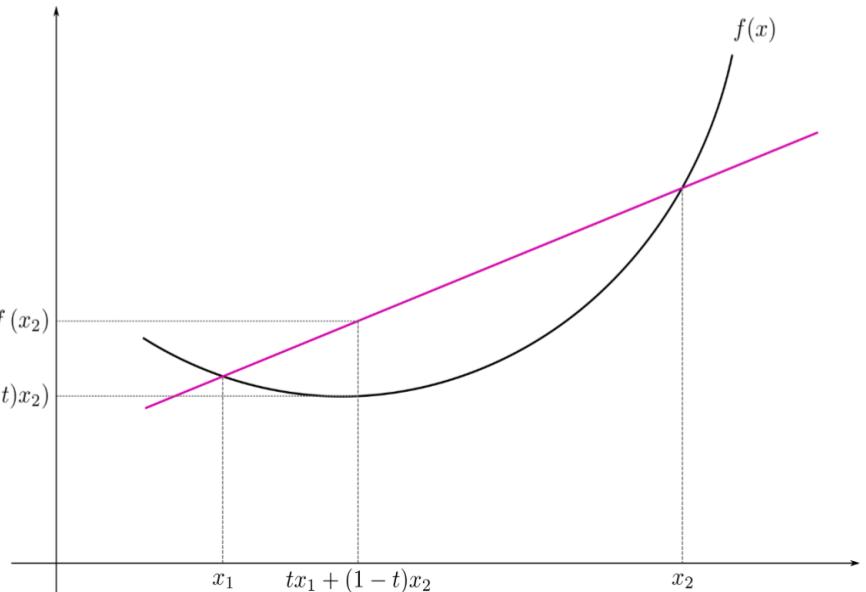


Figure 2: What convex functions look like

Probability and Statistics

Probability (Objective and Subjective)

The first approach is to define probability in terms of frequency of occurrence, as a percentage of successes in a moderately large number of similar situations.



Such an interpretation is often natural. For example, when we say that a perfectly manufactured coin lands on heads “with probability 50%,” we typically mean “roughly half of the time.”

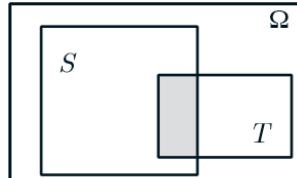
Consider, for example, a scholar who asserts that the Lliad and the Odyssey were composed by the same person, with probability 90%. Such an assertion conveys some information, but not in terms of frequencies, since the subject is a one-time event. Rather, it is an expression of the scholar’s subjective belief.

A blue-toned portrait engraving of Thomas Bayes, an English statistician and Presbyterian minister. He is depicted from the chest up, wearing a dark coat over a white cravat and a bow tie. The background is a textured blue.
$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

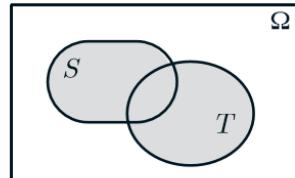
Set Operation

Examples of Venn diagrams.

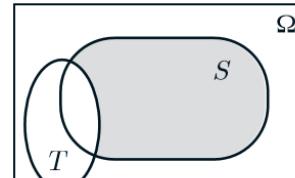
- (a) The shaded region is $S \cap T$.
- (b) The shaded region is $S \cup T$.
- (c) The shaded region is $S \cap c(T)$.
- (d) Here, $T \subset S$. The shaded region is the complement of S .
- (e) The sets S , T , and U are disjoint.
- (f) The sets S , T , and U form a partition of the set Ω



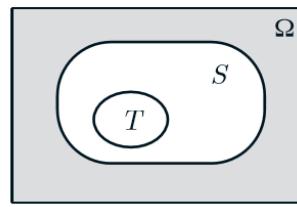
(a)



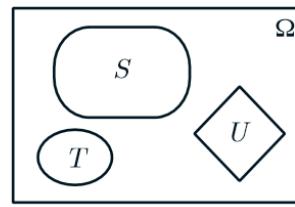
(b)



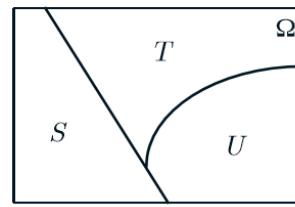
(c)



(d)



(e)



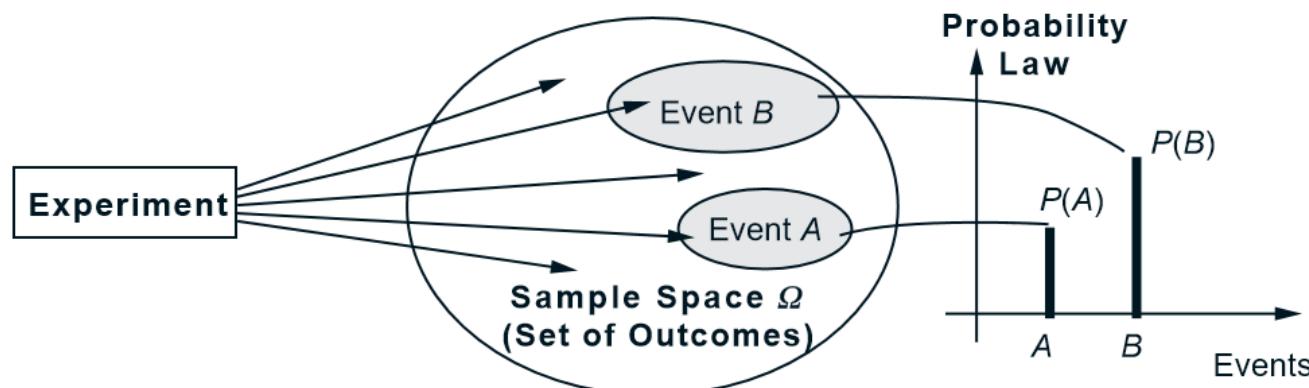
(f)

Probabilistic Models

Elements of a Probabilistic Model

The sample space Ω , which is the set of all possible outcomes of an experiment.

The **probability law**, which assigns to a set A of possible outcomes (also called an event) a nonnegative number $P(A)$ (called the probability of A) that encodes our knowledge or belief about the collective “likelihood” of the elements of A . The probability law must satisfy certain properties to be introduced shortly.



Probability Axioms

Probability Axioms

1. **(Nonnegativity)** $\mathbf{P}(A) \geq 0$, for every event A .
2. **(Additivity)** If A and B are two disjoint events, then the probability of their union satisfies

$$\mathbf{P}(A \cup B) = \mathbf{P}(A) + \mathbf{P}(B).$$

Furthermore, if the sample space has an infinite number of elements and A_1, A_2, \dots is a sequence of disjoint events, then the probability of their union satisfies

$$\mathbf{P}(A_1 \cup A_2 \cup \dots) = \mathbf{P}(A_1) + \mathbf{P}(A_2) + \dots$$

3. **(Normalization)** The probability of the entire sample space Ω is equal to 1, that is, $\mathbf{P}(\Omega) = 1$.

Conditional Probability

Properties of Conditional Probability

- The conditional probability of an event A , given an event B with $\mathbf{P}(B) > 0$, is defined by

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

and specifies a new (conditional) probability law on the same sample space Ω . In particular, all known properties of probability laws remain valid for conditional probability laws.

- Conditional probabilities can also be viewed as a probability law on a new universe B , because all of the conditional probability is concentrated on B .
- In the case where the possible outcomes are finitely many and equally likely, we have

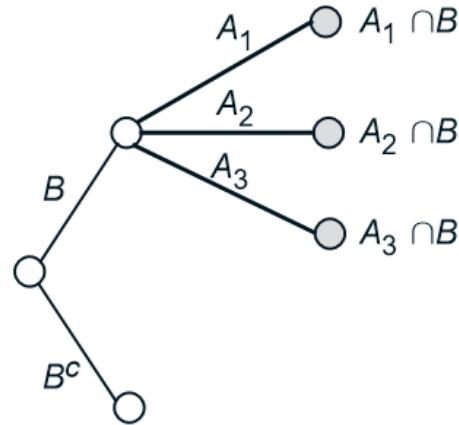
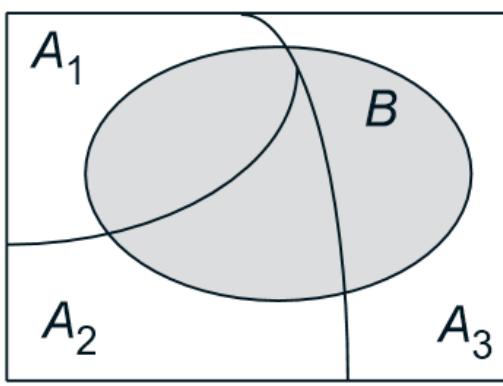
$$\mathbf{P}(A | B) = \frac{\text{number of elements of } A \cap B}{\text{number of elements of } B}.$$

Let's consider a problem of **conditional probability**:

My neighbor John has two kids.

1. He told me that one of his two kids is a boy, what is the probability that the other one is a girl.
2. If I saw one's kids is playing outside, that is a boy, what is the probability that the other one is a girl.

Total Probability Theorem



Total Probability Theorem

Let A_1, \dots, A_n be disjoint events that form a partition of the sample space (each possible outcome is included in one and only one of the events A_1, \dots, A_n) and assume that $\mathbf{P}(A_i) > 0$, for all $i = 1, \dots, n$. Then, for any event B , we have

$$\begin{aligned}\mathbf{P}(B) &= \mathbf{P}(A_1 \cap B) + \cdots + \mathbf{P}(A_n \cap B) \\ &= \mathbf{P}(A_1)\mathbf{P}(B | A_1) + \cdots + \mathbf{P}(A_n)\mathbf{P}(B | A_n).\end{aligned}$$

Independence

Bayes' Rule

Let A_1, A_2, \dots, A_n be disjoint events that form a partition of space, and assume that $\mathbf{P}(A_i) > 0$, for all i . Then, for any event B that $\mathbf{P}(B) > 0$, we have

$$\begin{aligned}\mathbf{P}(A_i | B) &= \frac{\mathbf{P}(A_i)\mathbf{P}(B | A_i)}{\mathbf{P}(B)} \\ &= \frac{\mathbf{P}(A_i)\mathbf{P}(B | A_i)}{\mathbf{P}(A_1)\mathbf{P}(B | A_1) + \dots + \mathbf{P}(A_n)\mathbf{P}(B | A_n)}\end{aligned}$$



REV. T. BAYES

Independence

- Two events A and B are said to be independent if

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

If in addition, $\mathbf{P}(B) > 0$, independence is equivalent to the condition

$$\mathbf{P}(A | B) = \mathbf{P}(A).$$

- If A and B are independent, so are A and B^c .
- Two events A and B are said to be conditionally independent, given another event C with $\mathbf{P}(C) > 0$, if

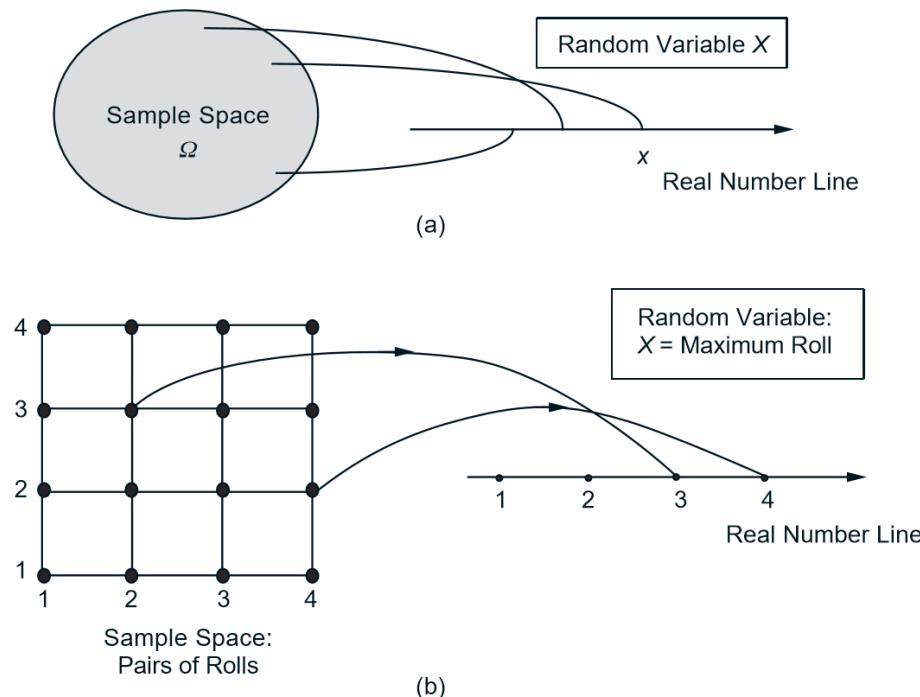
$$\mathbf{P}(A \cap B | C) = \mathbf{P}(A | C)\mathbf{P}(B | C).$$

If in addition, $\mathbf{P}(B \cap C) > 0$, conditional independence is equivalent to the condition

$$\mathbf{P}(A | B \cap C) = \mathbf{P}(A | C).$$

- Independence does not imply conditional independence, and vice versa.

Random Variable



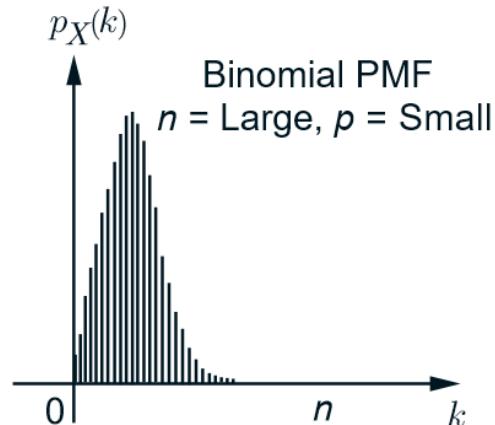
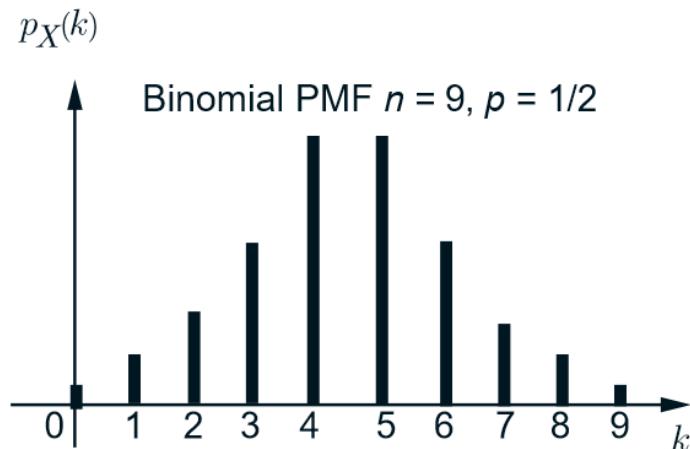
(a) Visualization of a random variable. It is a function that assigns a numerical value to each possible outcome of the experiment.

(b) An example of a random variable. The experiment consists of two rolls of a 4-sided die, and the random variable is the maximum of the two rolls. If the outcome of the experiment is (4,2), the experimental value of this random variable is 4.

Binomial Random Variable

At each toss, the coin comes up a head with probability p , and a tail with probability $1-p$, independently of prior tosses. Let X be the number of heads in the n -toss sequence. We refer to X as a binomial random variable with parameters n and p .

$$p_X(k) = \mathbf{P}(X = k) = \binom{n}{k} p^k (1-p)^{n-k}, \quad k = 0, 1, \dots, n. \quad \sum_{k=0}^n \binom{n}{k} p^k (1-p)^{n-k} = 1.$$

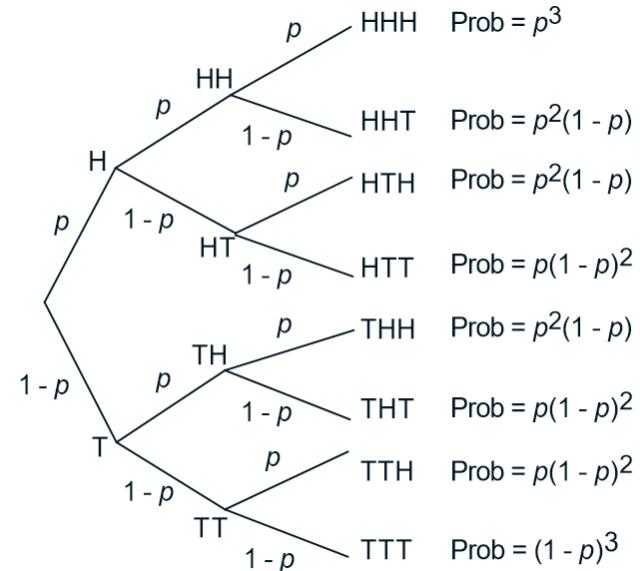


Binomial Probabilities

We showed above that the probability of any given sequence that contains k heads is $p^k(1 - p)^{n-k}$, so we have

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

where $\binom{n}{k}$ = number of distinct n -toss sequences that contain k heads.



The numbers $\binom{n}{k}$ (called “ n choose k ”) are known as the **binomial coefficients**, while the probabilities $p(k)$ are known as the **binomial probabilities**. Using a counting argument, to be given in Section 1.6, one finds that

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}, \quad k = 0, 1, \dots, n,$$

Expectation of Random Variable and Function of Random Variables

Expectation

We define the **expected value** (also called the **expectation** or the **mean**) of a random variable X , with PMF $p_X(x)$, by[†]

$$\mathbf{E}[X] = \sum_x x p_X(x).$$

Expected Value Rule for Functions of Random Variables

Let X be a random variable with PMF $p_X(x)$, and let $g(X)$ be a real-valued function of X . Then, the expected value of the random variable $g(X)$ is given by

$$\mathbf{E}[g(X)] = \sum_x g(x) p_X(x).$$

Variance

Variance

The variance $\text{var}(X)$ of a random variable X is defined by

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

and can be calculated as

$$\text{var}(X) = \sum_x (x - \mathbf{E}[X])^2 p_X(x).$$

It is always nonnegative. Its square root is denoted by σ_X and is called the **standard deviation**.

Covariance

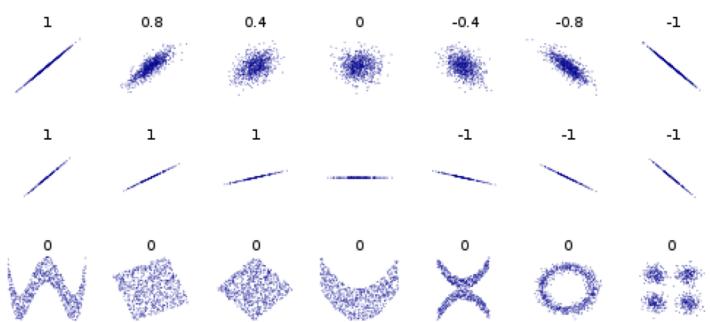
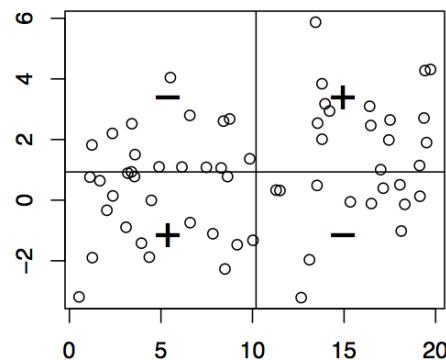
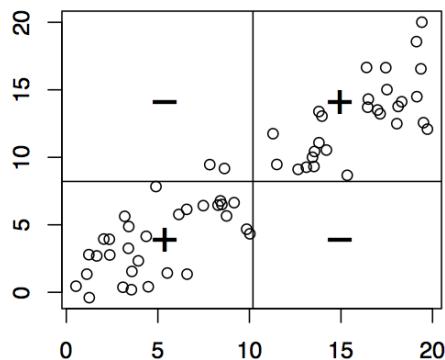
Covariance is a measure of the linear relationship between two random variables. We denote the covariance between X and Y as $\text{Cov}(X, Y)$, and it is defined to be

$$\text{Cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])]$$

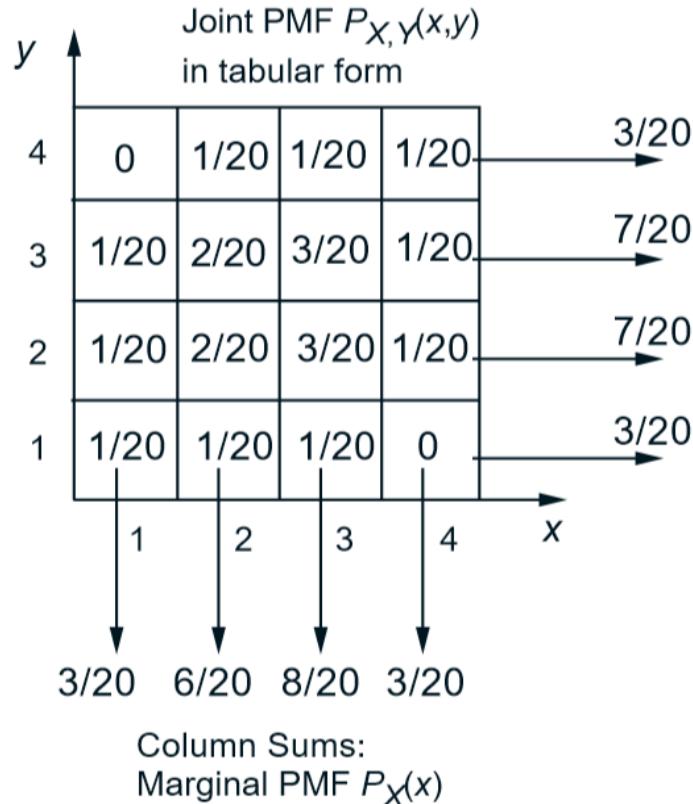
Note that the outer expectation must be taken over the joint distribution of X and Y .

Again, the linearity of expectation allows us to rewrite this as

$$\text{Cov}(X, Y) = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]$$



Joint Probability of More than One Random Variables



$$p_{X,Y,Z}(x, y, z) = \mathbf{P}(X = x, Y = y, Z = z),$$

$$p_X(x) = \sum_y \sum_z p_{X,Y,Z}(x, y, z).$$

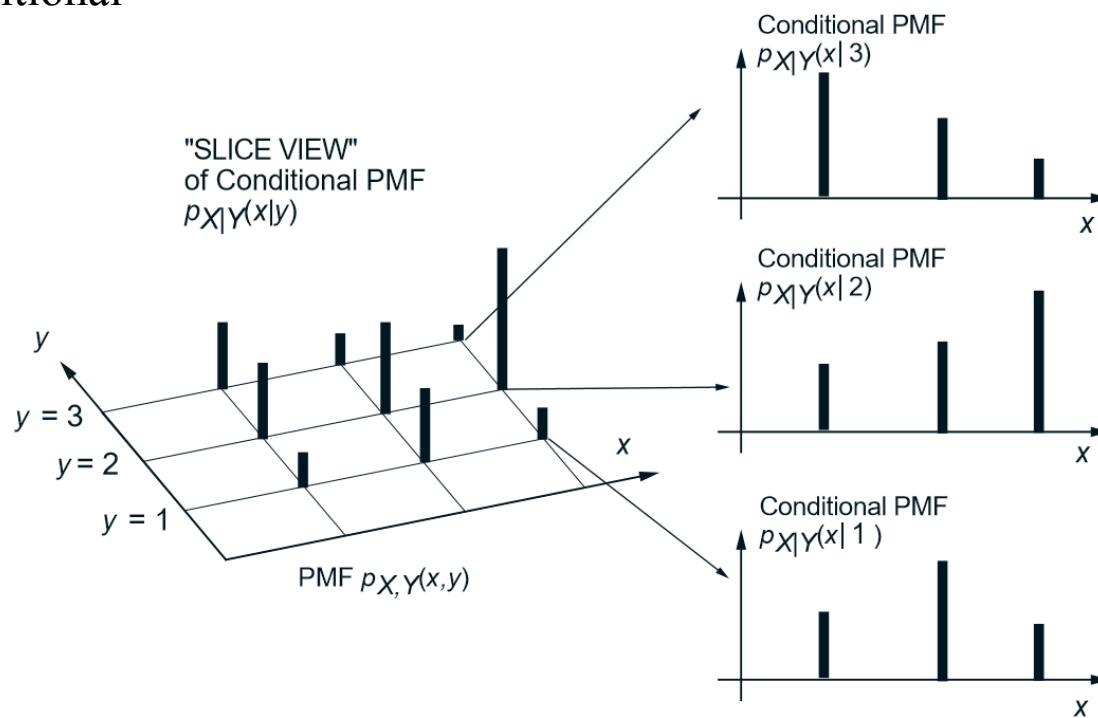
$$p_{X,Y}(x, y) = \sum_z p_{X,Y,Z}(x, y, z),$$

$$\mathbf{E}[g(X, Y, Z)] = \sum_{x,y,z} g(x, y, z) p_{X,Y,Z}(x, y, z),$$

Conditioning

The world is full of conditions, when we say independent, it usually implies “conditional independent”.

Independent or dependent?



Visualization of the conditional PMF $p_{X|Y}(x|y)$. For each y , we view the joint PMF along the slice $Y = y$ and renormalize so that

$$\sum_x p_{X|Y}(x|y) = 1.$$

Continuous PDF

a **probability distribution** is a mathematical function that, stated in simple terms, can be thought of as providing the probabilities of occurrence of different possible outcomes in an experiment. From frequency to a continuous function.

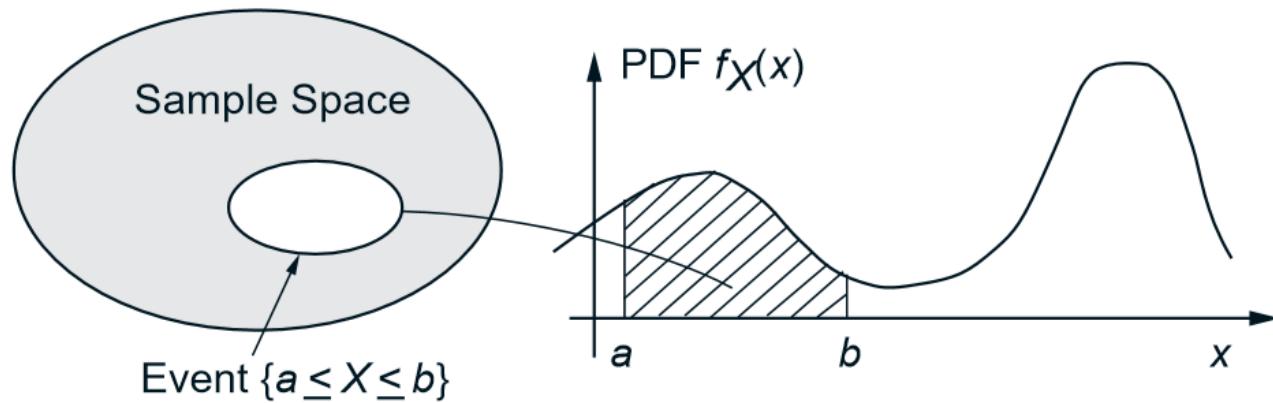
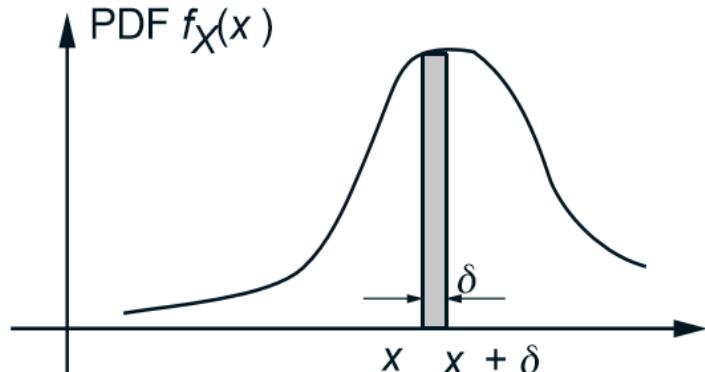
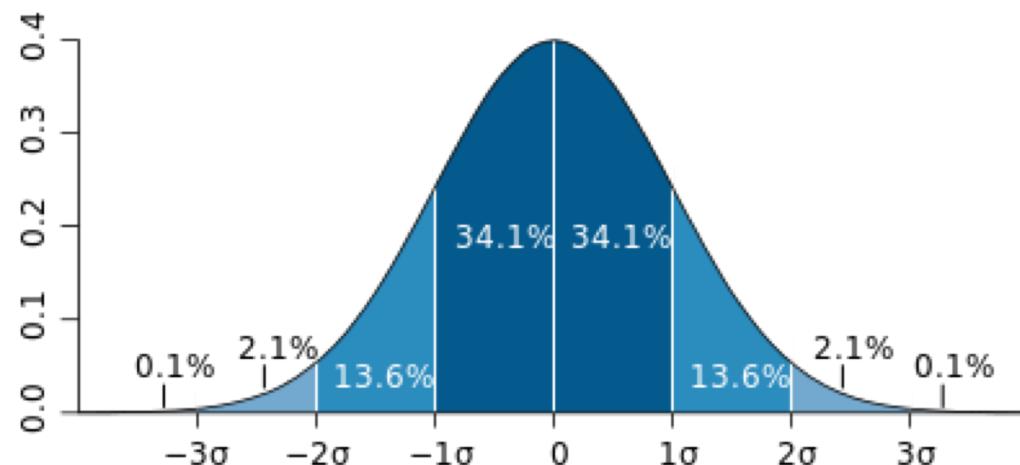


Illustration of a PDF. The probability that X takes value in an interval $[a, b]$ is $\int_a^b f_X(x) dx$, which is the shaded area in the figure.

Probability Distributions



Interpretation of the PDF $f_X(x)$ as “probability mass per unit length” around x . If δ is very small, the probability that X takes value in the interval $[x, x + \delta]$ is the shaded area in the figure, which is approximately equal to $f_X(x) \cdot \delta$.



Bernoulli Distribution

The probability distribution of any single experiment that asks a yes–no question; the question results in a Boolean valued outcome, a single bit of information whose value is success/yes/true/one with probability p and failure/no/false/zero with probability q .

If X is a random variable with this distribution, we have:

$$\Pr(X = 1) = p = 1 - \Pr(X = 0) = 1 - q.$$

The [probability mass function](#) f of this distribution, over possible outcomes k , is

$$f(k; p) = \begin{cases} p & \text{if } k = 1, \\ 1 - p & \text{if } k = 0. \end{cases}$$

This can also be expressed as

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



Jacob Bernoulli

$$\mathbb{E}(X) = p \quad \mathbb{E}[X^2] = \Pr(X = 1) \cdot 1^2 + \Pr(X = 0) \cdot 0^2 = p \cdot 1^2 + q \cdot 0^2 = p$$

$$\text{Var}[X] = \mathbb{E}[X^2] - \mathbb{E}[X]^2 = p - p^2 = p(1 - p) = pq$$

Binomial Distribution

The binomial distribution with parameters n and p is the discrete probability distribution of the number of successes in a sequence of n independent Bernoulli trials of yes–no questions.

$$(a+b)^1 = a + b$$

$$(a+b)^2 = a^2 + 2ab + b^2$$

$$(a+b)^3 = a^3 + 3a^2b + 3ab^2 + b^3$$

$$(x+y)^n = \sum_{k=0}^n \binom{n}{k} x^{n-k} y^k = \sum_{k=0}^n \binom{n}{k} x^k y^{n-k}$$

Pascal's triangle

Multinomial Distribution

$$f(x_1, \dots, x_k; n, p_1, \dots, p_k) = \Pr(X_1 = x_1 \text{ and } \dots \text{ and } X_k = x_k)$$
$$= \begin{cases} \frac{n!}{x_1! \cdots x_k!} p_1^{x_1} \times \cdots \times p_k^{x_k}, & \text{when } \sum_{i=1}^k x_i = n \\ 0 & \text{otherwise,} \end{cases}$$

Suppose one does an experiment of extracting n balls of k different colours from a bag, replacing the extracted ball after each draw. Balls from the same colour are equivalent.

The probability mass function can be expressed using the gamma function as:

$$f(x_1, \dots, x_k; p_1, \dots, p_k) = \frac{\Gamma(\sum_i x_i + 1)}{\prod_i \Gamma(x_i + 1)} \prod_{i=1}^k p_i^{x_i}$$

Where gamma function is an extension of factorial:
 $\Gamma(n) = (n - 1)!$
 $\Gamma(z) = \int_0^\infty x^{z-1} e^{-x} dx$

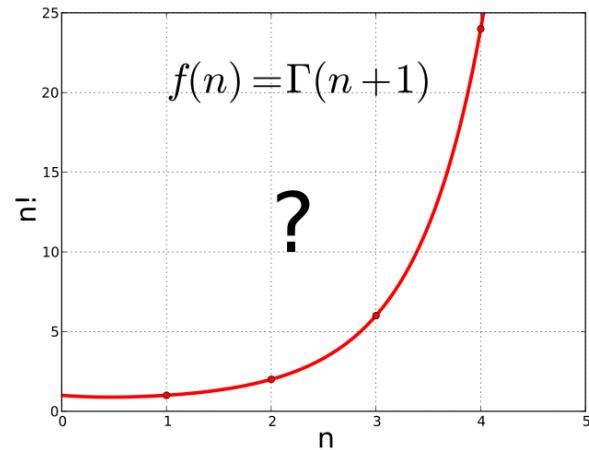
More on Gamma Function

It is easy graphically to interpolate the factorial function to non-integer values, but is there a formula that describes the resulting curve?

$$\begin{aligned}\Gamma(z+1) &= \int_0^\infty x^z e^{-x} dx \\ &= [-x^z e^{-x}]_0^\infty + \int_0^\infty z x^{z-1} e^{-x} dx \\ &= \lim_{x \rightarrow \infty} (-x^z e^{-x}) - (0e^{-0}) + z \int_0^\infty x^{z-1} e^{-x} dx\end{aligned}$$

Recognizing that as $x \rightarrow \infty, -x^z e^{-x} \rightarrow 0$,

$$\Gamma(z+1) = z \int_0^\infty x^{z-1} e^{-x} dx = z\Gamma(z)$$



$$\begin{aligned}\Gamma(1) &= \int_0^\infty x^{1-1} e^{-x} dx = [-e^{-x}]_0^\infty \\ &= \lim_{x \rightarrow \infty} (-e^{-x}) - (-e^{-0}) = 0 - (-1) = 1\end{aligned}$$

Given that $\Gamma(1) = 1$ and $\Gamma(n+1) = n\Gamma(n)$

$$\Gamma(n) = 1 \cdot 2 \cdot 3 \cdots (n-1) = (n-1)!$$

Gamma Distribution

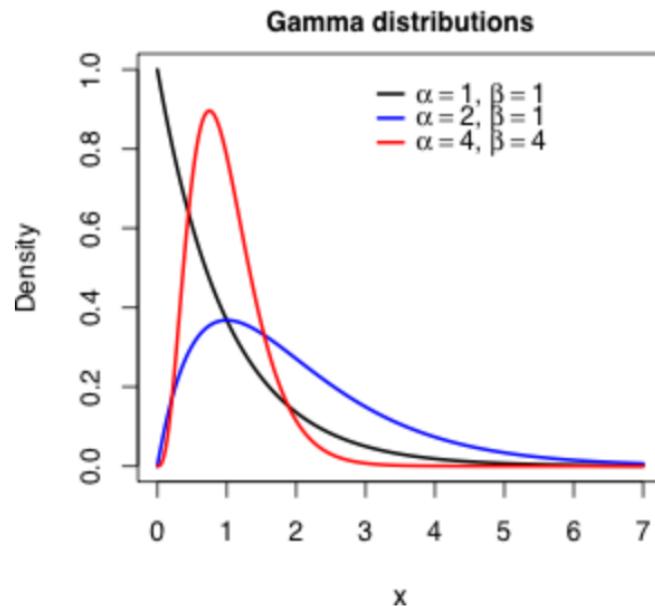
Like the lognormal the gamma distribution is unbounded on the right, defined for only positive X , and tends to yield skewed distributions.

$$X \sim \Gamma(\alpha, \beta) \equiv \text{Gamma}(\alpha, \beta)$$

The corresponding **probability density function** in the shape-rate parametrization is

$$f(x; \alpha, \beta) = \frac{\beta^\alpha x^{\alpha-1} e^{-\beta x}}{\Gamma(\alpha)} \quad \text{for } x > 0 \text{ and } \alpha, \beta > 0,$$

where $\Gamma(\alpha)$ is a complete gamma function.

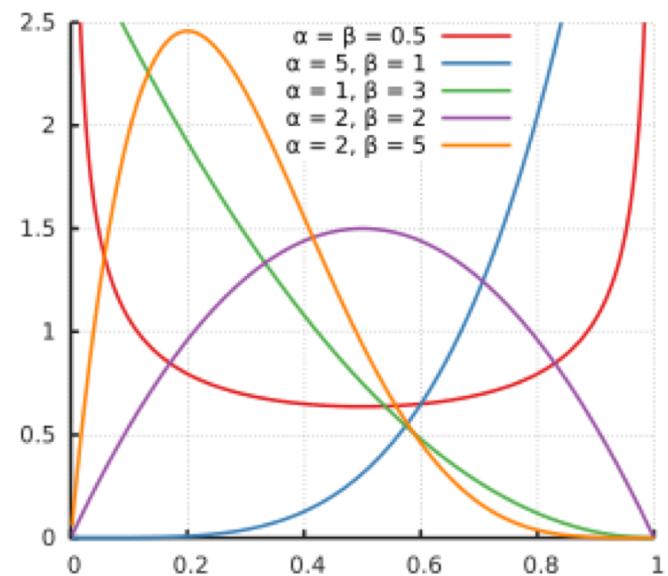


The gamma distribution is widely used as a conjugate prior in Bayesian statistics. It is the conjugate prior for the precision of a normal distribution. It is also the conjugate prior for the exponential distribution.

Beta Distribution

It is bounded on both sides. In this respect it resembles the binomial distribution. The standard beta distribution is constrained so that its domain is the interval $(0, 1)$.

$$\begin{aligned}f(x; \alpha, \beta) &= \text{constant} \cdot x^{\alpha-1}(1-x)^{\beta-1} \\&= \frac{x^{\alpha-1}(1-x)^{\beta-1}}{\int_0^1 u^{\alpha-1}(1-u)^{\beta-1} du} = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1}(1-x)^{\beta-1} \\&= \frac{1}{B(\alpha, \beta)} x^{\alpha-1}(1-x)^{\beta-1}\end{aligned}$$



The beta function, B is a normalization constant to ensure that the total probability integrates to 1.

Poisson Distribution

A discrete random variable X is said to have a Poisson distribution with parameter $\lambda > 0$, if, for $k = 0, 1, 2, \dots$,

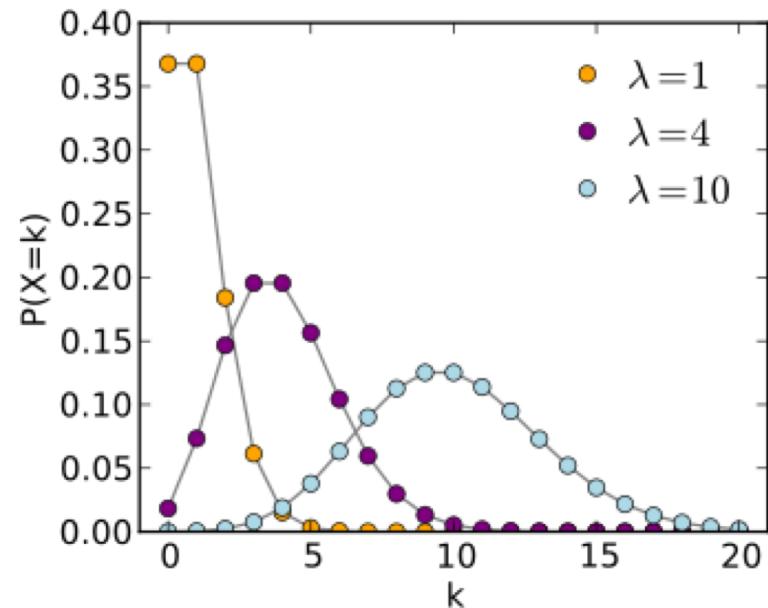
For example, on a particular river, overflow floods occur once every 100 years on average. Calculate the probability of $k = 0, 1, 2, 3, 4, 5$, or 6 overflow floods in a 100-year interval, assuming the Poisson model is appropriate. Because the average event rate is one overflow flood per 100 years, $\lambda = 1$, so that:

$$P(k \text{ overflow floods in 100 years}) = \frac{\lambda^k e^{-\lambda}}{k!} = \frac{1^k e^{-1}}{k!}$$

$$P(k = 0 \text{ overflow floods in 100 years}) = \frac{1^0 e^{-1}}{0!} = \frac{e^{-1}}{1} = 0.368$$

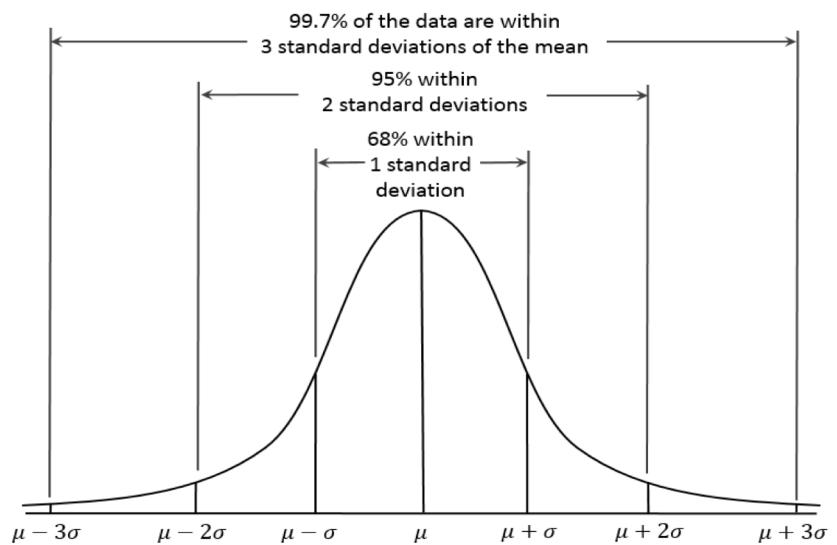
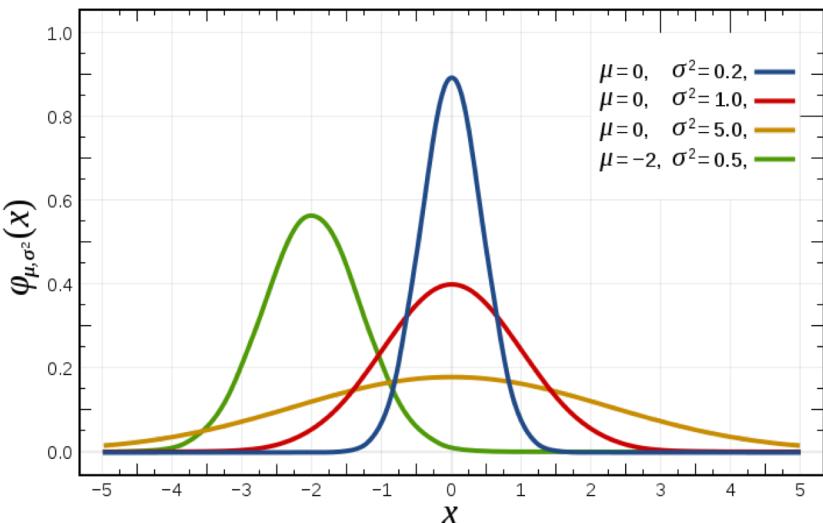
$$P(k = 1 \text{ overflow flood in 100 years}) = \frac{1^1 e^{-1}}{1!} = \frac{e^{-1}}{1} = 0.368$$

$$P(k = 2 \text{ overflow floods in 100 years}) = \frac{1^2 e^{-1}}{2!} = \frac{e^{-1}}{2} = 0.184$$



$$f(k; \lambda) = \Pr(X = k) = \frac{\lambda^k e^{-\lambda}}{k!}$$

Gaussian (Normal) Distribution



The probability density of the normal distribution is:

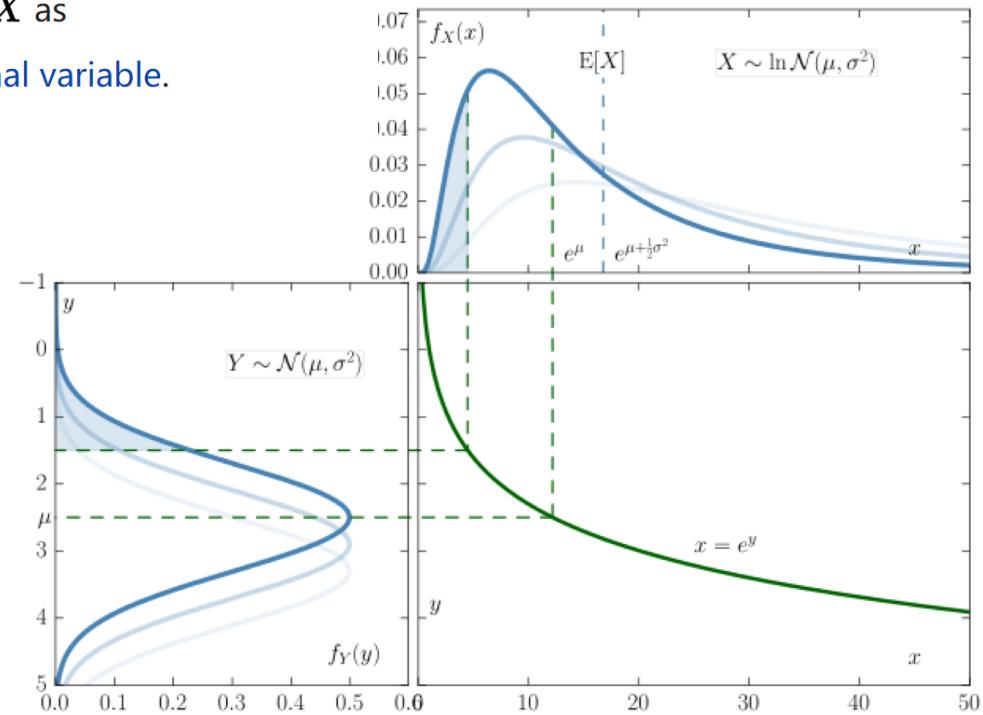
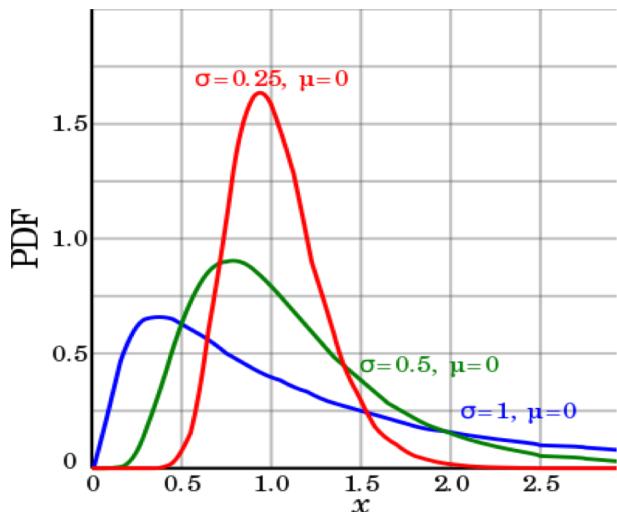
$$f(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

- μ is the mean or expectation of the distribution
- σ is the standard deviation
- σ^2 is the variance

Log Normal Distribution

Given a log-normally distributed random variable X and two parameters μ and σ that are, respectively, the mean and standard deviation of the variable's natural logarithm, then the logarithm of X is normally distributed, and we can write X as

$$X = e^{\mu + \sigma Z} \quad \text{with } Z \text{ a standard normal variable.}$$



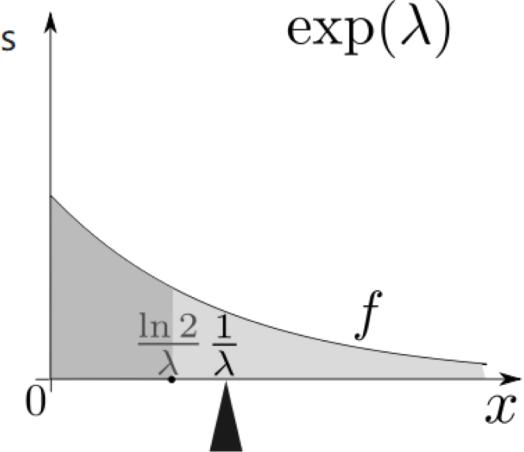
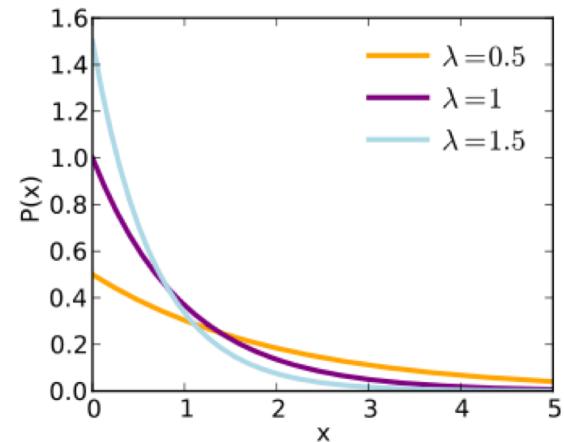
Exponential Distribution

The exponential distribution (also known as negative exponential distribution) is the probability distribution that describes the time between events in a Poisson process, i.e. a process in which events occur continuously and independently at a constant average rate. It is a particular case of the gamma distribution. It is the continuous analogue of the geometric distribution, and it has the key property of being memoryless.

The probability density function (pdf) of an exponential distribution is

$$f(x; \lambda) = \begin{cases} \lambda e^{-\lambda x} & x \geq 0, \\ 0 & x < 0. \end{cases} \quad E[X] = \frac{1}{\lambda} \quad \text{Var}[X] = \frac{1}{\lambda^2}$$

The exponential distribution occurs naturally when describing the lengths of the inter-arrival times in a homogeneous Poisson process.



Bayesian Examples

Marie is getting married tomorrow at an outdoor ceremony in the desert. In recent years, it has rained only **5 days each year**.

Unfortunately, the weatherman is forecasting rain for tomorrow. When it actually rains, the weatherman has forecast rain **90% of the time**. When it doesn't rain, he has forecast rain **10% of the time**. What is the probability it will rain on the day of Marie's wedding?

Event A: The weatherman has forecast rain.

Event B: It rains.

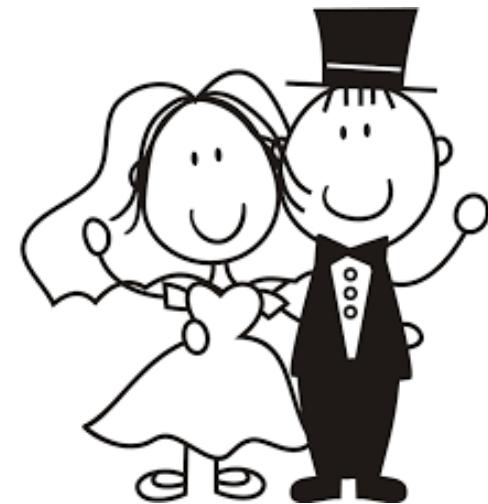
We know:

1. $p(B) = 5 / 365 = 0.0137$ [It rains 5 days out of the year.]
2. $p(\text{not } B) = 360 / 365 = 0.9863$
3. $p(A | B) = 0.9$ [When it rains, the weatherman has forecast rain 90% of the time.]
4. $p(A | \text{not } B) = 0.1$ [When it does not rain, the weatherman has forecast rain 10% of the time.]

Bayesian Example

We want to know $p(B | A)$, the probability it will rain on the day of Marie's wedding, given a forecast for rain by the weatherman. The answer can be determined from Bayes Rule:

1. $p(B | A) = p(A | B) \cdot p(B) / p(A)$
2. $p(A) = p(A | B) \cdot p(B) + p(A | \text{not } B) \cdot p(\text{not } B)$
 $= (0.9)(0.014) + (0.1)(0.986) = 0.111$
1. $p(B | A) = (0.9)(0.0137) / 0.111 = 0.111$



The **special** coin problem.

Simpson's Paradox

Simpson's paradox, or the Yule–Simpson effect, is a phenomenon in probability and statistics, in which a trend appears in different groups of data but disappears or reverses when these groups are combined.

Department	Female Applicants	Female Admitted	%	Male Applicants	Male Admitted	%	All Applicants	All Admitted	Overall %
Business School	100	49	49%	20	15	75%	120	64	53.3%
Law School	20	1	5%	100	10	10%	120	11	9.2%
Both	120	50	42%	120	25	21%	240	75	31.3%

Suppose two people, Lisa and Bart, each edit articles for two weeks. In the first week, Lisa fails to improve the only article she edited, and Bart improves 1 of the 4 articles he edited. In the second week, Lisa improves 3 of 4 articles she edited, while Bart improves the only article he edited.

	Week 1	Week 2	Total
Lisa	0/1	3/4	3/5
Bart	1/4	1/1	2/5

Acknowledgement

This slide of this class is modified from Lecture Notes of Dimitri P. Bertsekas and John N. Tsitsiklis – Introduction to Probability, MIT, 2000. & Wikipedia.

UNC Lecture Notes on Ecological Stats:

<https://www.unc.edu/courses/2008fall/ecol/563/001/docs/lectures/lecture3.htm>

Jeff Howbert Introduction to Machine Learning Winter 2012

Mathematics for Machine Learning Garrett Thomas

<http://gwthomas.github.io/docs/math4ml.pdf>

<https://rorasa.wordpress.com/2012/05/13/l0-norm-l1-norm-l2-norm-l-infinity-norm/>

<https://www.bilibili.com/video/av10852829/>