



# PROBABILITY

Damian Klimke

# Outline

- What is probability? And Why?
- Random Variables
- Probability Distributions
- Marginal Probability ( $\sigma$ )
- Conditional Probability
- Chain Rule of CP
- Independence and Conditional Independence
- Expectation, Variance and Covariance

- Common Probability Distributions
- Useful Properties of Common Functions
- Bayes' Rules
- Information Theory
- Structured Probabilistic Models
- Monte Carlo
- Markov Chains
- What for?

# What is probability? And Why?

- Life is probability
- Why do we study probability theory?
  - An effective model of uncertainty
  - Decision Making under uncertainty
- Examples:
  - Measurement sensors
  - Waiting time at a Bank's teller.
  - Value of a stock at a given day.
  - Outcome of a medical procedure.
  - A customer buying behavior.
  - **ChatBot, music etc.**
- One Decision Making Process: Collect Data, Model the Phenomenon, Extrapolate and make decisions.

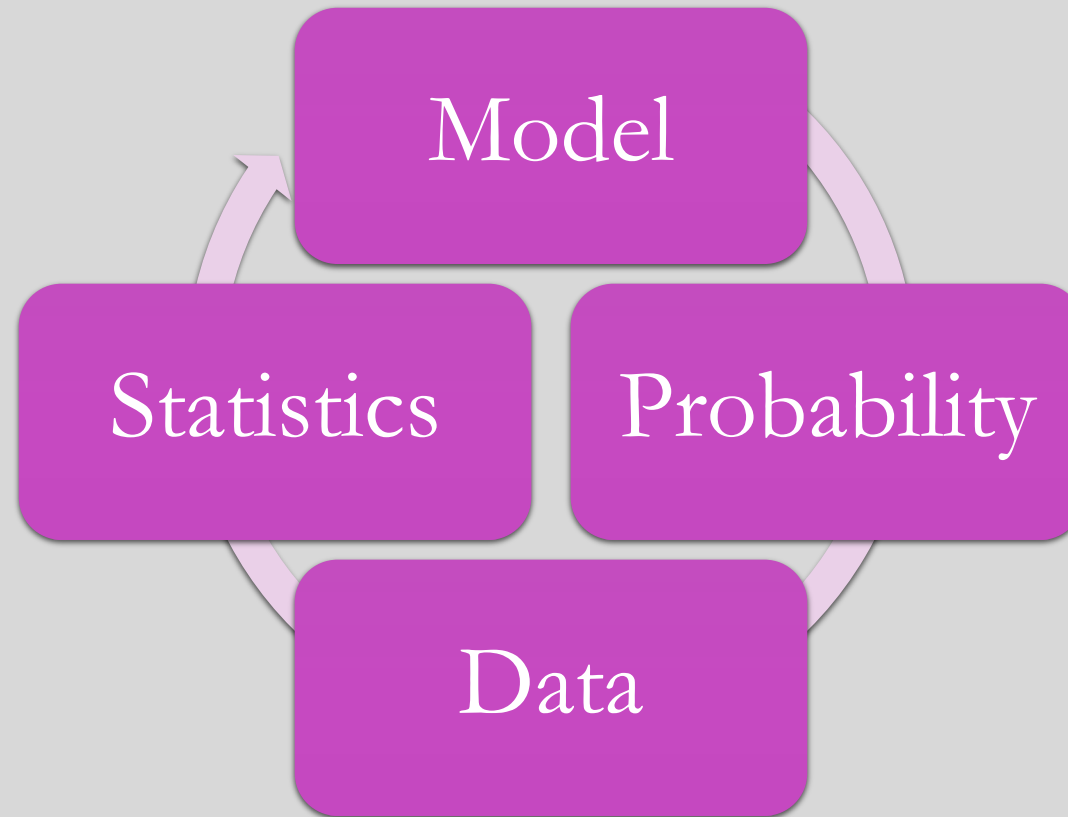
# What is probability? And Why?

- **Probability** is a measure quantifying the likelihood that events will occur. ... Probability quantifies as a number between 0 and 1, where, roughly speaking, 0 indicates impossibility and 1 indicates certainty. The higher the probability of an event, the more likely it is that the event will occur.
- A simple example is the tossing of a fair (unbiased) coin. Since the coin is fair, the two outcomes ("heads" and "tails") are both equally probable; the probability of "heads" equals the probability of "tails"; and since no other outcomes are possible, the probability of either "heads" or "tails" is  $1/2$  (which could also be written as 0.5 or 50%).
- These concepts have been given an axiomatic mathematical formalization in probability theory, which is used widely in such areas of study as mathematics, statistics, finance, gambling, science (in particular physics), artificial intelligence/machine learning, computer science, game theory, and philosophy to, for example, draw inferences about the expected frequency of events. Probability theory is also used to describe the underlying mechanics and regularities of complex systems.

# What is the difference? – Likelihood vs. Probability?

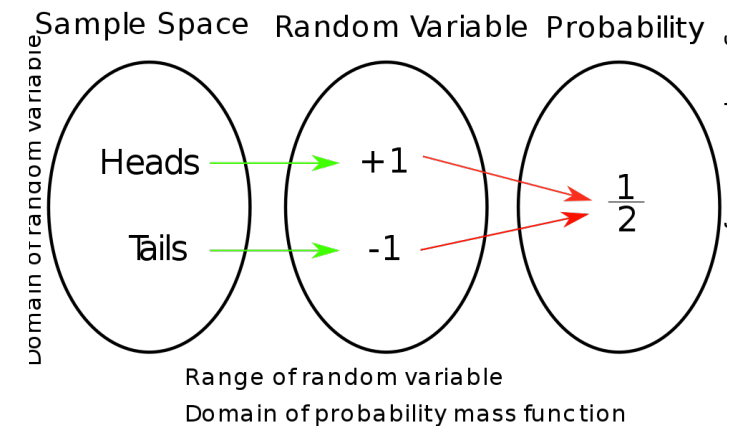
- Please search.
- <https://www.youtube.com/watch?v=pYxNSUDSFH4>

# What is probability? And Why?



# Random Variables

- Example:
- In probability and statistics, a **random variable**, **random quantity**, **aleatory variable**, or **stochastic variable** is described informally as a variable whose values depend on outcomes of a random phenomenon.<sup>[1]</sup> The formal mathematical treatment of random variables is a topic in probability theory. In that context, a random variable is understood as a measurable function defined on a probability space whose outcomes are typically real numbers.
- Help to express our thoughts mathematically





# Helpful links – Random Variable

- <https://towardsdatascience.com/understanding-random-variable-a618a2e99b93>

# Random Variables

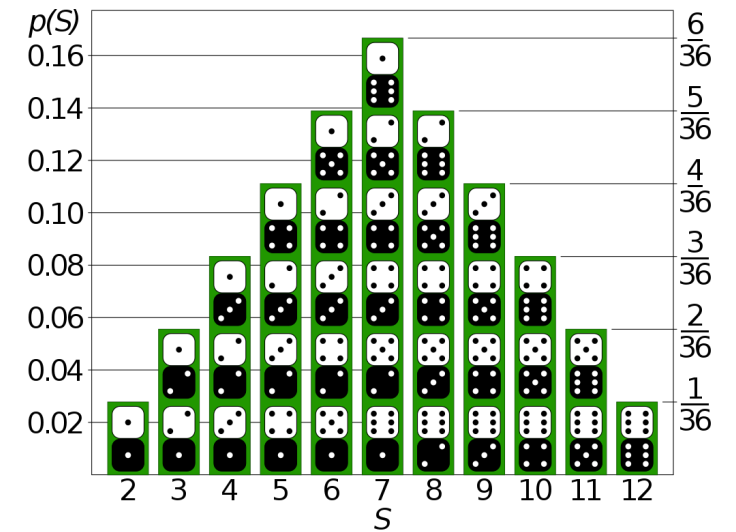
- Find an example (in your field)!

# Probability Distributions

- In probability theory and statistics, a **probability distribution** is a mathematical function that provides the probabilities of occurrence of different possible outcomes in an experiment. In more technical terms, the probability distribution is a description of a random phenomenon in terms of the probabilities of events. For instance, if the random variable  $X$  is used to denote the outcome of a coin toss ("the experiment"), then the probability distribution of  $X$  would take the value 0.5 for  $X = \text{heads}$ , and 0.5 for  $X = \text{tails}$  (assuming the coin is fair). Examples of random phenomena can include the results of an experiment or survey.

# Probability Distributions (Dice)

- The probability mass function (pmf)  $p(S)$  specifies the probability distribution for the sum  $S$  of counts from two dice. For example, the figure shows that  $p(11) = 2/36 = 1/18$ . The pmf allows the computation of probabilities of events such as  $P(S > 9) = 1/12 + 1/18 + 1/36 = 1/6$ , and all other probabilities in the distribution.



# Question: by Goldman Sachs

- You're given a fair coin. You flip the coin until either **Heads Heads Tails** (HHT) or **Heads Tails Tails** (HTT) appears. Is one more likely to appear first? If so, which one and with what probability?

# Probability Distributions – Discrete Variables & Mass Functions

- A **discrete probability distribution** is a probability distribution that can take on a countable number of values. For the probabilities to add up to 1, they have to decline to zero fast enough. For example, if  $P(X = n) = \frac{1}{2^n}$  for  $n = 1, 2, \dots$  the sum of probabilities would be  $\frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \dots = 1$ .
- In probability and statistics, a **probability mass function (PMF)** is a function that gives the probability that a discrete random variable is exactly equal to some value. The probability mass function is often the primary means of defining a discrete probability distribution, and such functions exist for either scalar or multivariate random variables whose domain is discrete.
- Definition:
- Suppose that  $X: S \mapsto A$  for  $A \subseteq \mathbb{R}$  is a discrete random variable for the sample space  $S$ . The **PMF**  $f_X: A \mapsto [0,1]$  for  $X$  is defined as:
  - $f_X(x) = P(X = x) = P(\{s \in S: X(s) = x\})$
- Think of it as mass!  $\sum_{x \in A} f_X(x) = 1$

# Probability Distributions – Continuous Variables & Density Functions

- A **continuous probability distribution** is a probability distribution with a cumulative distribution function that is absolutely continuous. Equivalently, it is a probability distribution on the real numbers that is absolutely continuous with respect to Lebesgue measure. Such distributions can be represented by their probability density functions. If the distribution of  $X$  is continuous, then  $X$  is called a **continuous random variable**. There are many examples of continuous probability distributions: normal, uniform, chi-squared, and others.
- Formally, if  $X$  is a continuous random variable, then it has a probability density function  $f(x)$ , and therefore its probability of falling into a given interval, say  $[a, b]$ , is given by the integral  $P[a \leq X \leq b] = \int_a^b f(x)dx$

# Question: by Linkedin

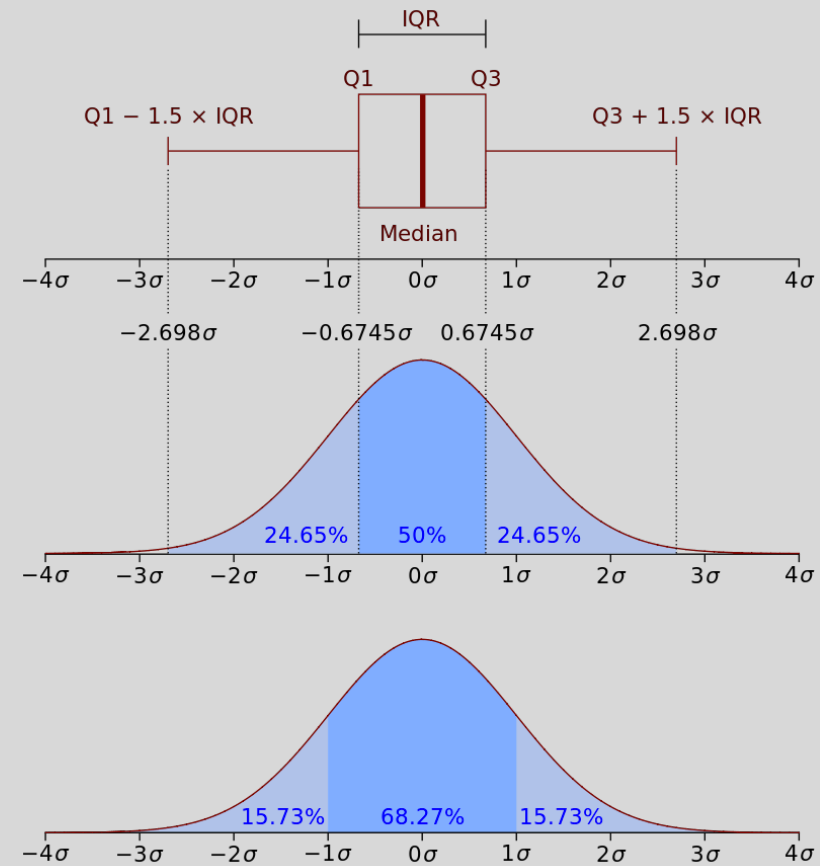
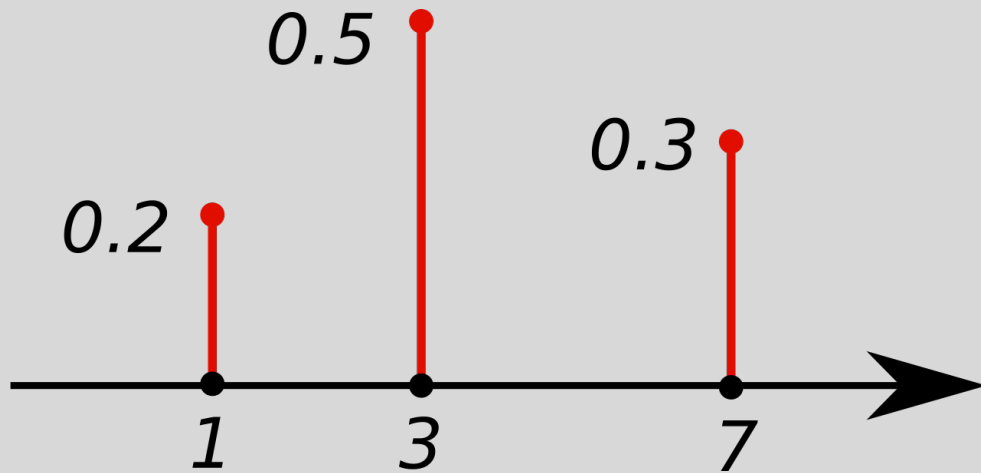
- Suppose we have two coins. One is fair and the other biased where the probability of it coming up heads is  $3/4$ .
- Let's say we select a coin at random and flip it two times. What is the probability that both flips result in the same side?



# Question by LinkedIn

- Imagine a deck of 500 cards numbered from 1 to 500. If all the cards are shuffled randomly and you are asked to pick three cards, one at a time, what's the probability of each subsequent card being larger than the previous drawn card?

# Probability Distributions – What is what?



# Probability Distributions

- Find examples and present them shortly to the class?
- [https://en.wikipedia.org/wiki/List\\_of\\_probability\\_distributions](https://en.wikipedia.org/wiki/List_of_probability_distributions)

# Marginal Probability

- Sometimes we know the probability distribution over a set of variables and we want to know the probability distribution over just a subset of them. The probability distribution over the subset is known as the marginal probability distribution.

(<https://www.deeplearningbook.org/contents/prob.html>)

- For example, suppose we have discrete random variables  $x$  and  $y$ , and we know  $P(x, y)$ . We can find  $P(x)$  with the sum rule:
- $\forall x \in \mathcal{X}, P(x = x) = \sum_y P(x = x, y = y)$
- For continuous variables we need the integration:
- $p(x) = \int p(x, y) dy$

| $y \backslash X$     | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $p_Y(y) \downarrow$ |
|----------------------|-------|-------|-------|-------|---------------------|
| $y_1$                | 4/32  | 2/32  | 1/32  | 1/32  | 8/32                |
| $y_2$                | 3/32  | 6/32  | 3/32  | 3/32  | 15/32               |
| $y_3$                | 9/32  | 0     | 0     | 0     | 9/32                |
| $p_X(x) \rightarrow$ | 16/32 | 8/32  | 4/32  | 4/32  | 32/32               |

Joint and marginal distributions of a pair of discrete random variables,  $X$  and  $Y$ , having nonzero **mutual information**  $I(X; Y)$ . The values of the joint distribution are in the 3x4 rectangle; the values of the marginal distributions are along the right and bottom margins.

# Conditional Probability

- In many cases, we are interested in the probability of some event, given that some other event has happened. This is called a conditional probability. We denote the conditional probability that  $y=\boldsymbol{y}$  given  $x=\boldsymbol{x}$  as  $P(y=\boldsymbol{y}|\boldsymbol{x}=\boldsymbol{x})$ . This conditional probability can be computed with the formula

$$P(y = \boldsymbol{y}|\boldsymbol{x} = \boldsymbol{x}) = \frac{P(y = \boldsymbol{y}, \boldsymbol{x} = \boldsymbol{x})}{P(\boldsymbol{x} = \boldsymbol{x})}$$

For:  $P(\boldsymbol{x} = \boldsymbol{x}) > 0$

# In contrast to:

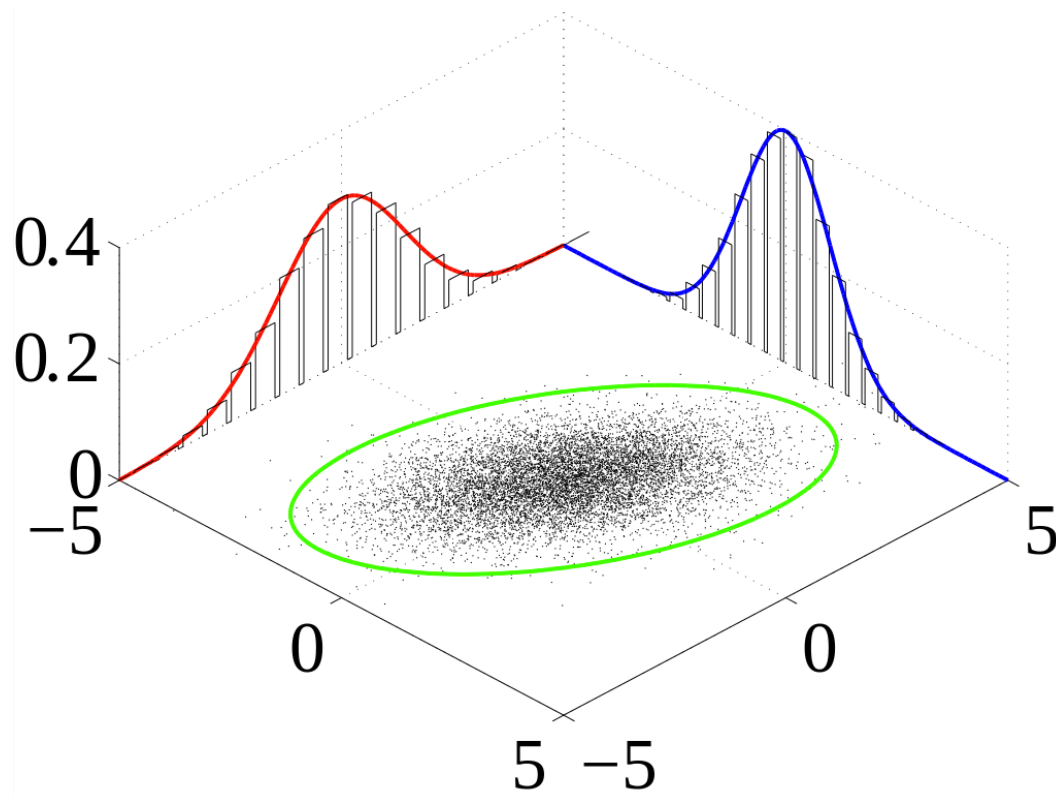
- In probability theory and statistics, the **marginal distribution** of a subset of a collection of random variables is the probability distribution of the variables contained in the subset. It gives the probabilities of various values of the variables in the subset without reference to the values of the other variables. This contrasts with a conditional distribution, which gives the probabilities contingent upon the values of the other variables.

# Example – Conditional Probability

- We consider the roll of a fair die and let  $X=1$  if the number is even (e.g. 2 or 4) and  $X=0$  otherwise. Further, we let  $Y=1$  if the number is prime (e.g. 2 or 3) and  $Y=0$  otherwise.

|   | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| X | 0 | 1 | 0 | 1 | 0 | 1 |
| Y | 0 | 1 | 1 | 0 | 1 | 0 |

- Then the unconditional probability the  $X=1$  is  $3/6=1/2$  (since there are six possible rolls of the die, of which three are even). But the probability that  $X=1$  conditional to  $Y=1$  is  **$1/3$**  (since there are three possible prime number rolls (2,3,5) of which one is even).
- **WHY?**



## JOINT DISTRIBUTION

Given random variables  $X, Y, \dots$ , that are defined on a probability space, the joint probability distribution for  $X, Y, \dots$  is a probability distribution that gives the probability that each of  $X, Y, \dots$  falls in any particular range or discrete set of values specified for that variable. In the case of only two random variables, this is called a **bivariate distribution**, but the concept generalizes to any number of random variables, giving a **multivariate distribution**.



# Chain Rule of CP

- For two events: A & B
  - $P(A \cap B) = P(A|B) * P(B)$
- Find an example?
- More than two events  $A_1, \dots, A_n$  :
  - $P(A_n \cap \dots \cap A_1) = P(A_n | A_{n-1} \cap \dots \cap A_1) * P(A_{n-1} \cap \dots \cap A_1)$
  - Also:

$$P(A_n \cap \dots \cap A_1) = \prod_{k=1}^n P(A_k | \cap_{j=1}^{k-1} A_j)$$

# Exercise – Chain Rule:

- With 4 events ( $n=4$ ), the chain rule is???
- $P(A_4 \cap A_3 \cap A_2 \cap A_1)$
- $= P(A_4 | A_3 \cap A_2 \cap A_1) * P(A_3 \cap A_2 \cap A_1)$
- $= P(A_4 | A_3 \cap A_2 \cap A_1) * P(A_3 | A_2 \cap A_1) * P(A_2 \cap A_1)$
- $= P(A_4 | A_3 \cap A_2 \cap A_1) * P(A_3 | A_2 \cap A_1) * P(A_2 | A_1) * P(A_1)$

# Independence

- Independence:
  - Two events are independent( $A \perp B$ ), **if and only if** their joint probability equals the product of their probabilities:
  - $P(A \cap B) = P(A)P(B)$
- Short: the events do not effect each other.

# Independence

- Find an example?

# Conditional Independence

- Conditional independence:
  - A and B are conditionally independent given C, **if and if only**,  $P(A \cap B|C) = P(A|C)P(B|C)$
  - $(A \perp B)|C \Leftrightarrow P(A \cap B|C) = P(A|C)P(B|C)$
  - Find an example in your field:

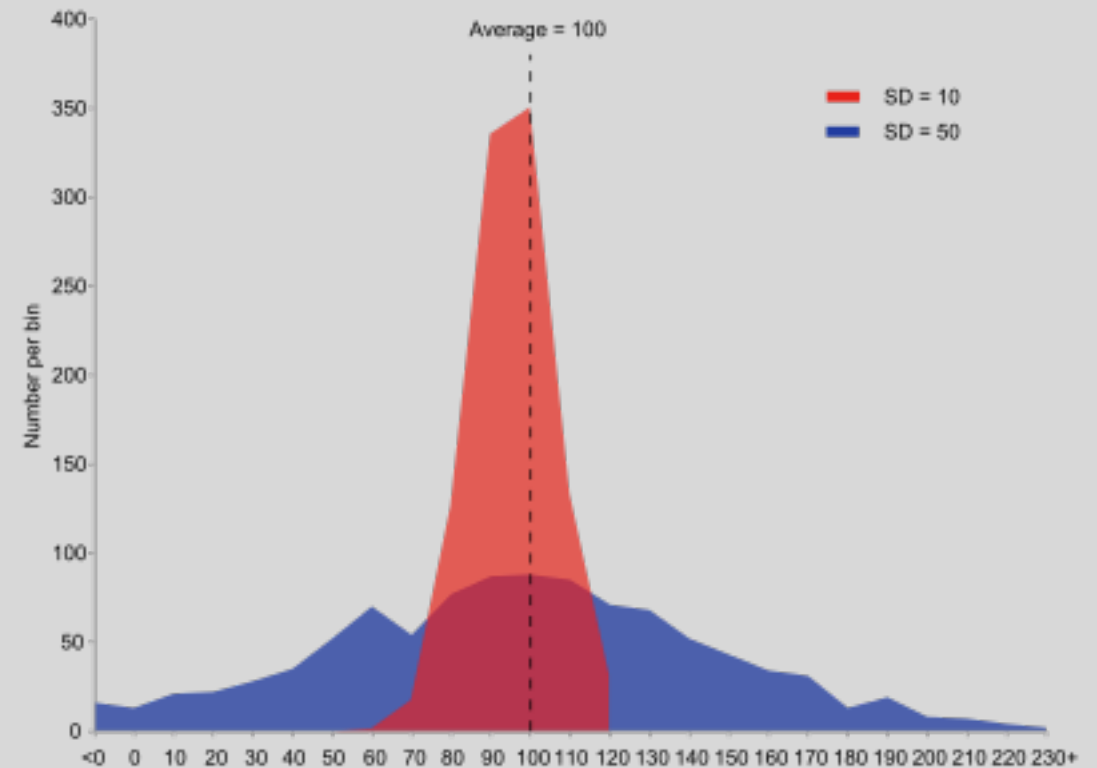
# Expectation

- Expectation or Expected Value
  - Intuitively, a random variable's expected value represents the average of a large number of independent realizations of the random variable.
- $X$  is a random variable with a finite number of outcomes  $(x_1, x_2, \dots, x_k)$  occurring with probabilities  $(p_1, p_2, \dots, p_k)$ .
  - $E[X] = \sum_{i=1}^k x_i p_i = x_1 p_1 + x_2 p_2 + \dots + x_k p_k$  for discrete variables
  - For continuous variables we compute it with an integral.

Find an example.

# Variance

- Variance: The variance of a standard  $X$  (function) is the expected value of the squared deviation from the mean  $X, \mu = E[X]$ :
- $Var(X) = E[(X - \mu)^2]$
- Example?? E.g.: Fair Die
- The  $\sqrt{Var(X)}$  is called the “standard deviation”



# Covariance

- The covariance gives some sense of how much two values are linearly related to each other, as well as the scale of these variables:
- X and Y (random variables (jointly distributed)): (If the covariance is positive, there is a positive correlation. In fact, correlation coefficients can simply be understood as a normalized version of covariance.)
  - $cov(X, Y) = E[(X - E[X])(Y - E[Y])]$
- Example or need?
- What is a covariance matrix and what for?
- <https://datascienceplus.com/understanding-the-covariance-matrix/>



# Common Probability Distributions

# 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 experiments, each asking a yes–no question, : success/yes/true/one (with probability  $p$ ) or failure/no/false/zero (with probability  $q = 1 - p$ ). A single success/failure experiment is also called a Bernoulli trial or Bernoulli experiment and a sequence of outcomes is called a Bernoulli process; for a single trial, i.e.,  $n = 1$ , the binomial distribution is a Bernoulli distribution. The binomial distribution is the basis for the popular binomial test of statistical significance.
- The binomial distribution is frequently used to model the number of successes in a sample of size  $n$  drawn with replacement from a population of size  $N$ . If the sampling is carried out without replacement, the draws are not independent and so the resulting distribution is a hypergeometric distribution, not a binomial one. However, for  $N$  much larger than  $n$ , the binomial distribution remains a good approximation, and is widely used.

# Binomial Distribution

- $n \in \{0,1,2, \dots\}$  – *number of trials*
- $p \in [0,1]$  – *success probability for each trial*
- $k \in \{0,1,2, \dots, n\}$  – *number of successes*
- $f(k, n, p) = P(k; n, p) = P(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}$

For  $k=0,1,2,\dots,n$ , where:

- $\binom{n}{k} = \frac{n!}{k!(n-k)!}$

# Question: by Google

- Let's say we are given a biased coin that comes up with heads 30% of the time when tossed.
- What is the probability of the coin landing as heads exactly 5 times out of 6 tosses?

# Bernoulli Distribution

- (Yes/No – distribution)
  - $P(X = 1) = p; P(X = 0) = 1 - p = q$

The PMF  $f$  of distribution, over possible outcomes  $k$ , is:

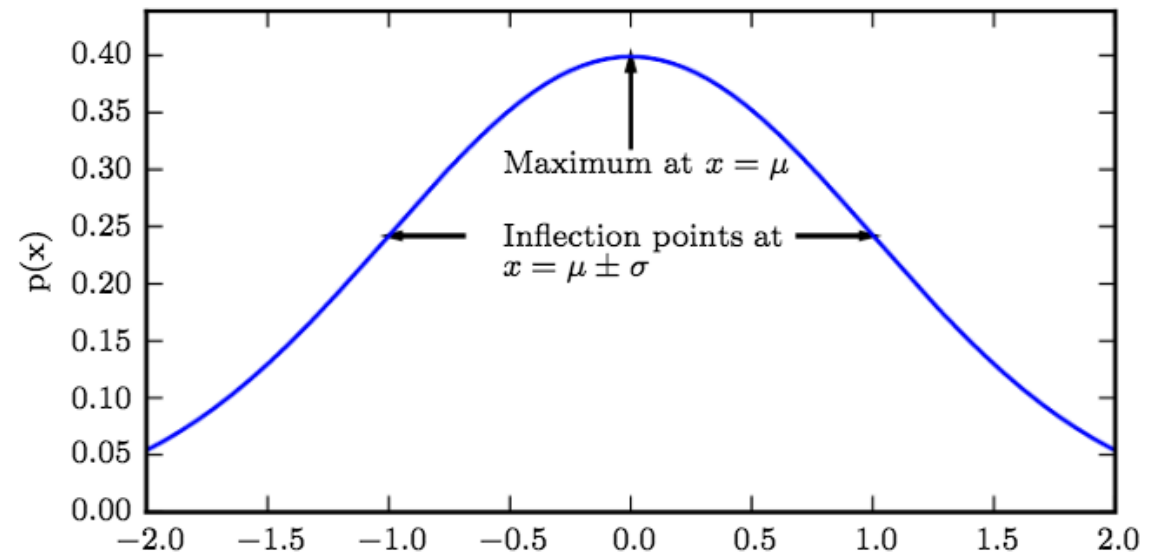
- $f(k; p) = \begin{cases} p & \text{if } k = 1, \\ q = 1 - p & \text{if } k = 0 \end{cases}$ 
  - Or:
- $f(k; p) = p^k (1 - p)^{1-k}$  for  $k \in \{0, 1\}$ 
  - $E[X] = p$
  - $Var_x(X) = p(1 - p)$
- What is it for? Find examples.

# Multinoulli Distribution & Categorical dis.

- Special case of the “multinomial distribution”
- The multinoulli, or categorical, distribution is a distribution over a single discrete variable with  $k$  different states, where  $k$  is finite.
- Vector  $\mathbf{p} \in [0,1]^{k-1}$ ,  $p_i$  gives the probability
- $k > 0$ , number of categories
- Often over “objects” (no need for the var, or  $E[x]$ )
- <https://www.statisticshowto.datasciencecentral.com/multinomial-distribution/>

# Gaussian Distribution

- Also known as: normal distribution.
- $N(x; \mu, \sigma^2) = \sqrt{\frac{1}{2\pi\sigma^2}} \exp(-\frac{1}{2\sigma^2} (x - \mu)^2)$
- $\mu = E[X], \sigma^2 = Var(X), \sigma = SD$
- central limit theorem
- multivariate normal distribution
- What for?



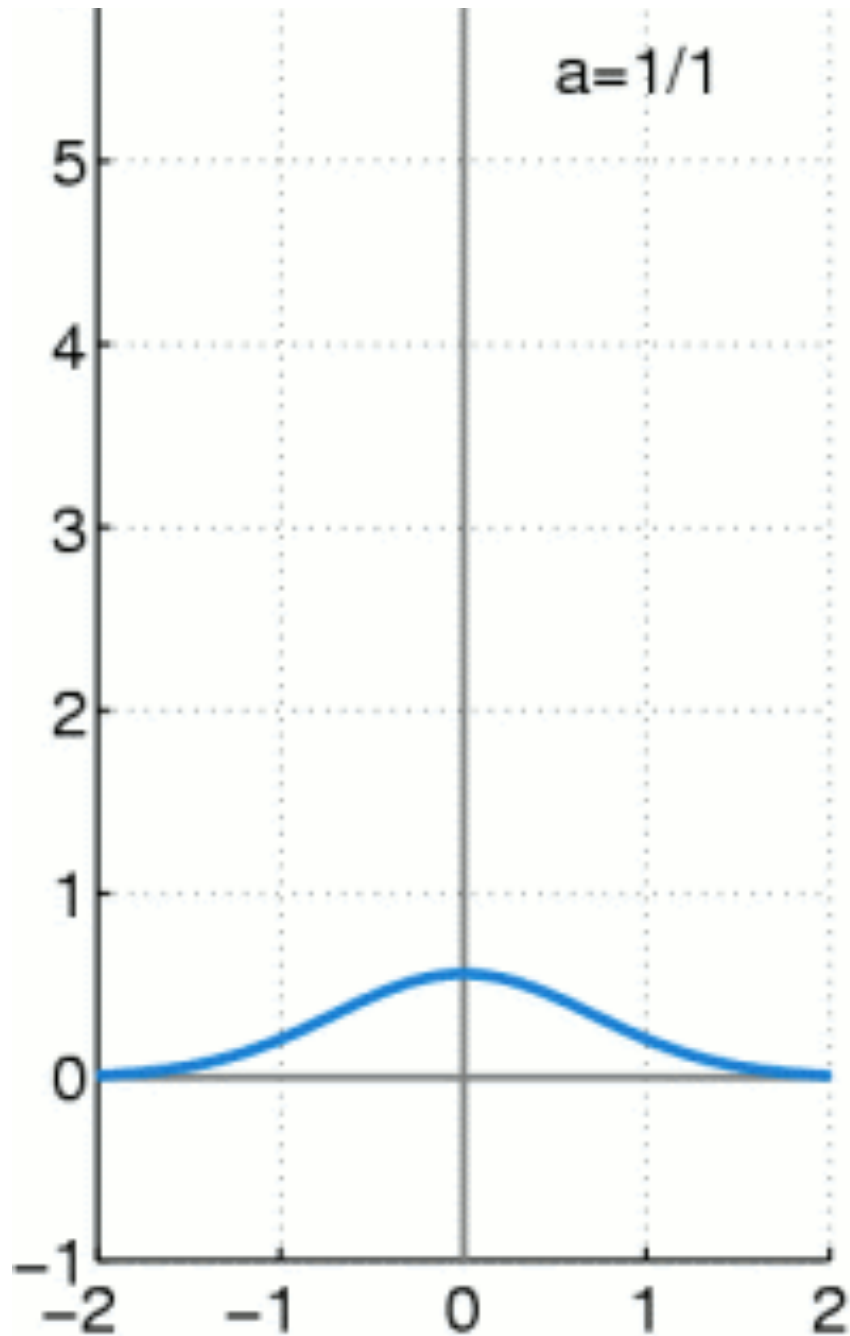
# Question by Amazon

- Given an unfair coin with the probability of heads and tails not equal to 50/50, what algorithm could generate a list of random ones and zeros?



# Laplace Distribution aka. Exponential Dis.

- In the context of deep learning, we often want to have a probability distribution with a sharp point at  $x = 0$
- $\mu = \text{Mean}; b > 0$  (*scale*)
- What for???
- A random variable has a Laplace  $(\mu, b)$  distribution, if its PMF is:
- $f(x; \mu, b) = \frac{1}{2b} \exp(-\frac{|x-\mu|}{b})$



# Dirac Dis.

- Dirac Delta Function as a “Limit”
- to specify that all the mass in a probability distribution clusters around a single point.
- $p(x) = \delta(x - \mu)$
- Aka: generalized function
- We can think of the Dirac delta function as being the limit point of a series of functions that put less and less density on all points other than zero.

By defining  $p(x)$  to be  $\delta$  shifted by  $-\mu$  we obtain an infinitely narrow and infinitely high peak of probability density where  $x = \mu$ .

A common use of the Dirac delta distribution is as a component of an **empirical distribution**,

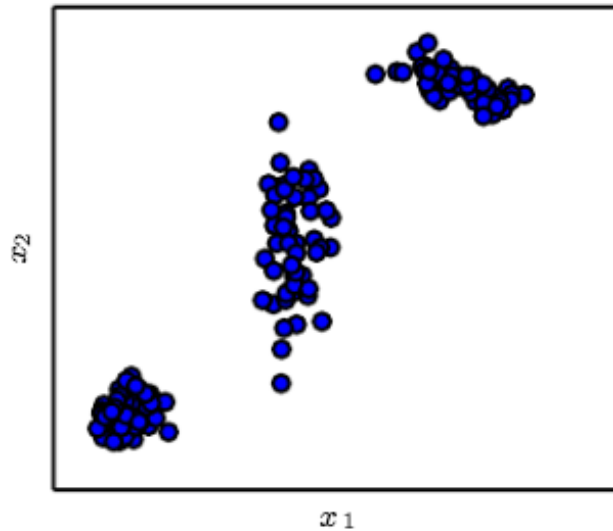
$$\hat{p}(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m \delta(\mathbf{x} - \mathbf{x}^{(i)}) \quad (3.28)$$

which puts probability mass  $\frac{1}{m}$  on each of the  $m$  points  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}$ , forming a given data set or collection of samples. The Dirac delta distribution is only necessary to define the empirical distribution over continuous variables. For discrete variables, the situation is simpler: an empirical distribution can be conceptualized as a multinoulli distribution, with a probability associated with each possible input value that is simply equal to the **empirical frequency** of that value in the training set.

## Empirical Dis. & Dirac

- What for?

# Mixtures of Dis.



- $P(x) = \sum_i P(c = i)P(x|c = i)$
- $P(c)$  is the multinoulli distribution over component identities
- E.g.: Empirical Dis. & Gaussian Mixture Model

A very powerful and common type of mixture model is the **Gaussian mixture model**, in which the components  $p(\mathbf{x} | c = i)$  are Gaussians. Each component has a separately parametrized mean  $\mu^{(i)}$  and covariance  $\Sigma^{(i)}$ . Some mixtures can have more constraints. For example, the covariances could be shared across components via the constraint  $\Sigma^{(i)} = \Sigma, \forall i$ . As with a single Gaussian distribution, the mixture of Gaussians might constrain the covariance matrix for each component to be diagonal or isotropic.

# Question by Postmates

- There are four people on the ground floor of a building that has five levels not including the ground floor. They all get into the same elevator.
- If each person is equally likely to get on any floor and they leave independently of each other, what is the probability that no two passengers will get off at the same floor?

# Bayes' Rule

- $P(A|B) = \frac{P(B|A)P(A)}{P(B)}$
- Often we know  $P(B|A)$ , but want to find  $P(A|B)$

# Question: by Facebook

- You are about to get on a plane to Seattle. You want to know if you should bring an umbrella. You call 3 random friends of yours who live there and ask each independently if it's raining. Each of your friends has a  $\frac{2}{3}$  chance of telling you the truth and a  $\frac{1}{3}$  chance of messing with you by lying. All 3 friends tell you that "Yes" it is raining.
- What is the probability that it's actually raining in Seattle?

# Question: by Google

- A jar holds 1000 coins. Out of all of the coins, 999 are fair and one is double-sided with two heads. Picking a coin at random, you toss the coin ten times.
- Given that you see 10 heads, what is the probability that the coin is double headed and the probability that the next toss of the coin is also a head?
- Give your answer to 3 significant figures.



# Question by Microsoft

- Amy and Brad take turns in rolling a fair six-sided die. Whoever rolls a "6" first wins the game. Amy starts by rolling first.
- What's the probability that Amy wins?

The basic intuition behind information theory is that learning that an unlikely event has occurred is more informative than learning that a likely event has occurred. A message saying “the sun rose this morning” is so uninformative as to be unnecessary to send, but a message saying “there was a solar eclipse this morning” is very informative.

We would like to quantify information in a way that formalizes this intuition.

- Likely events should have low information content, and in the extreme case, events that are guaranteed to happen should have no information content whatsoever.
- Less likely events should have higher information content.
- Independent events should have additive information. For example, finding out that a tossed coin has come up as heads twice should convey twice as much information as finding out that a tossed coin has come up as heads once.

## Information Theory

- <https://www.deeplearningbook.org/contents/prob.html#pf11>

# Structured Probabilistic Models

- $P(a, b, c) = p(a)p(b|a, c)p(c)$
- Suppose that a influences the value of b, and c influences the value of b, but that a and c are independent given b.
- structured probabilistic model, or graphical model
  - Directed & Undirected

# Directed:

**Directed** models use graphs with directed edges, and they represent factorizations into conditional probability distributions, as in the example above. Specifically, a directed model contains one factor for every random variable  $\mathbf{x}_i$  in the distribution, and that factor consists of the conditional distribution over  $\mathbf{x}_i$  given the parents of  $\mathbf{x}_i$ , denoted  $Pa_{\mathcal{G}}(\mathbf{x}_i)$ :

$$p(\mathbf{x}) = \prod_i p(\mathbf{x}_i \mid Pa_{\mathcal{G}}(\mathbf{x}_i)) . \quad (3.53)$$

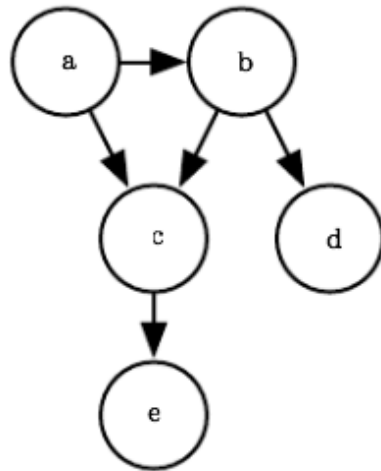


Figure 3.7: A directed graphical model over random variables  $a$ ,  $b$ ,  $c$ ,  $d$  and  $e$ . This graph corresponds to probability distributions that can be factored as

$$(3.54)$$

HOW  
WOULD THE  
FUNCTION  
LOOK LIKE?

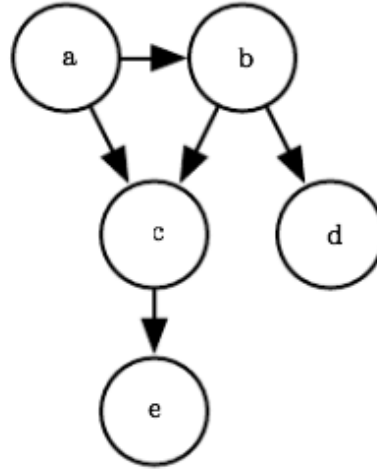


Figure 3.7: A directed graphical model over random variables  $a$ ,  $b$ ,  $c$ ,  $d$  and  $e$ . This graph corresponds to probability distributions that can be factored as

(3.54)

$$p(a, b, c, d, e) = p(a)p(b|a)p(c|a, b) p(d|b)p(e|c)$$

# Undirected:

- Clique:  $\mathcal{C}^i$

$$p(\mathbf{x}) = \frac{1}{Z} \prod_i \phi^{(i)}(\mathcal{C}^i).$$

- $\phi$  are factors
- (just functions, no probability dis.)

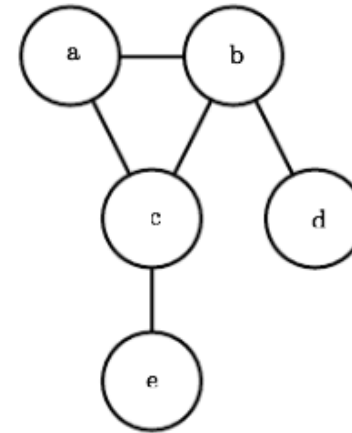


Figure 3.8: An undirected graphical model over random variables  $a$ ,  $b$ ,  $c$ ,  $d$  and  $e$ . This graph corresponds to probability distributions that can be factored as

$$p(a, b, c, d, e) = \frac{1}{Z} \phi^{(1)}(a, b, c) \phi^{(2)}(b, d) \phi^{(3)}(c, e). \quad (3.56)$$

Keep in mind that these graphical representations of factorizations are a language for describing probability distributions. They are not mutually exclusive families of probability distributions. Being directed or undirected is not a property of a probability distribution; it is a property of a particular **description** of a probability distribution, but any probability distribution may be described in both ways.



# Monte Carlo

- **Monte Carlo methods**, or **Monte Carlo experiments**, are a broad class of computational algorithms that rely on repeated random sampling to obtain numerical results. The underlying concept is to use randomness to solve problems that might be deterministic in principle. They are often used in physical and mathematical problems and are most useful when it is difficult or impossible to use other approaches. Monte Carlo methods are mainly used in three problem classes:<sup>[1]</sup> optimization, numerical integration, and generating draws from a probability distribution.
- Monte Carlo methods vary, but tend to follow a particular pattern:
  - Define a domain of possible inputs
  - Generate inputs randomly from a probability distribution over the domain
  - Perform a deterministic computation on the inputs
  - Aggregate the results

# Monte Carlo

- $\langle A \rangle = \sum_{x \in \Omega} P(x)A(x)$ , (Discrete)
- $\int_{x \in \Omega} P(x)A(x)d^n x$ ,
- Find an example.

### Definition of a Markov chain

A **stochastic process** in discrete-time is a family,  $(X(n))_{n \in \mathbb{N}_0}$ , of random variables indexed by the numbers  $\mathbb{N}_0 = \{0, 1, 2, \dots\}$ . The possible values,  $S$ , of  $X(n)$  are referred to as the **state space** of the process. In this course we consider only stochastic processes with values in a finite or countable state space. The mathematician may then think of a **random variable**,  $X$ , on  $S$  as a measurable map <sup>1</sup>

$$X: (\Omega, \mathcal{F}) \rightarrow (S, \mathcal{P}(S))$$

where  $\mathcal{P}(S)$  is the family of all subsets of  $S$ .

The distribution of a discrete-time stochastic process <sup>2</sup> with at most countable state space,  $S$ , is characterised by the point probabilities

$$P(X(n) = i_n, X(n-1) = i_{n-1}, \dots, X(0) = i_0)$$

for  $i_n, i_{n-1}, \dots, i_0 \in S$  and  $n \in \mathbb{N}_0$ . From the definition of elementary conditional probabilities it follows that

$$\begin{aligned} & P(X(n) = i_n, \dots, X(0) = i_0) \\ &= P(X(n) = i_n | X(n-1) = i_{n-1}, \dots, X(0) = i_0) \\ &\times P(X(n-1) = i_{n-1} | X(n-2) = i_{n-2}, \dots, X(0) = i_0) \\ &\times \dots \\ &\times P(X(1) = i_1 | X(0) = i_0) \times P(X(0) = i_0). \end{aligned}$$

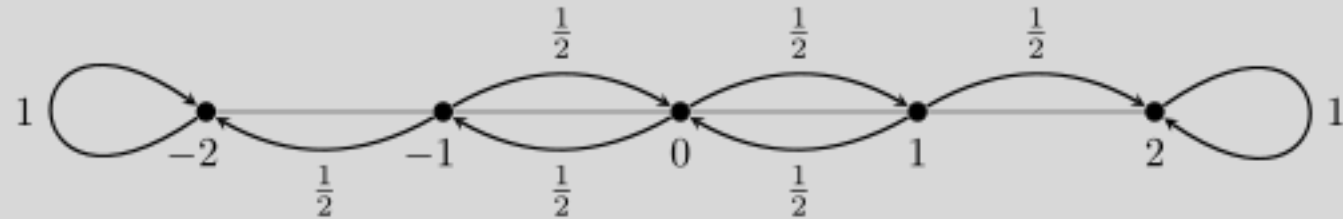
This is a general identity that holds for *any* discrete-time stochastic process on a countable state space. In these lecture notes we are only going to discuss the class of Markov chains to be defined below.

# MARKOV CHAINS

A sequence of random variables:

# Markov Chains/Models

- A **Markov chain** is a stochastic model describing a sequence of possible events in which the probability of each event depends only on the state attained in the previous event.
- Find an example.

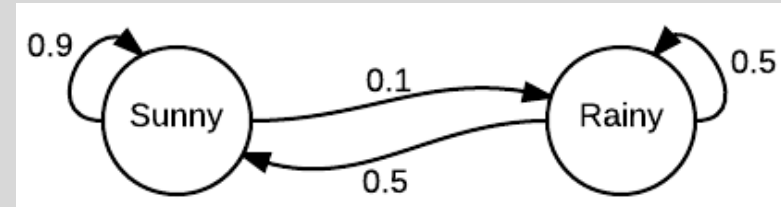


# Markov Chain

- The Markov Chain describes a sequence of possible events in which the probability of each event depends only on the state attained in the previous event. It is used for simulating sampling complex probability distributions.
- The Markov Chain (State Space of a chain) can be shown in a matrix.

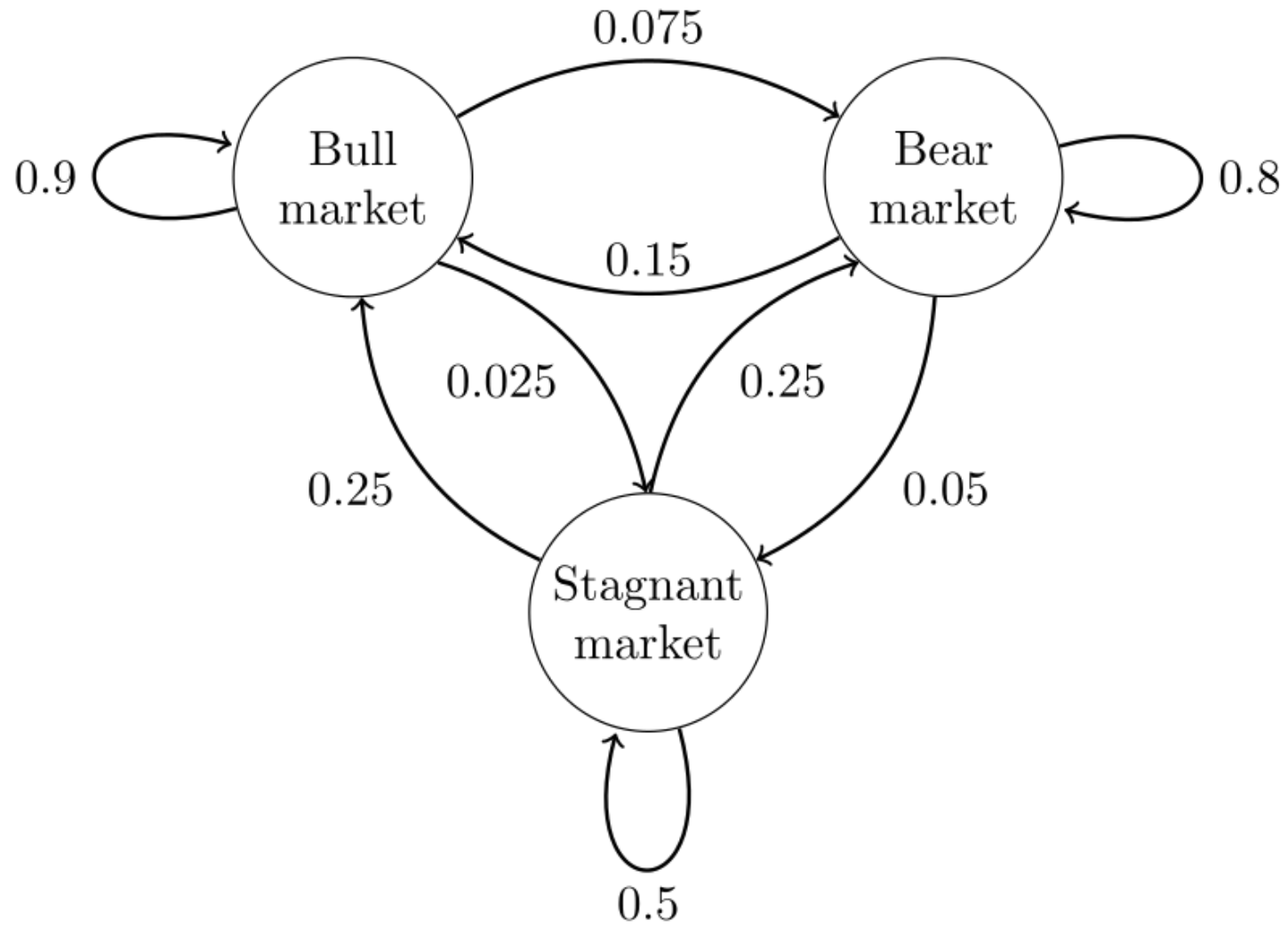
# Markov Chain: Example

- The probabilities can be displayed as a transition matrix:
  - $P = \begin{bmatrix} 0,9 & 0,1 \\ 0,5 & 0,5 \end{bmatrix}$  (The columns can be labeled „sunny“ & „rainy“)
- Predicting the weather:
- Day 0 (today) is sunny. (We know that, by looking outside)
- $x^0 = [1 \ 0]$ ,
- Tomorrow will be:
- $x^1 = x^0 P = [1 \ 0] \begin{bmatrix} 0,9 & 0,1 \\ 0,5 & 0,5 \end{bmatrix} = [0,9 \ 0,1]$ ,
- The day after tomorrow:
- $x^2 = x^1 P = x^0 P * P = x^0 p^2 = [1 \ 0] \begin{bmatrix} 0,9 & 0,1 \\ 0,5 & 0,5 \end{bmatrix}^2 = [0,86 \ 0,14]$ ,



# Remeber:

$$\circ x^n = x^{n-1}P \quad \text{or} \quad x^n = x^0 P^n$$



## Now You: Markov Chains

- What will the likelihood be in 3 “Times”, of the events, if the market is now bearish?



# Hashing Algorithms

- What does probability has to do with hashing???

# Links:

- <https://www.statisticshowto.datasciencecentral.com/probability-and-statistics/probability-main-index/>
  - Like a dictionary for DS
- <https://ocw.mit.edu/courses/electrical-engineering-and-computer-science/6-041-probabilistic-systems-analysis-and-applied-probability-spring-2006/lecture-notes/>
  - MIT Lecture notes
- <https://towardsdatascience.com/introduction-to-markov-chains-50da3645a50d>
  - DS
- <http://web.math.ku.dk/noter/filer/stoknoter.pdf>
  - Markov Chains

# Bibliography

- <https://en.wikipedia.org/wiki/Probability> (2019)
- [https://en.wikipedia.org/wiki/Random\\_variable](https://en.wikipedia.org/wiki/Random_variable)
- [https://en.wikipedia.org/wiki/Probability\\_distribution](https://en.wikipedia.org/wiki/Probability_distribution)
- [https://en.wikipedia.org/wiki/Probability\\_distribution](https://en.wikipedia.org/wiki/Probability_distribution)
- [https://en.wikipedia.org/wiki/Probability\\_distribution](https://en.wikipedia.org/wiki/Probability_distribution)
- [https://en.wikipedia.org/wiki/Conditional\\_independence](https://en.wikipedia.org/wiki/Conditional_independence)
- [https://en.wikipedia.org/wiki/Bernoulli\\_process](https://en.wikipedia.org/wiki/Bernoulli_process)
- [https://en.wikipedia.org/wiki/Monte\\_Carlo\\_method](https://en.wikipedia.org/wiki/Monte_Carlo_method)
- [https://en.wikipedia.org/wiki/Markov\\_chain](https://en.wikipedia.org/wiki/Markov_chain)
- <https://towardsdatascience.com/introduction-to-markov-chains-50da3645a50d>
- <http://web.math.ku.dk/noter/filer/stoknoter.pdf>