

Probability &

Statistics For

Machine learning.

Probability

The measure of how likely an event will occur

$$\rightarrow 10 \text{ kids } 3 \text{ play soccer } 7 \text{ don't } P(S) = 3/10 = 0.3$$

$$\rightarrow 2 \text{ coins flipped } \{HH, HT, TH, TT\} \quad P(HH) = 1/4$$

$$\rightarrow 3 \text{ coins flipped } \{HHH, HHT, HTH, HTT, THH, TTH, THT, TTT\} = 8 = 2^3 \quad P(HHH) = 1/8$$

$$\rightarrow \text{Dice rolled } P(6) = 1/6$$

$$\rightarrow 2 \text{ dice rolled } 6^2 = 36 \text{ outcomes } P(6, 6) = 1/36$$

Complement of Probability

$$\rightarrow 10 \text{ kids soccer, 3 play } P(\text{soccer}) = 1 - P(\text{not soccer}) = 1 - 0.3 = 0.7$$

$$\rightarrow \text{coin tossed 3 times } P(HHH) = 1 - P(HTH) = 1 - 1/8 = 7/8$$

Sum of Probabilities. Disjoint Events

$$\rightarrow 1 \text{ school kids play 1 sport } P(\text{soccer}) = 0.3 \quad P(\text{Basketball}) = 0.4$$

$$P(\text{soccer} \cup \text{Basketball}) = 0.3 + 0.4 = 0.7$$

$$P(A \cup B) = P(A) + P(B)$$

$$\rightarrow \text{when throwing a die } P(\text{even}) \text{ or } P(\text{odd})$$

$$\rightarrow 4/6 = 2/3$$

$$\rightarrow \text{roll two fair dice } P(\text{difference of 2}) \text{ or } P(\text{difference of 1})$$

$$\rightarrow \{(6,1), (5,3), (4,2), (3,1), (6,5), (5,4), (4,3), (3,2), (2,1)\} \times 2 \text{ to roll dice} \rightarrow 18/36 = 1/2$$

Sum of Probabilities (Joint Events)

$$\rightarrow P(\text{rainy}) \quad P(\text{windy}) \quad P(\text{cloudy}) \quad P(\text{cloudy or rainy}) = 150/1000 \text{ wtf!} \\ 60\% \quad 70\% \quad \text{rainy and windy can happen together} \\ \text{they are not independent}$$

$$\rightarrow \text{School kids can play as many sports as they want} \quad P(S) = 0.6 \quad P(B) = 0.5 \\ P(S \cup B) = P(S) + P(B) - P(S \cap B) = 0.6 + 0.5 - 0.6 \cdot 0.5 = 1/1 - \text{cannot infer}$$

Now 6 kids Play Soccer 5 basketball 3 both.

$$P(S) = 4/11 \quad P(B) = 5/11 \quad P(S \cap B) = 3/11$$

8

Disjoint vs Joint events

Mutually exclusive

non-mutually-exclusive

$$P(S \cup B) = P(S) + P(B)$$

$$P(S \cap B) = P(S) \cdot P(B)$$

$$\text{Q1: Sum of 7 or a difference of 1 rolling 2 fair die}$$

$$\text{sum } (6,1), (5,2), (4,3), (3,4), (2,5), (1,6) = 6$$

$$\text{diff } (6,5), (5,4), (4,3), (3,2), (2,1), (1,2), (2,3), (3,4), (4,5), (5,6) = 10 \\ (3,3), (4,4) \text{ repeat } \frac{10+6}{36} = 2/36 = [14/36]$$

Independent Event

When one event does not affect another event

\rightarrow School \rightarrow Go soccer \rightarrow not soccer

$$\text{Product Rule} \rightarrow P(A \cap B) = P(A) \cdot P(B)$$

$$\text{Q2: coin tossed 5 times heads all 5 times}$$

$$P(HHHHHHHHHH) = P(H) \cdot P(H) \cdot P(H) \cdot P(H) \cdot P(H) \\ = \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{32}$$

$$\text{Q3: if you roll a fair die 5 times. } P(5 \text{ sides})?$$

$$\rightarrow P(5) = \left(\frac{1}{6}\right)^5$$

$$\text{Q4: 30 friends at a party. What is more likely? +} \text{more people with same birthday or} \\ \text{- Assume 365 days. No Feb 29.} \quad \text{+ that no two of them have the same birthday}$$

With

$$\begin{array}{l} \text{Person 1} \quad \text{Person 2} \quad 3 \quad 4 \quad 5 \quad 6 \quad 7 \quad 8 \quad 9 \\ \frac{365}{365} \times \frac{364}{365} \times \frac{363}{365} \times \frac{362}{365} \times \frac{361}{365} \times \frac{360}{365} \times \frac{359}{365} \times \frac{358}{365} \times \frac{357}{365} = 0.905 \\ \text{which is pretty high.} \end{array}$$

\rightarrow This is what the graph looks like

(9)



30 friends 0.5 at 23

Conditional Probability

$$P(\text{today sunny}) = 0.6$$

$$\text{wt yesterday rainy to today } P(\text{rainy}) \text{ changes}$$

$$\text{Q5: Probability of landing heads twice if the first coin flipped is heads?} \\ \frac{1}{2} \quad P(HH) \mid \text{1st is H}$$

Product Rule for Independent events

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

Product rule for Conditional Probability

$$P(A \cap B) = P(A) \cdot P(B/A) \quad \text{single capital}$$

Lemma explain:

Rule 2 dice, probability that the sum is 10 & the first is A.
so let's take A first + first is 6 now we have to find the possibilities to
sum is 10 for 6H so $P(A) = \frac{1}{6}$, $P(B/A) = \frac{1}{6}/\frac{1}{6} = \frac{1}{6} \cdot \frac{1}{6} = \frac{1}{36}$

If you go opposite way it works also

$$P(B) \rightarrow \text{sum 10 + } \{6, 4, 5, 5, 4, 6\} \rightarrow 6 \text{ total possibilities. } \rightarrow 6/36$$

$$\begin{aligned} P(A|B) &\rightarrow 1/6 \\ &\rightarrow \frac{1}{6} \cdot \frac{1}{6} = 1/36 \end{aligned}$$

$$\text{Q6: } P(S) = 0.4 \quad P(CS) = 0.6 \quad P(\text{Running shoes}) = 0.8$$

$$P(S \cap C) = P(S) \cdot P(C/S) = 0.4 \cdot 0.8 \\ = 0.32 = 32 \text{ kids}$$

$$\text{Q7: } P(C|NS) = 0.5 \rightarrow P(NS \cap C) = ?$$

$$\begin{aligned} P(NS) \cdot P(C/NS) \\ \rightarrow 0.6 \cdot 0.5 = 0.3 = 30 \end{aligned}$$

no play football
% of students
that play football
wear running shoes

ML Use

You wanna calculate the probability of something based on other factors.

Spam detection → guess prob of spam based on features.

Image recognitions calculate prob that there is a cat in the image based on pixels of an image.

Image generation → you want to maximize the prob that a bunch of pixels form a face.

Text generation → bunch of characters form a word & words from a sentence.

Image recognition

- What is the probability that there is a cat in the image?
- $P(\text{cat image}) = P(\text{cat} | \text{pixels}, \text{pixels}, \dots, \text{pixels})$

Sentiment analysis

- Is this a happy sentence?
- Calculate $P(\text{happy} | \text{words in the sentence})$

$$\text{Picture} \rightarrow \text{model} \rightarrow P(\text{cat} | \text{pixels}) = 0.9$$

$$\text{Picture} \rightarrow \text{model} \rightarrow P(\text{cat} | \text{pixels}) = 0.1$$



Image generation

Generate a group of pixels such that the resulting image looks like human face.

Goal: generate images such that $P(\text{face} | \text{pixels})$ is high.

Week 1 Quiz

1) HT, TH, HH $\rightarrow P(H) = 1/2$ 2) $P(\text{no tail}) = P(HH, HT, TH, TT) = 3/4 = 75\%$

3) $P(\text{no tail}) = (1 + \text{out times})$ 4) 100 patients 50 headache so fever

$$10 - 1 - 1/10 = 6/10 = 60\% \quad P(\text{F} | \text{U}) = P(\text{F} \cap \text{U}) / P(\text{U})$$

5) A: 1000, B: 2000 log: 2000

$$P(B|A) = 1500$$

$$P(C|B|A) = \frac{P(A|B|A)}{P(B|A)} \Rightarrow \frac{1500}{2000} = 1/2 = 50\%$$

$$P(X|Y) = \frac{P(X \cap Y)}{P(Y)}$$

Random Variables

They can take many values: eg Temp, no. of heads on 10 coins

X = No. of heads

Flips coin

1 heads 0 heads
 $x=1$ $x=0$
 $P(x=1)=0.5$ $P(x=0)=0.5$

D is discrete
0 1 2 3 4

Countable
but can be infinite

Continuous
0.1, 0.00005

Infinite
takes value on an entire range

Probability Distributions

e.g. X = No. of heads in 3 coins.

$\Omega^3 = 8$ - total outcomes

$x=0$	- 1	1	3
$x=1$	- 3	2	
$x=2$	- 3	1	
$x=3$	- 1		

Flip 4 coins $2^4 = 16$:

$x=0 - 1$

$x=1 - 4$

$x=2 - 6$

$x=3 - 4$

$x=4 - 1$



$X_3 : 5 \text{ coin flips}$

No. tails = $2^5 = 32$ have to calculate.

So how $P(X_3 = ?)$? we can represent it as $P(X_3 = x)$, $x = 0, 1, 2, 3, 4, 5$

This is PMF - Probability Mass Function.

$$p_x(x) = P(X_3 = x), x = 1, 2, 3, 4, 5.$$

PMF must satisfy $\sum_x p_x(x) \geq 0$.

X_1, X_2, X_3, X_4, X_5 are very similar

They all represent No. of heads in n experiments

The way the probability distributes along the possible outcomes seems to have a similar pattern.

Could there be a single model to represent this pattern?

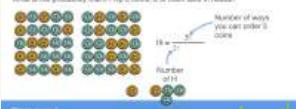
Binomial distribution

Let's say you toss a coin 5 times, 2 land on heads?

is there some way to find this?

Binomial Distribution: Example

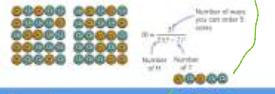
What is the probability that if I flip 5 coins, 2 of them land in heads?



L to count for no. of ways between some positions

Binomial Distribution: Example

What is the probability that if I flip 5 coins, 3 of them land in heads?



L to count for no. of ways between some positions

Binomial Distribution: Example

What is the probability that if I flip 5 coins, 0 of them land in heads?



$$\binom{n}{k} = \text{count all the combinations for landing } k \text{ heads in } n \text{ coin tosses}$$

$$\text{Property: } \binom{n}{k} = \binom{n}{n-k}$$

Probability of getting k heads in n coin tosses is the same as n-k tails.

General PMF for X: no. of heads in 5 coin tosses?

Coin $P(H) = p$ (usually 0.5)

Event $X = x$: x heads in 5 tosses

$$P_x^x(1-p)^{5-x}$$

probability of getting 5-x tails
of getting x heads

This is w/o care just one order

Now we need no. of ways \rightarrow Binomial coefficient

$$\binom{5}{x} p^x (1-p)^{5-x}, x = 0, 1, 2, 3, 4, 5$$

X follows a binomial distribution

for n coins

$$P_x(x) = \binom{n}{x} p^x (1-p)^{n-x}, x = 0, 1, 2, \dots, n$$

if Throw dice 5 times for 3 ones

$$P(1) = 1/6$$

$$P(3) = \binom{5}{3} \cdot \left(\frac{1}{6}\right)^3 \left(\frac{5}{6}\right)^2$$

$$= \frac{5!}{3! 2!} \cdot \left(\frac{1}{6}\right)^3 \left(\frac{5}{6}\right)^2$$

$$= \frac{5 \times 4 \times 3}{2 \times 1} \cdot \left(\frac{1}{6}\right)^3 \left(\frac{5}{6}\right)^2$$

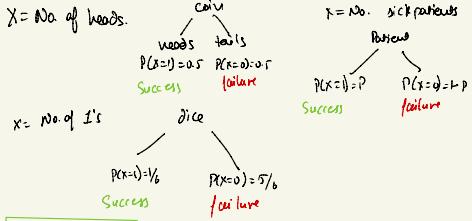
Binomial Coefficient

now this works more generally when you need to pick k out of n numbers
 Pick n items
 Pick m items
 n options $\leq n-1 \text{ options} \leq n-2 \dots \leq n-(k-1) \text{ options}$
 total length $n, n-1, n-2, \dots, k, k-1, \dots, 1$
 $\binom{n}{k}$

As these can be taken as independent events.



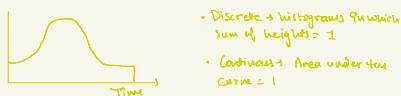
Bernoulli Distribution



Bernoulli	Binomial
Single trial	multiple (n) independent trials
2 outcomes (success/failure)	n outcomes (success/failure)
1 parameter (P success)	two parameters (n & P success)
0 or 1	0 to n
taking a coin once	taking a coin n times & getting no. heads

Probability Distributions (continuous)

No. of heads → can be listed ($0, 1, 2, 3, \dots$)
 amount of time on plane or on the bus ($0, 1, 2, 1, 2.13$)
 this cannot be listed.



Probability Density Function



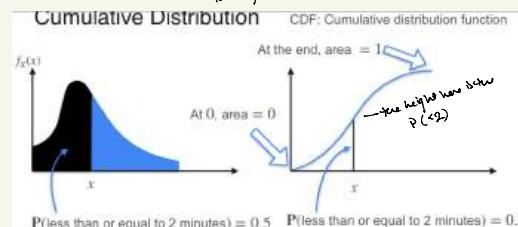
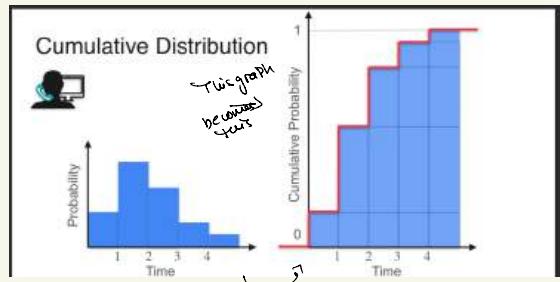
PDF: $f_X(x) \rightarrow$ tells you the rate you accumulate probability around each point.

- $P(a < X < b) = \text{area under } f_X(x)$
- $f_X(x)$ needs to satisfy:
 - It is defined for all numbers
 - $f_X(x) \geq 0$
 - Area under $f_X(x) = 1$

Cumulative Distribution Function

We don't always want to calculate areas.

The CDF tells you cumulative P from 0 to a specific n .



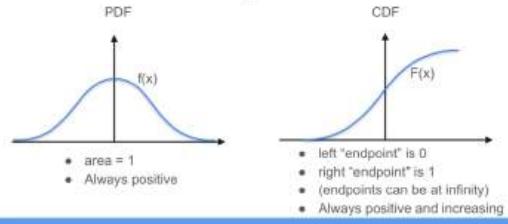
Cumulative Distribution Function: Formal Definition

The CDF shows how much probability the variable has accumulated until a certain value.

That means that:



PDF and CDF Summary



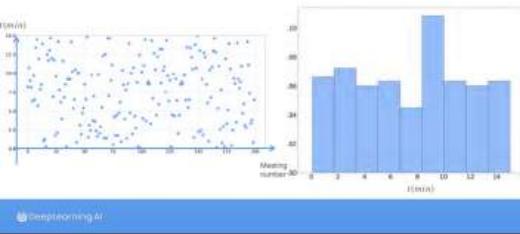
Uniform Distribution



You're calling a tech support line. They can answer any time between zero and 15 minutes and if they don't answer in this time, the line is disconnected.

Uniform Distribution: Motivation

wait time over 200 calls.



Uniform Distribution: Motivation

T: time (in minutes) you have to wait

Any value between 0 and 15 minutes must have the same frequency of occurrence.

The pdf must be constant for all values in the interval (0,15)

$$\text{Which constant?} \rightarrow 15 \times h = 1 \rightarrow h = \frac{1}{15} = 0.06$$

A Continuous random variable can be modeled with a uniform distribution if all possible values lie in an interval & have same frequency of occurrence

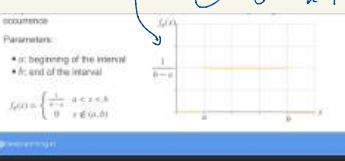
A Continuous random variable can be modeled with a uniform distribution if all possible values lie in an interval & have same frequency of occurrence

Parameters:

a: beginning of interval

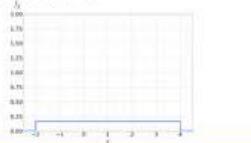
b: end of interval

$$f_x(x) = \begin{cases} \frac{1}{b-a} & a < x < b \\ 0 & x \notin (a,b) \end{cases}$$

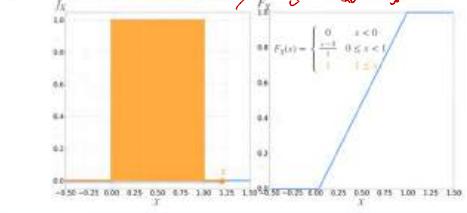


This is what the PDF would look like.

Uniform Distribution: PDF



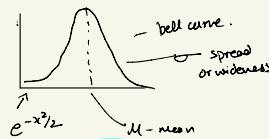
Uniform Distribution: CDF



For CDF with zero mode our cumulative distribution function is 0 for x < 0 and 1 for x > 15.

Normal / Gaussian Distribution

When n is very large the binomial distribution can be approximated by the Gaussian distribution very well.



$$\frac{1}{\sigma\sqrt{2\pi}} \cdot e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

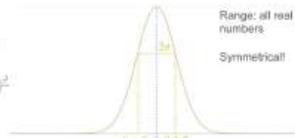
This is to regulate the regular area. — The idea is $\sigma\sqrt{2\pi}$ is the area under the original curve but the probability the area should add up to 1 so we divide it by $(\sigma\sqrt{2\pi})$. σ is the standard deviation or spread (from the mean $(x-\mu)$) → This tells you to bring x to 0 as the bell curve is usually at 0. σ we divide by to regulate the std deviation or spread.

Normal Distribution

Parameters:

- * μ : center of the bell
- * σ : spread of the bell

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



Range: all real numbers
Symmetrical

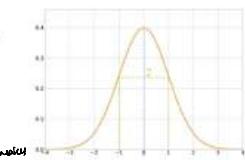
Standard Normal Distribution

Parameters:

- * $\mu = 0$
- * $\sigma = 1$

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

This's much easier this function much easier



The $\frac{(x-\mu)}{\sigma}$ → I explained earlier is standardization.

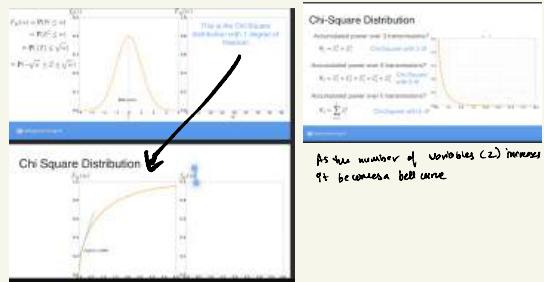
Chi Squared Distribution χ^2

Chi-Square Distribution: Motivation

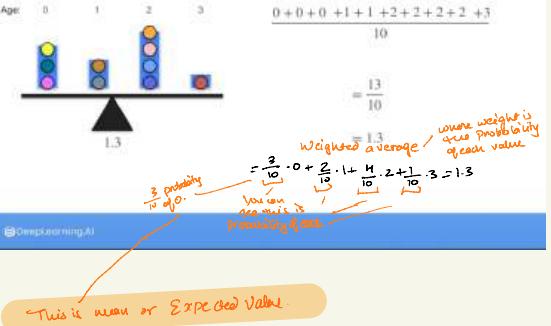


One common assumption is that the communication channel has noise with a standard normal distribution.

What is the power of the noise \rightarrow associated variance & dispersion of noise.
 $W = Z^2$



Mean: Example



Expected Value: Motivation Example 1

You play a game with a friend



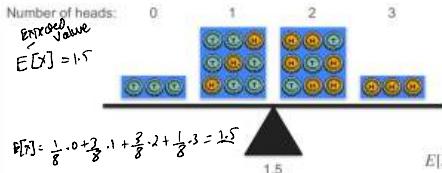
Game cost:

\$5

Long term: $0.5 \cdot \$10 + 0.5 \cdot \$0 = \$5 \Rightarrow$ You expect to win \$5 on average
 half of the time
 half of the time

Representation of mean or Expected Value

Flip a coin 8 times. Expected Value: Motivation Example 2



$$E[X] = 1.5$$

Expected Variable

X is a discrete random variable

PMF of X

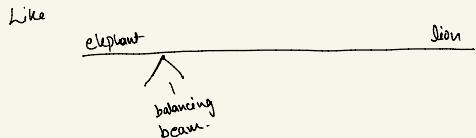
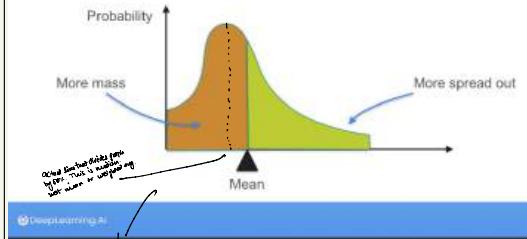
$$P_x(x_i) = P(X=x_i)$$

$$E[X] = \sum_i x_i p(x_i)$$

X is over continuous random variables

$$\int_{-\infty}^{+\infty} x_i f_X(x_i) dx_i \rightarrow \text{Summing the area of bars treated as rectangles.}$$

Expected Value: Common Misconception



Expected Value:

- $E[X]$
- Mean / Balancing point
- Defined for discrete & continuous random variables
- Weighted Average of the PMF / PDF