

Prior belief = belief before seeing the evidence
 Posterior belief = belief after seeing the evidence

Bayes Theorem With Discrete

Let M, N be discrete random variables

$$P(M=z | N=z) = \frac{P(N=z | M=z) P(M=z)}{P(N=z)}$$

More generally:

$$P(M=m_n | N=n) = \frac{P(N=n | M=m) P(M=m)}{P(N=n)}$$

Ex:

Let x be the change in gaze (in degrees) over 3 seconds after a sound is played.

$$P(\text{com here the sound}) = \frac{3}{5} \rightarrow Y \sim \text{Beta}(1, 1)$$

I observe $x=0$, what is the prob the baby can hear the sound

$$\begin{aligned} P(Y|x) &= \frac{P(x=0|y=1) P(y=1)}{P(x=0)} = \frac{\text{Since we don't have } P(x=0)}{\text{ }} \\ &= \frac{P(x=0|y=1) P(y=1)}{P(x=0|y=1) P(y=1) + P(x=0|y=0) P(y=0)} = \frac{3}{8} \end{aligned}$$

Inference With Continuous

Ex:

Q: At birth, girl elephant weights are distributed as a Gaussian with mean = 160kg, std = 7kg. At birth, boy elephant weights are distributed as a Gaussian with mean = 165kg, std = 3kg. All you know about a newborn elephant is that it is 163kg. What is the probability that it is a girl?



$$P(S|W) = \frac{P(W|S) P(S)}{P(W)} = \frac{P(W|S) P(S)}{P(W|S) P(S) + P(W|S^c) P(S^c)}$$

We have to say that the probability $P(W=163 | G=1)$ is 0 because we are speaking about continuous random variables.
But we can use ϵ so:

$$P(S|W) = \frac{\epsilon \cdot f(W|S) \cdot P(S)}{\epsilon \cdot f(W)}$$

Given that we use ϵ for this case to take the "instant" probability we can use the PDF of a normal distribution instead of the Φ for the CDF

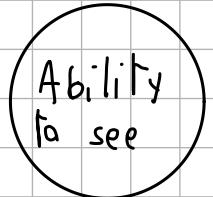
From here on when there is a question like how good or how bad something is i don't have to give a precise number but a probability distribution.

I can give an expectation but never a precise number

Dictionaries can also be used to represent a discretization of a continuous random var

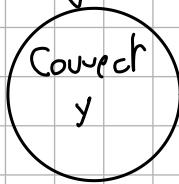
Ex:

As a Causal Model



A is a non-parametric distribution (dictionary)

Form size is fixed



$P(X=1 | A=x)$ is a Bernoulli random variable where the parameter p is a known function of x and form size s .

Q: What is $P(A=0 | x=0)$

Likelihood: Item response theory

If student i attempts problem j , the likelihood they answer it correctly is

$$P_{i,j} = \sigma(d_i - d_j)$$

ability of student i

difficulty of problem j

prob correct

squashing function

The ability - difficulty isn't necessary a not but also a number, then inside the squashing function it becomes a probability.

Modeling

Real world scenario with bunch of random variables.

Bayesian Networks

I have 4 health-related problems

Flu

Unconscious

In this case I have 4 joint probabilities.

Fever

Tired

I need to describe these using causality. For example flu causes fever. Even if it isn't deterministic it changes a lot. The probability

Constructing a Bayesian Network

Almost all the models that define causality are wrong but if I am able to build something with the causality I would be able to reconstruct the reality.

So the first step is calculating the joint distribution using causality.

The more causality relations we introduce the more the model grows exponentially.

The second step is giving the conditional probability between each event.

If I link everything in my network I'll finish to have a joint probability table.

The third step is to implicitly assume independences.

This means that I don't need to calculate relations that don't exist.

The model/graph must be acyclic.

$$P(\text{Joint}) = \prod_i P(X_i = x_i \mid \text{parents of } X_i)$$

Ex:

$$P(F_{lu} = 0, U = 1, F_{ev} = 0, T = 1)?$$

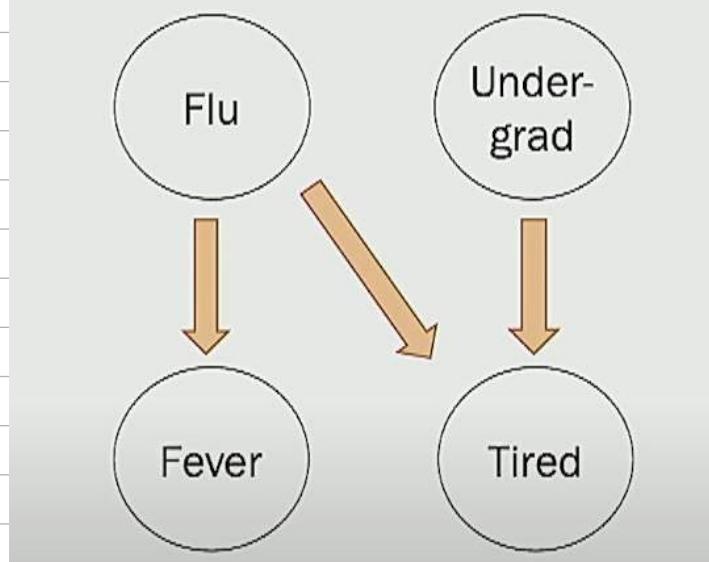
$$P(U = 1) = 0,8$$

$$P(F_{lu} = 0) = 0,9$$

$$P(F_{ev} = 0 \mid F_{lu} = 0) = 1 - 0,05$$

$$P(T = 1 \mid F_{lu} = 0, U = 1) = 0,8$$

$$P(F_{lu} = 1) = 0,1 \quad P(U = 1) = 0,8$$



$$\begin{aligned} P(F_{ev} = 1 \mid F_{lu} = 1) &= 0.9 \\ P(F_{ev} = 1 \mid F_{lu} = 0) &= 0.05 \end{aligned}$$

$$\begin{aligned} P(T = 1 \mid F_{lu} = 0, U = 0) &= 0.1 \\ P(T = 1 \mid F_{lu} = 0, U = 1) &= 0.8 \\ P(T = 1 \mid F_{lu} = 1, U = 0) &= 0.9 \\ P(T = 1 \mid F_{lu} = 1, U = 1) &= 1.0 \end{aligned}$$

Independence in RV

Recalling the independence of RV:

$$P(E \cap F) = P(E) P(F)$$

2 discrete random variables x and y are independent if:

For all x, y :

$$P(X=x, Y=y) = P(X=x) P(Y=y)$$

Intuitively: Knowing the value of X tells us nothing about the distribution of Y (and viceversa)

If 2 variables are not independent, they are called dependent

How do we discover independence from data?

The independence can give me lot of information about conditional probabilities

Covariance

Say X and Y are arbitrary random variables.

Covariance of X and Y :

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

It's the measure on how much 2 random variables vary together

- X and Y independent, $E[XY] = E[X]E[Y] \rightarrow \text{Cov}(X, Y) = 0$
- $\text{Cov}(X, Y) = 0$ does not imply X and Y independent.

An other way to calculate covariance:

$$\begin{aligned} \text{Cov}(X, Y) &= E[XY - E[X]Y - XE[Y] + E[X]E[Y]] \\ &= E[XY] - E[X]E[Y] - E[X]E[Y] + E[X]E[Y] \\ &= E[XY] - E[X]E[Y] \end{aligned}$$

I can actually discover the independence between random variables by covarying them.

The memory for a joint table is $O(n^2)$

Make a generative model

A good probabilistic model is generative. It explains the process through which the data is created.

Covariance of 0 doesn't mean independence

X and Y are random variables:

X is -1, 0 or 1 with equal probability:

$$y = \begin{cases} 0 & \text{if } X \neq 0 \\ 1 & \text{otherwise} \end{cases}$$

X and Y are random variables with PMF:

		-1	0	1	$p_Y(y)$
		0	1/3	0	1/3
		1	0	1/3	0
$p_X(x)$		1/3	1/3	1/3	1

$$Y = \begin{cases} 0 & \text{if } X \neq 0 \\ 1 & \text{otherwise} \end{cases}$$

$$E[X] = -1\left(\frac{1}{3}\right) + 0\left(\frac{1}{3}\right) + 1\left(\frac{1}{3}\right) = 0$$

$$E[Y] = 0\left(\frac{2}{3}\right) + 1\left(\frac{1}{3}\right) = \frac{1}{3}$$

$$\text{Cov}(x, y) = E[XY] - E[X]E[Y] = 0$$

The value of the covariance is 0 but the 2 random variables are not independent

Cauchy-Schwarz inequality.

$$-\text{std}(x)\text{std}(y) \leq \text{Cov}(x, y) \leq \text{std}(x)\text{std}(y)$$

Covariation is just normalized covariance

$$P(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \text{Var}(Y)}}$$

- We divide by this number because given the inequality before we always know that the covariance will be less.

Ex:

$$P(F_{lu} = 1 \mid U = 1, T = 1) = 0,1$$

1 - Compute all the joint probabilities

2 - Definition of conditional probabilities

$$\frac{\sum_y P(F_{lu} = 1, U = 1, F_{eu} = y, T = 1)}{\sum_y \sum_x P(F_{lu} = x, U = 1, F_{eu} = y, T = 1)} = 0,122$$

$$\sum_x \sum_y P(F_{lu} = x, U = 1, F_{eu} = y, T = 1)$$

I need to check all the possible combinations. I can do this by doing marginalisation

With thousands of R.V. is not possible to do like this.

It's not scalable doing like this.

Rejection Sampling Algorithm

Step 0: Define precisely a Bayesian Network

In this case instead of make all the calculations. i create lot of fake samples and then i count all the scenarios in which the characteristics i want are satisfied.

In order to do this i need to give all the needed probabilities from the graph.

Sampling could be a good idea. When i run a method in python all the values that they give me back are coming from sampling

Step 1: Rejecting

After the first pass I have a lot of samples.

Here I go through all the sample that I have and I throw away all the samples that are not consistent with the observation that I make.

From the smaller dataset that I have now I can't all the samples that respect a certain condition.

Considering the example used before:

$$P(Fev=1 \mid U=1, T=1) = \frac{10000}{30000}$$

At the beginning I have 100000 samples, from 100000 I have 30000 with $U=1$ and $T=1$ and from 30000 the samples with the fever are 10000

Why this approximation makes sense?

What is $P(flu=1 \mid U=1, T=1)$?

$$\text{Probability} \approx \frac{\# \text{Samples with } (flu=1, U=1, T=1)}{\# \text{Samples with } (U=1, T=1)}$$

Recalling our definition of probability as \rightarrow : $P(E) = \lim_{n \rightarrow \infty} \frac{N(E)}{n}$

$N(E) = \text{Fvials with } E \text{ occurred}$

When I'm sampling and I have random variables with different distributions I just need to sample from different distributions

Markov Chain Monte Carlo

It's a way to sample with conditioned variables fixed
We just fix one random variable during sampling

Flip a coin with unknown probability

Flip a coin $(n+m)$ times, comes up with n heads.

We don't know the probability X that coin comes up heads.

Frequentist

$$X = \lim_{n+m \rightarrow +\infty} \frac{n}{n+m} \approx \frac{n}{n+m}$$

X is often a single value

Bayesian (Prior is great!)

$$f_{x|n}(x|n) = \frac{P(N=n | X=x) f_x(x)}{P(N=n)}$$

X is a random variable. Leads to a belief distribution which captures confidence

Our belief before flipping coins is that: $X \sim \text{Unif}(0,1)$

Let N = number of heads.

Given $X = x$, coin flips independent: $(N|x) \sim \text{Bin}(n+m, x)$

$$f_{x|n}(x|n) = \frac{P(N=n | X=x) f_x(x)}{P(N=n)}$$

$$= \frac{\binom{n+m}{n} x^n (1-x)^m}{P(N=n)}$$

↙ P(N=n)

↙ To normalize the result

In this case we put $f(x) = 1$ because we choose a Beta with all the possibilities equally likely.

If you start with a $X \sim \text{Uni}(0,1)$ prior over probability and observe:

- n successes
- m failures

My new belief about the probability is:

$$f_X(x) = \frac{1}{c} \cdot x^n (1-x)^m$$

where:

$$c = \int_0^1 x^n (1-x)^m dx$$

Equivalently

If you start with a $X \sim \text{Uni}(0,1)$ prior over probability and observe:

$$\alpha = \text{num of successes} + 1$$

$$\beta = \text{num of failures} + 1$$

My new belief about probability is:

$$F_X(x) = \frac{1}{c} x^{\alpha-1} (1-x)^{\beta-1}$$

Where c :

$$c = \int_0^1 x^{\alpha-1} (1-x)^{\beta-1} dx$$

Beta Random Variable

X is a beta random variable: $X \sim \text{Beta}(\alpha, \beta)$
PDF where $\alpha, \beta > 0$:

$$f(x) = \begin{cases} \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1} \\ 0 \quad \text{otherwise} \end{cases}$$

$$B(\alpha, \beta) = \int_0^1 x^{\alpha-1} (1-x)^{\beta-1} dx$$

Basically with a Beta variable I have a certain prior

and i update it based on the observed data.
For example if i have:

$$P(p | H=3, T=1) \propto p^3 (1-p)^1$$

(p = probability of getting heads)

$$p \sim \text{Beta}(\alpha+H, \beta+T) = \text{Beta}(\alpha+H, \beta+T)$$

If i assume a uniform prior $\text{Beta}(1,1)$ Then the posterior would be $\text{Beta}(10, 2)$

In Bayesian statistics, before we see any data, we express our initial belief about a probability using a **prior distribution**. For a probability parameter p (e.g., the probability of getting heads in a coin flip), a common choice is the **Beta distribution**:

$$p \sim \text{Beta}(\alpha, \beta)$$

where:

- α is the number of **pseudo-observed** successes (e.g., heads).
- β is the number of **pseudo-observed** failures (e.g., tails).

The **Beta(1,1) distribution** is special because it is a **uniform prior**, meaning that before seeing any coin flips, we assume all values of p (from 0 to 1) are equally likely. This corresponds to saying, "We don't favor any particular probability of heads."

It's symmetric when $\alpha = \beta$:

$$\begin{aligned} E[X] &= \frac{\alpha}{\alpha+\beta} & \text{Var}(X) &= \frac{\alpha\beta}{(\alpha+\beta)^2 (\alpha+\beta+1)} \end{aligned}$$

Beta variables can come from experience

The Beta variable express how unconfident i am with my probability

What if the prior is Beta?

What if our prior belief about X was Beta:

$$F(X=x) = \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}$$

What is our posterior belief about x after observing n heads

Because we put our prior belief as Uni before, now we put in the posterior part as a Beta random variable.

$$f(X=x \mid N=n) = ?$$

so in this case i will not have $f(x)=1$ any more

$$\begin{aligned} f(X=x \mid N=n) &= \frac{P(N=n \mid X=x) f(X=x)}{P(N=n)} \\ &= \frac{\binom{n+m}{n} x^n (1-x)^m f(X=x)}{P(N=n)} \\ &= \frac{\binom{n+m}{n} x^n (1-x)^m \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}}{P(N=n)} \\ &= K_3 x^n (1-x)^m x^{\alpha-1} (1-x)^{\beta-1} \\ &= K_3 x^{n+\alpha-1} (1-x)^{m+\beta-1} \end{aligned}$$

If my prior is a Beta also my belief is a Beta

$$= X \mid N \sim \text{Beta}(\alpha+n, \beta+m)$$



This is called conjugate prior.

A beta understanding

If the prior distribution X (before seeing the data) is Beta

Then also the posterior is Beta

Beta is a conjugate distribution for Beta

- mathematically, conjugate means easy update
 \rightarrow add number of heads and tails seen to Beta parameters.

Difference between Bayesian and Frequentist.

Approach	Frequentist 🎲	Bayesian 🎲
Definition of Probability	Probability is the long-run frequency of an event occurring.	Probability is a measure of belief or degree of certainty.
What is p ? (e.g., coin flip probability)	p is a fixed but unknown number .	p is a random variable that follows a probability distribution.
How to Estimate p ?	Use sample proportions : $\hat{p} = \frac{\text{heads}}{\text{total flips}}$.	Use Bayes' theorem to update a probability distribution for p .
Uncertainty Representation	Uses confidence intervals based on sample data.	Uses a posterior distribution to represent uncertainty.
Prior Knowledge	Assumes no prior knowledge—only data matters.	Incorporates prior beliefs , which are updated with data.
Example: Is a Coin Fair?	Flip the coin many times. If the fraction of heads is 0.5, conclude it's fair.	Start with a belief (e.g., Beta(1,1)), then update with coin flip data.
Use Case	A/B testing, clinical trials, classic hypothesis testing.	Machine learning, Bayesian inference, small-sample learning.

Can set $X \sim \text{Beta}(a, b)$ as prior to reflect how biased you think coin is a priori.

- This is subjective probability (aka Bayesian)!
- The prior probability X based on seeing $(a+b-2)$ 'imaginary' flips where
 - $(a-1)$ of them were head
 - $(b-1)$ of them were tails
- Update to get posterior probability:
 $\rightarrow X | (n \text{ head and } m \text{ tails}) \sim \text{Beta}(a+n, b+m)$

Laplace Smoothing (Simple prior)

Prior: $X \sim \text{Beta}(\alpha = 2, \beta = 2)$

One imagined
head

One imagined
tail

Beta(2,2) is a popular prior to start with

If the expectation is the weighted average the point with the highest probability is called Mode

if $X \sim \text{Beta}(\alpha = 19, \beta = 8)$

$$E[X] = \frac{\alpha}{\alpha + \beta} = \frac{19}{19 + 8} \approx 0.70$$

$$\text{Mode}(x) = \frac{\alpha - 1}{\alpha + \beta - 2}$$

In the denominator I have the amount of success and it's divided by the number of total trials.

Adding Random Variables

IID Random Variables

Consider n random variables x_1, x_2, \dots, x_n

x_i are all independently and identically distributed (IID)

All have the same PMF/PDF

All have the same expectations

All have the same variance

ex:

2 players

$P(X=x)$ = amount of points scored by player 1

$P(Y=y)$ = amount of points scored by player 2

What is the probability of a tie?

$$P(\text{tie}) = \sum_{i=0}^{100} P(X=i, Y=i) \quad \text{Because they are independent}$$

The insight of convolution proof

What is the probability that $x+y=n$ ($P(x+y=n)$)

x	y	ind
0	n	0
1	$n-1$	1
2	$n-2$	2
.	:	.
:	:	:
n	0	n

They are mutually exclusive and countable
and cover all possible ways

In this case the 2 random variables are
dependent so I need to think about the joint
probability together

$$P(x+y) = \sum_{i=0}^n P(x=i, y=n-i)$$

If they are independent i need to do the summation
over the product of the probabilities.

$$P(x+y=n) = \sum_{i=0}^n P(x=i) P(y=n-i)$$

if we speak about continuous random variable
we would use an integral.

Sum of independent binomials

In some cases adding random variables is easy:

Let X and Y be independent binomials with the same value for p .

$$\therefore X \sim \text{Bin}(n_1, p) \text{ and } Y \sim \text{Bin}(n_2, p)$$

If we sum these 2 random variables we get $Z \sim \text{Bin}(n_1+n_2, p)$

When we sum we obtain a random variable because the result
is not deterministic.

Intuition:

X has n_1 trials and Y has n_2 trials

Each trial has some success probability p

Define Z to be n_1+n_2 trials, each with success probability p .

$$Z \sim \text{Bin}(n_1+n_2, p) \text{ and so } Z = X+Y$$

Generally, have n independent random variables

$$X_i \sim N(\mu_i, \sigma_i^2) \text{ for } i=1, 2, \dots, n$$

$$\sum_{i=1}^n X_i \sim N\left(\sum_{i=1}^n \mu_i, \sum_{i=1}^n \sigma_i^2\right)$$

Linear Transformation

γ is a linear transformation

$$X \sim N(\mu, \sigma^2)$$

$$\gamma = X + \lambda = 2 \cdot X$$

$$\gamma \sim N(2\mu, 2\sigma^2)$$

ex:

Gaussian Sampling and Elo Ratings

Each team has an elo score s , calculated on its past performance
each game a team has ability $A \sim N(s, 200^2)$
the team with the higher sampled ability wins
What is the probability that warios win a game?

$$P(A_w > A_o) ?$$

$$\text{Warios': } A_w \sim N(s=1657, 200^2)$$

$$\text{Opponent's } A_o \sim N(s=1670, 200^2)$$

Probability of winning a game

$$P(\text{Warios win}) = P(A_w > A_o)$$

|

$$= P(A_w - A_o > 0)$$

I apply a linear transformation to A_o

$$-1 \cdot A_o = -A_o \sim (-1670, (-1)^2 200^2)$$

The difference then would be a sum

$$\Delta = A_w + (-A_o) = (1657 - 1670, 200^2 + 200^2) = (X, 2 \cdot 200^2)$$

Then i can ask about what is the probability that this is > 0

$$P(\Delta > 0) = 1 - F_\Delta(0) \approx \gamma$$

We don't need anymore to complete we can solve like this.

Discrete vs Continuous

$$P(X+Y = \alpha) = \sum_{y=-\infty}^{+\infty} P(X=\alpha-y)P(Y=y)$$

$$f(X+Y = \alpha) = \int_{-\infty}^{+\infty} f(X=\alpha-y)f(Y=y) dy$$

Sum of Independent uni forms

Let x and y be independent variables

$$X \sim \text{Uni}(0,1) \text{ and } Y \sim \text{Uni}(0,1) \rightarrow f(x) = 1 \text{ for } 0 \leq x \leq 1$$

If you think about it: $(1 \leq \alpha \leq 2)$

$$f(X+Y = \alpha) ?$$

$$f(X+Y = \alpha) = \int_{-\infty}^{+\infty} F(X=\alpha-y)f(Y=y) dy$$

$$f(X+Y = \alpha) = \begin{cases} \alpha & 0 < x < 1 \\ 2-\alpha & 1 < x < 2 \\ 0 & otherwise \end{cases}$$

Since X and y are independent
is using the convolution
formula by multiplying the CDF
of one result by another
probability.

The sum of random variables isn't always a clean answer

Just look at the result coming from the sum of 2 Uni

Central Limit Theorem

Consider n independent and identically distributed variables X_1, X_2, \dots, X_n with $E[X_i] = \mu$ and $\text{Var}(X_i) = \sigma^2$

$$\sum_{i=1}^n X_i \sim N(n\mu, n\sigma^2) \text{ As } n \rightarrow \infty$$

With a huge number of i.i.d. RV, if I sum all of them, no matter what the result will be a normal distribution.

Proof of CLT

The Fourier transform of a PDF is called a characteristic function.
Take the characteristic function of the probability mass of the sample
normalize

Shows that this approaches an exponential function in the limit
as $n \rightarrow \infty$: $f(x) = e^{-x^2/2}$

This function is in turn the characteristic function of a Normal
Distribution.

Average of IID Variables

Let x_i be n i.i.d variables and \bar{x} be the average.

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

↓ Gaussian by CLT
 $N(\mu, \sigma^2)$

By the central limit theorem the mean of IID variables are distributed
normally. As $n \rightarrow \infty$

$$\bar{x} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

Coming from the variance
scaling property:

$$\text{Var}\left(\frac{1}{n} \sum_{i=1}^n x_i\right) = \frac{1}{n^2} \text{Var}\left(\sum_{i=1}^n x_i\right)$$

Sampling Definitions

I want to know the true mean and variance of happiness in Bhutan.

- Obviously I can't ask to all
- I ask to 200 people

$$\text{Happiness} = \{u_1, u_2, \dots, u_{200}\}$$

$$\mu = 83$$

Is it the true mean of happiness?

Mathematically, what is a sample?

Consider n random variables:

The sequence of random variables is a sample from the distribution F if:

- All the R.V. are iid
- All the variables have the same distribution function F where $E[X_i] = \mu, \text{Var}(X_i) = \sigma^2$

Estimating Cov Statistics (μ, σ^2)

How to estimate the cov statistics giving few samples.

How good is our estimation of the mean for example

$$E[\bar{X}] = E\left[\frac{1}{n} \sum_{i=1}^n X_i\right] = \frac{1}{n} E\left[\sum_{i=1}^n X_i\right] = \frac{1}{n} \sum_{i=1}^n E[X_i] = \frac{1}{n} n\mu = \mu$$

In this case this mean is called unbiased estimator of the population mean μ .

Intuition About the sample variance σ^2

For example population variance:

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$

↓
population mean

Now let's see the sample variance:

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

↓
sample mean

\bar{x} and μ surely differ because \bar{x} is the mean of the subset of the real population.

We put $n-1$ in the denominator because the variance we are calculating is an estimate using an estimate so it needs to be scaled a little bit more

Proof that s^2 is unbiased.

$$\begin{aligned} E[S^2] &= E\left[\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2\right] \Rightarrow (n-1)E[S^2] = E\left[\sum_{i=1}^n (X_i - \bar{X})^2\right] \\ (n-1)E[S^2] &= E\left[\sum_{i=1}^n ((X_i - \mu) + (\mu - \bar{X}))^2\right] \quad (\text{introduce } \mu - \mu) \\ &= E\left[\sum_{i=1}^n (X_i - \mu)^2 + \sum_{i=1}^n (\mu - \bar{X})^2 + 2 \sum_{i=1}^n (X_i - \mu)(\mu - \bar{X})\right] \\ &= E\left[\sum_{i=1}^n (X_i - \mu)^2 + n(\mu - \bar{X})^2 - 2n(\mu - \bar{X})^2\right] \\ &= E\left[\sum_{i=1}^n (X_i - \mu)^2 - n(\mu - \bar{X})^2\right] = \sum_{i=1}^n E[(X_i - \mu)^2] - nE[(\bar{X} - \mu)^2] \\ &= n\sigma^2 - n\text{Var}(\bar{X}) = n\sigma^2 - n\frac{\sigma^2}{n} = n\sigma^2 - \sigma^2 = (n-1)\sigma^2 \quad \text{Therefore } E[S^2] = \sigma^2 \end{aligned}$$

Standard error of the mean

$$E[\bar{X}] = \mu$$

We want to estimate $\frac{\sigma^2}{n}$

Def: The standard error of the mean is an estimate of the standard deviation of \bar{X}

$$SE = \sqrt{\frac{s^2}{n}} = \frac{s}{\sqrt{n}}$$

This measures how we think our mean estimate is.

Intuition:

s^2 is an unbiased estimate of σ^2

s^2/n is an unbiased estimate of $\text{Var}(\bar{X})$
 $\sqrt{s^2/n}$ is an estimate of $\sqrt{\text{Var}(\bar{X})}$

The estimated standard error is a measure of how much the sample mean varies across different samples (random samples = subset of the population)

The value of confidence in my claim is called the p-value

Bootstrap: Randomize algorithm

This algorithm allows you to:

- Know the distribution of statistics
- Calculate p values.
- Using computers.

How can i know how bad my estimated variance is if i have the possibility to know the true value coming from the real life?

Since we are using a subset of the entire dataset. The variance that i calculate is the one from the subset and it's called sample variance.

And i want to know how much this number change if i choose a different subset of other 200 people (This is basically the variance of the variance)

The sample variance is a random variable so it has a distribution.

The first half of the bootstrap is basically to choose 10000 samples, obtain the PUF and then take all needed conclusions.

Bootstrapping Assumption

$$F \approx F'$$

↓ → the sample distribution
the underlying distribution

The distribution i get from the histogram of my samples is a good approximation of the true distribution

Algorithm

- 1- Estimate the PUF using the sample
- 2- Repeat 10000 times
 - a- Resample k m (sample) from PUF
 - b- Recalculate the stat on the resample
- 3- You now have a distribution of your stat

In an easier way:

We take 10000 subsets from the main set

For every subset we calculate the variance, the mean, the median.

With the 10000 values we create a distribution.

At the end I have different PMFs for different metrics.

Considering the central limit theorem after the calculation of all the means if i sum all the values i will obtain a normal distribution

Then after bootstrapping i could've known the probability that the true mean is between a range of values

To calculate the probability of this range i can also simply count the number of values in the list that are within the range.

Bootstrap tells me how wrong I might be about a value

If i get a p-value < 0,05 is accepted

How can we use bootstrapping with the null hypothesis?

Null hypothesis: There is no true difference between random samples groups. The only differences that we observe are due to random chance / sampling error.

The bootstrapping in this case is useful due to the fact that we can check if the difference coming from the 2 groups is significant or just random noise

Basically the p value is the probability that our error / difference comes from noise and not due to errors.

Algorithmic Analysis

ex: Coupon Collector Problem.

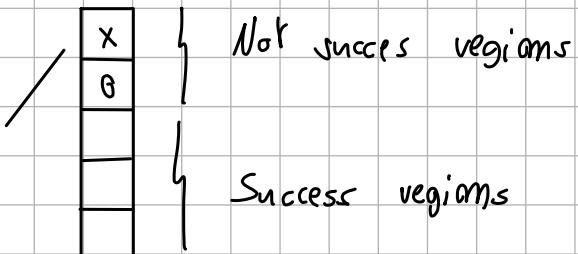
Consider a hash table with n buckets.

Each string equally likely to get hashed into any bucket.

Let $X = \# \text{ strings to hash until each bucket} \geq 1 \text{ string}$

What is $E[X]?$

string \rightarrow Hash function



5 buckets

In the first attempt we need 7 rounds to fill the hash table

Let $X_i = \#$ of trials to get success after i -th success
where success is hashing string to previously empty bucket
After i buckets have ≥ 1 string, probability of hashing a
string to an empty bucket is $P = \frac{(n-1)}{n}$

$$P(X_i = k) = \frac{n-i}{n} \left(\frac{i}{n}\right)^{k-i} \text{ equivalently } X_i \sim \text{Geo}\left(\frac{n-1}{n}\right)$$

$$E[X_i] = \frac{1}{P} = \frac{n}{n-i}$$

$$X = X_0 + X_1 + \dots + X_n \rightarrow E[X] = E[X_0] + E[X_1] + \dots + E[X_{n-1}]$$

Every bucket could be considered as a random variable so it'll
help to sum n different random variable.

$$E[X] = \frac{n}{n} + \frac{n}{n-1} + \frac{n}{n-2} + \dots + \frac{n}{1} = n \left(\frac{1}{n} + \frac{1}{n-1} + \dots + 1 \right) = O(n \log n)$$

This sum is called
harmonic sum and the
result is constant $\log n$

Conditional Expectation

X and Y are jointly discrete random variables:

We define the conditional expectation of x given $y = y$

$$E[X | Y = y] = \sum_x x P(X=x | Y=y) = \sum_x x p_{x|y}(x|y)$$

For continuous random variable:

$$F_{x|y}(x|y) = \frac{f_{x,y}(x,y)}{f_y(y)}$$

$$E[X|Y] = \int_{-\infty}^{+\infty} x f_{x|y}(x,y) dx$$

ex:

Roll 2 dices

$$X = \text{value of } D_1 + D_2 \quad Y = \text{Value of } D_2$$

$$\text{What is } E[X | Y=6]$$

$$P(X|Y) = \frac{P(XY)}{P(Y)} = \frac{1}{36} = \frac{1}{6}$$

$$E[X|Y] = \sum_{x=7}^{12} P(X=x | Y=6) = \frac{1}{6} (7+8+9+10+11+12) = \frac{57}{6}$$

Considering now $g(Y) = E[X|Y]$, it's a function of y .

In our case $E[X]$ is a value while $E[X|Y=y]$ is a function of y . Writing $E[X=5]$ doesn't make sense.

Law of total expectation

$$\begin{aligned} E[E[X|Y]] &= \sum_y E[X|Y=y] P(Y) & g(Y) &= E[X|Y] \\ &= \sum_y \sum_x x P(X=x | Y=y) P(Y) \\ &= \sum_y \sum_x x p(x,y) \\ &= \sum_x \sum_y x p(x,y) \end{aligned}$$

$$= \sum_x \sum_y p(x, y)$$

$$= \sum_x p(x) = E[X]$$

In this passage I use the definition of marginal probability by summing over all the values of y .

Analyzing Recursive Code

```
int Recurse() {
    int x = randomInt(1, 3); // Equally likely values
    if (x == 1) return 3;
    else if (x == 2) return (5 + Recurse());
    else return (7 + Recurse());
}
```

Let Y = value returned by $\text{Recurse}()$.

What is $E[Y]$?

We can see x as a distribution so it can be considered as a random value.

$$E[Y] = p(x=1)E[Y|X=1] + p(x=2)E[Y|X=2] + p(x=3)E[Y|X=3]$$

$$E[Y|X=1] = 3$$

$$E[Y|X=2] = E[5+Y] = 5 + E[Y]$$

$$E[Y|X=3] = E[7+Y] = 7 + E[Y]$$

$$E[Y] = 3\left(\frac{1}{3}\right) + (5+E[Y])\left(\frac{1}{3}\right) + (7+E[Y])\left(\frac{1}{3}\right) = \left(\frac{1}{3}\right)(15 + 2E[Y])$$

Machine Learning Estimation

While we have parameters

Consider different probability distributions:

Every one of them have special numbers for example the mean
The std ... etc

We call these prob distributions with these numbers parametric models

I don't have to describe all the points in the curve I just need to give the value of the parameters.

Given model, parameters yield actual distribution

- Usually refer to parameters of distribution as θ
- The θ can also be a vector of parameters.

With bootstrapping we already have a method to estimate some parameters for example the variance and the mean

Likelihood Given The Parameters

$$\text{Likelihood} = \prod_{i=1}^N f(x_i \mid \mu=5, \sigma=1)$$

$N =$ all the dataset

It's the PDF of all the points given some values of μ and σ in this case but it's in general

The likelihood is the product between all the PDF of all the points given a certain value of mean and standard deviation

Difference between max and argmax

If I have the function $f(x) = -x^2 + 5$

$\max_x -x^2 + 5$ will return the maximum value that the function can return

$\text{argmax}_x -x^2 + 5$ will return 0 because it return the maximum input used by the function to return the maximum value.

The argument that maximizes a function is also the argument that minimizes the log of the function

\rightarrow This is because log is monotonic so the larger the input the larger the output

We use the logarithm because computers are not very good at representing very very small numbers but with the log i am able to represent numbers with a little bit of sanity.

Maximum Likelihood Estimation

1. Decide on a model for the distribution of your samples. Define the PMF / PDF for your sample.

2. Write out the log likelihood function.

3. State that the optimal parameters are the argmax of the log likelihood function.

4. Use an optimization algorithm to calculate argmax

The likelihood function

$$L(\theta) = \prod_{i=1}^n f(x_i | \theta)$$

↓ ↓
n i.i.d. parameters
data points distribution

So i'm calculating the probability that certain R.V. all together have certain values.

In order to have all the values for the random variables at the same time i need to multiply their probability

The result will be how likely the data are generated given a certain model and certain parameters.

Maximum likelihood

$$L(\theta) = \prod_{i=1}^n f(x_i | \theta)$$

$$LL(\theta) = \sum_{i=1}^n \log(f(x_i | \theta))$$

$$\hat{\theta} = \arg \max_x LL(\theta)$$

MLE For Poisson

$$x \sim \text{Poi}(\lambda)$$

Now : need to find The maximum likelihood

$$\text{The PMF can be written as: } f(x_i | \lambda) = \frac{e^{-\lambda} \lambda^{x_i}}{x_i!}$$

$$\text{The likelihood can be written: } L(\theta) = f(x_1, \dots, x_n | \lambda) = \prod_{i=1}^n f(x_i | \lambda) = \prod_{i=1}^n \frac{e^{-\lambda} \lambda^{x_i}}{x_i!}$$

$$\begin{aligned} \text{The log likelihood} &= \sum_{i=1}^n \log \left(\frac{e^{-\lambda} \lambda^{x_i}}{x_i!} \right) \\ &= \sum_{i=1}^n -\lambda + x_i \log \lambda - \log(x_i!) \end{aligned}$$

Now ; just need to do the derivative to find The point in which The value of The function is at max.

$$\text{If I put that equal to 0 I'll have } \lambda = \frac{1}{n} \sum_{i=1}^n x_i$$

MLE For Bernoulli

$$x \sim \text{Bern}(p)$$

$$LL(\theta) = \sum_{i=1}^n \log f(x_i | p)$$

$$f(x_i | p) = \begin{cases} p & \text{if } x_i = 1 \\ 1-p & \text{if } x_i = 0 \end{cases}$$

Given that This function isn't differentiable we need to find a way

Differentiable PDF For Bernoulli

Consider i.i.d random variables x_1, x_2, \dots, x_n

$$x_i \sim \text{Bern}(p)$$

Given that the continuous form gives the same result we can write $f(x_i | p) = p^{x_i} (1-p)^{1-x_i}$

But this form happens to have sense just for 0 and 1. That's what we need in case of Bernoulli and it's also differentiable.

$$\text{Likelihood} = \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i}$$

$$LL(\theta) = \sum_{i=1}^n \log(p^{x_i} (1-p)^{1-x_i}) = \sum_{i=1}^n [x_i (\log p) + (1-x_i) \log(1-p)]$$

If we differentiate and set to 0:

$$y = \sum_{i=1}^n x_i$$

$$\frac{\partial LL(\theta)}{\partial p} = y \frac{1}{p} + (n-y) \frac{-1}{1-p} = 0 \rightarrow \boxed{p_{MLE} = \frac{y}{n} = \frac{1}{n} \sum_{i=1}^n x_i}$$

MLE For Gaussian

$$X \sim N(\mu, \sigma^2)$$

$$\text{Let } x_i \sim N(\mu, \sigma^2) \quad f(x_i | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$$

$$LL(\theta) = \sum_{i=1}^n \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} \right) = \sum_{i=1}^n \left[-\log \left(\frac{\sqrt{2\pi\sigma^2}}{2\sigma^2} \right) - \frac{(x_i - \mu)^2}{2\sigma^2} \right]$$

$$= -\sum_{i=1}^n \log \left(\sqrt{2\pi\sigma^2} + \frac{\sum_{i=1}^n [(x_i - \mu)^2]}{2\sigma^2} \right)$$

Let's derive with respect to μ

$$\frac{\partial LL(\theta)}{\partial \mu} = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) = 0$$

I can derive with respect to σ :

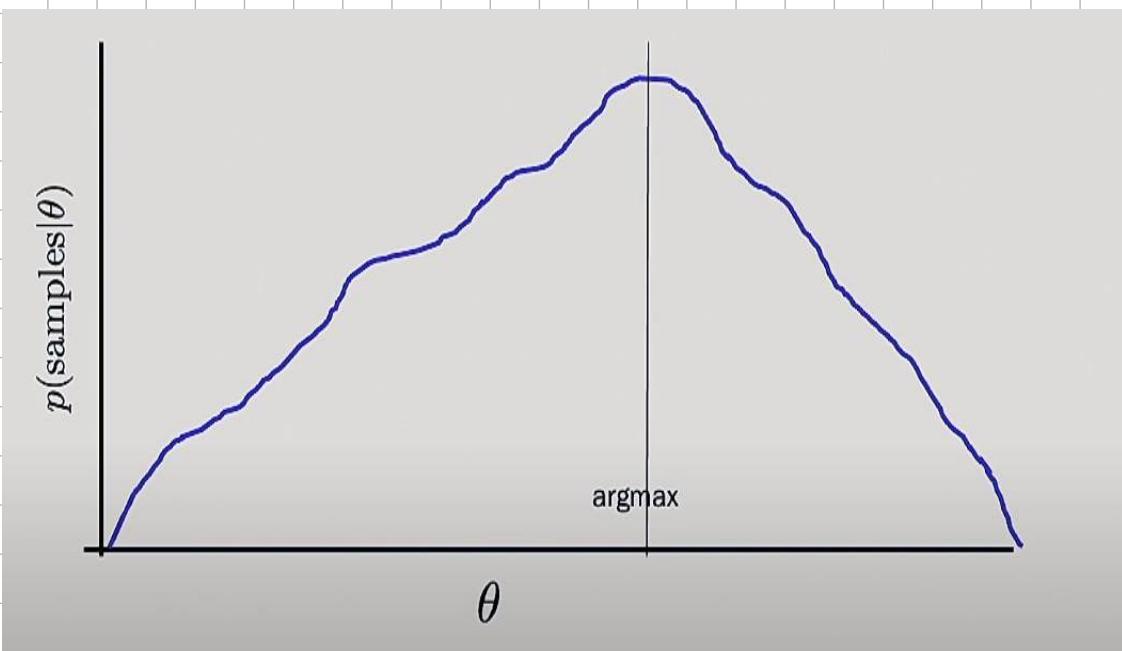
$$\frac{\partial \text{LL}(\sigma)}{\partial \sigma} = -\frac{n}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^n (x_i - \mu)^2 = 0$$

In this case I have 2 unknowns and 2 equations

Solving for μ_{MLE} : $\frac{1}{n} \sum_{i=1}^n x_i$

Solving for σ_{MLE} : $\sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \mu_{MLE})^2}$

Gradient Ascent



If my step size is small enough I can walk uphill and find a local maxima

With argmax I will have lot of imports and it can become computationally super expensive

I need to take the θ curve, taking the derivative and then follow the path where the derivative is positive

If I repeat this operation I will reach the top of the curve, at that point the values of the derivative in all the direction becomes negative.

Small step size guarantees you to reach the top but it will take more time

Repeat many times:

$$\theta_j^{\text{new}} = \theta_j^{\text{old}} + \eta \frac{\partial \text{LL}(\theta^{\text{old}})}{\partial \theta_j^{\text{old}}}$$

At the beginning
this is my initial guess.
↓
step size constant

In every optimisation algorithm i want to use gradient descent because i don't want to find the max but i want to find the minimum of a loss. In our case here we are maximizing the likelihood but you can also minimize the negative likelihood.

The problem with MLE is the fact that it overfit the data it sees it doesn't generalize the information. it doesn't think about data it hasn't seen.

Maxima A Posteriori

Basically the same as the Beta's variable but it doesn't apply only Bernoulli variable but also for other variables.

MLE Vs MAP

Data: x_1, \dots, x_n

Maximum likelihood estimation:

$$\hat{\theta}_{\text{MLE}} = \arg \max f(x_1 = x_1, \dots, x_n = x_n | \theta)$$

Maximum a posteriori:

$$\hat{\theta}_{\text{MAP}} = \arg \max f(\theta = \theta | x_1 = x_1, \dots, x_n = x_n)$$

I will obtain the parameters that are most likely given the values of the data.

Given that the parameter is the unknown and the data are the observable we can say that this is a more natural approach.

Notation Shorthand

$$\theta = \Theta = \theta$$

$$x^i = X^i = x^i$$

→ The i is in the apex because in all the data points could have multiple values for the same data point

$$\hat{\theta}_{MAP} = \arg \max f(\theta | x_1, \dots, x_n)$$

Basically it returns whatever point of the θ with the highest probability.

So basically this is used to avoid the overfitting of the data

Conjugate Distribution

MAP estimation:

$$\hat{\theta}_{MAP} = \arg \max f(\theta | x_1, x_2, \dots, x_n)$$

The mode of the posterior distribution of θ

The conjugate allow us to simplify the math behind the data

$$\arg \max_{\theta} f(\theta = \theta | \text{Data}) = \arg \max f(\text{Data} | \theta) f(\theta)$$

/

prior belief of
my parameters
How I think my
parameters are
before any data is
seen

Multinomial

Dirichlet ($\alpha_1, \alpha_2, \dots, \alpha_m$) is a conjugate for mult. nomial.

- It generalizes in the same way multinomial generalizes Bernoulli and binomial

Prior: Suppose $\sum_{i=1}^m \alpha_i = m$ imaginary trials, with α_i of outcome i

Naive Bayes

$$\text{MLE: } \underset{\theta}{\operatorname{argmax}} \left(\sum_i \log f(x_i | \theta) \right)$$

$$\text{MAP: } \underset{\theta}{\operatorname{argmax}} \left(\log(g(\theta)) + \sum_{i=1}^n \log(f(x_i | \theta)) \right)$$

Training Data

Assume iid data:

$$(x^1, y^1), (x^2, y^2) \dots$$

$$m = |x^i|$$

each data point has m features and a single output.

The output is basically if the user likes a particular element