

Distributions in probability

We are doing some experiment and collecting data. Now we want to see what is the distribution of probability of getting given observation. Distribution gives you theoretical model to define chance of observation value or range of value for given experiment(collection of data).

1. Bernauli distribution : Bernauli(p)

Simplest distribution of experiment whose outcomes are binary. Either 0 or 1.

Gender : Male(1) or Female(0)

coin toss : Head(1) or Tail(0)

Parameter p is probability of getting value 1 or true probability.

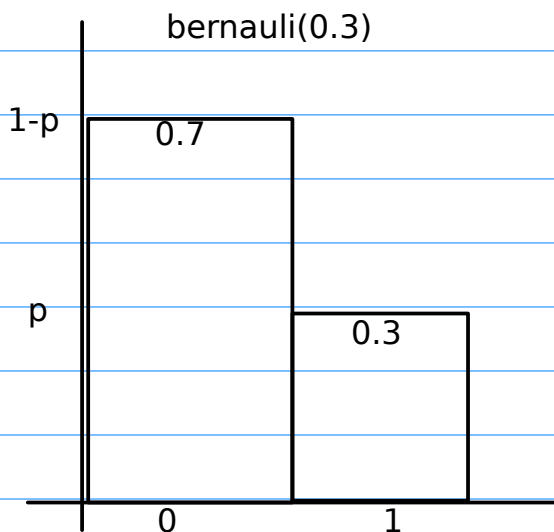
False probability (probability of getting value 0) $q = 1 - p$

$$\text{PMF} = P(X = x) = \begin{cases} p & \text{if } x = 1 \\ 1 - p & \text{if } x = 0 \end{cases}$$

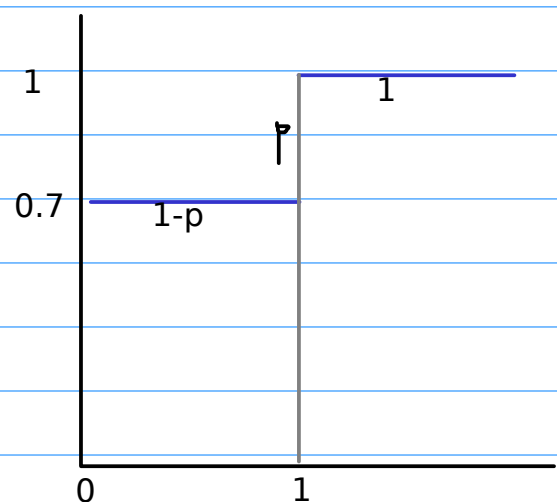
$E(x) = \text{mean} = p$

$\text{Var}(x) = pq = p(1-p)$

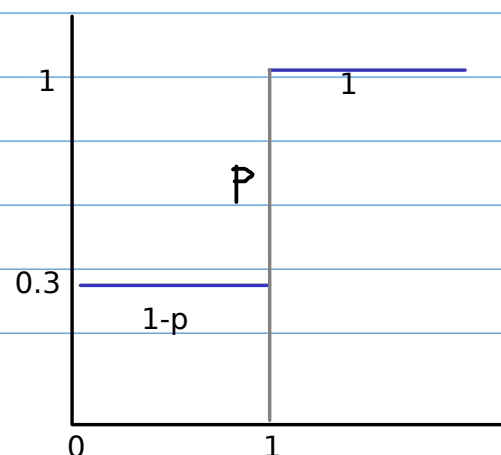
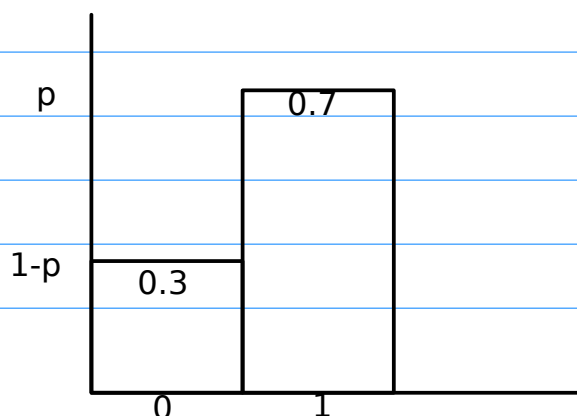
PMF



CDF



bernauli(0.7)



2. Binomial Distribution: Binomial(n, p)

Probability of getting value 1 for k times out of n times is distributed as binomial.

p = 0.5 probability of getting head

What is the probability of getting 7 heads(1) if you toss a coin 10 times.

p = 0.1 probability of getting faulty bulb

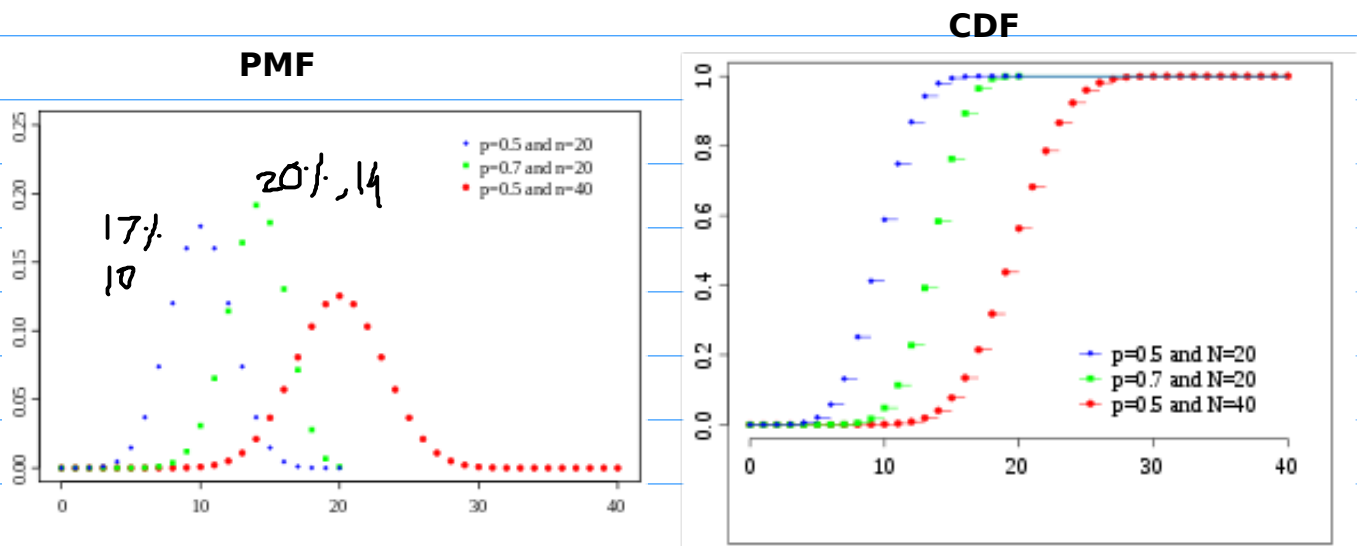
What is the probability of getting 2 faulty bulbs(1) in box of 20.

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

CDF
$$P(X \leq k) = \sum_{i=0}^k \binom{n}{i} p^i (1 - p)^{n-i}$$

$E(x) = \text{mean} = np$

$\text{Var}(x) = npq = np(1-p)$



for n = 20 we have two values of p 0.5 and 0.7

blue

When p = 0.5, it is high probability to have 10 points to be equal to 1

Green

When p = 0.7, it is high probability to have 14 points to be equal to 1

3. Uniform (Continuous) Distribution: $\text{uniform}(a, b)$

Probability of getting any value is similar across range a to b .

uniform colour, weight of cricket ball, random value generator

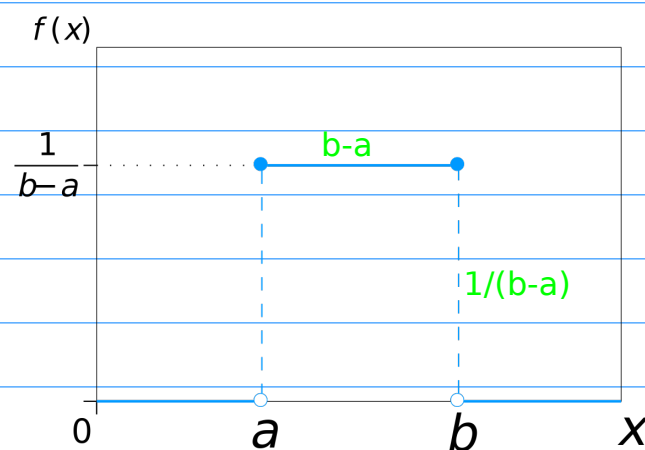
$$\text{PDF} = \begin{cases} = \frac{1}{b-a} & \text{if } x \in [a, b] \\ = 0 & \text{otherwise} \end{cases}$$

$$\text{CDF} = \begin{cases} = 0 & \text{if } x < a \\ = \frac{x-a}{b-a} & \text{if } x \in [a, b] \\ = 1 & \text{if } x > b \end{cases}$$

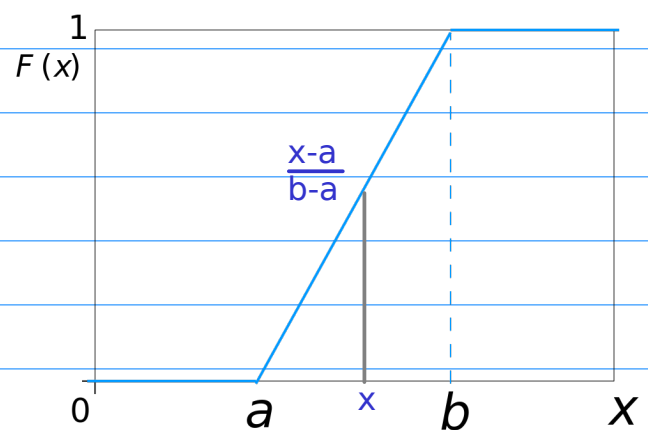
$$E(x) = \text{mean} = (a + b)/2$$

$$\text{Var}(x) = (b - a)^2/12$$

PDF



CDF



4. Poisson Distribution: $\text{poisson}(\lambda)$

Some event is occurring λ times per given time interval (sec, min, hr, day,...).
What is the probability that it will occur x times in given time.

Let's say $\lambda=10$, means we are getting average 10 calls per hour. Now we want to calculate what is the probability that you will get 10, 9, 11, 15, 6 calls in hour?

This probability distribution for given λ is called poisson distribution.

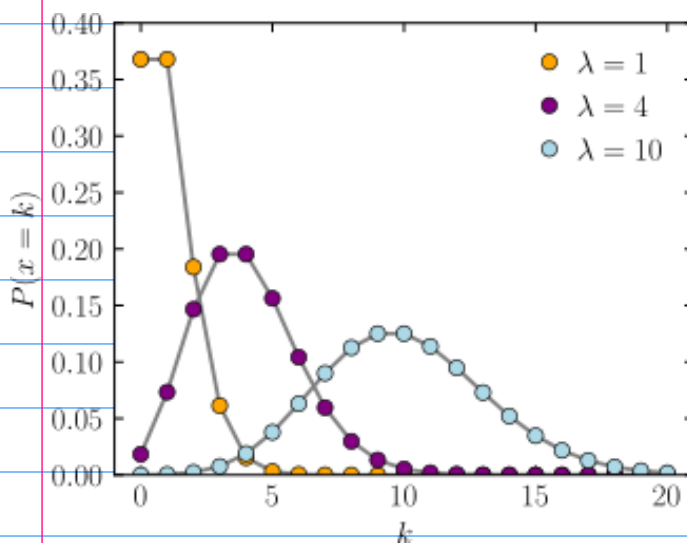
$$\text{PMF} = P(X=k) = f(k) = \frac{\lambda^k e^{-\lambda}}{k!}$$

$$\text{CDF} = P(X \leq k) = F(k) = e^{-\lambda} \sum_{i=0}^{\lfloor k \rfloor} \frac{\lambda^i}{i!}$$

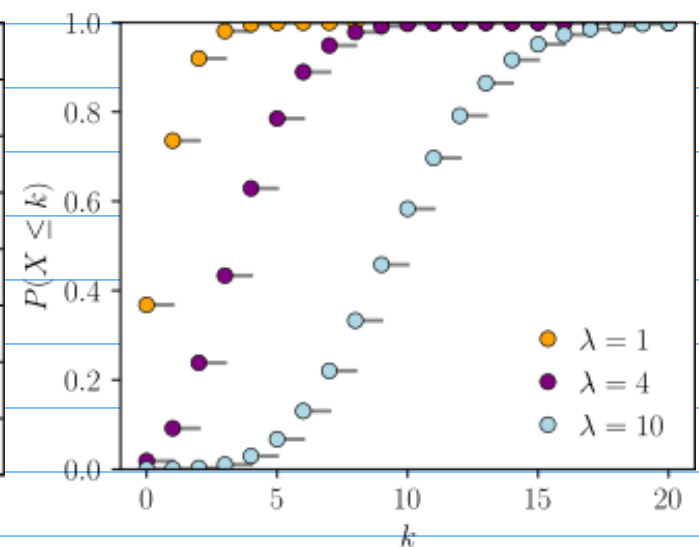
$$E(x) = \text{mean} = \lambda$$

$$\text{Var}(x) = \lambda$$

PMF



CDF



If X_1 and X_2 are two independent random variable with poisson distribution with rates λ_1 and λ_2 then we can say, (X_1+X_2) is poisson distribution of rate $(\lambda_1+\lambda_2)$

5. Exponential Distribution: $\text{exponential}(\lambda)$, $\beta=1/\lambda$

Some event is occurring λ times per given time interval (sec, min, hr, day,...).
What is the probability that event will occur after t time.

Let's say $\lambda=10$, means we are getting average 10 calls per hour. Now we want to calculate what is the probability that you will get call in 10 min, 30 min, 1hr, 1.5hr?

This probability distribution for given $\beta=1/\lambda$ is called exponential distribution which is inverse of poisson distribution. Or it is probability of time distribution between poisson events (with rate λ).

$$\text{PMF} = f(x) = \begin{cases} \lambda e^{-\lambda x} & \text{if } x \geq 0 \\ 0 & \text{if } x = 0 \end{cases}$$

$$\text{CDF} = F(x) = \begin{cases} 1 - e^{-\lambda x} & \text{if } x \geq 0 \\ 0 & \text{if } x = 0 \end{cases}$$

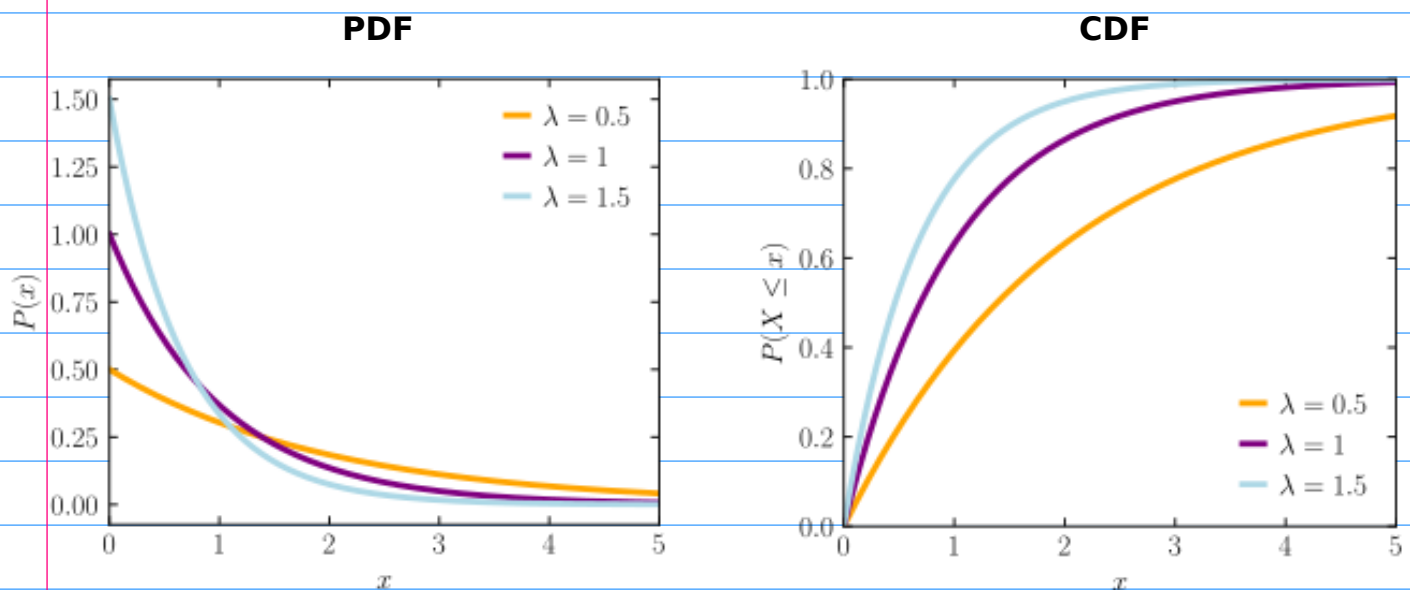
$$E(x) = \text{mean} = \beta = 1/\lambda$$

$$\text{Var}(x) = \beta^2 = 1/\lambda^2$$

Exponential distribution is memory less distribution. Probability that event will occur in at least t time is same as it occur after waiting for s time.

$$\Pr(X > s+t \mid X > s) = \Pr(X > t)$$

$$\Pr(\text{event occur in 40 min} \mid \text{waited for 30 min}) = \Pr(\text{event occur in 10 min})$$



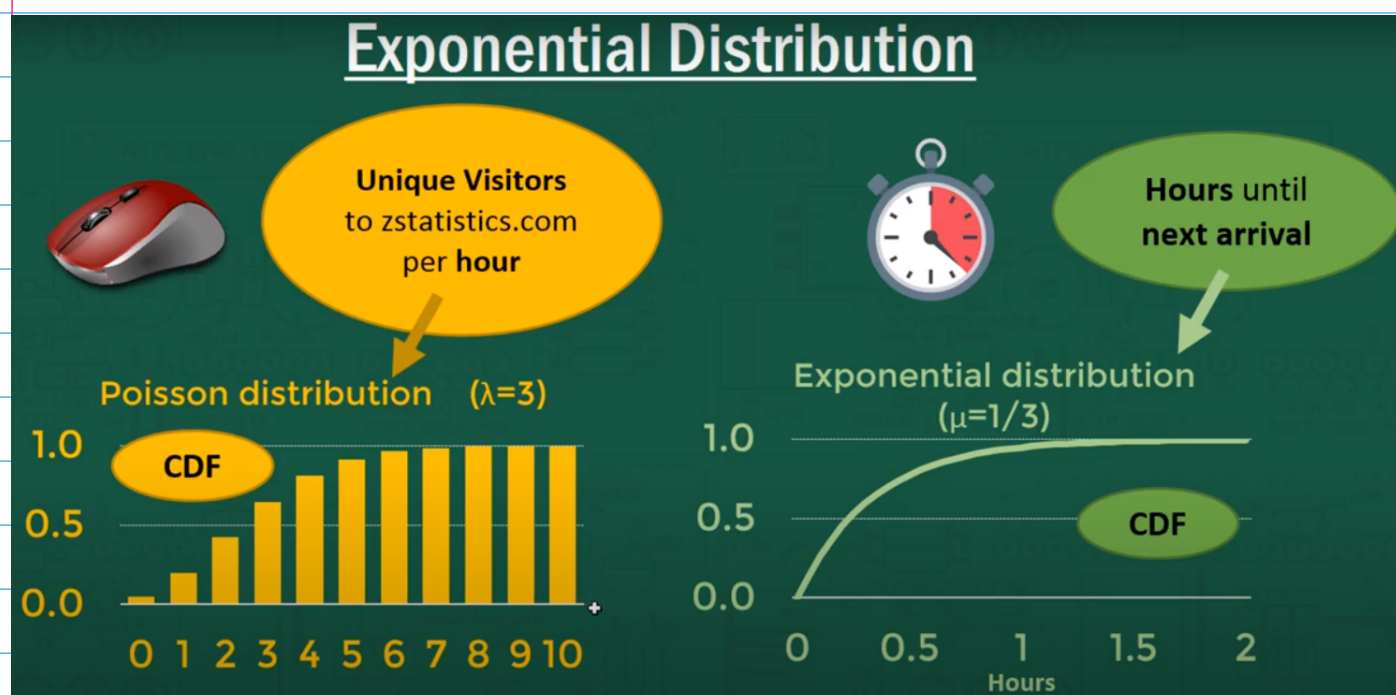
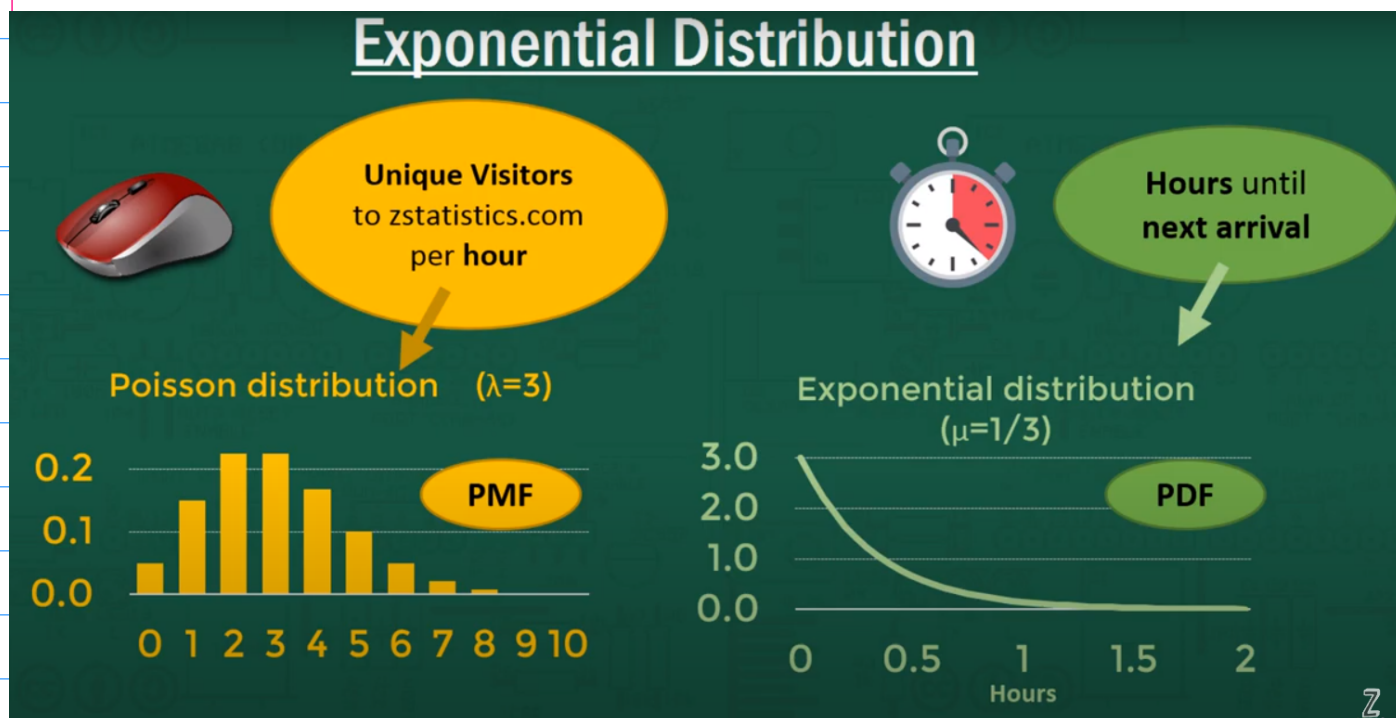
Let's take an example of unique visitors on website. Let's say number of unique visitors on webpage is poisson distributed with $\lambda = 3$ /hour. So you get average 3 unique visitors per hour or $\beta = 1/3$ hours per visitor. That means on an average you get visitor at 20 minutes interval.

Note= In below pictures we have taken β as $\mu = 1/\lambda$

PDF: We can observe that it is really unlikely to get visitor at 2 hours interval.

CDF: It is really likely (~100%) that we get visitor in 2 hours.

It is more likely to be interval of 0 to 30 min with $\beta = \mu = 20\text{min}$



6. Normal/Gaussian Distribution: $N(\mu, \sigma^2)$

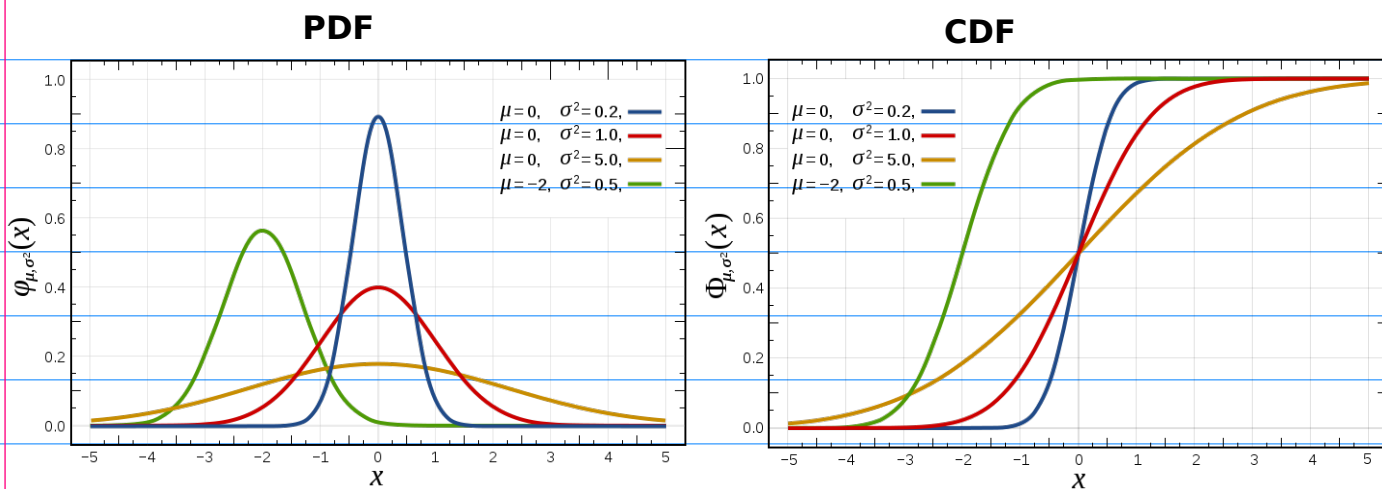
There are so many natural phenomenon that follow normal distribution, which is also known as Gaussian distribution. Normal distribution shows that large part of population tends to have value near to mean value μ and as you move away from the mean value probability start decreasing according to variance or σ^2 standard deviation σ . It forms bell curve of probability.

Height weight and length of living things follows Normal Distribution.

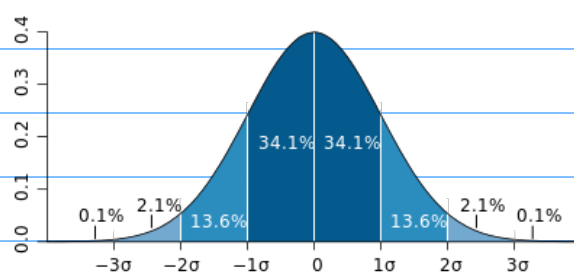
When we draw large number of samples from population of unknown distribution it follows Normal distribution.

$$\text{PDF}=f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} \quad \text{CDF}=F(x) = P(X \leq x) = \int_0^x f(x)dx$$

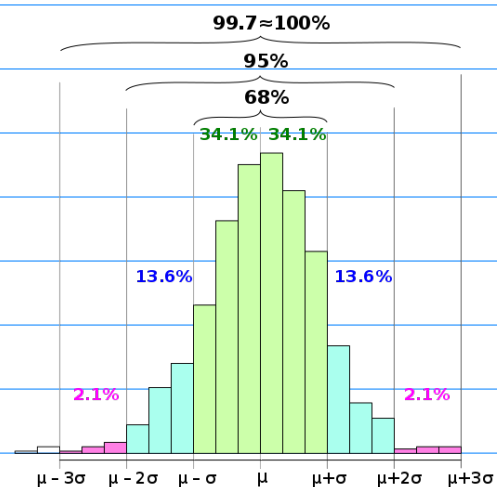
$$E(x) = \text{mean} = \mu \quad \text{Var}(x) = \sigma^2 \quad \text{standard deviation} = \sigma$$



Standard Deviation and Coverage: 68-95-99.7 rule



$$\begin{aligned} \Pr(\mu - 1\sigma \leq X \leq \mu + 1\sigma) &\approx 68.27\% \\ \Pr(\mu - 2\sigma \leq X \leq \mu + 2\sigma) &\approx 95.45\% \\ \Pr(\mu - 3\sigma \leq X \leq \mu + 3\sigma) &\approx 99.73\% \end{aligned}$$



7. Log Normal Distribution: lognormal(μ, σ^2)

There are so many phenomenon whose logarithm follow normal distribution, which is known as Log Normal Distribution. That means if random variable X is log normally distributed than $\ln(X)$ follows normal/gaussian distribution.

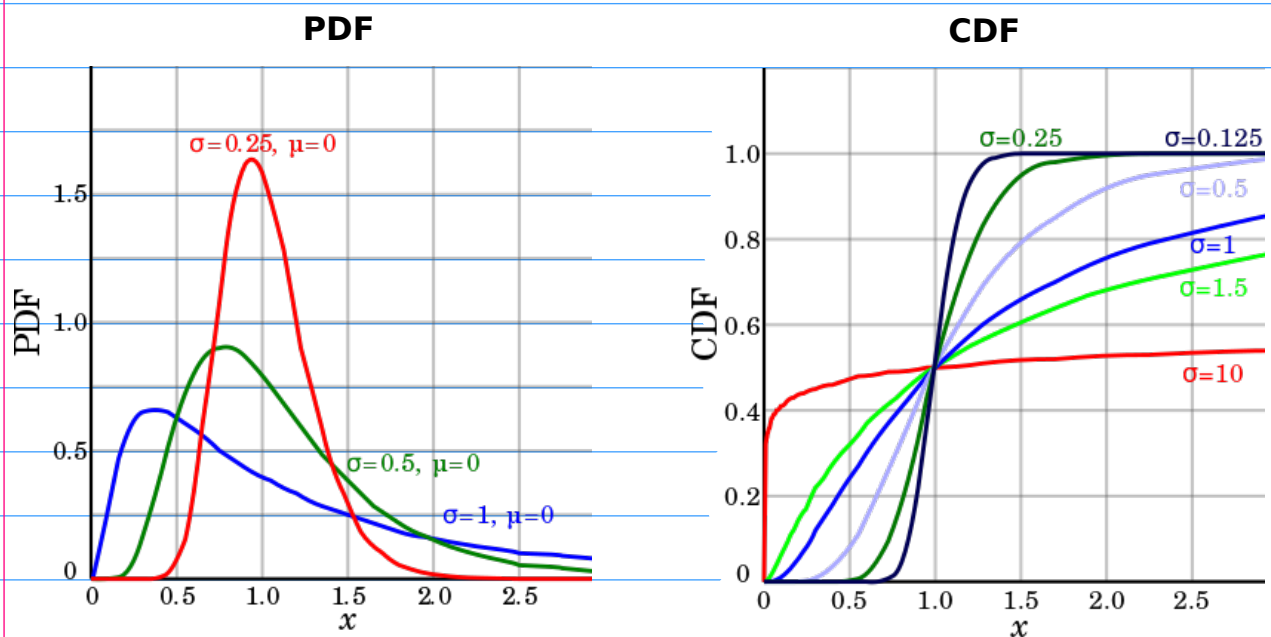
If something is more likely to happen in short interval and then probability of that happening decrease slowly over time is generally Log normally distributed.

Length of Comments posted in internet discussion forums

User's dwell time on online articles. Length of chess games.

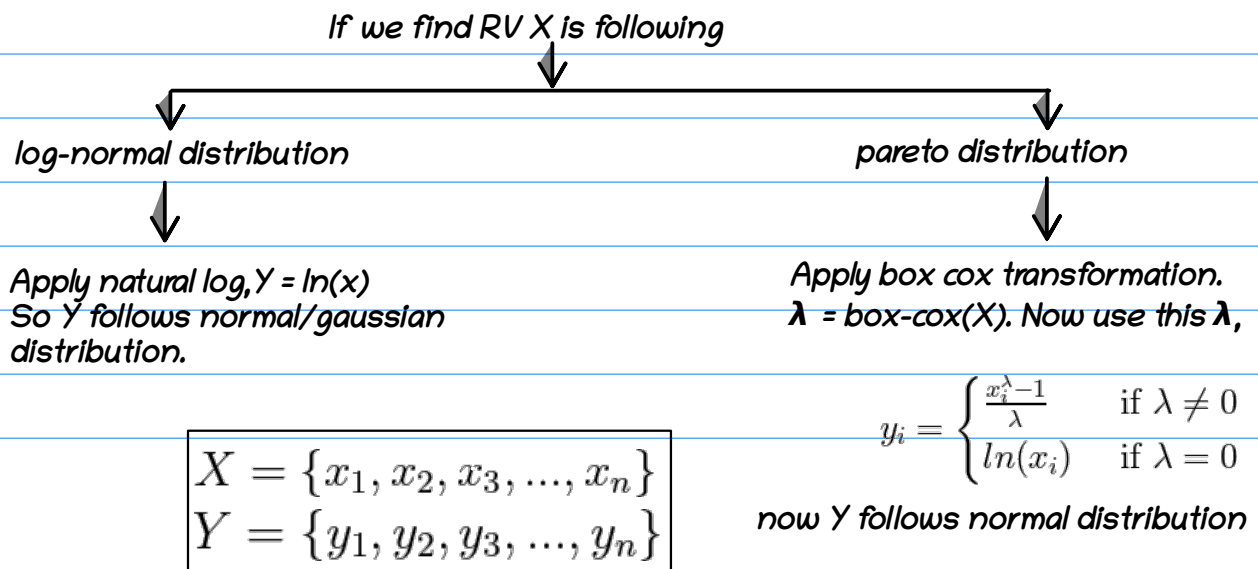
Income of 97-99% of the population. (top 2-3% follows pareto distribution)

$$\text{PDF}=f(x) = \frac{1}{x\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{\ln x - \mu}{\sigma}\right)^2} \quad \text{CDF}=F(x) = P(X \leq x) = \int_0^x f(x)dx$$



Practically if we find any distribution which is following log normal distribution

we take log of all data points so that new RV follows normal distribution.



8. Pareto distribution : $\text{pareto}(x_m, \alpha)$ infinite distribution

There are so many phenomenon whose logarithm follow normal distribution, which is known as Log Normal Distribution. That means if random variable X is log normally distributed than $\ln(X)$ follows normal/gaussian distribution.

If most of the part of RV (80%) is distributed in short interval and then rest 20% part's probability of happening decrease slowly over range is pareto distributed.

The sizes of human settlements (few cities, many villages)

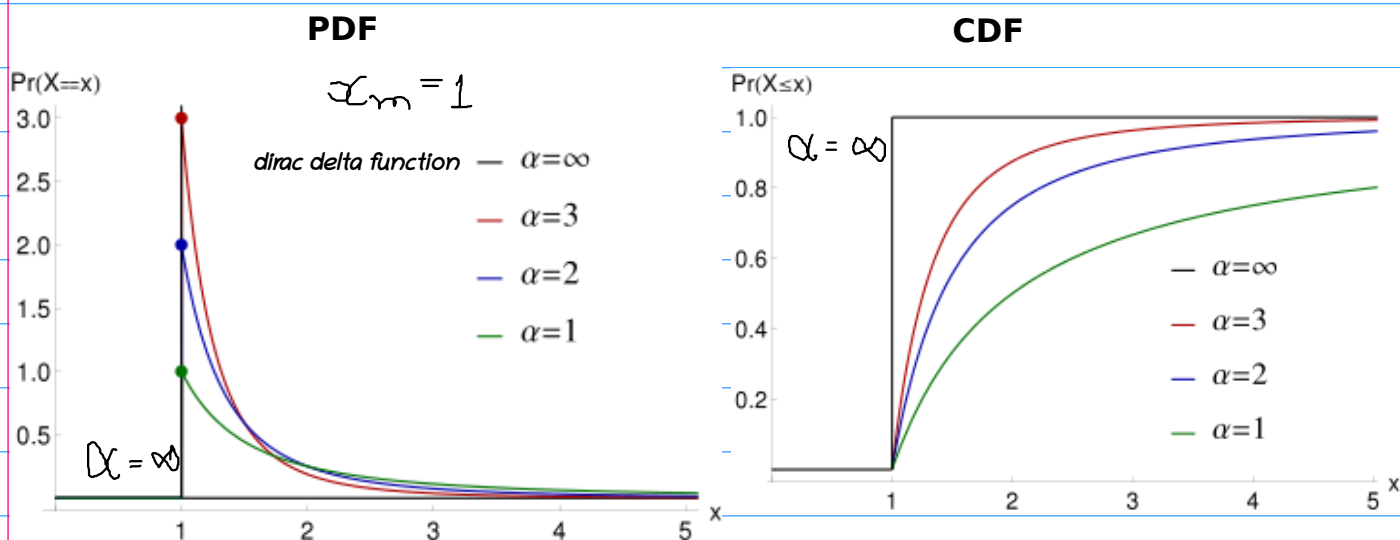
Oil reserves in oil fields (a few large fields, many small fields)

Size of meteorites.

parameters

$x_m > 0$, scale parameter - which decides the scale of distribution

$\alpha > 0$, shape parameter - which decided the tail of distribution (tail index)



Q: How to check if RV is following pareto distribution? What to do if it is Pareto?

A: Plot log of probability of number $p(x)$ versus log of that number, y axis = $\log(\text{prob}(x))$ and x axis = $\log(x)$. If it happens to be straight line than we can say that it is Pareto distribution(power law distribution).

We can also plot pareto-QQ plot also to check if RV is following pareto dist.

limitation: Pareto QQ will not work good for very small values of α is very small.

Log-log plot is necessary but insufficient evidence for a power law relationship, as many non power-law distributions will appear as straight lines.

If you find the distribution is pareto you can apply box-cox transformation to convert it into normal/gaussian distribution.