

Variable Transformation

Transformation of your variable through mathematical formulation

• Why do we require these mathematical formulation for transformation in Feature Engineering?

→ The distribution of our data through this gets converted into normal distribution.

The end goal is to normally distribute data.

Why normal distribution is required?

→ Problem solving around normal distribution becomes easy.

Linear / logistic regression assume that data is normal distributed, if not then we have to normalize it.

Few machine learning algo works better on normally distributed data, and these transformer help us normally distribute the data.

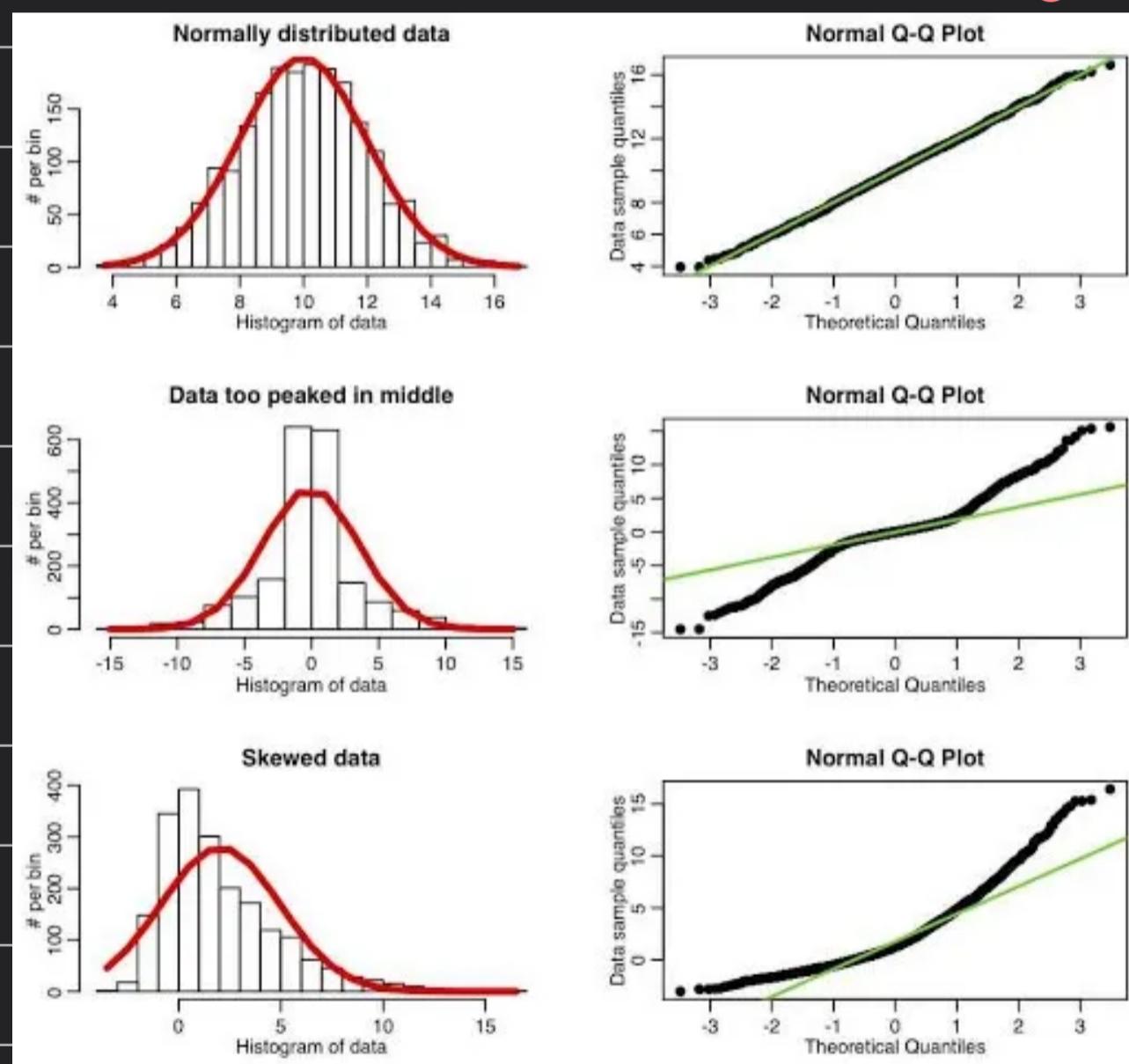
→ How to find if data is normally distributed?

There are few ways

1.) sns. distplot \rightarrow gives a distribution looking at which we can see if data is normally distributed

a.) pd. skew() \rightarrow pandas skew function
if o/p $\Rightarrow 0$ \rightarrow then Normally dist
o/p \Rightarrow -ve/+ve \rightarrow then data is skewed

b.) QQ Plots \rightarrow The best way & reliable way



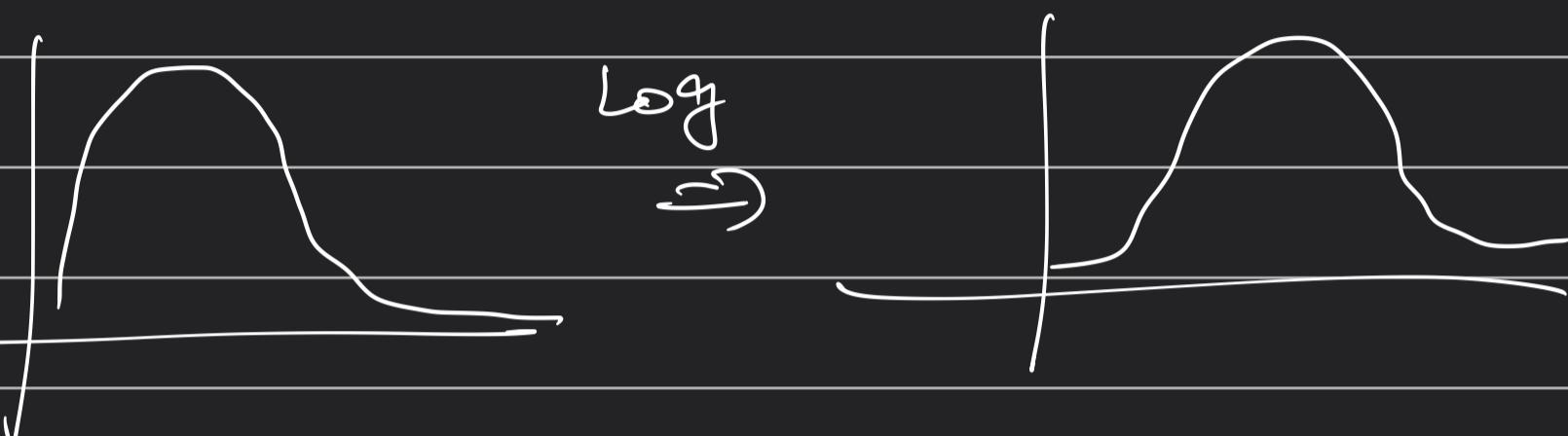
This is one way of reading QQ plots, other forms exist as well

▲ Log Transform

We take log of all the values present

in our variable. Base 2 or 10 depends, the data Normal distributed not completely.

we can't apply this on -ve data & can only be applied on right skewed data.



But why Log?

log bring values such as 1, b, 100, 1000 comes to an equivalent scale, so small values & very large values comes to an equivalent scale.

This is helpful in algo such as linear regression.

Additive scale gets converted to multiplicative scale.

★ Reciprocal Transform

$\frac{1}{X} \Rightarrow$ small value becomes big & big becomes small through this value of transformation variable

★ Sq Transform

Done for left skewed data, where you square all values in your variable

$$x'_1 \rightarrow x''_1$$

Box - Cox Transform

cannot be used on -ve values, or 0

$$(1) \quad x'_i = \begin{cases} x_i^\lambda - 1 & \text{if } \lambda \neq 0 \\ \ln(x_i) & \text{if } \lambda = 0 \end{cases}$$

The exponent λ is called λ (bhabda) that varies over range of -5 to 5 & in the process of examining all values to find appropriate or optimal value for λ . This help in best approximation for our distribution

(2)
 x ka power λ kitna shi zhega. we calculate first. we try all values in range of -5 to 5, 1.25, -1, 1, 2 etc.

we use two techniques for calculating

i) max likelihood

ii) Bayesian statistics

Yeo - Johnson

An adjustment of Box - Cox transform,

which can be applied to -ve numbers.

$$x_i^0 \underset{\text{def}}{=} \begin{cases} \left[(n_i^0 + 1)^{\lambda} - 1 \right] / \lambda & \text{if } \lambda \neq 0, n_i^0 \geq 0 \\ \ln(n_i^0 + 1) & \text{if } \lambda = 0, n_i^0 \geq 0 \\ -\left[(-n_i^0 + 1)^{2-\lambda} - 1 \right] / (2-\lambda) & \text{if } \lambda \neq 2, n_i^0 < 0 \\ -\ln(-n_i^0 + 1) & \text{if } \lambda = 2, n_i^0 < 0 \end{cases}$$

