



# Feature Engineering

Ayush

- 1) FE
- 2) Skewness, Kurtosis
- 3) Univariate Analysis
  - ↳ Converting Cont. data to discrete data
- 4) Creating new feature

⇒ Feature Engineering →

Q What is feature ?

Variable, "Attributes", "dimension" "Characteristics"

Aerofit  
case  
study

Rows

Columns →

	Education	Gender	Income	Skills	Product
1 <sup>st</sup> Row					

1 person  
and basically  
their data pts (characteristics).

Target  
↓

Column which  
you want  
to predict or  
analyze  
or generate insight  
about

Rows → Records

Columns → features in Data Science  
fields

features are used to predict Target.

Q I want to know whether a person is fit or not?

① Understanding fitness

⇒ Small Survey → weight & Height.

100 people	Weight	Height	Fit	Domain Expert
			0	are generally asked to label the data
			1	
			0	
			1	
			0	
			1	

Q Is only weight enough?? partially

Q Is only height enough??

⇒ Combination of feature → can be a good indicator.  
Some single feature → can still be a good indicator

Some single feature → might just be useless

$\Rightarrow \underline{\text{BMI}} \Rightarrow \text{Body Mass Index}$

$$\Rightarrow \frac{w}{h^2}$$

New feature are created to simplify improve the prediction

$w$	$h$	BMI	Fit
x	x	x	
x	x	x	
x	x	x	
x	x	x	

$\Rightarrow$  Domain knowledge  $\rightarrow$  Required to create certain new features.

$\Rightarrow$  Can our Machine Create relevant feature itself.  
"Machine learning".

$\Rightarrow$  Loan Worthiness Chart  $\rightarrow$

$\Rightarrow$  Gender, Income, Married, ...  $\rightarrow$  feature  
Loan Status  $\rightarrow$  Approved or Not  $\Rightarrow$  Target

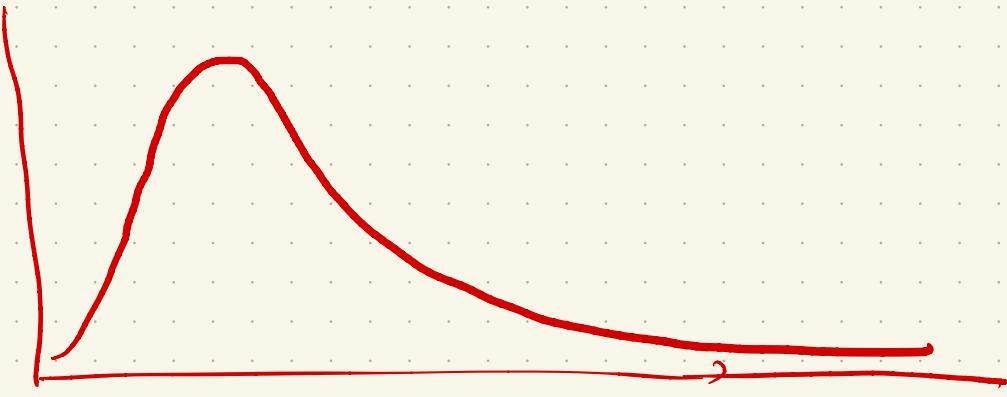
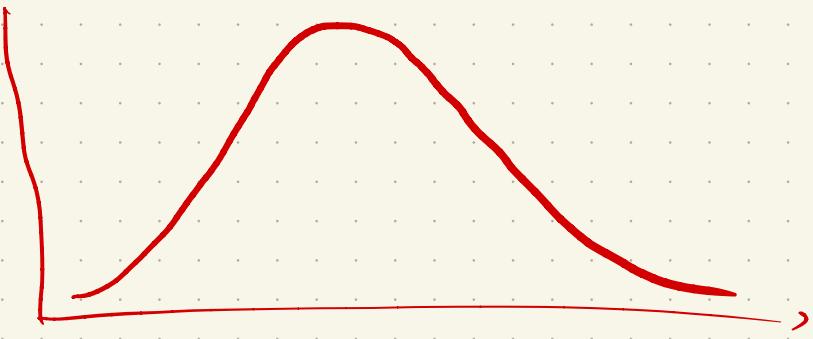
Variable Name	Description	Sample Data
Loan_ID	Loan reference number (unique ID)	LP001002; LP001003; ...
Gender	Applicant gender (Male or Female)	Male; Female
Married	Applicant marital status (Married or not married)	Married; Not Married
Dependents	Number of family members	0; 1; 2; 3+
Education	Applicant education/qualification (graduate or not graduate)	Graduate; Under Graduate
Self_Employed	Applicant employment status (yes for self-employed, no for employed/others)	Yes; No
ApplicantIncome	Applicant's monthly salary/income	5849; 4583; ...
CoapplicantIncome	Additional applicant's monthly salary/income	1508; 2358; ...
LoanAmount	Loan amount	128; 66; ...
Loan_Amount_Term	The loan's repayment period (in days)	360; 120; ...
Credit_History	Records of previous credit history (0: bad credit history, 1: good credit history)	0; 1
Property_Area	The location of property (Rural/Semiurban/Urban)	Rural; Semiurban; Urban
Loan_Status	Status of loan (Y: accepted, N: not accepted)	Y; N

→ Skewed → log normal.

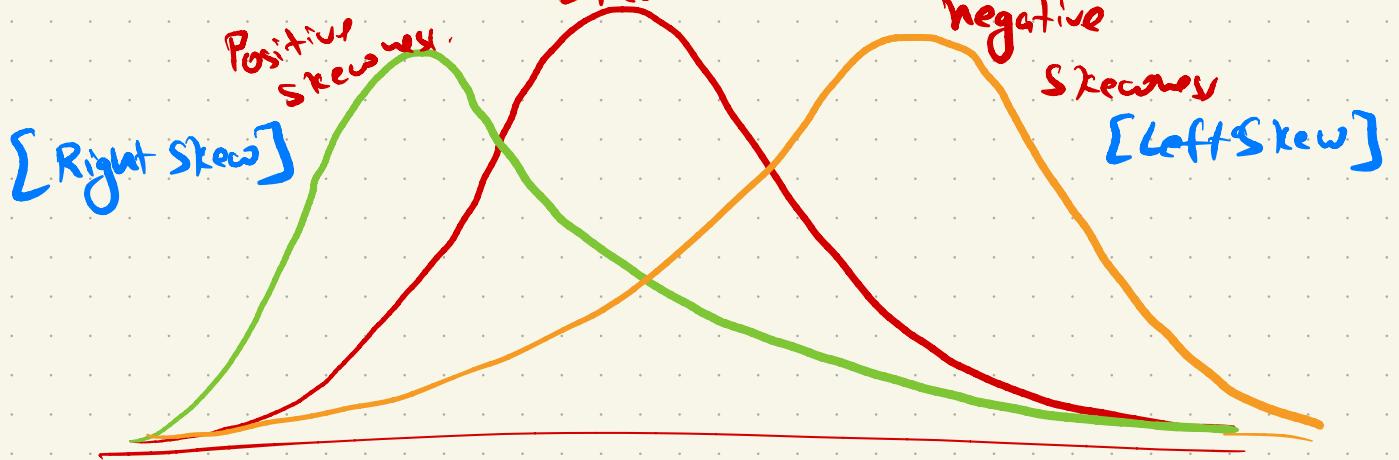
Is the data Skewed ??

Why this is important !

- ① Understand the distribution →
- ② Skewness impact loan approval (target)
- ③ Outliers understanding.

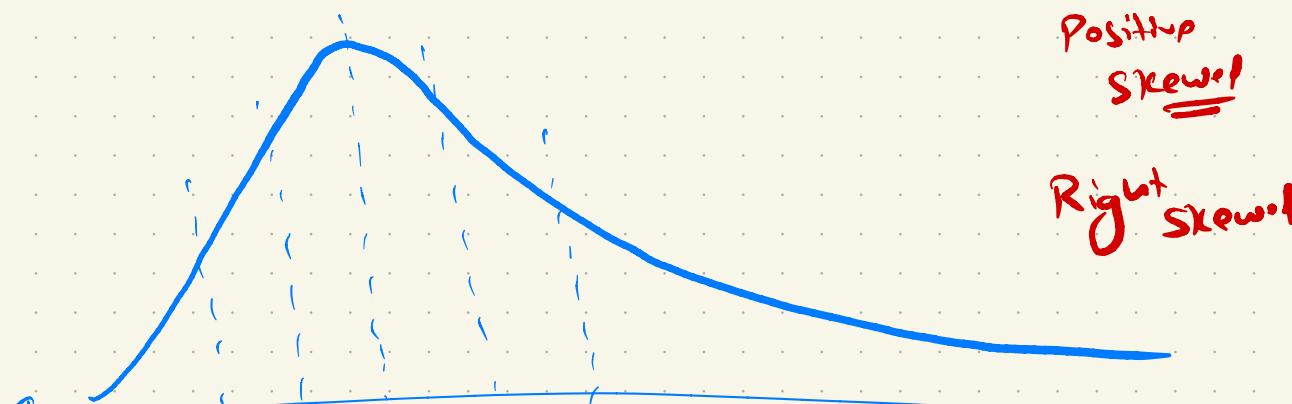


Type of Skewness



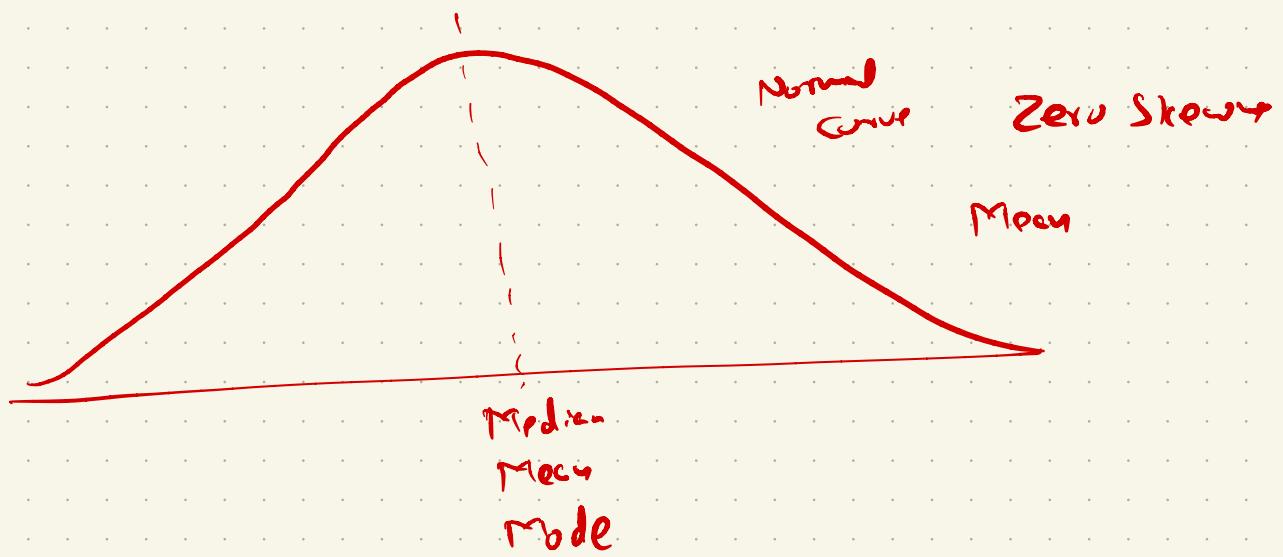
There are multiple maths formulae for skewness →

$$\text{Skew} = \frac{3(\text{mean} - \text{med})}{\sigma}$$



Mode  
Median  
Mean

Right Skew  $\rightarrow$  Mean is on right of Median.



Median  
Mean  
Mode



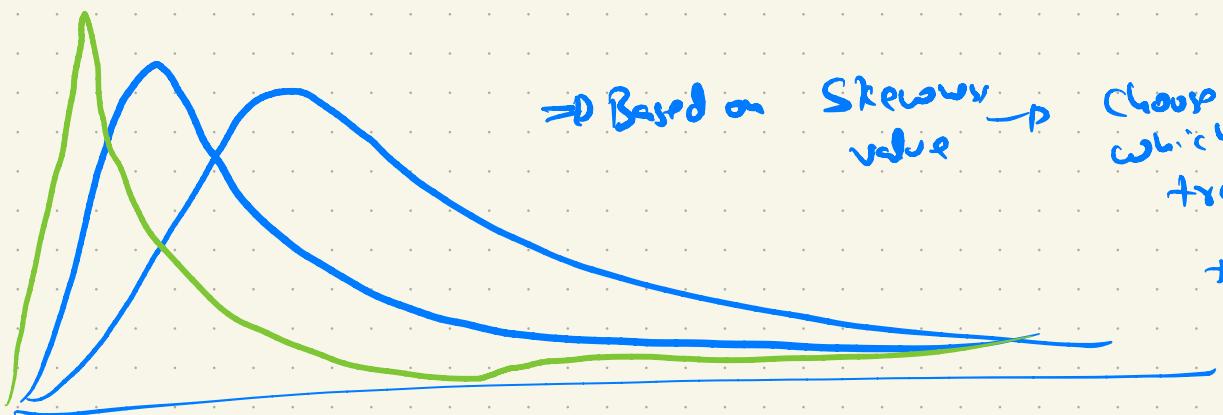
Mean Median Mode

Skewness  $\rightarrow$  tells us the shape of the distribution

$\mu \rightarrow$  mean

$\sigma \rightarrow$  Spread of dist'

Skew  $\rightarrow$  Shape of distribution



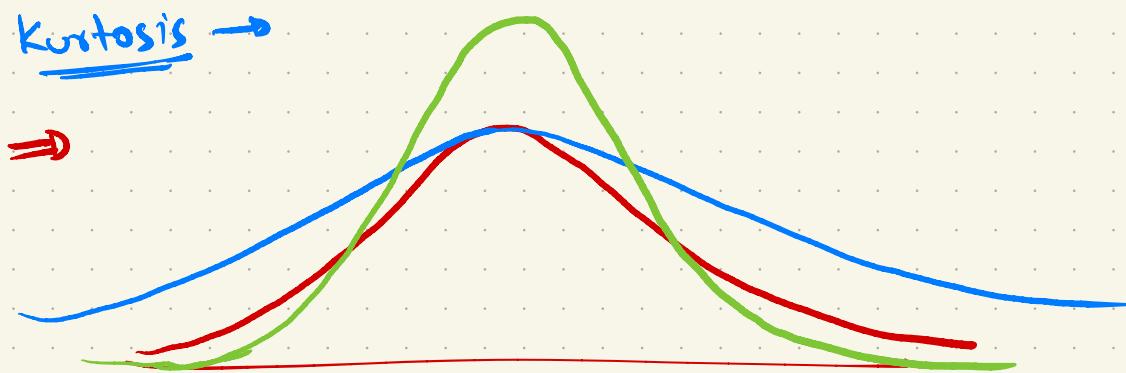
- Skewness  $\rightarrow$
- ① Data preparation  $\rightarrow$  transformation
  - ② Decision making can be improved
  - ③ Understand data & patterns

$\mu$   $\rightarrow$  mean  $\rightarrow$  central tendency

$\sigma \rightarrow$  Spread

Skew  $\rightarrow$  Shape of data.





$\Rightarrow \sigma \rightarrow$  Spread of data is.

$\Rightarrow$  Kurtosis  $\rightarrow$  Tells us how peaked the data is.

- ↳ Helps understanding data
- Outliers detection

$\Rightarrow$  Pointy  $\rightarrow$  More concentrated

Normal  $\rightarrow$

Flat  $\rightarrow$

- $\rightarrow$  Lepto kurtic
- $\rightarrow$  Meso kurtic
- $\rightarrow$  Platy kurtic

$\Rightarrow$  Skewness + kurtosis  $\rightarrow$  Interview Topics.

1

Still used  
if it is very extreme.

⇒ Univariate Analysis →

→ Analyse impact of 1 variable against target.

Total income → Loan

A = Salary 20 L

Loan 50 L →

Loan Ten 1 yr

B = Salary 10 L      Loan 50 L →

Term 20 yrs

⇒ Applicant should be able to pay EMI

⇒ EMI → Total loan Amount

loan Year × 12

- $\Rightarrow$  Domain knowledge  $\rightarrow$  income impact loan status  
 $\rightarrow$  we checked income  $\rightarrow$  didn't see result X  
 $\rightarrow$  we convert it to total income  $\rightarrow$  X  
 $\rightarrow$  Is income enough to pay the loan  $\rightarrow$  EMI, able to pay  
 $\rightarrow$  able to pay to loan status  $\rightarrow$  ✓

You are working on a dataset for predicting house prices. You have a feature 'YearBuilt' representing the year a house was built.

It does not contain any missing values. To improve the model's performance, what feature engineering technique can you apply to 'YearBuilt'?

- 
- $\Rightarrow$  • Feature Engineering  
 • Skewness + kurtosis  $\rightarrow$  to undervited data  
 • Univariate analysis  
 ↳ • Charts  
 ↳ • T-test to validate  
 ↳ • Binning  
 ↳ • Create new feature



- FE is an art →
- The more you practice → the better your model becomes

Collab Link : [https://colab.research.google.com/drive/1kLzjBZWRJS\\_drg98Piv8RZ4Hvd7JBkBv?usp=sharing](https://colab.research.google.com/drive/1kLzjBZWRJS_drg98Piv8RZ4Hvd7JBkBv?usp=sharing)

This blog has pretty simple but good starter problems and their solutions

<https://kindsonthegenius.com/blog/hypothesis-testing-solved-examplesquestions-and-solutions/>

You can also refer to Ch- 9 -11 of this book

[https://assets.openstax.org/oscms-prodcms/media/documents/IntroductoryStatistics-OP\\_i6tAI7e.pdf?](https://assets.openstax.org/oscms-prodcms/media/documents/IntroductoryStatistics-OP_i6tAI7e.pdf?gl=1*1nfheaw*ga*MTc0ODY0ODkwNy4xNzAyMzE2ODk4*ga_T746F8B0QC*MTcwMjMxNjg5OC4xLjEuMTcwMjMxNzM4Ni41NS4wLjA.)

[gl=1\\*1nfheaw\\*ga\\*MTc0ODY0ODkwNy4xNzAyMzE2ODk4\\*ga\\_T746F8B0QC\\*MTcwMjMxNjg5OC4xLjEuMTcwMjMxNzM4Ni41NS4wLjA.](https://assets.openstax.org/oscms-prodcms/media/documents/IntroductoryStatistics-OP_i6tAI7e.pdf?gl=1*1nfheaw*ga*MTc0ODY0ODkwNy4xNzAyMzE2ODk4*ga_T746F8B0QC*MTcwMjMxNjg5OC4xLjEuMTcwMjMxNzM4Ni41NS4wLjA.)