

Feature Engineering

①

Feature Engineering is the process of using domain knowledge to extract features from raw data. These features can be used to improve the performance of machine learning algorithms.

What is Feature Scaling?

Feature scaling is a technique to standardize the independent features present in the data in a fixed range.

age	salary	purchase
70	10000	
40	36000	
30	40000	

10 scale ← (age)
↳ 1000 scale (salary)
different scale

Why do we need Feature scaling?

age	Salary	Purchase
50	83600	1
27	48000	0

if we use KNN and find the distance between

$$\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2} \rightarrow \text{age} \rightarrow 529$$

↳ salary → 1225000000

↳ dominant

Types of Feature Scaling

Feature Scaling

1. Standardization
2. Normalization
 - a) Min-Max
 - b) Robust scaler

Standardization - Intuition

Also called as Z-score Normalization

<u>Age</u>	<u>salary</u>
27	-
15	-
33	-
63	-
90	-
05	-
⋮	⋮
500 value	-

$$X_i' = \frac{X_i - \bar{X}}{\sigma}$$

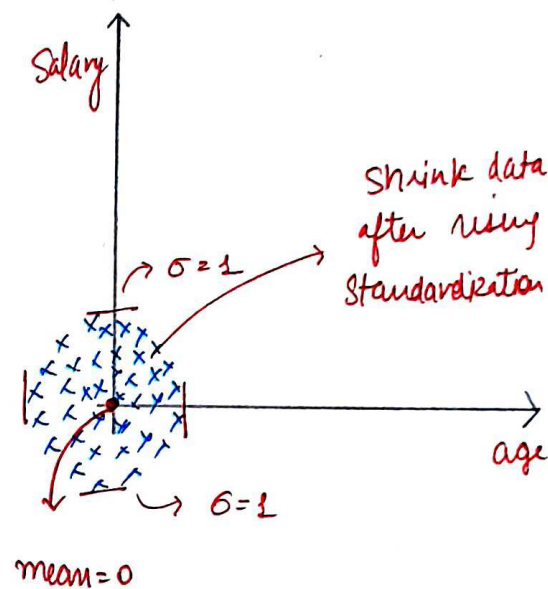
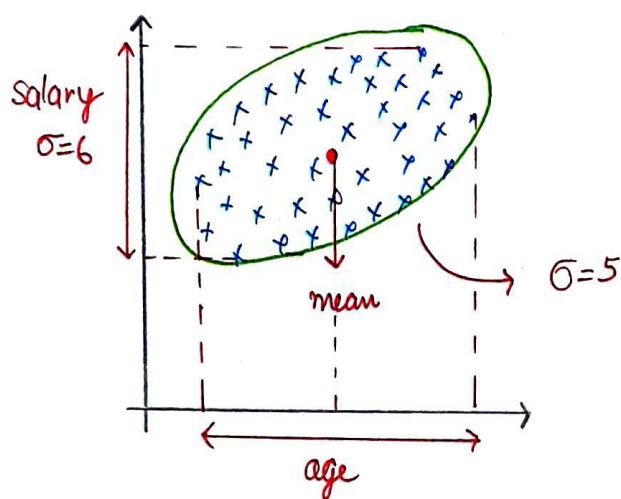
mean

Standard Deviation

* After scale the column and new scaled column

$$\mu = 0, \sigma = 1$$

Always (scaled colⁿ)



Impact of Outlier

There is not any impact of outlier on data.

When to use standardization?

1. k-Mean
2. k-Nearest Neighbours
3. PCA
4. Artificial Neural Network
5. Gradient Descent

Normalization

Normalization is a technique often applied as part of data preparation for machine learning. The goal of normalization is to change the values of numeric columns in the datasets to use a common scale, without distorting differences in the ranges of value or losing information.

1. Min Max scaling
2. Mean Normalization
3. Max
4. Robust scaling

MinMax Scaling

Weight (in kg)

130

67

81

61

32

.

.

.

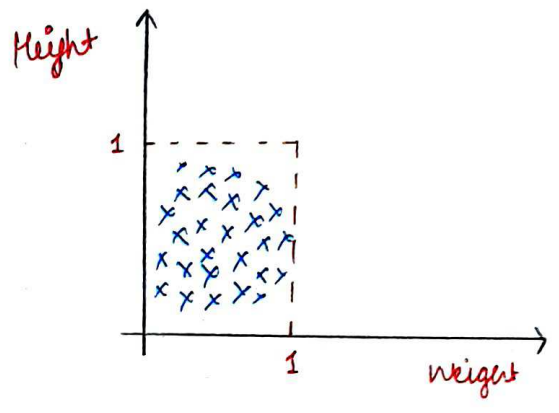
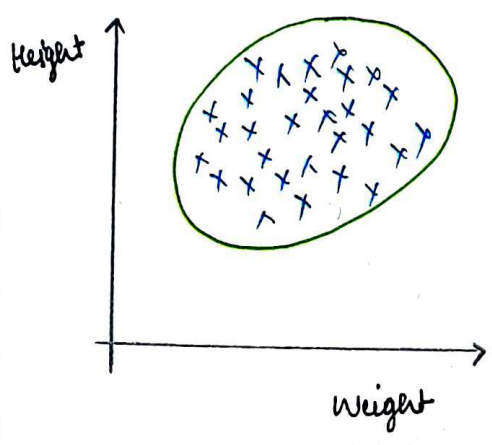
.

Normalize
→ Min Max Scaling →

$$X'_i = \frac{X_i - X_{\min}}{X_{\max} - X_{\min}}$$

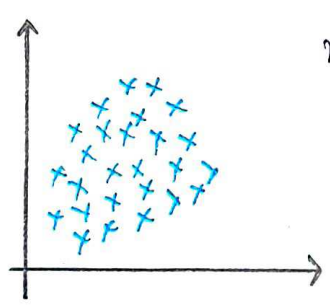
$$= \frac{130 - 32}{130 - 32} = 1$$

Height and weight

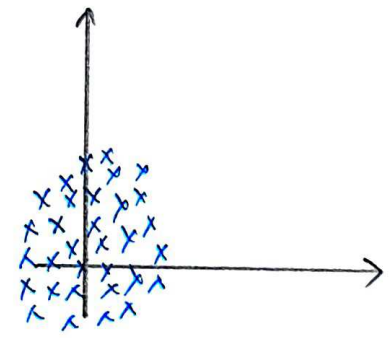


Mean Normalization

$$X'_i = \frac{X_i - X_{\text{mean}}}{X_{\text{max}} - X_{\text{min}}} \rightarrow [-1 \text{ to } 1]$$



mean centered



Max Absolute Scaling

Weight
200
100
300

$$X'_i = \frac{X_i}{|X_{\text{max}}|}$$

scikit learn
↑
Max Abs scale

use when data has many 0.

Robust Scaling

weights

200

300

100

$$X_i' = \frac{X_i - X_{\text{median}}}{\text{IQR} (75^{\text{th}} \text{ per} - 25^{\text{th}} \text{ per})}$$

* If Data have lots of outlier. (Try)

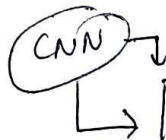
Normalization Vs standardization

1. > Is feature scaling required

2. > Min Max



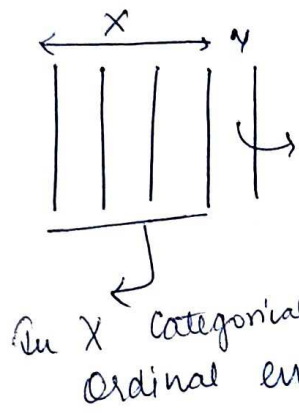
image →



→ [0 - 225]

→ min max

Encoding categorical variables



Ordinal Encoding

Education

- HS $\rightarrow 0$
- UG $\rightarrow 1$
- PG $\rightarrow 2$
- PG $\rightarrow 2$
- UG $\rightarrow 1$
- HS $\rightarrow 0$
- UG $\rightarrow 2$

PG > UG > HS

* There is a order in data so we use ordinal encoding otherwise one hot encoding can be use.

$\left\{ \begin{array}{l} PG \rightarrow 2 \\ UG \rightarrow 1 \\ HS \rightarrow 0 \end{array} \right\} \rightarrow$ we have define this order.

One Hot Encoding

\rightarrow dummy variable trap

color
 Yellow
 blue
 Red

<u>color-Y</u>	<u>color-B</u>	<u>color-R</u>
1	0	0
0	1	0
0	0	1

$[M-1 \text{ cols}]$
 \rightarrow accept

X

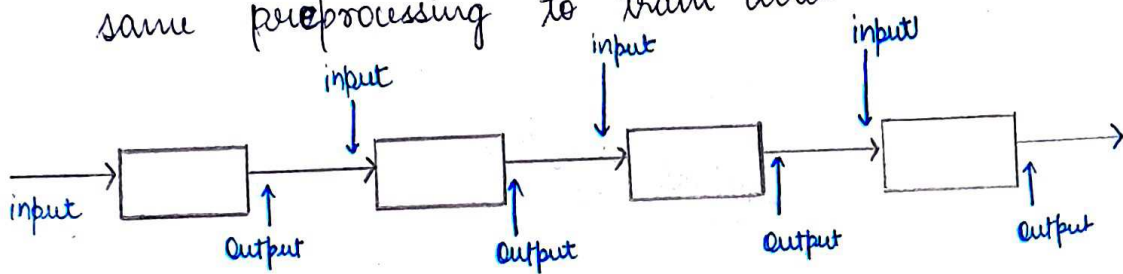
Y \rightarrow	1	0	0
B \rightarrow	0	1	0
R \rightarrow	0	0	1

Multicollinearity

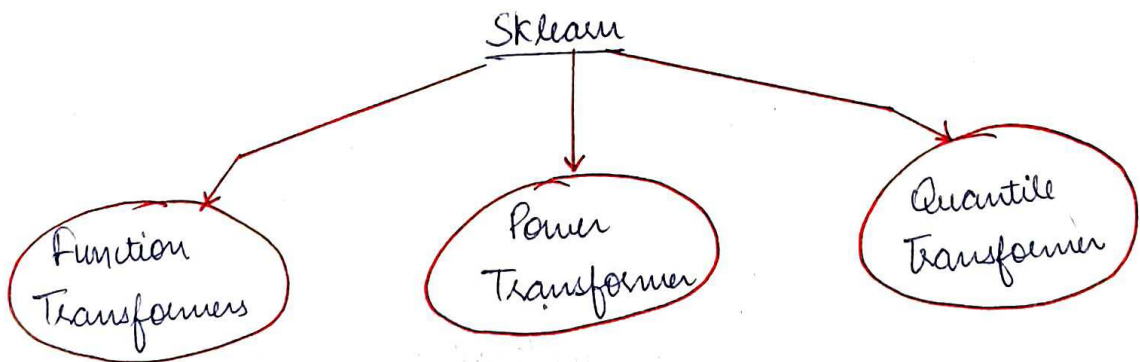
Scikit learn Pipelines

Pipelines chains together multiple steps so that the output of each step is used as input to the next step.

Pipelines makes it easy to apply the same preprocessing to train and test.



Function Transformer



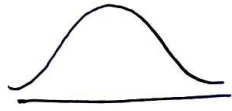
- | | |
|-------------------------|----------------|
| 1. Log Transformer | 1. Box Cox |
| 2. Reciprocal | 2. Yeo-Johnson |
| 3. Square Root / Square | |
| 4. Custom | |

How to find if data is normal?

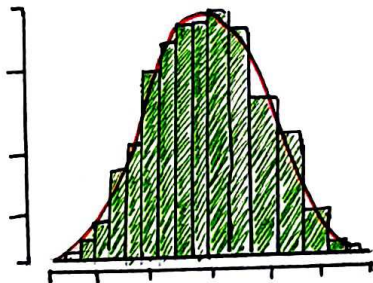
1.) sns.distplot

2.) pd.skew() = 0

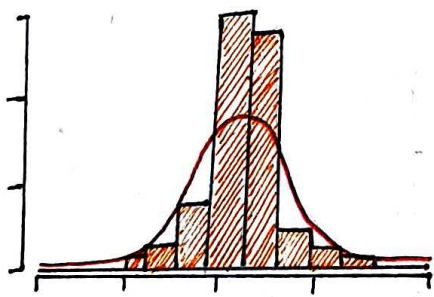
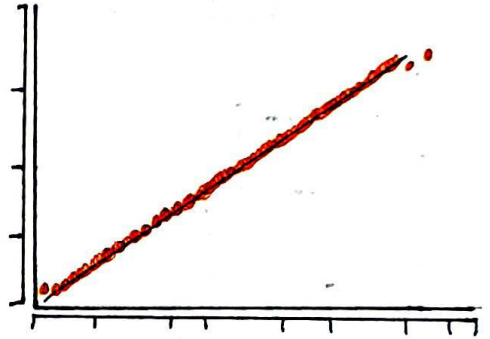
3.) Q-Q plot



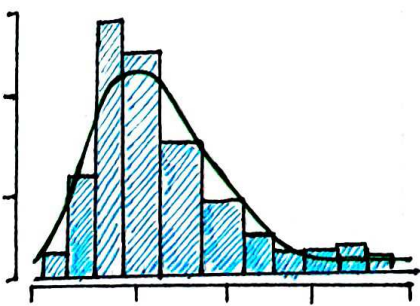
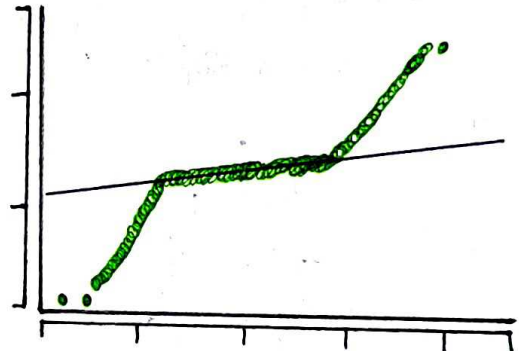
QQ plots



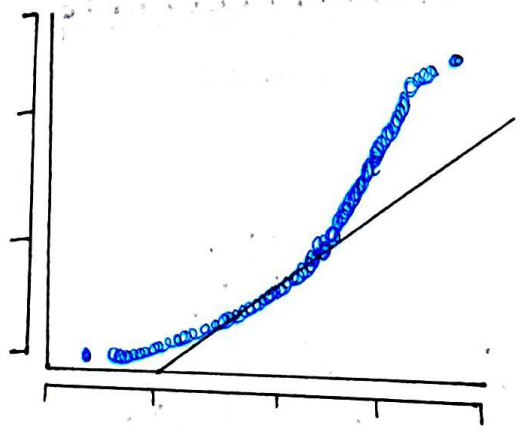
Normally distributed data

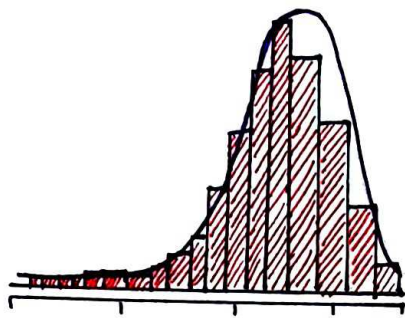


~~Skewed~~ Data
Data too peaked in middle

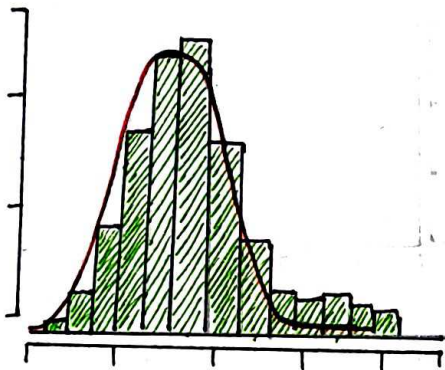
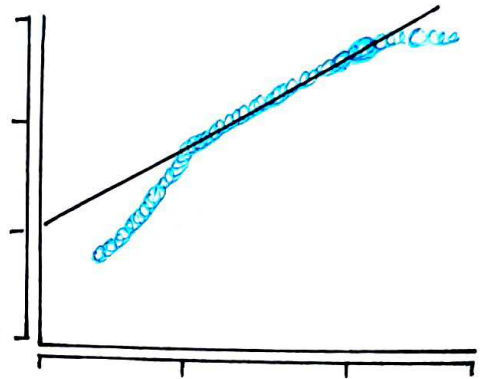


Skewed Data

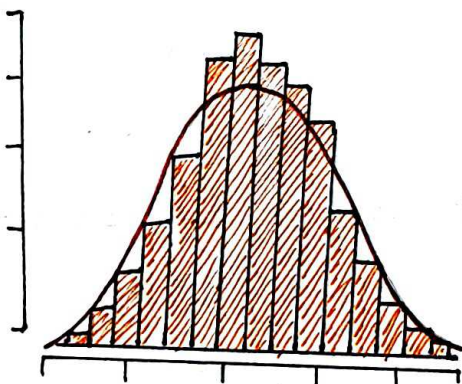
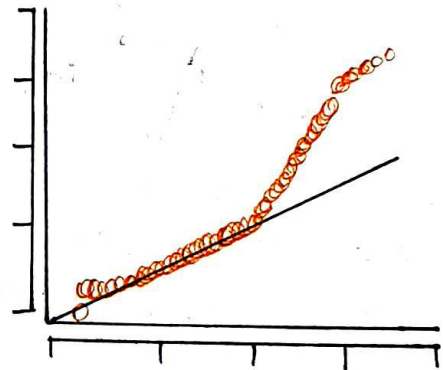




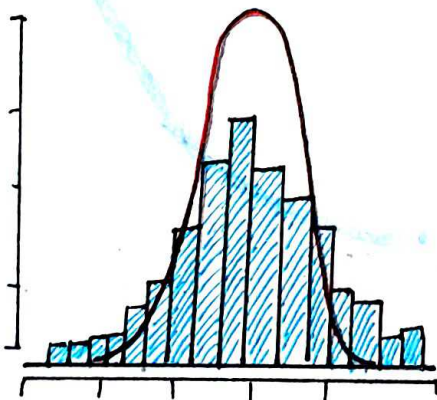
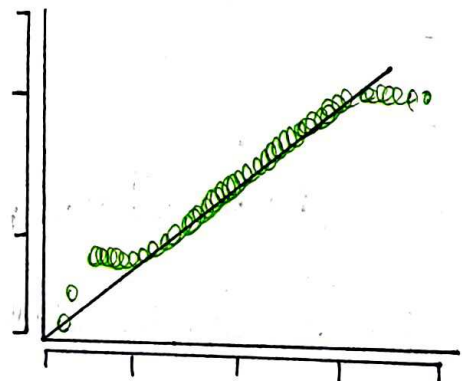
Skewed left



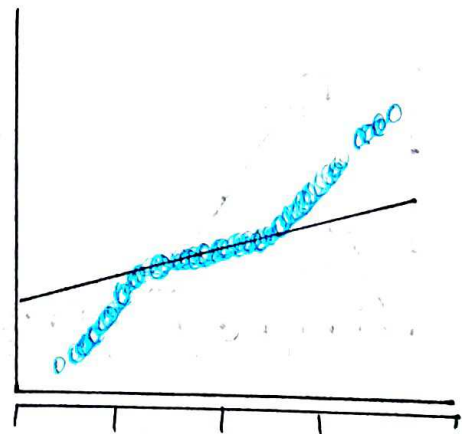
Skewed Right



Thin tails



Fat Tails



Log Transform

6

Age

21

→ (log)

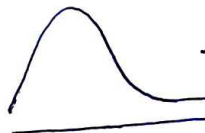
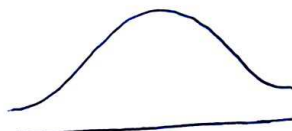
35

45

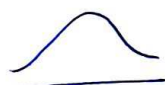
⋮

`np.log1p`

↳ * If data have 0 then always add 1 then perform log.



→



* log → not (-ve data)

* Right skewed $\xrightarrow{\text{log}}$ Center

Other Transform

Reciprocal $1/x$

Sq (x^2)

Sqrt (\sqrt{x})

↓

left skewed

Power Transform

Box Cox Transform

The exponent here is a variable called lambda (λ) that varies over the range of -5 to 5, and in the process of searching, we examine all values of λ . Finally we choose the optimal value (resulting in the best approximation to a normal distribution) for your variable.

$$x_i^{(\lambda)} = \begin{cases} \frac{x_i^\lambda - 1}{\lambda} & \text{if } \lambda \neq 0, \\ \ln(x_i) & \text{if } \lambda = 0, \end{cases}$$

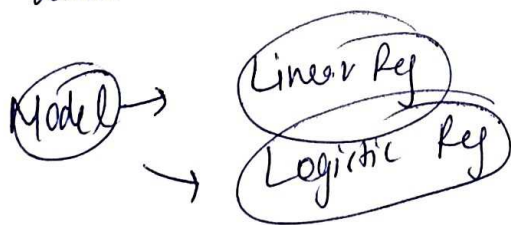
* Box Cox \rightarrow only applicable for $n > 0$

Yeo - Johnson Transform

This transformation is somewhat of an adjustment to the Box-Cox transformation, by which we can apply it to negative numbers.

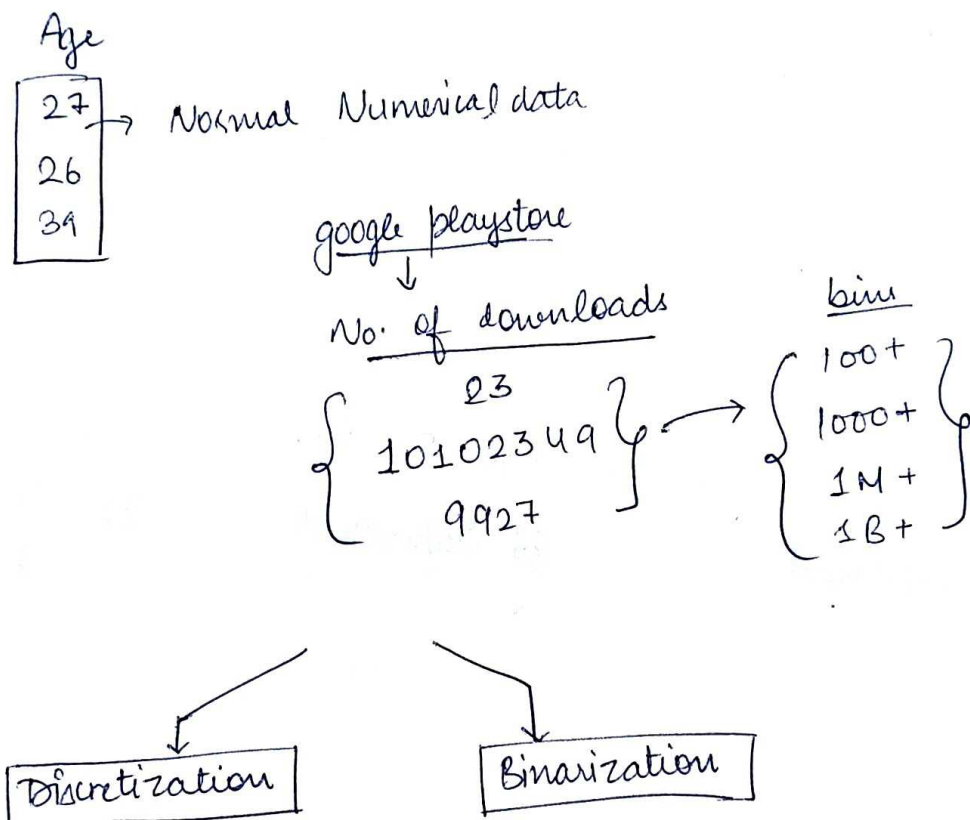
$$x_i^{(\lambda)} = \begin{cases} [(x_i + 1)^\lambda - 1] / \lambda & \text{if } \lambda \neq 0, x_i \geq 0; \\ \ln(x_i + 1) & \text{if } \lambda = 0, x_i \geq 0; \\ -[(-x_i + 1)^{2-\lambda} - 1] / (2-\lambda) & \text{if } \lambda \neq 2, x_i < 0; \\ -\ln(-x_i + 1) & \text{if } \lambda = 2, x_i < 0 \end{cases}$$

* This transform can work with -ve and 0 and positive data.



Encoding Numerical Features

7



Discretization

Discretization is the process of transforming continuous variables into discrete variables by creating a set of contiguous intervals that span the range of the variable's values. Discretization is also called binning, where bin is an alternative name for interval.

Why use Discretization:

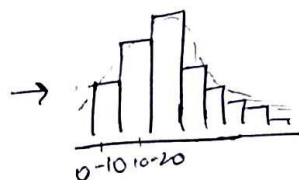
1. To handle outliers
2. To improve the value spread

Age

23, 42, 57, 81, - - - - - / 100

↓

0-10, 10-20, 20-30
 T_{bins} $T_{(6)}$ $T_{(10)}$
 check how many
 value betⁿ 0-10
 (5)



Types of Discretization

1.) Unsupervised

- a) Equal width
(uniform)

- b) Equal frequency (quantile)

- c) K-means
binning

2.1 Supervised

- a) Decision Tree (Binning)

3.] Custom Binning

Equal width / uniform Binning

Age

Age
27, 32, 34, 56, ... - max 100

$Bins = 10$

$$\frac{\text{max} - \text{min}}{\text{bins}} = \frac{100 - 0}{10} = 10$$

bins
 $(0-10)$, $(10-20)$, $(20-30)$ $(90-100)$
175

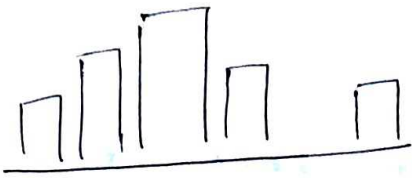
frequency,

5

16

17

5



→ equal binning

Advantage

- 1.) Outlier handle
- 2.) No change spread

age	age-trf	age-labels
43.0	5.0	(40.21, 48.16)
44.0	5.0	(40.21, 48.16)
15.0	1.0	(8.3, 16.3)
30.0	3.0	(24.9, 32.2)
35.0	4.0	(32.2, 40.21)
35.0	0.0	(0.4, 8.3)
2.0	2.0	(16.3, 24.9)
18.0	0.0	(0.4, 8.3)

Equal Frequency / Quantile Binning

Intervals = 10

Each interval contains 10% of total observations

Intervals:

0-16

10%
data

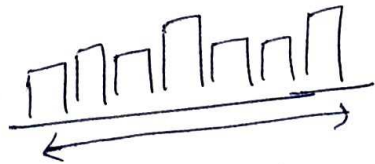
16-20

20%
data

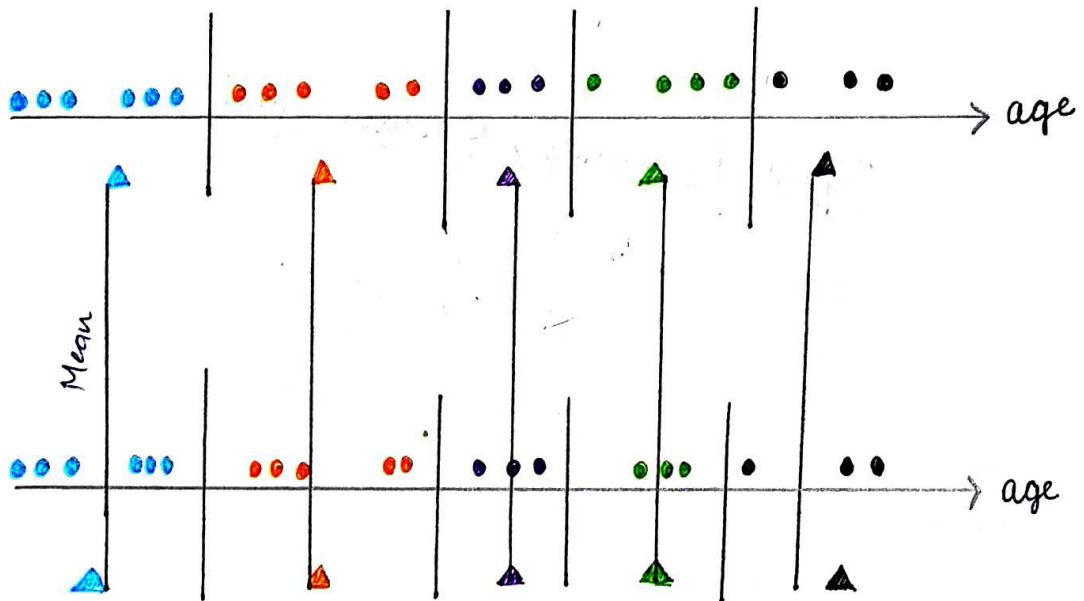
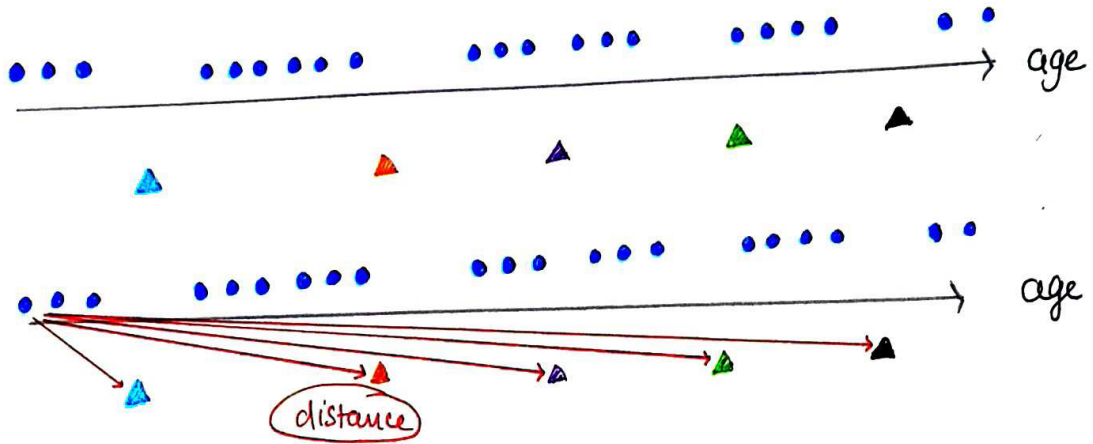
20-22

30%
data

- Handle outlier
- value spread uniform



KMeans Binning



Encoding the discretized variable

Sk learn



K Bins Discretized ()

bins = ?

Strategy

- uniform
- quantile
- kmeans

encoding

- ordinal
- OHE

Custom / Domain Based Binning

9

{ [0-18] → kids
[18-60] → adult
[60-80] → old } → pandas

Mixed Data

Mixed

Cabin Data

{ B5
D41
B6 } →

Categorical	Numerical
B	5
D	41
B	6

Data

7
3
1
A
C
4
D

Num
7
3
1
NA
NA
4
NA

Data (Categorical)
NA
NA
NA
A
C
NA
D