**"1"** → Feature Transformation → Only For **Continous Data**

Feature processing : Transformation.
↓
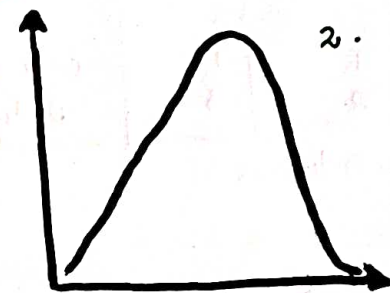columns / variables / Features (same)

EX: | x |
|---|
| 2 |
| 4 |
| 9 |
| 12 |
| 25 |

$x^{-0.5}$ √
$x^{1/2}$
$x^{1/3}$
$x^{1/4}$

1. $[x]^{1/2}, [x]^{1/3}, [x]^{1/4}$
2. np. log

**Right Skewed**

$n^{th}$ root (or) $\overline{\log(x)}$

+1.5
+2
+6.9 .....

−1 to +1 (normal) Distribution

1. $[x]^{2}, [x]^{3}, [x]^{4}$
2. np. exp

**Left Skewed**

$n^{th}$ power (or) exp

−1 to 1

# import Liberaris / packages.
# import numpy as np
# Import pandas as pd
# Import matplotlib. Pyplot as plt
# % matplotlib inline
# Import Seaborn as Sns

1. Root Transformations:
2. log Transformations
3. Reciprocal Transformat
4. Boxcox Transformatio

# titanic = pd. read_csv ( "titanic_csv")
<u>Location of csv</u>

titanic

| | Passenger Id | Survied | Pclass | Name | Sex | Age | Sibsp | Parch | Ticket | Fare |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 3 | Owen Harris | male | 22.0 | 1 | 0 | A/52771 | 7.2500 |
| 1 | 2 | 1 | 1 | cumrig | Femal | 38.0 | 1 | 0 | . | 71.2833 |
| 2 | 3 | 1 | 3 | Laina | Femal | 26.0 | 0 | 0 | . | 7.9250 |
| 3 | 4 | 1 | 1 | Futrelle | Femal | 35.0 | 1 | 0 | . | 53.1000 |
| : | : | | | | | | | | | 30,0000 |
| 88 | 890 | 1 | 1 | Behr | male | 26.0 | 0 | 0 | | 7.7500 |
| 890 | 891 | 0 | 3 | Dooley | male | 32.0 | 0 | 0 | | |

891 rows X 12 columns.

| Cabin | Embarked |
|---|---|
| NAN | S |
| C85 | C |
| NAN | S |
| C123 | S |
| : | C |
| 048 | |

# titanic = pd. read_csv ("titanic .csv", Usecols = <u>["Fare"]</u> )

titanic_head()

Out

| | Fare |
|---|---|
| 0 | 7.2500 |
| 1 | 71.2833 |
| 2 | 7.9250 |
| 3 | 53.1000 |
| 4 | 8.0500 |

Particular column

Usecols = [ "Fare"]
<u>col. Name.</u>

# Sns. distplot ( titanic ["Fare"])

plt. show()

titanic ["Fare"]. Skew()

Out:   + 4.7873



low
class
people
Having more
*Comparing
to
Highclass
People

# another way

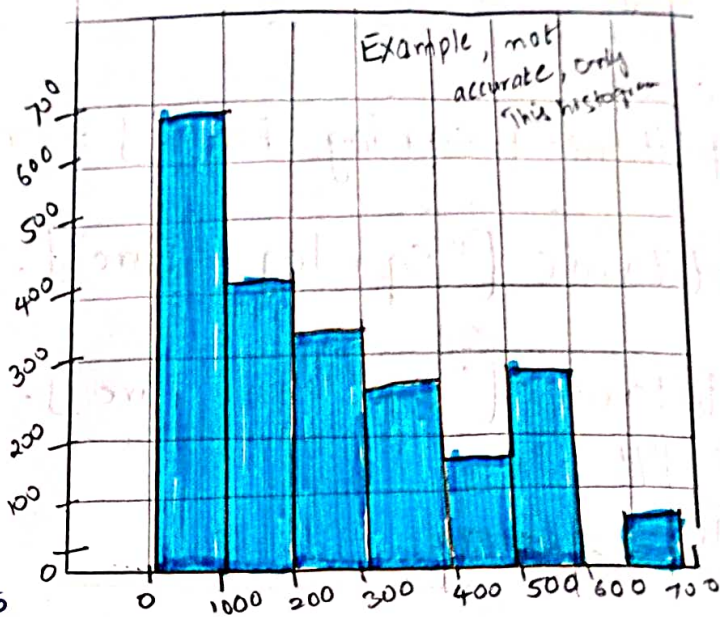    t - ["Fare"].hist()
    plt.show()

    titanic ["Fare"].skew()

    Out :- +4.7873

**Values**

Skew = +2.085 [ 2th root ]
Skew = +0.5    ( 4th root)
Skew = +0.21   (5th root)
Skew = -0.95   (6th root)

Consider,
This
which is
close To
"0"

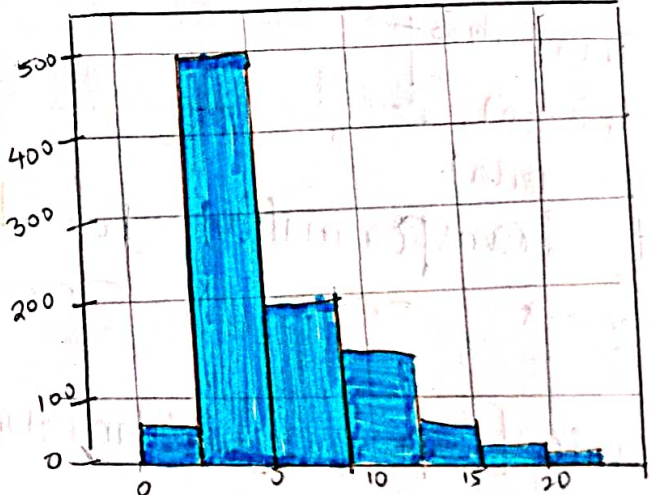* Root - Transformation # For Right Skewed

    # titanic ["sqr -Fare"] = titanic["Fare"] ** (1/2)
                                              → For Right Skewed
    # titanic ["sqr-Fare"].skew()              We use $n^{th}$ root
    # titanic ["sqr-Fare"].hist()                 √ , ∛, ∜ ∛
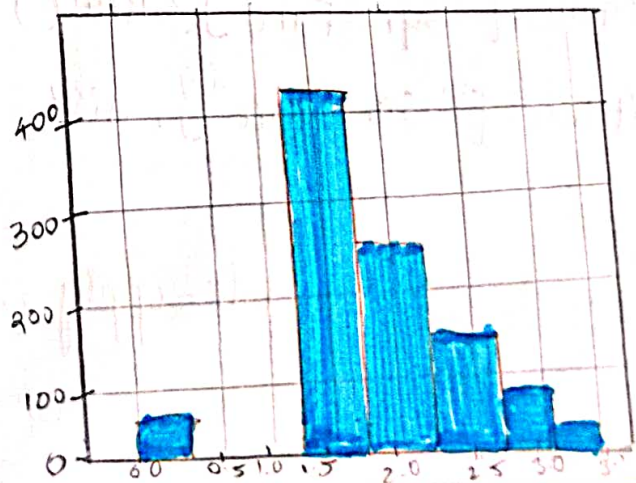
Out :- +2.085   **More Than +1
                or -1**

* Try with [ 1/5], [1/6].

    # titanic ["Sqr-Fare"] = titanic
                   ["Fare"] ** (1/5)

    # titanic ["Sqr-Fare"].hist()

    # titanic ["Sqr-Fare"].skew()

                        **less Than
                        -1 to +1**

    Out :- -0.212676
          **Final**

→ **log Transformation**

\#titanic ["sqr log_"Fare"] = **np.log** (titanic ["Fare"] **+0.01**)

\#titanic [" sqr log_" Fare"] .skew()

\#titanic ["sqr log_ "Fare"]. hist()
\# plt. show()

why?

Because In Data "Fare" column, we have some records "o" values. So,

Out **-2.41004**

Here, The Value, in This Have more, so we don't use log Transform For This col( Data)

If we calculate with log. it give α (infinite) value, so, we add +0.01 to "o" column.



we use " n " power

**Root** Transformation For (**Left Skewed**) $n^{th}$ power→ $(x)^2, (x)^3, (x)^4$ → left skewed

Fare data is not Left skewed, But To understand, Left skewed

$-n^{th}$ root $(x)^{1/2}, (x)^{1/3}, (x)^{1/4}$ → Right skewed

titanic ["sqr_Fare"] = titanic [Fare]**(2)

titanic ["sqr_ Fare"].skew ()

titanic ["sqr_ Fare"]. hist ()

it can be any value

(2) (3) (4) . . . .

* when The Data is left skewed.

Out :/

* We apply $n^{th}$ power To The Data (or) column.

## * Exponential

titanic ["sqr-exp Fare"] = np. exp (titanic ["Fare"])

titanic ["sqr-exp Fare"] .skew ()

titanic ["sqr-exp Fare"]. hist ()

plt. show

→ another method
For Left skewed
Distribution

* Evaluates $e^1x$ For Each
Element in The given
Input.

## * ==Reciprocal== ==Transformation==

titanic ["Rec - r"] = ==1/ [titanic ["Fare"]+0.01)==

titanic ["Rec - Fare"] . skew ()

titanic ["Rec - Fare"]. hist ()

plt. show ()

"Reciprocal" is that
We divide Every Value
with ①/ [Column] +0.0

of the
data

to Elimi
"0" value
In The Data.

## * ==Box COX== Transformation

| Fare | |
|------|------|
| 7.250 | $\frac{7.250-1}{\lambda}$ |
| 71.283 | $\frac{71.283^{\lambda}-1}{\lambda}$ |
| — | |
| — | |
| — | |

$$T[x] = \frac{x^{\lambda} - 1}{\lambda}$$

∴ $\lambda$ [lamba] Varies From -5 to +

∴ Where X is The response Variable and

⇒ In This Transformation, all Values o

"$\lambda$" are Considered and The optimal

Value for a given Variables are selected.

from scipy import stats

# stats.boxcox (titanic ["Fare"] $\frac{+0.01}{+1}$) → +0.01 (or) +1 (user defined)

Out :- (array ([2.3844558, 6.43357655, 2.512739737.
2.6079914, 5.7648427, 4.06763781..
]))

Which ever gives best result

From help

we see Two return. values

1. array values (boxcox)

2. optimal Lamda parameter (max log)

# Fitting data, Fitting Lambda

# titanic ["Fare_boxcox"], param = stats.boxcox
(titanic ["Fare"]+1)Out :- a,b = ↑ ↑

1

2

Multiple variable assignmen

stores

# print ("λ=, param) (or) param

Out :- -0.09

# titanic ["Fare_boxcox"]. skew ()

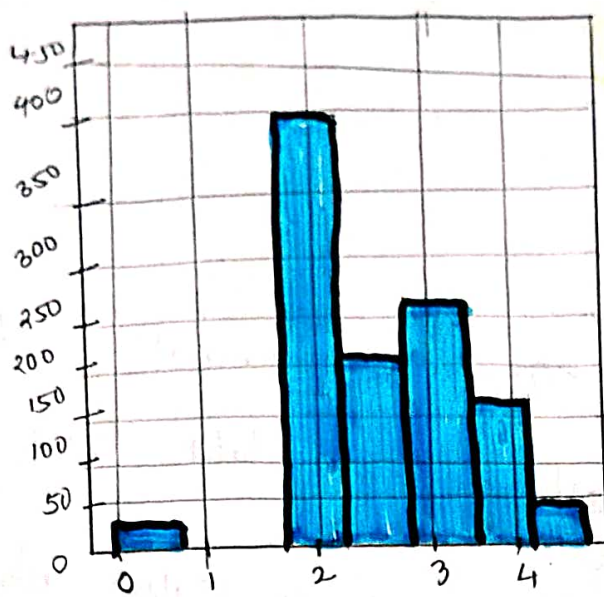Out :- -0.04 → lowest values while compare with other Transformations

## Code

# titanic ["Fare_boxcox], param = stats.boxcox (titanic.Fare +1

# titanic ["Fare_boxcox]. hist ()

# titanic ["Fare_boxcox]. skew()

Out :- -0.04

4/4/22
5:30pm.