5/4/22 10:30pm

**3** * **Discretization** *

↓

Converting Continuous Data To Discrete Data

↓
also called as "Binning"

code :-

Import pandas as pd

# stroke = pd.read_csv("Stroke prediction.csv")

# stroke.head()

Out

| id | gend | age | hyper Tension | heart-disease | Ever_married | Work type | Residance Type | avg-glucose Level | bmi | Smok status | Stroke |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 30669 | Male | 3.0 | 0 | 0 | No | children | Rural | 95.12 | 18.0 | NaN | 0 |
| 30468 | Male | 58.0 | 1 | 0 | Yes | Private | urban | 87.96 | 39.2 | Never | 0 |
| 16523 | Femal | 8.0 | 0 | 0 | No | Private | urban | 110.89 | 17.6 | NAN | 0 |
| 56543 | Femal | 70.0 | 0 | 0 | Yes | Private | Rural | 69.04 | 35.9 | Formaly smoker | 0 |
| 46136 | male | 14.0 | 0 | 0 | No | Never_worked | Rural | 161.28 | 19.1 | NAN | 0 |

discrete (cat) / contiou / discrete (cont) / discrete (cont) / discrete / discrete / discreu / contiou / contin / discr / discre

# stroke.Shape

Out : (43400, 12)

# stroke.info()

Out: information about stroke data.

⇒ Create Bins

\# intervals = [0, 12, 19, 30, 60, 90] → \#bins ⇒ Each interval

\# Categories = ["child", "teenager", "young adult", "middle aged", "Senior_citizen"]

**$ Creating New column & storing The data of "Age**

Divide.

\# Storke ["Age_category"] = pd. cut [X =

X = Stroke ['Age'], bins = intervals,

Labels = categories]

\#the Age column is Continous Data. so, We converting into Discrete Variable

intervals → categories
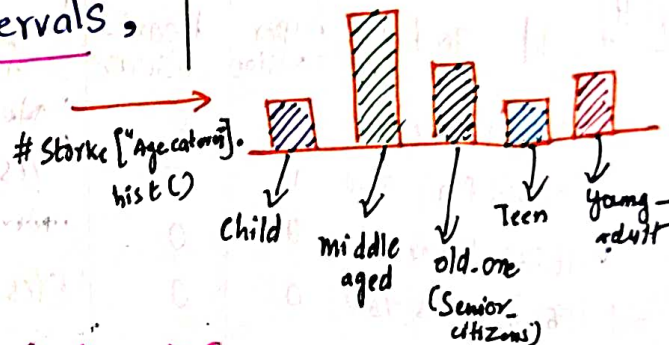* 0-12 → child
* 12-19 → teenager
* 19-30 → young adult
* 30-60 → Middle aged
* 60-90 → Senior citizen



\# Storke ["Age catorg]. hist ()

Child
middle aged
old.one (Senior citizens)
Teen
young adult

\# stroke.head ( )

Out

changed Continous Data To Discrete Categorical data

Newly-added

| Id | gender | age | hyper Tension | Heart disease | Ever married | work type | Residence Type | Avg. glucose Level | bmi | Smoke | Stroke | Age_Category |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 30669 | Male | 3.0 | - | - | - | - | - | - | - | - | - | child |
| 30468 | Male | 58.0 | - | - | - | - | - | - | - | - | - | middle aged |
| 16523 | Female | 8.0 | - | - | - | - | - | - | - | - | - | child |
| 56543 | Female | 70.0 | - | - | - | - | - | - | - | - | - | Senior citizen |
| 46143 | Male | 14.0 | - | - | - | - | - | - | - | - | - | teen |

**4** * **Encoding** *

Discrete

**Applicable For Categorical Variables**

* Machine Can't Understand Text Data. So, We Convert the "discrete categorical" to "discrete count" Data

* There are two types of Categorical Data

    1. Nominal (No Natural Sequence)
    2. Ordinal (Natural Sequence)

Ex:- ↓ Grades

| | |
|---|---|
| O | 10 |
| A+ | 9 |
| A | 8 |
| B+ | 7 |
| B | 6 |
| F | 0 |

⇒ O > A+
1079

# we should Treat Every variable Equally.

# odinal Data

cityname

| | |
|---|---|
| Hyd | 0 |
| Delhi | 1 |
| Mumbai | 2 |
| Kerala | 3 |
| Gujrat | 4 |

In this case, we can't Decide by numbering given To state.

which is greater

⇒ To convert Categorical Data
To Numeric :

# Nominal Data → (dummies)

    1. get Dummies (Pandas)
    2. One hot Encoding (sklearn)

# Odinal Data

    1. map (pandas)
    2. Label Encoder (sklearn)

| # One Hot Encoding | # Label Encoding |
|---|---|
| * Nominal | * Odinal |

```
import Pandas as pd
import Numpy as np
import matplotlib.pyplot as plt

% matplotlib inline
```

# df = pd.read_csv("homeprices.csv")

# df.head()

Out :-

| | Town | area | Price |
|---|---|---|---|
| 0 | chennai | 2600 | 5500000 |
| 1 | chennai | 3000 | 5650000 |
| 2 | chennai | 3200 | 6100000 |
| 3 | chennai | 3600 | 6800000 |
| 4 | Banglore | 2600 | 5850000 |

5pt

# Town

1. Chennai
2. Banglore
3. hyderabad    Discrete
   → Categorical Data

Nominal
Data

⇒ **Pd. get_dummies** [ Nominal Variable Encoding using pandas ]

# dummies = Pd. get_dummies ( df ["Town"])

# dummies

Out :-

| | Banglore | Chennai | Hyderabad |
|---|---|---|---|
| 0 | 0 | 1 | 0 |
| 1 | 0 | 1 | 0 |
| 2 | 0 | 1 | 0 |
| 3 | 0 | 1 | 0 |
| 4 | 1 | 0 | 0 |
| 5 | 1 | 0 | 0 |
| 6 | 1 | 0 | 0 |
| 7 | 1 | 0 | 0 |
| 8 | 0 | 0 | 1 |
| 9 | 0 | 0 | 1 |
| 10 | 0 | 0 | 1 |

**Out** :-

| | Banglore | Chennai | Hyderabd |
|---|---|---|---|
| 0 | 0 | 1 | 0 |
| 1 | 0 | 1 | 0 |
| 2 | 0 | 1 | 0 |
| 3 | 0 | 1 | 0 |
| 4 | 1 | 0 | 0 |
| 5 | 1 | 0 | 0 |
| 6 | 1 | 0 | 0 |
| 7 | 1 | 0 | 0 |
| 8 | 0 | 0 | 1 |
| 9 | 0 | 0 | 1 |
| 10 | 0 | 0 | 1 |
| 11 | 0 | 0 | 1 |

→ Dummies

# which Ever Record is there
it shows = 1

* For other columns it
shows = 0

→ Adding    → original    → converted one    by

# df_dummies = pd. Concat ([ df , dummies ] , axis = "columns"

df _ dumies

**Out** :-

| | Town | Area | Price | Banglore | Chennai | Hyderabad |
|---|---|---|---|---|---|---|
| 0 | Chennai | 2600 | 5500000 | 0 | 1 | 0 |
| 1 | Chennai | 3000 | 5650000 | 0 | 1 | 0 |
| 2 | Chennai | 3200 | 6100000 | 0 | 1 | 0 |
| 3 | Chennai | 3600 | 6800000 | 0 | 1 | 0 |
| 4 | Banglore | 3600 | 5850000 | 1 | 0 | 0 |
| 5 | Banglore | 2600 | 6150000 | 1 | 0 | 0 |
| 6 | Banglore | 2800 | 6500000 | 1 | 0 | 0 |
| 7 | Banglore | 3300 | 6100000 | 1 | 0 | 0 |
| | | | 7100000 | 0 | 0 | 1 |
| 8 | Hyderabad | 3600 | 5750000 | 0 | 0 | 1 |
| 9 | Hyderabad | 2600 | 6000000 | 0 | 0 | 1 |
| 10 | Hyderabad | 2900 | | 0 | 0 | 1 |
| 11 | Hyderabad | 3100 | 6200000 | 0 | 0 | 1 |

# Delete The Original column

# df_dummies . drop ("town", axis = "columns", inplace = True)

df_dummies

Out :

| | Area | Price | Banglore | chennai | Hyderabad |
|---|---|---|---|---|---|
| 0 | 2600 | 5500000 | 0 | 1 | 0 |
| 1 | 3000 | 5650000 | 0 | 1 | 0 |
| 2 | 3200 | 6100000 | 0 | 1 | 0 |
| 3 | 3800 | 6800000 | 0 | 1 | 0 |
| 4 | 2600 | 5850000 | 1 | 0 | 0 |
| 5 | 2800 | 6150000 | 1 | 0 | 0 |
| 6 | 3300 | 7100000 | 1 | 0 | 0 |
| 7 | 3600 | 5750000 | 1 | 0 | 0 |
| 8 | 2600 | 6000000 | 0 | 0 | 1 |
| 9 | 2900 | 6200000 | 0 | 0 | 1 |
| 10 | 3100 | 6900000 | 0 | 0 | 1 |
| 11 | 3600 | 6500000 | 0 | 0 | 1 |

# after deleting original column in New Data (ie. df_dummies)
Old_data (df) is still remains Same.

∅ **Dummy Variable Trap**

⇒
| Bang | chenn | Hyder |
|---|---|---|
| 0 | 1 | 0 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |
| 1 | 0 | 0 |
| 0 | 0 | 1 |
| 0 | 0 | 1 |

→

| Bang | chenn |
|---|---|
| 0 | 1 |
| 0 | 1 |
| 1 | 0 |
| 1 | 0 |
| 0 | 0 |
| 0 | 0 |

→ hyderabad

If we Remove One
of The column
also, still we can
Identify, col name by
(0,0)

# Town — 3 caterogies
↓ Dummies
3 columns [ Multi-collinearity/colleration Problem ]

# [If column] → ["n" categories]

We convert

↓

[(n-1) column]

# df_dummies . drop ("chennai", axis="columns", inplace=True

# df_dummies

out

|   | Area | Price | Banglore | Hyderabad |
|---|------|-------|----------|-----------|
| 0 | 2,600 | 5500000 | 0 | 0 |
| 1 | — | — | 0 | 0 |
| 2 | — | — | 0 | 0 |
| 3 | — | — | 0 | 0 |
| 4 | — | — | 0 | 0 |
| 5 | — | — | 1 | 0 |
| 6 | — | — | 1 | 0 |
| 7 | — | — | 1 | 0 |
| 8 | — | — | 1 | 0 |
| 9 | — | — | 0 | 1 |
| 10 | — | — | 0 | 1 |
| 11 | 36 000 | 6500000 | 0 | 1 |

(Same) as old data

# where ever Both column Having Equal Zero's, it i another column' data.

Chennai = 0 0

# df_dum = pd. get_dummies (df, drop_first = True)

# df_dum

Out :

| Area | Price | town_chennai | town_hyderabad |
|------|-------|--------------|----------------|
| 0 | | 1 | 0 |
| 1 | | 1 | 0 |
| 2 | | 1 | 0 |
| ... | | 1 | 0 |
| 10 | | | |
| 11 | | | |

↓
#This Diagram
is Exact of
(previous) Last page diagram.
But deleted banglore Replaced chennai'

town_banglore [deleted] → Because, B ✓ (Alphabet order)
chennai'
Replaced with ← h.yo
(0 0)

⇒ |One Hot Encoding| ———→ |Nominal using sklearn|

From sklearn. Preprocessing **import** One Hot Encoder

# enc = One Hot Encoder ( drop = "First")
↑
Stores in enc

# enc_df = pd. DataFrame ( enc.fit_transform,
↳ ( df [[ "town" ]]).
↳ to array ()

Converting

Output

# enc-df

Out

| | 0 | 1 |
|---|---|---|
| 0 | 1.0 | 0.0 |
| 1 | 1.0 | 0.0 |
| 2 | 1.0 | 0.0 |
| 3 | 1.0 | 0.0 |
| 4 | 0.0 | 0.0 |
| 5 | 0.0 | 0.0 |
| 6 | 0.0 | 0.0 |
| 7 | 0.0 | 0.0 |
| 8 | 0.0 | 1.0 |
| 9 | 0.0 | 1.0 |
| 10 | 0.0 | 1.0 |
| 11 | 0.0 | 1.0 |

# merge with main df

# df_ohe = df.join (enc-df)

# df_ohe.drop ("town", axis = "Columns", inplace = True)

# df_ohe

Out

| | area | Price | 0 | 1 |
|---|---|---|---|---|
| 0 | 2600 | 5500000 | 1.0 | 0.0 |
| 1 | 3000 | 5650000 | 1.0 | 0.0 |
| 2 | — | — | 1.0 | 0.0 |
| 3 | — | — | 1.0 | 0.0 |
| 4 | — | — | 0.0 | 0.0 |
| 5 | — | — | 0.0 | 0.0 |
| 6 | — | — | 0.0 | 0.0 |
| 7 | — | — | 0.0 | 0.0 |
| 8 | — | — | 0.0 | 0.0 |
| 9 | — | — | 0.0 | 1.0 |
| 10 | — | — | 0.0 | 1.0 |
| 11 | — | — | 0.0 | 1.0 |
| | | | 0.0 | 1.0 |

$\Rightarrow$ <span style="color:pink">Alphabetical Order</span>

**Label Encoder** $\rightarrow$ **Ordinal Variable** [sklearn]

# Used for binary category Variable

# dfnew = df.copy() $\rightarrow$ Creating "copy" From original Data.

# dfnew

Out

| | town | Area | Price |
|---|---|---|---|
| 0 | Chennai | 2600 | 5500000 |
| 1 | chennai | 3000 | 5650000 |
| 2 | Chennai | 3250 | 61000000 |
| . | : | : | : |
| . | : | : | |
| 11 | Hyderabad | 3600 | 6950000 |

From sklearn.preprocessing import Label Encoder

# Le = Label Encoder()

# dfnew.town = Le.fit_transform (dfnew.town)

# dfnew

Out :-

| | town | Area | Price |
|---|---|---|---|
| 0 | 1 | | |
| 1 | 1 | | |
| 2 | 1 | | |
| 3 | 1 | | |
| 4 | 0 | | |
| 5 | 0 | | |
| 6 | 0 | | |
| 7 | 0 | | |
| 8 | 2 | | |
| 9 | 2 | | |
| 10 | 2 | | |
| 11 | 2 | | |

<span style="color:red">Alphabetical Order</span>
$\downarrow$
Banglore [B] — 0
chennai [C] — 1
Hyderabd [H] — 2

√ User Desired values we use this.

⇒ map() → Ordinal Variable [Pandas]

# map function.

# df_m = df.copy()

# df_m ["town"] = df_m ["town"] . map ({ "chennai" :0, "Bangalore" :2 , "hyderabad": 1})

# df_m

Out :-

| | town | Area | Price |
|---|---|---|---|
| 0 | 0 | | |
| 1 | 0 | | |
| 2 | 0 | | |
| 3 | 0 | | |
| 4 | 2 | | |
| 5 | 2 | | |
| 6 | 2 | | |
| 7 | 2 | | |
| 8 | 1 | | |
| 9 | 1 | | |
| 10 | 1 | | |
| 11 | 1 | | |

# We can Replace with any Value

# Taken by our choices Here
chennai = 0 , 20
Banglore = 2 , 15
hyderabad = 1 , 25

Which Ever Order we write it Takes in the Order

⇒ Ordinal Encoder → Ordinal Variable (Ascending Order)

From Sklearn. Pre processing import Ordinal Encoder

# df_new1 = df.copy()

# Oe = Ordinal Encoder (categories = [[ "Banglore", "Hyderabad", "chennai"]])

# df_new1.town = Oe.fit_transform (df_new1 [[ "town"]])

# df_new1

Out :-

| town |
|---|
| 2.0 |
| 2.0 |
| 2.0 |
| 2.0 |
| 0.0 |
| 0.0 |
| 0.0 |
| 0.0 |
| 1.0 |
| 1.0 |
| 1.0 |
| 1.0 |

6/4/22 3:00AM