

last
class
=

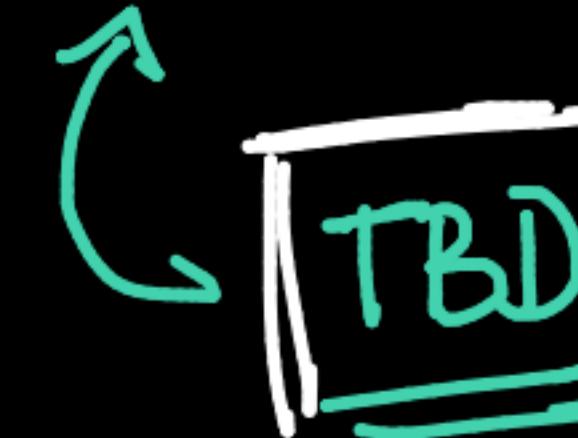
- Topics:
1. Handling missing values → many more
↳ basic
 2. Categorical to Numeric conversion } → Loan approval =
 3. Find most "useful" variables in Loan approval } → ...
 4. Pre-processing Datasets → Normalization & Standardization =

Ops:

→ [① 1 vs 2 vs 3 class break]
11.4% = 8.9%

19.7%
revise
solve problems & case-study

② Problem-solving sessions] least freq ans-
questions





EDA_FE.ipynb - Colaboratory x markov x +

colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=f5n7flegGQEe

Connect ▾  

+ Code + Text

↳ 11 cells hidden

↑ ↓ ↻ ⌂ ⌂ ⌂ ⌂ ⌂ ⌂ ⌂ ⌂ ⌂ ⌂ ⌂

Basic Data Exploration

[] ↳ 9 cells hidden

Basic Data visualization: Univariate

[] ↳ 11 cells hidden

Simple Feature Engineering ↗ abi[li]y-to-pay-EMI

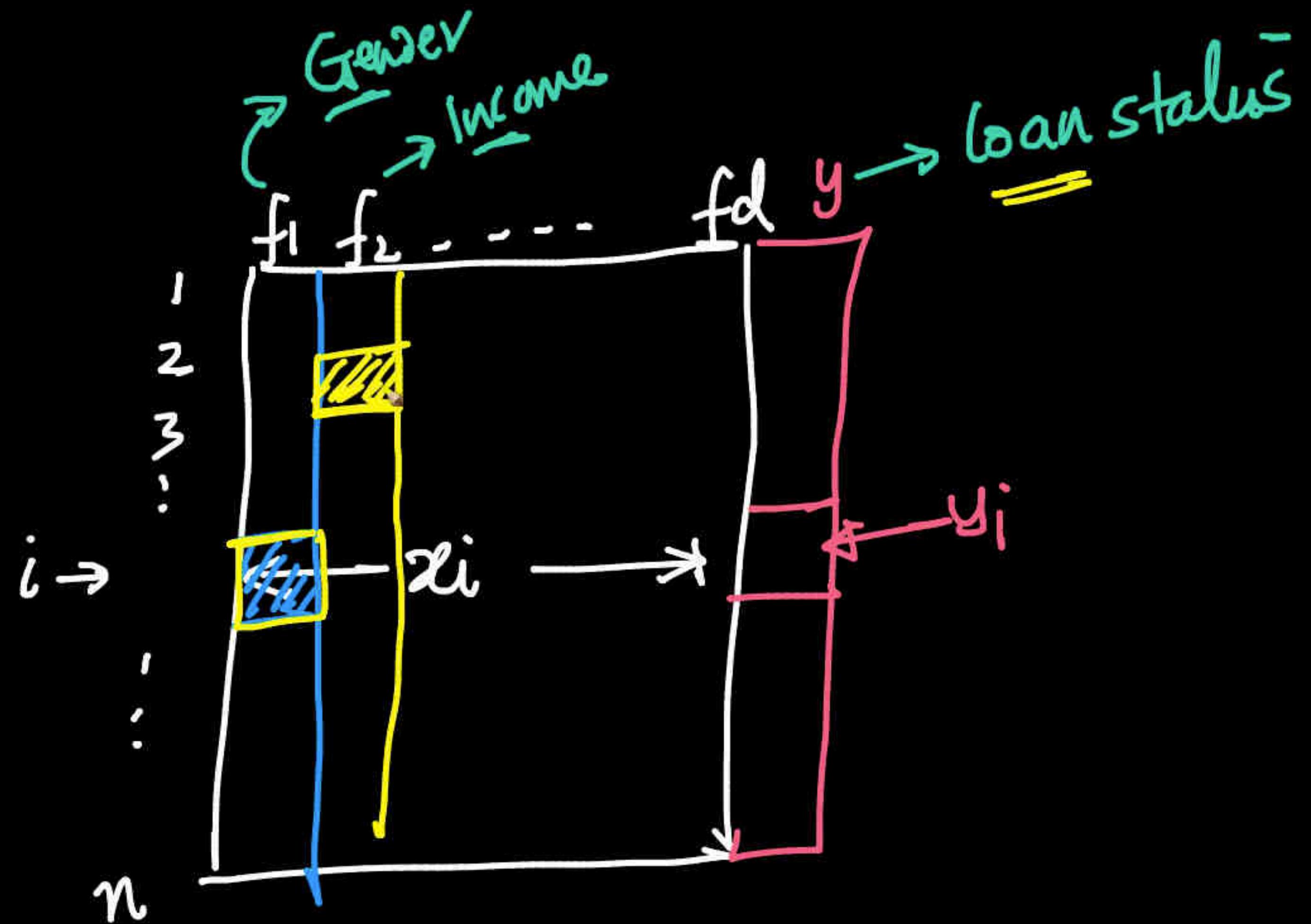
[] ↳ 37 cells hidden

Missing Values & Data Cleaning

[] ↳ 9 cells hidden

4 / 4

P



Missing Data

Q categorical
f_i: Gender → M or F (binary) discrete

missing-value-imputation



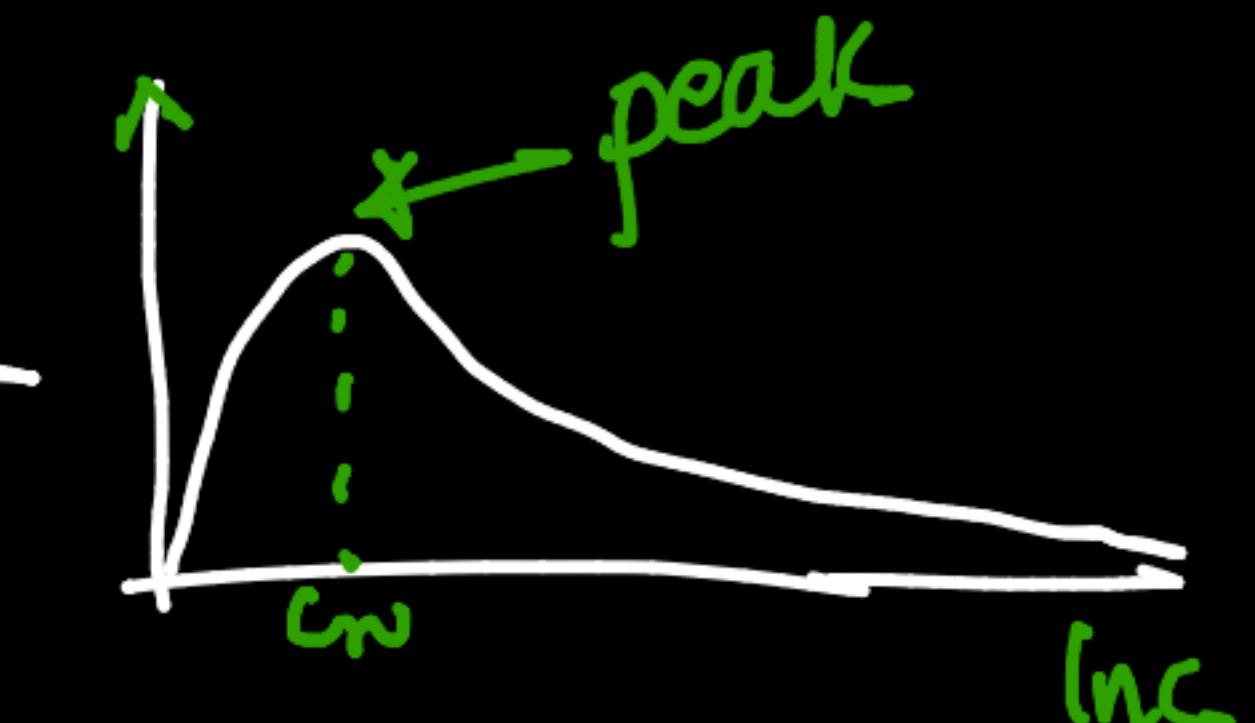
- ↳ mode = most freq. occurring value
- ↳ random
- "Other" ✓

↑ there is some info in the fact
that the data
is missing

Q

Income: Numeric

mean or median



614 2015

[options]

median income

→ mode

using Gender, age, Piaget
Model-based imputation
(later)

EDA_FE.ipynb - Colaboratory

markov

https://colab.research.google.com/drive/1_P-nKswaKnx5v7zIOx3pCxyJye73MvnN#scrollTo=SkB8xczbRc

Update :

+ Code + Text

Connect ▾

V

▼ Missionary Values & Dada Vitality



`data.isna().sum()`

```
Gender          13  
Married         13  
Dependents     13  
Education       13  
Self_Employed  32  
ApplicantIncome 13  
CoapplicantIncome 13  
LoanAmount      22  
Loan_Amount_Term 14  
Credit_History  50  
Property_Area   13  
Loan_Status      13  
TotalIncome      13  
Loan_Amount_per_year 30  
EMI              30  
Able_to_pay_EMI 13  
dtype: int64
```

EDA_FE.ipynb - Colaboratory

markov

colab.research.google.com/drive/1_P-nKswaKnx5y7ZIOx3pCxyJye73MvnN#scrollTo=SkB8xczb8c

Update :

+ Code + Text

Connect ▾

V

▼ Missing values & Data Cleaning

A set of small, light-gray navigation icons located at the bottom right of the screen. From left to right, they include: a double arrow for search, a magnifying glass for search, a circular arrow for refresh, a double arrow for history, a gear for settings, a square with rounded corners for full screen, a trash can for delete, and three vertical dots for more options.

▶ `data.isna().sum()`

```
Gender          13  
Married         5  
Dependents     15  
Education        0  
Self_Employed   32  
ApplicantIncome  0  
CoapplicantIncome 0  
LoanAmount      22  
Loan_Amount_Term 14  
Credit_History   50  
Property_Area    0  
Loan_Status       0  
TotalIncome       0  
Loan_Amount_per_year 36  
EMI              36  
Able_to_pay_EMI  0  
dtype: int64
```

```
[ ] # Function to create a data frame with number and percentage of missing data in a data frame
```

```
def missing_to_df(df):
```

EDA_FE.ipynb - Colaboratory

markov

```
Code + Text  
Loan_Amount_per_year      36  
EMI                         36  
Able_to_pay_EMI                0  
dtype: int64
```

✓ [32] # Function to create a data frame with number and percentage of missing data in a data frame

```
def missing_to_df(df):
    #Number and percentage of missing data in training data set for each column
    total_missing_df = df.isnull().sum().sort_values(ascending=False)
    percent_missing_df = (df.isnull().sum()/df.isnull().count()*100).sort_values(ascending=False)
    missing_data_df = pd.concat([total_missing_df, percent_missing_df], axis=1, keys=['Total', 'Percent'])
    return missing_data_df
```

```
[33] missing_df = missing_to_df(data)
      missing_df[missing_df['Total'] > 0]
```

	Total	Percent
Credit_History	50	8.143322
Loan_Amount_per_year	36	5.863192
EMI	36	5.863192
Self_Employed	32	5.211726

EDA_FE.ipynb - Colaboratory

markov

1

Code + Text

- ✓ RAM
- Disk

A small icon representing user settings or preferences.

▼

```
[33] missing_df = missing_to_df(data)
      missing_df[missing_df['Total'] > 0]
```

	Total	Percent
Credit_History	50	8.143322
Loan_Amount_per_year	36	5.863192
EMI	36	5.863192
Self_Employed	32	5.211726
LoanAmount	22	3.583062
Dependents	15	2.442997
Loan_Amount_Term	14	2.280130
Gender	13	2.117264
Married	3	0.488599

```
# Credit History=2 for nan/missing values.  
data[ 'Credit_History' ] = data[ 'Credit_History' ].fillna(2)
```

EDA_FE.ipynb - Colaboratory x markov x +

colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=bQzPBBynqbx2h

+ Code + Text RAM Disk

EMI 36

Able_to_pay_EMI 0

dtype: int64

{x}

[32] # Function to create a data frame with number and percentage of missing data in a data frame

```
def missing_to_df(df):
    #Number and percentage of missing data in training data set for each column
    total_missing_df = df.isnull().sum().sort_values(ascending=False)
    percent_missing_df = (df.isnull().sum()/df.isnull().count()*100).sort_values(ascending=False)
    missing_data_df = pd.concat([total_missing_df, percent_missing_df], axis=1, keys=[ 'Total', 'Percent'])
    return missing_data_df
```

[33] missing_df = missing_to_df(data)

missing_df[missing_df['Total'] > 0]

	Total	Percent
Credit_History	50	8.143322
Loan_Amount_per_year	36	5.863192
EMI	36	5.863192
Self_Employed	32	5.211726

RAM Disk

12 / 12

EDA_FE.ipynb - Colaboratory x markov x | + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=bQzPBBynqbx2h Update :

+ Code + Text

TotalIncome 0
Loan_Amount_per_year 36
EMI 36
Able_to_pay_EMI 0
dtype: int64

[32] # Function to create a data frame with number and percentage of missing data in a data frame

```
{ def missing_to_df(df):
    #Number and percentage of missing data in training data set for each column
    total_missing_df = df.isnull().sum().sort_values(ascending=False)
    percent_missing_df = (df.isnull().sum()/df.isnull().count()*100).sort_values(ascending=False)
    missing_data_df = pd.concat([total_missing_df, percent_missing_df], axis=1, keys=['Total', 'Percent'])
    return missing_data_df }
```

[33] missing_df = missing_to_df(data)
missing_df[missing_df['Total'] > 0]

	Total	Percent
Credit_History	50	8.143322
Loan_Amount_per_year	36	5.863192
EMI	36	5.863192

EDA_FE.ipynb - Colaboratory x markov x + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=bQzPBBynqbx2h Update :

+ Code + Text

Credit History vs Loan Approval

is CIBI Score

```
[ ] data['Credit_History'].value_counts()
```

1.0	475
0.0	89

```
[ ] Name: Credit_History, dtype: int64
```

```
[ ] sns.countplot(data = data, x = 'Credit_History', hue = 'Loan_Status')  
#Observation:  
## We can clearly see that the approval rate is 80% if your credit history is aligned with the guidelines.  
## Hence this is the most important question that can be considered.
```

```
[ ]
```

Missing Values & Data Cleaning

RAM Disk

14 / 14

EDA_FE.ipynb - Colaboratory x markov x +

colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=bQzPBBynqbx2h

+ Code + Text RAM Disk

Credit History vs Loan Approval

```
{x} [ ] data['Credit_History'].value_counts()
```

```
1.0 475
```

```
0.0 89
```

```
Name: Credit_History, dtype: int64
```



```
[ ] sns.countplot(data = data, x = 'Credit_History', hue = 'Loan_Status')  
#Observation:  
## We can clearly see that the approval rate is 80% if your credit history is aligned with the guidelines.  
## Hence this is the most important question that can be considered.
```

```
[ ]
```

Missing Values & Data Cleaning

EDA_FE.ipynb - Colaboratory x markov x + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=mezoRK5waMTr

+ Code + Text RAM Disk

Name: Credit_History, dtype: int64

Q {x} D :

0s

sns.countplot(data = data, x = 'Credit_History', hue = 'Loan_Status')

#Observation:

We can clearly see that the approval rate is 80% if your credit history is aligned with the guidelines.

Hence this is the most important question that can be considered.

<matplotlib.axes._subplots.AxesSubplot at 0x7fac71b2cc50>

Count

Credit_History

Loan_Status

Y

N

0.0

1.0

2.0

NO

Yes

350

300

250

200

150

100

50

0

16 / 16

EDA_FE.ipynb - Colaboratory x markov x | + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=mezoRK5waMTr

+ Code + Text RAM Disk

[33] missing_df = missing_to_df(data)
missing_df[missing_df['Total'] > 0]

	Total	Percent
Credit_History	50	8.143322
Loan_Amount_per_year	36	5.863192
EMI	36	5.863192
Self_Employed	32	5.211726
LoanAmount	22	3.583062
Dependents	15	2.442997
Loan_Amount_Term	14	2.280130
Gender	13	2.117264
Married	3	0.488599

[34] # Credit History=2 for nan/missing values.
data['Credit_History'] = data['Credit_History'].fillna(2)

GA - UH

2

1

0

Credit history

UID → Gibit

missing

if loan given → 1
else → 0

17 / 18

EDA_FE.ipynb - Colaboratory x markov x | + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=mezoRK5waMTr

+ Code + Text RAM Disk

EMI 36 5.863192

Self_Employed 32 5.211726

LoanAmount 22 3.583062

Dependents 15 2.442997

Loan_Amount_Term 14 2.280130

Gender 13 2.117264

Married 3 0.488599

[34] # Credit History=2 for nan/missing values.
data['Credit_History'] = data['Credit_History'].fillna(2)

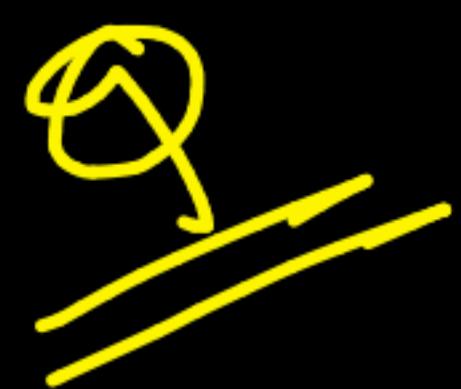
[] # Self_Employed = 'Other' for nan
data.Self_Employed.unique()

array(['No', 'Yes', nan], dtype=object)

[] data['Self_Employed'] = data['Self_Employed'].fillna('Other')

RAM Disk

18 / 19



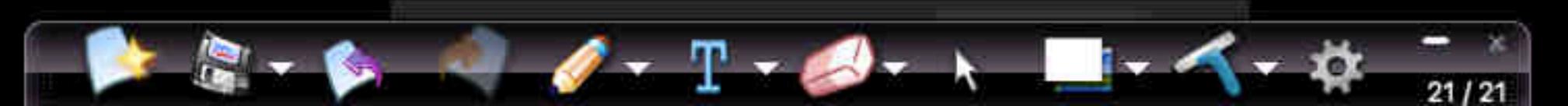
imputation f_i

$$P(\tilde{f}_i |$$

Model-based imputation
 $f_1, f_2, \dots, f_{i-1}, f_{i+1}, \dots, f_d$

Simple Bayes
(Gumberson)

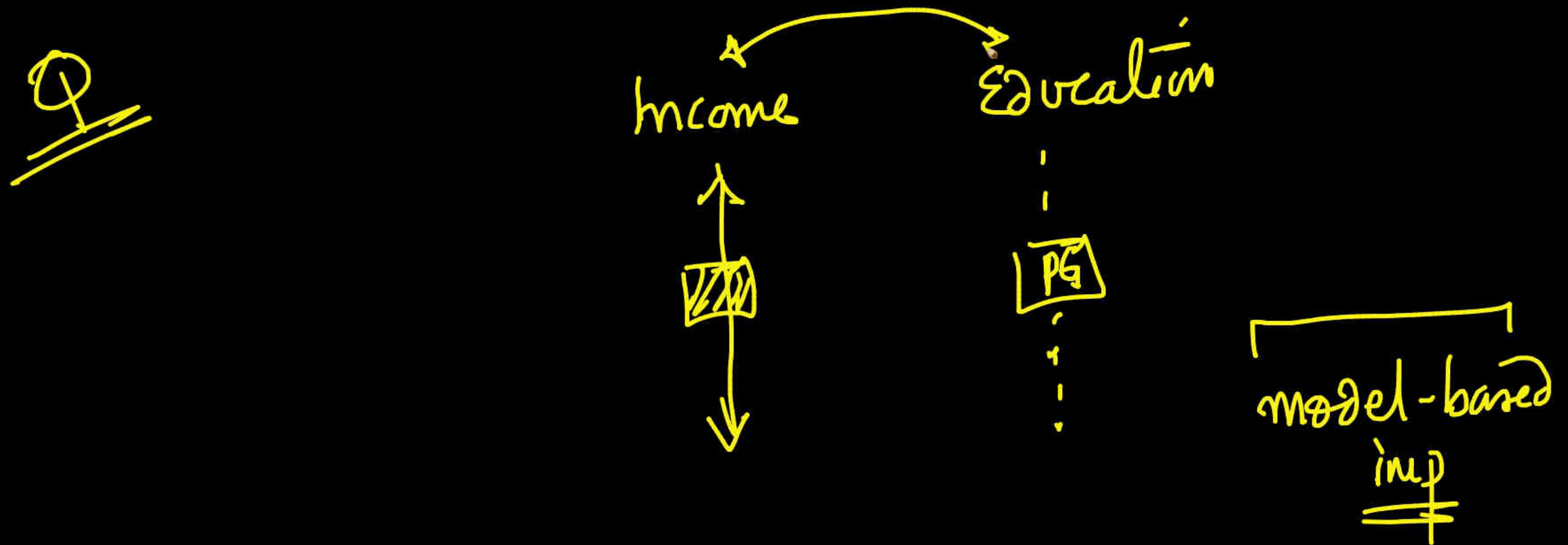
ML-models



Q

$$n = \# \text{rows} = \overbrace{1 \dots M}$$

1 rows → missing → OK to copy



EDA_FE.ipynb - Colaboratory x markov x | +

colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=mezoRK5waMTr

+ Code + Text RAM Disk

Total Percent

	Total	Percent
Credit_History	50	8.143322
Loan_Amount_per_year	36	5.863192
EMI	36	5.863192
Self_Employed	32	5.211726
LoanAmount	22	3.583062
Dependents	15	2.442997
Loan_Amount_Term	14	2.280130
Gender	13	2.117264
Married	3	0.488599

{x}

✓ Credit_History 50 8.143322

✓ Loan_Amount_per_year 36 5.863192

✓ EMI 36 5.863192

✓ Self_Employed 32 5.211726

✓ LoanAmount 22 3.583062

✓ Dependents 15 2.442997

✓ Loan_Amount_Term 14 2.280130

✓ Gender 13 2.117264

✓ Married 3 0.488599

[34] # Credit History=2 for nan/missing values.
data['Credit_History'] = data['Credit_History'].fillna(2)

[] # Self_Employed = 'Other' for nan
data.Self_Employed.unique()

RAM Disk

23 / 23

EDA_FE.ipynb - Colaboratory x markov x | + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=mezoRK5waMTr

+ Code + Text RAM Disk

Gender 13 2.117264

Married 3 0.488599

{x} [34] # Credit History=2 for nan/missing values.
data['Credit_History'] = data['Credit_History'].fillna(2)

[] # Self_Employed = 'Other' for nan
data.Self_Employed.unique()

array(['No', 'Yes', nan], dtype=object)

{ } data['Self_Employed'] = data['Self_Employed'].fillna('Other')

[] # median imputation for numerical columns.
from sklearn.impute import SimpleImputer

<> num_missing = ['EMI', 'Loan_Amount_per_year', 'LoanAmount', 'Loan_Amount_Term']

median_imputer = SimpleImputer(strategy = 'median')
for col in num_missing:

24 / 24

EDA_FE.ipynb - Colaboratory x markov x | + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=mezoRK5waMTr

+ Code + Text

RAM Disk

Gender 13 2.117264
Married 3 0.488599

[34] # Credit History=2 for nan/missing values.
data['Credit_History'] = data['Credit_History'].fillna(2)

[] # Self_Employed = 'Other' for nan
data.Self_Employed.unique()

array(['No', 'Yes', nan], dtype=object)

[] data['Self_Employed'] = data['Self_Employed'].fillna('Other')

[] # median imputation for numerical columns.
from sklearn.impute import SimpleImputer

num_missing = ['EMI', 'Loan_Amount_per_year', 'LoanAmount', 'Loan_Amount_Term']

median_imputer = SimpleImputer(strategy = 'median')
for col in num_missing:

$y_i = \text{loan-app}$

f_j

\vdots

\vdots

\vdots

\vdots

\checkmark more

N X

EDA_FE.ipynb - Colaboratory x markov x | + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=mezoRK5waMTr

+ Code + Text RAM Disk ✓

	Gender	13	2.117264
Q	Married	3	0.488599

{x} [34] # Credit History=2 for nan/missing values.
data['Credit_History'] = data['Credit_History'].fillna(2)

[] # Self_Employed = 'Other' for nan
data.Self_Employed.unique()

array(['No', 'Yes', nan], dtype=object)

[] data['Self_Employed'] = data['Self_Employed'].fillna('Other')

[] # median imputation for numerical columns.
from sklearn.impute import SimpleImputer

<> num_missing = ['EMI', 'Loan_Amount_per_year', 'LoanAmount', 'Loan_Amount_Term']

median_imputer = SimpleImputer(strategy = 'median')
for col in num_missing:

Chrome File Edit View History Bookmarks Profiles Tab Window Help

EDA_FE.ipynb - Colaboratory x markov x | +

colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=mezoRK5waMTr

+ Code + Text

RAM Disk

ML-lib

su'kut-Learn }

```
data.Self_Employed.unique()
array(['No', 'Yes', nan], dtype=object)

[ ] data['Self_Employed'] = data['Self_Employed'].fillna('Other')

[ ] # median imputation for numerical columns
from sklearn.impute import SimpleImputer
num_missing = ['EMI', 'Loan_Amount_per_year', 'LoanAmount', 'Loan_Amount_Term']

median_imputer = SimpleImputer(strategy = 'median')
for col in num_missing:
    data[col] = pd.DataFrame(median_imputer.fit_transform(pd.DataFrame(data[col])))

[ ] # Highest Freq imputation for some categorical columns.
cat_missing = ['Gender', 'Married', 'Dependents']

freq_imputer = SimpleImputer(strategy = 'most_frequent')
for col in cat_missing:
    data[col] = pd.DataFrame(freq_imputer.fit_transform(pd.DataFrame(data[col])))
```

Chrome File Edit View History Bookmarks Profiles Tab Window Help

EDA_FE.ipynb - Colaboratory x markov x | +

colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=mezoRK5waMTr

RAM Disk ✓

+ Code + Text

```
data.Self_Employed.unique()

array(['No', 'Yes', nan], dtype=object)

{x}
[ ] data['Self_Employed'] = data['Self_Employed'].fillna('Other')

[ ]
# median imputation for numerical columns.
from sklearn.impute import SimpleImputer

num_missing = ['EMI', 'Loan_Amount_per_year', 'LoanAmount', 'Loan_Amount_Term']

median_imputer = SimpleImputer(strategy = 'median')
for col in num_missing:
    data[col] = pd.DataFrame(median_imputer.fit_transform(pd.DataFrame(data[col])))

[ ]
# Highest Freq imputation for some categorical columns.
cat_missing = ['Gender', 'Married', 'Dependents']

freq_imputer = SimpleImputer(strategy = 'most_frequent')
for col in cat_missing:
    data[col] = pd.DataFrame(freq_imputer.fit_transform(pd.DataFrame(data[col])))
```

EDA_FE.ipynb - Colaboratory | markov | scikit-learn.org/stable/modules/generated/sklearn.impute.SimpleImputer.html

Toggle Menu

scikit learn

Install User Guide API Examples Community More

Prev Up Next

scikit-learn 1.1.1

Other versions

Please cite us if you use the software.

sklearn.impute.SimpleImputer

Examples using sklearn.impute.SimpleImputer

sklearn.impute.SimpleImputer

```
class sklearn.impute.SimpleImputer(*, missing_values=nan, strategy='mean', fill_value=None, verbose='deprecated', copy=True, add_indicator=False)
```

[source]

Imputation transformer for completing missing values.

Read more in the [User Guide](#).

New in version 0.20: SimpleImputer replaces the previous `sklearn.preprocessing.Imputer` estimator which is now removed.

Parameters:

missing_values : int, float, str, np.nan, None or pandas.NA, default=np.nan
The placeholder for the missing values. All occurrences of `missing_values` will be imputed. For pandas' dataframes with nullable integer dtypes with missing values, `missing_values` can be set to either `np.nan` or `pd.NA`.

strategy : str, default='mean'
The imputation strategy.

If "mean", then replace missing values using the mean along each column. Can only

EDA_FE.ipynb - Colaboratory | markov | sklearn.impute.SimpleImputer | + | scikit-learn.org/stable/modules/generated/sklearn.impute.SimpleImputer.html | Update | : |

scikit learn

Prev Up Next

scikit-learn 1.1.1

Other versions

Please cite us if you use the software.

sklearn.impute.SimpleImputer

Examples using sklearn.impute.SimpleImputer

Parameters: **missing_values : int, float, str, np.nan, None or pandas.NA, default=np.nan**

The placeholder for the missing values. All occurrences of `missing_values` will be imputed. For pandas' dataframes with nullable integer dtypes with missing values, `missing_values` can be set to either `np.nan` or `pd.NA`.

strategy : str, default='mean'

The imputation strategy.

- If "mean", then replace missing values using the mean along each column. Can only be used with numeric data.
- If "median", then replace missing values using the median along each column. Can only be used with numeric data.
- If "most_frequent", then replace missing using the most frequent value along each column. Can be used with strings or numeric data. If there is more than one such value, only the smallest is returned.
- If "constant", then replace missing values with `fill_value`. Can be used with strings or numeric data.

New in version 0.20: `strategy="constant"` for fixed value imputation.

fill_value : str or numerical value, default=None

When `strategy == "constant"`, `fill_value` is used to replace all occurrences of `fill_value`. If `fill_value` is not specified, the default `fill_value` will be 0 when imputing numerical data types.

EDA_FE.ipynb - Colaboratory

markov

X | sklearn.impute.SimpleImputer | X

```
+ Code + Text
```

```
data.Self_Employed.unique()
```

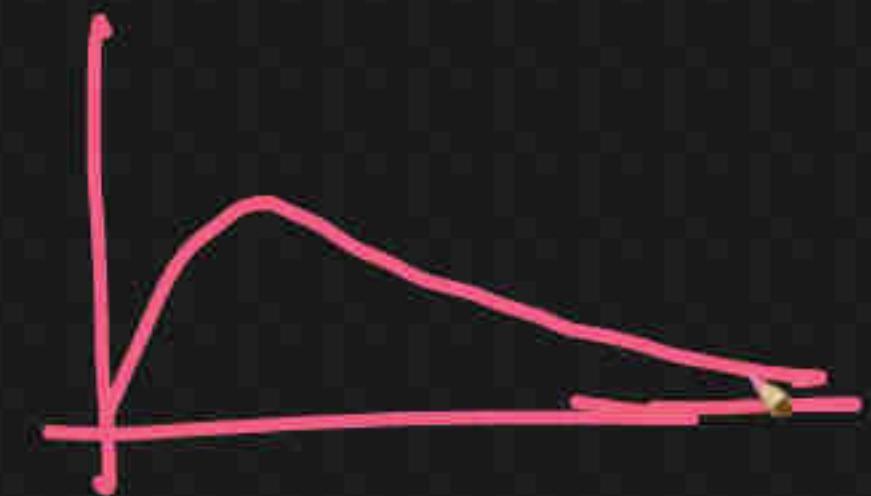
```
array(['No', 'Yes', nan], dtype=object)
```

```
[ ] data[ 'Self_Employed' ] = data[ 'Self_Employed' ].fillna( 'Other' )
```

```
# median imputation for numerical columns.  
from sklearn.impute import SimpleImputer  
  
num_missing = ['EMI', 'Loan_Amount_per_year', 'LoanAmount', 'Loan_Amount_Term']  
  
median_imputer = SimpleImputer(strategy = 'median')  
for col in num_missing:  
    data[col] = pd.DataFrame(median_imputer.fit_transform(pd.DataFrame(data[col])))
```

```
[ ] # Highest Freq imputation for some categorical columns  
cat_missing = ['Gender', 'Married', 'Dependents']
```

```
freq_imputer = SimpleImputer(strategy = 'most_frequent')
for col in cat_missing:
    data[col] = pd.DataFrame(freq_imputer.fit_transform(
```



EDA_FE.ipynb - Colaboratory x markov x | sklearn.impute.SimpleImputer x + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=NeFywieNcFYI

+ Code + Text RAM Disk ✓

```
array(['No', 'Yes', nan], dtype=object)
```

[] data['Self_Employed'] = data['Self_Employed'].fillna('Other')

{x}

[] # median imputation for numerical columns.
from sklearn.impute import SimpleImputer

num_missing = ['EMI', 'Loan_Amount_per_year', 'LoanAmount', 'Loan_Amount_Term']

median_imputer = SimpleImputer(strategy = 'median')
for col in num_missing:
 data[col] = pd.DataFrame(median_imputer.fit_transform(pd.DataFrame(data[col])))

[] # Highest Freq imputation for some categorical columns.
cat_missing = ['Gender', 'Married', 'Dependents']

freq_imputer = SimpleImputer(strategy = 'most_frequent')
for col in cat_missing:
 data[col] = pd.DataFrame(freq_imputer.fit_transform(pd.DataFrame(data[col])))

[] missing_df = missing_to_df(data)
missing_df[missing_df['Total'] > 0]

RAM Disk ✓

Up Down Reload Settings Copy Copy All Delete More

32 / 32

EDA_FE.ipynb - Colaboratory x markov x | sklearn.impute.SimpleImputer x + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=NeFywieNcFYI

+ Code + Text RAM Disk

```
num_missing = ['EMI', 'Loan_Amount_per_year', 'LoanAmount', 'Loan_Amount_Term']

median_imputer = SimpleImputer(strategy = 'median')
for col in num_missing:
    data[col] = pd.DataFrame(median_imputer.fit_transform(pd.DataFrame(data[col])))

[ ] # Highest Freq imputation for some categorical columns.
cat_missing = ['Gender', 'Married', 'Dependents']

freq_imputer = SimpleImputer(strategy = 'most_frequent')
for col in cat_missing:
    data[col] = pd.DataFrame(freq_imputer.fit_transform(pd.DataFrame(data[col])))

[ ] missing_df = missing_to_df(data)
missing_df[missing_df['Total'] > 0]
```

► Categorical to Numerical encoding

Nominal vs Ordinal variables

1. One Hot Encoding

```
EDA_FE.ipynb - Colaboratory x markov x | sklearn.impute.SimpleImputer x + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=_ZaDNwb8bpQf
```

+ Code + Text

[39] # Highest Freq imputation for some categorical columns.
cat_missing = ['Gender', 'Married', 'Dependents']

freq_imputer = SimpleImputer(strategy = 'most_frequent')
for col in cat_missing:
 data[col] = pd.DataFrame(freq_imputer.fit_transform(pd.DataFrame(data[col])))

[40] missing_df = missing_to_df(data)
missing_df[missing_df['Total'] > 0]



► Categorical to Numerical encoding

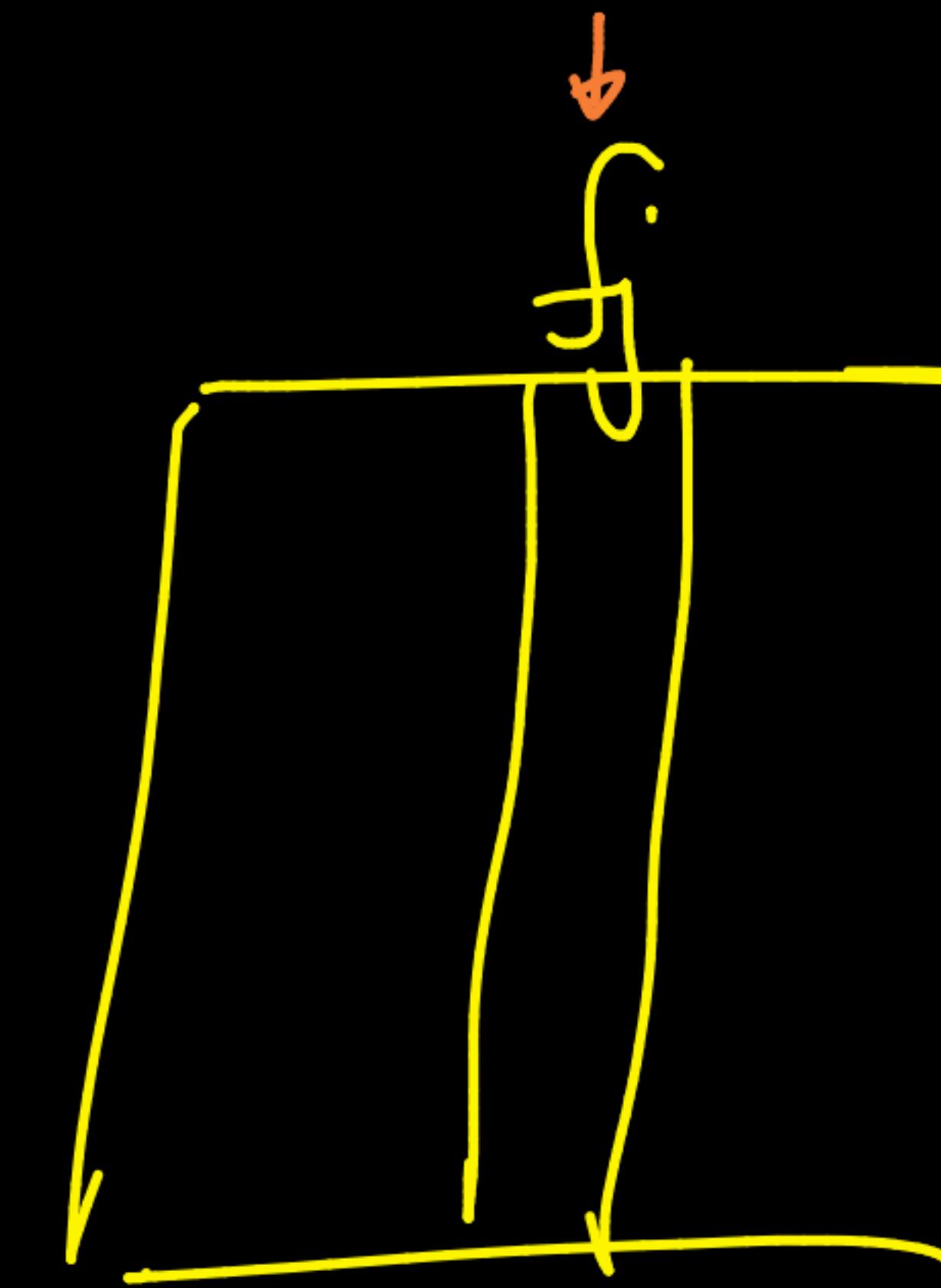
Nominal vs Ordinal variables

1. One Hot Encoding
2. Label encoding
3. Target Encoding

Appropriate encoding depends on what our task is (and) what we do next?

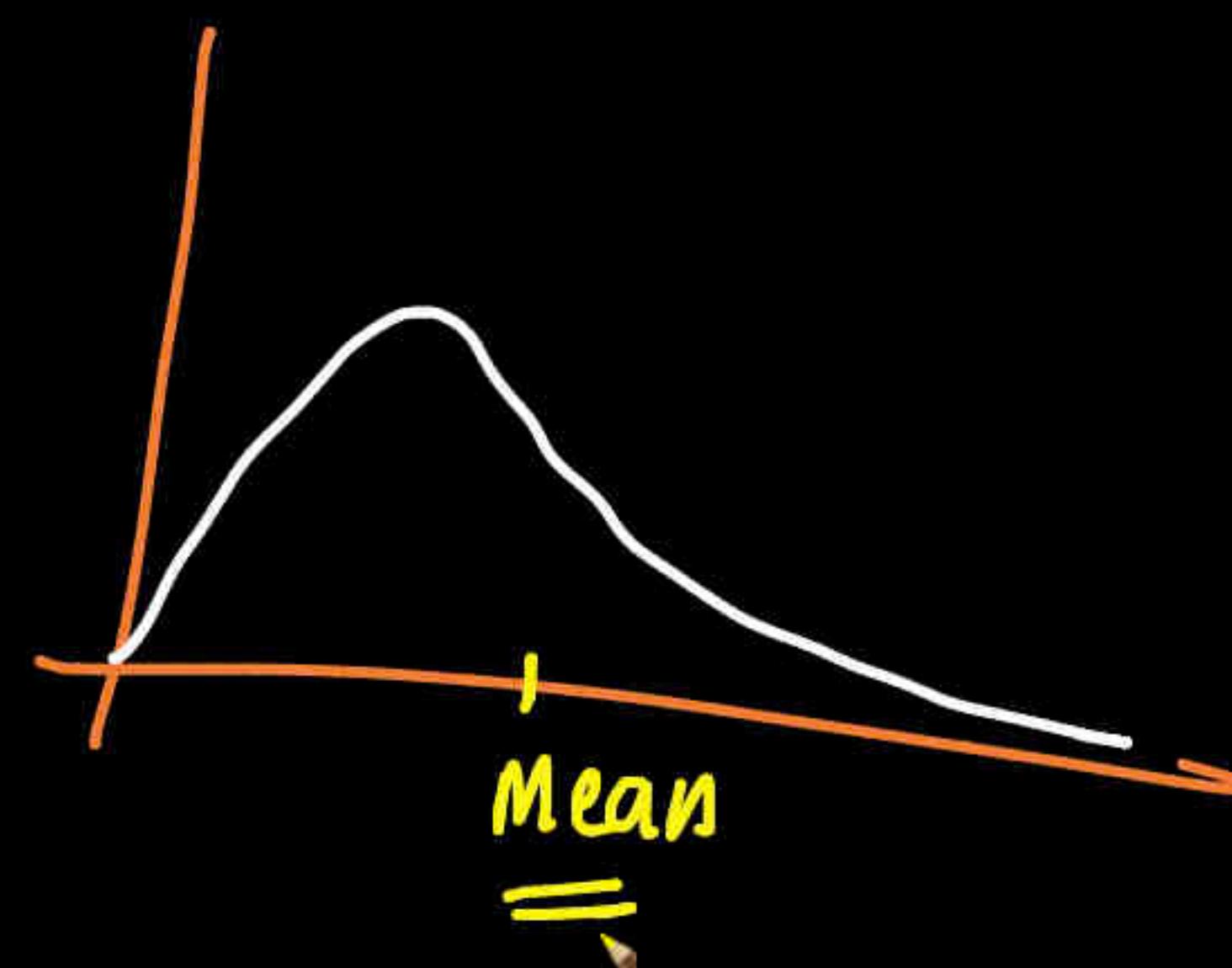
Tools: Compute Correlation (PCC and SPCC) between each feature and the Loan Status

Q

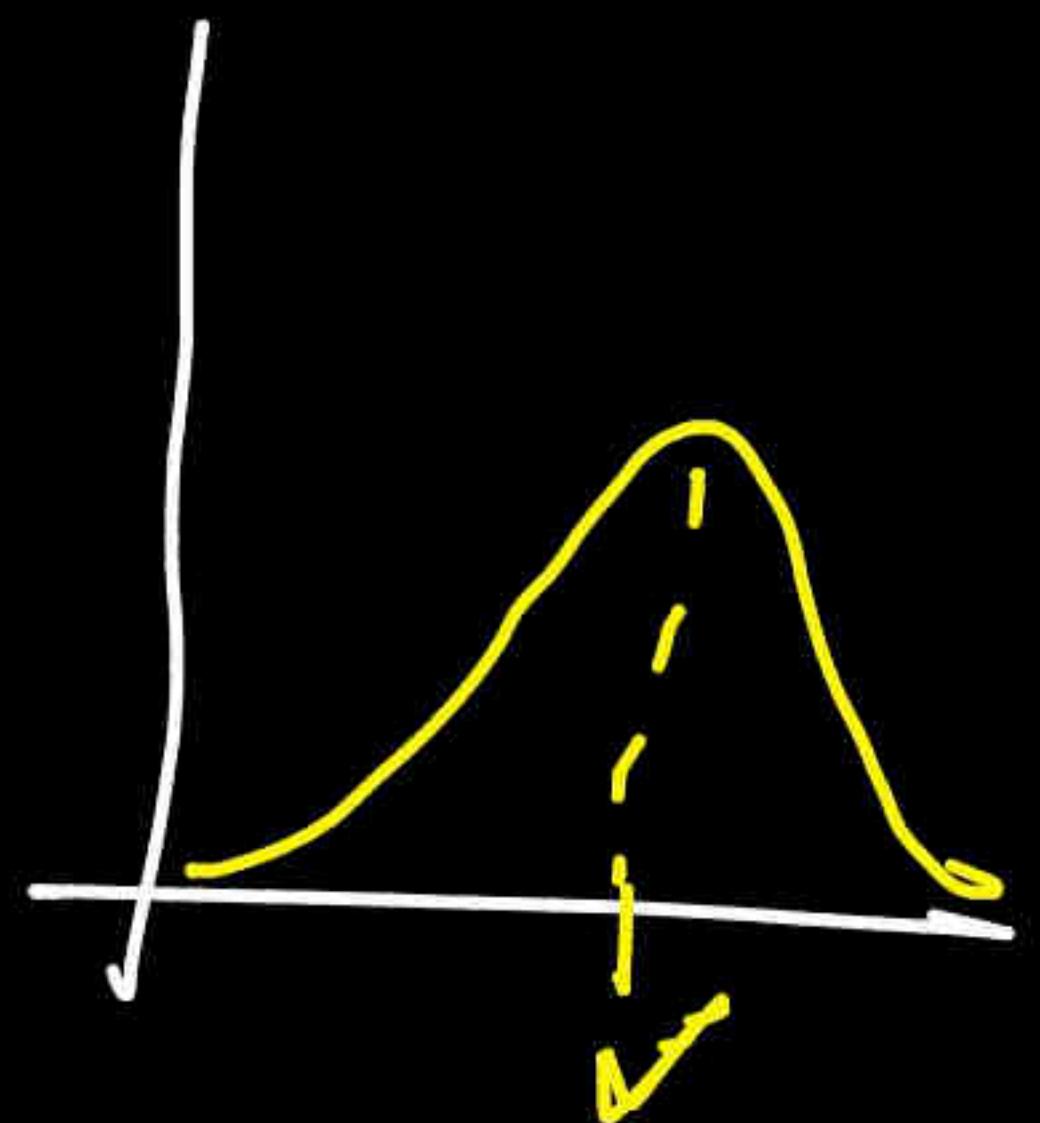


imputation
→ deep column ←

↑ >50% values are missing



(log)
boxcox



Chrome File Edit View History Bookmarks Profiles Tab Window Help

EDA_FE.ipynb - Colaboratory x markov x | sklearn.impute.SimpleImputer x + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=_ZaDNwb8bpQf

+ Code + Text

RAM Disk ✓

Update :

[38] # median imputation for numerical columns.

```
from sklearn.impute import SimpleImputer
```

{x}

num_missing = ['EMI' , 'Loan_Amount_per_year' , 'LoanAmount' , 'Loan_Amount_Term']

✓ median_imputer = SimpleImputer(strategy = 'median')

```
for col in num_missing:
```

```
    data[col] = pd.DataFrame(median_imputer.fit_transform(pd.DataFrame(data[col])))
```

[39] # Highest Freq imputation for some categorical columns.

```
cat_missing = ['Gender' , 'Married' , 'Dependents' ]
```

freq_imputer = SimpleImputer(strategy = 'most_frequent')

```
for col in cat_missing:
```

```
    data[col] = pd.DataFrame(freq_imputer.fit_transform(pd.DataFrame(data[col])))
```

[40] missing_df = missing_to_df(data)

```
missing_df[missing_df['Total'] > 0]
```

Total Percent

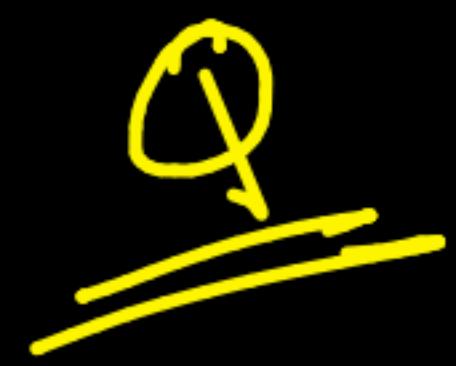
Categorical to Numerical encoding

37 / 37



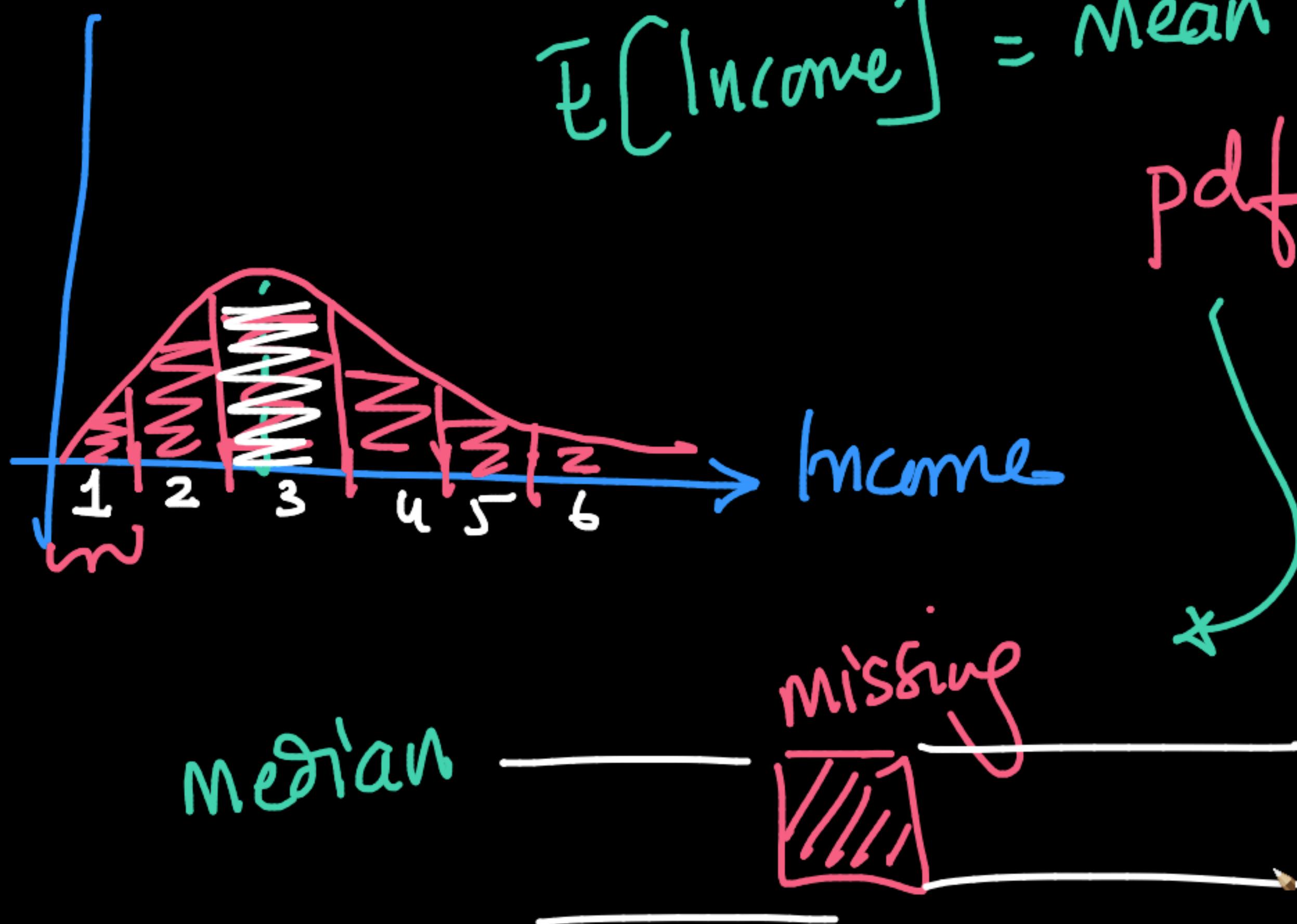
(MM 2008)

[00 2008 ← Missing → 2009]



$E[\text{Income}] = \text{Mean}$

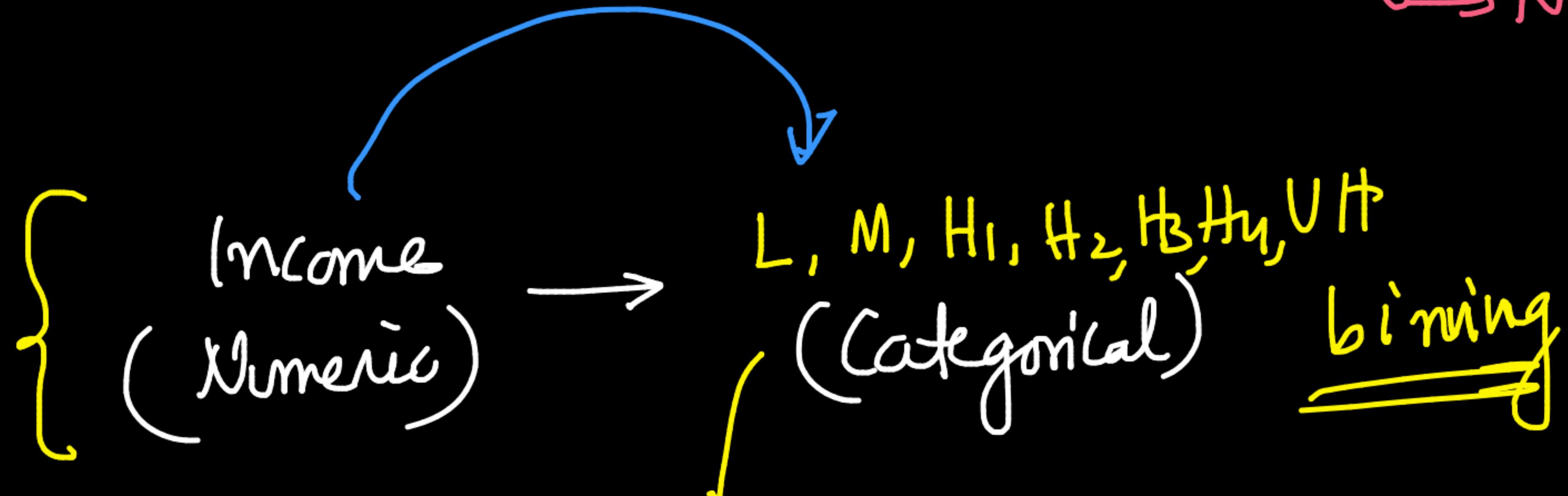
pdf



median

missing

Categorical
↳ Numeric



plots

χ^2 test

Income-bin vs
loan-approval

monotonic linear
SRCC & PCC → numeric

Gender
non-ordinal feature
 $\left\{ \begin{array}{l} M \rightarrow \boxed{+} \quad 1 \ 2 \ | \\ F \rightarrow \boxed{-} \quad 2 \ 1 \ 3 \\ O \rightarrow \boxed{0} \quad 3 \ 3 \ 2 \end{array} \right.$
numeric

loan status
↓
0 → NO
1 → YES
numeric

Edurealm

original features → natural ordering

HS	→ 0	4
UG	→ 1	5
PG	→ 2	10
PhD	→ 3	20
		new

Categorical

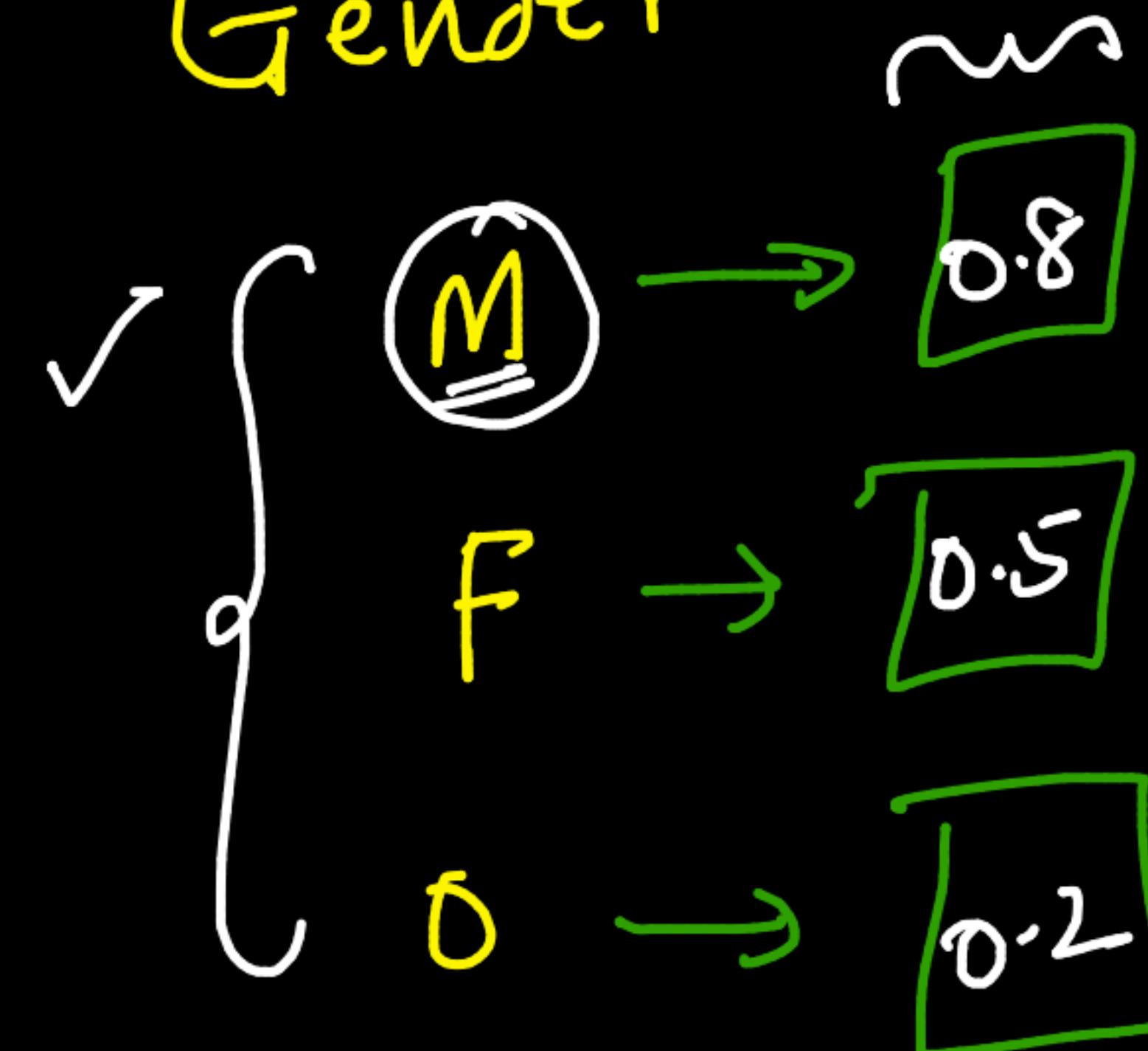
→ nominal → ?

→ ordinal → monotonic numbers

Gender

Data
613

new-cust
M
F



response coding ↗

$$p(\text{loanappn}=1 \mid G=M)$$

$$p(\text{loanApp}=1 \mid G=F)$$

$$p(\text{IA}=1 \mid G=O)$$

613 often

Obj / Task ↗

CC ↗

between

Gender & ↗

loanAppn ↗

loan approval

M → 0.8

F → 0.5

B → 0.2



response coding
=====

if G has no corr to loan approval

$$\begin{aligned} M &\rightarrow 0.6 \\ F &\rightarrow 0.61 \\ O &\rightarrow 0.59 \end{aligned}$$

EDA_FE.ipynb - Colaboratory x markov x | sklearn.impute.SimpleImputer x + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=sp3SJkJmKalm

+ Code + Text RAM Disk

```
data['Credit_History'] = data['Credit_History'].fillna(2)
```

{x} [36] # Self_Employed = 'Other' for nan
data.Self_Employed.unique()
array(['No', 'Yes', nan], dtype=object)

[37] data['Self_Employed'] = data['Self_Employed'].fillna('Other')

[38] # median imputation for numerical columns.
from sklearn.impute import SimpleImputer

num_missing = ['EMI', 'Loan_Amount_per_year', 'LoanAmount', 'Loan_Amount_Term']

median_imputer = SimpleImputer(strategy = 'median')
for col in num_missing:
 data[col] = pd.DataFrame(median_imputer.fit_transform(pd.DataFrame(data[col])))

[39] # Highest Freq imputation for some categorical columns.
cat_missing = ['Gender', 'Married', 'Dependents']

freq_imputer = SimpleImputer(strategy = 'most_frequent')



EDA_FE.ipynb - Colaboratory x markov x | sklearn.impute.SimpleImputer x + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=sp3SJkJmKalm

+ Code + Text RAM Disk

[38] # median imputation for numerical columns.

from sklearn.impute import SimpleImputer

✓ num_missing = ['EMI', 'Loan_Amount_per_year', 'LoanAmount', 'Loan_Amount_Term'] → list

median_imputer = SimpleImputer(strategy = 'median')

for col in num_missing:

data[col] = pd.DataFrame(median_imputer.fit_transform(pd.DataFrame(data[col])))

[39] # Highest Freq imputation for some categorical columns.

cat_missing = ['Gender', 'Married', 'Dependents']

freq_imputer = SimpleImputer(strategy = 'most_frequent')

for col in cat_missing:

data[col] = pd.DataFrame(freq_imputer.fit_transform(pd.DataFrame(data[col])))

[40] missing_df = missing_to_df(data)

missing_df[missing_df['Total'] > 0]

Total Percent

48 / 48

EDA_FE.ipynb - Colaboratory x markov x | sklearn.impute.SimpleImputer x + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=sp3SJkJmKalm

+ Code + Text RAM Disk

[38] # median imputation for numerical columns.
from sklearn.impute import SimpleImputer

num_missing = ['EMI', 'Loan_Amount_per_year', 'LoanAmount', 'Loan_Amount_Term']

median_imputer = SimpleImputer(strategy = 'median')
for col in num_missing:
 data[col] = pd.DataFrame(median_imputer.fit_transform(pd.DataFrame(data[col])))

[39] # Highest Freq imputation for some categorical columns.
cat_missing = ['Gender', 'Married', 'Dependents']

freq_imputer = SimpleImputer(strategy = 'most_frequent')
for col in cat_missing:
 data[col] = pd.DataFrame(freq_imputer.fit_transform(pd.DataFrame(data[col])))

[40] missing_df = missing_to_df(data)
missing_df[missing_df['Total'] > 0]

Total Percent

49 / 49

```
EDA_FE.ipynb - Colaboratory x markov x | sklearn.impute.SimpleImputer x + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=sp3SJkJmKalm + Code + Text RAM Disk ✓
```

[38] # median imputation for numerical columns.

```
from sklearn.impute import SimpleImputer
```

{ num_missing = ['EMI', 'Loan_Amount_per_year', 'LoanAmount', 'Loan_Amount_Term']

```
median_imputer = SimpleImputer(strategy = 'median')
```

```
for col in num_missing:
```

```
    data[col] = pd.DataFrame(median_imputer.fit_transform(pd.DataFrame(data[col])))
```

[39] # Highest Freq imputation for some categorical columns.

```
cat_missing = ['Gender', 'Married', 'Dependents']
```

{ freq_imputer = SimpleImputer(strategy = 'most_frequent')

```
for col in cat_missing:
```

```
    data[col] = pd.DataFrame(freq_imputer.fit_transform(pd.DataFrame(data[col])))
```

[40] missing_df = missing_to_df(data)

```
missing_df[missing_df['Total'] > 0]
```

Total Percent



Chrome File Edit View History Bookmarks Profiles Tab Window Help

EDA_FE.ipynb - Colaboratory x markov x sklearn.impute.SimpleImputer x + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=sp3SJkJmKalm

+ Code + Text RAM Disk ✓

[38] # median imputation for numerical columns.
from sklearn.impute import SimpleImputer

num_missing = ['EMI', 'Loan_Amount_per_year', 'LoanAmount', 'Loan_Amount_Term']

median_imputer = SimpleImputer(strategy = 'median')
for col in num_missing:
 data[col] = pd.DataFrame(median_imputer.fit_transform(pd.DataFrame(data[col])))

[39] # Highest Freq imputation for some categorical columns.
cat_missing = ['Gender', 'Married', 'Dependents']

freq_imputer = SimpleImputer(strategy = 'most_frequent')
for col in cat_missing:
 data[col] = pd.DataFrame(freq_imputer.fit_transform(pd.DataFrame(data[col])))

[40] missing_df = missing_to_df(data)
missing_df[missing_df['Total'] > 0]

Total Percent

RAM Disk ✓

51 / 51

Chrome File Edit View History Bookmarks Profiles Tab Window Help

EDA_FE.ipynb - Colaboratory x markov x | sklearn.impute.SimpleImputer x + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=sp3SJkJmKalm

+ Code + Text RAM Disk

[38] # median imputation for numerical columns.

```
from sklearn.impute import SimpleImputer

num_missing = ['EMI', 'Loan_Amount_per_year', 'LoanAmount', 'Loan_Amount_Term']

median_imputer = SimpleImputer(strategy = 'median')
for col in num_missing:
    data[col] = pd.DataFrame(median_imputer.fit_transform(pd.DataFrame(data[col])))
```

[39] # Highest Freq imputation for some categorical columns.

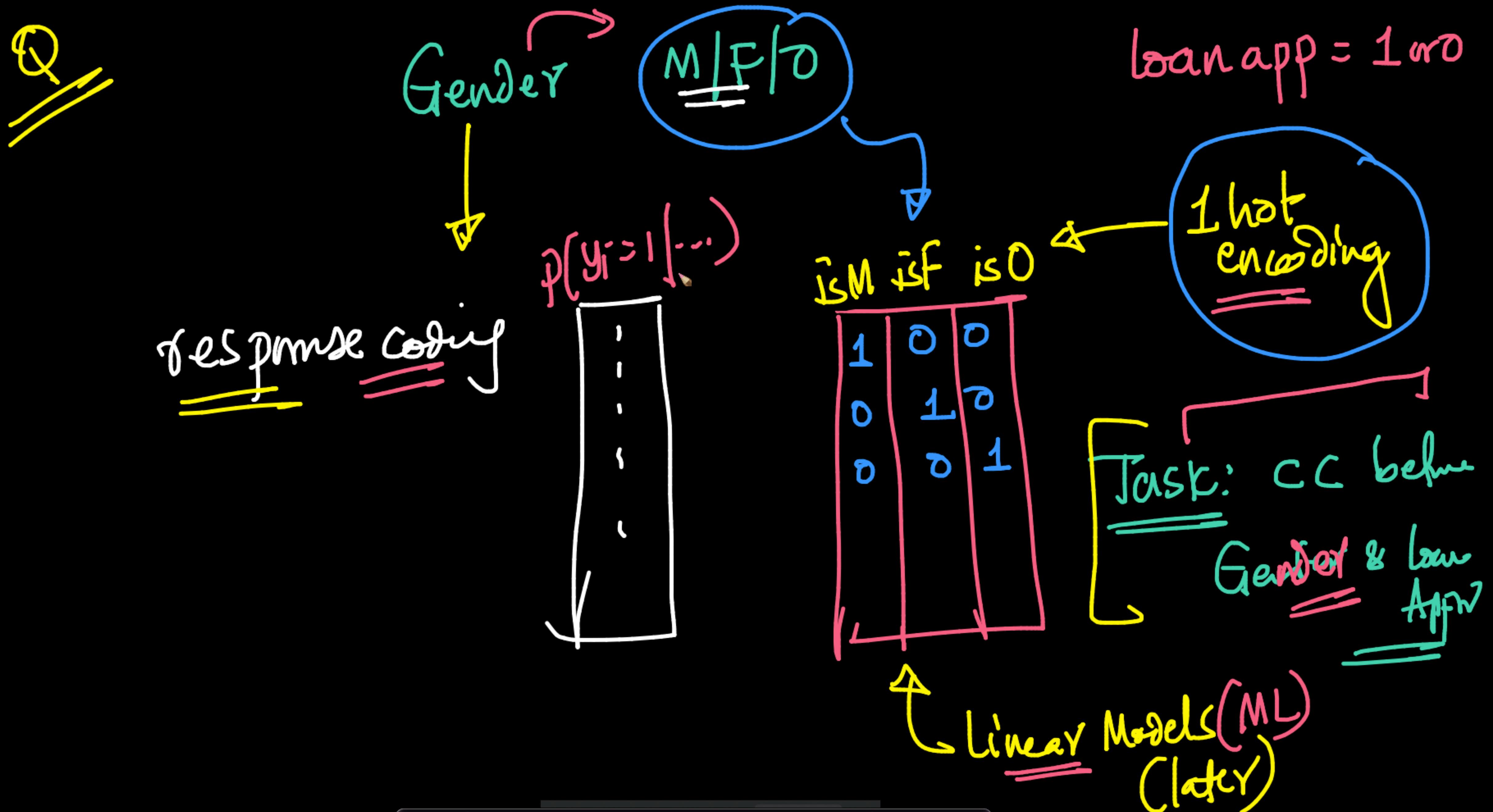
```
cat_missing = ['Gender', 'Married', 'Dependents']

freq_imputer = SimpleImputer(strategy = 'most_frequent')
for col in cat_missing:
    data[col] = pd.DataFrame(freq_imputer.fit_transform(pd.DataFrame(data[col])))
```

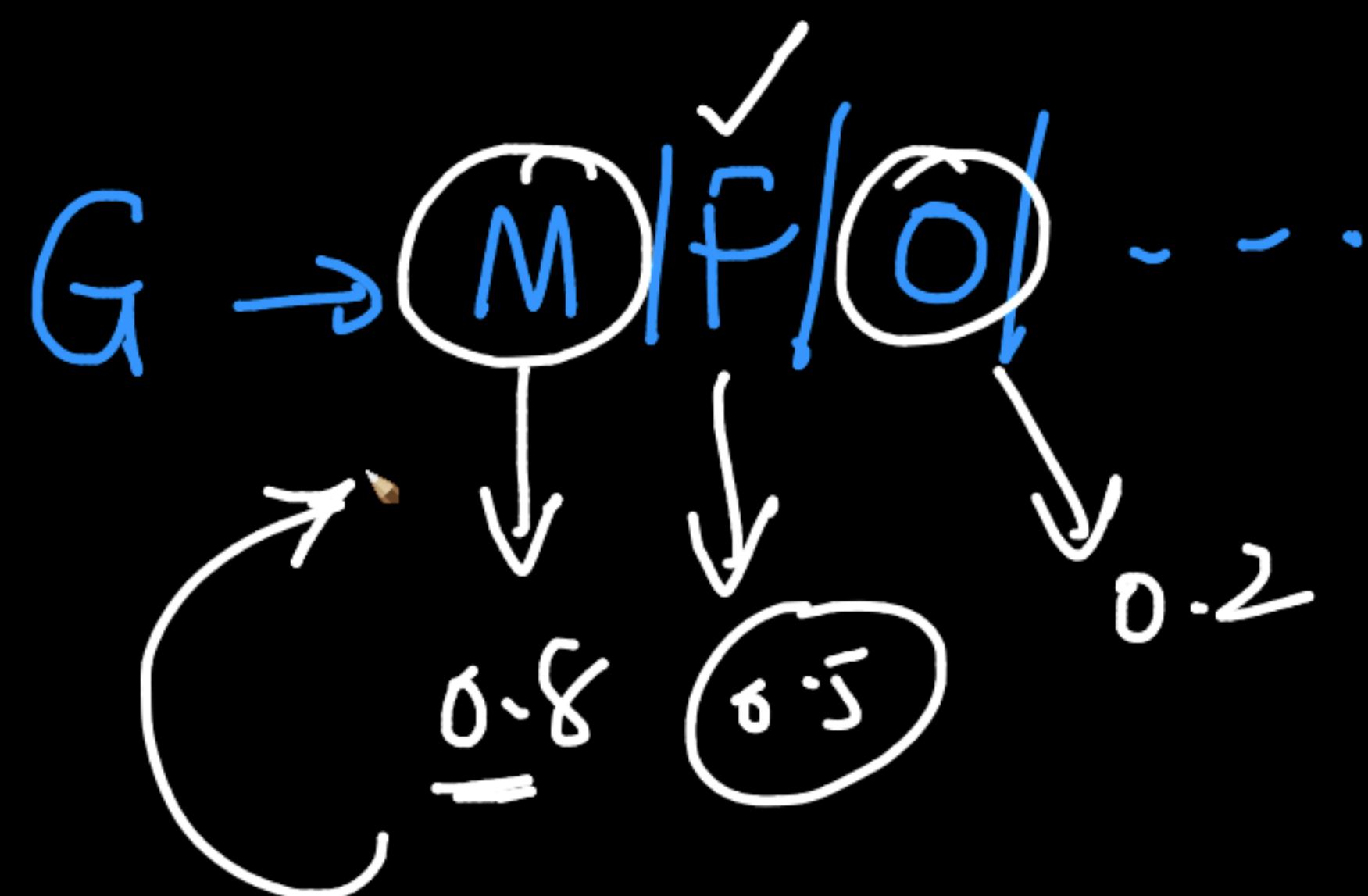
[40] missing_df = missing_to_df(data)
missing_df[missing_df['Total'] > 0]

Total Percent

52 / 52



is M Correlate with loan Appr





Educator

originality

responding
creativity

0.3

0.6

0.8

0.4

better

HS

UG

PG

PhD

1

2

3

4

1

2

3

4

Q

Evaluation

UG

HS

PhD

PG



10	2
----	---

$$\frac{2}{10}$$

data sparsity



not reliable

EDA_FE.ipynb - Colaboratory x markov x | sklearn.impute.SimpleImputer x + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=sp3SJkJmKalm

+ Code + Text RAM Disk

▼ Categorical to Numerical encoding

{ x } Nominal vs Ordinal variables

- 1. One Hot Encoding ✓
- 2. Label encoding ↗
- 3. Target Encoding ✓

Appropriate encoding depends on what our task is (and) what we do next?

Task: Compute Correlation (PCC and SRCC) between each feature and the Loan_Status

```
[ ] s = (data.dtypes == 'object')
object_cols = list(s[s].index)
object_cols
```

```
['Gender',
 'Married',
 'Education',
 'Self_Employed',
 'Property_Area',
 'Loan_Status']
```

+ Code + Text

✓ RAM Disk



▼ Categorical to Numerical encoding

Nominal vs Ordinal variables

1. One Hot Encoding
2. Label encoding
3. Target Encoding

Appropriate encoding depends on what our task is (and) what we do next?

Task: Compute Correlation (PCC and SRCC) between each feature and the Loan-Status

```
[ ] s = (data.dtypes == 'object')
object_cols = list(s[s].index)
object_cols
```

```
'Gender',
'Married',
'Education',
'Self_Employed',
'Property_Area',
'Loan_Status']
```

loan approved
0 , 1
2

+ Code + Text

✓ RAM Disk



▼ Categorical to Numerical encoding

Nominal vs Ordinal variables

1. One Hot Encoding
2. Label encoding
3. Target Encoding

Appropriate encoding depends on what our task is (and) what we do next?

Task: Compute Correlation (PCC and SRCC) between each feature and the Loan_Status

```
[ ] s = (data.dtypes == 'object')
object_cols = list(s[s].index)
object_cols
```

```
'Gender',
'Married',
'Education',
'Self_Employed',
'Property_Area',
'Loan_Status']
```

Gender
M, F, O
0, 2, J

EDA_FE.ipynb - Colaboratory x markov x | sklearn.impute.SimpleImputer x +

colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=ffGjSIJtgvPW

+ Code + Text RAM Disk

RAM Disk

Loan Status

{x} ✓

0s

Y 422
N 192
Name: Loan_Status, dtype: int64

[] from sklearn.preprocessing import LabelEncoder

label_encoder = LabelEncoder()
col='Loan_Status'
data[col] = label_encoder.fit_transform(data[col])

[] data[col].value_counts()

1 → 422
0 192
Name: Loan_Status, dtype: int64

EDA_FE.ipynb - Colaboratory x markov x | sklearn.impute.SimpleImputer x + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=U-iTHyZPiAek

+1 ✓ RAM Disk

- J SRCC

Gender

{x} [] #Gender
data['Gender'].value_counts()

Male 502
Female 112
Name: Gender, dtype: int64

[] from sklearn.preprocessing import LabelEncoder

label_encoder = LabelEncoder()
col='Gender'
data[col] = label_encoder.fit_transform(data[col])

[] data[col].value_counts()

1 502
0 112
Name: Gender, dtype: int64

EDA_FE.ipynb - Colaboratory x markov x | sklearn.impute.SimpleImputer x +

colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=_K4bmriUIDco

+ Code + Text

RAM Disk

Main menu

2 cells hidden

{x}

Property Area

```
[ ] # col='Property_Area'  
col='Property_Area'  
data[col].value_counts()
```

Property_Area	Count
Semiurban	233
Urban	202
Rural	179

```
[ ] !pip install category_encoders
```

Looking in indexes: <https://pypi.org/simple>, <https://us-python.pkg.dev/colab-wheels/public/simple/>

Requirement already satisfied: category_encoders in /usr/local/lib/python3.7/dist-packages (2.4.1)

Requirement already satisfied: scikit-learn>=0.20.0 in /usr/local/lib/python3.7/dist-packages (from category_encoder)

Requirement already satisfied: patsy>=0.5.1 in /usr/local/lib/python3.7/dist-packages (from category_encoders) (0.5.2)

Requirement already satisfied: pandas>=0.21.1 in /usr/local/lib/python3.7/dist-packages (from category_encoders) (1.1.5)

Requirement already satisfied: scipy>=1.0.0 in /usr/local/lib/python3.7/dist-packages (from category_encoders) (1.4.6)

EDA_FE.ipynb - Colaboratory x markov x | sklearn.impute.SimpleImputer x + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=_K4bmriUIDco

+ Code + Text RAM Disk

Requirement already satisfied: numpy>=1.14.0 in /usr/local/lib/python3.7/dist-packages (from category_encoders) (1.2)

Requirement already satisfied: pytz>=2017.3 in /usr/local/lib/python3.7/dist-packages (from pandas>=0.21.1->category_encoders) (2018.3)

Requirement already satisfied: python-dateutil>=2.7.3 in /usr/local/lib/python3.7/dist-packages (from pandas>=0.21.1->category_encoders) (2.7.5)

Requirement already satisfied: six in /usr/local/lib/python3.7/dist-packages (from patsy>=0.5.1->category_encoders) (1.12.0)

Requirement already satisfied: threadpoolctl>=2.0.0 in /usr/local/lib/python3.7/dist-packages (from scikit-learn>=0.20.0->category_encoders) (2.0.0)

Requirement already satisfied: joblib>=0.11 in /usr/local/lib/python3.7/dist-packages (from scikit-learn>=0.20.0->category_encoders) (0.14.0)

```
[ ] from category_encoders import TargetEncoder
[ ] te = TargetEncoder()
[ ] data[col] = te.fit_transform(data[col], data['Loan_Status'])

[ ] col='Property_Area'
[ ] data[col].value_counts()

0.768240    233
0.658416    202
0.614525    179
Name: Property_Area, dtype: int64
```

<>

☰ Education

63/63

EDA_FE.ipynb - Colaboratory x markov x | sklearn.impute.SimpleImputer x + colab.research.google.com/drive/1_P-nKswaKnx5v77lOx3pCxvJve73MvnN#scrollTo=_K4bmriUIDco

+ Code + Text ✓ RAM Disk

Requirement already satisfied: python-dateutil>=2.7.3 in /usr/local/lib/python3.7/dist-packages (from pandas>=0.21.1)

Requirement already satisfied: six in /usr/local/lib/python3.7/dist-packages (from patsy>=0.5.1->category_encoders)

Requirement already satisfied: threadpoolctl>=2.0.0 in /usr/local/lib/python3.7/dist-packages (from scikit-learn>=0.20.0->ca

[x]

```
[ ] from category_encoders import TargetEncoder
```

te = TargetEncoder()
data[col] = te.fit_transform data[col], data['Loan_Status'])

1 vs 0

```
[ ] col='Property_Area'  
data[col].value_counts()
```

{
0.768240 233
0.658416 202
0.614525 179
Name: Property_Area, dtype: int64

<> ▾ Education

```
[ ] col='Education'  
data[col].value_counts()
```

EDA_FE.ipynb - Colaboratory x markov x | sklearn.impute.SimpleImputer x + colab.research.google.com/drive/1_P-nKswaKnx5v77lOx3pCxvJve73MvnN#scrollTo=_K4bmriUIDco

+ Code + Text RAM Disk ✓

```
[ ] from category_encoders import TargetEncoder
```

```
{x} te = TargetEncoder()
data[col] = te.fit_transform(data[col], data['Loan_Status'])
```

```
[ ] col='Property_Area'
data[col].value_counts()
```

✓ { 0.768240 233 → SU
0.658416 202 → U
0.614525 179 → P
Name: Property_Area, dtype: int64

▼ Education

```
<> [ ] col='Education'
data[col].value_counts()
```

```
>- Graduate 480
Not Graduate 134
Name: Education, dtype: int64
```

EDA_FE.ipynb - Colaboratory x markov x | sklearn.impute.SimpleImputer x + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=_K4bmriUIDco

+ Code + Text RAM Disk

Self Employed

```
{x} [ ] col='Self_Employed'  
      data[col].value_counts()
```

No	500
Yes	82
Other	32

```
Name: Self_Employed, dtype: int64
```

```
[ ] from category_encoders import TargetEncoder  
  
te = TargetEncoder()  
data[col] = te.fit_transform(data[col], data['Loan_Status'])  
data[col].value_counts()
```

✓ 0.686000	500 No
✓ 0.682927	82 Yes
✓ 0.718750	32 Other

```
<> Name: Self_Employed, dtype: int64
```

```
[ ] s = (data.dtypes == 'object')  
object_cols = list(s[s].index)  
object_cols
```

EDA_FE.ipynb - Colaboratory x markov x | sklearn.impute.SimpleImputer x +

colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=_K4bmriUIDco

+ Code + Text RAM Disk ✓

Self Employed

{x} [] col='Self_Employed'
data[col].value_counts()

No	500
Yes	82
Other	32

Name: Self_Employed, dtype: int64

[] from category_encoders import TargetEncoder

te = TargetEncoder()
data[col] = te.fit_transform(data[col], data['Loan_Status'])
data[col].value_counts()

0.686000	500
0.682927	82
0.718750	32

Name: Self_Employed, dtype: int64

[] s = (data.dtypes == 'object')
object_cols = list(s[s].index)
object_cols

File Home Recent Help

67 / 67

EDA_FE.ipynb - Colaboratory x markov x | sklearn.impute.SimpleImputer x +

colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=_K4bmriUIDco

+ Code + Text RAM Disk

Self Employed

{x}

```
[ ] col='Self_Employed'  
      data[col].value_counts()
```

No 500
Yes 82
Other 32
Name: Self_Employed, dtype: int64

[] from category_encoders import TargetEncoder

```
te = TargetEncoder()  
data[col] = te.fit_transform(data[col], data['Loan_Status'])  
data[col].value_counts()
```

0.686000 500
0.682927 82
0.718750 32
Name: Self_Employed, dtype: int64

[] s = (data.dtypes == 'object')
object_cols = list(s[s].index)
object_cols

EDA_FE.ipynb - Colaboratory x markov x | sklearn.impute.SimpleImputer x +

colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=_K4bmriUIDco

+ Code + Text

NO MORE NON NUMERIC COLS.

[]

{x} [] data.dtypes

Gender	int64
Married	int64
Dependents	float64
Education	int64
Self_Employed	float64
ApplicantIncome	int64
CoapplicantIncome	float64
LoanAmount	float64
Loan_Amount_Term	float64
Credit_History	float64
Property_Area	float64
Loan_Status	int64
TotalIncome	float64
Loan_Amount_per_year	float64
EMI	float64
Able_to_pay_EMI	int64
dtype: object	

<>

Correlation Coefficients

RAM Disk

Update

EDA_FE.ipynb - Colaboratory x markov x | sklearn.impute.SimpleImputer x + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=_K4bmriUIDco

+ Code + Text

RAM Disk ✓

Correlation Coefficients

{x}

[] #PCC

```
plt.figure(figsize=(15, 15))
sns.heatmap(data.corr(method='pearson'), square=True, annot=True)
```

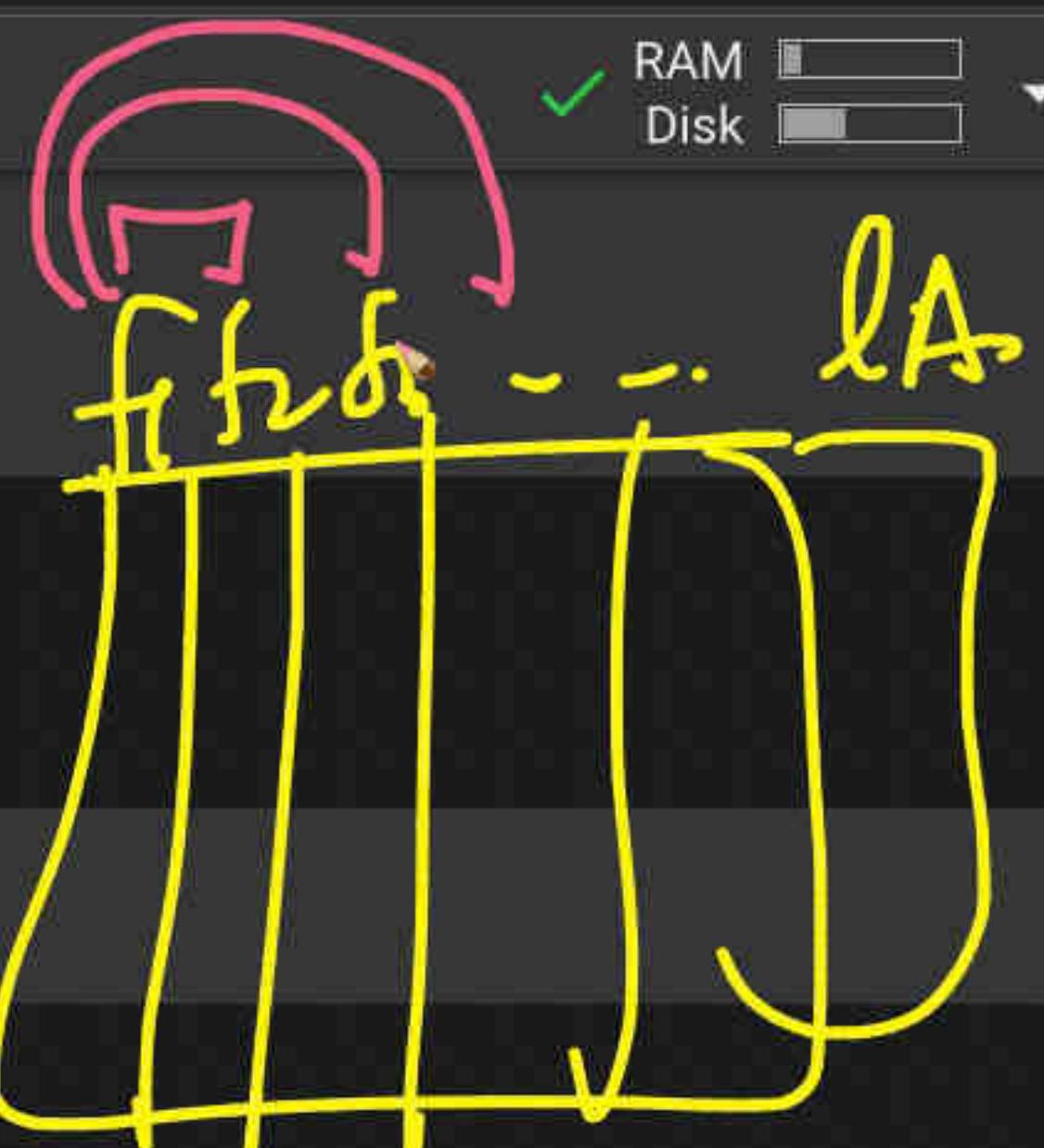
Annotations: A red bracket highlights the 'method' parameter in the sns.heatmap call, and a pink curly brace encloses the entire plt.figure and sns.heatmap lines.

[] #SRCC

```
plt.figure(figsize=(15, 15))
sns.heatmap(data.corr(method='spearman'), square=True, annot=True)
```

[] plt.scatter(data['Credit_History'], data['Loan_Status'])
plt.show()
#sometimes scatter plots can be misleading do to catgeorical nature of the data

[] sns.countplot(data = data, x = 'Credit_History', hue = 'Loan_Status')



EDA_FE.ipynb - Colaboratory x markov x | sklearn.impute.SimpleImputer x +

colab.research.google.com/drive/1_P-nKswaKnx5v77lOx3pCxvJve73MvnN#scrollTo=_K4bmriUIDco

+ Code + Text RAM Disk

RAM Disk

Correlation Coefficients

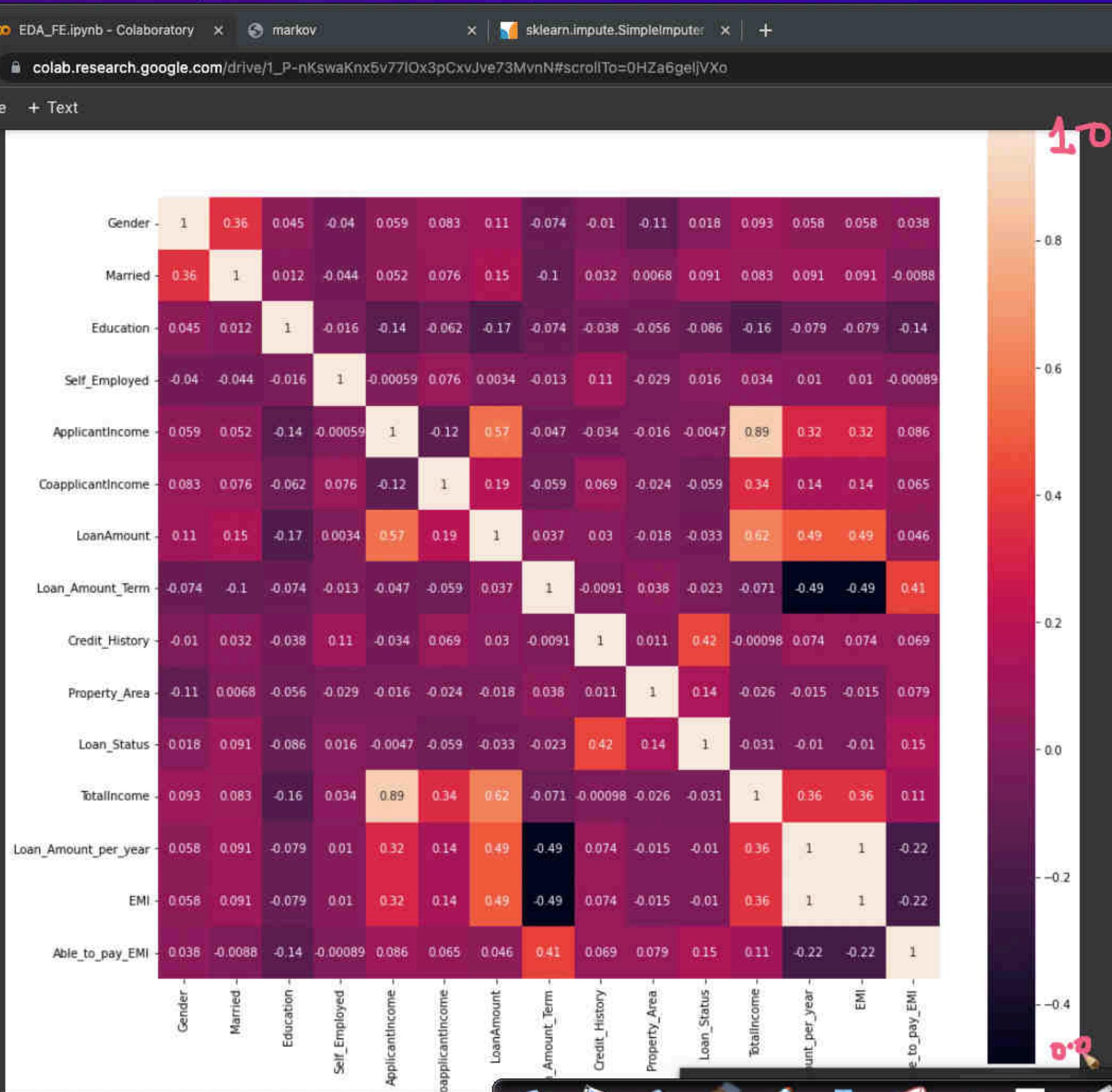
```
{x} [ ] #PCC
    plt.figure(figsize=(15, 15))
    sns.heatmap(data.corr(method='pearson'), square=True, annot=True)
```

=====

```
[ ] #SRCC
    plt.figure(figsize=(15, 15))
    sns.heatmap(data.corr(method='spearman'), square=True, annot=True)
```

```
[ ] plt.scatter(data['Credit_History'], data['Loan_Status'])
    plt.show()
    #sometimes scatter plots can be misleading do to catgeorical nature of the data
```

```
[ ] sns.countplot(data = data, x = 'Credit_History', hue = 'Loan_Status')
```



EDA_FE.ipynb - Colaboratory x markov x | sklearn.impute.SimpleImputer x +

colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=OHZa6geljVXo

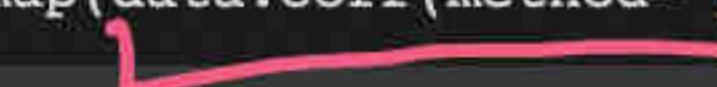
+ Code + Text RAM Disk

RAM Disk

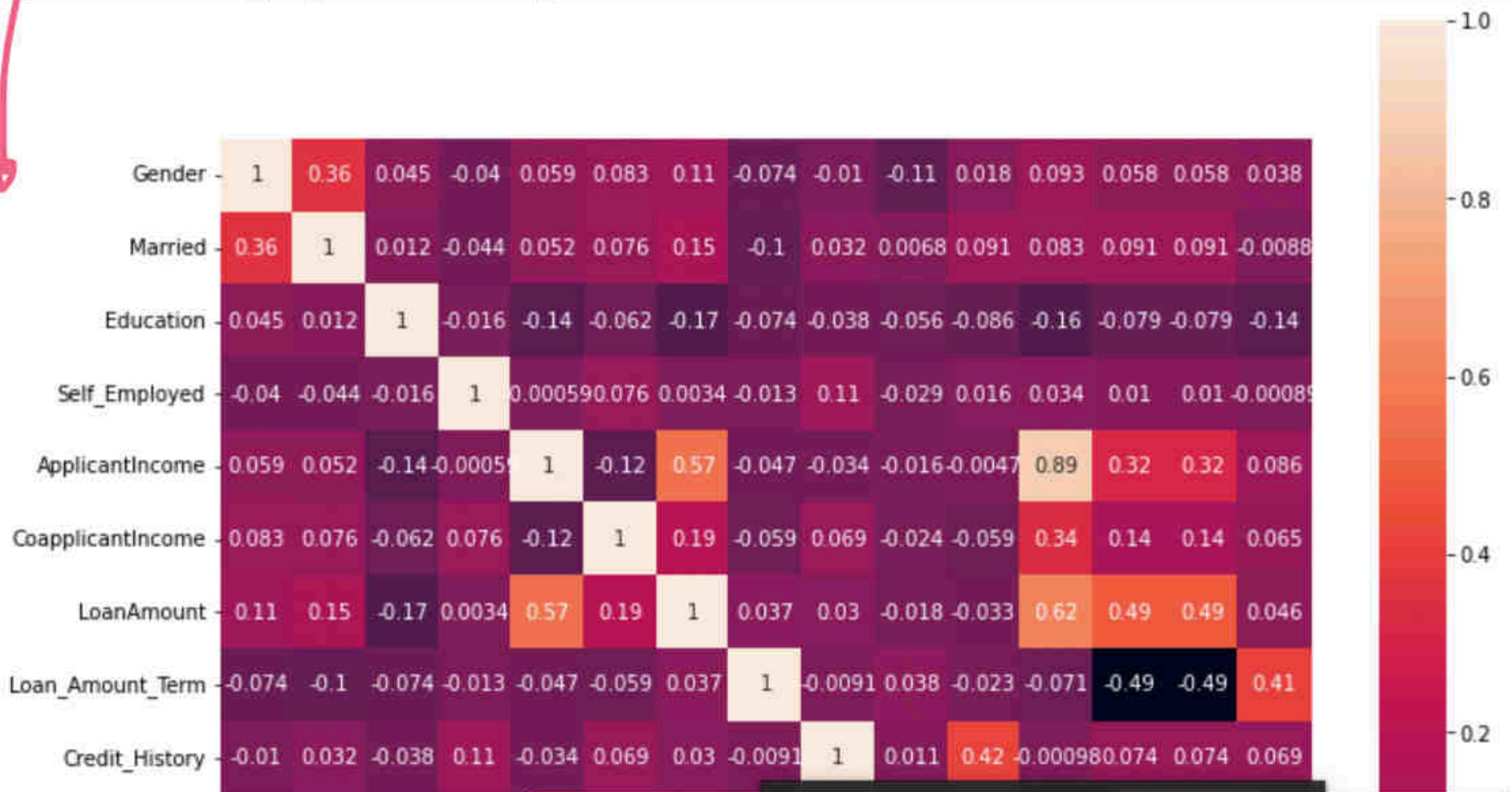
Correlation Coefficients

```
#PCC
plt.figure(figsize=(12, 12))
sns.heatmap(data.corr(method='pearson'), square=True, annot=True)
```



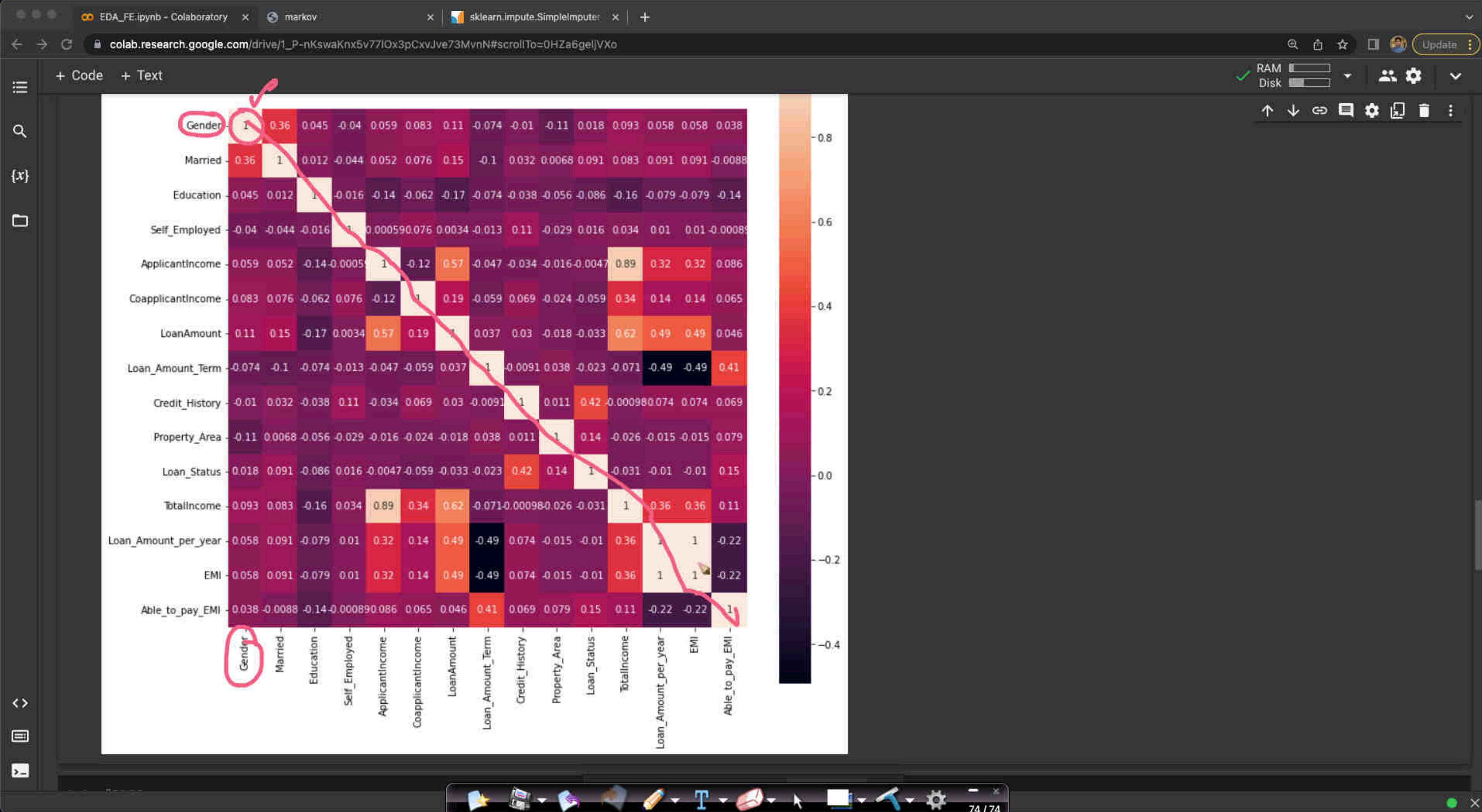


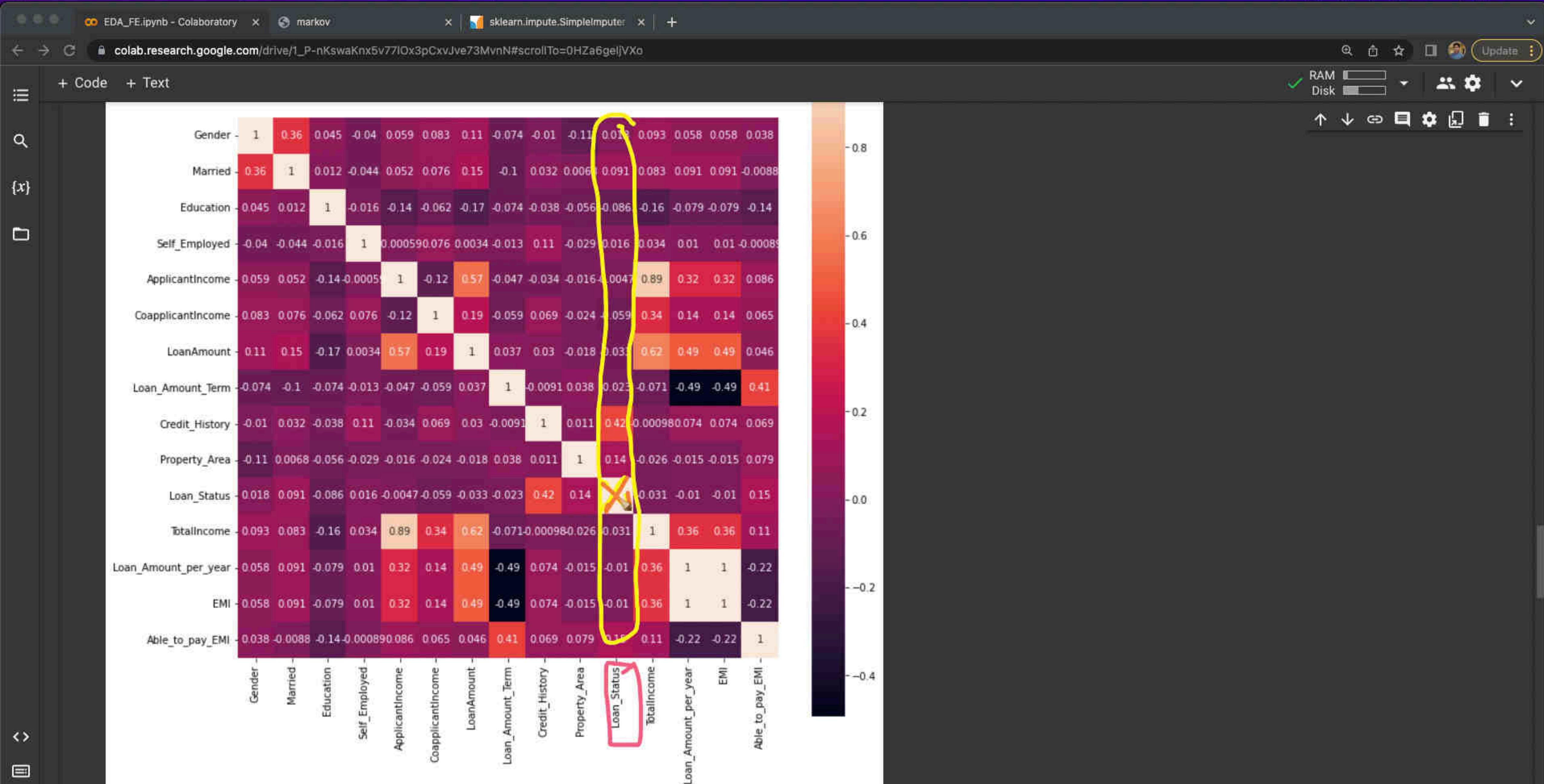
`<matplotlib.axes._subplots.AxesSubplot at 0x7fac6315f450>`

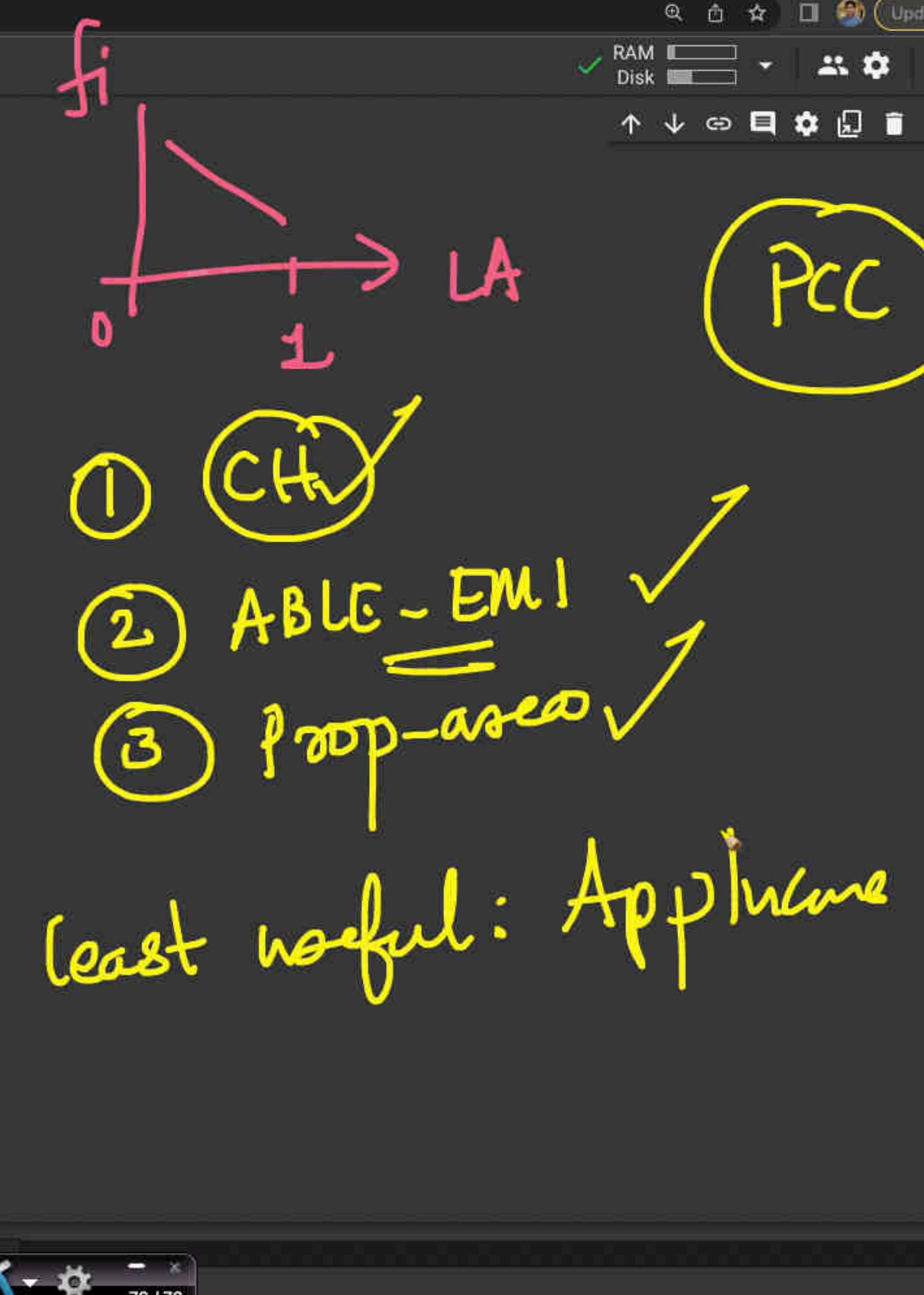
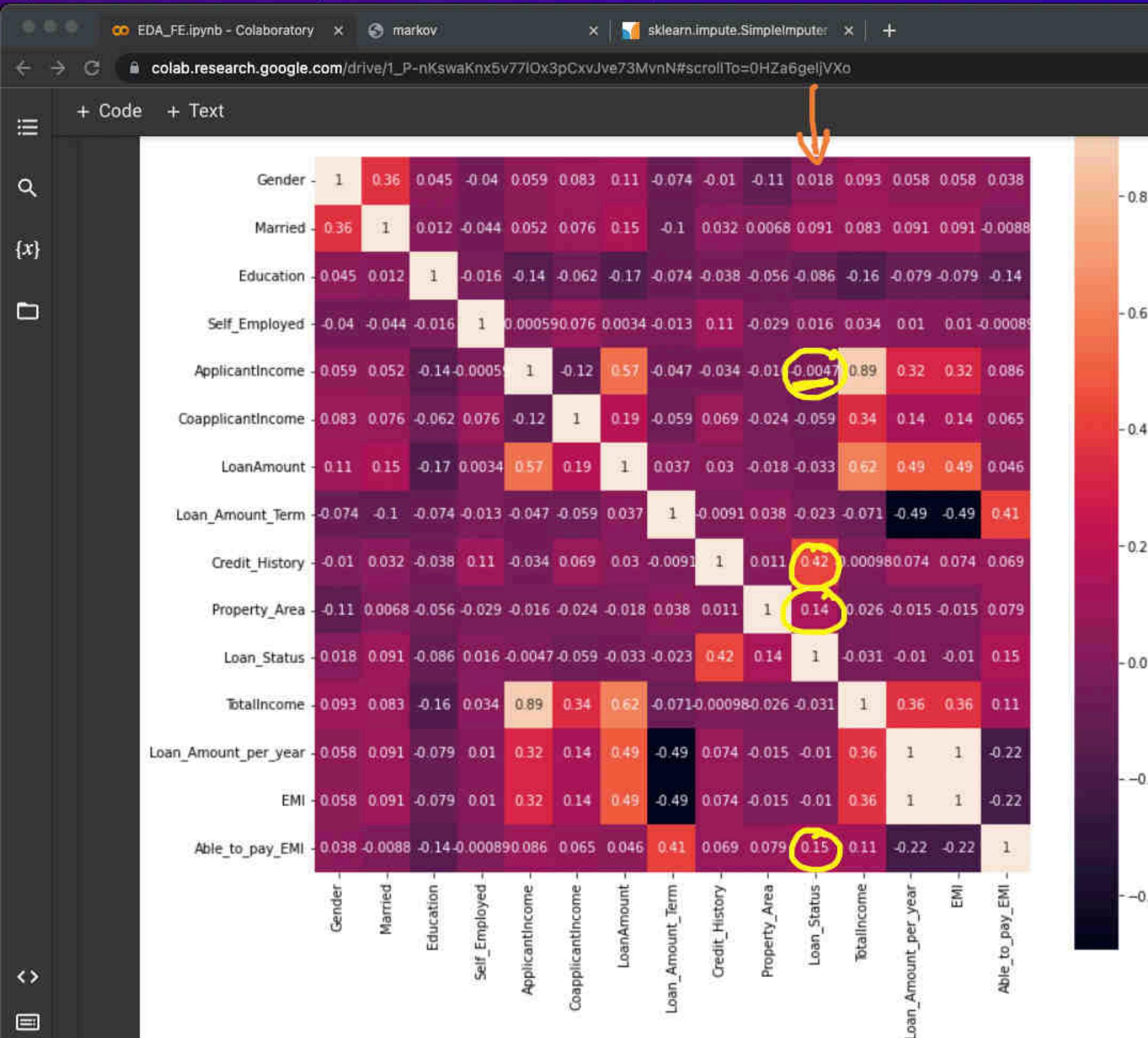


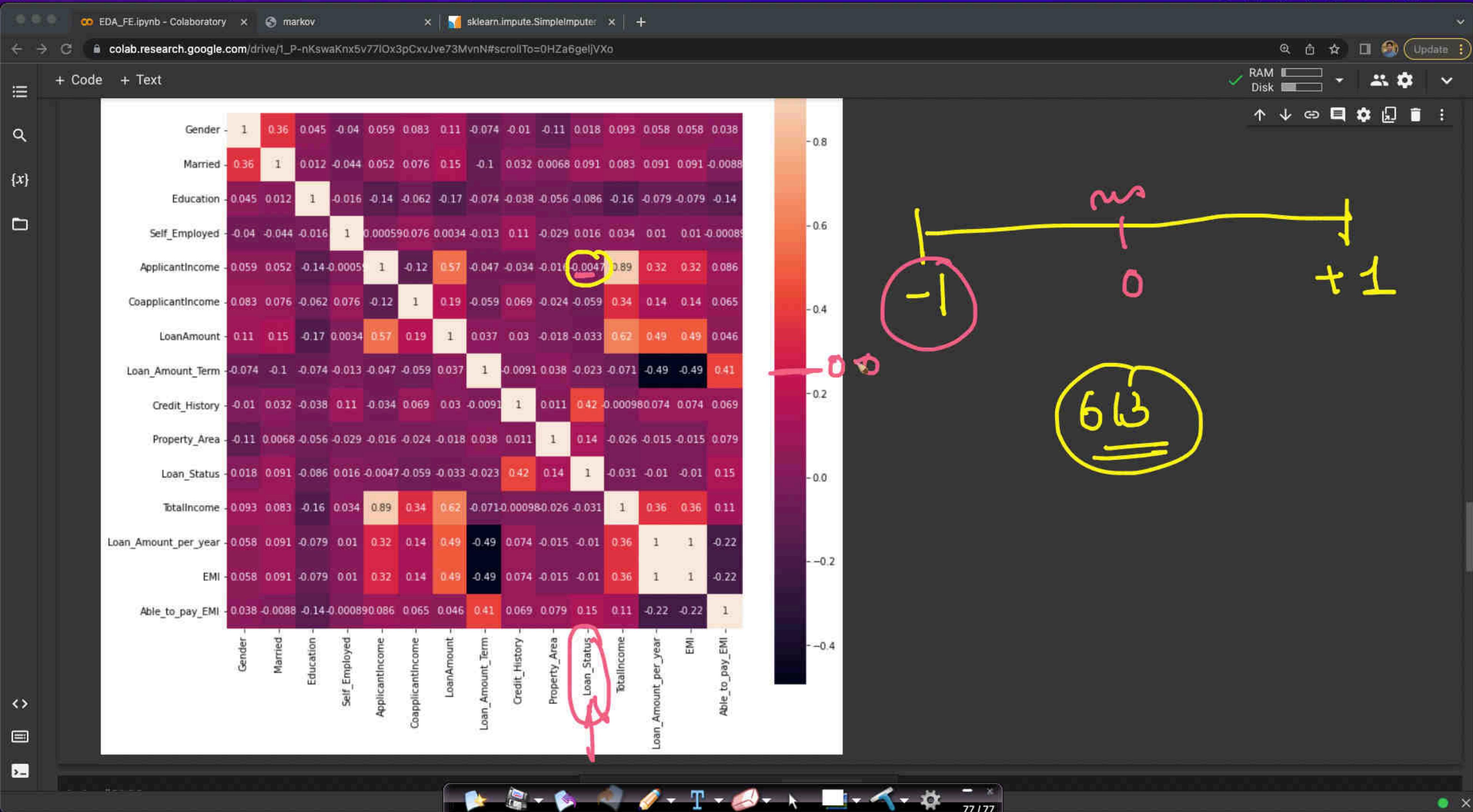
	Gender	Married	Education	Self_Employed	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount_Term	Credit_History						
Gender	1	0.36	0.045	-0.04	0.059	0.083	0.11	-0.074	-0.01	-0.11	0.018	0.093	0.058	0.058	0.038
Married	0.36	1	0.012	-0.044	0.052	0.076	0.15	-0.1	0.032	0.0068	0.091	0.083	0.091	0.091	-0.0088
Education	0.045	0.012	1	-0.016	-0.14	-0.062	-0.17	-0.074	-0.038	-0.056	-0.086	-0.16	-0.079	-0.079	-0.14
Self_Employed	-0.04	-0.044	-0.016	1	0.000590	0.076	0.0034	-0.013	0.11	-0.029	0.016	0.034	0.01	0.01	-0.00089
ApplicantIncome	0.059	0.052	-0.14	-0.000590	1	-0.12	0.57	-0.047	-0.034	-0.016	-0.0047	0.89	0.32	0.32	0.086
CoapplicantIncome	0.083	0.076	-0.062	0.076	-0.12	1	0.19	-0.059	0.069	-0.024	-0.059	0.34	0.14	0.14	0.065
LoanAmount	0.11	0.15	-0.17	0.0034	0.57	0.19	1	0.037	0.03	-0.018	-0.033	0.62	0.49	0.49	0.046
Loan_Amount_Term	-0.074	-0.1	-0.074	-0.013	-0.047	-0.059	0.037	1	-0.0091	0.038	-0.023	-0.071	-0.49	-0.49	0.41
Credit_History	-0.01	0.032	-0.038	0.11	-0.034	0.069	0.03	-0.0091	1	0.011	0.42	-0.000980	0.074	0.074	0.069

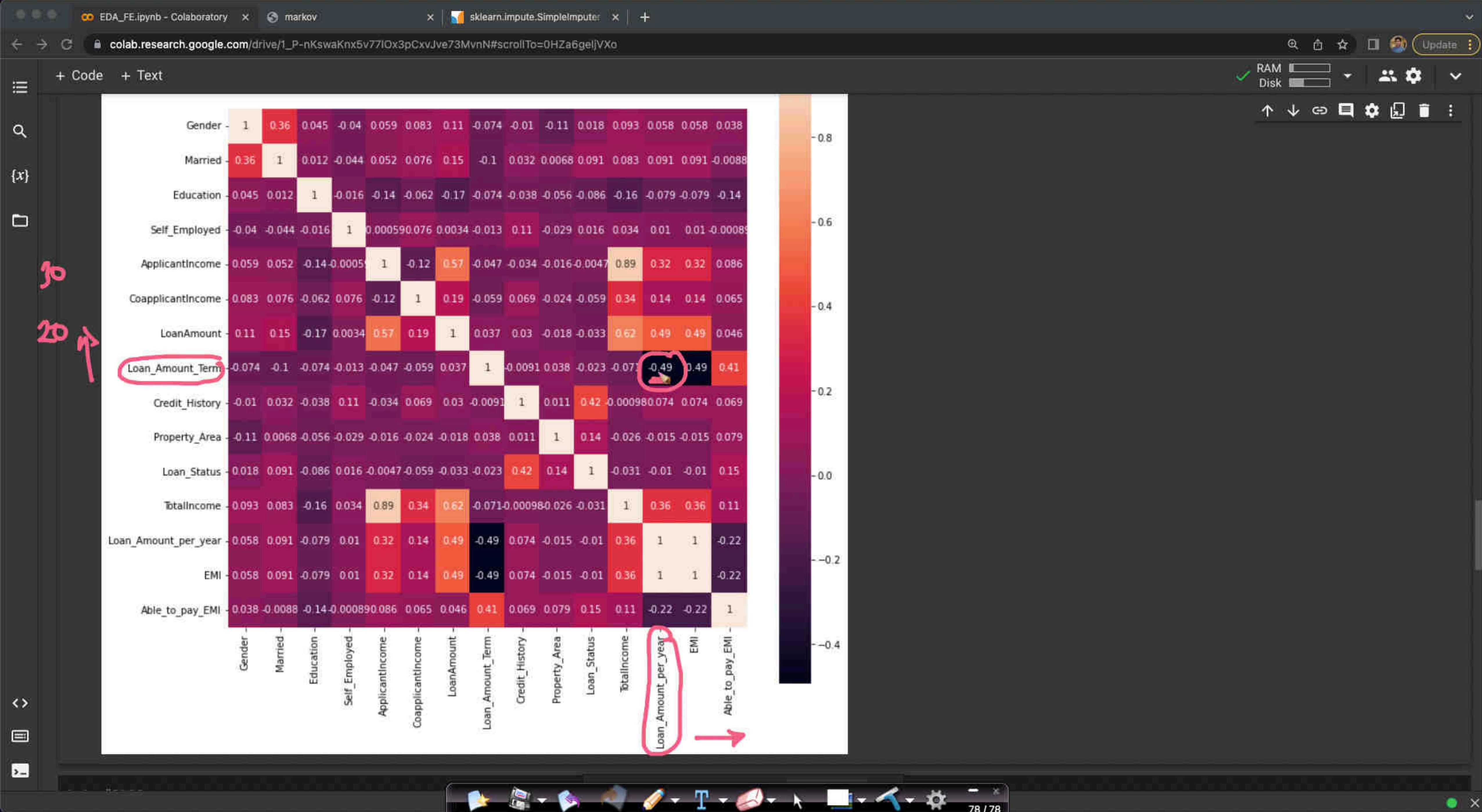
73 / 73

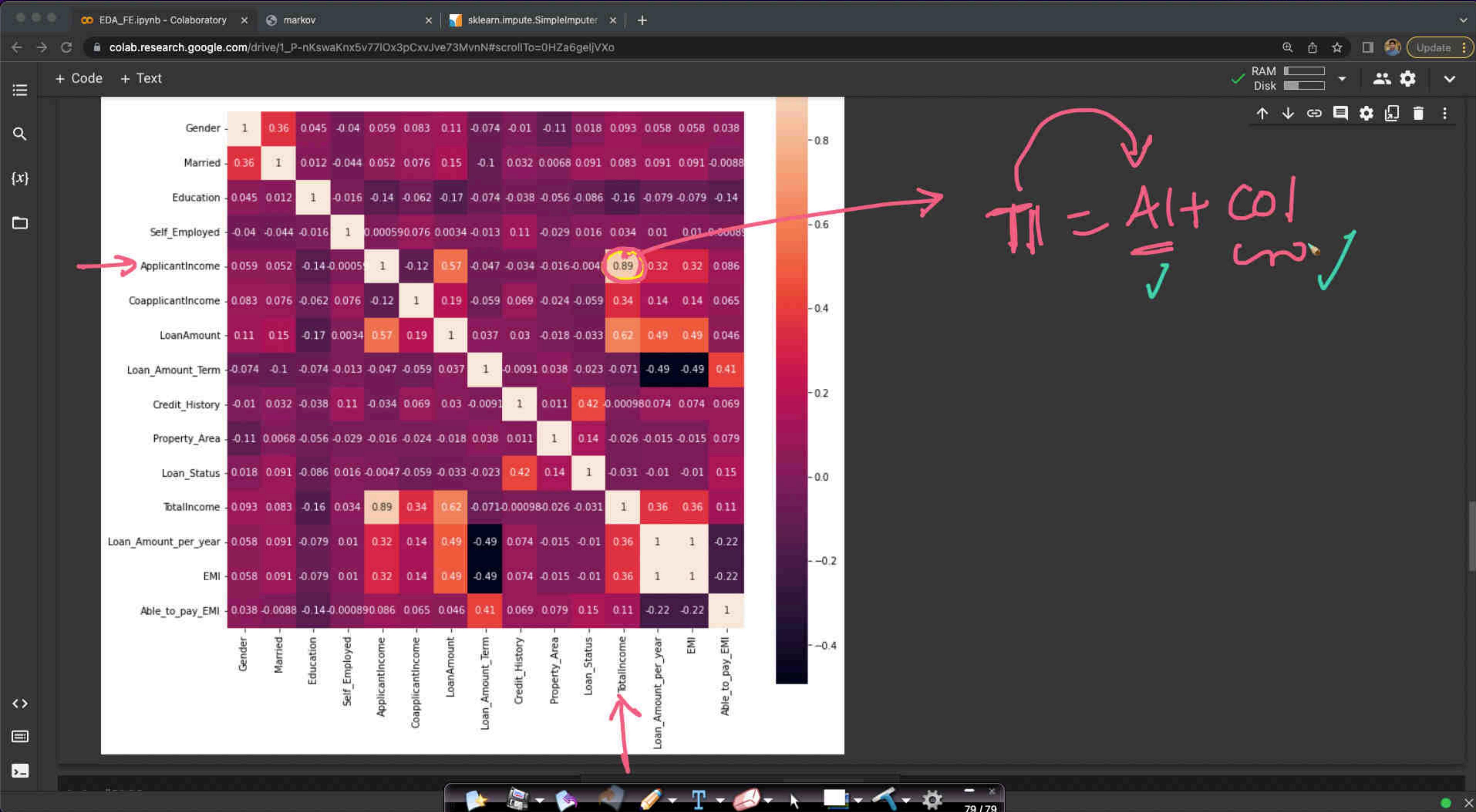


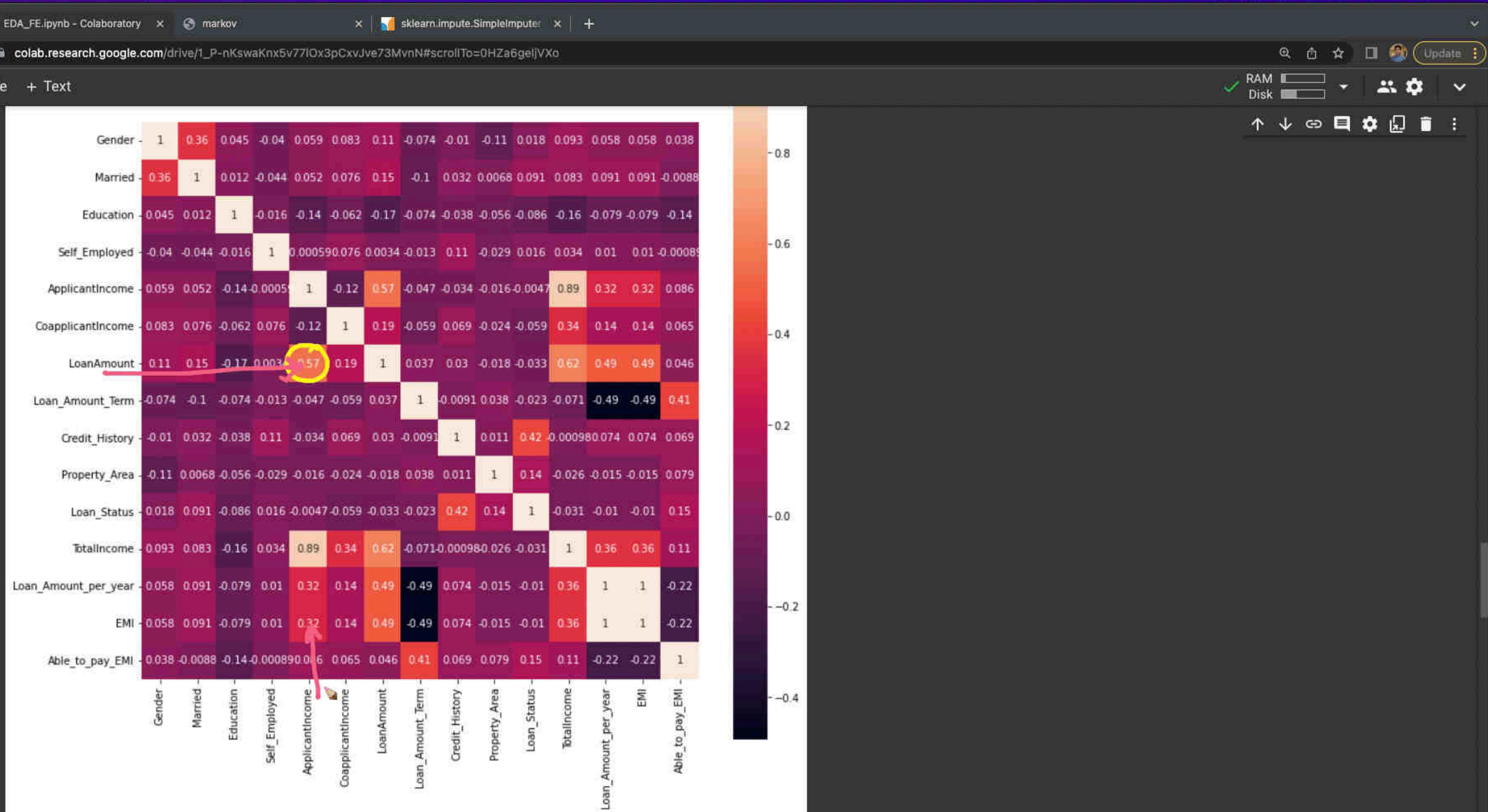


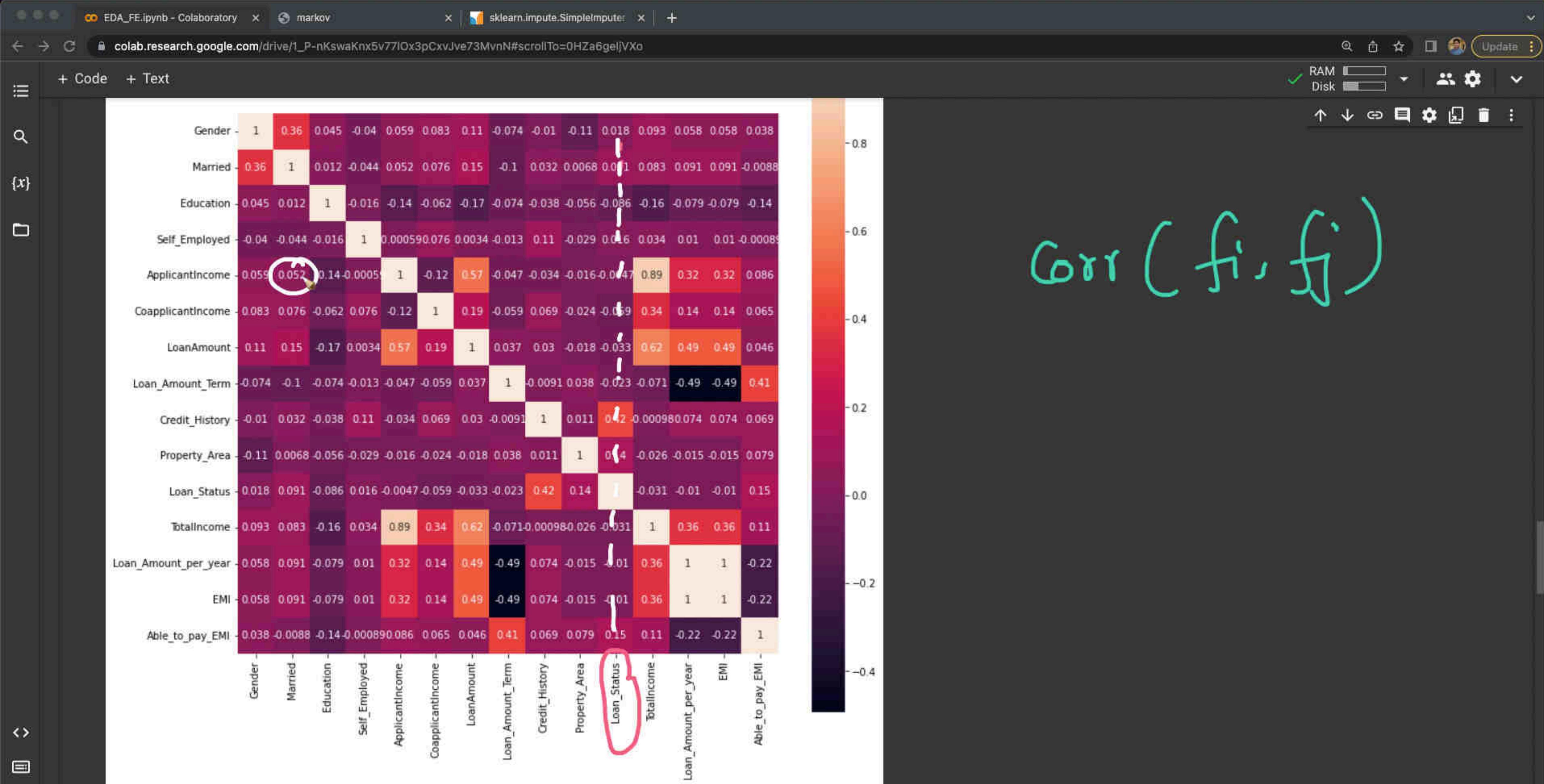




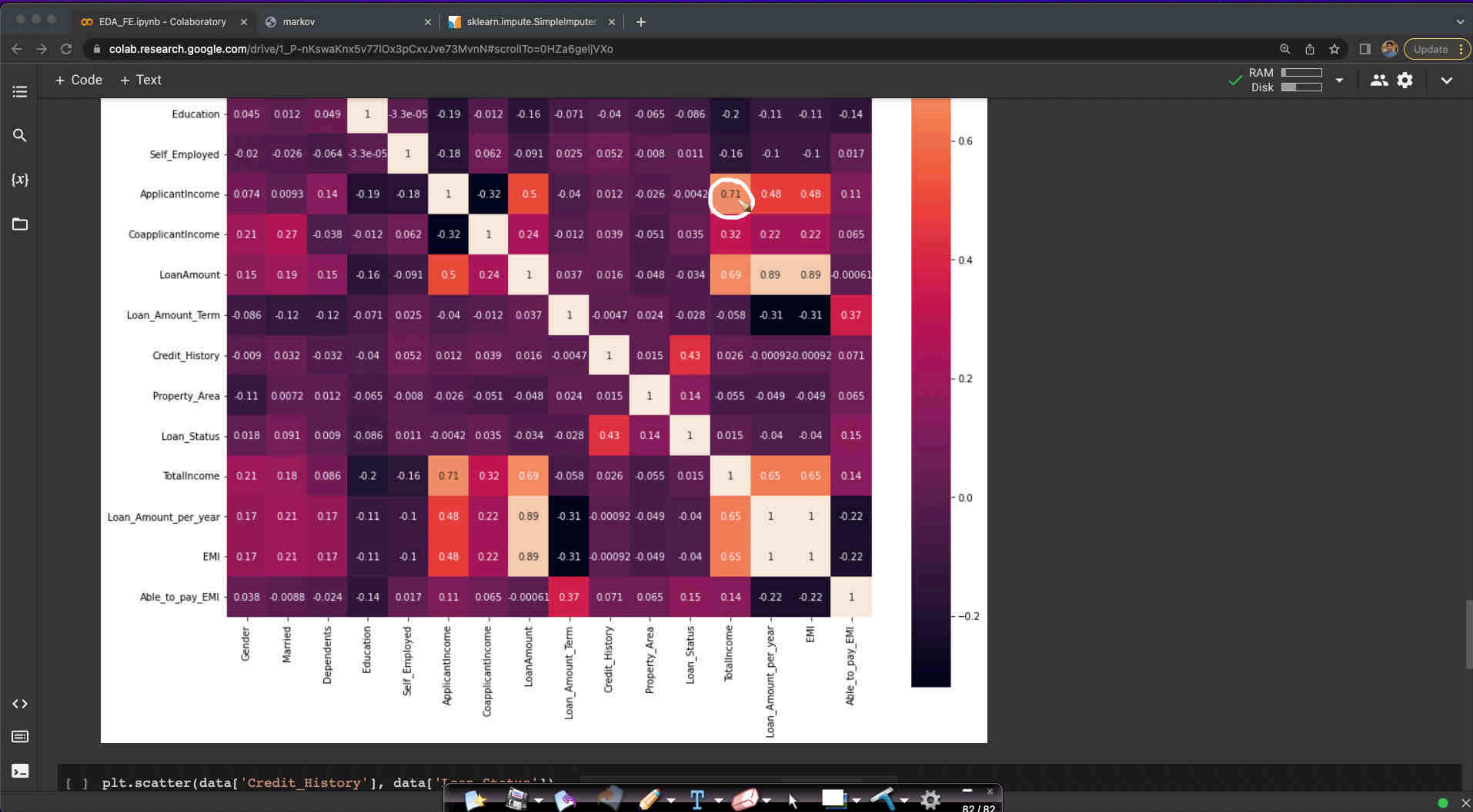


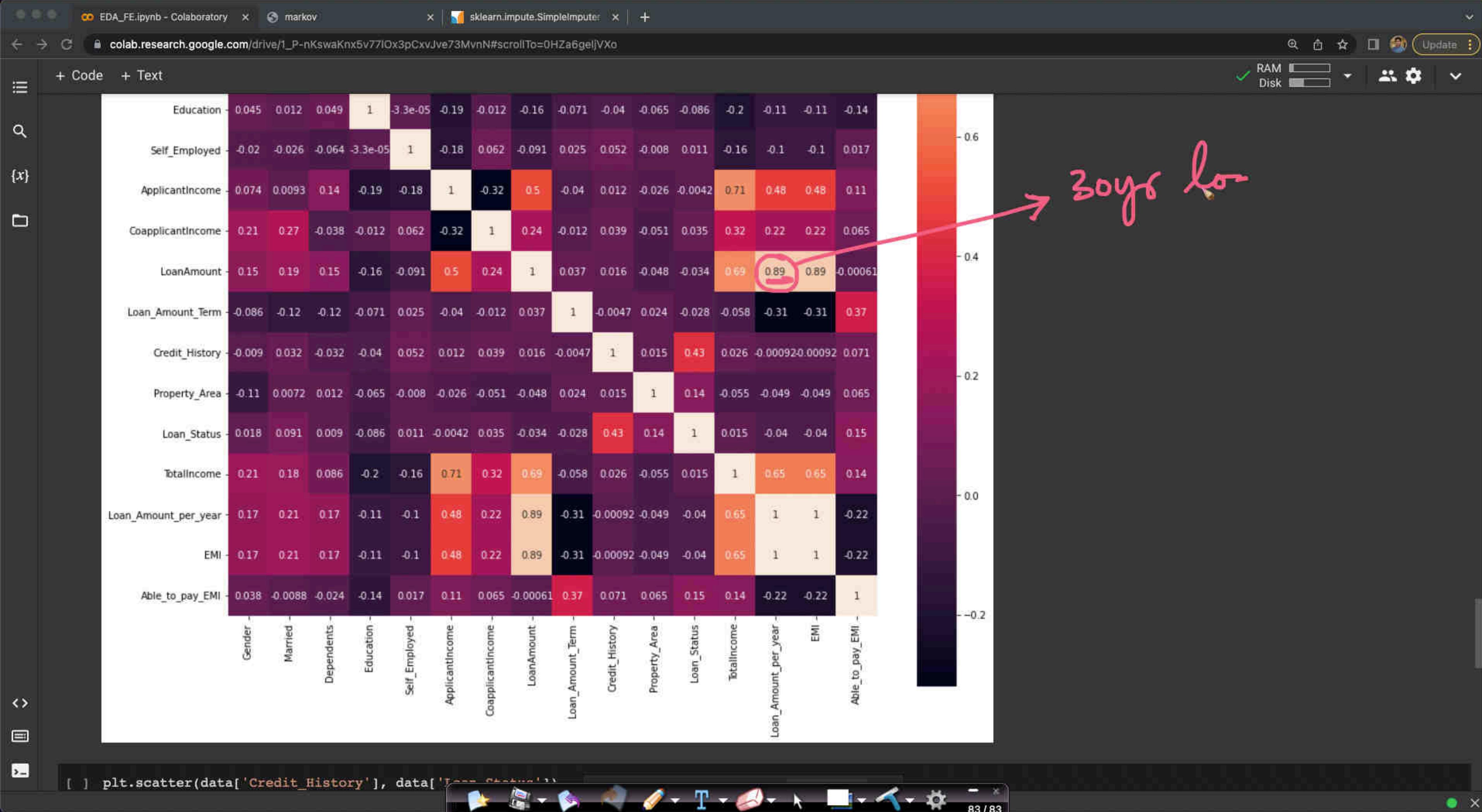


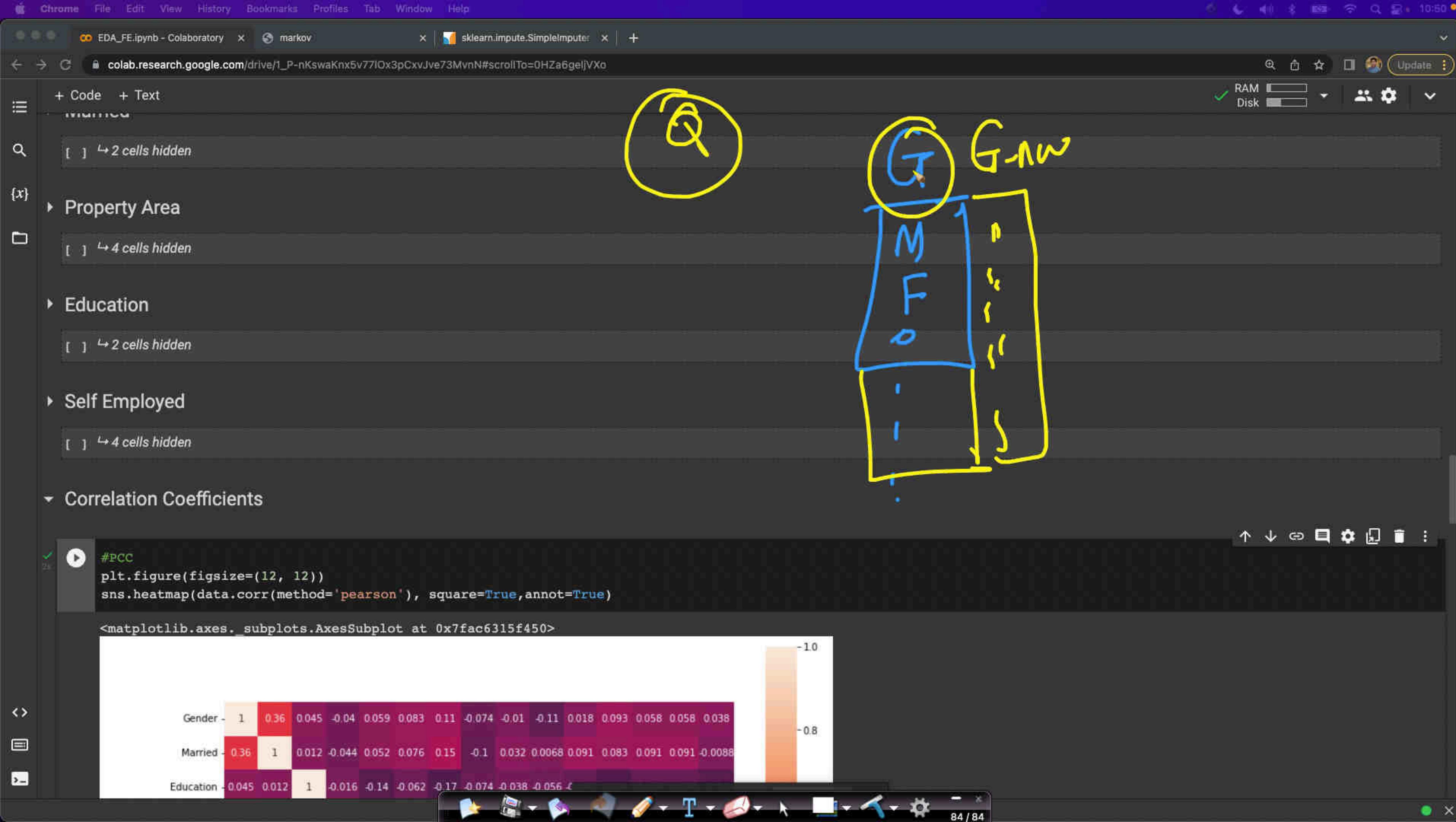


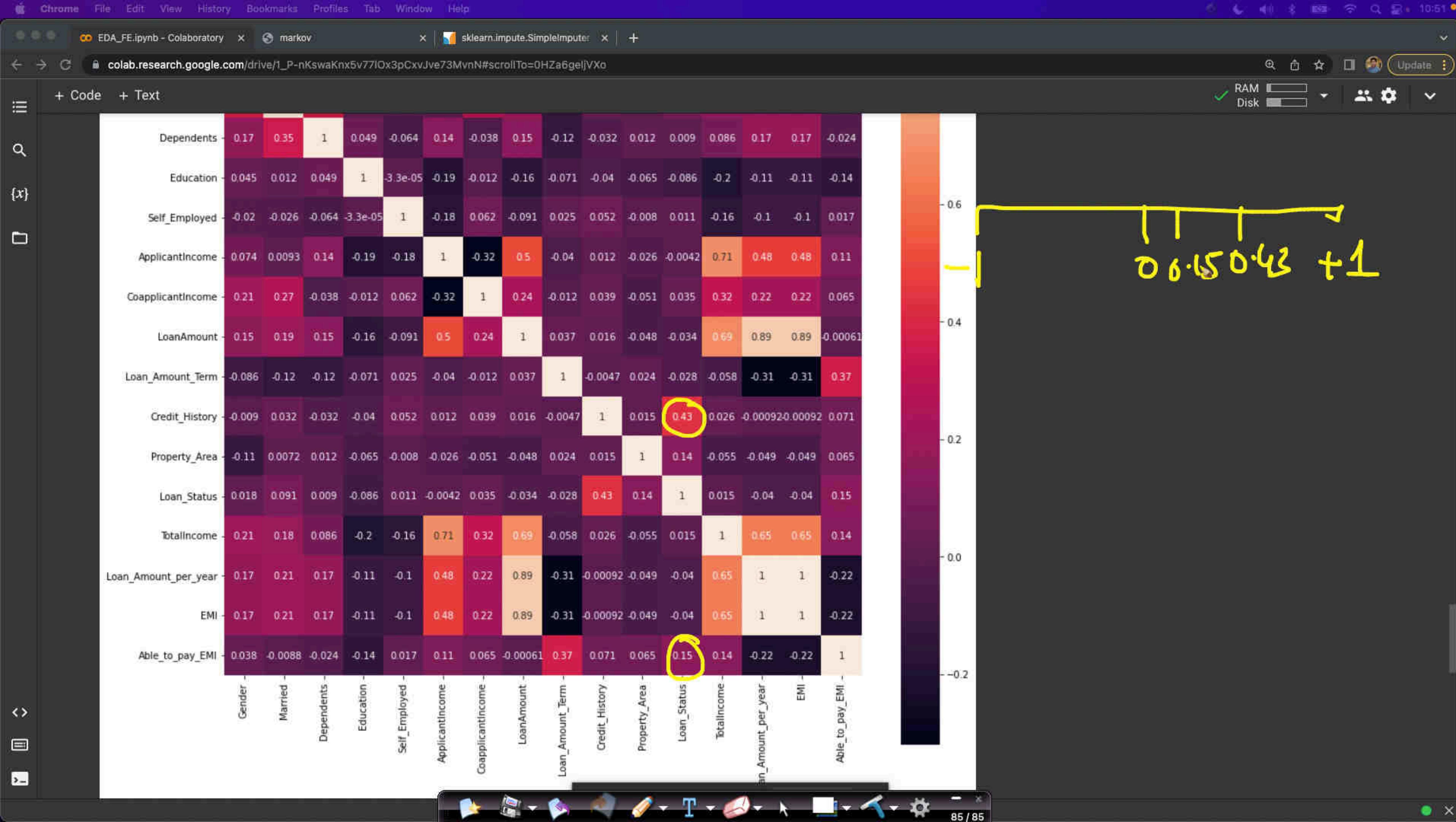


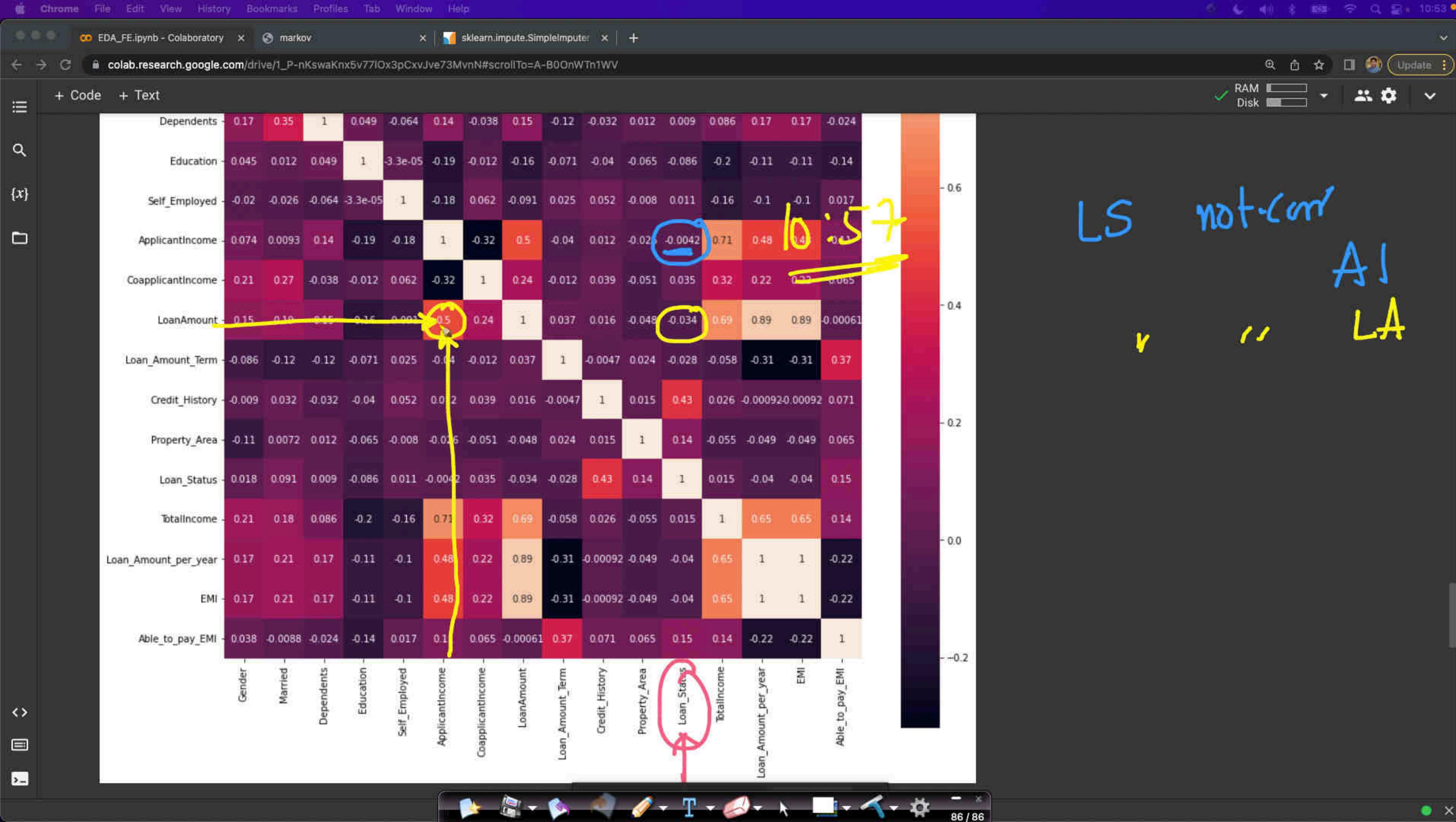
Corr (f_i, f_j)

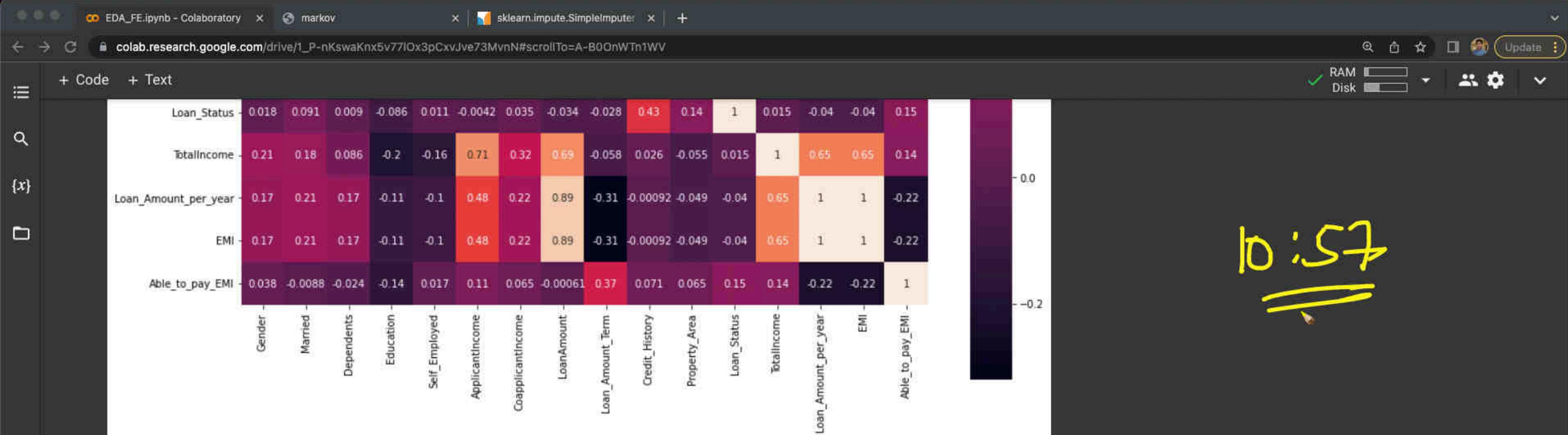




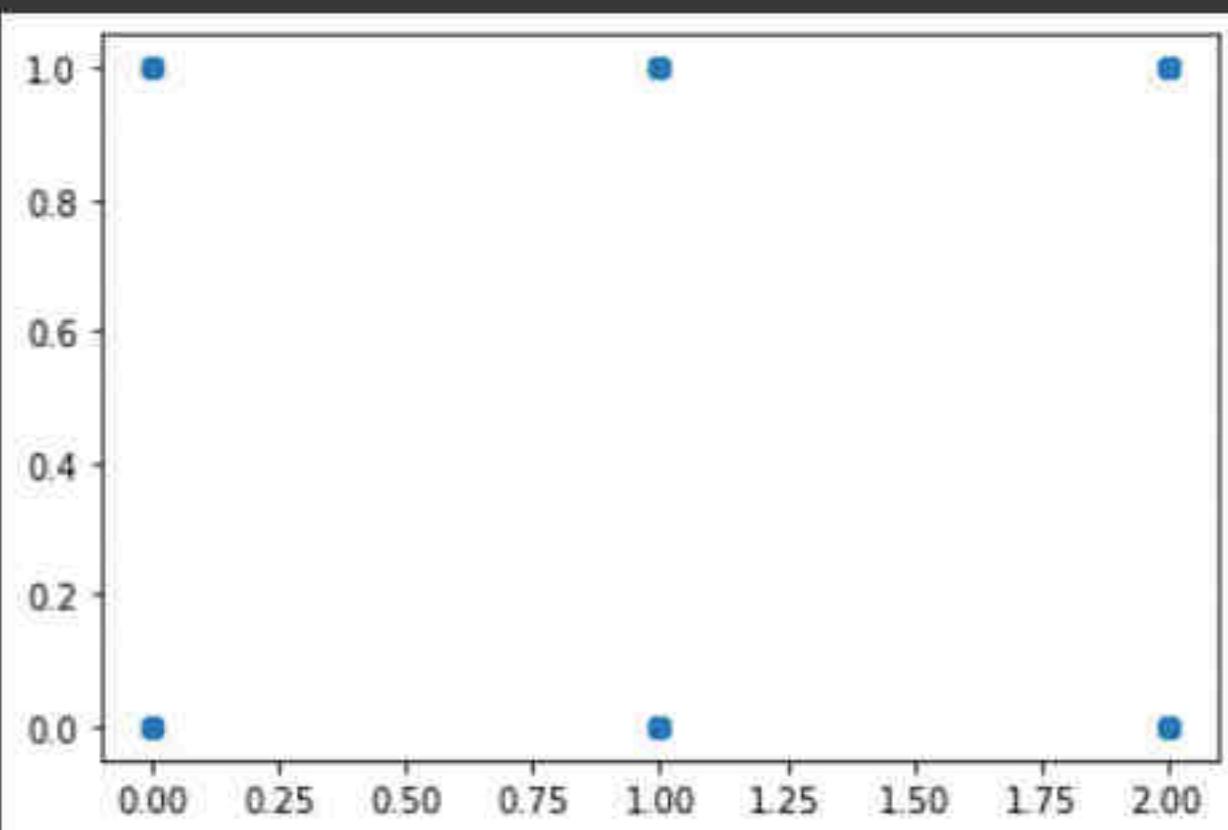








```
[ ] plt.scatter(data['Credit_History'], data['Loan_Status'])
plt.show()
#sometimes scatter plots can be misleading do to catgeorical nature of the data
```





Multi-class

LS
 $y_i \rightarrow \underbrace{S}_{\sim} \underbrace{M}_{\sim} \underbrace{L}_{\sim} \underbrace{k}_{\sim}$

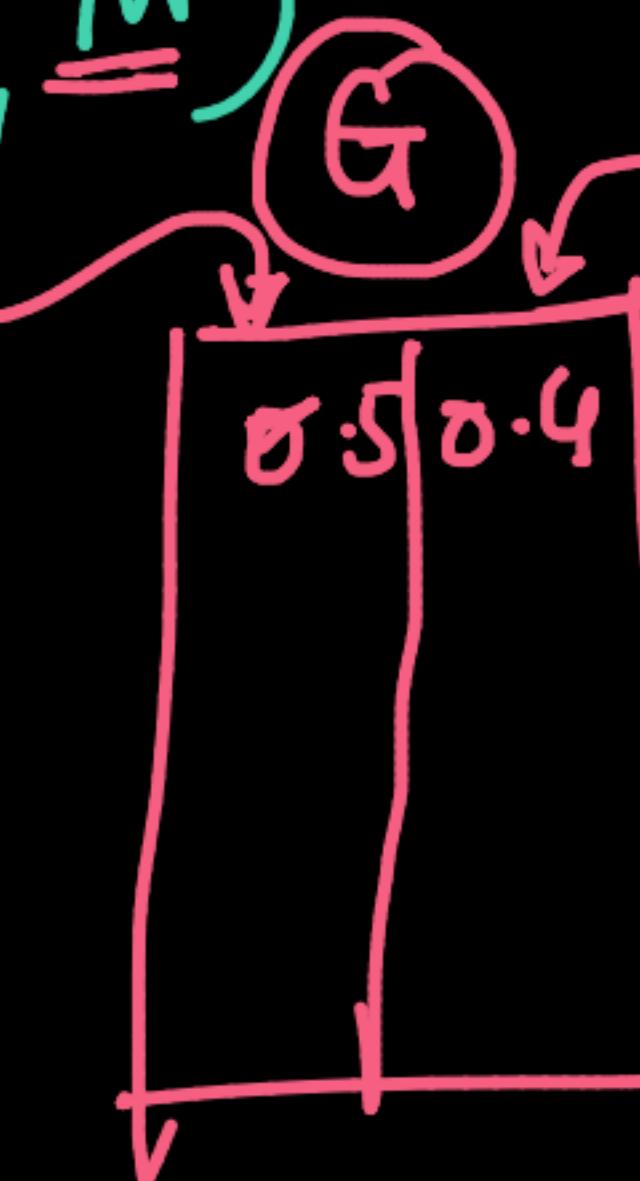
$\{ G = M | F | o$



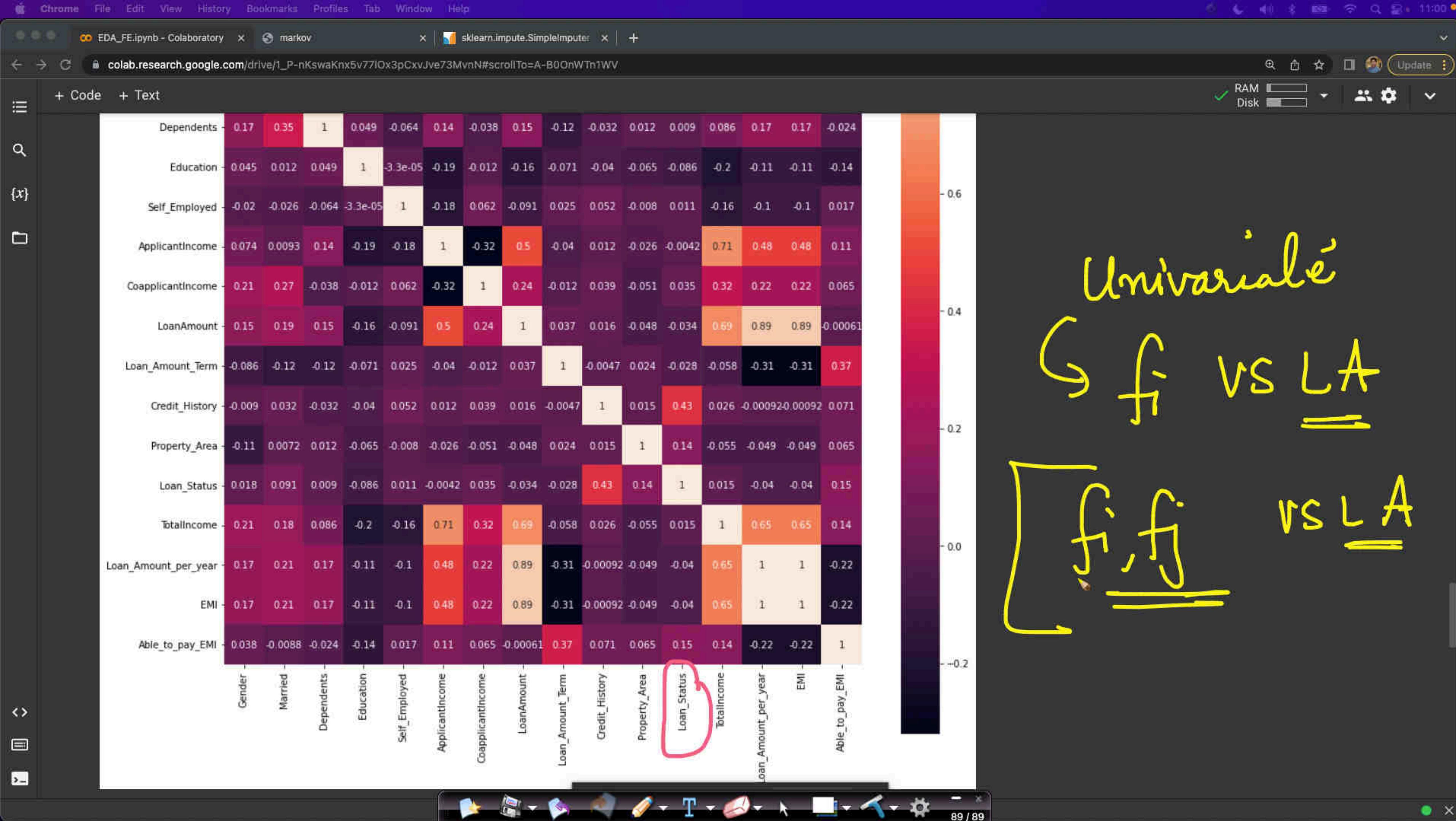
(later)
 feature
 imp
 (ML)

$P(y_i = S | \underbrace{G = M}_{\sim})$

$P(y_i = M | \underbrace{G = M}_{\sim})$



$\chi^2 - \text{est}$



EDA_FE.ipynb - Colaboratory x markov x sklearn.impute.SimpleImputer x + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=A-B0OnWTn1WV

+ Code + Text

RAM Disk

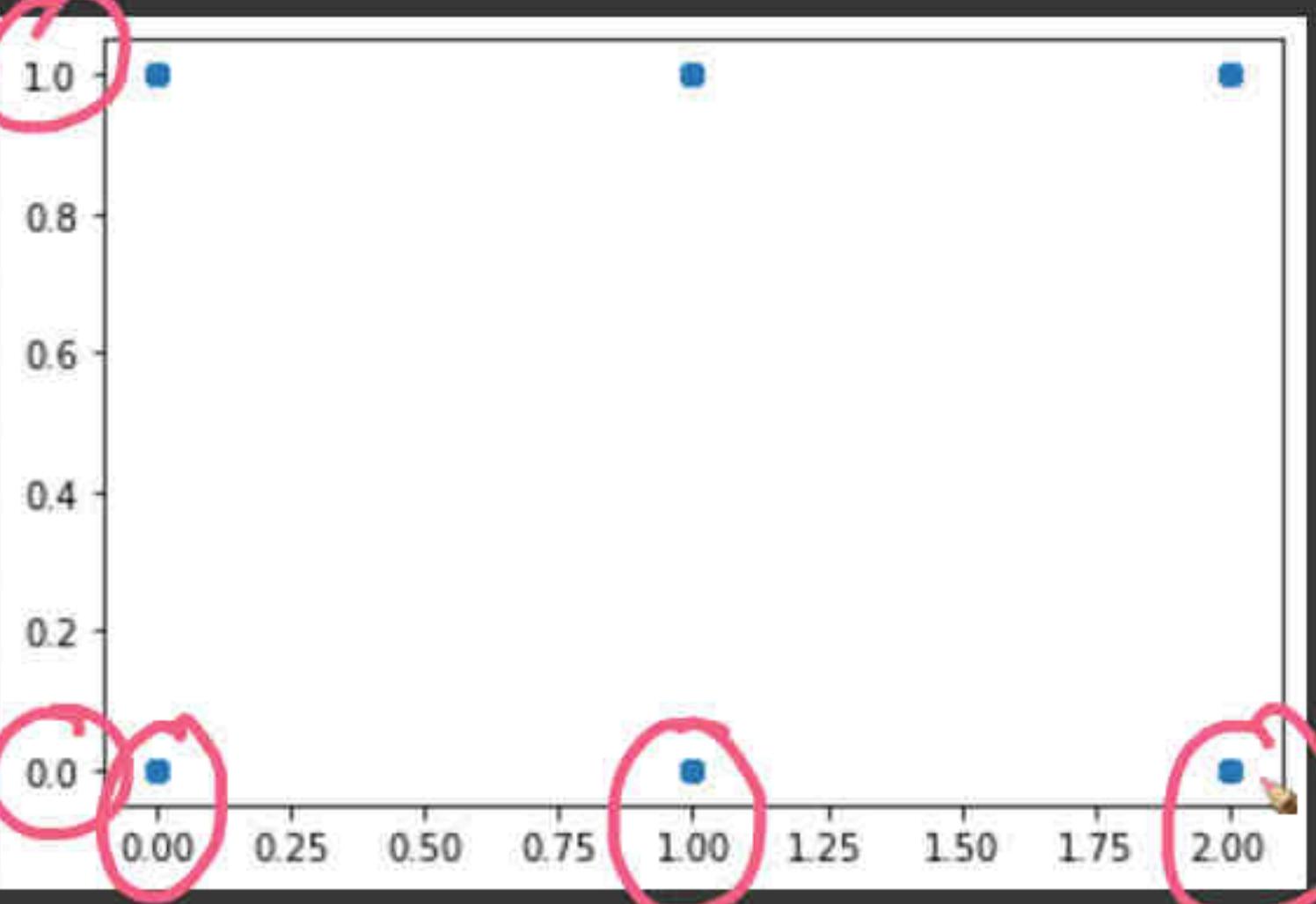
Dependents Education Self_Employed Applicant Coapplicant LoanAmount Credit_History Property_Loan Total_Amount Loan_Amount_perCapita Able_to_pay

{x}

LA

CH

```
plt.scatter(data['Credit_History'], data['Loan_Status'])
plt.show()
#sometimes scatter plots can be misleading do to catgeorical nature of the data
```



```
sns.countplot(data = data, x = 'Credit_History', hue = 'Loan_Status')
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f3e2eb7d950>
```

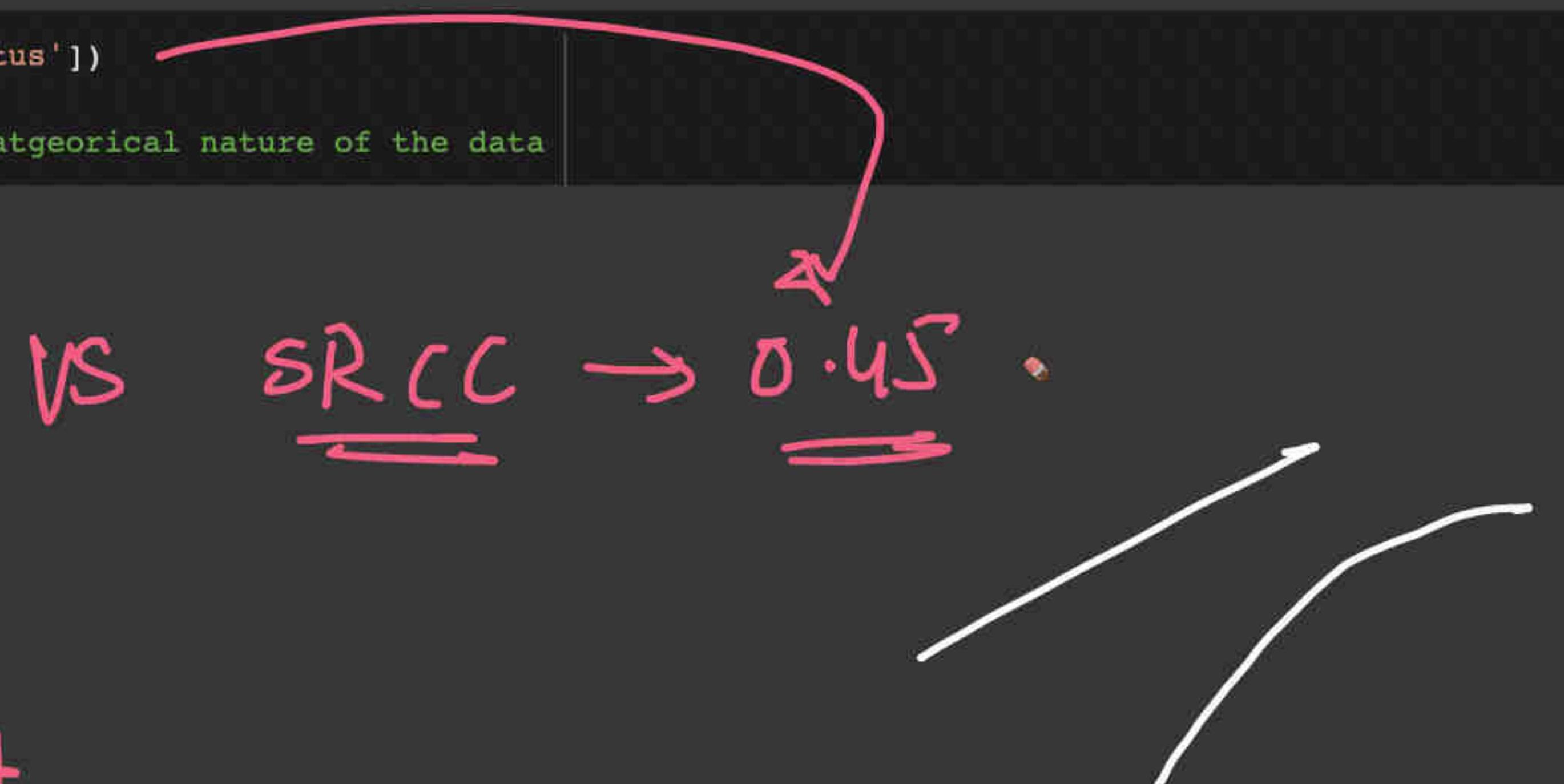
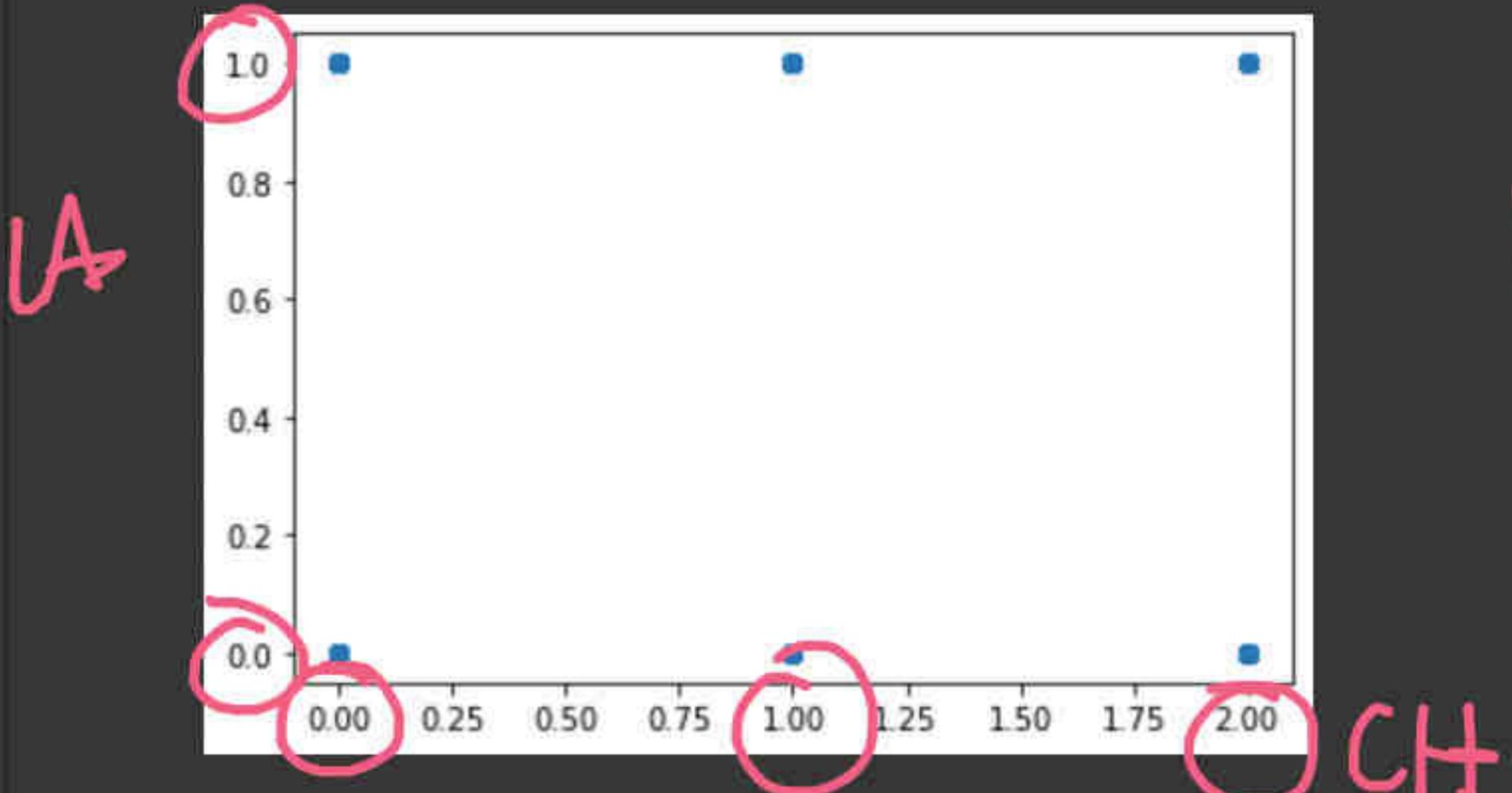
EDA_FE.ipynb - Colaboratory x markov x sklearn.impute.SimpleImputer x + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=A-B0OnWTn1WV

+ Code + Text

RAM Disk

Dependents Education Self_Employed Applicant Coapplicant LoanAmount Loan_Amount_perCapita Credit_History Property_Area Total_Amount_Applyed_for_Loan Able_to_pay

```
plt.scatter(data['Credit_History'], data['Loan_Status'])
plt.show()
#sometimes scatter plots can be misleading do to catgeorical nature of the data
```



```
sns.countplot(data = data, x = 'Credit_History', hue = 'Loan_Status')
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f3e2eb7d950>
```

EDA_FE.ipynb - Colaboratory x markov x sklearn.impute.SimpleImputer x +

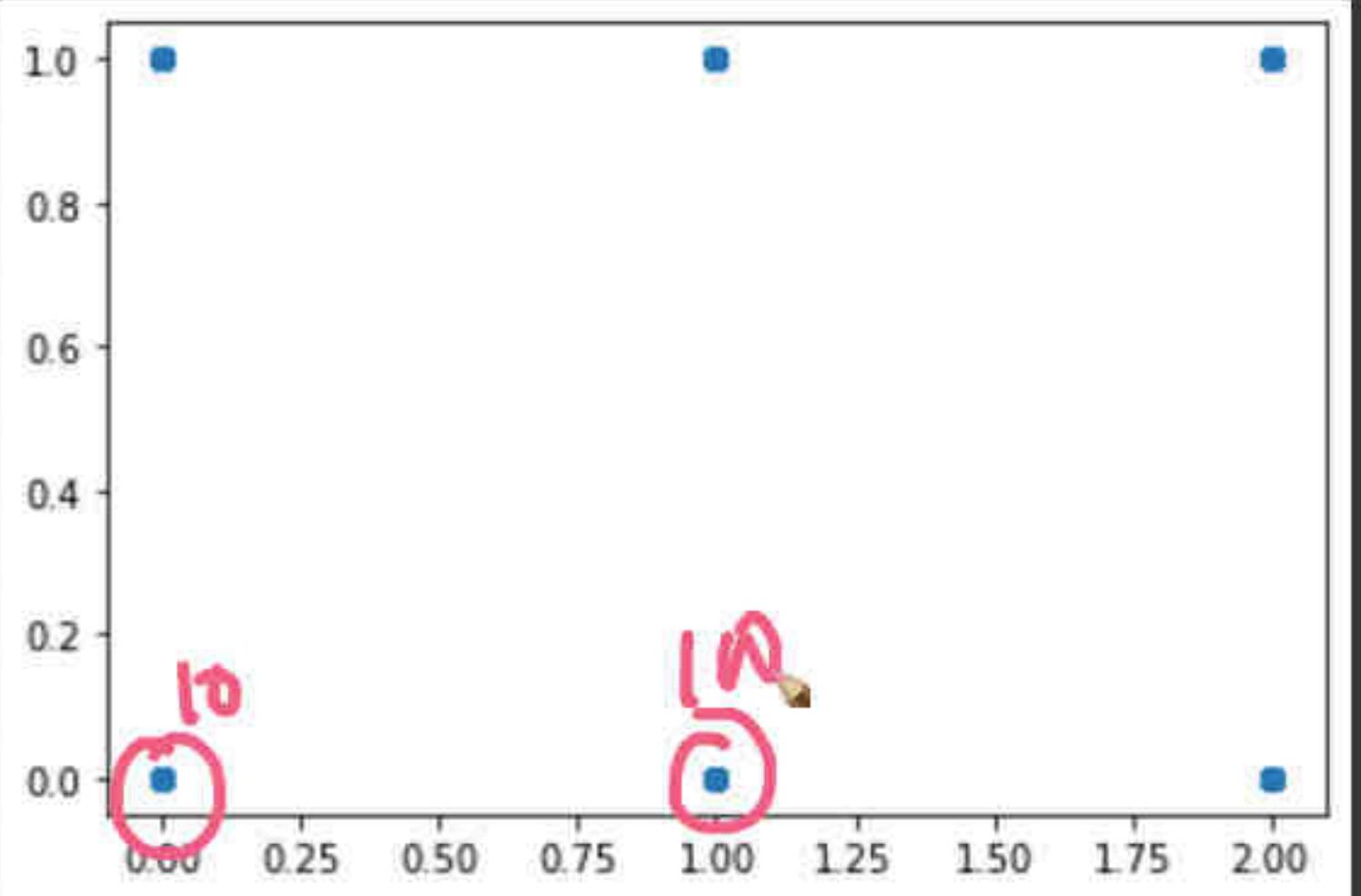
colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=A-B0OnWTn1WV

+ Code + Text

RAM Disk

Dependents Education Self_Employed Applicant Coapplicant LoanAmount Credit_History Property_Area Total Loan_Amount_perCapita Able_to_pay

```
plt.scatter(data['Credit_History'], data['Loan_Status'])
plt.show()
#sometimes scatter plots can be misleading do to catgeorical nature of the data
```



```
sns.countplot(data =data, x = 'Credit_History', hue = 'Loan_Status')
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f3e2eb7d950>
```

EDA_FE.ipynb - Colaboratory x markov x sklearn.impute.SimpleImputer x + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=A-B0OnWTn1WV

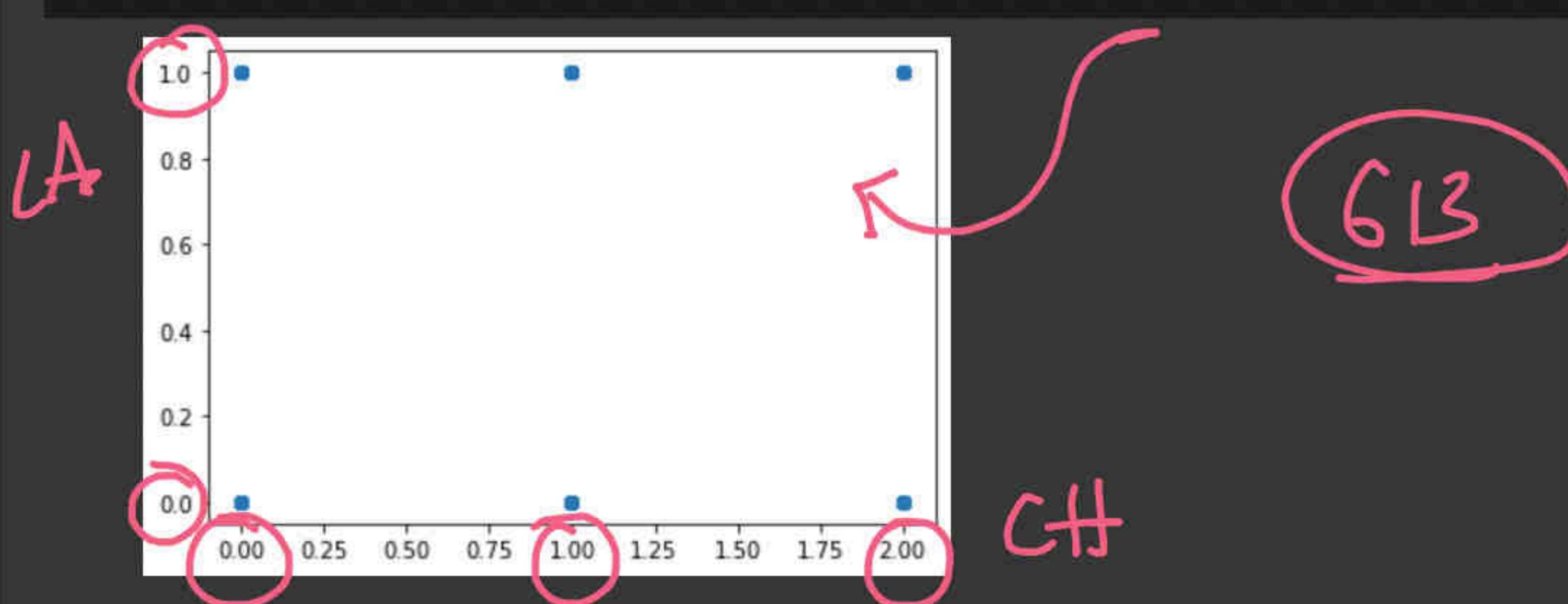
+ Code + Text

RAM Disk

Dependents Education Self_Employed Applicant Coapplicant LoanAmount Loan_Amount_perCapita Credit_History Property_Area Total_Amount_Applyed_for_Loan Able_to_pay_Loan

{x}

plt.scatter(data['Credit_History'], data['Loan_Status'])
plt.show()
#sometimes scatter plots can be misleading do to catgeorical nature of the data



sns.countplot(data = data, x = 'Credit_History', hue = 'Loan_Status')

<matplotlib.axes._subplots.AxesSubplot at 0x7f3e2eb7d950>

EDA_FE.ipynb - Colaboratory x markov x sklearn.impute.SimpleImputer x + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=A-B0OnWTn1WV

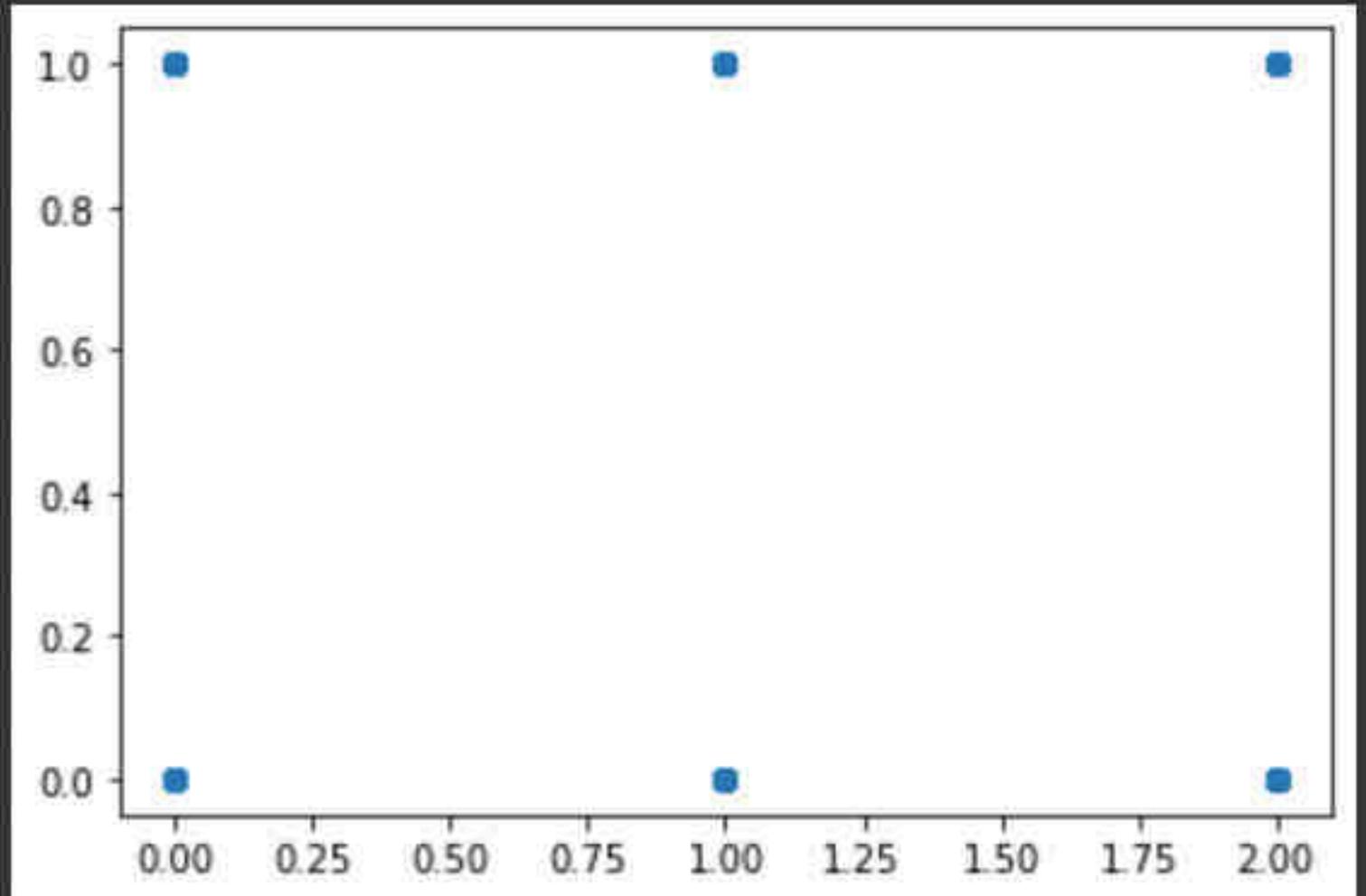
+ Code + Text

RAM Disk

Dependents Education Self_Employed Applicant Coapplicant LoanAmount Credit_History Property_Area Total_Loan_Amount_per_Capita Able_to_pay

```
plt.scatter(data['Credit_History'], data['Loan_Status'])
plt.show()
#sometimes scatter plots can be misleading do to catgeorical nature of the data
```

PCC = 



```
sns.countplot(data = data, x = 'Credit_History', hue = 'Loan_Status')
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f3e2eb7d950>
```

EDA_FE.ipynb - Colaboratory x markov x | sklearn.impute.SimpleImputer x +

colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=A-B0OnWTn1WV

+ Code + Text

RAM Disk

{x}

[] sns.countplot(data =data, x = 'Credit_History', hue = 'Loan_Status')

<matplotlib.axes._subplots.AxesSubplot at 0x7f3e2eb7d950>

Credit_History	Loan_Status	Count
0.0	0	85
0.0	1	10
1.0	0	100
1.0	1	370
2.0	0	10
2.0	1	40

Column Standarization and Normalization

- Mean centering and Variance scaling (Standard Scaling)
- MinMax Scaling

95 / 95

col. std &
normalization

$X = x_1 \dots x_n$

$\frac{x_i - m}{s} \Rightarrow \text{standardization}$

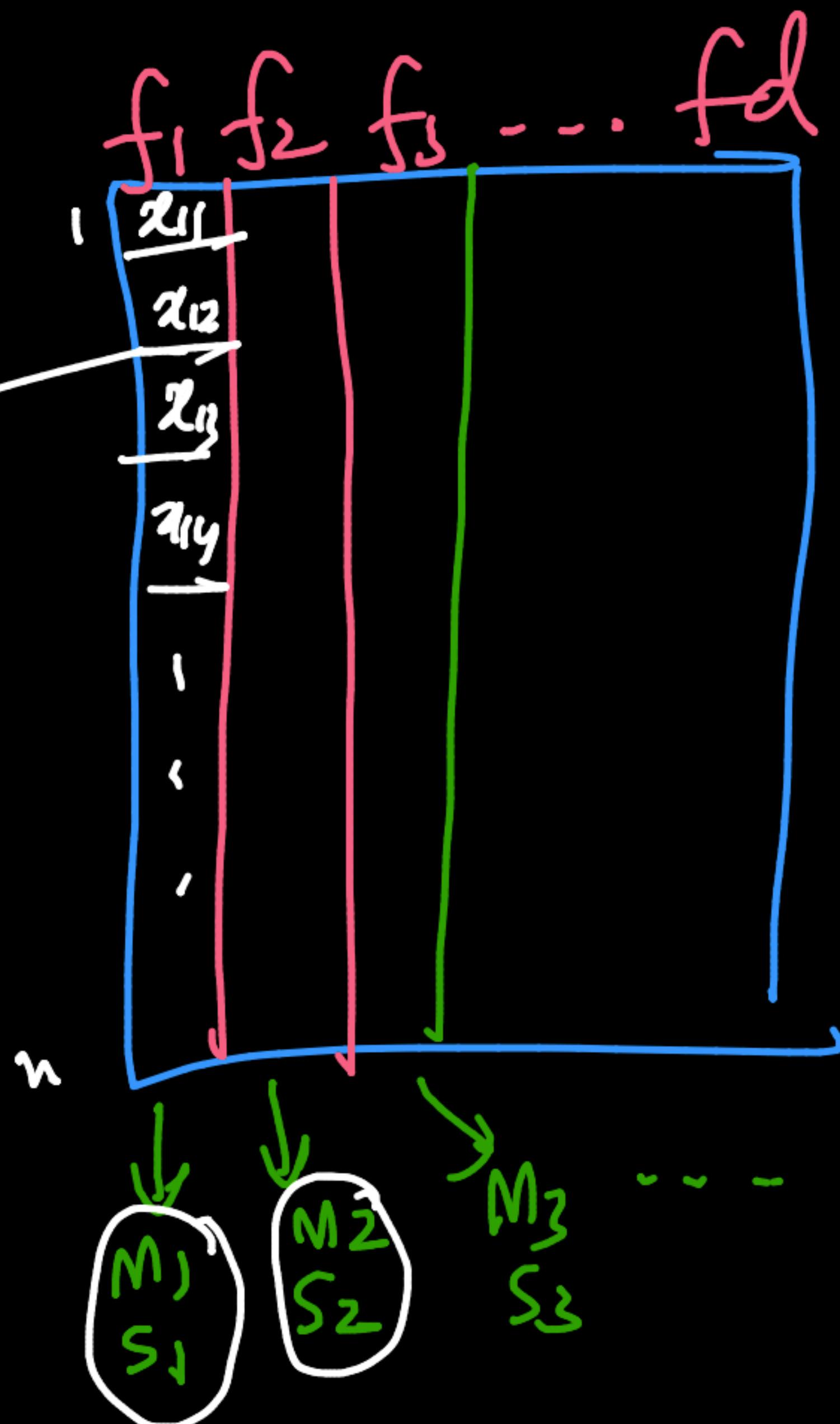
Poe-Poole's t_{ij}



Linear seg

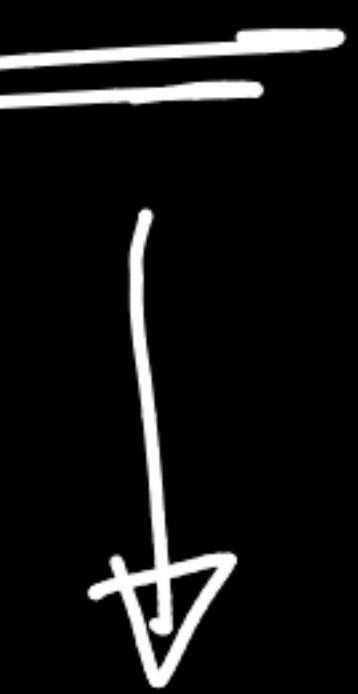
DL

(later)



$$\frac{x_{1j} - M_j}{S_1}$$

column · SF

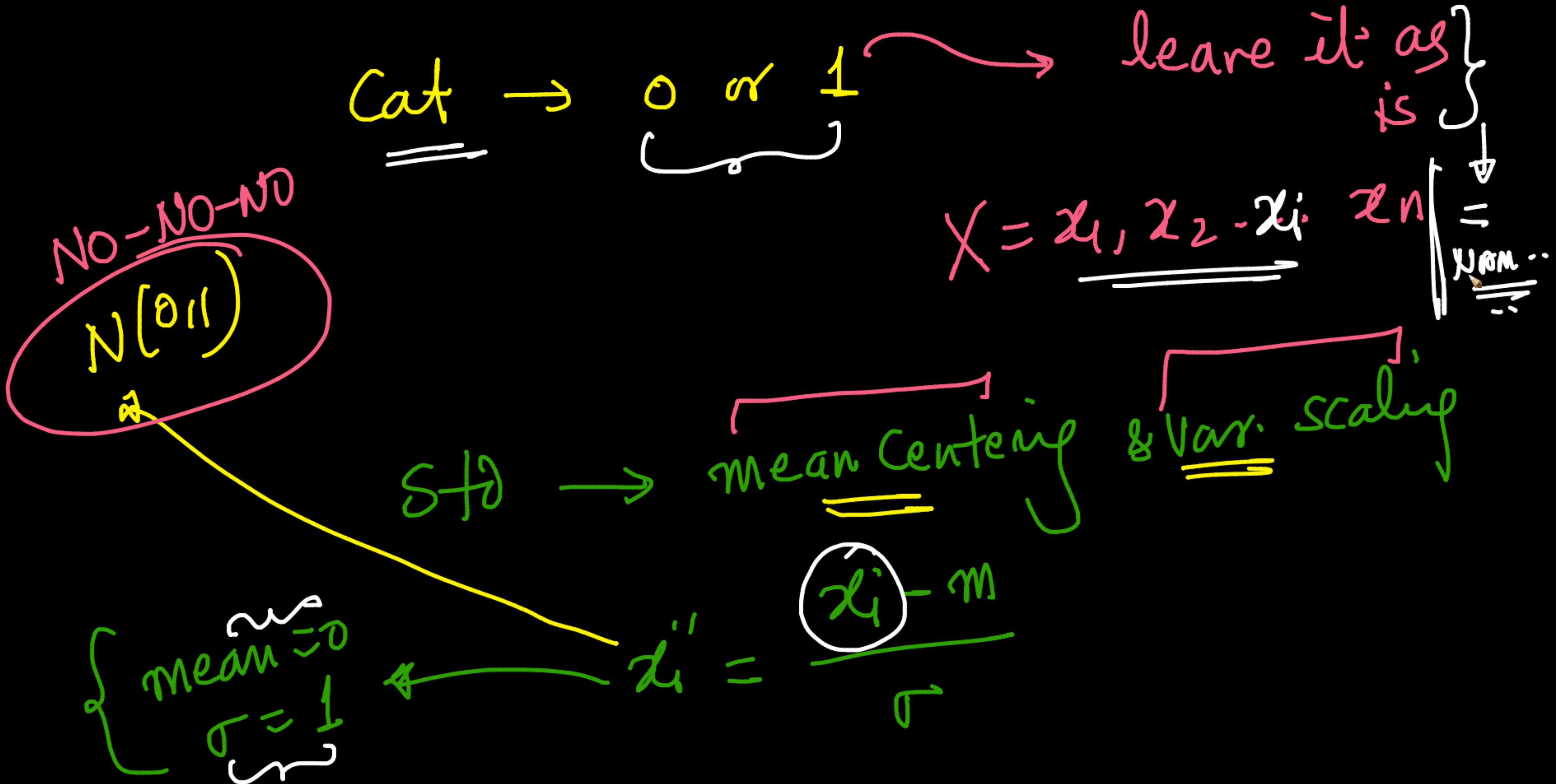


each col
is SF ..
independently

{ everything has to be numbers

Categorical → Numeric → ML

✓ [Why → aids in optimization problems (ML)
 ↓
 f) are sensible
 ↑
 g) labeled ...]



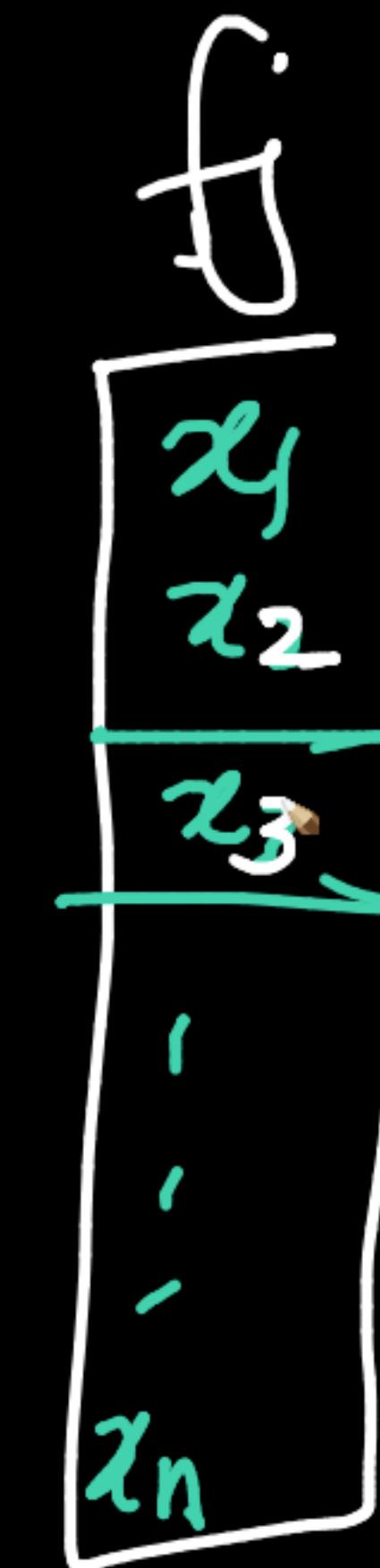
Normalization
=

→ Image-processing & CV
for each col

Min-Max
Scaling

min
Max

$$\frac{x_i - \text{min}}{\text{max} - \text{min}} = \tilde{x}_i \xrightarrow{\text{range}} [0, 1]$$



EDA_FE.ipynb - Colaboratory x markov x | sklearn.impute.SimpleImputer x + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=-6tEgElmk0a5

+ Code + Text RAM Disk ✓

Column Standardization and Normalization

{x} • Mean centering and Variance scaling (Standard Scaling) 

• MinMax Scaling 

from sklearn.preprocessing import StandardScaler, MinMaxScaler

scaler = StandardScaler()
std_data = scaler.fit_transform(data)
std_data = pd.DataFrame(std_data, columns=data.columns)
std_data.head()

	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount
0	0.472343	-1.372089	-0.737806	-0.528362	-0.174052	0.072991	-0.554487	-0.211241	0.211241
1	0.472343	0.728816	0.253470	-0.528362	-0.174052	-0.134412	-0.038732	-0.211241	0.211241
2	0.472343	0.728816	-0.737806	-0.528362	-0.586643	-0.393747	-0.554487	-0.948996	0.948996
3	0.472343	0.728816	-0.737806	1.892641	-0.174052	-0.462062	0.251980	-0.306435	0.306435
4	0.472343	-1.372089	-0.737806	-0.528362	-0.174052	0.097728	-0.554487	-0.056551	0.056551

EDA_FE.ipynb - Colaboratory x markov x | sklearn.impute.SimpleImputer x + colab.research.google.com/drive/1_P-nKswaKnx5v77IOx3pCxvJve73MvnN#scrollTo=-6tEgElmk0a5

+ Code + Text RAM Disk

Column Standardization and Normalization

{x}

- Mean centering and Variance scaling (Standard Scaling)
- MinMax Scaling

```
from sklearn.preprocessing import StandardScaler, MinMaxScaler  
scaler = StandardScaler()  
std_data = scaler.fit_transform(data)  
std_data = pd.DataFrame(std_data, columns=data.columns)  
std_data.head()
```

	Gender	Married	Dependents	Education	Self_Employed	ApplicantIncome	CoapplicantIncome	LoanAmount	Loan_Amount
0	0.472343	-1.372089	-0.737806	-0.528362	-0.174052	0.072991	-0.554487	-0.211241	0.211241
1	0.472343	0.728816	0.253470	-0.528362	-0.174052	-0.134412	-0.038732	-0.211241	0.211241
2	0.472343	0.728816	-0.737806	-0.528362	-0.586643	-0.393747	-0.554487	-0.948996	0.948996
3	0.472343	0.728816	-0.737806	1.892641	-0.174052	-0.462062	0.251980	-0.306435	0.306435
4	0.472343	-1.372089	-0.737806	-0.528362	-0.174052	0.097728	-0.554487	-0.056551	0.056551

====

Gaussian



"3"

Kurtosis

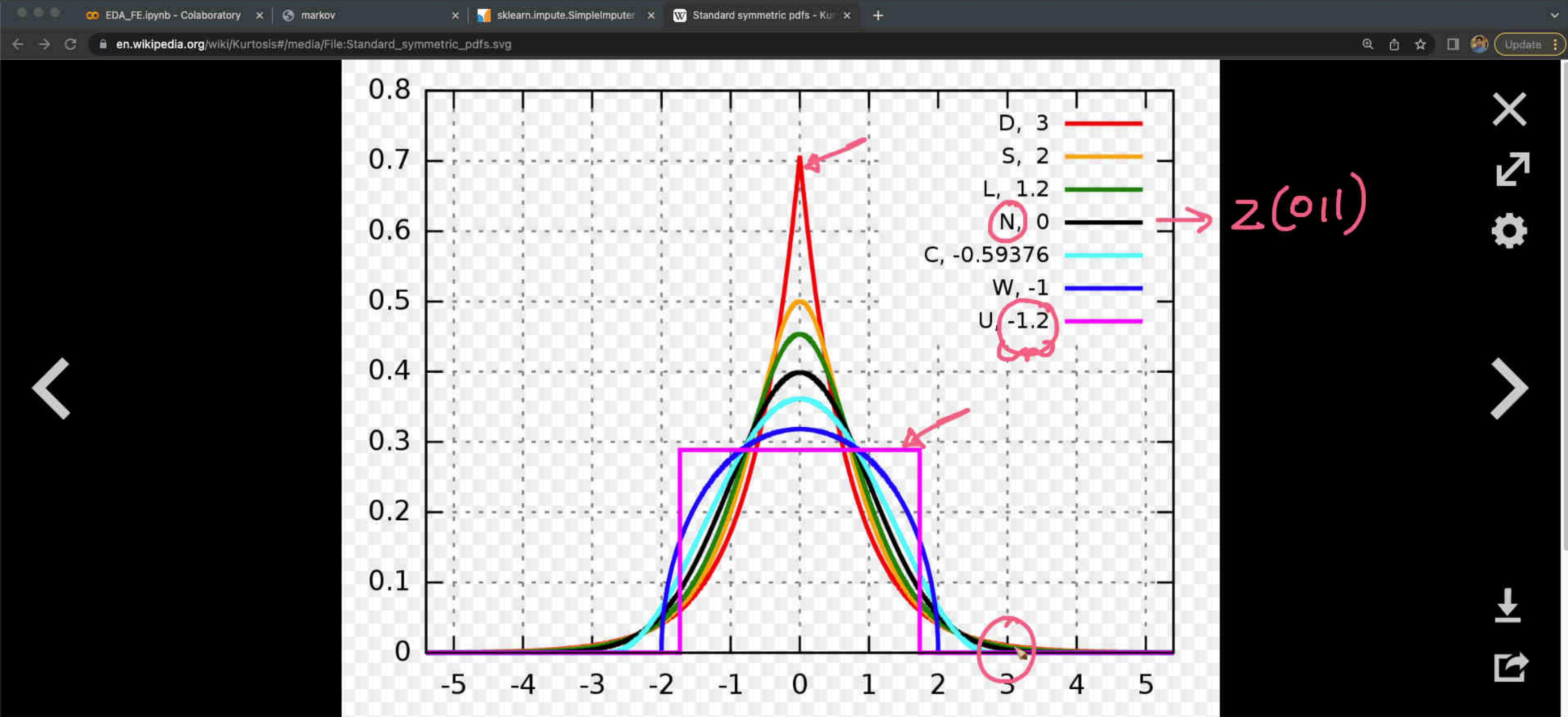


how tailed is your
dist

Excess Kurtosis

====

$$= K - 3.$$



Probability density functions for selected distributions with mean 0, variance 1 and different excess kurtosis

 More details



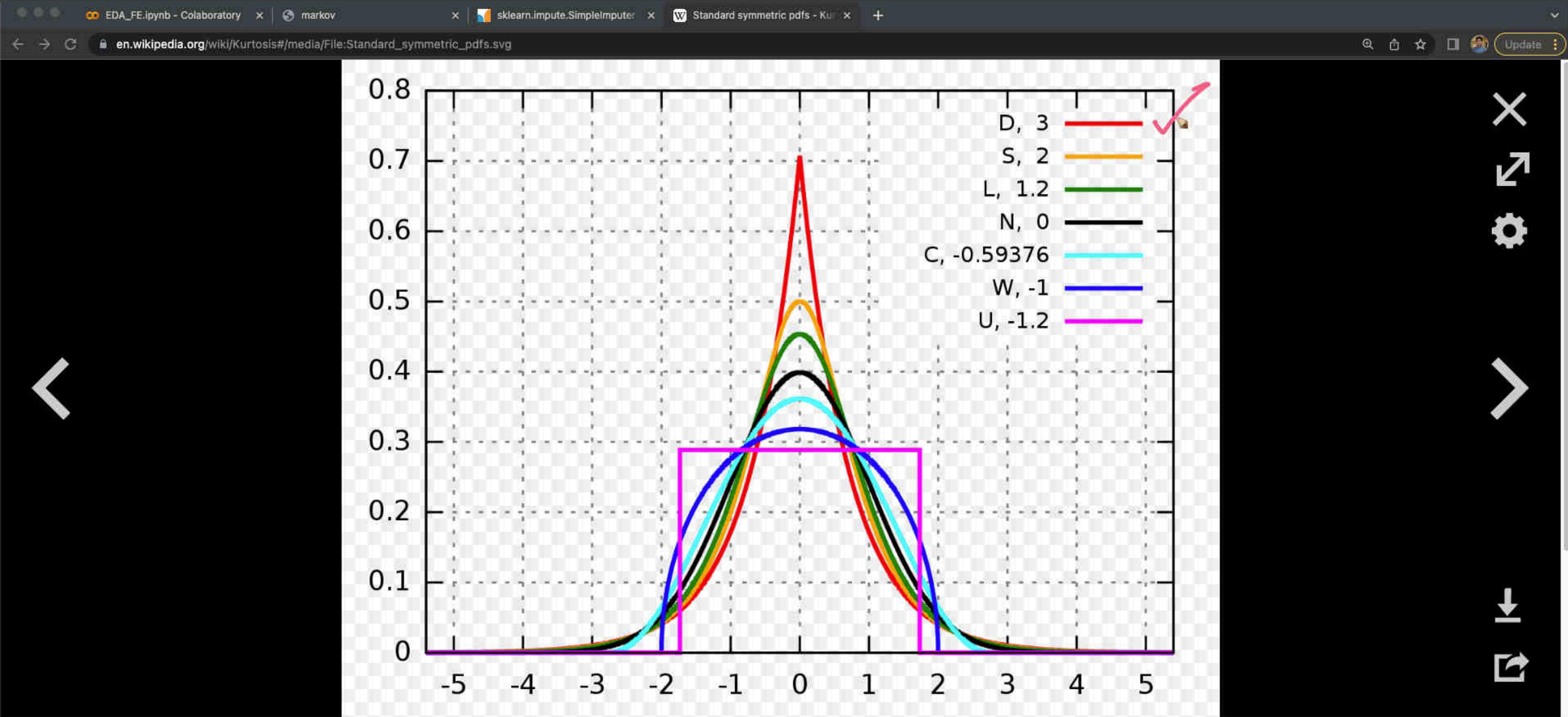
Original: MarkSweep

Vector: Andel - M



Public Domain

[view terms](#)

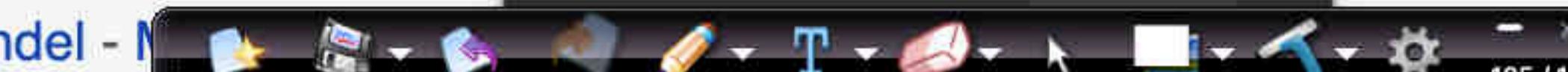


Probability density functions for selected distributions with mean 0, variance 1 and different excess kurtosis

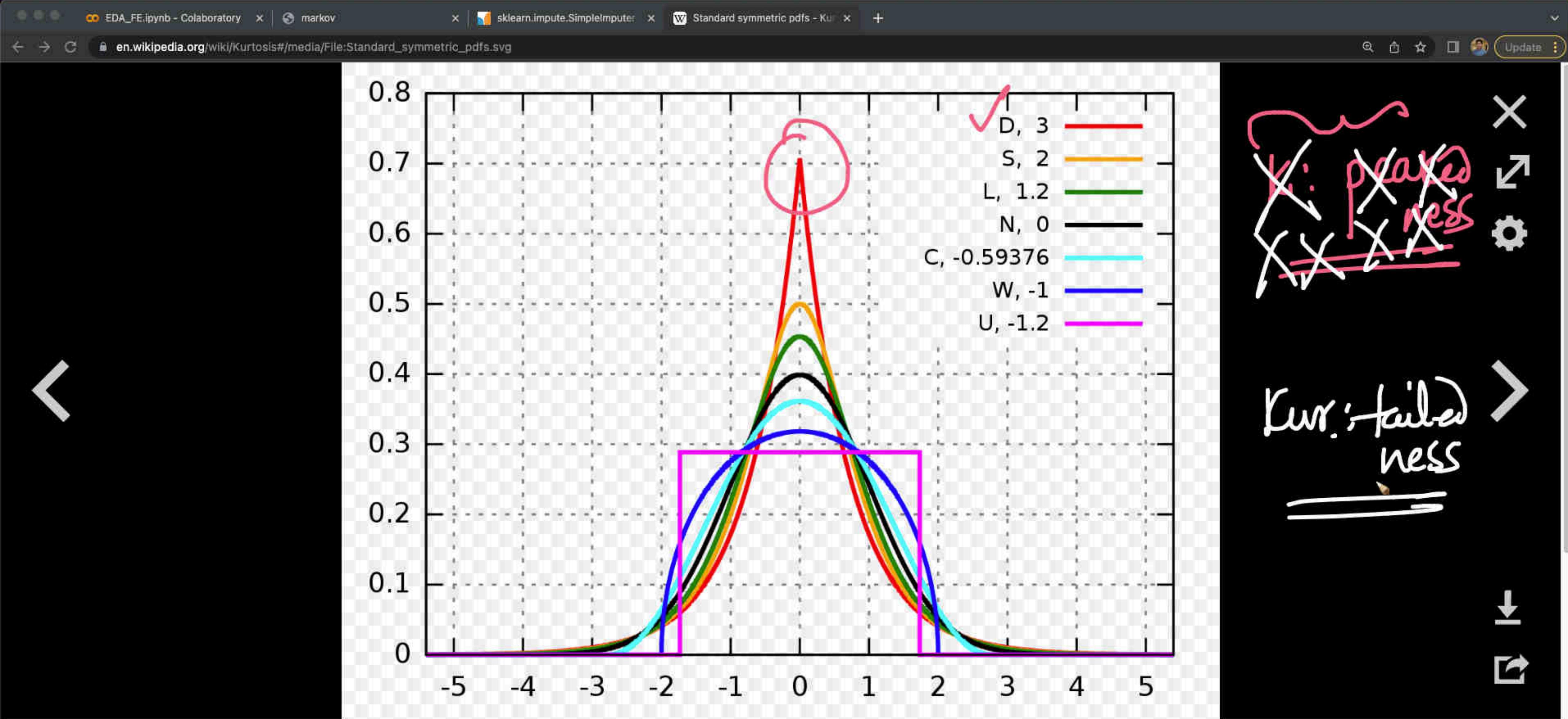
 More details



Original: MarkSweep Vector: Andel - M



Public Domain view terms



More details



Original: MarkSweep Vector: Andel - M



Public Domain view terms

EDA_FE.ipynb - Colaboratory | markov | sklearn.impute.SimpleImputer | Kurtosis - Wikipedia | en.wikipedia.org/wiki/Kurtosis | Update

10 External links

Pearson moments [edit]

The kurtosis is the fourth [standardized moment](#), defined as

$$\text{Kurt}[X] = E\left[\left(\frac{X - \mu}{\sigma}\right)^4\right] = \frac{E[(X - \mu)^4]}{(E[(X - \mu)^2])^2} = \frac{\mu_4}{\sigma^4},$$

where μ_4 is the fourth [central moment](#) and σ is the [standard deviation](#). Several letters are used in the literature to denote the kurtosis. A very common choice is κ , which is fine as long as it is clear that it does not refer to a [cumulant](#). Other choices include γ_2 , to be similar to the notation for skewness, although sometimes this is instead reserved for the excess kurtosis.

The kurtosis is bounded below by the squared [skewness](#) plus 1:^{[4]:432}

$$\frac{\mu_4}{\sigma^4} \geq \left(\frac{\mu_3}{\sigma^3}\right)^2 + 1,$$

where μ_3 is the third [central moment](#). The lower bound is realized by the [Bernoulli distribution](#). There is no upper limit to the kurtosis of a general probability distribution, and it may be infinite.

A reason why some authors favor the excess kurtosis is that cumulants are [extensive](#). Formulas related to the extensive property are more naturally expressed in terms of the excess kurtosis. For example, let X_1, \dots, X_n be independent random variables for which the fourth moment is finite. Then the fourth moment of the sum variable defined by the sum of the X_i . The excess kurtosis of X is

$X: x_1 x_2 \dots x_n$

$E\left[\left(\frac{X - \mu}{\sigma}\right)^4\right]$

EDA_FE.ipynb - Colaboratory x markov x sklearn.impute.SimpleImputer x Kurtosis - Wikipedia x +

en.wikipedia.org/wiki/Kurtosis

10 External links

Pearson moments [edit]

The kurtosis is the fourth [standardized moment](#), defined as

$$\text{Kurt}[X] = E\left[\left(\frac{X - \mu}{\sigma}\right)^4\right] = \frac{E[(X - \mu)^4]}{(E[(X - \mu)^2])^2} = \frac{\mu_4}{\sigma^4},$$

where μ_4 is the fourth [central moment](#) and σ is the [standard deviation](#). Several letters are used in the literature to denote the kurtosis. A very common choice is κ , which is fine as long as it is clear that it does not refer to a [cumulant](#). Other choices include γ_2 , to be similar to the notation for skewness, although sometimes this is instead reserved for the excess kurtosis.

The kurtosis is bounded below by the squared [skewness](#) plus 1:^{[4]:432}

$$\frac{\mu_4}{\sigma^4} \geq \left(\frac{\mu_3}{\sigma^3}\right)^2 + 1,$$

where μ_3 is the third [central moment](#). The lower bound is realized by the [Bernoulli distribution](#). There is no upper limit to the kurtosis of a general probability distribution, and it may be infinite.

A reason why some authors favor the excess kurtosis is that cumulants are [extensive](#). Formulas related to the extensive property are more naturally expressed in terms of the excess kurtosis. For example, let X_1, \dots, X_n be independent random variables for which the fourth moment is finite. Then the kurtosis of the sum of the X_i is equal to the sum of the kurtoses of the X_i plus $n(n-1)(n-2)(n-3)$ times the square of the standard deviation of the X_i .

EDA_FE.ipynb - Colaboratory | markov | sklearn.impute.SimpleImputer | Kurtosis - Wikipedia

en.wikipedia.org/wiki/Kurtosis

Further reading

10 External links

$$X: x_1 \ x_2 \ \dots \ x_n \rightarrow$$

$$\frac{1}{n} \sum_{i=1}^n \left(\frac{x_i - \mu}{s} \right)^4 = \text{Kurt}$$

$$Ek = \text{Kurt} - 3$$

Pearson moments [edit]

The kurtosis is the fourth [standardized moment](#), defined as

$$\text{Kurt}[X] = E \left[\left(\frac{X - \mu}{\sigma} \right)^4 \right] = \frac{E[(X - \mu)^4]}{(E[(X - \mu)^2])^2} = \frac{\mu_4}{\sigma^4},$$

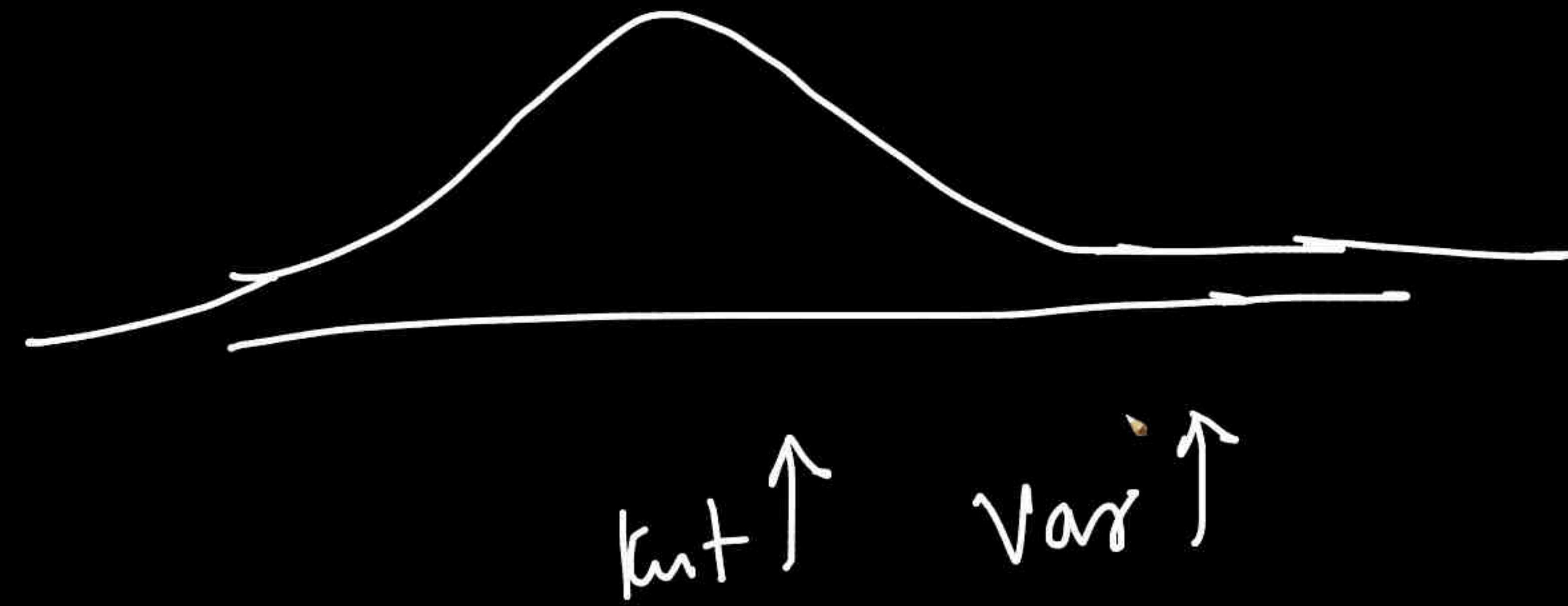
where μ_4 is the fourth [central moment](#) and σ is the [standard deviation](#). Several letters are used in the literature to denote the kurtosis. A very common choice is κ , which is fine as long as it is clear that it does not refer to a [cumulant](#). Other choices include γ_2 , to be similar to the notation for skewness, although sometimes this is instead reserved for the excess kurtosis.

The kurtosis is bounded below by the squared [skewness](#) plus 1:^{[4]:432}

$$\frac{\mu_4}{\sigma^4} \geq \left(\frac{\mu_3}{\sigma^3} \right)^2 + 1,$$

where μ_3 is the third [central moment](#). The lower bound is realized by the [Bernoulli distribution](#). There is no upper limit to the kurtosis of a general probability distribution, and it may be infinite.

A reason why some authors favor the excess kurtosis is that cumulants are [extensive](#). Formulas related to the extensive property are more naturally expressed in terms of the excess kurtosis. For example, let X_1, \dots, X_n be independent random variables for which



EDA_FE.ipynb - Colaboratory x | markov x | sklearn.impute.SimpleImputer x | Kurtosis - Wikipedia x +

en.wikipedia.org/wiki/Kurtosis

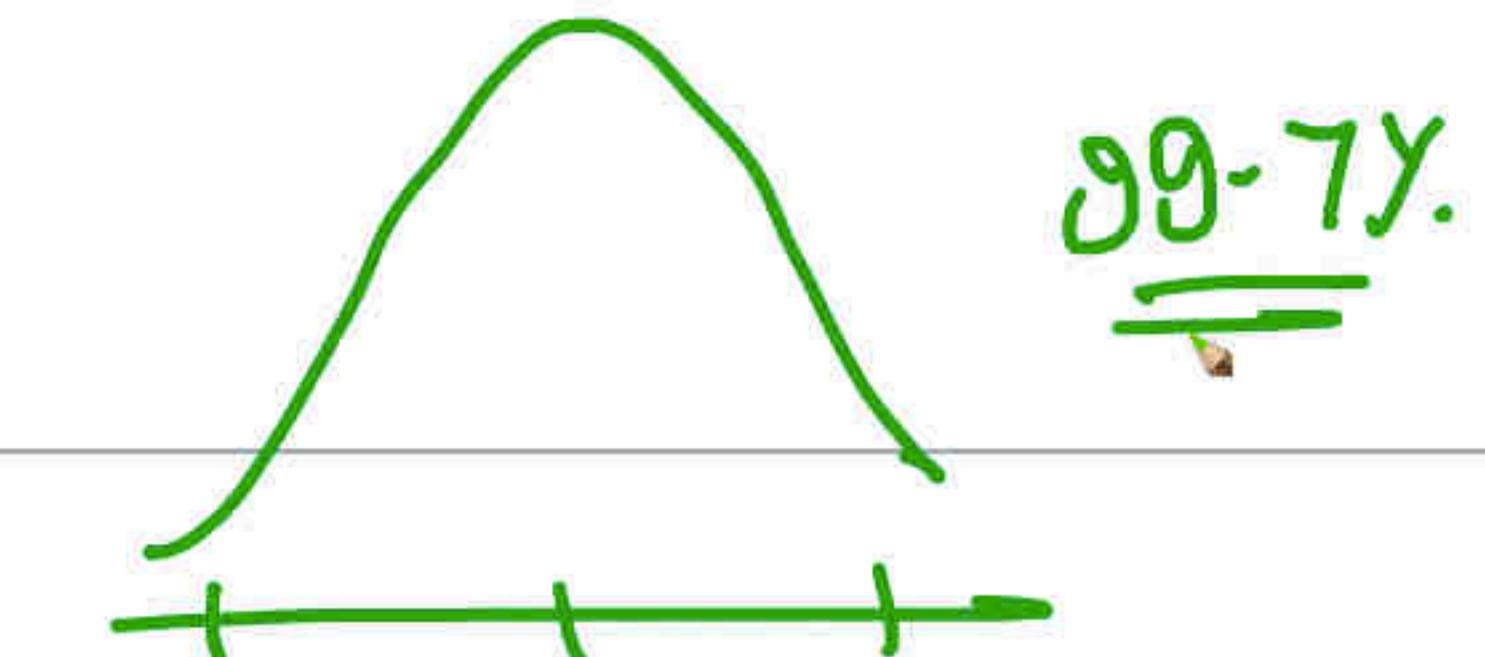
Further reading

10 External links

Pearson moments [edit]

The kurtosis is the fourth [standardized moment](#), defined as

$$\text{Kurt}[X] = E\left[\left(\frac{X - \mu}{\sigma}\right)^4\right] = \frac{E[(X - \mu)^4]}{(E[(X - \mu)^2])^2} = \frac{\mu_4}{\sigma^4},$$



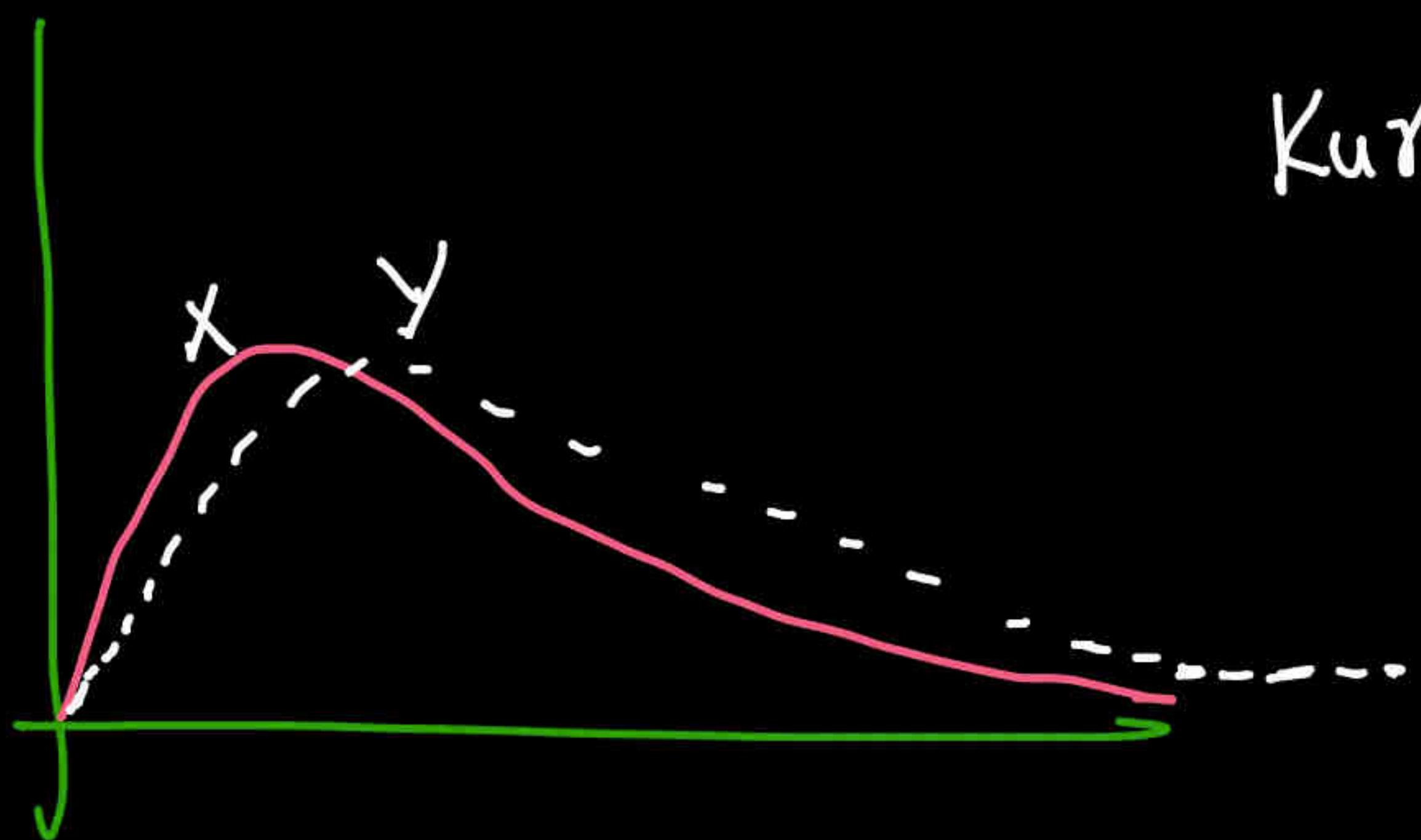
where μ_4 is the fourth [central moment](#) and σ is the [standard deviation](#). Several letters are used in the literature to denote the kurtosis. A very common choice is κ , which is fine as long as it is clear that it does not refer to a [cumulant](#). Other choices include γ_2 , to be similar to the notation for skewness, although sometimes this is instead reserved for the excess kurtosis.

The kurtosis is bounded below by the squared [skewness](#) plus 1:^{[4]:432}

$$\frac{\mu_4}{\sigma^4} \geq \left(\frac{\mu_3}{\sigma^3}\right)^2 + 1,$$

where μ_3 is the third [central moment](#). The lower bound is realized by the [Bernoulli distribution](#). There is no upper limit to the kurtosis of a general probability distribution, and it may be infinite.

A reason why some authors favor the excess kurtosis is that cumulants are [extensive](#). Formulas related to the extensive property are more naturally expressed in terms of the [excess kurtosis](#). For example, let X_1, \dots, X_n be independent random variables for which the sum $X = \sum_{i=1}^n X_i$ is defined by the sum of the X_i . The excess kurtosis of X is



$$\text{Kurt}(X) < \text{Kurt}(Y)$$

≡

Kurt

Normal →

3

X 3

1

Ex $\text{funt} = \underline{\text{list}}^{\sim 3}$

三

2

4

Q

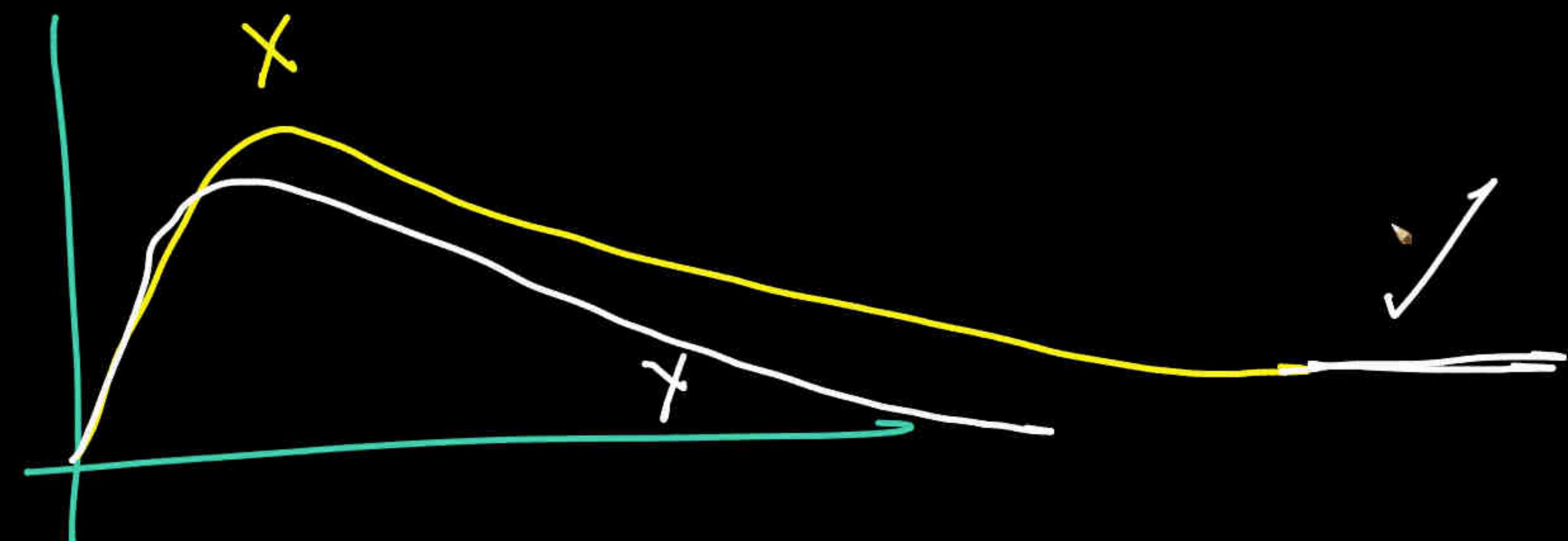
X

Y

Why?

Extrn: 3

2



Q

\bar{X} : mean $[lb \quad ub] \leftarrow$ sample of size n

95% C.I. of pop-mean is $[lb \quad ub]$

95% C.I.
OK

Using sample sizes of n , if we repeatedly compute the the C.I.
there is a 95% prob of the μ lying in these C.I's

Q

X: $x_1 \ x_2 \ \dots \ x_n$

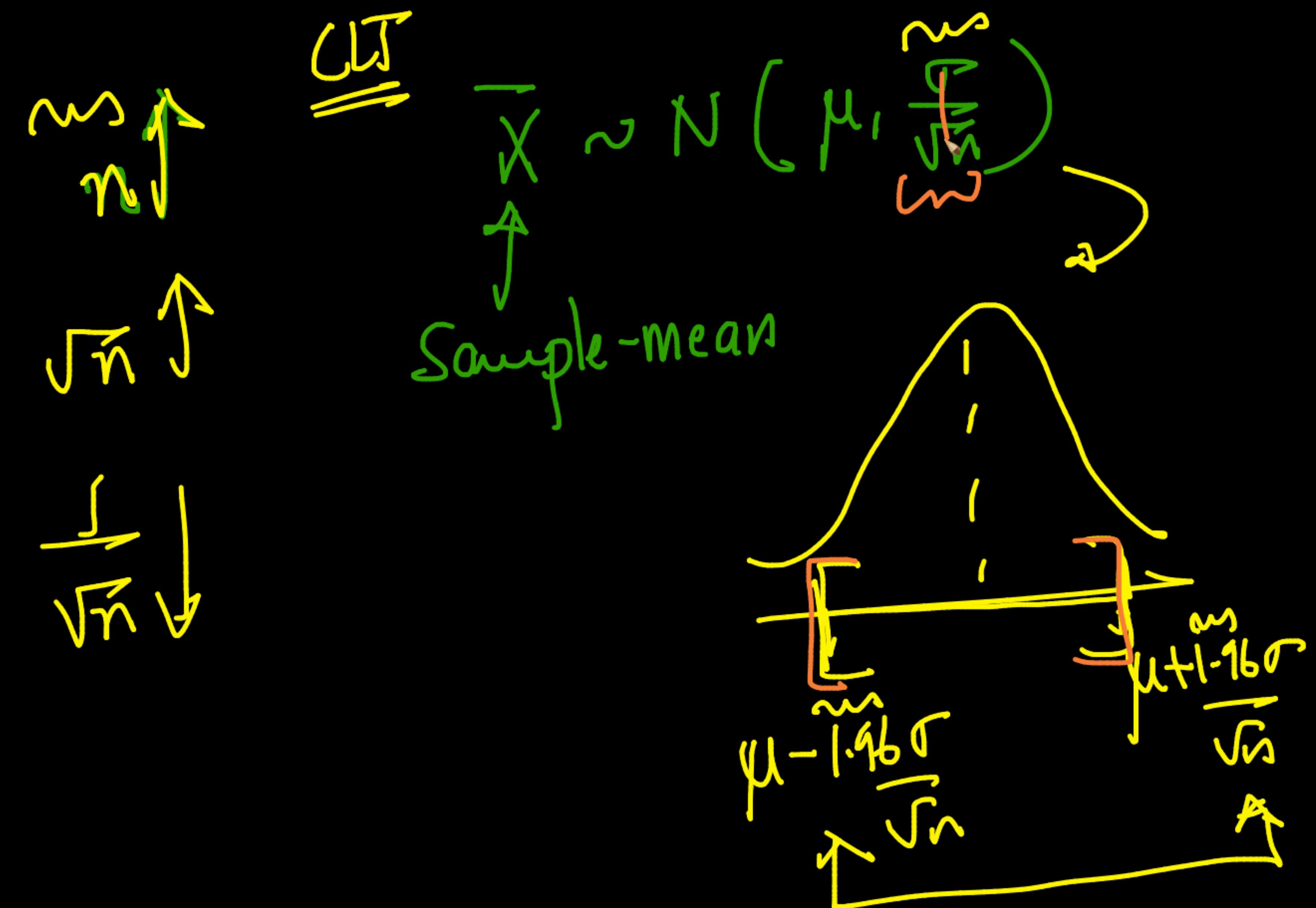
CDF (empirical)

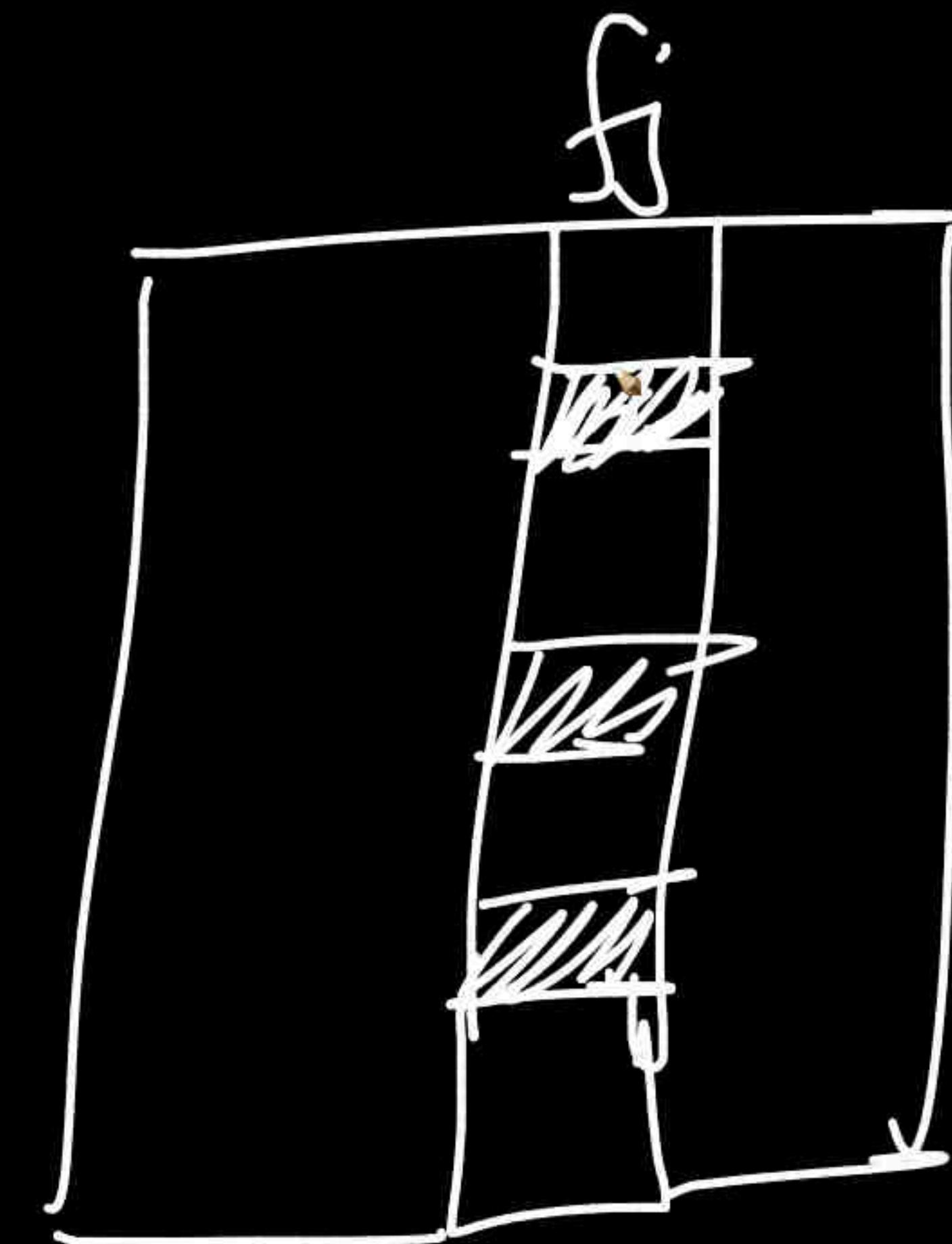
PDF (empirical)

$X \sim N(\mu=0, \sigma=2)$

Th-CDF

Th-PDF

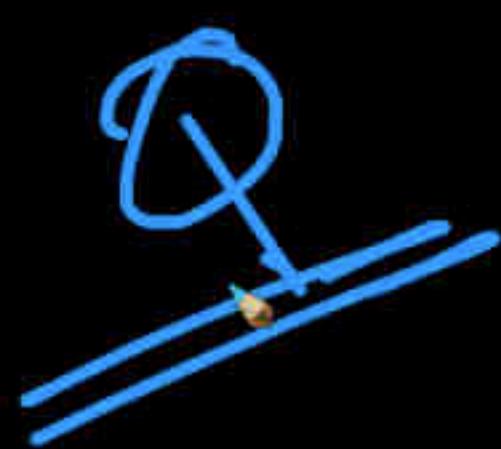




Q
—

f-disb(?)

$$N(\theta_{11}) = \underline{\underline{z-disb}}$$



Kur 2

Definition [edit]

Fisher's moment coefficient of skewness [edit]

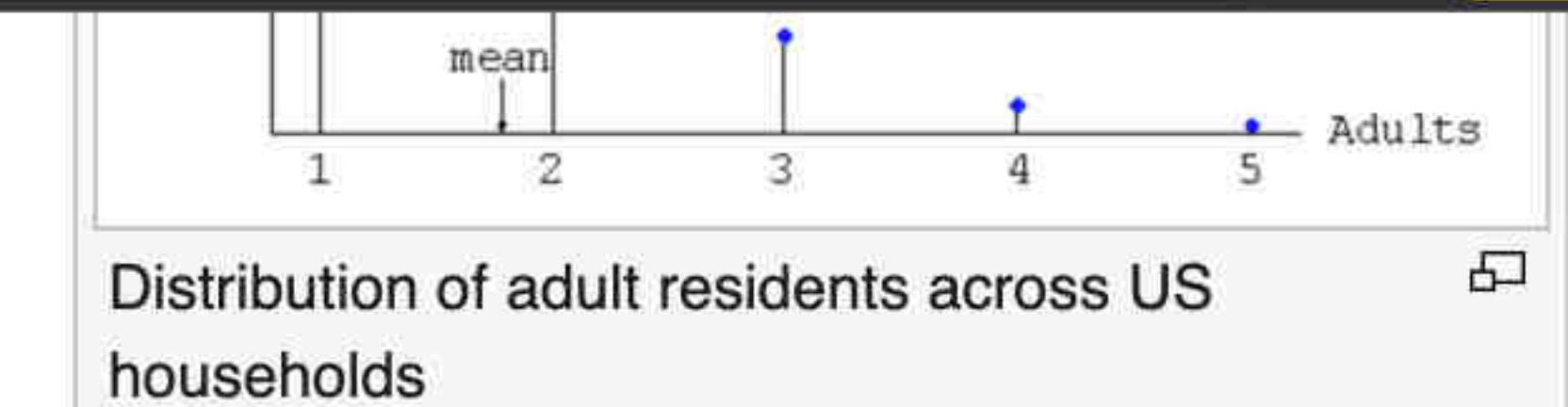
The skewness of a random variable X is the third standardized moment $\tilde{\mu}_3$, defined as:^{[4][5]}

$$\tilde{\mu}_3 = E\left[\left(\frac{X - \mu}{\sigma}\right)^3\right] = \frac{\mu_3}{\sigma^3} = \frac{E[(X - \mu)^3]}{(E[(X - \mu)^2])^{3/2}} = \frac{\kappa_3}{\kappa_2^{3/2}}$$

where μ is the mean, σ is the standard deviation, E is the expectation operator, μ_3 is the third central moment, and κ_t are the t -th cumulants. It is sometimes referred to as Pearson's moment coefficient of skewness,^[5] or simply the moment coefficient of skewness,^[4] but should not be confused with Pearson's other skewness statistics (see below). The last equality expresses skewness in terms of the ratio of the third cumulant κ_3 to the 1.5th power of the second cumulant κ_2 . This is analogous to the definition of kurtosis as the fourth cumulant normalized by the square of the second cumulant. The skewness is also sometimes denoted $\text{Skew}[X]$.

If σ is finite, μ is finite too and skewness can be expressed in terms of the non-central moment $E[X^3]$ by expanding the previous formula,

$$\begin{aligned}\tilde{\mu}_3 &= E\left[\left(\frac{X - \mu}{\sigma}\right)^3\right] \\ &= \frac{E[X^3] - 3\mu E[X^2] + 3\mu^2 E[X] - \mu^3}{\sigma^3}\end{aligned}$$



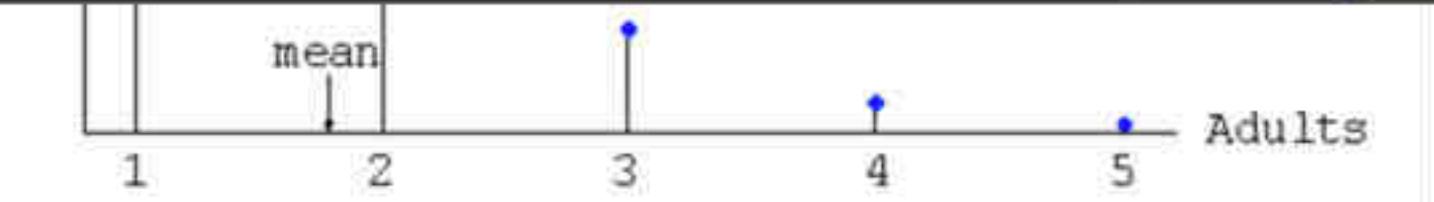
EDA_FE.ipynb - Colaboratory | markov | sklearn.impute.SimpleImputer | Kurtosis - Wikipedia | Skewness - Wikipedia +

en.wikipedia.org/wiki/Skewness#Definition

Definition [edit]

Fisher's moment coefficient of skewness [edit]

The skewness of a random variable X is the third standardized moment $\tilde{\mu}_3$, defined as:^{[4][5]}



Distribution of adult residents across US households

$$\tilde{\mu}_3 = E\left[\left(\frac{X - \mu}{\sigma}\right)^3\right] = \frac{\mu_3}{\sigma^3} = \frac{E[(X - \mu)^3]}{(E[(X - \mu)^2])^{3/2}} = \frac{\kappa_3}{\kappa_2^{3/2}}$$

where μ is the mean, σ is the standard deviation, E is the expectation operator, μ_3 is the third central moment, and κ_t are the t -th cumulants. It is sometimes referred to as Pearson's moment coefficient of skewness,^[5] or simply the moment coefficient of skewness,^[4] but should not be confused with Pearson's other skewness statistics (see below). The last equality expresses skewness in terms of the ratio of the third cumulant κ_3 to the 1.5th power of the second cumulant κ_2 . This is analogous to the definition of kurtosis as the fourth cumulant normalized by the square of the second cumulant. The skewness is also sometimes denoted $\text{Skew}[X]$.

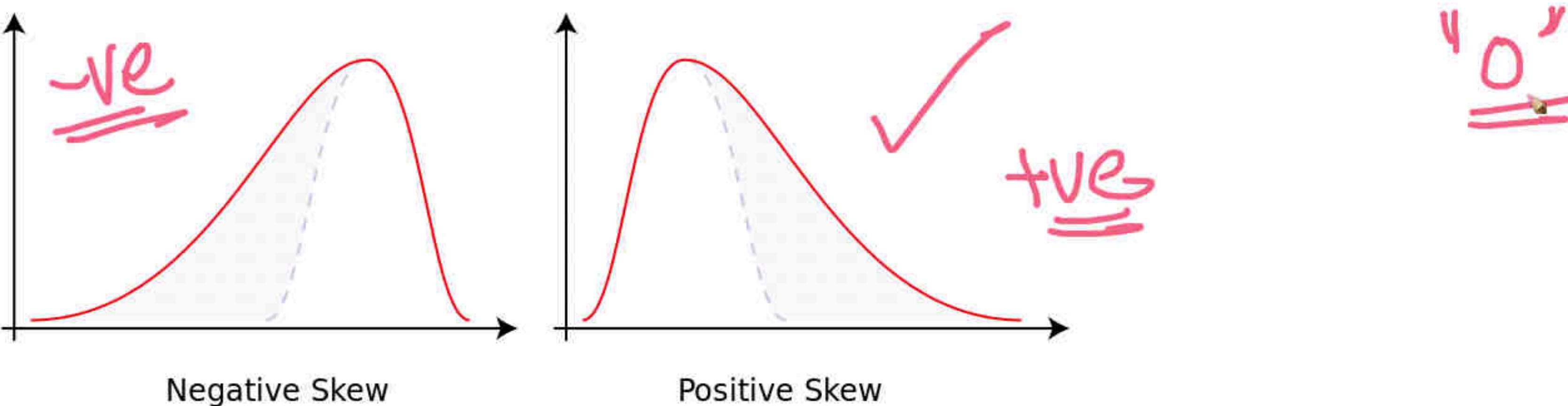
If σ is finite, μ is finite too and skewness can be expressed in terms of the non-central moment $E[X^3]$ by expanding the previous formula,

$$\begin{aligned}\tilde{\mu}_3 &= E\left[\left(\frac{X - \mu}{\sigma}\right)^3\right] \\ &= \frac{E[X^3] - 3\mu E[X^2] + 3\mu^2 E[X] - \mu^3}{\sigma^3}\end{aligned}$$

EDA_FE.ipynb - Colaboratory | markov | sklearn.impute.SimpleImputer | Kurtosis - Wikipedia | Skewness - Wikipedia + en.wikipedia.org/wiki/Skewness#Definition Update

be skewed or leaning to the right; *left* instead refers to the left tail being drawn out and, often, the mean being skewed to the left of a typical center of the data. A left-skewed distribution usually appears as a *right-leaning* curve.^[1]

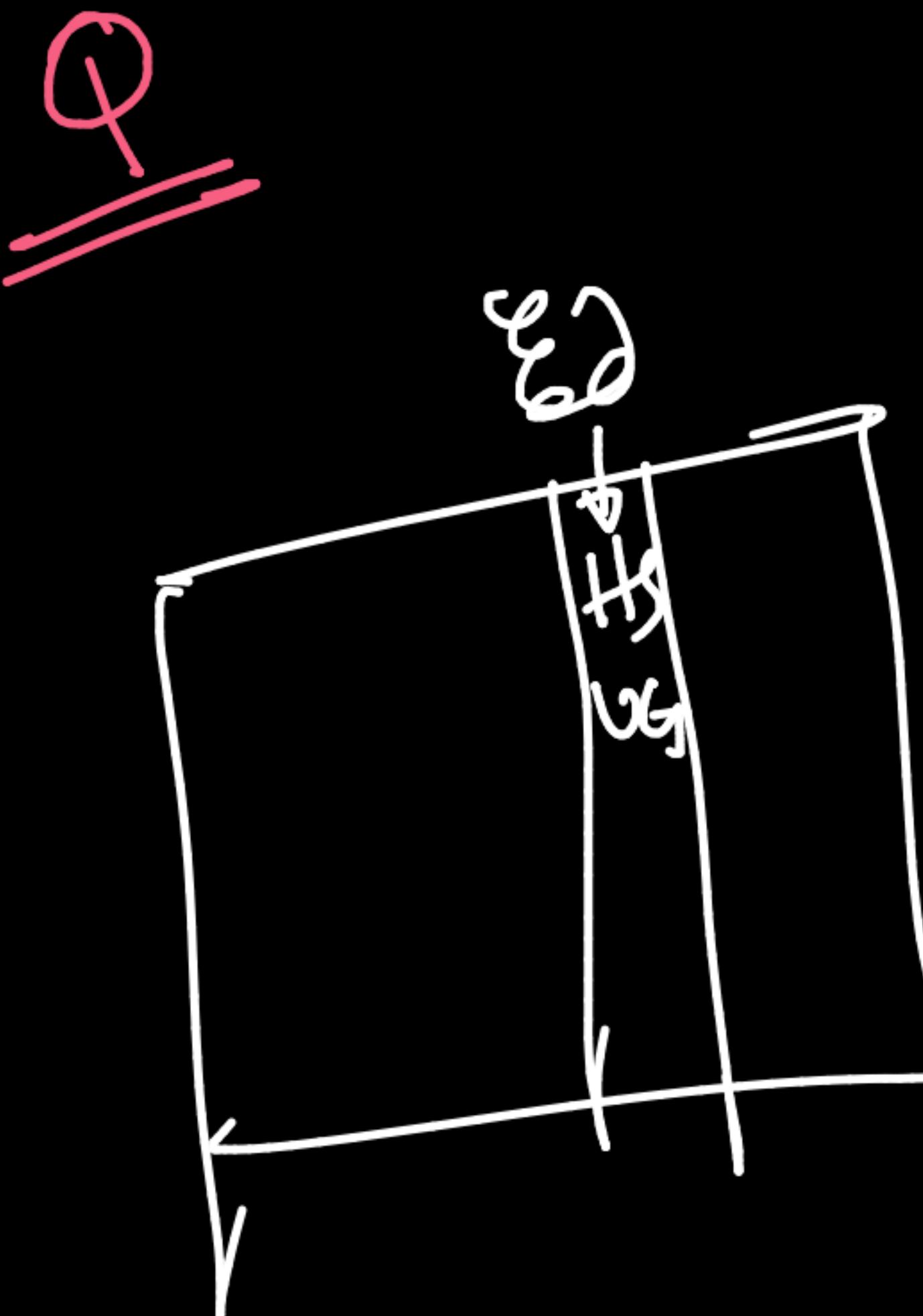
2. *positive skew*: The right tail is longer; the mass of the distribution is concentrated on the left of the figure. The distribution is said to be *right-skewed*, *right-tailed*, or *skewed to the right*, despite the fact that the curve itself appears to be skewed or leaning to the left; *right* instead refers to the right tail being drawn out and, often, the mean being skewed to the right of a typical center of the data. A right-skewed distribution usually appears as a *left-leaning* curve.^[1]



Negative Skew Positive Skew

Skewness in a data series may sometimes be observed not only graphically but by simple inspection of the values. For instance, consider the numeric sequence (49, 50, 51), whose values are evenly distributed around a central value of 50. We can transform this sequence into a negatively skewed distribution by adding a value far below the mean, which is probably a negative **outlier**, e.g. (40, 49, 50, 51). Therefore, the mean of the sequence becomes 47.5, and the median is 49.5. Based on the formula of **nonparametric skew**, defined as $(\mu - \nu)/\sigma$, the skew is negative. Similarly, we can make the sequence pos

is probably a positive outlier, e.g. (49, 50,



Ed: HS, UG, PG, Phd

1 → HS

2 → UG

3 → PG

4 → Phd

5 → HS

1	0	0	0
0	1	0	0
0	0	1	0
0	0	0	1
1	0	0	0



fi

:



:

:

:



fi

$$E(X) = \frac{1}{n} \sum_{i=1}^n x_i$$

$x_1 x_2 \dots x_n$

PMF: $P(X=k) \quad \forall k$

$$E(X) = \sum_k k \cdot P(X=k)$$

$$\int_F k \cdot p(k) dk$$

EDA_FE.ipynb - Colaboratory x markov x sklearn.impute.SimpleImputer x Kurtosis - Wikipedia x Skewness - Wikipedia x Moment (mathematics) - Wikipedia x New Tab x +

math.dartmouth.edu/~m20x18/markov

markov

1 / 4 | - 175% + | ☰ ⚡

MATH 20 – INEQUALITIES OF MARKOV AND CHEBYSHEV

Often, given a random variable X whose distribution is unknown but whose expected value μ is known, we may want to ask how likely it is for X to be ‘far’ from μ , or how likely it is for this random variable to be ‘very large.’ This would give us some idea of the spread of the distribution, though perhaps not a complete picture.

Proposition 1 (Markov’s Inequality). *Let X be a random variable that takes only nonnegative values. Then for any positive real number a ,*

$$P(X \geq a) \leq \frac{E(X)}{a}$$

provided $E(X)$ exists.

For example, Markov’s inequality tells us that as long as X doesn’t take negative values, the probability of X being at least k times its expected value is at most $\frac{1}{k}$. Notice that the only things we assumed about this random variable are that it can’t be negative and has finite mean; we don’t need to know