

Topics:

- Imbalanced data (Code + Case-study) → Fraud-detection

- Decision Trees (intuition & Math)

✓ [- ROC & AU-ROC
= = -]

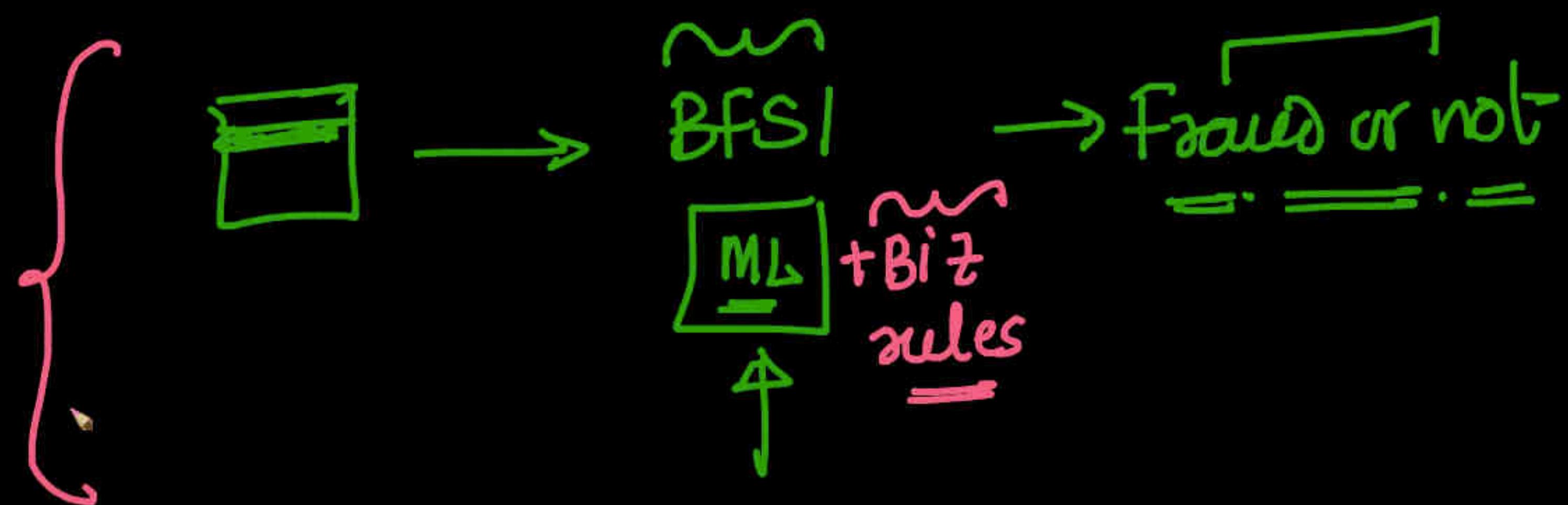
→ GBDT → Xgboost }

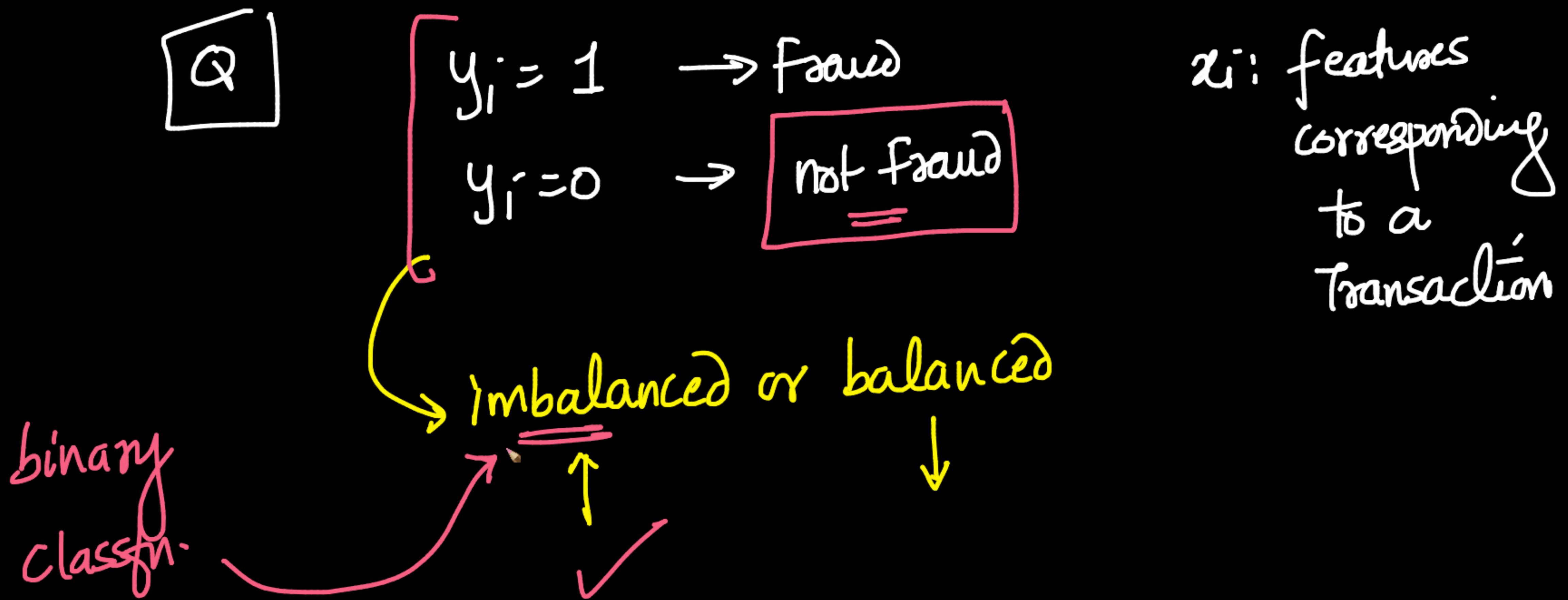
Next
= :
;

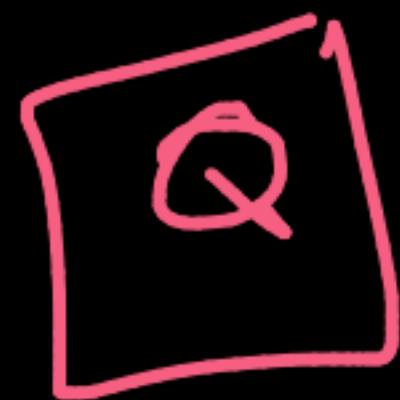
- Algo.
- Various
scenarios

✓ { - RF & GBDT

credit card







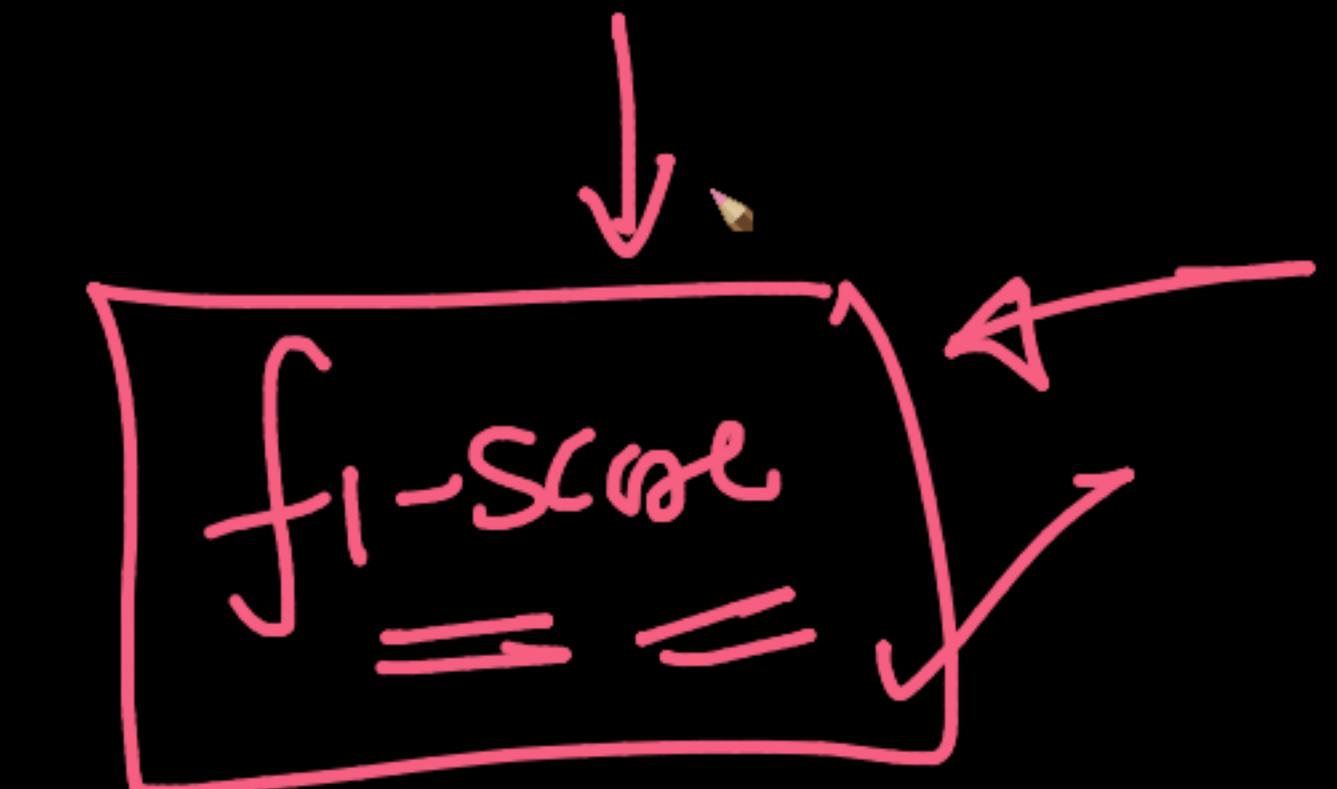
Business metric to measure the model perf-

on:

✓ Recall ← what %age of fraudulent-
xns did our model flag
as fraudulent

Problem: $y_i = 1 \text{ if}$
 $y_i = 0 \text{ else}$
recall = 100%.

✓ Precision: high



∞ imbalanced Data.ipynb - Colab | plot (-x*log(x)-(1-x)*log(1-x)) | Category Encoders — Category Encoders | +

contrib.scikit-learn.org/category_encoders/

Category Encoders

2.4.0

Search docs

Backward Difference Coding

BaseN

Binary

CatBoost Encoder

Count Encoder

Generalized Linear Mixed Model Encoder

Hashing

Helmet Coding

James-Stein Encoder

Leave One Out

M-estimate

One Hot

Ordinal

Polynomial Coding

Sum Coding

Target Encoder

Weight of Evidence

Wrappers

Quantile Encoder

Summary Encoder

» Category Encoders

OHE

View page source

Category Encoders

A set of scikit-learn-style transformers for encoding categorical variables into numeric with different techniques. While ordinal, one-hot, and hashing encoders have similar equivalents in the existing scikit-learn version, the transformers in this library all share a few useful properties:

- First-class support for pandas dataframes as an input (and optionally as output)
- Can explicitly configure which columns in the data are encoded by name or index, or infer non-numeric columns regardless of input type
- Can drop any columns with very low variance based on training set optionally
- Portability: train a transformer on data, pickle it, reuse it later and get the same thing out.
- Full compatibility with sklearn pipelines, input an array-like dataset like any other transformer

Usage

install as:

```
pip install category_encoders
```

or

```
conda install -c conda-forge category_encoders
```

To use:

5 / 5

co_imbalanced Data.ipynb - Colab | plot (-x*log(x)-(1-x)*log(1-x)) | Category Encoders - Catego | +

colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=dp3Jn34iJCBE

Update

+ Code + Text

RAM Disk

Up Down Left Right Home End Page Up Page Down

```
[ ] import pandas as pd
import numpy as np
from numpy import argmax

from datetime import date, time, timedelta
import pendulum # for time formatting

import matplotlib.pyplot as plt
import seaborn as sns

import category_encoders as ce # for categorical encoding
from sklearn.preprocessing import LabelEncoder, StandardScaler

from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.metrics import confusion_matrix, accuracy_score, precision_score, recall_score, f1_score, precision_re
```

```
[ ] !wget "https://drive.google.com/uc?export=download&id=1HC7MfMyg4Kph051bc21JKIxsubmc18C" -O orders.csv
```

```
--2022-05-04 07:24:45-- https://drive.google.com/uc?export=download&id=1HC7MfMyg4Kph051bc21JKIxsubmc18C
Resolving drive.google.com (drive.google.com)... 172.253.122.138, 172.253.122.102, 172.253.122.100, ...
Connecting to drive.google.com|||172.253.122.138|:443... connected.
```



∞ imbalanced Data.ipynb - Colab × plot (-x*log(x)-(1-x)*log(1-x)) × Category Encoders - Category Encoders × + colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=dp3Jn34iJCBE

+ Code + Text RAM Disk

Requirement already satisfied: threadpoolctl>=2.0.0 in /usr/local/lib/python3.7/dist-packages
Installing collected packages: category-encoders
Successfully installed category-encoders-2.4.0

6s

{x}

[] import pandas as pd
import numpy as np
from numpy import argmax

from datetime import date, time, timedelta
import pendulum # for time formatting

import matplotlib.pyplot as plt
import seaborn as sns

import category_encoders as ce # for categorical encoding
from sklearn.preprocessing import LabelEncoder, StandardScaler

from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.metrics import confusion_matrix, accuracy_score, precision_score, recall_score, f1_score, precision_re

<>

[] !wget "https://drive.google.com/uc?export=download&id=1HC7MfMyg4Kph051bc21JKIxsubmc18C" -O orders.csv

--2022-05-04 07:24:45-- https://drive.google.com/uc?export=download&id=1HC7MfMyg4Kph051bc21JKIxsubmc18C

7/7

[2,4,5,8,1,5,12]

∞ imbalanced Data.ipynb - Colab × plot (-x*log(x)-(1-x)*log(1-x)) × Category Encoders - Category Encoders × + colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=dp3Jn34iJCBE

+ Code + Text RAM Disk

6s

Requirement already satisfied: threadpoolctl>=2.0.0 in /usr/local/lib/python3.7/dist-packages
Installing collected packages: category-encoders
Successfully installed category-encoders-2.4.0

{x}

```
[ ] import pandas as pd
import numpy as np
from numpy import argmax

✓ from datetime import date, time, timedelta
import pendulum # for time formatting

✓ import matplotlib.pyplot as plt
import seaborn as sns

import category_encoders as ce # for categorical encoding
from sklearn.preprocessing import LabelEncoder, StandardScaler

from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.metrics import confusion_matrix, accuracy_score, precision_score, recall_score, f1_score, precision_re
```

<>

```
[ ] !wget "https://drive.google.com/uc?export=download&id=1HC7MfMyg4Kph051bc21JKIxsubmc18C" -O orders.csv
```

--2022-05-04 07:24:45-- https://drive.google.com/uc?export=download&id=1HC7MfMyg4Kph051bc21JKIxsubmc18C

8/8

∞ imbalanced Data.ipynb - Colab × plot (-x*log(x)-(1-x)*log(1-x)) × Category Encoders - Category Encoders × + colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=dp3Jn34iJCBE

+ Code + Text RAM Disk

6s

Requirement already satisfied: threadpoolctl>=2.0.0 in /usr/local/lib/python3.7/dist-packages
Installing collected packages: category-encoders
Successfully installed category-encoders-2.4.0

{x}

```
[ ] import pandas as pd
import numpy as np
from numpy import argmax

from datetime import date, time, timedelta
import pendulum # for time formatting

import matplotlib.pyplot as plt
import seaborn as sns

import category_encoders as ce # for categorical encoding
from sklearn.preprocessing import LabelEncoder, StandardScaler
from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.metrics import confusion_matrix, accuracy_score, precision_score, recall_score, f1_score, precision_re
```

<>

```
[ ] !wget "https://drive.google.com/uc?export=download&id=1HC7MfMyg4Kph051bc21JKIxsubmc18C" -O orders.csv
```

--2022-05-04 07:24:45-- https://drive.google.com/uc?export=download&id=1HC7MfMyg4Kph051bc21JKIxsubmc18C

```
[ ] import pandas as pd
import numpy as np
from numpy import argmax

{x}
from datetime import date, time, timedelta
import pendulum # for time formatting

import matplotlib.pyplot as plt
import seaborn as sns

import category_encoders as ce # for categorical encoding
from sklearn.preprocessing import LabelEncoder, StandardScaler
from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.metrics import confusion_matrix, accuracy_score, precision_score, recall_score, f1_score, precision_re
```

```
[ ] !wget "https://drive.google.com/uc?export=download&id=1HC7MfMyg4Kph051bc21JKIxebmc18C" -O orders.csv
```

```
<> --2022-05-04 07:24:45-- https://drive.google.com/uc?export=download&id=1HC7MfMyg4Kph051bc21JKIxebmc18C_
Resolving drive.google.com (drive.google.com)... 172.253.122.138, 172.253.122.102, 172.253.122.100, ...
Connecting to drive.google.com (drive.google.com)|172.253.122.138|:443... connected.
HTTP request sent, awaiting response... 303 See Other
Location: https://doc-0k-14-docs.googleusercontent.com/docs/securesc/ha0ro937gcuc717deffksulhg5h7mbpl/8bsmgmlggd4mtl
```



∞ imbalanced Data.ipynb - Colab × plot (-x*log(x)-(1-x)*log(1-x)) × Category Encoders - Categorical Data × + colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=wN4r7EK88w10

+ Code + Text

RAM Disk ✓

[11] []

```
'amount', 'device', 'payment_id', 'customer_ip', 'customer_id',
'payment_method', 'payment_method_provider', 'payment_method_bin',
'payment_method_type', 'payment_method_product',
'payment_method_card_category', 'payment_method_issuer_bank',
'payment_method_issuer_country', 'is_fraudulent', 'time_diff',
'created_at_hod_sin', 'created_at_hod_cos', 'created_dom_sin',
'created_dom_cos', 'created_dow_sin', 'created_dow_cos',
'created_wom_sin', 'created_wom_cos', 'experience_dom_sin',
'experience_dom_cos', 'experience_dow_sin', 'experience_dow_cos',
'experience_wom_sin', 'experience_wom_cos'],
dtype='object')
```

df["is_fraudulent"].value_counts()

1 0
75817 1012

Name: is_fraudulent, dtype: int64

1 0
1:75

```
[ ] df.unique() # unique values per feature
```

	order_id	city	category_name	product_id	product_name
1	76829	2	96	353	340

order_id 76829
city 2
category_name 96
product_id 353
product_name 340

imbalanced Data.ipynb - Colab plot (-x*log(x)-(1-x)*log(1-x)) Category Encoders - Catego + colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=wN4r7EK88w10 Update

+ Code + Text

RAM Disk

[11] []

```
'amount', 'device', 'payment_id', 'customer_ip', 'customer_id',
'payment_method', 'payment_method_provider', 'payment_method_bin',
'payment_method_type', 'payment_method_product',
'payment_method_card_category', 'payment_method_issuer_bank',
'payment_method_issuer_country', 'is_fraudulent', 'time_diff',
'created_at_hod_sin', 'created_at_hod_cos', 'created_dom_sin',
'created_dom_cos', 'created_dow_sin', 'created_dow_cos',
'created_wom_sin', 'created_wom_cos', 'experience_dom_sin',
'experience_dom_cos', 'experience_dow_sin', 'experience_dow_cos',
'experience_wom_sin', 'experience_wom_cos'],
dtype='object')
```

df["is_fraudulent"].value_counts()

False 75817
True 1012
Name: is_fraudulent, dtype: int64

[] df.unique() # unique values per feature

order_id	76829
city	2
category_name	96
product_id	353
product_name	340

12/12

imbalanced Data.ipynb - Colab plot (-x*log(x)-(1-x)*log(1-x)) Category Encoders - Categorical News - Latest News, Breaking + colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=wN4r7EK88w10 Update

+ Code + Text RAM Disk

[11] [✓] [RAM] [Disk]

```
'amount', 'device', 'payment_id', 'customer_ip', 'customer_id',
'payment_method', 'payment_method_provider', 'payment_method_bin',
'payment_method_type', 'payment_method_product',
'payment_method_card_category', 'payment_method_issuer_bank',
'payment_method_issuer_country', 'is_fraudulent', 'time_diff',
'created_at_hod_sin', 'created_at_hod_cos', 'created_dom_sin',
'created_dom_cos', 'created_dow_sin', 'created_dow_cos',
'created_wom_sin', 'created_wom_cos', 'experience_dom_sin',
'experience_dom_cos', 'experience_dow_sin', 'experience_dow_cos',
'experience_wom_sin', 'experience_wom_cos'],
dtype='object')
```

df["is_fraudulent"].value_counts()

False 75817
True 1012
Name: is_fraudulent, dtype: int64

[] df.unique() # unique values per feature

	order_id	city	category_name	product_id	product_name
1	76829	2	96	353	340

1: 1000 ✓

13 / 13

imbalanced Data.ipynb - Colab plot (-x*log(x)-(1-x)*log(1-x)) Category Encoders - Categorical News - Latest News, Breaking + colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=wN4r7EK88w10 Update :+ Code + Text ✓ RAM Disk

[11] [+] amount, 'device', 'payment_id', 'customer_ip', 'customer_id',
payment_method, 'payment_method_provider', 'payment_method_bin',
'payment_method_type', 'payment_method_product',
'payment_method_card_category', 'payment_method_issuer_bank',
'payment_method_issuer_country', 'is_fraudulent', 'time_diff',
'created_at_hod_sin', 'created_at_hod_cos', 'created_dom_sin',
'created_dom_cos', 'created_dow_sin', 'created_dow_cos',
'created_wom_sin', 'created_wom_cos', 'experience_dom_sin',
'experience_dom_cos', 'experience_dow_sin', 'experience_dow_cos',
'experience_wom_sin', 'experience_wom_cos'),
dtype='object')

df["is_fraudulent"].value_counts()

False 75817
True 1012
Name: is_fraudulent, dtype: int64

[] df.unique() # unique values per feature

order_id 76829
city 2
category_name 96
product_id 353
product_name 340

∞ imbalanced Data.ipynb - Colab × plot (-x*log(x)-(1-x)*log(1-x)) × Category Encoders - Categorical × News - Latest News, Breaking × +

colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=3tU0JQJ8KKsC

+ Code + Text RAM Disk

[] df["is_fraudulent"].value_counts()

{x}

False 75817
True 1012
Name: is_fraudulent, dtype: int64

df.unique() # unique values per feature

order_id	76829
city	2
category_name	96
product_id	353
product_name	340
amount	8627
device	5
payment_id	76829
customer_ip	44208
customer_id	51109
payment_method	1
payment_method_provider	5
payment_method_bin	6991
payment_method_type	8
payment_method_product	144
payment_method_card_category	2

order_id

Y_i = 1 or 0

✓

15 / 15

∞ imbalanced Data.ipynb - Colab × Google plot (-x*log(x)-(1-x)*log(1-x)) × Category Encoders - Categorical × News - Latest News, Breaking × +

colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=3tU0JQJ8KKsC

+ Code + Text RAM Disk  

```
[ ] df["is_fraudulent"].value_counts()
```

{x} False 75817
True 1012
Name: is_fraudulent, dtype: int64

df.unique() # unique values per feature

card	76829
city	2
category_name	96
product_id	353
product_name	340
amount	8627
device	5
payment_id	76829
customer_ip	44208
customer_id	51109
payment_method	1
payment_method_provider	5
payment_method_bin	6991
payment_method_type	8
payment_method_product	144
payment_method_card_category	2
payment_method_issuer_bank	2046

RAM Disk  

Up Down Reload Settings Copy Paste Delete More

16 / 16

∞ imbalanced Data.ipynb - Colab × Google plot (-x*log(x)-(1-x)*log(1-x)) × Category Encoders - Categorical Data × News - Latest News, Breaking × +

colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=3tU0JQJ8KKsC

+ Code + Text

[] df["is_fraudulent"].value_counts()

False 75817
True 1012
Name: is_fraudulent, dtype: int64

{x}

df.unique() # unique values per feature

order_id 76829
city 2
category_name 96
product_id 353
product_name 340
amount 8627
device 5
payment_id 76829
customer_ip 44208
customer_id 51109
payment_method 1
payment_method_provider 5
payment_method_bin 6991
payment_method_type 8
payment_method_product 144
payment_method_card_category 2
payment_method_issuer_bank 2046

RAM Disk

Update

Up Down Reload Settings Copy Copy All Delete More

0s

17 / 17

+ Code + Text

✓ RAM
Disk



```
[ ] df["is_fraudulent"].value_counts()
```

```
False    75817
True     1012
Name: is_fraudulent, dtype: int64
```



✓
Us

```
df.nunique() # unique values per feature
```

```
order_id                76829
city                     2
category_name             96
product_id                353
product_name               340
amount                   8627
device                    5
payment_id                76829
customer_ip               44208
customer_id                51109
payment_method              14
payment_method_provider           5
payment_method_bin                 6991
payment_method_type                  8
payment_method_product                144
payment_method_card_category          2
payment_method_issuer_bank            2046
```

∞ imbalanced Data.ipynb - Colab × plot (-x*log(x)-(1-x)*log(1-x)) × Category Encoders - Categorical Data × News - Latest News, Breaking × +

colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=3tU0JQJ8KKsC

+ Code + Text

RAM Disk

True 1012
Name: is_fraudulent, dtype: int64

↑ ↓ ⌂ ⚙️ 📁 🗑️ :

{x} ✓

df.nunique() # unique values per feature

order_id	76829
city	2
category_name	96
product_id	353
product_name	340
amount	8627
device	5
payment_id	76829
customer_ip	44208
customer_id	51109
payment_method	1
payment_method_provider	5
payment_method_bin	6991
payment_method_type	8
payment_method_product	144
payment_method_card_category	2
payment_method_issuer_bank	2046
payment_method_issuer_country	151
is_fraudulent	2
time_diff	117

353

340

19 / 19

colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=3tU0JQJ8KKsC

+ Code + Text

RAM Disk

order_id 76829
city 2
category_name 96
product_id 353
product_name 340
amount 8627
device 5
payment_id 76829
customer_ip 44208
customer_id 51109
payment_method 1
payment_method_provider 5
payment_method_bin 6991
payment_method_type 8
payment_method_product 144
payment_method_card_category 2
payment_method_issuer_bank 2046
payment_method_issuer_country 151
is_fraudulent 2
time_diff 117
created_at_hod_sin 24
created_at_hod_cos 13
created_dom_sin 31
created_dom_cos 18
created_dow_sin 6
created_dow_cos 4
created_wom_sin 5
created_wom_cos 1

encrypted- P

 $P_i \rightarrow 20 \rightarrow 5$

imbalanced Data.ipynb - Colab plot (-x*log(x)-(1-x)*log(1-x)) Category Encoders - Categorical News - Latest News, Breaking

colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=3tU0JQJ8KKsC

+ Code + Text

RAM Disk

order_id 76829
city 2
category_name 96
product_id 353
product_name 340
amount 8627
device 5
payment_id 76829
customer_ip 44208
customer_id 51109
payment_method 1
payment_method_provider 5
payment_method_bin 6991
payment_method_type 8
payment_method_product 144
payment_method_card_category 2
payment_method_issuer_bank 2046
payment_method_issuer_country 151
is_fraudulent 2
time_diff 117
created_at_hod_sin 24
created_at_hod_cos 13
created_dom_sin 31
created_dom_cos 18
created_dow_sin 6
created_dow_cos 4
created_wom_sin 5
created_wom_cos 1

ip by freq > 10
ip = other

Chrome File Edit View History Bookmarks Profiles Tab Window Help

imbalanced Data.ipynb - Colab plot (-x*log(x)-(1-x)*log(1-x)) Category Encoders - Categorical News - Latest News, Breaking + colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=3tU0JQJ8KKsC

+ Code + Text RAM Disk Update

order_id 76829
city 2
category_name 96
product_id 353
product_name 340
amount 8627
device 5
payment_id 76829
customer_ip 44208
customer_id 51109
payment_method 1
payment_method_provider 5
payment_method_bin 6991
payment_method_type 8
payment_method_product 144
payment_method_card_category 2
payment_method_issuer_bank 2046
payment_method_issuer_country 151
is_fraudulent 2
time_diff 117
created_at_hod_sin 24
created_at_hod_cos 13
created_dom_sin 31
created_dom_cos 18
created_dow_sin 6
created_dow_cos 4
created_wom_sin 5
created_wom_cos 1

regularize $W \rightarrow 0$

RAM Disk

Up Down Left Right Home Stop Refresh Back Forward

22 / 22

Chrome File Edit View History Bookmarks Profiles Tab Window Help

colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=3tU0JQJ8KKsC

+ Code + Text

RAM Disk

order_id 76829
city 2
category_name 96
product_id 353
{x} product_name 340
amount 8627
device 5
payment_id 76829
customer_ip 44208 → IP
customer_id 51109
payment_method 1
payment_method_provider 5
payment_method_bin 6991
payment_method_type 8
payment_method_product 144
payment_method_card_category 2
payment_method_issuer_bank 2046
payment_method_issuer_country 151
is_fraudulent 2
time_diff 117
created_at_hod_sin 24
created_at_hod_cos 13
created_dom_sin 31
created_dom_cos 18
created_dow_sin 6
created_dow_cos 4
created_wom_sin 5
created_wom_cos 1

Country:
State
City
Zipcode

23 / 23

imbalanced Data.ipynb - Colab plot (-x*log(x)-(1-x)*log(1-x)) Category Encoders - Catego... News - Latest News, Breaking...

colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=R79B_o5dLPOR

+ Code + Text

RAM Disk

[14]

{x}

[15] df.shape
(76829, 30)

drop duplicates
if df.shape[0] == df.drop_duplicates().shape[0] :
 print('No duplicates Found')
else:
 duplicates = df.shape[0] - df.drop_duplicates().shape[0]
 print('{} duplicates found'.format(duplicates))

6821 duplicates found

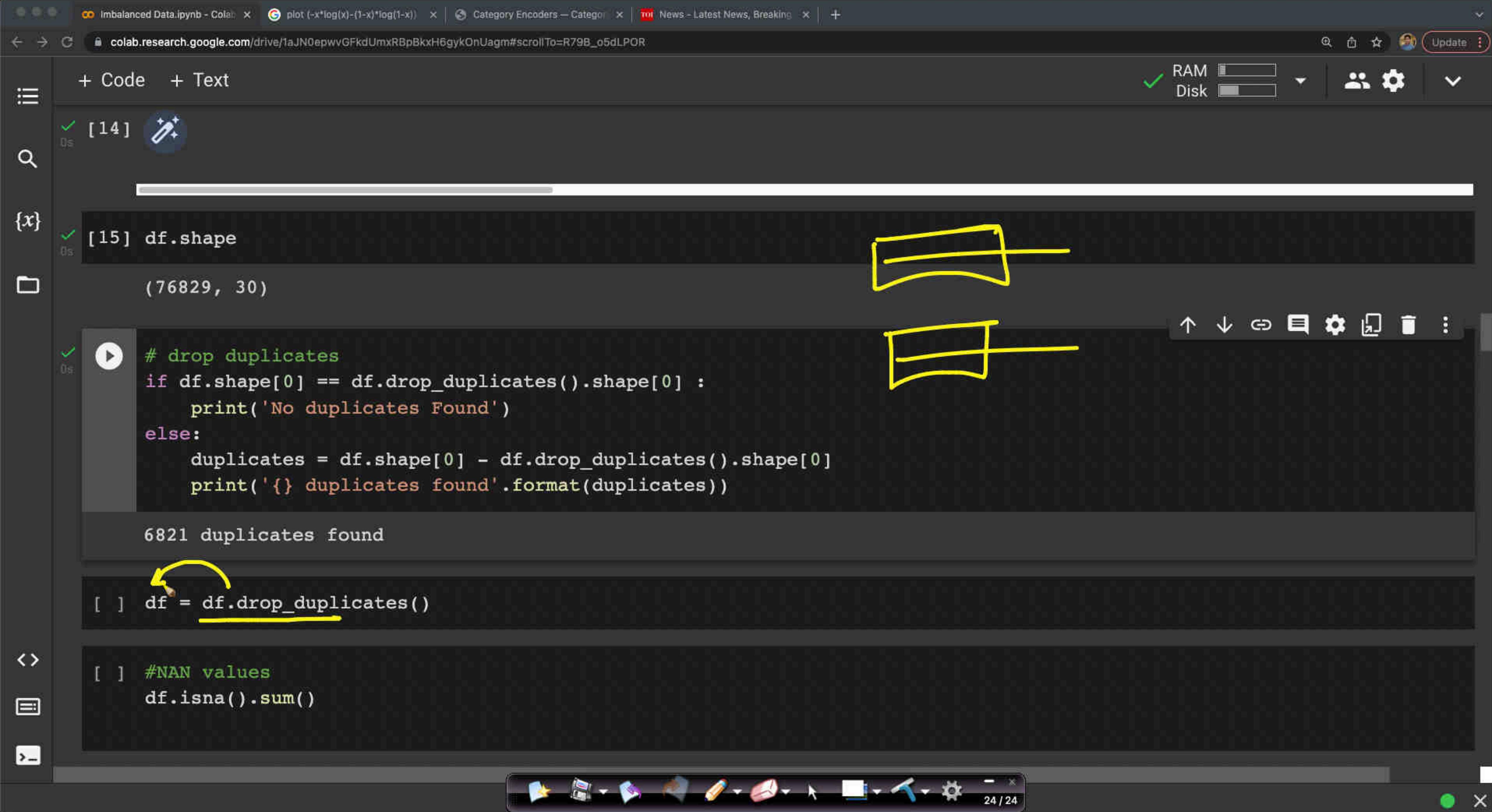
[] df = df.drop_duplicates()

[] #NAN values
df.isna().sum()

RAM Disk

Up Down Reload Settings Copy Clear More

24 / 24



∞ imbalanced Data.ipynb - Colab × Google plot (-x*log(x)-(1-x)*log(1-x)) × Category Encoders - Categorical × News - Latest News, Breaking × +

colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=R79B_o5dLPOR: Update :

+ Code + Text RAM Disk

```
[ ] #NAN values
df.isna().sum()
```

{x}

city	0
category_name	83
product_id	0
amount	0
device	0
customer_ip	0
customer_id	0
payment_method_provider	0
payment_method_bin	46
payment_method_type	61
payment_method_product	1886
payment_method_card_category	2413
payment_method_issuer_bank	2838
payment_method_issuer_country	66
is_fraudulent	0
time_diff	0
created_at_hod_sin	0
created_at_hod_cos	0
created_dom_sin	0
created_dom_cos	0
created_dow_sin	0
created_dow_cos	0
created_wom_sin	0

<>

25 / 25

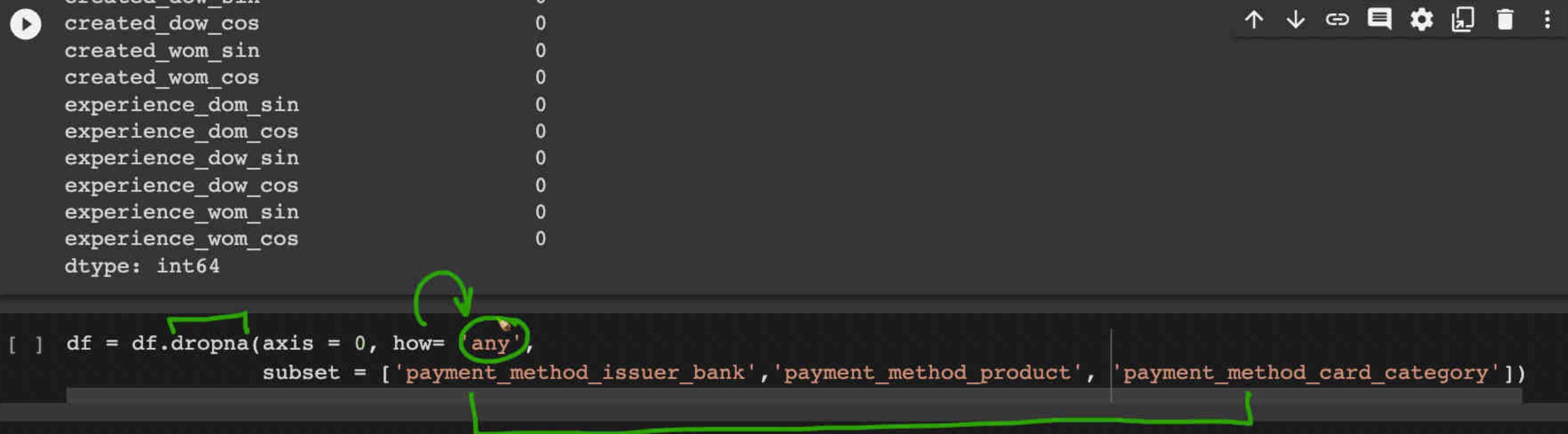
Imbalanced Data.ipynb - Colab | plot (-x*log(x)-(1-x)*log(1-x)) | Category Encoders — Categorical Data Handling | News - Latest News, Breaking

Code + Text



```
created_aow_sin
created_dow_cos
created_wom_sin
created_wom_cos
experience_dom_sin
experience_dom_cos
experience_dow_sin
experience_dow_cos
experience_wom_sin
experience_wom_cos
dtype: int64
```

RAM Disk



city
category_name
product_id
amount
device
customer_ip
customer_id
payment method provider

Imbalanced Data.ipynb - Colab | plot (-x*log(x)-(1-x)*log(1-x)) | Category Encoders — Categorical Data Handling | News - Latest News, Breaking

Code + Text

```
subset = ['payment_method_issuer_bank', 'payment_method_product', 'payment_method_card_category'])
```

df.isna().sum()

city	0
category_name	75
product_id	0
amount	0
device	0
customer_ip	0
customer_id	0
payment_method_provider	0
payment_method_bin	0
payment_method_type	0
payment_method_product	0
payment_method_card_category	0
payment_method_issuer_bank	0
payment_method_issuer_country	11
is_fraudulent	0
time_diff	0
created_at_hod_sin	0
created_at_hod_cos	0
created_dom_sin	0
created_dom_cos	0
created_dow_sin	0
created_dow_cos	0



 Imbalanced Data.ipynb - Colab | plot (-x*log(x)-(1-x)*log(1-x)) | Category Encoders – Category | News - Latest News, Breaking | +

Code + Text

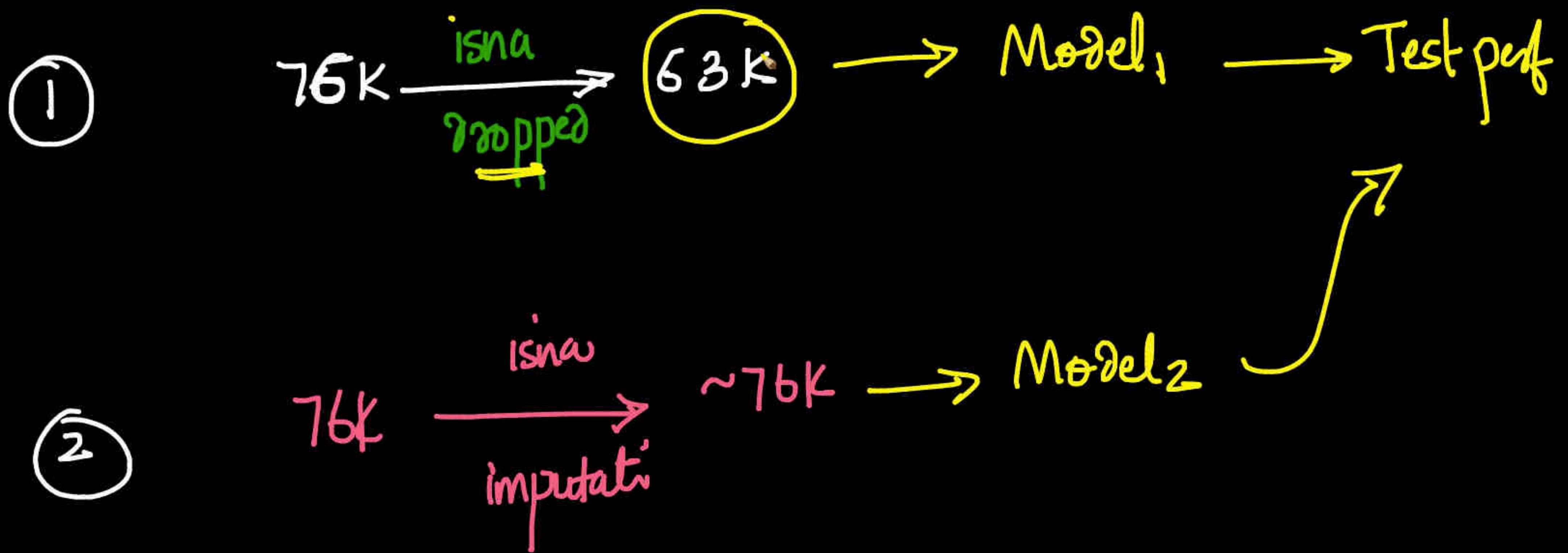
[22]

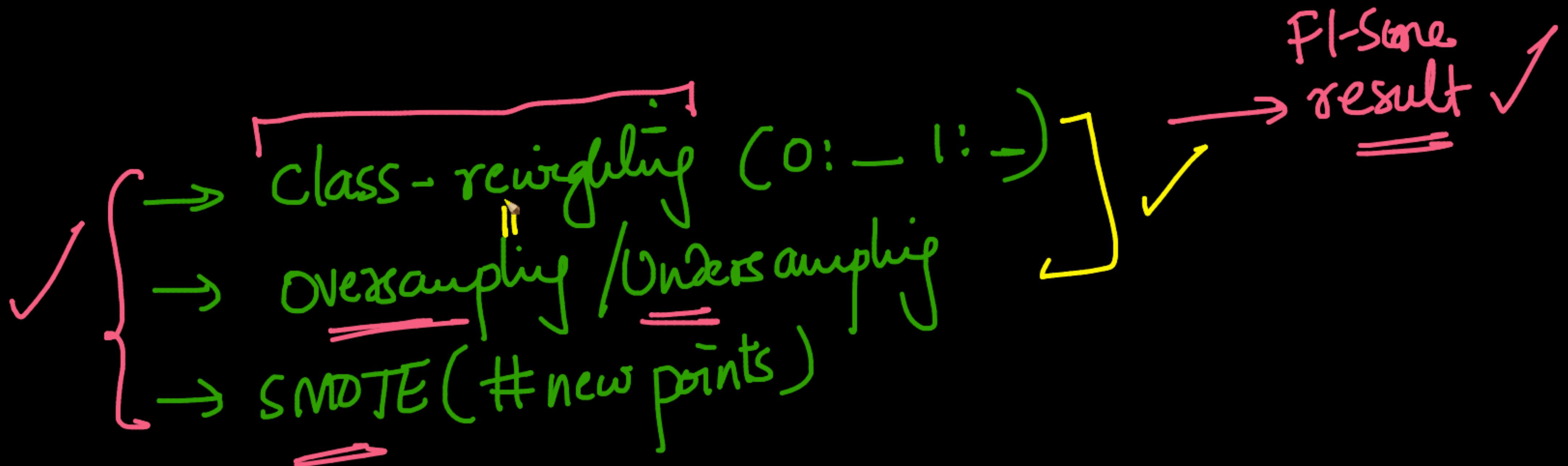
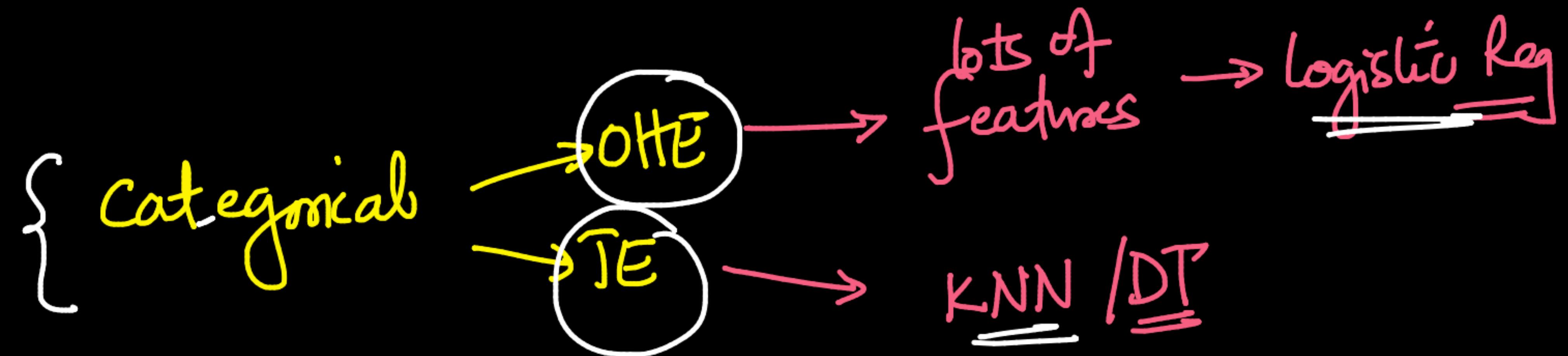
✓ ⏪ df.isna().sum()

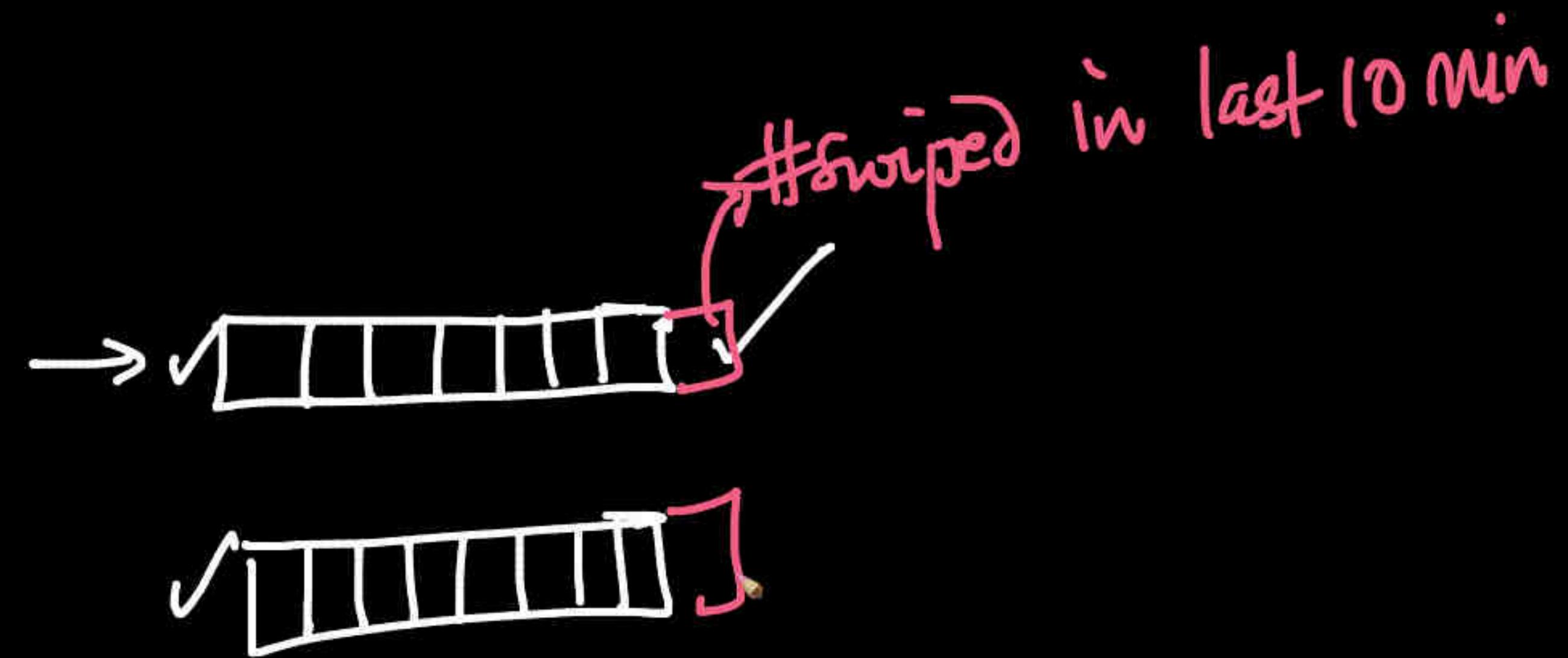
city 0
category_name 0
product_id 0
amount 0
device 0
customer_ip 0
customer_id 0
payment_method_provider 0
payment_method_bin 0
payment_method_type 0
payment_method_product 0
payment_method_card_category 0
payment_method_issuer_bank 0
payment_method_issuer_country 0
is_fraudulent 0
time_diff 0
created_at_hod_sin 0
created_at_hod_cos 0
created_dom_sin 0
created_dom_cos 0
created_dow_sin 0
created_dow_cos 0
created_wom_sin 0

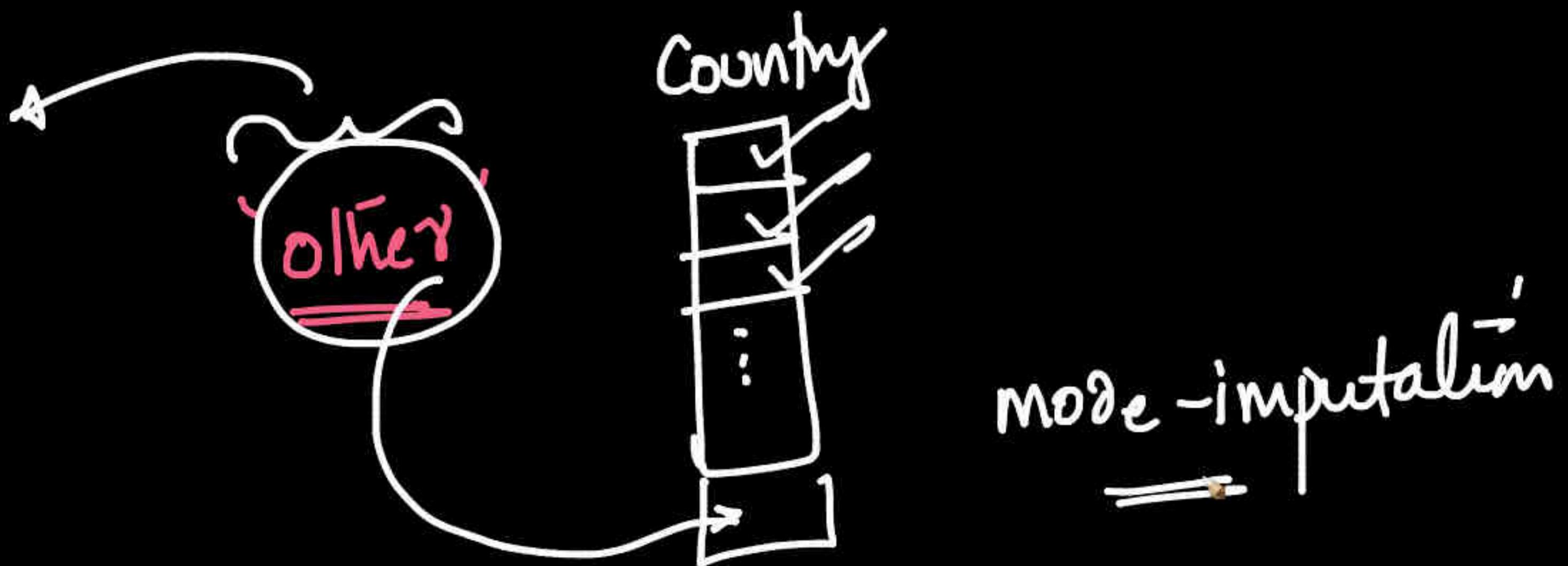
✓ *dop*
✓ *isna* → *imputalím*

75K → 63K







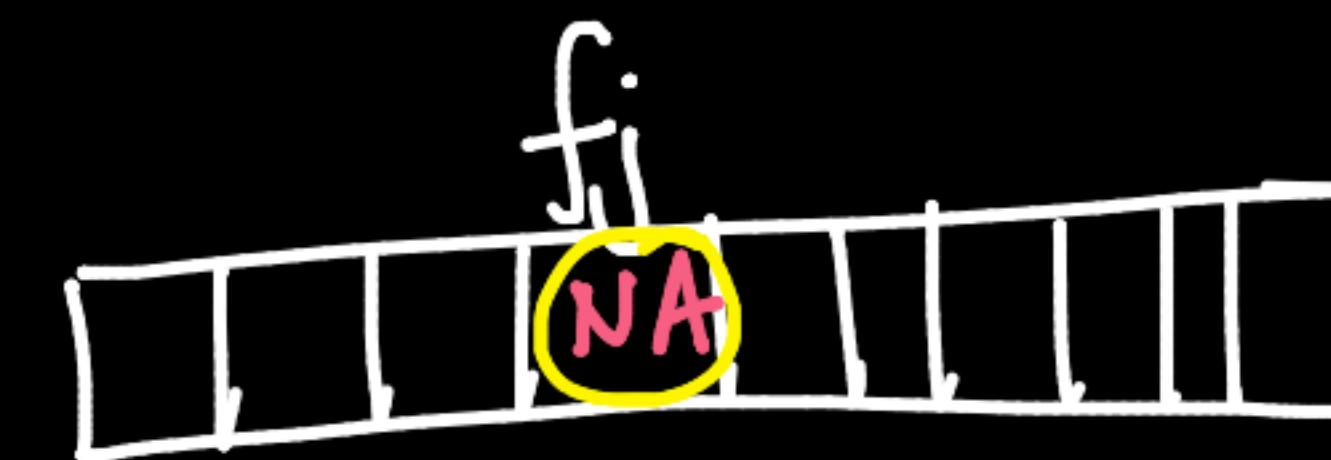




Test-time

x_{qj} :

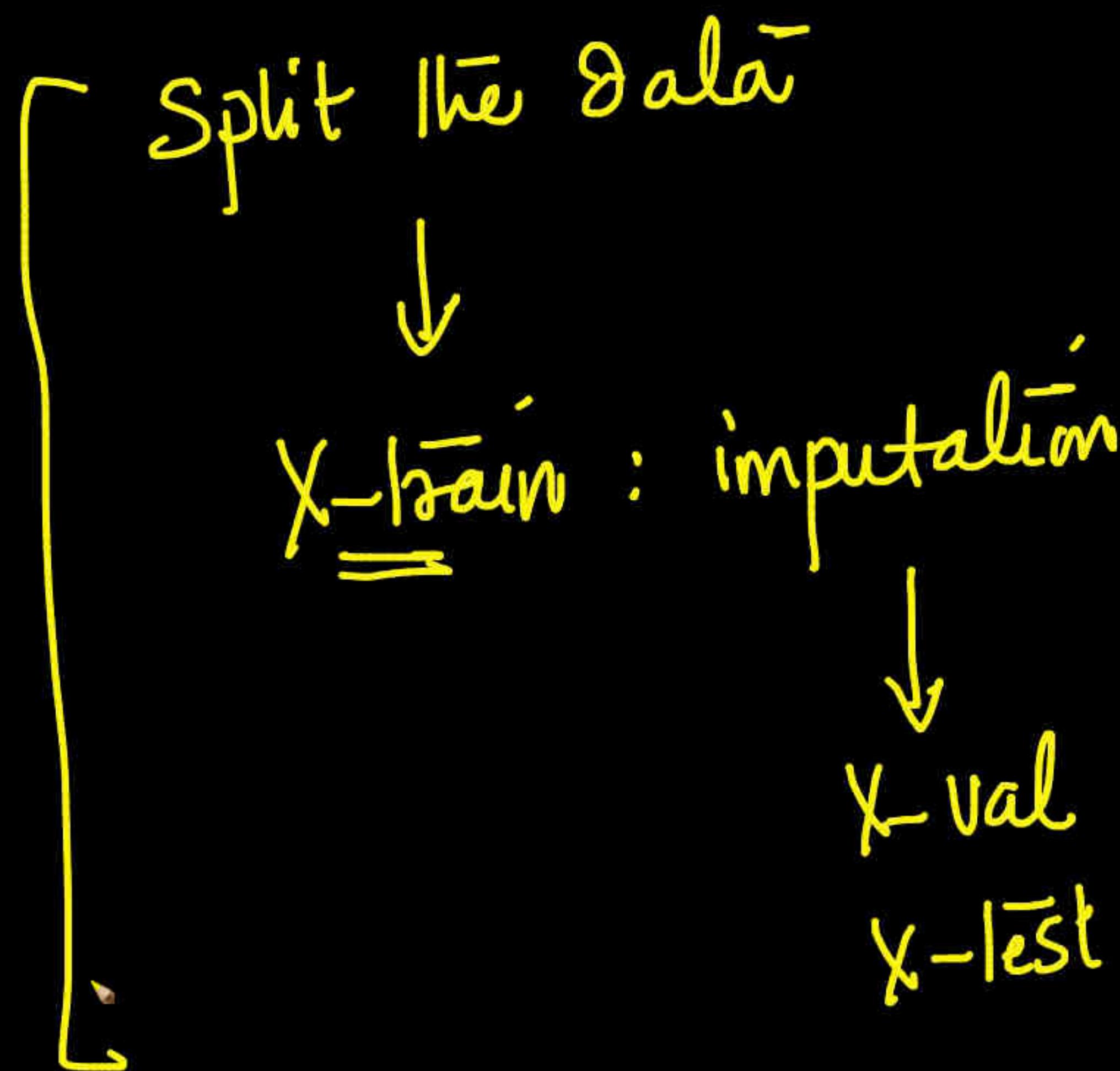
$y_{qj}:$

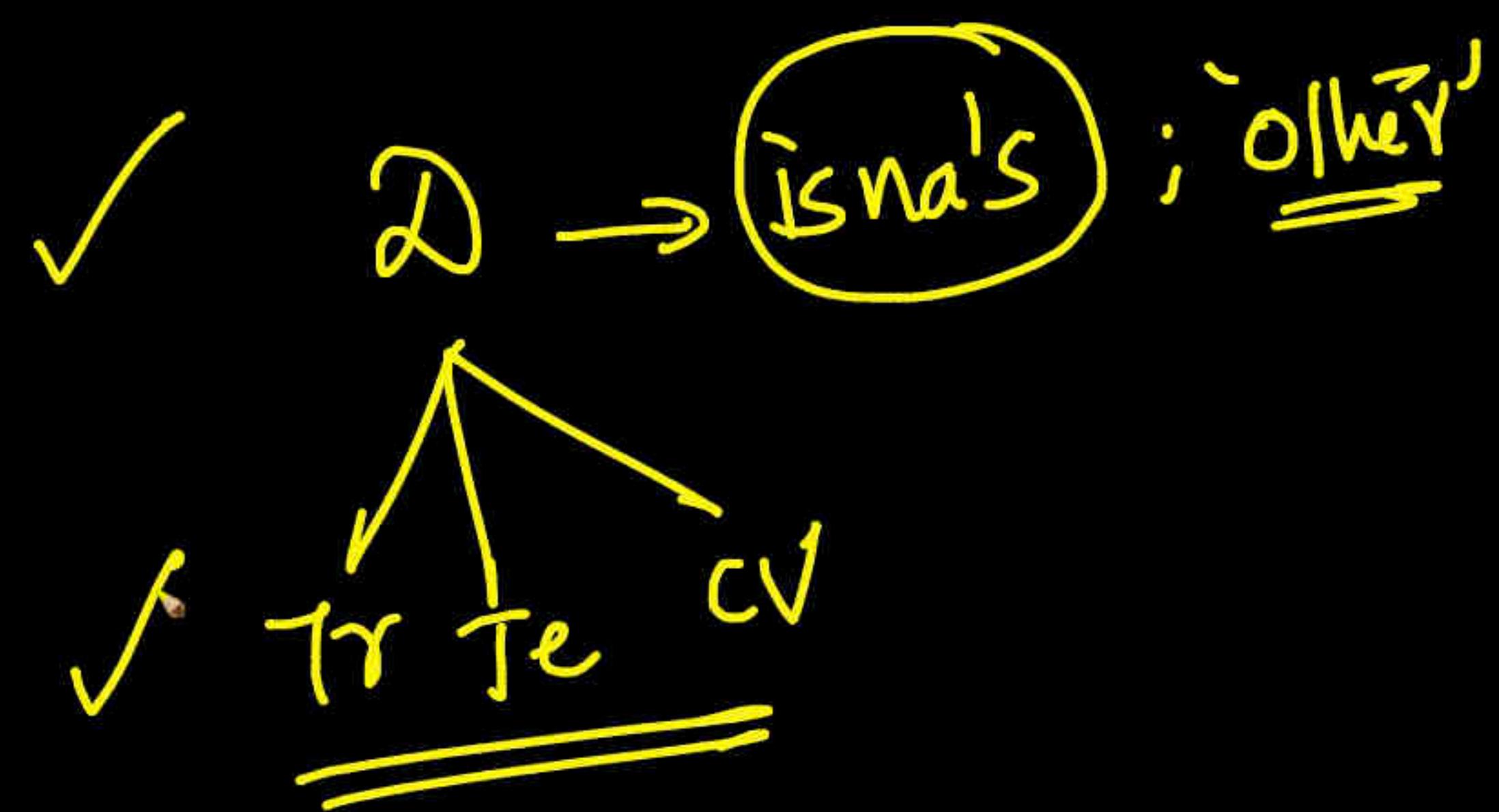


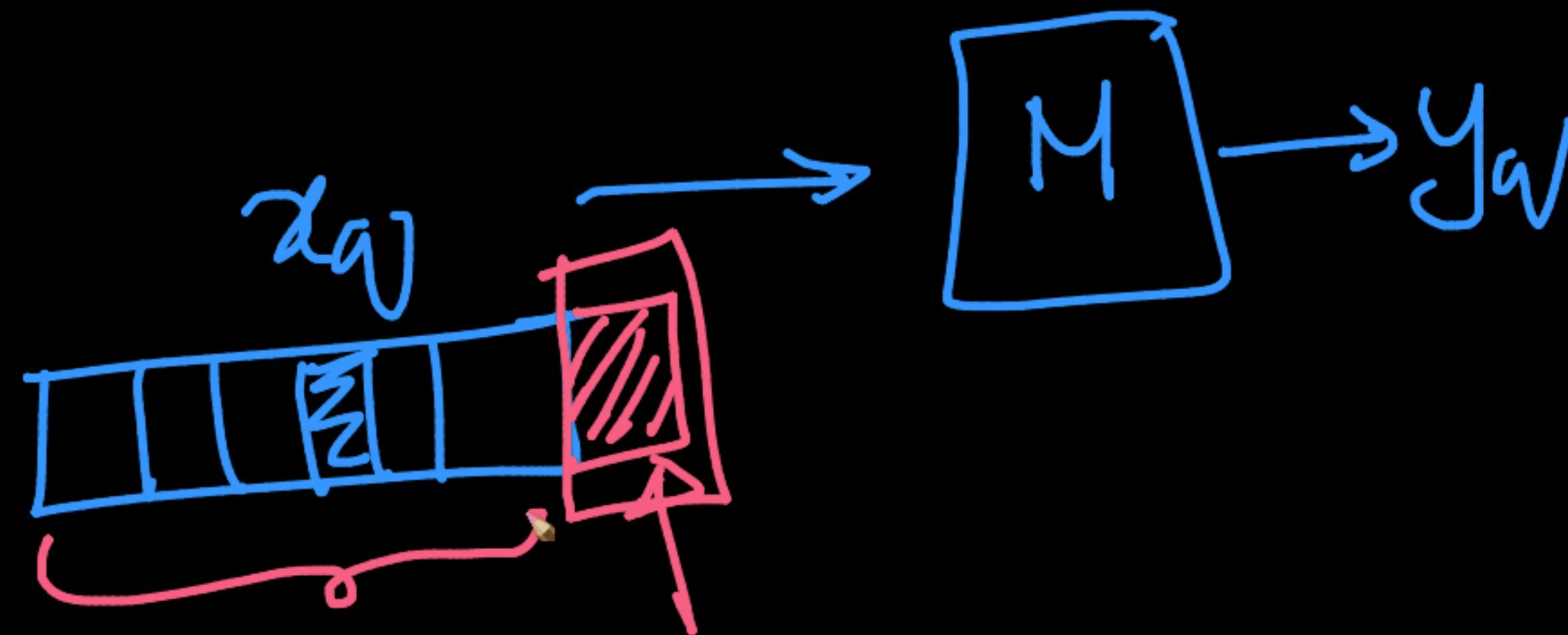
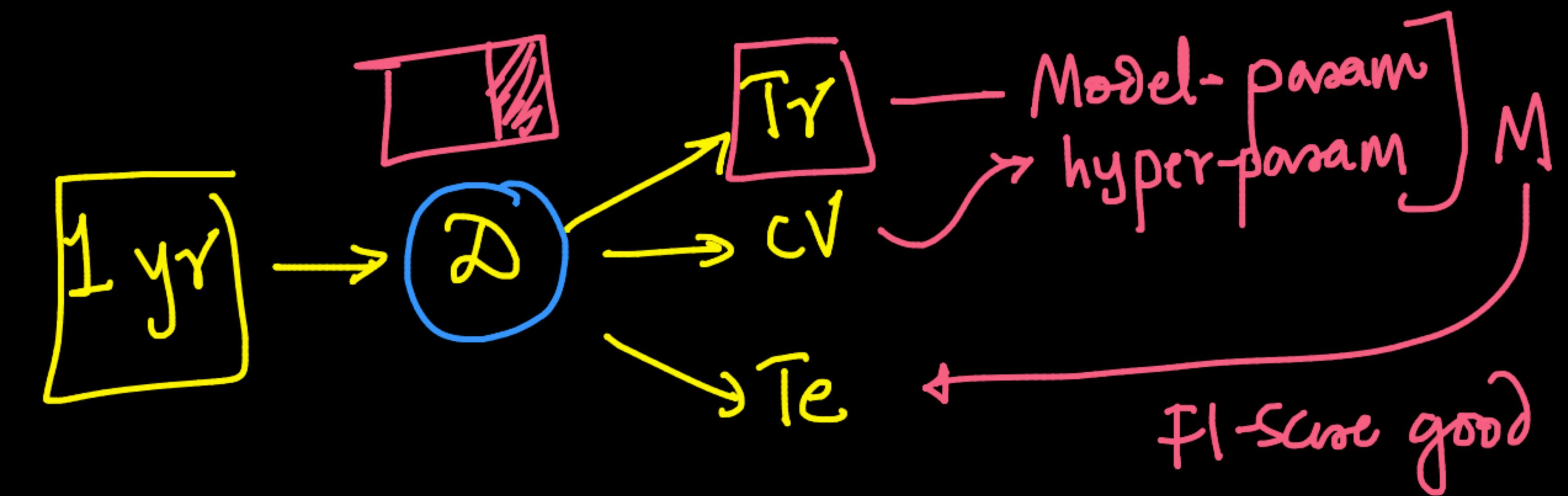
X
ISNA
same, imputation

f_j : numeric \rightarrow mean/
median

f_j : Categorical \rightarrow Mode







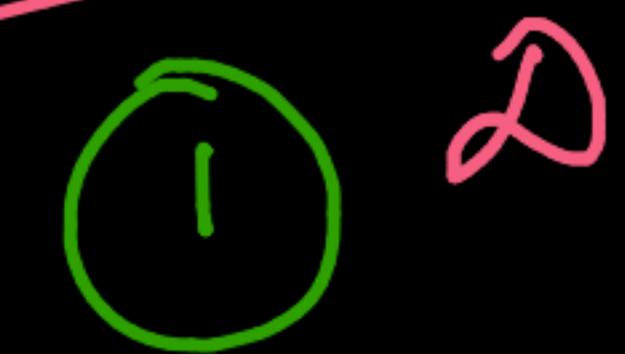
1

don't  in production
→ isn't drop

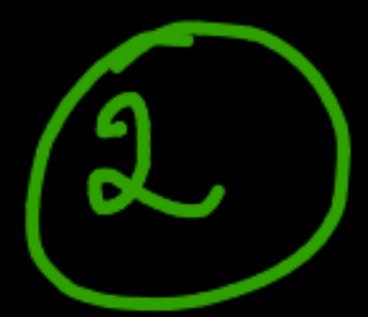
2

NANs in production
→ imputation (Mean/Mode
Median)

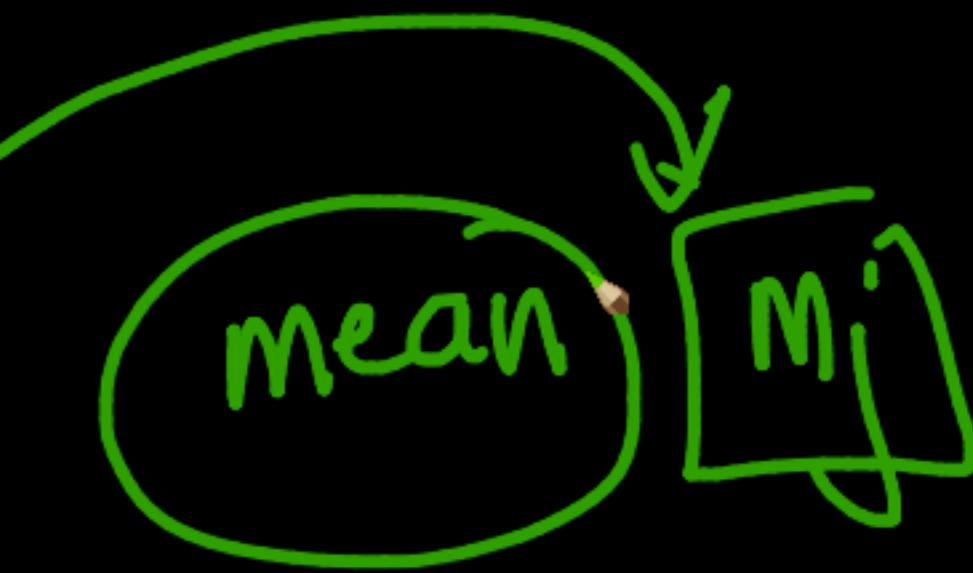
~~Imputation:~~



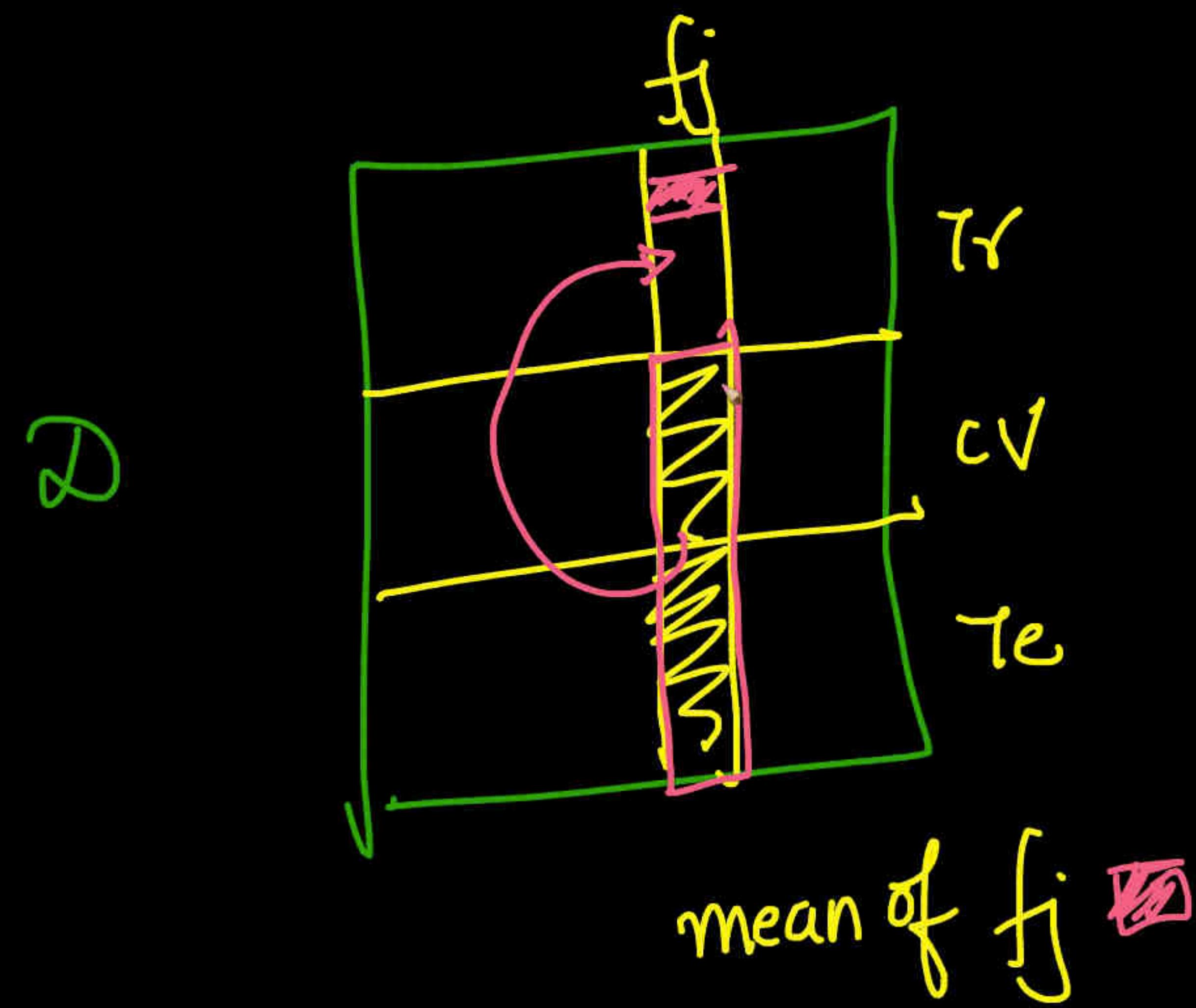
D_{TR}, D_{LE}, D_{CV}



Impulé using D_{TR}



CV: { m/s }
Test



∞ imbalanced Data.ipynb - Colab | plot (-x*log(x)-(1-x)*log(1-x)) | Category Encoders - Categorical Data | News - Latest News, Breaking | +

colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=nJlxUT34N9F

+ Code + Text

RAM Disk

Update

S -> Sn

{x}

```
[ ] set(X.columns) - set(X.select_dtypes(['number']).columns)
```

category_name
city
customer_id
customer_ip
device
payment_method_card_category
payment_method_issuer_bank
payment_method_issuer_country
payment_method_product
payment_method_provider
payment_method_type

[] # Target encoding using mean

ce_target = ce.TargetEncoder(cols = ['category_name', 'city', 'customer_id', 'customer_ip', 'device', 'payment_method_ca
X_train = ce_target.fit_transform(X_train, y_train)

[] X_val = ce_target.transform(X_val)
X test = ce target.transform(X test)

+ Code + Text

✓ RAM
Disk



```
[ ] 'customer_id',
[ ] 'customer_ip',
[ ] 'device',
[ ] 'payment_method_card_category',
[ ] 'payment_method_issuer_bank',
[ ] 'payment_method_issuer_country',
[ ] 'payment_method_product',
[ ] 'payment_method_provider',
[ ] 'payment_method_type'}
```

```
[ ] # Target encoding using mean
[ ]
ce_target = ce.TargetEncoder(cols = ['category_name', 'city', 'customer_id', 'customer_ip', 'device', 'payment_method_ca
x_train = ce_target.fit_transform(X_train, y_train)
```

```
[ ] X_val = ce_target.transform(X_val)
[ ] X_test = ce_target.transform(X_test)
```

```
[ ] # Hyper-param tuning
[ ] from sklearn.linear_model import LogisticRegression
[ ] from sklearn.pipeline import make_pipeline
[ ] from sklearn.metrics import f1_score
```

```
[ ] train_scores = []
```

∞ imbalanced Data.ipynb - Colab × plot (-x*log(x)-(1-x)*log(1-x)) × Category Encoders - Categorical × News - Latest News, Breaking × + colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=nJlxUT34N9F

+ Code + Text

RAM Disk

RAM Disk

customer_id',
customer_ip',
device',
payment_method_card_category',
payment_method_issuer_bank',
payment_method_issuer_country',
payment_method_product',
payment_method_provider',
payment_method_type'}

Non-numeric

Target encoding using mean

ce_target = ce.TargetEncoder(cols = ['category_name', 'city', 'customer_id', 'customer_ip', 'device', 'payment_method_card_category'])

x_train = ce_target.fit_transform(X_train, y_train)

X_val = ce_target.transform(X_val)
X_test = ce_target.transform(X_test)

Hyper-param tuning

from sklearn.linear_model import LogisticRegression
from sklearn.pipeline import make_pipeline
from sklearn.metrics import f1_score

train_scores = []

+ Code + Text

```
[ ] 'customer_id',
'customer_ip',
'device',
'payment_method_card_category',
'payment_method_issuer_bank',
'payment_method_issuer_country',
'payment_method_product',
'payment_method_provider',
'payment_method_type'
```

$$P(y_i=1 \mid \text{city} = \text{Singapore}) = \frac{1000}{20,000}$$

city : Singapore : 1000

$$\text{city} : \cancel{\text{Singapore}} : \cancel{1000}$$

$$\text{Dubai} : \cancel{30,000} = P(y_i=1 \mid \text{city} = \text{Dubai})$$

Dubai : 30,000

```
[ ] # Target encoding using mean
#
ce_target = ce.TargetEncoder(cols = ['category_name', 'city', 'customer_id', 'customer_ip', 'device', 'payment_method_ca
x_train = ce_target.fit_transform(X_train, y_train)
```

```
[ ] X_val = ce_target.transform(X_val)
X_test = ce_target.transform(X_test)
```

```
[ ] # Hyper-param tuning
from sklearn.linear_model import LogisticRegression
from sklearn.pipeline import make_pipeline
from sklearn.metrics import f1_score
```

```
train_scores = []
```

f_j : categorical features

$$c_1 = \frac{300}{1000} = 0.3$$

c-1



$$P(y=1 | f_j=c_2) = c_2 = \frac{4+\alpha}{10+\beta} = 0.4$$



:

c_{10}

$$\frac{1}{75} \sim P(y=1)$$

Smoothing
strategies

$$C_1: \frac{\overbrace{300+\alpha}^1}{\overbrace{1000+\beta}^{75}} = \frac{301}{1075} \approx \frac{300}{1000}$$

$$C_2: \frac{\overbrace{4+\alpha}^1}{\overbrace{10+\beta}^{75}} = \frac{5}{85} \neq \frac{4}{10}$$

contrib.scikit-learn.org/category_encoders/targetencoder.html

- Backward Difference Coding
- BaseN
- Binary
- CatBoost Encoder
- Count Encoder
- Generalized Linear Mixed Model Encoder
- Hashing
- Helmert Coding
- James-Stein Encoder
- Leave One Out
- M-estimate
- One Hot
- Ordinal
- Polynomial Coding
- Sum Coding

Target Encoder

- Weight of Evidence
- Wrappers
- Quantile Encoder
- Summary Encoder

return_df: bool

boolean for whether to return a pandas DataFrame from transform (otherwise it will be a numpy array).

handle_missing: str

options are 'error', 'return_nan' and 'value', defaults to 'value', which returns the target mean.

handle_unknown: str

options are 'error', 'return_nan' and 'value', defaults to 'value', which returns the target mean.

min_samples_leaf: int

minimum samples to take category average into account.

smoothing: float

smoothing effect to balance categorical average vs prior. Higher value means stronger regularization. The value must be strictly bigger than 0.

References

1

A Preprocessing Scheme for High-Cardinality Categorical Attributes in Classification and Prediction Problems, from

[Backward Difference Coding](#)[BaseN](#)[Binary](#)[CatBoost Encoder](#)[Count Encoder](#)[Generalized Linear Mixed Model Encoder](#)[Hashing](#)[Helmert Coding](#)[James-Stein Encoder](#)[Leave One Out](#)[M-estimate](#)[One Hot](#)[Ordinal](#)[Polynomial Coding](#)[Sum Coding](#)[Target Encoder](#)[Weight of Evidence](#)[Wrappers](#)[Quantile Encoder](#)[Summary Encoder](#)**return_df: bool**

boolean for whether to return a pandas DataFrame from transform (otherwise it will be a numpy array).

handle_missing: str

options are 'error', 'return_nan' and 'value', defaults to 'value', which returns the target mean.

handle_unknown: str

options are 'error', 'return_nan' and 'value', defaults to 'value', which returns the target mean.

min_samples_leaf: int

minimum samples to take category average into account.

smoothing: float

smoothing effect to balance categorical average vs prior. Higher value means stronger regularization. The value must be strictly bigger than 0.

References

1

A Preprocessing Scheme for High-Cardinality Categorical Attributes in Classification and Prediction Problems, from

∞ imbalanced Data.ipynb - Colab x plot (-x*log(x)-(1-x)*log(1-x)) x Category Encoders - Catego x Target Encoder - Category En x +

colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=EdTxI8iPNcXQ

+ Code + Text

RAM Disk

RAM Disk

[27]

```
device',
'payment_method_card_category',
'payment_method_issuer_bank',
'payment_method_issuer_country',
'payment_method_product',
'payment_method_provider',
'payment_method_type'}
```

Target encoding using mean

ce_target = ce.TargetEncoder(cols = ['category_name','city','customer_id','customer_ip','device','payment_method_ca
x_train = ce_target.fit_transform(X_train, y_train)

[] X_val = ce_target.transform(X_val)
X_test = ce_target.transform(X_test)

[] # Hyper-param tuning
from sklearn.linear_model import LogisticRegression
from sklearn.pipeline import make_pipeline
from sklearn.metrics import f1_score

train_scores = []
val_scores = []
scaler = StandardScaler()

48 / 48

∞ imbalanced Data.ipynb - Colab plot (-x*log(x)-(1-x)*log(1-x)) Category Encoders - Category En... Target Encoder - Category En... +

colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=EdTxI8iPNcXQ

+ Code + Text

RAM Disk

device',
[27] 'payment_method_card_category',
'payment_method_issuer_bank',
'payment_method_issuer_country',
'payment_method_product',
'payment_method_provider',
'payment_method_type'}

Target encoding using mean

ce_target = ce.TargetEncoder(cols = ['category_name','city','customer_id','customer_ip','device','payment_method_ca
X_train = ce_target.fit_transform(X_train, y_train)

[] X_val = ce_target.transform(X_val) ✓
X_test = ce_target.transform(X_test)

[] # Hyper-param tuning
from sklearn.linear_model import LogisticRegression
from sklearn.pipeline import make_pipeline
from sklearn.metrics import f1_score

train_scores = []
val_scores = []
scaler = StandardScaler()
1 2 3

fj: $\{C_1, C_2, \dots, C_D\}$ train data

49 / 49

∞ imbalanced Data.ipynb - Colab plot (-x*log(x)-(1-x)*log(1-x)) Category Encoders - Category Encoder Target Encoder - Category Encoder

colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=xGwDDIXTNq8W

Update

+ Code + Text

✓ RAM
Disk

Users Settings

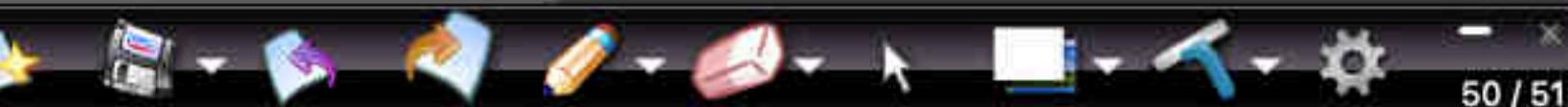
Up Down Left Right Home End

```
[ ] # Hyper-param tuning
from sklearn.linear_model import LogisticRegression ✓
from sklearn.pipeline import make_pipeline
from sklearn.metrics import f1_score

[ ] train_scores = []
[ ] val_scores = []
scaler = StandardScaler()
✓ l=0.01
✓ h= 1000.0
✓ d=50.0

[ ] for la in np.arange(l,h,d):
    scaled_lr = make_pipeline( scaler, LogisticRegression(C=1/la))
    scaled_lr.fit(X_train, y_train)
    train_y_pred = scaled_lr.predict(X_train)
    val_y_pred = scaled_lr.predict(X_val)
    train_score = f1_score(y_train, train_y_pred)
    val_score = f1_score(y_val, val_y_pred)
    train_scores.append(train_score)
    val_scores.append(val_score)

[ ] plt.figure()
```



∞ imbalanced Data.ipynb - Colab plot (-x*log(x)-(1-x)*log(1-x)) × Category Encoders - Category En × Target Encoder - Category En × +

colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=xGwDDIXTNq8W

Update :

+ Code + Text

✓ RAM Disk

👤⚙️

↑ ↓ ⌂ ⚙️ 📁 🗑️ :

```
[ ] # Hyper-param tuning
from sklearn.linear_model import LogisticRegression
from sklearn.pipeline import make_pipeline
from sklearn.metrics import f1_score

train_scores = []
val_scores = []
scaler = StandardScaler()
l=0.01
h= 1000.0
d=50.0

for la in np.arange(l,h,d): ✓
    scaled_lr = make_pipeline(scaler, LogisticRegression(C=1/la))
    scaled_lr.fit(X_train, y_train)
    train_y_pred = scaled_lr.predict(X_train)
    val_y_pred = scaled_lr.predict(X_val)
    train_score = f1_score(y_train, train_y_pred)
    val_score = f1_score(y_val, val_y_pred)
    train_scores.append(train_score)
    val_scores.append(val_score)

[ ] plt.figure()
```

* No class
rewighting

* No class
rewighting

∞ imbalanced Data.ipynb - Colab | plot (-x*log(x)-(1-x)*log(1-x)) | Category Encoders - Category En | Target Encoder - Category En | +

colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=xGwDDIXTNq8W

Update :

+ Code + Text

✓ RAM
Disk

Users Settings

Up Down Left Right Home End Page Up Page Down

```
[ ] # Hyper-param tuning
from sklearn.linear_model import LogisticRegression
from sklearn.pipeline import make_pipeline
from sklearn.metrics import f1_score

train_scores = []
val_scores = []
scaler = StandardScaler()
l=0.01
h= 1000.0
d=50.0

for la in np.arange(l,h,d):
    scaled_lr = make_pipeline(scaler, LogisticRegression(C=1/la))
    scaled_lr.fit(X_train, y_train)
    train_y_pred = scaled_lr.predict(X_train)
    val_y_pred = scaled_lr.predict(X_val)
    train_score = f1_score(y_train, train_y_pred)
    val_score = f1_score(y_val, val_y_pred)
    train_scores.append(train_score)
    val_scores.append(val_score)

[ ] plt.figure()
```



∞ imbalanced Data.ipynb - Colab × plot (-x*log(x)-(1-x)*log(1-x)) × Category Encoders - Category Encoder × Target Encoder - Category Encoder × +

colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=xGwDDIXTNq8W

+ Code + Text

RAM Disk

Update

[] # Hyper-param tuning

{x} from sklearn.linear_model import LogisticRegression

{} from sklearn.pipeline import make_pipeline

{} from sklearn.metrics import f1_score

{} train_scores = []

{} val_scores = []

{} scaler = StandardScaler()

{} l=0.01

{} h= 1000.0

{} d=50.0

{} for la in np.arange(l,h,d):

{} scaled_lr = make_pipeline(scaler, LogisticRegression(C=1/la))

{} scaled_lr.fit(X_train, y_train)

{} train_y_pred = scaled_lr.predict(X_train)

{} val_y_pred = scaled_lr.predict(X_val)

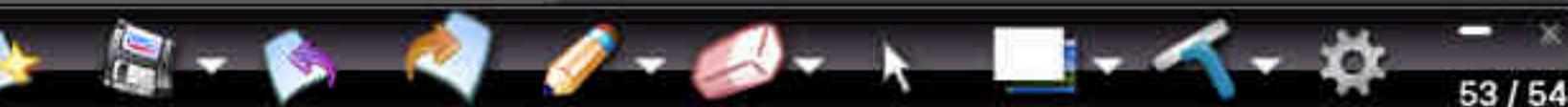
{} train_score = f1_score(y_train, train_y_pred)

{} val_score = f1_score(y_val, val_y_pred)

<> { train_scores.append(train_score)

{} { val_scores.append(val_score)

[] plt.figure()



∞ imbalanced Data.ipynb - Colab × plot (-x*log(x)-(1-x)*log(1-x)) × Category Encoders - Category Encoder × Target Encoder - Category Encoder × +

colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=xGwDDIXTNq8W

+ Code + Text RAM Disk

Update

[] # Hyper-param tuning

{x} from sklearn.linear_model import LogisticRegression

train_scores = []

val_scores = []

scaler = StandardScaler()

l=0.01

h= 1000.0

d=50.0

for la in np.arange(l,h,d):

scaled_lr = make_pipeline(scaler, LogisticRegression(C=1/la))

scaled_lr.fit(X_train, y_train)

train_y_pred = scaled_lr.predict(X_train)

val_y_pred = scaled_lr.predict(X_val)

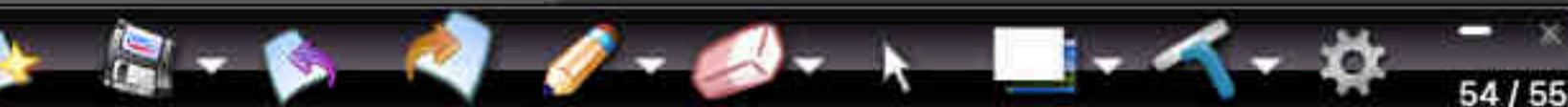
train_score = f1_score(y_train, train_y_pred)

val_score = f1_score(y_val, val_y_pred)

train_scores.append(train_score)

val_scores.append(val_score)

[] plt.figure()



∞ imbalanced Data.ipynb - Colab plot (-x*log(x)-(1-x)*log(1-x)) × Category Encoders - Category En × Target Encoder - Category En × +

colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=xGwDDIXTNq8W

Update :

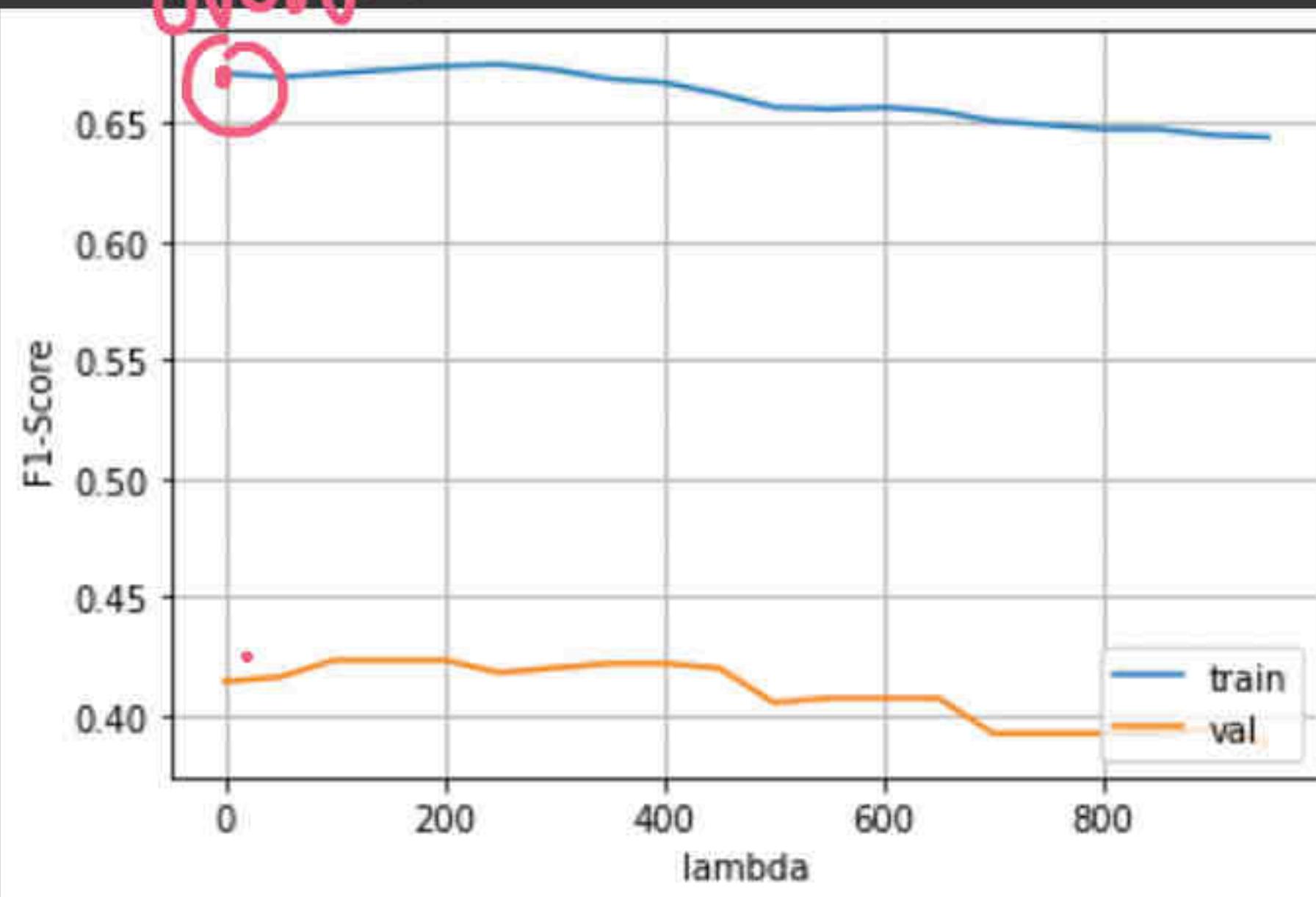
+ Code + Text

✓ RAM Disk



```
[ ] plt.plot(list(np.arange(l,h,d)), val_scores, label="val")
[ ] plt.legend(loc='lower right')
[ ] plt.xlabel("lambda")
[ ] plt.ylabel("F1-Score")
[ ] plt.grid()
[ ] plt.show()
```

overfit



0.425 ← no class rebalancing

```
▶ # minority class needs more re-weighting
▶ # Hyper-param tuning
```

∞ imbalanced Data.ipynb - Colab × plot (-x*log(x)-(1-x)*log(1-x)) × Category Encoders - Categorical Encoding × Target Encoder - Category En × +

colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=xGwDDIXTNq8W

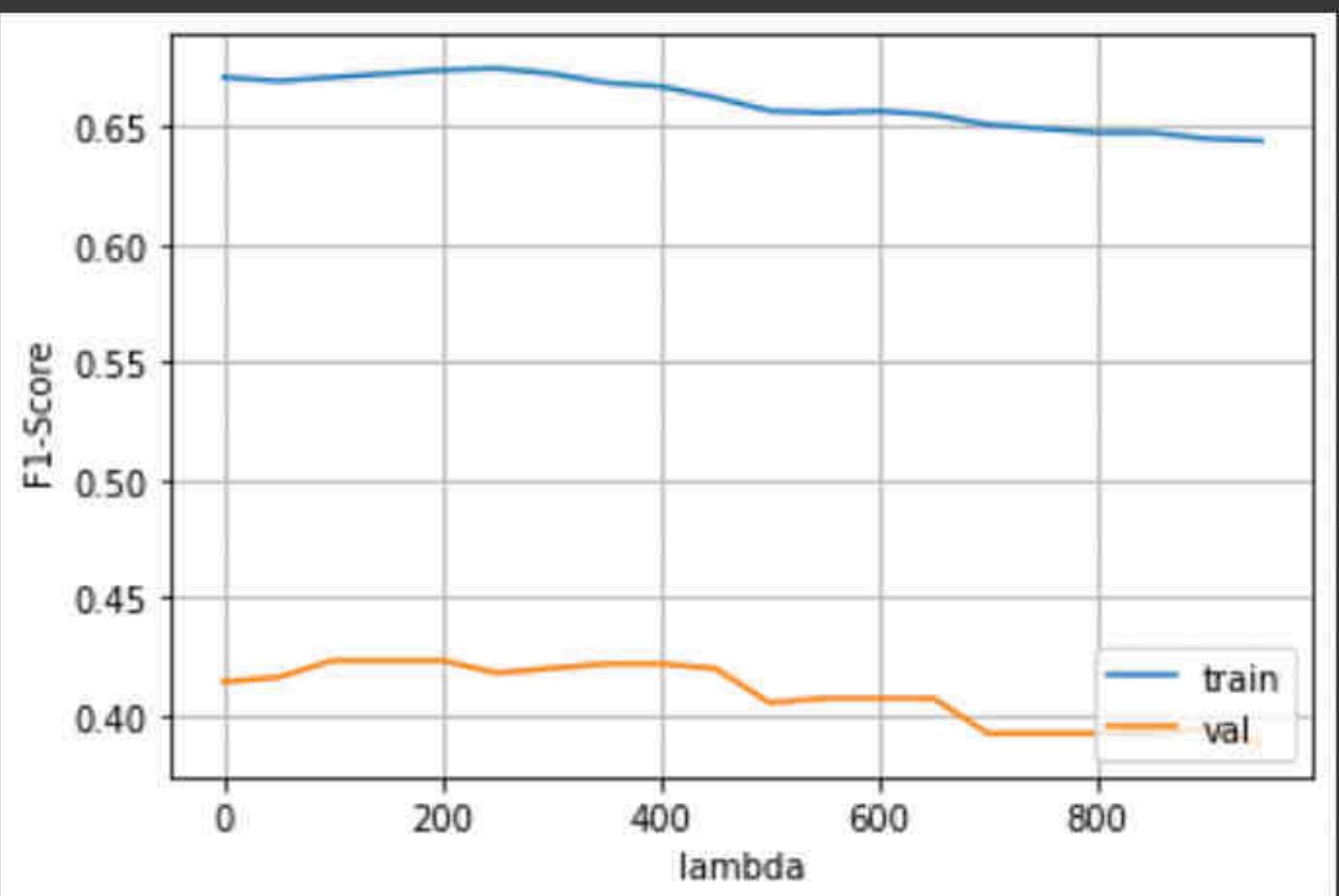
Update :

+ Code + Text

✓ RAM
Disk



```
[ ] plt.plot(list(np.arange(l,h,d)), val_scores, label="val")
[ ] plt.legend(loc='lower right')
[ ] plt.xlabel("lambda")
[ ] plt.ylabel("F1-Score")
[ ] plt.grid()
[ ] plt.show()
```



▶ # minority class needs more re-weighting

▶ # Hyper-param tuning

imbalanced Data.ipynb - Colab plot (-x*log(x)-(1-x)*log(1-x)) Category Encoders - Category Encoder Target Encoder - Category Encoder

colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=xGwDDIXTNq8W

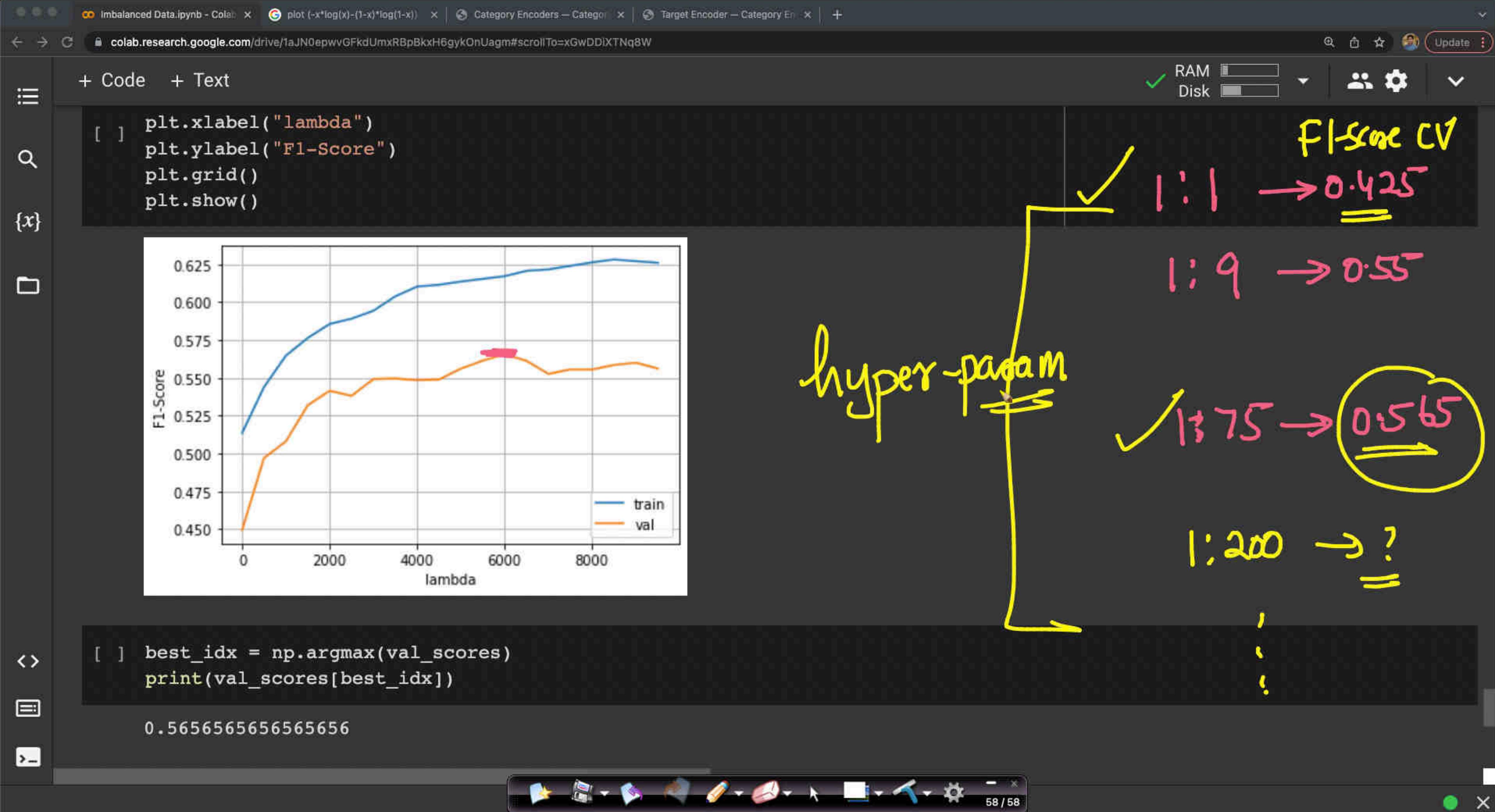
+ Code + Text RAM Disk

[] train_scores = []
val_scores = []
scaler = StandardScaler()
 $\{x\}$ l=0.01
h= 1000.0
d=50.0

for la in np.arange(l,h,d):
 scaled_lr = make_pipeline(scaler, LogisticRegression(C=1/la, class_weight={ 0:0.1, 1:0.9 }))
 scaled_lr.fit(X_train, y_train)
 train_y_pred = scaled_lr.predict(X_train)
 val_y_pred = scaled_lr.predict(X_val)
 train_score = f1_score(y_train, train_y_pred)
 val_score = f1_score(y_val, val_y_pred)
 train_scores.append(train_score)
 val_scores.append(val_score)

#plotting
plt.figure()
plt.plot(list(np.arange(l,h,d)), train_scores, label="train")
plt.plot(list(np.arange(l,h,d)), val_scores, label="val")
plt.legend(loc='lower right')
plt.xlabel("lambda")
plt.ylabel("F1-Score")
plt.grid()

1
q



∞ imbalanced Data.ipynb - Colab | plot (-x*log(x)-(1-x)*log(1-x)) | Category Encoders - Category En | Target Encoder - Category En | +

colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=MUgn2VXV-aFV

+ Code + Text

RAM Disk

0 2000 4000 6000 8000 lambda

↑ ↓ ⌂ ⚙ ⚙ ⚙ ⚙ ⚙ ⚙

{x} {y}

0.5656565656565656

0.565 : Val-Dala F1

[] # Model with lambda_best
best_idx = np.argmax(val_scores)
l_best = l+d*best_idx
scaled_lr = make_pipeline(scaler, LogisticRegression(C=1/l_best, class_weight={0:0.01, 1:0.75}))
scaled_lr.fit(X_train, y_train)

y_pred_test = scaled_lr.predict(X_test)
test_score = f1_score(y_test, y_pred_test)
print(test_score)

0.5721925133689839

test

test

[] # Alternatives metrics for imbalanced data:
Brier-Score: $1/n \sum [Y_i - Y_{i_hat}]^2 = \text{AVG Squared Error}$

59 / 59

∞ imbalanced Data.ipynb - Colab × plot (-x*log(x)-(1-x)*log(1-x)) × Category Encoders - Category En × Target Encoder - Category En × +

colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=_s4jB_N49xYb

Update :

+ Code + Text

✓ RAM
Disk



0.5656565656565656

```
[ ] # Model with lambda_best
best_idx = np.argmax(val_scores)
l_best = 1+best_idx
scaled_lr = make_pipeline( scaler, LogisticRegression(C=1/l_best, class_weight={ 0:0.01, 1:0.75 }))  
scaled_lr.fit(X_train, y_train)

y_pred_test = scaled_lr.predict(X_test)
test_score = f1_score(y_test, y_pred_test)

print(test_score)
```

0.5721925133689839

```
[30] # minority class needs more weighting

# Hyper-param tuning
from sklearn.linear_model import LogisticRegression
from sklearn.pipeline import make_pipeline
from sklearn.metrics import f1_score

train_scores = []
```

```
∞ imbalanced Data.ipynb - Colab × plot (-x*log(x)-(1-x)*log(1-x)) × Category Encoders - Categorical Encoding × Target Encoder - Category Encoding × +  
colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=JM5J906_gLj  
+ Code + Text RAM Disk ✓    
[33] # Model with lambda_best  
best_idx = np.argmax(val_scores)  
l_best = l+d*best_idx  
scaled_lr = make_pipeline( scaler, LogisticRegression(C=1/l_best, class_weight={ 0:0.01, 1:0.75 }))  
scaled_lr.fit(X_train, y_train)  
  
y_pred_test = scaled_lr.predict(X_test)  
test_score = f1_score(y_test, y_pred_test)  
  
print(test_score)
```

0.6007751937984496



```
✓ confusion_matrix(y_test, y_pred_test)  
array([[13549,    175],  
       [    31,   155]])  
↑ 0  
↓ 1
```

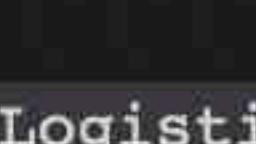
```
<> [30] # minority class needs more weighting  
      # Hyper-param tuning  
      from sklearn.linear_model import LogisticRegression  
      from sklearn.pipeline import make_pipeline
```

```
∞ imbalanced Data.ipynb - Colab × plot (-x*log(x)-(1-x)*log(1-x)) × Category Encoders - Categorical Encoding × Target Encoder - Category Encoding × +  
colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=JM5J906_gLj  
+ Code + Text RAM Disk    
[33] # Model with lambda_best  
best_idx = np.argmax(val_scores)  
l_best = l+d*best_idx  
scaled_lr = make_pipeline( scaler, LogisticRegression(C=1/l_best, class_weight={ 0:0.01, 1:0.75 }))  
scaled_lr.fit(X_train, y_train)  
  
y_pred_test = scaled_lr.predict(X_test)  
test_score = f1_score(y_test, y_pred_test)  
  
print(test_score)
```

0.6007751937984496

$$\frac{155}{186}$$

```
✓ confusion_matrix(y_test, y_pred_test)  
array([[13549,    175],  
       [    31,   155]])  
  
<> [30] # minority class needs more weighting  
# Hyper-param tuning  
from sklearn.linear_model import LogisticRegression  
from sklearn.pipeline import make_pipeline  
  
[31] # Using Pipeline
```



```
∞ imbalanced Data.ipynb - Colab × plot (-x*log(x)-(1-x)*log(1-x)) × Category Encoders - Categorical Encoding × Target Encoder - Category Encoding × +  
colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=JM5J906_gLj  
+ Code + Text RAM Disk ✓    
[33] # Model with lambda_best  
best_idx = np.argmax(val_scores)  
l_best = l+d*best_idx  
scaled_lr = make_pipeline( scaler, LogisticRegression(C=1/l_best, class_weight={ 0:0.01, 1:0.75 }))  
scaled_lr.fit(X_train, y_train)  
  
y_pred_test = scaled_lr.predict(X_test)  
test_score = f1_score(y_test, y_pred_test)  
  
print(test_score)
```

0.6007751937984496



```
✓ confusion_matrix(y_test, y_pred_test)  
array([[13549, 175],  
       [ 31, 155]])
```

```
<> [30] # minority class needs more weighting  
      # Hyper-param tuning  
      from sklearn.linear_model import LogisticRegression  
      from sklearn.pipeline import make_pipeline
```



imbalanced Data.ipynb - Colab plot (-x*log(x)-(1-x)*log(1-x)) Category Encoders - Category Encoder Target Encoder - Category Encoder

colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=JM5J906_gLj

+ Code + Text RAM Disk

[33] # Model with lambda_best
best_idx = np.argmax(val_scores)
l_best = l+d*best_idx
scaled_lr = make_pipeline(scaler, LogisticRegression(C=1/l_best, class_weight={ 0:0.01, 1:0.75 }))
scaled_lr.fit(X_train, y_train)

y_pred_test = scaled_lr.predict(X_test)
test_score = f1_score(y_test, y_pred_test)

print(test_score)

0.6007751937984496



confusion_matrix(y_test, y_pred_test)

array([[13549, 175],
 [31, 155]])

[30] # minority class needs more weighting

Hyper-param tuning
from sklearn.linear_model import LogisticRegression
from sklearn.pipeline import make_pipeline

imbalanced Data.ipynb - Colab plot (-x*log(x)-(1-x)*log(1-x)) Category Encoders - Category Encoder Target Encoder - Category Encoder

colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=JM5J906_gLj

+ Code + Text RAM Disk

[33] # Model with lambda_best

```
best_idx = np.argmax(val_scores)
l_best = l+d*best_idx
scaled_lr = make_pipeline( scaler, LogisticRegression(C=1/l_best, class_weight={ 0:0.01, 1:0.75 }))
```

scaled_lr.fit(X_train, y_train)

y_pred_test = scaled_lr.predict(X_test)
test_score = f1_score(y_test, y_pred_test)

print(test_score)

0.6007751937984496

confusion_matrix(y_test, y_pred_test)

	predicted	
actual	0	1
0	13549	175
1	31	155

Recall = $\frac{175}{186} \approx 80\%$

Pr: $\frac{155}{330} \approx 48\%$

[30] # minority class needs more weighting

Hyper-param tuning

```
from sklearn.linear_model import LogisticRegression
from sklearn.pipeline import make_pipeline
```

+ Code + Text

✓ RAM
Disk

```
# Alternatives metrics for imbalanced data:  
[ ] # Brier-Score: 1/n*SUM[Y_i - Y_i_hat]^2 = AVG Squared Error  
[ ] # G-mean: SQRT(Sensitivity * Specificity)=SQRT(TPR*TNR) ----Seldom used
```

F1-Score

- AUC (ROC) - later
=

$$\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

on 1
pools

avg squared error with probabilities

$\sqrt{\text{TPR} \times \text{TNR}}$

∞ imbalanced Data.ipynb - Colab × Google plot (-x*log(x)-(1-x)*log(1-x)) × Category Encoders - Category Encoder × Target Encoder - Category Encoder × +

colab.research.google.com/drive/1aJN0epwvGFkdUmxRBpBkxH6gykOnUagm#scrollTo=-BsklhG4-zo4

Update :

+ Code + Text

✓ RAM
Disk



```
# Alternatives metrics for imbalanced data:  
[ ] # Brier-Score: 1/n*SUM[Y_i - Y_i_hat]^2 = AVG Squared Error  
# G-mean: SQRT(Sensitivity * Specificity)=SQRT(TPR*TNR) ----Seldom used
```

{x}

[]

□

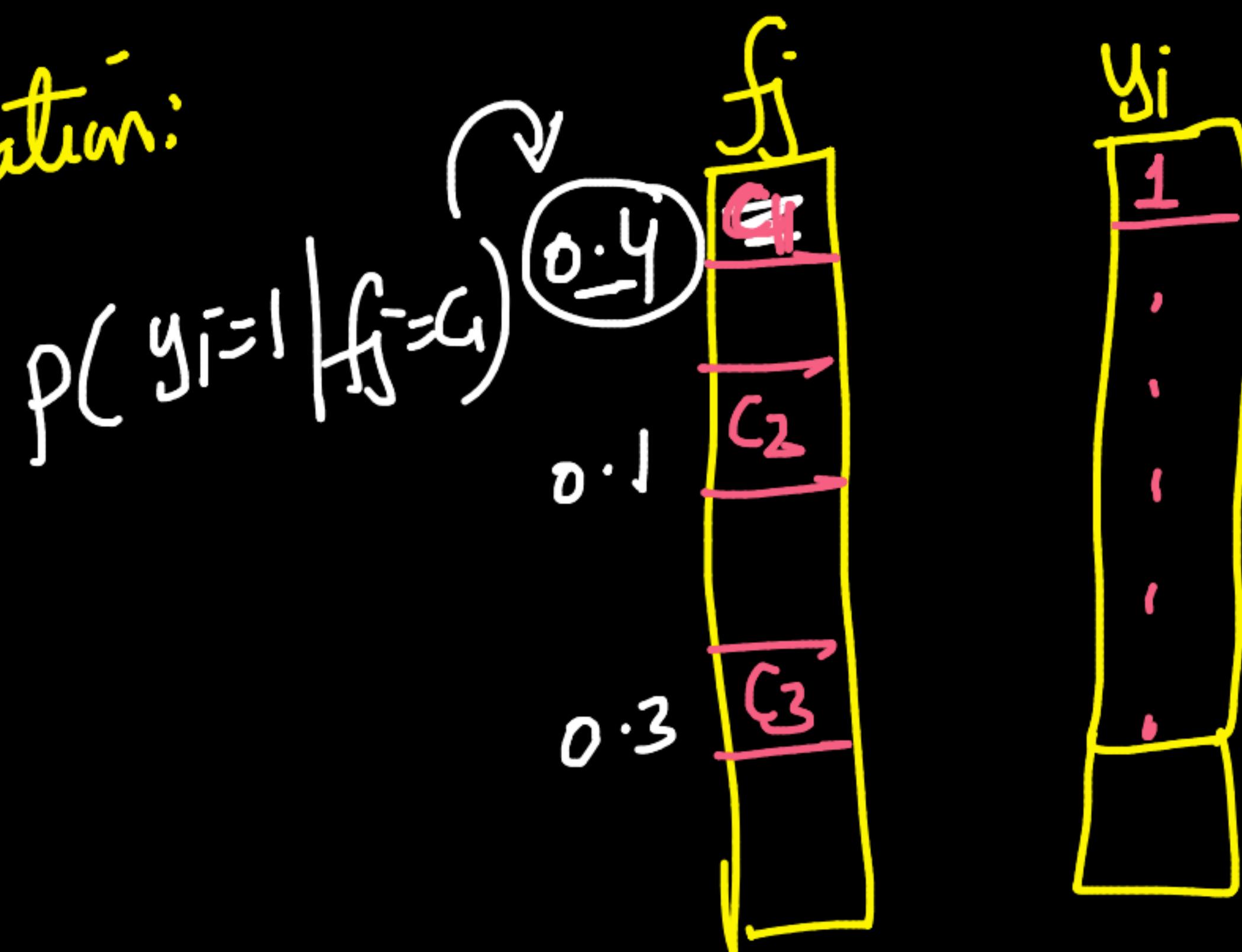
<>

≡

>-



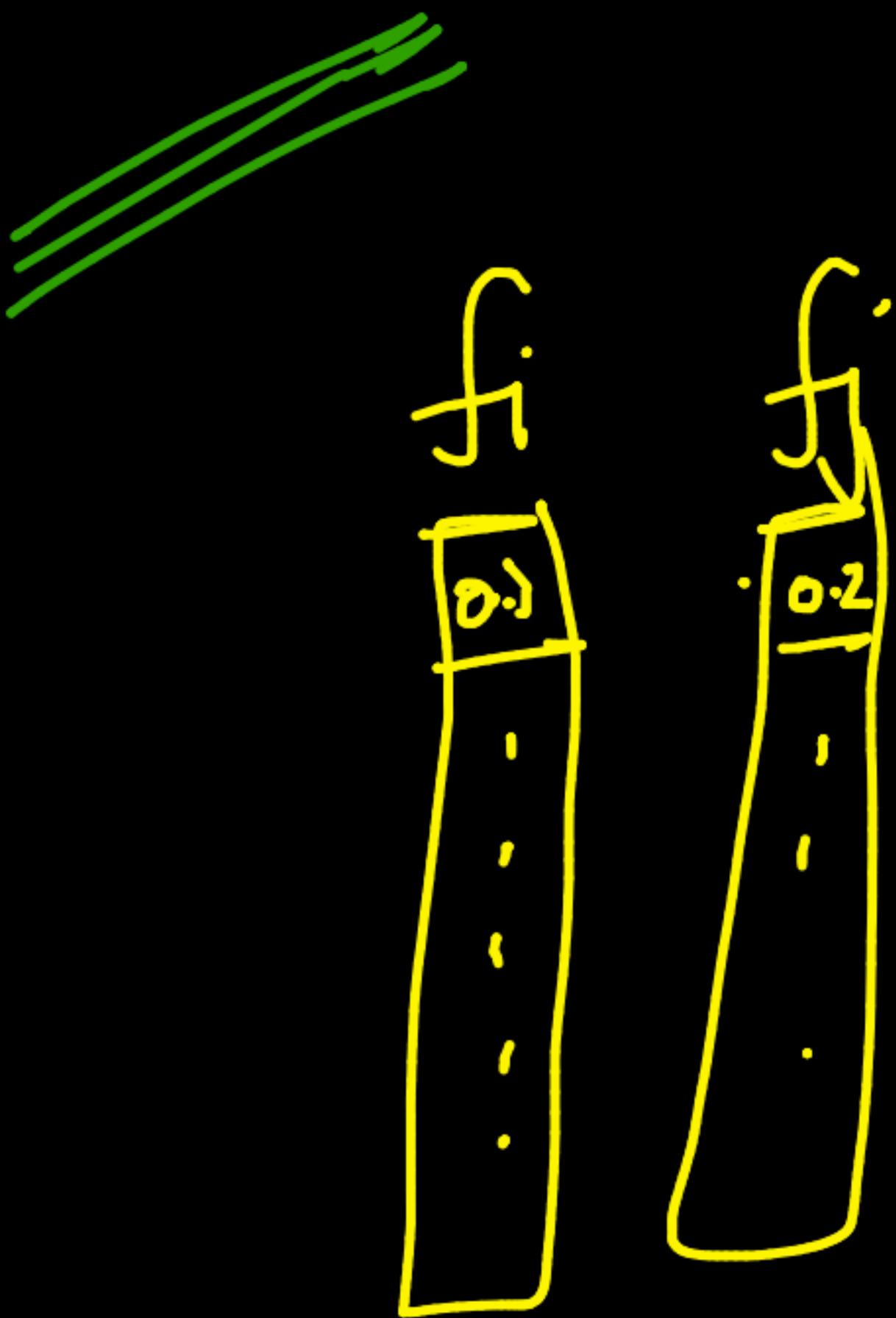
Correlation:



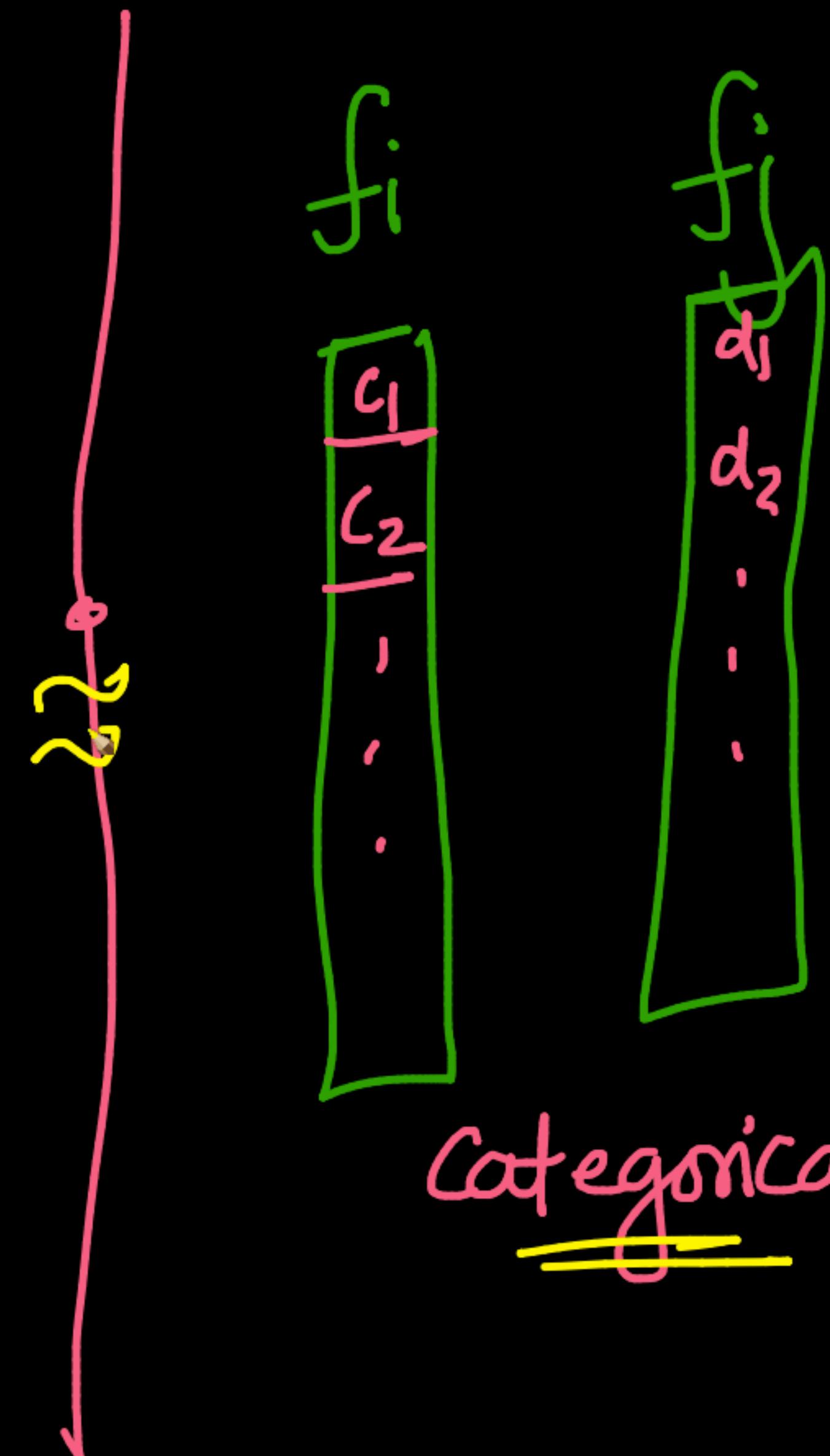
Target encoding

$c_1 \xrightarrow[\text{corr}]{\text{highly}} y_i = 1$

$\{ p(y_i=1 | c_1) \rightarrow y_i = 1$



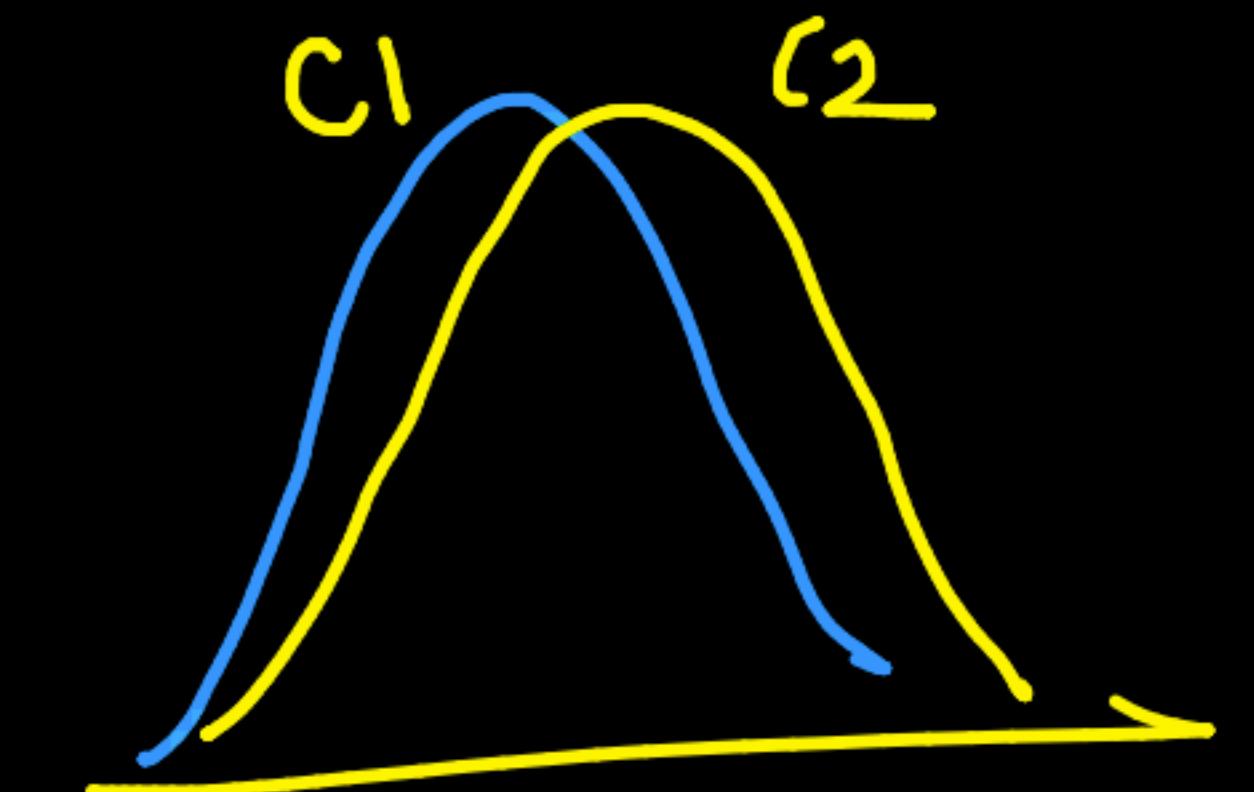
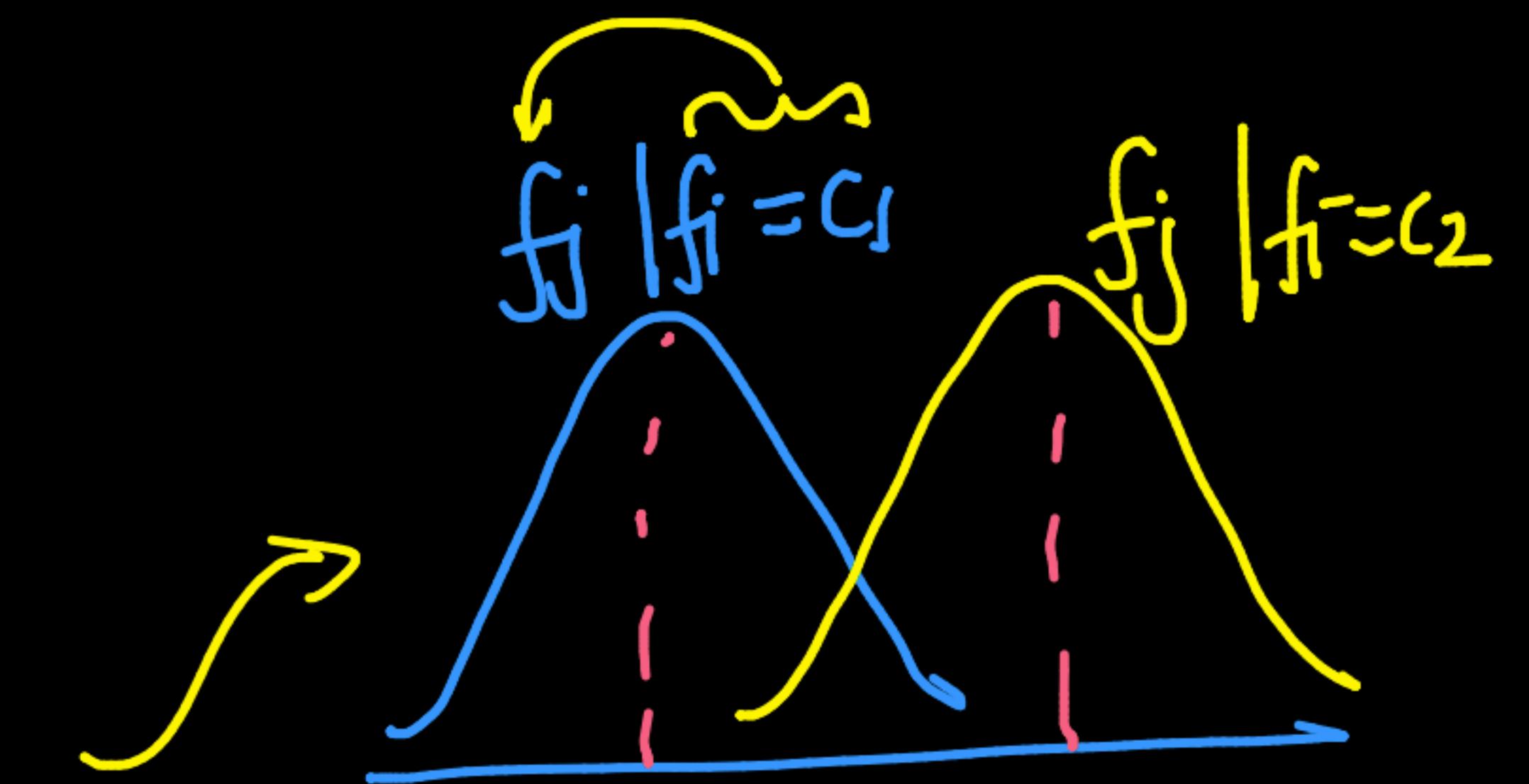
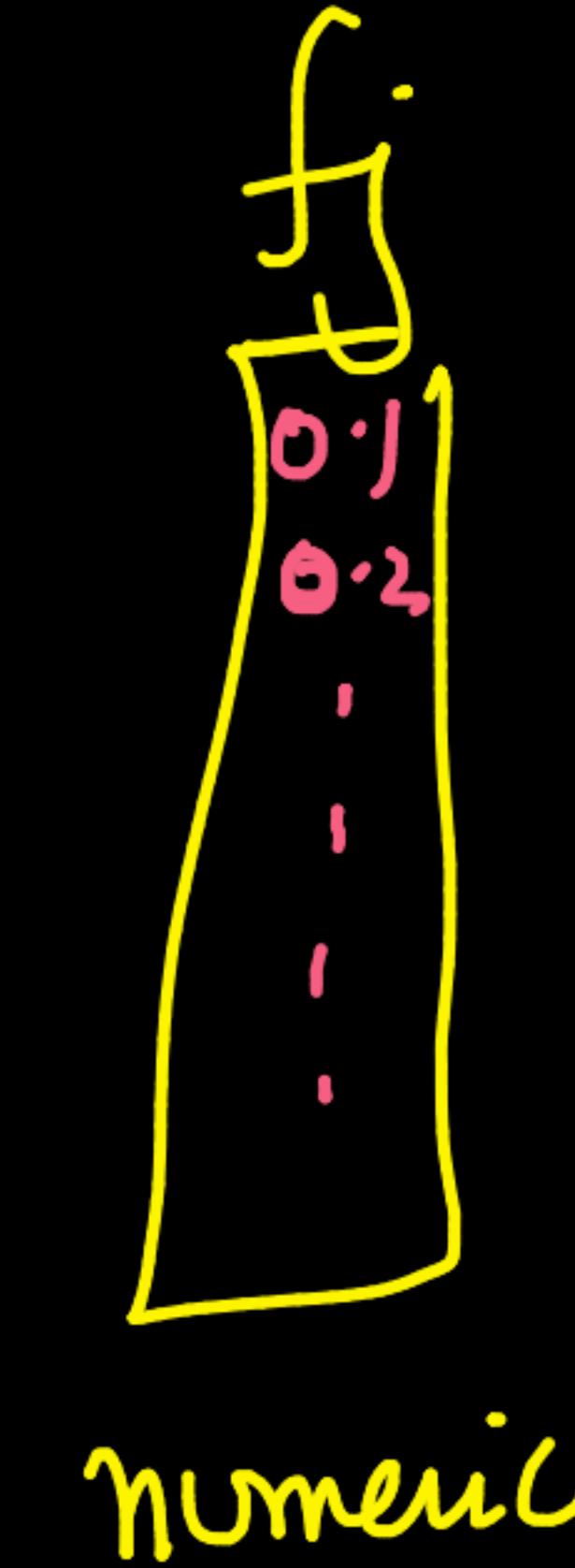
{ PCC or SRCC



χ^2 -test

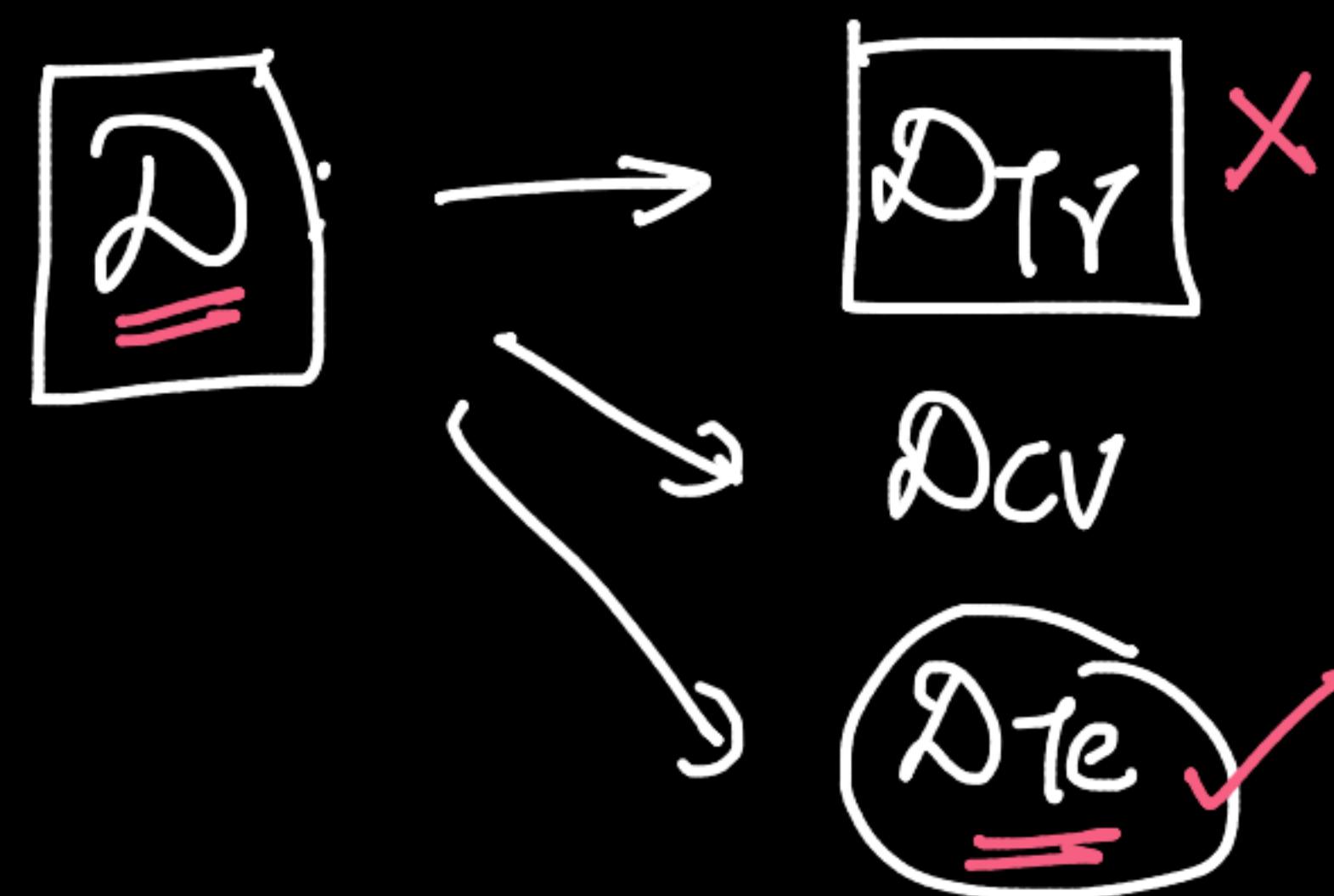
✓ {
z-test
nT-Test

✓ { ANOVA



✓ KL-Divergence

best
Dtrain



time-series : TBSplitting



$$P(Y_i=1 \mid f_j = c_1) = \frac{4\alpha + \beta}{10\alpha + \beta}$$

↑ brain
↑ test

The equation shows the probability of $Y_i=1$ given $f_j = c_1$. The numerator is $4\alpha + \beta$ and the denominator is $10\alpha + \beta$. Red arrows point from the terms α and β to the words "brain" and "test" respectively.

- Pr, recall
- $\text{F1-Score} \checkmark$
- G-mean: $\sqrt{\frac{\text{Sens} + \text{Spec}}{2}}$

en.wikipedia.org/wiki/Confusion_matrix

Community portal Recent changes Upload file Tools What links here Related changes Special pages Permanent link Page information Cite this page Wikidata item Print/export Download as PDF Printable version Languages العربية Deutsch Español Français 한국어 Italiano 日本語 Português 中文 文 5 more Edit links

Identical sets of classes in both dimensions (each combination of dimension and class is a variable in the contingency table).

Contents [hide]

- 1 Example
- 2 Table of confusion
- 3 Confusion matrices with more than two categories
- 4 See also
- 5 References

Example [edit]

Given a sample of 12 individuals, 8 that have been diagnosed with cancer and 4 that are cancer-free, where individuals with cancer belong to class 1 (positive) and non-cancer individuals belong to class 0 (negative), we can display that data as follows:

Individual Number	1	2	3	4	5	6	7	8	9	10	11	12
Actual Classification	1	1	1	1	1	1	1	1	0	0	0	0

Assume that we have a classifier that distinguishes between individuals with and without cancer in some way, we can take the 12 individuals and run them through the classifier. The classifier then makes 9 accurate predictions and misses 3: 2 individuals with cancer wrongly predicted as being cancer-free (sample 1 and 2), and 1 person without cancer that is wrongly predicted to have cancer (sample 9).

Individual Number	1	2	3	4	5	6	7	8	9	10	11	12
Actual Classification	1	1	1	1	1	1	1	1	0	0	0	0
Predicted Classification	0	0	1	1	1	1	1	1	1	0	0	0

Notice that if we compare the actual classification set to the predicted classification set, we see that samples 1 and 2 were predicted as negative (0), while samples 9 and 10 were predicted as positive (1). This indicates a misclassification rate of 3 out of 12, or 25%.

true negative (TN)
A test result that correctly indicates the absence of a condition or characteristic

false positive (FP)
A test result which wrongly indicates that a particular condition or attribute is present

false negative (FN)
A test result which wrongly indicates that a particular condition or attribute is absent

sensitivity, recall, hit rate, or true positive rate (TPR)
 $\text{TPR} = \frac{\text{TP}}{\text{P}} = \frac{\text{TP}}{\text{TP} + \text{FN}} = 1 - \text{FNR}$

specificity, selectivity or true negative rate (TNR)
 $\text{TNR} = \frac{\text{TN}}{\text{N}} = \frac{\text{TN}}{\text{TN} + \text{FP}} = 1 - \text{FPR}$

precision or positive predictive value (PPV)
 $\text{PPV} = \frac{\text{TP}}{\text{TP} + \text{FP}} = 1 - \text{FDR}$

negative predictive value (NPV)
 $\text{NPV} = \frac{\text{TN}}{\text{TN} + \text{FN}} = 1 - \text{FOR}$

miss rate or false negative rate (FNR)
 $\text{FNR} = \frac{\text{FN}}{\text{P}} = \frac{\text{FN}}{\text{FN} + \text{TP}} = 1 - \text{TPR}$

false-out or false positive rate (FPR)
 $\text{FPR} = \frac{\text{FP}}{\text{N}} = \frac{\text{FP}}{\text{FP} + \text{TN}} = 1 - \text{TNR}$

false discovery rate (FDR)
 $\text{FDR} = \frac{\text{FP}}{\text{FP} + \text{TP}} = 1 - \text{PPV}$

The template for any binary confusion matrix uses the four kinds of results discussed above (true positives, false negatives, false positives, and true negatives) along with the positive and negative classifications. The four outcomes can be formulated in a 2×2 *confusion matrix*, as follows:

		Predicted condition	
		Positive (PP)	Negative (PN)
Actual condition	Total population = P + N	Positive (P)	Negative (N)
	Positive (P)	True positive (TP)	False negative (FN)
Negative (N)		False positive (FP)	True negative (TN)

Sources: [11][12][13][14][15][16][17][18]

view · talk · edit

TN
—
TN+FP

The color convention of the three data tables above were picked to match this confusion matrix, in order to easily differentiate the data.

Now, we can simply total up each type of result, substitute into the template, and create a confusion matrix that will concisely summarize the results of testing the classifier:

		Predicted condition	
		Cancer	Non-cancer
Actual condition	Total 8 + 4 = 12	Cancer	Non-cancer
Cancer	6	2	

$$P + N \quad TP + TN + FP + FN$$

balanced accuracy (BA)

$$BA = \frac{TPR + TNR}{2}$$

F1 score

is the harmonic mean of precision and sensitivity:

$$F_1 = 2 \times \frac{PPV \times TPR}{PPV + TPR} = \frac{2TP}{2TP + FP + FN}$$

phi coefficient (ϕ or r_ϕ) or Matthews correlation coefficient (MCC)

$$MCC = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$

Fowlkes–Mallows index (FM)

$$FM = \sqrt{\frac{TP}{TP + FP} \times \frac{TP}{TP + FN}} = \sqrt{PPV \times TPR}$$

informedness or bookmaker informedness (BM)

$$BM = TPR + TNR - 1$$

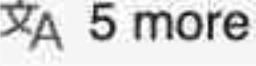
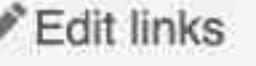
markedness (MK) or deltaP (Δp)

$$MK = PPV + NPV - 1$$

Diagnostic odds ratio (DOR)

$$DOR = \frac{LR+}{LR-}$$

Sources: Fawcett (2006),^[1] Piryonesi and El-Diraby (2020),^[2] Powers (2011),^[3] Ting (2011),^[4] CAWCR,^[5] D. Chicco & G. Jurman (2020, 2021),^[6] Tharwat (2018).^[8]

Recent changes
Upload file
Tools
What links here
Related changes
Special pages
Permanent link
Page information
Cite this page
Wikidata item
Print/export
Download as PDF
Printable version
Languages 
العربية
Deutsch
Español
Français
한국어
Italiano
日本語
Português
中文
 5 more
 Edit links

variable in the contingency table).

Contents [hide]

- 1 Example
- 2 Table of confusion
- 3 Confusion matrices with more than two categories
- 4 See also
- 5 References

Example [edit]

Given a sample of 12 individuals, 8 that have been diagnosed with cancer and 4 that are cancer-free, where individuals with cancer belong to class 1 (positive) and non-cancer individuals belong to class 0 (negative), we can display that data as follows:

Individual Number	1	2	3	4	5	6	7	8	9	10	11	12
Actual Classification	1	1	1	1	1	1	1	1	0	0	0	0

Assume that we have a classifier that distinguishes between individuals with and without cancer in some way, we can take the 12 individuals and run them through the classifier. The classifier then makes 9 accurate predictions and misses 3: 2 individuals with cancer wrongly predicted as being cancer-free (sample 1 and 2), and 1 person without cancer that is wrongly predicted to have cancer (sample 9).

Individual Number	1	2	3	4	5	6	7	8	9	10	11	12
Actual Classification	1	1	1	1	1	1	1	1	0	0	0	0
Predicted Classification	0	0	1	1	1	1	1	1	1	0	0	0

Notice, that if we compare the actual there are 4 different outcomes that could result in any combination of true/false positive/negative.

G-mean

A test result that correctly indicates the absence of a condition or characteristic

false positive (FP)

A test result which wrongly indicates that a particular condition or attribute is present

false negative (FN)

A test result which wrongly indicates that a particular condition or attribute is absent

sensitivity, recall, hit rate, or true positive rate (TPR)

$$\text{TPR} = \frac{\text{TP}}{\text{P}} = \frac{\text{TP}}{\text{TP} + \text{FN}} = 1 - \text{FNR}$$

specificity, selectivity or true negative rate (TNR)

$$\text{TNR} = \frac{\text{TN}}{\text{N}} = \frac{\text{TN}}{\text{TN} + \text{FP}} = 1 - \text{FPR}$$

precision or positive predictive value (PPV)

$$\text{PPV} = \frac{\text{TP}}{\text{TP} + \text{FP}} = 1 - \text{FDR}$$

negative predictive value (NPV)

$$\text{NPV} = \frac{\text{TN}}{\text{TN} + \text{FN}} = 1 - \text{FOR}$$

miss rate or false negative rate (FNR)

$$\text{FNR} = \frac{\text{FN}}{\text{P}} = \frac{\text{FN}}{\text{FN} + \text{TP}} = 1 - \text{TPR}$$

false-out or false positive rate (FPR)

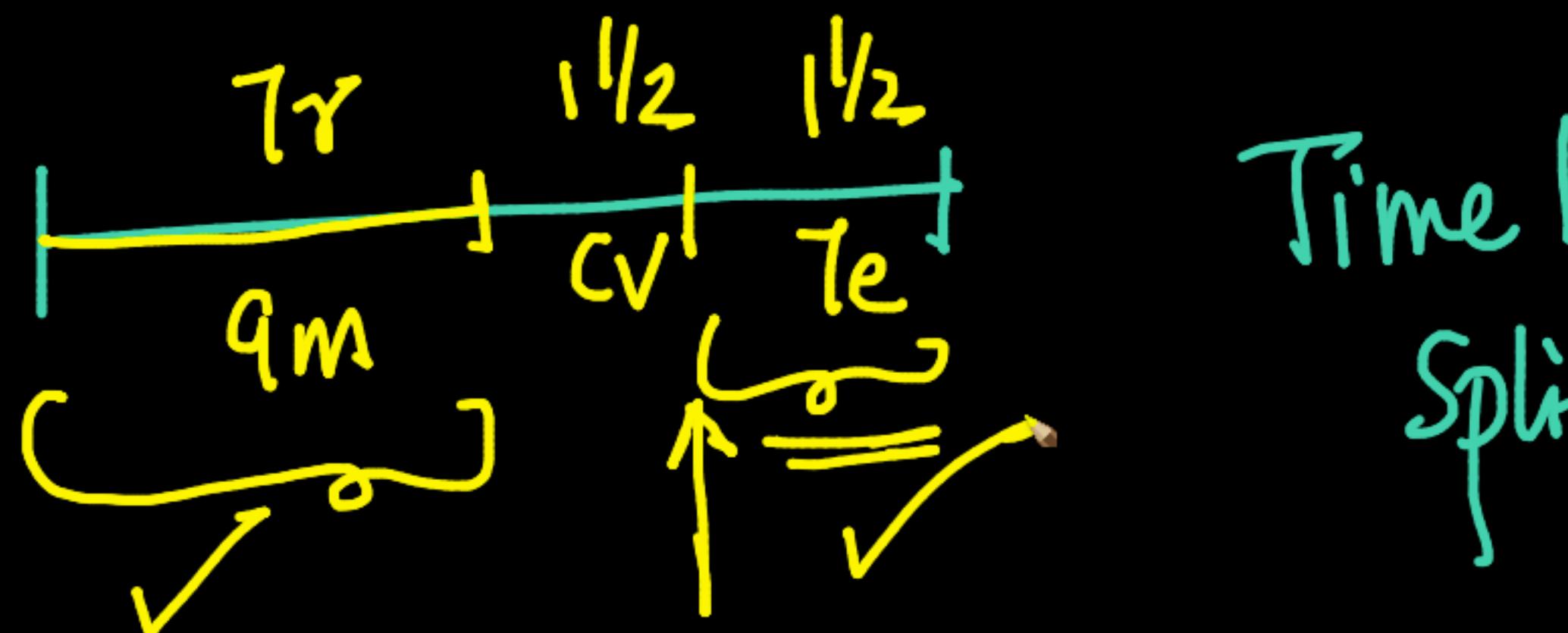
$$\text{FPR} = \frac{\text{FP}}{\text{N}} = \frac{\text{FP}}{\text{FP} + \text{TN}} = 1 - \text{TNR}$$

false discovery rate (FDR)

$$\text{FDR} = \frac{\text{FP}}{\text{FP} + \text{TP}} = 1 - \text{PPV}$$

false omission rate (FOR)

Domain → File, Gamean



← things changes over time significantly



Search the docs ...

Under-sampling methods

Over-sampling methods

RandomOverSampler

SMOTE

SMOTENC

SMOTEN

ADASYN

BorderlineSMOTE

KMeansSMOTE

SVMSMOTE

Combination of over- and under-sampling methods

Ensemble methods

Batch generator for Keras

Batch generator for TensorFlow

Miscellaneous

SMOTE

```
{ class imblearn.over_sampling.SMOTE(*, sampling_strategy='auto',  
random_state=None, k_neighbors=5, n_jobs=None)
```

[source]

Class to perform over-sampling using SMOTE.

On this page

Examples using
imblearn.over_sampling.SMO

Edit this page

This object is an implementation of SMOTE - Synthetic Minority Over-sampling Technique as presented in [1].

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

Parameters: `sampling_strategy : float, str, dict or callable, default='auto'`

Sampling information to resample the data set.

- When `float`, it corresponds to the desired ratio of the number of samples in the minority class over the number of samples in the majority class after resampling. Therefore, the ratio is expressed as $\alpha_{os} = N_{rm}/N_M$ where N_{rm} is the number of samples in the minority class after resampling and N_M is the number of samples in the majority class.



Search the docs ...

Under-sampling methods



Over-sampling methods



RandomOverSampler

SMOTE

SMOTENC

SMOTEN

ADASYN

BorderlineSMOTE

KMeansSMOTE

SVMSMOTE

Combination of over- and under-sampling methods



Ensemble methods



Batch generator for Keras



Batch generator for TensorFlow



Miscellaneous

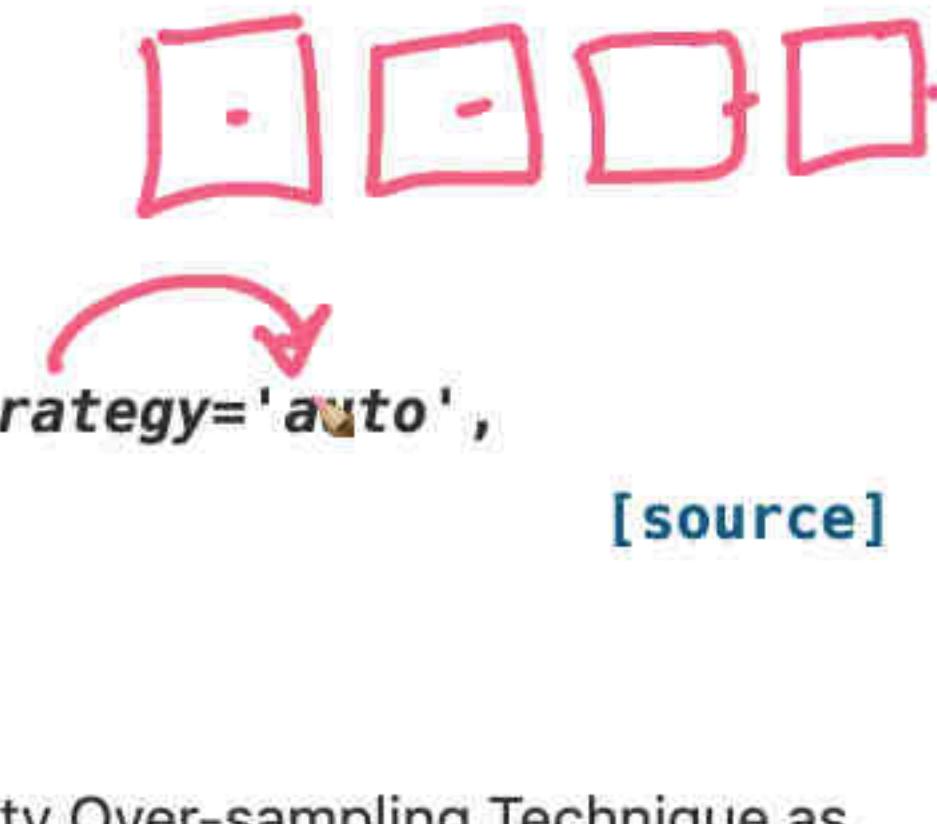


SMOTE

```
class imblearn.over_sampling.SMOTE(*, sampling_strategy='auto',
random_state=None, k_neighbors=5, n_jobs=None)
```

[source]

Class to perform over-sampling using SMOTE.



This object is an implementation of SMOTE - Synthetic Minority Over-sampling Technique as presented in [1].

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

Parameters: `sampling_strategy : float, str, dict or callable, default='auto'`

Sampling information to resample the data set.

- When `float`, it corresponds to the desired ratio of the number of samples in the minority class over the number of samples in the majority class after resampling. Therefore, the ratio is expressed as $\alpha_{os} = N_{rm}/N_M$ where N_{rm} is the number of samples in the minority class after resampling and N_M is the number of samples in the majority class.

On this page

Examples using
imblearn.over_sampling.SMO

Edit this page



Search the docs ...

Under-sampling methods



Over-sampling methods



RandomOverSampler

SMOTE

SMOTENC

SMOTEN

ADASYN

BorderlineSMOTE

KMeansSMOTE

SVMSMOTE

Combination of over- and under-sampling methods



Ensemble methods



Batch generator for Keras



Batch generator for TensorFlow



Miscellaneous

**Warning**

`float` is only available for **binary** classification. An error is raised for multi-class classification.

- When `str`, specify the class targeted by the resampling. The number of samples in the different classes will be equalized. Possible choices are:
`'minority'`: resample only the minority class;
`'not minority'`: resample all classes but the minority class;
`'not majority'`: resample all classes but the majority class;
`'all'`: resample all classes;
`'auto'`: equivalent to `'not majority'`.
- When `dict`, the keys correspond to the targeted classes. The values correspond to the desired number of samples for each targeted class.
- When callable, function taking `y` and returns a `dict`. The keys correspond to the targeted classes. The values correspond to the desired number of samples for each class.

random_state : int, RandomState instance, default=None

Control the randomization of the algorithm.

- If int, `random_state` is the seed used by the random number generator;

If `None`, the random state is the `RandomState` instance used by `np.random`.

On this page

Examples using
imblearn.over_sampling.SMO

Edit this page

imbalanced Data.ipynb - Colab | plot (-x*log(x)-(1-x)*log(1-x)) | SMOTE — Version 0.9.0

imbalanced-learn.org/stable/references/generated/imblearn.over_sampling.SMOTE.html

Getting Started User Guide API reference Examples Release history About us

Sampling technique," Journal of artificial intelligence research, 321-357, 2002.

Search the docs ...

Under-sampling methods

Over-sampling methods

- RandomOverSampler
- SMOTE**
- SMOTENC
- SMOTEN
- ADASYN

BorderlineSMOTE

KMeansSMOTE

SVMSMOTE

Combination of over- and under-sampling methods

Ensemble methods

Examples

```
>>> from collections import Counter
>>> from sklearn.datasets import make_classification
>>> from imblearn.over_sampling import SMOTE
>>> X, y = make_classification(n_classes=2, class_sep=2,
... weights=[0.1, 0.9], n_informative=3, n_redundant=1, flip_y=0,
... n_features=20, n_clusters_per_class=1, n_samples=1000, random_state=10)
>>> print('Original dataset shape %s' % Counter(y))
Original dataset shape Counter({1: 900 0: 100})
>>> sm = SMOTE(random_state=42)
>>> X_res, y_res = sm.fit_resample(X, y)
>>> print('Resampled dataset shape %s' % Counter(y_res))
Resampled dataset shape Counter({0: 900, 1: 900})
```

Methods

fit(X, y) Check inputs and statistics of the sampler.

fit(X, y)

82 / 83



 Search the docs ...

Under-sampling methods

Over-sampling methods

RandomOverSampler

SMOTE

SMOTENC

SMOTEN

ADASYN

BorderlineSMOTE

KMeansSMOTE

SVMSMOTE

Combination of over- and under-sampling methods

Ensemble methods

Examples

```
>>> from collections import Counter
>>> from sklearn.datasets import make_classification
>>> from imblearn.over_sampling import SMOTE
>>> X, y = make_classification(n_classes=2, class_sep=2,
...    weights=[0.1, 0.9], n_informative=3, n_redundant=1, flip_y=0,
...    n_features=20, n_clusters_per_class=1, n_samples=1000, random_state=10)
>>> print('Original dataset shape %s' % Counter(y))
Original dataset shape Counter({1: 900, 0: 100})
>>> sm = SMOTE(random_state=42)
>>> X_res, y_res = sm.fit_resample(X, y)
>>> print('Resampled dataset shape %s' % Counter(y_res))
Resampled dataset shape Counter({0: 900, 1: 900})
```

1:900

018

Methods

fit(X, y)

Check inputs and statistics of the sampler.

learner, the weights are recalculated such that a higher weight is assigned to samples that the current weak learner misclassified. This weight determines the probability that the sample will appear in the training of the next weak learner. For this reason, boosting algorithms like AdaBoost are particularly useful for class imbalance problems because higher weight is given to the minority class at each successive iteration as data from this class is often misclassified.

Data sampling methods combined with boosting can be an effective way to deal with class imbalance problems. To demonstrate, let's use Lending Club loan data (available [here](#)) to try to predict whether someone will default on a loan using two interesting sampling approaches that have been suggested and combined with AdaBoost: [SMOTEBoost](#) and [RUSBoost](#).

DT
RF
GBDT
Boosting

bagging

SMOTEBoost is an oversampling method based on the SMOTE algorithm (Synthetic Minority Oversampling Technique). SMOTE uses k-nearest neighbors to create synthetic examples of the minority class. SMOTEBoost then injects the SMOTE method at each boosting iteration. The advantage of this approach is that while standard boosting gives equal weights to all misclassified data, SMOTE gives more examples of the minority class at each boosting step. Similarly, RUSBoost achieves the same goal by performing random undersampling (RUS) at each boosting iteration instead of SMOTE.

Trying it out: data preparation

First, I import the data and get rid of all entries that do not correspond to "Fully Paid" or "Charged Off". I set the "Fully Paid" target to 0 and "Charged Off" to 1.

Get started Sign in to Medium with Google

Srikanth Varma Chekuri
cvarma@scaler.com

Continue as Srikanth Varma

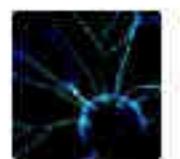


Anna Vasilyeva
34 Followers

Follow

More from Medium

Celia Sagastume in MLearning.ai



Wavelet Feature Engineering for EEG Seizure Prediction

rajni singh



Why should you deploy your ML model in shadow mode?

Edwin Shao in Terminal 1

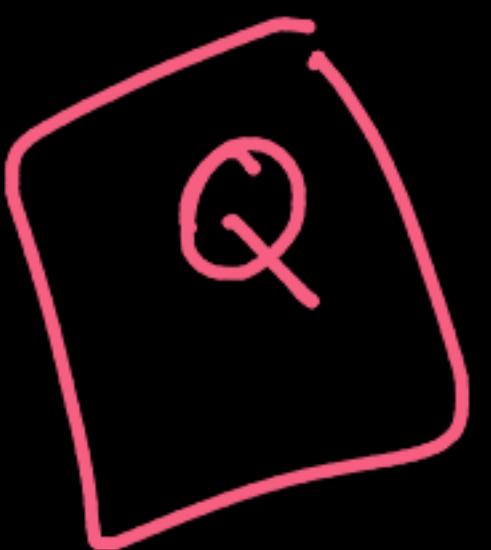


T1 & Three Revolutions

Pádraig Cunningham in Towards Data S...



Ensembles in Machine Learning



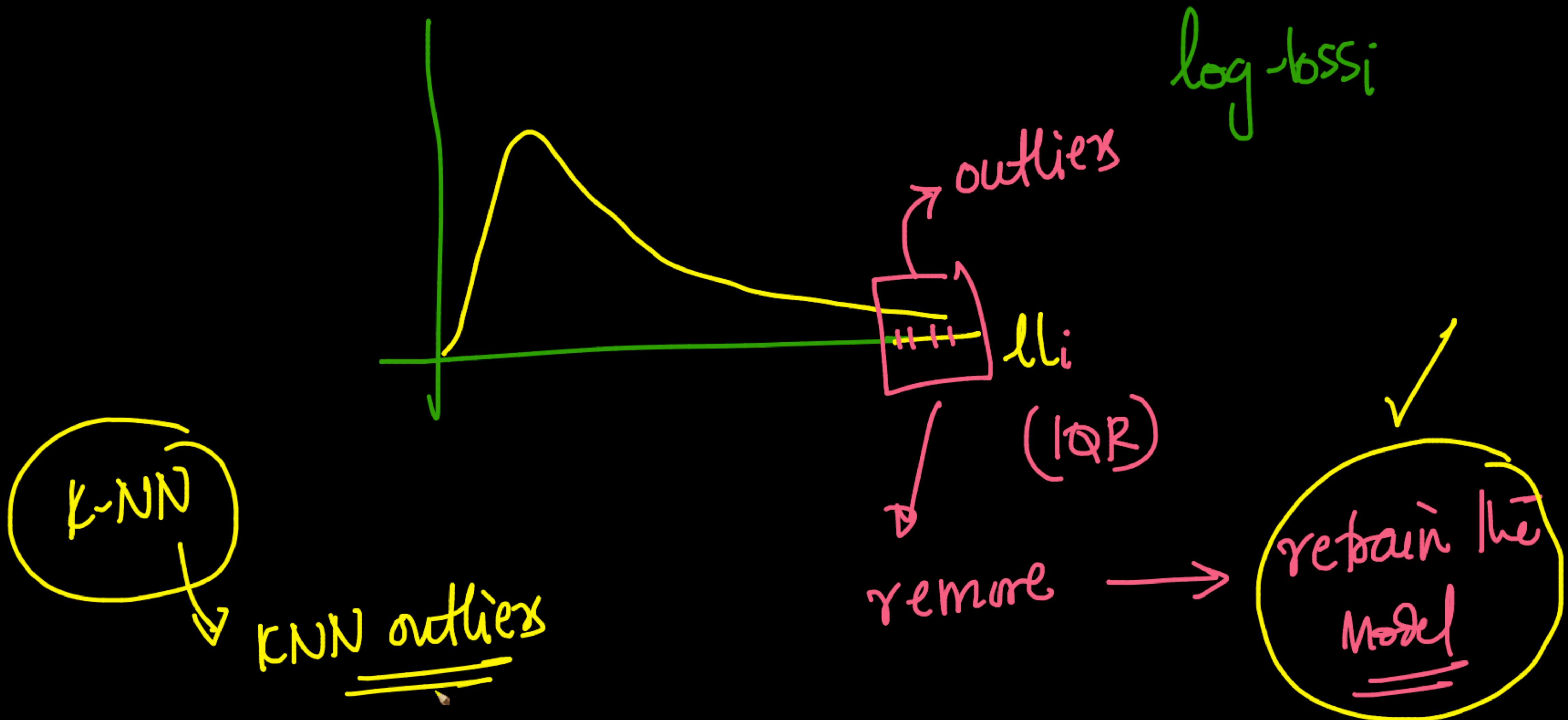
Logistic-regression

hint: outlier impact logistic regression

✓ IQR-based
PCR-column



model \rightarrow extreme outliers
 log-loss_i
 $\text{or } \underline{\Gamma(\omega^T x_i + b)} \forall i$



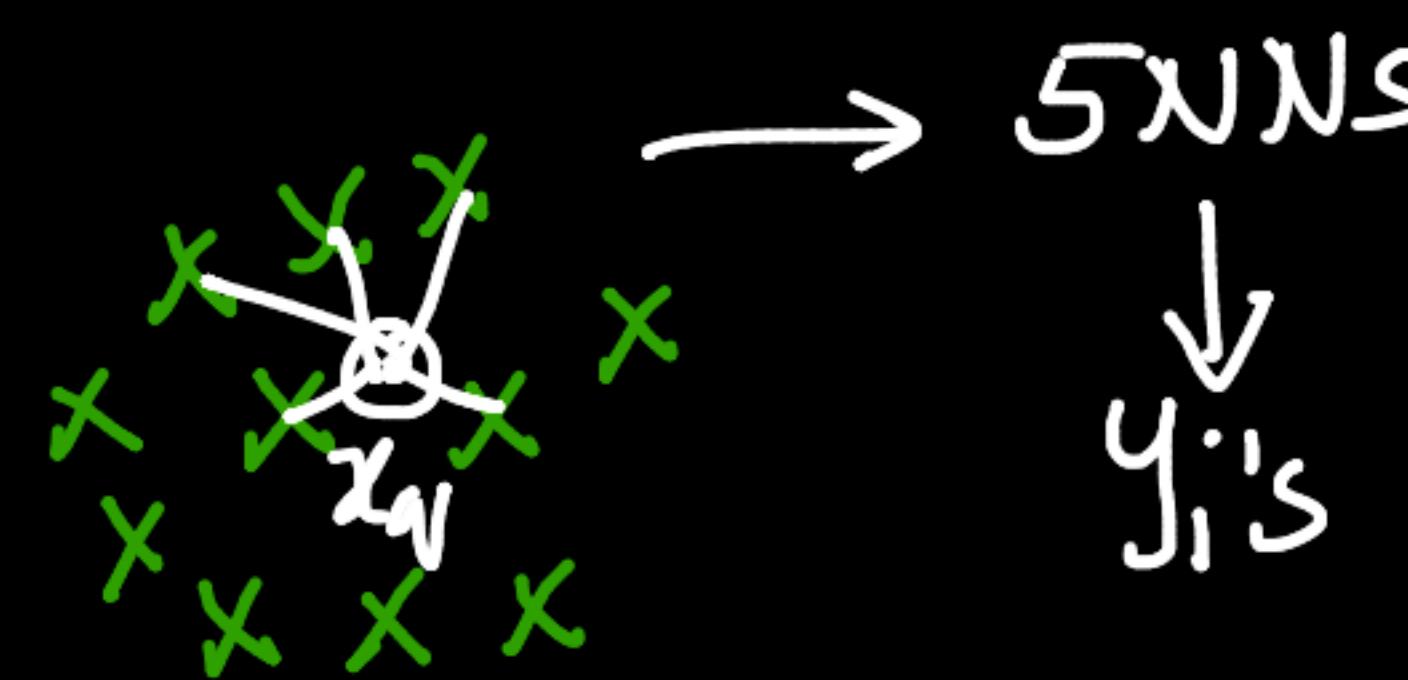
Task



outliers w.r.t. logistic regr- model



$k=5$



x_i 's

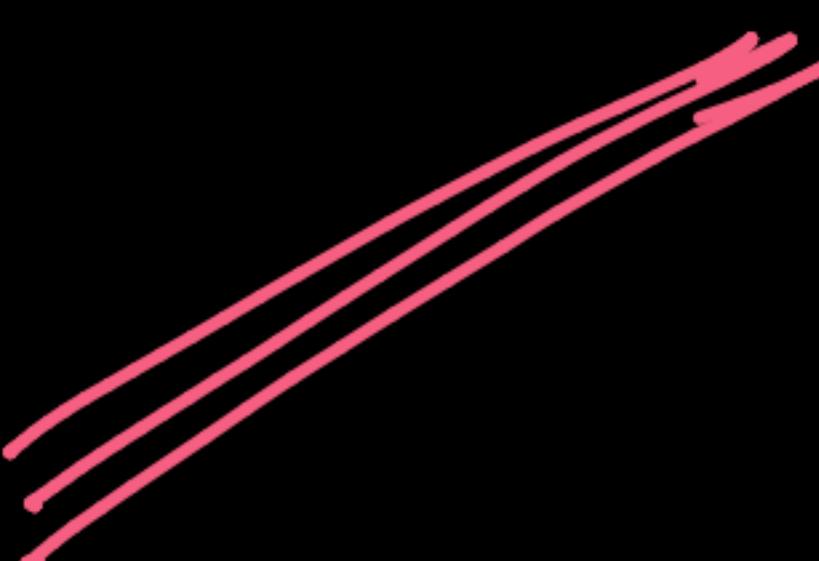
$y_i \in \mathbb{R}$

y_i 's \rightarrow mean/
median
 $= \underline{\underline{y_q}}$

KNN for
regression ✓
~~regression~~

logistic reg \rightarrow classfn
Linear reg \rightarrow reg

Decision Trees:



K-NN → locality

c: logistic reg] → Linear planes/fns
reg: linear Reg]

AUC (next
later)

binary-classfn
best: AUC → 1.0
random → 0.5

∞ imbalanced Data.ipynb - Colab | plot (-x*log(x)-(1-x)*log(1-x)) | sklearn.metrics.auc → scikit-learn

scikit-learn.org/stable/modules/generated/sklearn.metrics.auc.html

roc_auc_score

Compute the area under the ROC curve.

average_precision_score

Compute average precision from prediction scores.

precision_recall_curve

Compute precision-recall pairs for different probability thresholds.

Prev Up Next

scikit-learn 1.0.2

Other versions

Please [cite us](#) if you use the software.

[sklearn.metrics.auc](#)

Examples using

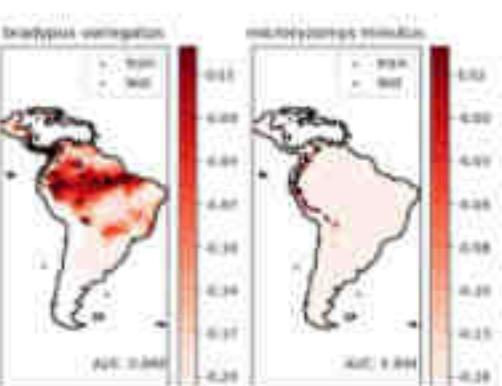
[sklearn.metrics.auc](#)

Examples

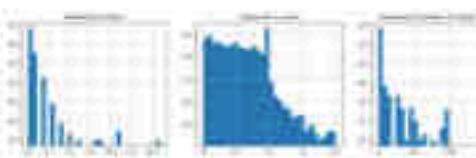
```
>>> import numpy as np
>>> from sklearn import metrics
>>> y = np.array([1, 1, 2, 2])
>>> pred = np.array([0.1, 0.4, 0.35, 0.8])
>>> fpr, tpr, thresholds = metrics.roc_curve(y, pred, pos_label=2)
>>> metrics.auc(fpr, tpr)
0.75
```

>>>

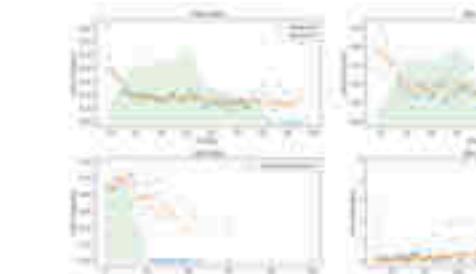
Examples using `sklearn.metrics.auc`



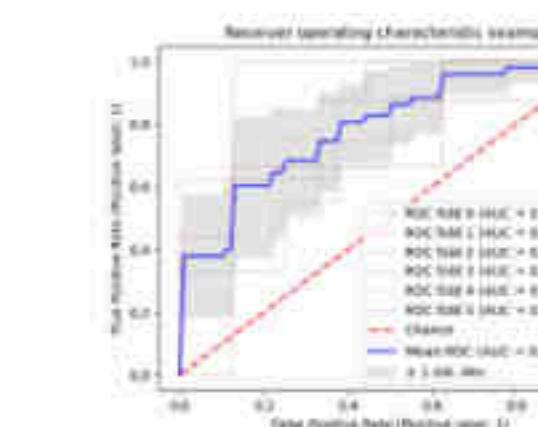
Species distribution modeling



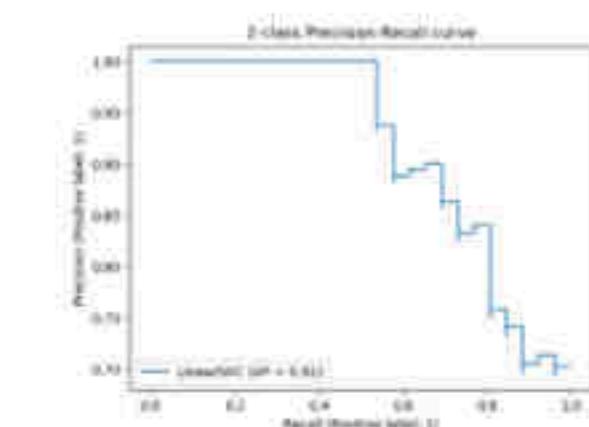
Poisson regression and non-normal loss



Tweedie regression on insurance claims



Receiver Operating Characteristic (ROC)



Precision-Recall

∞ imbalanced Data.ipynb - Colab | plot (-x*log(x)-(1-x)*log(1-x)) | sklearn.metrics.auc → scikit-learn

scikit-learn.org/stable/modules/generated/sklearn.metrics.auc.html

roc_auc_score

Compute the area under the ROC curve.

average_precision_score

Compute average precision from prediction scores.

precision_recall_curve

Compute precision-recall pairs for different probability thresholds.

Prev Up Next

scikit-learn 1.0.2

Other versions

Please [cite us](#) if you use the software.

[sklearn.metrics.auc](#)

Examples using

[sklearn.metrics.auc](#)

Examples

```
>>> import numpy as np
>>> from sklearn import metrics
>>> y = np.array([1, 1, 2, 2])
>>> pred = np.array([0.1, 0.4, 0.35, 0.8])
>>> fpr, tpr, thresholds = metrics.roc_curve(y, pred, pos_label=2)
>>> metrics.auc(fpr, tpr)
```

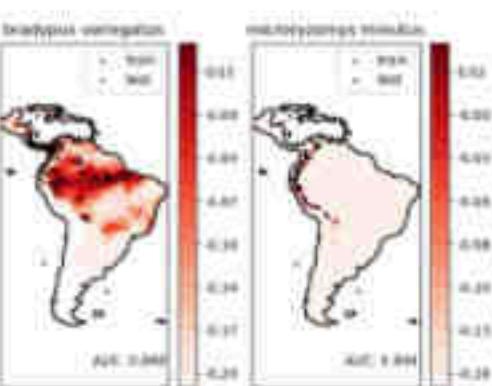
0.75

1.0

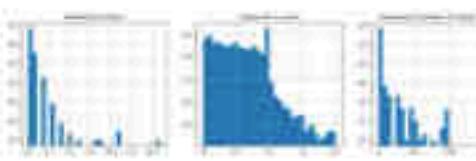
0.5

X
O

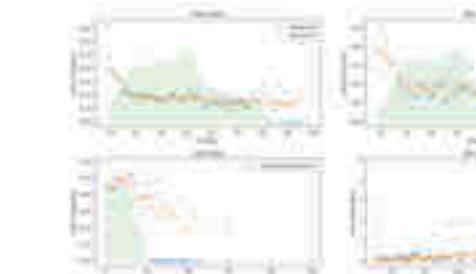
Examples using `sklearn.metrics.auc`



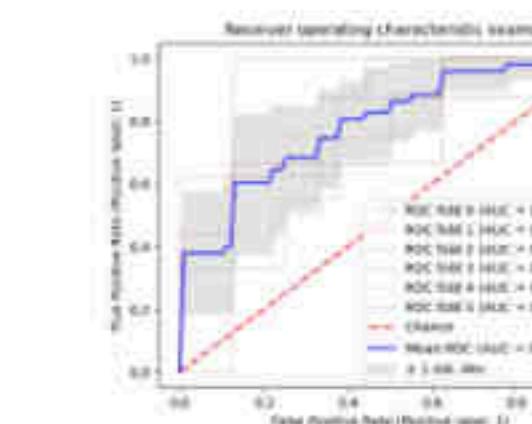
Species distribution modeling



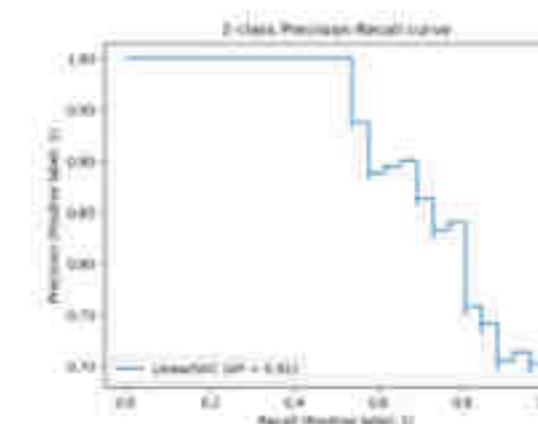
Poisson regression and non-normal loss



Tweedie regression on insurance claims



Receiver Operating Characteristic (ROC)



Precision-Recall

Decision Tree

if-else

Root

$f_1 < a$

Y

$f_2 < c$

Y

+ve

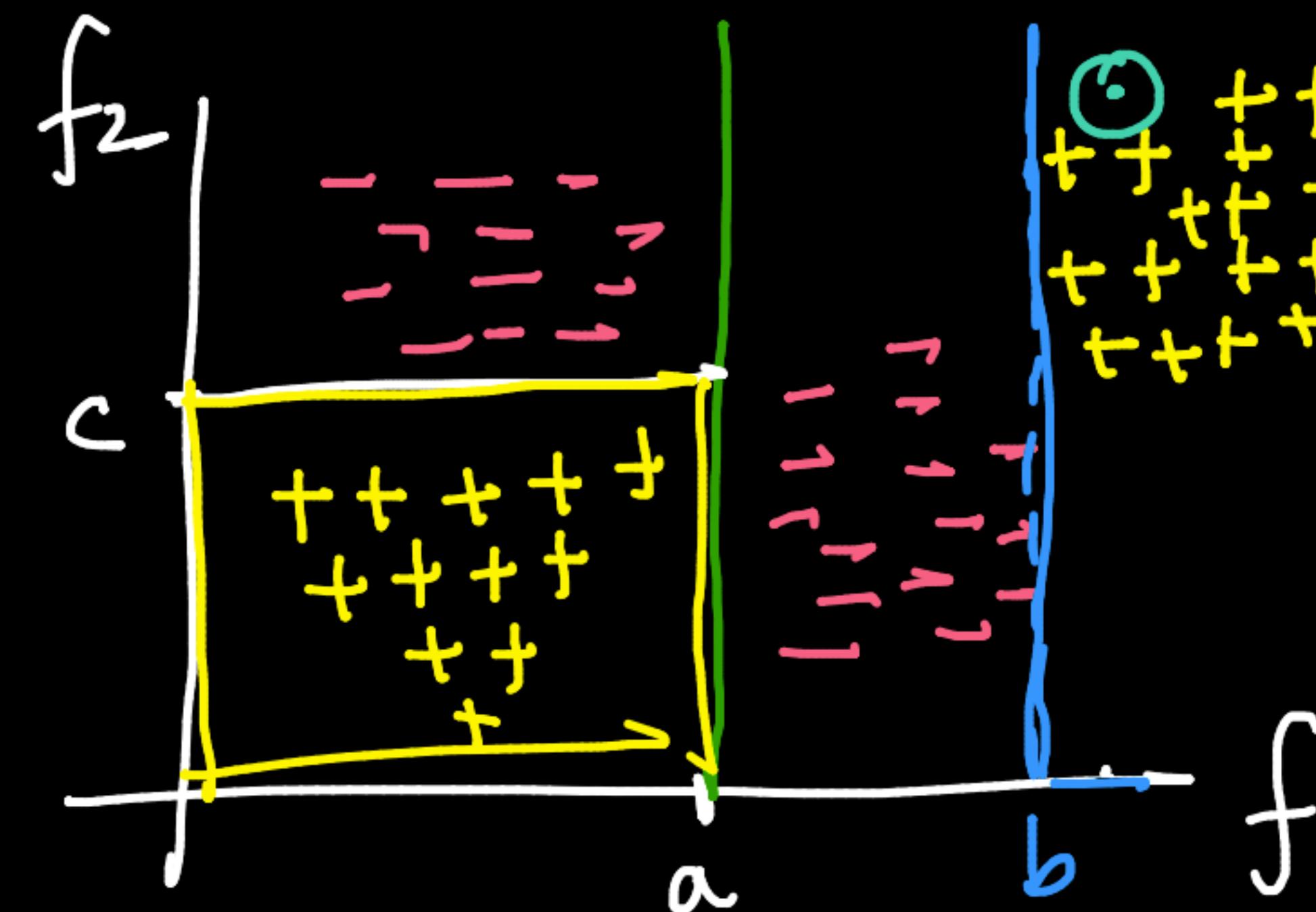
$f_1 < b$

N

-ve

internal-nodes

+ve → leaf-nodes



Nested if-else

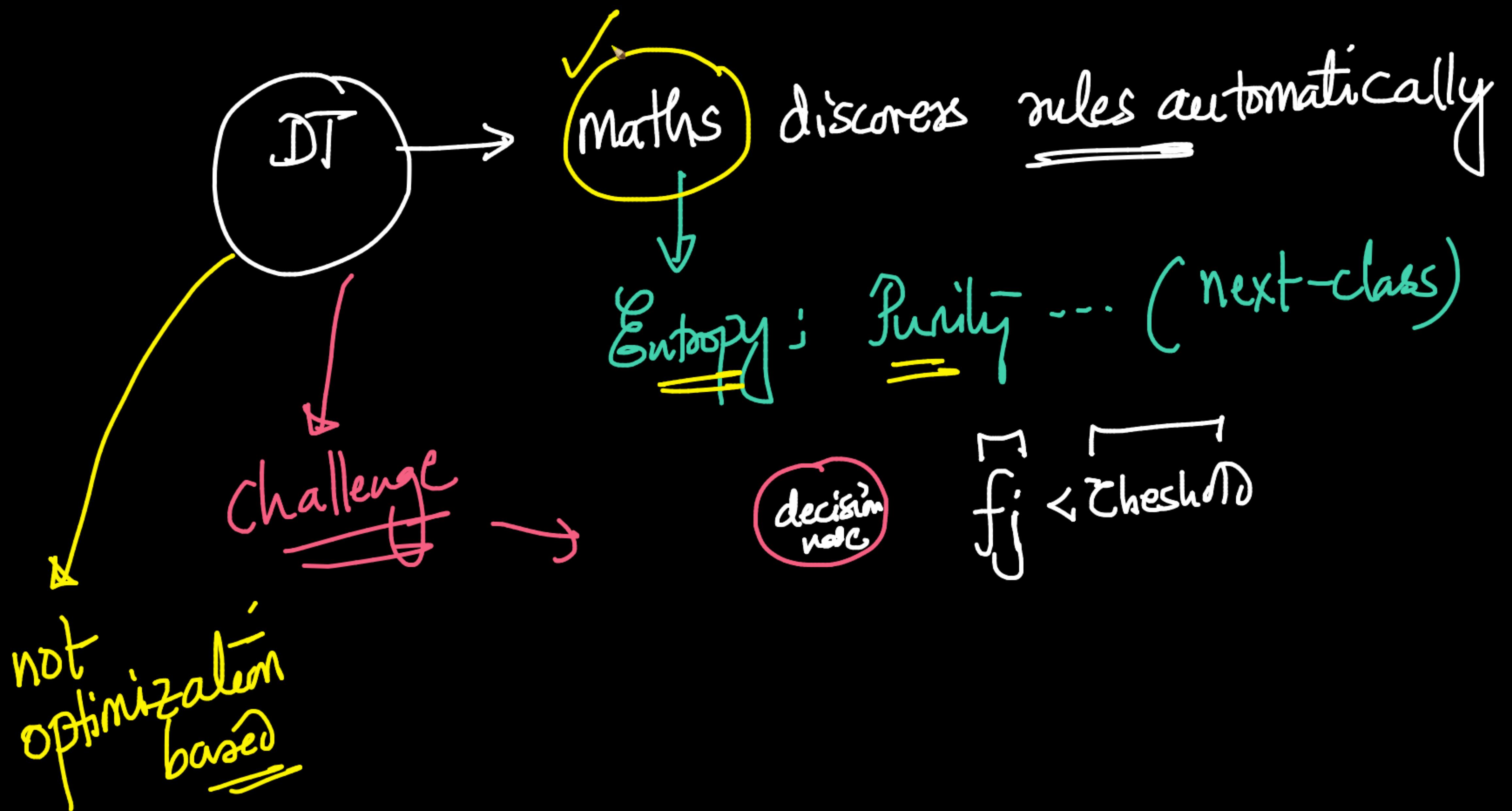


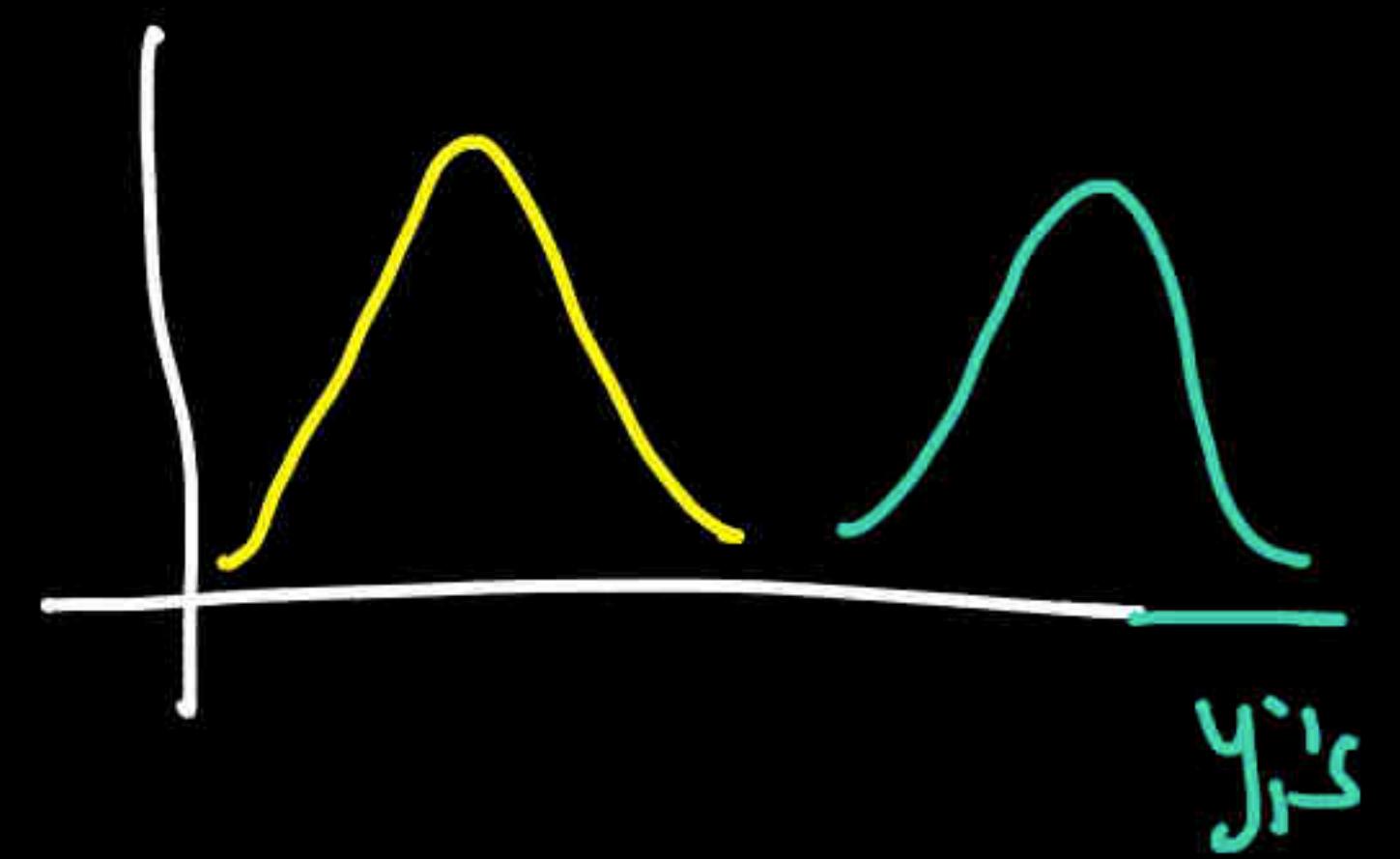
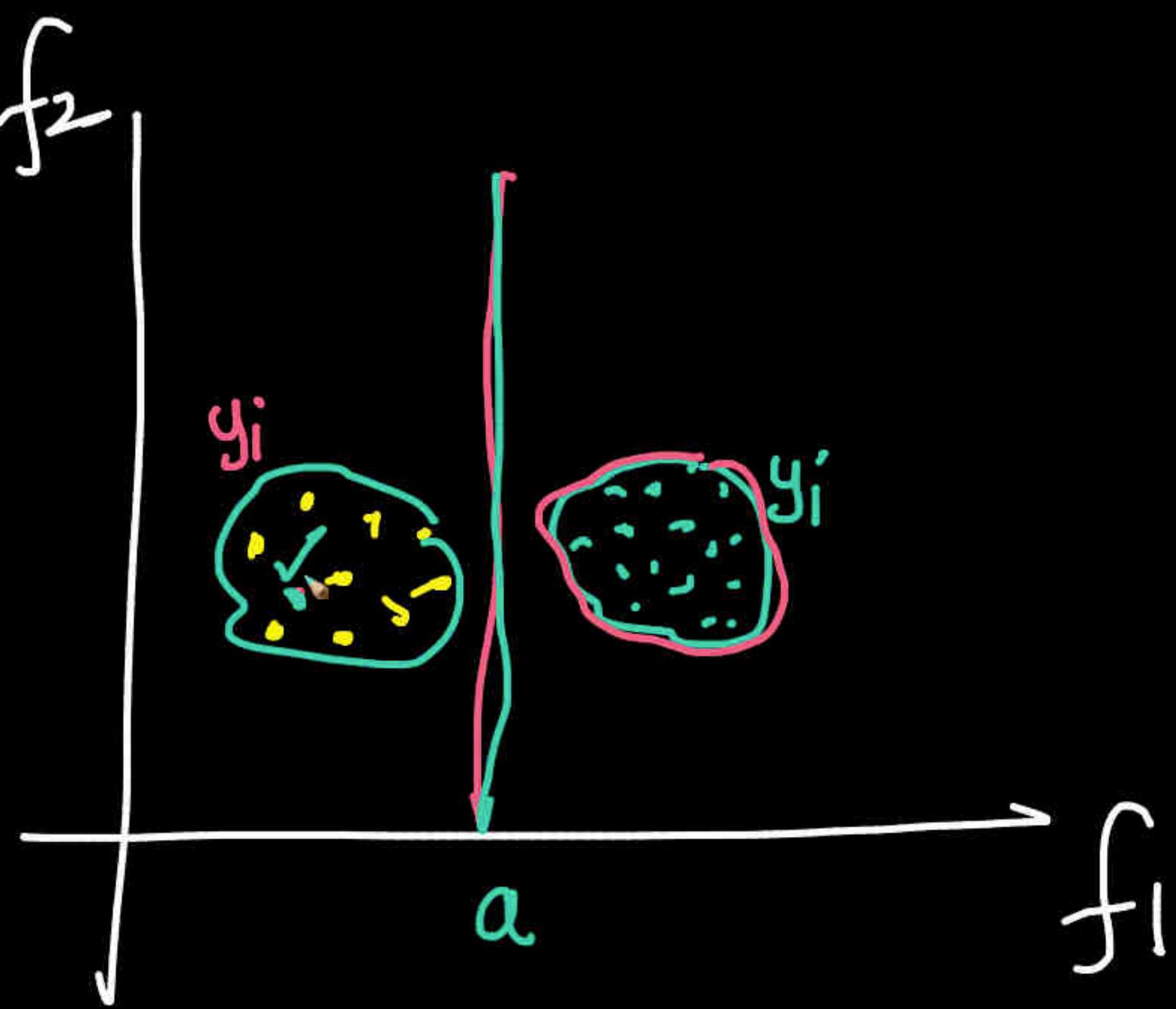
- KNN
- Logistic-Reg (poly-feature)

Rule-based system → nested if ... else
(Prog)

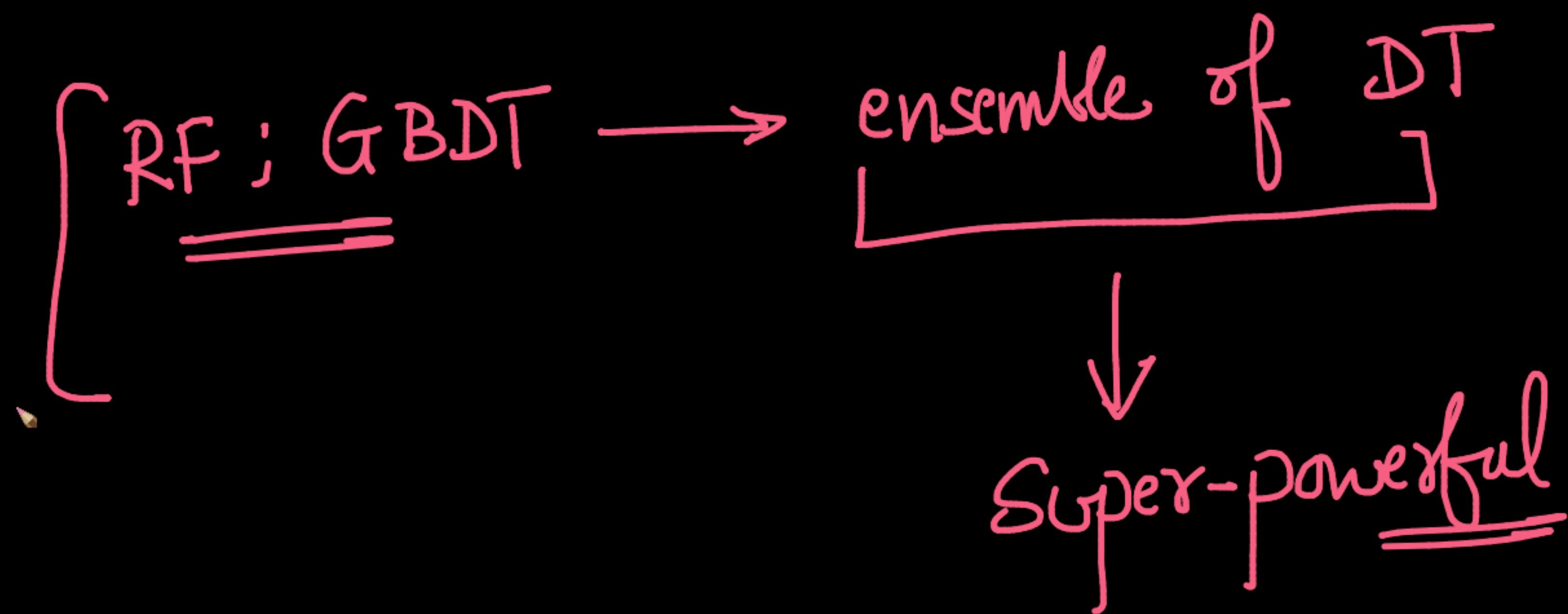


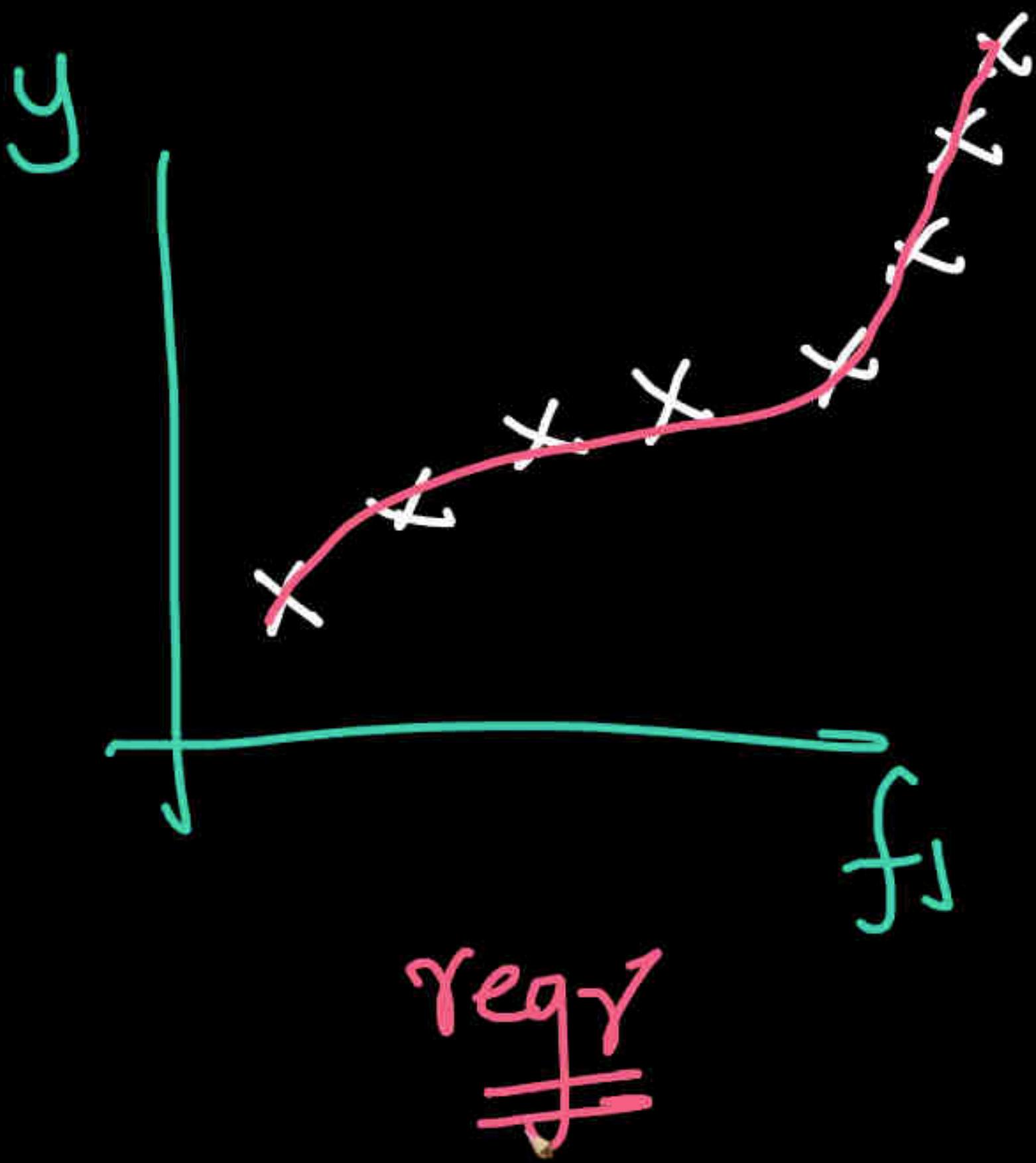
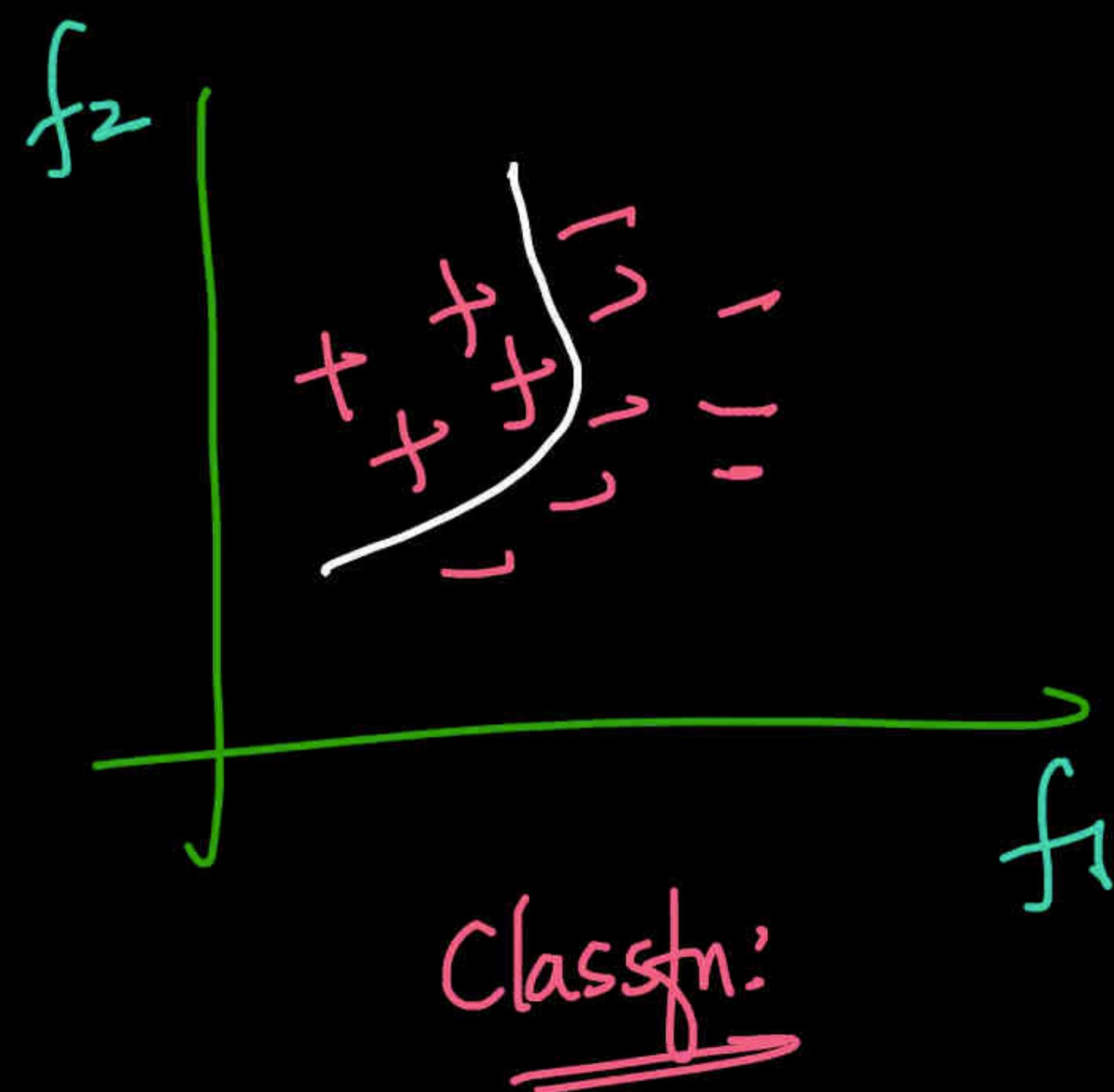
geom: DT is a collection of axis ||
hyperplanes





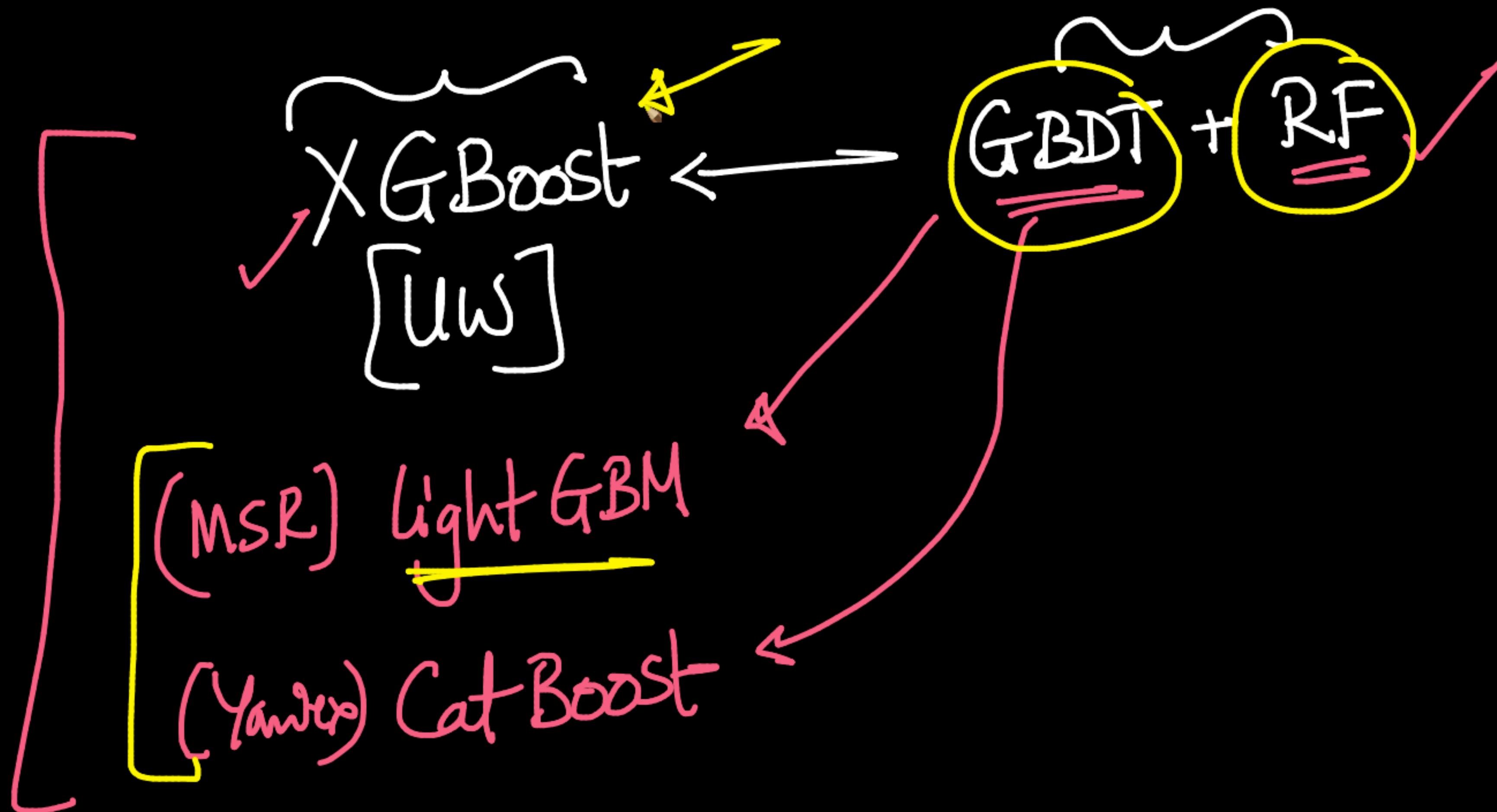
① DT → not very powerful





Next class:

- Entropy; build a DT
- AUC / ROC
- ensembles: RF



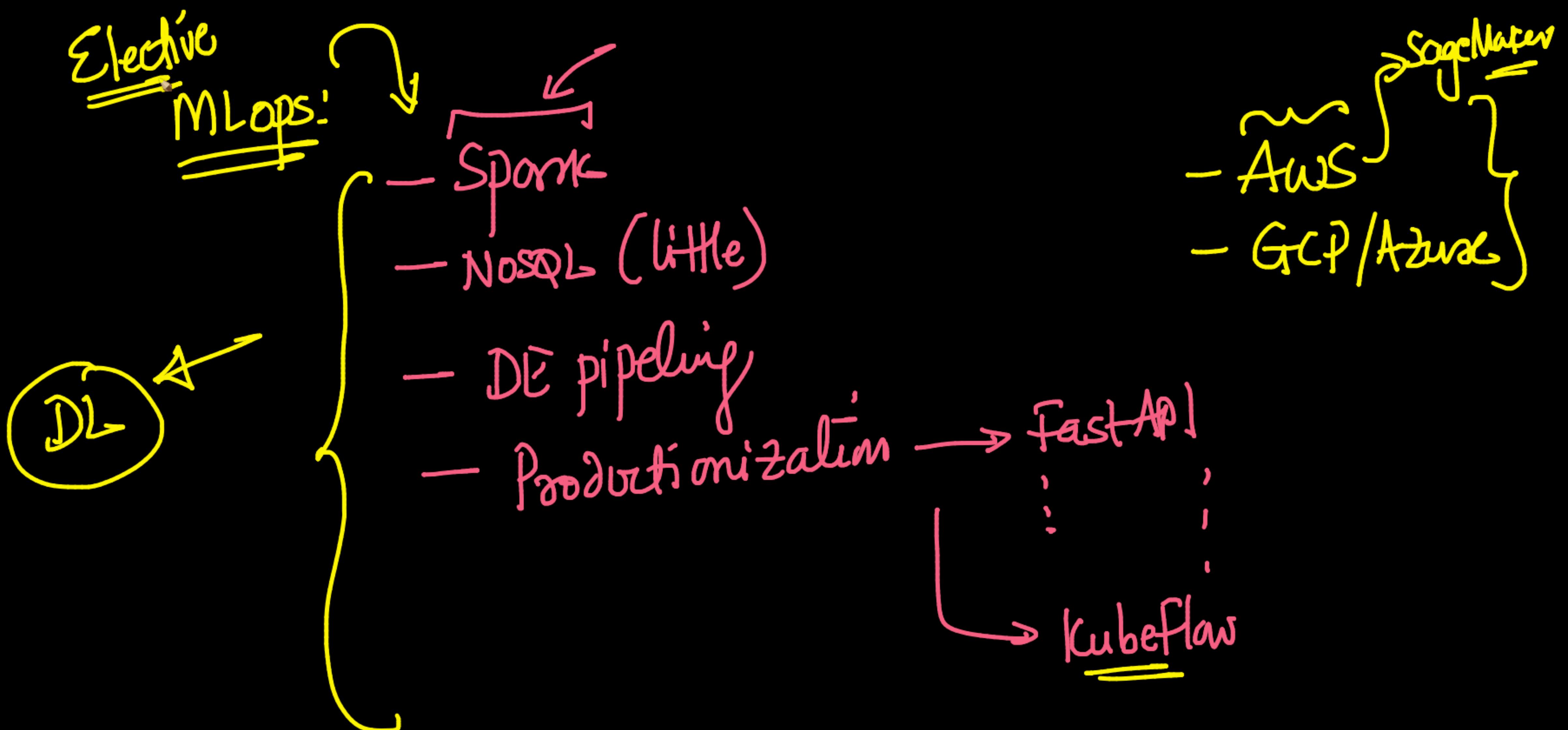


Tableau ; SQL
== ==

DA → :

imbalanced Data.ipynb - Colab plot (-x*log(x)-(1-x)*log(1-x)) LightGBM: A Highly Efficient Gradient Boosting Decision Tree

proceedings.neurips.cc/paper/2017/file/6449f44a102fde848669bdd9eb6b76fa-Paper.pdf

1 / 9 | - 150% + | ☰ ⌂

LightGBM: A Highly Efficient Gradient Boosting Decision Tree

1

2

3

4

LightGBM: A Highly Efficient Gradient Boosting Decision Tree

Guolin Ke¹, Qi Meng², Thomas Finley³, Taifeng Wang¹,
Wei Chen¹, Weidong Ma¹, Qiwei Ye¹, Tie-Yan Liu¹

¹Microsoft Research ²Peking University ³ Microsoft Redmond

¹{guolin.ke, taifengw, wche, weima, qiwye, tie-yan.liu}@microsoft.com;
²qimeng13@pku.edu.cn; ³tfinely@microsoft.com;

Abstract

Gradient Boosting Decision Tree (GBDT) is a popular machine learning algorithm, and has quite a few effective implementations such as XGBoost and pGBRT. Although many engineering optimizations have been adopted in these implementations, the efficiency and scalability are still unsatisfactory when the feature dimension is high and data size is large. A major reason is that for each feature, they need to scan all the data instances to estimate the information gain of all possible split points, which is very time consuming. To tackle this problem, we propose two novel techniques: *Gradient-based One-Side Sampling* (GOSS) and *Exclusive Feature Bundling* (EFB). With GOSS, we exclude a significant proportion of features from the computation of information gain. This is because most features play a more important role in the computation of information gain, GOSS can obtain quite accurate estimation of the information gain with a much smaller data size. With EFB, we bundle mutually exclusive features (i.e., they rarely take nonzero values simultaneously), to reduce the number of features. We prove that finding the optimal bundling of exclusive features is NP-hard, but a greedy algorithm



103 / 103