

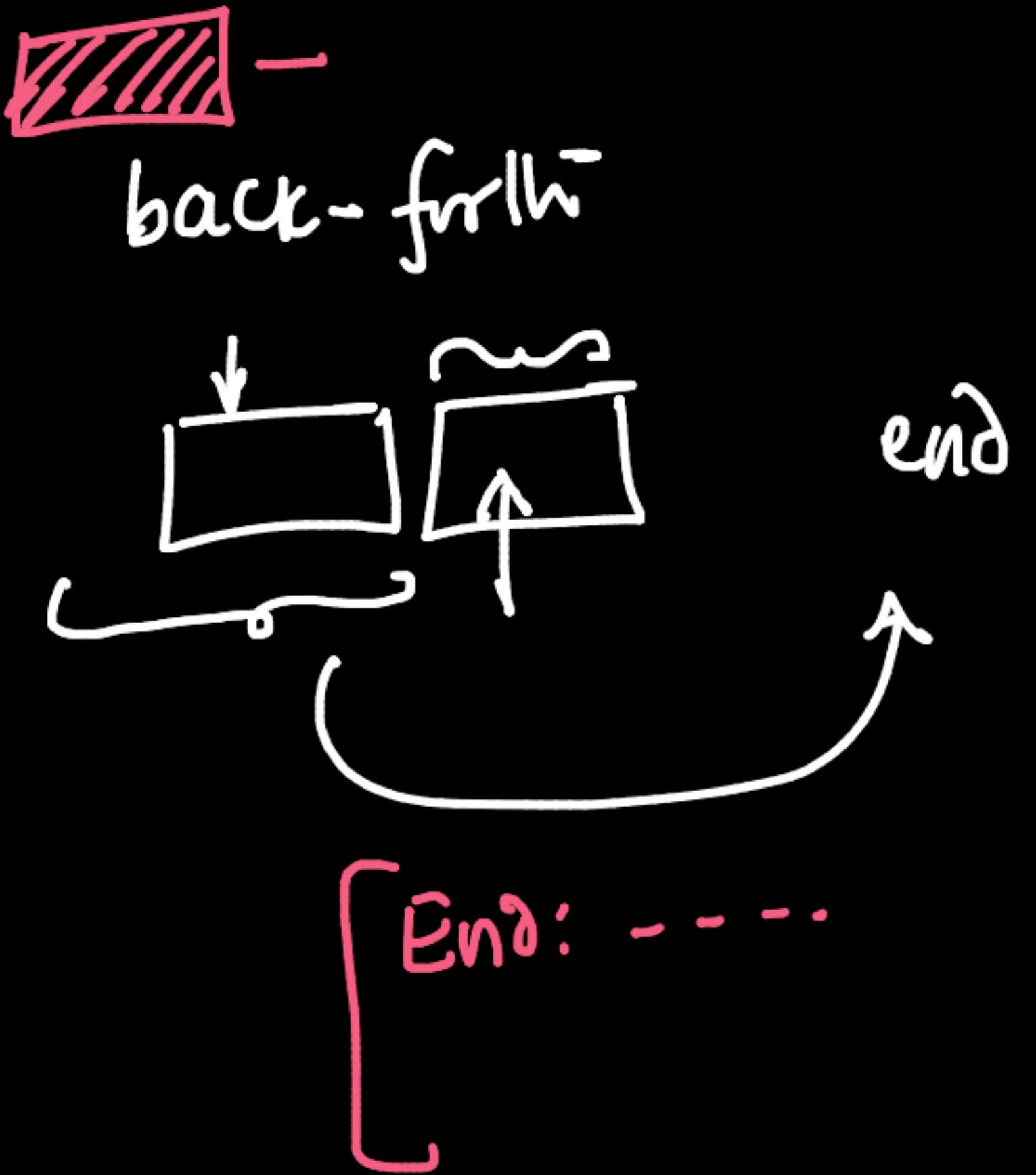
Topics:

- Box-plot; Bar-plot & scatter plot (retail-data) → P2005. 2000+
- Probability distributions
 - { - PMF, CDF
 - PDF
 - Bernoulli & Binomial dist
 - Air tickets overbooking

~~Ops:~~

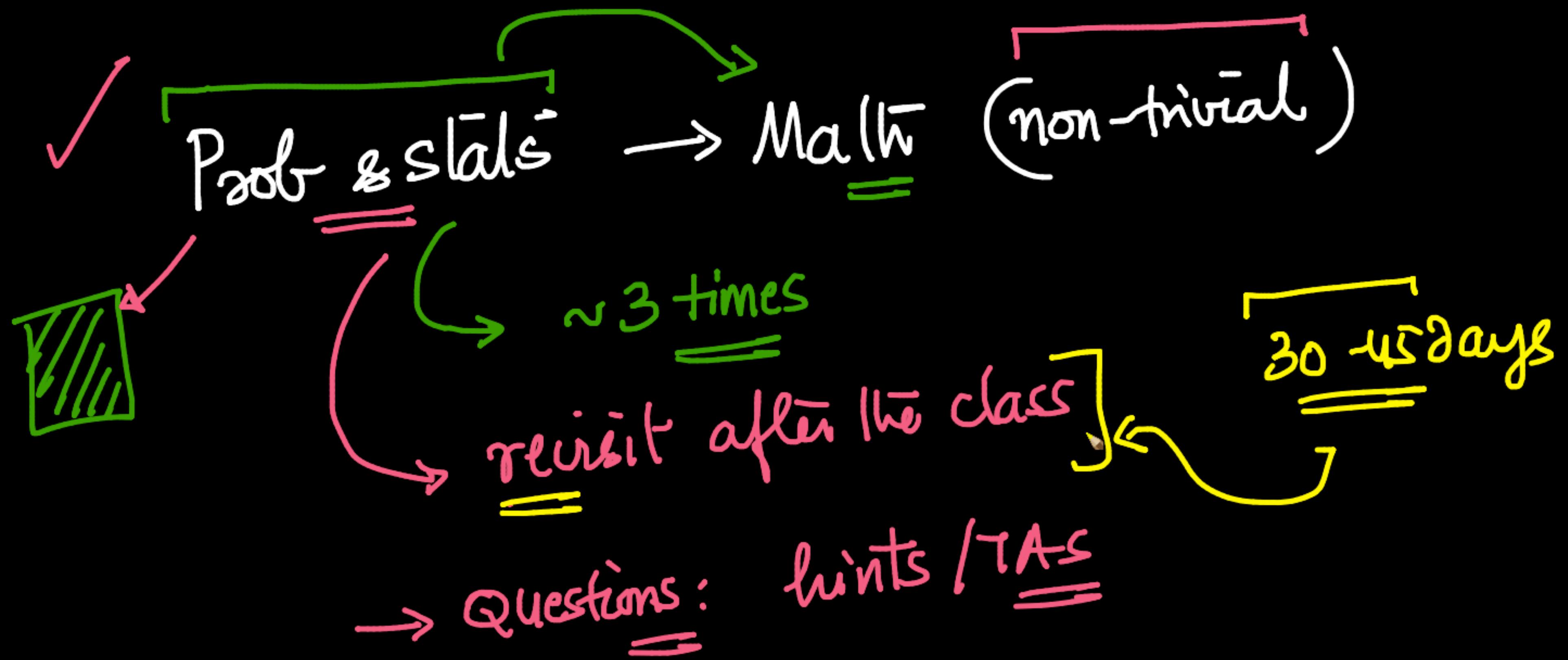
{ - Questions : Question - tab
- chat : communication; Yes/No

Psds: go-slow ; pick-up-speed



- scratch-pad; pot of code }
- future: detailed text-notes

→ TA (n) whatsapp/slack



Box-plot

Probability2.ipynb - Colaboratory

ProbabilityDisb_1.ipynb - Colaboratory

New Tab

RAM Disk

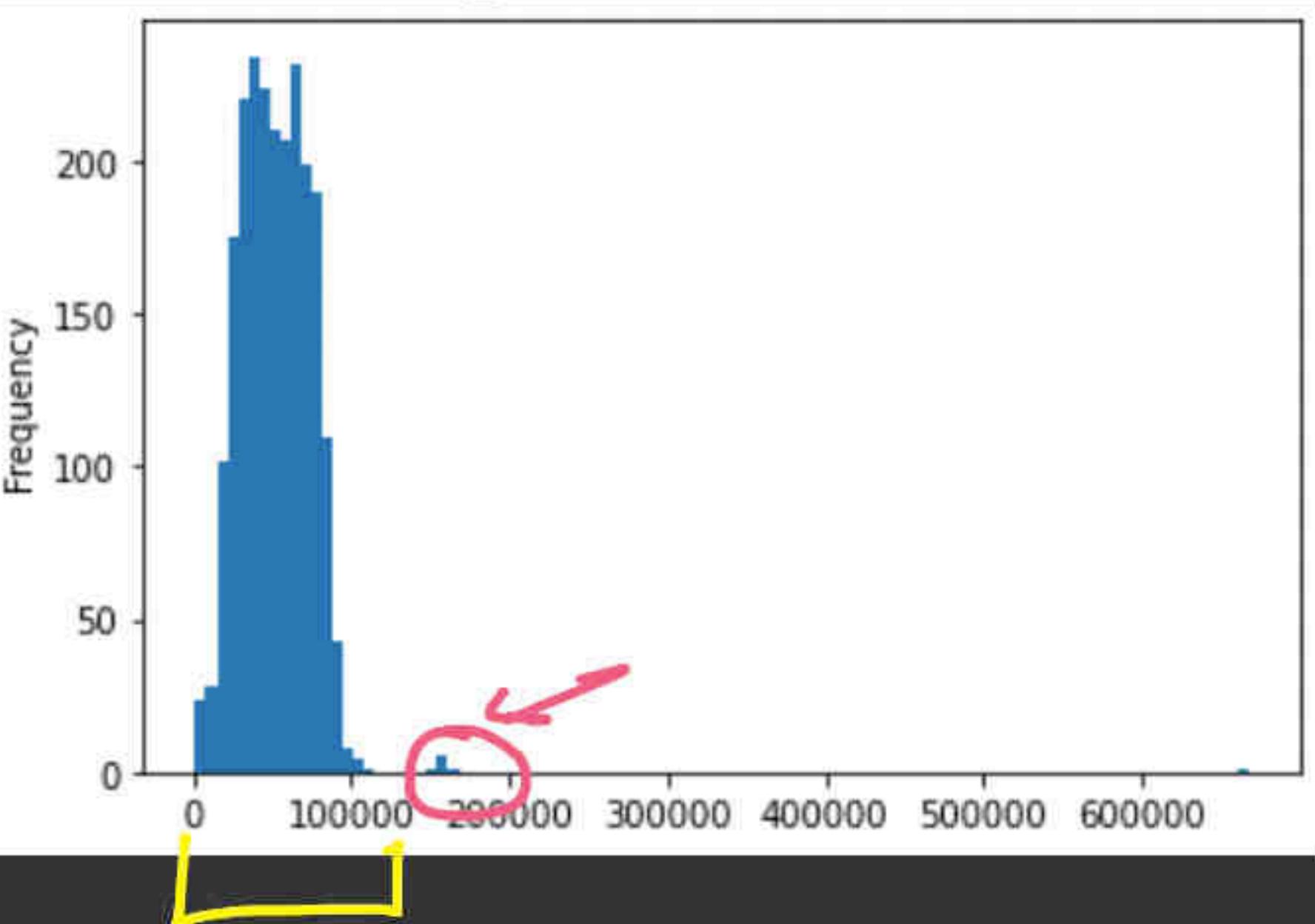
Update

colab.research.google.com/drive/15T7TVIAASw9Tdi4JxqJxuDqRgFBGOW3#scrollTo=SGqzqcJbVA9

+ Code + Text

```
[ ] df['Income'].plot.hist(bins=100)  
# few extreme values to the right.
```

{x} <matplotlib.axes._subplots.AxesSubplot at 0x7feb674d2450>

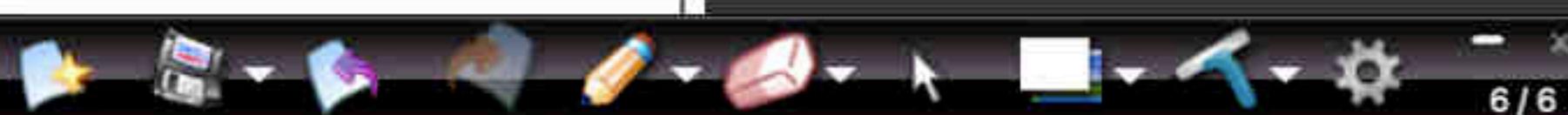


histogram
KDE

#box plot

```
<> ax = df['Income'].plot.box()  
ax.set_ylabel("Income in Dollars")
```

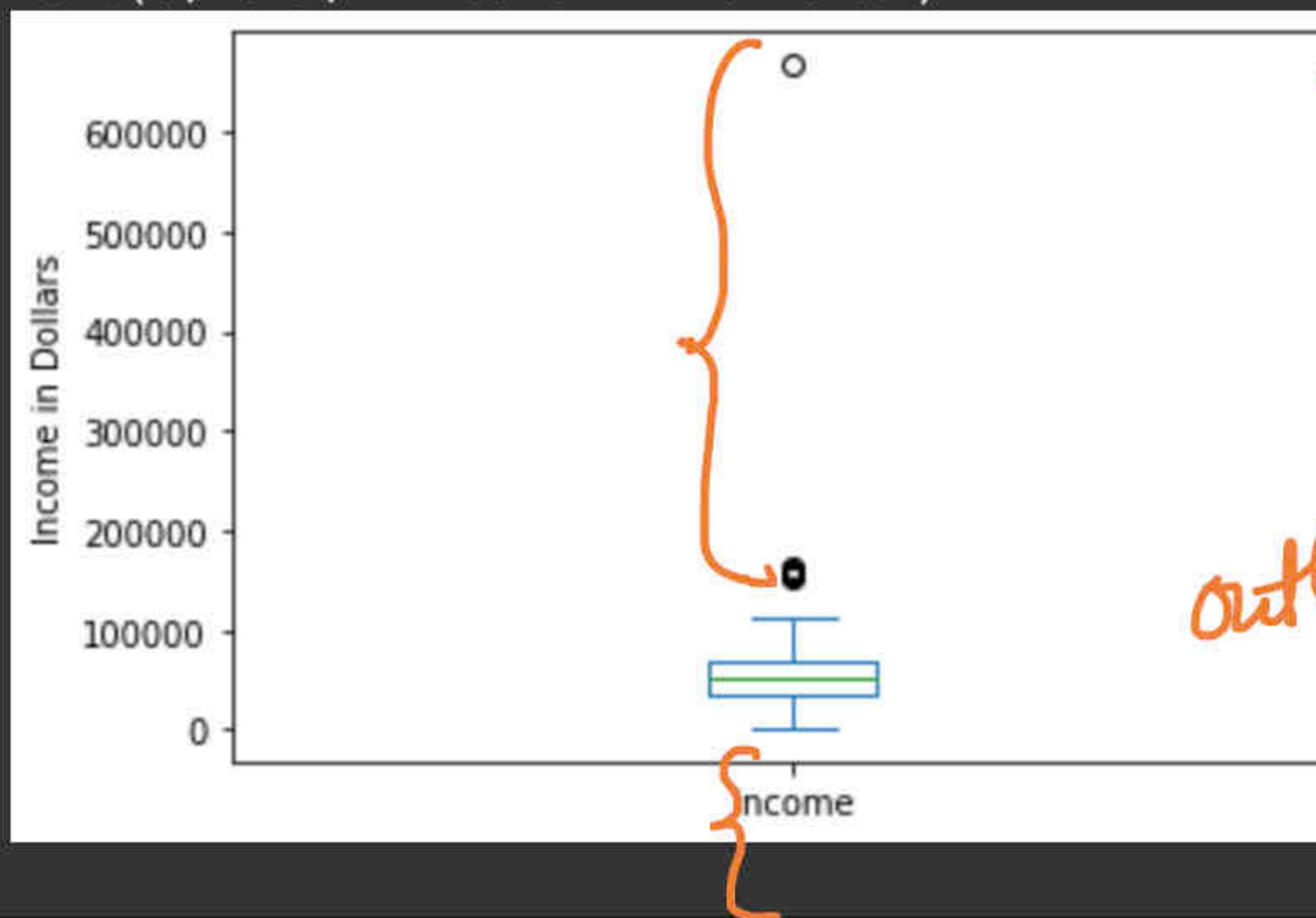
```
>- Text(0, 0.5, 'Income in Dollars')
```



+ Code + Text

```
[ ] #box plot  
ax = df['Income'].plot.box()  
ax.set_ylabel("Income in Dollars")
```

```
Text(0, 0.5, 'Income in Dollars')
```



outliers + {
outliers → $1.5 \text{ IQR} + Q_3$

whiskers ← $Q_3 = P_{75}$
box ← $Q_2 = P_{50}$
← $Q_1 = P_{25}$

outliers → $Q_1 - 1.5 \text{ IQR}$
Income

```
[ ] #box plot
```

```
ax = df['MntWInow'].plot.box()
```



Probability2.ipynb - Colaboratory

ProbabilityDisb_1.ipynb - Colaboratory

New Tab

colab.research.google.com/drive/1l5T7TVIAASw9Tdi4JxqJxuDqRgFBGOW3#scrollTo=SGqzqcJbVA9

Update

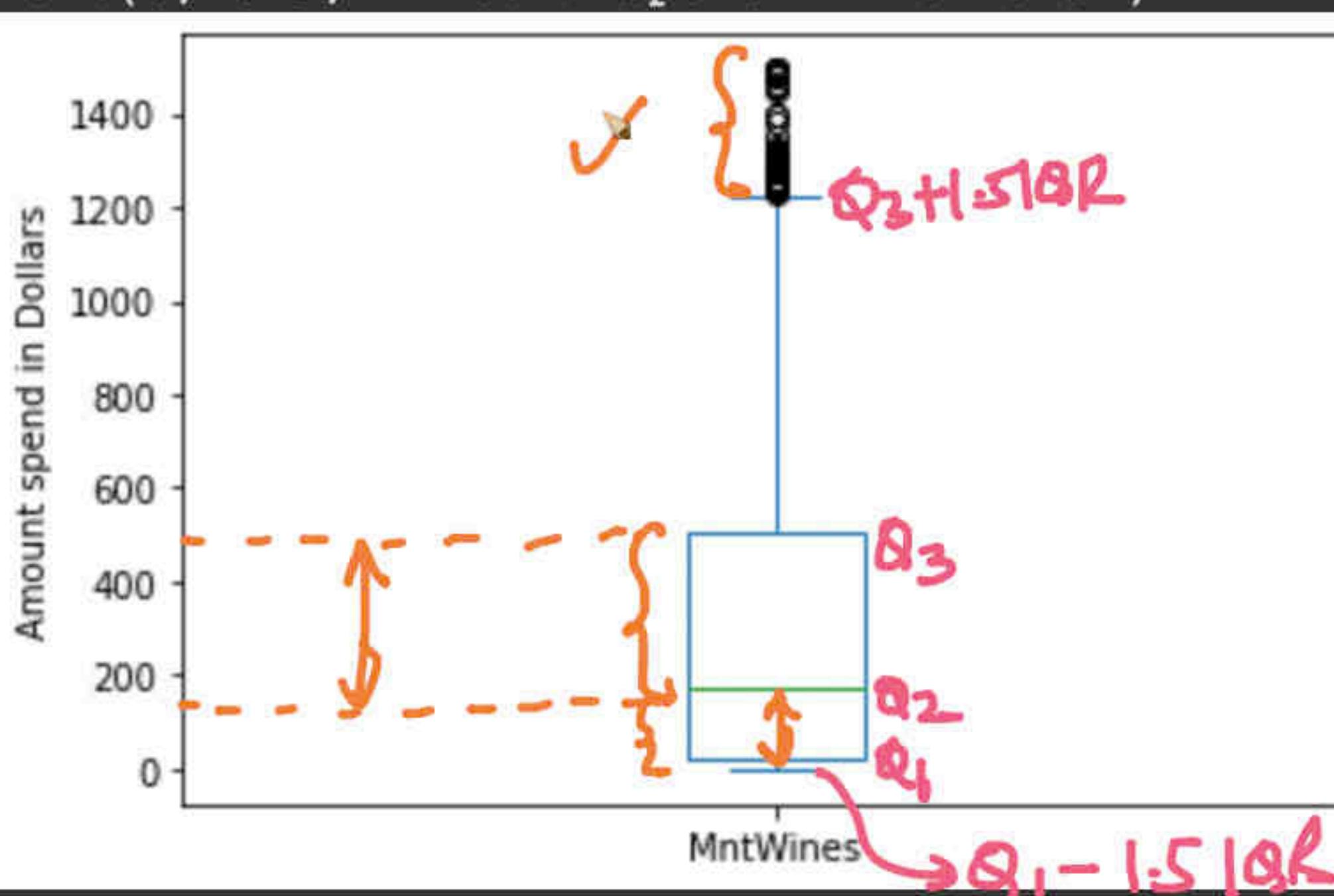
+ Code

+ Text

RAM
Disk

```
[ ] #box plot
ax = df['MntWines'].plot.box()
ax.set_ylabel("Amount spend in Dollars")
```

```
Text(0, 0.5, 'Amount spend in Dollars')
```



```
[ ] #bar plot
```

```
df['Education'].value_counts().plot.bar()
```



Probability2.ipynb - Colaboratory

ProbabilityDisb_1.ipynb - Colaboratory

New Tab

colab.research.google.com/drive/1l5T7TVIAASw9TdI4JxqJxuDqRgFBGOW3#scrollTo=SGqzqcJbVA9

Update

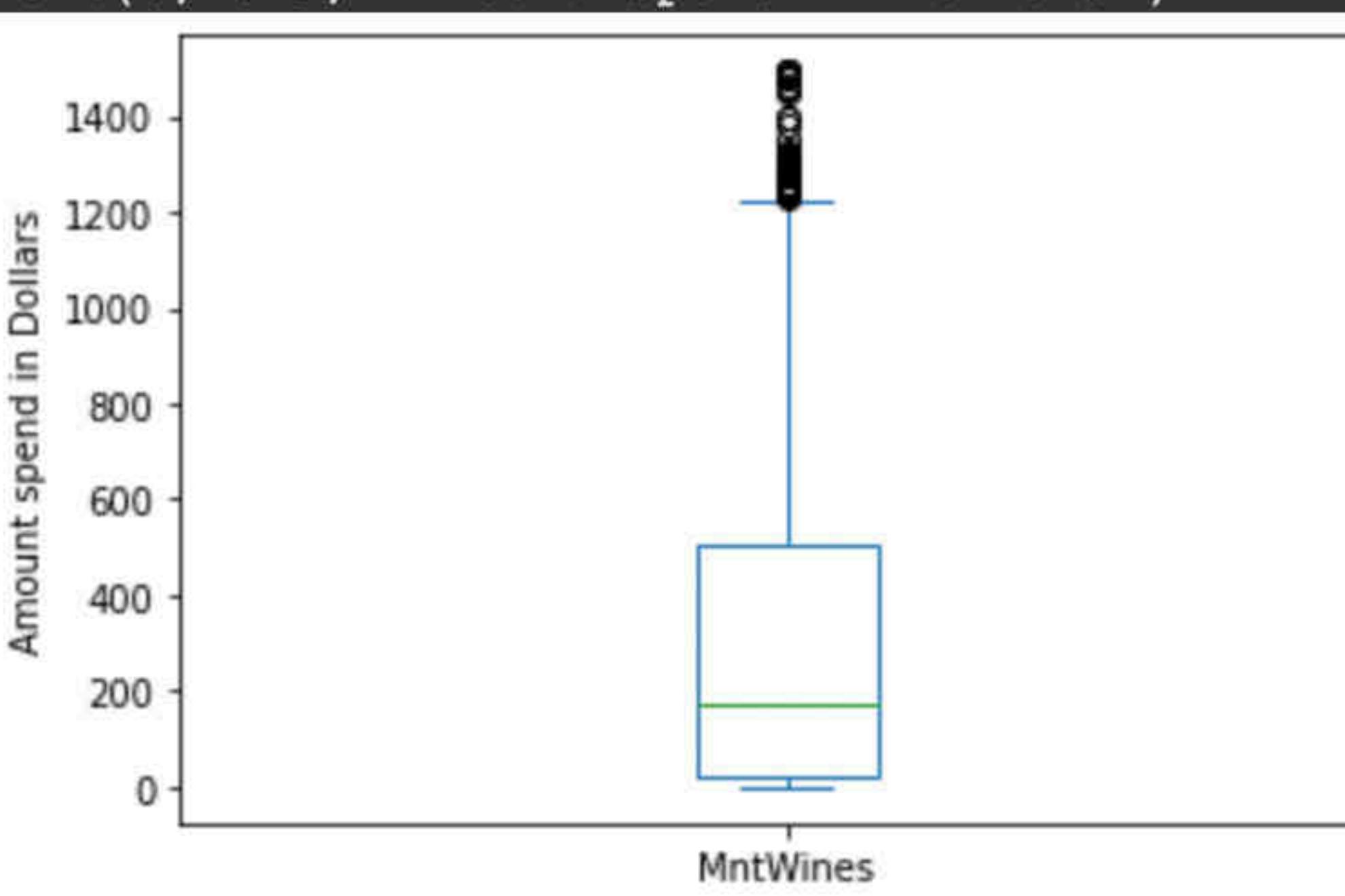
+ Code

+ Text

RAM
Disk

```
[ ] #box plot
ax = df['MntWines'].plot.box()
ax.set_ylabel("Amount spend in Dollars")
```

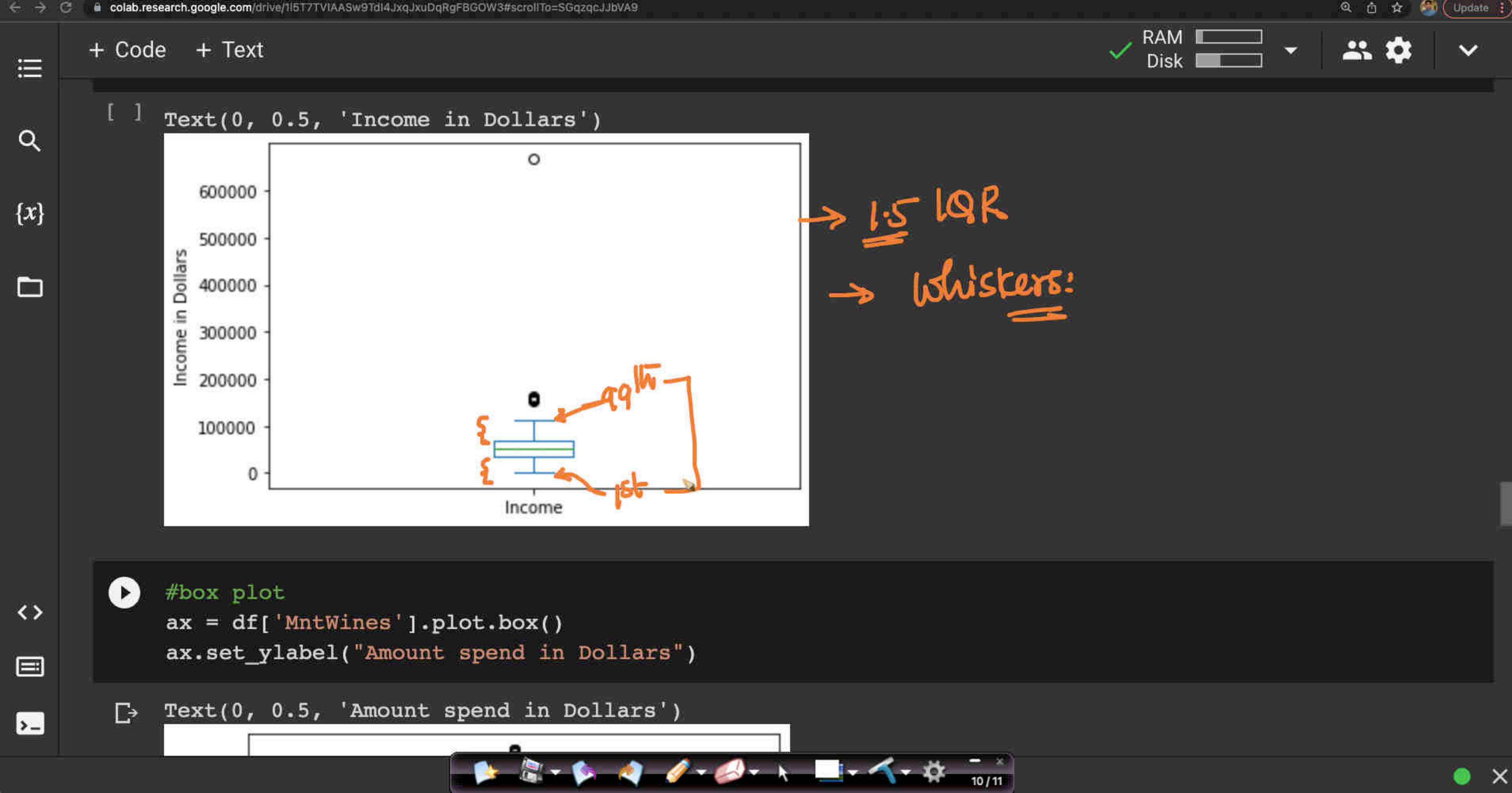
```
Text(0, 0.5, 'Amount spend in Dollars')
```



```
[ ] #bar plot
```

```
df['Education'].value_counts().plot.bar()
```





Probability2.ipynb - Colaboratory | ProbabilityDisb_1.ipynb - Colaboratory | New Tab | pandas.DataFrame.boxplot - Colaboratory

pandas.pydata.org/docs/reference/api/pandas.DataFrame.boxplot.html

Getting started User Guide API reference Development Release notes 1.4.2 Update

pandas

pandas.DataFrame.plot.hexbin
pandas.DataFrame.plot.hist
pandas.DataFrame.plot.kde
pandas.DataFrame.plot.line
pandas.DataFrame.plot.pie
pandas.DataFrame.plot.scatter
pandas.DataFrame.boxplot
pandas.DataFrame.hist
pandas.DataFrame.sparse.density
pandas.DataFrame.sparse.from_spmatrix
pandas.DataFrame.sparse.to_coo
pandas.DataFrame.sparse.to_dense
pandas.DataFrame.from_dict
pandas.DataFrame.from_records
pandas.DataFrame.to_parquet
pandas.DataFrame.to_pickle
pandas.DataFrame.to_csv
pandas.DataFrame.to_hdf
pandas.DataFrame.to_sql
pandas.DataFrame.to_dict
pandas.DataFrame.to_excel
pandas.DataFrame.to_json
pandas.DataFrame.to_html
pandas.DataFrame.to_feather

pandas.DataFrame.boxplot

`DataFrame.boxplot(column=None, by=None, ax=None, fontsize=None, rot=0, grid=True, figsize=None, layout=None, return_type=None, backend=None, **kwargs)` [source]

Make a box plot from DataFrame columns.

Make a box-and-whisker plot from DataFrame columns, optionally grouped by some other columns. A box plot is a method for graphically depicting groups of numerical data through their quartiles. The box extends from the Q1 to Q3 quartile values of the data, with a line at the median (Q2). The whiskers extend from the edges of box to show the range of the data. By default, they extend no more than 1.5 * IQR (IQR = Q3 - Q1) from the edges of the box, ending at the farthest data point within that interval. Outliers are plotted as separate dots.

For further details see Wikipedia's entry for [boxplot](#).

Parameters: `column : str or list of str, optional`
Column name or list of names, or vector. Can be any valid input to `pandas.DataFrame.groupby()`.

`by : str or array-like, optional`
Column in the DataFrame to `pandas.DataFrame.groupby()`. One box-plot will be done per value of columns in `by`.

`ax : object of class matplotlib.axes.Axes, optional`
The matplotlib axes to be used by boxplot.



11 / 12

Probability2.ipynb - Colaboratory ProbabilityDisb_1.ipynb - Colaboratory New Tab

colab.research.google.com/drive/1l5T7TVIAASw9Tdi4JxqJxuDqRgFBGOW3#scrollTo=qhN6-BWVpojt

Update

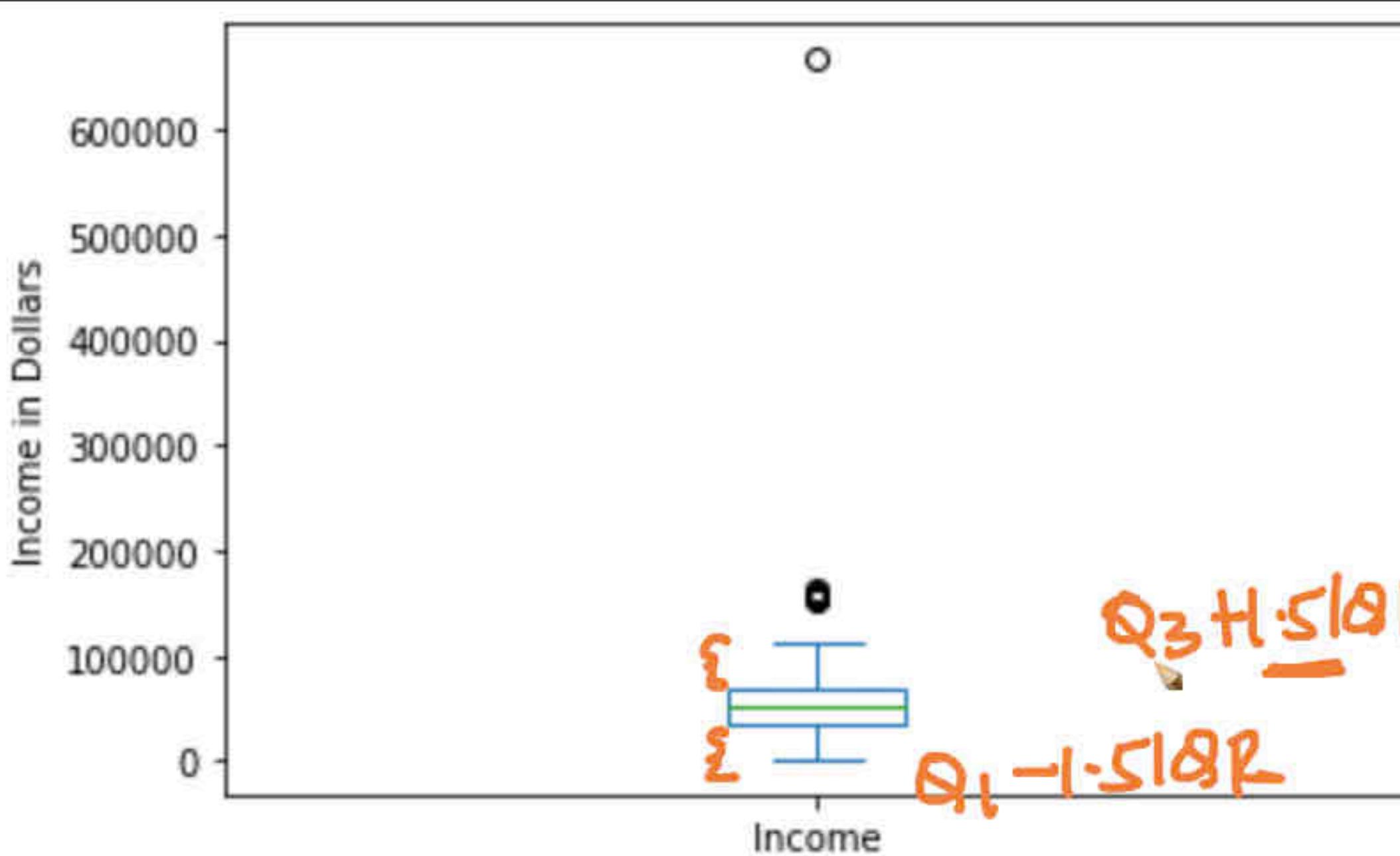
+ Code + Text

RAM
Disk



```
ax = df['Income'].plot.box()  
ax.set_ylabel("Income in Dollars")
```

{x} Text(0, 0.5, 'Income in Dollars')



```
[ ] #box plot  
ax = df['MntWines'].plot.box()  
ax.set_ylabel("Amount spend in Dollars")
```



Probability2.ipynb - Colaboratory

ProbabilityDisb_1.ipynb - Colaboratory

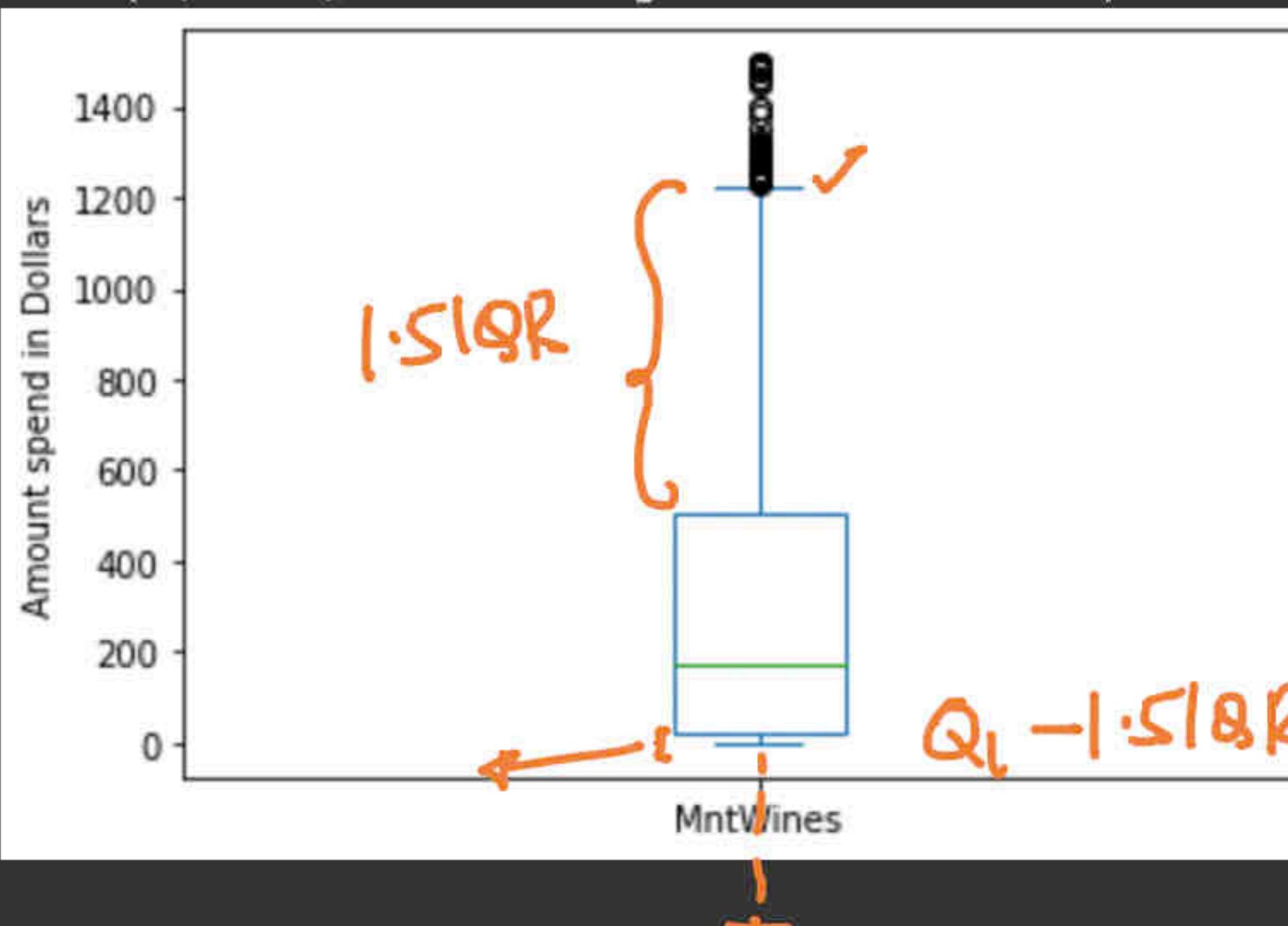
New Tab

pandas.DataFrame.boxplot - Colaboratory

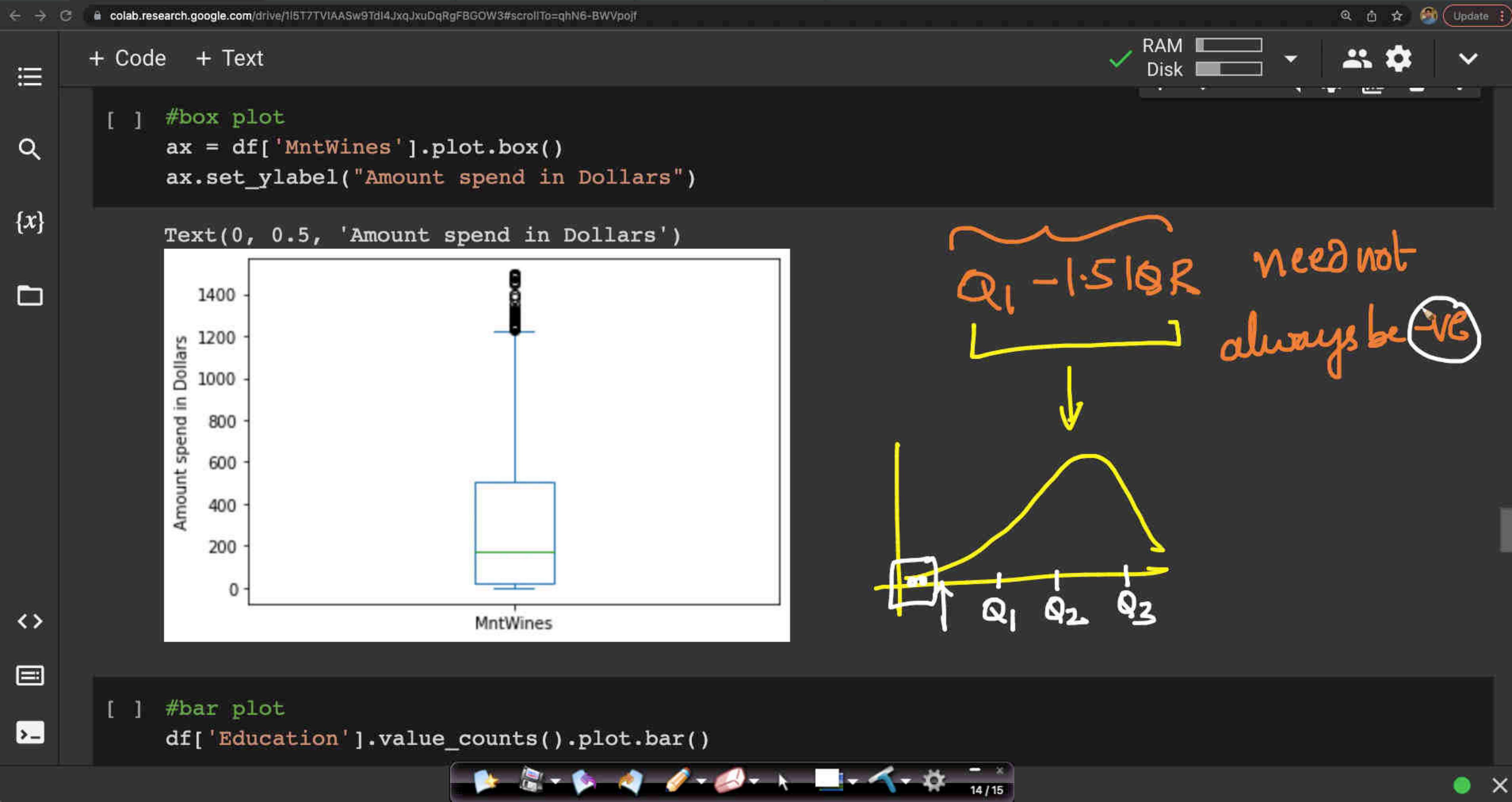
+

```
[ ] #box plot
ax = df['MntWines'].plot.box()
ax.set_ylabel("Amount spend in Dollars")
```

Text(0, 0.5, 'Amount spend in Dollars')

pandas: I.5IQR (Q₃) Max(min)Q_L - I.5IQR

```
[ ] #bar plot
df['Education'].value_counts().plot.bar()
```



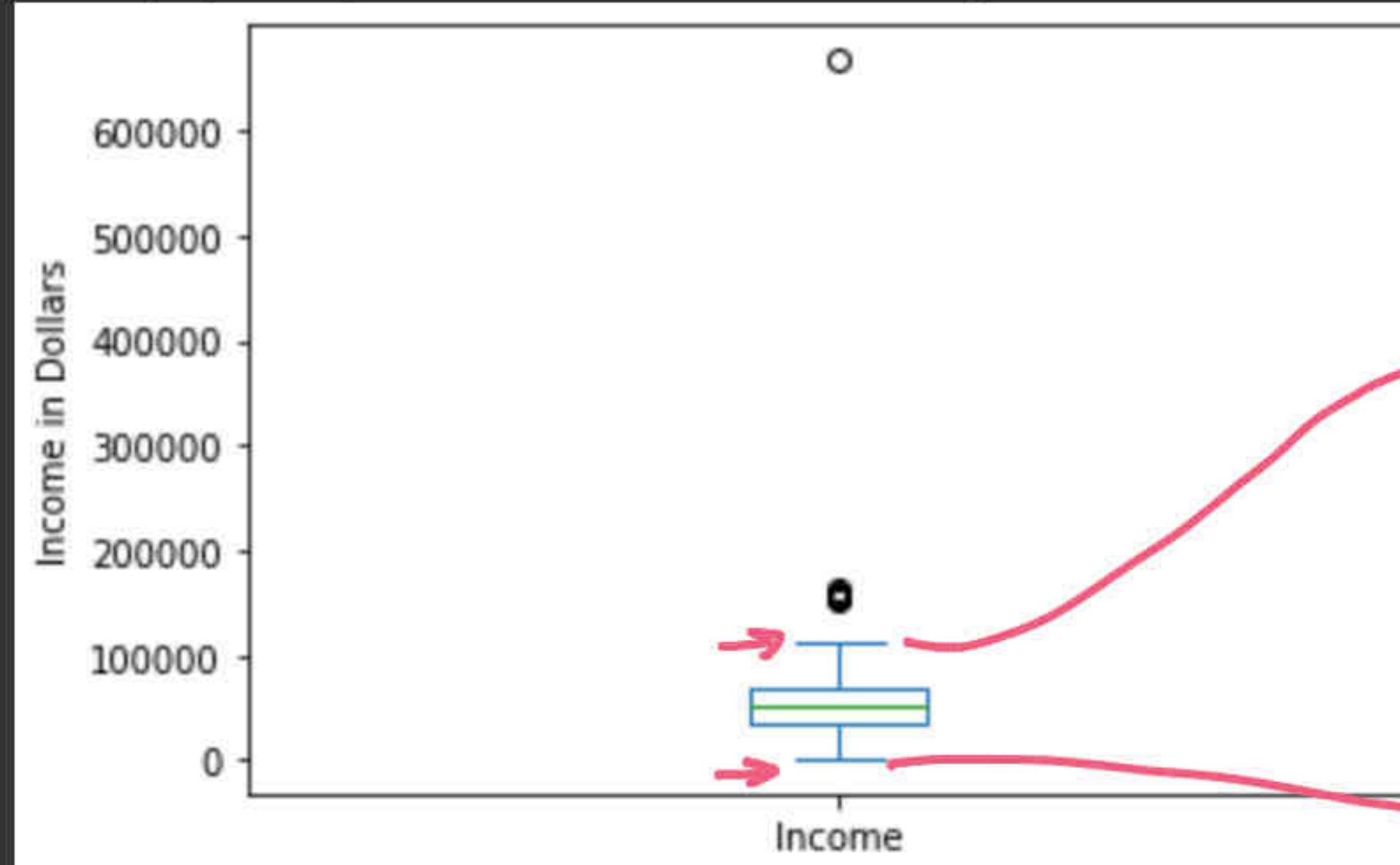
Probability2.ipynb - Colaboratory X ProbabilityDisb_1.ipynb - Colaboratory X New Tab

colab.research.google.com/drive/15T7TVIAASw9Tdi4JxqJxuDqRgFBGOW3#scrollTo=qhN6-BWVpojt

RAM Disk

+ Code + Text

Text(0, 0.5, 'Income in Dollars')



pandas:

min

UW = $(Q_3 + 1.5 \cdot IQR, \text{Max r.v.})$

LW = $\max(Q_1 - 1.5 \cdot IQR, \text{Min r.v.})$

min
r.v.

#box plot
ax = df['MntWines'].plot.box()
ax.set_ylabel("Amount spend in Dollars")

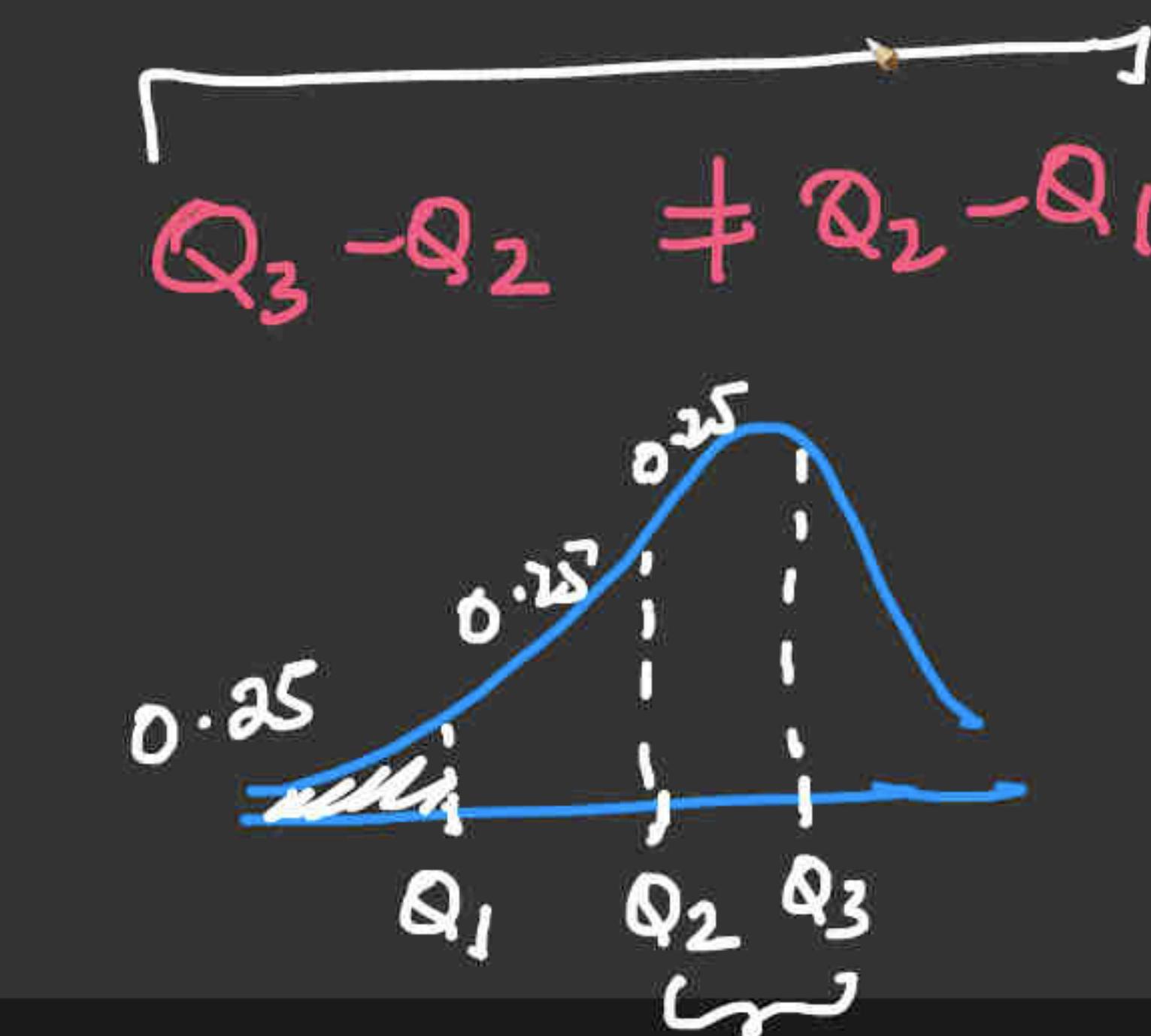
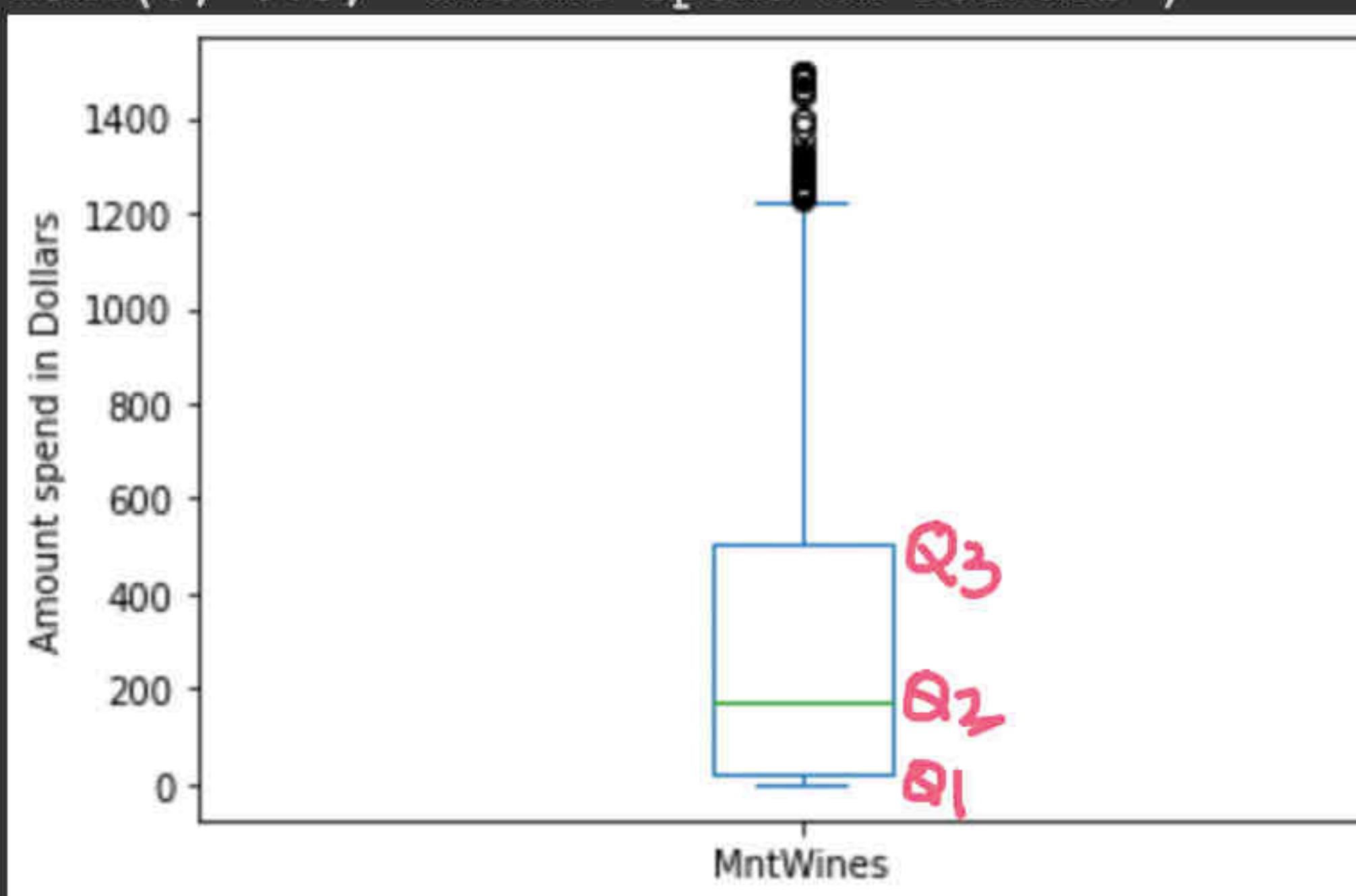
Text(0, 0.5, 'Amount spend in Dollars')

+ Code + Text

RAM
Disk



```
[ ] Text(0, 0.5, 'Amount spend in Dollars')
```



```
<cell> #bar plot  
df['Education'].value_counts().plot.bar()
```

```
<cell> <matplotlib.axes._subplots.AxesSubplot at 0x7feb53e3f9d0>
```

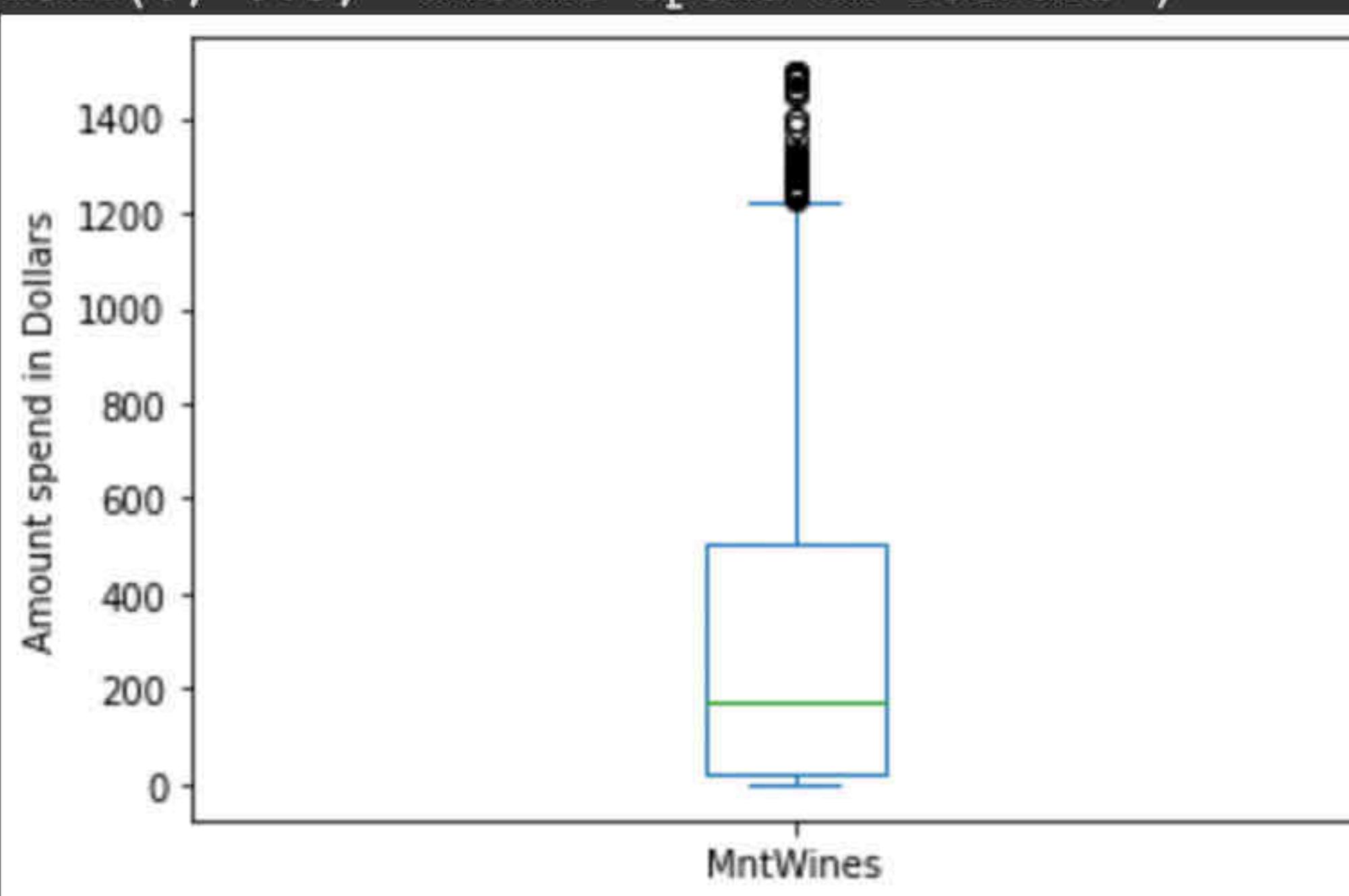


+ Code + Text

RAM
Disk



```
[ ] Text(0, 0.5, 'Amount spend in Dollars')
```



box-plot → Ques ...

```
[ ] #bar plot  
df['Education'].value_counts().plot.bar()
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7feb53e3f9d0>
```



RAM Disk

+ Code + Text

[] #bar plot

df['Education'].value_counts().plot.bar()

bar-plot

{x}

<matplotlib.axes._subplots.AxesSubplot at 0x7feb53e3f9d0>

freq

Education Level	Frequency (approx.)
Graduation	1000
PhD	500
Master	380
2n Cycle	220
Basic	60

values discrete r.v

18 / 19

Probability2.ipynb - Colaboratory | ProbabilityDisb_1.ipynb - Colaboratory | New Tab | pandas.DataFrame.boxplot - Colaboratory | +

colab.research.google.com/drive/1l5T7TVIAASw9TdI4JxqJxuDqRgFBGOW3#scrollTo=0j3hZrmRp-34

+ Code + Text RAM Disk

RAM Disk

Code Cell

Run

df['Country'].value_counts().plot.bar()

{x}

fees

<matplotlib.axes._subplots.AxesSubplot at 0x7feb53db4f90>

Country	Value
BR	~1000
AUS	~150
CA	~250
IND	~150
GER	~100
US	~100
ME	~10

2000+

country

[] df['Marital_Status'].value_counts()

File Edit View Insert Cell Kernel Help

19 / 20

Probability2.ipynb - Colaboratory | ProbabilityDisb_1.ipynb - Colaboratory | New Tab

colab.research.google.com/drive/1l5T7TVIAASw9TdI4JxqJxuDqRgFBGOW3#scrollTo=0j3hZrmRp-34

Update

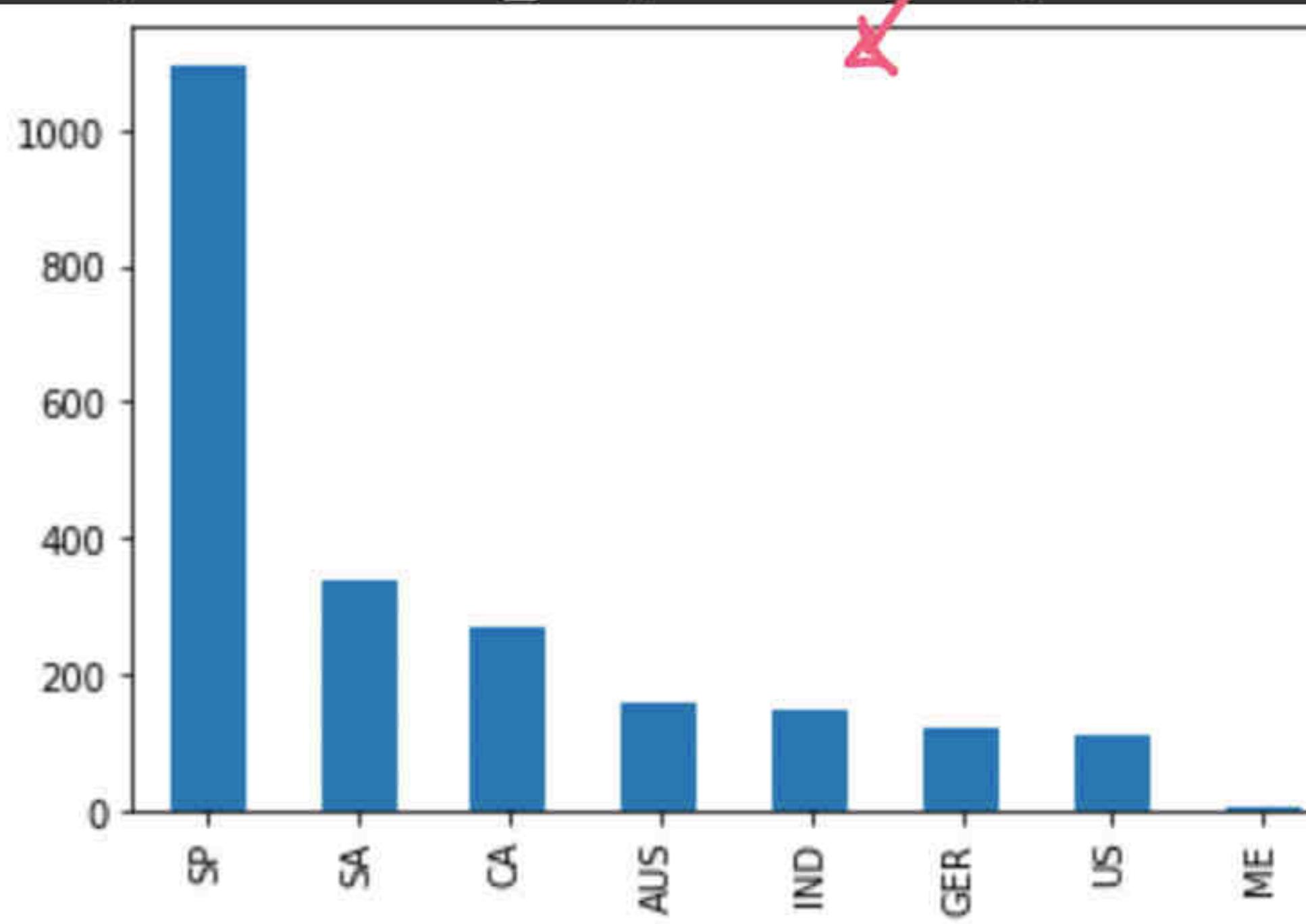
+ Code + Text

RAM
Disk



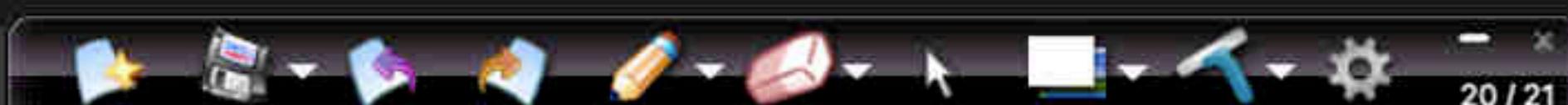
```
df['Country'].value_counts().plot.bar()
```

{x} <matplotlib.axes._subplots.AxesSubplot at 0x7feb53db4f90>



Categorical discrete RV =

```
[ ] df['Marital_Status'].value_counts()
```



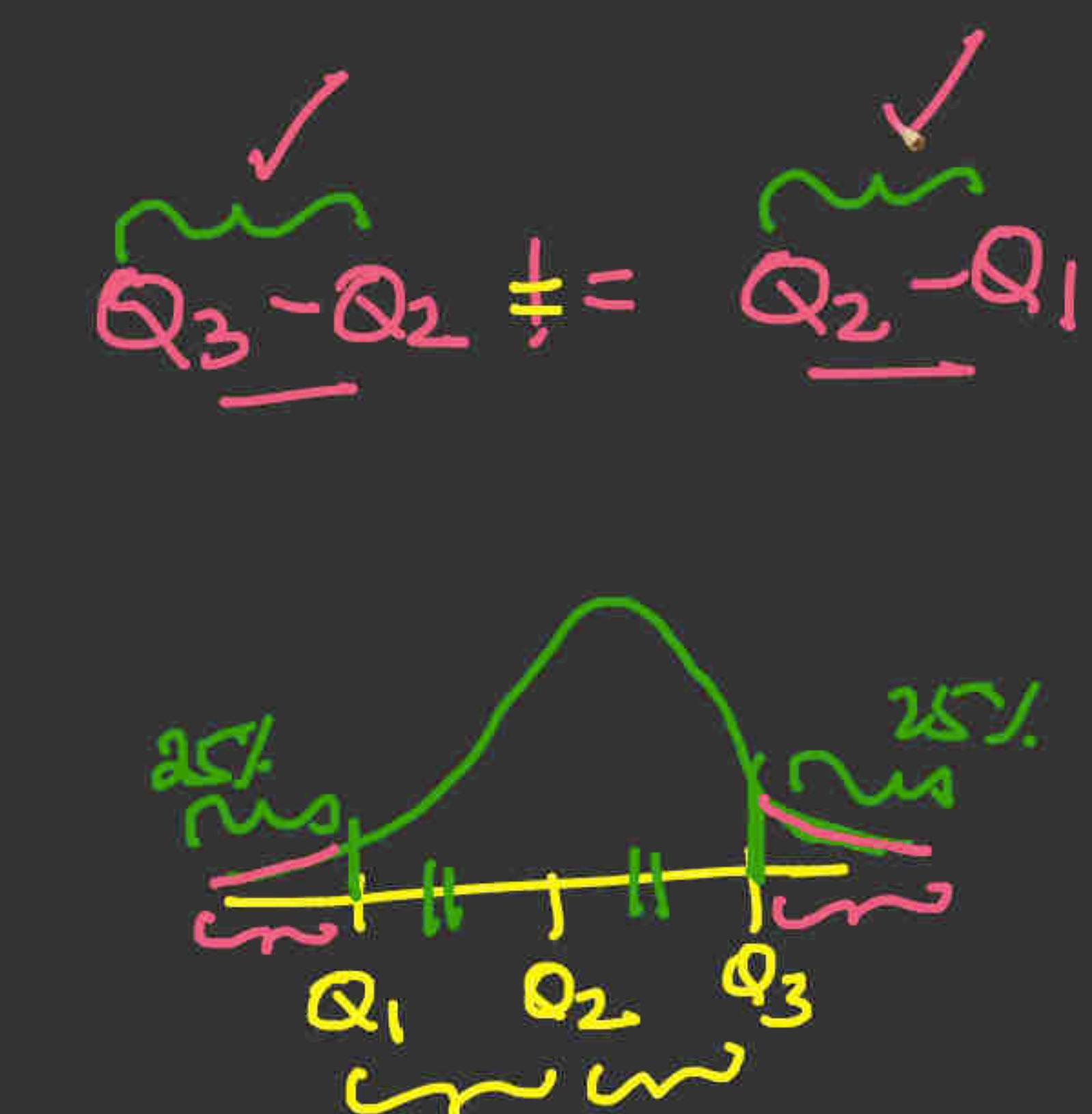
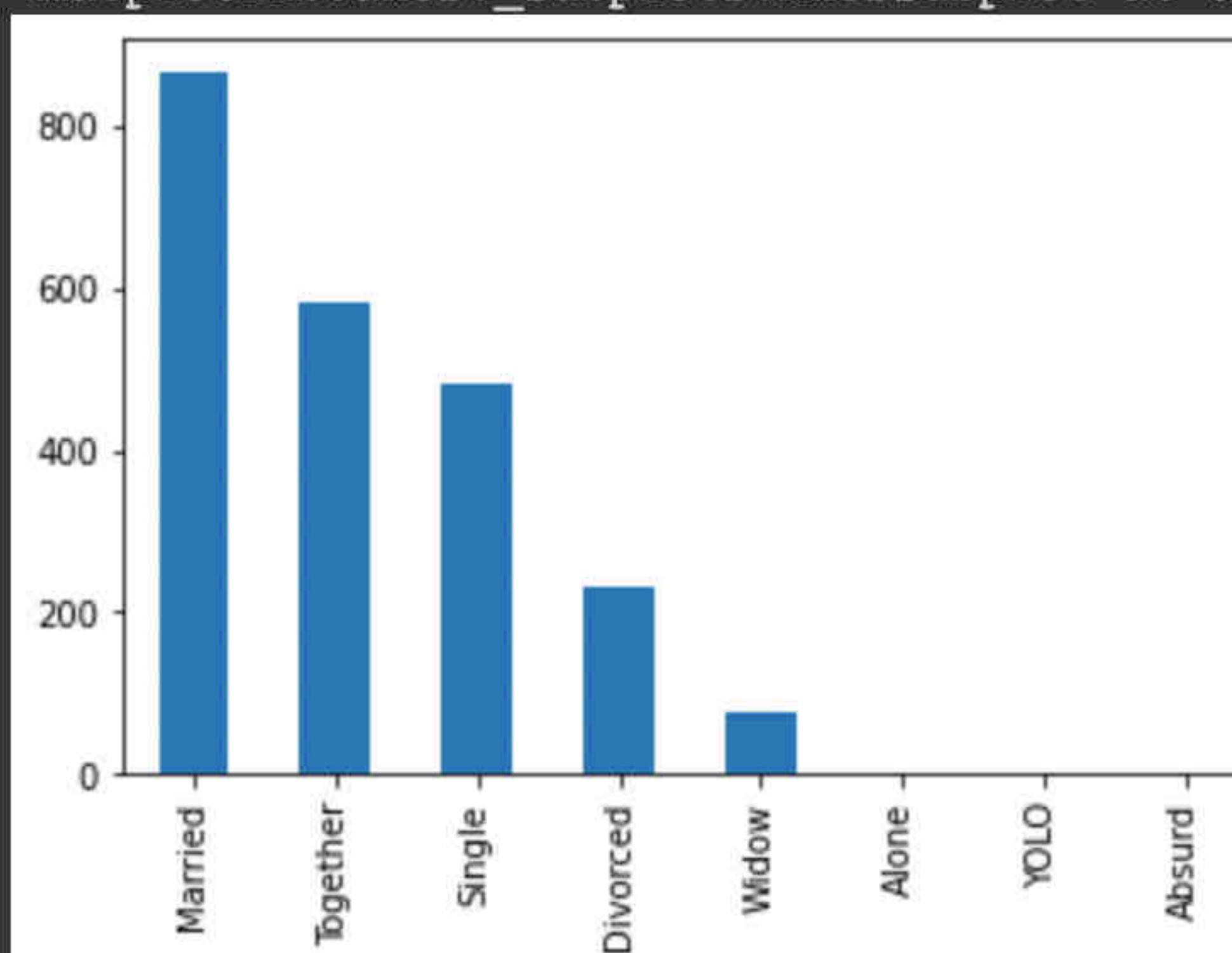
+ Code + Text

RAM
Disk



```
df['Marital_Status'].value_counts().plot.bar()
```

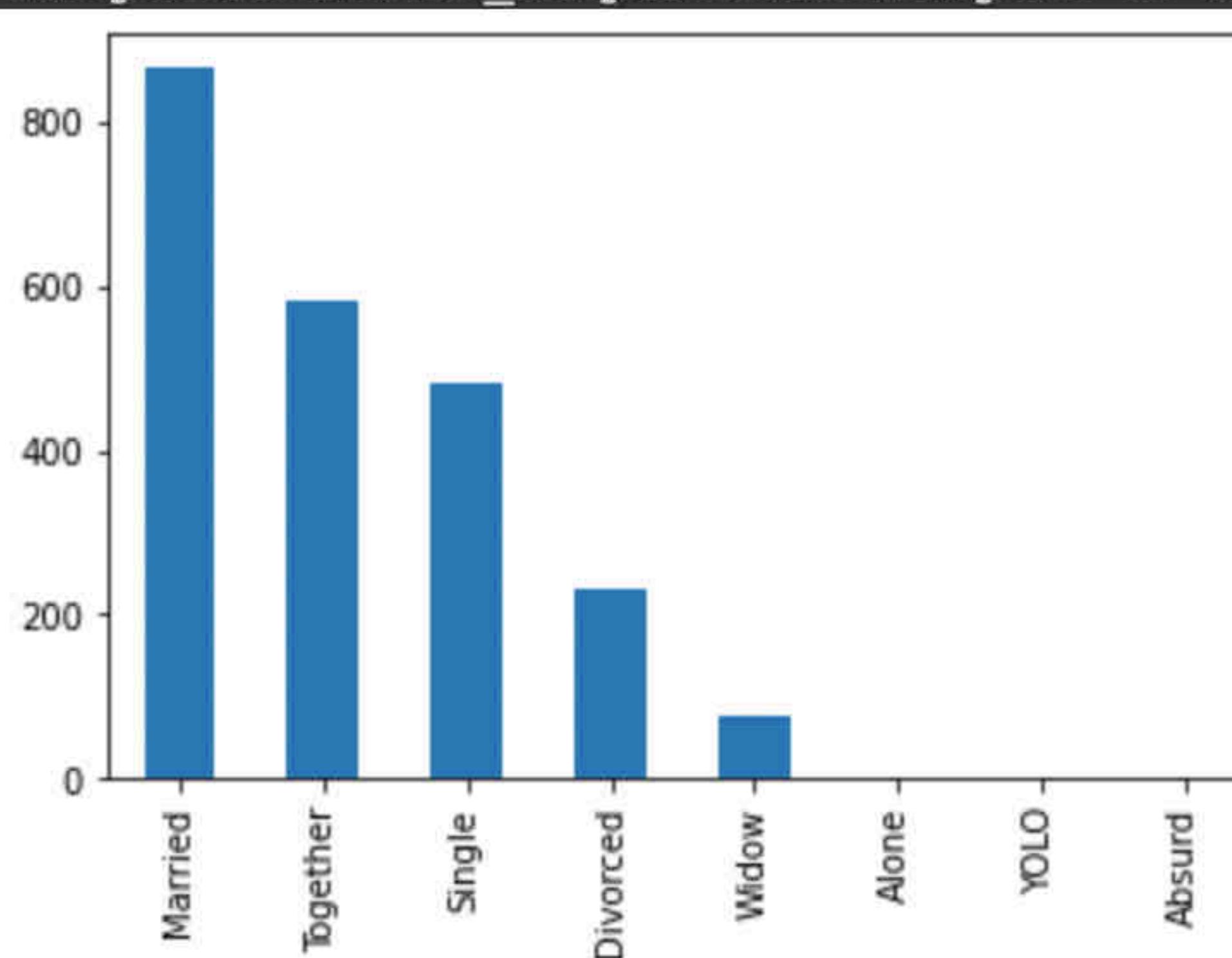
<matplotlib.axes._subplots.AxesSubplot at 0x7ff8e90e8490>





```
df[ 'Marital Status' ].value_counts().plot.bar()
```

```
<matplotlib.axes. subplots.AxesSubplot at 0x7ff8e90e8490>
```

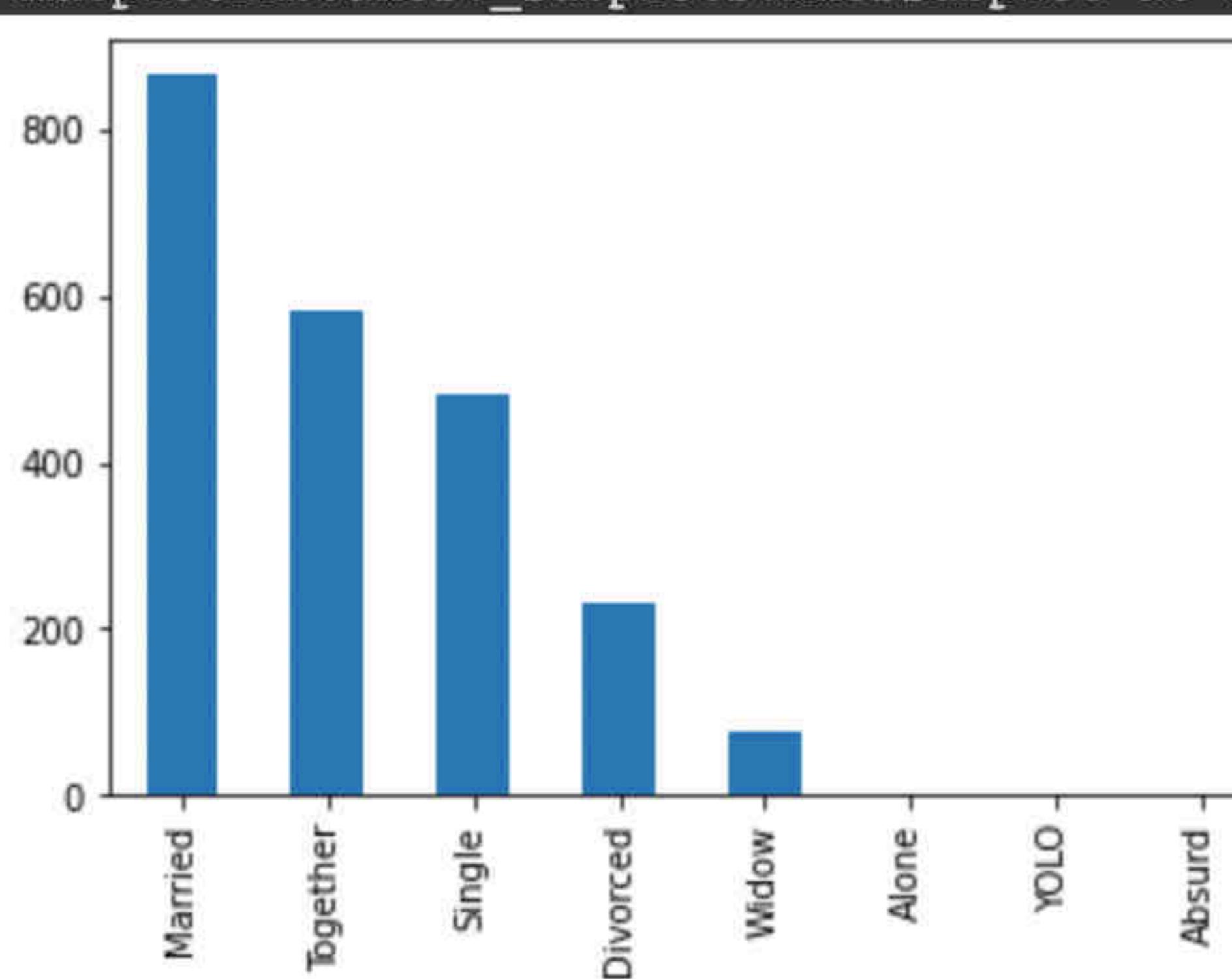


Hand-drawn diagram illustrating the relationships between Python data science libraries:

- Data Science** (Central Box):
 - NumPy
 - Pandas
 - Matplotlib
 - SciPy
- NumPy → Pandas
- SciPy → Pandas
- Pandas → Seaborn
- Pandas → Matplotlib
- Matplotlib (with a bracket and a pencil icon)

+ Code + Text

RAM Disk



```
df[ 'Marital Status' ].value_counts().plot.bar()
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7ff8e90e849>
```

3490>

integer

age

numerical
continuous open

categorical
(ordinal)
(rare)

4, 5, 6, ... 90

{ 10 | 20 | ... }

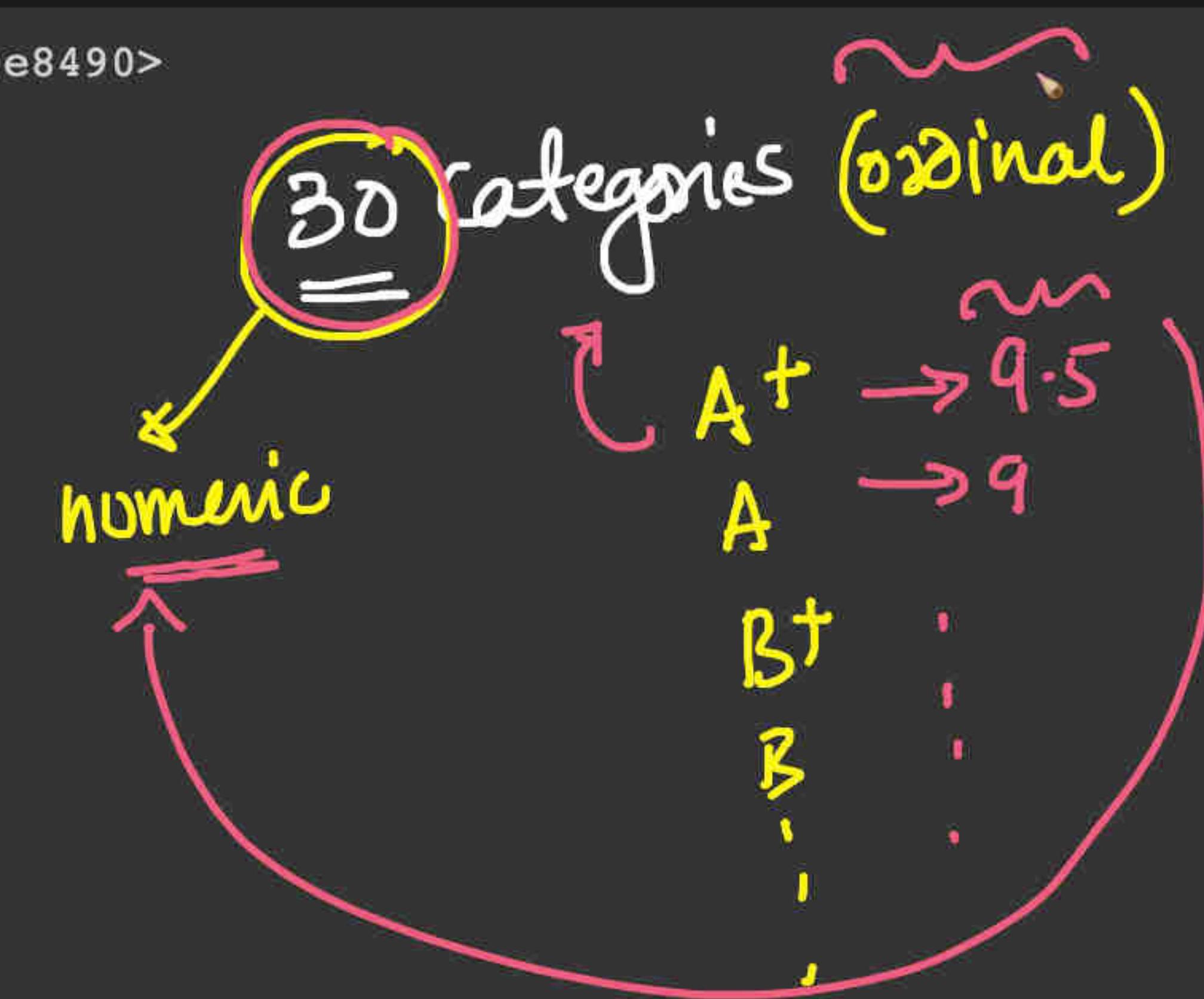
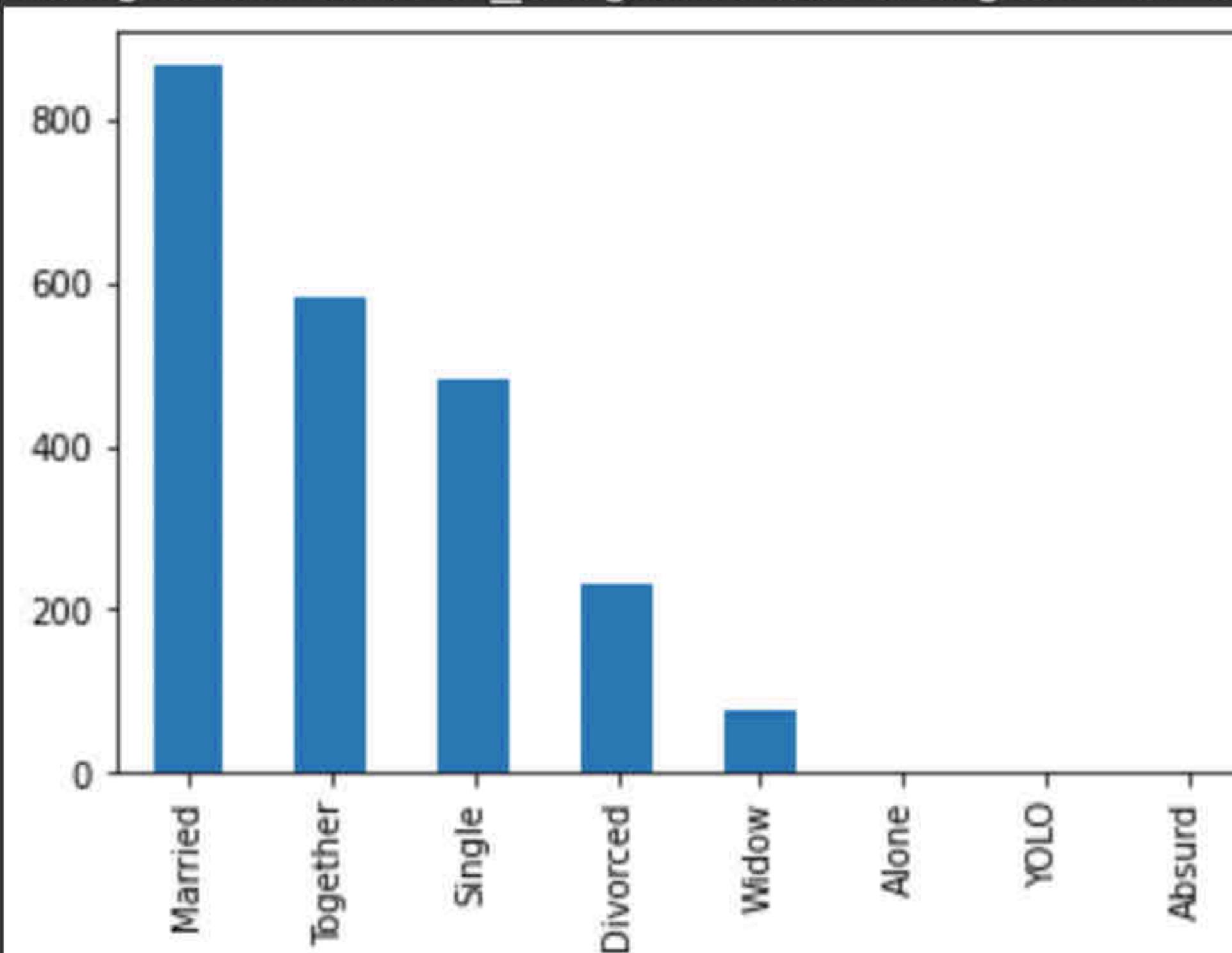
+ Code + Text

RAM
Disk



```
df['Marital_Status'].value_counts().plot.bar()
```

<matplotlib.axes._subplots.AxesSubplot at 0x7ff8e90e8490>



Probability2.ipynb - Colaboratory | ProbabilityDisb_1.ipynb - Colaboratory | New Tab | pandas.DataFrame.boxplot - Colaboratory | +

colab.research.google.com/drive/1l5T7TVIAASw9Tdi4JxqJxuDqRgFBGOW3#scrollTo=A2JIMTNsnIMy

+ Code + Text RAM Disk

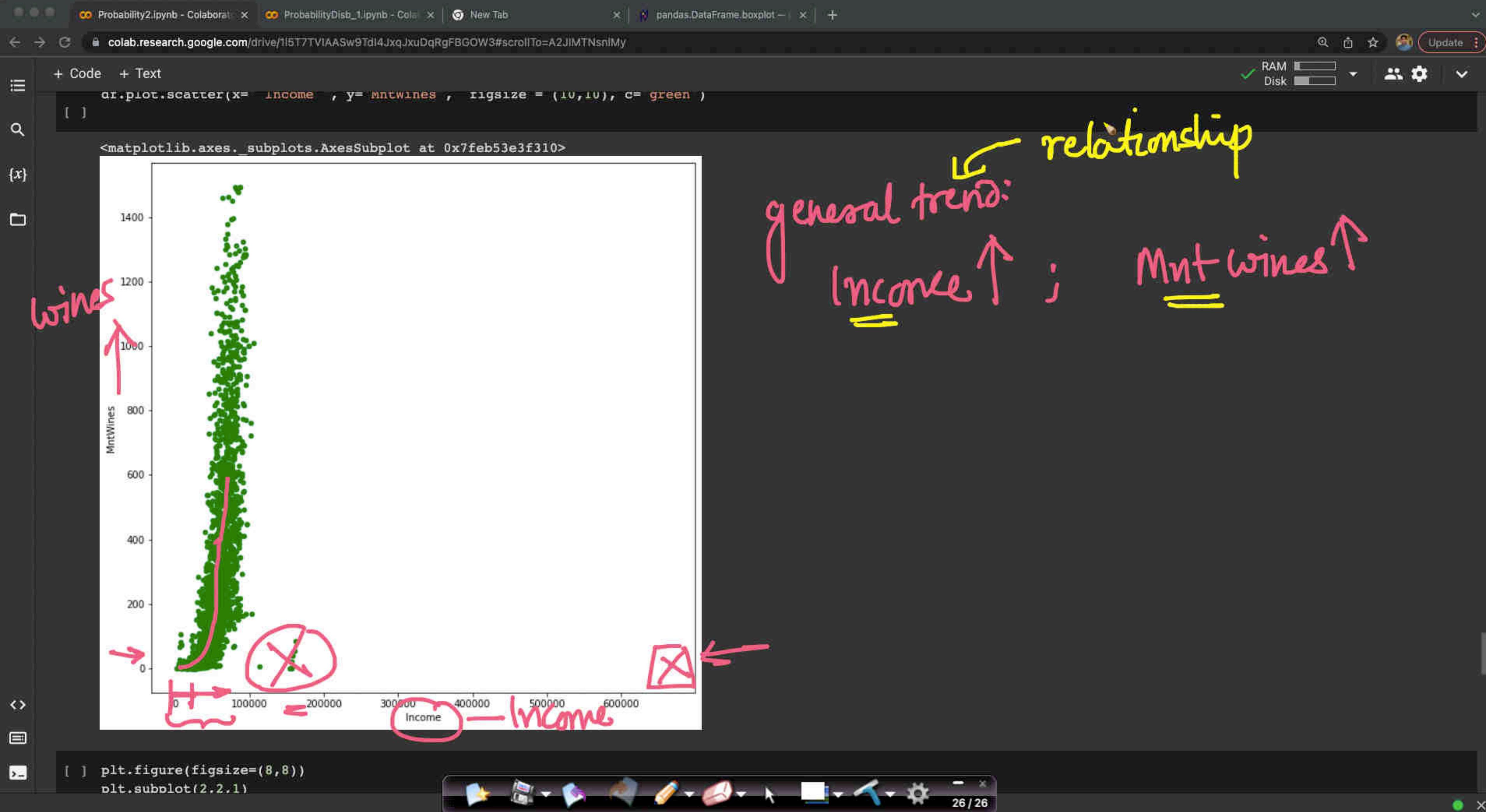
#scatter plot

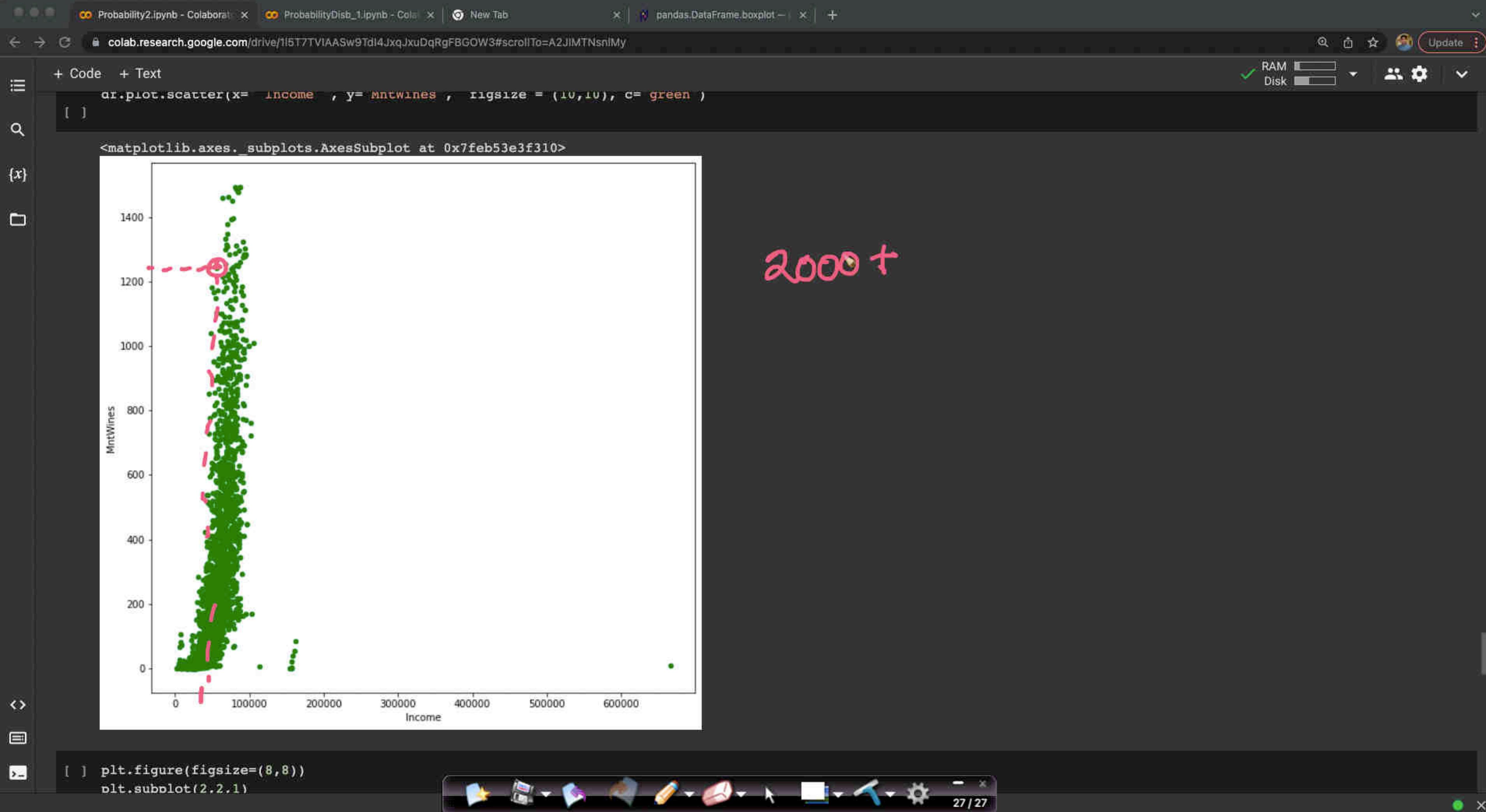
```
df.plot.scatter(x=' Income ', y='MntWines', figsize = (10,10), c='green')
```

{x} <matplotlib.axes._subplots.AxesSubplot at 0x7feb53e3f310>

MntWines

25 / 25





colab.research.google.com/drive/1l5T7TVIAASw9Tdi4JxqJxuDqRgFBGOW3#scrollTo=A2JIMTNsnIMy

+ Code + Text

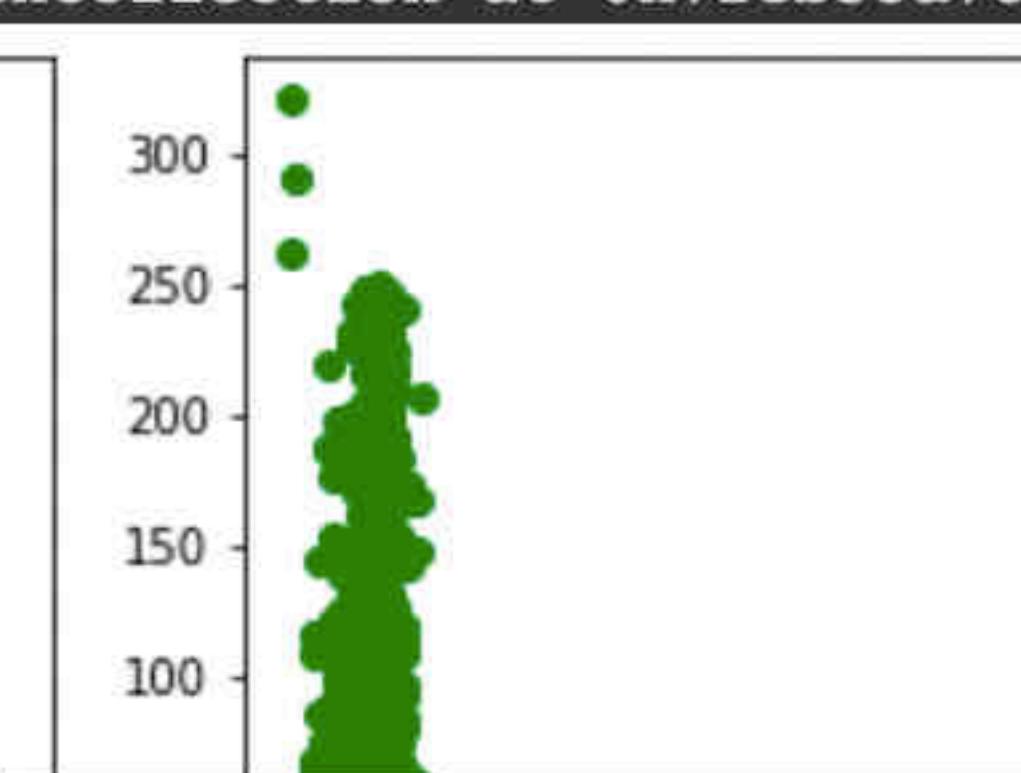
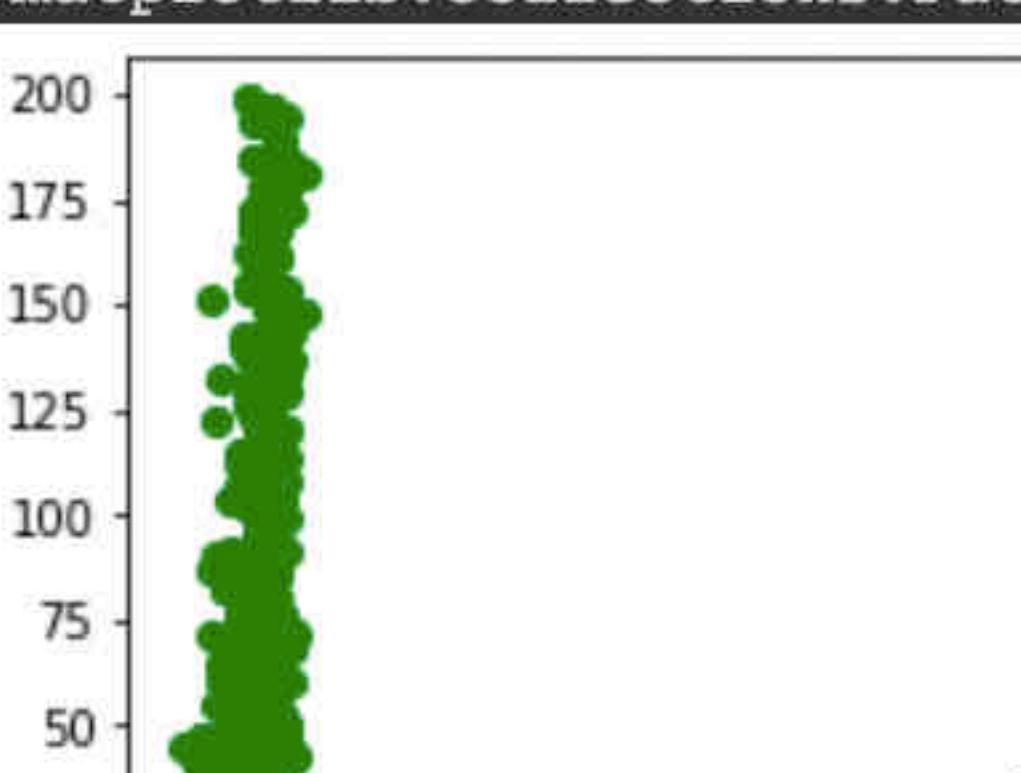
RAM Disk

Income

{x}

plt.figure(figsize=(8,8))
plt.subplot(2,2,1)
plt.scatter(x=df['Income'], y=df['MntFruits'], c='green')
plt.subplot(2,2,2)
plt.scatter(x=df['Income'], y=df['MntGoldProds'], c='green')
plt.subplot(2,2,3)
plt.scatter(x=df['Income'], y=df['MntSweetProducts'], c='green')
plt.subplot(2,2,4)
plt.scatter(x=df['Income'], y=df['MntWines'], c='green')

<matplotlib.collections.PathCollection at 0x7feb53a78d50>



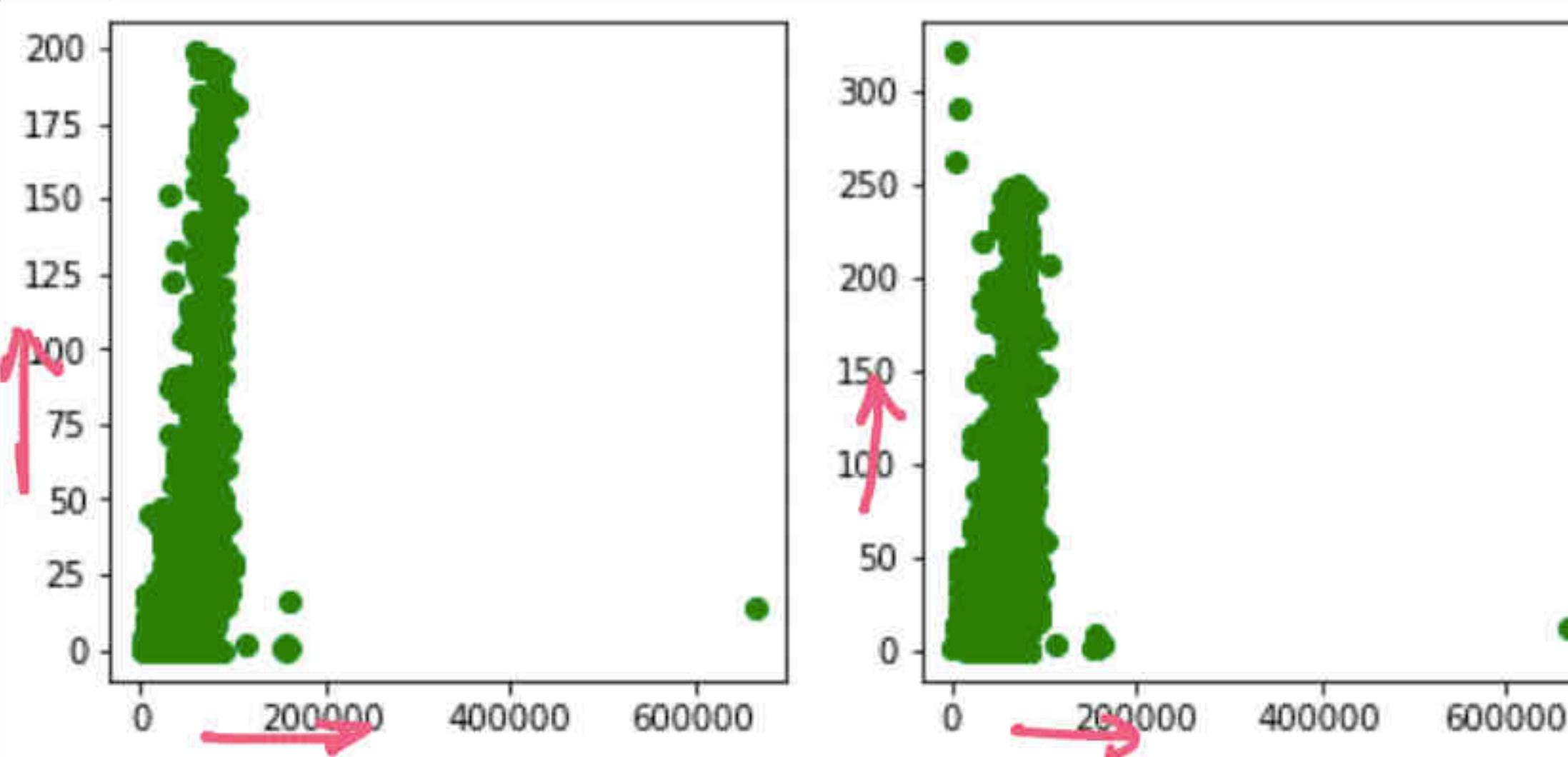
28 / 28

+ Code + Text

RAM Disk

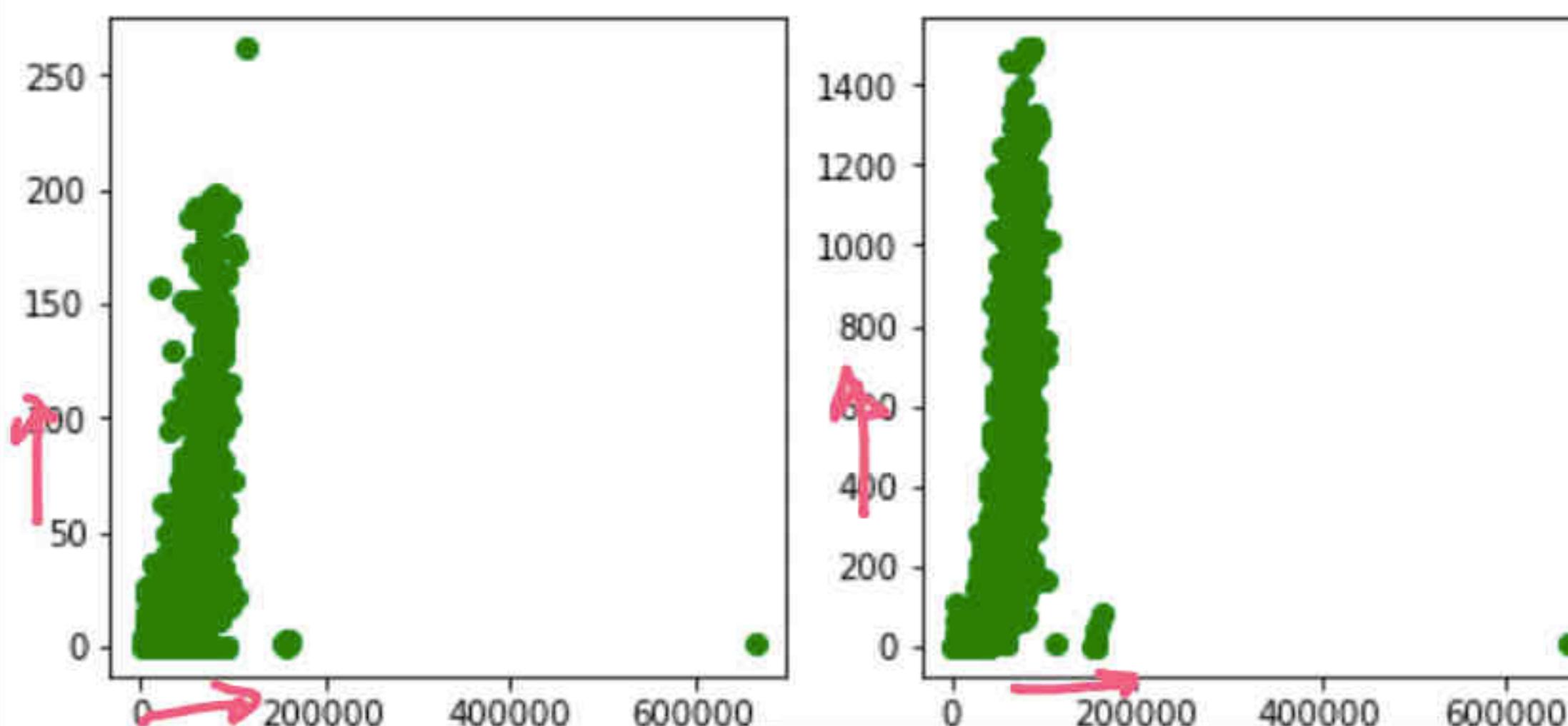


<matplotlib.collections.PathCollection at 0x7feb53a78d50>



In general
Income ↑

Mut- ... ↑



In general
Income ↑

Mut- ... ↑

+ Code + Text

✓ RAM Disk



```
# Have they responded to the marketing survey by our company  
df['Response']
```

```
0    1  
1    1  
2    0  
3    0  
4    1  
..  
2235   0  
2236   0  
2237   0  
2238   0  
2239   1  
Name: Response, Length: 2240, dtype: int64
```

Prime - Sub

```
df.boxplot(by='Response', column='Income', figsize=(10,6))  
plt.show()
```

```
/usr/local/lib/python3.7/dist-packages/matplotlib/cbook/__init__.py:1376: VisibleDeprecationWarning: Creating an nda  
X = np.atleast_1d(X.T if isinstance(X, np.ndarray) else np.asarray(X))
```

Boxplot grouped by Response

Income

Probability2.ipynb - Colaboratory · ProbabilityDisb_1.ipynb - Colaboratory · New Tab

colab.research.google.com/drive/1l5T7TVIAASw9Tdi4JxqJxuDqRgFBGOW3#scrollTo=pUstlAPqkDB

Update

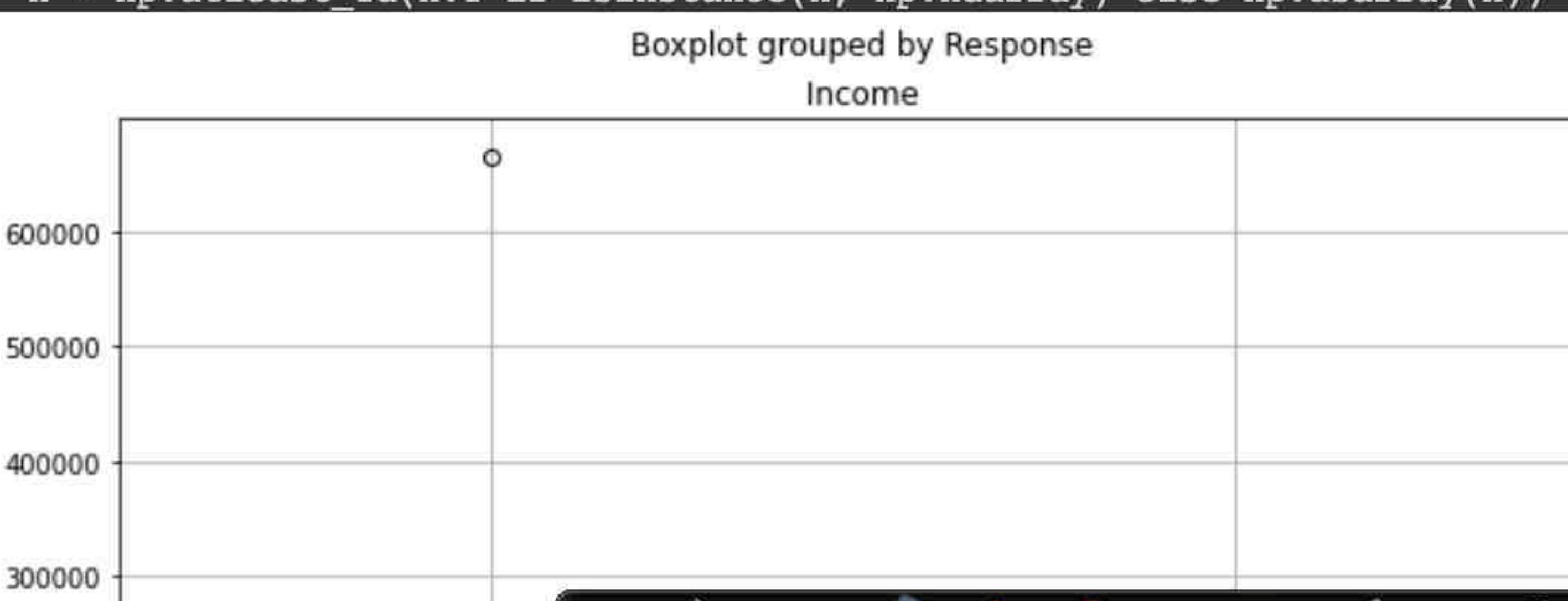
+ Code + Text

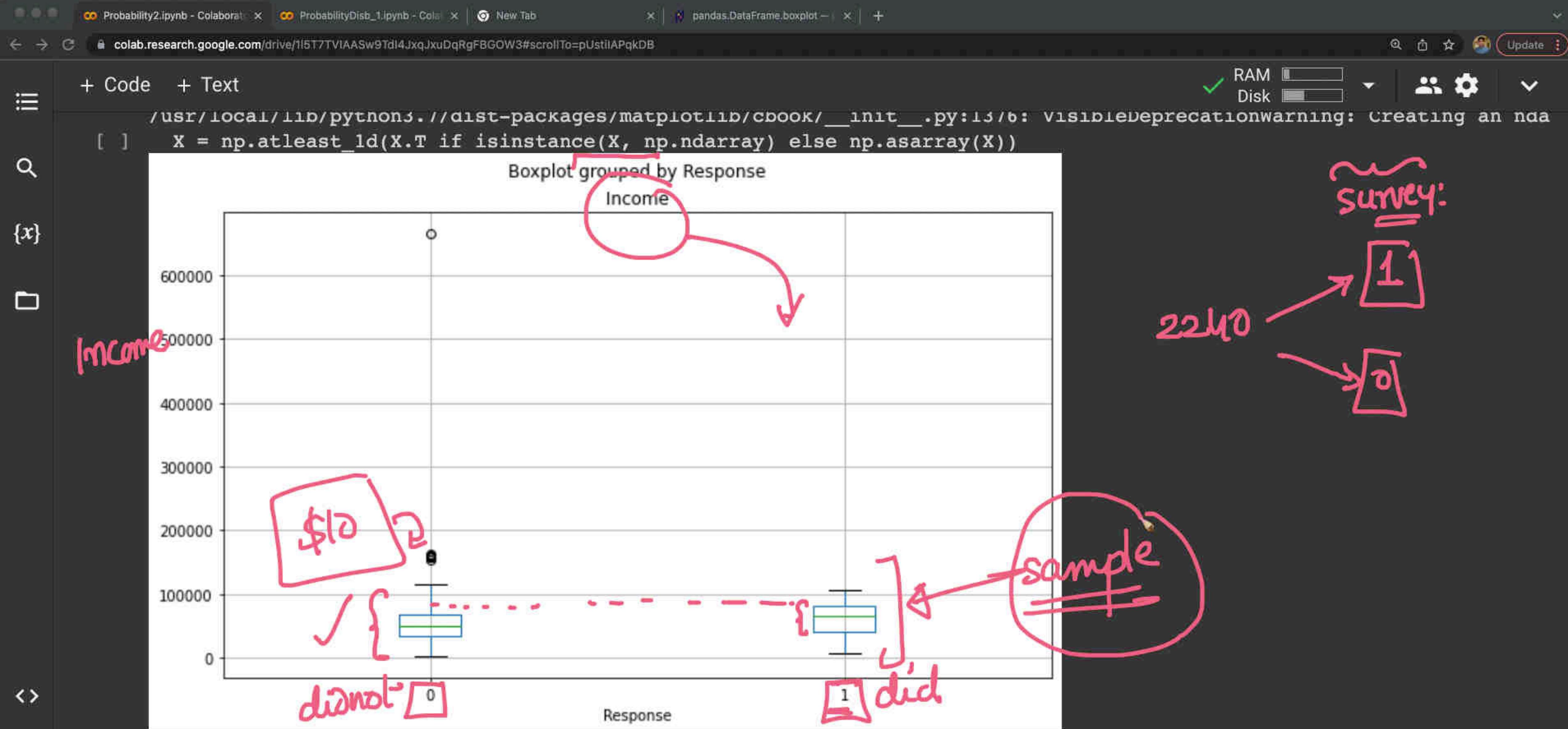
RAM Disk

```
4      1
 ..
2235    0
2236    0
{x} 2237    0
2238    0
2239    1
Name: Response, Length: 2240, dtype: int64
```

```
[ ] df.boxplot(by='Response', column='Income', figsize=(10,6))
plt.show()
```

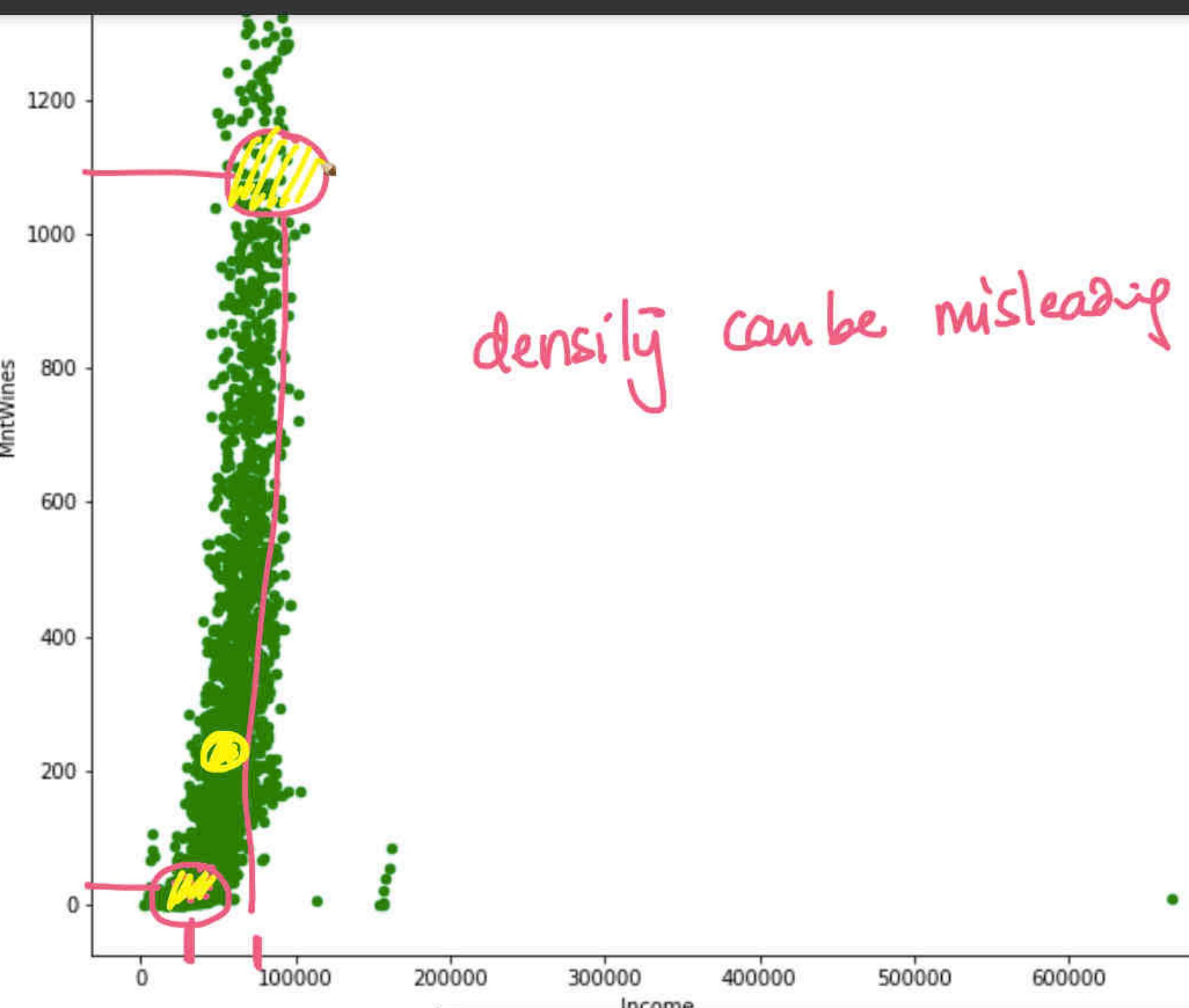
```
/usr/local/lib/python3.7/dist-packages/matplotlib/cbook/__init__.py:1376: VisibleDeprecationWarning: Creating an nda
  X = np.atleast_1d(X.T if isinstance(X, np.ndarray) else np.asarray(X))
```





+ Code + Text

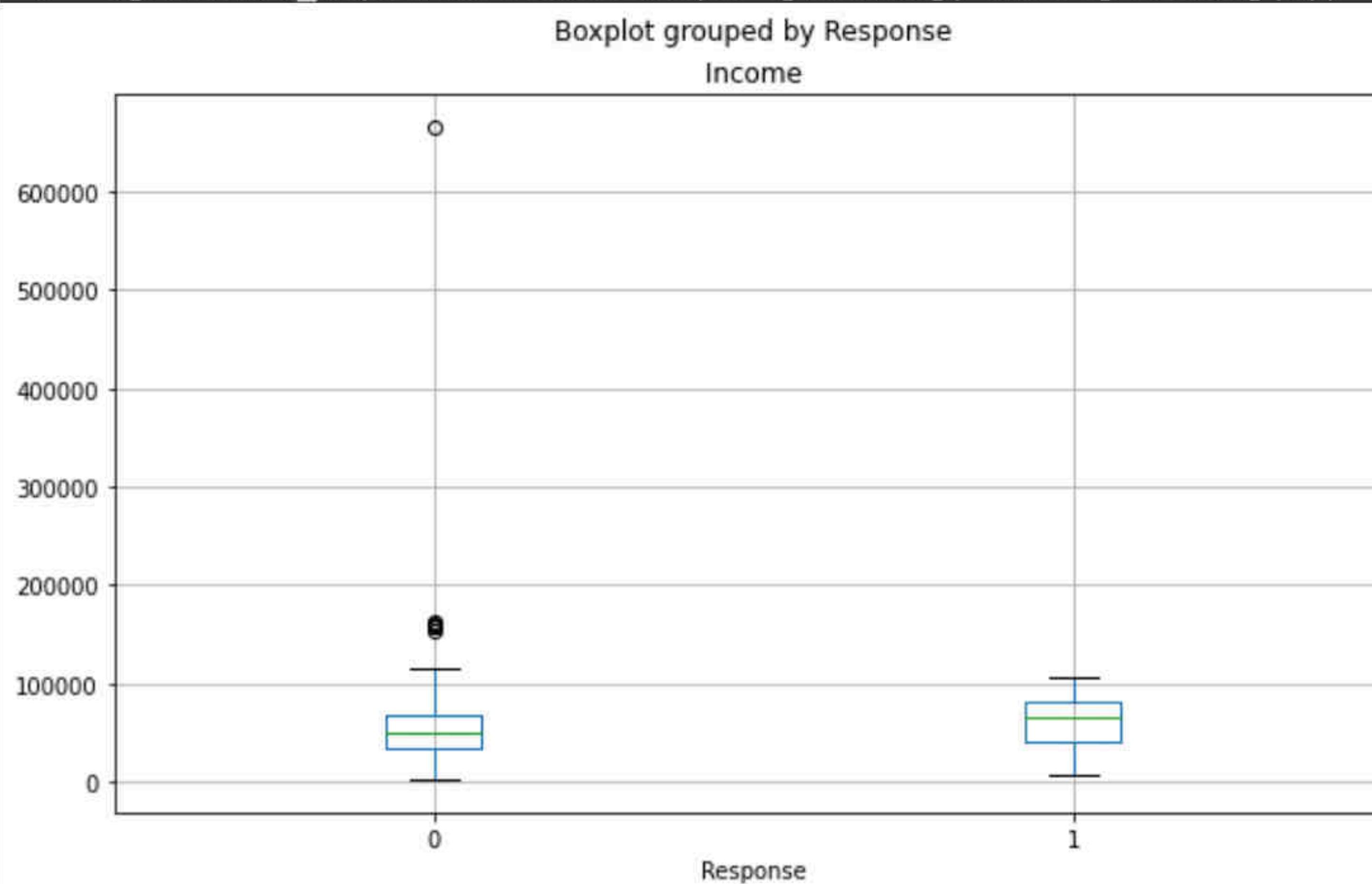
✓ RAM Disk



+ Code + Text

RAM Disk

```
[ ] /usr/local/lib/python3.7/dist-packages/matplotlib/cbook/__init__.py:1376: VisibleDeprecationWarning: Creating an nda  
X = np.atleast_1d(X.T if isinstance(X, np.ndarray) else np.asarray(X))
```



Response

$\rightarrow Y = 1$

$\rightarrow N = 0$

\downarrow

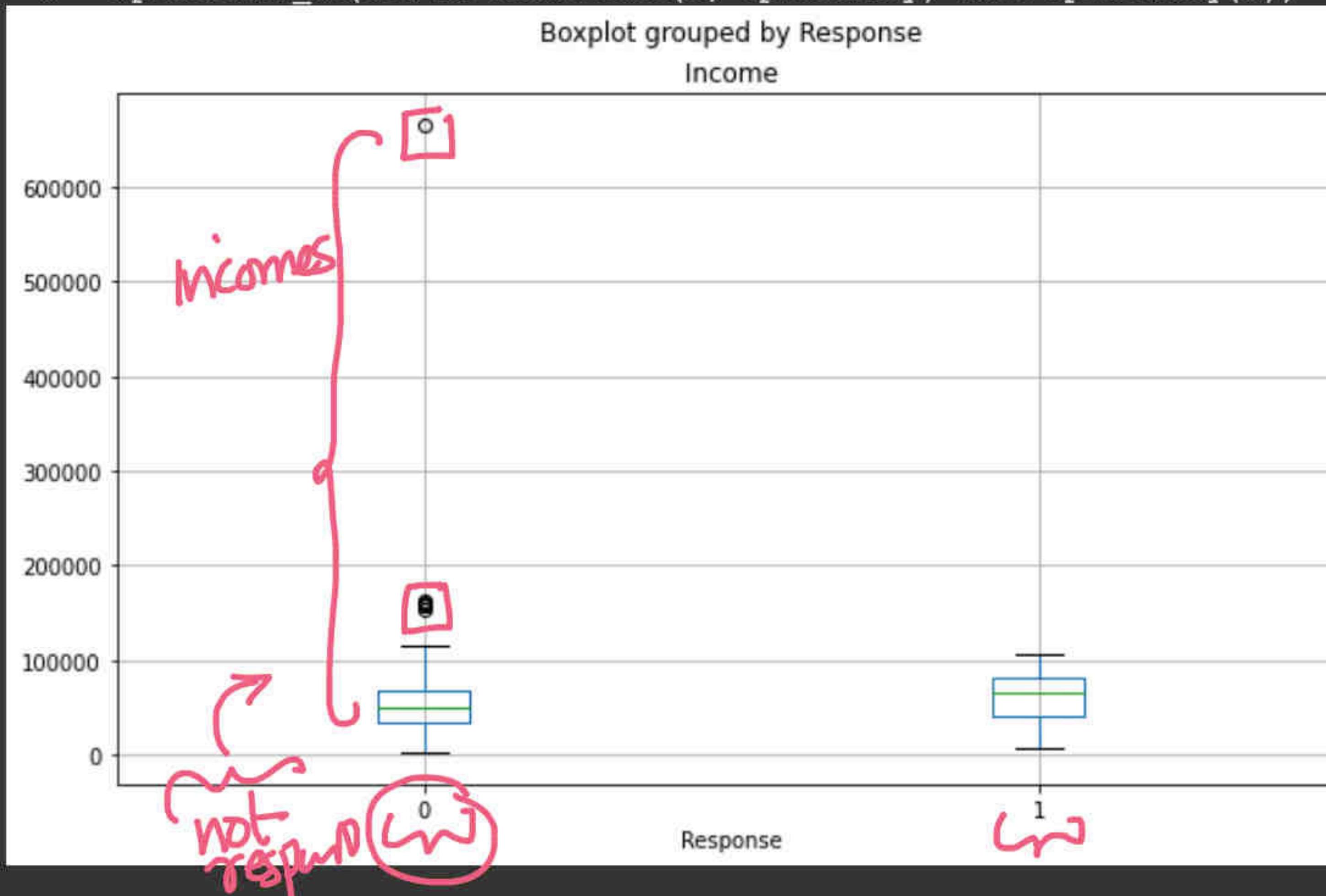
Categorical

+ Code + Text

✓ RAM []
Disk []

V

```
[ ] /usr/local/lib/python3.7/dist-packages/matplotlib/cbook/__init__.py:1376: VisibleDeprecationWarning: Creating an ndarray  
x = np.atleast_1d(x.T if isinstance(x, np.ndarray) else np.asarray(x))
```

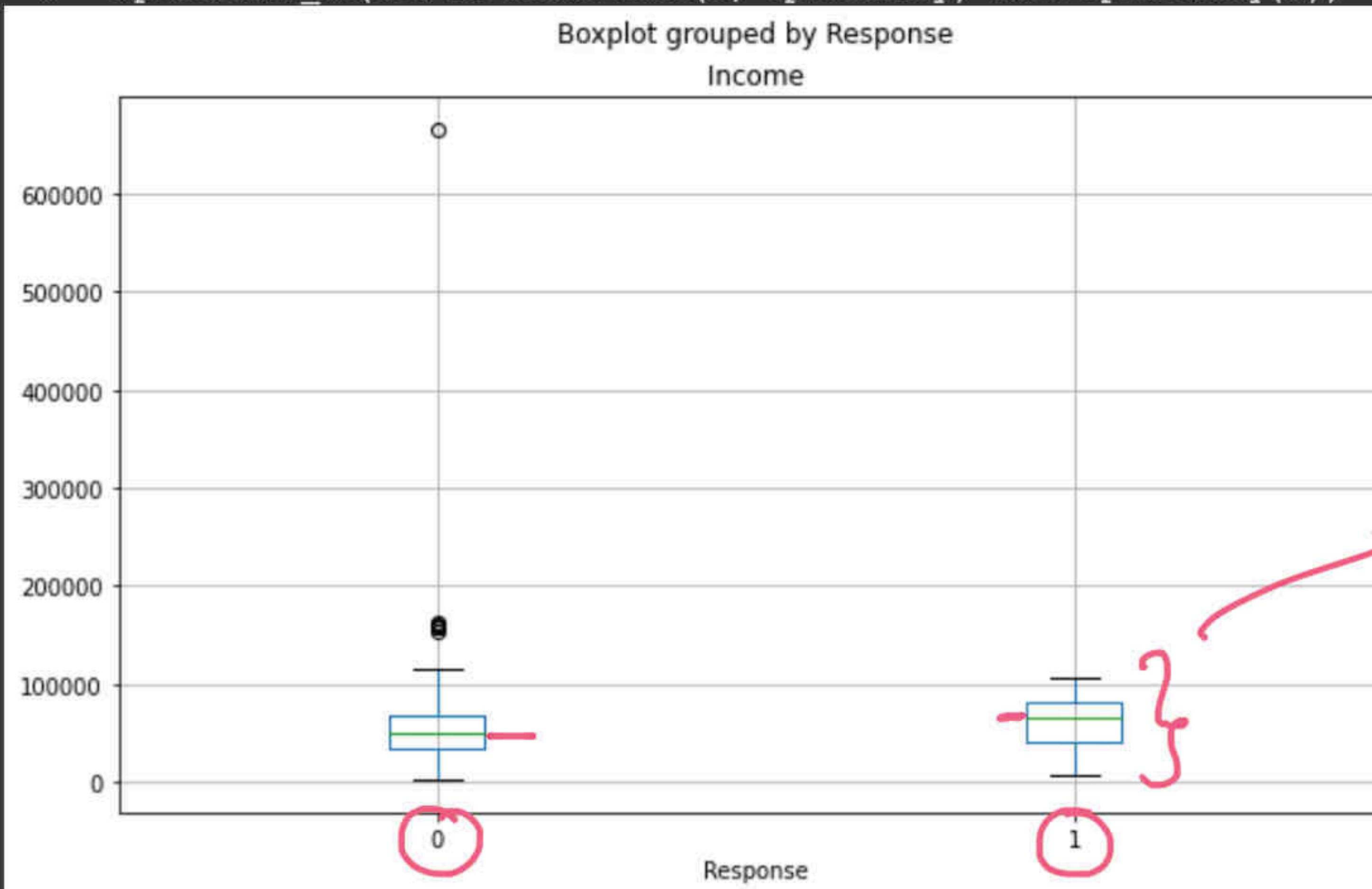


R-1

+ Code + Text

RAM Disk

```
[ ] /usr/local/lib/python3.7/dist-packages/matplotlib/cbook/__init__.py:1376: VisibleDeprecationWarning: Creating an nda  
X = np.atleast_1d(X.T if isinstance(X, np.ndarray) else np.asarray(X))
```

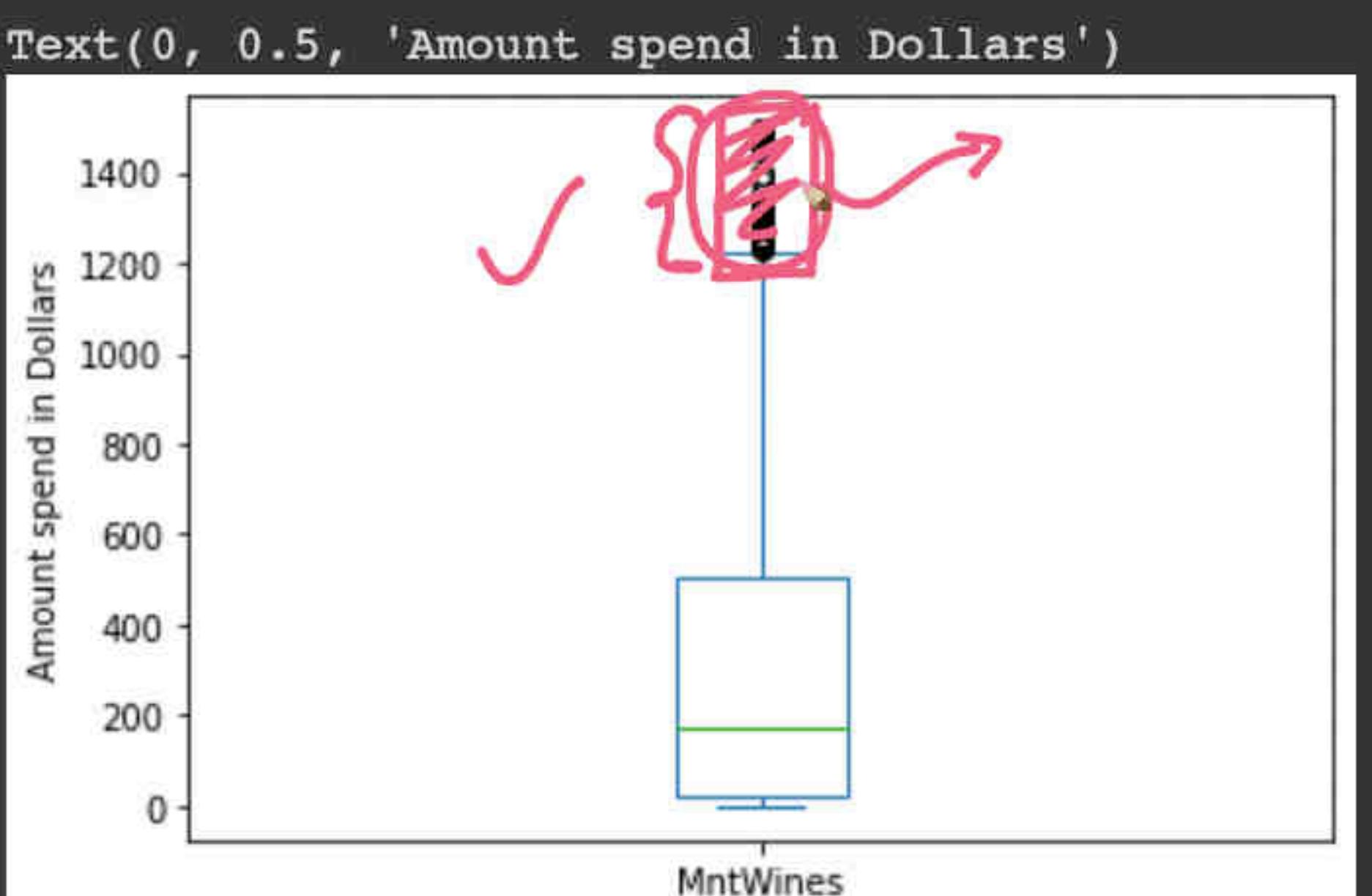


RAM Disk

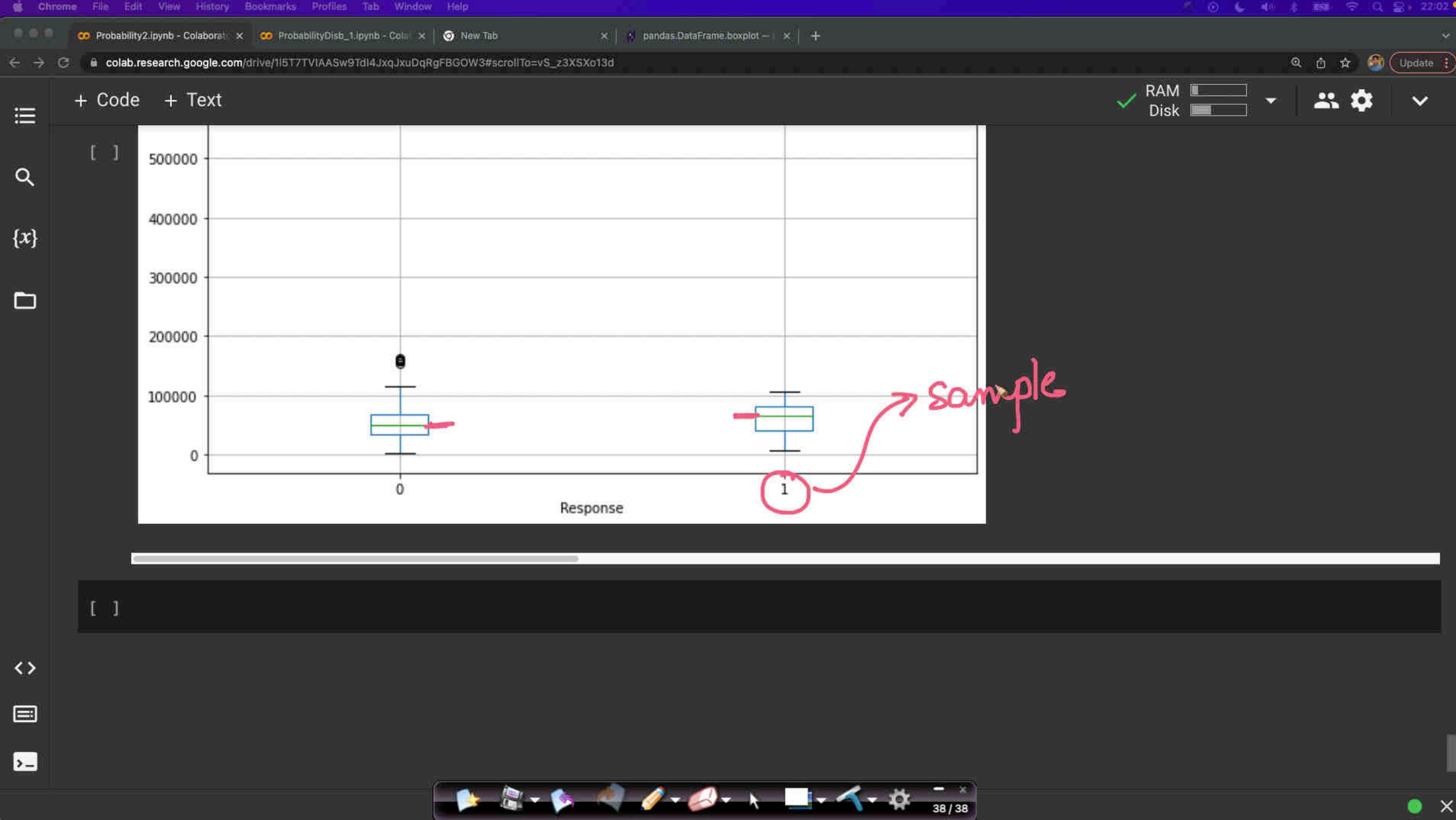
Code



```
[ ] #box plot  
ax = df['MntWines'].plot.box()  
ax.set_ylabel("Amount spend in Dollars")
```

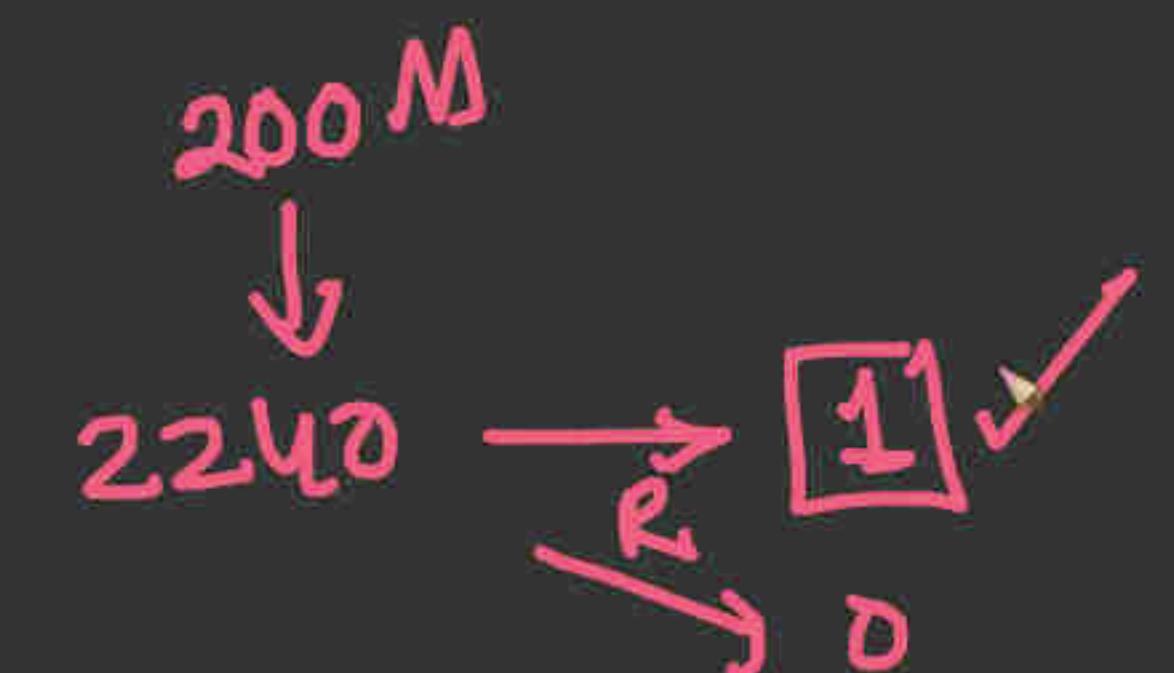
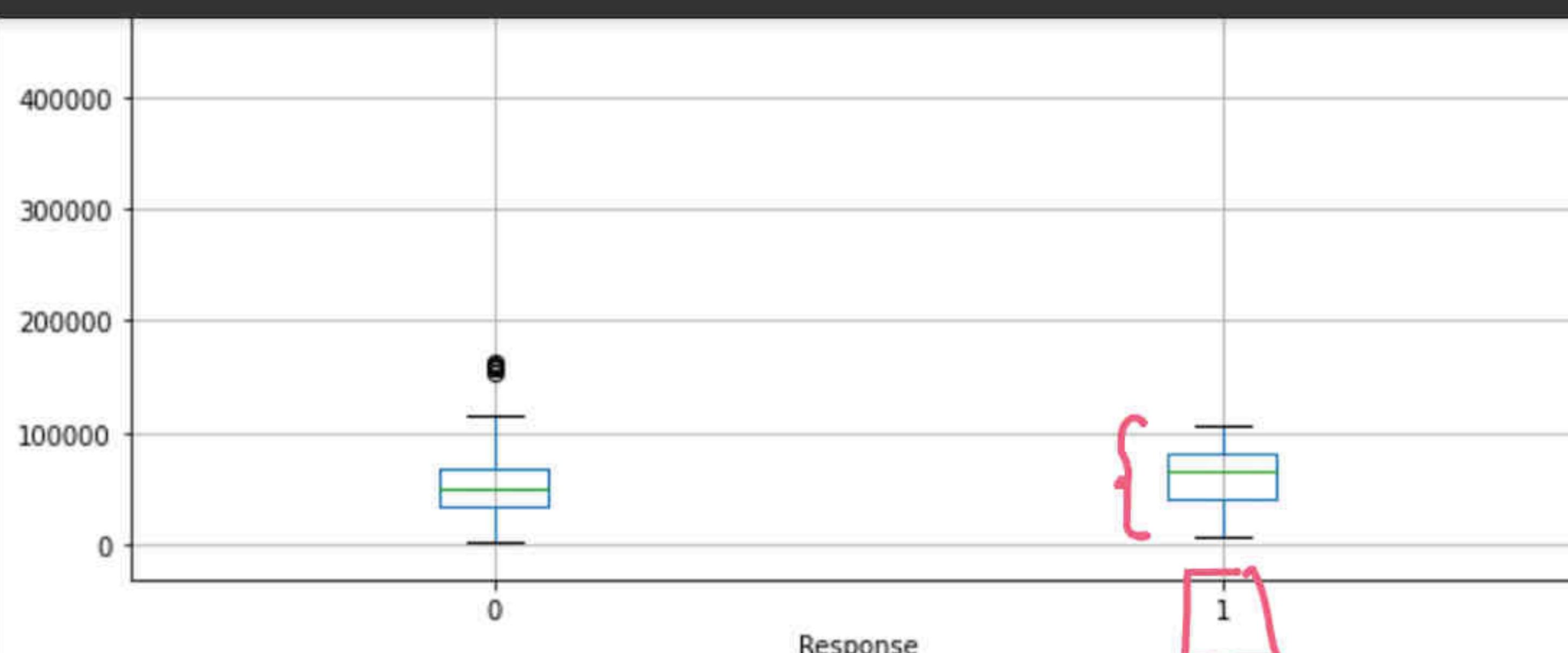


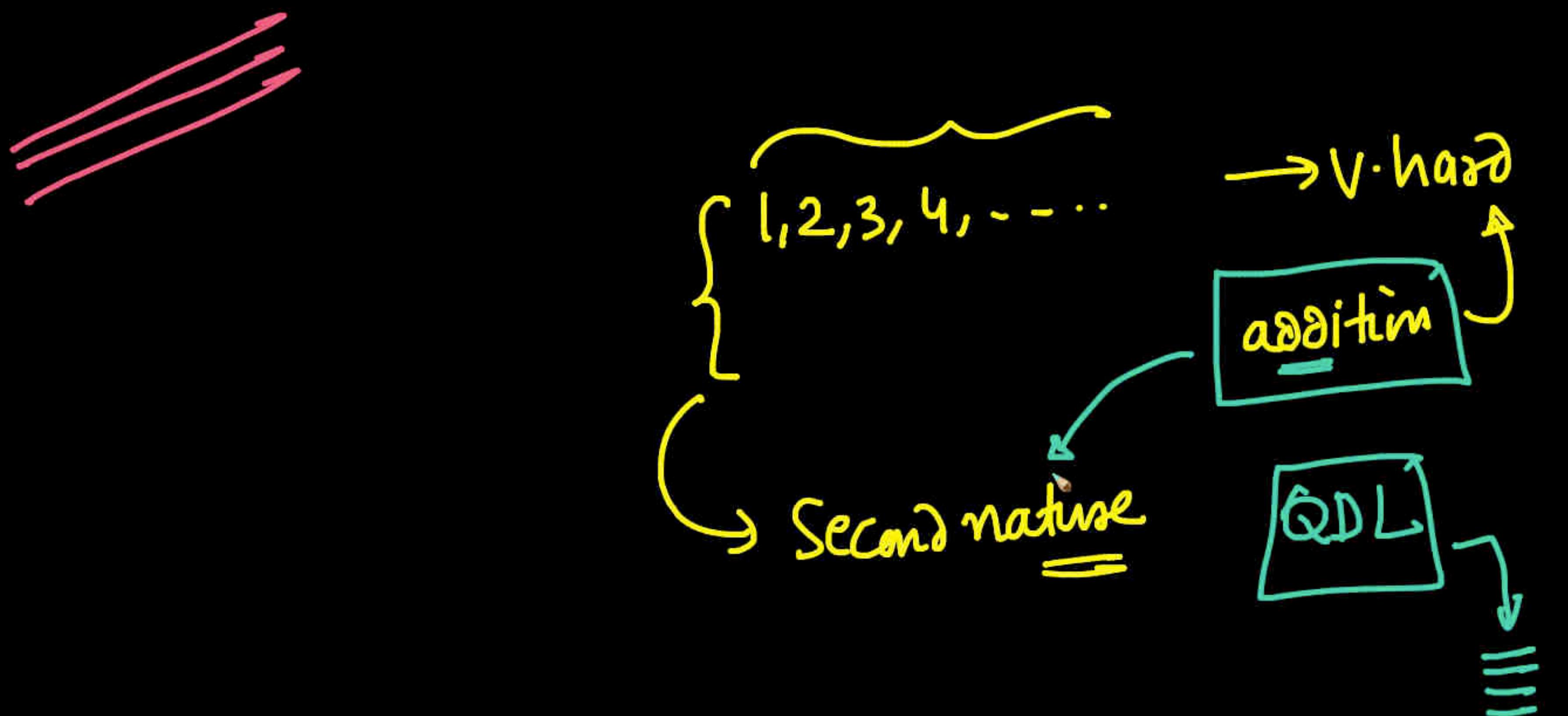
```
[ ] #bar plot  
df['Education'].value_counts().plot.bar()
```



+ Code + Text

✓ RAM Disk





which plot to use
where

→ Continuous r.v: work-exp

bin-width
histogram
KDE; density plots

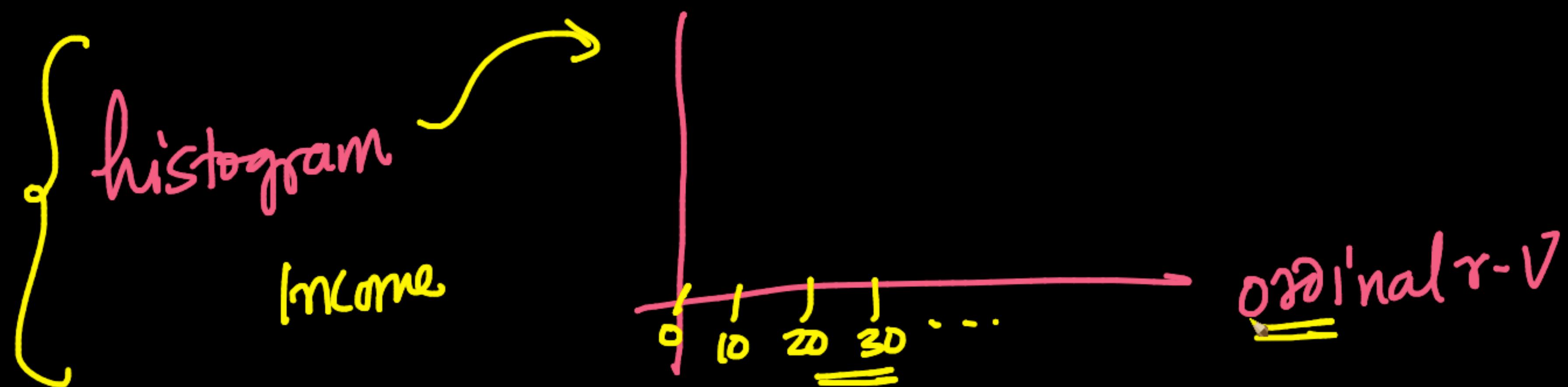
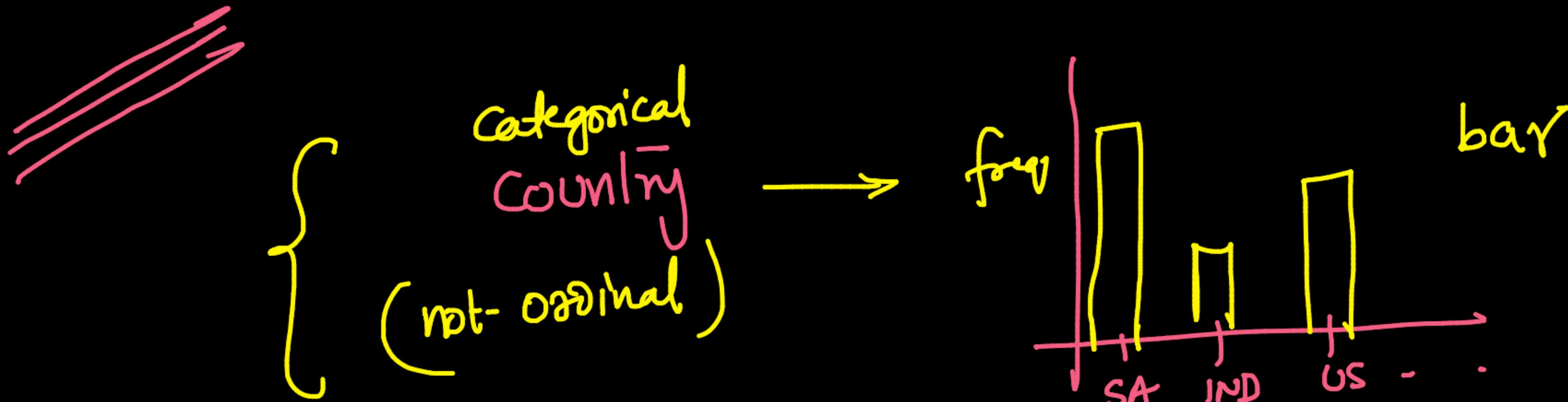
→ bar-plot:

0-2 2-4 4-6

6-8 8+

]

binned

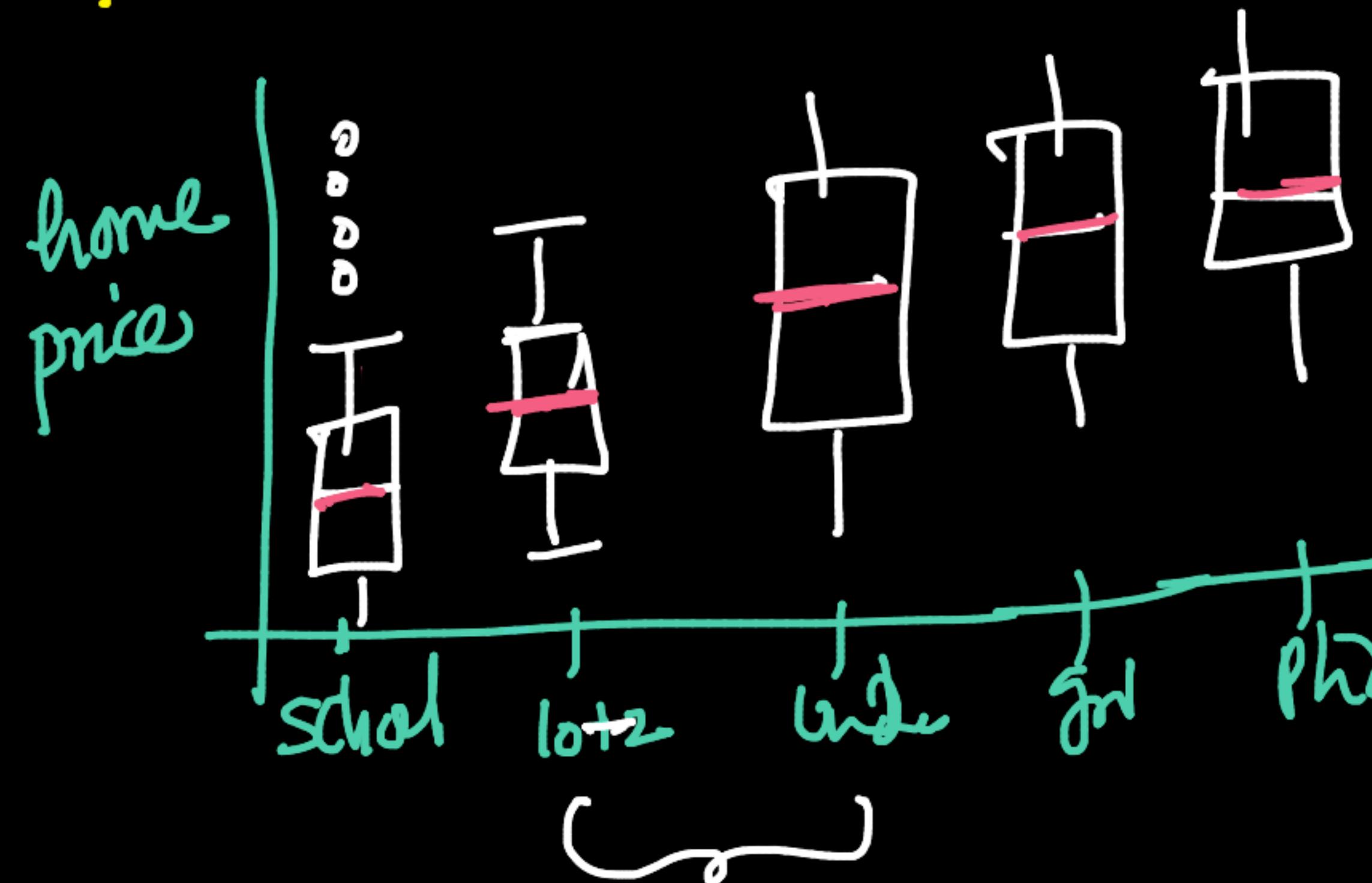


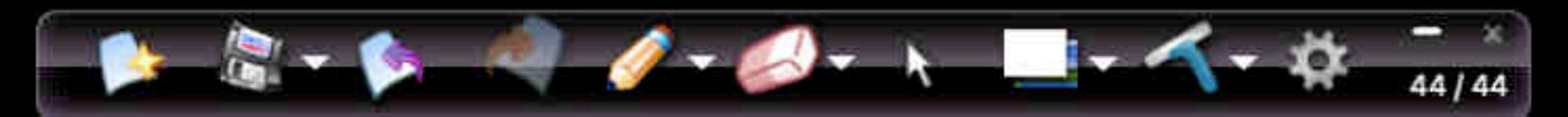
→ home-prices : outliers → boxplot

→ ^{original} home-price & education-level relationship

scatter?

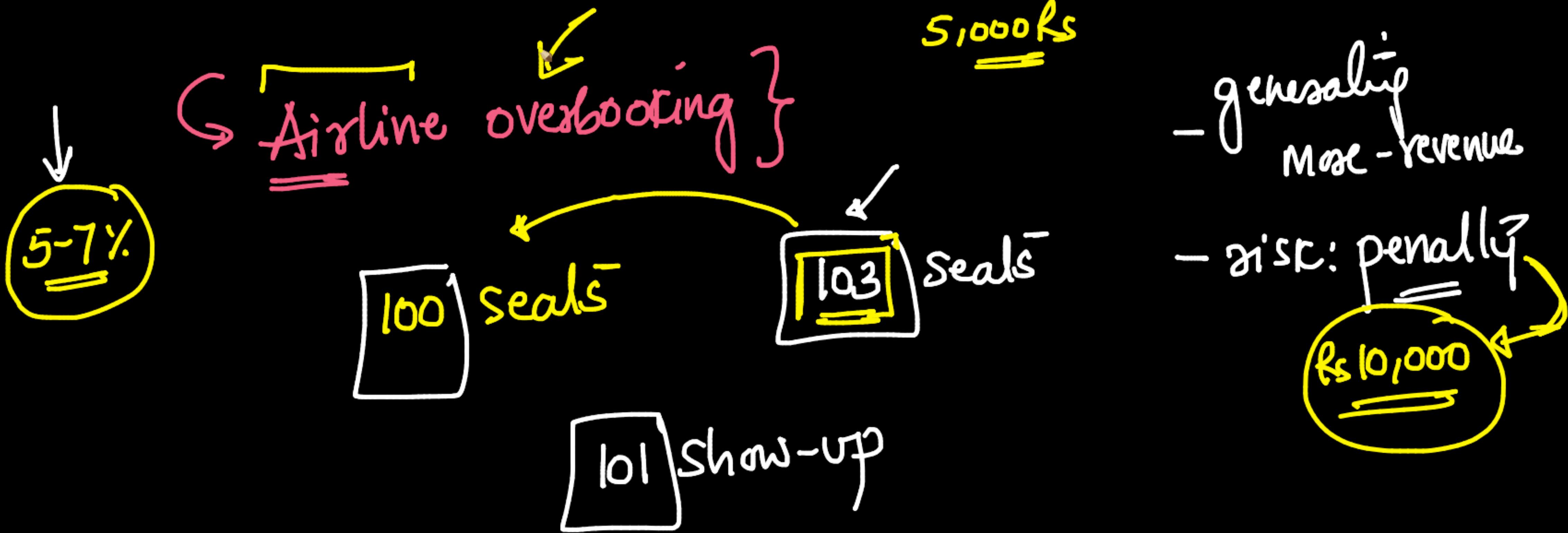
bar plot





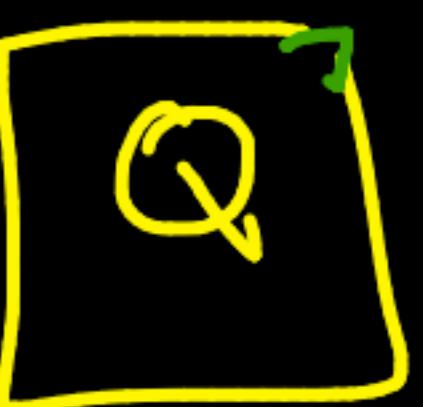
Probability - dist $\frac{2}{2}$

- what ; where
- why
- PMF, PDF ; CDF ..
- real-world ex
- intuition



- generating more revenue

- risk: penalty



Airline

-data [showing-up]

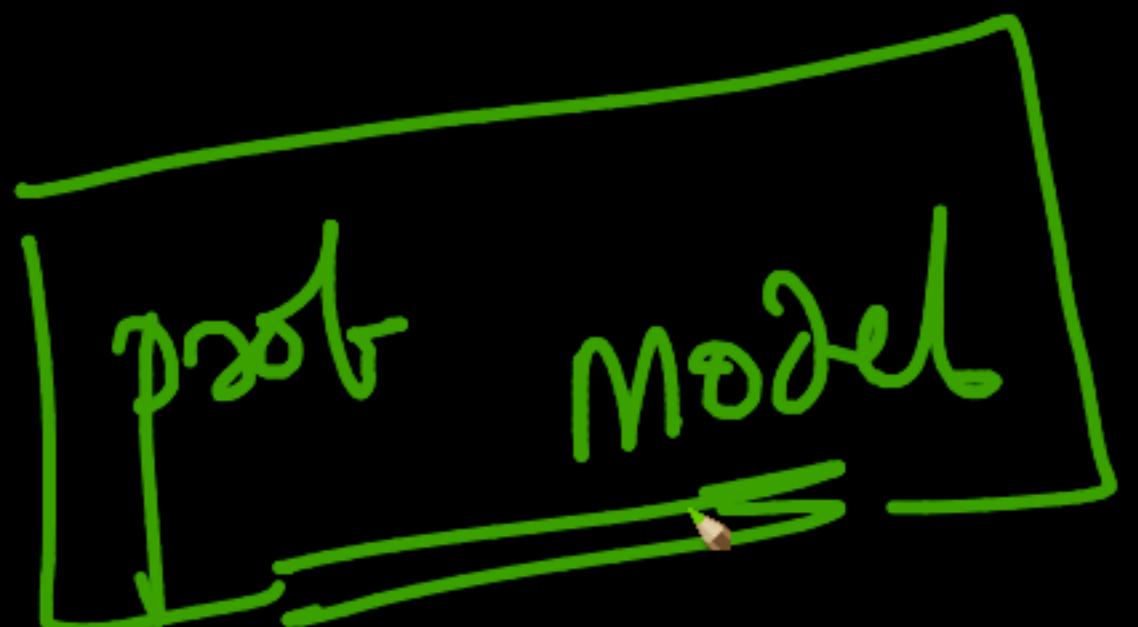
100-seats

- Max net rev by over booking

↳ revenue - penalty

more tickets

↙ 100 show-up



Probability2.ipynb - Colaboratory | ProbabilityDistr_1.ipynb - Colaboratory | New Tab | pandas.DataFrame.boxplot - Colaboratory | +

colab.research.google.com/drive/1JZ-TbmRIY8mmD018C_grijQHluBGNQxZ#scrollTo=SxnxYGBUXIQP

+ Code + Text

RAM Disk

[3] flight.csv 100%[=====] 11.89K --.KB/s in 0s

2022-04-26 16:58:15 (57.7 MB/s) - 'flight.csv' saved [12175/12175]

{x}

flights = pd.read_csv('./flight.csv')
flights.info()

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 1000 entries, 0 to 999  
Data columns (total 3 columns):  
 #   Column      Non-Null Count  Dtype    
---    
 0   Passenger_ID  1000 non-null   int64   
 1   Flight_ID     1000 non-null   object   
 2   Arrived       1000 non-null   int64   
dtypes: int64(2), object(1)  
memory usage: 23.6+ KB
```

✓ flights.head()

	Passenger_ID	Flight_ID	Arrived
0	1811	A320	1

48 / 48

Probability2.ipynb - Colaboratory | ProbabilityDistr_1.ipynb - Colaboratory | New Tab | pandas.DataFrame.boxplot - Colaboratory

colab.research.google.com/drive/1JZ-TbmRIY8mmD018C_grijQHluBGNQxZ#scrollTo=46VlhWzHWhRZ

+ Code + Text RAM Disk  

[1] `import pandas as pd
import numpy as np
import matplotlib.pyplot as plt

import seaborn as sns
from scipy import stats`

[2] `id = "1PazlhissU63pozk0jJckyjIJuM_u-0JF"
print("https://drive.google.com/uc?export=download&id=" + id)`

https://drive.google.com/uc?export=download&id=1PazlhissU63pozk0jJckyjIJuM_u-0JF

[3] `!wget "https://drive.google.com/uc?export=download&id=1PazlhissU63pozk0jJckyjIJuM_u-0JF" -O flight.csv`

--2022-04-26 16:58:14-- https://drive.google.com/uc?export=download&id=1PazlhissU63pozk0jJckyjIJuM_u-0JF
Resolving drive.google.com (drive.google.com)... 74.125.26.101, 74.125.26.138, 74.125.26.100, ...
Connecting to drive.google.com (drive.google.com)|74.125.26.101|:443... connected.
HTTP request sent, awaiting response... 303 See Other
Location: <https://doc-14-14-docs.googleusercontent.com/docs/securesc/ha0ro937gcuc717deffksulhg5h7mbp1/4n3oik3dp7qkgo>
Warning: wildcards not supported in HTTP.
--2022-04-26 16:58:15-- <https://doc-14-14-docs.googleusercontent.com/docs/securesc/ha0ro937gcuc717deffksulhg5h7mbp1>
Resolving doc-14-14-docs.googleusercontent.com (doc-14-14-docs.googleusercontent.com)... 173.194.218.132, 2607:f8b0:
Connecting to doc-14-14-docs.googleusercontent.com (doc-14-14-docs.googleusercontent.com)|173.194.218.132|:443... co
HTTP request sent, awaiting response... 200 OK

49 / 49

Probability2.ipynb - Colaboratory | ProbabilityDistr_1.ipynb - Colaboratory | New Tab | pandas.DataFrame.boxplot - Colaboratory | +

colab.research.google.com/drive/1JZ-TbmRIY8mmD018C_grijQHluBGNQxZ#scrollTo=46VlhWzHWhRZ

+ Code + Text

https://drive.google.com/uc?export=download&id=1PazlhissU63pozk0jJckyjIJuM_u-0JF

[3] !wget "https://drive.google.com/uc?export=download&id=1PazlhissU63pozk0jJckyjIJuM_u-0JF" -O flight.csv

{x} --2022-04-26 16:58:14-- https://drive.google.com/uc?export=download&id=1PazlhissU63pozk0jJckyjIJuM_u-0JF
Resolving drive.google.com (drive.google.com)... 74.125.26.101, 74.125.26.138, 74.125.26.100, ...
Connecting to drive.google.com (drive.google.com)|74.125.26.101|:443... connected.
HTTP request sent, awaiting response... 303 See Other
Location: https://doc-14-14-docs.googleusercontent.com/docs/securesc/ha0ro937gcuc717deffksulhg5h7mbpl/4n3oik3dp7qkgo
Warning: wildcards not supported in HTTP.
--2022-04-26 16:58:15-- https://doc-14-14-docs.googleusercontent.com/docs/securesc/ha0ro937gcuc717deffksulhg5h7mbpl
Resolving doc-14-14-docs.googleusercontent.com (doc-14-14-docs.googleusercontent.com)... 173.194.218.132, 2607:f8b0:
Connecting to doc-14-14-docs.googleusercontent.com (doc-14-14-docs.googleusercontent.com)|173.194.218.132|:443... co
HTTP request sent, awaiting response... 200 OK
Length: 12175 (12K) [text/csv]
Saving to: 'flight.csv'

flight.csv 100%[=====>] 11.89K --- KB/s in 0s

2022-04-26 16:58:15 (57.7 MB/s) - 'flight.csv' saved [12175/12175]

[4] flights = pd.read_csv('./flight.csv')
flights.info()

<class 'pandas.core.frame.DataFrame'>

Done (Total rows: 1000 entries)

RAM Disk

Up Down Refresh Settings

Linux Server

Terminal / Shell / Command Line

colab.research.google.com/drive/1JZ-TbmRIY8mmD018C_grijQHluBGNQxZ#scrollTo=46VlhWzHWhRZ

+ Code + Text

RAM Disk

https://drive.google.com/uc?export=download&id=1PazlhissU63pozk0jJckyjIJuM_u-0JF

[3] !wget "https://drive.google.com/uc?export=download&id=1PazlhissU63pozk0jJckyjIJuM_u-0JF" -O flight.csv

--2022-04-26 16:58:14-- https://drive.google.com/uc?export=download&id=1PazlhissU63pozk0jJckyjIJuM_u-0JF
Resolving drive.google.com (drive.google.com)... 74.125.26.101, 74.125.26.138, 74.125.26.100, ...
Connecting to drive.google.com (drive.google.com)|74.125.26.101|:443... connected.
HTTP request sent, awaiting response... 303 See Other
Location: <https://doc-14-14-docs.googleusercontent.com/docs/securesc/ha0ro937gcuc717deffksulhg5h7mbpl/4n3oik3dp7qkgo>
Warning: wildcards not supported in HTTP.
--2022-04-26 16:58:15-- <https://doc-14-14-docs.googleusercontent.com/docs/securesc/ha0ro937gcuc717deffksulhg5h7mbpl>
Resolving doc-14-14-docs.googleusercontent.com (doc-14-14-docs.googleusercontent.com)... 173.194.218.132, 2607:f8b0:
Connecting to doc-14-14-docs.googleusercontent.com (doc-14-14-docs.googleusercontent.com)|173.194.218.132|:443... co
HTTP request sent, awaiting response... 200 OK
Length: 12175 (12K) [text/csv]
Saving to: 'flight.csv'

flight.csv 100%[=====>] 11.89K --.KB/s in 0s

2022-04-26 16:58:15 (57.7 MB/s) - 'flight.csv' saved [12175/12175]

[4] flights = pd.read_csv('./flight.csv')
flights.info()

<class 'pandas.core.frame.DataFrame'>

Done (Total rows: 1000 entries)

Probability2.ipynb - Colaboratory | ProbabilityDistr_1.ipynb - Colaboratory | New Tab

colab.research.google.com/drive/1JZ-TbmRIY8mmD018C_grijyQHluBGNQxZ#scrollTo=46VlnWzHWhRZ

Update

+ Code + Text

✓ RAM Disk



```
[4] flights = pd.read_csv('./flight.csv')
flights.info()
```

{x} <class 'pandas.core.frame.DataFrame'>

RangeIndex: 1000 entries, 0 to 999

Data columns (total 3 columns):

#	Column	Non-Null Count	Dtype
0	Passenger_ID	1000 non-null	int64
1	Flight_ID	1000 non-null	object
2	Arrived	1000 non-null	int64

dtypes: int64(2), object(1)

memory usage: 23.6+ KB

▶ flights.head()

Passenger_ID Flight_ID Arrived

0	1811	A320	1
---	------	------	---

1	1812	A320	1
---	------	------	---

2	1813	B777	1
---	------	------	---

3	1814	B737	1
---	------	------	---

4	1815	B737	1
---	------	------	---



```
Probability2.ipynb - Colaboratory | ProbabilityDistr_1.ipynb - Colaboratory | New Tab
```

```
colab.research.google.com/drive/1JZ-TbmRIY8mmD018C_grijQHluBGNQxZ#scrollTo=46VlnWzHWhRZ
```

```
Update
```

+ Code + Text

✓ RAM Disk



memory usage: 23.6+ KB

```
[ ] flights.head()
```

{x}

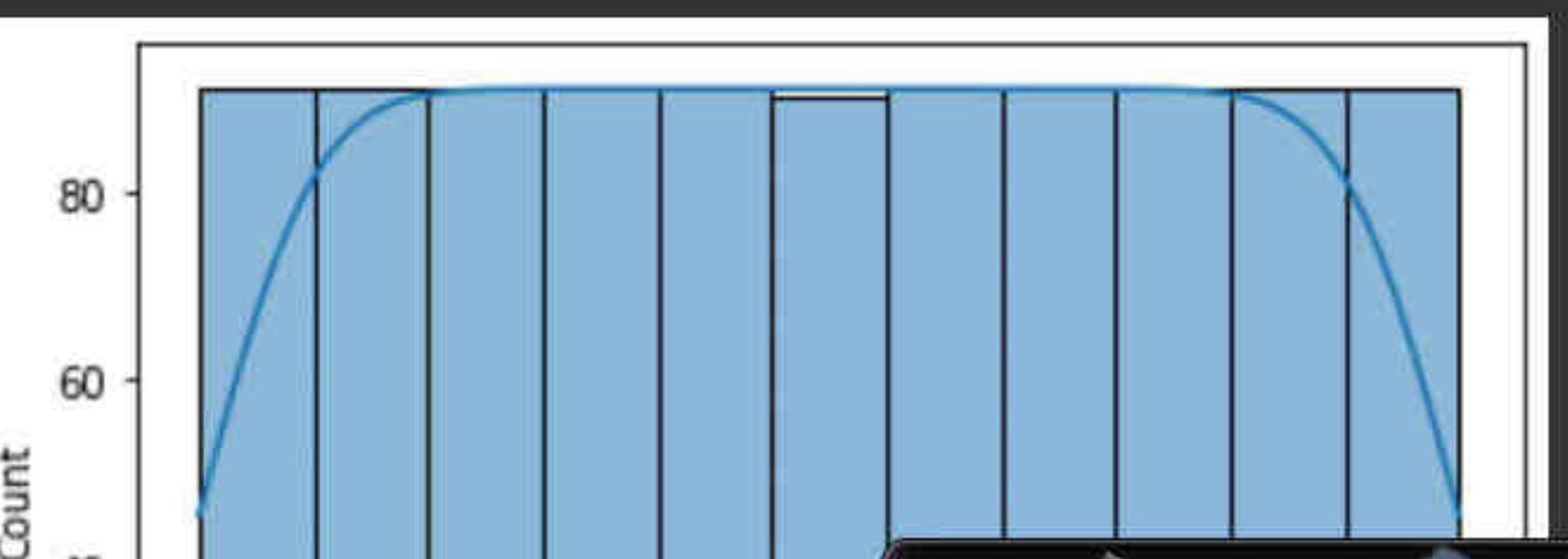
□

{ } { }

	Passenger_ID	Flight_ID	Arrived
0	1811	A320	1
1	1812	A320	1
2	1813	B777	1
3	1814	B737	1
4	1815	B737	1

1811 A320 → 1/0

```
[ ] sns.histplot(flights["Passenger_ID"], kde=True)  
plt.show()
```



colab.research.google.com/drive/1JZ-TbmRIY8mmD018C_grijQHluBGNQxZ#scrollTo=zeaTTNEFWzxG

+ Code + Text

[4] flights.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 3 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   Passenger_ID  1000 non-null   int64  
 1   Flight_ID     1000 non-null   object  
 2   Arrived       1000 non-null   int64  
dtypes: int64(2), object(1)
memory usage: 23.6+ KB
```

int64
object
int64

[] flights.head()

	Passenger_ID	Flight_ID	Arrived
0	1811	A320	1
1	1812	A320	1
2	1813	B777	1
3	1814	B737	1
4	1815	B737	1

Probability2.ipynb - Colaboratory | ProbabilityDsb_1.ipynb - Colaboratory | New Tab

colab.research.google.com/drive/1JZ-TbmRIY8mmD018C_grijQHluBGNQxZ#scrollTo=zeaTTNEFWzxG

Update

+ Code + Text

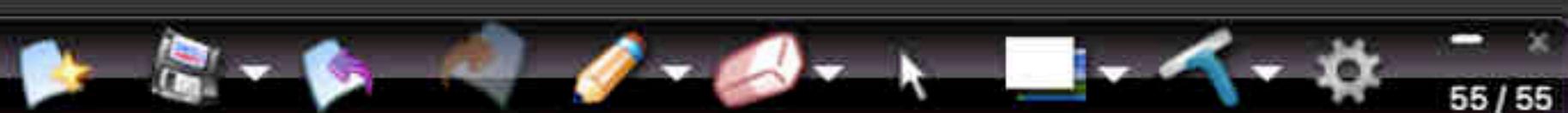
flights = pd.read_csv('./flights.csv')

[4] flights.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 3 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   Passenger_ID  1000 non-null   int64  
 1   Flight_ID     1000 non-null   object  
 2   Arrived       1000 non-null   int64  
dtypes: int64(2), object(1)
memory usage: 23.6+ KB
```

[] flights.head()

	PID	FID	A
	Passenger_ID	Flight_ID	Arrived
0	1811	A320	1
1	1812	A320	1
2	1813	B777	1
3	1814	B737	1
4	1815	B737	1



Probability2.ipynb - Colaboratory | ProbabilityDistr_1.ipynb - Colaboratory | New Tab

colab.research.google.com/drive/1JZ-TbmRIY8mmD018C_grijQHluBGNQxZ#scrollTo=zeaTTNEFWzxG

Update

+ Code + Text

flights = pd.read_csv('~/flights.csv')

[4] flights.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 3 columns):
 # Column Non-Null Count Dtype

 0 Passenger_ID 1000 non-null int64
 1 Flight_ID 1000 non-null object
 2 Arrived 1000 non-null int64
 dtypes: int64(2), object(1)
 memory usage: 23.6+ KB

sample

[] flights.head()

	Passenger_ID	Flight_ID	Arrived
0	1811	A320	1
1	1812	A320	1
2	1813	B777	1
3	1814	B737	1
4	1815	B737	1

RAM Disk





+ Code + Text

[4] Data columns (total 3 columns):

#	Column	Non-Null Count	Dtype
0	Passenger_ID	1000 non-null	int64
1	Flight_ID	1000 non-null	object
2	Arrived	1000 non-null	int64

dtypes: int64(2), object(1)
memory usage: 23.6+ KB

[] flights.head()

	Passenger_ID	Flight_ID	Arrived
0	1811	A320	1
1	1812	A320	1
2	1813	B777	1
3	1814	B737	1
4	1815	B737	1

sns.histplot(flights["Passenger_ID"], kde=True)
plt.show()

random-var

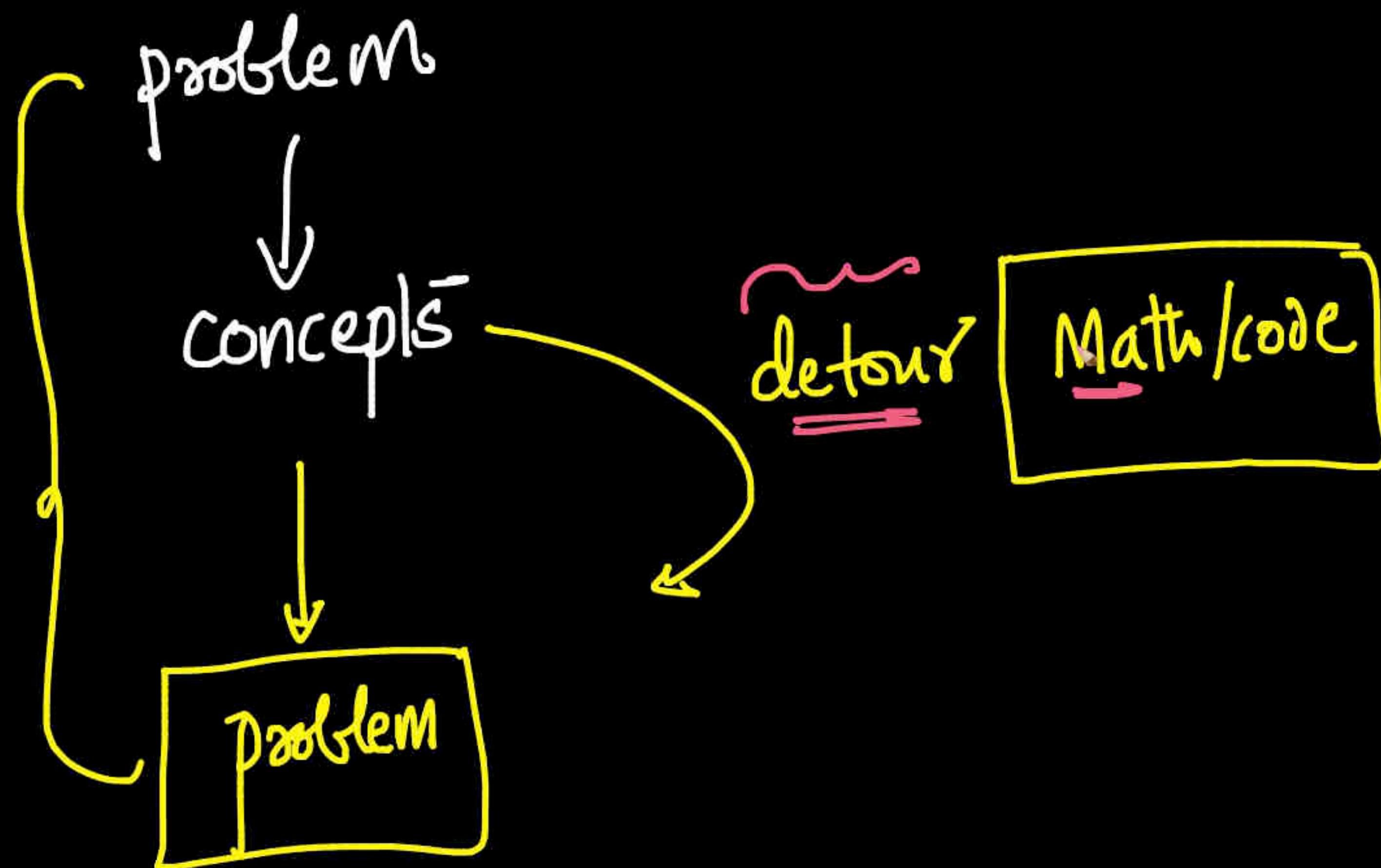
A = 1 or 0

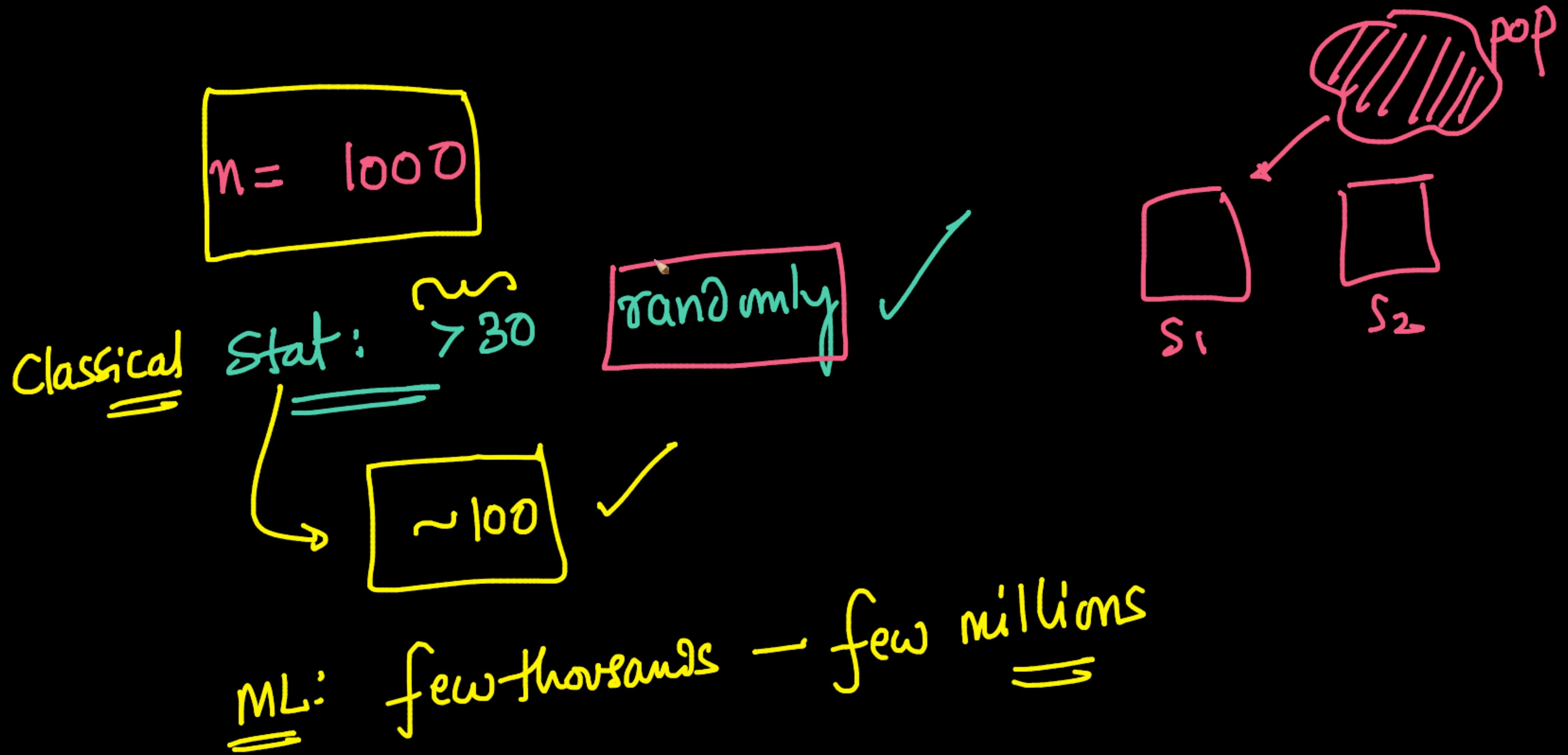
↳ categorical / binary

pID: r.v. numeric → continuous

1000 - 9999

C



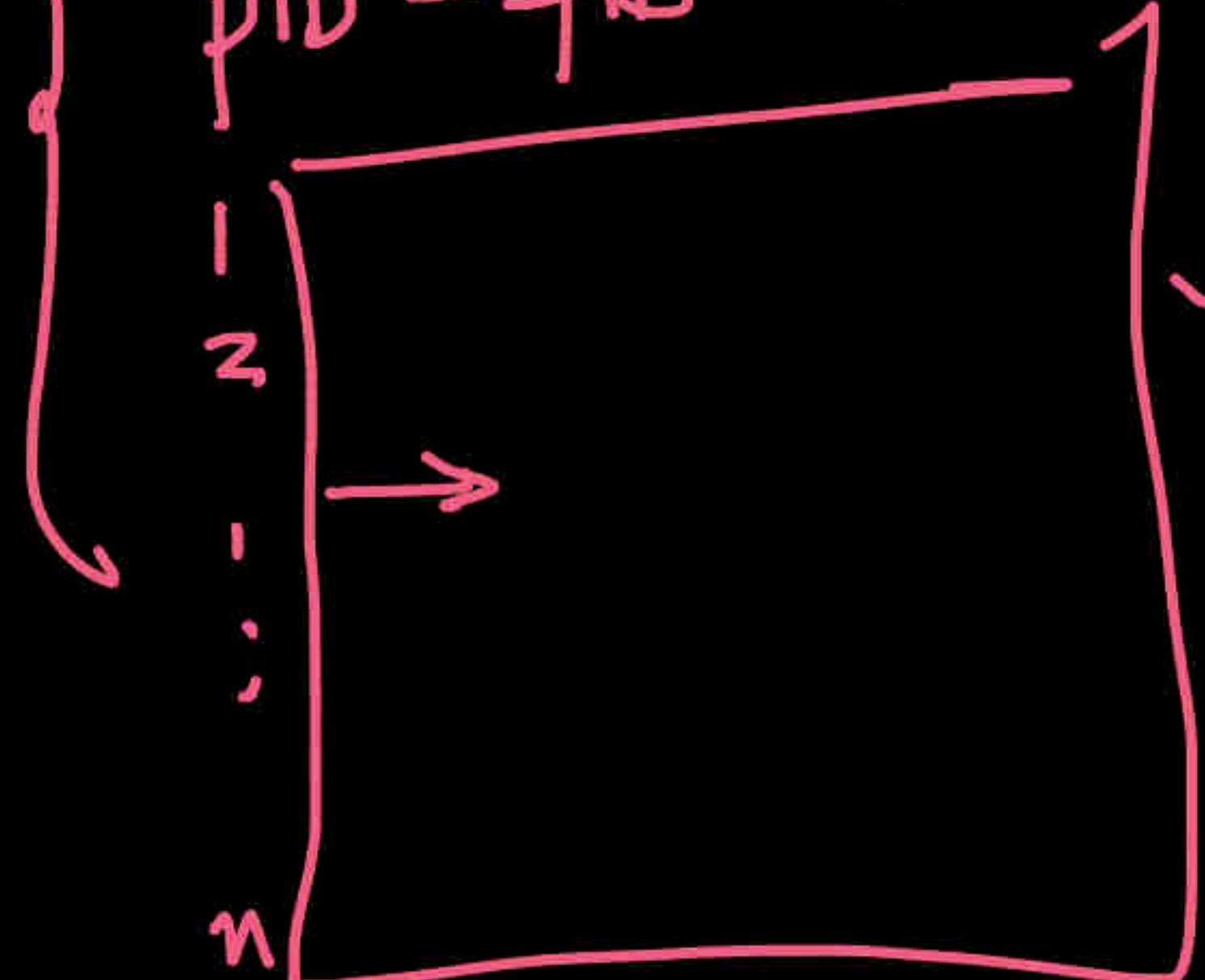


airline

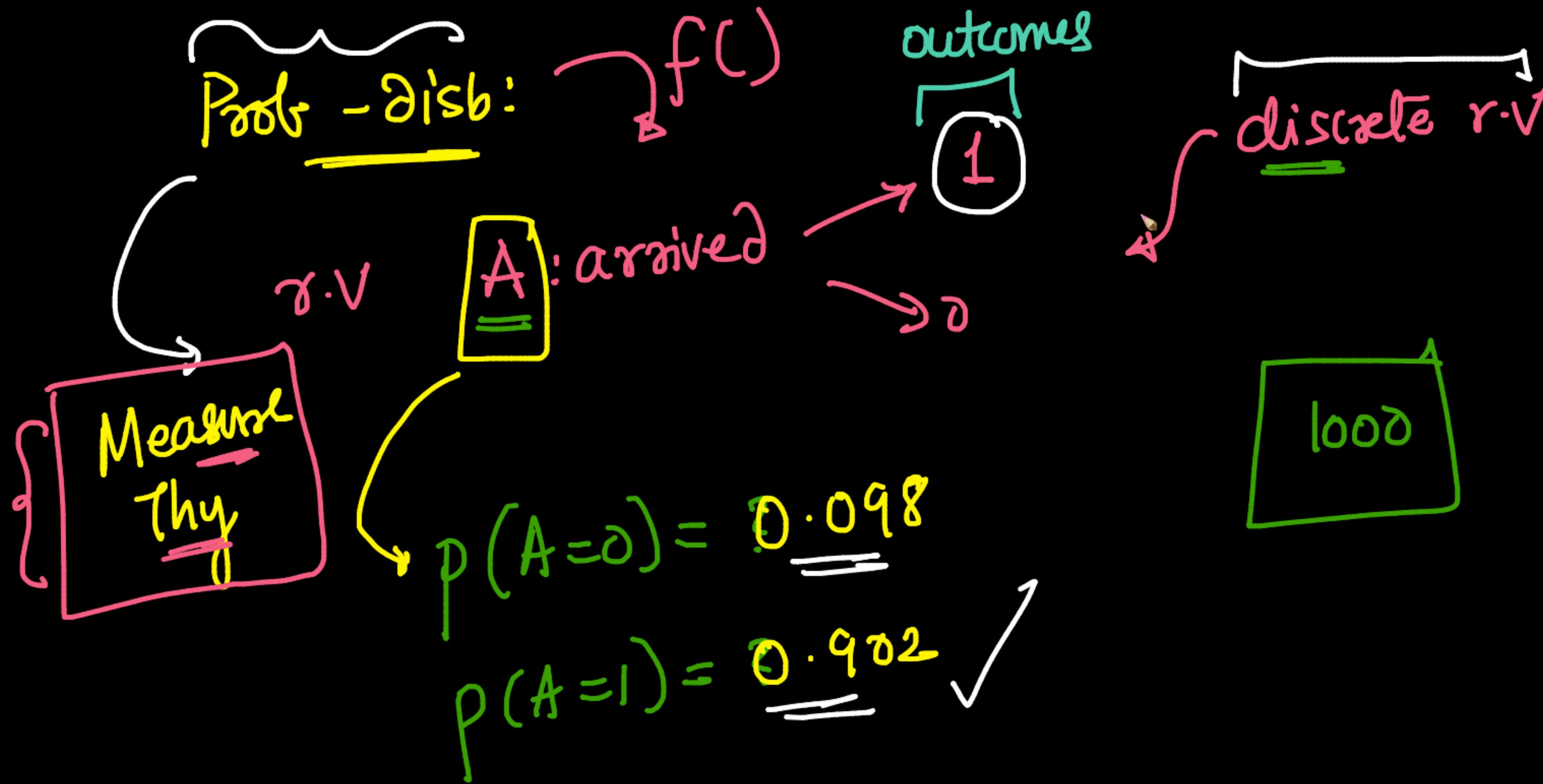
IM → 1000 flights

~~5428~~

{ PID - fid - isArrived

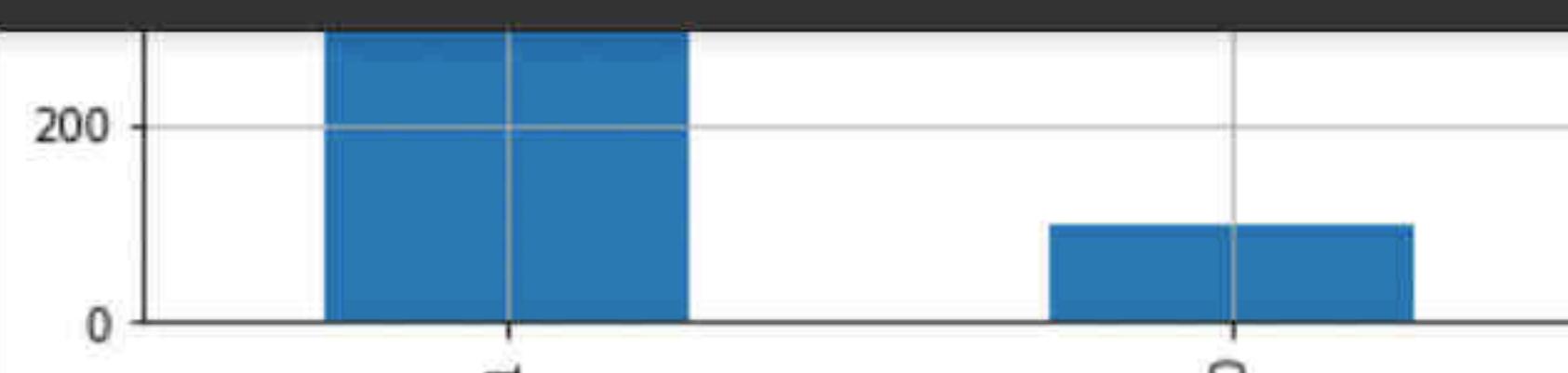


1000 dataset
~~=~~



+ Code + Text

[5]



{x}

```
shows_up_probability = flights['Arrived'].value_counts(normalize=True)[1]
print(shows_up_probability)
```

0.902

```
[ ] flights['Arrived'].value_counts()
  1    902
  0    98
Name: Arrived, dtype: int64
```

$$\frac{902}{1000} \rightarrow 0.902$$

$$\frac{902}{1000} = P(A=1) = 0.902$$

```
[ ] import math

PENALTY = 10000
def comb(n, r):

    num1 = math.factorial(n)
    num2 = math.factorial(r)
    num3 = math.factorial(n-r)
```

$$\frac{98}{1000} = P(A=0)$$



Fair dice → 6 sides

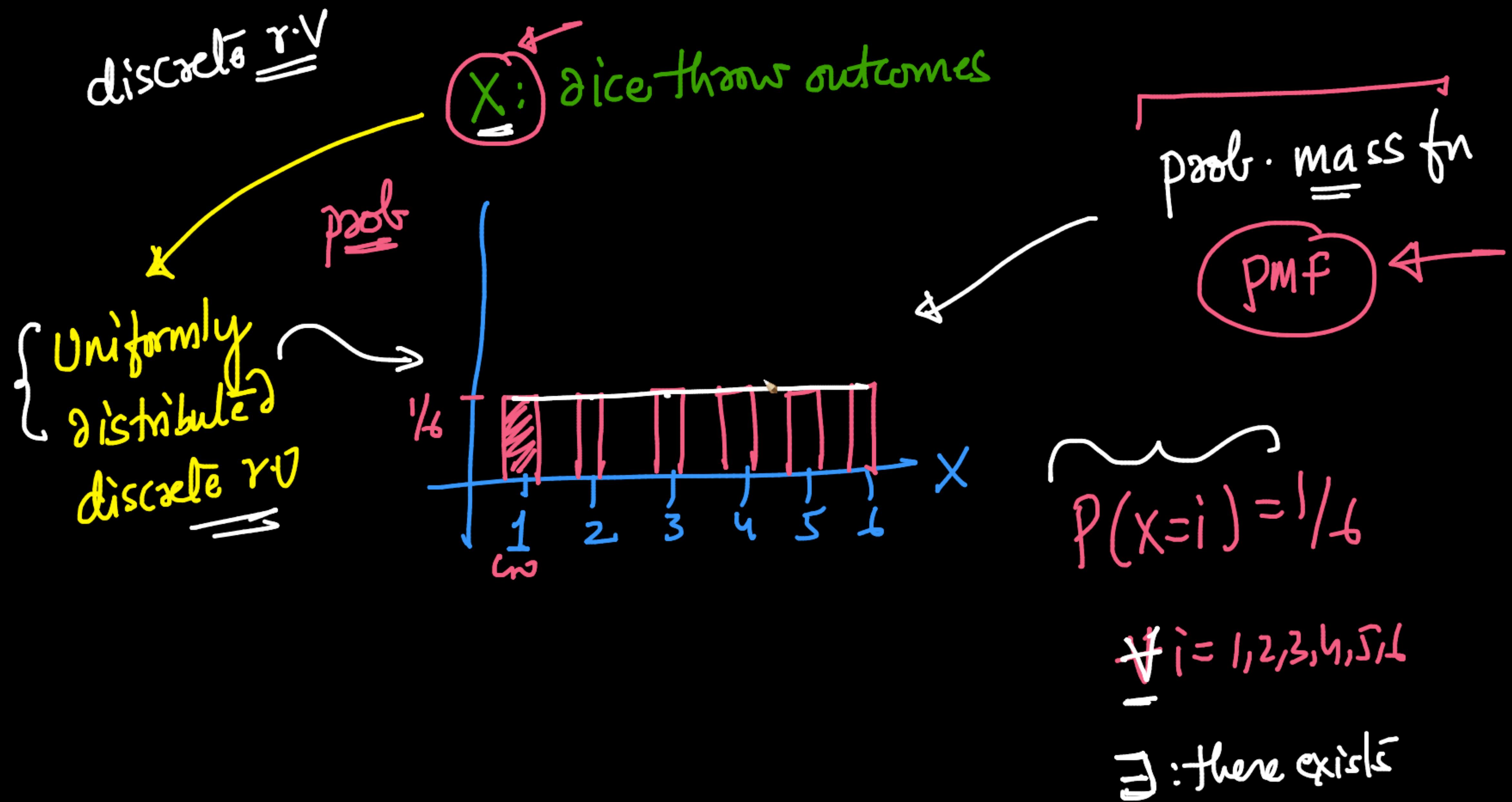
1, 2, 3, 4, 5

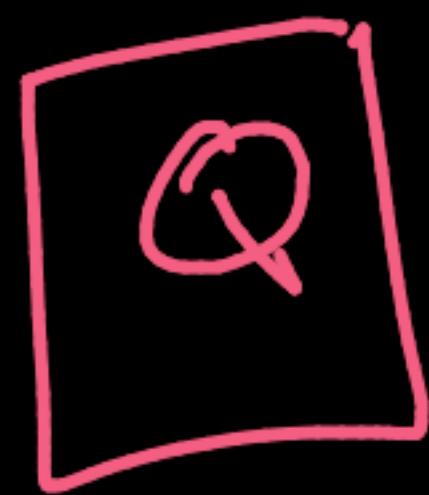
4.5

g.v X^1 : represent the outcome of a throw
→ $\{1, 2, 3, 4, 5, 6\}$ Categorical

$$P(X=1) = \frac{1}{6} = P(X=2) \dots = P(X=6)$$

Prob. dist





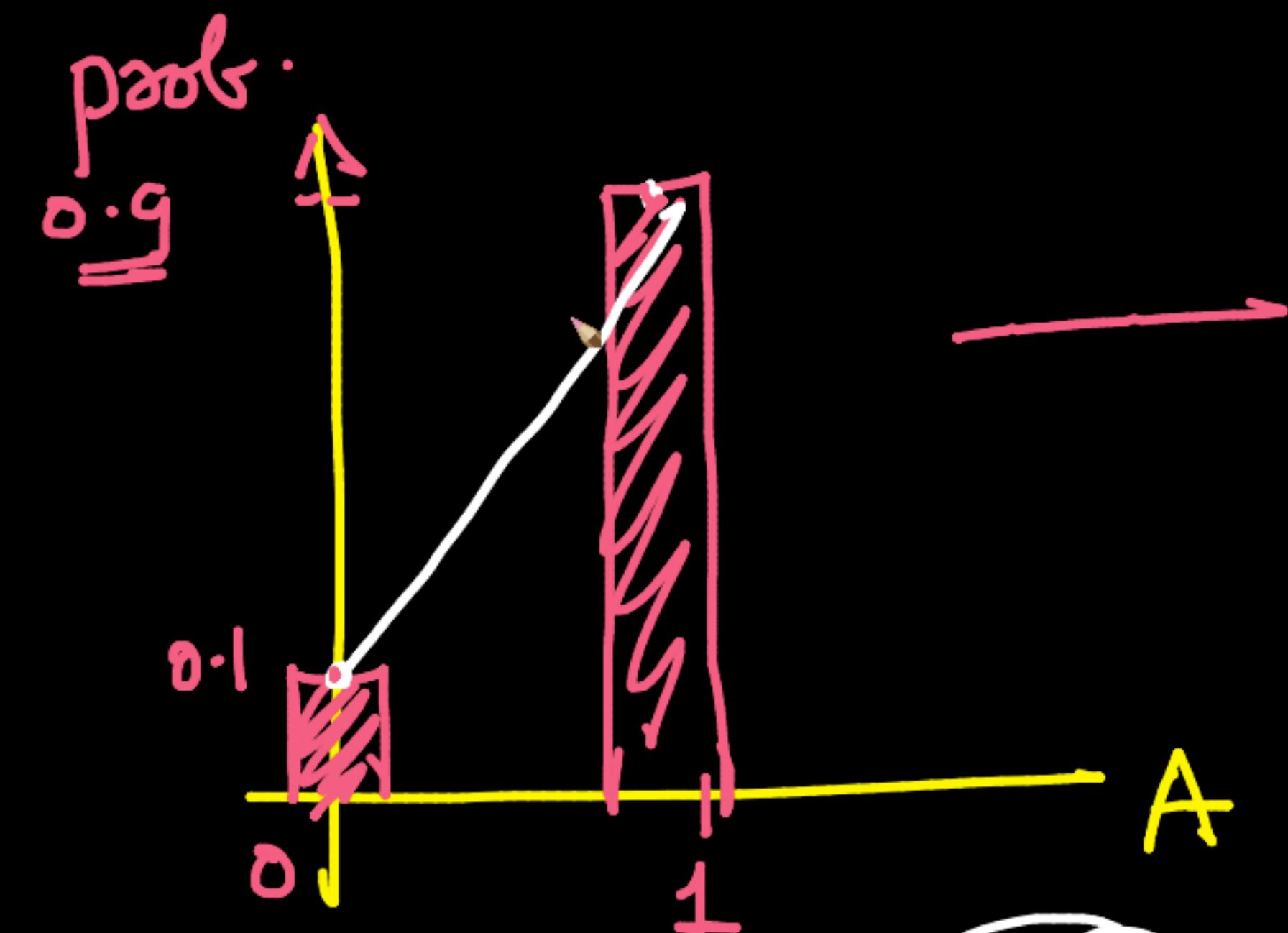
discrete $A = \{0 \text{ or } 1\}$

$\checkmark P(A=0) = 0.1$

$\checkmark P(A=1) = 0.9$

0	1
---	---

1.0



Extremely

viz ✓
intuition

Terminology → overwhelming]
→ contextualize] - revision


{ 80 mul-pool

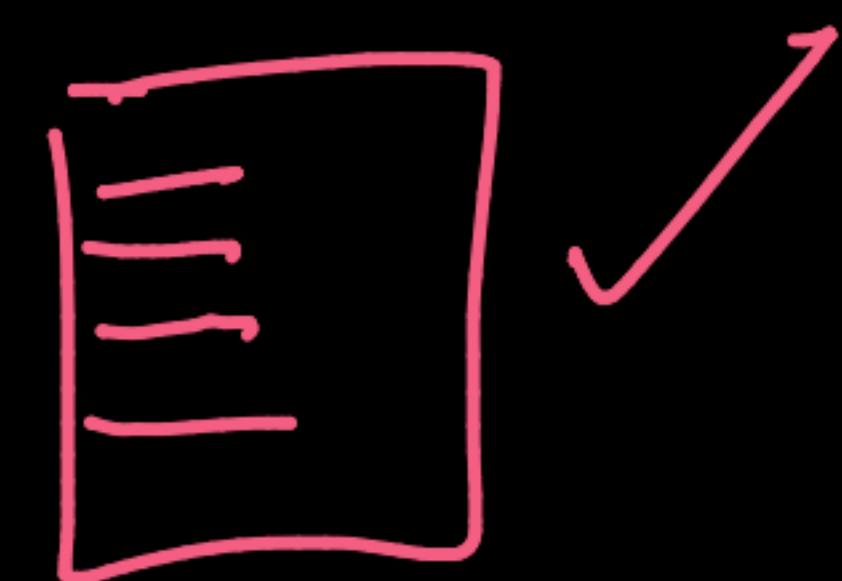
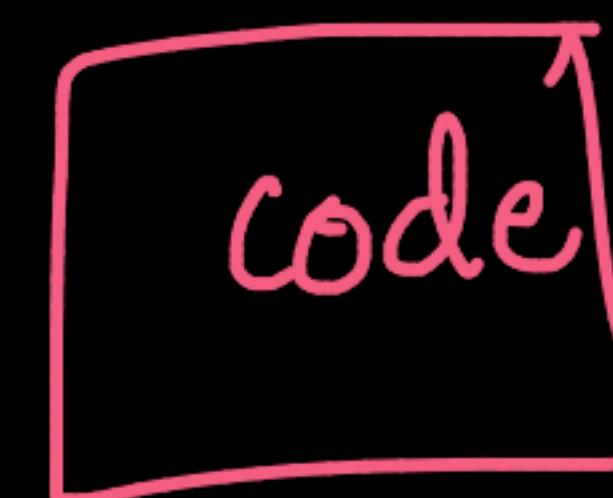
[PMF : discrete r.v]

PDF (later) : continuous r.v

{ → notes ✓

→ work-out the concepts in your own [test]

Feynman



freq:
— week → thorough
— month
— 3-month

Q

A: 0 or 1

$$P(A=0) = \boxed{0 \cdot 1} \quad P(A=1) = \boxed{0 \cdot 9}$$

Uniform r.v? \rightarrow No

Q

fair coin
X: T or H

$$P(X=T) = P(X=H) = \underline{\underline{\frac{1}{2}}}$$

Uniform r.v

CDF: Cumulative disb. fn

discrete r.v

fair dice $X \in \{1, 2, 3, 4, 5, 6\}$

$P(X \leq i)$ $\forall i = 1, 2, 3, 4, 5, 6$

$$P(X \leq 1) = 1/6$$

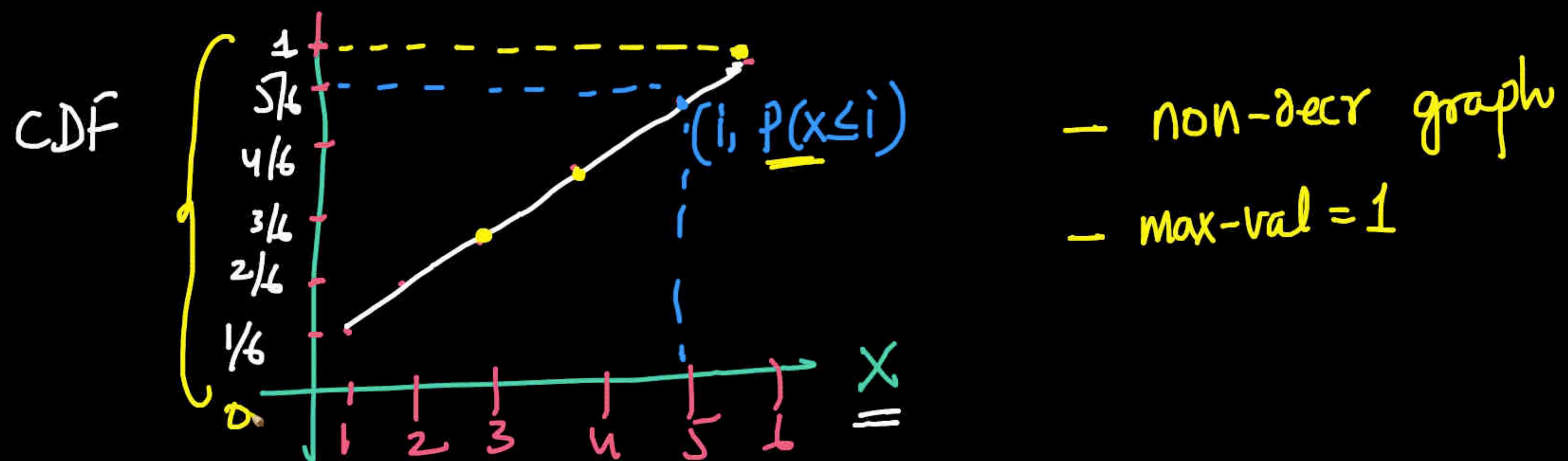
$$P(X \leq 2) = 2/6 = P(X=1) + P(X=2)$$

:

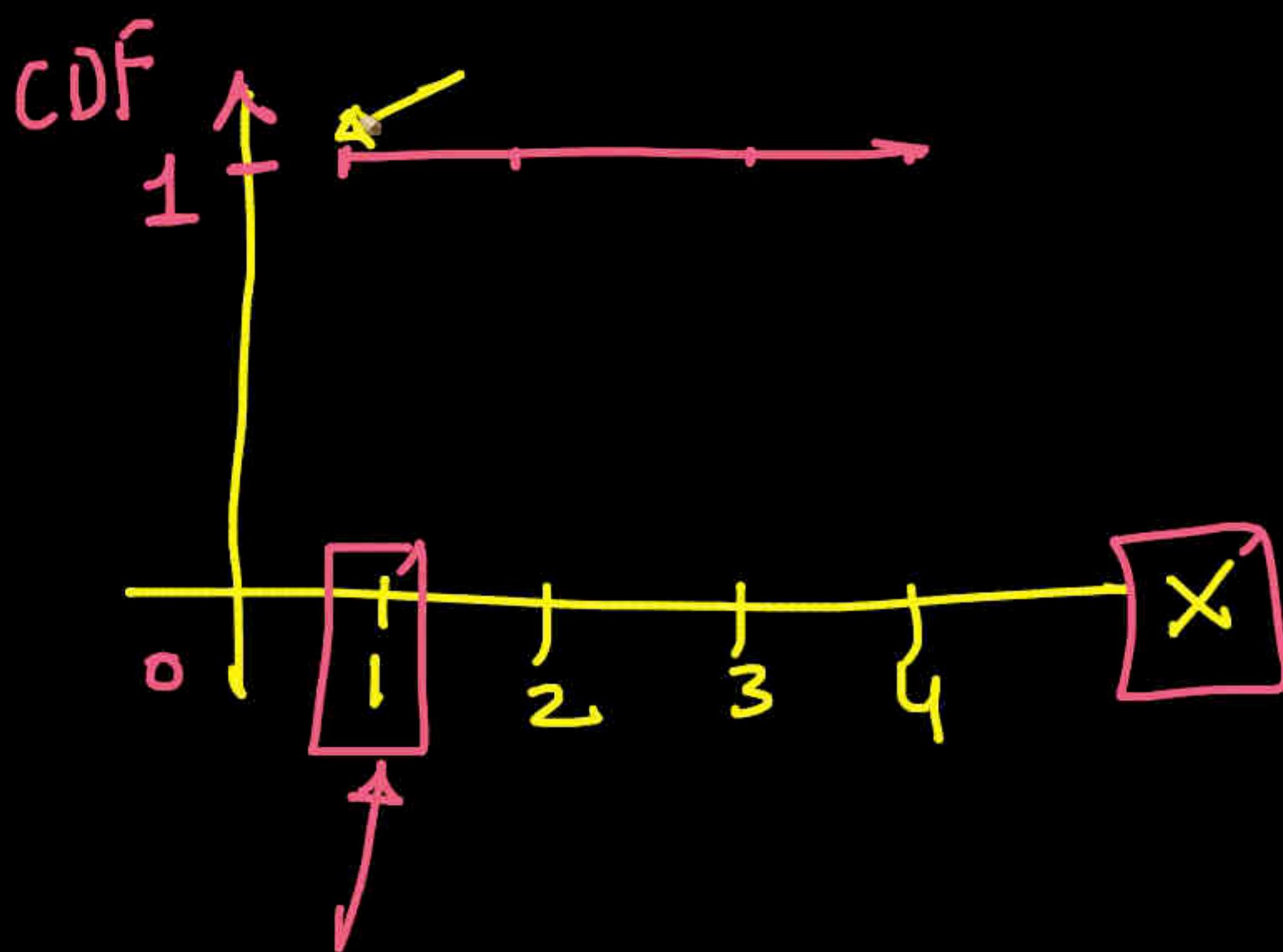
$$P(X \leq 6) = 6/6$$

PMF: $P(X=i)$

$\forall i$



$$P(X \leq u) = P(X \leq 3) + \sum_{x > 3}$$



$$P(X \leq 1) = 1$$

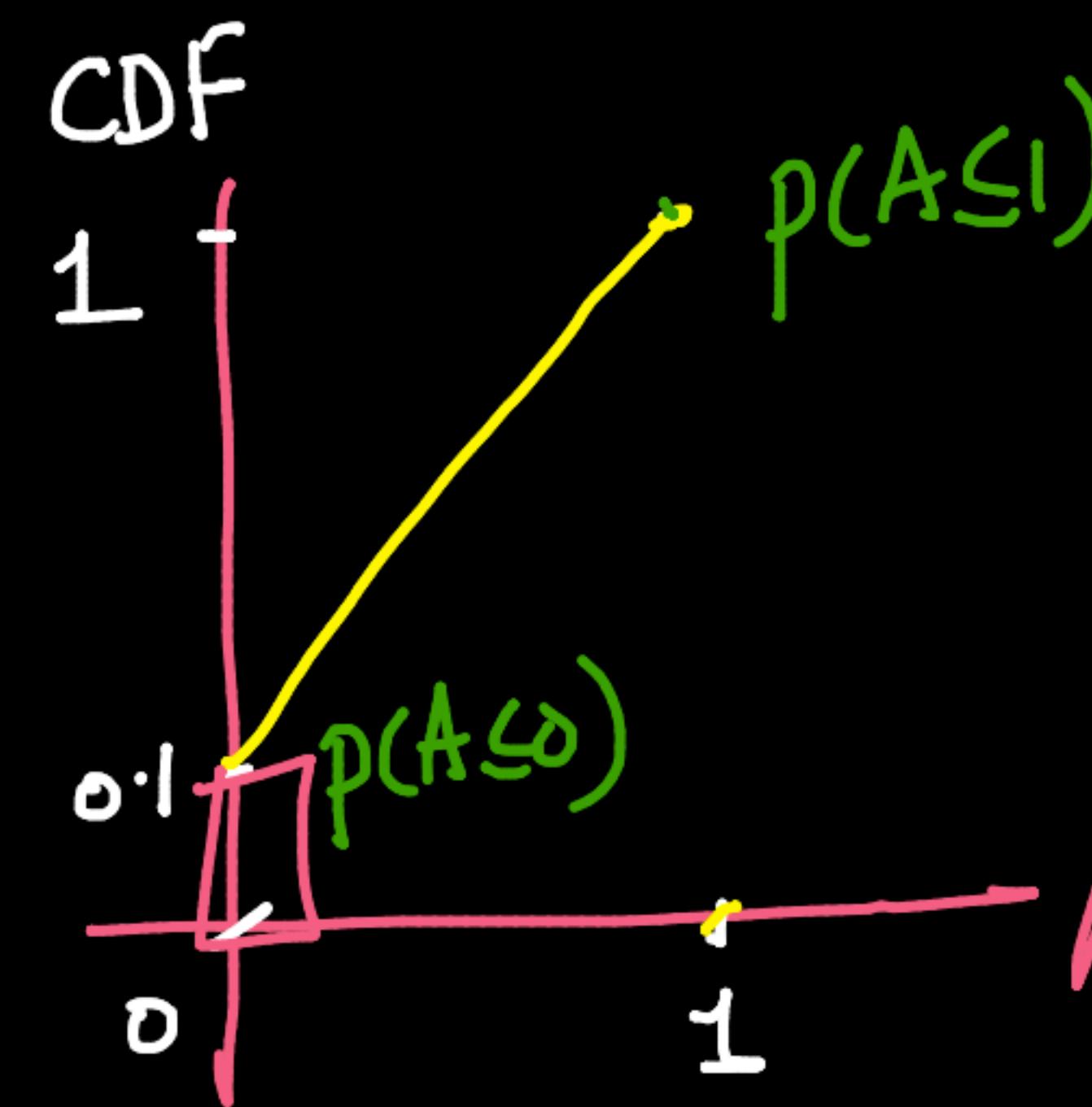
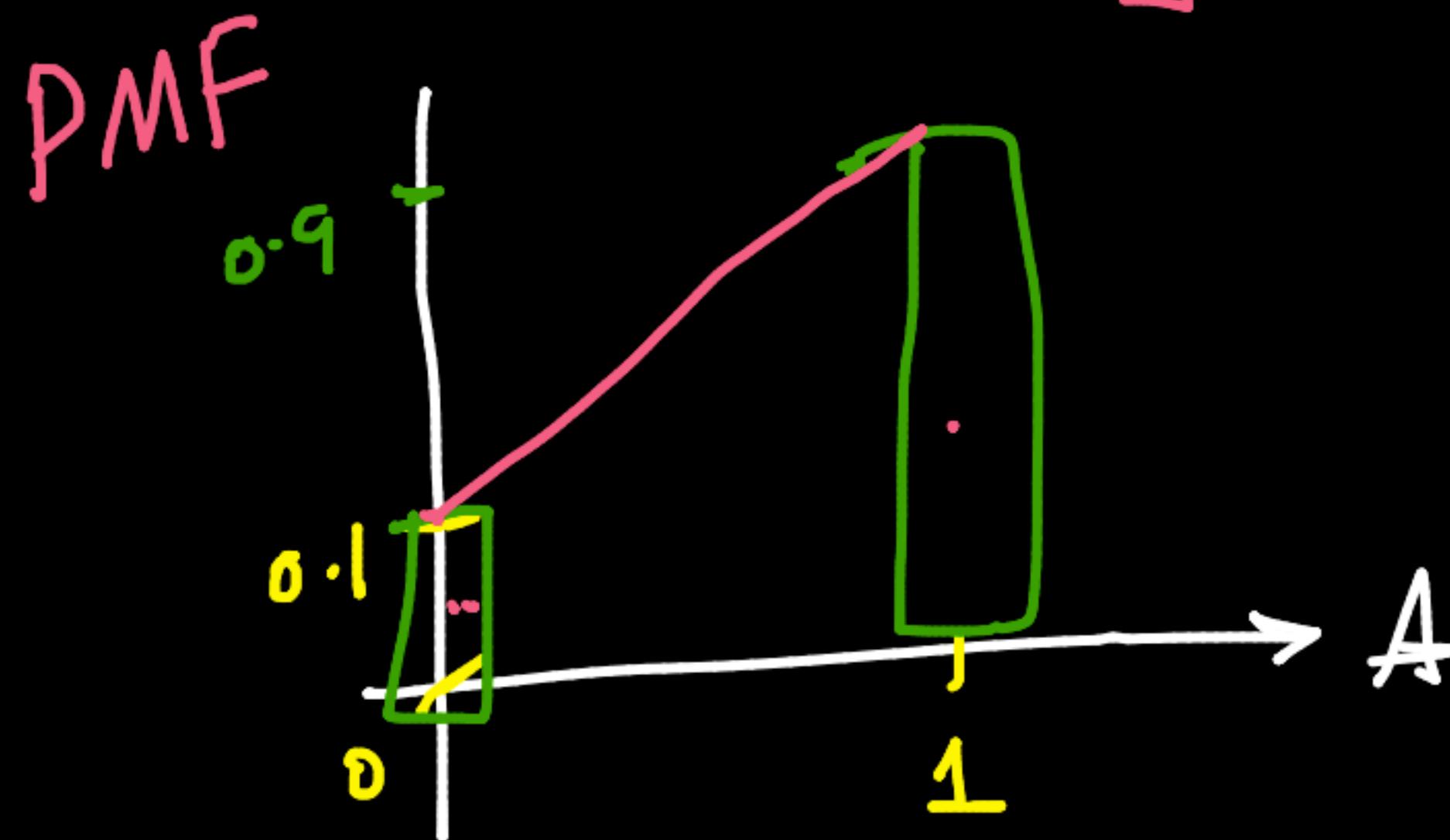
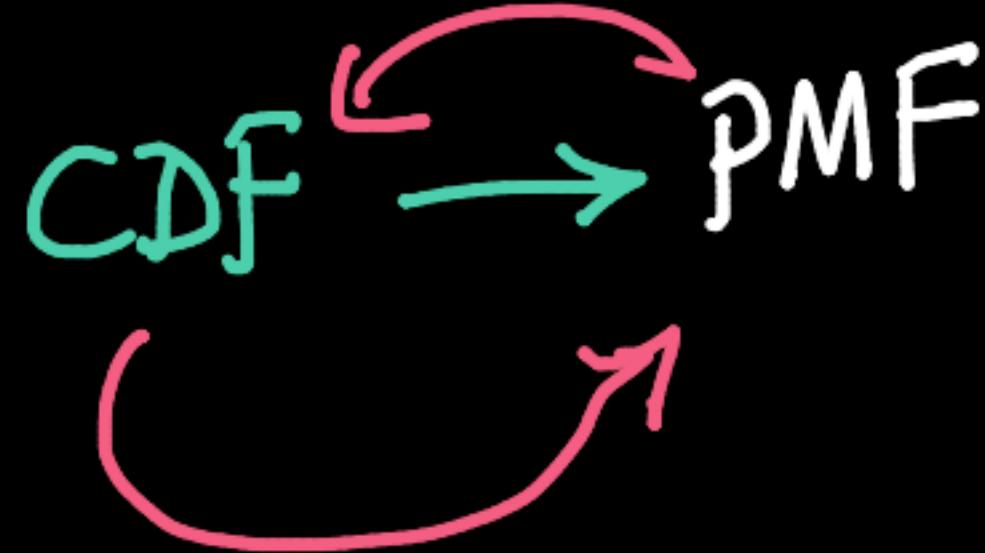
$$P(X \leq 2) = 1$$

$$P(X \leq 3) = 1$$

$$P(X \leq 4) = 1$$

$$X \in \{1, 2, 3, 4\}$$

= discrete r.v

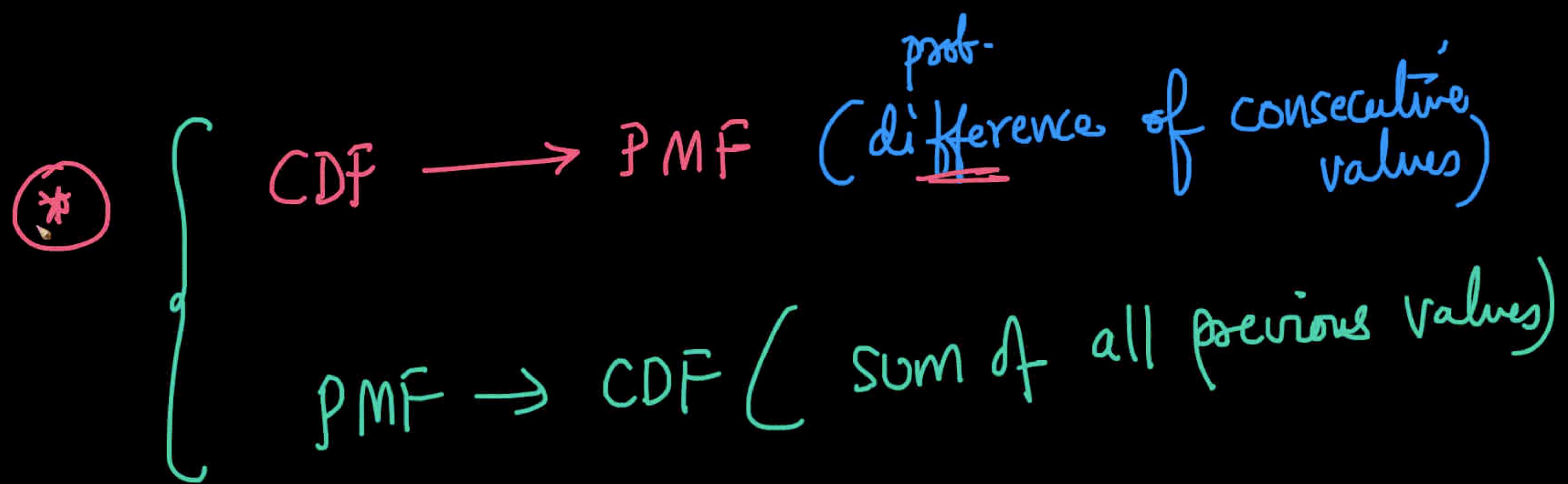


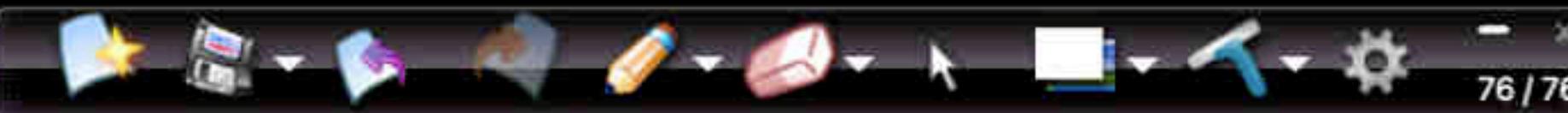
$$P(A \leq 0) = 0.1$$

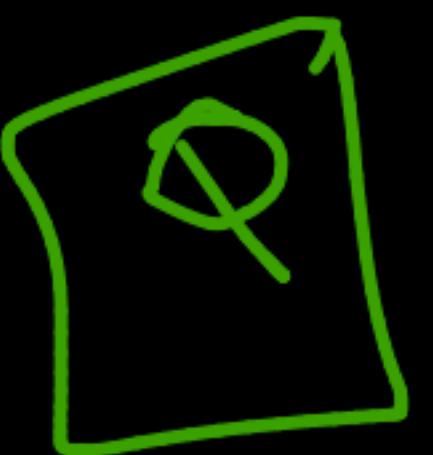
$$P(A \leq 1) = 1.0$$

$$A \in \{0, 1\}$$

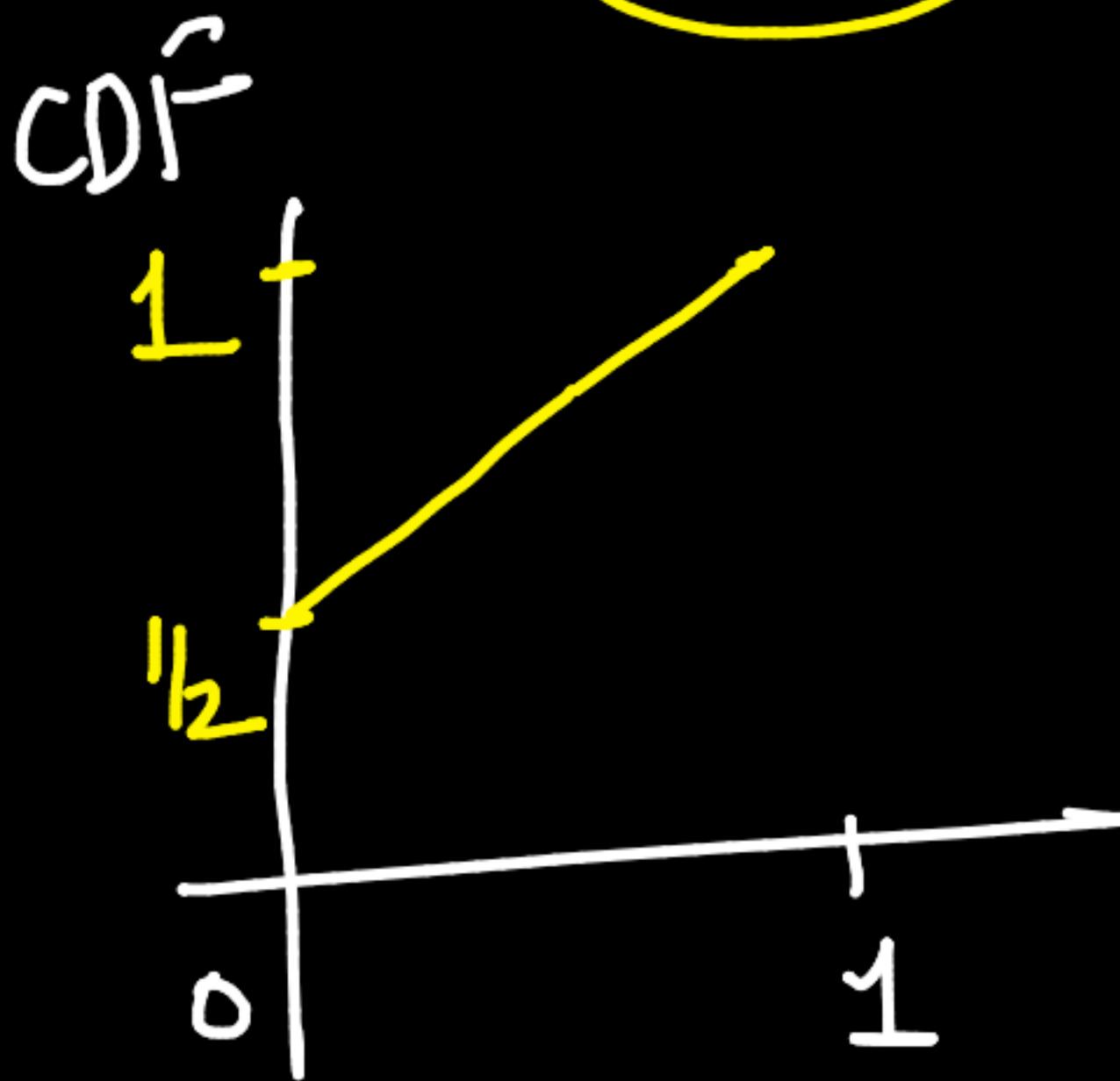
$$P(A=1)$$







Y.V
X. discacelē
H&T
1 0
↓
create
matrix

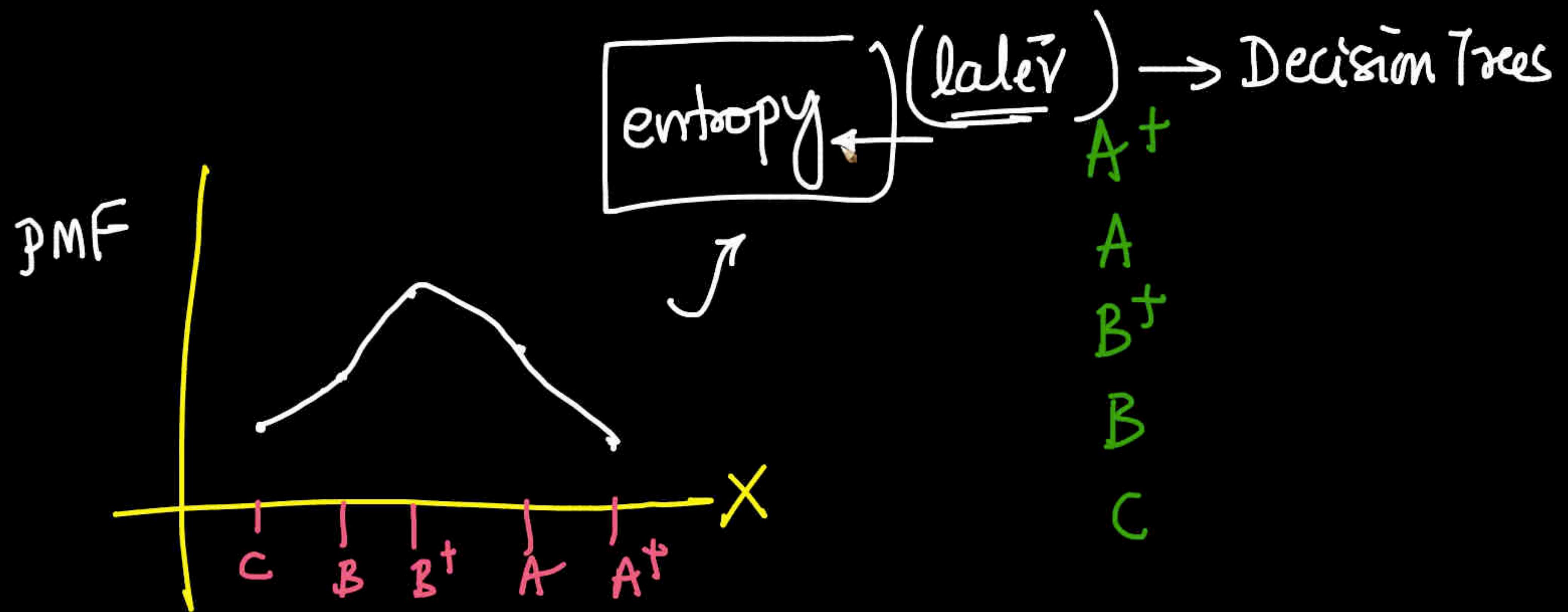


CDF
✓
 $P(X \leq i)$

PMF

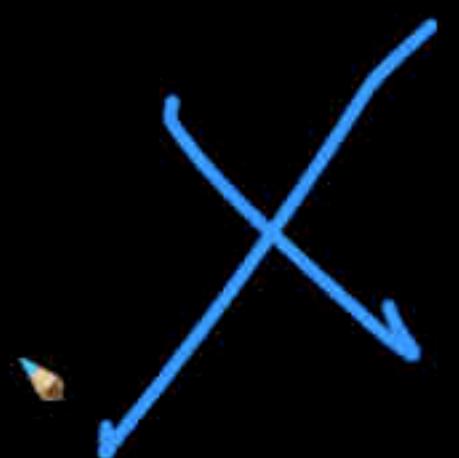
$p(x=i)$

$$p(x=1) = \frac{1}{2} = p(x=H)$$



{ q - 10:30 - class
10:30 - 10:40 break
10:40 - 12 - class
Q&A

Slow



✓ { 9:00 - 10:30
10:30 - 10:35
10:35 - 11:30
11:30 } — stick to this

revise

→ speed up
[5 classes later]

{ CDF , PDF of continuous r.v
|
|
|
Bernoulli r.v
Binomial r.v

decent - pace

= Simulation - code





Search Google or type a URL

Update

Gmail Images



Google

Search Google or type a URL



Learning



Colaboratory



GitHub



My Drive



InterviewBit S...



YouTube



Scaler Academ...



InterviewBit

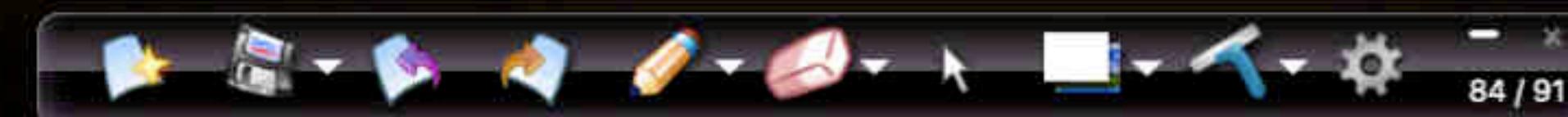


(99+) Feed

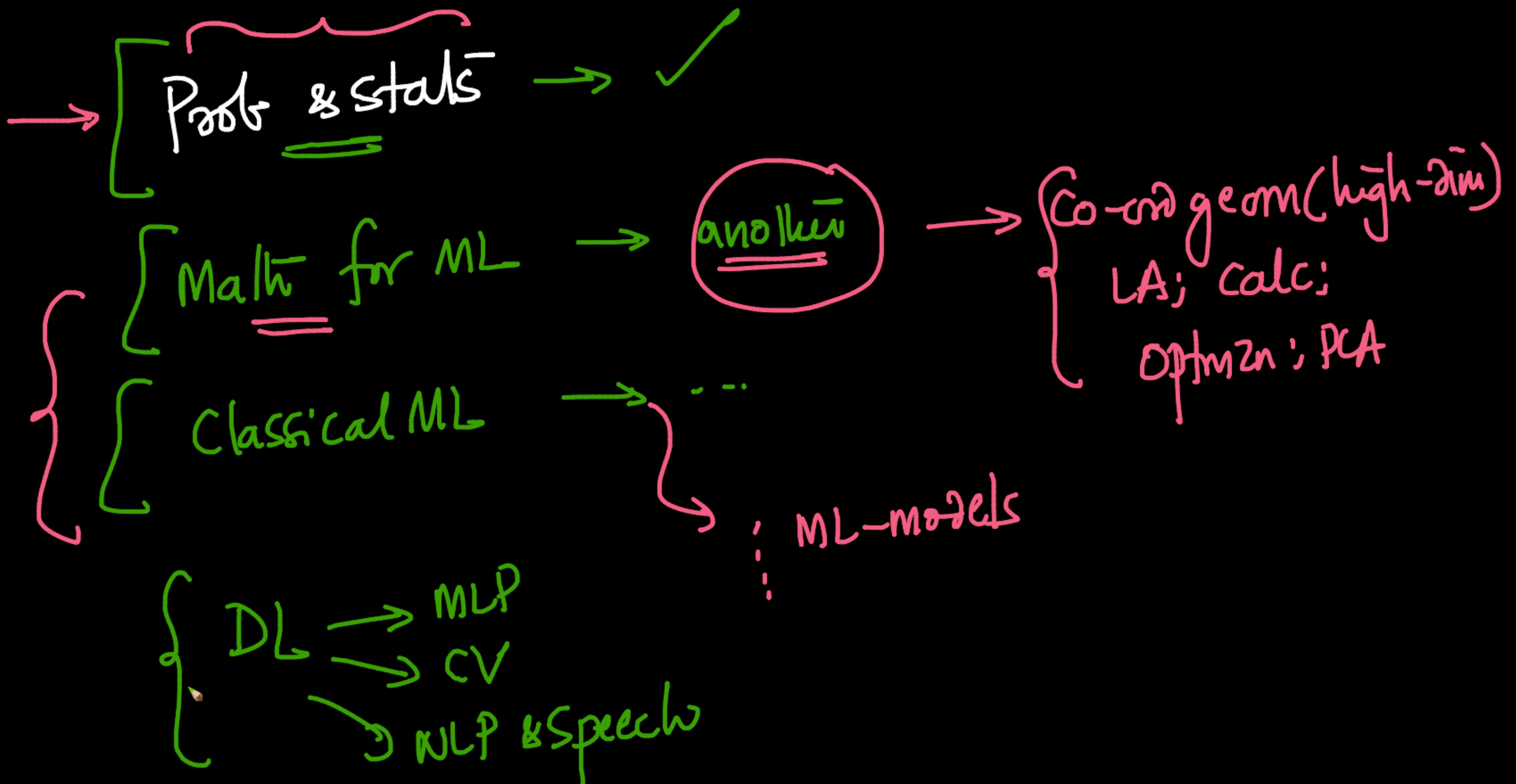


Add shortcut

Stay in control with easy-to-use privacy settings on Google

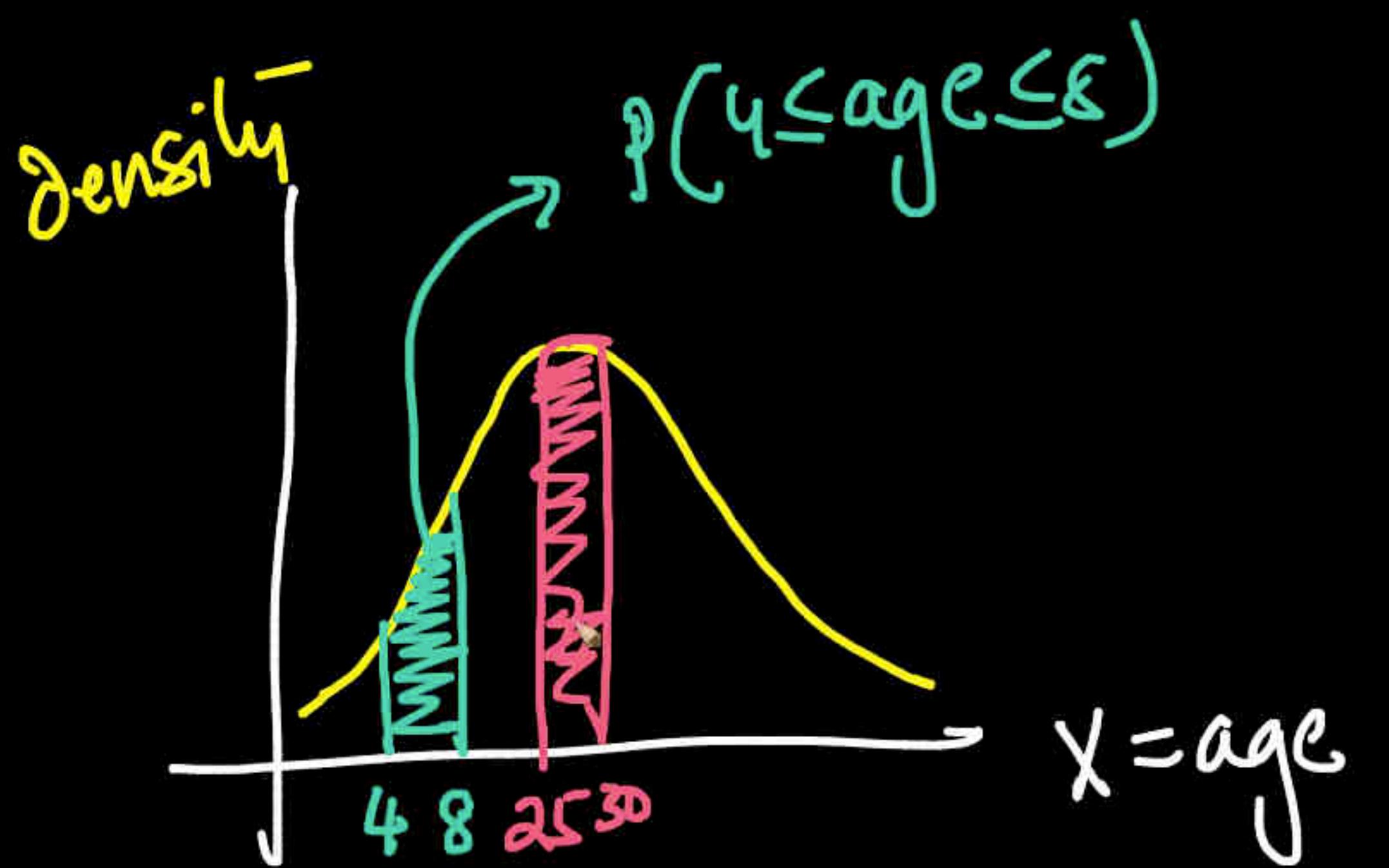


84 / 91



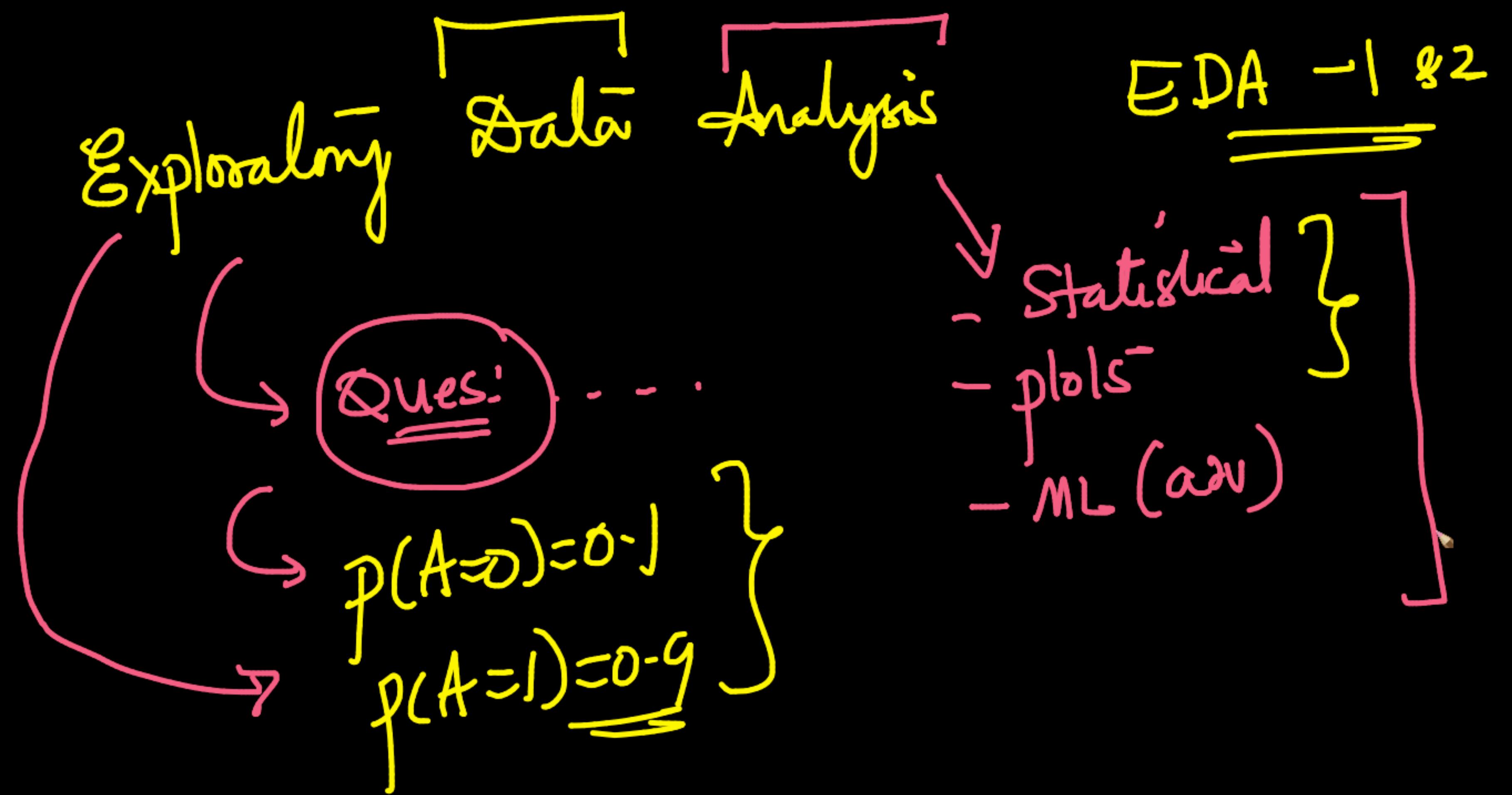
Outliers → remove
transformation (log-normal;
box-cox)
imputation - ML

B.R.
stats



PDF & CDF (next-class)

EDA:



Outlier

- IQR →
- ML
- DL

