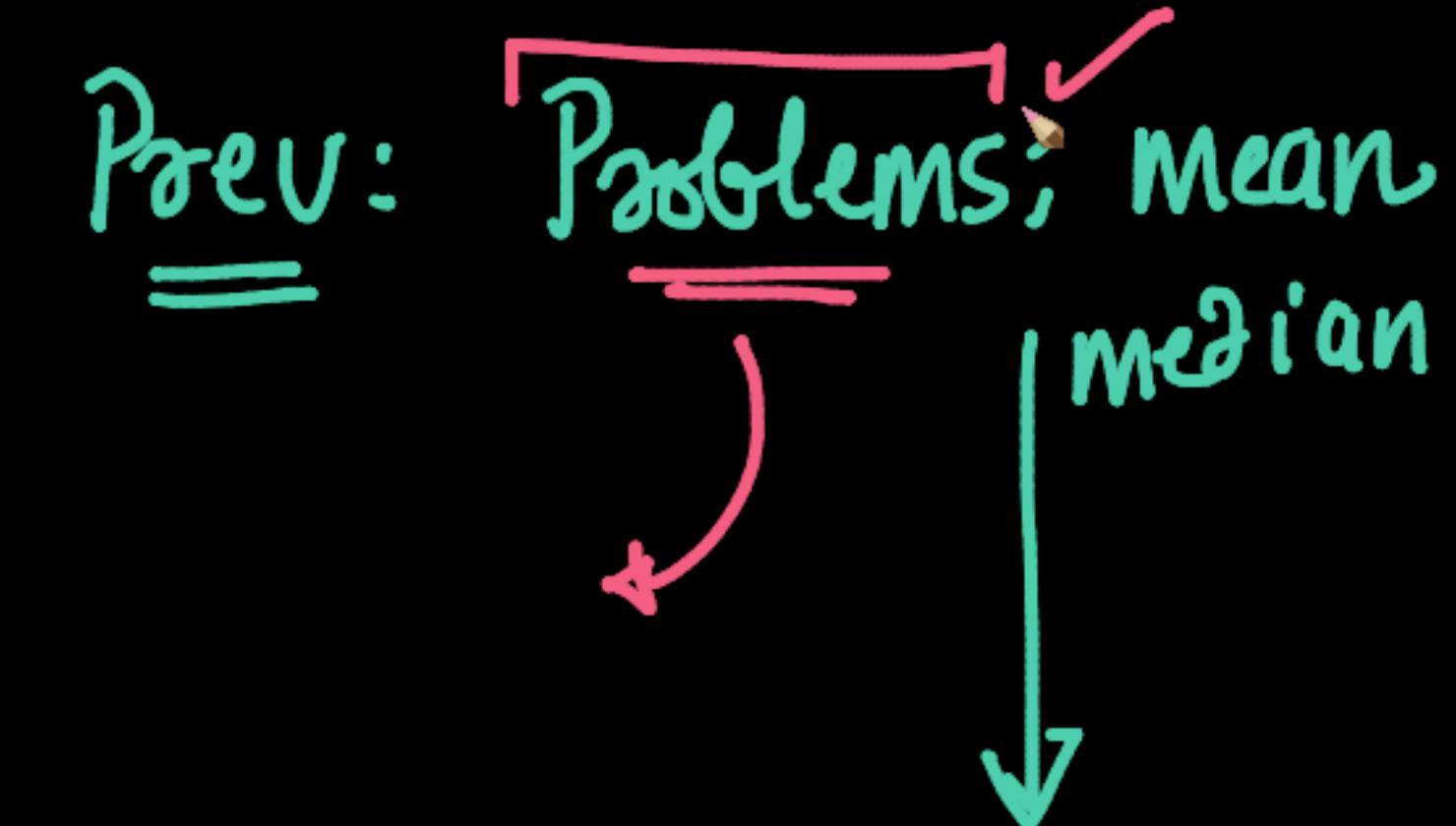


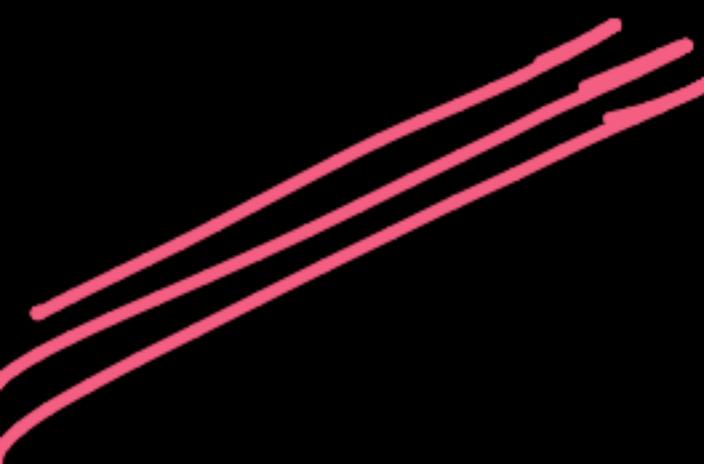
Topics:

- mode; weighted-mean
- std.; variance; MAD
- percentiles; Quantiles; IQR & outliers
- histogram; density-plots
- box plot; bar plot; scatterplot

Retail
data →



{ Airlines, { booting } → - Probability distributions →



mean: average

median: middle-value

-Code-

Mode:



{ most frequently
occurring value
=

discrete r.v

Probability2.ipynb - Colaboratory x ProbabilityDisb_1.ipynb - Colaboratory x +

colab.research.google.com/drive/1l5T7TVIAASw9Tdl4JxqJxuDqRgFBGOW3#scrollTo=SGqzqcJJbVA9

Update

+ Code + Text

RAM Disk

```
[ ] #mode  
df["Education"].value_counts()
```

```
{x}  
Graduation    1127  
PhD           486  
Master         370  
2n Cycle      203  
Basic          54  
Name: Education, dtype: int64
```

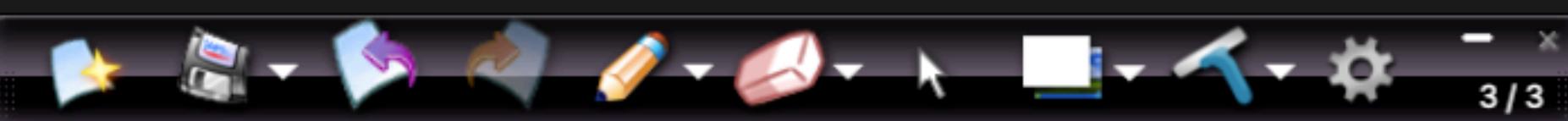
Education: discrete-r.v (categorical)

→ ordinal

```
[ ] #variance  
print("Gold:", df['MntGoldProds'].std())  
print("Fruits:", df['MntFruits'].std())  
print("Sweets:", df['MntSweetProducts'].std())  
print("Wine:", df['MntWines'].std())  
print("Meat:", df['MntMeatProducts'].std())  
print("Fish:", df['MntFishProducts'].std())
```

```
Gold: 52.167438914997064  
Fruits: 39.77343376457871  
Sweets: 41.2804984878548  
Wine: 336.5973926053717  
Meat: 225.71537251175445  
Fish: 54.62897940287769
```

```
[10] from scipy import stats  
print(stats.median_absolute_deviation(df['MntGoldProds']))
```



Probability2.ipynb - Colaboratory | ProbabilityDisb_1.ipynb - Colaboratory | +

colab.research.google.com/drive/1l5T7TVIAASw9Tdl4JxqJxuDqRgFBGOW3#scrollTo=SGqzqcJJbVA9

+ Code + Text RAM Disk Update

[] #mode

df["Education"].value_counts()

Graduation 1127 ✓

PhD 486

Master 370

2n Cycle 203

Basic 54

Name: Education, dtype: int64

[] #variance

print("Gold:", df['MntGoldProds'].std())

print("Fruits:", df['MntFruits'].std())

print("Sweets:", df['MntSweetProducts'].std())

print("Wine:", df['MntWines'].std())

print("Meat:", df['MntMeatProducts'].std())

print("Fish:", df['MntFishProducts'].std())

Gold: 52.167438914997064

Fruits: 39.77343376457871

Sweets: 41.2804984878548

Wine: 336.5973926053717

Meat: 225.71537251175445

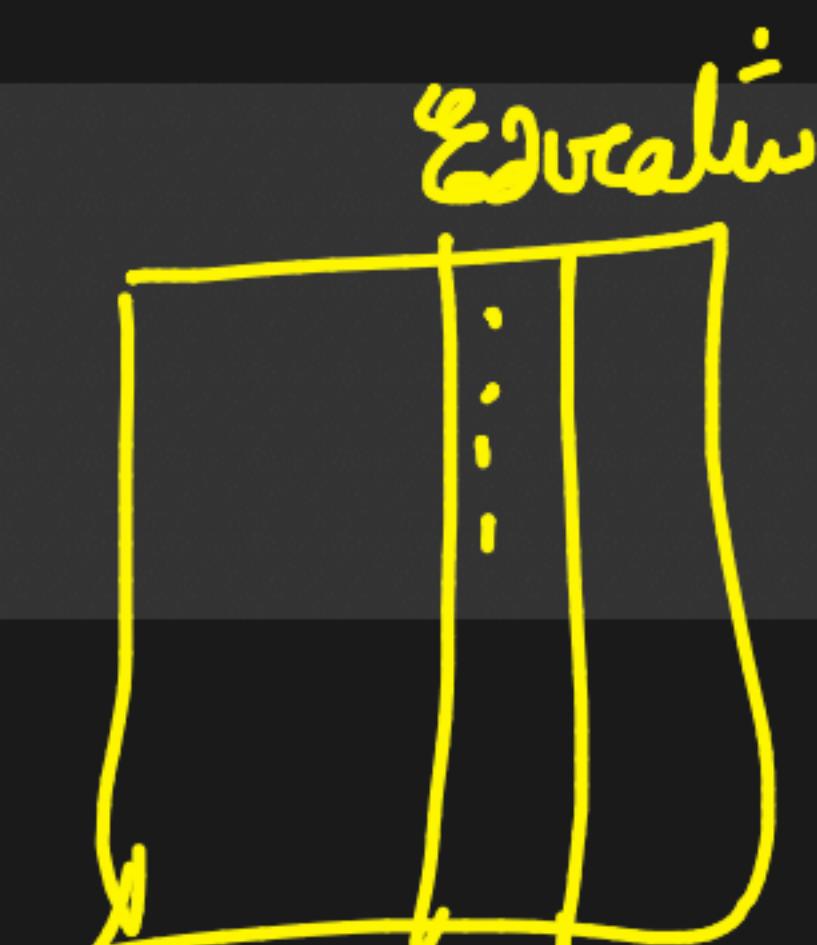
Fish: 54.62897940287769

[10] from scipy import stats

print(stats.median_absolute_deviation(df['MntGoldProds']))

Mode

Education



Probability2.ipynb - Colaboratory x ProbabilityDisb_1.ipynb - Colaboratory x +

colab.research.google.com/drive/1l5T7TVIAASw9Tdl4JxqJxuDqRgFBGOW3#scrollTo=SGqzqcJJbVA9

Update

+ Code + Text

RAM Disk

```
[ ] #mode  
df["Education"].value_counts()
```

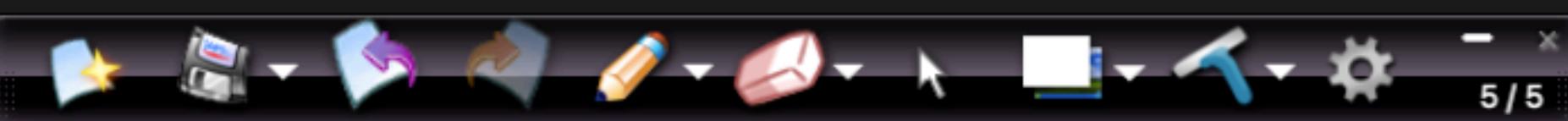
```
{ } Graduation    1127  
      PhD          486  
      Master        370  
      2n Cycle      203  
      Basic         54  
Name: Education, dtype: int64
```

Mode: Most frequently occurring value

```
[ ] #variance  
print("Gold:", df['MntGoldProds'].std())  
print("Fruits:", df['MntFruits'].std())  
print("Sweets:", df['MntSweetProducts'].std())  
print("Wine:", df['MntWines'].std())  
print("Meat:", df['MntMeatProducts'].std())  
print("Fish:", df['MntFishProducts'].std())
```

```
Gold: 52.167438914997064  
Fruits: 39.77343376457871  
Sweets: 41.2804984878548  
Wine: 336.5973926053717  
Meat: 225.71537251175445  
Fish: 54.62897940287769
```

```
[10] from scipy import stats  
print(stats.median_absolute_deviation(df['MntGoldProds']))
```



Income: continuous r.v.

mode: numerical
=

50,000\$

\$ 49-51K (Most)

probability dist (peaks)
= \rightarrow (later)

↓
↓

Ratings: 1
2
3
4
5

{ Ordinal → can be ordered
=

Education:-
(discrete)

high school
10+2
Grad
Masters
:

ML-model
=

Countries :- not ordinal

↳ categorial / discrete

{ mean → average
median → middle value

120

mode → continuous & discrete
↳ popular, shoe-sizes

6, 7, 8, 8.5, ...

average → mean

money spent on wines by Cust-1
to more recent monthly purchases

↓
by Cust-1 with more importance/weightage

$$\bar{x}_w = \frac{w_1x_1 + w_2x_2 + w_3x_3 + \dots + w_{12}x_{12}}{w_1 + w_2 + \dots + w_{12}}$$



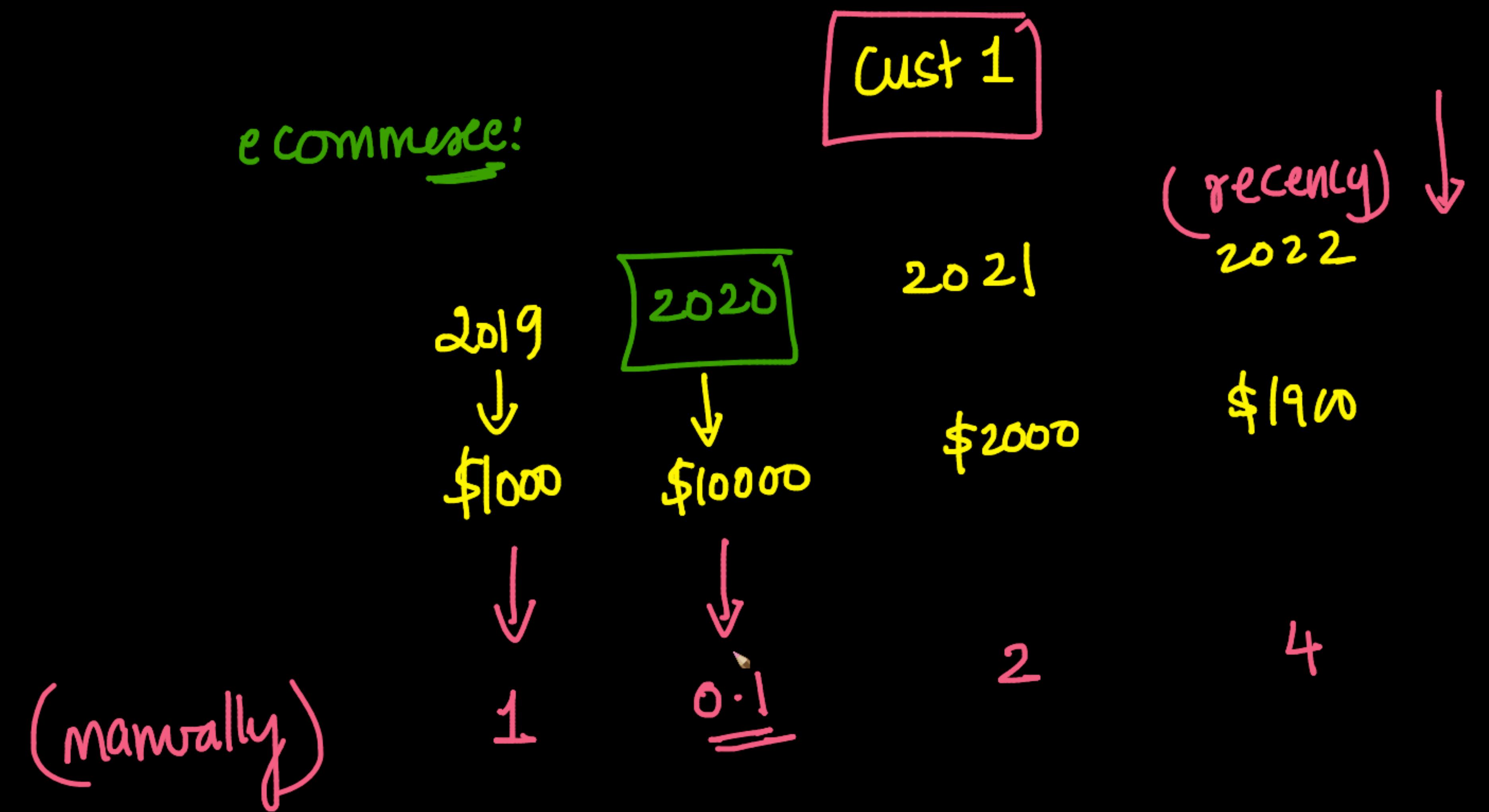
$$\left\{ \begin{array}{l} x_1 = 100; x_2 = 50; \\ x_3 = 150; x_4 = 200 \\ x_5 = 10; x_6 = 20 \end{array} \right.$$

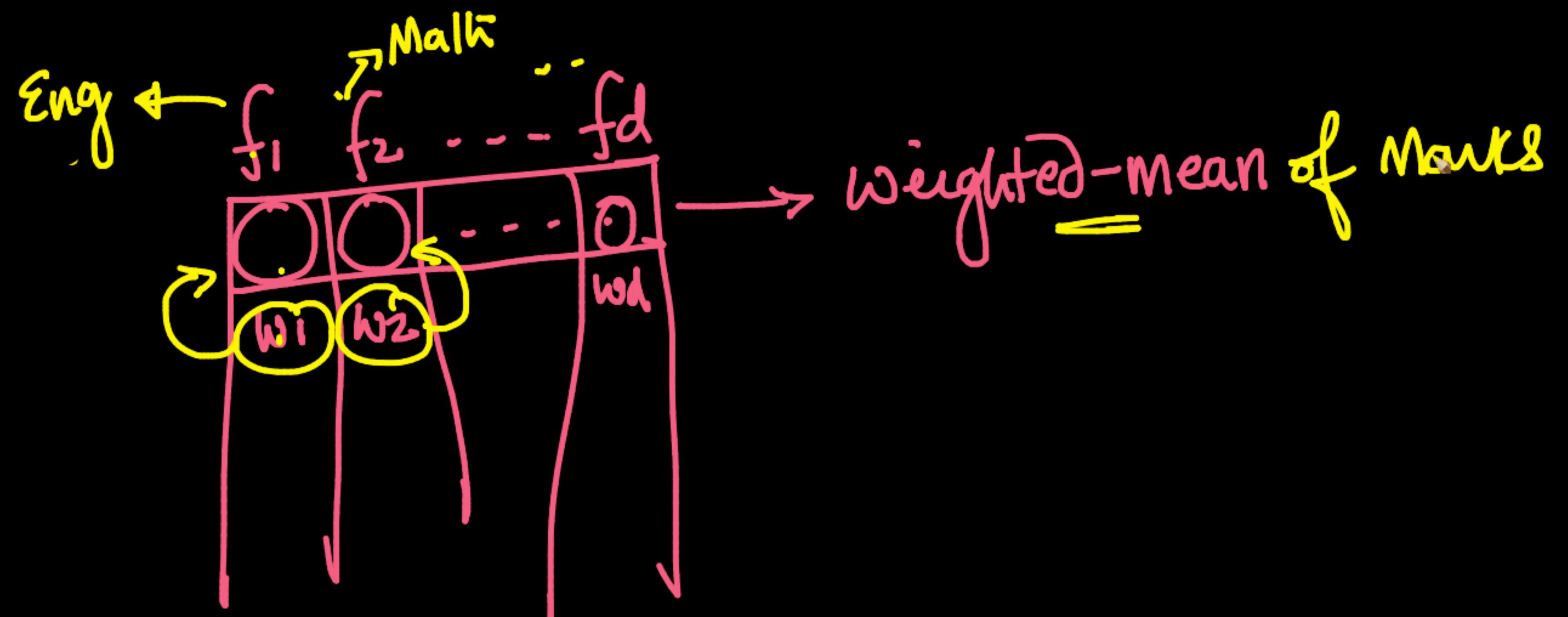
$$\bar{x}_w = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}$$

m = 1000 customers

Q { weighted mean with higher weightage to people
living in France than Germany than Spain . . .

$$\bar{x}_w = \frac{\sum_{i=1}^{1000} w_i x_i}{\sum_{i=1}^{1000} w_i}$$





Amazon:

mean
{ median
w-mean

→ 'central' value

\$10
=

mode → most 'frequent' value



$$\rightarrow \text{Variance} = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}$$

unbiased $\hat{\sigma}^2$

$$\frac{(n-1) \text{ in den}}{den}$$

$$\rightarrow \text{std-dev} = s = \sqrt{\text{Variance}}$$

\hookrightarrow units is same as x_i

$x_1, x_2, \dots, x_k, \dots, x_n$

\bar{x} : mean

→ median vs
mean

→ abs() vs
square

→ Variance = $\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}$

↑
= $\frac{10^8}{n}$

→ std-dev = $s = \sqrt{\text{Variance}}$

↳ Units is same as x_i

$\frac{\$16,000}{\sqrt{10}}$

unbiased $\hat{\sigma}$

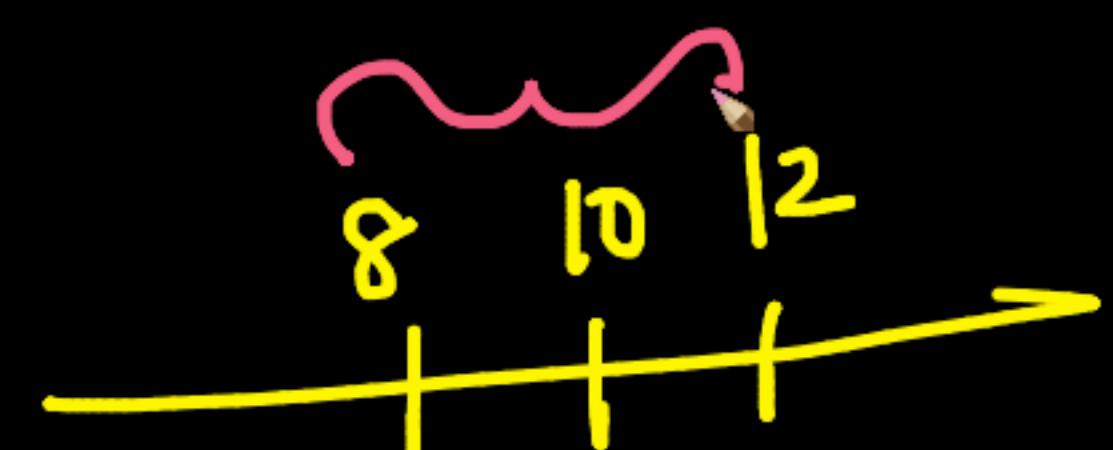
$$\frac{(n-1) \text{ in den}}{n}$$

$x_1, x_2, \dots, x_k, \dots, x_n$ $\$$

\bar{x} : mean $\frac{\$80,000}{10}$

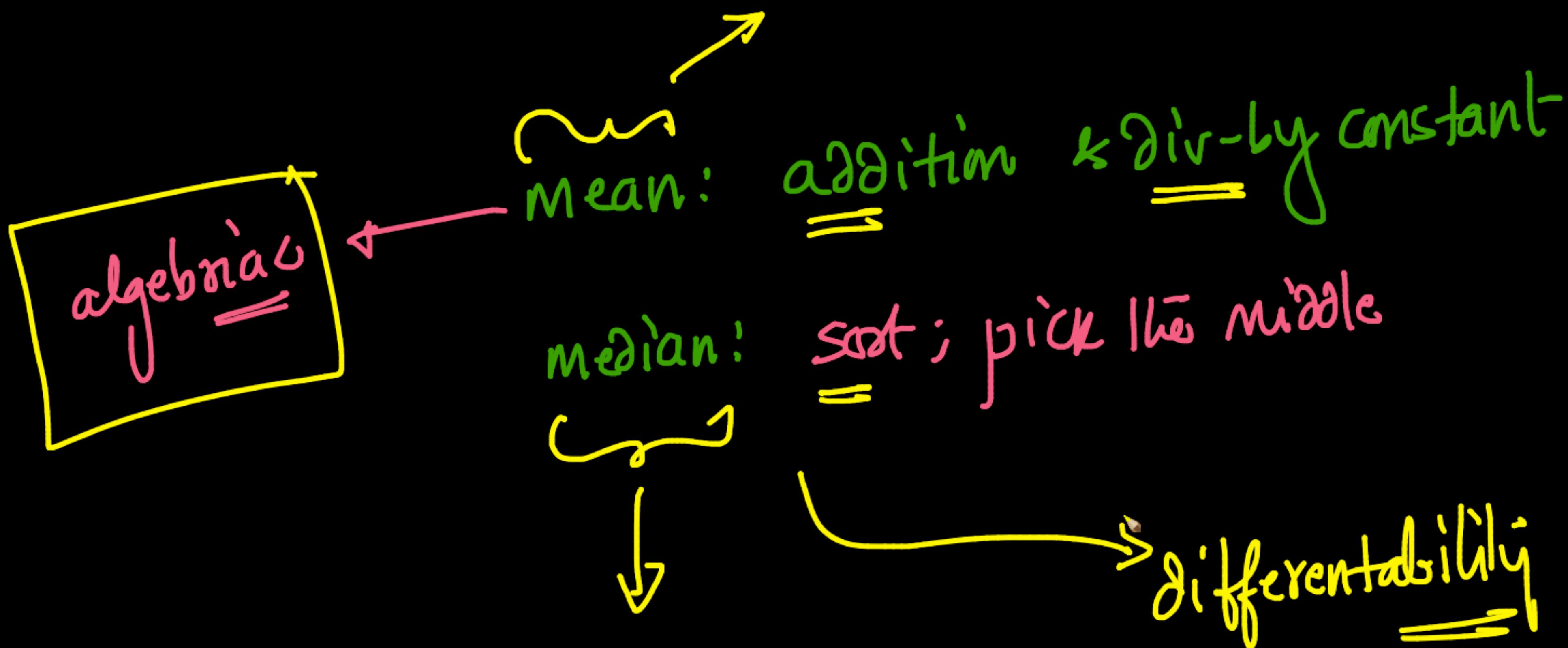
Median vs mean

→ abs() vs square



$$\frac{\sum_{i=1}^n (x_i - \bar{x})}{n}$$

$$=\frac{0}{3}=0$$

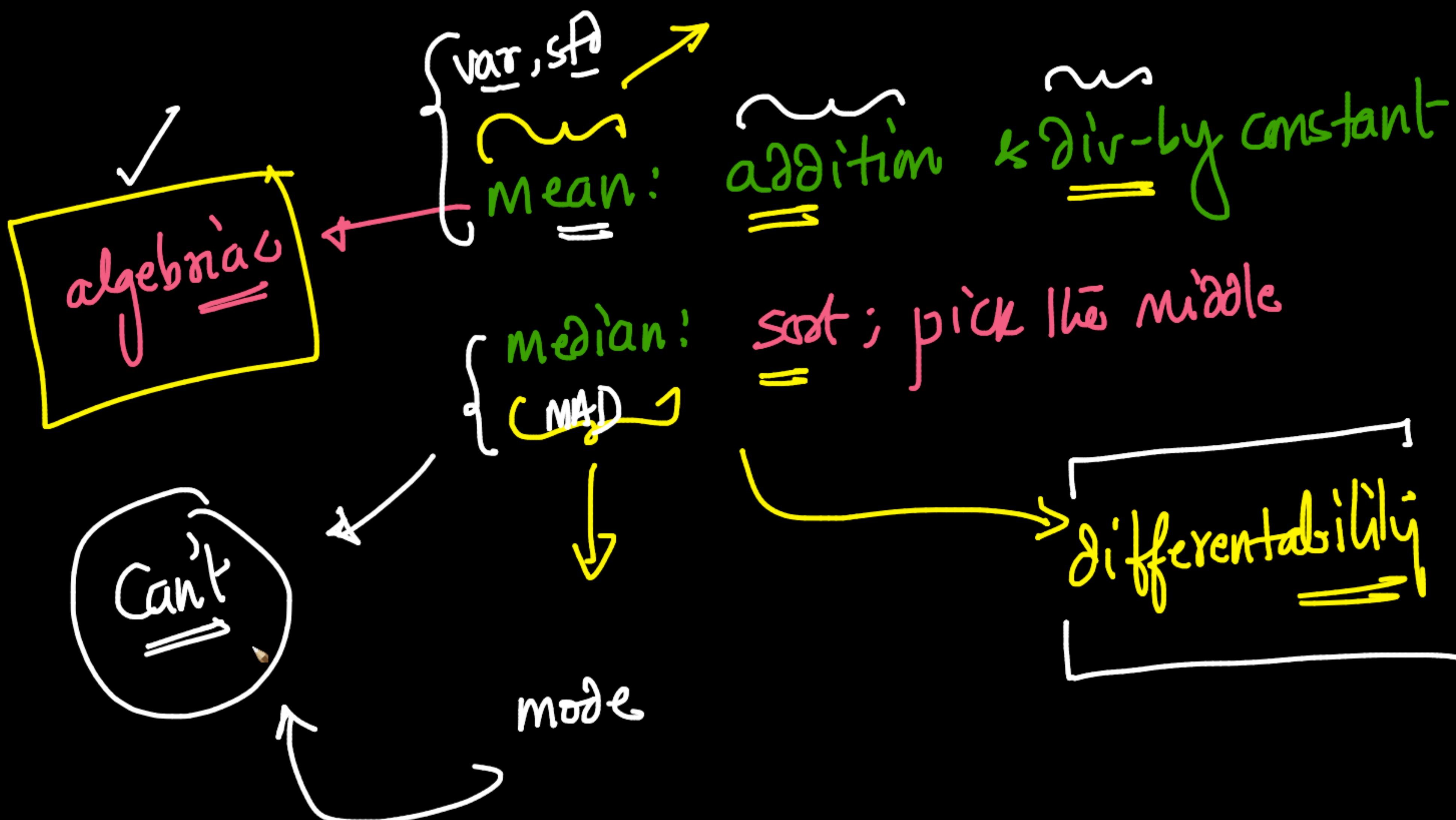


{ Median abs deviation
around the median
(MAD)

$$\text{median} = \text{median} \left(\left| x_i - \underline{x_M} \right| \right)$$

Median abe-dev
around the mean
↓
very seldom used

x_1, x_2, \dots, x_n : $\check{x_M}$ = median
(x_i 's)
Very very outlier proof

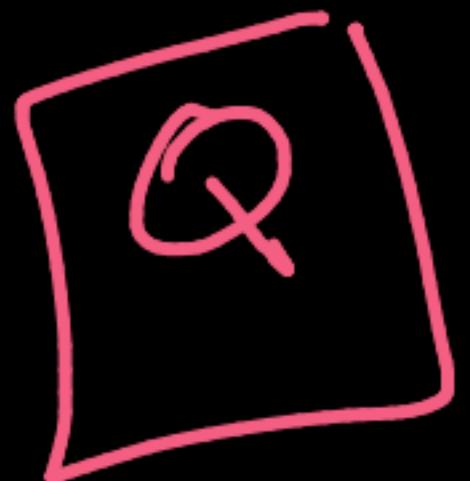


MAD vs STD

↓

robust
estimator

↓
impacted by outliers



when does
weighted mean behave like the mean?



$$\left\{ \bar{x}_w = \frac{\frac{1}{w_1}x_1 + \frac{1}{w_2}x_2 + \dots + \frac{1}{w_n}x_n}{\frac{1}{w_1} + \frac{1}{w_2} + \dots + \frac{1}{w_n}} = \frac{x_1 + x_2 + \dots + x_n}{n} = \bar{x} \right.$$

Percentiles:

x_1, x_2, \dots, x_n

{ one-def } { k^{th} percentile = p_k }

"at least" \underbrace{k} percentage of points $\leq p_k$

e.g.: $3, 2, 1, 5, 6$ $\xrightarrow{\text{sort}}$ $\underbrace{1, 2, 3, 5, 6}_{n=5}$

$p_{50} = 3$ $3/5 = 60\%$. 4 points \leq

n=5

2,3,1,6,5 → find p₁₀ ?

sort

1 2 3 4 5
2,3,5,6

p=10

~ 20% of values are ≤ 1

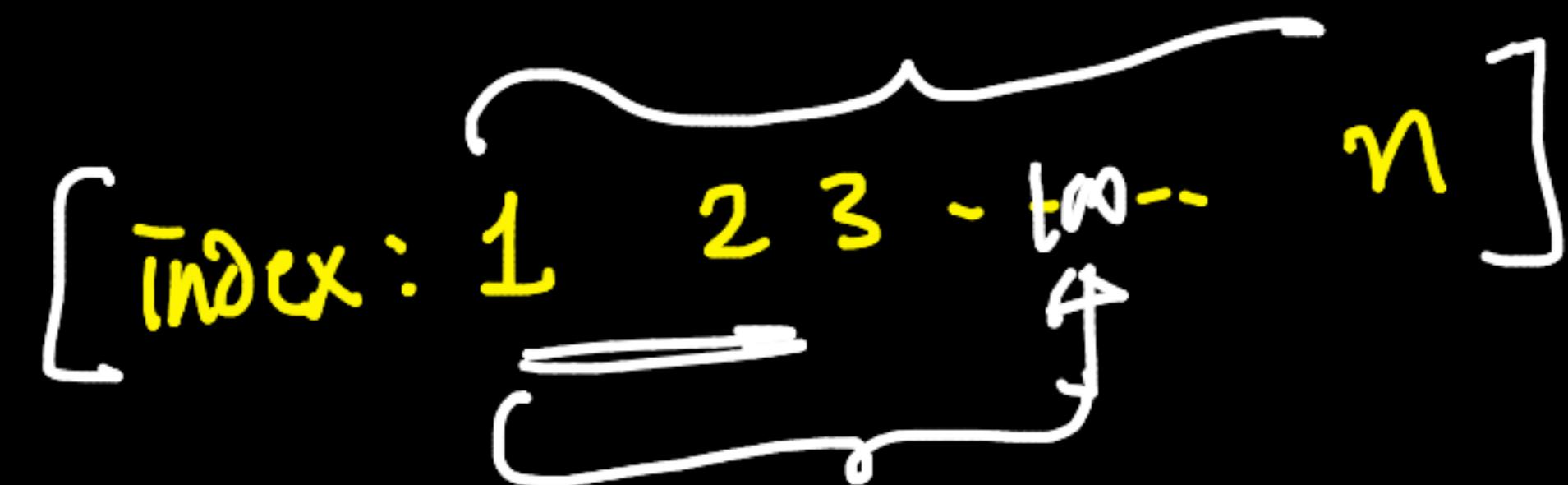
$$P_{10} = 1 \quad \checkmark$$

$$\left[\frac{n \times p}{100} \right] = \left[\frac{50}{100} \right] = \left[\underline{\underline{0.5}} \right] = 1$$

Q

n-values ; p^{th} - percentile value

→ sort them



$$\lceil 0.1 \rceil = 1$$

[index of p^{th} percentile = $\lceil \frac{n \times p}{100} \rceil$ in the sorted array.]

$$\frac{20}{100} \times 1$$

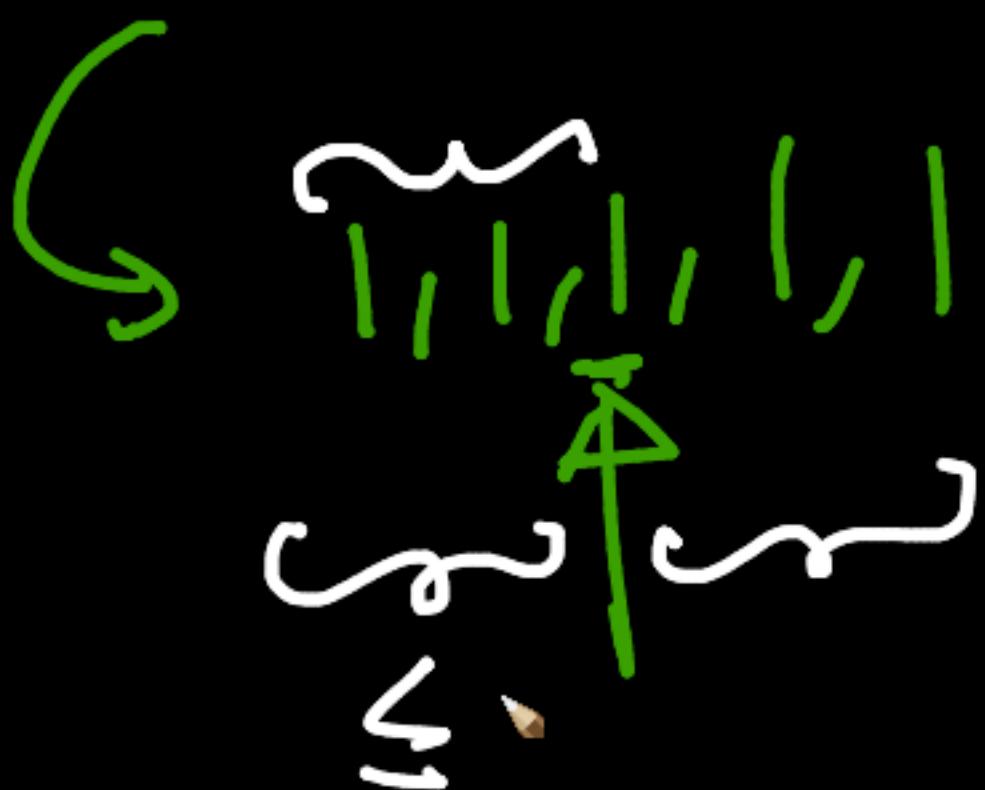
Median = P₅₀

median
≤ 50%
≥ 50%

Q

I, I, I, I, I

P₅₀



Quantiles:

$Q_0: P_0 \rightarrow \min$

$\checkmark \left\{ \begin{array}{l} Q_1 = P_{25} \\ Q_2 = P_{50} \rightarrow \text{median} \\ Q_3 = P_{75} \end{array} \right.$

$Q_4 = P_{100} \rightarrow \max$

$x_1, x_2, \dots, x_n \rightarrow \text{sort}$

$$\left\lceil \frac{25 \times n}{100} \right\rceil = Q_1$$

$$\left\lceil \frac{75 \times n}{100} \right\rceil = Q_3$$

Mut Wines



Uni - Variable

Variable

$P_0 \rightarrow P_{50}$

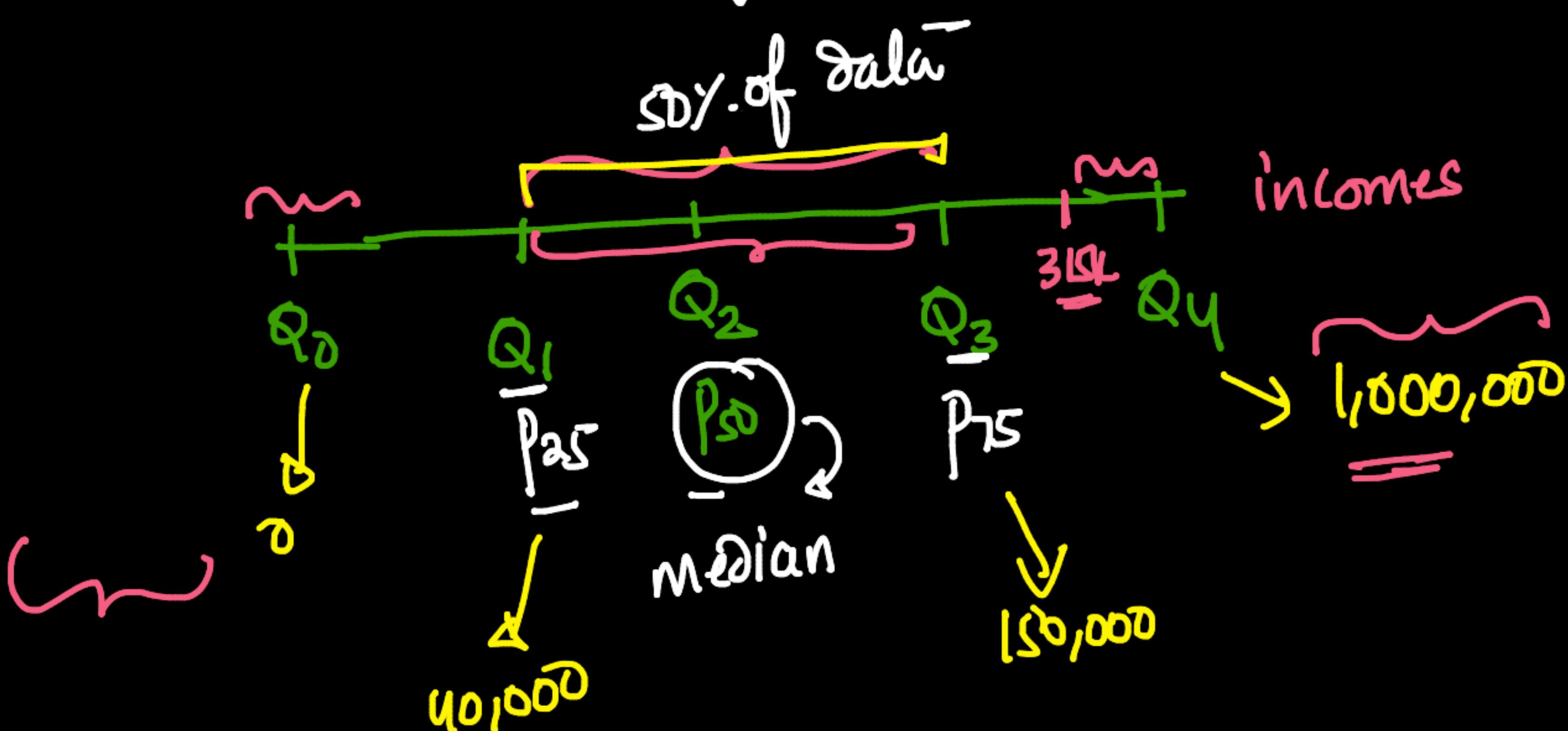
$P_{50} \rightarrow P_{100}$

IQR & outliers

$$\hookrightarrow Q_3 - Q_1$$

$$150,000 - 40,000$$

$$= \underline{\underline{110,000}} \text{ $}$$



I Rule-of-thumb:

detect outliers
subjective
domain specific

(OR)

$$\begin{aligned} & Q_1 - 1.5 \text{ IQR} = 40,000 - 1.5 \times 110,000 = -\text{ve} \\ & Q_3 + 1.5 \text{ IQR} = 150,000 + 1.5 \times 110,000 \\ & = 150k + 165k = \underline{\underline{315,000}} \end{aligned}$$

+ Code + Text

```
#outliers using IQR range
r = 1.5*stats.iqr(df['MntWines'])
lb = np.percentile(df['MntWines'], 25)-r
ub = np.percentile(df['MntWines'], 75)+r
print(lb)
print(ub)
```

✓ 77.0
1225.0

Q₁ - 1.5 IQR

ub
 $Q_3 + 1.5 \text{ IQR}$
PIN=Max
1493
1225

[14] df['MntWines'].max()

1493

[15] sum(df['MntWines'] > ub)

35

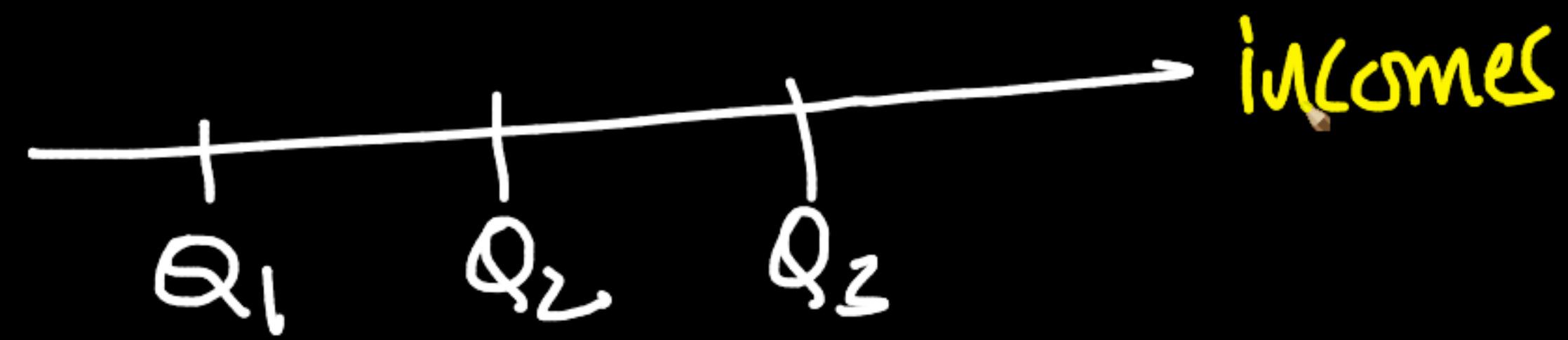
```
df['MntWines'].plot.hist()
#chooses nbins=10
```

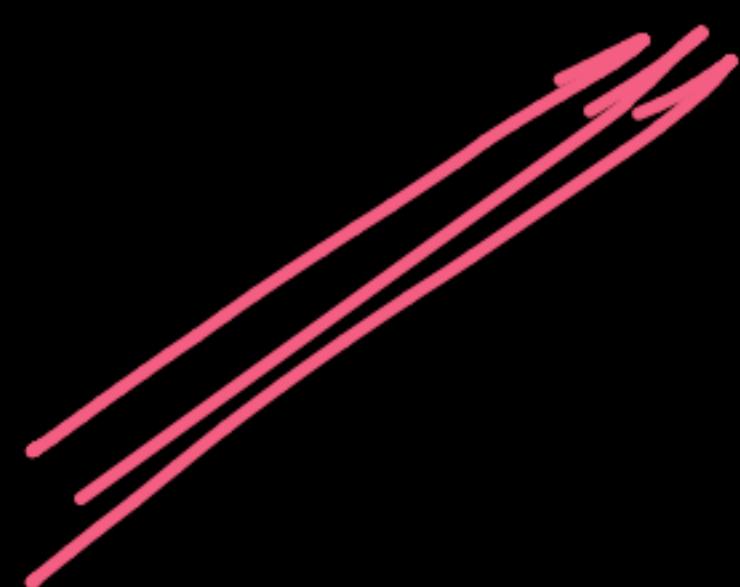
<matplotlib.axes._subplots.AxesSubplot at 0x7feb570fc8d0>



$$\overbrace{Q_3 - Q_2}^{\sim} \leq \overbrace{Q_2 - Q_1}^{\sim} =$$

not always
↑





Recap:

Uni-variate



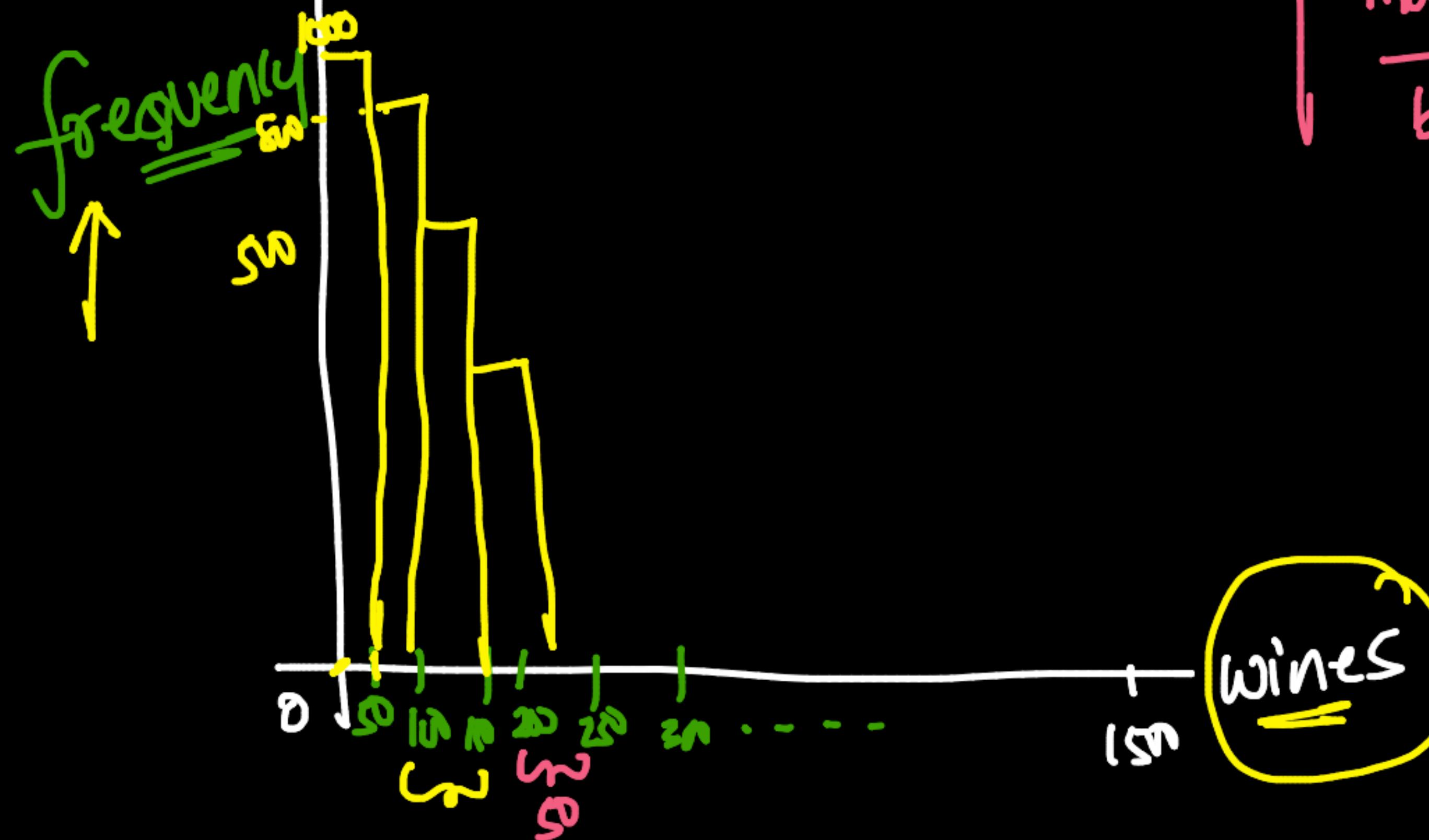
Mean
median
mode
percentiles
Quantiles
IQR

stdDev
MAD

Visualize

Visualize data disb

histogram



$\frac{\text{Max} - \text{Min}}{\text{binSize}}$ bins

min

Max

bins = 30 bins ✓

bin-size = 50

colab.research.google.com/drive/1I5T7TVIAASw9TdI4JxqJxuDgBgFBGOW3#scrollTo=IfNrEJzI2EX

Update

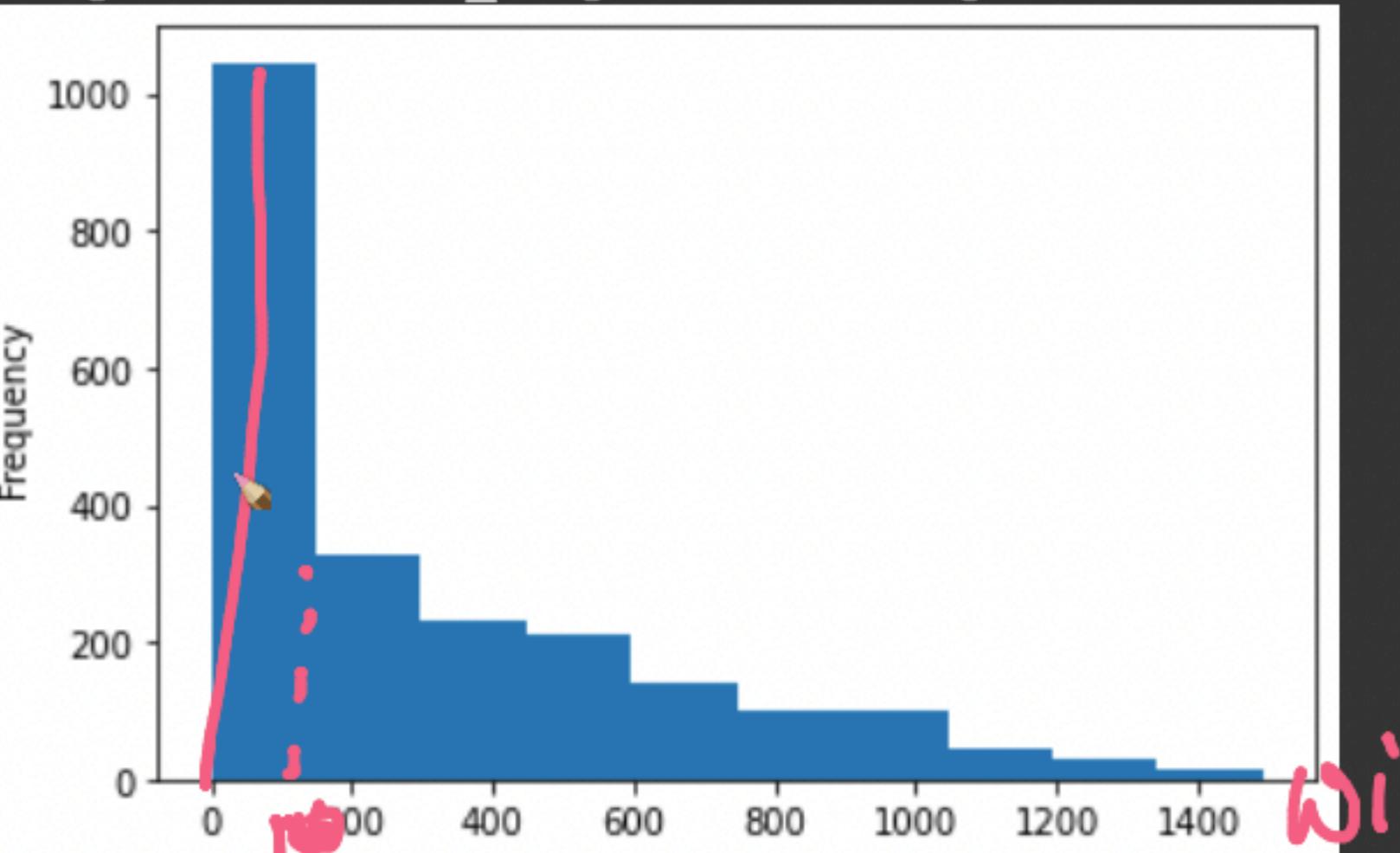
+ Code + Text

RAM Disk

```
16] df[ 'MntWines' ].plot.hist()  
    #chooses nbins=10
```

pandas

```
matplotlib.axes. subplots.AxesSubplot at 0x7feb570fc8d0>
```



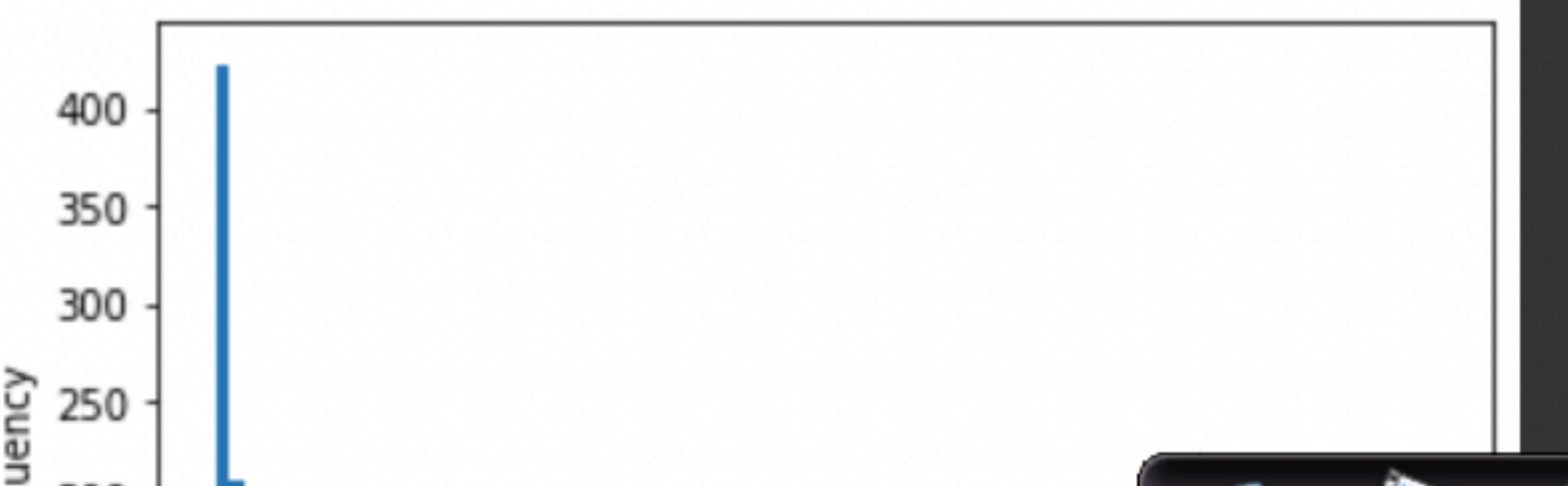
dines

↑ ↓ ⌂ ⏷ ⚙ ⏺ ⏻ ⏸

✓
18

```
df[ 'MntWines' ].plot.hist(bins=100)
```

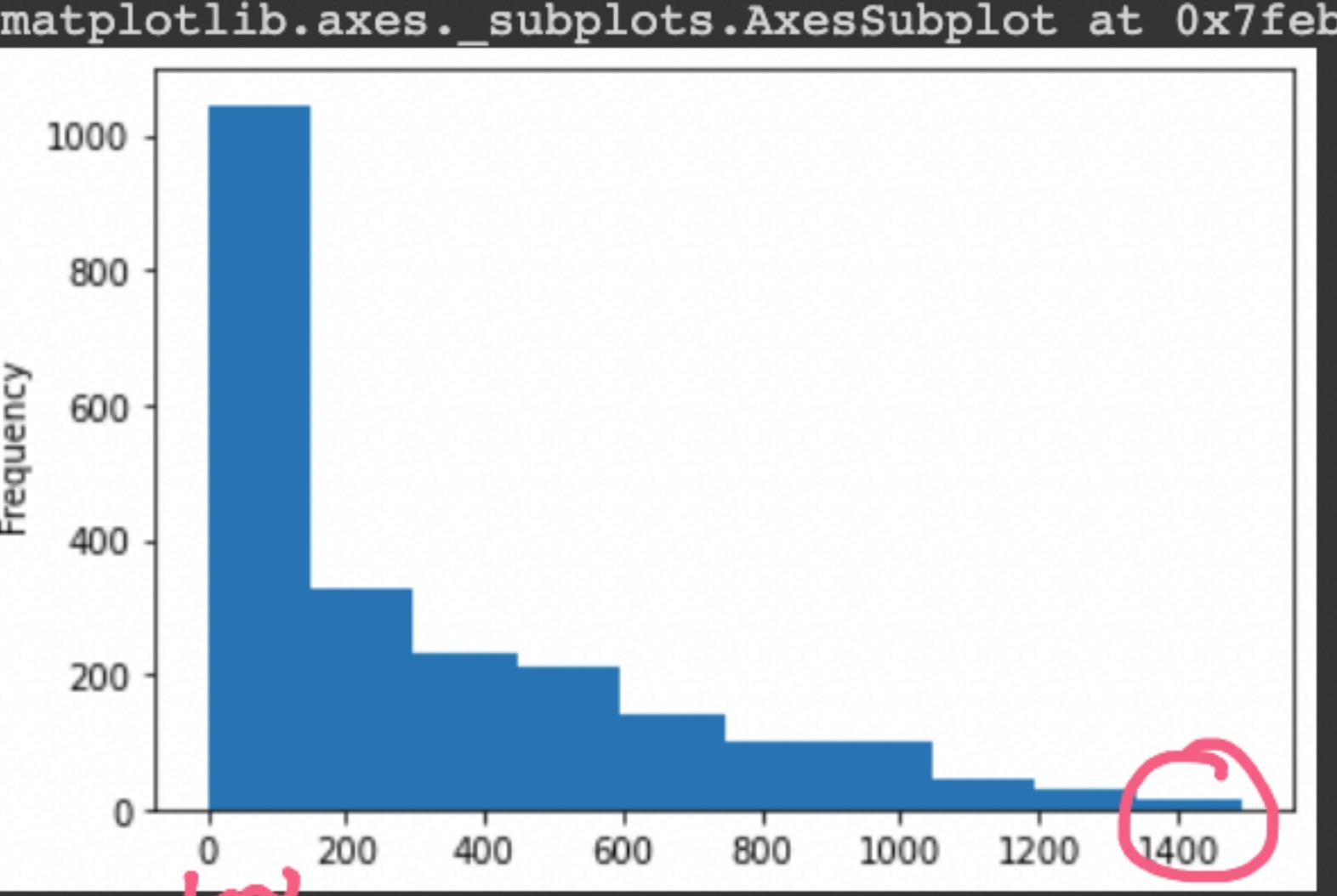
```
[1]: <matplotlib.axes._subplots.AxesSubplot at 0x7feb56ffbf50>
```



+ Code + Text

RAM Disk ✓

```
[16] df['MntWines'].plot.hist()  
#chooses nbins=10
```

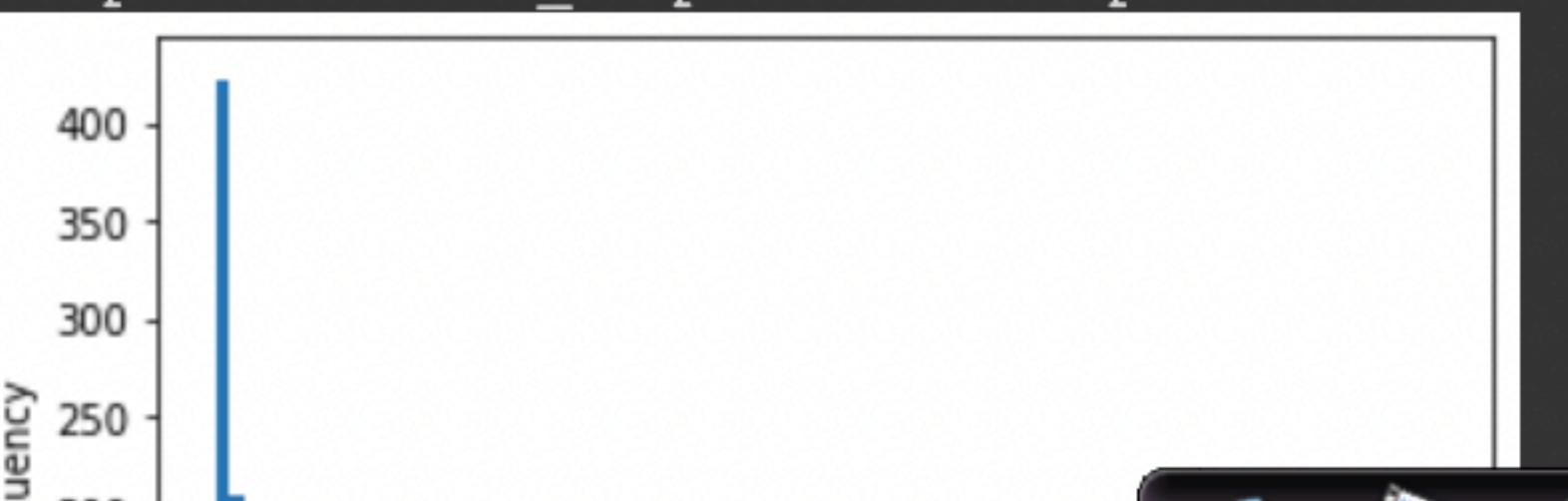


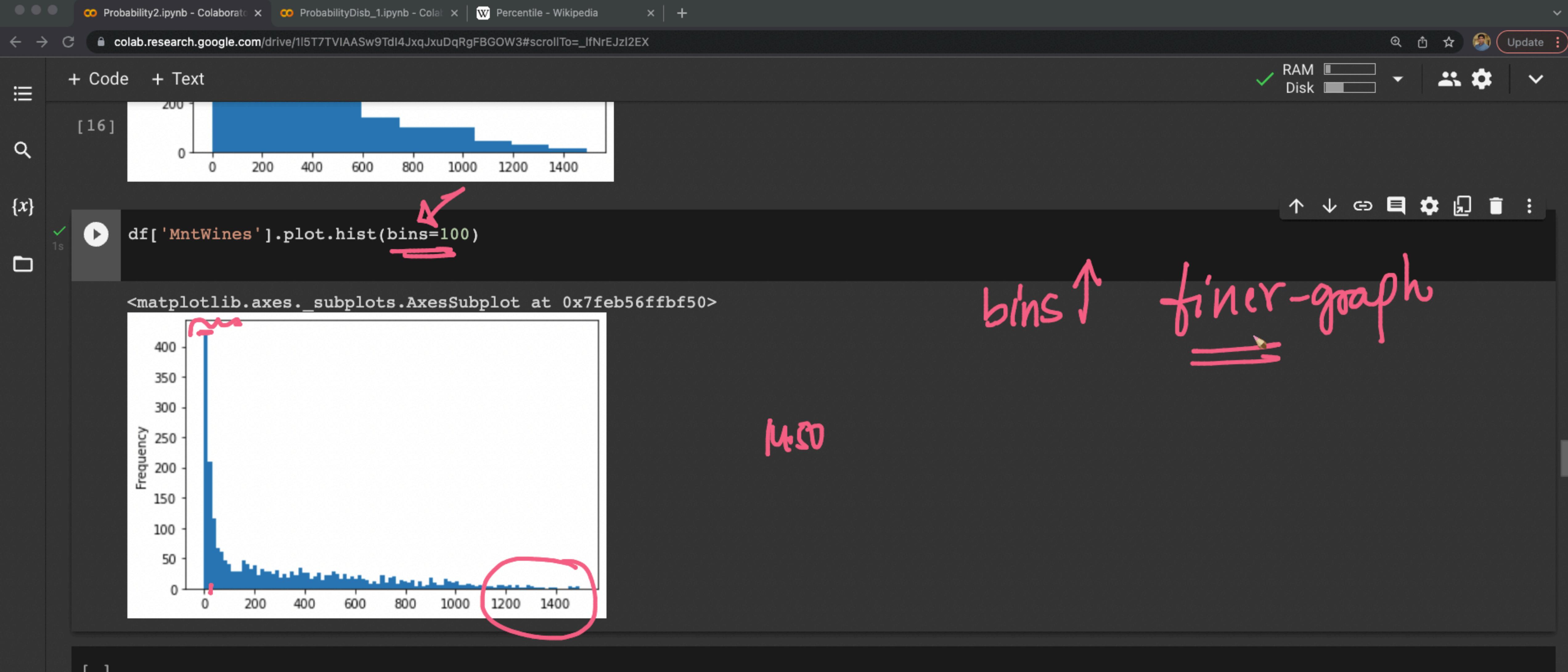
✓

```
df[ 'MntWines' ].plot.hist(bins=100)
```



```
[1]: <matplotlib.axes._subplots.AxesSubplot at 0x7feb56ffbf50>
```





[]

[22] import seaborn as sns
sns.kdeplot(data=df['MntWines'])

`<matplotlib.axes._subplots.AxesSubplot at 0x7feb568a5150>`



Smoothed histogram

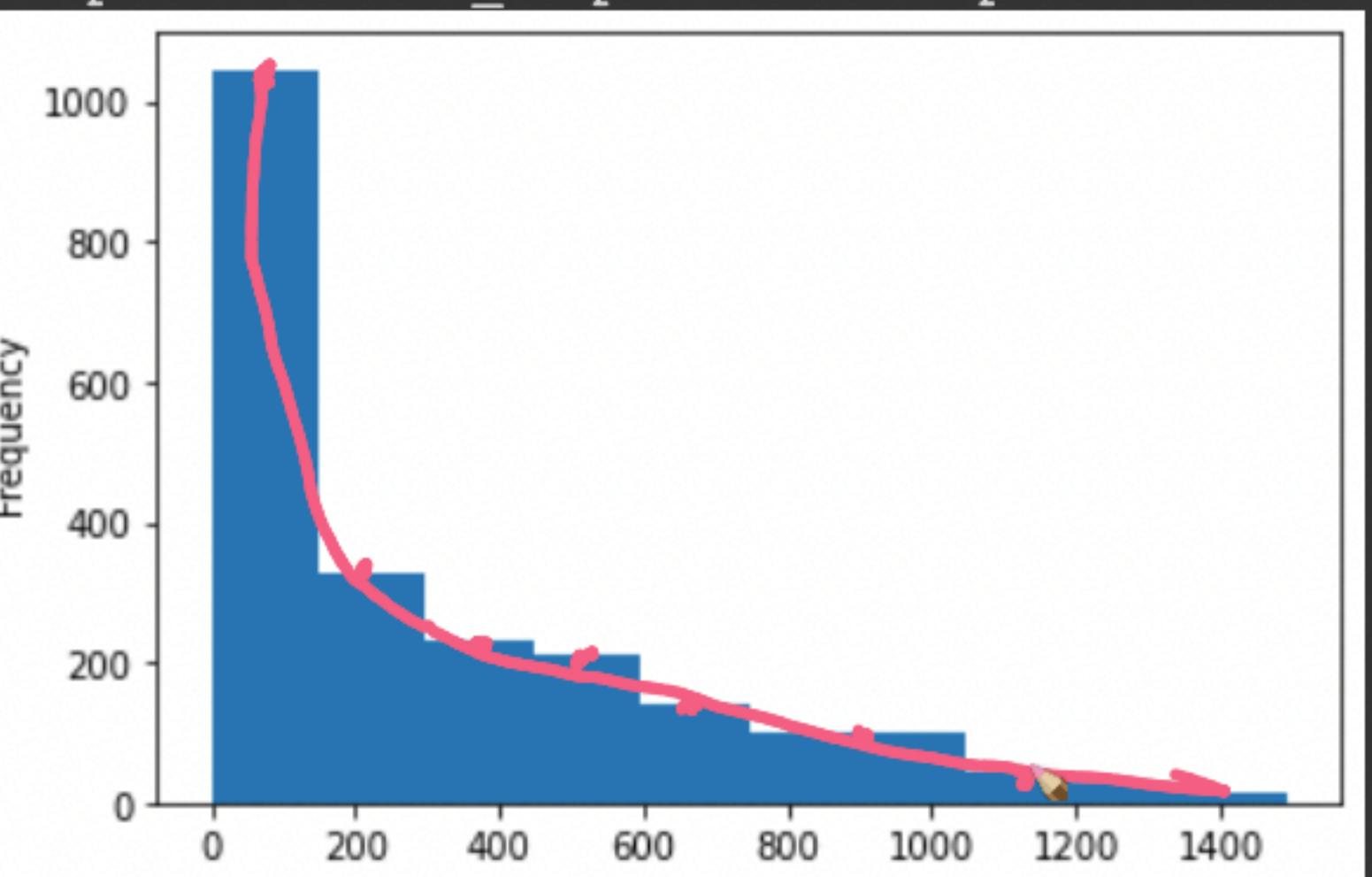
Probability2.ipynb - Colaboratory | ProbabilityDisb_1.ipynb - Colaboratory | Percentile - Wikipedia

Update

+ Code + Text

RAM Disk

▼

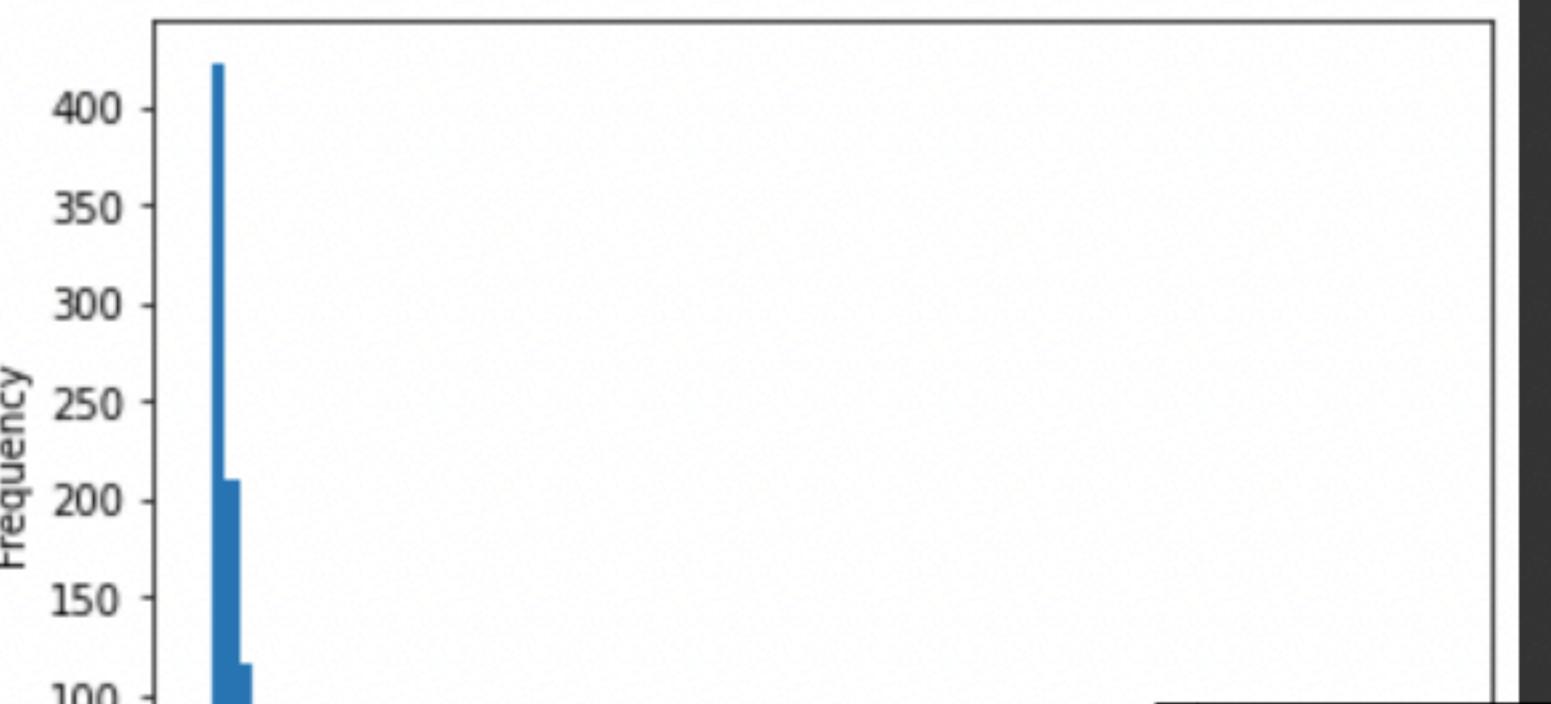


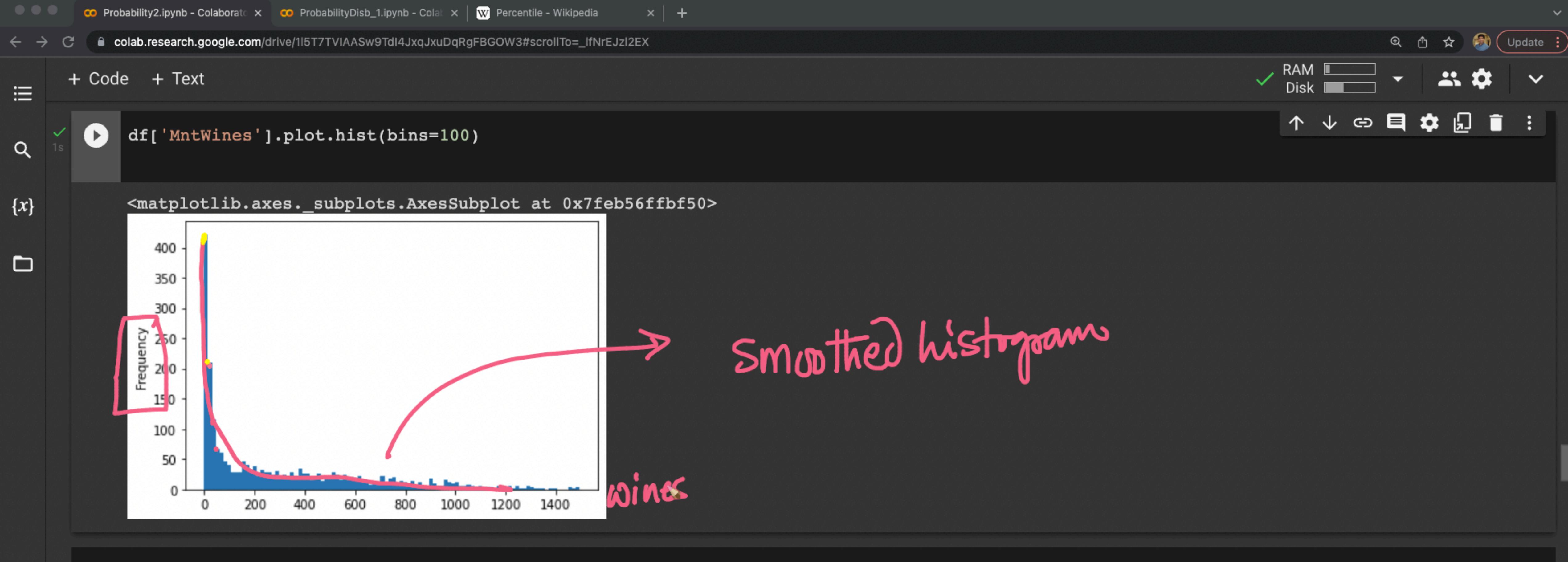
1s

```
f[ 'MntWines' ].plot.hist(bins=100)
```

↑ ↓ ⌂ ⌂ ⌂ ⌂ ⌂ ⌂ ⌂

```
<matplotlib.axes._subplots.AxesSubplot at 0x7feb56ffbf50>
```

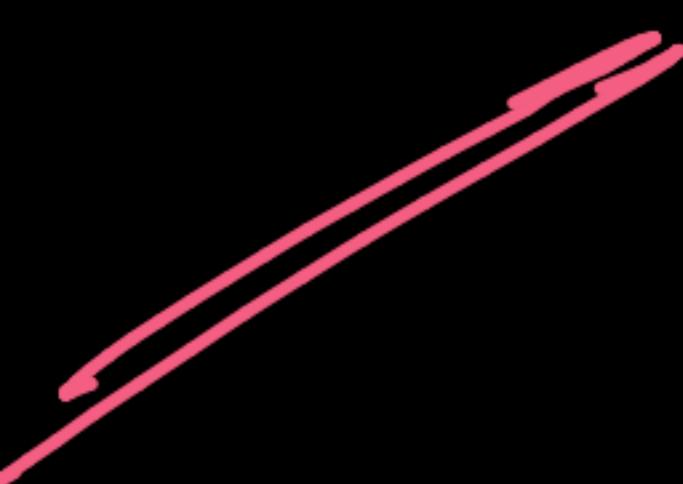




```
[22] import seaborn as sns  
sns.kdeplot(data=df['MntWines'])
```

<matplotlib.axes._subplots.AxesSubplot at 0x7feb568a5150>





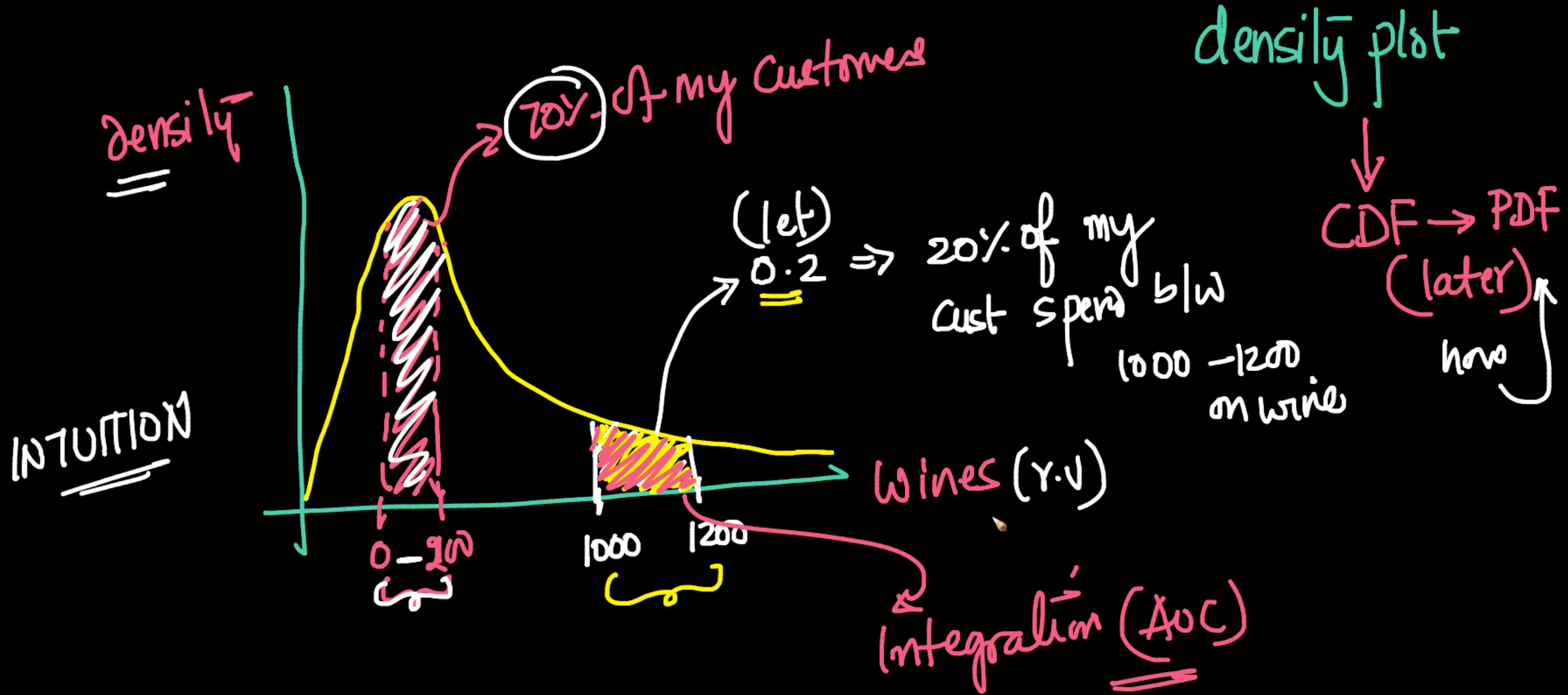
Kernel density plots

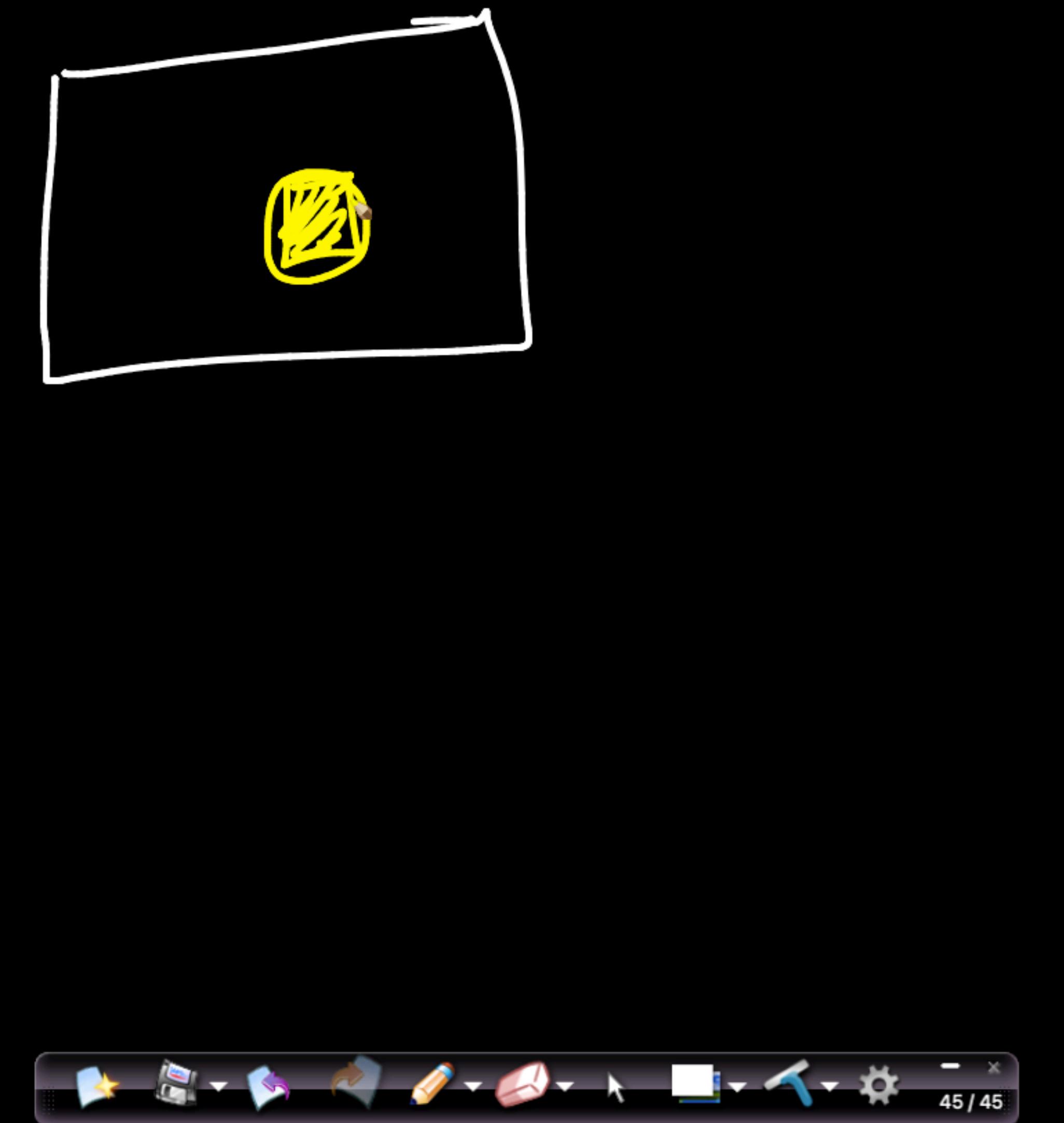
kde



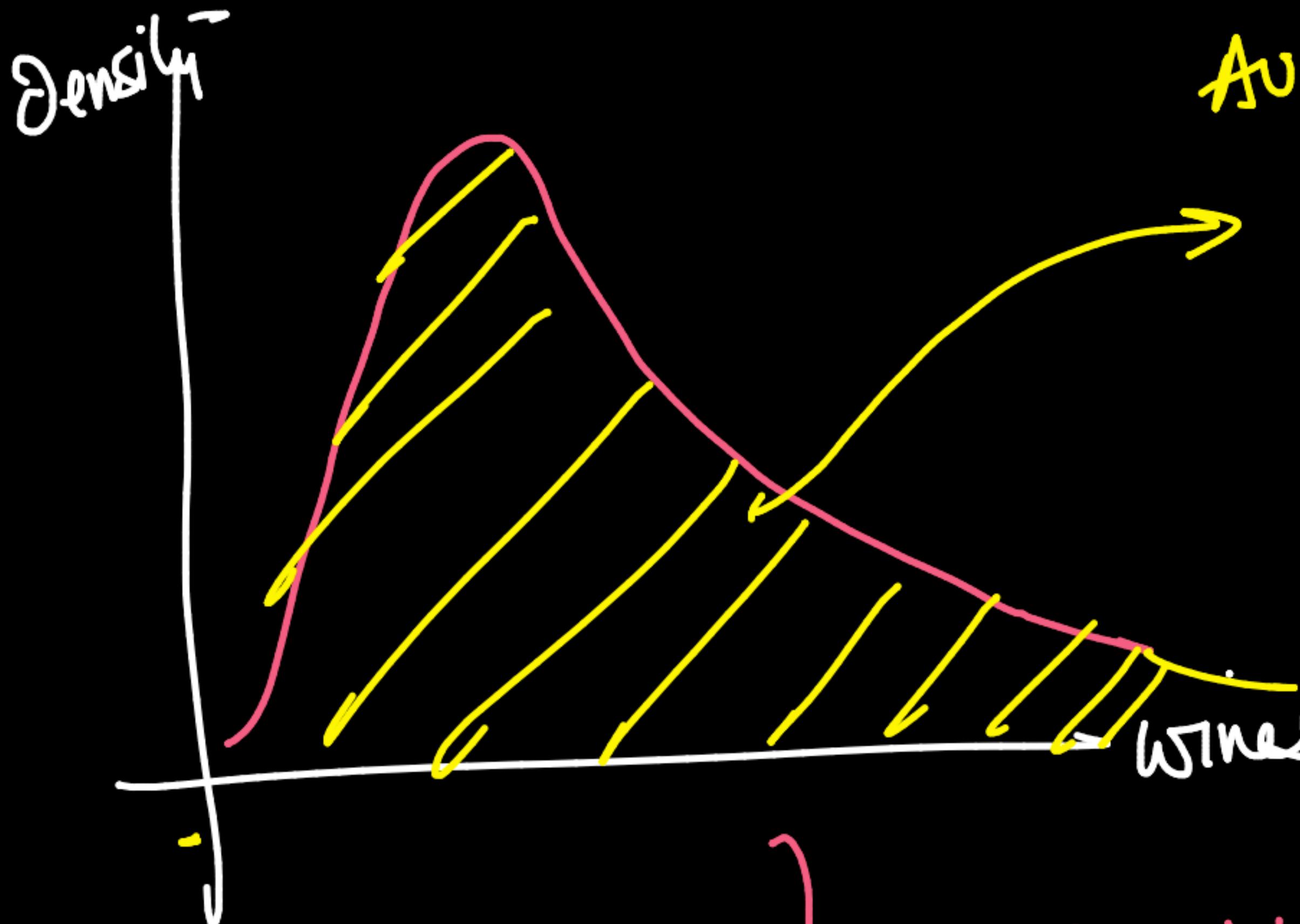
alternative form of smoothing

[
- Kernel
- Gaussian
Kernel]



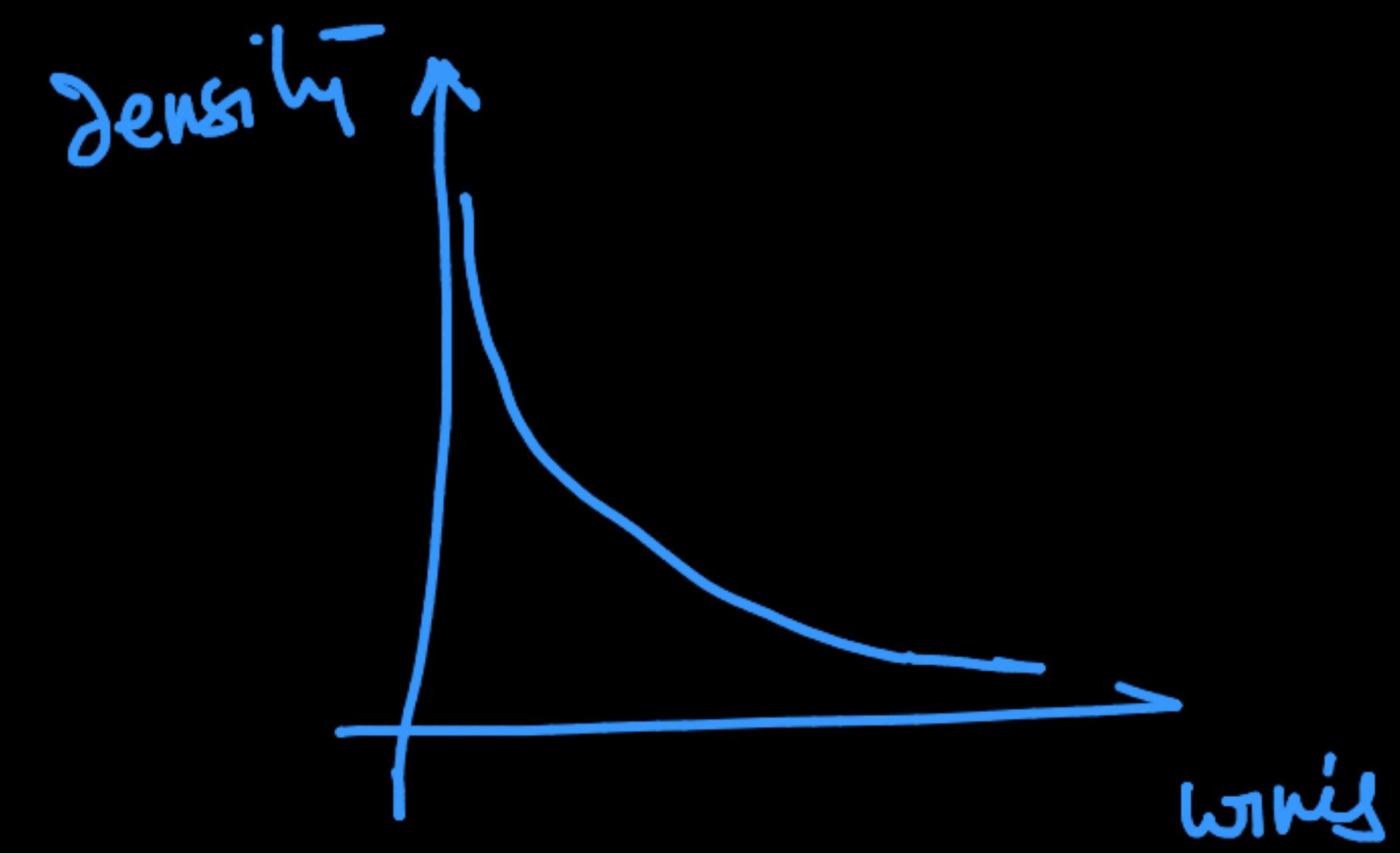
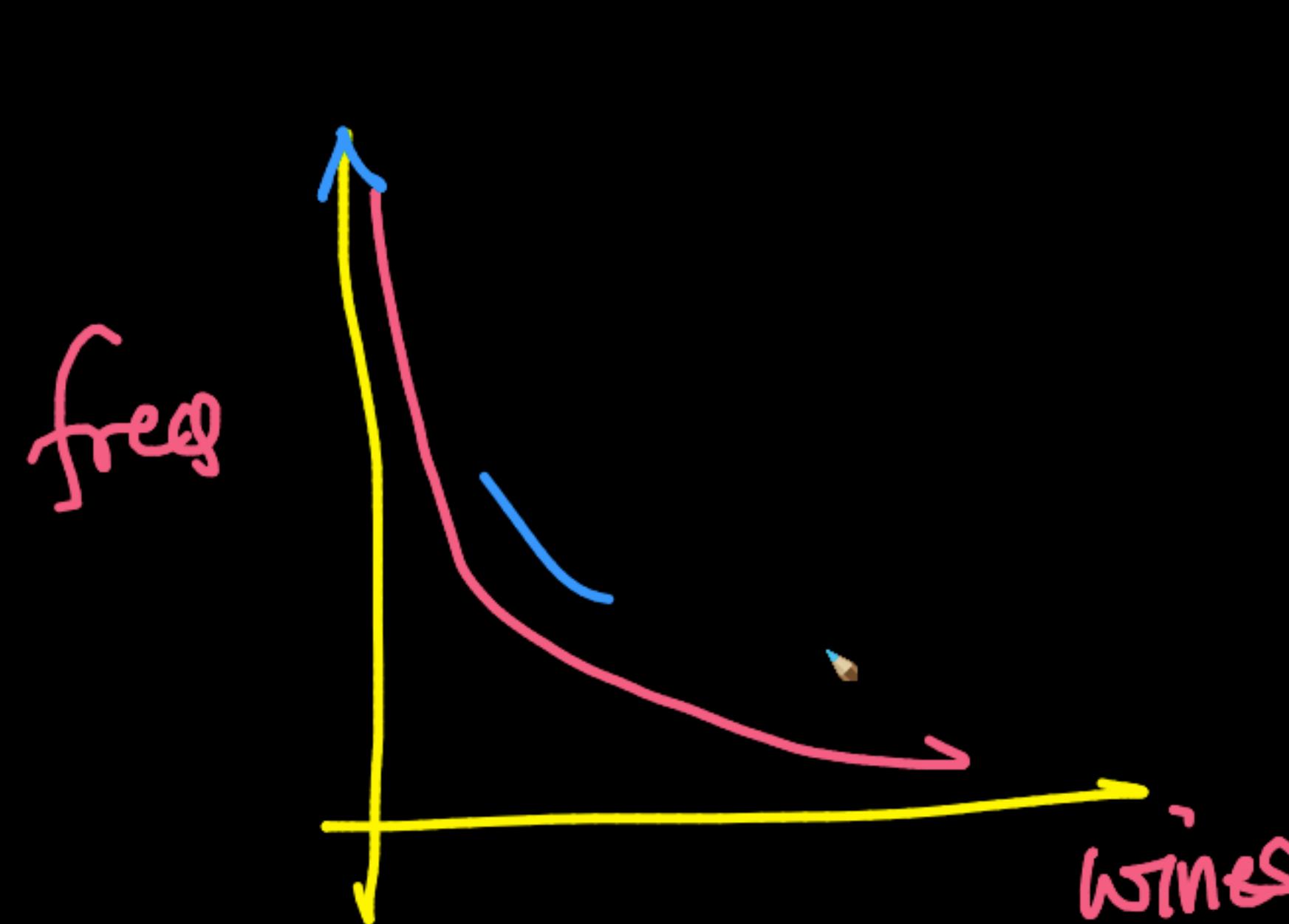


Q



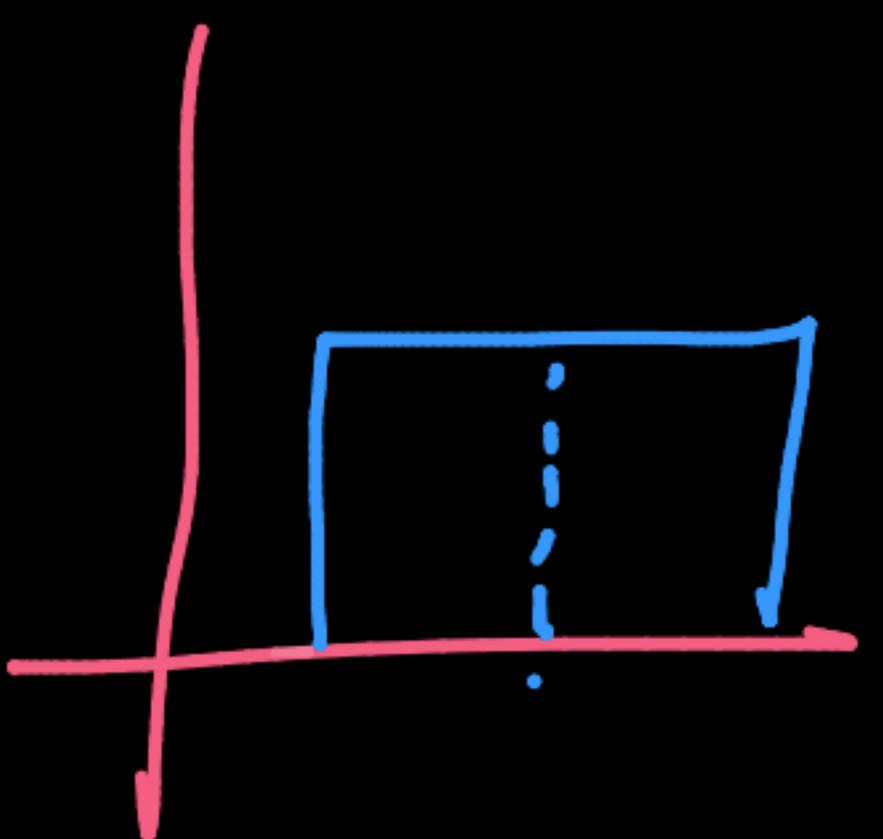
AUG
1 $\Rightarrow 100$ y.

↳ prob. dist. plots (Clatéy)

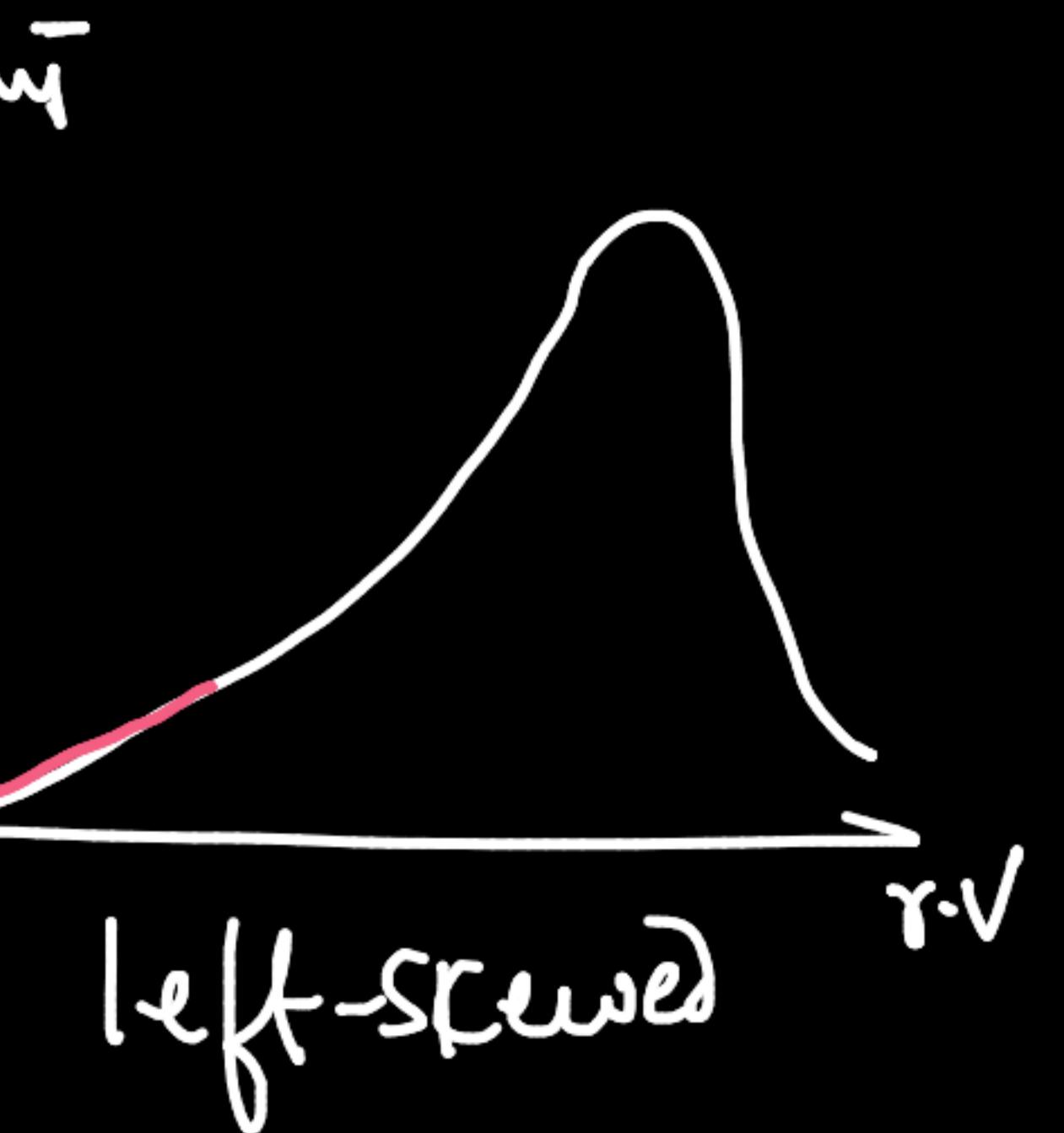
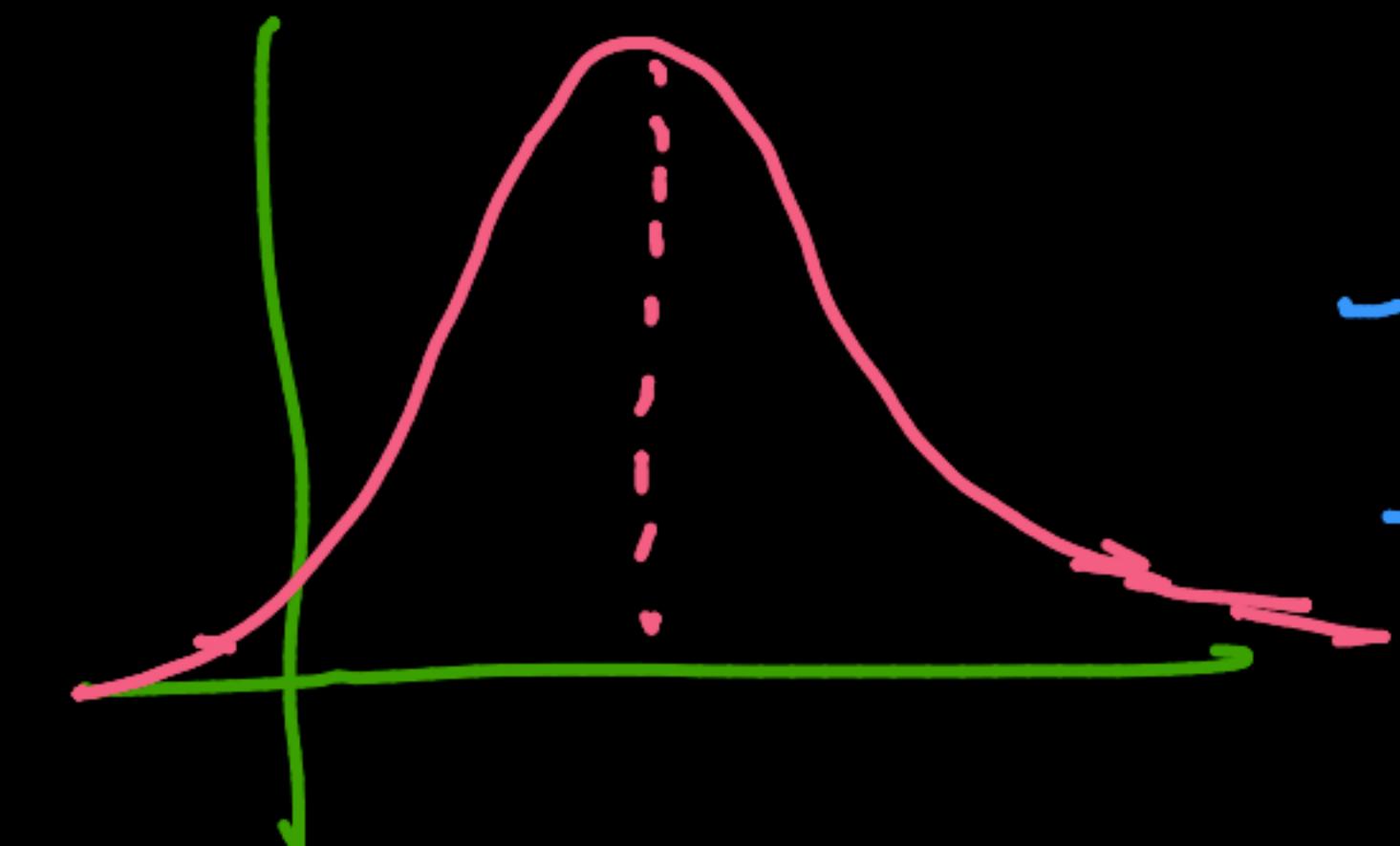




right-skewed



Symmetric



- Tails ?
 - Gaussian
 - log-normal
 - (futnre)

Probability2.ipynb - Colaboratory | ProbabilityDisb_1.ipynb - Colaboratory | Percentile - Wikipedia

colab.research.google.com/drive/1l5T7TVIAASw9Tdl4JxqJxuDqRgFBGOW3#scrollTo=_IfNrEJzl2EX

+ Code + Text

RAM Disk

1s

{x} []

[22] import seaborn as sns
sns.kdeplot(data=df['MntWines'])

<matplotlib.axes._subplots.AxesSubplot at 0x7feb568a5150>

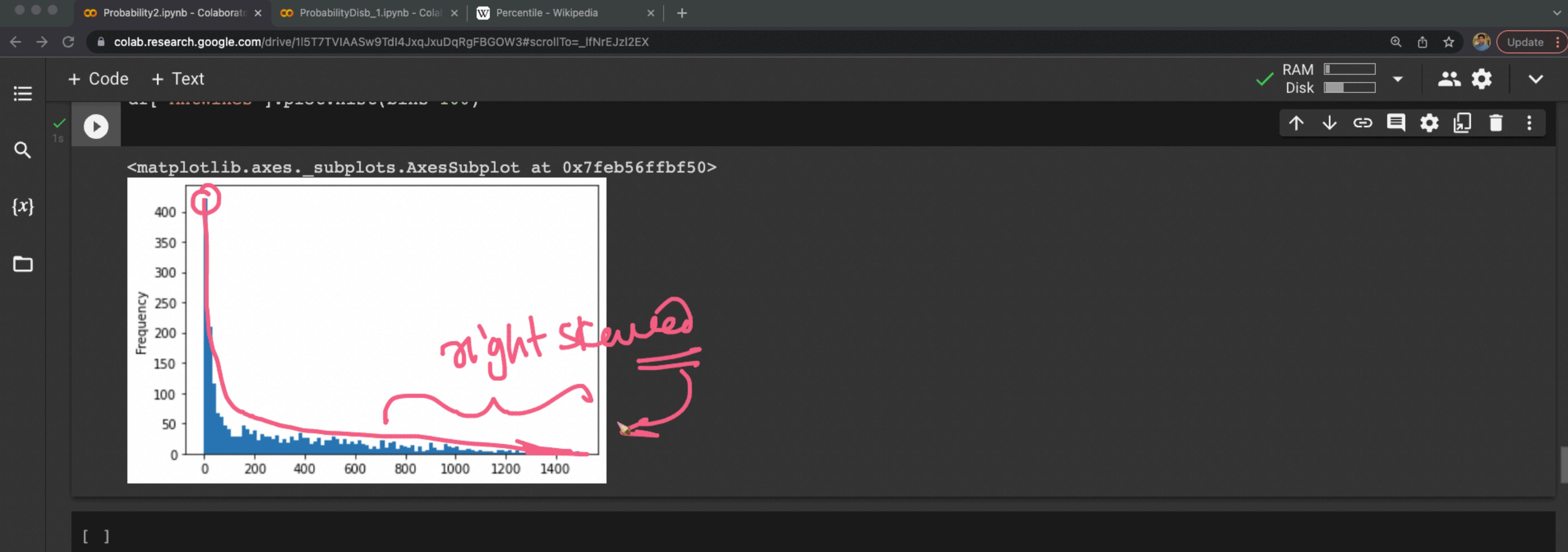
Density

MntWines

[24] df['Income']

	Income
0	\$84,835.00
1	\$57,091.00
2	\$67,267.00
3	\$32,474.00

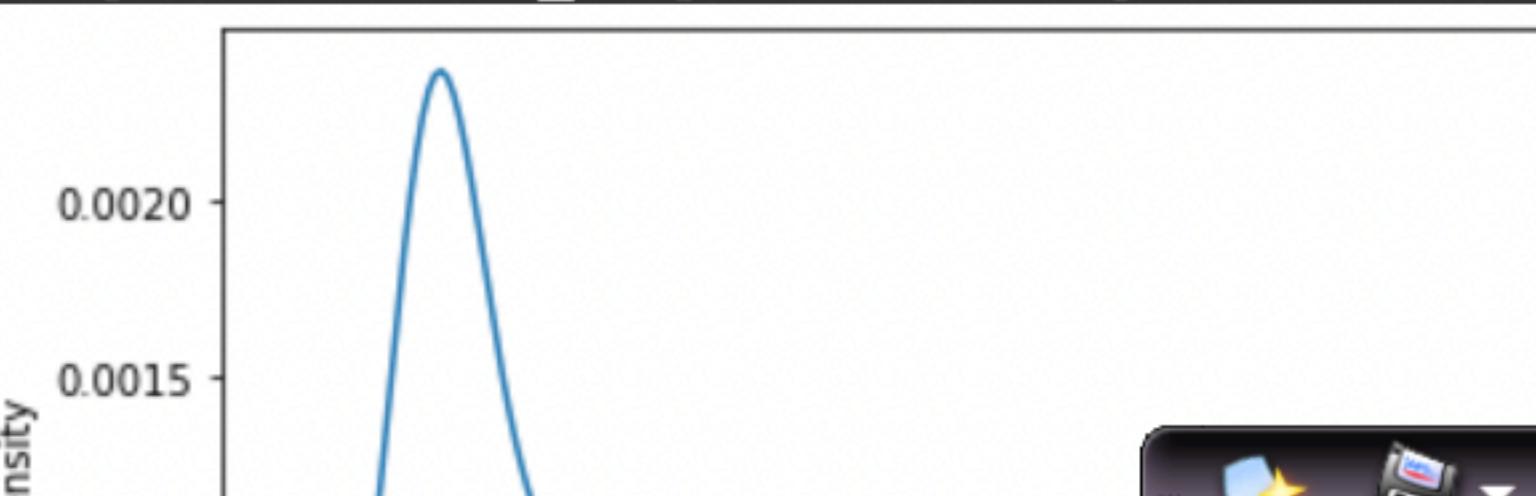
49 / 49

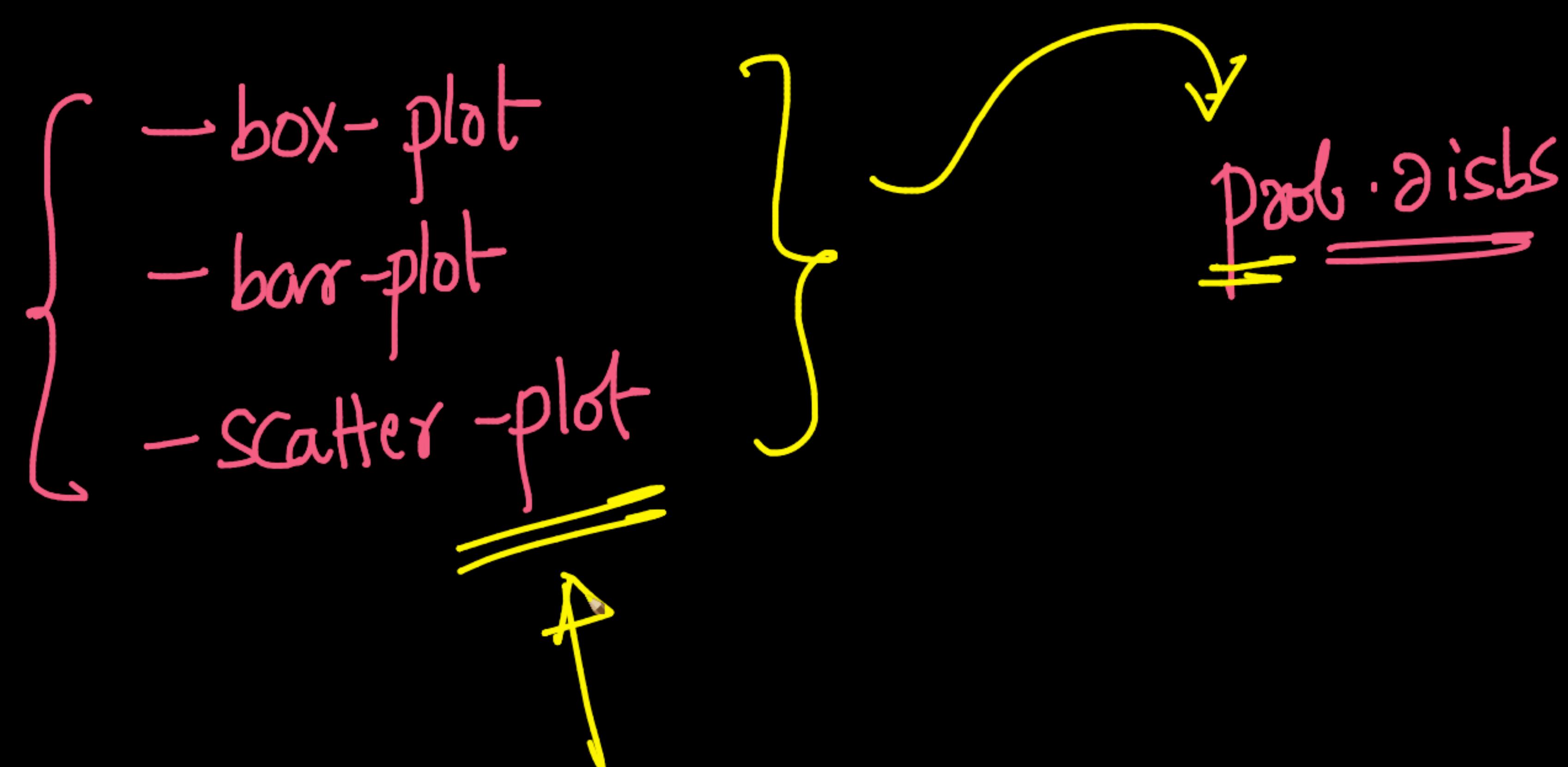


0s

import seaborn as sns
sns.kdeplot(data=df['MntWines'])

<matplotlib.axes._subplots.AxesSubplot at 0x7feb568a5150>





Probability2.ipynb - InterviewBit Software Services x Live | DSML Advanced : Pro x + en.wikipedia.org scaler.com/meetings/i/dsml-advanced-probability-distributions-1-3/live Update :

DSML Advanced : Probability Distributions 1 | Lecture

GEOMRTT

$$\bar{x}_w = \frac{\cancel{2x_1} + \cancel{2x_2} + \dots + \cancel{2x_n}}{2n}$$

You are sharing your screen now

Stop Sharing

Srikanth Varma Chekuri (You) (Screen)

02:44:06

Scott's normal reference rule [edit]

w_j's are equal to

Srikanth Varma Chekuri (You)

Chat

Notify me about Nothing

New a message +

Avijit Swain

Hi. Did you get the chance to look at the question that I asked in the last class. Independent events are not mutually exclusive and vice versa.

Already Answered Answer Now

11:43 pm 0

Harpreet Singh To: Everyone 11:44 pm

yes sir

To: Everyone Enable/Disable Chat

Type message

Doubt Session Ongoing

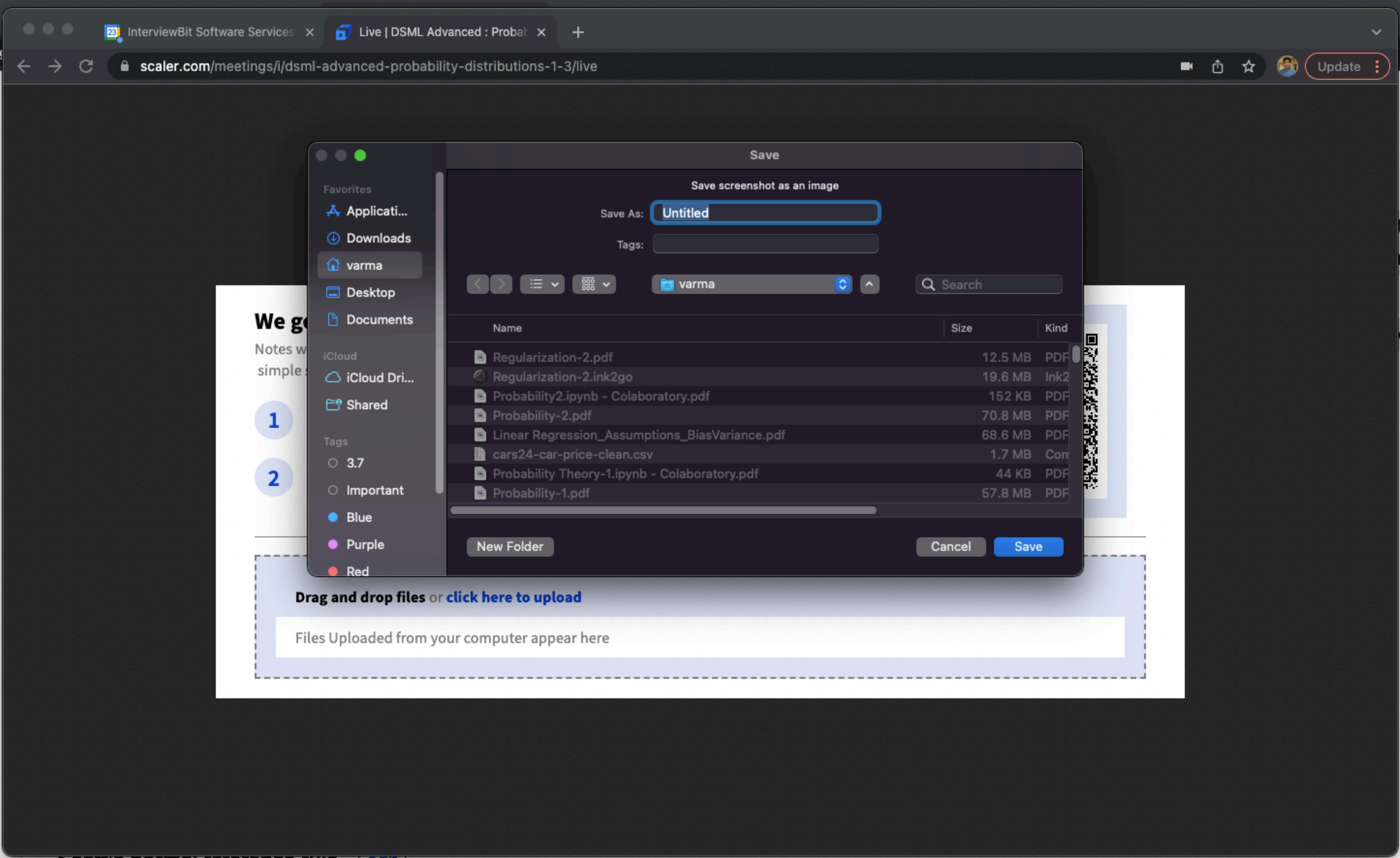
Bin width h is given by

$$h = \frac{3.49\hat{\sigma}}{\sqrt[3]{n}},$$

mutually excl $\not\Rightarrow$ independent

$$\left. \begin{array}{l} p(A) \neq 0 \\ p(B) \neq 0 \end{array} \right\} \checkmark$$

$$\text{Var} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$



Scott's normal reference rule [edit]

Bin width h is given by

$$h = \frac{3.49\hat{\sigma}}{\sqrt[3]{n}},$$

