

Machine Learning

$[5, 1, 1, 0, 1] = \beta$ (in)

Basics

i) Lists

i) Define a list:

$x = [1, 2, 3, 4, 5, 6]$

[] M

11: turtle

ii) To print the length of a list

`print(len(x))`

[1, 3, 8] = 8

() food . 8

iii) slicing

$x[:3]$

Output: [1, 2, 3]

[8, 6, 1] : turtle

iv) Negative-slicing

$x[-2:]$

Output: [5, 6]

[1, 6, 8]

v) Append a list to another list

`x.extend([7, 8])`

Output: [1, 2, 3, 4, 5, 6, 7, 8]

vi) Append an element to the existing list

`x.append(9)`

Output: [1, 2, 3, 4, 5, 6, 7, 8, 9]

vii.) $y = [10, 11, 12]$

list of lists = $[x, y]$

list of lists

$\left[[1, 2, 3, 4, 5, 6, 7, 8, 9], [10, 11, 12] \right]$ ↳ sort (i)

$y[1]$

$[3, 2, 1, 8, 9, 1] = k$

output: 11

viii.) #Sorting a list (Ascending order by default)

$z = [3, 2, 1]$

$z.sort()$

z

output: $[1, 2, 3]$

$[8, 9, 1] : \text{tuple}$

ix.) Sorting a list (reverse = True) (Descending order)

$[10, 9, 8, 7, 6, 5, 4, 3, 2, 1]$

$z.sort(reverse=True)$

$[9, 8, 7] : \text{tuple}$

z

$[3, 2, 1]$

few writings at will is wrong. (v)
 $([8, 7]) : \text{tuple}$

+ List items are ordered, changeable and allow duplicate values.

+ List items are indexed. $[2, 3, 4, 5, 6, 7] : \text{tuple}$

x.) $list1 = list(["abc", 1, "cherry"])$

list1

(P) wrong . b

Outputs the list. Here, the list is created using a list constructor.

$[P, 8, F, 3, 2, H, S, G, 7] : \text{tuple}$

Q) Tuples

A tuple is a collection which is ordered & unchangeable.

tuple = ("apple", "banana")

print(tuple)

"banana" = ["banana"]

i) To print the length of the tuple :

x = (1, 2, 3)

len(x)

Output : 3

ii) y = (4, 5, 6)

y[2]

Output : 5

iii) listoftuples = [x, y]

List of Tuples

Output: [(1, 2, 3), (4, 5, 6)]

iv) (age, income) = "32,12000".split(',')

print(age)

print(income)

Output:

32

12000

3.) Dictionary

Dictionary is a map / hash table in other languages.

i) captains = { }

Captains ["Enterprise"] = "Kirk"

Captains ["Enterprise D"] = "Picard"

Captains ["Deep space Nine"] = "Sisko"

Captains ["Voyager"] = "Janeway"

print (captains ["Voyager"])

Output:

Janeway

ii) print (captains.get ("Enterprise"))

Output:

Kirk

iii) print (captains.get ("NX-01"))

None

iv) for ship in captains:

print (ship + ":" + captains[ship])

v) for ship in Captains:

print (ship)

Output:

Enterprise

ENT D

DSN

Voyager

1

4) Functions

```
# def &gt;(x):  
    return x*x  
print(&gt;(2))
```

Output:

4

```
# def doSomething(p,x):  
    return p(x)  
print(doSomething(say, 4))
```

(k) trying

Output:

16

```
# print(doSomething(lambda x: x*x*x, 3))
```

Output:

27

5) Boolean Expressions

```
# print(1 == 3)  
False
```

if 1 != 3:
 print("How")

```
# print(True or False)  
True
```

elif 1 > 3:
 print("Yikes")

else:

print("All good")

```
# print(1 is 3)  
False
```

5.) Looping

```
# for x in range(10):  
    print(x) Output: 0 1 2 3 4 5 6 7 8 9
```

```
# for x in range(10):  
    if (x is 1):  
        continue  
    if (x > 5):  
        break  
    print(x) Output:  
0  
1  
2  
3  
4  
5
```

```
# x=0  
while (x<10):  
    print(x)  
    x += 1 Output:  
0 1 2 3 4 5 6 7 8 9
```

6.) Pandas:

> Pandas is a python library that makes handling tabular data easier.

> Introduces "data frames" and "series" that allow you to slice and dice rows & columns of information.

7.) Numpy:

> Numpy arrays are multi-dimensional array objects. It's easy to create a pandas dataframe from a numpy array and pandas Dataframes can be cast as numpy arrays.

* Scikit-Learn:

> Load, clean and manipulate your input data using pandas. Then convert pandas Dataframe into a Numpy array as it's being passed into some Scikit-Learn function. That conversion can happen automatically.

Exercise 1

import numpy as np
import pandas as pd

df = pd.read_csv("PastHires.csv")

df.head() # prints top 5 rows of the dataframe

df.head(10) # prints the top 10 rows

df.tail(4) # prints the last 4 rows

df.shape # (prints (rows, cols))

df.size # (prints rows * cols)

len(df) # prints the no. of rows

df.columns # prints the column names of the dataframe

df['Level of Education']

prints all the data present at that particular column.

df['Employed?'][:3]

prints (0-3) rows data of given column.

df['Hired'][4]

prints date present at 4th row of column Hired.

df[['Years Exp', 'Hired']]

prints data of both the columns given

df[['Years Exp', 'Hired']][,:3]

prints data of both the columns [0-3 rows]

df.sort_values(['Years Exp'])

prints data present at Years Exp column by sorting

df.sort_values(['Years Exp'], ascending=False)

prints data present at years Exp column sorting in descending.

Edcnt = df[['Level of Education']].value_counts()

Outputs each unique value present at column given and the frequency of that value.

Edcnt.plot(kind='bar')

outputs a bar graph of the given dataframe

prevEmplir = df[['prev', 'Hired']][5:10]

prevEmplir.plot(kind='hist')

Outputs a histogram of the given dataframe.

Types of Data

- i) Numerical
- ii) Categorical
- iii) Ordinal

Numerical

* Represents some sort of quantitative measurement

Ex: Heights of people, page load times, stock prices etc.

• Discrete Data

Integer based; often counts of some event.

- How many purchases did a customer make in a year?
- How many times did I flip "heads"?

• Continuous Data

Has an infinite number of possible values

- How much time did it take for a user to check out?
- How much rain fell on a given day?

Categorical

* Qualitative data that has no inherent mathematical meaning

- Gender, yes/no (binary data), Race, state of Residence, product category, political party etc..

* You can assign numbers to categories in order to represent them more compactly, but the numbers don't have mathematical meaning.

Ordinal

- * A mixture of numerical and categorical
- * Categorical data that has mathematical meaning.

Ex: Movie ratings on a 1-5 scale.

Ratings must be 1, 2, 3, 4 or 5

But these values have mathematical meaning
1 means it's a worse movie than a 2.

Quiz!!

* How much gas is in your gas tank?

> Numerical

* A rating of your overall health (choices - 1, 2, 3, 4) corresponding to ("poor", "moderate", "good", "excellent")

> Categorical

* The race of your classmates
categorical

* Ages in years

> Numerical

* Money spent in a store buying

> Numerical

Mean

- > AKA Average
- > Sum (no. of samples) / number of values
- > Ex:

No. of children in each house on the street:

$$\text{Mean} = (0+2+3+2+1+0+0+2+0) / 9 = 1.11$$

Median

- > Sort the values and take the value at the midpoint

Ex: 0, 2, 3, 2, 1, 0, 0, 2, 0

Sort: 0, 0, 0, 0, 1, 2, 2, 2, 3
↑
Median = 1

* If you have an even no. of samples, take the avg of the two in the middle.

- > Median is less susceptible to outliers than the mean.

Ex: the mean household income in the US is \$ 72,641
But the median is only \$ 51,939 - because the mean is skewed by a handful of billionaires.

** Median better represents the "typical" American in this example.

Mode

- > The most common value in a dataset.
Not relevant to continuous numerical data.
Works only for discrete numerical data.

0, 2, 3, 2, 1, 0, 0, 2, 0

- > How many of each value are there?

0: 4 1: 1 2: 3 3: 1

Mode = 0

Normal Distribution

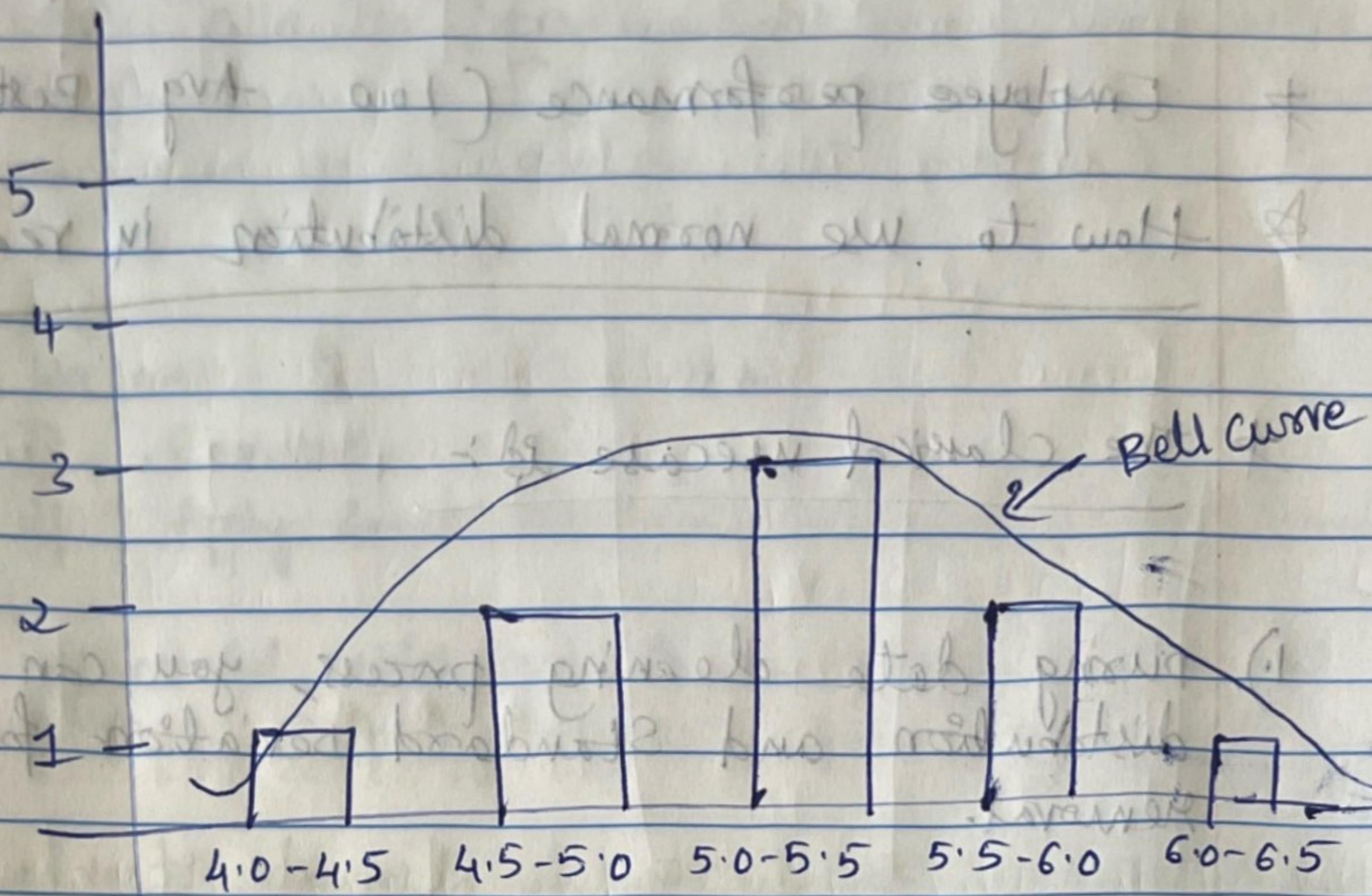
- * Let's say you want to do date analysis on people's height

Name	Height (ft)
Bob	6.2
Thomas	5.7
Nina	4.6
Mit	5.4
Sofia	5.9
Madhan	4.3
Tao	5.1
Deepika	5.2
Rafiq	4.9

- * When you plot this data on a histogram, it will look like:

Count

(Histogram is a freq distribution)



↗ When you draw a curve that passes through this histogram it looks like a bell.

↗ This is called a normal distribution. Most of your data samples are around avg value. Some far away from data samples on the left & right.

Ex:

2 bedroom apt price in Bangalore city

Count

300

200

100

35lakh

90lakh

3crre



→ we can even consider (Test score out of 100)

★ Employee performance (Low Avg Best performance)

★ How to use normal distribution in real life??

★ The classical usecase is :-

1) During data cleaning process, you can use normal distribution and Standard deviation for Outlier removal.

Ex:- Adding an outliers to our data

Outlier is a very different value from our data.

stab num	Name	height (ft)
Rob		6.2
Thomas		5.7
Nina		4.6
Mit		5.4
Sofia		5.9
Mohan		4.3
Tao		5.1
Deep		5.2
Rafiq		4.9
Smith		9.07 → outlier

→ This kind of outliers can occur because of an error in data collection / transmission process / sometimes might be valid [max in the history] [They stand out even on a histogram when plotted.]

→ We need to treat outliers

i) Remove them

ii) Apply some methods to treat them (Transformation)

→ ML model might get skewed if you don't treat outliers. So, it's better to treat them.

→ What formula do we use to remove Outliers?

1.) First, let's learn about:

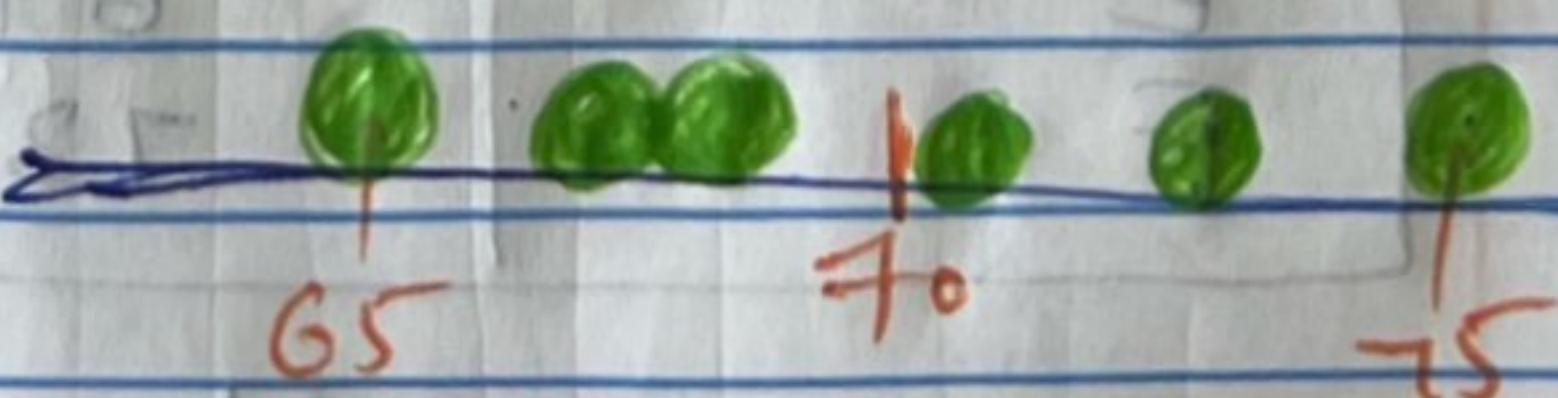
Mean Absolute Deviation

Standard Deviation

Ex:- History Test

Name	Score
A	75
B	72
C	68
D	65
E	67
F	73

Plotting on a Line chart



$$\text{Avg} = 70$$

[Nearer to average]

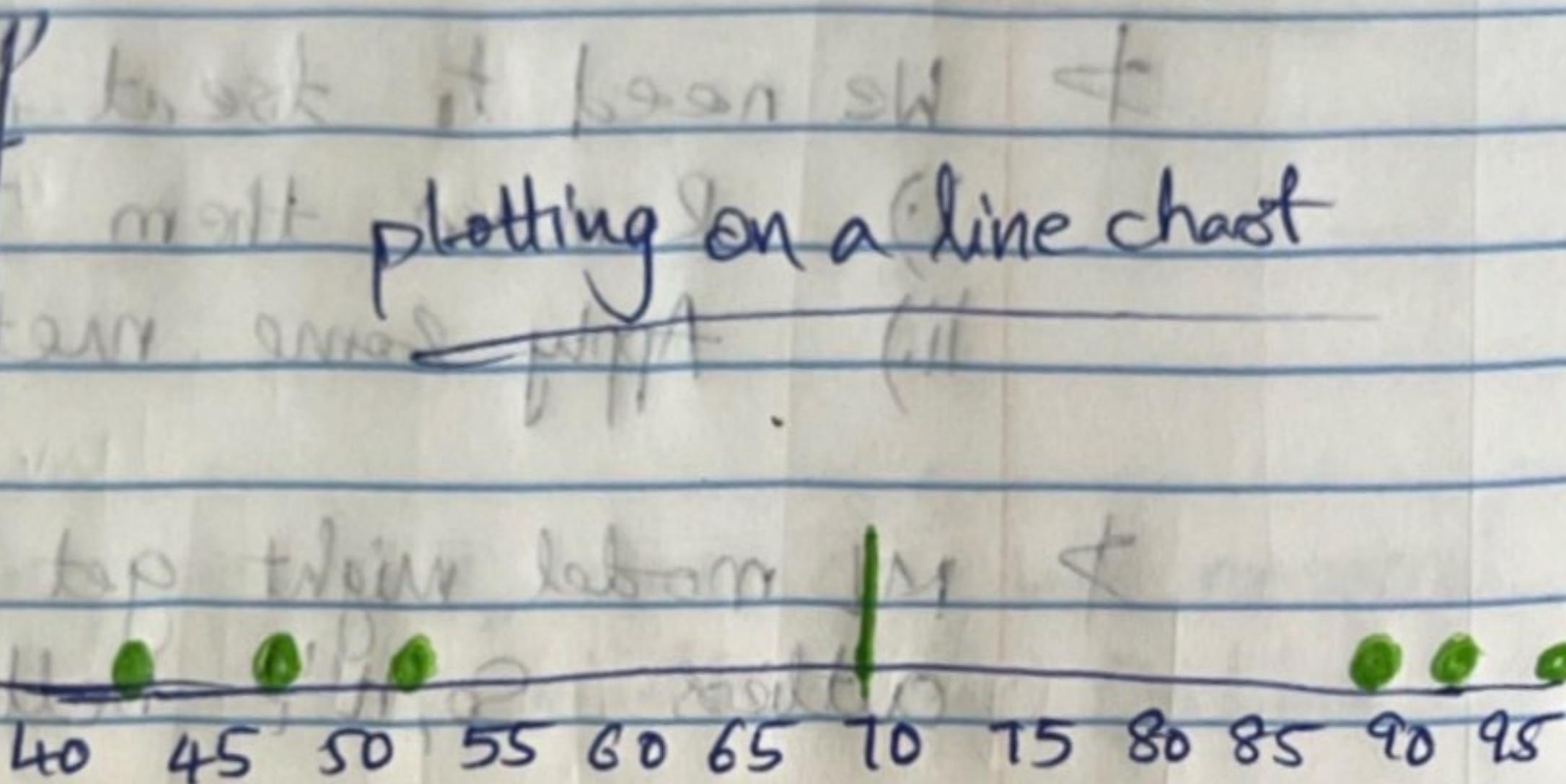
We are trying to know how far are our data points from the average [How spread out].

Ex:

Consider Other Sets

Math Test

Name	Score
A	93
B	96
C	43
D	47
E	51
F	90



[Quite spread out from average]

$$\text{Avg} = 70$$

* Let's Compute Absolute for each individual score

Ex:

History Test

Name	Score	Abs(Score - Avg)
A	75	5
B	72	2
C	68	2
D	65	5
E	67	3
F	73	2

$$\text{Mean} = \frac{19}{6} = 3.16$$

[Expression of absolute]

Math Test

Name	Score	Abs (score - Avg)
A	93	23
B	96	26
C	43	27
D	47	23
E	51	19
F	90	20
Mean		23

$$\text{Mean} = \frac{161}{6} = 23$$

* This mean is known as Mean Absolute Deviation (MAD).

- * So, this is a useful metric to represent how spread out the data points are.
- * Higher the MAD, more the data points are spread out.

* There could be a scenario where using MAD might not be enough.

Ex:

Name	Score	Abs (Score - Avg)
A	75	5
B	72	2
C	68	2
D	65	5
E	67	3
F	73	3
Mean		3.33

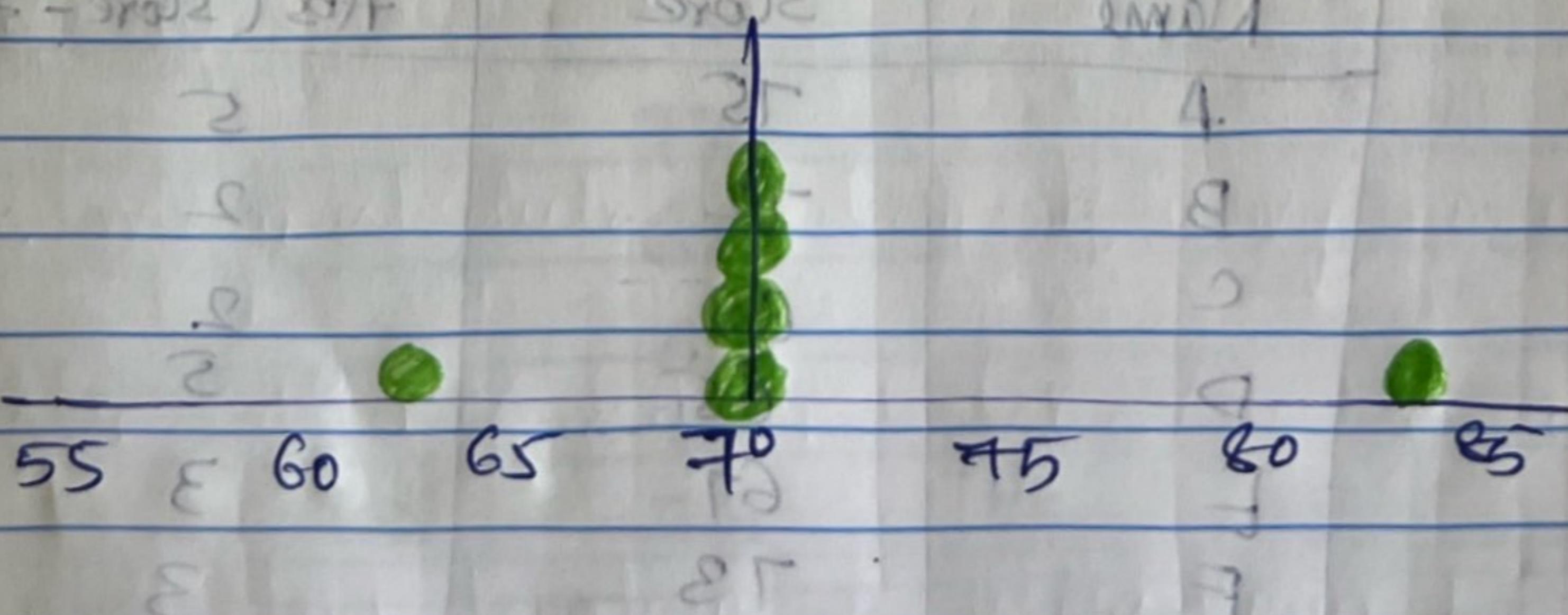
plotting:



Let's take another example:

Name	Score	Abs (Score - Avg)
A	88	13
B	70	0
C	70	0
D	63	7
E	70	0
F	70	0
MAD		3.33

plotting:



So, in this case MAD is not very useful. we need to come up with something else.

Name	Score	Abs(Score - Avg)	$(Score - Avg)^2$
A	75	5	25
B	72	2	4
C	68	2	4
D	65	5	25
E	67	3	9
F	73	3	9
		Avg	12.66
		\sqrt{Avg}	3.55

Name	Score	Abs(Score - Avg)	$(Score - Avg)^2$
A	83	13	169
B	70	0	0
C	70	0	0
D	63	7	49
E	70	0	0
F	70	0	0
		Avg	36.33
		\sqrt{Avg}	6.02

higher the (\sqrt{Avg}), more the data points are spread out.

→ this ($\sqrt{\text{Avg}}$) is known as Standard Deviation.

Formula

$$\sigma = \sqrt{\frac{\sum (x_i - \mu)^2}{N}}$$

σ = Population Standard deviation

N = Size of population

x_i = Each value from the population

μ = The population mean

Notes:

- 1.) MAD and SD (Both measure the dispersion of your data by computing the distance of the data to its mean).

→ The mean absolute deviation is using norm L_1

(Also called Manhattan distance / rectilinear distance)

→ The standard deviation is using norm L_2

(Also called Euclidean distance)

* The difference between the two norms is that the standard deviation is calculating the square of the differences whereas the mean absolute deviation is only looking at the absolute difference.

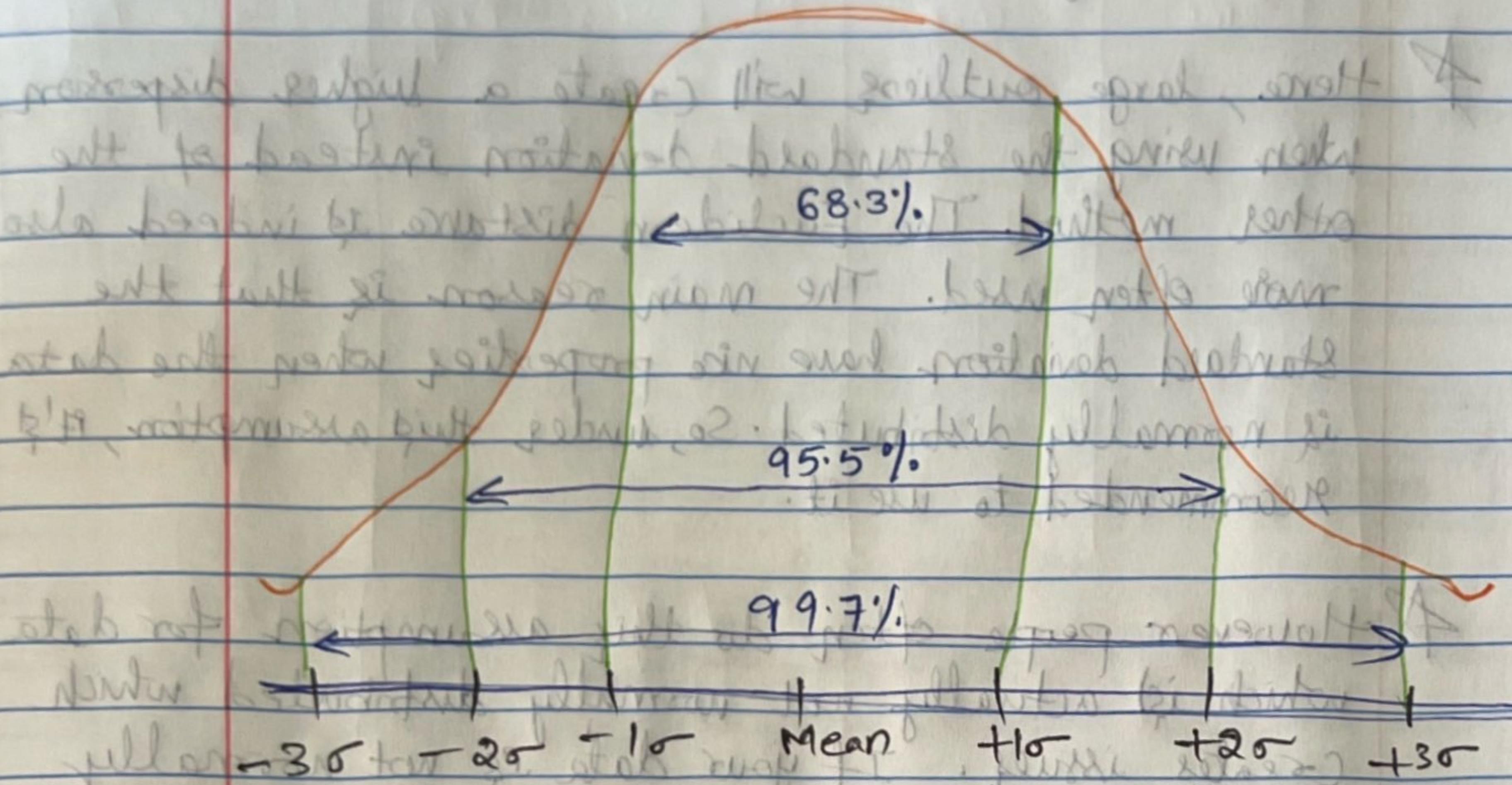
* Hence, large outliers will create a higher dispersion when using the standard deviation instead of the other method. The Euclidean distance is indeed also more often used. The main reason is that the standard deviation have nice properties when the data is normally distributed. So, under this assumption, it's recommended to use it.

* However people often do this assumption for data which is actually not normally distributed which creates issues. If your data is not normally distributed, you can still use the standard deviation, but you should be careful with the interpretation of the results.

* Identifying whether data is normally distributed is an essential step in many statistical analyses because many techniques assume normality. The methods to determine if your data are normally distributed:

- 1) Visual Inspection Methods (histogram, Q-plot)
- 2) Statistical Tests
- 3) Descriptive statistics

Using Standard Deviation to remove Outliers:



By conducting several tests on normally distributed data, mathematicians found that (68.3%) of data in any normal distribution comes in (-1σ to $+1\sigma$), 95% of data points fall under (-2σ to $+2\sigma$) and 99.7% of data under (-3σ to $+3\sigma$).

We can use this knowledge to remove outliers. If the data point is beyond ($+3\sigma$ or -3σ) can be treated as an outlier (This is the general guideline).

* Z Score :

* How many standard deviation away a datapoint is from mean.

Ex:	Name	Height (ft)	Z Score
	A	6.2	$(6.2 - 5.25) / 0.61 = 1.53$
	B	5.7	$(5.7 - 5.25) / 0.61 = 0.72$
	C	4.6	$(4.6 - 5.25) / 0.61 = -1.06$
	D	5.4	$(5.4 - 5.25) / 0.61 = 0.23$
	E	5.9	$(5.9 - 5.25) / 0.61 = 1.04$
	F	4.3	$(4.3 - 5.25) / 0.61 = -1.55$
	G	5.1	$(5.1 - 5.25) / 0.61 = -0.25$
	H	5.2	$(5.2 - 5.25) / 0.61 = -0.09$
	I	4.9	$(4.9 - 5.25) / 0.61 = -0.58$

$$\boxed{\text{Average} = 5.25 \quad | \quad \text{Standard deviation} = 0.61}$$



$$\therefore Z = \frac{x - \mu}{\sigma}$$

x = Each data point

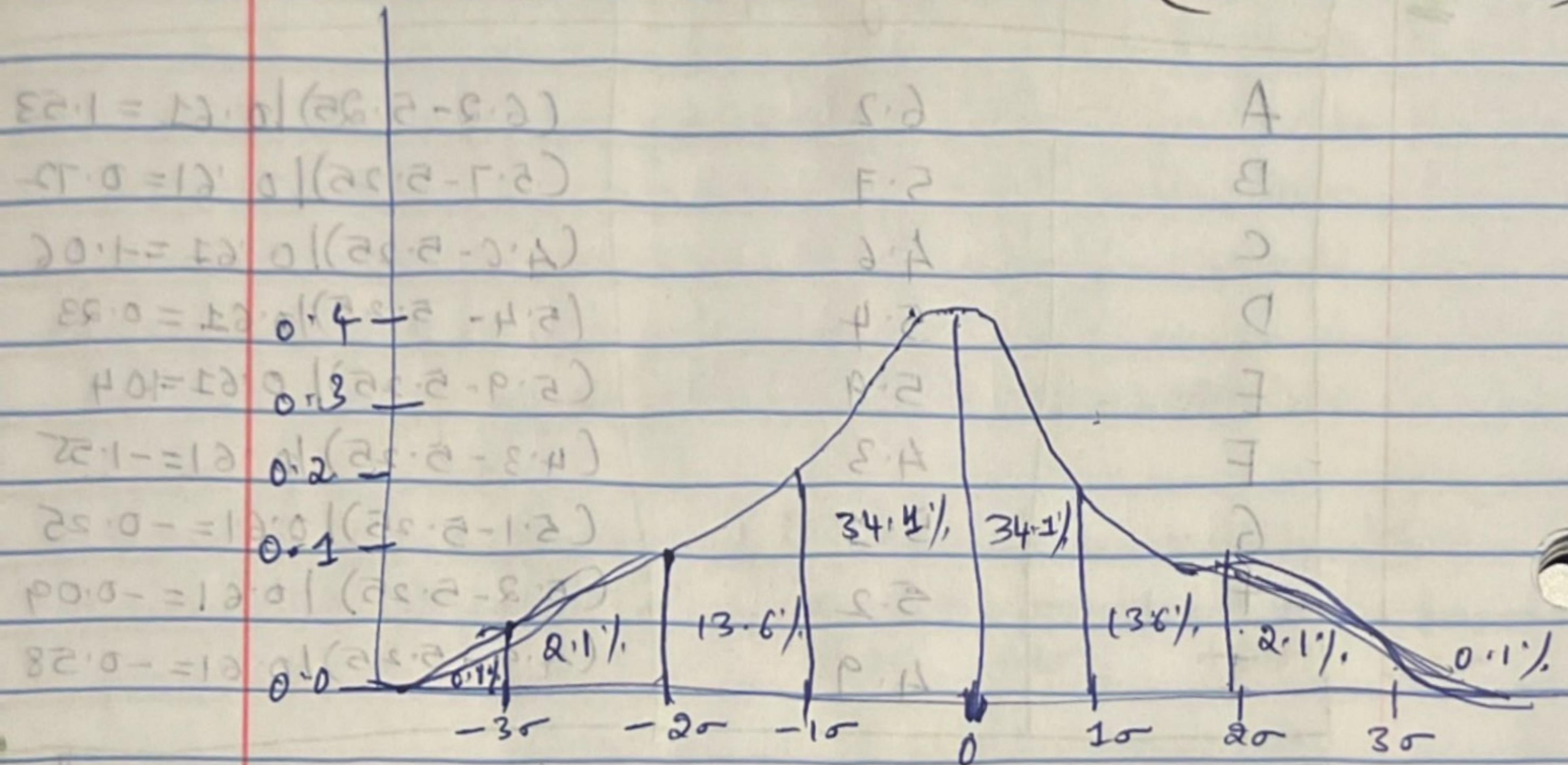
μ = Average (Mean)

σ = Standard deviation

Probability Density Functions

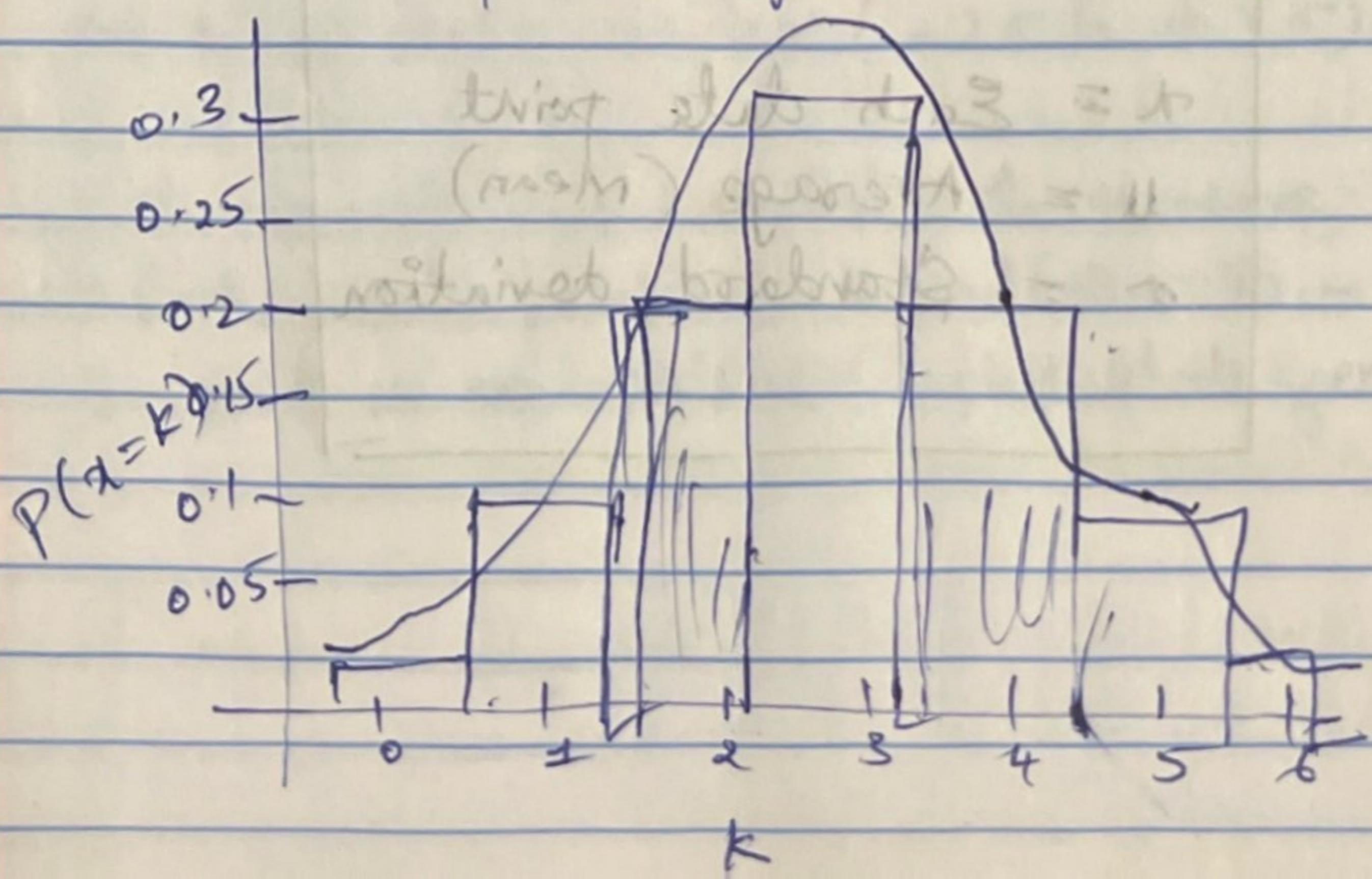
- Gives you the probability of a data point falling within some given range of a given value.

(Continuous Data)



Probability Mass Function

- Visualize the probability of discrete data.



Percentiles:

- # Percentile is a measure used to indicate the value below which a given percentage of observations fall.

$$P_x = \frac{x(n+1)}{100}$$

P_x = The value at which x percentage of data lie below that value
 n = Total no. of observations

How to calculate percentile?

- # Imagine we have a set that represents the marks of 20 students.

Marks Scored out of 100

89

97

→ To find the 90th percentile

78

45

i) Sort the data in the ascending order

94

50

66

69

50

73

43

94

$$P_{90} = \frac{90(20+1)}{100}$$

75

87

$$P_{90} = \frac{90(21)}{100}$$

31

77

53

45

$$P_{90} = 18.9 \approx 19$$

$(P_{90} = 94) \rightarrow 90\% \text{ of students got less than } 94 \text{ and } 10\% \text{ of students got more than } 94.$

$$P_{90} = 94$$

* Let's find 25th percentile

$$P_{25} = \frac{25(20+1)}{100}$$

$$P_{25} = 5.25 + 1 \approx 5$$

* Find the percentile for $\boxed{\text{marks} = 78}$

$$P = \frac{n}{N} \times 100$$

n = Ordinal rank of values

N = Total values in the dataset

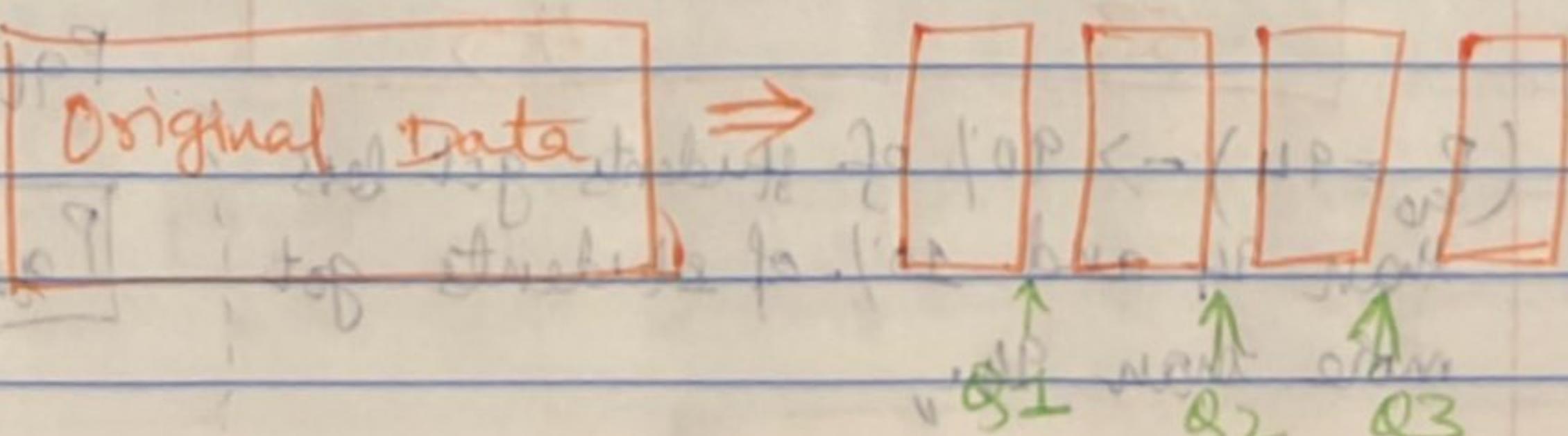
$$\text{Given, } P \neq 78 \Rightarrow 78 \neq \frac{n}{20} \times 100$$

$$P = \frac{12 \times 100}{20}$$

$$\boxed{P = 60}$$

Quartiles

* A quartile divides a sorted data set into 4 equal parts, so that each part represents 1/4 of the data set.



First quartile (Q_1) - 25% of the dataset

Second quartile (Q_2) - 50% of the dataset

Third quartile (Q_3) - 75% of the dataset

Q_1 - Lower quartile

Q_2 - Median

Q_3 - Upper quartile

25% of all the data has a value $\leq Q_1$

50% of all the data has a value $\leq M$

75% of all the data has a value $\leq Q_3$

50% of all the data lies between Q_1 and Q_3

If a measurement falls to the right of Q_3 of a set of data, we know it's in the top 25% of the data

If a measurement falls to the left of the lower quartile of a set of data, then we know it's in the bottom 25% of the data

How to Calculate Quartiles

To find Q_1, Q_2, Q_3 , we need to first sort the data in ascending order.

Find the Median = Middle Value in the dataset, which is the second Quartile.

Ex:

$$n = 15$$

$$\text{Median} = \frac{(n+1)}{2} \quad (n = \text{odd})$$

$$= 8$$

$n \geq 10$

$$\therefore \boxed{Q_2 = 8^{\text{th}} \text{ observation}}$$

$$Q_1 = (n+1) \times (1/4)$$

$$Q_1 = \frac{16}{4} = 4$$

$$\therefore \boxed{Q_1 = 4}$$

$$Q_3 = (n+1) \times (3/4)$$

$$Q_3 = \frac{16 \times 3}{4} = 12$$

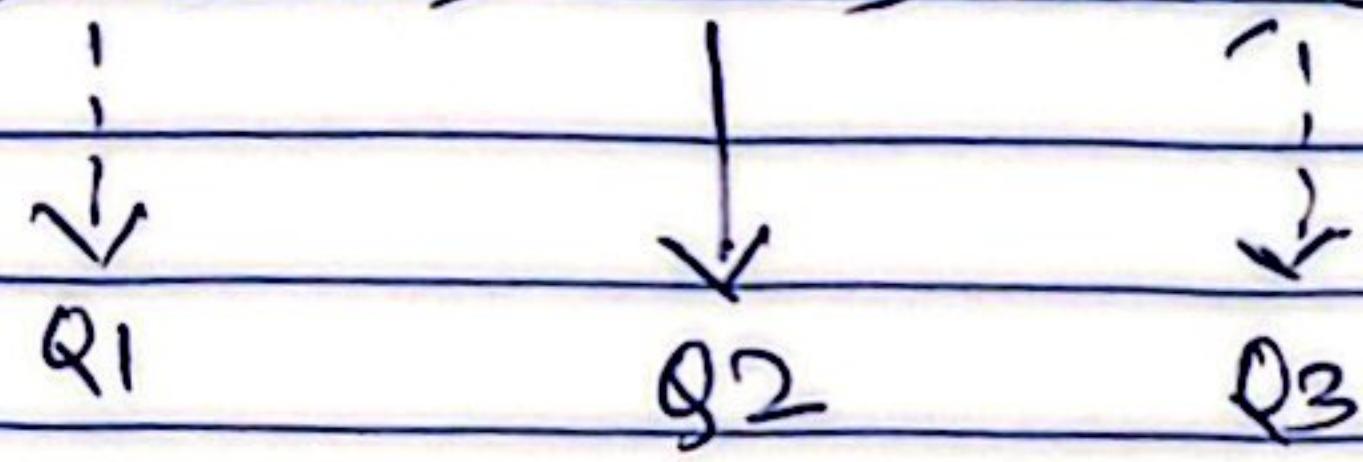
Ex

	Name	Income
0	Rob	5000
1	Rafiq	6000
2	Nina	4000
3	Sofia	7500
4	Mohan	8000
5	Tao	7000
6	Elon Musk	10000000

$$n = 7$$

+ Sorting the Income values:

[4000, 5000, 6000, 7000, 7500, 8000, 10000000]



$$Q_1 = \frac{n}{4} = \frac{7}{4} = 1.7 \approx 2$$

$$Q_2 = \text{Median} = \frac{(n+1)}{2} = \frac{8}{2} = 4$$

[$n = \text{odd} = 7$]

$$Q_3 = \frac{n \times 3}{4} = \frac{21}{4} = 5.25 \approx 5$$