

## Introduction to pandas

### Measures of Central Tendency

#### 1) mean (Arithmetic Average)

- The sum of all values divided by the total number of values
- It's the common measure but very **sensitive** to Outliers

المتوسط : المجموع على العدد

- It's used when the data is normal distribution

- There is no "outliers"

حساس جداً للقيم الشاذة

In pandas: `df['column'].mean()`

#### 2) median (middle value)

- The middle point of the data when it is arranged in ascending order.
- It's not affected by outliers
- If you have a few people with massive salaries, the median gives a more "honest" picture of the average person's income.

الوسيط : الفصير الأوسط بعد ترتيب البيانات

لا يتأثر بالقيم الشاذة

In pandas: `df['column'].median()`

#### 3) mode (most frequent)

- The value that appears most often in the dataset.
- Mainly used for Categorical Data (words/labels)

الأنز تكراراً

يفضل استخدامه مع بيانات النصية

In pandas: `df['column'].mode()`

Note: **Strategy: Skewness, Distribution:**

before you choose which one to use for Imputation (filling missing values)

You look at the distribution of your data:

- 1- Symmetrical (normal distribution) عندما تكون البيانات منطوية "لا يوجد بيزاناس شاذة"

use the **Mean**

- 2- Skewed (Asymmetrical Distribution)

"بيزاناس شاذة" الانحراف عالية جداً

The mean is pulled toward the tail" Use the **Median**.

### 3. Data Categorical:

Use the **Mode**

`df['Column'].describe()`

↳ Calculate all of them.

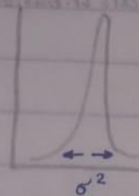
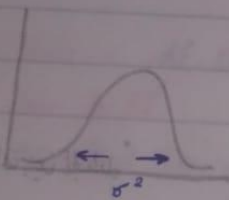
### 4. Variance ( $\sigma^2$ )

• It measures how far each number of the set from the mean.

• It is average of the squared difference from the mean.

• High variance means the data is spread out

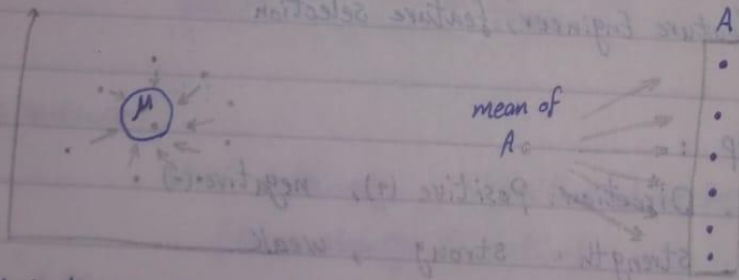
• Low variance means the data points are clustered closely around the average.



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

$$\sum_{i=1}^n \frac{(x_i - \mu)^2}{n}$$

In AI High Variance often refers to overfitting (where the model is too sensitive to small fluctuation).



if  $\mu$  is large, it means there is Variance.

if  $\mu$  is small, it means there is no variance.

$\sigma^2$

In pandas `df['Column'].var()`

means here the difference between elements of A and there mean ( $\mu$ )

## 5) standard Deviation ( $\sigma$ )

• The square root of the variance.

$$\sigma = \sqrt{\text{Variance}}$$

• we use it to identify Outliers

The difference :

Variance ( $\sigma^2$ )

standard deviation ( $\sigma$ )

• The average of squared difference from the mean.

• The square root of the variance

• squared units (it makes it hard to visualize)

• Original Units (This is intuitive)

• Used in mathematical optimization and loss functions.

• Use to detect Outliers and for data scaling.

There are used for Calculate the distribution for data.

In pandas: `df['Column'].std()`

## 6) Covariance

• It means how two variables move together.

positive Covariance: Both variables move together

negative Covariance: One increases, the other decreases

• Used for "feature Engineer",

• "PCA"

• Relation ship :

• Direction: positive (+), negative (-)

• strength: strong, weak

### Four Quadrants:

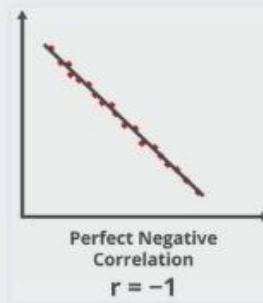
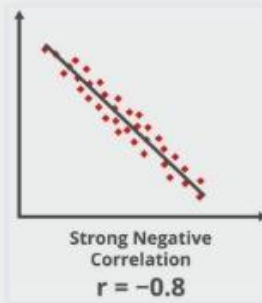
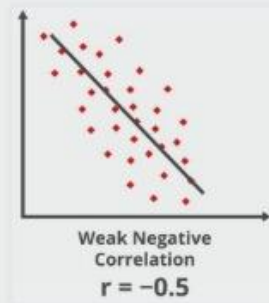
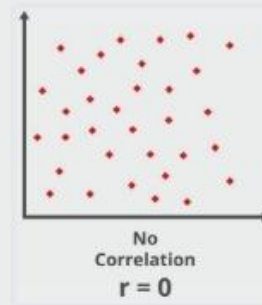
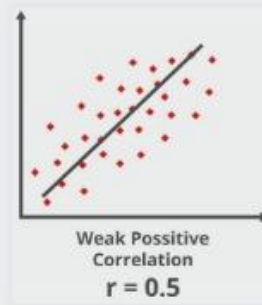
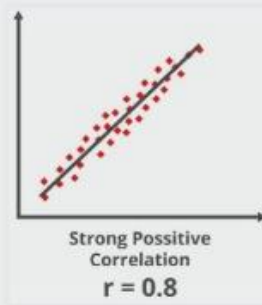
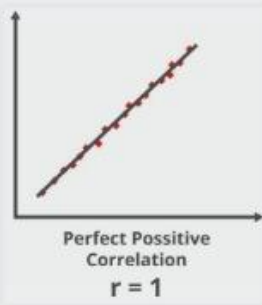
strong positive: when Column A goes up, Column B goes up significantly

strong negative: when Column A goes up, Column B goes down significantly

weak positive: There is a slight upward trend, but it is 'noisy'

weak negative: There is a slight downward trend,

# Correlation

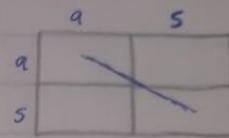




## 7) Correlation

- A standardized version of Covariance . It ranges from -1 to +1
- +1 perfect positive relationship.
- -1 perfect negative relationship.
- 0 No relationship at all.
- It's vital for "Feature selection". if two columns are 99% correlated you can drop one of them because they provide the same information for AI model.
- Used in

- Confusion matrix.
- Correlation matrix
- heatmap



"main diagonal" means relationship with 'itself'.

In Pandas: `df.corr()`

Use Case In AI:

Standard Deviation:

Data scaling and Outlier detection.

Variance:

Understanding model error

Covariance:

Used internally in algorithms like PCA  
(Dimensionality Reduction)

Correlation:

selecting which features (columns) are most important for the model.

**Feature Engineering:** Creating new information for existing columns. If we have "length", "width", you might create a new feature like 'Area'.

**Feature selection:** The process of picking the most important variable. if two columns are (highly related), we can drop one of them to make model faster and more accurate.