

# variability

November 14, 2024

Variability: - Variability refers to how spread out or dispersed the values in a dataset are from each other and from the central value (often the mean). It provides insights into the consistency and predictability of data points within a set. High variability indicates that data points are more spread out, whereas low variability suggests that data points are closer to the mean or central value.

Variability helps in understanding the distribution of data and identifying patterns, extremes

These metrics are crucial in fields like statistics, finance, engineering, and research, as the

## 1. Deviation

- Definition: Deviation measures the difference between each data point and the mean of the dataset.
- Formula:
  - Deviation =  $x_i - \text{mean}$
  - Where:
    - \*  $x_i$  = data points
    - \* mean = the average of the dataset
- Example:
  - For the data set: [4, 7, 10, 5, 12] with mean = 7.6:
  - Deviation for 4 =  $4 - 7.6 = -3.6$
  - Deviation for 7 =  $7 - 7.6 = -0.6$
  - Deviation for 10 =  $10 - 7.6 = 2.4$
  - Deviation for 5 =  $5 - 7.6 = -2.6$
  - Deviation for 12 =  $12 - 7.6 = 4.4$
- Practical Usage:
  - Measures how spread out the data points are from the mean.
  - Helps in assessing the consistency of data in business operations.
- Additional Information:
  - Deviation is often used in early data exploration to identify how spread out values are from the average.
  - Positive deviations indicate values above the mean, while negative deviations are below the mean.
- Scenario:
  - A logistics company wants to monitor the average delivery time across its various locations.
- Problem Statement:
  - The company has set a standard delivery time of 2 hours but notices that some deliveries deviate significantly from this target. They need to measure these deviations to improve their delivery accuracy.

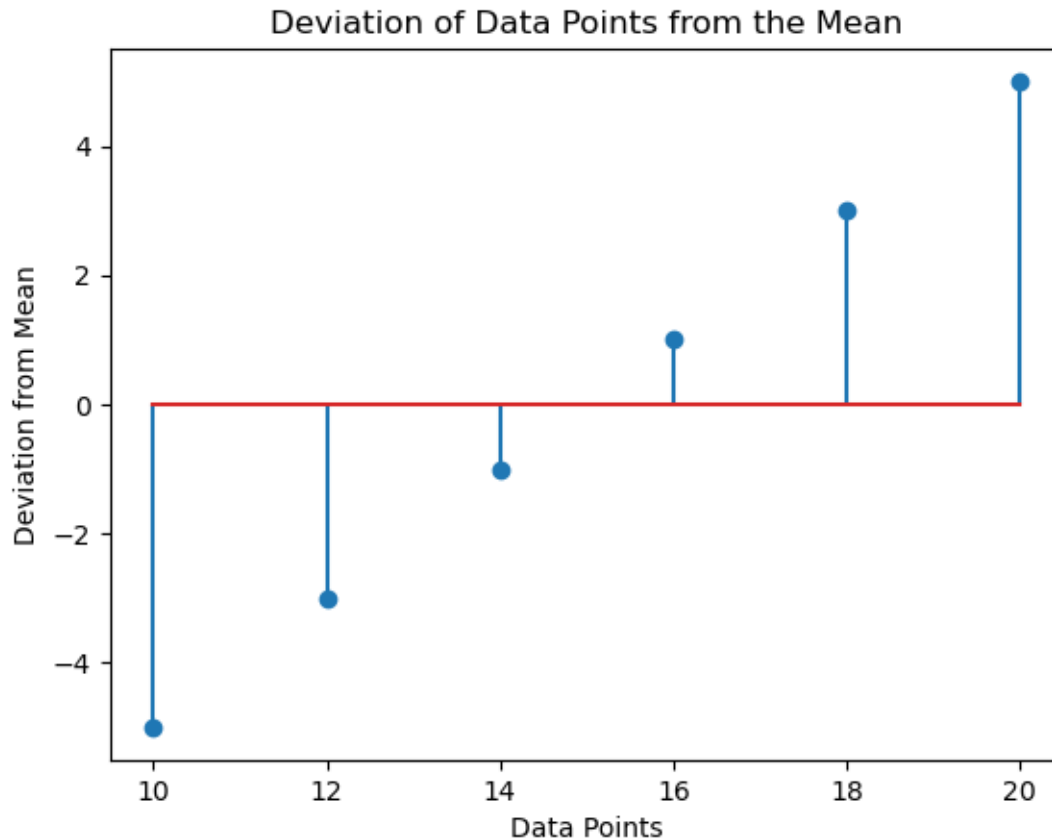
- Solution:
  - Using deviation, the company can calculate how each delivery time differs from the average delivery time. This allows them to identify locations where delivery times deviate too much from the norm and take corrective action. Deviation is chosen here because it captures each individual delivery's performance relative to the average, highlighting specific issues rather than an overall variability.
- Alternate Solutions:
  - Standard Deviation: Provides a broader view of delivery time consistency across locations, but may not pinpoint specific deviations.
  - Variance: Like standard deviation, this could be used but might make the deviation figures less interpretable for direct comparison.

```
[8]: import matplotlib.pyplot as plt

# Sample data
data = [10, 12, 14, 16, 18, 20]
mean = sum(data) / len(data)

# Calculate deviations
deviations = [x - mean for x in data]

# Plotting without 'use_line_collection'
plt.stem(data, deviations)
plt.xlabel("Data Points")
plt.ylabel("Deviation from Mean")
plt.title("Deviation of Data Points from the Mean")
plt.show()
```



## 2. Variance

- Definition: Variance measures the average of the squared deviations from the mean. It provides an indication of how spread out the values are.
  - Formula:
    - $\text{Variance} = (\Sigma(x_i - \text{mean})^2) / n$
    - Where:
      - $x_i$  = data points
      - mean = the average of the dataset
      - n = number of data points
  - Example:
    - For the data set: [4, 7, 10, 5, 12] with mean = 7.6
    - $\text{Variance} = ((4-7.6)^2 + (7-7.6)^2 + (10-7.6)^2 + (5-7.6)^2 + (12-7.6)^2) / 5 = 9.2$
  - Practical Usage:
    - Commonly used in finance to measure the risk of an asset or portfolio.
    - Helps in understanding the variability in production or sales data.
  - Additional Information:
    - Variance is an essential statistic in finance, often used to measure risk. A high variance indicates that data points are spread out widely, suggesting higher volatility.
- Scenario: A hedge fund needs to assess the risk associated with different stocks in their portfolio.

- Problem Statement:
  - The fund wants to understand how much each stock's price fluctuates over time to identify high-risk and low-risk investments.
- Solution:
  - Variance is used here to quantify the average squared deviation of stock prices from the mean price. Stocks with higher variance are more volatile, indicating higher risk, while stocks with lower variance are considered stable. Variance is ideal because it provides a clear mathematical way to quantify fluctuation and is standard in financial risk assessment.
- Alternate Solutions:
  - Mean Absolute Deviation (MAD): MAD could be used to measure fluctuation in a simpler way but lacks the weight variance places on larger deviations, which is critical in finance.
  - Standard Deviation: Provides similar information but is easier to interpret due to being on the same scale as the original data.

```
[9]: import numpy as np

# Sample data
data = [4, 7, 10, 5, 12]
mean = np.mean(data)
variance = np.var(data)

print("Variance:", variance)
```

Variance: 9.040000000000001

### 3. Standard Deviation

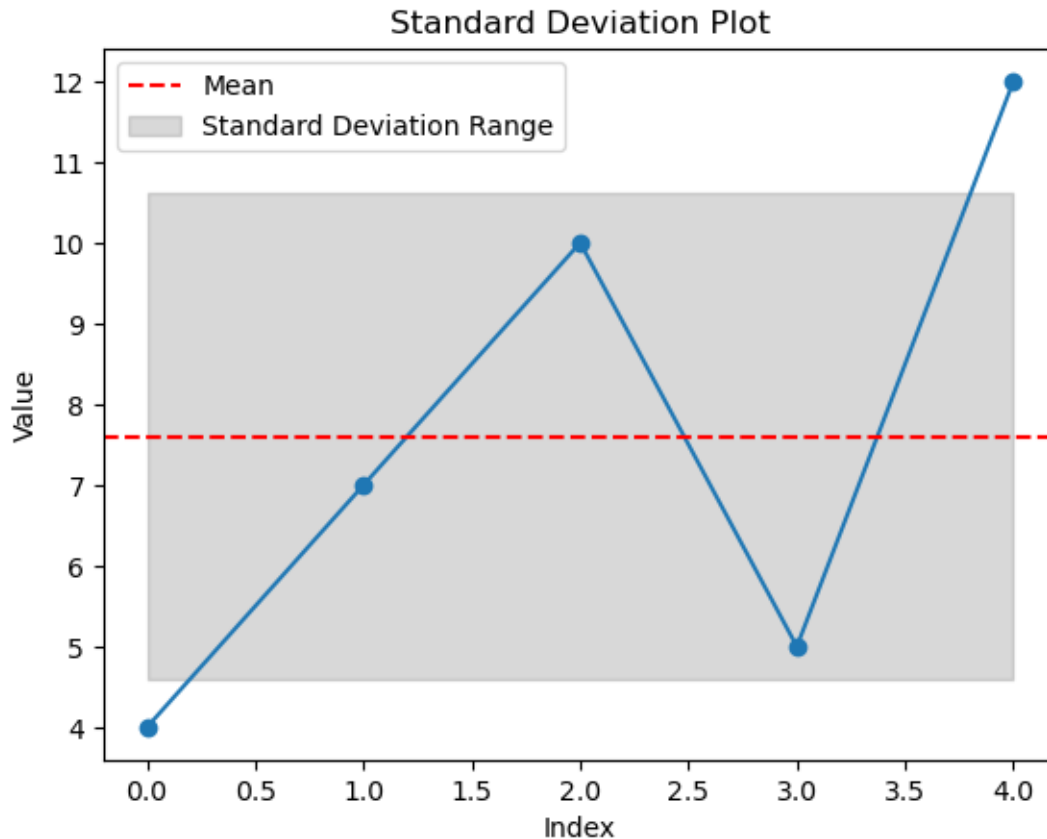
- Definition: Standard deviation is the square root of the variance and represents the spread of data points around the mean.
- Formula:
  - Standard Deviation =  $\sqrt{\text{Variance}}$
- Example:
  - For the data set: [4, 7, 10, 5, 12], variance = 9.2
  - Standard Deviation =  $\sqrt{9.2} = 3.03$
- Practical Usage:
  - Provides a measure of consistency in data (e.g., measuring the consistency of product quality).
  - Used in risk management in finance to assess volatility.
- Additional Information:
  - Standard deviation is widely used in fields such as statistics, business, and finance for measuring data consistency and risk.
- Scenario:
  - An e-commerce company is analyzing the consistency of product ratings across different categories.
- Problem Statement:
  - The company wants to determine which product categories have the most consistent ratings to better understand customer satisfaction across categories.
- Solution:

- Standard deviation is used to measure the spread of ratings for each category. Categories with a lower standard deviation in ratings indicate that customers have similar experiences with products, suggesting consistency. Higher standard deviation in ratings indicates varying customer experiences. Standard deviation is selected here because it's widely interpretable and provides insight into the consistency of feedback on the same scale as the data.
- Alternate Solutions:
  - Interquartile Range (IQR): Useful for understanding the middle 50% of ratings, which can help with outlier-heavy data but might miss overall spread.
  - Mean Absolute Deviation: Provides a straightforward way to view variability in ratings but lacks the ability to weigh outliers as heavily.

```
[10]: import numpy as np
import matplotlib.pyplot as plt

# Sample data
data = [4, 7, 10, 5, 12]
mean = np.mean(data)
std_dev = np.std(data)

# Plotting Standard Deviation Range
plt.plot(data, marker='o', linestyle='-')
plt.axhline(mean, color='r', linestyle="--", label="Mean")
plt.fill_between(range(len(data)), mean-std_dev, mean+std_dev, color="gray",
                 alpha=0.3, label="Standard Deviation Range")
plt.xlabel("Index")
plt.ylabel("Value")
plt.legend()
plt.title("Standard Deviation Plot")
plt.show()
```



#### 4. Mean Absolute Deviation (MAD)

- Definition: MAD is the average of the absolute deviations from the mean, providing a measure of the spread without squaring the deviations.
- Formula:
  - $MAD = (\sum |x_i - \text{mean}|) / n$
  - Where:
    - $x_i$  = data points
    - mean = the average of the dataset
    - n = number of data points
- Example:
  - For the data set: [4, 7, 10, 5, 12] with mean = 7.6
  - $MAD = (|4-7.6| + |7-7.6| + |10-7.6| + |5-7.6| + |12-7.6|) / 5 = 2.24$
- Practical Usage:
  - Used in fields where outliers are problematic, as it does not square the deviations.
  - Applied in manufacturing to monitor quality control.
- Additional Information:
  - MAD is preferred over variance or standard deviation in situations where a dataset contains outliers, as it doesn't square deviations and, therefore, is less affected by extreme values.
- Scenario:

- A customer service team wants to track response times and maintain consistency for inquiries.
- Problem Statement:
  - The company has noticed variability in response times across different days and wants a simple way to measure these changes to maintain customer satisfaction.
- Solution:
  - Mean Absolute Deviation (MAD) helps the team to calculate the average deviation of response times from the mean. MAD provides an easy-to-understand measure of consistency without overemphasizing extreme response times. It's useful here because the company is more interested in an intuitive, straightforward measurement than in squaring deviations (as in variance).
- Alternate Solutions:
  - Standard Deviation: Could provide a more statistically conventional measure but may overemphasize extreme deviations.
  - Range: Shows the absolute span of response times but doesn't account for consistency around the average.

```
[11]: import numpy as np

# Sample data
data = [4, 7, 10, 5, 12]
mean = np.mean(data)
mad = np.mean([abs(x - mean) for x in data])

print("Mean Absolute Deviation:", mad)
```

Mean Absolute Deviation: 2.7199999999999998

## 5. Range

- Definition: The range is the difference between the highest and lowest values in a dataset.
- Formula:  $\text{Range} = \max(x_i) - \min(x_i)$  Where:
  - $x_i$  = data points
- Example: For the data set: [4, 7, 10, 5, 12]:  $\text{Range} = 12 - 4 = 8$
- Practical Usage:
  - Provides a quick measure of the data spread.
  - Used in weather forecasts to show temperature variation.
- Additional Information:
- Range is often used as a quick way to measure data spread, though it is sensitive to outliers and doesn't indicate internal data structure.
- Scenario:
  - A manufacturing plant is monitoring the temperature of machinery to ensure operational safety.
- Problem Statement:
  - Machinery should operate within a specific temperature range. Sudden spikes or drops may indicate issues that require maintenance.
- Solution:
  - Range is ideal for this scenario because it quickly shows the highest and lowest recorded temperatures, providing a simple check for extreme variations that may indicate potential machine issues. Range is useful here due to its simplicity and

direct insight into maximum deviation from desired temperature limits.

- Alternate Solutions:
  - Interquartile Range (IQR): Would help if outliers skew the data, focusing instead on the middle 50% of values.
  - Variance or Standard Deviation: Good for understanding spread but unnecessary for simple threshold-based monitoring.

```
[12]: data = [4, 7, 10, 5, 12]
      range_value = max(data) - min(data)

      print("Range:", range_value)
```

Range: 8

## 6. Percentile

- Definition: A percentile indicates the relative standing of a value in a dataset. The p-th percentile is the value below which p% of the data fall.
- Formula:
  - $\text{Percentile} = (p/100) * (n + 1)$
  - Where:
    - \* p = percentile
    - \* n = number of data points
- Example:
  - For the data set: [4, 7, 10, 5, 12] and p = 50 (50th percentile)
  - The 50th percentile (median) is 7
- Practical Usage:
  - Used to understand the distribution of data (e.g., income distribution, test scores).
  - Helps in decision-making, such as salary benchmarking or grading.
- Additional Information:
  - Percentiles are common in standardized testing and grading. They help establish a relative standing in comparison to other values in the dataset.
- Scenario:
  - A hospital wants to assess patient wait times in the emergency room.
- Problem Statement:
  - The hospital aims to understand how their wait times compare to benchmarks (e.g., 80th percentile of wait times within acceptable limits) to improve patient experience.
- Solution:
  - Percentiles allow the hospital to assess where the bulk of patient wait times fall, focusing on specific benchmarks like the 80th or 90th percentile to identify prolonged wait times. This approach is useful because it provides actionable insights about typical versus extreme wait times, making it easy to set goals based on real performance metrics.
- Alternate Solutions:
  - Mean and Median Wait Times: Useful for general understanding but may not highlight extreme wait times as effectively as percentiles.
  - Range or IQR: Could help track variation but doesn't directly offer percentile-based goals.



```
[13]: import numpy as np

# Sample data
data = [4, 7, 10, 5, 12]
percentile_50 = np.percentile(data, 50)

print("50th Percentile (Median):", percentile_50)
```

50th Percentile (Median): 7.0

## 7. Interquartile Range (IQR)

- Definition: The IQR measures the range within which the middle 50% of the data lies, defined as the difference between the 75th percentile (Q3) and 25th percentile (Q1).
- Formula:
  - $IQR = Q3 - Q1$
  - Where:
    - \*  $Q1 = 25\text{th percentile}$
    - \*  $Q3 = 75\text{th percentile}$
- Example:
  - For the data set: [4, 7, 10, 5, 12] with  $Q1 = 5$  and  $Q3 = 10$
  - $IQR = 10 - 5 = 5$
- Practical Usage:
  - Used to identify outliers (values that fall outside  $1.5 \times IQR$  above  $Q3$  or below  $Q1$ ).
  - Applied in data visualization, particularly in box plots, to show data spread.
- Additional Information:
  - IQR is used to identify outliers and understand data spread. Data points that lie beyond 1.5 times the IQR above  $Q3$  or below  $Q1$  are often considered outliers.
- Scenario:
  - An educational testing company wants to assess score distributions and detect outliers for student performance on a standardized test.
- Problem Statement:
  - The company needs a measure that helps in identifying students with exceptionally high or low scores for specialized assistance programs or additional testing.
- Solution:
  - Interquartile Range (IQR) focuses on the middle 50% of the scores, allowing the company to understand the spread without being skewed by outliers. IQR is effective here as it allows them to detect extreme scores (outside the typical range), making it ideal for identifying potential outliers for further evaluation.
- Alternate Solutions:
  - Standard Deviation: Useful to measure overall spread but might overemphasize extreme values.
  - Range: Provides the absolute range of scores but doesn't reflect internal distribution patterns, limiting outlier identification.

```
[14]: import numpy as np
import matplotlib.pyplot as plt

# Sample data
```

```
data = [4, 7, 10, 5, 12]
q1 = np.percentile(data, 25)
q3 = np.percentile(data, 75)
iqr = q3 - q1

# Box plot to illustrate IQR
plt.boxplot(data, vert=False)
plt.xlabel("Value")
plt.title(f"Interquartile Range (IQR) = {iqr}")
plt.show()
```

