

# Blog Tuần 3 – Module 2

## *Từ Thống Kê Đến Firebase và Nghệ Thuật Trình Bày*

**Tác giả:** GRID034

---

Tuần thứ ba của Module 2 đưa bạn vào một hành trình đặc biệt: từ nền tảng thống kê xác suất đến việc vận dụng nó trong xây dựng hệ thống trí tuệ nhân tạo với Firebase, và kết thúc bằng kỹ năng trình bày thông tin chuyên nghiệp. Đây là sự kết hợp giữa tư duy định lượng và nghệ thuật kể chuyện dữ liệu.

Các chủ đề nổi bật bao gồm:

### **1. Thống kê mô tả cơ bản**

Làm quen với các khái niệm nền tảng như biến ngẫu nhiên, kỳ vọng, phương sai, độ lệch chuẩn. Giải thích qua ví dụ trực quan và ứng dụng phân tích dữ liệu học sinh, thu nhập, hoặc kết quả thi.

### **2. Hiệp phương sai và Hệ số tương quan**

Khám phá mối quan hệ giữa các biến số: khi nào chúng đồng biến, khi nào nghịch biến? Cách tính Covariance và Correlation, minh họa qua dữ liệu điểm thi và ứng dụng trong phát hiện xu hướng trong AI.

### **3. Firebase và Firestore trong AI**

Giới thiệu toàn diện về nền tảng Firebase, cách cấu hình Realtime Database và Firestore, lưu trữ dữ liệu không cần server backend, kết nối Python để thực hiện các thao tác CRUD và đồng bộ hóa dữ liệu.

### **4. Huấn luyện mô hình Machine Learning trên Firebase**

Sử dụng dữ liệu thật (Iris dataset), lưu trữ, truy xuất và huấn luyện các mô hình như Decision Tree, SVM, Logistic Regression. Tích hợp Streamlit để tạo web app đơn giản dự đoán kết quả và ghi lại dữ liệu người dùng.

### **5. Kỹ năng trình bày và kể chuyện bằng dữ liệu**

Trang bị kỹ năng thiết kế slide chuyên nghiệp và truyền đạt thông điệp hiệu quả: nguyên tắc 5-5-5, quy tắc màu 60-30-10, biểu đồ chọn lọc (bar, line, scatter). Sử dụng kỹ thuật SCQA và sơ đồ Logical Tree để kể chuyện dữ liệu logic và thuyết phục.

Tất cả được trình bày một cách logic, minh họa rõ ràng và áp dụng thực tiễn, giúp bạn xây dựng kỹ năng tổng hợp: từ phân tích – mô hình hóa – đến truyền đạt hiệu quả.

# Basic Statistic

Dao Lam Hoang

## Mở đầu

Thống kê là nền tảng không thể thiếu trong khoa học dữ liệu, trí tuệ nhân tạo và nhiều lĩnh vực khác. Tài liệu này nhằm giới thiệu các khái niệm cơ bản trong thống kê như: biến ngẫu nhiên, các hàm phân phối xác suất, kỳ vọng, phương sai và tương quan. Mỗi phần đều đi kèm ví dụ minh họa rõ ràng, giúp người đọc dễ tiếp cận và hiểu sâu vấn đề.

## 1. Biến ngẫu nhiên

**Biến ngẫu nhiên**  $X$  là một hàm số  $X : \Omega \rightarrow \mathbb{R}$ , ánh xạ một kết quả  $s \in \Omega$  tới một số thực trên trục số thực, tức là  $X(s) \in \mathbb{R}$ .

### Biến ngẫu nhiên liên tục $X$

$X(s) : \Omega \rightarrow \mathbb{R}$ , ánh xạ mỗi  $s$  trong không gian mẫu (không đếm được) đến một giá trị thực.

**Ví dụ:** Thời gian chờ đợi xe buýt  $T$  (phút), nhận mọi giá trị thực không âm.

$$f_T(t) = \lambda e^{-\lambda t}, \quad t \geq 0, \lambda > 0.$$

### Biến ngẫu nhiên rời rạc $X$

$X(s) : \Omega \rightarrow \mathbb{R}$ , ánh xạ mỗi  $s$  trong không gian mẫu hữu hạn hoặc đếm được đến một giá trị thực.

**Ví dụ:** Số chấm trên mặt xúc xắc  $X$ :

$$X \in \{1, 2, 3, 4, 5, 6\}, \quad P(X = x) = \frac{1}{6}$$

## 2. Các Hàm Phân Phối Xác Suất

### 2.1 Hàm phân phối xác suất rời rạc (Probability Mass Function - PMF)

Cho biến ngẫu nhiên rời rạc  $X$  có tập giá trị  $S = \{x_1, x_2, \dots\}$ . Hàm phân phối xác suất của  $X$  được định nghĩa bởi:

$$p_X(x) = P(X = x), \quad \forall x \in S,$$

thỏa mãn:

$$p_X(x) \geq 0, \quad \sum_{x \in S} p_X(x) = 1.$$

**Ví dụ:** Gieo một con xúc xắc đều, biến ngẫu nhiên  $X$  biểu diễn số mặt xúc xắc xuất hiện. Tập

giá trị:  $S = \{1, 2, 3, 4, 5, 6\}$  và

$$p_X(x) = P(X = x) = \frac{1}{6}, \quad x = 1, 2, 3, 4, 5, 6.$$

## 2.2 Hàm phân phối xác suất tích lũy (Cumulative Distribution Function - CDF)

Hàm phân phối tích lũy của biến ngẫu nhiên  $X$  được định nghĩa bởi:

$$F_X(x) = P(X \leq x).$$

Với biến ngẫu nhiên rời rạc,

$$F_X(x) = \sum_{t \leq x} p_X(t).$$

**Ví dụ:** Với  $X$  là số mặt xúc xắc khi gieo kẻ trên,

$$F_X(3) = P(X \leq 3) = p_X(1) + p_X(2) + p_X(3) = \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = \frac{3}{6} = 0.5.$$

Ngoài ra,  $F_X$  có dạng bậc thang, như sau:

$$F_X(x) = \begin{cases} 0, & x < 1, \\ \frac{1}{6}, & 1 \leq x < 2, \\ \frac{2}{6}, & 2 \leq x < 3, \\ \frac{3}{6}, & 3 \leq x < 4, \\ \frac{4}{6}, & 4 \leq x < 5, \\ \frac{5}{6}, & 5 \leq x < 6, \\ 1, & x \geq 6. \end{cases}$$

## 2.3 Hàm mật độ xác suất (Probability Density Function - PDF)

Với biến ngẫu nhiên liên tục  $X$ , hàm mật độ xác suất  $f_X(x)$  thỏa mãn:

$$P(a \leq X \leq b) = \int_a^b f_X(x) dx,$$

với mọi  $a, b \in \mathbb{R}$ , và

$$f_X(x) \geq 0, \quad \int_{-\infty}^{+\infty} f_X(x) dx = 1.$$

Hàm phân phối tích lũy tương ứng:

$$F_X(x) = \int_{-\infty}^x f_X(t) dt.$$

**Ví dụ:** Giả sử biến ngẫu nhiên  $X$  biểu diễn thời gian chờ đợi (tính bằng phút) cho một xe buýt đến, phân phối theo phân phối mũ với tham số  $\lambda > 0$ :

$$f_X(x) = \begin{cases} \lambda e^{-\lambda x}, & x \geq 0, \\ 0, & \text{ngược lại.} \end{cases}$$

Khi đó, hàm phân phối tích lũy là:

$$F_X(x) = P(X \leq x) = \int_0^x \lambda e^{-\lambda t} dt = 1 - e^{-\lambda x}, \quad x \geq 0.$$

### 3. Kỳ vọng, trung bình, phương sai và độ lệch chuẩn

#### 3.1 Kỳ vọng (Expected Value), Trung bình (Mean)

**Kỳ Vọng:** Là giá trị trung bình lý thuyết mà biến ngẫu nhiên có thể nhận được nếu phép thử được lặp lại vô số lần.

Cho biến ngẫu nhiên rời rạc  $X$  với giá trị  $x_1, x_2, \dots, x_n$  và xác suất tương ứng  $p_1, p_2, \dots, p_n$ :

$$\mathbb{E}[X] = \mu = \sum_{i=1}^n x_i p_i = \mathbb{E}[X^2] - \mu^2$$

Trong trường hợp liên tục:

$$\mathbb{E}[X] = \int_{-\infty}^{+\infty} x f_X(x) dx$$

#### 3.2 Phương sai (Variance) và Độ lệch chuẩn (Standard Deviation)

**Phương sai:** Đo mức độ phân tán hoặc độ biến động của các giá trị biến ngẫu nhiên quanh giá trị kỳ vọng.

$$\text{Var}(X) = \mathbb{E}[(X - \mu)^2] = \sum_{i=1}^n (x_i - \mu)^2 p_i$$

**Độ lệch chuẩn:** giúp đo độ phân tán theo cách trực quan và dễ hiểu hơn

$$\sigma_X = \sqrt{\text{Var}(X)}$$

### 3.3 Ví dụ: Tung một đồng xu hai lần

#### Ví dụ

- Gọi  $X$  là số lần xuất hiện mặt ngửa khi tung hai đồng xu liên tiếp.

- Các giá trị của  $X$ : 0, 1, 2.
- Xác suất:
  - $P(X = 0) = \frac{1}{4}$  (cả hai lần đều là mặt sấp)
  - $P(X = 1) = \frac{2}{4} = \frac{1}{2}$  (một lần ngửa, một lần sấp)
  - $P(X = 2) = \frac{1}{4}$  (cả hai lần đều là mặt ngửa)

- Tính toán

$$\mathbb{E}[X] = 0 \cdot \frac{1}{4} + 1 \cdot \frac{1}{2} + 2 \cdot \frac{1}{4} = 0 + \frac{1}{2} + \frac{2}{4} = 1$$

$$\text{Var}(X) = (0 - 1)^2 \cdot \frac{1}{4} + (1 - 1)^2 \cdot \frac{1}{2} + (2 - 1)^2 \cdot \frac{1}{4} = 1 \cdot \frac{1}{4} + 0 + 1 \cdot \frac{1}{4} = \frac{1}{2}$$

$$\sigma_X = \sqrt{\text{Var}(X)} = \sqrt{\frac{1}{2}} \approx 0.707$$

## 4. Hiệp phương sai và Hệ số tương quan

### 4.1 Hiệp phương sai (Covariance)

**Hiệp phương sai** giữa  $X$  và  $Y$  là biểu thị hướng và mức độ liên hệ tuyến tính giữa chúng.

Cho hai biến ngẫu nhiên  $X$  và  $Y$  với kỳ vọng  $\mu_X = \mathbb{E}[X]$  và  $\mu_Y = \mathbb{E}[Y]$ .

$$\text{Cov}(X, Y) = \mathbb{E}[(X - \mu_X)(Y - \mu_Y)] = \mathbb{E}[XY] - \mu_X \mu_Y = \frac{\sum (x - \mu_x)(y - \mu_y)}{n}$$

- Nếu  $\text{Cov}(X, Y) > 0$ ,  $X$  và  $Y$  có xu hướng tăng cùng nhau.
- Nếu  $\text{Cov}(X, Y) < 0$ ,  $X$  và  $Y$  có xu hướng ngược chiều.

### 4.2 Hệ số tương quan (Correlation coefficient)

Hệ số tương quan Là hiệp phương sai đã được chuẩn hóa để nằm trong khoảng  $[-1, 1]$ .

Với  $\sigma_X = \sqrt{\text{Var}(X)}$  và  $\sigma_Y = \sqrt{\text{Var}(Y)}$  là độ lệch chuẩn của  $X$  và  $Y$ , ta có hệ số tương quan:

$$\rho_{X,Y} = \text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y},$$

- $\rho_{X,Y} = 1$ : tuyến tính tăng hoàn hảo.
- $\rho_{X,Y} = -1$ : tuyến tính giảm hoàn hảo.
- $\rho_{X,Y} = 0$ : không có tương quan tuyến tính.

### 4.3 Ví dụ: Hiệp phương sai và hệ số tương quan

#### Ví dụ

Giả sử ta có hai biến ngẫu nhiên  $X$  và  $Y$  biểu diễn điểm Toán và Lý của 5 học sinh:

Học sinh	$X$ (Toán)	$Y$ (Lý)
1	6	7
2	7	6
3	8	9
4	9	10
5	10	9

Tính kỳ vọng:

$$\mu_X = \frac{6 + 7 + 8 + 9 + 10}{5} = 8, \quad \mu_Y = \frac{7 + 6 + 9 + 10 + 9}{5} = 8.2$$

Tính hiệp phương sai:

$$\begin{aligned} \text{Cov}(X, Y) &= \frac{1}{5} \sum (x_i - \mu_X)(y_i - \mu_Y) \\ &= \frac{1}{5} [(-2)(-1.2) + (-1)(-2.2) + 0(0.8) + 1(1.8) + 2(0.8)] \\ &= 1.72 \end{aligned}$$

Tính độ lệch chuẩn:

$$\sigma_X = \sqrt{\frac{1}{5} \sum (x_i - \mu_X)^2} = \sqrt{2}, \quad \sigma_Y = \sqrt{\frac{1}{5} \sum (y_i - \mu_Y)^2} \approx 1.326$$

Hệ số tương quan:

$$\rho_{X,Y} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} \approx \frac{1.72}{\sqrt{2} \cdot 1.326} \approx 0.918$$

**Kết luận:**  $\rho_{X,Y} \approx 0.918$  cho thấy mối tương quan tuyến tính mạnh giữa điểm Toán và điểm Lý.

### Kết luận

Qua bài viết này, chúng ta đã tìm hiểu những khái niệm cốt lõi trong thống kê cơ bản: từ biến ngẫu nhiên, các hàm phân phối xác suất, đến các đại lượng như kỳ vọng, phương sai, tương quan. Những khái niệm này không chỉ là nền tảng lý thuyết mà còn có vai trò thiết yếu trong việc phân tích và hiểu dữ liệu trong thực tiễn. Việc nắm vững những nội dung này sẽ giúp bạn tự tin hơn khi bước vào các chủ đề nâng cao như suy luận thống kê, học máy hoặc phân tích dữ liệu lớn.

# Hiểu Đúng Về Thống Kê và Ứng Dụng Trong Machine Learning: Từ Mean, Median đến Variance

*Vũ Thái Sơn*

Hành trình khám phá ý nghĩa thực sự của các chỉ số thống kê và vai trò quan trọng trong AI

## 1. Giới thiệu: Tại sao thống kê quan trọng trong Machine Learning?

Bạn có bao giờ tự hỏi tại sao các thuật toán AI có thể "hiểu" dữ liệu và dự đoán chính xác như thần? Bí mật nằm ở những khái niệm thống kê đơn giản như Mean, Median, và Variance – những thứ gần gũi như tính điểm trung bình lớp học hay lý giải tại sao Netflix gợi ý đúng phim bạn thích [1]. Hãy cùng bắt đầu chuyến phiêu lưu thống kê để khám phá sức mạnh của Machine Learning!

### 1.1 Vì sao thống kê là "linh hồn" của Machine Learning?

Machine Learning không phải phép màu – nó là thống kê được áp dụng một cách thông minh. Theo Nalisnick [2], ngay cả các mô hình deep learning hiện đại cũng bắt rễ từ những nguyên lý thống kê cơ bản. Trong bài blog này, chúng ta sẽ:

- Hiểu rõ ý nghĩa của Mean, Median, và Variance trong AI.
- Thực hành với Python để cảm nhận sức mạnh của thống kê.
- Áp dụng vào các bài toán thực tế, từ phân tích dữ liệu đến xử lý ảnh.

## 2. Mean - Không chỉ là "trung bình" đơn thuần

### 2.1 Mean là gì?

Mean (trung bình) là "trọng tâm" của dữ liệu, được tính bằng:

$$\mu = \frac{1}{n} \sum_{i=1}^n x_i$$

**Ví dụ đời thường:** Hãy tưởng tượng bạn là giáo viên chấm điểm cho 30 học sinh. Điểm trung bình 8.5 không chỉ nói lên mức độ "trung bình" mà còn cho biết: nếu điểm dao động từ 5 đến 10, lớp học có sự phân hóa; nếu điểm đều quanh 8–9, lớp đồng đều hơn [3].

### 2.2 Mean trong Machine Learning: Chuẩn hóa dữ liệu

Mean là công cụ quan trọng trong chuẩn hóa dữ liệu, giúp các thuật toán so sánh dữ liệu trên cùng một thang đo:



```

1 import numpy as np
2 from sklearn.preprocessing import StandardScaler
3
4 # Sample dataset: student math scores
5 scores = [85, 90, 78, 92, 88, 76, 94, 82]
6
7 # Calculate mean
8 mean_score = np.mean(scores)
9 print(f"Mean score: {mean_score:.2f}")
10
11 # Standardize scores using StandardScaler
12 scaler = StandardScaler()
13 standardized_scores = scaler.fit_transform(np.array(scores).reshape(-1, 1))
14
15 print("Original vs Standardized scores:")
16 for i in range(len(scores)):
17     print(f"Score: {scores[i]} -> Standardized: {standardized_scores[i][0]:.2f}")

```

## 2.3 Hạn chế của Mean

Mean rất nhạy với outliers (giá trị ngoại lai). Ví dụ: thu nhập trung bình của một khu phố có thể bị "kéo lệch" bởi một vài tỷ phú, hay thời gian phản hồi của website tăng vọt vì vài request chậm chạp.

## 3. Median - "Người trọng tài công bằng"

### 3.1 Tại sao Median quan trọng?

Median là giá trị ở giữa khi sắp xếp dữ liệu, không bị ảnh hưởng bởi outliers, giống như một trọng tài công bằng không bị lung lay bởi những giá trị cực đoan [4]:

$$\text{Median} = \begin{cases} x_{(n+1)/2} & \text{nếu } n \text{ lẻ} \\ \frac{x_{n/2} + x_{n/2+1}}{2} & \text{nếu } n \text{ chẵn} \end{cases}$$

### 3.2 Ví dụ thực tế: Thu nhập

Người	Thu nhập (triệu VNĐ/tháng)
Anh A	8
Chị B	12
Anh C	15
Chị D	5
CEO E	200
<b>Mean</b>	48 triệu
<b>Median</b>	12 triệu

Bảng 1: So sánh Mean vs Median trong thu nhập

Thu nhập 200 triệu của CEO kéo Mean lên 48 triệu, nhưng Median 12 triệu phản ánh đúng hơn mức thu nhập phổ biến của nhóm.

### 3.3 Median trong xử lý ảnh: Giảm nhiễu

Median rất hiệu quả trong việc loại bỏ nhiễu ảnh mà không làm mờ chi tiết:

```

1 import cv2
2 import numpy as np
3 import matplotlib.pyplot as plt
4
5 # Load a noisy image in grayscale
6 image = cv2.imread('noisy_image.ôn định, dễ dự đo', cv2.IMREAD_GRAYSCALE)
7
8 # Apply median filter to reduce noise
9 denoised_image = cv2.medianBlur(image, 5) # 5x5 kernel
10
11 # Display results
12 plt.figure(figsize=(10, 5))
13 plt.subplot(1, 2, 1)
14 plt.imshow(image, cmap='gray')
15 plt.title('Original Noisy Image')
16 plt.axis('off')
17
18 plt.subplot(1, 2, 2)
19 plt.imshow(denoised_image, cmap='gray')
20 plt.title('After Median Filter')
21 plt.axis('off')
22
23 plt.tight_layout()
24 plt.show()

```

**Giải thích:** Bộ lọc Median thay thế các pixel nhiễu bằng giá trị trung vị của các pixel lân cận, giữ được các cạnh sắc nét [5].

## 4. Variance và Standard Deviation - "Cảm biến" sự biến động

### 4.1 Variance và Standard Deviation là gì?

Variance đo mức độ "rải rác" của dữ liệu quanh Mean:

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

Standard Deviation là căn bậc hai của Variance, cho cảm giác trực quan hơn về độ phân tán:

$$\text{Std}(X) = \sqrt{\text{Var}(X)}$$

### 4.2 Ví dụ sinh động: Hai lớp học khác nhau

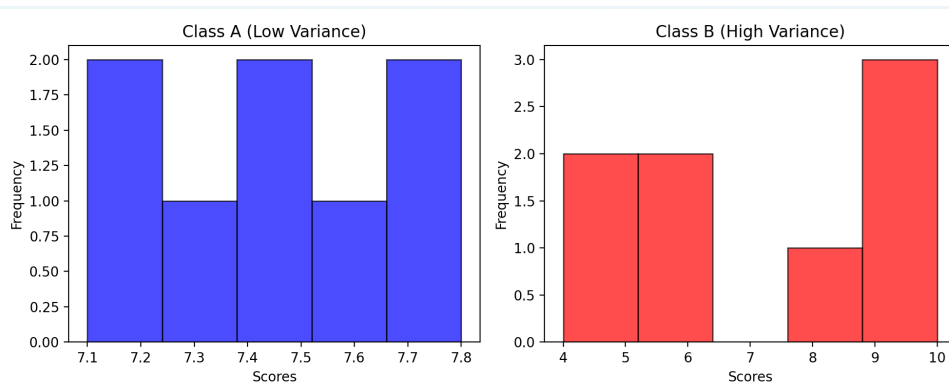
Hãy tưởng tượng hai lớp học có cùng điểm trung bình 7.5:

- **Lớp A:** Điểm từ 7.0–8.0 (Std = 0.3)
- **Lớp B:** Điểm từ 4.0–10.0 (Std = 1.8)

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 # Scores for two classes
5 class_A = [7.2, 7.5, 7.8, 7.3, 7.6, 7.4, 7.7, 7.1]
6 class_B = [4.5, 8.2, 6.0, 9.5, 5.5, 8.8, 10.0, 4.0]
7
8 # Calculate Mean and Standard Deviation
9 print(f'Class A - Mean: {np.mean(class_A):.2f}, Std: {np.std(class_A):.2f}')
10 print(f'Class B - Mean: {np.mean(class_B):.2f}, Std: {np.std(class_B):.2f}')
11
12 # Visualize distributions
13 plt.figure(figsize=(10, 4))
14 plt.subplot(1, 2, 1)
15 plt.hist(class_A, bins=5, color='blue', alpha=0.7, edgecolor='black')
16 plt.title('Class A (Low Variance)')
17 plt.xlabel('Scores')
18 plt.ylabel('Frequency')
19
20 plt.subplot(1, 2, 2)
21 plt.hist(class_B, bins=5, color='red', alpha=0.7, edgecolor='black')
22 plt.title('Class B (High Variance)')
23 plt.xlabel('Scores')
24 plt.ylabel('Frequency')
25
26 plt.tight_layout()
27 plt.show()

```



Hình 1: Phân bố điểm của 2 lớp A và B .

**Kết luận:** Lớp A ổn định, dễ dự đoán, trong khi lớp B có sự chênh lệch lớn, cần tìm hiểu thêm nguyên nhân.

## 5. So sánh tổng quan và khi nào dùng gì

### 5.1 Bảng tham khảo nhanh

Chỉ số	Khi nào dùng	Nhược điểm
<b>Mean</b>	Dữ liệu phân bố chuẩn, ít outliers	Nhạy cảm với outliers
<b>Median</b>	Dữ liệu có outliers, phân bố lệch	Bỏ qua thông tin của giá trị cực đại
<b>Variance/Std</b>	Đo lường sự biến động	Nhạy cảm với outliers

Bảng 2: So sánh các chỉ số thống kê [6]

### 5.2 Hướng dẫn chọn chỉ số

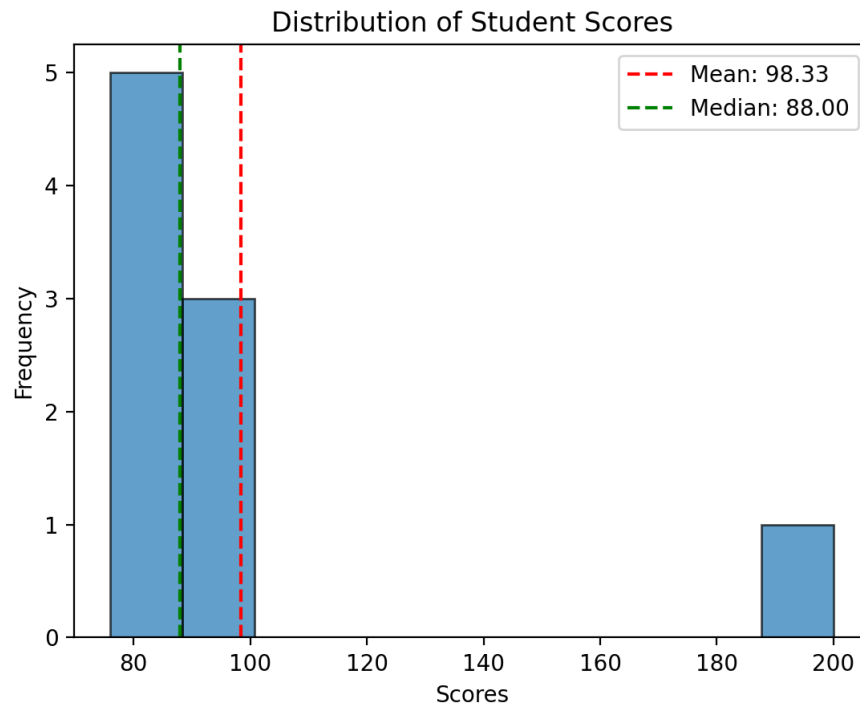
Cách chọn chỉ số phù hợp:

1. Vẽ histogram để xem phân bố dữ liệu.
2. Kiểm tra outliers bằng boxplot.
3. Quyết định:
  - Phân bố chuẩn, ít outliers? Dùng **Mean**.
  - Nhiều outliers hoặc dữ liệu lệch? Dùng **Median**.
  - Cần đo sự biến động? Dùng **Variance/Std**.

## 6. Thực hành: Phân tích dữ liệu với Python

Hãy cùng thử phân tích dữ liệu đơn giản để hiểu rõ hơn:

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4
5 # Sample dataset: student scores with an outlier
6 data = {
7     'scores': [85, 90, 78, 92, 88, 76, 94, 82, 200]
8 }
9 df = pd.DataFrame(data)
10
11 # Calculate statistics
12 mean = df['scores'].mean()
13 median = df['scores'].median()
14 std = df['scores'].std()
15
16 print(f"Mean: {mean:.2f}")
17 print(f"Median: {median:.2f}")
18 print(f"Standard Deviation: {std:.2f}")
19
20 # Visualize distribution
21 plt.hist(df['scores'], bins=10, alpha=0.7, edgecolor='black')
22 plt.axvline(mean, color='red', linestyle='--', label=f"Mean: {mean:.2f}")
23 plt.axvline(median, color='green', linestyle='--', label=f"Median: {median:.2f}")
24 plt.title('Distribution of Student Scores')
25 plt.xlabel('Scores')
26 plt.ylabel('Frequency')
27 plt.legend()
28 plt.show()
```



Hình 2: Ảnh hưởng của outlier đến giá trị Mean và Median

Đoạn code mẫu sử dụng biểu đồ hiển thị và cho thấy giá trị ngoại lai (200) làm Mean bị lệch, nhưng Median vẫn phản ánh đúng xu hướng chung.

## 7. Kết luận: Thống kê mang lại lợi ích gì?

Thống kê không chỉ là phép tính – nó là ngôn ngữ để hiểu dữ liệu và là nền tảng của AI. Mean, Median, và Variance giúp bạn:

- **Hiểu sâu dữ liệu:** Không chỉ nhìn con số mà còn thấy câu chuyện phía sau.
- **Đưa ra quyết định thông minh:** Chọn đúng chỉ số cho từng tình huống.
- **Xây dựng AI tốt hơn:** Nền tảng vững chắc dẫn đến mô hình hiệu quả.

### 7.1 Bước tiếp theo

Muốn đi xa hơn? Hãy:

1. Thực hành với dữ liệu thực tế bằng module statistics của Python [7].
2. Đọc *Pattern Recognition and Machine Learning* của Christopher Bishop [8].
3. Tìm hiểu về xác suất để làm chủ các kỹ thuật AI nâng cao.

**Hãy nhớ:** Mọi đột phá trong AI đều bắt đầu từ thống kê. Nắm vững những kiến thức cơ bản này, bạn sẽ sẵn sàng chinh phục thế giới Machine Learning!

## Tài liệu

- [1] DataCamp, “Unveiling the magic of statistical machine learning,” <https://www.datacamp.com/tutorial/unveiling-the-magic-of-statistical-machine-learning>, 2024, accessed: 2025-07-19.
- [2] E. Nalisnick, “A brief tour of deep learning from a statistical perspective,” *Annual Review of Statistics and Its Applications*, vol. 10, pp. 1–24, 2023.
- [3] GeeksforGeeks, “Statistics for machine learning,” <https://www.geeksforgeeks.org/machine-learning/statistics-for-machine-learning/>, 2024, accessed: 2025-07-19.
- [4] DataCamp, “Mean vs. median: Knowing the difference,” <https://www.datacamp.com/tutorial/mean-vs-median>, 2025, accessed: 2025-07-19.
- [5] Tutorialspoint, “Mean, median, and mode in machine learning,” [https://www.tutorialspoint.com/machine\\_learning/machine\\_learning\\_mean\\_median\\_mode.htm](https://www.tutorialspoint.com/machine_learning/machine_learning_mean_median_mode.htm), 2025, accessed: 2025-07-19.
- [6] Simplilearn, “Statistics for machine learning: A complete guide,” <https://www.simplilearn.com/tutorials/machine-learning-tutorial/statistics-for-machine-learning>, 2025, accessed: 2025-07-19.
- [7] W3Schools, “Python statistics module,” [https://www.w3schools.com/python/module\\_statistics.asp](https://www.w3schools.com/python/module_statistics.asp), 2025, accessed: 2025-07-19.
- [8] C. M. Bishop, *Pattern Recognition and Machine Learning*. Springer, 2006.

# Thống Kê Cơ Bản và Hệ Số Tương Quan: Khám Phá Mỗi Quan Hệ Dữ Liệu

*Bùi Đức Xuân*

## 1. Thống Kê Mô Tả Cơ Bản

Thống kê mô tả là nền tảng để tóm tắt và mô tả các đặc điểm chính của một tập dữ liệu.

### Giá Trị Trung Bình (Mean)

- **Trung bình tổng thể (Population Mean -  $\mu$ ):** Đây là giá trị trung bình của tất cả các phần tử trong một tập dữ liệu tổng thể. Công thức tính là  $\mu = \frac{1}{n} \sum_{i=1}^n x_i$ . Ví dụ, để tính cân nặng trung bình của toàn bộ 200,000 con chuột, chúng ta sẽ cần thu thập dữ liệu từ tất cả chúng.
- **Trung bình mẫu (Sample Mean -  $\bar{x}$ ):** Trong thực tế, việc thu thập dữ liệu từ toàn bộ tổng thể thường tốn kém và mất thời gian. Do đó, chúng ta thường ước tính trung bình tổng thể bằng cách sử dụng một mẫu nhỏ hơn. Trung bình mẫu được tính bằng  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ . Ví dụ, từ 5 mẫu cân nặng chuột (3, 13, 19, 24, 29g), trung bình mẫu là 17.6g.

### Phương Sai (Variance) & Độ Lệch Chuẩn (Standard Deviation)

- **Phương sai tổng thể (Population Variance):** Được tính bằng  $\text{var}(X) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$ . Độ lệch chuẩn ( $\sigma$ ) là căn bậc hai của phương sai tổng thể.
- **Ước tính phương sai mẫu (Sample Variance):** Khi ước tính phương sai của tổng thể từ một mẫu, nếu chúng ta chia cho 'n', chúng ta sẽ **liên tục đánh giá thấp phương sai thực tế** xung quanh giá trị trung bình của tổng thể. Để khắc phục điều này, công thức ước tính phương sai mẫu chính xác hơn là  $\text{var}(X) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$ . Độ lệch chuẩn mẫu cũng được ước tính bằng căn bậc hai của phương sai mẫu.

## 2. Bậc Tự Do (Degrees of Freedom)

Bậc tự do là **số lượng giá trị trong phép tính cuối cùng của một thống kê được tự do thay đổi** (theo Wikipedia).

- **Ví dụ minh họa:**
  - Nếu bạn tung đồng xu 100 lần và muốn biết có bao nhiêu mặt sấp và mặt ngửa, bạn chỉ cần hỏi tôi một câu: "Bạn có bao nhiêu mặt sấp?" Khi biết số mặt sấp, bạn tự động biết số mặt ngửa (100 - số mặt sấp). Ở đây, chỉ có 1 bậc tự do.
  - Nếu bạn hỏi về màu của đèn giao thông và được cho biết "nó không phải màu vàng cũng không phải màu đỏ", bạn ngay lập tức biết nó là màu xanh lá. Có 3 kết quả có thể (đỏ, vàng, xanh) nhưng chỉ cần 2 thông tin để xác định (không đỏ, không vàng), do đó có 2 bậc tự do.



- **Liên hệ với Trung bình mẫu:** Khi tính phương sai mẫu, nếu bạn đã biết giá trị trung bình của mẫu, giá trị cuối cùng trong mẫu không còn độc lập nữa. Điều này giải thích lý do tại sao chúng ta chia cho 'n-1' thay vì 'n' trong công thức phương sai mẫu, vì một giá trị đã bị "khóa" bởi trung bình mẫu.

### 3. Hiệp Phương Sai (Covariance)

Hiệp phương sai là một phép đo thống kê giúp chúng ta hiểu mối quan hệ giữa hai biến.

- **Ý tưởng chính:** Hiệp phương sai giúp nhận biết ba loại mối liên hệ chính giữa hai biến: xu hướng đồng biến (positive trend), xu hướng nghịch biến (negative trend), hoặc không có xu hướng nào.
- **Minh họa bằng ví dụ:** Hãy tưởng tượng chúng ta đếm số lượng táo xanh và táo đỏ ở các siêu thị khác nhau.
  - Nếu số lượng táo xanh và táo đỏ cùng thấp (ví dụ: Hà Nội, Hồ Chí Minh) hoặc cùng cao (ví dụ: Cần Thơ, Đà Nẵng, Bình Định) so với giá trị trung bình của chúng, điều này cho thấy một **xu hướng đồng biến**. Khi đó, Hiệp phương sai (Cov)  $> 0$ .
  - Nếu một biến tăng trong khi biến kia giảm, điều này chỉ ra một **xu hướng nghịch biến**. Khi đó, Cov  $< 0$ .
  - Nếu không có mối liên hệ rõ ràng nào, Cov sẽ gần bằng 0.
- **Hạn chế của Hiệp phương sai:** Mặc dù Hiệp phương sai cho biết hướng của mối quan hệ (đồng biến hay nghịch biến), nhưng nó **khó diễn giải và phụ thuộc vào thang đo của dữ liệu**. Giá trị Hiệp phương sai không cho biết liệu độ dốc của đường biểu diễn mối quan hệ có dốc hay không, hoặc các điểm dữ liệu có gần đường xu hướng hay không.

### 4. Hệ Số Tương Quan (Correlation Coefficient)

Hệ số tương quan khắc phục những hạn chế của Hiệp phương sai, cung cấp một thước đo chuẩn hóa và dễ diễn giải hơn về mối quan hệ giữa hai biến. Khi làm việc với các biến liên tục, hệ số tương quan thường dùng là **Pearson's r**.

- **Định nghĩa:** Hệ số tương quan, thường được biểu thị là 'r', là một phép đo hướng và độ mạnh của mối quan hệ giữa hai biến. Nó là một "bước đệm tính toán" từ Hiệp phương sai.
- **Giá trị và Ý nghĩa:** Hệ số tương quan có giá trị trong khoảng từ **-1 đến +1**.
  - **+1:** Cho thấy một **tương quan dương hoàn hảo** (perfect positive correlation), nghĩa là khi một biến tăng, biến kia cũng tăng theo một đường thẳng.
  - **-1:** Cho thấy một **tương quan âm hoàn hảo** (perfect negative correlation), nghĩa là khi một biến tăng, biến kia giảm theo một đường thẳng.
  - **0:** Cho thấy **không có mối quan hệ tuyến tính** nào giữa hai biến.
  - **Giá trị càng gần +1 hoặc -1:** Cho thấy mối quan hệ tuyến tính giữa hai biến càng mạnh. Ví dụ,  $0 < \text{Correlation} < 1$  cho thấy mối quan hệ dương nhưng không hoàn hảo.

- **Công thức (Pearson's r):**

$$\rho_{xy} = \frac{E[(x - \mu_x)(y - \mu_y)]}{\sqrt{\text{var}(x)\text{var}(y)}} = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{\sqrt{(n \sum x_i^2 - (\sum x_i)^2)(n \sum y_i^2 - (\sum y_i)^2)}}$$

- **Tính chất:**

- Phạm vi từ -1 đến +1.
- Tính đối xứng:  $\rho_{xy} = \rho_{yx}$ .
- **Không nhạy cảm với thang đo dữ liệu:** Hệ số tương quan không thay đổi khi dữ liệu được nhân với một hằng số dương ( $\rho_{x,y} = \rho_{ax,by}$  với  $a, b > 0$ ).
- **Không nhạy cảm với phép dịch chuyển:** Hệ số tương quan không thay đổi khi một hằng số được cộng vào dữ liệu ( $\rho_{x,y} = \rho_{x+c,y+d}$ ).
- **Tương quan (Correlation) và Hồi quy (Regression):** Sự khác biệt chính là tương quan đo lường **mức độ của mối quan hệ** giữa hai biến độc lập (x và y), trong khi hồi quy là cách một biến **ảnh hưởng đến** biến khác.
- **Tương quan (Correlation) và Nhân quả (Causation):** Một nguyên tắc quan trọng là **”Tương quan không phải là nhân quả”**. Điều này có nghĩa là chỉ vì hai biến có mối liên hệ với nhau không nhất thiết có nghĩa là biến này gây ra biến kia. Ví dụ, sẽ là phi đạo đức nếu thực hiện một thí nghiệm để xác định liệu hút thuốc có gây ung thư phổi hay không; tuy nhiên, mối tương quan mạnh giữa chúng có thể được quan sát một cách tự nhiên.
- **Ưu điểm của Tương quan:**
  - Cho phép nhà nghiên cứu khảo sát các biến tự nhiên mà có thể phi đạo đức hoặc không thực tế nếu thử nghiệm bằng thực nghiệm.
  - Cho phép nhà nghiên cứu thấy rõ ràng và dễ dàng liệu có mối quan hệ giữa các biến hay không, và có thể hiển thị dưới dạng đồ họa.

## 5. Ứng Dụng của Hệ Số Tương Quan

Hệ số tương quan được ứng dụng rộng rãi trong nhiều lĩnh vực, đặc biệt là trong thị giác máy tính.

- **Tìm kiếm khuôn mẫu (Template Matching) & Phát hiện vật thể (Object Detection):**
  - Trong ứng dụng này, hệ số tương quan được sử dụng để tìm một khuôn mẫu (template) cụ thể trong một hình ảnh lớn hơn.
  - Nó hoạt động bằng cách tính toán hệ số tương quan giữa khuôn mẫu và các phần khác nhau của hình ảnh. Vị trí có hệ số tương quan cao nhất (gần +1) cho thấy sự trùng khớp tốt nhất.
  - Khả năng của  $\rho$  hoạt động tốt dưới các biến đổi như dịch chuyển và thay đổi tỷ lệ (ví dụ:  $P2 = 1.2P1 + 10$  và  $P1$  và  $P2$  vẫn có  $\rho$  cao) khiến nó rất hiệu quả trong việc nhận diện các đối tượng ngay cả khi chúng xuất hiện ở các kích thước hoặc vị trí khác nhau.

## Kết Luận

Từ trung bình và phương sai mô tả dữ liệu đơn lẻ, đến hiệp phương sai và hệ số tương quan khám phá mối quan hệ giữa các biến, những công cụ thống kê này là chìa khóa để phân tích và hiểu dữ liệu. Đặc biệt, hệ số tương quan cung cấp một cách diễn giải chuẩn hóa về hướng và độ mạnh của mối quan hệ, mở ra nhiều ứng dụng thực tế.

Hãy tưởng tượng dữ liệu của bạn là một thành phố rộng lớn. Các thống kê mô tả cơ bản như trung bình hay phương sai giống như việc bạn mô tả một tòa nhà cụ thể trong thành phố đó – chiều cao của nó, diện tích của nó. Hiệp phương sai giống như việc bạn cố gắng xem liệu các tòa nhà cao tầng có xu hướng ở gần các con sông lớn hay không. Còn hệ số tương quan thì giống như việc bạn tạo ra một bản đồ chi tiết với các con đường được đánh dấu rõ ràng, cho biết không chỉ hướng đi mà còn độ chắc chắn của mỗi con đường (mối quan hệ). Điều này giúp bạn dễ dàng ”điều hướng” và đưa ra dự đoán về cách các phần khác nhau của thành phố liên kết với nhau.

# Firestore Toàn Tập: Từ NoSQL Realtime đến Machine Learning với Iris Dataset

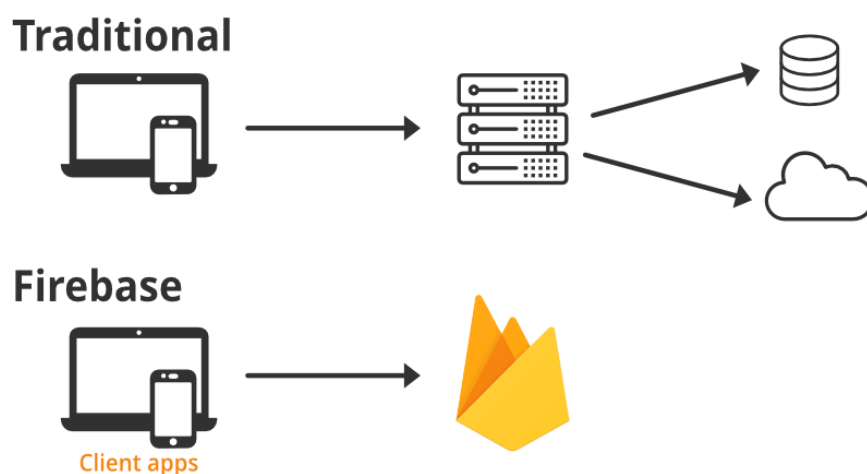
Đàm Nguyễn Khánh

## 1. Firebase là gì?

**Firebase** là một nền tảng **Backend-as-a-Service (BaaS)** do Google phát triển. Nói đơn giản, nó cung cấp cho bạn một “**hạ tầng phía sau**” gồm cơ sở dữ liệu, xác thực người dùng, lưu trữ file và các công cụ vận hành khác để bạn chỉ cần tập trung viết ứng dụng (frontend) mà không phải cài đặt server hay cấu hình phức tạp.

### Giải thích thuật ngữ

**Backend-as-a-Service (BaaS):** là dịch vụ cung cấp sẵn các chức năng backend như cơ sở dữ liệu, xác thực, lưu trữ file, API... mà không cần bạn tự triển khai hoặc bảo trì server.



Hình 3: So sánh giữa kiến trúc truyền thống (Traditional) và kiến trúc sử dụng Firebase.

## 2. Hai lựa chọn NoSQL trong Firebase: Realtime Database và Firestore

Firebase cung cấp 2 hệ quản trị cơ sở dữ liệu dạng **NoSQL**:

- **Realtime Database:** Dữ liệu được lưu dưới dạng JSON tree, phù hợp các ứng dụng realtime nhẹ.
- **Cloud Firestore:** Cấu trúc dạng **document** (tài liệu) – **collection** (bộ sưu tập), tương tự MongoDB, mạnh mẽ hơn trong truy vấn và mở rộng.

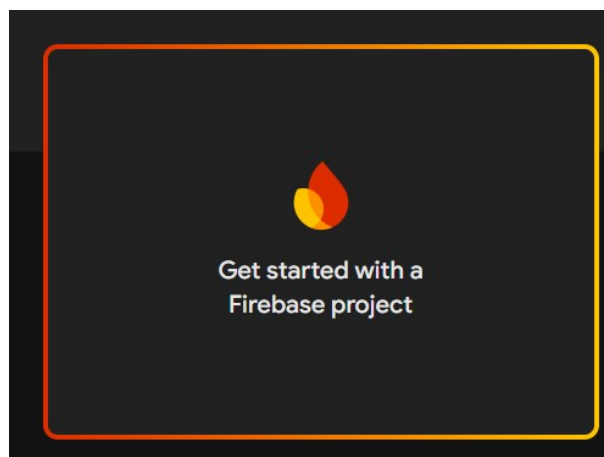
#### Giải thích thuật ngữ

**Realtime:** dữ liệu thay đổi sẽ được đồng bộ ngay lập tức tới người dùng mà không cần reload.

### 3. Hướng dẫn thiết lập Firebase Project

#### Bước 1: Tạo Project

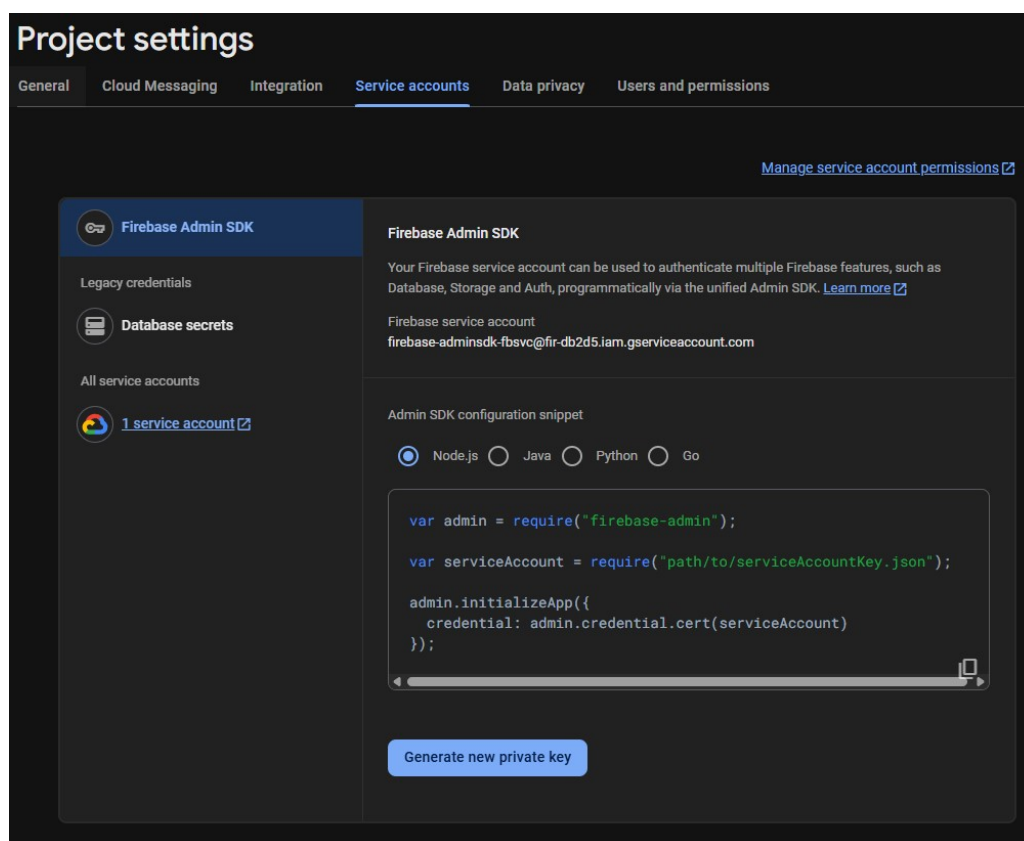
Truy cập <https://console.firebase.google.com>, chọn Add project, đặt tên và khởi tạo.



Hình 4: Giao diện khởi tạo dự án Firebase lần đầu tiên. Người dùng cần bấm vào ô **“Get started with a Firebase project”** để bắt đầu quy trình tạo mới một project trên Firebase Console. Việc tạo project này là bước khởi đầu để có thể sử dụng các dịch vụ như Firestore, Realtime Database, Authentication, Functions,... từ Firebase.

#### Bước 2: Tạo Service Account

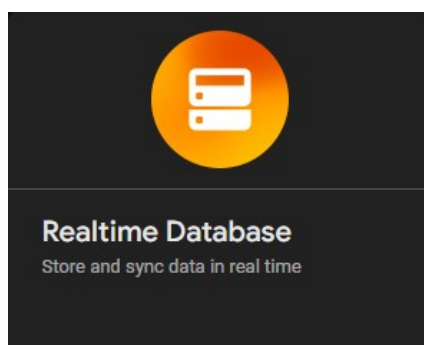
Vào Project Settings → Service Accounts → Generate new private key. Lưu file JSON – bạn sẽ dùng file này để xác thực khi kết nối Python tới Firebase.



Hình 5: Giao diện tạo **Service Account** trong mục Project Settings → Service accounts của Firebase Console. Người dùng cần bấm vào nút Generate new private key để tải về một file .json chứa thông tin xác thực. File này sẽ được dùng trong mã Python (hoặc Node.js, Java, Go) để kết nối tới Firebase thông qua SDK.

### Bước 3: Bật Realtime Database

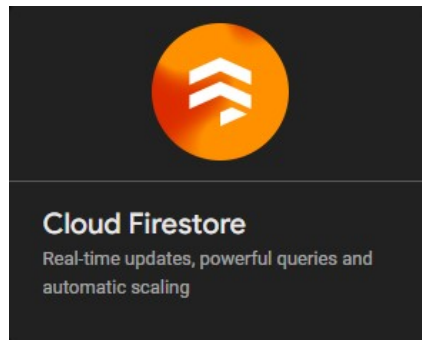
Chọn chế độ Test Mode để cho phép đọc/ghi tự do trong quá trình thử nghiệm.



Hình 6: Biểu tượng và tùy chọn kích hoạt **Realtime Database** trong Firebase. Đây là một hệ quản trị cơ sở dữ liệu NoSQL dạng cây JSON, cho phép đồng bộ dữ liệu theo thời gian thực giữa nhiều client. Khi bật tính năng này lần đầu, Firebase sẽ yêu cầu chọn khu vực (region) và chính sách bảo mật (Security Rules). Trong giai đoạn thử nghiệm, bạn có thể chọn test mode để dễ dàng đọc/ghi mà không bị giới hạn quyền. *Lưu ý:* test mode chỉ nên dùng cho mục đích phát triển, không nên dùng trong môi trường thật.

## Bước 4: Bật Cloud Firestore

Chọn vùng (region) phù hợp, nên trùng với Realtime DB để giảm độ trễ và chi phí.



Hình 7: Biểu tượng và tùy chọn kích hoạt **Cloud Firestore** trong Firebase. Firestore sử dụng mô hình Document – Collection rất giống MongoDB, phù hợp cho các ứng dụng web, mobile và cả backend AI.

Firestore là hệ cơ sở dữ liệu NoSQL hiện đại của Google, hỗ trợ:

- **Realtime updates:** dữ liệu cập nhật theo thời gian thực.
- **Powerful queries:** hỗ trợ truy vấn theo điều kiện, sắp xếp, phân trang.
- **Automatic scaling:** tự mở rộng để phục vụ lượng người dùng lớn.

## 4. Kết nối Firestore từ Python

### Khởi tạo kết nối

```
1 import firebase_admin
2 from firebase_admin import credentials, firestore
3
4 cred = credentials.Certificate("serviceAccountKey.json")
5 firebase_admin.initialize_app(cred)
6 db = firestore.client()
```

### Các thao tác CRUD cơ bản

CRUD viết tắt của:

- **Create:** thêm dữ liệu
- **Read:** đọc dữ liệu
- **Update:** cập nhật dữ liệu
- **Delete:** xóa dữ liệu

Các ví dụ:

- Thêm document: `db.collection("users").add({...})`
- Truy vấn: `.where("age", ">", 18)`
- Cập nhật: `doc_ref.update({...})`
- Xóa: `doc_ref.delete()`

## 5. Bộ Dữ liệu Iris và Cách Lưu vào Firestore

### Giới thiệu về bộ dữ liệu Iris

Bộ dữ liệu **Iris** là một trong những dataset kinh điển trong Machine Learning. Mỗi mẫu là thông tin của một bông hoa thuộc 1 trong 3 loài:

- Iris Setosa
- Iris Versicolor
- Iris Virginica

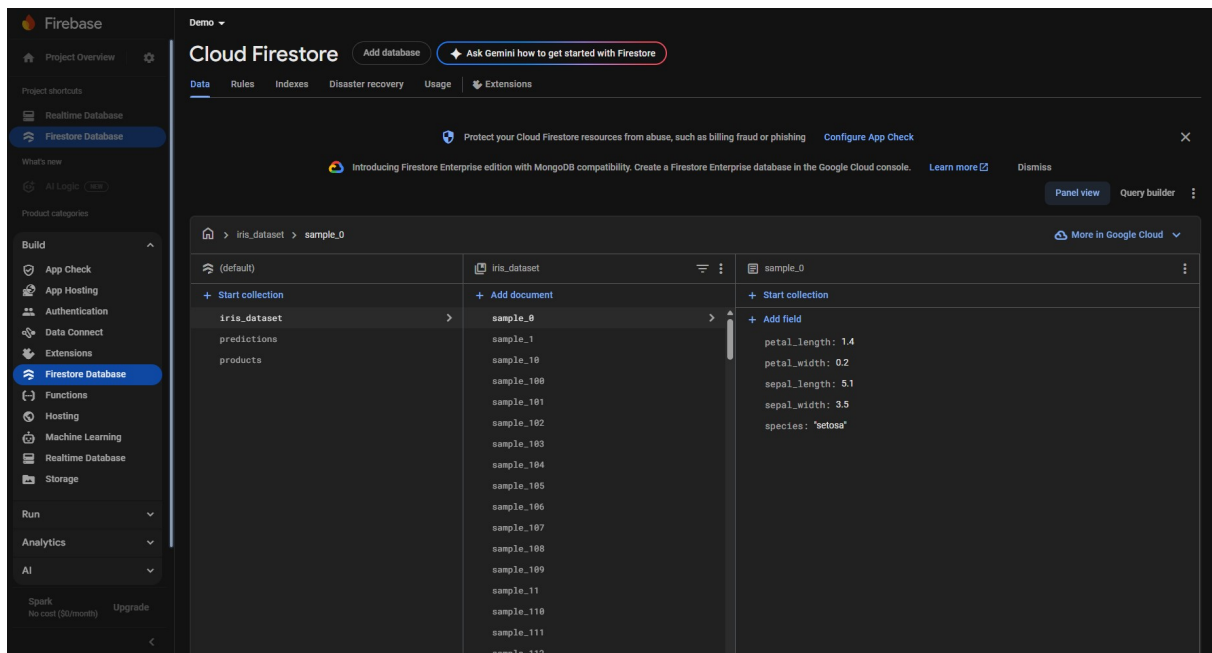
Các đặc trưng gồm:

- Sepal length (chiều dài đài hoa)
- Sepal width (chiều rộng đài hoa)
- Petal length (chiều dài cánh hoa)
- Petal width (chiều rộng cánh hoa)

### Lưu dữ liệu vào Firestore

Mỗi mẫu sẽ là một **document** trong **collection** tên `iris_samples`. Tạo script Python để duyệt từng dòng của `pandas.DataFrame` và thêm vào Firestore bằng `.add()` hoặc `.set()`.





Hình 8: Giao diện quản lý dữ liệu của **Cloud Firestore** trong Firebase.

Hình ảnh hiển thị:

- Một **collection** có tên là `iris_dataset`, chứa các **documents** dạng `sample_0`, `sample_1`, ...
- Mỗi document đại diện cho một mẫu hoa trong bộ dữ liệu Iris.
- Bên phải, các trường dữ liệu (fields) bao gồm:
  - `petal_length`, `petal_width`, `sepal_length`, `sepal_width`: các đặc trưng đo lường của hoa.
  - `species`: tên loài hoa, ví dụ "setosa".

Đây là cách dữ liệu dạng document được tổ chức trong Firestore: không cần bảng, mỗi document là một bản ghi độc lập, các fields linh hoạt về kiểu và số lượng.

## 6. Machine Learning với Iris trên Firebase

### Tải dữ liệu về để huấn luyện

Sử dụng Firestore API để đọc toàn bộ dữ liệu, convert sang `pandas.DataFrame` để huấn luyện mô hình.

### Huấn luyện mô hình

Các mô hình bạn có thể thử:

- `DecisionTreeClassifier`
- `RandomForestClassifier`
- `SVC` (Support Vector Machine)

- LogisticRegression

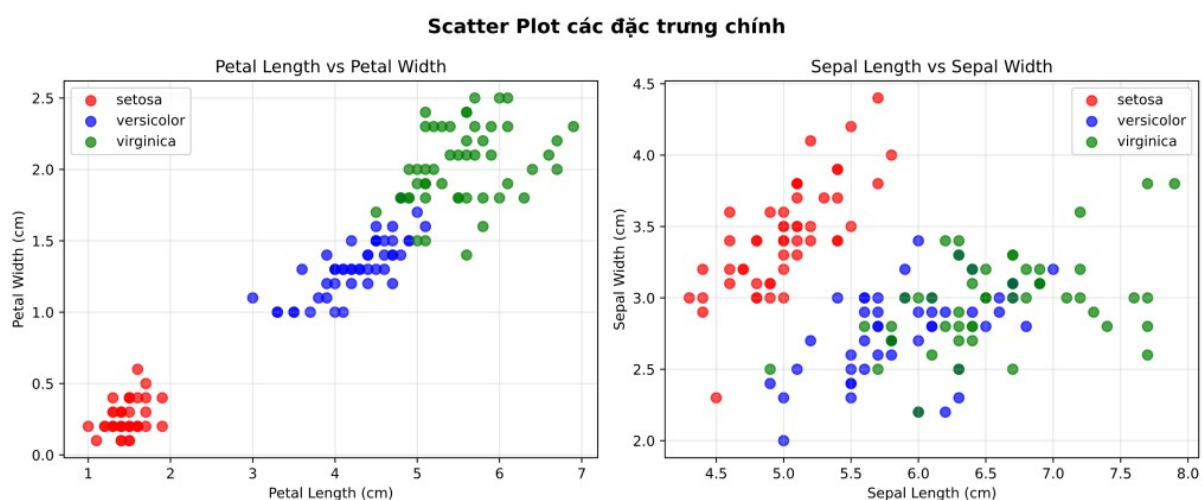
Chia tập train/test, đánh giá độ chính xác, chọn mô hình tốt nhất để sử dụng.

Gợi ý chèn hình: bảng so sánh độ chính xác các mô hình.

## Trực quan hóa dữ liệu

Vẽ biểu đồ **scatter plot** để thấy sự phân tách giữa các loài.

- Petal Length vs Petal Width
- Sepal Length vs Sepal Width



Hình 9: Biểu đồ **Scatter Plot** thể hiện sự phân bố các đặc trưng hình thái trong bộ dữ liệu **Iris**, phân loại theo ba loài hoa: setosa, versicolor, và virginica.

**Biểu đồ bên trái** so sánh **Petal Length** và **Petal Width** (chiều dài và chiều rộng cánh hoa):

- Cho thấy sự phân tách rõ ràng giữa ba loài.
- Setosa (đỏ) nằm hoàn toàn tách biệt.
- Versicolor (xanh dương) và Virginica (xanh lá) có một số vùng chồng lấn.

**Biểu đồ bên phải** so sánh **Sepal Length** và **Sepal Width** (chiều dài và chiều rộng đài hoa):

- Sự phân chia giữa các loài kém rõ ràng hơn.
- Các cụm màu không hoàn toàn tách biệt, đặc biệt giữa versicolor và virginica.

Biểu đồ này minh họa tầm quan trọng của việc chọn đúng đặc trưng khi huấn luyện mô hình phân loại.

## 7. Xây dựng Giao diện với Streamlit

**Streamlit** là framework siêu đơn giản để tạo web app bằng Python, rất thích hợp để demo AI/ML.

Chức năng chính:

- Form nhập 4 thông số hoa
- Dự đoán loài
- Ghi kết quả vào Firestore

## 8. Kết luận

Firebase kết hợp với Python và Streamlit giúp bạn xây dựng hệ thống AI nhỏ gọn, dễ triển khai và không cần lo về hạ tầng.

Với kiến thức trong bài viết, bạn có thể:

- Hiểu rõ về hai hệ NoSQL của Firebase
- Lưu trữ dữ liệu AI lên Firestore
- Huấn luyện và đánh giá mô hình Machine Learning
- Tạo app web dự đoán và kết nối cơ sở dữ liệu thực tế

# Thiết Kế Slide & Kể Chuyện Hiệu Quả

*Đàm Nguyên Khánh*

## Giới thiệu

Trong lĩnh vực khoa học dữ liệu và trí tuệ nhân tạo (AI), khả năng trình bày dữ liệu và kể chuyện (storytelling) là kỹ năng quan trọng không kém việc xây dựng các mô hình phân tích. Bài viết này sẽ cung cấp cho bạn những nguyên tắc thiết kế slide chuyên nghiệp và kỹ thuật kể chuyện dữ liệu để trình bày thông tin một cách hiệu quả, thuyết phục và hấp dẫn.

## 1. Tầm Quan Trọng của Kỹ Năng Thuyết Trình

- **Học tập:** Nổi bật khi báo cáo, tạo cơ hội thực tập hoặc nghiên cứu.
- **Tìm việc:** Gây ấn tượng trong phỏng vấn, thể hiện sự tự tin.
- **Công việc:** Tạo tác động, truyền đạt ý tưởng hiệu quả, thuyết phục đầu tư.



Hình 10: Một buổi thuyết trình dữ liệu hiệu quả: Diễn giả cần sử dụng biểu đồ rõ ràng, slide gọn gàng, tương tác tốt với khán giả.

## 2. Những Lỗi Thường Gặp trong Thuyết Trình

### 2.1. Lỗi Thiết Kế Slide

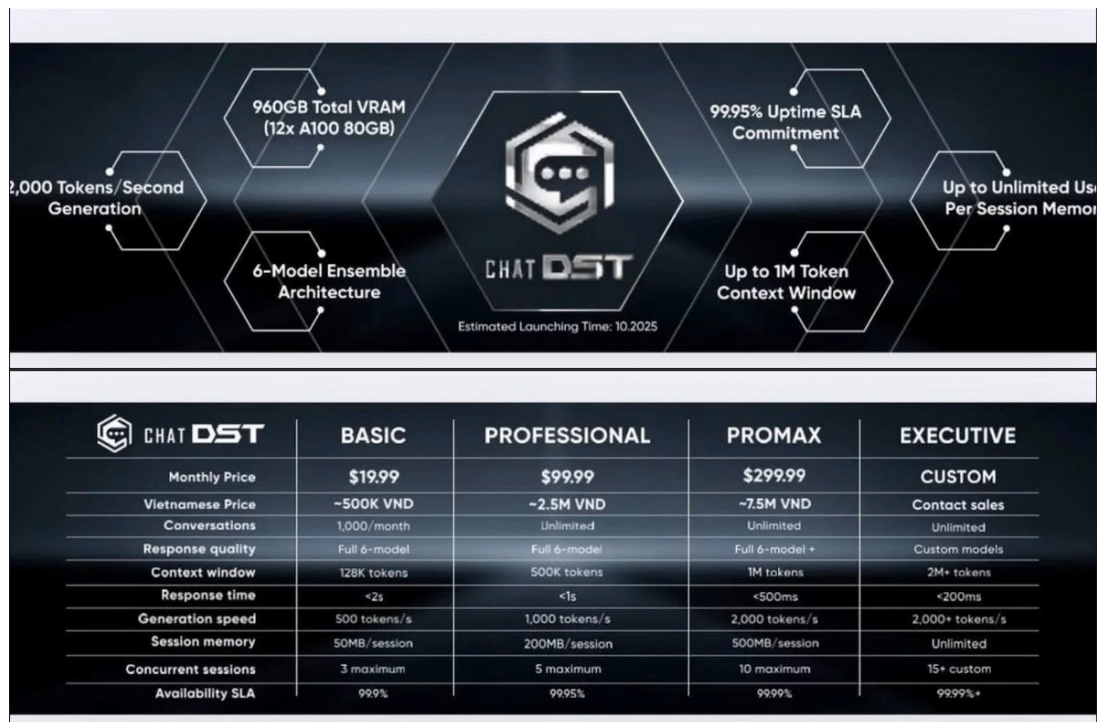
- Quá nhiều chữ (do copy-paste từ tài liệu).
- Màu sắc lộn xộn, thiếu tính nhất quán.
- Biểu đồ phức tạp, thiếu nhãn và đơn vị.



Hình 11: Ví dụ về một slide không hiệu quả: Quá nhiều chữ nhỏ, biểu đồ phức tạp và dày đặc, màu sắc chưa có định hướng rõ ràng khiến khán giả khó tập trung và ghi nhớ thông tin chính.

## 2.2. Lỗi Kể Chuyện

- Thiếu cấu trúc rõ ràng.
- Sử dụng quá nhiều thuật ngữ kỹ thuật không giải thích.
- Thiếu tương tác với người nghe.

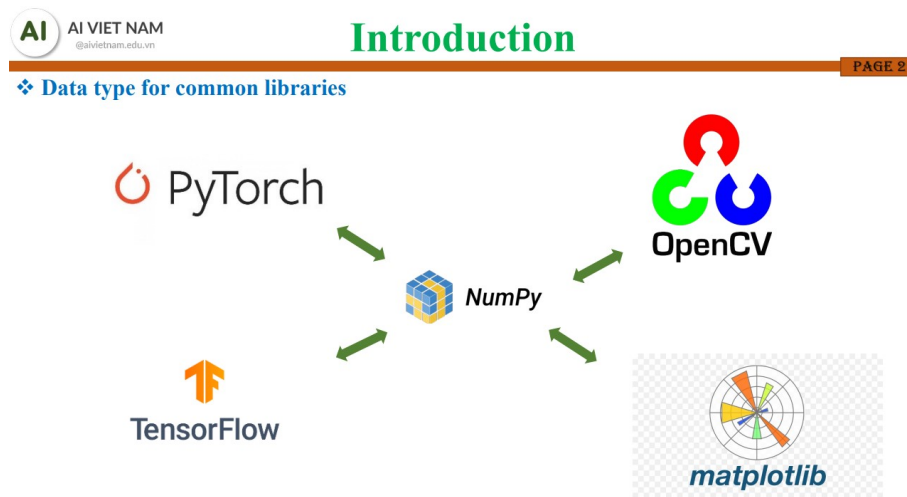


Hình 12: Ví dụ về lỗi trình bày: Slide chứa quá nhiều thuật ngữ kỹ thuật chuyên sâu như ”token”, ”context window”, ”6-model ensemble” khiến người nghe không chuyên khó tiếp cận và không hiểu được giá trị thật sự của sản phẩm.

### 3. Thiết Kế Slide Chuyên Nghiệp

#### 3.1. Bốn Nguyên Tắc Vàng

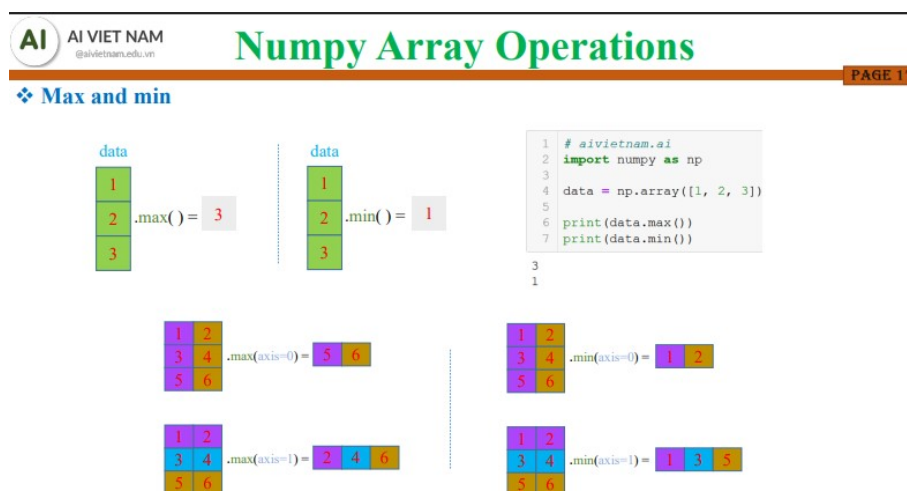
- **Đơn giản (Simplicity):** Mỗi slide chỉ chứa một thông điệp chính.
- **Nhất quán (Consistency):** Thống nhất phong cách thiết kế.
- **Tương phản (Contrast):** Nhấn mạnh nội dung chính.
- **Cân bằng (Balance):** Bố cục hài hòa, dễ nhìn.



Hình 13: Slide minh họa tốt mối quan hệ giữa NumPy và các thư viện phổ biến trong lĩnh vực AI như PyTorch, TensorFlow, OpenCV và Matplotlib. Thiết kế đơn giản (chỉ truyền đạt một ý chính), sử dụng phong cách nhất quán (màu sắc và font chữ), có tương phản rõ ràng (biểu tượng nổi bật trên nền trắng), và bố cục cân bằng (NumPy ở trung tâm, các thư viện phân bố đều xung quanh).

### 3.2. Kỹ Thuật Cụ Thể

- Quy tắc 5-5-5: 5 dòng, 5 từ/dòng, tối đa 5 slide text liên tiếp.
- Typography: Ưu tiên font sans-serif, kích thước hợp lý.
- Màu sắc: Quy tắc 60-30-10 cho nền, phụ, và nhấn.
- White Space: Ít nhất 20

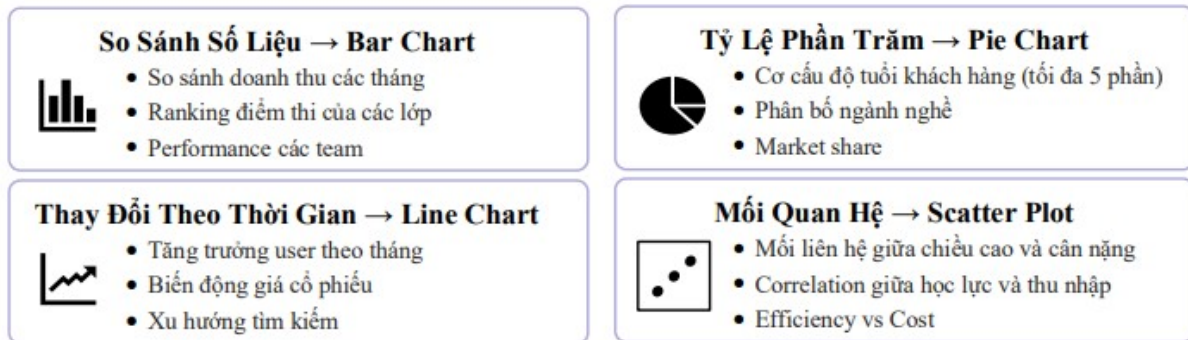


Hình 14: Ví dụ áp dụng quy tắc 60-30-10 hiệu quả: 60% màu nền trắng tạo cảm giác dễ chịu và chuyên nghiệp, 30% màu phụ (xanh dương và xanh lá) dùng cho nhãn và tiêu đề phụ để phân nhóm rõ ràng, và 10% màu nhấn (cam/đỏ) giúp làm nổi bật giá trị kết quả quan trọng. Thiết kế này giúp người học nhanh chóng nắm được nội dung chính mà không bị quá tải thông tin.



### 3.3. Thiết Kế Biểu Đồ Hiệu Quả

- Bar chart: So sánh.
- Line chart: Xu hướng thời gian.
- Pie chart: Phân bổ (tối đa 5 phần).
- Scatter plot: Mối quan hệ giữa hai biến.



Hình 15: Tổng hợp 4 loại biểu đồ dữ liệu thường dùng: Bar Chart (so sánh số liệu), Line Chart (xu hướng theo thời gian), Pie Chart (phân bổ phần trăm) và Scatter Plot (mối quan hệ giữa các biến). Mỗi biểu đồ được minh họa bằng biểu tượng rõ ràng và kèm ví dụ thực tế, giúp người dùng lựa chọn đúng loại biểu đồ cho từng mục tiêu phân tích.

## 4. Kể Chuyện Bằng Dữ Liệu

### 4.1. Framework SCQA

- Situation: Mô tả hiện trạng.
- Complication: Nêu ra vấn đề.
- Question: Câu hỏi trung tâm.
- Answer: Giải pháp cụ thể.

### 4.2. Logical Tree

Cấu trúc phân cấp logic theo nguyên tắc MECE:

- Không trùng lặp, bao phủ toàn bộ vấn đề.
- Quy tắc 3: mỗi nhánh có 2-4 phân nhánh, ưu tiên 3 ý chính.

## 5. Kỹ Thuật Thuyết Trình

### 5.1. Kiểm Soát Giọng Nói

- Tốc độ hợp lý (140-160 từ/phút).
- Ngắt nghỉ đúng lúc.





Hình 16: Sơ đồ cây quyết định (Decision Tree) minh họa cách tổ chức logic giữa các lựa chọn: “Mở cửa hàng” hoặc “Không mở cửa hàng”, và các kết quả tương ứng trong từng tình huống (bùng nổ, suy thoái, hoặc giữ nguyên). Đây là một ví dụ điển hình của Logical Tree giúp người ra quyết định phân tích rủi ro - lợi ích một cách có hệ thống.

- Âm lượng, ngữ điệu thay đổi để tạo điểm nhấn.

## 5.2. Ngôn Ngữ Cơ Thể

- Tư thế thẳng, tay tự nhiên.
- Eye contact với khán giả.
- Biểu cảm khuôn mặt linh hoạt.

## 5.3. Kết Nối Với Khán Giả

- Mở đầu ấn tượng bằng câu hỏi, số liệu, câu chuyện.
- Đặt câu hỏi tương tác.
- Xử lý câu hỏi chuyên nghiệp.

## 5.4. Xử Lý Lo Lắng Khi Thuyết Trình

- Chuẩn bị kỹ, luyện tập thường xuyên.
- Thở sâu, tưởng tượng tích cực.
- Sai sót nhỏ: bình tĩnh xử lý, tiếp tục trình bày.



Hình 17: Hình ảnh minh họa một diễn giả thuyết trình đầy tự tin trước đám đông, áp dụng hiệu quả các kỹ thuật trình bày: tư thế vững vàng, cử chỉ tay mở rộng thể hiện sự làm chủ, eye contact mạnh mẽ và phong thái quyết đoán. Đây là ví dụ điển hình cho việc kết hợp tốt ngôn ngữ cơ thể, giọng nói và kết nối với khán giả để tạo ảnh hưởng mạnh mẽ.

## Tổng Kết

Kỹ năng trình bày và kể chuyện không chỉ giúp nhà khoa học dữ liệu trình bày phân tích rõ ràng mà còn tạo sức ảnh hưởng đến quyết định kinh doanh. Hãy bắt đầu từ những nguyên tắc đơn giản và thực hành thường xuyên để nâng cao khả năng thuyết trình của bạn.

*"Your ideas are only as good as your ability to communicate them."*