# Tool by Nguyễn Vạn Phúc Huy

## Matplotlib

### 📊 1. Vẽ Đồ Thị Cơ Bản

- import matplotlib.pyplot as plt: Nhâp thư viên pyplot.
- plt.style.use("style name"): Chon kiểu hiển thi (vd: "ggplot", ("seaborn-v0 8"), ("dark background").
- plt.plot(y): Vẽ đồ thi với y là tung đô, hoành đô là chỉ số (index).
- plt.plot(x, y): Vẽ đổ thị với x là hoành độ, y là tung độ.
- plt.show(): Hiển thị đồ thị vừa vẽ.

### 2. Tùy Chỉnh Trục & Tiêu Đề

- plt.xlim(min, max): Đặt giới han cho truc hoành.
- [plt.ylim(min, max)]: Đặt giới hạn cho trục tung (có thể dùng None cho một giới hạn).
- plt.title("Tiêu đề"): Thêm tiêu để cho đồ thi.
- plt.xlabel("Nhãn trục X"): Thêm nhãn cho truc hoành.
- plt.ylabel("Nhãn trục Y"): Thêm nhãn cho trục tung.
- plt.grid(True): Hiển thi lưới (hoặc False để ẩn).
- [plt.xticks(ticks\_list, labels\_list, rotation=45)]: Tùy chỉnh các vạch chia (ticks) và nhãn (labels) trên trục X, có thể xoav nhãn.
- plt.yticks(ticks\_list, labels\_list): Tương tự cho trục Y.

### 🤲 3. Định Dạng Đường, Màu & Marker

- plt.plot(x, y, format string): Đinh dang nhanh (ít dùng).
  - ∘ Ví du: 'b--' (blue dashed line xanh dương, nét đứt).
  - ∘ Ví du: 'r:' (red dotted line đỏ, nét chấm).
  - o Ví dụ: 'g^' (green triangles xanh lá, marker hình tam giác).
- · Cách dùng rõ ràng (nên dùng):
  - o color='blue' (hoặc '#0000FF', 'b'): Chỉ định màu sắc.
  - linestyle='--' (hoặc ls): Kiểu đường ('-', '--', ':', '-.').
  - linewidth=2.5 (hoặc lw): Đô dày của đường.
  - marker='o': Kiểu điểm đánh dấu ('o', '^', 's', '+', 'x').
  - o (markersize=8) (hoặc (ms)): Kích thước của marker.
- Ví du đầy đủ: plt.plot(x, y, color='green', linestyle=':', marker='o', lw=2)

# 🛂 4. Kích Thước & Đồ Thị Con (Cách Cơ Bản)

- plt.figure(figsize=(width, height)): Tạo một Figure mới với kích thước tùy chỉnh (đơn vị inch). (Lưu ý: plt.subplots() ở muc 8 là cách làm gon hơn).
- plt.subplot(rows, cols, index): Chia figure thành lưới rows x (cols) và chọn đổ thị con thứ index (bắt đầu từ
   1) để vẽ lênh plt.plot() tiếp theo.
  - o Ví dụ: plt.subplot(2, 2, 1) (Lưới 2x2, chọn ô thứ 1 trái trên).
  - Ví dụ: plt.subplot(2, 1, 2) (Lưới 2x1, chọn ô thứ 2 dưới).
- [plt.tight\_layout()]: Tự động căn chỉnh các đồ thị con cho vừa vặn, tránh chồng chéo.

## 5. Thêm Chữ, Chú Thích & Nhãn

- plt.text(x, y, "Nội dung"): Thêm chữ tại tọa độ (x, y).
  - Tham số tùy chon: fontsize, color, horizontalalignment='center', rotation.
- plt.plot(..., label="Tên đường"): Đặt nhãn cho một đường vẽ (dùng cho legend).
- plt.legend(loc="vi trí"): Hiển thi bảng chú giải các nhãn.
  - o Ví dụ: loc="lower right", loc="best" (tự động chọn vị trí tốt nhất).
- [plt.annotate("Chú thích", xy=(x\_point, y\_point), xytext=(x\_text, y\_text), arrowprops=dict(facecolor='black')): Thêm chú thích có mũi tên chỉ từ chữ (xytext) đến điểm (xy).

## 6. Các Loại Đồ Thị Khác

- plt.hist(data, bins=n): Vẽ biểu đồ tần suất (histogram) với n khoảng (bins).
- plt.scatter(x, y): Vẽ biểu đồ phân tán (scatter plot).
  - o Tham số tùy chon: s (kích thước điểm), c (màu sắc, có thể là 1 mảng giá trì), alpha (đô trong suốt).
- plt.boxplot(data): Vẽ biểu đồ hộp (boxplot), (data) có thể là list các mảng.
- plt.bar(categories, values): Vẽ biểu đồ côt đứng.
- (plt.barh(categories, values): Vẽ biểu đồ cột ngang.
- plt.pie(sizes, labels=labels\_list, autopct='%1.1f%%'): Ve biểu đổ tròn.
- [plt.imshow(matrix\_data, cmap='viridis')]: Hiển thị ma trận 2D (heatmap, ảnh).
- (plt.colorbar()): Thêm thanh màu (thường dùng với (imshow) hoặc (scatter) khi (c) là một mảng giá trị).

### 💾 7. Lưu Hình Vẽ

- (plt.savefig("tên\_file.png"): Lưu đồ thị thành file ảnh.
  - o Tham số tùy chọn:
  - o dpi=300: (Dots Per Inch) Tăng độ phân giải cho ảnh (chuẩn in ấn).
  - o transparent=True : Nền trong suốt.
  - bbox inches='tight': Cắt bỏ các khoảng trắng thừa xung quanh hình.
- Vi du: plt.savefig("my figure.png", dpi=300, bbox inches='tight')

## 👑 8. API Hướng Đối Tượng (OO) - Cách Dùng Chuẩn

Đây là cách dùng linh hoạt và được khuyên dùng, đặc biệt khi xử lý nhiều đổ thị con. Thay vì (plt.plot()), (plt.title()), ta gọi phương thức từ đối tượng (ax) (Axes).

- Figure (fig): Toàn bộ cửa sổ chứa hình.
- Axes (ax): Từng ô đồ thi con (nơi thực sư vẽ).

### A. Một đồ thi con (Cách dùng chuẩn)

```
# Tạo 1 Figure và 1 Axes

fig, ax = plt.subplots(figsize=(10, 6))

# Về trên ax

ax.plot(x, y, label="Dữ liệu", color='red')

# Tùy chính trên ax (thường có 'set_' ở đầu)

ax.set_title("Tiêu đề (cách 00)")

ax.set_xlabel("Nhân X")

ax.set_ylabel("Nhân Y")

ax.set_xlim(0, 10)

ax.grid(True)

ax.legend()

plt.show() # Vẫn dùng plt.show() để hiển thị
```

### **Pandas**

## 🛓 1. Nhập & Xuất Dữ Liệu

- pd.read csv('file.csv'): Doc file CSV.
  - Tham số: (sep=','), (header=0), (names=['a', 'b']), (index\_col='col\_name')
- pd.read excel('file.xlsx', sheet name='Sheet1'): Doc file Excel.
- pd.read json('file.json'): Đoc file JSON.
- pd.read sql(query, connection): Đoc từ database SQL.
- pd.DataFrame(data\_dict): Tao DataFrame từ dictionary (ví dụ: ({'col1': [1,2], 'col2': [3,4]}).
- df.to csv('file.csv', index=False): Lưu DataFrame thành file CSV (không lưu côt index).
- df.to excel('file.xlsx', sheet\_name='Data'): Lưu DataFrame thành file Excel.
- df.to json('file.json', orient='records'): Luu DataFrame thành file JSON.

### 2. Khám Phá Nhanh (Inspection)

- df.head(n): Xem n dòng đầu tiên (mặc định là 5).
- df.tail(n): Xem n dòng cuối cùng (mặc định là 5).
- df.sample(n): Lấy ngẫu nhiên (n) dòng.
- df. shape: Trả về (số dòng, số cột).
- df.info(): Thông tin tổng quan (kiểu dữ liệu, số lượng non-null, bộ nhớ).
- df.describe(): Thống kê mô tả cho các cột số (count, mean, std, min, max, quartiles).
- df.describe(include='object'): Thống kê mô tả cho các côt object/categorical (count, unique, top, freq).
- df.columns : Hiển thị danh sách các tên cột.
- df.index: Hiển thị thông tin về index (hàng)
- df.dtypes: Hiển thị kiểu dữ liệu (dtype) của mỗi cột
- df['col'].value\_counts(): Đếm các giá trị duy nhất trong một cột (Series)

## of 3. Chọn & Lọc Dữ Liệu (Selection & Filtering)

#### · Chọn Cột:

- o (df['col']) hoặc (df.col): Chọn một cột (trả về Series).
- o df[['col1', 'col2']]: Chon nhiều côt (trả về DataFrame).

#### · Chọn Dòng (theo nhãn/vị trí):

- df.loc[label]: Chọn một dòng theo nhãn index (ví dụ: df.loc[0]).
- o df.loc['label1':'label5']: Chon một khoảng dòng theo nhãn (slice, bao gồm cả điểm cuối)
- df.iloc[vi\_trí]: Chọn một dòng theo vị trí số nguyên (ví dụ: df.iloc[0]).
- o df.iloc[0:5]: Chon 5 dòng đầu tiên (slice, không bao gồm điểm cuối, giống Python list).

#### · Chọn Dòng & Cột:

- df.loc[row\_labels, col\_labels]: Chọn theo nhãn.
  - *Ví dụ:* (df.loc[0:5, ['name', 'age']]
- df.iloc[row\_positions, col\_positions]: Chọn theo vị trí.
  - Ví du: df.iloc[0:5, [0, 2]]

#### · Loc theo điều kiện (Boolean Indexing):

- df[df['age'] > 30]: Lọc các dòng thỏa mãn điều kiện.
- o df[df['country'].isin(['USA', 'Canada'])]: Loc các dòng có giá trị nằm trong một danh sách.
- o df[(df['age'] > 30) & (df['gender'] == 'Male')]: Lọc với nhiều điều kiện (dùng & (AND), ∫ (OR), ⊘ (NOT))
- o (df.query('age > 30 and gender == "Male"'): Cách lọc khác dùng chuỗi truy vấn (query string).

### 🖴 4. Làm Sach & Biến Đổi Dữ Liêu (Cleaning & Transformation)

#### • Xử lý NaN (Missing Values):

- o df.isna().sum(): Đếm số lượng giá tri NaN trong mỗi côt.
- o df.dropna(): Xóa tất cả các dòng có chứa NaN.
- o df.dropna(axis=1): Xóa tất cả các cột có chứa NaN.
- df.fillna(value): Điền giá trị (ví dụ: 0, df['col'].mean()) cho các ô bị thiếu (NaN).

#### · Thay Đổi Dữ Liệu:

- o df.rename(columns={'old name': 'new name'}): Đổi tên côt.
- o df.replace(old\_value, new\_value): Thay thế giá trị.
- df['col'].astype('new type'): Thay đổi kiểu dữ liêu của côt (ví du: 'int', 'float', 'category').

#### Xóa Dòng/Côt:

- df.drop(columns=['col\_name']): Xóa cột.
- df.drop(labels=[index\_label]): Xóa dòng theo index

#### Duplicates

- o df.duplicated(): Kiểm tra các dòng trùng lặp (trả về Boolean Series).
- (df.drop\_duplicates(subset=['col\_name'], keep='first'): Xóa các dòng trùng lặp (chỉ giữ lại dòng đầu tiên)

#### • Áp dụng hàm (Apply/Map):

- o df['col'].map(function or dict): Áp dung hàm hoặc dictionary cho từng phần tử của một Cột (Series).
- o df.apply(function, axis=0): Áp dung hàm lên từng côt (trả về 1 giá tri cho mỗi côt).
- o df.apply(function, axis=1): Áp dung hàm lên từng dòng (trả về 1 giá trị cho mỗi dòng).
  - Ví du (lambda): (df['new\_col'] = df['old\_col'].apply(lambda x: x \* 2)

### 🔀 5. Sắp Xếp & Sửa Index

- df.sort values(by='col name', ascending=True): Sắp xếp DataFrame theo giá trị của một cột.
- (df.sort\_values(by=['col1', 'col2']): Sắp xếp theo nhiều cột.
- df.sort index(): Sắp xếp DataFrame theo index.
- (df.reset\_index()): Chuyển index hiện tại thành một cột mới và dùng index số nguyên mặc định.
- df.set index('col name'): Đặt một cột làm index của DataFrame.

## 📊 6. Thống Kê & Groupby

#### Thống Kê Cơ Bản:

- .mean(), .std(), .median(): Trung bình, đô lệch chuẩn, trung vi.
- (.sum(), .count(): Tổng, số lượng (không tính NaN).
- o .max(), .min(): Giá tri lớn nhất, nhỏ nhất.
- o .nunique(): Đếm số lượng giá trị duy nhất.
- o .corr(): Tính ma trận tương quan giữa các cột số
- (quantile(x)): Tính tứ phân vị ở x%

#### · Grouping:

- df.groupby('col'): Nhóm dữ liêu theo các giá tri của một cột.
- o df.groupby(['col1', 'col2']): Nhóm theo nhiều cột.
- o df.groupby('col').sum(): Tính tổng cho mỗi nhóm.
- df.groupby('col')['data\_col'].mean(): Tính trung bình của cột 'data\_col' cho mỗi nhóm.
- o df.groupby('col').agg(func\_dict); Áp dụng nhiều hàm thống kê cùng lúc.
  - Vi du: df.groupby('group').agg({'data1': 'mean', 'data2': 'sum'})

## 7. Hợp Nhất & Kết Nối (Merge & Concat)

- pd.concat([df1, df2]): Nối các DataFrame theo chiều doc (stack rows, axis=0).
- pd.concat([df1, df2], axis=1): Nối các DataFrame theo chiều ngang (ghép côt, axis=1).
- pd.merge(df left, df right, on='key col'): Kết nối (join) hai DataFrame dưa trên một cột chung (key).
  - o Tham số how (cách join):
    - (how='inner') (mặc định): Chỉ giữ lại các hàng có key chung ở cả hai bảng.
    - (how='left'): Giữ tất cả các hàng của bảng trái, điền NaN nếu không có key ở bảng phải.
    - (how='right': Giữ tất cả các hàng của bảng phải.
    - how='outer': Giữ tất cả các hàng của cả hai bảng.

### 8. Pivot Tables

- (df.pivot\_table(values='data\_col', index='row\_col', columns='col\_col', aggfunc='mean')):
  - Tao môt bảng tổng hợp (giống Excel Pivot Table).
  - o values : Côt dữ liêu để tính toán.
  - index: Cột để dùng làm các hàng của bảng pivot.
  - o columns: Cột để dùng làm các cột của bảng pivot.
  - aggfunc: Hàm thống kê (ví dụ: 'mean', 'sum', 'count').

### 9. Dữ Liệu Thời Gian (Time Series)

- pd.to\_datetime(df['date\_col']): Chuyển đổi cột sang kiểu dữ liệu datetime.
- Sau khi đặt cột datetime làm index (dùng df.set\_index('date\_col')):
- df['2025']: Chọn tất cả dữ liệu trong năm 2025.
- df['2025-10']: Chọn tất cả dữ liệu trong tháng 10/2025.
- df.resample('M').mean(): Tính trung bình dữ liêu theo mỗi tháng ('M'), ('D' ngày, 'W' tuần, 'Q' quý).

### 10. Trưc Quan Hóa (Tích hợp)

- df.plot(kind='line'): Vẽ biểu đồ đường (mặc định)
- df.plot(kind='bar'): Vẽ biểu đồ cột đứng.
- df.plot(kind='barh'): Vẽ biểu đổ côt ngang
- df.plot(kind='scatter', x='col1', y='col2'): Vẽ biểu đổ phân tán
- df['col'].hist(bins=20): Vẽ biểu đồ histogram cho một cột.
- df.boxplot(): Vẽ biểu đồ hộp cho các cột số.

# Phân Phối Xác Suất (Scipy Stats)

### 1. 📖 Giới Thiêu

- Xác Suất: Đại lượng dùng để đo khả năng xảy ra của một sư kiện (event), có giá tri từ 0 đến 1.
- Phân Phối Xác Suất (PPXS): Thể hiện xác suất xuất hiện của các giá trị của một biến ngẫu nhiên. Giúp hiểu quy luật của dữ liêu, dư đoán giá tri, so sánh dữ liêu với các PPXS đã biết.
- Biến Ngẫu Nhiên (BNN): Các đặc trưng dữ liêu được giả định là biến ngẫu nhiên.
  - Rởi rạc (Discrete): Nhận các giá trị đếm được (vd: số lần tung xúc xắc, số câu đúng, số khách hàng). Mô tả bằng Hàm khối xác suất (PMF).
  - Liên tục (Continuous): Nhận các giá trị trong một khoảng (vd: chiếu cao, nhiệt độ, thời gian). Mô tả bằng Hàm mật độ xác suất (PDF).

## 2. In Thư Viện scipy.stats

- Thư viên chính trong Python để làm việc với các PPXS.
- Import thư viện chung và các phân phối cụ thể:

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from scipy.stats import binom  # Nhị thức
from scipy.stats import poisson  # Poisson
from scipy.stats import uniform  # Đều
from scipy.stats import norm  # Chuán
from scipy.stats import expon  # Mū
from scipy.stats import t  # Student
from scipy.stats import chi2  # Chi-binh phương
```

· Một số phân phối thông dụng:

Scipy Package	Distribution Name	Tên Phân Phối	Ký hiệu
binom	Binomial Distribution	Phân Phối Nhị Thức	B(n, p)
poisson	Poisson Distribution	Phân Phối Poisson	$P(\lambda)$
uniform	Uniform Distribution	Phân Phối Đều	U(a, b)
norm	Normal/Gauss Distribution	Phân Phối Chuẩn	$N(\mu, \ \sigma^2)$
expon	Exponential Distribution	Phân Phối Mũ	$Exp(\lambda)$
t	Student t-Distribution	Phân Phối Student	t(df)
chi2	Chi-Square Distribution	Phân Phối Chi Bình Phương	$\chi^2(df)$

### 3. Các Hàm & Tham Số Chung

### A. Các Hàm Tính Xác Suất & Lấy Mẫu (Function Details)

- (dist>.rvs(...): Tao mẫu ngẫu nhiên (Random Variates).
  - o size: Số lương mẫu cần tạo.
  - Các tham số đặc trưng của phân phối (vd: (n), (p) cho (binom).
- <dist>.pmf(k, ...): Hàm khối xác suất (Probability Mass Function).
  - Dùng cho BNN Rời rac.
  - $\circ$  Tính xác suất P(X=k). (Hàm độ lớn f(x) cho phân phối rời rạc).
- <dist>.pdf(x, ...): Hàm mật đô xác suất (Probability Density Function).
  - o Dùng cho BNN Liên tuc.
  - $\circ f(x)$  (Hàm mật độ f(x) cho phân phối liên tục. Lưu ý:  $f(x) \ge 0$  và tích phân toàn miền bằng 1, f(x) không phải là xác suất tại điểm x).
- (dist>.cdf(x, ...): Hàm phân phối tích lũy (Cumulative Distribution Function).
  - $\circ$  Tính xác suất  $P(X \leq x)$  (Hàm F(x)).
- <dist>.ppf(q, ...): Hàm điểm phần trăm (Percent Point Function).
  - Hàm ngược của CDF (Inverse of CDF), còn gọi là Quantile Function.
  - $\circ$  Tìm x sao cho  $P(X \le x) = q$  (với q là xác suất). (Tìm điểm phân vi).
- <dist>.sf(x, ...): Hàm sống sót (Survival Function).
  - Tính xác suất P(X > x) (bằng 1 cdf(x)). (Hàm Survival).
- (dist>.isf(q, ...): Hàm ngược của hàm sống sót (Inverse Survival Function).
  - Tìm x sao cho P(X > x) = q. (Hàm ngược của hàm Survival).

### B. Các Tham Số Chung (loc, scale)

Hầu hết các phân phối liên tục đều có 2 tham số chuẩn hóa:

• loc : (Location) Dich chuyển vi trí phân phối.

```
\circ Với norm, loc chính là trung bình \mu.
```

- Với uniform, (loc là điểm bắt đầu a.
- scale: (Scale) Co dãn phân phối (quy định độ rộng)
  - $\circ$  Với norm, scale chính là đô lệch chuẩn  $\sigma$ .
  - $\circ$  Với (uniform), (scale) là độ rộng b-a.
  - $\circ$  Với expon, scale là  $1/\lambda$  (nghịch đảo của rate  $\lambda$ , tức là giá trị trung bình).

### C. Các Hàm Thống Kê

- (dist>.mean(...): Trả về kỳ vọng (trung bình) của phân phối.
- (<dist>.median(...): Trả về trung vi.
- $\langle \text{dist} \rangle$ .var $(\dots)$ : Trả về phương sai  $(\sigma^2)$ .
- $\langle \text{dist} \rangle$ .std $(\dots)$ : Trả về độ lệch chuẩn  $(\sigma)$ .

### 4. Phân Phối Nhị Thức (binom) - Rời Rạc

- Ký hiệu: B(n,p)
  - on: Số lần thực hiện phép thử Bernoulli độc lập.
  - op: Xác suất thành công trong một lần thử.
- Biến ngẫu nhiên X: Số lần thành công trong n phép thử.
- Ví dụ: Một bài thi trắc nghiệm gồm 10 câu, mỗi câu có 4 lựa chọn (p=0.25). Bạn đánh lụi tất cả. X là số câu trả lời đúng.  $X \sim B(10, 0.25)$ .

```
from scipy.stats import binom
import matplotlib.pyplot as plt
import numpy as no
n, p = 10, 0.25
# 1. Phát sinh mẫu ngẫu nhiên (rvs)
# Phát sinh 1 giá trị (1 bạn làm bài)
mau_1 = binom.rvs(n, p)
print(f"Một kết quả ngẫu nhiên: {mau_1}")
# Phát sinh 1000 giá tri (1000 ban làm bài)
mau 1000 = binom.rvs(n, p, size=1000)
# Vẽ histogram của mẫu
plt.figure(figsize=(10, 5))
dorongcot = 0.4
offset = dorongcot / 2
plt.hist(mau 1000, bins=np.arange(n + 2) - offset, width=dorongcot, density=True, label='Histogram từ mẫu RVS')
plt.xticks(range(n + 1))
plt.xlabel('Số câu đúng (k)')
plt.ylabel('Tần suất tương đối / Mật độ')
plt.title(f'Histogram Phân phối Nhị thức (n={n}, p={p}) từ 1000 mẫu')
# 2. Tính và vẽ hàm PMF (xác suất lý thuyết P(X=k))
k_{values} = np.arange(0, n + 1)
pmf_values = binom.pmf(k_values, n, p)
plt.stem(k values, pmf values, linefmt='r-', markerfmt='ro', basefmt=' ', label='PMF Lý thuyết')
plt.legend()
nlt.show()
# 3. Tính và vẽ hàm CDF (xác suất tích lũy P(X <= k))
plt.figure(figsize=(10, 5))
cdf values = binom.cdf(k values, n, p)
plt.plot(k_values, cdf_values, marker='o', drawstyle='steps-post', label='CDF') # steps-post cho đồ thi bâc thang
plt.title(f'Hàm Phân phối Tích lũy (CDF) Nhị thức (n={n}, p={p})')
plt.xlabel('Số câu đúng (k)')
plt.ylabel('P(X <= k)')</pre>
plt.xticks(k values)
plt.grid(True, which='both', linestyle='--', linewidth=0.5)
plt.legend()
plt.show()
```

```
# 4. Tính toán xác suất cụ thể

# P(X = 3) - Xác suất đúng chính xác 3 câu
p_eq_3 = binom.pmf(3, n, p)
print(f"P(X = 3) = {p_eq_3:.4f}")

# P(X <= 2) - Xác suất đúng tối đa 2 câu
p_le_2 = binom.cdf(2, n, p)
print(f"P(X <= 2) = {p_le_2:.4f}")

# P(X > 5) - Xác suất đúng nhiều hơn 5 câu
p_gt_5 = binom.sf(5, n, p) # <=> 1 - binom.cdf(5, n, p)
print(f"P(X > 5) = {p_gt_5:.4f}")

# Tìm k nhỏ nhất sao cho P(X <= k) >= 0.9 (Phân vị thứ 90)
k_q90 = binom.ppf(0.9, n, p)
print(f"Phân vị thứ 90 (k sao cho P(X<=k)>=0.9): {k_q90}")
```

### 5. Phân Phối Poisson (poisson) - Rời Rạc

- Ký hiệu:  $P(\lambda)$ 
  - lambda (λ): Trung binh số lần xảy ra sự kiện trong một khoảng (thời gian, không gian). Trong scipy tham số này tên là mu.
- Biến ngẫu nhiên X: Số lần sư kiên xảy ra.
- **Ví dụ**: Số khách hàng đến siêu thị trong 1 giờ (trung bình  $\lambda=20$ ).

```
from scipy.stats import poisson
# Trung bình 20 khách/giờ
mu lambda = 20
# P(X = 20) - Xác suất có đúng 20 khách
p eq 20 = poisson.pmf(20, mu=mu lambda)
print(f"P(X = 20) = \{p_eq_20:.4f\}")
# P(X <= 15) - Xác suất có tối đa 15 khách
p le 15 = poisson.cdf(15, mu=mu lambda)
print(f"P(X <= 15) = \{p_le_15:.4f\}")
# P(X > 30) - Xác suất có hơn 30 khách
p gt 30 = poisson.sf(30, mu=mu lambda) # <=> 1 - poisson.cdf(30, mu=mu lambda)
print(f"P(X > 30) = \{p \ gt \ 30:.4f\}")
# Vẽ PMF và CDF tương tự như Binomial
k values poi = np.arange(0, 41) # Chon khoảng giá tri hợp lý quanh lambda
pmf_poi = poisson.pmf(k_values_poi, mu=mu_lambda)
cdf_poi = poisson.cdf(k_values_poi, mu=mu_lambda)
fig, axes = plt.subplots(1, 2, figsize=(12, 4))
axes[0].stem(k_values_poi, pmf_poi, linefmt='b-', markerfmt='bo', basefmt=' ')
axes[0].set(xlabel='Số khách (k)', title=f'Poisson PMF ($\lambda$={mu lambda})')
axes[1].plot(k values poi, cdf poi, marker='.', drawstyle='steps-post')
axes[1].set(xlabel='Số khách (k)', title=f'Poisson CDF ($\lambda$={mu_lambda})')
plt.tight_layout()
plt.show()
```

### 6. Phân Phối Chuẩn (norm) - Liên Tuc

- Ký hiệu:  $N(\mu,\sigma^2)$ 
  - $1oc (\mu)$ : Kỳ vong (trung bình).
  - $\circ$  (scale) ( $\sigma$ ): Độ lệch chuẩn. (Lưu ý: tham số là  $\sigma$ , không phải  $\sigma^2$ ).
- Biến ngẫu nhiên X: Biến liên tục, đồ thị hình chuông đối xứng qua  $\mu$ .
- **Ví du**: Chiều cao nam thanh niên trưởng thành ( $\mu=170$  cm,  $\sigma=3$  cm),  $X\sim N(170,3^2)$ .

```
from scipy.stats import norm
import numpy as no
import matplotlib.pyplot as plt
mu, sigma = 170, 3
# 1. Phát sinh mẫu ngẫu nhiên (rvs) và vẽ histogram
plt.figure(figsize=(10, 5))
mau 10k = norm.rvs(loc=mu, scale=sigma, size=10000)
plt.hist(mau 10k, bins=100, density=True, alpha=0.6, color='g', label='Histogram từ mẫu RVS')
# 2. Tính và vẽ hàm PDF (hàm mật đô xác suất) lý thuyết
x values = np.linspace(mu - 4*sigma, mu + 4*sigma, 200) # Khoảng giá tri quanh mu (~99.99%)
pdf values = norm.pdf(x values, loc=mu, scale=sigma)
plt.plot(x values, pdf values, color='red', lw=2, label=f'PDF Lý thuyết $N({mu}, {sigma}^2)$')
plt.title(f'Phân phối Chuẩn (Normal) - PDF và Histogram Mẫu')
plt.xlabel('Chiêu cao (cm)')
plt.ylabel('Mật độ')
plt.legend()
plt.grid(True, linestyle='--', alpha=0.7)
plt.show()
# 3. Tính toán với Phân phối Chuẩn Tắc (Z ~ N(0, 1))
print("\n--- Tính toán với Phân phối Chuẩn Tắc Z ~ N(0, 1) ---")
\# P(7 <= -1)
p le neg1 = norm.cdf(-1) # Măc đinh loc=0, scale=1
print(f"P(Z \leftarrow -1) = \{p \text{ le neg1:.5f}\}")
\# P(Z > 1) = 1 - P(Z \le 1)
p_gt_1 = 1 - norm.cdf(1)
# Hoăc dùng hàm sf: P(Z > 1)
p gt 1 sf = norm.sf(1)
print(f"P(Z > 1) = \{p \text{ gt 1:.5f}\} \text{ (dùng 1-cdf)} = \{p \text{ gt 1 sf:.5f}\} \text{ (dùng sf)"}\}
\# P(-1 \le 7 \le 1)
p between = norm.cdf(1) - norm.cdf(-1)
print(f"P(-1 <= Z <= 1) = {p_between:.5f}")</pre>
# 4. Tìm điểm phân vi (Ouantile) với Phân phối Chuẩn Tắc
print("\n--- Tìm điểm phân vi (Ouantile) với Z ~ N(0, 1) ---")
# Tim z sao cho P(Z \le z) = 0.02275
z_val_ppf = norm.ppf(0.02275)
print(f"z sao cho P(Z \le z) = 0.02275 là: {z_val_ppf:.3f}") # Gần bằng -2
# Tim z sao cho P(Z > z) = 0.02275
z val isf = norm.isf(0.02275)
print(f"z sao cho P(Z > z) = 0.02275 là: {z_val_isf:.3f}") # Gần bằng 2
# 5. Tính toán với phân phối chuẩn bất kỳ (ví du X ~ N(16, 4^2))
print("\n--- Tính toán với Phân phối Chuẩn X ~ N(16, 4^2) ---")
mu_x, sigma_x = 16, 4
# P(X <= 8)
p_le_8 = norm.cdf(8, loc=mu_x, scale=sigma_x)
print(f"P(X <= 8) = \{p_1e_8:.5f\}")
# Tim x sao cho P(X \le x) = 0.02275
x_val_ppf = norm.ppf(0.02275, loc=mu_x, scale=sigma_x)
print(f"x sao cho P(X <= x) = 0.02275 là: {x_val_ppf:.3f}") # = 16 + (-2)*4 = 8
# Vẽ minh họa P(-1 <= Z <= 1)
plt.figure(figsize=(8, 4))
x_z = np.linspace(-4, 4, 200)
pdf z = norm.pdf(x z)
plt.plot(x_z, pdf_z, label='PDF của Z ~ N(0,1)')
x_{fill} = np.linspace(-1, 1, 100)
plt.fill_between(x_fill, norm.pdf(x_fill), color='blue', alpha=0.35, label='P(-1 <= Z <= 1)')
plt.title('Minh hoa diên tích dưới đường cong PDF chuẩn tắc')
plt.text(-0.4, 0.1, f'{p_between:.3f}', fontsize=12)
plt.legend()
plt.grid(True, linestyle='--', alpha=0.7)
plt.show()
```

7. Phân Phối Đều (uniform) - Liên Tục

- Ký hiệu: U(a,b)
  - o a : Giá trị nhỏ nhất.
  - b : Giá tri lớn nhất.
- Ý nghĩa: Mọi giá trị trong khoảng [a, b] đều có cùng khả năng xảy ra (hàm mật đô là hằng số trong khoảng này).
- Tham số scipy: loc = a, scale = b a.

```
from scipy.stats import uniform
a, b = 10, 30
# loc = 10, scale = 30 - 10 = 20
dist_u = uniform(loc=a, scale=(b-a))
# P(X <= 15)
p u le 15 = dist u.cdf(15)
print(f"P(X \le 15) = \{p_u_le_15:.2f\}") # (15-10) / (30-10) = 0.25
# P(20 <= X <= 25)
p u between = dist u.cdf(25) - dist u.cdf(20)
print(f''P(20 \le X \le 25) = \{p \text{ u between: } .2f\}'') \# (25-20) / (30-10) = 0.25
# Lấy mẫu và vẽ histogram
mau deu = dist u.rvs(10000)
plt.hist(mau_deu, bins=20, density=True, alpha=0.6, label='Histogram Mau')
# Vẽ PDF lý thuyết
x_u = np.linspace(a - 5, b + 5, 100) # Vē rộng hơn khoảng [a,b]
pdf_u = dist_u.pdf(x_u)
plt.plot(x_u, pdf_u, 'r-', lw=2, label='PDF Lý thuyết')
plt.title(f'Phân phối Đều U({a},{b})')
plt.legend()
plt.show()
```

### 8. Phân Phối Mũ (expon) - Liên Tục

- Ký hiệu:  $Exp(\lambda)$ 
  - λ: Tốc đô (rate) trung bình của sư kiên (ví du: số sư kiên / đơn vi thời gian).
- Ý nghĩa: Mô tả thời gian chờ cho đến khi sư kiên tiếp theo xảy ra (trong quá trình Poisson).
- Tham số scipy: scale = 1 / \lambda.
  - Chú ý: scale chính là thời gian chờ trung bình giữa các sự kiện. Nếu trung bình thời gian chờ là 5 phút ( $1/\lambda = 5$ ), thì scale=5.

```
from scipy.stats import expon
# Ví dụ: Trung bình thời gian chờ xe bus là 10 phút.
# 1/lambda = 10 (thời gian chờ trung bình)
# Ta dùng scale = 10
dist_e = expon(scale=avg_wait_time) # loc=0 là mặc định (thời gian không âm)
# P(X <= 5) - Xác suất chờ tối đa 5 phút
n = 1e = 5 = dist e.cdf(5)
print(f"P(X \le 5) = \{p e le 5:.4f\}")
# P(X > 10) - Xác suất chờ hơn 10 phút (hơn thời gian trung bình)
p_e_gt_10 = dist_e.sf(10)
print(f"P(X > 10) = \{p e gt 10:.4f\}") # Khoảng 0.368 (1/e)
# Vẽ PDE và CDE
x_e = np.linspace(0, avg_wait_time * 4, 200) # Ve đến khoảng 4 lần trung bình
pdf e = dist e.pdf(x e)
cdf_e = dist_e.cdf(x_e)
fig, axes = plt.subplots(1, 2, figsize=(12, 4))
axes[0].plot(x e, pdf e, label=f'PDF ($1/\lambda$={avg wait time})')
axes[0].set(xlabel='Thời gian chờ (phút)', title='Exponential PDF')
axes[0].legend()
axes[0].grid(True, linestyle='--', alpha=0.7)
```

```
axes[1].plot(x_e, cdf_e, label=f'CDF ($1/\lambda$={avg_wait_time})')
axes[1].set(xlabel='Thời gian chờ (phút)', title='Exponential CDF')
axes[1].legend()
axes[1].grid(True, linestyle='--', alpha=0.7)
plt.tight_layout()
plt.show()
```

### 9. Phân Phối Student (t) - Liên Tục

- Ký hiệu: t(df)
  - o df: Bâc tư do (Degrees of Freedom).
- Ý nghĩa: Tương tự phân phối Chuẩn nhưng có "đuôi nặng" (heavy tails), nghĩa là các giá trị cực đoan (xa trung bình) dễ xảy ra hơn. đặc biệt khi df nhỏ.
- Sử dụng: Khi ước lượng trung bình tổng thể từ mẫu nhỏ (thường n < 30) và chưa biết phương sai  $\sigma^2$  tổng thể. Khi df lớn (ví du df > 30), phân phối t xấp xì phân phối t norm.

```
from scipy.stats import t
df_low = 3 # Bậc tự do thấp (đuôi nặng)
df medium = 10
df_high = 50 # Bậc tự do cao (gần giống normal)
x t = np.linspace(-5, 5, 200)
plt.figure(figsize=(10, 6))
plt.plot(x_t, norm.pdf(x_t), 'k-', lw=2.5, alpha=0.8, label='Normal (Z)')
plt.plot(x t, t.pdf(x t, df=df low), 'r--', lw=1.5, label=f't (df={df low})')
plt.plot(x t, t.pdf(x t, df=df medium), 'g-.', lw=1.5, label=f't (df={df medium})')
plt.plot(x t, t.pdf(x t, df=df high), 'b:', lw=1.5, label=f't (df={df high})')
plt.legend()
plt.title('So sánh Phân phối t với Phân phối Chuẩn tắc')
plt.xlabel('Giá tri')
plt.ylabel('Mật độ')
plt.grid(True, linestyle='--', alpha=0.7)
nlt.show()
```

## v 10. 🔬 Ước Lượng Tham Số (Fitting)

- Mục đích: Tìm các tham số của một phân phối lý thuyết (ví dụ: μ, σ của norm) sao cho phân phối đó khớp (fit) nhất với dữ liêu mẫu ta có.
- Hàm: <dist>.fit(data) Trả về các tham số ước lương được.

```
# 1. Tao dữ liêu mẫu (ví du: từ N(15, 2^2))
np.random.seed(42) # Để kết quả có thể tái lặp
mu_true, sigma_true = 15, 2
data_sample = norm.rvs(loc=mu_true, scale=sigma_true, size=1000)
# 2. Giả sử không biết mu, sigma. Dùng .fit() để ước lượng từ data_sample
mu fit, sigma fit = norm.fit(data sample)
print(f"Tham số gốc: mu = {mu_true:.4f}, sigma = {sigma_true:.4f}")
print(f"Tham số ước lượng: mu = {mu_fit:.4f}, sigma = {sigma_fit:.4f}")
# 3. Vẽ biểu đồ so sánh histogram của dữ liêu và PDF với tham số ước lương
plt.figure(figsize=(10, 6))
plt.hist(data_sample, bins=30, density=True, alpha=0.6, label='Histogram Dữ liệu Mẫu')
x_fit = np.linspace(data_sample.min(), data_sample.max(), 200)
pdf fit = norm.pdf(x fit, loc=mu fit, scale=sigma fit)
plt.plot(x_fit, pdf_fit, 'r-', lw=2, label=f'PDF Uớc lượng N({mu_fit:.2f}, {sigma_fit:.2f}^2)')
plt.title('Ước lượng tham số Phân phối Chuẩn từ Dữ liệu Mẫu')
plt.xlabel('Giá trị')
plt.ylabel('Mật độ')
plt.legend()
```

```
plt.grid(True, linestyle='--', alpha=0.7)
plt.show()
```

# Kiểm tra giữa kỳ

Dưới đây là tóm tắt khung bài làm và các hàm chính được sử dụng trong bài kiểm tra, dựa trên để bài ([TKMT] GK.docx) và đáp án (DapAn.ipynb).

# 連 Các Thư Viện Cần Thiết

Khung bài làm tiêu chuẩn luôn bắt đầu bằng việc nhập các thư viện cốt lõi:

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from scipy.stats import poisson, binom, norm
```

### 🗸 🚜 Bài 1: Phân Tích Dữ Liêu (Dataset: mpg.csv)

Bài này tập trung vào kỹ năng thao tác, phân tích và trực quan hóa dữ liêu bằng Pandas và Seaborn.

- 1. Đọc và Làm Sạch Dữ Liệu
  - [cite start] Muc tiêu: Tải dữ liêu và xử lý các giá tri thiếu hoặc trùng lặp[cite: 47, 48, 49].
  - Cách làm:
    - 1. Doc file: mpg\_df = pd.read\_csv('mpg.csv')
    - 2. Kiểm tra thiếu: mpg\_df.isnull().sum() (Đáp án cho thấy cột 'horsepower' có 6 giá trị thiếu).
    - 3. Xử lý thiếu: mpg df.dropna(inplace=True) (Xóa các dòng có giá tri thiếu).
    - 4. Xử lý trùng lặp: mpg\_df.drop\_duplicates(inplace=True) (Xóa các dòng bị trùng).

### 2. Tính Thống Kê Cơ Bản

- [cite start] Muc tiêu: Tính các chỉ số thống kê mô tả cho các côt cu thể[cite: 50, 51].
- · Cách làm:
  - c [cite\_start]Trung binh (Mean): mean\_values = mpg\_df[['cylinders', 'horsepower', 'weight']].mean()
    [cite: 50]
  - c [cite\_start]Trung vi (Median): median\_values = mpg\_df[['cylinders', 'horsepower',
     'weight']].median() [cite: 50]
  - [cite start]Miền giá trị (Range): mpg range = (mpg df['mpg'].min(), mpg df['mpg'].max()) [cite: 51]
  - o [cite start]Phân vi (Quartiles): [cite: 51]
    - Q1: mpg q1 = mpg df['mpg'].quantile(0.25)
    - Q3: mpg\_q3 = mpg\_df['mpg'].quantile(0.75)
  - o [cite\_start]Nhận xét: Dựa vào miền giá trị và Q1, Q3 để xem có giá trị nào quá xa không (outlier)[cite: 52].
- 3. Phân Tích Theo Nhóm (Groupby)
  - [cite\_start]Muc tiêu: Tính giá trị trung bình của 'mpg' cho từng nhóm 'origin' [cite: 53] [cite\_start]và tìm nhóm tiêu hao ít nhiên liệu nhất (MPG cao nhất) hoặc nhiều nhất (MPG thấp nhất)[cite: 54].
  - · Cách làm:
    - 1. **Tính trung bình theo nhóm:** (mpg by origin = mpg df.groupby('origin')['mpg'].mean())
    - 2. Tìm khu vực tiêu hao ít nhiên liệu nhất (MPG cao nhất): mpg\_by\_origin.idxmax()
    - 3. [cite\_start]Tim khu vực tiêu hao nhiều nhiên liệu nhất (MPG thấp nhất): [mpg\_by\_origin.idxmin()] (Đáp án dùng [idxmin] để tìm xe tiêu hao ít nhiên liệu nhất, điểu này có vẻ ngược với để[cite: 37, 54]. Lưu ý: [mpg] (Miles/Gallon) càng cao thì càng ít tốn nhiên liệu. Do đó, [idxmax()] mới là xe ít tiêu hao nhiên liệu nhất).

### 4. Trưc Quan Hóa (Phân Bố & So Sánh)

- [cite\_start]Muc tiêu: Vẽ Histogram của 'weight' [cite: 55] [cite\_start]và Boxplot so sánh 'horsepower' giữa các thập niên[cite: 56].
- · Cách làm:
  - o Histogram (với Seaborn):

```
sns.histplot(mpg_df['weight'], bins=20, kde=True)
plt.title("Phân bố trọng lượng xe")
plt.xlabel("Trọng lượng (pound)")
plt.ylabel("Số lượng xe")
plt.show()
```

• Boxplot (với Seaborn): \* Mẹo: Để chia 'model\_year' (70, 71, ..., 80, 81, 82) thành thập niên 70 và 80, ta dùng phép chia lấy nguyên // 10 \* 10.

```
# Tạo cột thập niên (70 hoặc 80)
decade = mpg_df['model_year'] // 10 * 10
sns.boxplot(x=decade, y=mpg_df['horsepower'])
plt.title("So sánh công suất giữa thập niên 70 và 80")
plt.xlabel("Thập niên sản xuất")
plt.ylabel("Công suất (hp)")
plt.show()
```

- [cite\_start]Nhận xét: Nhìn vào boxplot để so sánh trung vị, Q1, Q3 và các điểm ngoại lai (outlier) giữa hai nhóm[cite: 57].
- 5. Trực Quan Hóa (Mối Quan Hệ)
  - [cite start] Muc tiêu: Vẽ 4 biểu đồ Scatter Plot để xem mối quan hệ qiữa 'mpq' và các biến khác[cite: 58].
  - · Cách làm:
    - Sử dụng plt.subplots để tạo một lưới 2x2.
    - Vẽ từng biểu đổ scatter vào mỗi ax (truc) tương ứng.

```
fig, axes = plt.subplots(2, 2, figsize=(12, 10)) # Tạo lưới 2x2

# Biểu đô 1: MPG vs Cylinders
sns.scatterplot(x=mpg_df['cylinders'], y=mpg_df['mpg'], ax=axes[0, 0])
axes[0, 0].set_title("MPG vs Cylinders")

# Biểu đô 2: MPG vs Horsepower
sns.scatterplot(x=mpg_df['horsepower'], y=mpg_df['mpg'], ax=axes[0, 1])
axes[0, 1].set_title("MPG vs Horsepower")

# Biểu đô 3: MPG vs Acceleration
sns.scatterplot(x=mpg_df['acceleration'], y=mpg_df['mpg'], ax=axes[1, 0])
axes[1, 0].set_title("MPG vs Acceleration")

# Biểu đô 4: MPG vs Model Year
sns.scatterplot(x=mpg_df['model_year'], y=mpg_df['mpg'], ax=axes[1, 1])
axes[1, 1].set_title("MPG vs Model Year")

plt.tight_layout() # Tự động căn chính
plt.show()
```

- [cite\_start]**Nhận xét:** Đánh giá yếu tố ảnh hưởng lớn nhất [cite: 59] bằng cách xem biểu đổ nào có xu hướng (tuyến tính. loa....) rõ ràng nhất (ví dụ: 'horsepower' và 'cylinders' có vẻ tương quan nghịch manh với 'mpg').
- La Bài 2: Phân Phối Xác Suất (Scipy.stats)

Bài này tập trung vào việc áp dụng 3 phân phối xác suất cơ bản để giải quyết bài toán thực tế.

- (a) Phân Phối Poisson (Số lần xuất hiện trung bình)
- [cite\_start] Mục tiêu: Tính xác suất có nhiều hơn 3 lỗi (P(X>3)), biết trung bình là 2 lỗi  $(\lambda=2)$  [cite: 62, 63].
- Phân tích: P(X > 3) = 1 P(X < 3).
- Hàm: Sử dụng poisson.cdf (Hàm phân phối tích lũy).
  - k = 3 (số lỗi)
  - mu = 2 (tham số  $\lambda$  trong scipy.stats là mu)
- · Cách làm:

```
# P(X > 3) = 1 - P(X <= 3)
prob_more_than_3 = 1 - poisson.cdf(3, mu=2)
print(f"Xác suất có hơn 3 lỗi: {prob_more_than_3:.4f}") # Kết quả: 0.1429</pre>
```

- (b) Phân Phối Nhị Thức (Số lần thành công trong n phép thủ)
  - [cite start] Muc tiêu: Tính xác suất có chính xác 18 vi mạch đạt chuẩn (P(X=18))[cite: 65].
  - · Phân tích:
    - $\circ$  [cite start]Số phép thử (tổng số vi mạch): n=20 [cite: 64]
    - $\circ$  [cite\_start]Xác suất thành công (đạt chuẩn): p=0.95 [cite: 64]
    - $\circ$  [cite start]Số lần thành công mong muốn: k=18 [cite: 65]
  - Hàm: Sử dụng binom.pmf (Hàm khối xác suất cho giá trị chính xác).
  - · Cách làm:

```
prob_18_success = binom.pmf(18, n=20, p=0.95)
print(f"Xác suất có đúng 18 vi mạch đạt: {prob_18_success:.4f}") # Kết quả: 0.1887
```

- (c) Phân Phối Chuẩn (Biến liên tục)
  - [cite start] Muc tiêu 1: Tính xác suất hoàn thành trong vòng 45 phút  $(P(X \le 45))$  [cite: 67].
  - · Phân tích:
    - [cite start]Trung binh:  $\mu = 50$  phút [cite: 66]
    - $\circ$  [cite\_start]Độ lệch chuẩn:  $\sigma=5$  phút [cite: 66]
    - $\circ$  [cite start]Giá tri cần tính: x=45 [cite: 67]
  - Hàm: Sử dụng norm. cdf (Hàm phân phối tích lũy tính diên tích bên trái của x).
  - · Cách làm (Muc tiêu 1):

```
# P(X <= 45) với mu=50, sigma=5
prob_under_45 = norm.cdf(45, loc=50, scale=5)
print(f"Xác suất hoàn thành trong 45 phút: {prob_under_45:.4f}") # Kết quả: 0.1587
```

- [cite\_start] Mục tiêu 2: Tìm thời gian x để nằm trong top 5% nhanh nhất ( $P(X \le x) = 0.05$ )[cite: 68].
- Phân tích: Đây là bài toán tìm điểm phân vi (ngược của CDF).
  - $\circ~$  Xác suất (quantile): q=0.05 (vì là 5% nhanh nhất, tức là 5% ở đuôi bên trái)
- Hàm: Sử dụng norm.ppf (Hàm điểm phần trăm Percent Point Function).
- Cách làm (Muc tiêu 2):

```
# Tim x sao cho P(X <= x) = 0.05
thresh_top_5 = norm.ppf(0.05, loc=50, scale=5)
print(f"Thời gian cần để vào top 5% nhanh nhất: {thresh_top_5:.2f} phút") # Kết quả: 41.78 phút
```