# Python for Data Analytics

Xi Li (李曦) Professor of Marketing Professor of Information and Innovation Management Director, Asia Case Research Centre Associate Director, Institute of Digital Economy and Innovation

Ph.D. (Management), University of TorontoM.Phil. (Operations Research), HKUSTB.E. (Computer Science), Tsinghua University

Email: xili@hku.hk Office hour: by appointment

# Course Website

Here: https://ximarketing.github.io/\_pages/teaching/ Password: ATC All course materials will be posted on this website.

# Python

You should already know about these...

- Install modules in Python
- pip install
- Running Python codes
- Branch and Loop
- Functions

Let Vincent (or me) know if you don't know about these topics.

# Python

In the next a few classes, we are learning more things!

- Statistical Tests
- Statistical Analysis and Basic Machine Learning: Linear and Logistic Regression, Multinomial Logit Models, Artificial Neural Networks, and Clustering
- Advanced Algorithms: Probabilistic, Dynamic Programming, Divide and Conquer
- Unstructured Data: Text and Images
- API and Data Scraping

# Python

Our plan is too ambiguous and we only have 15 hours left. How can we learn so many topics in 5 sessions?

# Schedule

Classes 6 and 7: Python for Data Analytics Class 8: Advanced Algorithms Class 9: Unstructured Data Analysis Class 10: API, Data Scraping and Class Review

# Statistical Analyses and Visualization with GPT

# Reading a document online

I have prepared an online CSV file, which is available here. Alternatively, its URL can be found here: https://ximarketing.github.io/data/python\_data1.csv

Now, let's ask GPT about how to load this CSV file to python.

# Which LLM should I use?

HKU has subscribed to a few GPT models: Login here to use ChatGPT: https://chatgpt.hku.hk/

• GPT 4o > GPT 4 > GPT 3.5.

• GPT 40 can handle images while other versions cannot. You can also use POE assistant here: https://poe.com/

- POE assistant can also handle images as input
   You can also use DeepSeek here: https://www.deepseek.com
  - DeepSeek can also be used to process images as input

# I want to read the data https://ximarketing.github.io/data/python\_data1.csv to python. Show me the code

 $\sim$ 



#### GPT-3.5 (Tokens used 174)

You can use the following Python code to read data from a CSV file hosted on a website:

Xi Li

```
import pandas as pd
url = "https://ximarketing.github.io/data/python data1.csv"
df = pd.read csv(url)
print(df)
```

Make sure you have the pandas library installed in your Python environment to run this code. You can install it using pip install pandas . This code will read the data from the CSV file at the specified URL and store it in a pandas DataFrame called df .

# From your disk



# From your disk



#### 

1 import pandas as pd
2 url = "https://ximarketing.github.io/data/python\_data1.csv"
3 df = pd.read\_csv(url)
4 print(df)

## Plotting the relationship between age and credit\_score

#### Xi Li

Plotting the relationship between age and credit\_score

17/02/2025, 09:23:56

#### •••

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 url = "https://ximarketing.github.io/data/python_datal.csv"
5 df = pd.read_csv(url)
6
7 plt.figure(figsize=(10, 6))
8 plt.scatter(df['age'], df['credit_score'], color='green', alpha=0.5)
9 plt.title('Relationship between Age and Credit Score')
10 plt.xlabel('Relationship between Age and Credit Score')
10 plt.xlabel('Age')
11 plt.ylabel('Credit Score')
12 plt.grid(True)
13 plt.show()
```



#### Relationship between Age and Credit Score

### You can also ask for a density heatmap.



### Code for generating density heatmap:

#### 

```
1 import pandas as pd
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4
5 url = "https://ximarketing.github.io/data/python_datal.csv"
6 df = pd.read_csv(url)
7
8 plt.figure(figsize=(10, 6))
9 sns.kdeplot(x=df['age'], y=df['credit_score'], cmap='Blues',
shade=True, cbar=True)
10 plt.title('Density Heatmap of Age and Credit Score')
11 plt.xlabel('Age')
12 plt.ylabel('Credit Score')
13 plt.show()
```

# Testing the correlation between age and credit score:

#### Xi Li

I want to test the correlation between age and credit\_score. Show me the code

17/02/2025, 09:35:55

#### Xi Li

I also need the significance of the test

17/02/2025, 09:36:52

#### •••

- 1 import pandas as pd
- 2 url = "https://ximarketing.github.io/data/python\_data1.csv"
- 3 df = pd.read\_csv(url)
- 4 correlation = df['age'].corr(df['credit\_score'])
- 5 print("Pearson correlation coefficient between age and credit\_score:", correlation)

### Adding a Regression Line

#### •••

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 url = "https://ximarketing.github.io/data/python_datal.csv"
5 df = pd.read_csv(url)
6 plt.figure(figsize=(10, 6))
7 plt.scatter(df['age'], df['credit_score'], color='green', alpha=0.5)
8 plt.title('Relationship between Age and Credit Score')
9 plt.xlabel('Age')
10 plt.ylabel('Credit Score')
11 plt.grid(True)
12 sns.regplot(x='age', y='credit_score', data=df, scatter=False, color='blue'
13 plt.show()
```



# Adding regression output

Method: Date: Mo Time: No. Observations Df Residuals: Df Model:	Least Squares F-statistic: n, 17 Feb 2025 Prob (F-statistic): 09:50:13 Log-Likelihood: : 1000 AIC: 998 BIC: 1	93.84 2.85e-21 -5288.6 1.058e+04 1.059e+04	
 coe const 623.9154 age 2.9000	f std err t $P> t $ [0.02 10.567 59.042 0.000 60 0 0.299 9.687 0.000 2	25 0.975] 03.179 644.652 2.313 3.487	
 Omnibus: Prob(Omnibus): Skew: Kurtosis:	2.290 Durbin-Watson: 0.318 Jarque-Bera (JB): -0.110 Prob(JB): 3.063 Cond. No.	2.067 2.174 0.337 246.	

Notes: [1] Standard Errors assume that the covariance matrix of the errors is correctly specified. Let 's Try Quadratic Regression Now!

$$ext{credit score}_i = lpha + eta_1 imes ext{age}_i + eta_2 imes ext{age}_i^2$$

#### •••

```
1 import pandas as pd
 2 import numpy as np
 3 import matplotlib.pyplot as plt
 4 import statsmodels.api as sm
 5 url = "https://ximarketing.github.io/data/python data1.csv"
 6 data = pd.read csv(url)
 7 data['age squared'] = data['age'] ** 2
 8 X = data[['age', 'age squared']]
9 X = sm.add constant(X) # Add a constant term for the intercept
10 y = data['credit score']
11 model = sm.OLS(y, X).fit()
12 print(model.summary())
13 plt.figure(figsize=(10, 6))
14 X pred = np.linspace(data['age'].min(), data['age'].max(), 100)
15 X pred squared = X pred ** 2
16 X pred combined = np.column stack((np.ones like(X pred), X pred, X pred squared))
17 y pred = model.predict(X pred combined)
18 plt.plot(X pred, y pred, color='red', label='Regression Line')
19 plt.scatter(data['age'], data['credit score'], label='Data', s=10)
20 equation = f'Credit Score = {model.params[0]:.2f} + {model.params[1]:.2f}*Age +
   {model.params[2]:.2f}*Age^2'
21 r squared = f'R-squared = {model.rsquared:.2f}'
22 plt.text(0.5, 1.05, equation, ha='center', va='center', transform=plt.gca().transAxes,
   fontsize=10)
23 plt.text(0.5, 1.01, r squared, ha='center', va='center', transform=plt.gca().transAxes,
   fontsize=10)
24 plt.xlabel('Age')
25 plt.ylabel('Credit Score')
```



### Exercise 1:

# Use GPT to design a desktop APP to visualize your scatter plot and regression!

## Exercise 1:

Use GPT to design a desktop APP to visualize your scatter plot and regression!

Note: Creating a mobile APP is possible but much more complex. You also need different APPs for iOS and Android. When creating the desktop APP, it would be easy to use the tkinter module (you can inform GPT).

# This is my APP!



### Exercise 2:

Use GPT to design a desktop APP to allow users to click on the screen to input the points and generate a regression line fitting the points.

```
1 import tkinter as tk
 2 import numpy as np
 3 from sklearn.linear model import LinearRegression
4 \text{ points} = []
5 def generate point(event):
       x = event.x
       y = event.y
       points.append((x, y))
       canvas.create oval(x - 2, y - 2, x + 2, y + 2, fill = "red")
10 def perform regression():
       X = [point[0] for point in points]
11
       y = [point[1] for point in points]
12
13
       X = np.array(X).reshape(-1, 1)
       y = np.array(y)
14
       model = LinearRegression()
15
       model.fit(X, y)
16
       y pred = model.predict(X)
17
18
       for i in range(len(X) - 1):
19
           x1, y1 = X[i][0], y pred[i]
           x_{2}, y_{2} = X[i + 1][0], y pred[i + 1]
20
           canvas.create line(x1, y1, x2, y2, fill="blue")
21
       equation = "y = {:.2f}x + {:.2f}".format(model.coef [0], model.intercept )
22
23
       canvas.create text(10, 10, anchor='nw', text=equation, fill='black', font=('Arial', 12))
24 def clear canvas():
25
       canvas.delete("all")
26
       del points[:]
27 \text{ app} = \text{tk.Tk()}
28 app.title("Regression App")
29 canvas = tk.Canvas(app, width=800, height=500, bg="white")
30 canvas.grid(row=0, column=0, columnspan=2, padx=10, pady=10)
31 canvas.bind("<Button-1>", generate point) # Bind left mouse click to generate point function
32 regression button = tk.Button(app, text="Regression", command=perform regression)
33 regression button.grid(row=1, column=0, padx=10, pady=10)
34 clear button = tk.Button(app, text="Clear", command=clear canvas)
35 clear button.grid(row=1, column=1, padx=10, pady=10)
```

```
36 app.mainloop()
```

## *t*-test: What is a *t*-test?

In a *t*-test, we try to compare whether two groups have the same mean. For example, suppose that we want to know whether HKU graduates, on average, make higher salaries than CUHK graduates do. In this case, we would like to compare the means of HKU graduates and CUHK graduates.

The dataset is here:

https://ximarketing.github.io/data/python\_data2.csv

## Let's view the histogram first (with the help of GPT)!

#### 

- 1 import pandas as pd
  2 import matplotlib.pyplot as plt
  - 3 data = pd.read\_csv("https://ximarketing.github.io/data/python\_data2.csv")
  - 4 salary\_hku = data[data['School'] == 'HKU']['Salary']
  - 5 salary\_cuhk = data[data['School'] == 'CUHK']['Salary']
  - 6 plt.hist(salary\_hku, bins=10, alpha=0.5, label='HKU', color='lightblue')
  - 7 plt.hist(salary\_cuhk, bins=10, alpha=0.5, label='CUHK', color='lightcoral')
  - 8 plt.xlabel('Salary')
  - 9 plt.ylabel('Frequency')
- 10 plt.title('Salary Histogram of HKU and CUHK')
- 11 plt.legend()
- 12 plt.show()



Ask GPT to help you generate the code for running a *t*-test to compare the salaries of the two schools!

# The code I got (with a bit modification):

#### •••

```
1 import pandas as pd
 2 from scipy import stats
 3 # Load the data from the CSV file
 4 data = pd.read csv("https://ximarketing.github.io/data/python data2.csv")
 5 # Separate the salary data for HKU and CUHK graduates
 6 hku salaries = data[data['School'] == 'HKU']['Salary']
 7 cuhk salaries = data[data['School'] == 'CUHK']['Salary']
 8 # Calculate the means
 9 hku mean = hku salaries.mean()
10 cuhk mean = cuhk salaries.mean()
11 # Perform t-test
12 t statistic, p value = stats.ttest ind(hku salaries, cuhk salaries)
  print("Mean salary for HKU graduates:", hku mean)
13
   print("Mean salary for CUHK graduates:", cuhk mean)
14
15 print("T-Statistic:", t statistic)
16 print("P-Value:", p value)
```

# Asking for visualization (including standard error)

#### •••

```
1 import pandas as pd
 2 from scipy import stats
 3 import matplotlib.pyplot as plt
 4 import numpy as np
 5 data = pd.read csv("https://ximarketing.github.io/data/python data2.csv")
 6 hku salaries = data[data['School'] == 'HKU']['Salary']
 7 cuhk salaries = data[data['School'] == 'CUHK']['Salary']
 8 hku mean = hku salaries.mean()
 9 cuhk mean = cuhk salaries.mean()
10 hku std = hku salaries.std()
11 cuhk std = cuhk salaries.std()
12 hku se = hku std / np.sqrt(len(hku salaries))
  cuhk se = cuhk std / np.sqrt(len(cuhk salaries))
13
14 t statistic, p value = stats.ttest ind(hku salaries, cuhk salaries)
15 plt.bar(['HKU', 'CUHK'], [hku mean, cuhk mean], yerr=[hku se, cuhk se],
   capsize=5, color=['blue', 'green'])
16 plt.ylabel('Mean Salary')
17 plt.title('Mean Salary for HKU and CUHK Graduates with Standard Errors')
18 plt.show()
```


#### Exercise:

Use GPT to design a desktop APP. The APP takes input from the user (the values of two groups), runs the *t*-test and outputs the result to the user.

 $\chi^2$  test: What is a  $\chi^2$  test? Can you propose an example of a  $\chi^2$  test?

#### Example of $\chi$ -Square Test

Suppose that we are interested in whether education affects marriage status, and we have collected the following data:

	Primary School	Secondary School	Bachelor	Master	Doctor
Single	12	18	36	19	4
Married	45	59	107	47	8
Separed	11	12	18	2	0
Divorced	12	20	23	5	2

#### Example of $\chi$ -Square Test

#### What conclusions can you draw from the table?

	Primary School	Secondary School	Bachelor	Master	Doctor
Single	12	18	36	19	4
Married	45	59	107	47	8
Separed	11	12	18	2	0
Divorced	12	20	23	5	2

### My code produced by GPT

#### •••

```
1 import numpy as np
   from scipy.stats import chi2 contingency
 2
   # Create a 2x2 contingency table (replace the values with your own data)
 4
   observed table = np.array([
 5
                    [12,18,36,19,4],
 6
                    [45,59,107,47,8],
                    [11,12,18,2,0],
 9
                    [12, 20, 23, 5, 2]])
10
11
   # Perform the Chi-square test
   chi2 stat, p val, dof, expected table = chi2 contingency(observed table)
12
13
14
   print("Chi-square Statistic:", chi2 stat)
   print("P-value:", p val)
15
  print("Degrees of Freedom:", dof)
16
  print("Expected Table:")
17
  print(expected table)
18
```

#### •••

```
import numpy as np
   from scipy.stats import chi2 contingency
 2
 3
   # Create a 2x2 contingency table (replace the values with your own data)
 4
   observed table = np.array([
 5
                    [12,18,36,19,4],
 6
                    [45,59,107,47,8],
                    [11,12,18,2,0],
                    [12, 20, 23, 5, 2]])
 9
10
11
   # Perform the Chi-square test
12
   chi2 stat, p val, dof, expected table = chi2 contingency(observed table)
13
   print("Chi-square Statistic:", chi2 stat)
14
   print("P-value:", p val)
15
16 print("Degrees of Freedom:", dof)
   print("Expected Table:")
17
   print(expected table)
18
```

The significance value is p = 0.206 > 5%, implying that "there is *no* evidence that education affects your marriage status."

#### $\chi$ -Square Test: From Documents

- 1 import numpy as np
- 2 import pandas as pd
- 3 from scipy import stats
- 4 import matplotlib.pyplot as plt
- 5 import numpy as np
- 6 data = pd.read\_csv('https://ximarketing.github.io/data/python\_data3.csv')
- 7 table = pd.crosstab(data.Gender, columns=data.Purchase)
- 8 print(table)
- 9 print(stats.chi2\_contingency(table, correction=False))

Logistic Regression What's it?

#### In short, a logistic regression specifies the following function:

$$\Pr[ ext{Y=1}] = rac{\exp(eta X)}{1+\exp(eta X)}$$

Why do we choose this "weird" function?

Exercise: Consider the following dataset https://ximarketing.github.io/data/loan.csv

Age	Income	LoanAmount	CreditScore	MonthsEmployed	InterestRate	LoanTerm	Education	EmploymentType	MaritalStatus	HasMortgage	LoanPurpose	Default
56	85994	50587	520	80	15.23	36	Bachelors	Full-time	Divorced	Yes	Other	0
69	50432	124440	458	15	4.81	60	Masters	Full-time	Married	No	Other	0
46	84208	129188	451	26	21.17	24	Masters	Unemployed	Divorced	Yes	Auto	1
32	31713	44799	743	0	7.07	24	High School	Full-time	Married	No	Business	0
60	20437	9139	633	8	6.51	48	Bachelors	Unemployed	Divorced	No	Auto	0
25	90298	90448	720	18	22.72	24	High School	Unemployed	Single	Yes	Business	1

The data is about the loan default information, where the outcome is Default (1 = default, 0 = no default).

Education includes: high school, masters, bachelors, and PhD

Employment type includes: Full-time, part-time, unemployed, and selfemployed

Marital status includes: single, married, and divorced

Loan purpose includes: auto, business, education, home, and other.

Age	Income	LoanAmount	CreditScore	MonthsEmployed	InterestRate	LoanTerm	Education	EmploymentType	MaritalStatus	HasMortgage	LoanPurpose	Default
56	85994	50587	520	80	15.23	36	Bachelors	Full-time	Divorced	Yes	Other	0
69	50432	124440	458	15	4.81	60	Masters	Full-time	Married	No	Other	0
46	84208	129188	451	26	21.17	24	Masters	Unemployed	Divorced	Yes	Auto	1
32	31713	44799	743	0	7.07	24	High School	Full-time	Married	No	Business	0
60	20437	9139	633	8	6.51	48	Bachelors	Unemployed	Divorced	No	Auto	0
25	90298	90448	720	18	22.72	24	High School	Unemployed	Single	Yes	Business	1

We would like to run a logistic regression to predict which type of borrowers are more likely to default in the future.

#### •••

```
1 import pandas as pd
 2 import statsmodels.api as sm
  url = 'https://ximarketing.github.io/data/loan.csv'
  data = pd.read csv(url)
  categorical vars = ['Education', 'EmploymentType', 'MaritalStatus',
   'HasMortgage', 'LoanPurpose']
  data = pd.get dummies(data, columns=categorical vars, drop first=True)
 8
 9
10 numerical vars = ['Age', 'Income', 'LoanAmount', 'CreditScore',
   'MonthsEmployed', 'InterestRate', 'LoanTerm']
11 X = data[numerical vars + [col for col in data.columns if
   col.startswith(tuple(categorical vars))]]
12 y = data['Default']
13 X = sm.add constant(X)
14 model = sm.Logit(y, X.astype(float))
15 result = model.fit()
16
17 print(result.summary())
```

### How to interpret this result?

Age	-0.0392	0.000	-84.774	0.000	-0.040	-0.038
Income	-8.753e-06	1.7e-07	-51.488	0.000	-9.09e-06	-8.42e-06
LoanAmount	4.229e-06	9.31e-08	45.398	0.000	4.05e-06	4.41e-06
CreditScore	-0.0008	4.09e-05	-18.343	0.000	-0.001	-0.001
MonthsEmployed	-0.0097	0.000	-50.927	0.000	-0.010	-0.009
InterestRate	0.0687	0.001	67.565	0.000	0.067	0.071
LoanTerm	9.458e-05	0.000	0.248	0.804	-0.001	0.001
Education_High School	0.0783	0.018	4.404	0.000	0.043	0.113
Education_Masters	-0.1301	0.018	-7.059	0.000	-0.166	-0.094
Education_PhD	-0.1759	0.019	-9.490	0.000	-0.212	-0.140
EmploymentType_Part-time	0.2816	0.019	14.790	0.000	0.244	0.319
<pre>EmploymentType_Self-employed</pre>	0.2362	0.019	12.293	0.000	0.199	0.274
EmploymentType_Unemployed	0.4416	0.019	23.672	0.000	0.405	0.478
MaritalStatus_Married	-0.2267	0.016	-14.155	0.000	-0.258	-0.195
MaritalStatus_Single	-0.0644	0.016	-4.139	0.000	-0.095	-0.034
HasMortgage_Yes	-0.1560	0.013	-12.007	0.000	-0.181	-0.130
LoanPurpose_Business	0.0442	0.020	2.193	0.028	0.005	0.084
LoanPurpose_Education	-0.0169	0.020	-0.828	0.408	-0.057	0.023
LoanPurpose_Home	-0.1935	0.021	-9.234	0.000	-0.235	-0.152
LoanPurpose_Other	-0.0071	0.020	-0.346	0.729	-0.047	0.033

Multinomial Logit Model What's it?

### Multinomial Logit Model

In a linear regression, our dependent variable is a number. In a logistic regression, our dependent variable is a probability, e.g., the probability of success. In a multinomial logit regression, our dependent variable is a choice, usually a choice among different brands. In the above example, HKU president's choice is one of the four brands, BMW, Tesla, Audi, and Mercedes Benz.

#### Suppose that consumers have three choices, A, B, C.

Now, given  $X_i$ , we would like to come up with three functions  $f_A(X_i)$ ,  $f_B(X_i)$  and  $f_C(X_i)$ , such that

 $egin{aligned} & \Pr[Y_i = A] pprox f_A(X_i), \ & \Pr[Y_i = B] pprox f_B(X_i), \ & \Pr[Y_i = C] pprox f_C(X_i). \end{aligned}$ 

#### These functions are:

$$f_A(X_i) = rac{\exp(lpha_A+eta_A X_i)}{\exp(lpha_A+eta_A X_i)+\exp(lpha_B+eta_B X_i)+\exp(lpha_C+eta_C X_i)}$$

$$f_B(X_i) = rac{\exp(lpha_B+eta_BX_i)}{\exp(lpha_A+eta_AX_i)+\exp(lpha_B+eta_BX_i)+\exp(lpha_C+eta_CX_i)}$$

$$f_C(X_i) = rac{\exp(lpha_C+eta_C X_i)}{\exp(lpha_A+eta_A X_i)+\exp(lpha_B+eta_B X_i)+\exp(lpha_C+eta_C X_i)}$$

## Let's analyze the following data: https://ximarketing.github.io/data/bankchoice.csv

	Choice	Age	Female	Income	Education	Job
1	CCB	62	0	7	2	Industry
2	CCB	34	1	4	5	Retired
3	CCB	68	0	5	2	Industry
4	CCB	60	0	3	2	Education
5	ICBC	18	0	6	2	Industry
6	CCB	18	0	4	3	Student
7	CCB	51	0	7	3	Education
8	CCB	25	1	3	2	Unemployed
9	BOC	42	1	4	4	Education
10	CCB	71	1	5	3	Service
11	BOC	23	0	4	5	Student
12	BOC	30	0	2	5	Retired

### The code I got from GPT:

#### 

```
1 import pandas as pd
 2 import statsmodels.api as sm
 3 from patsy import dmatrices
 4 data = pd.read csv('https://ximarketing.github.io/data/bankchoice.csv')
 5 y, X = dmatrices('Choice ~ Age + Income + Education + C(Job, Treatment)',
   data, return type='dataframe')
 6 model = sm.MNLogit(y, X).fit()
 7 summary = model.summary2()
 8 pd.set option( 'display.max columns', None)
  pd.set option('display.width', None)
 9
   choices = y.design info.column names
10
   choices.pop()
11
   for i, choice in enumerate(choices, start=1):
12
13
       print(f"Choice: {choice}")
14
       print(summary.tables[i])
```

```
15 print('\n')
```

## Do you understand the output?

#### Choice: Choice[ABC]

		y = 0	Coef.	Std.Err.	t	P> t	[0.025	0.975]
Interce	ept	Intercept	0.473001	0.195306	2.421850	1.544173e-02	0.090209	0.855793
C(Job,	Treatment)[T.Finance]	C(Job, Treatment)[T.Finance]	4.291024	1.009358	4.251242	2.125879e-05	2.312719	6.269329
C(Job,	Treatment)[T.Government]	C(Job, Treatment)[T.Government]	0.520523	0.246418	2.112360	3.465559e-02	0.037553	1.003493
C(Job,	Treatment)[T.Industry]	C(Job, Treatment)[T.Industry]	2.442084	0.518532	4.709615	2.481853e-06	1.425781	3.458387
C(Job,	Treatment)[T.Retired]	C(Job, Treatment)[T.Retired]	-2.563242	0.145518	-17.614645	1.901798e-69	-2.848451	-2.278032
C(Job,	Treatment)[T.Service]	C(Job, Treatment)[T.Service]	-0.042547	0.189207	-0.224869	8.220814e-01	-0.413386	0.328292
C(Job,	Treatment)[T.Student]	C(Job, Treatment)[T.Student]	-1.018968	0.161727	-6.300563	2.965671e-10	-1.335947	-0.701990
C(Job,	Treatment)[T.Unemployed]	C(Job, Treatment)[T.Unemployed]	0.226454	0.182337	1.241956	2.142528e-01	-0.130919	0.583828
Age		Age	-0.032753	0.002334	-14.031712	9.971443e-45	-0.037328	-0.028178
Income		Income	0.613154	0.028165	21.769850	4.480917e-105	0.557951	0.668357
Educati	ion	Education	1.272112	0.037794	33.658878	2.312053e-248	1.198036	1.346187

Choice	Choice[BOC]							
		y = 1	Coef.	Std.Err.	t	P> t	[0.025	0.975]
Interce	ept	Intercept	0.836615	0.193845	4.315904	1.589513e-05	0.456686	1.216543
C(Job,	Treatment)[T.Finance]	C(Job, Treatment)[T.Finance]	4.021791	1.009072	3.985632	6.730076e-05	2.044046	5.999537
C(Job,	Treatment)[T.Government]	C(Job, Treatment)[T.Government]	0.843469	0.244775	3.445892	5.691779e-04	0.363718	1.323220
C(Job,	Treatment)[T.Industry]	C(Job, Treatment)[T.Industry]	2.159684	0.517983	4.169410	3.053894e-05	1.144456	3.174913
C(Job,	Treatment)[T.Retired]	C(Job, Treatment)[T.Retired]	-3.167509	0.144007	-21.995543	3.177106e-107	-3.449758	-2.885261
C(Job,	Treatment)[T.Service]	C(Job, Treatment)[T.Service]	-0.328630	0.187726	-1.750587	8.001701e-02	-0.696565	0.039305
C(Job,	Treatment)[T.Student]	C(Job, Treatment)[T.Student]	-1.415400	0.160125	-8.839354	9.627386e-19	-1.729239	-1.101561
C(Job,	Treatment)[T.Unemployed]	C(Job, Treatment)[T.Unemployed]	-0.858125	0.181573	-4.726057	2.289217e-06	-1.214002	-0.502248
Age		Age	-0.017986	0.002319	-7.756240	8.748421e-15	-0.022530	-0.013441
Income		Income	0.882860	0.028079	31.441465	5.491352e-217	0.827825	0.937895
Educati	ion	Education	0.809948	0.037629	21.524823	9.116791e-103	0.736197	0.883699

Choice: Choice[CCB]

		y = 2	Coef.	Std.Err.	t	P> t	[0.025	0.975]
Intercep	ot	Intercept	-0.709650	0.216728	-3.274379	1.058945e-03	-1.134429	-0.284871
C(Job, T	reatment)[T.Finance]	C(Job, Treatment)[T.Finance]	5.027076	1.013953	4.957901	7.125898e-07	3.039766	7.014387
C(Job, T	[reatment)[T.Government]	C(Job, Treatment)[T.Government]	2.005531	0.263119	7.622133	2.495179e-14	1.489827	2.521236
C(Job, T	[reatment)[T.Industry]	C(Job, Treatment)[T.Industry]	4.737702	0.525177	9.021162	1.861047e-19	3.708375	5.767029
C(Job, T	reatment)[T.Retired]	C(Job, Treatment)[T.Retired]	-0.591359	0.168584	-3.507807	4.518174e-04	-0.921776	-0.260941
C(Job, T	reatment)[T.Service]	C(Job, Treatment)[T.Service]	1.289505	0.209591	6.152477	7.628194e-10	0.878714	1.700296
C(Job, T	reatment)[T.Student]	C(Job, Treatment)[T.Student]	0.707630	0.183555	3.855134	1.156664e-04	0.347868	1.067391
C(Job, T	<pre>reatment)[T.Unemployed]</pre>	C(Job, Treatment)[T.Unemployed]	1.207999	0.203195	5.945033	2.764005e-09	0.809745	1.606253
Age		Age	-0.003158	0.002413	-1.308492	1.907064e-01	-0.007888	0.001572
Income		Income	0.517294	0.028628	18.069586	5.532719e-73	0.461184	0.573404
Educatio	n	Education	0.209854	0.039220	5.350756	8.758775e-08	0.132985	0.286723

## Handwriting Digit Prediction

We can further build a handwriting digit prediction model based on multinomial logit regression.

Input:  $x_i \in \{0, 1\}$ , the color of pixel *i*, black vs. white. Output:  $Y \in \{0, 1, \dots, 9\}$ .

## Handwriting Digit Prediction

Let's try the following one, which is created by GPT.

🦸 Han	_		×
	Predi	ict	
	Clea	er	
Draw a	digit and	I click Pre	edict

#### Code for data training

#### •••

```
1 from sklearn.linear model import LogisticRegression
 2 from sklearn.datasets import fetch openml
 3 from sklearn.model selection import train test split
 4 from sklearn.metrics import accuracy score
 5 import joblib
 7 mnist = fetch openml('mnist 784', version=1)
9 X = mnist.data / 255.0
10 y = mnist.target.astype(int)
11 # Split the data into training and testing sets
12 X train, X test, y train, y test = train test split(X, y, test size=0.25, random state=42)
13 # Train the logistic regression model
14 model = LogisticRegression(solver='lbfgs', max iter=1000, multi class='multinomial')
15 model.fit(X train, y train)
16 # Evaluate the model
17 y pred = model.predict(X test)
18 accuracy = accuracy score(y test, y pred)
19 print(f'Accuracy: {accuracy * 100:.2f}%')
20 # Save the model
21 joblib.dump(model, 'logistic regression model.pkl')
22 import tkinter as tk
23 from tkinter import *
24 import numpy as np
25 from PIL import Image, ImageDraw, ImageOps
26 import joblib
27 # Load the trained model (make sure you have a trained model saved as 'logistic regression mode
28 model = joblib.load('logistic regression model.pkl')
```

- Logistic regression is the extension of linear regression.
- Multinomial logit model is the extension of logistic regression.
- Artificial neural network (and deep neural network) is the extension of multinomial logit model.



https://www.youtube.com/embed/bfmFfD2RIcg?enablejsapi=1

An artificial neural network has a few layers:

- Input layer: Which takes input from data  $x_i$
- Output layer: Which produces an out  $\Pr[y_i = 1]$
- Hidden layers: There are one or more hidden layers which are responsible for making calculations (i.e., logistic function).

Why do we create artificial neural networks like this? This is because human brains also work in similar ways! We also have cells in our brain which calculates logistic function! So, an artifical neural network simply mimics our human brain.

In mathematics, there is a famous "universal approximation theorem," which states that, if you have an artificial neural work that is large enough, theoretically, you should be able to use it to approximate any functions... What can we do?

- Autonomous driving
- Play GO
- Digit recognition...

Again, let us build an artificial neural network in python for digit recognition with the help of GPT!

Again, let us build an artificial neural network in python for digit recognition with the help of GPT!

Note: You will need to use the module tensorflow for your ANN. The installation of tensorflow is a bit different.

# Installing TensorFlow



On your Anaconda, lunch anaconda prompt
# Installing TensorFlow

C:\WINDOWS\system32\cmd.exe

Microsoft Windows [Version 10.0.19045.5487] (c) Microsoft Corporation. All rights reserved.

(base) C:\Users\FBE>pip install tensorflow

Input "pip install tensorflow", and enter the command

# Pie Chart

# Can you create a pie chart that visualize the sales of different EV brands?



#### Code generated by GPT

#### 

```
1 import matplotlib.pyplot as plt
2 # Data for car sales
3 brands = ['BYD', 'Tesla', 'Zeekr', 'BMW', 'MG']
4 sales = [848, 202, 195, 144, 92]
5 # Create a pie chart
6 plt.figure(figsize=(8, 8))
7 plt.pie(sales, labels=brands, autopct='%1.1f%%', startangle=140)
8 plt.axis('equal') # Equal aspect ratio ensures that pie is drawn as a circle.
9 # Add a title
10 plt.title('Car Sales Distribution')
11 # Display the pie chart
12 plt.show()
```



#### I am making it more interactive (plotly needed)

#### 

- 1 import plotly.graph\_objects as go
- 2 # Data for car sales
- 3 brands = ['BYD', 'Tesla', 'Zeekr', 'BMW', 'MG']
- 4 sales = [848, 202, 195, 144, 92]
- 5 # Create a pie chart trace
- 7 # Update layout for interactivity
- 8 fig.update\_layout(title='Car Sales Distribution', showlegend=False)
- 9 # Show the figure
- 10 fig.show()

# Dumbbell chart

#### Try to create a dumbbell chart like this one!



The data is available here:

http://ximarketing.github.io/data/cancer.csv

Try to use GPT to create the dumbbell chart for you!

#### •••

```
1 import matplotlib.pyplot as plt
 2 import pandas as pd
 3 data = pd.read csv('https://ximarketing.github.io/data/cancer.csv')
 4 # Sort the data based on 2010 survival rate
 5 data sorted = data.sort values(by='2010', ascending=False)
 6 data sorted = data sorted.iloc[::-1]
 7 data sorted.reset index(drop=True, inplace=True)
 8 # Plotting
 9 plt.figure(figsize=(10, 6))
10 plt.subplots adjust(left=0.22) # Adjust the left margin
11 for i, row in data sorted.iterrows():
       plt.plot([row['2010'], row['1970']], [i, i], marker='o', color='b') #
12
   Dumbbell circles
       plt.annotate('', xy=(row['2010'], i), xytext=(row['1970'], i),
13
14
                    arrowprops=dict(arrowstyle='->', connectionstyle='arc3',
   color='gray')) # Arrows
15 # Adding labels and title
16 plt.yticks(range(len(data sorted)), data sorted['Cancer'])
17 plt.xlabel('Survival Rate (%)')
18 plt.ylabel('Cancers')
19 plt.title('Cancer Survival Rate Comparison (1970 vs. 2010)')
20 # Showing the plot
21 plt.grid(axis='x', linestyle='--', linewidth=0.5)
22 plt.show()
```

Generating Visualization by Uploading Images

### Example: Generate a Donut Chart Image can be downloaded here



	Xi Li	
	Give me the python code to produce this donut chart	
	JPG 14.71 KB	19/02/2025, 08:51:42
GPT-40 (Tokens used 685)		

matplotlib library.

#### My Output from GPT



#### Example: Generate a Population Pyramid Image can be downloaded here



example code:

<u>---</u>

 $\circ$ 

#### •••

```
1 import matplotlib.pyplot as plt
 2 import numpy as np
 3 # Sample data similar to the provided population pyramid
 4 age groups = ['0-4', '5-9', '10-14', '15-19', '20-24', '25-29', '30-34', '35-39']
                 '40-44', '45-49', '50-54', '55-59', '60-64', '65-69', '70-74',
                 '75-79', '80-84', '85-89', '90-94', '95-99', '100+']
 7 male population = [20000, 18000, 17000, 16000, 15000, 14000, 13000, 12000,
                      11000, 10000, 9000, 8000, 7000, 6000, 5000, 4000,
                      3000, 2000, 1000, 500, 1001
10 female population = [19000, 17000, 16000, 15000, 14000, 13000, 12000,
11
                        11000, 10000, 9000, 8000, 7000, 6000, 5000, 4000,
12
                        3000, 2000, 1000, 500, 300, 100]
13 # Convert the populations to numpy arrays for easier manipulation
14 male population = np.array(male population)
15 female population = np.array(female population)
16 # Create the figure and axis
17 fig, ax = plt.subplots(figsize=(10, 8))
18 # Plot the male and female bars
19 ax.barh(age groups, male population, color='blue', label='Males')
20 ax.barh(age groups, -female population, color='red', label='Females')
21 # Add labels and title
22 ax.set xlabel('Population')
23 ax.set title('Population Pyramid for Hong Kong, December 2020')
24 ax.legend()
25 # Customize the x-axis to show positive numbers for female population
26 ax.set xticks(np.arange(-20000, 21000, 5000))
27 ax.set xticklabels([str(abs(x)) for x in np.arange(-20000, 21000, 5000)])
28 # Show the plot
29 plt.show()
```

#### How about an animated GIF? Image can be downloaded here

#### Xi Li

I want a dynamic GIF like this showing how population pyramid changes over time. Show me python code





#### GPT-40 (Tokens used 1219)

Creating a dynamic GIF like the one you've shown requires several steps, including data preparation, plotting the population pyramids for each year, and then combining these plots into an animated GIF. Here's an example of how you can achieve this using Python with libraries such as matplotlib, pandas, and imageio.

#### Generate a Figure like this -- I don't even know its name! Image can be downloaded here



#### Not too bad...



## Generate a Figure like this -- I don't even know its name! Image can be downloaded here



#### After a few trials, I got it here:



### Generate a Figure like this -- It's called a Chord Diagram. Image can be downloaded here



## It's better than the original one. (It may not run on Spyder; try Jupyter Notebook on Anaconda)



In this assignment, you need to develop a Python desktop APP using the "tkinter" module with the help of GPT. The main function of the APP is to load, analyze and visualize data. The requirements are as follows.

1. Data loading: The APP should be able to

- load data itself from an online source (e.g., URL) OR,
- allow the user to input data herself manually OR,
- allow the user to upload a data file to the APP from her device.

In this assignment, you need to develop a Python desktop APP using the "tkinter" module with the help of GPT. The main function of the APP is to load, analyze and visualize data. The requirements are as follows.

2. Data analysis: The APP has a user-interface which allows the user to simply (e.g., by clicking a button) perform some data analyses (chosen by the developer).

In this assignment, you need to develop a Python desktop APP using the "tkinter" module with the help of GPT. The main function of the APP is to load, analyze and visualize data. The requirements are as follows.

3. Data visualization: The APP should offer users of the APP a convenient way to visualize the output of data analysis.

# Example

An APP that loads data from a URL. Then, the user can choose the dependent variable and independent variable and click a button to run a regression. The APP demonstrates the regression equation, visualizes the result and explains the meaning.

# Example

An APP that asks the user for data and runs different types of *t*-tests. The user chooses the type of *t*-test and inputs the data himself/herself. Then, the APP visualizes the result and explains the meaning to the user.

## Deliverables

- Your code;
- Your prompt when using GPT or other LLMs, which explains how you created the APP;
- If the APP needs any external data files, also submit your data files (size limit is 5 MB).
- Make sure your code runs well.

#### Deadline

#### April 5, 2025, 23:59 PM.

This is a group assignment. We may select and demonstrate some good projects in class.