

60



ChatGPT Prompts

for
Data
Science



By Travis Tang
<https://www.linkedin.com/in/travistang>

1. Train Classification Model

I want you to act as a data scientist and code for me. I have a dataset of [describe dataset]. Please build a machine learning model that predict [target variable].



I want you to act as a data scientist and code for me. I have a dataset of customer characteristics and churn. Please build a machine learning model that predict if a customer will churn. Please do not write code for data visualisation and exploration.



To build a machine learning model for predicting customer churn, we will first need to import the necessary libraries and load the data into our environment.

Copy code

```
# Import libraries
import pandas as pd
import numpy as np

# Load data
data = pd.read_csv('data.csv')
```

Next, we will handle missing values and encode categorical variables, if necessary. We will also need to ensure that the data is in the correct format for modeling. This may involve tasks such as scaling numerical columns or one-hot encoding categorical variables.

Copy code

```
# Handle missing values
data = data.fillna(data.mean())

# One-hot encode categorical variables
data = pd.get_dummies(data)

# Scale numerical columns
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
data[['num_column_1', 'num_column_2', 'num_column_3']] =
scaler.fit_transform(data[['num_column_1', 'num_column_2',
'num_column_3']])
```

Now that the data is cleaned and prepared, we can split it into training and testing sets and build a machine learning model.

 Copy code

```
# Split data into training and testing sets
X = data.drop('churn', axis=1)
y = data['churn']
X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2, random_state=42)

# Build model
from sklearn.ensemble import RandomForestClassifier

model = RandomForestClassifier()
model.fit(X_train, y_train)
```

Finally, we can evaluate the model on the testing set.

 Copy code

```
# Evaluate model
from sklearn.metrics import accuracy_score

predictions = model.predict(X_test)
accuracy = accuracy_score(y_test, predictions)
print('Accuracy:', accuracy)
```

2. Automatic Machine Learning

I want you to act as an automatic machine learning (AutoML) bot using TPOT for me. I am working on a model that predicts [...]. Please write python code to find the best classification model with the highest AUC score on the test set.

The image shows a screenshot of a ChatGPT interface. The user's message is at the top, followed by a response from the AI. The AI has provided a Python script to use TPOT for building a classification model.

```
import pandas as pd
from sklearn.model_selection import train_test_split
from tpot import TPOTClassifier

# Load the data
df = pd.read_csv("credit_card_default_data.csv")

# Extract the target column
y = df["default"]

# Extract the features
X = df.drop("default", axis=1)

# Split the data into training and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2)

# Create a TPOT classifier
tpot = TPOTClassifier(scoring="roc_auc", verbosity=2)

# Fit the classifier to the training data
tpot.fit(X_train, y_train)

# Score the classifier on the test set
print(tpot.score(X_test, y_test))

# Export the best model as a Python script
tpot.export("best_model.py")
```

3. Tune Hyperparameter

I want you to act as a data scientist and code for me. I have trained a [model name]. Please write the code to tune the hyper parameters.



I want you to act as a data scientist and code for me. I have trained a scikit-learn decision tree. Please write the code to tune the hyper parameters.



To tune the hyperparameters of a decision tree model using scikit-learn, we can use the `GridSearchCV` function from the `sklearn.model_selection` module. This function allows us to specify a range of values for each hyperparameter, and it will train and evaluate a model for each combination of hyperparameter values.

Copy code

```
# Import libraries
import pandas as pd
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import GridSearchCV

# Load data
data = pd.read_csv('data.csv')
X = data.drop('target', axis=1)
y = data['target']

# Create a decision tree classifier
model = DecisionTreeClassifier()

# Define the hyperparameter grid
param_grid = {'max_depth': [2, 4, 6, 8, 10],
              'min_samples_leaf': [2, 4, 6, 8, 10]}

# Create the grid search object
grid_search = GridSearchCV(model, param_grid, cv=5,
                           scoring='accuracy')

# Fit the grid search object to the training data
grid_search.fit(X, y)

# Get the best hyperparameters
print('Best max_depth:', grid_search.best_params_['max_depth'])
print('Best min_samples_leaf:', grid_search.best_params_['min_samples_leaf'])
```

4. Explore Data

I want you to act as a data scientist and code for me. I have a dataset of [describe dataset]. Please write code for data visualisation and exploration.



I want you to act as a data scientist and code for me. I have a dataset of customer characteristics and churn. Please write code for data visualisation and exploration.



To begin exploring and visualizing the customer churn data, we will first need to import the necessary libraries and load the data into our environment.

Copy code

```
# Import libraries
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

# Load data
data = pd.read_csv('data.csv')

# Summary statistics of the numerical columns
data.describe()

# Count plots of the categorical columns
for col in data.select_dtypes(include='object').columns:
    sns.countplot(x=col, data=data)
    plt.show()

# Correlation matrix of the numerical columns
sns.heatmap(data.select_dtypes(exclude='object').corr(), annot=True)

# Box plots of the numerical columns by target
for col in data.select_dtypes(exclude='object').columns:
    sns.boxplot(x='churn', y=col, data=data)
    plt.show()

# Bar plots of the categorical columns by target
for col in data.select_dtypes(include='object').columns:
    sns.barplot(x='churn', y=col, data=data)
    plt.show()
```

5. Generate Data

I want you to act as a fake data generator. I need a dataset that has x rows and y columns: [insert column names]

The screenshot shows a ChatGPT interface. The user (represented by a person icon) asks for a dataset with 1,000 rows and 3 columns: "customer_id", "product", and "revenue". The AI (represented by a green swirl icon) responds with a Python code snippet to generate such a dataset.

```
import random
import pandas as pd

# Set the product names
products = ["Product A", "Product B", "Product C", "Product D",
"Product E", "Product F"]

# Set the customer names
customers = ["Customer 1", "Customer 2", "Customer 3", "Customer 4",
"Customer 5"]

# Create an empty dataframe
df = pd.DataFrame(columns=["customer_id", "product", "revenue"])

# Generate 1,000 rows of fake data
for i in range(1000):
    # Select a random product
    product = random.choice(products)

    # Select a random customer
    customer = random.choice(customers)

    # Generate a random revenue amount
    revenue = random.randint(100, 10000)

    # Add the data to the dataframe
    df = df.append({"customer_id": customer, "product": product,
"revenue": revenue}, ignore_index=True)

# Preview the data
df.head()
```

6. Write Regex

I want you to act as a coder. Please write me a regex in python that [describe regex]



I want you to act as a coder. Please write me a regex in python that would match any strings starts with an alphanumeric character, contains the word "data", contains the word "science", and ends with an exclamation point.



Sure! Here is a regex in python that will match strings that start with an alphanumeric character, contain the word "data", contain the word "science", and end with an exclamation point:

Copy code

```
import re

def match_string(string):
    regex = r"^[a-zA-Z0-9].*data.*science.*!$"
    return bool(re.match(regex, string))
```

This regex uses the following elements:

- `^`: This matches the start of the string.
- `[a-zA-Z0-9]`: This matches any alphanumeric character.
- `.*`: This matches any number of any characters (except a newline), including none.
- `data`: This matches the string "data".
- `science`: This matches the string "science".
- `!`: This matches the exclamation point.
- `\$`: This matches the end of the string.

The `match_string` function will return `True` if the string matches the regex, and `False` if it does not.

7. Train Time Series

I want you to act as a data scientist and code for me. I have a time series dataset [describe dataset]. Please build a machine learning model that predict [target variable]. Please use [time range] as train and [time range] as validation.



I want you to act as a data scientist and code for me. I have a time series dataset of sales revenue from 2019 to 2022. Please build a machine learning model that predict sales revenue for 12 months in 2023. Please use 2019 to 2021 data as train and 2022 data as validation set. Please do not write code for data visualisation and exploration.



To build a machine learning model for forecasting future sales revenue, we will first need to import the necessary libraries and load the data into our environment.

Copy code

```
# Import libraries
import pandas as pd
import numpy as np

# Load data
data = pd.read_csv('data.csv')
```

Next, we will handle missing values and encode categorical variables, if necessary. We will also need to ensure that the data is in the correct format for modeling. This may involve tasks such as scaling numerical columns and creating lag variables.

Copy code

```
# Handle missing values
data = data.fillna(data.mean())

# One-hot encode categorical variables
data = pd.get_dummies(data)

# Scale numerical columns
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
data[['sales']] = scaler.fit_transform(data[['sales']])

# Create lag variables
for i in range(1, 13):
    data['lag_{}'.format(i)] = data.sales.shift(i)
data = data.dropna()
```

Now that the data is cleaned and prepared, we can split it into training and validation sets and build a machine learning model.

 Copy code

```
# Split data into training and validation sets
train = data[data.date < '2022']
valid = data[data.date >= '2022']

X_train = train.drop('sales', axis=1)
y_train = train['sales']
X_valid = valid.drop('sales', axis=1)
y_valid = valid['sales']
```

To build the model, we will use a random forest regressor from the `sklearn` library.

 Copy code

```
# Build model
from sklearn.ensemble import RandomForestRegressor

model = RandomForestRegressor()
model.fit(X_train, y_train)
```

Finally, we can evaluate the model on the validation set.

 Copy code

```
# Evaluate model
from sklearn.metrics import mean_absolute_error

predictions = model.predict(X_valid)
mae = mean_absolute_error(y_valid, predictions)
print('Mean Absolute Error:', mae)
```

8. Address Imbalance Data

I want you to act as a coder. I have trained a machine learning model on an imbalanced dataset. The predictor variable is the column [Insert column name]. In python, how do I oversample and/or undersample my data?



I want you to act as a coder. I have trained a machine learning model on an imbalanced dataset. The predictor variable is the column "scam". The dataset has 1% scam dataset and 99% non-scam dataset. In python, how do I oversample and/or undersample my data?



To oversample the minority class in your dataset, you can use the `SMOTE` (Synthetic Minority Oversampling Technique) method, which is available in the `imblearn` library. Here is an example of how to use it:

Copy code

```
from imblearn.over_sampling import SMOTE

# Split the dataset into features and labels
X = df.drop('scam', axis=1)
y = df['scam']

# Oversample the minority class using SMOTE
smote = SMOTE()
X_over, y_over = smote.fit_resample(X, y)
```

To undersample the majority class, you can use the `RandomUnderSampler` method from the `imblearn` library. Here is an example of how to use it:

Copy code

```
from imblearn.under_sampling import RandomUnderSampler

# Split the dataset into features and labels
X = df.drop('scam', axis=1)
y = df['scam']

# Undersample the majority class using RandomUnderSampler
rus = RandomUnderSampler()
X_under, y_under = rus.fit_resample(X, y)
```

9. Get Feature Importance

I want you to act as a data scientist and explain the model's results. I have trained a decision tree model and I would like to find the most important features. Please write the code.

 I want you to act as a data scientist and explain the model's results. I have trained a decision tree model and I would like to find the most important features. Please write the code.

 Certainly! Here is some code that you can use to find the most important features in a decision tree model:

Copy code

```
import numpy as np

# Load your model and the data that you used to train it
model = # Load your model
X_train = # Load your training data

# Get the feature importances from the model
importances = model.feature_importances_

# Sort the feature importances in descending order
indices = np.argsort(importances)[::-1]

# Print the feature ranking
print("Feature ranking:")

for f in range(X_train.shape[1]):
    print("%d. feature %d (%f)" % (f + 1, indices[f],
                                    importances[indices[f]]))
```

This code will use the `feature_importances_` attribute of the decision tree model to compute the importance of each feature in the model. It will then sort the features in descending order by their importance and print out the ranking.

Note that the feature importances returned by the model are relative to each other, so they can be used to compare the importance of different features, but they are not directly comparable to the values of the features themselves.

10. Visualize Data with Matplotlib

I want you to act as a coder in python. I have a dataset [name] with columns [name]. [Describe graph requirements]



I want you to act as a coder in python. I have a dataset called iris_df. It has the columns Petal Length, Petal Width, Sepal Length, Sepal Width and Species Type. On the left, I want a graph that depicts the relationship between Petal Width and Petal Length. On the right, I want a graph that depicts the relationship between Sepal Width and Sepal Length. For both graphs, please differentiate between the different species. Please make it visually pleasing and colorblind-friendly. Please use matplotlib, ensure that the graph is large enough for a presentation, and save the output to a file called "output.png"



Certainly! Here is some code that should accomplish what you've described:

Copy code

```
import matplotlib.pyplot as plt

# Extract data for Petal Width and Petal Length
petal_width = iris_df['Petal Width']
petal_length = iris_df['Petal Length']
species = iris_df['Species Type']

# Create a figure with two subplots
fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(10, 5))

# Plot the relationship between Petal Width and Petal Length on the left
for spec in set(species):
    spec_df = iris_df[iris_df['Species Type'] == spec]
    ax1.scatter(spec_df['Petal Width'], spec_df['Petal Length'],
               label=spec)
    ax1.set_xlabel('Petal Width')
    ax1.set_ylabel('Petal Length')
    ax1.legend()

# Plot the relationship between Sepal Width and Sepal Length on the right
for spec in set(species):
    spec_df = iris_df[iris_df['Species Type'] == spec]
    ax2.scatter(spec_df['Sepal Width'], spec_df['Sepal Length'],
               label=spec)
    ax2.set_xlabel('Sepal Width')
    ax2.set_ylabel('Sepal Length')
    ax2.legend()

# Save the figure to a file
fig.savefig('output.png')

# Show the plot
plt.show()
```

11. Visualize Image Grid Matplotlib

I want you to act as a coder. I have a folder of images.

[Describe how files are organised in directory] [Describe how you want images to be printed]



I want you to act as a coder. I have a folder of images. The subdirectories are "TRAIN" and "TEST". Within "Train" and "Test" are 10 folders from "0" to "9". Within each folder, there are JPG, JPEG and PNG images. Write code that randomly selects 20 images from the TRAIN subdirectory and display them in a 4 by 5 grid. Repeat the same for the test dataset.



Sure! Here is some code that should accomplish what you've described:

Copy code

```
import os
import random
import matplotlib.pyplot as plt
from matplotlib.image import imread

# Set the base directory
base_dir = 'path/to/your/folder'

# Set the subdirectories
train_dir = os.path.join(base_dir, 'TRAIN')
test_dir = os.path.join(base_dir, 'TEST')

# Set the number of images to display
num_images = 20
```

```
# Select random images from the train directory
train_images = []
for i in range(10):
    subdir = os.path.join(train_dir, str(i))
    files = [os.path.join(subdir, f) for f in os.listdir(subdir) if
f.endswith('.jpg', '.jpeg', '.png')]]
    train_images += random.sample(files, num_images // 10)

# Select random images from the test directory
test_images = []
for i in range(10):
    subdir = os.path.join(test_dir, str(i))
    files = [os.path.join(subdir, f) for f in os.listdir(subdir) if
f.endswith('.jpg', '.jpeg', '.png')]]
    test_images += random.sample(files, num_images // 10)

# Display the train images
fig, ax = plt.subplots(4, 5, figsize=(20, 20))
for i, image_path in enumerate(train_images):
    ax[i//5][i%5].imshow(imread(image_path))
plt.show()

# Display the test images
fig, ax = plt.subplots(4, 5, figsize=(20, 20))
for i, image_path in enumerate(test_images):
    ax[i//5][i%5].imshow(imread(image_path))
plt.show()
```

12. Explain Model with Lime

I want you to act as a data scientist and explain the model's results. I have trained a [library name] model and I would like to explain the output using LIME. Please write the code.



I want you to act as a data scientist and explain the model's results. I have trained a scikit-learn K-nearest-neighbour model and I would like to explain the output using LIME. Please write the code.



Certainly! Here is some code that you can use to explain the results of a scikit-learn K-nearest-neighbors model using LIME (Local Interpretable Model-agnostic Explanations):

Copy code

```
import lime
import lime.lime_tabular
import numpy as np

# Load your model and the data that you used to train it
model = # Load your model
X_train = # Load your training data
X_test = # Load the test data that you want to explain

# Create a LIME explainer object
explainer = lime.lime_tabular.LimeTabularExplainer(X_train,
feature_names=feature_names, class_names=class_names,
kernel_width=3)

# Choose an instance to explain
i = # Choose an index of the test instance that you want to explain
x = X_test[i]

# Explain the prediction for this instance
exp = explainer.explain_instance(x, model.predict_proba,
num_features=5, top_labels=1)

# Print the explanation
print(exp.as_list())
```

This code will create a LIME explainer object, which you can use to generate an explanation for a single instance in your test dataset. The explanation will consist of a list of the most important features that contributed to the model's prediction, along with the effect of each feature on the prediction.

13. Explain Model with Shap

I want you to act as a data scientist and explain the model's results. I have trained a scikit-learn XGBoost model and I would like to explain the output using a series of plots with Shap. Please write the code.

 I want you to act as a data scientist and explain the model's results. I have trained a scikit-learn XGBoost model and I would like to explain the output using a series of plots with Shap. Please write the code.

 Certainly! Here is some code that you can use to create a series of plots using Shap (SHapley Additive exPlanations) to explain the results of a scikit-learn XGBoost model:

Copy code

```
import shap
import matplotlib.pyplot as plt

# Load your model and the data that you used to train it
model = # Load your model
X_train = # Load your training data

# Create a SHAP explainer object
explainer = shap.TreeExplainer(model)

# Explain the model's predictions on the training data
shap_values = explainer.shap_values(X_train)

# Loop through the features and create a plot for each one
for i in range(X_train.shape[1]):
    # Create a subplot
    plt.subplot(X_train.shape[1], 1, i+1)

    # Create a SHAP summary plot
    shap.summary_plot(shap_values[:,i], X_train, plot_type="bar",
show=False)

    # Show the plot
    plt.show()
```

14. Write Multithreaded Functions

I want you to act as a coder. Can you help me parallelize this code across threads in python?



I want you to act as a coder. Can you help me parallelize this function in python?

```
import time
def print_time(delay):
    count = 0
    while count < 5:
        time.sleep(delay)
        count += 1
```



Sure! Here is an example of how you can parallelize the `print_time` function using the `threading` module in Python:

Copy code

```
import threading
import time

# Define a function for the thread
def print_time(threadName, delay):
    count = 0
    while count < 5:
        time.sleep(delay)
        count += 1

# Create two threads as follows
thread1 = threading.Thread(target=print_time, args=("Thread-1", 2,))
thread2 = threading.Thread(target=print_time, args=("Thread-2", 4,))

# Start the threads
thread1.start()
thread2.start()
```

This will create two threads, `thread1` and `thread2`, and start them simultaneously. The `print_time` function will be executed by both threads, with the `threadName` and `delay` arguments being passed to it from the `args` parameter in the `Thread` constructor.

15. Compare Function Speed

I want you to act as a software developer. I would like to compare the efficiency of two algorithms that performs the same thing in python. Please write code that helps me run an experiment that can be repeated for 5 times. Please output the runtime and other summary statistics of the experiment.

[Insert functions]



I want you to act as a software developer. I would like to compare the efficiency of two algorithms that performs the same thing in python. Please write code that helps me run an experiment that can be repeated for 5 times. Please output the runtime and other summary statistics of the experiment.

Function 1:

```
def compute1(a, b):
    return (a + b) / (a * b)
```

Function 2:

```
import numpy as np
def compute2(a, b):
    return np.divide(np.sum([a,b]),np.multiply(a,b))
```



As a software developer, I would suggest the following code to run an experiment comparing the efficiency of the two algorithms, which can be repeated for 5 times and outputs the runtime and summary statistics of the experiment:

Copy code

```
import time
import statistics

def measure_time(func, a, b):
    start = time.perf_counter()
    func(a, b)
    end = time.perf_counter()
    return end - start
```

CHATGPT PROMPTS FOR DATA SCIENCE

```
def run_experiment(func1, func2, a, b, num_repeats=5):
    times1 = []
    times2 = []

    for i in range(num_repeats):
        times1.append(measure_time(func1, a, b))
        times2.append(measure_time(func2, a, b))

    mean_time1 = sum(times1) / len(times1)
    mean_time2 = sum(times2) / len(times2)
    stdev1 = statistics.stdev(times1)
    stdev2 = statistics.stdev(times2)

    print(f"{func1.__name__}: mean time = {mean_time1:.6f}s, stdev = {stdev1:.6f}s")
    print(f"{func2.__name__}: mean time = {mean_time2:.6f}s, stdev = {stdev2:.6f}s")
    if mean_time1 < mean_time2:
        print(f"{func1.__name__} is faster than {func2.__name__}")
    elif mean_time1 > mean_time2:
        print(f"{func2.__name__} is faster than {func1.__name__}")
    else:
        print(f"{func1.__name__} and {func2.__name__} have the same performance")

# Example usage
run_experiment(compute1, compute2, 10, 20)
```

This code defines a `measure_time` function that measures the time taken to execute a given function with a set of arguments. It then defines a `run_experiment` function that uses the `measure_time` function to compare the efficiency of two functions.

16. Create NumPy Array

I want you to act as a data scientist. I need to create a numpy array. This numpy array should have the shape of (x,y,z). Please initialize the numpy array with random values.



I want you to act as a data scientist. I need to create a numpy array. This numpy array should have the shape of (10, 3, 4, 5). Please initialize the numpy array with random values between 1 to 5.



Sure! Here's how you can create a numpy array with the desired shape and initialized with random values between 1 and 5:

Copy code

```
import numpy as np

# Create a numpy array with the desired shape
arr = np.random.randint(1, 6, size=(10, 3, 4, 5))

# Print the array to verify its contents
print(arr)
```

This will create a numpy array with a shape of (10, 3, 4, 5) and fill it with random integers between 1 and 5 (inclusive).

I hope this helps! Let me know if you have any questions.



17. Write Unit Test

I want you to act as a software developer. Please write unit tests for the function [Insert function]. The test cases are:
[Insert test cases]

Credit: @svpino

 I want you to act as a software developer. Please write unit tests for the function def compute(a, b)
The test cases are:
Input: -2, 2 | Output: 0
Input: 2,2 | Output:1
Input: None, 2 | Output: TypeError
Input: 0,2 | Output: ZeroDivisionError

 Certainly! Here is some sample code in python that defines the function `compute` and writes unit tests for it using the `unittest` module:

Copy code

```
import unittest

def compute(a, b):
    return (a + b) / (a * b)

class TestCompute(unittest.TestCase):
    def test_compute_1(self):
        self.assertEqual(compute(-2, 2), 0)

    def test_compute_2(self):
        self.assertEqual(compute(2, 2), 1)

    def test_compute_3(self):
        with self.assertRaises(TypeError):
            compute(None, 2)

    def test_compute_4(self):
        with self.assertRaises(ZeroDivisionError):
            compute(0, 2)

if __name__ == '__main__':
    unittest.main()
```

18. Validate Column

I want you to act as a data scientist. Please write code to test if that my pandas Dataframe [insert requirements here]

 I want you to act as a data scientist. Please write code to test if that my pandas Dataframe "df" has any duplicate rows and has only 3 unique values in the "val" column.

 Certainly! Here is some example code that uses Pandas to test if the DataFrame "df" has any duplicate rows and has only 3 unique values in the "val" column:

 Copy code

```
import pandas as pd

# Test for duplicates
if df.duplicated().any():
    print("DataFrame has duplicate rows.")
else:
    print("DataFrame does not have duplicate rows.")

# Test for unique values
unique_vals = df["val"].unique()
if len(unique_vals) == 3:
    print("DataFrame has 3 unique values in the 'val' column.")
else:
    print("DataFrame does not have 3 unique values in the 'val' column.")
```

19. Explain Python

I want you to act as a code explainer. What is this code doing?

[Insert code]

Credit: @svpino

 I want you to act as a code explainer. What is this code doing? reconstructed =
np.clip(reconstructed, 0, 1)
plt.imshow(np.transpose(reconstructed, (1, 2, 0)))
plt.show()

 Certainly! This code appears to be displaying an image using the `imshow()` function from the Matplotlib library.

- `np.clip(reconstructed, 0, 1)`: This function clips the values in the `reconstructed` array so that all values less than 0 are set to 0, and all values greater than 1 are set to 1. This is likely being done to ensure that the values in the array are within the valid range for pixel intensity values (which are usually between 0 and 1).
- `np.transpose(reconstructed, (1, 2, 0))`: This function transposes the dimensions of the `reconstructed` array. The dimensions of an image array are usually `(height, width, channels)`, but the `imshow()` function expects the dimensions to be `(width, height, channels)`. The `transpose()` function reorders the dimensions of the array so that it has the correct shape for `imshow()`.
- `plt.imshow(...)` : This function displays the image stored in the transposed `reconstructed` array.
- `plt.show()`: This function displays the figure created by `imshow()`.

20. Explain SQL

I want you to act as a data science instructor. Can you please explain to me what this SQL code is doing? [Insert SQL code]



I want you to act as a data science instructor. Can you please explain to me what this SQL code is doing? `SELECT year, revenue, SUM(revenue) OVER (ROWS BETWEEN UNBOUNDED PRECEDING AND CURRENT ROW) FROM table`



Certainly! Here is an explanation of what this SQL code is doing:

- `SELECT year, revenue`: This clause specifies the columns that you want to include in the output. In this case, the output will include the "year" and "revenue" columns.
- `SUM(revenue) OVER (ROWS BETWEEN UNBOUNDED PRECEDING AND CURRENT ROW)`: This is a window function that calculates the sum of the "revenue" column for each row in the table. The `OVER()` clause defines the window over which the function will operate, in this case the `ROWS BETWEEN UNBOUNDED PRECEDING AND CURRENT ROW` window. This window includes all rows from the start of the table (`UNBOUNDED PRECEDING`) up to the current row (`CURRENT ROW`).
- `FROM table`: This clause specifies the table that you want to retrieve data from. In this case, the table is simply named "table".

Overall, this SQL code is selecting the "year" and "revenue" columns from the "table" table, and calculating the sum of the "revenue" column for each row in the table. The output will include a column with the sum of the "revenue" column for each row, as well as the original "year" and "revenue" columns.

21. Explain Google Sheets Formula

I want you to act as a Google Sheets formula explainer. Explain the following Google Sheets command. [Insert formula]



I want you to act as a Google Sheets formula explainer. Explain the following Google Sheets command. =filter(C2:C12, A2:A12=F2)



The `FILTER` function in Google Sheets is used to filter a range of cells based on specified criteria and return the resulting rows. The function has the following syntax:

Copy code

```
=FILTER(range, condition1, [condition2], ...)
```

Here's an explanation of each of the arguments:

- `range`: This is the range of cells that you want to filter. In this case, the range is `C2:C12`.
- `condition1`, `condition2`, etc.: These are the conditions that you want to use to filter the range. In this case, the condition is `A2:A12=F2`, which specifies that you want to include only the rows where the value in column A is equal to the value in cell F2.

So, in this case, the `FILTER` function is being used to filter the range `C2:C12` based on the condition `A2:A12=F2`, which means that only the rows where the value in column A is equal to the value in cell F2 will be included in the resulting dataset. The resulting dataset will include all of the values from column C that meet this criteria.

Note that the `FILTER` function is an array function, which means that you'll need to enter it using Ctrl + Shift + Enter (on Windows) or Command + Shift + Enter (on Mac). This will cause the formula to be enclosed in curly braces {} to indicate that it's an array formula.



22. Improve Code Speed

I want you to act as a software developer. Please help me improve the time complexity of the code below. [Insert code]



I want you to act as a software developer. Please help me improve the time complexity of the code below. sum=0.

```
for i in range(y+1):
    sum+=i
return sum
```



To improve the time complexity of the code you provided, you could use a mathematical formula to calculate the sum of the integers from 0 to `y`, rather than using a loop. This would eliminate the need to loop through the integers and add them up one by one, and would allow you to calculate the sum in a single step.

Here is how you could rewrite the code to use a mathematical formula to calculate the sum:

Copy code

```
def sum_range(y):
    return y * (y + 1) // 2
```

This function calculates the sum of the integers from 0 to `y` using the formula $\frac{y * (y + 1)}{2}$, which is the formula for the sum of an arithmetic series. This formula has a time complexity of O(1), because it performs a fixed number of operations regardless of the value of `y`.

23. Optimize Pandas

I want you to act as a code optimizer. Can you point out what's wrong with the following Pandas code and optimize it?
[Insert code here]



I want you to act as a code optimizer. Can you point out what's wrong with the following Pandas code and optimize it? df = pd.read_csv("titles.csv")

```
# Mutation - DON'T DO THIS
df_bad = df.query("runtime > 30 & type == 'SHOW'")
df_bad["score"] = df_bad[["imdb_score", "tmdb_score"]].sum(axis=1)
df_bad = df_bad[["seasons", "score"]]
df_bad = df_bad.groupby("seasons").agg(["count", "mean"])
df_bad = df_bad.droplevel(axis=1, level=0)
df_bad = df_bad.query("count > 10")
```



There's nothing inherently wrong with the code you've provided, but there are a few suggestions I can make to optimize it:

1. Consider using the `isin()` method to filter the dataframe instead of using the query method. This can be faster and more readable, especially if you have multiple conditions.

For example:

Copy code

```
df_filtered = df[df["runtime"].isin([30, 40, 50]) &
df["type"].isin(["SHOW", "MOVIE"])]
```

2. Instead of using `df_bad["score"] = df_bad[["imdb_score", "tmdb_score"]].sum(axis=1)`, you can create a new column using the `assign()` method, which can be more efficient and readable:

Copy code

```
df_bad = df_bad.assign(score=df_bad["imdb_score"] +
df_bad["tmdb_score"])
```

3. Instead of using `df_bad = df_bad[["seasons", "score"]]` to select specific columns, you can use the `filter()` method to filter the dataframe based on a list of column names:

Copy code

```
df_bad = df_bad.filter(items=["seasons", "score"])
```

24. Optimize Pandas Again

I want you to act as a code optimizer. Can you point out what's wrong with the following Pandas code and optimize it?
[Insert code here]



I want you to act as a code optimizer. Can you optimize this Pandas code for me?

```
vals = []
for ind, row in df.iterrows():
    country = row["prod_country1"]
    if country in top3:
        vals.append("top3")
    elif country in top10:
        vals.append("top10")
    elif country in top20:
        vals.append("top20")
    else:
        vals.append("other")
df["prod_country_rank"] = vals
```



Here are a few suggestions for optimizing the code you provided:

1. Instead of using the ``iterrows()`` method to loop through the rows of the dataframe, you can use the ``apply()`` method to apply a function to each row. This can be more efficient and easier to read. For example:

Copy code

```
def classify_country(row):
    country = row["prod_country1"]
    if country in top3:
        return "top3"
    elif country in top10:
        return "top10"
    elif country in top20:
        return "top20"
    else:
        return "other"

df["prod_country_rank"] = df.apply(classify_country, axis=1)
```

25. Optimize Python

I want you to act as a code optimizer. The code is poorly written. How do I correct it? [Insert code here]



I want you to act as a code optimizer. The code is poorly written. How do I correct it?

```
numbers = [1, 2, 3]
letters = ["A", "B", "C"]
```

```
for index in range(len(numbers)):
    print(numbers[index], letters[index])
```



The code you've provided is functional, but there are a few ways you could make it more concise and efficient:

1. Instead of using the `range` function and indexing into the lists, you can use the `zip` function to iterate over the two lists simultaneously:

Copy code

```
for number, letter in zip(numbers, letters):
    print(number, letter)
```

2. If you don't need the index values, you can use the `enumerate` function to iterate over the lists:

Copy code

```
for index, (number, letter) in enumerate(zip(numbers, letters)):
    print(index, number, letter)
```

3. You can also use list comprehensions to create a new list that combines the elements of `numbers` and `letters`:

Copy code

```
combined = [(number, letter) for number, letter in zip(numbers,
letters)]
```

Which of these approaches is the best choice will depend on your specific needs and the context in which the code is being used.



26. Optimize SQL

I want you to act as an SQL code optimizer. The following code is slow. Can you help me speed it up? [Insert SQL]



I want you to act as an SQL code optimizer. The following code is slow. Can you help me speed it up? SELECT
*
FROM TestTable
WHERE DATEPART(YEAR, SomeMyDate) = '2021';



Here are some suggestions to optimize the performance of the above SQL code:

1. Use an index on the `SomeMyDate` column: An index on the `SomeMyDate` column will allow the database to quickly find rows that match the condition in the `WHERE` clause, rather than having to scan the entire table.
2. Use a covered index: If the `TestTable` table has a large number of columns and you only need to retrieve a few of them, you can create a "covered index" that includes all the columns you need. This will allow the database to retrieve the data from the index itself, rather than having to look up the rows in the table.
3. Use the `EXTRACT()` function instead of `DATEPART()`: The `EXTRACT()` function may be slightly faster than `DATEPART()` for extracting the year from a date. Here is how you can rewrite the query using `EXTRACT()`:

Copy code

```
SELECT *  
FROM TestTable  
WHERE EXTRACT(YEAR FROM SomeMyDate) = 2021;
```

27. Simplify Python

I want you to act as a code simplifier. Can you simplify the following code?

 I want you to act as a code simplifier. Can you simplify the following code? `elems = [4, 7, 18]`

```
elem0 = elems[0]
elem1 = elems[1]
elem2 = elems[2]
```

 To simplify the code you've provided, you could use a tuple or a list to store the elements of `elems` together, like this:

Copy code

```
elems = [4, 7, 18]
elem0, elem1, elem2 = elems
```

Alternatively, you could eliminate the separate variables entirely and use `elems` directly in your code.

Which of these approaches is the best choice will depend on your specific needs and the context in which the code is being used.



28. Write Documentation

I want you to act as a software developer. Please provide documentation for func1 below. [Insert function]

Credit: @svpino

The image shows a dark-themed chat interface. At the top, a user icon and the text "I want you to act as a software developer. Please provide documentation for func1 below." are visible. Below this, the AI's response is shown in a code block:

```
def func1(y):
    sum=0.
    for i in range(y+1):
        sum+=i
    return sum
```

Below the code, the AI provides the documentation:

Here is the documentation for the `func1` function:

Copy code

```
def func1(y):
    """
    This function calculates the sum of the integers from 0 to y
    (inclusive).

    Parameters:
        y (int): The upper bound of the range of integers to be
        summed.

    Returns:
        float: The sum of the integers from 0 to y.
    """
    sum=0.
    for i in range(y+1):
        sum+=i
    return sum
```

29. Improve Readability

I want you to act as a code analyzer. Can you improve the following code for readability and maintainability? [Insert code]



I want you to act as a code analyzer. Can you improve the following code for readability and maintainability?

```
def list (a,b):
    randomCamelCase = a
    random_snake_case=b
    return [randomCamelCase, random_snake_case]
```



Certainly! Here is the improved code for readability and maintainability:

Copy code

```
def get_list(a, b):
    random_camel_case = a
    random_snake_case = b
    return [random_camel_case, random_snake_case]
```

30. Format SQL

I want you to act as a SQL formatter. Please format the following SQL code. Please convert all reserved keywords to uppercase [Insert requirements]. [Insert Code]

The screenshot shows a ChatGPT interface. The user (me) has sent a message asking for the SQL code to be formatted with uppercase reserved keywords. The AI (ChatGPT) has responded with the formatted SQL code. A 'Copy code' button is visible above the code block.

I want you to act as a SQL formatter for me. Can you please format the following SQL code for me. Please convert all reserved keywords to uppercase, keep everything else at lowercase, provide proper indentations, and keep indentations at 3 spaces per indentation. select c.year, t.* from revenue_table c cross join lateral (values (c.q1, 'Q1'), (c.q2, 'Q2'),(c.q3, 'Q3')) as t(sales, quarter) order by year, quarter;

Here is the formatted SQL code for you:

```
SELECT c.year, t.*  
FROM revenue_table c  
CROSS JOIN LATERAL (  
    VALUES (c.q1, 'Q1'),  
    (c.q2, 'Q2'),  
    (c.q3, 'Q3')  
) AS t(sales, quarter)  
ORDER BY year, quarter;
```

31. Translate Between DBMS

I want you to act as a coder and write SQL code for MySQL. What is the equivalent of PostgreSQL's DATE_TRUNC for MySQL?

A screenshot of a ChatGPT conversation. The user asks: "I want you to act as a coder and write SQL code for MySQL. What is the equivalent of DATE_TRUNC for MySQL?" The AI responds: "MySQL does not have a built-in function equivalent to PostgreSQL's `DATE_TRUNC` function, which allows you to truncate a date or timestamp to a specified precision. However, you can achieve similar results using the `DATE_FORMAT` function, which allows you to format a date or timestamp as a string using a variety of format specifiers." Below this, it says: "For example, to truncate a timestamp to the nearest hour, you can use the following query:" followed by a code block containing the SQL query.

```
SELECT DATE_FORMAT(timestamp_column, '%Y-%m-%d %H:00:00') AS truncated_timestamp
FROM table_name;
```

32. Translate Python to R

I want you to act as a code translator. Can you please convert the following code from python to R? [Insert code]

Credit: @svpino

 I want you to act as a code translator. Can you please convert the following code from python to R?

```
pd.concat([df_1, ..., df_n], axis=0)
```

 Certainly! In R, you can use the `bind_rows` function from the `dplyr` library to concatenate multiple data frames along the rows (axis=0). Here is the equivalent code in R:

```
library(dplyr)  
  
bind_rows(df_1, ..., df_n)
```

Copy code

This code first loads the `dplyr` library, and then uses the `bind_rows` function to concatenate the data frames `df_1`, `df_2`, ..., `df_n` along the rows. The resulting data frame will have all the rows from the input data frames, in the order they were specified.

33. Translate R to Python

I want you to act as a code translator. Can you please convert the following code from R to python? [Insert code]

Credit: @svpino



I want you to act as a code translator. Can you please convert the following code from R to python?
spread(
 df, key = 'key_name'
 value = 'value_name'
)



Certainly! In python, you can use the `pivot` method of a pandas DataFrame to spread a column into multiple columns. Here is the equivalent code in python using pandas:

Copy code

```
df.pivot(index='key_name', columns='value_name', values='value')
```

This code uses the `pivot` method to reshape the data frame by spreading the `key_name` column into multiple columns and using the `value_name` column as the values in the new columns. The resulting data frame will have one column for each unique value in the `value_name` column, and the values in the new columns will come from the `value` column.

34. Explain to Five-Year-Old

I want you to act as a data science instructor. Explain [concept] to a five-year-old.



I want you to act as a data science instructor. Explain what is neural network to a five-year-old.



Sure! A neural network is a type of computer program that is designed to recognize patterns. It does this by looking at lots of examples and trying to find things that they have in common.

Imagine you have a bunch of pictures of different animals, and you want to know which ones are dogs. You could show the pictures to a person, and they would be able to tell you which ones are dogs by looking for certain features, like four legs, a wagging tail, and a furry face. A neural network does something similar, but it does it all by itself, without any help from a person.

To do this, a neural network has a lot of little parts called "neurons," which are kind of like tiny brains. These neurons are arranged in layers, and they work together to analyze the patterns in the pictures. When they see a pattern that looks like a dog, they can say "Hey, this looks like a dog!"

Neural networks are really good at finding patterns, and they can be used for all sorts of things, like recognizing faces in photos, understanding speech, and even playing games like chess. They're one of the most powerful tools that data scientists have to solve all kinds of problems.



35. Explain to Undergraduate

I want you to act as a data science instructor. Explain [concept] to an undergraduate.



I want you to act as a data science instructor. Explain what is neural network to an undergraduate student.



A neural network is a type of machine learning model that is inspired by the structure and function of the brain. It consists of layers of interconnected "neurons," which process and transmit information.

Each neuron receives input from other neurons or from external sources, and combines these inputs using weights that reflect the importance of each input. The neuron then applies an activation function to the combined input to produce an output, which is passed on to other neurons in the next layer.

The weights and activation functions of the neurons are learned from training data, allowing the neural network to learn complex patterns and relationships.

Neural networks can be used for a wide variety of tasks, such as image classification, language translation, and time series prediction. They are particularly well-suited to tasks that involve large and complex datasets, and are able to learn and adapt as they are exposed to more data.



36. Explain to Professor

I want you to act as a data science instructor. Explain [concept] to a professor.



I want you to act as a data science instructor. Explain what is neural network to a professor.



Sure! A neural network is a type of machine learning algorithm that is inspired by the structure and function of the brain. It is composed of a large number of interconnected processing nodes, called "neurons," which are organized into layers. Each neuron receives input from other neurons, and uses this input to make a decision or prediction based on the patterns it sees in the data.

Neural networks are typically used to solve problems that involve pattern recognition or classification. For example, a neural network might be trained to recognize handwritten digits based on a large dataset of labeled examples, or to classify emails as spam or non-spam based on their content.

One of the key advantages of neural networks is that they are able to learn and adapt to new data over time, without the need for explicit programming. This makes them particularly useful for tasks that involve complex, non-linear relationships, or for situations where the desired output is not easily defined in advance.

Neural networks are widely used in a variety of fields, including image and speech recognition, natural language processing, and finance. They have also been used to achieve state-of-the-art results on a number of benchmarks and real-world tasks.



37. Explain to Business Stakeholder

I want you to act as a data science instructor. Explain [concept] to a business stakeholder.



I want you to act as a data science instructor. Explain what is neural network to a non-technical business stakeholder in a company.



Certainly! A neural network is a type of computer program that is designed to recognize patterns in data. It's like a "virtual brain" that can learn from examples and make decisions or predictions based on what it has learned.

For example, let's say you have a company that sells products online, and you want to figure out which of your customers are most likely to buy a new product that you just released. You could use a neural network to analyze data about your customers' past purchases, their demographic information, and other factors that might be relevant. Based on this analysis, the neural network could identify patterns that seem to be associated with a higher likelihood of making a purchase, and it could use these patterns to predict which customers are most likely to buy the new product.

This is just one example of how neural networks can be used to solve problems in business. They can also be used for things like predicting which employees are most likely to leave the company, identifying fraudulent activity, and optimizing supply chain operations, among many other applications.

Overall, neural networks are a powerful tool that can help businesses make better decisions and improve their operations by finding patterns and insights in data that would be difficult or impossible for a human to identify.



38. Explain Like Stackoverflow

I want you to act as an answerer on StackOverflow. You can provide code snippets, sample tables and outputs to support your answer. [Insert technical question]



I want you to act as an answerer on StackOverflow. You can provide code snippets, sample tables and outputs to support your answer. What is the use of constraints in SQL?



In SQL, constraints are used to specify rules for the data in a table. If a constraint is violated, the database will return an error and the transaction will be rolled back. Constraints are used to ensure the integrity and correctness of the data in a database.

There are several types of constraints that can be defined in SQL:

- NOT NULL: This constraint specifies that a column cannot contain a NULL value.
- UNIQUE: This constraint specifies that all the values in a column must be unique.
- PRIMARY KEY: This constraint specifies a column or a set of columns that uniquely identifies each row in a table. A table can have only one primary key.
- FOREIGN KEY: This constraint specifies that the values in a column must match the values in a column of a different table (the parent table). This is used to enforce referential integrity.
- CHECK: This constraint specifies a condition that the values in a column must satisfy.
- DEFAULT: This constraint specifies a default value for a column. If no value is specified for the column when inserting a row, the default value will be used.

Here is an example of how constraints can be defined in SQL:

Copy code

```
CREATE TABLE Customers (
    CustomerId INT PRIMARY KEY,
    Name VARCHAR(255) NOT NULL,
    Email VARCHAR(255) UNIQUE,
    Age INT CHECK (Age > 0)
);
```

In this example, the `Customers` table has four columns: `CustomerId`, `Name`, `Email`, and `Age`. The `CustomerId` column is defined as the primary key of the table, the `Name` column is defined as NOT NULL, the `Email` column is defined as UNIQUE, and the `Age` column is defined with a CHECK constraint that ensures that the values in the column are greater than zero.

39. Suggest Edge Cases

I want you to act as a software developer. Please help me catch edge cases for this function [insert function]



I want you to act as a software developer. Please help me catch edge cases for this function

```
def compute(a, b):
    return (a+b)/(a*b)
```



Certainly! Here are some edge cases you might want to consider when testing the `compute` function:

- When `a` and `b` are both zero. In this case, the function will raise a `ZeroDivisionError`, since you cannot divide by zero.
- When `a` is zero and `b` is non-zero. In this case, the function will also raise a `ZeroDivisionError`.
- When `a` and `b` are both negative. This might produce a different result than when `a` and `b` are both positive, depending on the values.
- When `a` is negative and `b` is positive, or vice versa. This might also produce a different result than when `a` and `b` are both positive or both negative.
- When `a` is very large and `b` is very small, or vice versa. This might produce a different result than when `a` and `b` are both large or both small.

These are just a few examples of edge cases you might want to consider.

Depending on the specific requirements of your function, there might be other edge cases to consider as well.

40. Suggest Dataset

I want you to act as a data science career coach. I want to build a predictive model for [...]. At the same time, I would like to showcase my knowledge in [...]. Can you please suggest the five most relevant datasets for my use case?



I want you to act as a data science career coach. I want to build a predictive model for movies. At the same time, I would like to showcase my knowledge in computer vision. Can you please suggest the five most relevant datasets for my use case?



Sure! Here are five potential datasets that you might find useful for building a predictive model for movies:

1. The MovieLens dataset: This dataset contains movie ratings and tag recommendations from the MovieLens website, and includes information about the movies, users, and ratings.
2. The IMDb dataset: This dataset contains information about movies and TV shows, including cast, crew, and plot summaries.
3. The Netflix Prize dataset: This dataset contains movie ratings and other data from the Netflix Prize competition, which challenged participants to build a recommendation system for movies.
4. The TMDB 5000 Movie dataset: This dataset contains information about movies, including their budget, revenue, and cast, as well as movie posters and other images.
5. The Film Popularity dataset from the UCI Machine Learning Repository: This dataset contains information about the popularity of movies, including their gross earnings and IMDb ratings.

41. Suggest Portfolio Ideas

I want you to act as a data science coach. My background is in [...] and I would like to [career goal]. I need to build a portfolio of data science projects that will help me land a role in [...] as a [...]. Can you suggest five specific portfolio projects that will showcase my expertise in [...] and are of relevance to [company]?



I want you to act as a data science coach. I have a degree in civil engineering and I would like to make a switch into data science. I know I need to build a portfolio of data science projects that will help me land a role in Tesla as a data scientist. Can you suggest five specific portfolio projects that will showcase my expertise in civil engineering and are of relevance to Tesla?



Sure! Here are five potential portfolio projects that could be relevant for a data science role at Tesla and showcase your expertise in civil engineering:

1. Infrastructure optimization: Tesla is building a network of charging stations and other infrastructure to support its electric vehicles. You could build a project that helps optimize the design and placement of these stations based on factors such as traffic patterns, population density, and the availability of renewable energy sources.
2. Predictive maintenance: Tesla's vehicles and charging stations are complex systems that require regular maintenance to operate efficiently. You could build a machine learning model that predicts when maintenance will be needed, so that it can be scheduled in advance and downtime can be minimized.
3. Traffic flow prediction: Tesla is working on developing self-driving vehicles, which will need to be able to navigate complex traffic environments. You could build a model that predicts traffic flow in different locations, which could be used to help autonomous vehicles navigate more efficiently.

42. Suggest Resources

I want you to act as a data science coach. I would like to learn about [topic]. Please suggest 3 best specific resources. You can include [specify resource type]



I want you to act as a data science coach. I would like to learn about data engineering.
Please suggest 3 best specific resources. You can include websites, courses and books.



Sure! Here are three resources that I would recommend for learning about data engineering:

1. **Online Courses:** There are many online courses available that cover data engineering concepts and best practices. Some options include "Data Engineering on Google Cloud Platform" on Coursera and "Data Engineering Fundamentals" on Pluralsight.
2. **Books:** There are also several books that are great for learning about data engineering. Some options include "Designing Data-Intensive Applications" by Martin Kleppmann and "Building Data Pipelines with Python" by Alec Watson.
3. **Websites:** There are many websites that offer tutorials, articles, and other resources on data engineering. Some options include the "Data Engineering" section of the Dataquest blog and the "Data Engineering" section of the Databricks blog.

43. Suggest Time Complexity

I want you to act as a software developer. Please compare the time complexity of the two algorithms below. [Insert two functions]



I want you to act as a software developer. Please compare the time complexity of the two algorithms below.

```
def func1(y):
    sum=0.
    for i in range(y+1):
        sum+=i
    return sum
def func2(y):
    return y/2*(1+y)
```



As a software developer, I would say that the time complexity of the `func1` function is $O(n)$, where n is the value of the input `y`. This is because the function performs a single loop that runs from 0 to `y`, and the time taken to execute the function increases linearly with the value of `y`.

The time complexity of the `func2` function is $O(1)$, because the function performs a fixed number of operations regardless of the value of the input `y`. The time taken to execute the function is constant, and does not depend on the size of the input.

Therefore, the `func2` function is more efficient than the `func1` function, because it has a lower time complexity. The `func2` function will always be faster than the `func1` function, regardless of the size of the input.



44. Suggest Feature Engineering

I want you to act as a data scientist and perform feature engineering. I am working on a model that predicts [insert feature name]. There are columns: [Describe columns]. Can you suggest features that we can engineer for this machine learning problem?



I want you to act as a data scientist and perform feature engineering. I am working on a model that predicts whether person will default on a credit card loan in October. There are columns: LIMIT_BAL (Credit card limit), SEX, EDUCATION, MARRIAGE, AGE, BILL_AMT[1-6] (amount of bill statement for April to September), PAY_AMT[1-6] payment amount for April to September. Can you suggest features that we can engineer for this machine learning problem?



Certainly! Here are a few ideas for features that you could engineer for this machine learning problem:

1. Payment history: You could create a new feature that captures the overall payment history of the borrower. For example, you could calculate the average payment amount over the past six months, or count the number of months in which the borrower paid the minimum amount due.
2. Credit utilization: You could calculate the credit utilization ratio for each borrower, which is the ratio of the total amount of credit being used to the total credit limit. This can be a useful feature because high credit utilization is often associated with higher risk of default.
3. Payment to bill ratio: You could create a new feature that captures the ratio of payment amount to bill amount for each month. This can give you an idea of how much of the bill the borrower is paying off each month, which can be a useful predictor of default risk.

45. Suggest A/B Testing Steps

I want you to act as a statistician. [Describe context] Please design an A/B test for this purpose. Please include the concrete steps on which statistical test I should run.



I want you to act as a statistician. My web designer has changed the landing page's background from blue to red. Please design an A/B test for this purpose. Please include the concrete steps on which statistical test I should run.



Certainly! Here are the steps for designing an A/B test to compare the effectiveness of a blue landing page background versus a red landing page background, including which statistical test to use:

1. Sampling:

- Choose a representative sample of visitors to your website. This could be a random sample, or you could stratify the sample based on factors such as age, gender, or geographic location to ensure that the sample is representative of your target audience.
- Divide the sample into two groups: the "treatment" group, which will see the red background, and the "control" group, which will see the blue background. It is important to ensure that the two groups are comparable in terms of the characteristics of the visitors.
- Set the sample size to be large enough to detect a statistically significant difference between the two groups, if one exists.

2. Running the experiment:

- Implement the changes to the landing page for the treatment group. Make sure to keep all other elements of the page the same for both groups, so that any differences in the outcomes can be attributed to the background color.
- Track the key metric that you are interested in, such as the conversion rate (the percentage of visitors who take a desired action, such as filling out a form or making a purchase). You will need to have a way to track this metric for both the treatment and control groups.
- Run the experiment for a sufficient length of time to allow for a sufficient number of visitors to be included in the sample. The length of time will depend on the size of the sample and the expected effect size.

3. Analyzing the results:

- Calculate the difference in the key metric between the treatment and control groups. This will be the "effect size" of the change in background color.
- Use a two-sample t-test to determine whether the observed difference is statistically significant. This test is appropriate because you are comparing the means of two independent samples (the treatment and control groups).
- If the p-value of the t-test is less than your predetermined significance level (usually 0.05), you can conclude that the difference is statistically significant, and that the change in background color had an effect on the key metric.

46. Career Coaching

I want you to act as a career advisor. I am looking for a role as a [role name]. My background is [...]. How do I land the role and with what resources exactly in 6 months?



I want you to act as a career advisor. I am looking for a role as a data engineer in Amazon as a fresh graduate. I have an undergraduate degree in chemical engineering and have self-studied data science through MOOCs. How do I land the role and with what resources exactly in 6 months?



Sure, I'd be happy to help you out. As a data engineer at Amazon, you'll be responsible for designing, building, maintaining, and troubleshooting data pipelines. It's a role that combines software engineering with data management and analysis skills.

To land a data engineering role at Amazon, you'll need to demonstrate your technical skills, problem-solving abilities, and attention to detail. Here are a few things you can do to improve your chances:

1. Continue to build your technical skills: As a chemical engineering graduate, you have a strong foundation in math and science, which will be useful in data engineering. However, you'll also need to be proficient in programming languages such as Python and SQL, as well as tools like Hadoop, Spark, and AWS. Consider taking online courses or earning certifications to enhance your skills.
2. Gain practical experience: While a degree and MOOCs are a good start, practical experience is crucial for landing a data engineering role. Consider working on data-related projects in your spare time, or interning at a company that uses data engineering. This will give you the opportunity to apply your skills and build a portfolio of work to show potential employers.
3. Network: Networking can be a powerful tool for finding job opportunities. Attend industry events and meetups, join online forums and groups, and reach out to professionals in the field to learn more about data engineering and make connections. You never know who might be able to help you find your next job.
4. Tailor your resume and cover letter: When applying for a data engineering role at Amazon, make sure to highlight your relevant skills and experience. Include any projects you've worked on, any relevant coursework or certifications, and any relevant work experience. Customize your resume and cover letter for each job you apply to, and make sure to highlight how your skills and experience make you a strong fit for the role.

47. Correct Own Chatgpt Code

Your above code is wrong. [Point out what is wrong]. Can you try again?

A screenshot of a ChatGPT conversation. The user asks for a command to copy files recursively from "Desktop/files1" to "Desktop/files2". ChatGPT responds with a command that includes "-rT" instead of the correct "-r" and "-T". It then corrects itself, suggesting the use of "mkdir -p" and "cp -r".

I want you to act as a Linux terminal. Can you provide me the command to recursively copy all files from "Desktop/files1" to a folder that has not been created "Desktop/files2"?

To recursively copy all files from "Desktop/files1" to a folder that has not been created "Desktop/files2" in the Linux terminal, you can use the `cp` command with the `-r` and `-T` options, like this:

`cp -rT Desktop/files1 Desktop/files2`

Copy code

Your above terminal command is wrong. There is no "-T" option for cp. Can you try again?

I apologize for the error in my previous answer. The `-T` option is not a valid option for the `cp` command. To recursively copy all files from "Desktop/files1" to a folder that has not been created "Desktop/files2" in the Linux terminal, you can use the `cp` command with the `-r` option and the `mkdir -p` command, like this:

`mkdir -p Desktop/files2
cp -r Desktop/files1/* Desktop/files2/`

Copy code

48. Correct Python Code

I want you to act as a software developer. This code is supposed to [expected function]. Please help me debug this python code that cannot be run. [Insert function]



I want you to act as a software developer. This code is supposed to add numbers from 1 to k. Please help me debug this python code that cannot be run.

```
def add_from_1_to_k(k):
    if(k>0):
        result = k+add_from_1_to_k(k-1)
    return result
```



There is an error in the code you provided because it is missing a return statement for the case where `k` is 0 or negative. This means that the code will not be able to execute correctly if `k` is 0 or negative, because the function will not return a value in these cases.

To fix this error, you can add a return statement for the case where `k` is 0 or negative, like this:

Copy code

```
def add_from_1_to_k(k):
    if k <= 0:
        return 0
    else:
        result = k + add_from_1_to_k(k-1)
    return result
```

49. Correct SQL Code

I want you to act as a SQL code corrector. This code does not run in [your DBMS, e.g. PostgreSQL]. Can you correct it for me? [SQL code here]



I want you to act as a SQL code corrector. This code does not run. Can you correct it for me?
SELECT name
FROM dish
WHERE name = 'Prawn Salad'
GROUP BY name
HAVING count(*) = 1
LIMIT 5 ORDER BY name ;



It looks like there is a problem with the syntax of your `SELECT` statement. The `LIMIT` clause should come before the `ORDER BY` clause. Try the following instead:

Copy code

```
SELECT name
FROM dish
WHERE name = 'Prawn Salad'
GROUP BY name
HAVING count(*) = 1
ORDER BY name
LIMIT 5;
```

This should correct the syntax error and allow the query to run successfully.



50. Troubleshoot PowerBI Model

I want you to act as a PowerBI modeler. Here is the details of my current project. [Insert details]. Do you see any problems with the table?

Credit: Mathias Halkjær Petersen



I want you to act as a PowerBI modeler. Here is the details of my current project.

Table: Emission

- Column: Timestamp (datetime)
- Column: Value (string)
- Column: Unit (string)
- Column: EmissionType (string)
- Column: Entity (string)

Do you see any problems with the table?



It is difficult to determine if there are any problems with the Emission table without more context. However, here are a few things to consider when working with this table in Power BI:

1. Data types: Make sure that the data types for each column are appropriate for the data they contain. For example, the "Timestamp" column should have a data type of "Date/Time", and the "Value" column should have a numeric data type such as "Decimal Number" or "Whole Number".
2. Data quality: Check for any errors or inconsistencies in the data. For example, you might want to ensure that all values in the "Unit" column are spelled correctly, or that the "EmissionType" column does not contain any null values.
4. Performance: Consider the size of the Emission table and the complexity of your measures and visuals. If the table is very large or your measures and visuals are very complex, it could impact the performance of your Power BI report.

Overall, the key is to understand the data in the Emission table and how you want to use it in your Power BI report. This will help you identify any potential problems and come up with solutions to address them.

51. Create Running Average

I want you to act as a data scientist and write SQL code for me. I have a table with two columns [Insert column names]. I would like to calculate a running average for [which value]. What is the SQL code that works for PostgreSQL 14?



I want you to act as a data scientist and write SQL code for me. I have a table with two columns "date", "product" and "revenue". I would like to calculate a running average for each product. What is the SQL code that works for PostgreSQL 14?



Here is the SQL code that you can use to calculate the running average for each product in PostgreSQL 14:

Copy code

```
SELECT
    date,
    product,
    revenue,
    AVG(revenue) OVER (
        PARTITION BY product
        ORDER BY date
        ROWS BETWEEN UNBOUNDED PRECEDING AND CURRENT ROW
    ) as running_avg
FROM sales_table
ORDER BY date;
```

This code will do the following:

1. Select the `date`, `product`, `revenue`, and `running_avg` columns from the `sales_table` table.
2. Calculate the running average of the `revenue` column for each product using a window function.
3. The `AVG()` function is a window function that calculates the average of a set of rows. The `OVER()` clause specifies the window for the window function. The `PARTITION BY` clause divides the rows of the table into partitions to which the window function is applied. In this case, the rows are partitioned by `product`. The `ORDER BY` clause specifies the order in which the rows are processed within each partition. In this case, the rows are ordered by `date`. The `ROWS` keyword specifies that the window is defined in terms of rows, and the `BETWEEN` keyword specifies the start and end of the window.
4. The `running_avg` column will contain the average of all `revenue` values from the start of the partition up to the current row, for each product.
5. The resulting table will be ordered by `date`.



52. Solve Leetcode Question

Assume you are given the tables... with the columns... Output the following...

Credit: Data Lemur



Assume that you are given the table below containing information on viewership by device type (where the three types are laptop, tablet, and phone). Define "mobile" as the sum of tablet and phone viewership numbers. Write a query to compare the viewership on laptops versus mobile devices.

Output the total viewership for laptop and mobile devices in the format of "laptop_views" and "mobile_views". The viewership table has the columns "user_id", "device_type", and "view_time"



To compare the viewership on laptops versus mobile devices, you can use the following SQL query:

Copy code

```
SELECT
    SUM(CASE WHEN device_type = 'laptop' THEN view_time ELSE 0 END)
AS laptop_views,
    SUM(CASE WHEN device_type IN ('tablet', 'phone') THEN view_time
ELSE 0 END) AS mobile_views
FROM viewership
```

This query uses a `CASE` expression to sum the "view_time" values for laptop and mobile devices separately. The `SUM` function is then used to calculate the total viewership for each device type.

53. Write Google Sheets Formula

I want you to act as a bot that generates Google Sheets formula. Please generate a formula that [describe requirements]



I want you to act as a bot that generates Google Sheets formula. Please generate a formula that puts data from another sheet into current sheet.



Here is a formula that puts data from another sheet into the current sheet:

Copy code

```
=IMPORTRANGE("sheet_id", "range")
```

This formula uses the `IMPORTRANGE` function, which imports data from another sheet. The first argument is the "sheet_id" of the sheet that you want to import data from. The sheet_id is a unique identifier for the sheet, and you can find it in the URL of the sheet when it's open in your browser. For example, if the URL of the sheet is `https://docs.google.com/spreadsheets/d/abc123/edit`, the sheet_id is `abc123`.

The second argument is the range of cells that you want to import. This can be a single cell (e.g., "A1"), a range of cells (e.g., "A1:B10"), or a named range.

For example, if you want to import the data from cells A1 through B10 in a sheet with the sheet_id `abc123`, you would use the following formula:

Copy code

```
=IMPORTRANGE("abc123", "A1:B10")
```

Note that to use the `IMPORTRANGE` function, you'll need to grant access to the sheet that you want to import data from. To do this, you can go to the "Data" menu and select "Import data." Then, choose the sheet that you want to import data from and click "Allow access."

Once you've granted access, you can use the `IMPORTRANGE` function to import data from the other sheet into the current sheet. The imported data will be automatically updated whenever the source data is changed.



54. Write R

I want you to act as a data scientist using R. Can you write an R script that [Insert requirement here]



I want you to act as a data scientist using R. Can you write an R script that imports the iris dataset and uses ggplot2 to perform data visualization. I would like to use a violinplot to compare the distribution of sepal length between flower species.



Certainly! Here is an R script that imports the iris dataset and uses ggplot2 to create a violinplot comparing the distribution of sepal length between flower species:

Copy code

```
# Load the ggplot2 library
library(ggplot2)

# Load the iris dataset
data("iris")

# Create a violinplot of sepal length by species
ggplot(data = iris, aes(x = Species, y = Sepal.Length)) +
  geom_violin()
```

This script loads the ggplot2 library and the iris dataset, and then creates a violinplot using the `geom_violin` function. The plot shows the distribution of sepal length for each flower species, with the width of the violins indicating the density of the data at each point.

To use this script, copy and paste it into an R script file or an R console, and then run it by pressing the "Run" button or by using the "Ctrl + Enter" keys. The code will create the violinplot and display it in the R graphics window. You can also save the plot to a file by using the `ggsave` function, like this:

Copy code

```
ggsave("violinplot.png", plot = p)
```

55. Write Shell

I want you to act as a Linux terminal expert. Please write the code to [describe requirements]



I want you to act as a Linux terminal expert. Please write the code to create a directory called "new folder", and create a shell file "txt.sh". It should contain the command to list all files in "new folder"



To create a directory called "new folder" and a shell file "txt.sh" that contains the command to list all files in "new folder", you can use the following commands in the Linux terminal:

Copy code

```
# Create the directory
mkdir new_folder

# Change into the new directory
cd new_folder

# Create the shell file
echo "ls" > txt.sh

# Make the shell file executable
chmod +x txt.sh
```

56. Write VBA

I want you to act as an Excel VBA developer. Can you write a VBA that [Insert function here]?

 I want you to act as an Excel VBA developer. Can you write a VBA that automatically print all sheet names?

 Yes, I can write a VBA code that will automatically print all sheet names in an Excel workbook. Here is a sample code that you can use:

Sub PrintSheetNames()

```
Dim ws As Worksheet
For Each ws In ThisWorkbook.Worksheets
    Debug.Print ws.Name
Next ws
End Sub
```

 Copy code

To use this code, copy and paste it into the VBA editor in Excel, and then run it by pressing the "Run" button or by using the F5 key. The code will loop through all the worksheets in the current workbook and print the name of each sheet to the debug window.

Note: You may need to enable the debug window in Excel to see the output. To do this, go to the "View" tab in the ribbon, click the "Immediate Window" button, and then press the "Alt + F11" keys to open the VBA editor. Then, go to the "View" menu, click the "Debug Windows" submenu, and click the "Immediate Window" option. The debug window will appear at the bottom of the screen, and you will be able to see the output of the code as it is executed.

57. Format Tables

I want you to act as a document formatter. Please format the following into a nice table for me to place in Google Docs?
[insert text table here]

 I want you to act as a document formatter for me. Can you please format the following into a nice table for me to place in Google Docs?

1 January 2021 Friday New Year's Day New Year's Day
12 February 2021 Friday Chinese New Year Chinese New Year
13 February 2021 Saturday Chinese New Year Chinese New Year
2 April 2021 Friday Good Friday Good Friday
1 May 2021 Saturday Labour Day Labour Day

 Here is the formatted table for you:

Date	Day	Holiday	Description
1 January	Friday	New Year's Day	New Year's Day
12 February	Friday	Chinese New Year	Chinese New Year
13 February	Saturday	Chinese New Year	Chinese New Year
2 April	Friday	Good Friday	Good Friday
1 May	Saturday	Labour Day	Labour Day



58. Summarize Book

I want you to act as a technical book summarizer. Can you please summarize the book [name] with 5 main points?



I want you to act as a technical book summarizer. Can you please summarize the book Storytelling With Data: A Data Visualization Guide for Business Professionals with 5 main points?



Certainly! Here are five main points from "Storytelling with Data: A Data Visualization Guide for Business Professionals" by Cole Nussbaumer Knaflic:

1. The purpose of data visualization is to communicate information clearly and effectively through graphical means.
2. It is important to consider the context in which the data will be presented and the audience that will be viewing it.
3. Good data visualization should include appropriate and effective use of charts, graphs, and other visual encodings.
4. The design and aesthetics of the visualization are important, but function should be the primary focus.
5. It is important to clearly label and annotate the visualization so that the information being conveyed is easily understood.



59. Summarize Paper

I want you to act as an academic. Please summarise the paper [...] in simple terms in one paragraph.



I want you to act as an academic. Please summarise the paper "Distilling the knowledge in a neural network" in simple terms in one paragraph.



"Distilling the knowledge in a neural network" is a paper that presents a method for compressing a large, complex neural network model into a smaller, more efficient model that can be deployed in real-world applications. The authors propose a method called "knowledge distillation" in which the large model is used to train the smaller model, effectively transferring the knowledge from the large model to the smaller one. The smaller model is able to achieve similar performance to the large model, but with fewer parameters and faster inference times. This method has been shown to be effective in a variety of tasks, including image classification and machine translation.



60. Provide Emotional Support

I want you to provide emotional support to me. [Explain problem here.]



I want you to provide emotional support to me. I am a data scientist who has been training my deep learning model for 3 months, and my model is still underperforming a linear regression model.



I'm here to help in any way I can. It can be frustrating when we put a lot of time and effort into a project and don't see the results we were hoping for, but it's important to remember that it's a normal part of the process and that setbacks can often lead to valuable learning experiences.

Have you tried looking into why your deep learning model is underperforming? It could be that the data you are using is not well-suited for a deep learning model, or that there are issues with the way you have designed and trained the model. It can be helpful to try to identify what might be causing the issue and then see if there are any changes you can make to address it. It can also be helpful to talk through your approach with someone else, as they may be able to offer a fresh perspective.

In the meantime, it's important to take care of yourself and try to find ways to manage any stress or frustration you may be feeling. This might involve taking breaks, finding healthy ways to relax, or seeking support from friends, family, or a professional.



Follow



Travis Tang

[www.linkedin.com
/in/travistang](https://www.linkedin.com/in/travistang)

for more data
science content