# Activity_Course 4 TikTok project lab

February 1, 2025

## 1 TikTok Project

**Course 4 - The Power of Statistics**

You are a data professional at TikTok. The current project is reaching its midpoint; a project proposal, Python coding work, and exploratory data analysis have all been completed.

The team has reviewed the results of the exploratory data analysis and the previous executive summary the team prepared. You received an email from Orion Rainier, Data Scientist at TikTok, with your next assignment: determine and conduct the necessary hypothesis tests and statistical analysis for the TikTok classification project.

A notebook was structured and prepared to help you in this project. Please complete the following questions.

## 2 Course 4 End-of-course project: Data exploration and hypothesis testing

In this activity, you will explore the data provided and conduct hypothesis testing.

**The purpose** of this project is to demostrate knowledge of how to prepare, create, and analyze hypothesis tests.

**The goal** is to apply descriptive and inferential statistics, probability distributions, and hypothesis testing in Python.

*This activity has three parts:*

**Part 1:** Imports and data loading

- What data packages will be necessary for hypothesis testing?

**Part 2:** Conduct hypothesis testing

- How will descriptive statistics help you analyze your data?

- How will you formulate your null hypothesis and alternative hypothesis?

**Part 3:** Communicate insights with stakeholders

- What key business insight(s) emerge from your hypothesis test?

- What business recommendations do you propose based on your results?

Be sure to complete this activity before moving on. The next course item will provide you with a completed exemplar to compare to your own work.

# 3  Data exploration and hypothesis testing

# 4  PACE stages

Throughout these project notebooks, you'll see references to the problem-solving framework PACE. The following notebook components are labeled with the respective PACE stage: Plan, Analyze, Construct, and Execute.

## 4.1  PACE: Plan

1. What is your research question for this data project? Later on, you will need to formulate the null and alternative hypotheses as the first step of your hypothesis test. Consider your research question now, at the start of this task.

We are researching whether the video_view_count variable is correlated with the verified_status variable.

*Complete the following steps to perform statistical analysis of your data:*

### 4.1.1  Task 1. Imports and Data Loading

Import packages and libraries needed to compute descriptive statistics and conduct a hypothesis test.

Hint:

Be sure to import `pandas`, `numpy`, `matplotlib.pyplot`, `seaborn`, and `scipy`.

```
[1]:  # Import packages for data manipulation
      import numpy as np
      import pandas as pd

      # Import packages for data visualization
      import matplotlib.pyplot as plt
      import seaborn as sns

      # Import packages for statistical analysis/hypothesis testing
      from scipy import stats
```

Load the dataset.

**Note:** As shown in this cell, the dataset has been automatically loaded in for you. You do not need to download the .csv file, or provide more code, in order to access the dataset and proceed with this lab. Please continue with this activity by completing the following instructions.

```
[2]:  # Load dataset into dataframe
      data = pd.read_csv("tiktok_dataset.csv")
```

## 4.2   PACE: Analyze and Construct

1. Data professionals use descriptive statistics for Exploratory Data Analysis. How can comput-
   ing descriptive statistics help you learn more about your data in this stage of your analysis?

Descriptive statistics allow us to quickly explore and understand large amounts of data. With a
couple of statistics, mean and variance, we can learn a lot about the distributions of variables we
have sampled. This information can help us gauge the types of relationships the different variables
may have.

### 4.2.1   Task 2. Data exploration

Use descriptive statistics to conduct Exploratory Data Analysis (EDA).

Hint:

Refer back to *Self Review Descriptive Statistics* for this step-by-step proccess.

Inspect the first five rows of the dataframe.

```
[3]:  # Display first few rows
      data.head(10)
```

```
[3]:      # claim_status      video_id  video_duration_sec  \
      0   1          claim  7017666017                  59
      1   2          claim  4014381136                  32
      2   3          claim  9859838091                  31
      3   4          claim  1866847991                  25
      4   5          claim  7105231098                  19
      5   6          claim  8972200955                  35
      6   7          claim  4958886992                  16
      7   8          claim  2270982263                  41
      8   9          claim  5235769692                  50
      9  10          claim  4660861094                  45


                              video_transcription_text verified_status  \
      0   someone shared with me that drone deliveries a…    not verified
      1   someone shared with me that there are more mic…    not verified
      2   someone shared with me that american industria…    not verified
      3   someone shared with me that the metro of st. p…    not verified
      4   someone shared with me that the number of busi…    not verified
      5   someone shared with me that gross domestic pro…    not verified
      6   someone shared with me that elvis presley has …    not verified
      7   someone shared with me that the best selling s…    not verified
      8   someone shared with me that about half of the …    not verified
```

```
9    someone shared with me that it would take a 50…            verified
```

```
    author_ban_status  video_view_count  video_like_count  video_share_count  \
0       under review          343296.0           19425.0              241.0
1             active          140877.0           77355.0            19034.0
2             active          902185.0           97690.0             2858.0
3             active          437506.0          239954.0            34812.0
4             active           56167.0           34987.0             4110.0
5       under review          336647.0          175546.0            62303.0
6             active          750345.0          486192.0           193911.0
7             active          547532.0            1072.0               50.0
8             active           24819.0           10160.0             1050.0
9             active          931587.0          171051.0            67739.0
```

```
    video_download_count  video_comment_count
0                    1.0                  0.0
1                 1161.0                684.0
2                  833.0                329.0
3                 1234.0                584.0
4                  547.0                152.0
5                 4293.0               1857.0
6                 8616.0               5446.0
7                   22.0                 11.0
8                   53.0                 27.0
9                 4104.0               2540.0
```

```
[4]:  # Generate a table of descriptive statistics about the data
      data.describe()
```

```
[4]:                    #        video_id  video_duration_sec  video_view_count  \
      count  19382.000000  1.938200e+04        19382.000000      19084.000000
      mean    9691.500000  5.627454e+09           32.421732     254708.558688
      std     5595.245794  2.536440e+09           16.229967     322893.280814
      min        1.000000  1.234959e+09            5.000000         20.000000
      25%     4846.250000  3.430417e+09           18.000000       4942.500000
      50%     9691.500000  5.618664e+09           32.000000       9954.500000
      75%    14536.750000  7.843960e+09           47.000000     504327.000000
      max    19382.000000  9.999873e+09           60.000000     999817.000000
```

```
             video_like_count  video_share_count  video_download_count  \
      count      19084.000000       19084.000000          19084.000000
      mean       84304.636030       16735.248323           1049.429627
      std       133420.546814       32036.174350           2004.299894
      min            0.000000           0.000000              0.000000
      25%          810.750000         115.000000              7.000000
      50%         3403.500000         717.000000             46.000000
      75%       125020.000000       18222.000000           1156.250000
```

```
max          657830.000000      256130.000000       14994.000000

       video_comment_count
count        19084.000000
mean           349.312146
std            799.638865
min              0.000000
25%              1.000000
50%              9.000000
75%            292.000000
max           9599.000000
```

Check for and handle missing values.

```
[5]: # Check for missing values
     data.isna().sum()
```

```
[5]: #                            0
     claim_status             298
     video_id                   0
     video_duration_sec         0
     video_transcription_text 298
     verified_status            0
     author_ban_status          0
     video_view_count         298
     video_like_count         298
     video_share_count        298
     video_download_count     298
     video_comment_count      298
     dtype: int64
```

```
[6]: # Drop rows with missing values
     data = data.dropna(axis=0)
```

```
[7]: # Display first few rows after handling missing values
     data.head(10)
```

```
[7]:    # claim_status    video_id  video_duration_sec  \
     0  1        claim  7017666017                  59
     1  2        claim  4014381136                  32
     2  3        claim  9859838091                  31
     3  4        claim  1866847991                  25
     4  5        claim  7105231098                  19
     5  6        claim  8972200955                  35
     6  7        claim  4958886992                  16
     7  8        claim  2270982263                  41
     8  9        claim  5235769692                  50
```

```
9  10      claim  4660861094                                45
```

```
                        video_transcription_text verified_status  \
0  someone shared with me that drone deliveries a…    not verified
1  someone shared with me that there are more mic…    not verified
2  someone shared with me that american industria…    not verified
3  someone shared with me that the metro of st. p…    not verified
4  someone shared with me that the number of busi…    not verified
5  someone shared with me that gross domestic pro…    not verified
6  someone shared with me that elvis presley has …    not verified
7  someone shared with me that the best selling s…    not verified
8  someone shared with me that about half of the …    not verified
9  someone shared with me that it would take a 50…        verified
```

```
  author_ban_status  video_view_count  video_like_count  video_share_count  \
0      under review          343296.0           19425.0              241.0
1            active          140877.0           77355.0            19034.0
2            active          902185.0           97690.0             2858.0
3            active          437506.0          239954.0            34812.0
4            active           56167.0           34987.0             4110.0
5      under review          336647.0          175546.0            62303.0
6            active          750345.0          486192.0           193911.0
7            active          547532.0            1072.0               50.0
8            active           24819.0           10160.0             1050.0
9            active          931587.0          171051.0            67739.0
```

```
   video_download_count  video_comment_count
0                   1.0                  0.0
1                1161.0                684.0
2                 833.0                329.0
3                1234.0                584.0
4                 547.0                152.0
5                4293.0               1857.0
6                8616.0               5446.0
7                  22.0                 11.0
8                  53.0                 27.0
9                4104.0               2540.0
```

You are interested in the relationship between `verified_status` and `video_view_count`. One approach is to examine the mean value of `video_view_count` for each group of `verified_status` in the sample data.

```python
[8]: # Compute the mean `video_view_count` for each group in `verified_status`
     data.groupby('verified_status')['video_view_count'].mean()
```

```
[8]: verified_status
     not verified    265663.785339
```

```
verified          91439.164167
Name: video_view_count, dtype: float64
```

### 4.2.2   Task 3. Hypothesis testing

Before you conduct your hypothesis test, consider the following questions where applicable to complete your code response:

1. Recall the difference between the null hypothesis and the alternative hypotheses. What are your hypotheses for this data project?

The null hypothesis represents what is expected about a sample statistic based on the understood status quo of a population parameter, whereas the alternative hypothesis will be a suprising, contradictory hypothesis about the population parameter.

We test the accuracy of the hypotheses with tests about sample statistics. In this case we want to see if the video_view_count is significantly impacted by the verified status of the author.

Your goal in this step is to conduct a two-sample t-test. Recall the steps for conducting a hypothesis test:

1. State the null hypothesis and the alternative hypothesis
2. Choose a signficance level
3. Find the p-value
4. Reject or fail to reject the null hypothesis

Hypotheses:

Null: There is not significant difference in video view count between verified and unverified accounts.

Alternate: There will be a signficant difference in video view count between verified and unverified accounts.

You choose 5% as the significance level and proceed with a two-sample t-test.

```python
# Conduct a two-sample t-test to compare means
not_verified = data[data["verified_status"] == "not␣
 ↪verified"]["video_view_count"]
verified = data[data["verified_status"] == "verified"]["video_view_count"]

stats.ttest_ind(a=not_verified, b=verified, equal_var=False)
```

```
[9]: Ttest_indResult(statistic=25.499441780633777, pvalue=2.6088823687177823e-120)
```

**Question:** Based on the p-value you got above, do you reject or fail to reject the null hypothesis?

We reject the null hypothesis since our p-value is far smaller than the .05 significance level we were testing for.

## 4.3   PACE: Execute

## 4.4   Step 4: Communicate insights with stakeholders

*Ask yourself the following questions:*

1. What business insight(s) can you draw from the result of your hypothesis test?

We have learned that we can expect a significant deviation in video views from verified accounts vs unverified accounts, but the direction and magnitude of this difference was not learned from this hypothesis test.

**Congratulations!** You've completed this lab. However, you may not notice a green check mark next to this item on Coursera's platform. Please continue your progress regardless of the check mark. Just click on the "save" icon at the top of this notebook to ensure your work has been logged.