# PREDICTING
# THE AVERAGE YEARLY EARNING

## FOR THE TOP 1000 CHANNELS ON YOUTUBE

By : Manal AlQahtani, Nadia AlGhamdi

As Project 2 of SDAIA Data Science Bootcamp ( T5 )

SDAIA Academy

METIS

# OUTLINE

**01**
INTRODUCTION

**02**
WORKFLOW

**03**
TOOLS

**04**
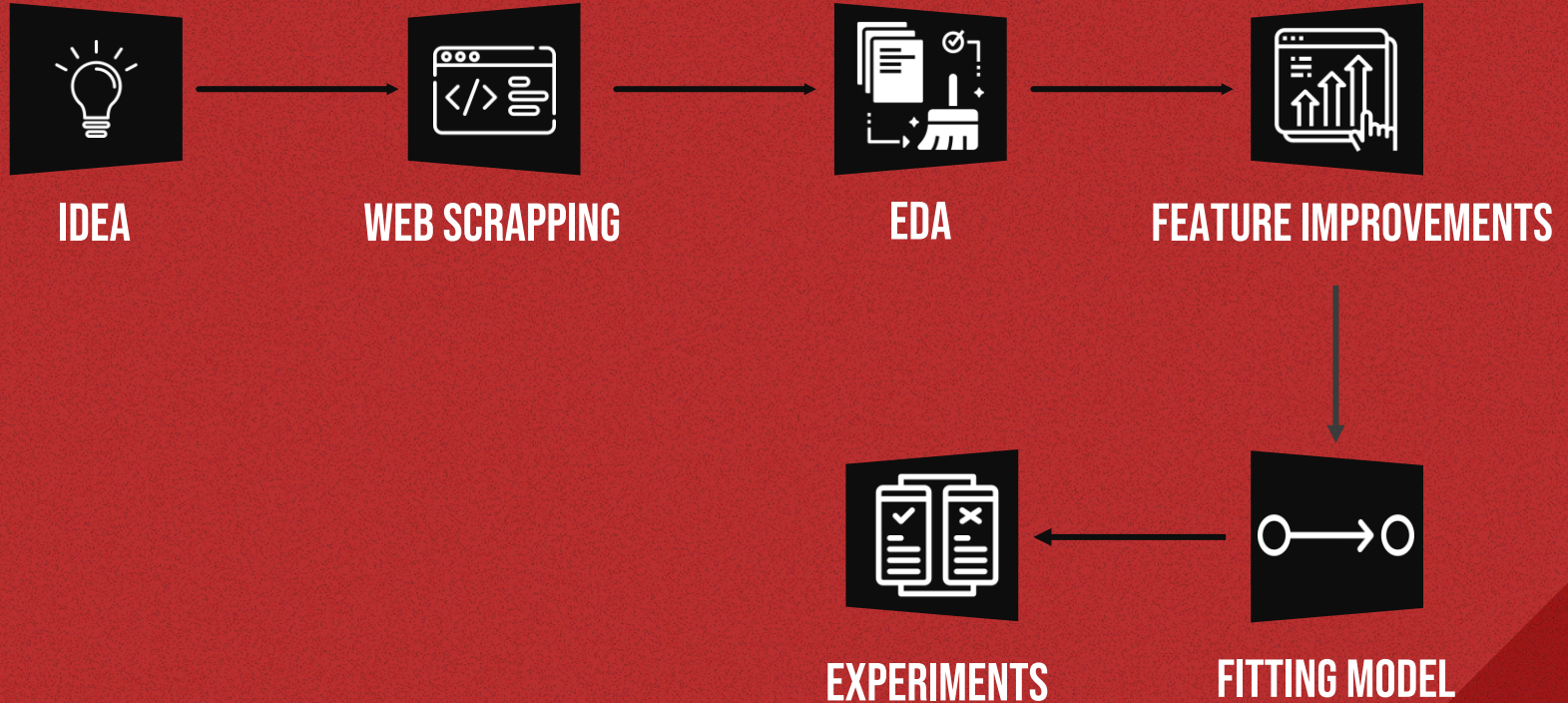IMPLEMENTATION

**05**
CONCLUSION

SDAIA Academy
أكاديمية سدايا

METIS

# INTRODUCTION

- Nowadays Social media is considered a wealth source, where everyone can make profit out of it .

- And the YouTube is considered half the internet, where 1.9 billion users logging in it .

# TOOLS

## PYTHON

Jupyter notebook     BeautifulSoup

NumPy, Pandas     Selenium

Matplotlib, Seaborn     Sklearn

## HTML

HTML

CSS

# DATASET



- Merged 2 datasets on Channel ID

- Create DataFrame out of them
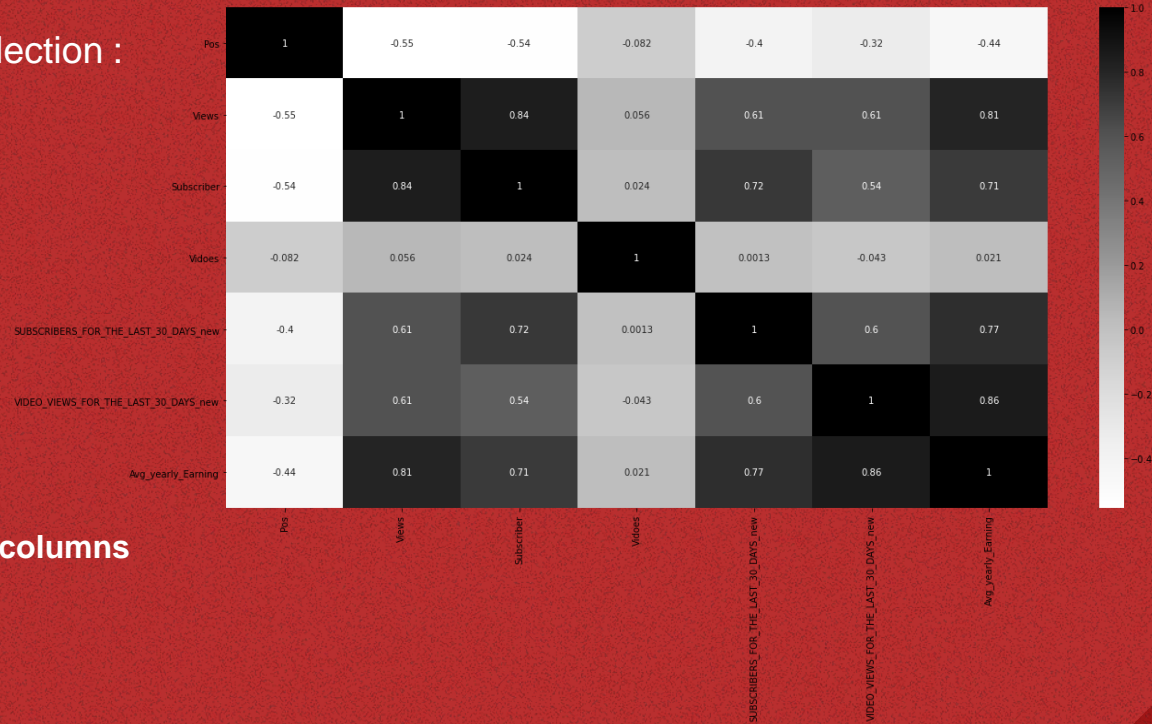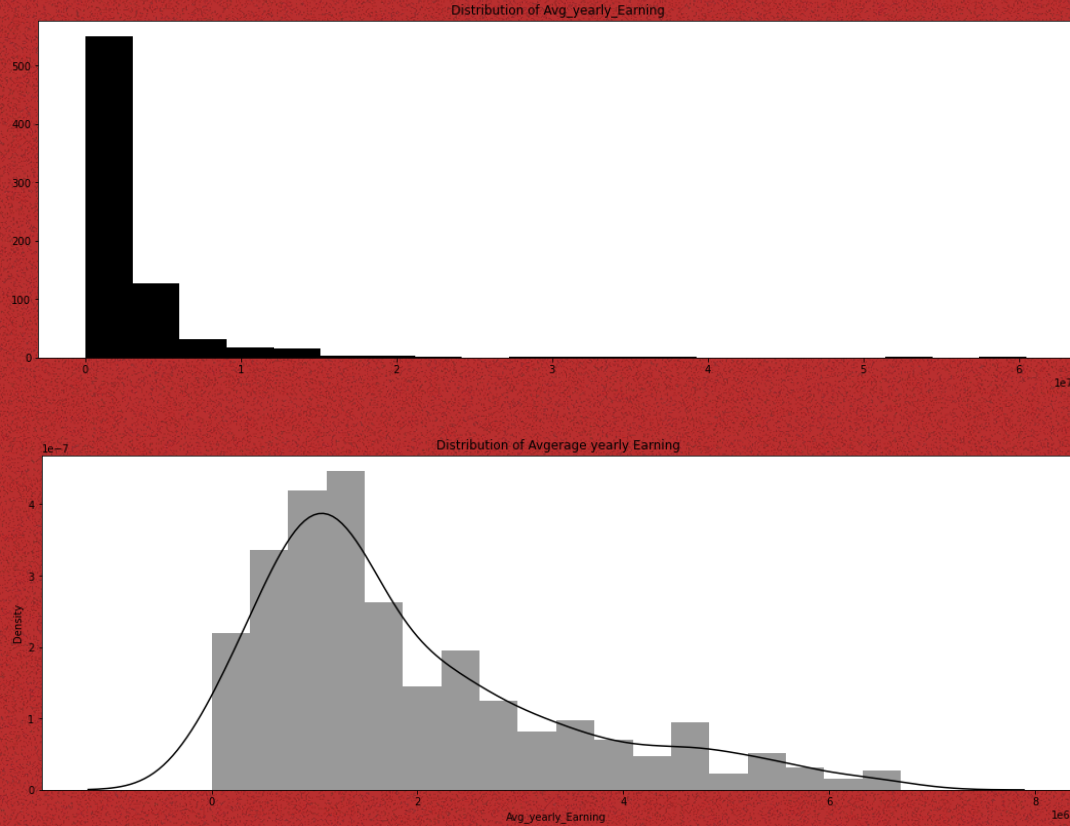
# EDA

- Features Engineering :
  Getting the Average monthly/ yearly earning
  out of min and max .

- Features Selection :



**687 rows × 13 columns**

# EDA

- Outliers :



Distribution of Avg_yearly_Earning

Distribution of Avgerage yearly Earning

- Testing Linearity between dependent and independents variables  :

- Log Transfromation to establish linear correlation between Y and features :

- Encoding categorical variables :
    encode with value between 0 and n_classes – 1 .

# MODELING

- **Linear regression result :**

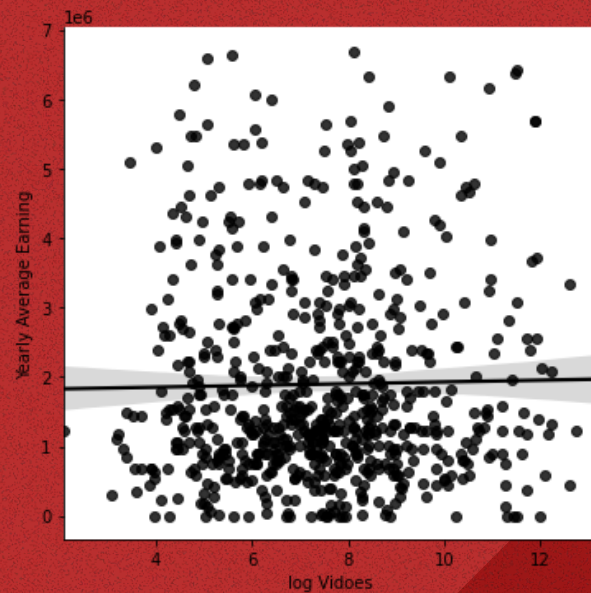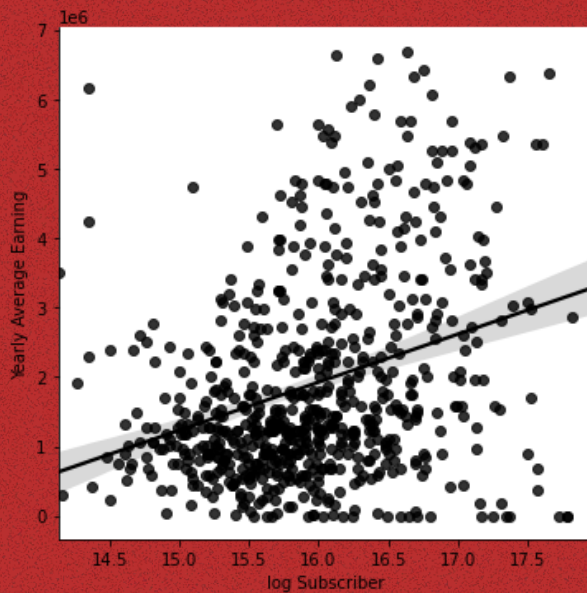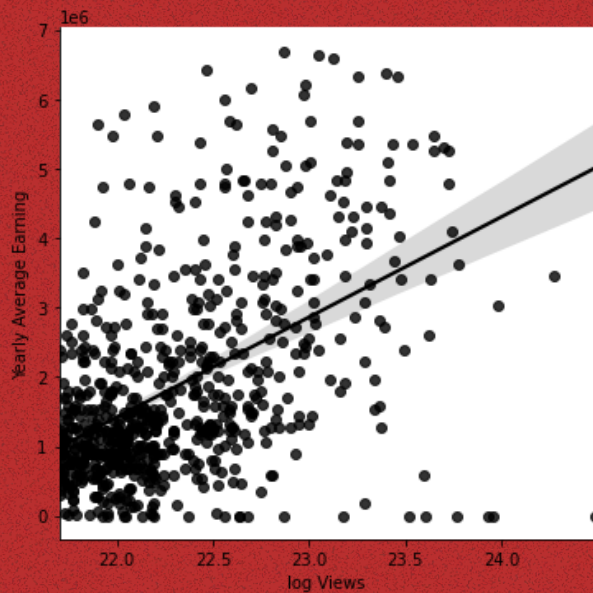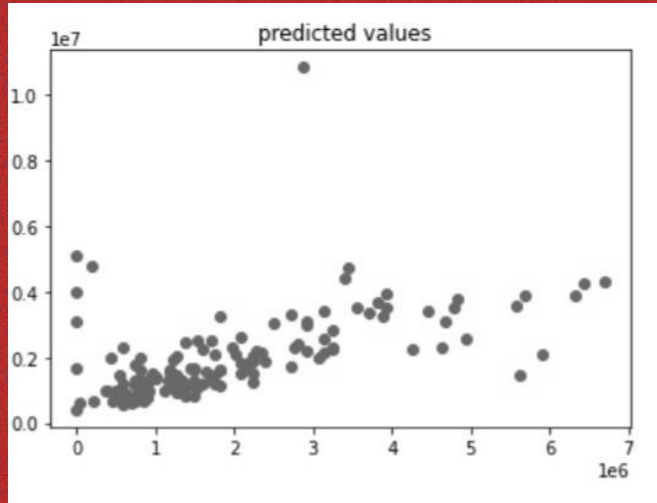| | r_squared | MAE | MSE | RMSE |
|---|---|---|---|---|
| **Linear Regression** | 0.634396 | 653861.543143 | 8.949487e+11 | 9.460173e+05 |

```
Index(['Pos', 'Channel', 'SOCIAL_BLADE_RANK', 'VIDEO_VIEWS_RANK',
       'COUNTRY_RANK', 'MUSIC_RANK', 'CHANNEL_TYPE',
       'SUBSCRIBERS_FOR_THE_LAST_30_DAYS_new',
       'VIDEO_VIEWS_FOR_THE_LAST_30_DAYS_new', 'Avg_yearly_Earning', 'Views',
       'Subscriber', 'Vidoes'],
      dtype='object')
```

```
1  r = pd.DataFrame(lr.coef_, columns = ['Coeffcients'])
2  r
```

| | Coeffcients |
|---|---|
| 0 | 1.574779e+03 |
| 1 | -1.184288e+02 |
| 2 | -6.511618e+01 |
| 3 | -1.164976e+02 |
| 4 | -4.005173e+02 |
| 5 | -6.594658e+02 |
| 6 | 2.075652e+04 |
| 7 | 6.866476e+00 |
| 8 | 3.876942e-03 |
| 9 | 2.492482e+06 |
| 10 | -5.830345e+05 |
| 11 | -2.199701e+03 |

- **Making prediction from LR**

- **Assess the performance of our model :**

## - Improving Linear Regression

|  | r_squared | MAE | MSE | RMSE |
|---|---|---|---|---|
| **Linear Regression** | 0.634396 | 653861.543143 | 8.949487e+11 | 9.460173e+05 |
| **Polynomial degree 4** | 0.951278 | 221565.225829 | 1.192647e+11 | 3.453473e+05 |
| **Polynomial degree 5** | 0.299057 | 787351.038821 | 1.715812e+12 | 1.309890e+06 |

**- Regulization :**

|  | R Squared | MAE | MSE | RMSE | Best Alpha |
|---|---|---|---|---|---|
| **Lasso** | 0.471861 | 571276.702456 | 9.130710e+11 | 955547.479501 | 0.0 |
| **Ridge** | 0.471861 | 571276.702456 | 9.130710e+11 | 955547.479501 | 0.0 |
| **Elastic Net** | 0.471861 | 571276.702456 | 9.130710e+11 | 955547.479501 | 0.0 |

# CONCLUSION

Even though polynomial showed low MSE, R squared is high to level that we reach overfitting . Thus, linear regression isn't the appropriate model to meet our goal !