# Deep Learning for Fraud Detection using GCP

Global IT Security Summit
5th Edition

# Table of Contents

# Fraud and it's Cost

According to ClearTax

     Financial fraud happens when someone deprives you of your money or otherwise harms your financial health through misleading, deceptive, or other illegal practices. This can be done through a variety of methods such as identity theft or investment fraud.

Jim Gee, Partner and National Head of Forensic Services at Crowe, said:

- "In the ten years since the first Financial Cost of Fraud report was published, the global economy has suffered rising losses each year, owing to a multitude of new and diverse threats. Losses in 2018 averaged 7.15% of expenditure, compared to 4.6% in 2007.

- "The figures quoted in the 2019 report are stark. Globally, fraud losses equate to a shocking US$5.127 trillion each year, which represents almost 70% of the $7.442 trillion which world spends on healthcare each year.

# Fighting Fraud with AI and Cloud

According to Svetlana Belyalova, head of operational risk management at Rosbank

"The use of AI to fight financial fraud, both internal and external, has become a hot topic. AI is the future of fraud management, irrespective of the system you are using…It brings a lot of value in both data management and decision-making."

To fight Financial Fraud we can use various variations of Classical Machine Learning Algorithms Like Isolation Forest, Reinforcement Learning techniques like MDP, or Deep Learning Family of Recurrent Neural Networks.

Moreover, putting these models into production adds another layer of complexity, therefore recuring to Cloud services such Google Cloud Platform is essential to the feasibility and efficiency of the  End-to-End Artificial Intelligence pipelines.

# The Datasets we have used

We have used two different datasets to benchmark our models:

- The first is :

Credit cards transactions by European cardholders anonymized dataset.

https://www.kaggle.com/mlg-ulb/creditcardfraud

- The second is :

Paysim synthetic dataset of mobile money transactions.

https://github.com/EdgarLopezPhD/PaySim

# Deep Learning models for Fraud Detection

In the Realm of Artificial Intelligence we have many sub fields, and our focus we'll be on Deep Learning, and specifically the family of Recurrent Neural Networks, that we have implemented using TensorFlow.
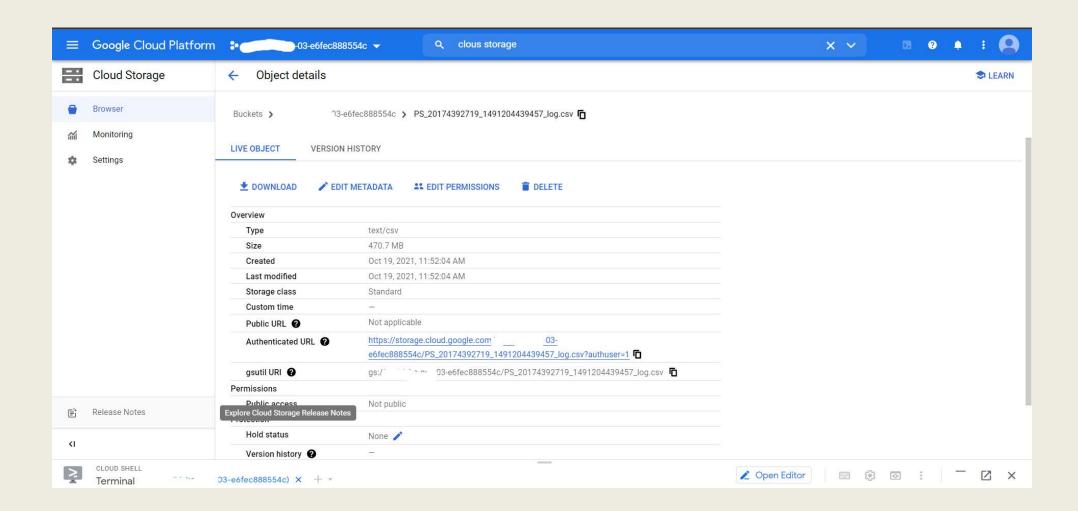
Through our research we have found that sequential models have better results if we introduce on top of them CNN layers(Convolutional Neural Networks) that extracts the best features and reduce the volume of the processed data.
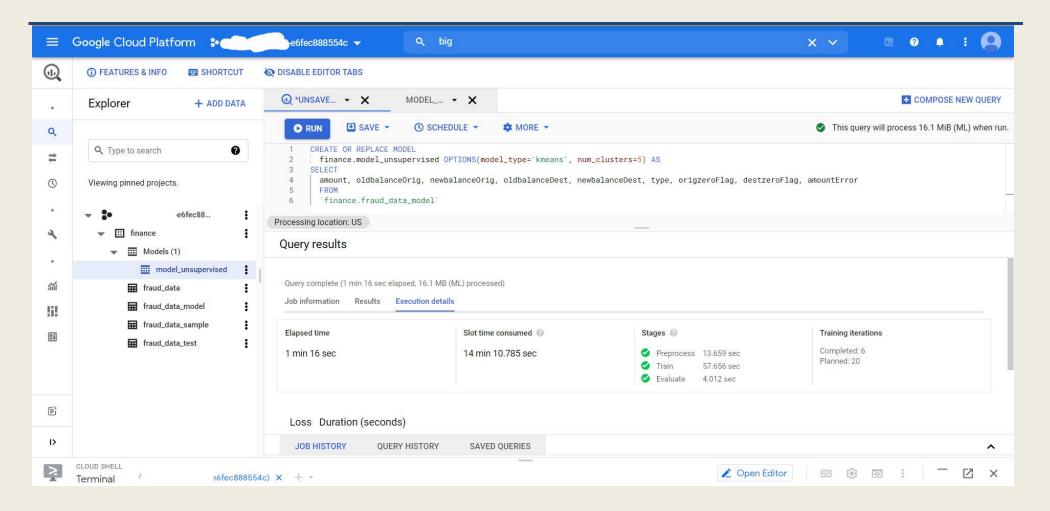
The benchmarking baseline have been established through the basic RNN(Recurrent Neural Networks), then we have used more advanced variation starting from LSTM(Long Short Term Memory) to BILSTM(Bidirectional LSTM).
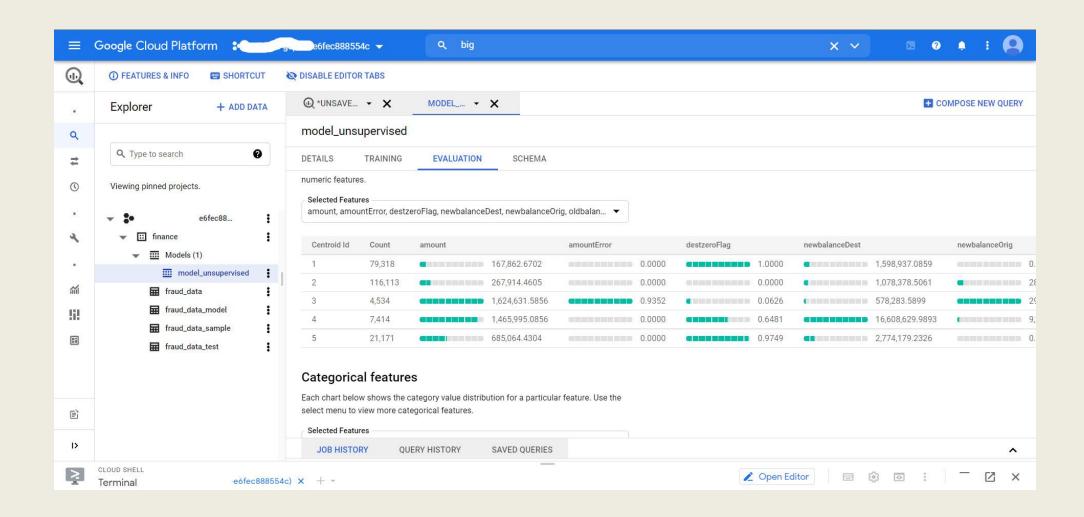
These algorithms have consumed a lot of computing resources and their not time efficient, so we have used Keras which have a great hyperparameter tuning, and we replaced BILSTM with BIGRU (Bidirectional Gated Recurrent Unit).
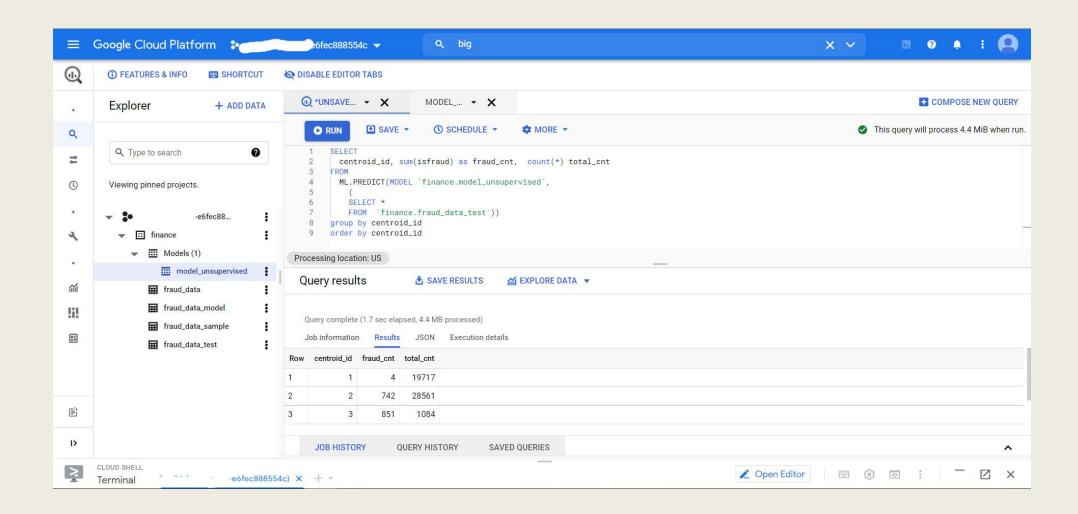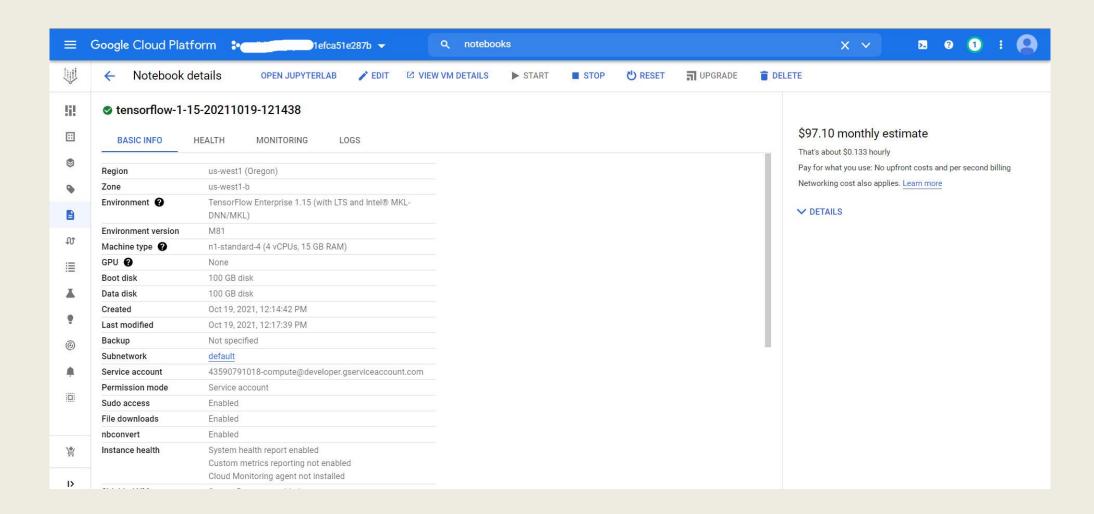
# Deep Learning models for Fraud Detection

```python
tf.keras.backend.clear_session()
tf.random.set_seed(51)
np.random.seed(51)

tf.keras.backend.clear_session()
dataset = windowed_dataset(x_train, window_size, batch_size, shuffle_buffer_size)


model = tf.keras.models.Sequential([
  tf.keras.layers.Lambda(lambda x: tf.expand_dims(x, axis=-1),
                         input_shape=[None]),

  tf.keras.layers.Conv1D(filters=32, kernel_size=5,
                         strides=1, padding="causal",
                         activation="relu",
                         input_shape=[None, 1]),

  tf.keras.layers.Bidirectional(tf.keras.layers.GRU(32, return_sequences=True)),
  tf.keras.layers.Bidirectional(tf.keras.layers.GRU(32, return_sequences=True)),

  tf.keras.layers.LSTM(64, return_sequences=True),
  tf.keras.layers.LSTM(64, return_sequences=True),

  tf.keras.layers.Dense(30, activation="relu"),
  tf.keras.layers.Dense(10, activation="relu"),
  tf.keras.layers.Dense(1),
  tf.keras.layers.Lambda(lambda x: x * 400.0)
])

optimizer = tf.keras.optimizers.SGD(lr=1e-5, momentum=0.9)
model.compile(loss=tf.keras.losses.Huber(),
              optimizer=optimizer,
              metrics=["mae"])
history_cnn_bigru = model.fit(dataset, epochs=150)
```

# Leveraging AI models with GCP

# Leveraging AI models with GCP

# Leveraging AI models with GCP

# Leveraging AI models with GCP

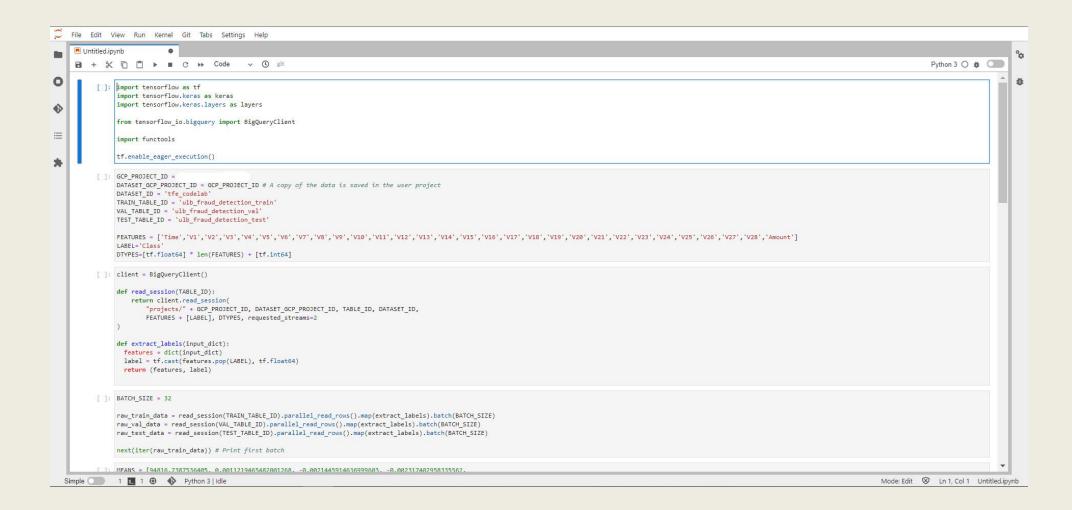# Leveraging AI models with GCP
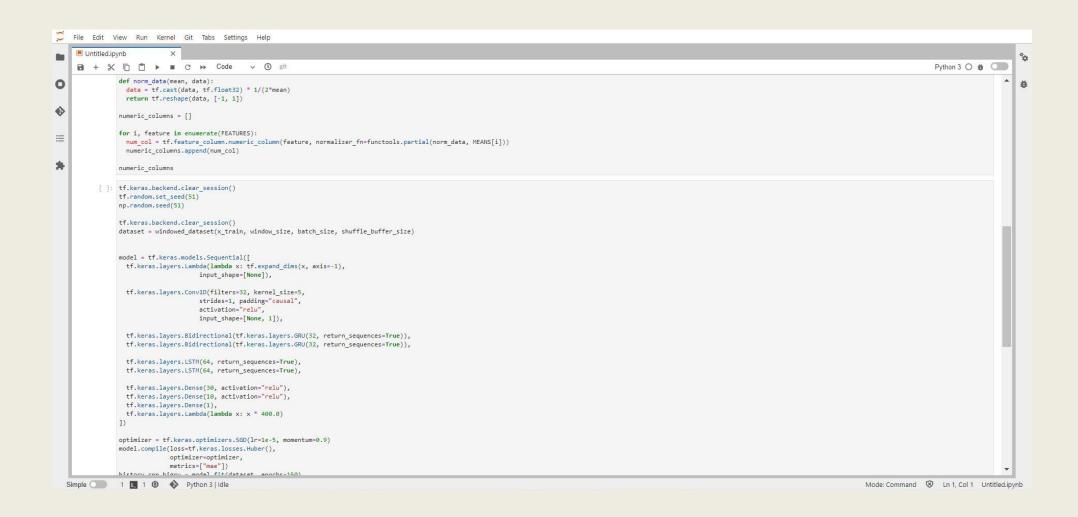
# Leveraging AI models with GCP

# Leveraging AI models with GCP

# Thank You

Conclusion & Perspectives

Q & A