**SPILL THE TECH WORKSHOP**

# Data Science Pipeline

Randy Galawana
Data Engineer

Telkomsel
by Telkom Indonesia

# Data Science Pipeline in General

**Data Retrieval**

```
query_database(query)
name_columns(mapping)
```

**Data Cleaning**

```
convert_to_double(column)
normalize_data(column)
```

**Model Training**

```
fit_model(hyperparams)
save_model(filename)
```

**Model Evaluation**

```
model_predict(test_data)
confusion_matrix(pred, true)
```

## Data Engineer Domain

- Big Data
- Different Tech in all use case
- Data Warehouse and Data Lake
- Advance Data Transformation
- Data Quality Checks

## Data Scientist Domain

- Sample Data
- Using most common tech / playground
- Statistic and Visualization
- Simple Data Cleaning
- Simple Data Transformation

Telkomsel

# Data Retrieval using Pandas

Retrieve data from known data sources

## Database (JDBC / ODBC)

Using pandas read_sql function with sqlalchemy connection, example in notebook.  Support most of ODBC and JDBC database

## IO File

Using pandas io read
https://pandas.pydata.org/docs/user_guide/io.html
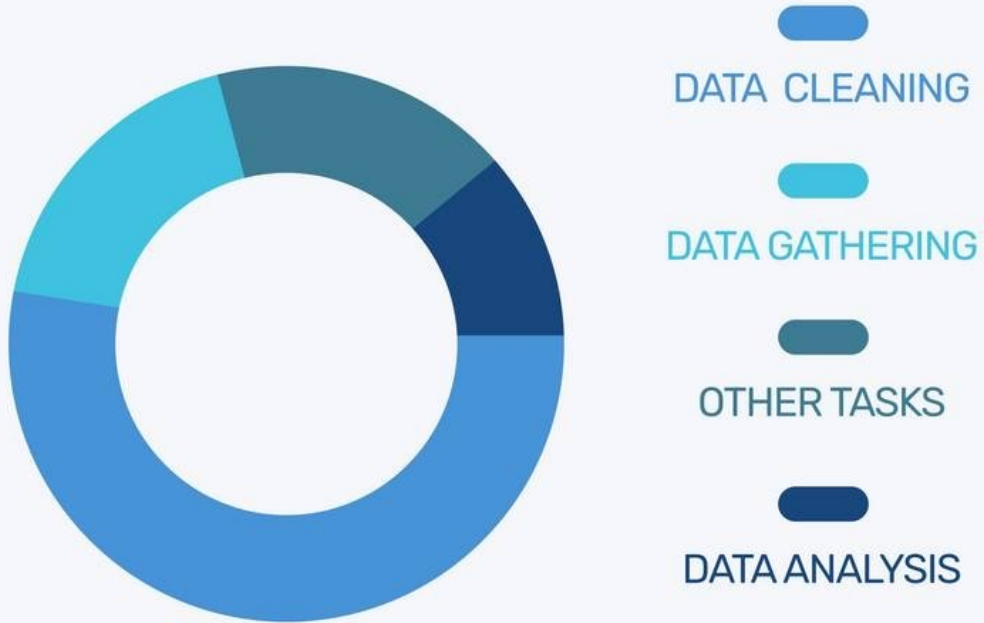
Support almost file extension and format

## API's

Using pandas io read
https://pandas.pydata.org/docs/user_guide/io.html

Support almost file extension and format

Telkomsel

# Data Cleaning

Why Data Cleaning need to include in Pipeline



DATA CLEANING

DATA GATHERING

OTHER TASKS

DATA ANALYSIS

Data scientists spend 60% of their time on cleaning data.

MonkeyLearn

Telkomsel

# Data Cleaning using Pandas

Most common Data cleaning using pandas and python

## Handling Null Values / Missing

Impu ter

- Simple Imputer
- Mode, Mean Values
- Hot Deck

Drop Na

- Thresholding

## Standardization and Normalization

### Category Encoder

- One hot encoding
- Label encoder, etc

### Scaler

- MinMaxScaler
- StandartScaler, etc

## De-Duplication



Original data — Deduplicated data

# Missing Value Treatment

# Data Transformation using pandas

## Pivot

df

| | foo | bar | baz | zoo |
|---|---|---|---|---|
| **0** | one | A | 1 | x |
| **1** | one | B | 2 | y |
| **2** | one | C | 3 | z |
| **3** | two | A | 4 | q |
| **4** | two | B | 5 | w |
| **5** | two | C | 6 | t |

```
df.pivot(index='foo',
        columns='bar',
        values='baz')
```

| bar | A | B | C |
|---|---|---|---|
| **foo** | | | |
| **one** | 1 | 2 | 3 |
| **two** | 4 | 5 | 6 |

## Melt

df3

| | first | last | height | weight |
|---|---|---|---|---|
| **0** | John | Doe | 5.5 | 130 |
| **1** | Mary | Bo | 6.0 | 150 |

```
df3.melt(id_vars=['first', 'last'])
```

| | first | last | variable | value |
|---|---|---|---|---|
| **0** | John | Doe | height | 5.5 |
| **1** | Mary | Bo | height | 6.0 |
| **2** | John | Doe | weight | 130 |
| **3** | Mary | Bo | weight | 150 |

## Stack

df2

| | | A | B |
|---|---|---|---|
| **first** | **second** | | |
| **bar** | one | 1 | 2 |
| | two | 3 | 4 |
| **baz** | one | 5 | 6 |
| | two | 7 | 8 |

MultiIndex

stacked = df2.stack()

| first | second | | |
|---|---|---|---|
| **bar** | one | A | 1 |
| | | B | 2 |
| | two | A | 3 |
| | | B | 4 |
| **baz** | one | A | 5 |
| | | B | 6 |
| | two | A | 7 |
| | | B | 8 |

MultiIndex

## Unstack(1)

stacked

| first | second | | |
|---|---|---|---|
| **bar** | one | A | 1 |
| | | B | 2 |
| | two | A | 3 |
| | | B | 4 |
| **baz** | one | A | 5 |
| | | B | 6 |
| | two | A | 7 |
| | | B | 8 |

MultiIndex

```
stacked.unstack(1)
        or
stacked.unstack('second')
```

| | second | one | two |
|---|---|---|---|
| **first** | | | |
| **bar** | A | 1 | 3 |
| | B | 2 | 4 |
| **baz** | A | 5 | 7 |
| | B | 6 | 8 |

MultiIndex

Telkomsel

# Data Transformation using pandas

# Data Transformation using pandas

## Transpose



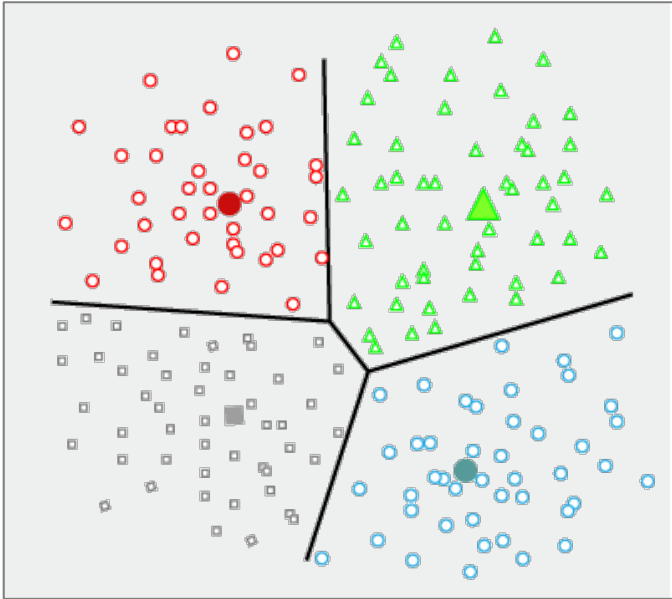## Group by / Aggregate



Telkomsel

# EDA and Data Viz using Seaborn

# Machine Learning



## Unsupervised Learning
### Clustering

## Supervised Learning
### Classification
### Regression

Telkomsel

# Machine Learning Using Scikit Learn



scikit-learn
algorithm cheat-sheet

**classification**

**regression**

**clustering**

**dimensionality reduction**

Telkomsel

Back

scikit learn

# Machine Learning Model Evaluation

## FOR REGRESSION

| | |
|---|---|
| **01** | MEAN ABSOLUTE ERROR |
| **02** | MEAN SQUARED ERROR |
| **03** | ROOT MEAN SQUARED ERROR |
| **04** | R – SQUARED |
| **05** | ADJUSTED R – SQUARED |

## FOR CLASSIFICATION

| | |
|---|---|
| **01** | ACCURACY SCORE |
| **02** | CONFUSION MATRIX |
| **03** | PRECISION & RECALL |
| **04** | F1 – SCORE |
| **05** | AUC – ROC CURVE |

# Sample Big Data ML Pipeline Stack



- Data Ingestion : Kafka, Apache NIFI, Sqoop
- Orchestrator : Apache NIFI, Airflow, Oozie
- Big Data Lake : Hadoop (HDFS, Hive)
- Data Transformation : Apache Spark (PySpark), Hive SQL, Apache Flink
- Data Visualization : Tablaeu Server, ReDash, PowerBI, Dashboard web
- Machine Learning : CDSW, SparkML, scikit-learn, kedro, MLFlow
- Data Quality : Apache Griffin

Telkomsel

# DATA & AI LANDSCAPE 2019

## INFRASTRUCTURE

### HADOOP ON-PREMISE
### HADOOP IN THE CLOUD
### STREAMING / IN-MEMORY
### NoSQL DATABASES
### NewSQL DATABASES
### GRAPH DBs
### MPP DBs
### CLOUD EDW
### SERVERLESS
### DATA TRANSFORMATION
### DATA INTEGRATION
### DATA GOVERNANCE
### MGMT / MONITORING
### STORAGE
### CLUSTER SVCS
### DATA GENERATION & LABELLING
### AI OPS
### GPU DBs & CLOUD
### HARDWARE

## ANALYTICS & MACHINE INTELLIGENCE

### DATA ANALYST PLATFORMS
### DATA SCIENCE PLATFORMS
### BI PLATFORMS
### VISUALIZATION
### MACHINE LEARNING
### COMPUTER VISION
### HORIZONTAL AI
### SPEECH & NLP
### SEARCH
### LOG ANALYTICS
### SOCIAL ANALYTICS
### WEB / MOBILE / COMMERCE ANALYTICS

## APPLICATIONS – ENTERPRISE

### SALES
### MARKETING - B2B
### MARKETING - B2C
### CUSTOMER EXPERIENCE / SERVICE
### ENTERPRISE PRODUCTIVITY
### HUMAN CAPITAL
### LEGAL
### REGTECH & COMPLIANCE
### FINANCE
### BACK OFFICE AUTOMATION & RPA
### SECURITY

## APPLICATIONS – INDUSTRY

### ADVERTISING
### EDUCATION
### REAL ESTATE
### GOV'T
### INTELLIGENCE
### FINANCE - INVESTING
### FINANCE - LENDING
### INSURANCE
### HEALTHCARE
### LIFE SCIENCES
### TRANSPORTATION
### AGRICULTURE
### COMMERCE
### INDUSTRIAL
### OTHER

## CROSS-INFRASTRUCTURE/ANALYTICS

## OPEN SOURCE

### FRAMEWORKS
### QUERY / DATA FLOW
### DATA ACCESS & DATABASES
### ORCHESTRATION & MGMT
### STREAMING & MESSAGING
### STAT TOOLS & LANGUAGES
### AI OPS & INFRA
### AI / MACHINE LEARNING / DEEP LEARNING
### SEARCH
### LOGGING & MONITORING
### VISUALIZATION
### COLLABORATION
### SECURITY

## DATA SOURCES & APIs

### HEALTH
### IOT
### FINANCIAL & ECONOMIC DATA
### AIR / SPACE / SEA
### PEOPLE / ENTITIES
### LOCATION INTELLIGENCE
### OTHER

## DATA RESOURCES

### DATA SERVICES
### INCUBATORS & SCHOOLS
### RESEARCH

FIRSTMARK
EARLY STAGE VENTURE CAPITAL

# Prerequisites

What need to be installed before Hands on

-Anaconda Jupyter
 Notebook
-Python v3.6 +

Or use **google collabs**

**https://colab.research.google.com/**

Telkomsel

# Get Sample Code

Download sample code from

**https://github.com/galawana/tselworkshop1**

(Download Zip and extract or git clone repository)

# Demo and Hands On

Upload extracted folder into **drive.google.com**

# Thank You

© Telkomsel 2022