# 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

# Prerequisites

**What need to be installed before the workshop**

-Anaconda Jupyter
 Notebook
-Python v3.6 +

Or use google collabs

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

# Get Sample Code

Download sample code from

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

(Download Zip and extract or git clone repository)
If using google collabs, upload to drive.google.com

# 01
—

# Data Retrieval / Ingestion

How to Ingest data from multiple source

# 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

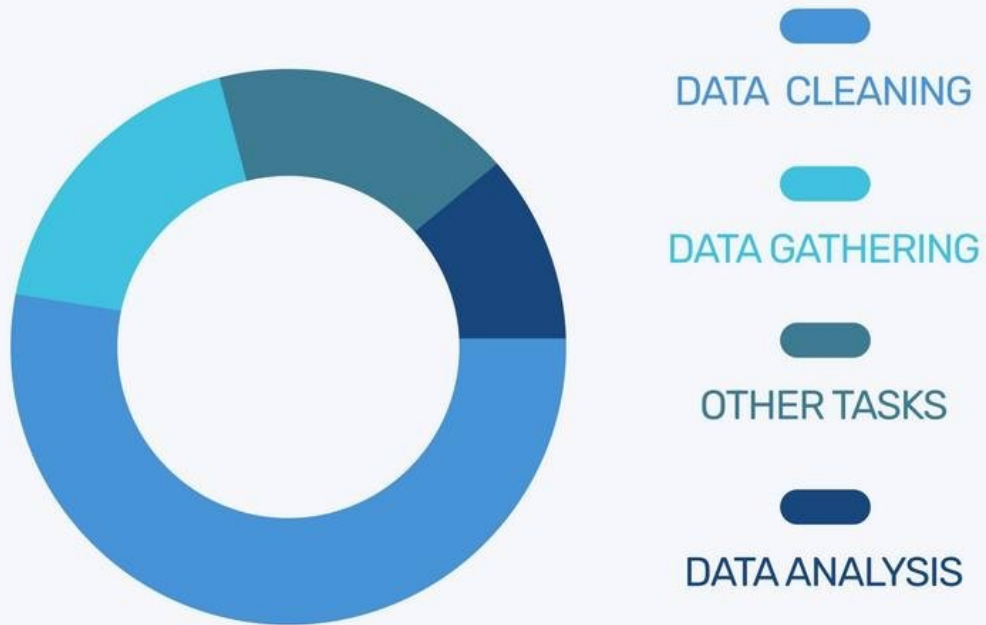Support almost file extension and format

### API's



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

Support almost file extension and format

Telkomsel

Open Sample code in
**Data_Retrieval.ipnyb**

**Demo and Hands on**

# 02
—

# Data Cleaning and Transformation

Clean the data, and do simple to advance transformation

Telkomsel

# 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
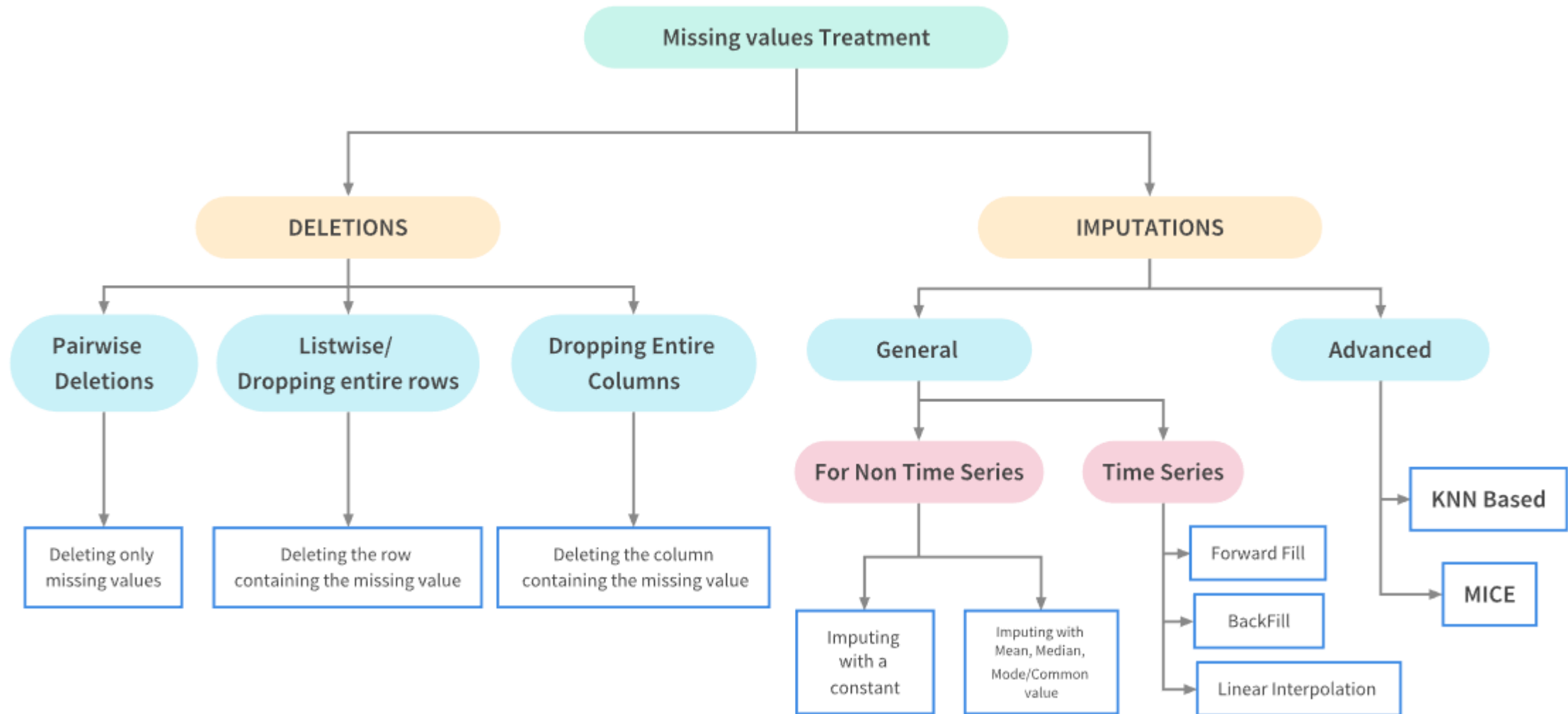
### Category Encoder

- One hot encoding
- Label encoder, etc

### Scaler

- MinMaxScaler
- StandartScaler, etc

## De-Duplication



Original data

Deduplicated data

Telkomsel

# Missing Value Treatment

Open Sample code in
**Data_Cleaning.ipnyb**

Demo and Hands on

# 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



## Join / Merge

SP04 | SP01 SP02 SP03 | SP05

LEFT

SP04 | SP01 SP02 SP03 | SP05

RIGHT

SP04 | SP01 SP02 SP03 | SP05

INNER

SP04 | SP01 SP02 SP03 | SP05

OUTER

## Concat

|   | Name | Age |
|---|------|-----|
| 0 | Sam  | 14  |
| 1 | Emma | 15  |

|   | Name  | Age |
|---|-------|-----|
| 0 | Karen | 10  |
| 1 | Rahul | 13  |

concat() →

|   | Name  | Age |
|---|-------|-----|
| 0 | Sam   | 14  |
| 1 | Emma  | 15  |
| 0 | Karen | 10  |
| 1 | Rahul | 13  |

Telkomsel

# Data Transformation using pandas

## Transpose

## Group by / Aggregate



©w3resource.com

Open Sample code in
**Data_Transformation.ipnyb**

**Demo and Hands on**

Telkomsel
by Telkom Indonesia

# 03
—

# EDA and Data Visualization with Seaborn

Explore and Tell the story from the data using
seaborn visualization library

Telkomsel

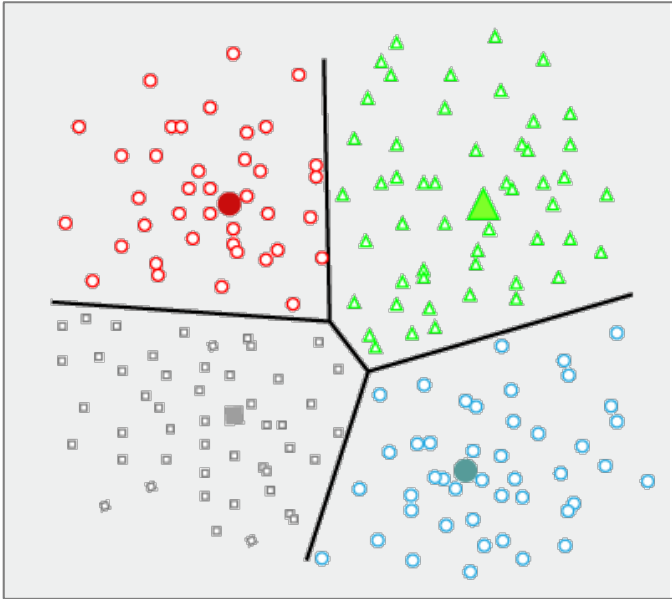# Seaborn chart library

Open Sample code in
**EDA.ipnyb**

**Demo and Hands on**

# 04

## Model Training dan Evaluation with scikit learn

Training ML model and evaluate the model

Telkomsel

# Type of Machine Learning

**Unsupervised Learning**

**Supervised Learning**
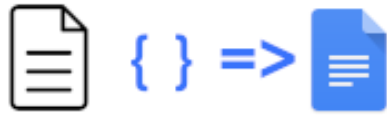
Clustering



Classification



Regression



Telkomsel

# Scikit Learn ML Library

# Scikit Learn Pipelines



Data Pipelines & ML Pipelines

Transformer → Estimator

Function that takes data and fit & transforms them into augmented data or feature

StandardScaler,TfidfVectorizer
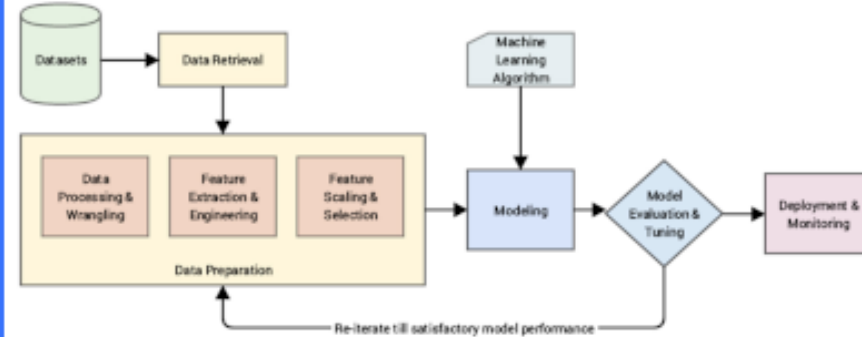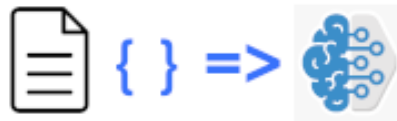
Function that takes data as input and fit the data and produces a model we can use to predict

LogisticRegression,KNN

Data To Data → Data To Model

Jesus Saves @JCharisTech

Telkomsel

# Model Evaluation

Open Sample code in
**Machine_Learning.ipnyb**

**Demo and Hands on**

# Thank You

© Telkomsel 2022