# Database Project (Fall 2023)
# HW3 (100pts, Due date: Nov 3)

**Student ID**: 2020315798

**Student Name**: Choi Jin Woo

**Instruction:** In this homework, we provide you with a jupyter notebook file (DBP_HW3.ipynb). You should follow the instructions in these documents carefully.

**Submit two files as follows:**

- DBP_HW3_StudentID.zip
    - DBP_HW3_StudentID.ipynb
    - DBP_HW3_StudentID.pdf

**[Spark SQL]**

**1. [20pts]** The following data is related to the manufacturing process of computer manufacturers.

*Data:*
- **Already in the manufacturing process → used for problem (a):**
[Row(phase='packing', model='book_pro', serial='book_pro1'),
Row(phase='packing', model='book_pro', serial='book_pro2'),
Row(phase='packing', model='plus', serial='plus1'),
Row(phase='packing', model='book_pro', serial='book_pro3'),
Row(phase='packing', model='plus', serial='plus2'),
Row(phase='inspection', model='book_pro', serial='book_pro4'),
Row(phase='inspection', model='plus', serial='plus3'),
Row(phase='inspection', model='book_pro', serial='book_pro5'),
Row(phase='inspection', model='book_pro', serial='book_pro6')]

- **The manufacturing process to add → used for problem (b):**
[Row(phase='assembly', model='book_pro', serial='book_pro7'),
Row(phase='assembly', model='plus', serial='plus4')]

(a) Create a DataFrame with the given data and display the generated DataFrame.

**[Answer]**
Enter your code and result here. You must show your result (captured image).

(a) Create a DataFrame with the given data and display the generated DataFrame.

```python
# ============= EDIT HERE =============

df = sc.parallelize([
    Row(phase='packing', model='book_pro', serial='book_pro1'),
    Row(phase='packing', model='book_pro', serial='book_pro2'),
    Row(phase='packing', model='plus', serial='plus1'),
    Row(phase='packing', model='book_pro', serial='book_pro3'),
    Row(phase='packing', model='plus', serial='plus2'),
    Row(phase='inspection', model='book_pro', serial='book_pro4'),
    Row(phase='inspection', model='plus', serial='plus3'),
    Row(phase='inspection', model='book_pro', serial='book_pro5'),
    Row(phase='inspection', model='book_pro', serial='book_pro6')
]).toDF()

# ===================================

df.show(truncate=False)
```

```
+----------+--------+---------+
|phase     |model   |serial   |
+----------+--------+---------+
|packing   |book_pro|book_pro1|
|packing   |book_pro|book_pro2|
|packing   |plus    |plus1    |
|packing   |book_pro|book_pro3|
|packing   |plus    |plus2    |
|inspection|book_pro|book_pro4|
|inspection|plus    |plus3    |
|inspection|book_pro|book_pro5|
|inspection|book_pro|book_pro6|
+----------+--------+---------+
```

(b) After adding **two laptops to the manufacturing process**, find the number of products for each model.

**[Answer]**
Enter your code and result here. You must show your result (captured image).

(b) After adding two laptops to the manufacturing process, find the number of products for each model.

```python
# ============= EDIT HERE =============

df2 = sc.parallelize([
    Row(phase='assembly', model='book_pro', serial='book_pro7'),
    Row(phase='assembly', model='plus', serial='plus4')
]).toDF()

new_df = df.unionByName(df2).groupBy("model").count()

# ===================================
new_df.show(truncate=False)
```

```
+--------+-----+
|model   |count|
+--------+-----+
|book_pro|7    |
|plus    |4    |
+--------+-----+
```

(c) Group the data in the joined DataFrame by 'phases' and count the number of data for each phase.

**[Answer]** Enter your code and result here. You have to show your snapshot result.

(c) Group the data in the joined DataFrame by 'phases' and count the number of data for each phase.

```python
# ============= EDIT HERE =============
group_df = df.unionByName(df2).groupBy("phase").count()

# ===================================
group_df.show(truncate=False)
```

```
+----------+-----+
|phase     |count|
+----------+-----+
|packing   |5    |
|inspection|4    |
|assembly  |2    |
+----------+-----+
```

**2. [20pts]** The following data are *manufacturing process* and *customer information* for computer manufacturers.
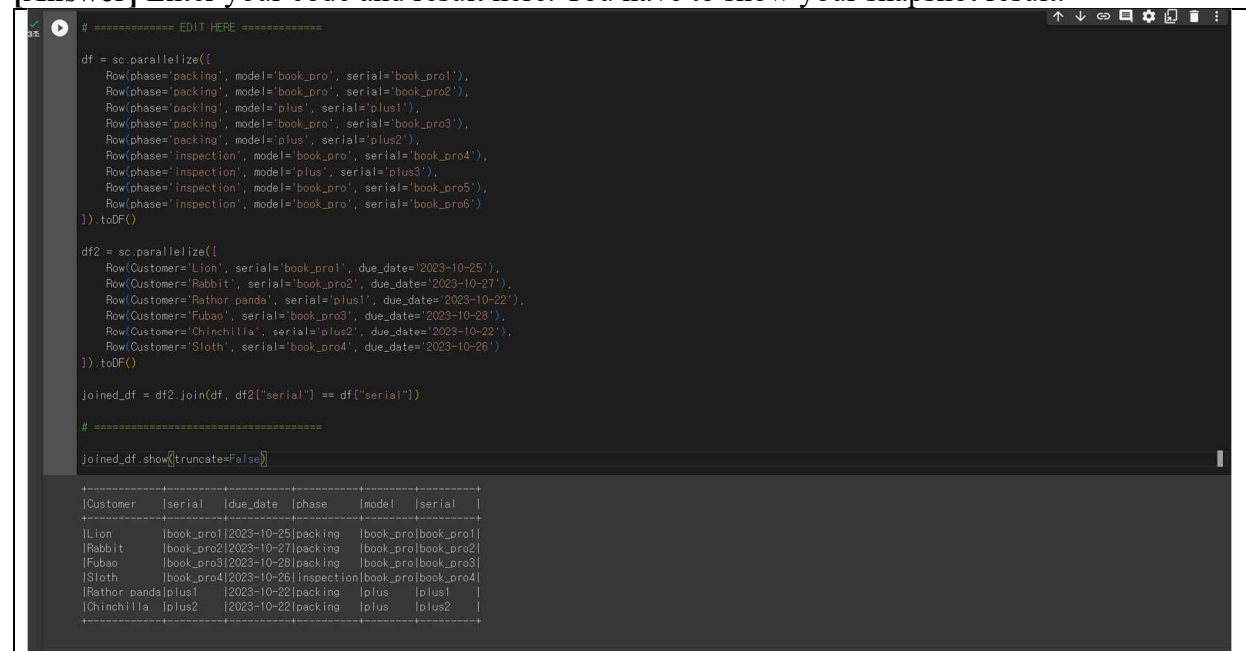
- ***Data 1(manufacturing process):***
[Row(phase='packing', model='book_pro', serial='book_pro1'),
Row(phase='packing', model='book_pro', serial='book_pro2'),
Row(phase='packing', model='plus', serial='plus1'),
Row(phase='packing', model='book_pro', serial='book_pro3'),
Row(phase='packing', model='plus', serial='plus2'),
Row(phase='inspection', model='book_pro', serial='book_pro4'),
Row(phase='inspection', model='plus', serial='plus3'),
Row(phase='inspection', model='book_pro', serial='book_pro5'),
Row(phase='inspection', model='book_pro', serial='book_pro6')]

- ***Data 2(customer information):***
[Row(Customer='Lion', serial='book_pro1', due_date='2023-10-25'),
Row(Customer='Rabbit', serial='book_pro2', due_date='2023-10-27'),
Row(Customer='Rathor panda', serial='plus1', due_date='2023-10-22'),
Row(Customer='Fubao', serial='book_pro3', due_date='2023-10-28'),
Row(Customer='Chinchilla', serial='plus2', due_date='2023-10-22'),
Row(Customer='Sloth', serial='book_pro4', due_date='2023-10-26')]

(a) Create a DataFrame for the two given data and join Data 1 with Data 2 using an inner join based on the 'serial' column. (left side: Data 2, right side: Data 1)
**[Answer]** Enter your code and result here. You have to show your snapshot result.



```python
# ============= EDIT HERE ==============

df = sc.parallelize([
    Row(phase='packing', model='book_pro', serial='book_pro1'),
    Row(phase='packing', model='book_pro', serial='book_pro2'),
    Row(phase='packing', model='plus', serial='plus1'),
    Row(phase='packing', model='book_pro', serial='book_pro3'),
    Row(phase='packing', model='plus', serial='plus2'),
    Row(phase='inspection', model='book_pro', serial='book_pro4'),
    Row(phase='inspection', model='plus', serial='plus3'),
    Row(phase='inspection', model='book_pro', serial='book_pro5'),
    Row(phase='inspection', model='book_pro', serial='book_pro6')
]).toDF()

df2 = sc.parallelize([
    Row(Customer='Lion', serial='book_pro1', due_date='2023-10-25'),
    Row(Customer='Rabbit', serial='book_pro2', due_date='2023-10-27'),
    Row(Customer='Rathor panda', serial='plus1', due_date='2023-10-22'),
    Row(Customer='Fubao', serial='book_pro3', due_date='2023-10-28'),
    Row(Customer='Chinchilla', serial='plus2', due_date='2023-10-22'),
    Row(Customer='Sloth', serial='book_pro4', due_date='2023-10-26')
]).toDF()

joined_df = df2.join(df, df2["serial"] == df["serial"])

# ==========================================

joined_df.show(truncate=False)
```

```
+------------+---------+----------+----------+--------+---------+
|Customer    |serial   |due_date  |phase     |model   |serial   |
+------------+---------+----------+----------+--------+---------+
|Lion        |book_pro1|2023-10-25|packing   |book_pro|book_pro1|
|Rabbit      |book_pro2|2023-10-27|packing   |book_pro|book_pro2|
|Fubao       |book_pro3|2023-10-28|packing   |book_pro|book_pro3|
|Sloth       |book_pro4|2023-10-26|inspection|book_pro|book_pro4|
|Rathor panda|plus1    |2023-10-22|packing   |plus    |plus1    |
|Chinchilla  |plus2    |2023-10-22|packing   |plus    |plus2    |
+------------+---------+----------+----------+--------+---------+
```

(b) Use an SQL query to select the data from the joined DataFrame where the **'due_date'** is on or after **'2023-10-25'**. And briefly explain the method you used.

**[Answer]** Enter your code and result here. You have to show your snapshot result.

```python
# ============ EDIT HERE ============
joined_df.createOrReplaceTempView("joined_df")

sql_query = """
    SELECT *
    FROM joined_df
    WHERE due_date >= '2023-10-25'
"""

sql_df = spark.sql(sql_query)

# ===================================

sql_df.show(truncate=False)
```

```
+--------+---------+----------+----------+--------+---------+
|Customer|serial   |due_date  |phase     |model   |serial   |
+--------+---------+----------+----------+--------+---------+
|Lion    |book_pro1|2023-10-25|packing   |book_pro|book_pro1|
|Rabbit  |book_pro2|2023-10-27|packing   |book_pro|book_pro2|
|Fubao   |book_pro3|2023-10-28|packing   |book_pro|book_pro3|
|Sloth   |book_pro4|2023-10-26|inspection|book_pro|book_pro4|
+--------+---------+----------+----------+--------+---------+
```

[createOrReplaceTempView]
To save Dataframe into temporary SQL table named "joined_df"

[sql_query]
SQL query that filters where "due_date" is on or after "2023-10-25"

[spark.sql()]
Returns new Dataframe with query result

**[Spark ML]**
**3. [60pts]** We provide you with a ***Fashion-MNIST*** dataset.

***Dataset Description:***

**Training set:** 60,000 examples
**Test set:** 10,000 examples
Each example is a 28x28 grayscale image associated with a label from 10 classes.

| Label | Description |
|-------|-------------|
| 0 | T-shirt/top |
| 1 | Trouser |
| 2 | Pullover |
| 3 | Dress |
| 4 | Coat |
| 5 | Sandal |
| 6 | Shirt |
| 7 | Sneaker |
| 8 | Bag |
| 9 | Ankle boot |

For more information, visit this website: https://github.com/zalandoresearch/fashion-mnist

(a) Load the provided dataset, convert it into a DataFrame, and show it. You should follow the following instructions.

   ***-Instructions 1:*** *Assemble the features into a vector column and name the column "features."*

*-Instructions 2: Rename the target column to "label."*

[**Answer**] Enter your code and result here. You have to show your snapshot result.



(b) Train models to classify the classes of the Fashion MNIST dataset and **report the results for test data**. The models used are **Logistic Regression, Decision Tree, and Random Forest.**

For detailed explanations of the usage of each model, please refer to the official documentation below.

- Logistic Regression: [Link]
- Decision Tree: [Link]
- Random Forest: [Link]

[**Answer**] Fill in the table below.

|  | **Logistic Regression** | **Decision Tree** | **Random Forest** |
|---|---|---|---|
| **Test accuracy** | 0.7561 | 0.7441 | 0.7724 |

[**Answer**] Enter your code and result here. You have to show your snapshot result.

## Logistic Regression

Reference: https://spark.apache.org/docs/latest/api/python/reference/api/pyspark.ml.classification.LogisticRegression.html

```python
from pyspark.ml.classification import LogisticRegression

# Training and Test
# ============= EDIT HERE =============
lr = LogisticRegression(maxIter=100, regParam=0.05, elasticNetParam=0.5)
lr_model = lr.fit(train_data)
lr_preds = lr_model.transform(test_data)
lr_accuracy = evaluator.evaluate(lr_preds)

# ===================================

print(f"Accuracy: {lr_accuracy}")
```

```
Accuracy: 0.7560803665844201
```

## Decision Tree

Reference: https://spark.apache.org/docs/latest/api/python/reference/api/pyspark.ml.classification.DecisionTreeClassifier.html

```python
from pyspark.ml.classification import DecisionTreeClassifier

# On model declaration, fix seed, maxDepth,  to the following values
seed = 2023

# ============= EDIT HERE =============

dt = DecisionTreeClassifier(seed=seed, maxDepth=15)
dt_model = dt.fit(train_data)
dt_preds = dt_model.transform(test_data)
dt_accuracy = evaluator.evaluate(dt_preds)

# ===================================

print(f"Accuracy: {dt_accuracy}")
```

```
Accuracy: 0.7440958759252732
```

## Random Forest

Reference: https://spark.apache.org/docs/latest/api/python/reference/api/pyspark.ml.classification.RandomForestClassifier.html

```python
[12] from pyspark.ml.classification import RandomForestClassifier

    # On model declaration, fix seed to the following values
    seed = 2023

    # ============= EDIT HERE =============
    rfc = RandomForestClassifier(seed=seed)
    rfc_model = rfc.fit(train_data)
    rfc_preds = rfc_model.transform(test_data)
    rfc_accuracy = evaluator.evaluate(rfc_preds)

    # ===================================

    print(f"Accuracy: {rfc_accuracy}")
```

```
Accuracy: 0.7724121724826695
```