**databricks**

# 14 DAYS

## AI CHALLENGE

# DAY 13

## Topic:

Model Comparison & Feature Engineering

## Challenge:

1. Train 3 different models

2. Compare metrics in MLflow

3. Build Spark ML pipeline

4. Select best model

```python
#Loading Data
# Prepare data
df = spark.table("ecommerce.gold.product_metrics").toPandas()
if df.shape[0] == 0:
    print("No data available in ecommerce.silver.product_metrics.
    Cannot proceed with train/test split.")
else:
    X = df[["total_events"]]
    y = df["purchases"]
    X_train, X_test, y_train, y_test = train_test_split(X, y,
    test_size=0.2, random_state=42)
```

> 📊 See performance (1)

> ▤  df:  pandas.core.frame.DataFrame = [product_id: object, brand: object ... 4 more
fields]
> ▤  X:  pandas.core.frame.DataFrame = [total_events: int64]
> ▤  X_test:  pandas.core.frame.DataFrame = [total_events: int64]
> ▤  X_train:  pandas.core.frame.DataFrame = [total_events: int64]

✓ 11:34 PM (<1s)                    4: Cell 4

```python
models = {
    "LinearRegression": LinearRegression(),
    "DecisionTreeRegressor": DecisionTreeRegressor(max_depth = 5),
    "RandomForestRegressor": RandomForestRegressor
    (n_estimators=100, random_state=42)
}
```

```python
from pyspark.ml import Pipeline
from pyspark.ml.feature import VectorAssembler
from pyspark.ml.regression import LinearRegression

assembler = vectorAssembler = VectorAssembler(inputCols=["total_events"], outputCol="features")
lr = LinearRegression(featuresCol="features", labelCol="purchases")
pipeline = Pipeline(stages=[assembler, lr])
spark_df =spark.table("ecommerce.gold.product_metrics").fillna(0)
train, test = spark_df.randomSplit([0.8, 0.2], seed=42)
pipeline_model = pipeline.fit(train)
# Evaluate the model
predictions = pipeline_model.transform(test)
display(predictions.select("product_id", "brand", "purchases", "prediction"))
```

> 📊 See performance (1)

> ▦ predictions: pyspark.sql.connect.dataframe.DataFrame
> ▦ spark_df: pyspark.sql.connect.dataframe.DataFrame = [product_id: string, brand: string ... 4 more fields]
> ▦ test: pyspark.sql.connect.dataframe.DataFrame = [product_id: string, brand: string ... 4 more fields]
> ▦ train: pyspark.sql.connect.dataframe.DataFrame = [product_id: string, brand: string ... 4 more fields]

Table ∨    +

|   | ᴬᴮᵪ product_id | ᴬᴮᵪ brand | ¹²₃ purchases | 1.2 prediction |
|---|---------------|-----------|---------------|----------------|
| 1 | 100000010 | eksmo | 0 | -19.7436973545211... |
| 2 | 100000014 | eksmo | 0 | -20.3976444950013... |
| 3 | 100000017 | null | 18 | 10.555853487728463 |

| ☰ | 📈 | ▦ | 🔍 metrics.rmse < 1 and params.model = "tree" | ⓘ |

⇅ Sort: Created ⌄        ▦ Group by ⌄

| ☐ | 👁 | **Run Name** | | 🔍 Search metric charts |
| ☐ | 👁 | 🟣 RandomForestRegressor | | |
| ☐ | 👁 | 🟢 DecisionTreeRegressor | | ⌄  **Model metrics (1)** |
| ☐ | 👁 | 🟠 LinearRegression | | |
| ☐ | 👁 | 🔴 RandomForestRegressor | | ⠿  r2                    ⌞⌟  ⋮ |
| ☐ | 👁 | 🟣 DecisionTreeRegressor | | |
| ☐ | 👁 | 🟢 LinearRegression | | |
| ☐ | 👁 | 🔴 LinearRegression | | |
| ☐ | 👁 | 🔵 LinearRegression | | |

Model metrics r2 chart:
- 🟣 0.92
- 🟢 0.93
- 🟠 0.92
- 🩷 0.92
- 🟣 0.93
- 🟢 0.92

Axis: 0    0.2    0.4    0.6    0.8