



Forecasting Retail Sales for Walmart at scale

IDS 561 – Analytics for Big Data



1

Problem Setting

Context & reason for our Analysis

Problem Setting

The Competitive Landscape & Thin Margins in the Retail Segment



Fierce Competition



Expansion of Stores



Retain Loyal Customers

Solution: Leverage Data



PySpark

2

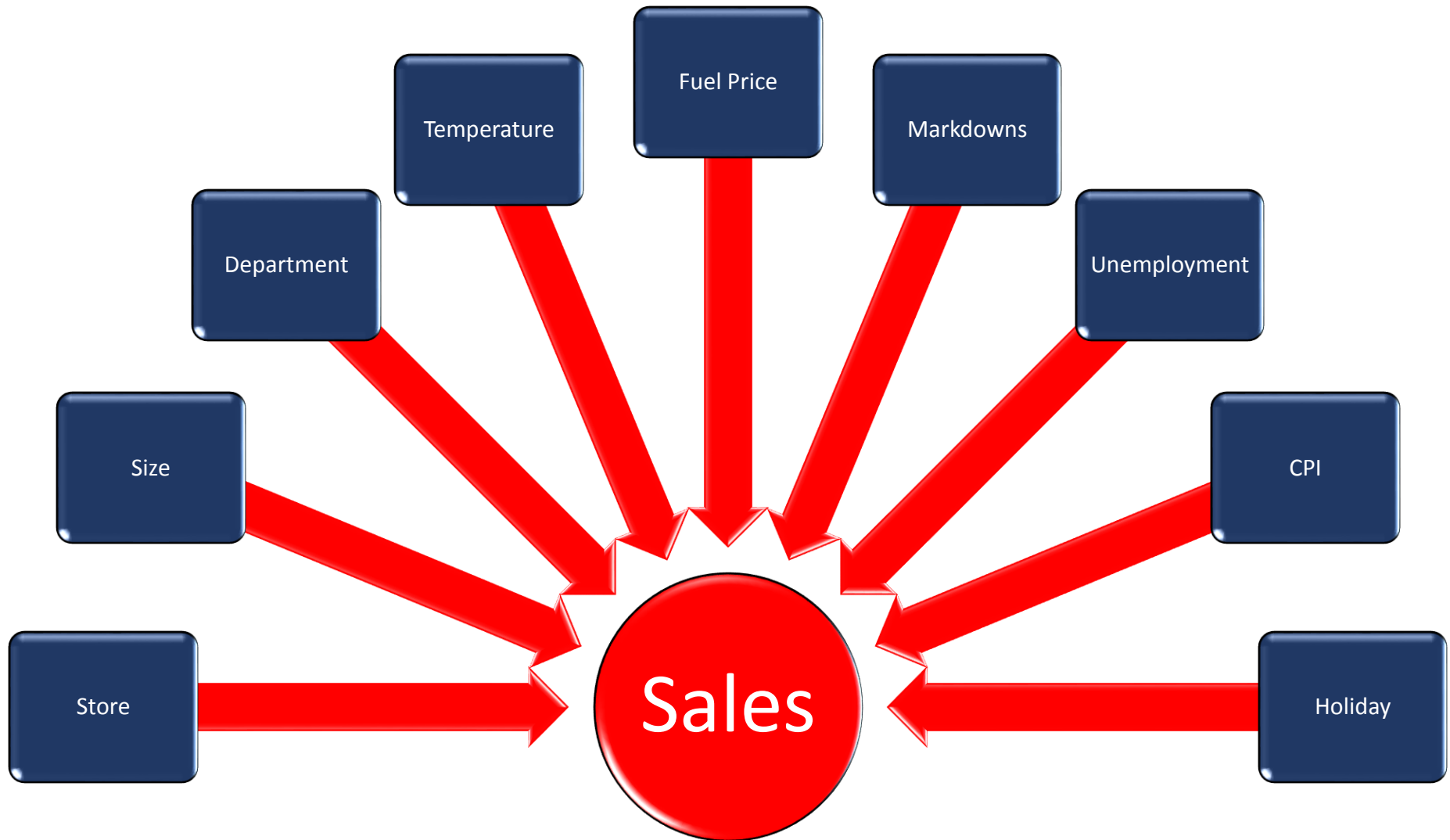
Data Description

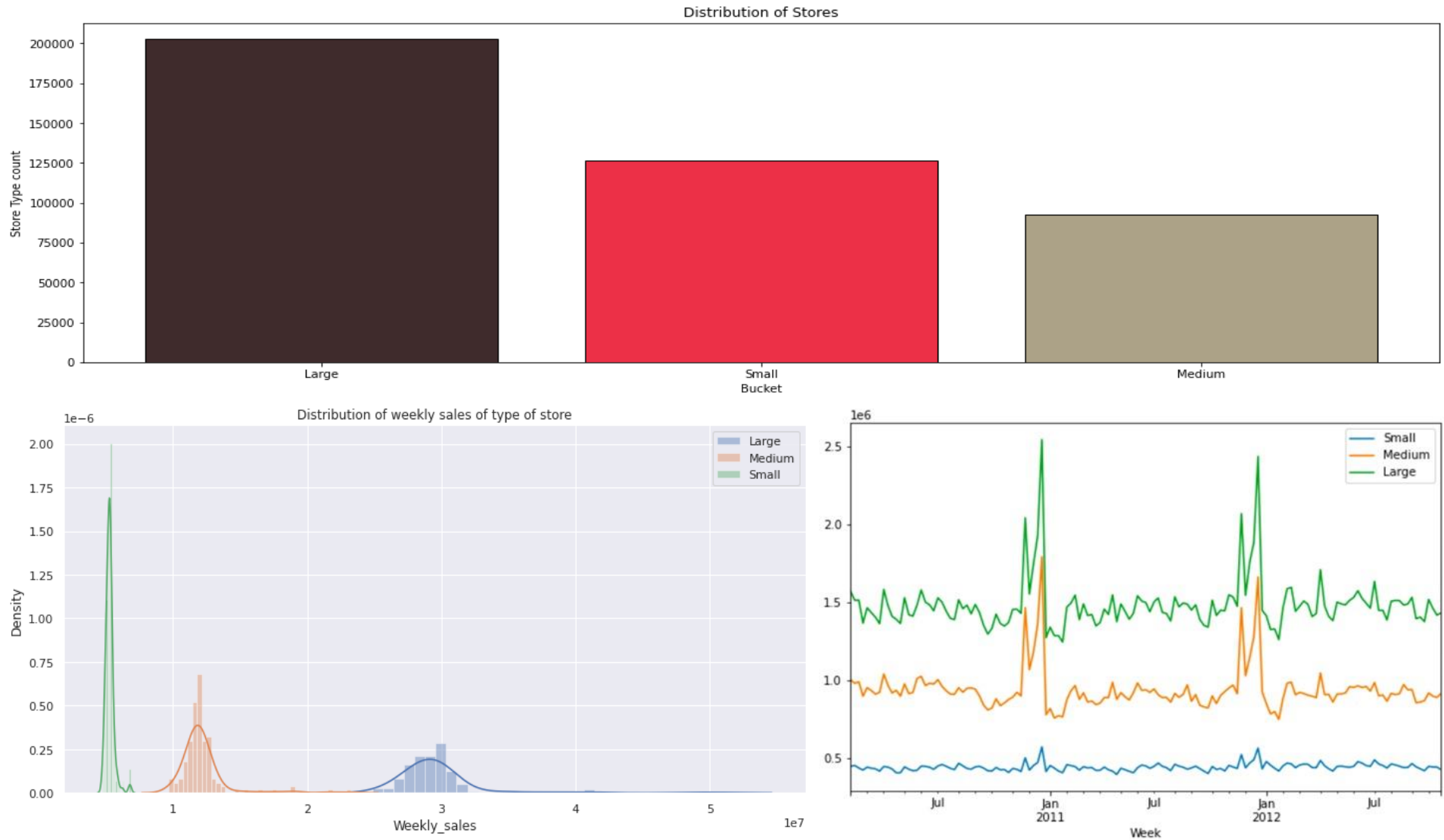
Details about the Dataset

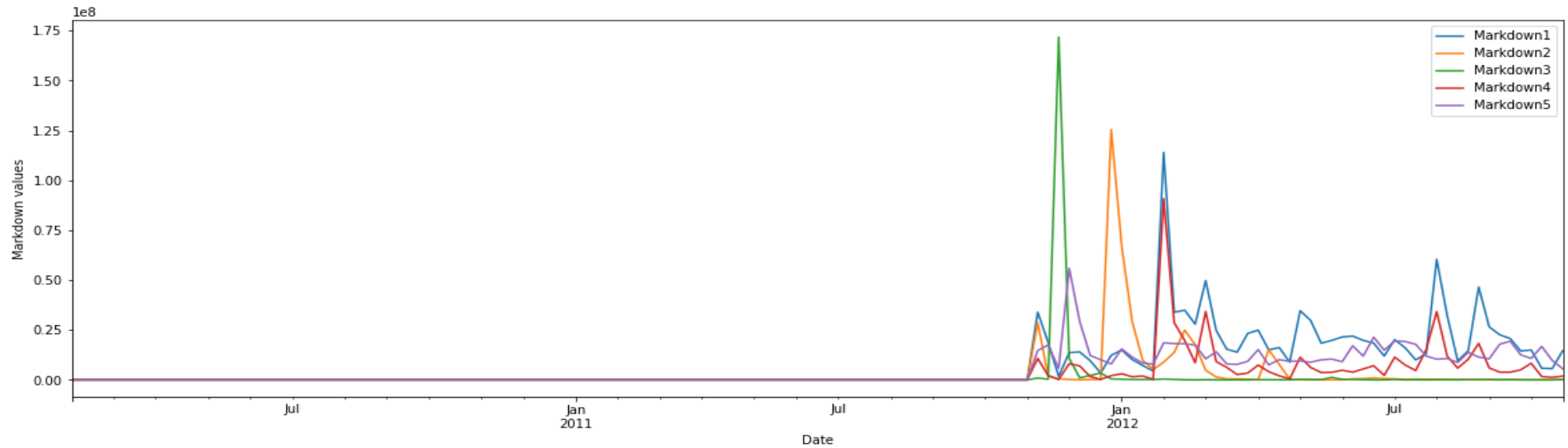
- We have used publicly available sales data for 45 Walmart stores across the United States. (Available on Kaggle)
- 420K+ rows for training and 100K+ rows of validation data.
- The primary data consists of weekly sales per store per department. Additionally, we have information on indicators like CPI, fuel prices, unemployment, temperature, holiday, etc.
- The entire dataset is split across five files, which is merged into a single dataset as a part of the data cleaning process.

```
[ ] 1 merged_data.show()
```

Dept	Weekly_Sales	Store	Date	Temperature	Fuel_Price	MarkDown1	MarkDown2	MarkDown3	MarkDown4	MarkDown5	CPI	Unemployment	IsHoliday	Type	Size	Bucket
1	24924.5	1	2010-02-05	42.31	2.572	NA	NA	NA	NA	NA	211.0963582	8.106	FALSE	A	151315	Large
1	46039.49	1	2010-02-12	38.51	2.548	NA	NA	NA	NA	NA	211.2421698	8.106	TRUE	A	151315	Large
1	41595.55	1	2010-02-19	39.93	2.514	NA	NA	NA	NA	NA	211.2891429	8.106	FALSE	A	151315	Large
1	19403.54	1	2010-02-26	46.63	2.561	NA	NA	NA	NA	NA	211.3196429	8.106	FALSE	A	151315	Large
1	21827.9	1	2010-03-05	46.5	2.625	NA	NA	NA	NA	NA	211.3501429	8.106	FALSE	A	151315	Large
1	21043.39	1	2010-03-12	57.79	2.667	NA	NA	NA	NA	NA	211.3806429	8.106	FALSE	A	151315	Large



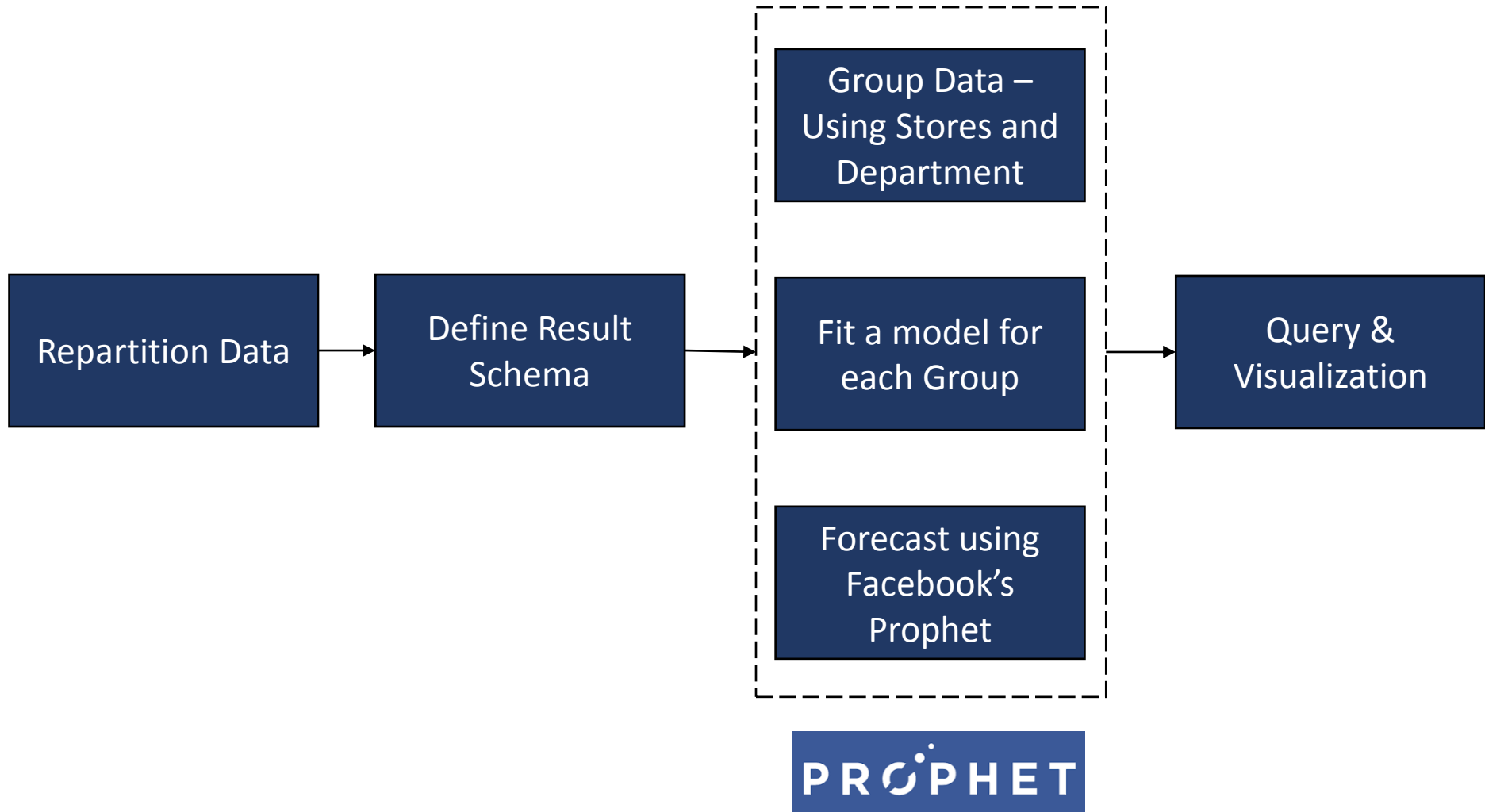




3

Methodology

Our Approach to Understand and Predict Sales



Multiple Time Series Forecasting



```
[ ] 1 train.explain()
```

```
== Physical Plan ==
```

```
FileScan csv [Store#244,Dept#245,Date#246,Weekly_Sales#247] Batched: false, DataFilters: [], Format: CSV, Location: InMemoryFileIndex(1 paths)
```



```
1 sql_statement = '''
2   SELECT
3     Store as store,
4     Dept as dept,
5     CAST(Date as date) as ds,
6     SUM(Weekly_Sales) as y
7   FROM train
8   GROUP BY store, dept, ds
9   ORDER BY store, dept, ds
10  '''
11
12 train = (
13   spark
14   .sql( sql_statement )
15   .repartition(sc.defaultParallelism, ['store', 'dept'])
16   ).cache()
```



```
[ ] 1 train.explain()
```

```
== Physical Plan ==
```

```
InMemoryTableScan [store#279, dept#280, ds#281, y#282]
```

```
+-- InMemoryRelation [store#279, dept#280, ds#281, y#282], StorageLevel(disk, memory, deserialized, 1 replicas)
```

```
+-- Exchange hashpartitioning(store#279, dept#280, 2), REPARTITION_BY_NUM, [id=#225]
```

```
+-- *(3) Sort [store#279 ASC NULLS FIRST, dept#280 ASC NULLS FIRST, ds#281 ASC NULLS FIRST], true, 0
```

```
+-- Exchange rangepartitioning(store#279 ASC NULLS FIRST, dept#280 ASC NULLS FIRST, ds#281 ASC NULLS FIRST), true, 0
```

```
+-- *(2) HashAggregate(keys=[store#244, dept#245, _groupingexpression#288], functions=[sum(Weekly_Sales)], _rowid_#289)
```

```
+-- Exchange hashpartitioning(store#244, dept#245, _groupingexpression#288, 200), ENSURE_REQUIREMENTS, [id=#226]
```

```
+-- *(1) HashAggregate(keys=[store#244, dept#245, _groupingexpression#288], functions=[partial_sum(Weekly_Sales)], _rowid_#290)
```

```
+-- *(1) Project [Store#244, Dept#245, Weekly_Sales#247, cast(Date#246 as date) AS _date_#291], false, [id=#227]
```

```
+-- FileScan csv [Store#244,Dept#245,Date#246,Weekly_Sales#247] Batched: false, DataFilters: [], Format: CSV, Location: InMemoryFileIndex(1 paths)
```

Multiple Time Series Forecasting



```
1 result_schema = StructType([
2   StructField('ds',DateType()),
3   StructField('store',IntegerType()),
4   StructField('dept',IntegerType()),
5   StructField('y',FloatType()),
6   StructField('yhat',FloatType()),
7   StructField('yhat_upper',FloatType()),
8   StructField('yhat_lower',FloatType())
9 ])
```



```
1 def forecast_store_dept( history_pd: pd.DataFrame ) -> pd.DataFrame:
2
3     model = Prophet(
4         interval_width=0.95,
5         daily_seasonality=False,
6         weekly_seasonality=True,
7         yearly_seasonality=True,
8         holidays = holidays
9     )
10
11     model.fit( history_pd )
12
13     future_pd = model.make_future_dataframe(
14         periods=52,
15         freq='w',
16         include_history=True
17     )
18
19     forecast_pd = model.predict( future_pd )
20
```



```
1 from pyspark.sql.functions import current_date
2
3 results = (
4     train
5     .groupBy('store', 'dept')
6     .applyInPandas(forecast_store_dept, schema=result_schema)
7     .withColumn('training_date', current_date() )
8 )
9
10 results.createOrReplaceTempView('forecasts')
```



```
1 results.coalesce(1)

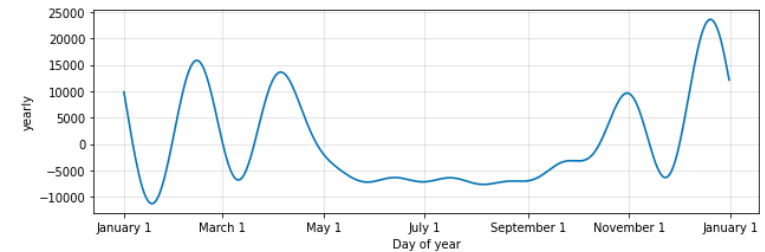
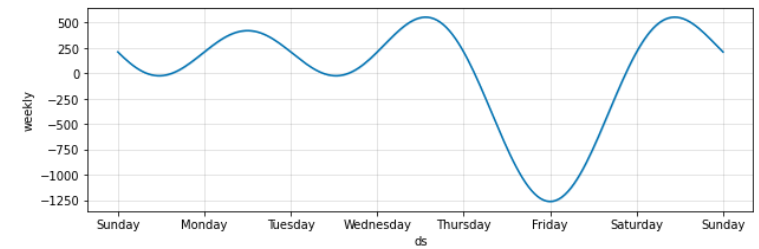
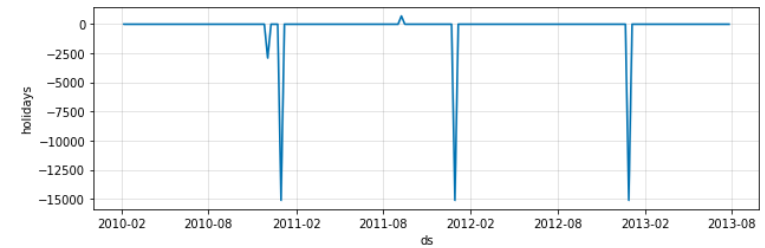
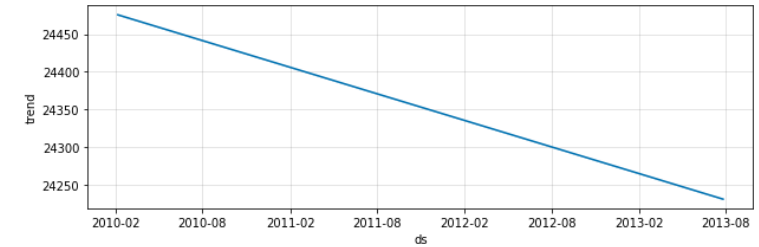
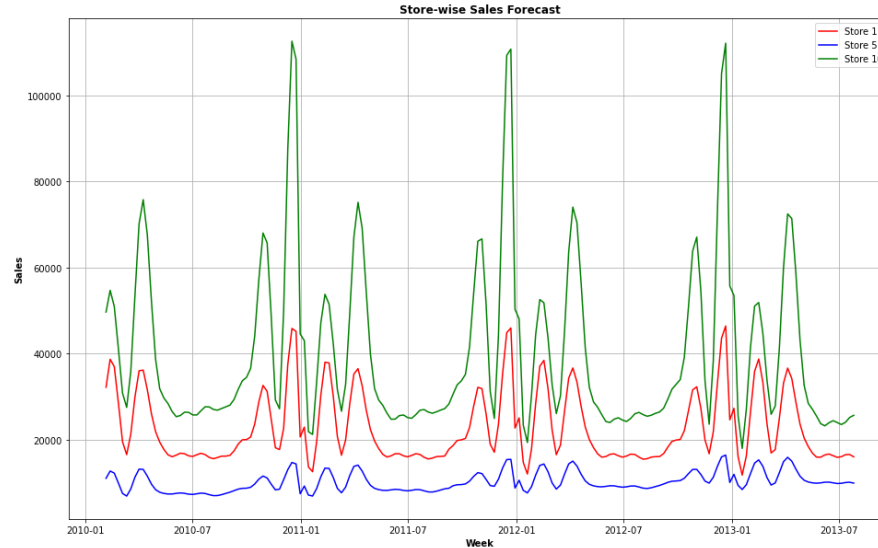
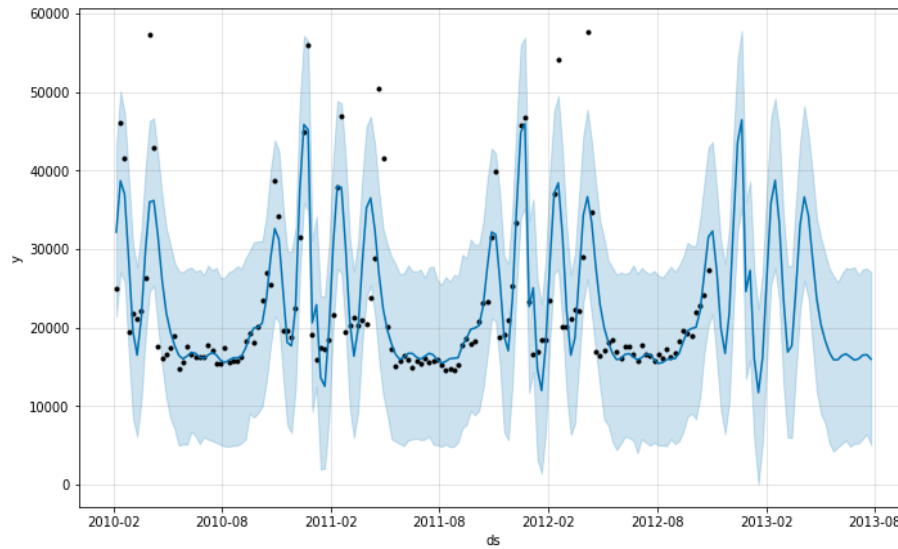
DataFrame[ds: date, store: int, dept: int, y: float, yhat: float, yhat_upper: float, yhat_lower: float]
```

```
1 test_sql = '''
2 SELECT
3   ds,
4   y,
5   yhat,
6   yhat_upper,
7   yhat_lower
8 FROM forecasts
9 WHERE store = 1 AND dept = 1
10 '''
11
12 test_result = spark.sql(test_sql)
```

4

Results

Insights derived from the analysis



Thank you

Made by – Anirudha Balkrishna, Sanjay Madesha & Shubham Khode