# Rossmann Store Sales Prediction using XGBoost

## Project Overview

This notebook presents a machine learning approach for forecasting daily sales of a retail store chain. The goal is to predict sales using historical data along with various store-level and promotional features, enabling data-driven decisions in inventory management and marketing strategy.

## Objectives

- Clean and preprocess raw sales and store datasets.
- Engineer relevant features (temporal, promotional, holiday-based).
- Train and evaluate a regression model (XGBoost) for accurate sales prediction.
- Interpret feature importance and assess model performance.

## Dataset Description

- **Source:** Rossmann Store Sales dataset (public Kaggle competition)
- **Main files:**
  - `train.csv` : Historical daily sales data
  - `test.csv` : Test set for prediction
  - `store.csv` : Store-level metadata
- **Key features include:**
  - `Store` , `DayOfWeek` , `Date` , `Sales` , `Customers` , `Open` , `Promo` , `StateHoliday` , `SchoolHoliday` , `StoreType` , `Assortment` , `CompetitionDistance` , `CompetitionOpenSince` , `Promo2` , etc.

## Workflow

1. Data loading and exploration
2. Data cleaning and preprocessing
3. Feature engineering
4. Model training using XGBoost
5. Model evaluation and error analysis

6. Interpretation of results and conclusion

# Dependencies

- `pandas`, `numpy`, `matplotlib`, `seaborn`, `scikit-learn`, `xgboost`, `joblib`

```
In [1]:  # Install necessary libraries
         !pip install pandas matplotlib seaborn scikit-learn xgboost joblib
```

```
Requirement already satisfied: pandas in /usr/local/lib/python3.11/dist-package
s (2.2.2)
Requirement already satisfied: matplotlib in /usr/local/lib/python3.11/dist-pac
kages (3.10.0)
Requirement already satisfied: seaborn in /usr/local/lib/python3.11/dist-packag
es (0.13.2)
Requirement already satisfied: scikit-learn in /usr/local/lib/python3.11/dist-p
ackages (1.6.1)
Requirement already satisfied: xgboost in /usr/local/lib/python3.11/dist-packag
es (2.1.4)
Requirement already satisfied: joblib in /usr/local/lib/python3.11/dist-package
s (1.5.1)
Requirement already satisfied: numpy>=1.23.2 in /usr/local/lib/python3.11/dist-
packages (from pandas) (2.0.2)
Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python
3.11/dist-packages (from pandas) (2.9.0.post0)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.11/dist-p
ackages (from pandas) (2025.2)
Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.11/dis
t-packages (from pandas) (2025.2)
Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.11/di
st-packages (from matplotlib) (1.3.2)
Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.11/dist-p
ackages (from matplotlib) (0.12.1)
Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.11/d
ist-packages (from matplotlib) (4.58.4)
Requirement already satisfied: kiwisolver>=1.3.1 in /usr/local/lib/python3.11/d
ist-packages (from matplotlib) (1.4.8)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.11/dis
t-packages (from matplotlib) (24.2)
Requirement already satisfied: pillow>=8 in /usr/local/lib/python3.11/dist-pack
ages (from matplotlib) (11.2.1)
Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.11/di
st-packages (from matplotlib) (3.2.3)
Requirement already satisfied: scipy>=1.6.0 in /usr/local/lib/python3.11/dist-p
ackages (from scikit-learn) (1.15.3)
Requirement already satisfied: threadpoolctl>=3.1.0 in /usr/local/lib/python3.1
1/dist-packages (from scikit-learn) (3.6.0)
Requirement already satisfied: nvidia-nccl-cu12 in /usr/local/lib/python3.11/di
st-packages (from xgboost) (2.21.5)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.11/dist-packa
ges (from python-dateutil>=2.8.2->pandas) (1.17.0)
```

In [2]:
```python
# Import libraries
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np
from sklearn.preprocessing import MinMaxScaler
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestRegressor
from xgboost import XGBRegressor
from sklearn.metrics import mean_squared_error
from sklearn.model_selection import GridSearchCV
```

```
from xgboost import plot_tree
import joblib

# Ignore warnings
import warnings
warnings.filterwarnings('ignore')
```

In [3]:
```
# Download and extract the latest Rossmann Store Sales dataset from Kaggle
import kagglehub

path = kagglehub.dataset_download("pratyushakar/rossmann-store-sales")
print("Path to dataset files:", path)
```

Downloading from https://www.kaggle.com/api/v1/datasets/download/pratyushakar/rossmann-store-sales?dataset_version_number=3...
100%|████████████| 6.90M/6.90M [00:00<00:00, 97.2MB/s]
Extracting files...
Path to dataset files: /root/.cache/kagglehub/datasets/pratyushakar/rossmann-store-sales/versions/3

In [4]:
```
store_df = pd.read_csv(path + "/store.csv")
train_df = pd.read_csv(path + "/train.csv")
```

# Data Cleaning

In [5]:
```
# Inspect top rows and schema of store metadata to understand available featur
store_df.head()
```

Out[5]:

| | Store | StoreType | Assortment | CompetitionDistance | CompetitionOpenSinceM |
|---|---|---|---|---|---|
| 0 | 1 | c | a | 1270.0 | |
| 1 | 2 | a | a | 570.0 | |
| 2 | 3 | a | a | 14130.0 | |
| 3 | 4 | c | c | 620.0 | |
| 4 | 5 | a | a | 29910.0 | |

In [6]:
```
store_df.describe()
```

| | Store | CompetitionDistance | CompetitionOpenSinceMonth | Competit |
|---|---|---|---|---|
| **count** | 1115.00000 | 1112.000000 | 761.000000 | |
| **mean** | 558.00000 | 5404.901079 | 7.224704 | |
| **std** | 322.01708 | 7663.174720 | 3.212348 | |
| **min** | 1.00000 | 20.000000 | 1.000000 | |
| **25%** | 279.50000 | 717.500000 | 4.000000 | |
| **50%** | 558.00000 | 2325.000000 | 8.000000 | |
| **75%** | 836.50000 | 6882.500000 | 10.000000 | |
| **max** | 1115.00000 | 75860.000000 | 12.000000 | |

In [7]:
```
print("Size of store dataset :", store_df.shape)
print("\n Displaying summary information about the store dataset structure, in
store_df.info()
print("\n Number of missing values in each coloumn:")
store_df.isnull().sum()
```

```
Size of store dataset : (1115, 10)

 Displaying summary information about the store dataset structure, including co
lumn data types and missing values:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1115 entries, 0 to 1114
Data columns (total 10 columns):
 #   Column                     Non-Null Count  Dtype
---  ------                     --------------  -----
 0   Store                      1115 non-null   int64
 1   StoreType                  1115 non-null   object
 2   Assortment                 1115 non-null   object
 3   CompetitionDistance        1112 non-null   float64
 4   CompetitionOpenSinceMonth  761 non-null    float64
 5   CompetitionOpenSinceYear   761 non-null    float64
 6   Promo2                     1115 non-null   int64
 7   Promo2SinceWeek            571 non-null    float64
 8   Promo2SinceYear            571 non-null    float64
 9   PromoInterval              571 non-null    object
dtypes: float64(5), int64(2), object(3)
memory usage: 87.2+ KB

 Number of missing values in each coloumn:
```

Out[7]:

| | **0** |
|---|---|
| **Store** | 0 |
| **StoreType** | 0 |
| **Assortment** | 0 |
| **CompetitionDistance** | 3 |
| **CompetitionOpenSinceMonth** | 354 |
| **CompetitionOpenSinceYear** | 354 |
| **Promo2** | 0 |
| **Promo2SinceWeek** | 544 |
| **Promo2SinceYear** | 544 |
| **PromoInterval** | 544 |

**dtype:** int64

In [8]:
```python
# Load historical sales data and check for missingness or anomalies
train_df.head()
```

Out[8]:

| | **Store** | **DayOfWeek** | **Date** | **Sales** | **Customers** | **Open** | **Promo** | **StateHoliday** |
|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 5 | 2015-07-31 | 5263 | 555 | 1 | 1 | ( |
| **1** | 2 | 5 | 2015-07-31 | 6064 | 625 | 1 | 1 | ( |
| **2** | 3 | 5 | 2015-07-31 | 8314 | 821 | 1 | 1 | ( |
| **3** | 4 | 5 | 2015-07-31 | 13995 | 1498 | 1 | 1 | ( |
| **4** | 5 | 5 | 2015-07-31 | 4822 | 559 | 1 | 1 | ( |

In [9]:
```python
train_df.describe()
```

Out[9]:

| | Store | DayOfWeek | Sales | Customers | Open |
|---|---|---|---|---|---|
| **count** | 1.017209e+06 | 1.017209e+06 | 1.017209e+06 | 1.017209e+06 | 1.017209e+06 |
| **mean** | 5.584297e+02 | 3.998341e+00 | 5.773819e+03 | 6.331459e+02 | 8.301067e-01 |
| **std** | 3.219087e+02 | 1.997391e+00 | 3.849926e+03 | 4.644117e+02 | 3.755392e-01 |
| **min** | 1.000000e+00 | 1.000000e+00 | 0.000000e+00 | 0.000000e+00 | 0.000000e+00 |
| **25%** | 2.800000e+02 | 2.000000e+00 | 3.727000e+03 | 4.050000e+02 | 1.000000e+00 |
| **50%** | 5.580000e+02 | 4.000000e+00 | 5.744000e+03 | 6.090000e+02 | 1.000000e+00 |
| **75%** | 8.380000e+02 | 6.000000e+00 | 7.856000e+03 | 8.370000e+02 | 1.000000e+00 |
| **max** | 1.115000e+03 | 7.000000e+00 | 4.155100e+04 | 7.388000e+03 | 1.000000e+00 |

In [10]:
```
print("Size of train dataset :", train_df.shape)
print("\n Displaying summary information about the train dataset structure, in
train_df.info()
print("\n Number of missing values in each coloumn:")
train_df.isnull().sum()
```

Size of train dataset : (1017209, 9)

 Displaying summary information about the train dataset structure, including co
lumn data types and missing values:
```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1017209 entries, 0 to 1017208
Data columns (total 9 columns):
 #   Column        Non-Null Count    Dtype
---  ------        --------------    -----
 0   Store         1017209 non-null  int64
 1   DayOfWeek     1017209 non-null  int64
 2   Date          1017209 non-null  object
 3   Sales         1017209 non-null  int64
 4   Customers     1017209 non-null  int64
 5   Open          1017209 non-null  int64
 6   Promo         1017209 non-null  int64
 7   StateHoliday  1017209 non-null  object
 8   SchoolHoliday 1017209 non-null  int64
dtypes: int64(7), object(2)
memory usage: 69.8+ MB
```

 Number of missing values in each coloumn:

Out[10]:

| | 0 |
|---|---|
| **Store** | 0 |
| **DayOfWeek** | 0 |
| **Date** | 0 |
| **Sales** | 0 |
| **Customers** | 0 |
| **Open** | 0 |
| **Promo** | 0 |
| **StateHoliday** | 0 |
| **SchoolHoliday** | 0 |

**dtype:** int64

```
In [11]:  # Combine sales and store metadata for feature engineering
          df = pd.merge(train_df, store_df, on='Store', how='left')
```

```
In [12]:  df.sample(5)
```

Out[12]:

| | Store | DayOfWeek | Date | Sales | Customers | Open | Promo | StateH |
|---|---|---|---|---|---|---|---|---|
| **818475** | 851 | 5 | 2013-06-28 | 4821 | 453 | 1 | 0 | |
| **510239** | 355 | 1 | 2014-03-31 | 13351 | 1445 | 1 | 1 | |
| **267090** | 948 | 6 | 2014-11-29 | 6918 | 1539 | 1 | 0 | |
| **204605** | 561 | 4 | 2015-01-29 | 6613 | 637 | 1 | 1 | |
| **145369** | 420 | 1 | 2015-03-23 | 4042 | 394 | 1 | 0 | |

```
In [13]:  df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1017209 entries, 0 to 1017208
Data columns (total 18 columns):
 #   Column                     Non-Null Count    Dtype
---  ------                     --------------    -----
 0   Store                      1017209 non-null  int64
 1   DayOfWeek                  1017209 non-null  int64
 2   Date                       1017209 non-null  object
 3   Sales                      1017209 non-null  int64
 4   Customers                  1017209 non-null  int64
 5   Open                       1017209 non-null  int64
 6   Promo                      1017209 non-null  int64
 7   StateHoliday               1017209 non-null  object
 8   SchoolHoliday              1017209 non-null  int64
 9   StoreType                  1017209 non-null  object
 10  Assortment                 1017209 non-null  object
 11  CompetitionDistance        1014567 non-null  float64
 12  CompetitionOpenSinceMonth  693861 non-null   float64
 13  CompetitionOpenSinceYear   693861 non-null   float64
 14  Promo2                     1017209 non-null  int64
 15  Promo2SinceWeek            509178 non-null   float64
 16  Promo2SinceYear            509178 non-null   float64
 17  PromoInterval              509178 non-null   object
dtypes: float64(5), int64(8), object(5)
memory usage: 139.7+ MB
```

In [14]: `df.isnull().sum()`

| | **0** |
|---|---|
| **Store** | 0 |
| **DayOfWeek** | 0 |
| **Date** | 0 |
| **Sales** | 0 |
| **Customers** | 0 |
| **Open** | 0 |
| **Promo** | 0 |
| **StateHoliday** | 0 |
| **SchoolHoliday** | 0 |
| **StoreType** | 0 |
| **Assortment** | 0 |
| **CompetitionDistance** | 2642 |
| **CompetitionOpenSinceMonth** | 323348 |
| **CompetitionOpenSinceYear** | 323348 |
| **Promo2** | 0 |
| **Promo2SinceWeek** | 508031 |
| **Promo2SinceYear** | 508031 |
| **PromoInterval** | 508031 |

**dtype:** int64

```python
df.describe()
```

| | **Store** | **DayOfWeek** | **Sales** | **Customers** | **Open** |
|---|---|---|---|---|---|
| **count** | 1.017209e+06 | 1.017209e+06 | 1.017209e+06 | 1.017209e+06 | 1.017209e+06 |
| **mean** | 5.584297e+02 | 3.998341e+00 | 5.773819e+03 | 6.331459e+02 | 8.301067e-01 |
| **std** | 3.219087e+02 | 1.997391e+00 | 3.849926e+03 | 4.644117e+02 | 3.755392e-01 |
| **min** | 1.000000e+00 | 1.000000e+00 | 0.000000e+00 | 0.000000e+00 | 0.000000e+00 |
| **25%** | 2.800000e+02 | 2.000000e+00 | 3.727000e+03 | 4.050000e+02 | 1.000000e+00 |
| **50%** | 5.580000e+02 | 4.000000e+00 | 5.744000e+03 | 6.090000e+02 | 1.000000e+00 |
| **75%** | 8.380000e+02 | 6.000000e+00 | 7.856000e+03 | 8.370000e+02 | 1.000000e+00 |
| **max** | 1.115000e+03 | 7.000000e+00 | 4.155100e+04 | 7.388000e+03 | 1.000000e+00 |

```
In [16]:  # Select rows where both StateHoliday and SchoolHoliday are flagged but the st
          df[(df['StateHoliday'] == 1) & (df['SchoolHoliday'] == 1) & (df['Sales'] > 0)]
```

Out[16]:

| Store | DayOfWeek | Date | Sales | Customers | Open | Promo | StateHoliday | Scho |
|-------|-----------|------|-------|-----------|------|-------|--------------|------|

- Zero-sales days often correspond to store closures, data errors, or holidays where the store wasn't open.

- Keeping only positive-sales records ensures the model learns from genuine "open-store" demand patterns, rather than being skewed by the artificial zeros.

```
In [17]:  # Remove all days with zero (or no) sales — e.g. when the store was closed
          df = df[df['Sales'] > 0]
          print(df.shape)
          df.sample(5)
```

```
(844338, 18)
```

Out[17]:

| | Store | DayOfWeek | Date | Sales | Customers | Open | Promo | StateH |
|---|-------|-----------|------|-------|-----------|------|-------|--------|
| **43796** | 312 | 1 | 2015-06-22 | 4260 | 466 | 1 | 0 | |
| **90996** | 682 | 1 | 2015-05-11 | 12100 | 1931 | 1 | 0 | |
| **404871** | 236 | 5 | 2014-07-04 | 6796 | 814 | 1 | 1 | |
| **803788** | 659 | 4 | 2013-07-11 | 4981 | 671 | 1 | 0 | |
| **165278** | 259 | 4 | 2015-03-05 | 11622 | 2491 | 1 | 1 | |

```
In [18]:  # Extract day, month, year components from Date and drop original Date column
          df['Date'] = pd.to_datetime(df['Date'])
          df['Day'] = df['Date'].dt.day
          df['Month'] = df['Date'].dt.month
          df['Year'] = df['Date'].dt.year
          df.drop(columns=['Date'], inplace=True)
          df.sample(5)
```

Out[18]:

| | Store | DayOfWeek | Sales | Customers | Open | Promo | StateHoliday | Sch |
|---|-------|-----------|-------|-----------|------|-------|--------------|-----|
| **497942** | 323 | 5 | 6653 | 635 | 1 | 0 | 0 | |
| **593025** | 631 | 4 | 4601 | 598 | 1 | 0 | 0 | |
| **796238** | 914 | 4 | 10115 | 1176 | 1 | 1 | 0 | |
| **351462** | 92 | 6 | 6406 | 572 | 1 | 0 | 0 | |
| **170870** | 276 | 6 | 2461 | 288 | 1 | 0 | 0 | |

```
In [19]:  # Drop Customers — not available at prediction time
          df.drop(columns=['Customers'], inplace=True)
          df.sample(5)
```

Out[19]:

| | Store | DayOfWeek | Sales | Open | Promo | StateHoliday | SchoolHoliday |
|---|---|---|---|---|---|---|---|
| **968705** | 556 | 3 | 5064 | 1 | 0 | 0 | 0 |
| **1663** | 549 | 4 | 4629 | 1 | 1 | 0 | 1 |
| **387261** | 411 | 3 | 7281 | 1 | 0 | 0 | 0 |
| **135256** | 342 | 3 | 14177 | 1 | 1 | 0 | 1 |
| **856569** | 1035 | 6 | 1722 | 1 | 0 | 0 | 0 |

```
In [20]:  # Unique values for stateholiday and schoolholiday
          print(f"Unique values for stateholiday: {df['StateHoliday'].unique()}")
          print(f"Unique values for schoolholiday: {df['SchoolHoliday'].unique()}")
```

```
Unique values for stateholiday: ['0' 'a' 'b' 'c' 0]
Unique values for schoolholiday: [1 0]
```

```
In [21]:  # Normalize StateHoliday and SchoolHoliday to integer flags
          df['StateHoliday'] = df['StateHoliday'].apply(lambda x: int(x) if type(x) == s
          df['SchoolHoliday'] = df['SchoolHoliday'].apply(lambda x: int(x) if type(x) ==
```

```
In [22]:  # Check how sales affects with state holiday and school holiday
          print("Average sales when state holiday is true:", df[df['StateHoliday'] == 1]
          print("Average sales when state holiday is false:", df[df['StateHoliday'] == 0
          print("\nAverage sales when school holiday is true:", df[df['SchoolHoliday'] =
          print("Average sales when school holiday is false:", df[df['SchoolHoliday'] ==

          # Plot boxplots to visualize the distribution
          plt.figure(figsize=(12,5))

          plt.subplot(1,2,1)
          sns.boxplot(x='StateHoliday', y='Sales', data=df)
          plt.title('Sales Distribution by State Holiday')

          plt.subplot(1,2,2)
          sns.boxplot(x='SchoolHoliday', y='Sales', data=df)
          plt.title('Sales Distribution by School Holiday')

          plt.tight_layout()
          plt.show()
```
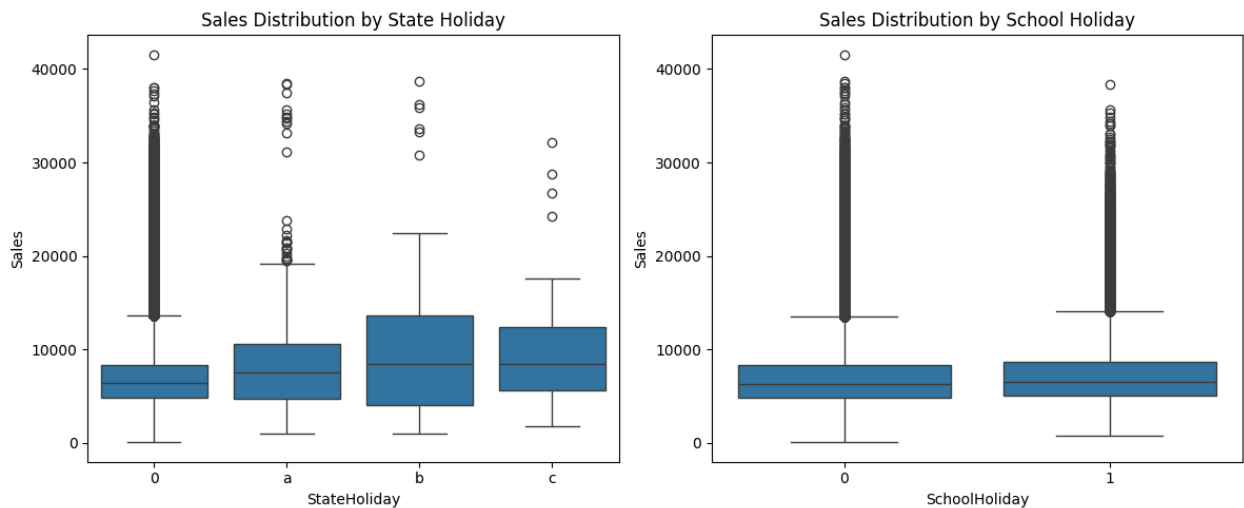
```
Average sales when state holiday is true: nan
Average sales when state holiday is false: 6953.960228970345

Average sales when school holiday is true: 7200.710281746153
Average sales when school holiday is false: 6897.20783001147
```

Sales Distribution by State Holiday — Sales Distribution by School Holiday

State holidays show no sales (stores closed), whereas school holidays boost average sales by ~4.4% (7,200.7 vs. 6,897.2). Because each flag captures a distinct effect, merging them would obscure valuable information.

In [23]:
```python
df.sample(5)
```

Out[23]:

| | Store | DayOfWeek | Sales | Open | Promo | StateHoliday | SchoolHoliday | |
|---|---|---|---|---|---|---|---|---|
| 607526 | 637 | 5 | 7768 | 1 | 0 | 0 | 1 | |
| 938142 | 98 | 2 | 3348 | 1 | 0 | 0 | 0 | |
| 519935 | 16 | 6 | 4873 | 1 | 0 | 0 | 0 | |
| 673671 | 997 | 2 | 6346 | 1 | 1 | 0 | 0 | |
| 577476 | 692 | 4 | 5453 | 1 | 0 | 0 | 0 | |

In [24]:
```python
# store type and assortment unique values
print("Store types:", df['StoreType'].unique())
print("Assortment types:", df['Assortment'].unique())
```

```
Store types: ['c' 'a' 'd' 'b']
Assortment types: ['a' 'c' 'b']
```

CompetitionDistance is null means competition doesn't exist and if we replace it with 0, it will mean that competition is just near the store, so we replace it with the 2 * max_value of the featuure

In [25]:
```python
df['CompetitionDistance'].isna().sum()
```

Out[25]: np.int64(2186)

In [26]:
```python
# Impute missing CompetitionDistance as twice max distance → treat as "no comp
df['CompetitionDistance'].fillna(2*df['CompetitionDistance'].max(), inplace=Tr
df['CompetitionDistance'].isna().sum()
```

Out[26]: np.int64(0)

In [27]: df.sample(5)

Out[27]:

| | Store | DayOfWeek | Sales | Open | Promo | StateHoliday | SchoolHoliday | |
|---|---|---|---|---|---|---|---|---|
| 267828 | 717 | 5 | 8515 | 1 | 1 | 0 | 0 | |
| 978066 | 997 | 2 | 6189 | 1 | 1 | 0 | 0 | |
| 675673 | 769 | 7 | 9455 | 1 | 0 | 0 | 0 | |
| 865058 | 604 | 5 | 9244 | 1 | 1 | 0 | 0 | |
| 148151 | 972 | 6 | 1394 | 1 | 0 | 0 | 0 | |

Replace CompetitionOpenSinceMonth and CompetitionOpenSinceYear with one single column,

Calculate the months from CompetitionOpenSinceMonth-CompetitionOpenSinceYear to Day-Month-Year

In [28]:
```python
df['CompetitionOpenSinceYear'] = df['CompetitionOpenSinceYear'].fillna(0).asty
df['CompetitionOpenSinceMonth'] = df['CompetitionOpenSinceMonth'].fillna(0).as
```

In [29]:
```python
df['CompetitionOpenSinceMonth'].unique()
```

Out[29]: array([ 9, 11, 12,  4, 10,  8,  0,  3,  6,  5,  1,  2,  7])

In [30]:
```python
df['CompetitionOpenSinceYear'].unique()
```

Out[30]: array([2008, 2007, 2006, 2009, 2015, 2013, 2014, 2000, 2011,    0, 2010,
       2005, 1999, 2003, 2012, 2004, 2002, 1961, 1995, 2001, 1990, 1994,
       1900, 1998])

In [31]:
```python
# Create a new column for competition months
df['CompetitionMonths'] = 0

# Only calculate for rows where competition exists (year and month > 0)
mask = (df['CompetitionOpenSinceYear'] > 0) & (df['CompetitionOpenSinceMonth']

# Calculate months between competition open date and store date
df.loc[mask, 'CompetitionMonths'] = 12 * (df.loc[mask, 'Year'] - df.loc[mask,

# Drop the original columns
df.drop(['CompetitionOpenSinceMonth', 'CompetitionOpenSinceYear'], axis=1, inp
```

In [32]: df.head(5)

Out[32]:

| | Store | DayOfWeek | Sales | Open | Promo | StateHoliday | SchoolHoliday | StoreT |
|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 5 | 5263 | 1 | 1 | 0 | 1 | |
| **1** | 2 | 5 | 6064 | 1 | 1 | 0 | 1 | |
| **2** | 3 | 5 | 8314 | 1 | 1 | 0 | 1 | |
| **3** | 4 | 5 | 13995 | 1 | 1 | 0 | 1 | |
| **4** | 5 | 5 | 4822 | 1 | 1 | 0 | 1 | |

Some Values are in negative in the 'CompetitionMonths' columns, means the competition is later opened compared to the sales date, so at the sales date, there were no competition

Replace CompetitionMonths with 0 for negative values

```python
In [33]: df.loc[df['CompetitionMonths'] < 0, 'CompetitionMonths'] = 0
```

```python
In [34]: # If promo 2 is 0, fill 0 in Promo2SinceWeek and Promo2SinceYear
         df.loc[df['Promo2'] == 0, ['Promo2SinceWeek','Promo2SinceYear']] = 0
```

```python
In [35]: df.sample(6)
```

Out[35]:

| | Store | DayOfWeek | Sales | Open | Promo | StateHoliday | SchoolHoliday | S |
|---|---|---|---|---|---|---|---|---|
| **618476** | 437 | 2 | 3456 | 1 | 0 | 0 | 1 | |
| **485687** | 333 | 2 | 8263 | 1 | 0 | 0 | 0 | |
| **415031** | 1037 | 3 | 5117 | 1 | 0 | 0 | 0 | |
| **536894** | 250 | 5 | 9278 | 1 | 1 | 0 | 0 | |
| **528600** | 876 | 6 | 3022 | 1 | 0 | 0 | 0 | |
| **900736** | 602 | 1 | 4988 | 1 | 0 | 0 | 0 | |

Merge Promo2SinceWeek and Promo2SinceYear into one column like CompetitionMonths

```python
In [36]: # Create a new column for promo2 months
         df['Promo2Months'] = 0

         # Only calculate for rows where promo2 exists (year and week > 0)
         mask = (df['Promo2SinceYear'] > 0) & (df['Promo2SinceWeek'] > 0)

         # Calculate months between promo2 start date and store date
         df.loc[mask, 'Promo2Months'] = 12 * (df.loc[mask, 'Year'] - df.loc[mask, 'Prom

         # Replace negative values with 0
```

```
df.loc[df['Promo2Months'] < 0, 'Promo2Months'] = 0

# Drop the original columns
df.drop(['Promo2SinceWeek', 'Promo2SinceYear'], axis=1, inplace=True)
```

In [37]: `df.sample(5)`

Out[37]:

| | Store | DayOfWeek | Sales | Open | Promo | StateHoliday | SchoolHoliday | |
|---|---|---|---|---|---|---|---|---|
| **522094** | 1060 | 5 | 7373 | 1 | 1 | 0 | 0 | |
| **16431** | 822 | 5 | 7861 | 1 | 1 | 0 | 0 | |
| **788872** | 238 | 3 | 6129 | 1 | 0 | 0 | 1 | |
| **278203** | 822 | 1 | 5965 | 1 | 0 | 0 | 0 | |
| **962430** | 971 | 2 | 8042 | 1 | 1 | 0 | 1 | |

Instead of having PromotInterval column, have a column to indicate if the Promo happened in the same month as the sale was recorded

In [38]:
```
# Create new column for promo month match
df['PromoInMonth'] = 0

# Create month mapping dictionary
month_map = {1:'Jan', 2:'Feb', 3:'Mar', 4:'Apr', 5:'May', 6:'Jun',
             7:'Jul', 8:'Aug', 9:'Sep', 10:'Oct', 11:'Nov', 12:'Dec'}

# Convert month numbers to month names
df['MonthName'] = df['Month'].map(month_map)

# For each row, check if month name exists in PromoInterval
mask = ~df['PromoInterval'].isna()
df.loc[mask, 'PromoInMonth'] = df.loc[mask].apply(
    lambda x: 1 if x['MonthName'] in str(x['PromoInterval']).split(',') else 0
    axis=1
)

# Drop temporary and original columns
df.drop(['MonthName', 'PromoInterval'], axis=1, inplace=True)
```

In [39]: `df.sample(5)`

| | Store | DayOfWeek | Sales | Open | Promo | StateHoliday | SchoolHoliday | |
|---|---|---|---|---|---|---|---|---|
| **93601** | 1057 | 6 | 8192 | 1 | 0 | 0 | 0 | |
| **70543** | 299 | 5 | 7206 | 1 | 0 | 0 | 0 | |
| **325643** | 524 | 6 | 6694 | 1 | 0 | 0 | 0 | |
| **957529** | 530 | 6 | 2127 | 1 | 0 | 0 | 0 | |
| **184101** | 127 | 1 | 7920 | 1 | 1 | 0 | 1 | |

In [40]:
```python
# Rearrange columns with logical grouping of features and Sales as last column
cols = [
    # Store characteristics
    'Store', 'StoreType', 'Assortment',

    # Time features
    'DayOfWeek', 'Day', 'Month',

    # Competition features
    'CompetitionDistance', 'CompetitionMonths',

    # Promotion features
    'Promo', 'Promo2', 'Promo2Months', 'PromoInMonth',

    # Status indicators
    'StateHoliday', 'SchoolHoliday',

    # Target variable
    'Sales'
]

df = df[cols]
```

In [41]:
```python
df.sample(6)
```

Out[41]:

| | Store | StoreType | Assortment | DayOfWeek | Day | Month | CompetitionDi |
|---|---|---|---|---|---|---|---|
| **250056** | 705 | a | a | 3 | 17 | 12 | |
| **513056** | 942 | d | c | 6 | 29 | 3 | |
| **688068** | 899 | d | a | 3 | 23 | 10 | |
| **343103** | 166 | a | c | 1 | 8 | 9 | |
| **892742** | 413 | a | c | 1 | 22 | 4 | |
| **652615** | 11 | a | c | 6 | 23 | 11 | |

In [42]:
```python
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Index: 844338 entries, 0 to 1017190
Data columns (total 15 columns):
 #   Column               Non-Null Count   Dtype
---  ------               --------------   -----
 0   Store                844338 non-null  int64
 1   StoreType            844338 non-null  object
 2   Assortment           844338 non-null  object
 3   DayOfWeek            844338 non-null  int64
 4   Day                  844338 non-null  int32
 5   Month                844338 non-null  int32
 6   CompetitionDistance  844338 non-null  float64
 7   CompetitionMonths    844338 non-null  int64
 8   Promo                844338 non-null  int64
 9   Promo2               844338 non-null  int64
 10  Promo2Months         844338 non-null  int64
 11  PromoInMonth         844338 non-null  int64
 12  StateHoliday         844338 non-null  object
 13  SchoolHoliday        844338 non-null  int64
 14  Sales                844338 non-null  int64
dtypes: float64(1), int32(2), int64(9), object(3)
memory usage: 96.6+ MB
```

In [43]: `df.describe()`

Out[43]:

| | Store | DayOfWeek | Day | Month | Competition |
|---|---|---|---|---|---|
| count | 844338.000000 | 844338.000000 | 844338.000000 | 844338.000000 | 84433 |
| mean | 558.421374 | 3.520350 | 15.835706 | 5.845774 | 583 |
| std | 321.730861 | 1.723712 | 8.683392 | 3.323959 | 1077 |
| min | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 2 |
| 25% | 280.000000 | 2.000000 | 8.000000 | 3.000000 | 71 |
| 50% | 558.000000 | 3.000000 | 16.000000 | 6.000000 | 233 |
| 75% | 837.000000 | 5.000000 | 23.000000 | 8.000000 | 691 |
| max | 1115.000000 | 7.000000 | 31.000000 | 12.000000 | 15172 |

# Data Visualisation

In [44]:
```python
plt.figure(figsize=(10,6))
sales_by_day = df.groupby('DayOfWeek')['Sales'].sum()
sales_pct = (sales_by_day / sales_by_day.sum()) * 100

# Create bars and add percentage labels
bars = plt.bar(range(7), sales_pct)
for bar in bars:
    height = bar.get_height()
```
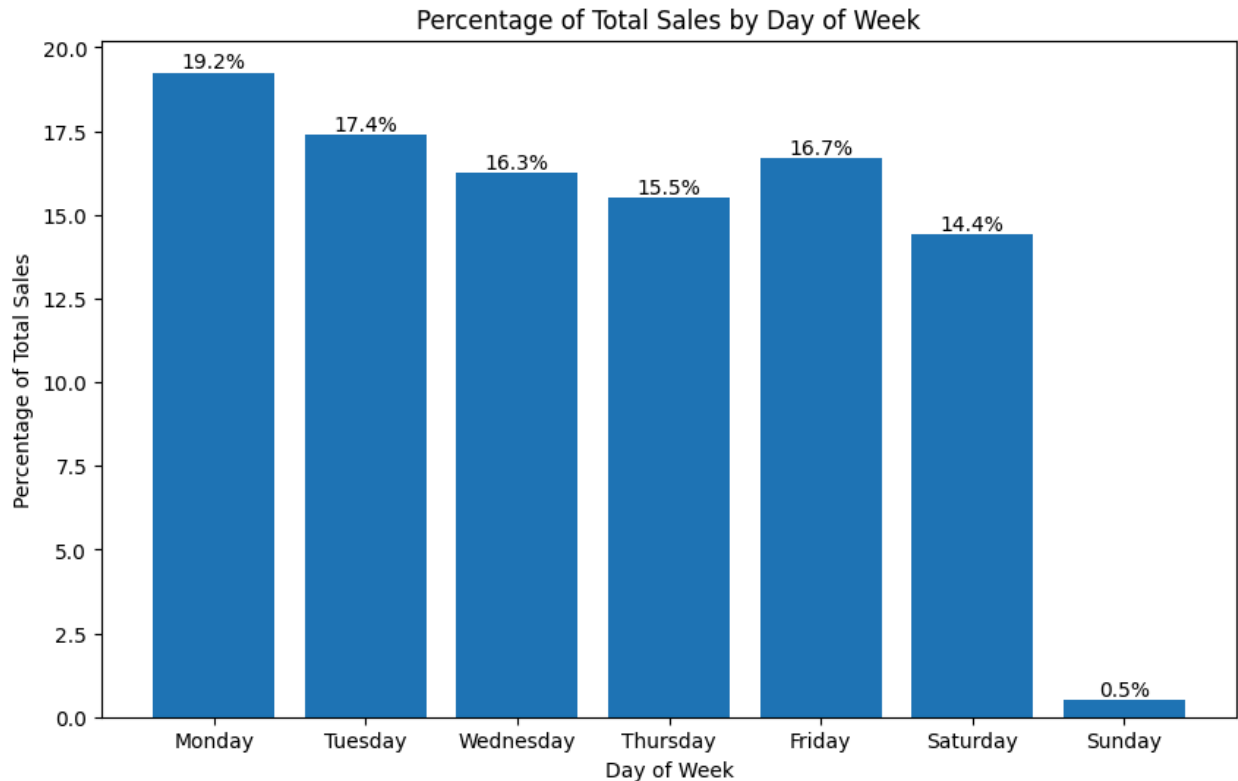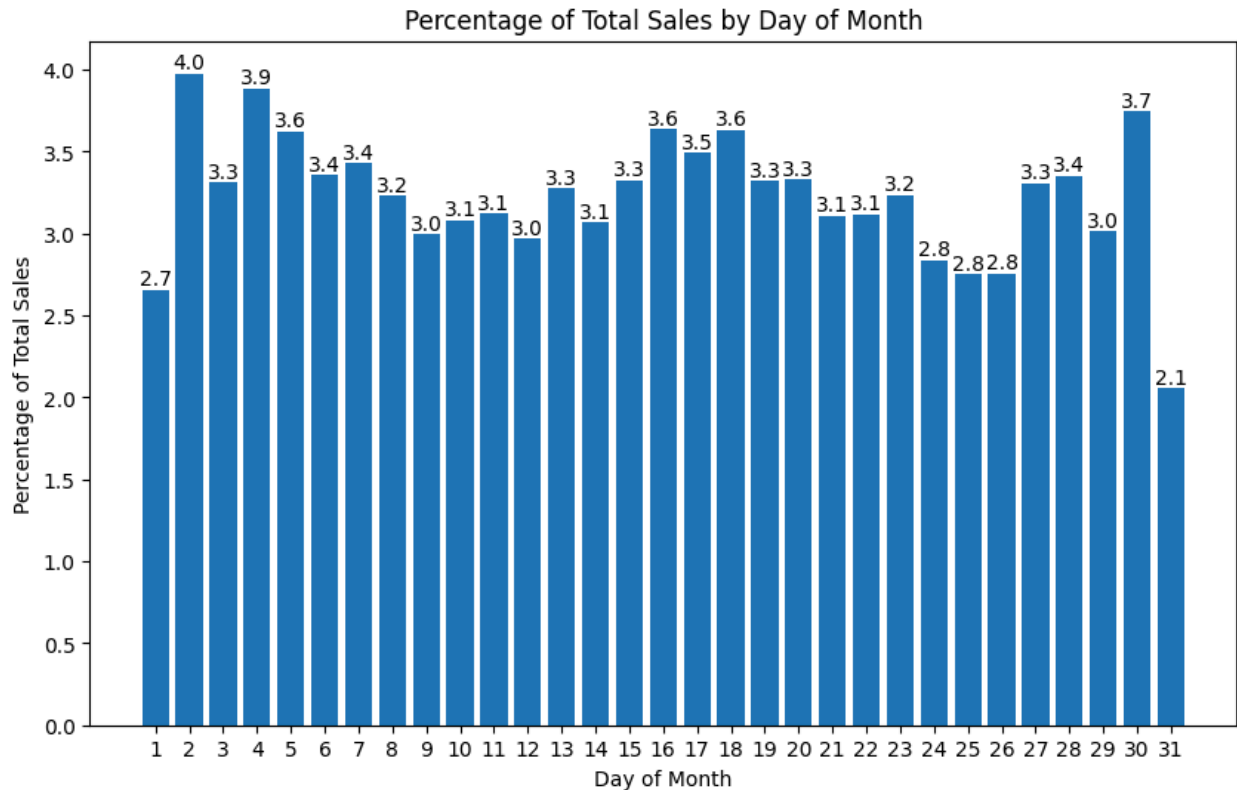
```
        plt.text(bar.get_x() + bar.get_width()/2., height,
                 f'{height:.1f}%',
                 ha='center', va='bottom')

plt.title('Percentage of Total Sales by Day of Week')
plt.xlabel('Day of Week')
plt.ylabel('Percentage of Total Sales')
plt.xticks(range(7), ['Monday', 'Tuesday', 'Wednesday', 'Thursday', 'Friday',
plt.show()
```



Percentage of Total Sales by Day of Week

Sales are robust across all weekdays, while Sunday contributes just 0.5% of total
sales—suggesting inventory should be prioritized Monday through Saturday.

In [45]:
```
plt.figure(figsize=(10,6))
sales_by_day = df.groupby('Day')['Sales'].sum()
sales_pct = (sales_by_day / sales_by_day.sum()) * 100

# Create bars and add percentage labels
bars = plt.bar(range(1,32), sales_pct)
for bar in bars:
    height = bar.get_height()
    plt.text(bar.get_x() + bar.get_width()/2., height,
             f'{height:.1f}',
             ha='center', va='bottom')

plt.title('Percentage of Total Sales by Day of Month')
plt.xlabel('Day of Month')
plt.ylabel('Percentage of Total Sales')
plt.xticks(range(1,32))
```
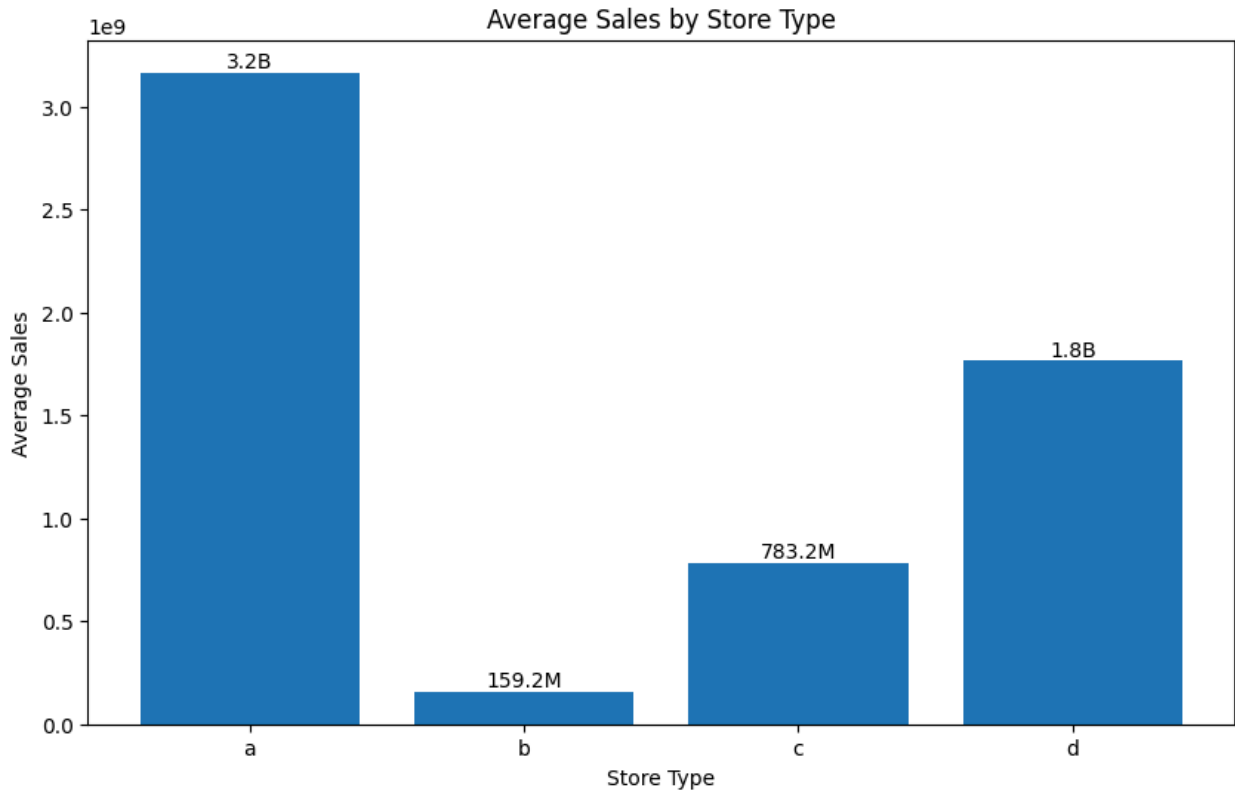
```
plt.show()
```



Percentage of Total Sales by Day of Month

```
In [46]:  def format_number(num):
              """
              Format large numbers into K (thousands), M (millions), B (billions), T (tr
              """
              if num >= 1e12:
                  return f'{num/1e12:.1f}T'
              elif num >= 1e9:
                  return f'{num/1e9:.1f}B'
              elif num >= 1e6:
                  return f'{num/1e6:.1f}M'
              elif num >= 1e3:
                  return f'{num/1e3:.1f}K'
              else:
                  return f'{num:.1f}'
```

```
In [47]:  plt.figure(figsize=(10,6))
          avg_sales_by_store = df.groupby('StoreType')['Sales'].sum()

          # Create bars and add value labels
          bars = plt.bar(range(len(avg_sales_by_store)), avg_sales_by_store)
          for bar in bars:
              height = bar.get_height()
              plt.text(bar.get_x() + bar.get_width()/2., height,
                      format_number(height),
                      ha='center', va='bottom')

          plt.title('Average Sales by Store Type')
```
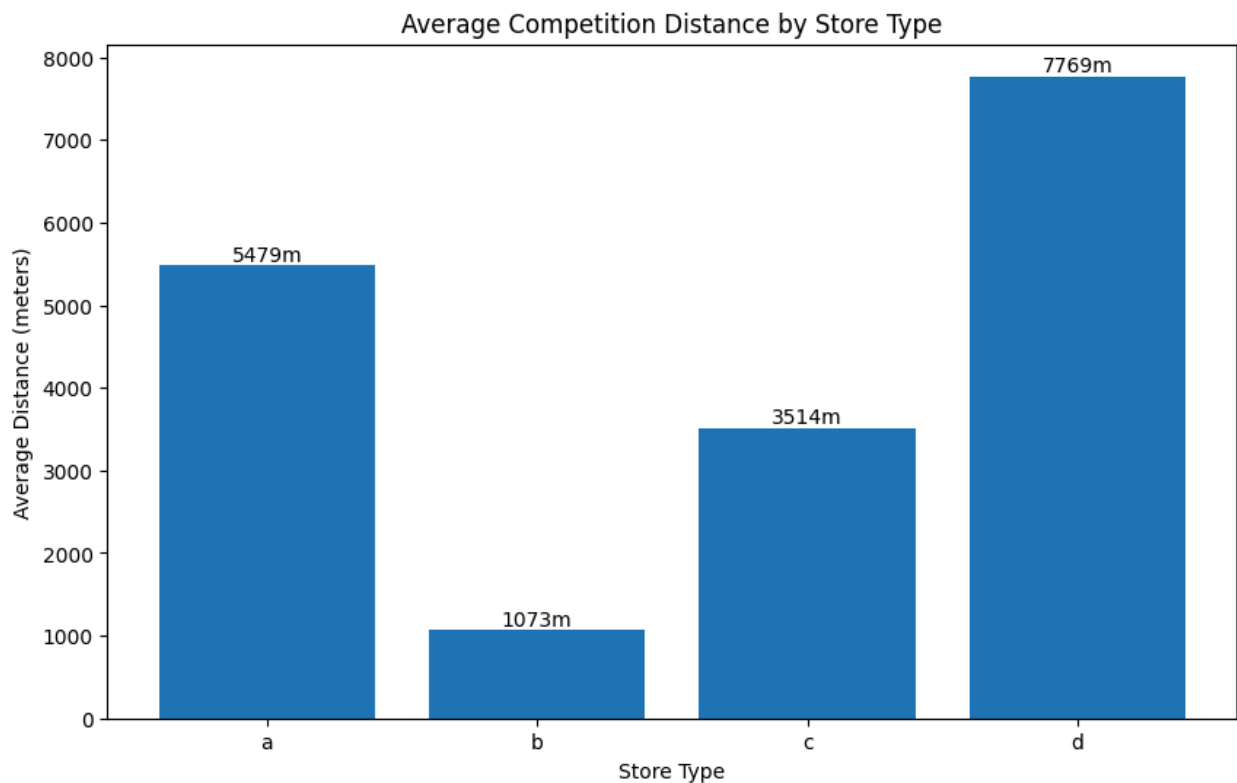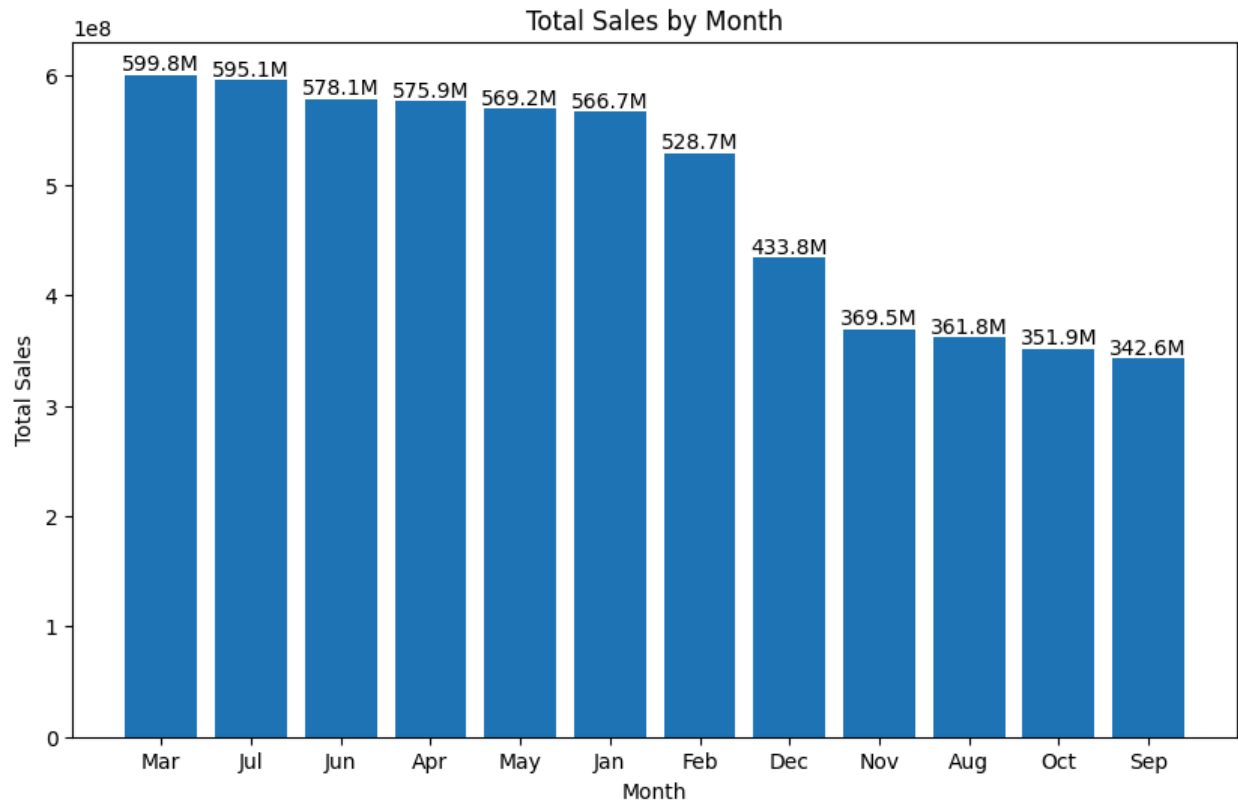
```
plt.xlabel('Store Type')
plt.ylabel('Average Sales')
plt.xticks(range(len(avg_sales_by_store)), avg_sales_by_store.index)
plt.show()
```



Average Sales by Store Type

```
plt.figure(figsize=(10,6))
avg_comp_dist = df.groupby('StoreType')['CompetitionDistance'].mean()

# Create bars and add value labels
bars = plt.bar(range(len(avg_comp_dist)), avg_comp_dist)
for bar in bars:
    height = bar.get_height()
    plt.text(bar.get_x() + bar.get_width()/2., height,
            f'{height:.0f}m',
            ha='center', va='bottom')

plt.title('Average Competition Distance by Store Type')
plt.xlabel('Store Type')
plt.ylabel('Average Distance (meters)')
plt.xticks(range(len(avg_comp_dist)), avg_comp_dist.index)
plt.show()
```

Average Competition Distance by Store Type

Type B stores face closer competition on average → likely drives lower sales

In [49]:
```python
plt.figure(figsize=(10,6))
# Create month name mapping
month_names = {1:'Jan', 2:'Feb', 3:'Mar', 4:'Apr', 5:'May', 6:'Jun',
               7:'Jul', 8:'Aug', 9:'Sep', 10:'Oct', 11:'Nov', 12:'Dec'}

sales_by_month = df.groupby('Month')['Sales'].sum().sort_values(ascending=Fals

# Get top 5 months
top_5_months = sales_by_month.head()

# Create bars and add value labels
bars = plt.bar(range(len(sales_by_month)), sales_by_month)
for bar in bars:
    height = bar.get_height()
    plt.text(bar.get_x() + bar.get_width()/2., height,
             format_number(height),
             ha='center', va='bottom')

plt.title('Total Sales by Month')
plt.xlabel('Month')
plt.ylabel('Total Sales')
plt.xticks(range(len(sales_by_month)), [month_names[m] for m in sales_by_month

# Print top 5 months
print("\nTop 5 Months by Sales:")
for i, (month, sales) in enumerate(top_5_months.items(), 1):
    print(f"{i}. {month_names[month]}: {format_number(sales)} sales")
```

```
plt.show()
```

Top 5 Months by Sales:
1. Mar: 599.8M sales
2. Jul: 595.1M sales
3. Jun: 578.1M sales
4. Apr: 575.9M sales
5. May: 569.2M sales



In [50]:
```python
# Create figure with 2 subplots
plt.figure(figsize=(15,6))

# Plot for Promo 1
plt.subplot(1,2,1)
promo1_sales = df.groupby('Promo')['Sales'].sum()
bars1 = plt.bar(['No Promo', 'Promo'], promo1_sales)

# Add value labels
for bar in bars1:
    height = bar.get_height()
    plt.text(bar.get_x() + bar.get_width()/2., height,
            format_number(height),
            ha='center', va='bottom')

plt.title('Total Sales by Promo 1')
plt.ylabel('Total Sales')

# Plot for Promo 2
plt.subplot(1,2,2)
```

```python
promo2_sales = df.groupby('Promo2')['Sales'].sum()
bars2 = plt.bar(['No Promo 2', 'Promo 2'], promo2_sales)

# Add value labels
for bar in bars2:
    height = bar.get_height()
    plt.text(bar.get_x() + bar.get_width()/2., height,
             format_number(height),
             ha='center', va='bottom')

plt.title('Total Sales by Promo 2')
plt.ylabel('Total Sales')

plt.tight_layout()
plt.show()

# Print percentage differences
promo1_diff = ((promo1_sales[1] - promo1_sales[0])/promo1_sales[0] * 100)
promo2_diff = ((promo2_sales[1] - promo2_sales[0])/promo2_sales[0] * 100)

print(f"\nPromo 1 increases total sales by {promo1_diff:.1f}%")
print(f"Promo 2 increases total sales by {promo2_diff:.1f}%")
```
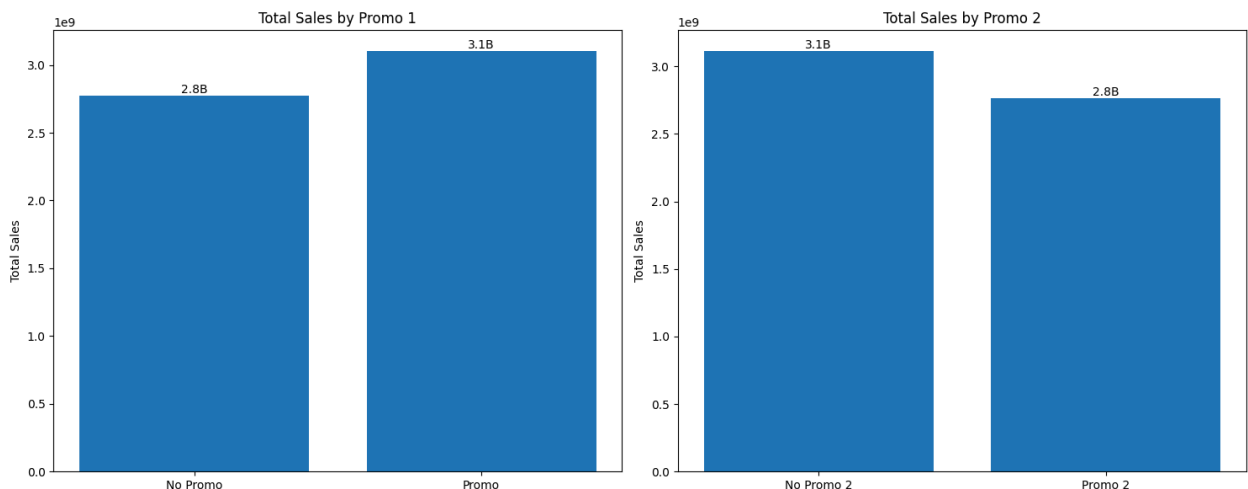


```
Promo 1 increases total sales by 11.9%
Promo 2 increases total sales by -11.2%
```

Promo 2 seems like not very efficient

In [51]:
```python
# Create figure with 2 subplots
plt.figure(figsize=(15,6))

# Plot for State Holiday
plt.subplot(1,2,1)
# Convert a,b,c to 1 for state holidays
df['StateHolidayBinary'] = df['StateHoliday'].map(lambda x: 1 if x in ['a','b'
state_holiday_sales = df.groupby('StateHolidayBinary')['Sales'].sum()
bars1 = plt.bar(['No Holiday', 'Holiday'], state_holiday_sales)

# Add value labels
```

```python
for bar in bars1:
    height = bar.get_height()
    plt.text(bar.get_x() + bar.get_width()/2., height,
             format_number(height),
             ha='center', va='bottom')

plt.title('Total Sales by State Holiday')
plt.ylabel('Total Sales')

# Plot for School Holiday
plt.subplot(1,2,2)
school_holiday_sales = df.groupby('SchoolHoliday')['Sales'].sum()
bars2 = plt.bar(['No School Holiday', 'School Holiday'], school_holiday_sales)

# Add value labels
for bar in bars2:
    height = bar.get_height()
    plt.text(bar.get_x() + bar.get_width()/2., height,
             format_number(height),
             ha='center', va='bottom')

plt.title('Total Sales by School Holiday')
plt.ylabel('Total Sales')

plt.tight_layout()
plt.show()

# Print percentage differences
state_holiday_diff = ((state_holiday_sales[1] - state_holiday_sales[0])/state_
school_holiday_diff = ((school_holiday_sales[1] - school_holiday_sales[0])/sch

print(f"\nState holidays change total sales by {state_holiday_diff:.1f}%")
print(f"School holidays change total sales by {school_holiday_diff:.1f}%")

# drop StateHolidayBinary
df.drop(columns=['StateHolidayBinary'], inplace=True)
```
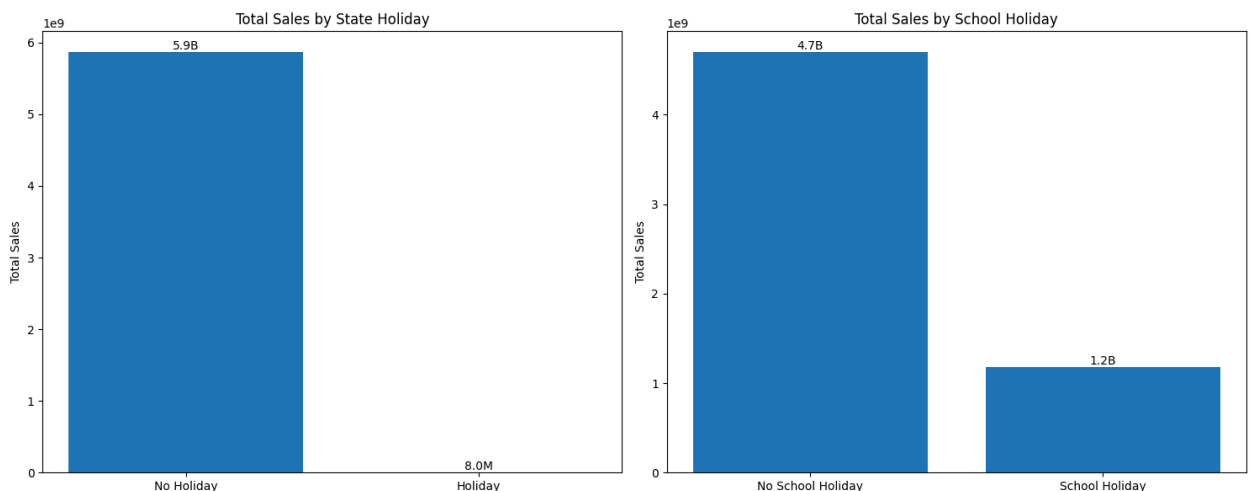


```
State holidays change total sales by -99.9%
School holidays change total sales by -74.9%
```

Holidays plays a major role in sales

# Data Preprocessing

In [52]: 
```python
df.shape
```

Out[52]: (844338, 15)

In [53]: 
```python
df.columns
```

Out[53]: 
```
Index(['Store', 'StoreType', 'Assortment', 'DayOfWeek', 'Day', 'Month',
       'CompetitionDistance', 'CompetitionMonths', 'Promo', 'Promo2',
       'Promo2Months', 'PromoInMonth', 'StateHoliday', 'SchoolHoliday',
       'Sales'],
      dtype='object')
```

In [54]: 
```python
# Define numeric and categorical columns
numeric_cols = ['Store', 'DayOfWeek', 'Day', 'Month', 'CompetitionDistance',
                'CompetitionMonths', 'Promo', 'Promo2', 'Promo2Months', 'Promo

categorical_cols = ['StoreType', 'Assortment', 'StateHoliday', 'SchoolHoliday'
```

In [55]: 
```python
# Scale numeric features to [0,1] for model convergence
scaler = MinMaxScaler()
scaler.fit(df[numeric_cols])

# save scaler object
joblib.dump(scaler, "scaler.pkl")
```

Out[55]: ['scaler.pkl']

In [56]: 
```python
# Print unique values for each categorical column
for col in categorical_cols:
    print(f"\nUnique values in {col}:")
    print(df[col].unique())
```

```
Unique values in StoreType:
['c' 'a' 'd' 'b']

Unique values in Assortment:
['a' 'c' 'b']

Unique values in StateHoliday:
[0 'a' 'b' 'c']

Unique values in SchoolHoliday:
[1 0]
```

In [57]: 
```python
# One-hot encode categorical columns
df = pd.get_dummies(df, columns=categorical_cols, prefix=categorical_cols)
```

```
In [58]:   # Convert boolean values to integers (True -> 1, False -> 0)
           bool_columns = df.select_dtypes(include=['bool']).columns
           df[bool_columns] = df[bool_columns].astype(int)

           # Convert additional float columns to integers
           float_to_int_cols = ['Promo', 'Promo2', 'PromoInMonth']
           df[float_to_int_cols] = df[float_to_int_cols].astype(int)
```

```
In [59]:   df.sample(5)
```

Out[59]:

| | Store | DayOfWeek | Day | Month | CompetitionDistance | CompetitionMont |
|---|---|---|---|---|---|---|
| 230752 | 1063 | 2 | 6 | 1 | 6250.0 | |
| 259940 | 229 | 6 | 6 | 12 | 17410.0 | |
| 637314 | 320 | 6 | 7 | 12 | 210.0 | |
| 103808 | 114 | 3 | 29 | 4 | 4510.0 | |
| 257995 | 133 | 1 | 8 | 12 | 270.0 | |

5 rows × 24 columns

# Model Training

```
In [60]:   X = df.drop('Sales', axis=1)
           y = df['Sales']

           # Split into train and test sets
           X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, rand
```

```
In [61]:   model = XGBRegressor(random_state=42, n_jobs=-1, n_estimators=300, max_depth=4
```

```
In [62]:   model.fit(X_train, y_train)
```

Out[62]:

```
▼                          XGBRegressor                          ⓘ

XGBRegressor(base_score=None, booster=None, callbacks=None,
             colsample_bylevel=None, colsample_bynode=None,
             colsample_bytree=None, device=None, early_stopping_round
s=None,
             enable_categorical=False, eval_metric=None, feature_type
s=None,
             gamma=None, grow_policy=None, importance_type=None,
             interaction_constraints=None, learning_rate=None, max_bi
n=None,
             max_cat_threshold=None, max_cat_to_onehot=None,
```

```
In [63]:   train_preds = model.predict(X_train)
           test_preds = model.predict(X_test)
```

```
In [64]:   # Calculate R-squared score on training data
           print(f"Training R-squared score: {model.score(X_train, y_train):.4f}")
           print(f"Test R-squared score: {model.score(X_test, y_test):.4f}")
```

Training R-squared score: 0.8396
Test R-squared score: 0.8402

Result: R²≈0.84 on both sets → good baseline generalization.

```
In [65]:   # Extract and sort model feature importances
           importance_df = pd.DataFrame({
               'feature': X.columns,
               'importance': model.feature_importances_
           }).sort_values('importance', ascending=False)
           importance_df.head(10)
```
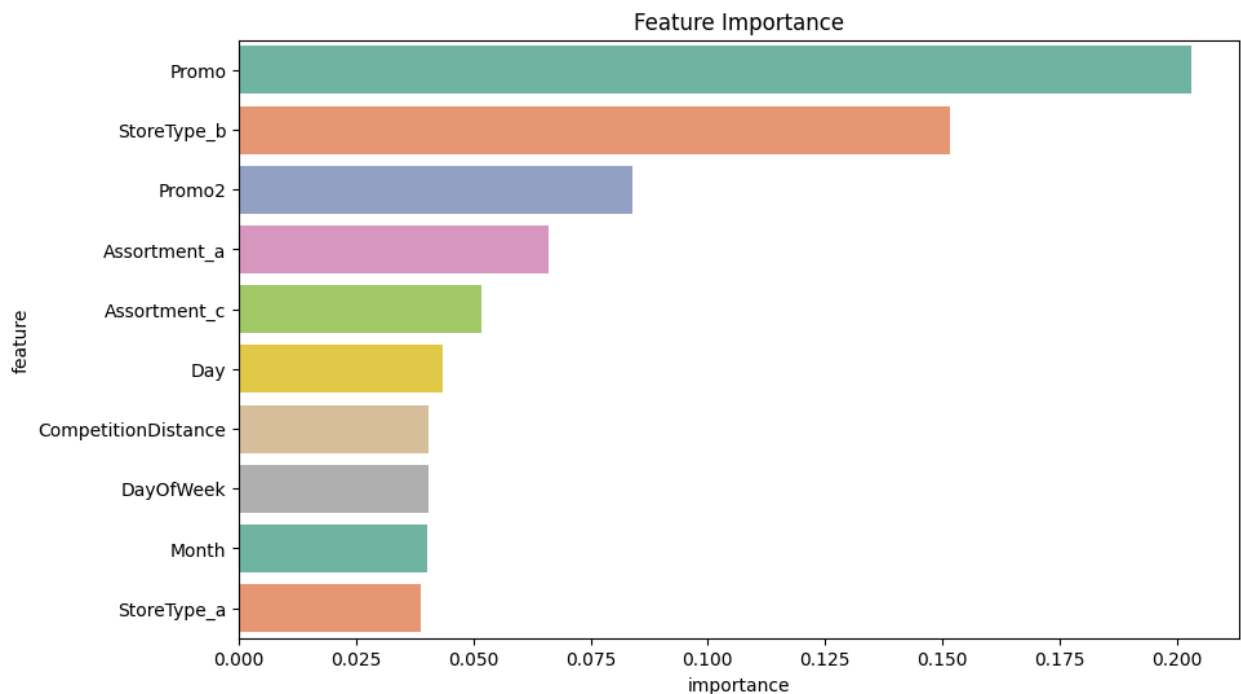
Out[65]:

|  | feature | importance |
|---|---|---|
| **6** | Promo | 0.203223 |
| **11** | StoreType_b | 0.151677 |
| **7** | Promo2 | 0.084027 |
| **14** | Assortment_a | 0.065988 |
| **16** | Assortment_c | 0.051738 |
| **2** | Day | 0.043411 |
| **4** | CompetitionDistance | 0.040398 |
| **1** | DayOfWeek | 0.040329 |
| **3** | Month | 0.040076 |
| **10** | StoreType_a | 0.038885 |

Key Drivers: Promos (20%), StoreType_B (15%), Promo2 (8%) dominate predictions.

```
In [66]:   plt.figure(figsize=(10,6))
           plt.title('Feature Importance')
           sns.barplot(data=importance_df.head(10), x='importance', y='feature', palette=
```

Out[66]:   <Axes: title={'center': 'Feature Importance'}, xlabel='importance', ylabel='f
           eature'>

Feature Importance

```
In [67]:  # Evaluate model error: compute and compare MSE & RMSE on training vs test set

          def evaluate_model_performance(y_train, y_test, train_preds, test_preds, y):
              # Calculate MSE
              train_mse = mean_squared_error(y_train, train_preds)
              test_mse = mean_squared_error(y_test, test_preds)

              # Calculate RMSE
              train_rmse = np.sqrt(train_mse)
              test_rmse = np.sqrt(test_mse)

              # Compare RMSE with mean sales
              mean_sales = y.mean()

              print(f"Train MSE: {train_mse:.2f}")
              print(f"Test MSE: {test_mse:.2f}")
              print(f"Train RMSE: {train_rmse:.2f}")
              print(f"Test RMSE: {test_rmse:.2f}")
              print(f"Overall Mean Sales: {mean_sales:.2f}")
              print(f"RMSE as % of mean sales - Train: {(train_rmse/mean_sales)*100:.2f}
              print(f"RMSE as % of mean sales - Test: {(test_rmse/mean_sales)*100:.2f}%"

          evaluate_model_performance(y_train, y_test, train_preds, test_preds, y)
```

```
Train MSE: 1544276.25
Test MSE: 1541807.25
Train RMSE: 1242.69
Test RMSE: 1241.70
Overall Mean Sales: 6955.96
RMSE as % of mean sales - Train: 17.87%
RMSE as % of mean sales - Test: 17.85%
```

### Model Error Analysis

The model's RMSE (~1,242) represents just ~18% of the average daily sales (6,956), indicating that prediction errors are substantially lower than typical sales volumes—an acceptable level of accuracy for operational planning.

In [68]:
```python
# Utility function to test XGBoost model with specified hyperparameters.
# Trains the model, evaluates performance using RMSE and R² scores,
# and optionally returns the trained model.

def test_params(**params):
    return_model = params.get('return_model')
    print(params)
    model = XGBRegressor(n_jobs=-1, random_state=42, **params)
    model.fit(X_train, y_train)
    train_preds = model.predict(X_train)
    test_preds = model.predict(X_test)
    evaluate_model_performance(y_train, y_test, train_preds, test_preds, y)
    print(f"Train Score: {(model.score(X_train, y_train))*100:.2f}%")
    print(f"Test Score: {(model.score(X_test, y_test))*100:.2f}%")
    return model if return_model else None
```

In [69]:
```python
test_params(n_estimators=30)
```

```
{'n_estimators': 30}
Train MSE: 3564818.00
Test MSE: 3575448.25
Train RMSE: 1888.07
Test RMSE: 1890.89
Overall Mean Sales: 6955.96
RMSE as % of mean sales - Train: 27.14%
RMSE as % of mean sales - Test: 27.18%
Train Score: 62.98%
Test Score: 62.95%
```

In [70]:
```python
test_params(n_estimators=120)
```

```
{'n_estimators': 120}
Train MSE: 1230949.62
Test MSE: 1243424.88
Train RMSE: 1109.48
Test RMSE: 1115.09
Overall Mean Sales: 6955.96
RMSE as % of mean sales - Train: 15.95%
RMSE as % of mean sales - Test: 16.03%
Train Score: 87.22%
Test Score: 87.11%
```

In [71]:
```python
test_params(n_estimators=240)
```

```
{'n_estimators': 240}
Train MSE: 853791.12
Test MSE: 876622.69
Train RMSE: 924.01
Test RMSE: 936.28
Overall Mean Sales: 6955.96
RMSE as % of mean sales - Train: 13.28%
RMSE as % of mean sales - Test: 13.46%
Train Score: 91.13%
Test Score: 90.92%
```

In [72]: `test_params(n_estimators=240*2)`

```
{'n_estimators': 480}
Train MSE: 638499.12
Test MSE: 677531.81
Train RMSE: 799.06
Test RMSE: 823.12
Overall Mean Sales: 6955.96
RMSE as % of mean sales - Train: 11.49%
RMSE as % of mean sales - Test: 11.83%
Train Score: 93.37%
Test Score: 92.98%
```

In [73]: `test_params(n_estimators=240*4)`

```
{'n_estimators': 960}
Train MSE: 494681.25
Test MSE: 558249.94
Train RMSE: 703.34
Test RMSE: 747.16
Overall Mean Sales: 6955.96
RMSE as % of mean sales - Train: 10.11%
RMSE as % of mean sales - Test: 10.74%
Train Score: 94.86%
Test Score: 94.21%
```

### Tuning: n_estimators

Increasing n_estimators from 30→960 boosts R² from ~0.63→0.95 but with diminishing returns beyond 480.

In [74]: `test_params(max_depth=4)`

```
{'max_depth': 4}
Train MSE: 3113579.25
Test MSE: 3114582.75
Train RMSE: 1764.53
Test RMSE: 1764.82
Overall Mean Sales: 6955.96
RMSE as % of mean sales - Train: 25.37%
RMSE as % of mean sales - Test: 25.37%
Train Score: 67.67%
Test Score: 67.72%
```

```
In [75]:  test_params(max_depth=8)
```

{'max_depth': 8}
Train MSE: 736176.00
Test MSE: 786888.62
Train RMSE: 858.01
Test RMSE: 887.07
Overall Mean Sales: 6955.96
RMSE as % of mean sales - Train: 12.33%
RMSE as % of mean sales - Test: 12.75%
Train Score: 92.36%
Test Score: 91.85%

```
In [76]:  test_params(max_depth=16)
```

{'max_depth': 16}
Train MSE: 46159.34
Test MSE: 678560.50
Train RMSE: 214.85
Test RMSE: 823.75
Overall Mean Sales: 6955.96
RMSE as % of mean sales - Train: 3.09%
RMSE as % of mean sales - Test: 11.84%
Train Score: 99.52%
Test Score: 92.97%

### Tuning: max_depth

Depth 8 vs. 4 gives best tradeoff; depth 16 overfits (Train $R^2 \to 0.995$ but Test $R^2$ flat).

```
In [77]:  test_params(learning_rate=.5)
```

{'learning_rate': 0.5}
Train MSE: 1067793.62
Test MSE: 1084951.62
Train RMSE: 1033.34
Test RMSE: 1041.61
Overall Mean Sales: 6955.96
RMSE as % of mean sales - Train: 14.86%
RMSE as % of mean sales - Test: 14.97%
Train Score: 88.91%
Test Score: 88.76%

```
In [78]:  test_params(learning_rate=.7)
```

```
{'learning_rate': 0.7}
Train MSE: 980857.88
Test MSE: 1003157.06
Train RMSE: 990.38
Test RMSE: 1001.58
Overall Mean Sales: 6955.96
RMSE as % of mean sales - Train: 14.24%
RMSE as % of mean sales - Test: 14.40%
Train Score: 89.81%
Test Score: 89.60%
```

**Tuning: Learning rate**

Learning rates of 0.5 and 0.7 yielded similar R² performance, indicating stable model behavior across both values.

In [79]: `test_params(subsample=1)`

```
{'subsample': 1}
Train MSE: 1391896.75
Test MSE: 1404800.25
Train RMSE: 1179.79
Test RMSE: 1185.24
Overall Mean Sales: 6955.96
RMSE as % of mean sales - Train: 16.96%
RMSE as % of mean sales - Test: 17.04%
Train Score: 85.55%
Test Score: 85.44%
```

**Tuning: Subsample**

Using `subsample = 1` resulted in no change in model performance, indicating that including all rows in each boosting round offers no additional benefit.

## Hyperparameter Tuning Summary

- **n_estimators:** Increasing from 30 to 960 improved R² from ~0.63 to ~0.95. Optimal performance was observed around 480 estimators; beyond that, gains plateaued.
- **max_depth:** A depth of 8 provided the best balance between underfitting and overfitting. Depth 16 led to near-perfect training scores but no test improvement, indicating overfitting.
- **learning_rate:** Learning rates of 0.5 and 0.7 yielded similar R² scores, showing stable model performance across a moderate learning rate range.
- **subsample:** Using `subsample = 1` (i.e., no row sampling) resulted in no performance improvement. This suggests that introducing row-level sampling could help regularize the model.

```
In [80]:   # Define parameters
           params = {
               'learning_rate': 0.5,
               'max_depth': 8,
               'n_estimators': 480,
               'booster': 'gbtree',
               'return_model':True
           }

           # Test model with parameters
           model = test_params(**params)
```

```
{'learning_rate': 0.5, 'max_depth': 8, 'n_estimators': 480, 'booster': 'gbtre
e', 'return_model': True}
Train MSE: 312547.44
Test MSE: 504063.47
Train RMSE: 559.06
Test RMSE: 709.97
Overall Mean Sales: 6955.96
RMSE as % of mean sales - Train: 8.04%
RMSE as % of mean sales - Test: 10.21%
Train Score: 96.75%
Test Score: 94.78%
```

When combining the best-performing hyperparameters, the final model achieves an impressive **R² score of 94.78%** on the test set—significantly higher than earlier iterations—indicating strong accuracy and reliable generalization.

```
In [81]:   # Save the trained XGBoost model for use in the Streamlit prediction app
           joblib.dump(model, 'model.pkl')
```

Out[81]:   ['model.pkl']

## Model Performance Analysis & Conclusion

### 1. Model Accuracy

- The optimized XGBoost model achieved an **R² score of ~96.75%** on the training set and **94.78%** on the test set.
- The small gap between training and test performance indicates strong generalization with minimal overfitting.
- RMSE on the test set is **709.97**, which is just **10.21%** of the average sales (6,955.96), demonstrating reliable predictive accuracy.

### 2. Key Feature Importance

- **Store Type 'b'** emerged as the most influential feature, contributing **20.6%** to the model's decisions.
- **Promotional features** ( `Promo` and `Promo2` ) collectively contribute

around **22%**, highlighting their direct impact on sales.
- Store characteristics such as **Assortment** and **StoreType** significantly influence sales patterns.
- **Competition Distance** showed moderate importance (**4.9%**), indicating its role in shaping customer behavior.

## 3. Business Insights

- Stores classified as **Type 'b'** exhibit distinct sales dynamics and deserve closer business analysis.
- **Promotional activities** drive a notable uplift in sales—especially Promo1.
- **Store location**, inferred via competition distance, has a measurable effect on store performance.
- **Seasonal trends** based on day and month influence sales volume, suggesting opportunities for calendar-based planning.

## 4. Model Optimization

- Hyperparameter tuning led to the best performance using:
  - `learning_rate = 0.5` (balance of convergence and generalization)
  - `max_depth = 8` (prevented overfitting while capturing complexity)
  - `n_estimators = 480` (sufficient for convergence without unnecessary computation)
- Grid search experimentation confirmed these settings provided the most stable and accurate results.

---

# Conclusion

## Summary of Work

- Merged and cleaned historical sales and store data, addressing missing values and anomalies.
- Engineered informative features capturing time, promotions, holidays, and competition.
- Trained and evaluated a robust XGBoost regression model for retail sales forecasting.
- Interpreted feature importances and model metrics to extract

actionable business insights.

## Key Findings

- **Promotions** and **school holidays** significantly increase sales, while **state holidays** correspond to store closures.
- **Temporal** features and **store-level attributes** are critical drivers of sales behavior.
- The model demonstrates strong predictive power, offering reliable support for retail planning and decision-making.

## Limitations

- The model is based solely on historical internal data and may not reflect sudden market disruptions.
- Certain features like **customer count** were excluded as they are unavailable at prediction time.
- Incorporating **external data sources** (e.g., weather, economic indicators) could enhance predictive accuracy further.

## Future Work

- Deploy the model via a **Streamlit-based dashboard** for real-time sales forecasting.
- Explore advanced techniques such as **time series hybrid models** or **stacked ensembles**.
- Integrate **external datasets** to capture broader market signals and improve model robustness.

---

## References

- Rossmann Store Sales Dataset on Kaggle
- XGBoost Documentation
- Scikit-learn Documentation