

WALMART PROJECT DOCUMENTATION

DESCRIPTION

One of the leading retail stores in the US, Walmart, would like to predict the sales and demand accurately. There are certain events and holidays which impact sales on each day. There are sales data available for 45 stores of Walmart. The business is facing a challenge due to unforeseen demands and runs out of stock some times, due to the inappropriate machine learning algorithm. An

ideal ML algorithm will predict demand accurately and ingest factors like economic conditions including CPI, Unemployment Index, etc.

Walmart runs several promotional markdown events throughout the year. These markdowns precede prominent holidays, the four largest of all, which are the Super Bowl, Labour Day, Thanksgiving, and Christmas. The weeks including these holidays are weighted five times higher in the evaluation than non-holiday weeks. Part of the challenge presented by this competition is modeling the effects of markdowns on these holiday weeks in the absence of complete/ideal historical data. Historical sales data for 45 Walmart stores located in different regions are available.

Dataset Description

This is the historical data which covers sales from 2010-02-05 to 2012-11-01, in the file Walmart_Store_sales. Within this file you will find the following fields:

- Store - the store number
- Date - the week of sales
- Weekly_Sales - sales for the given store
- Holiday_Flag - whether the week is a special holiday week 1 – Holiday week 0 – Non-holiday week
- Temperature - Temperature on the day of sale
- Fuel_Price - Cost of fuel in the region
- CPI – Prevailing consumer price index
- Unemployment - Prevailing unemployment rate

Holiday Events

Super Bowl: 12-Feb-10, 11-Feb-11, 10-Feb-12, 8-Feb-13

Labour Day: 10-Sep-10, 9-Sep-11, 7-Sep-12, 6-Sep-13

Thanksgiving: 26-Nov-10, 25-Nov-11, 23-Nov-12, 29-Nov-13

Christmas: 31-Dec-10, 30-Dec-11, 28-Dec-12, 27-Dec-13

Analysis Tasks

Basic Statistics tasks

- Q1 -> Which store has maximum sales
- Ans -> The Store 20 has maximum Sales. (By using aggregate and Arrange function)
- Q2 -> Which store has maximum standard deviation i.e., the sales vary a lot. Also, find out the coefficient of mean to standard deviation
- Ans -> The Store with Maximum Standard Deviation is Store 14. (By using summarise with sd function and arrange function.)
- Q3 -> Which store/s has good quarterly growth rate in Q3'2012
- Ans -> The Store with Highest Quarterly Growth rate is Store 7. (By using filter function, summarise with sum function, Transform function and arrange function.)
- Q4 -> Some holidays have a negative impact on sales. Find out holidays which have higher sales than the mean sales in non-holiday season for all stores together

The screenshot displays the RStudio environment with a data frame loaded. The data frame has columns: Store, Date, Weekly_Sales, Holiday_Flag, Temperature, Fuel_Price, CPI, Unemployment, Quarter, Year, Quarter_Year, and check_data. The console shows the output of a linear regression model:

```

Residuals:
    min       1q   median       3q      max
-287703.3  -85237.7  -22986.6  61308.8  87882.9

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) -3483483.4    3001974.9   -1.160  0.2479
store_data_s1_no 235.9      34356.9    0.165  0.8690
store_data_s1_noCPI 19855.8    13547.0    1.466  0.1450
store_data_s1_noUnemployment 324652.4    59278.7    2.120  0.0367
store_data_s1_noFuelPrice -67463.6    49616.7   -1.360  0.1761
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 131300 on 138 degrees of freedom
Multiple R-squared:  0.08517,    Adjusted R-squared:  0.05865
F-statistic: 3.722 on 4 and 138 Df, p-value: 0.02479

> holiday_sales = transform(holiday_sales, check_data = ifelse(weekly_sales > 1041256, "TRUE", "FALSE"))
> view(holiday_sales)
>
  
```

- Q5 -> Provide a monthly and semester view of sales in units and give insights

• Ans ->

RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

Environment History Connections Tutorial

Global Environment

- holiday_sales: 450 obs. of 12 variables
- mtcars_final: 32 obs. of 13 variables
- q2_2012_sales: 585 obs. of 11 variables
- q3_2012_sales: 585 obs. of 11 variables
- sales_by_month: 33 obs. of 3 variables
- sales_by_month_2010: 11 obs. of 3 variables
- sales_by_semester: 6 obs. of 3 variables
- sales_by_semester_2010: 2 obs. of 3 variables
- std_avg_sales: 45 obs. of 4 variables
- std_avg_sales_sort: 45 obs. of 3 variables
- store_sun_com: 45 obs. of 2 variables
- store_sun_q2_2012: 45 obs. of 2 variables

Files Plots Packages Help Viewer

Home . R files

- Demo_2_Perform Regression Analysis with multiple variables.csv: 36.8 KB, Aug 27, 2021, 11:01 AM
- Demo_3_Decision Tree Classification.csv: 48.7 KB, Aug 27, 2021, 10:13 AM
- Demo_4_K-Fold Cross validation.csv: 334.8 KB, Aug 27, 2021, 10:13 AM
- Diabetes.csv: 24.4 KB, Aug 29, 2021, 10:56 AM
- ECommerce.xlsx: 6.1 MB, Aug 23, 2021, 11:16 AM
- final output of cluster.csv: 574.1 KB, Sep 4, 2021, 1:22 PM
- ggplot.R: 764 B, Aug 19, 2021, 5:46 PM
- homework_class.R: 1.8 KB, Aug 29, 2021, 4:31 PM
- homework_class.R: 1.8 KB, Aug 29, 2021, 2:01 PM
- homework_class.R: 1.2 KB, Aug 31, 2021, 11:31 AM
- hbc.csv: 5.8 KB, Aug 15, 2021, 1:31 PM
- hypothesis.R: 238 B, Aug 22, 2021, 1:45 PM
- Lesson_3_Data Structures
- Lesson_7_Regression Analysis
- Lesson_8_Classification
- M2_Movie_Metadata.csv: 1.4 MB, Aug 21, 2021, 1:22 PM
- math distribution.jpeg: 17.3 KB, Aug 21, 2021, 12:24 PM
- reading_files.R: 762 B, Aug 19, 2021, 5:19 PM
- regression.R: 1.1 KB, Aug 27, 2021, 11:14 AM
- Text.csv: 4.4 KB, Aug 18, 2021, 10:21 AM
- walmart_project.R: 969 B, Sep 5, 2021, 11:42 AM
- Walmart_Store_sales.csv: 355.2 KB, Sep 5, 2021, 9:32 AM

Console

```
R 4.1.1 ~> fit <- lm(store_sun_q2_2012 ~ store_sun_com + holiday_sales + mtcars_final + q2_2012_sales + q3_2012_sales + sales_by_month + sales_by_month_2010 + sales_by_semester + sales_by_semester_2010 + std_avg_sales + std_avg_sales_sort + store_sun_com + store_sun_q2_2012)
```

Estimate Std. Error t value Pr(>|t|)

(Intercept) -3483483.4 3001974.9 -1.160 0.2479

store_data_sl_no\$sl_no 235.9 1426.9 0.165 0.8690

store_data_sl_no\$CP 19855.8 33547.0 1.466 0.1450

store_data_sl_no\$unemployment 124852.4 59178.7 2.110 0.0367 *

store_data_sl_no\$fuel_price -67463.6 49616.7 -1.360 0.1761

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 151300 on 138 degrees of freedom

Multiple R-squared: 0.08517, Adjusted R-squared: 0.05965

F-statistic: 3.212 on 4 and 138 Df, p-value: 0.01479

> holiday_sales = transform(holiday_sales, check_data = ifelse(weekly_sales > 1041256, "TRUE", "FALSE"))

> view(holiday_sales)

> sales_by_semester = summarise(group_by(walmart_data, Year, Semester), semester_sales = sum(weekly_sales))

'summarise()' has grouped output by 'Year'. You can override using the '.groups' argument.

> view(sales_by_semester)

> sales_by_month = summarise(group_by(walmart_data, Year, Month), month_sales = sum(weekly_sales))

'summarise()' has grouped output by 'Year'. You can override using the '.groups' argument.

> view(sales_by_month)

RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

Environment History Connections Tutorial

Global Environment

- holiday_sales: 450 obs. of 12 variables
- mtcars_final: 32 obs. of 13 variables
- q2_2012_sales: 585 obs. of 11 variables
- q3_2012_sales: 585 obs. of 11 variables
- sales_by_month: 33 obs. of 3 variables
- sales_by_month_2010: 11 obs. of 3 variables
- sales_by_semester: 6 obs. of 3 variables
- sales_by_semester_2010: 2 obs. of 3 variables
- std_avg_sales: 45 obs. of 4 variables
- std_avg_sales_sort: 45 obs. of 3 variables
- store_sun_com: 45 obs. of 2 variables
- store_sun_q2_2012: 45 obs. of 2 variables

Files Plots Packages Help Viewer

Home . R files

- Demo_2_Perform Regression Analysis with multiple variables.csv: 36.8 KB, Aug 27, 2021, 11:01 AM
- Demo_3_Decision Tree Classification.csv: 48.7 KB, Aug 27, 2021, 10:13 AM
- Demo_4_K-Fold Cross validation.csv: 334.8 KB, Aug 27, 2021, 10:13 AM
- Diabetes.csv: 24.4 KB, Aug 29, 2021, 10:56 AM
- ECommerce.xlsx: 6.1 MB, Aug 23, 2021, 11:16 AM
- final output of cluster.csv: 574.1 KB, Sep 4, 2021, 1:22 PM
- ggplot.R: 764 B, Aug 19, 2021, 5:46 PM
- homework_class.R: 1.8 KB, Aug 29, 2021, 4:31 PM
- homework_class.R: 1.8 KB, Aug 29, 2021, 2:01 PM
- homework_class.R: 1.2 KB, Aug 31, 2021, 11:31 AM
- hbc.csv: 5.8 KB, Aug 15, 2021, 1:31 PM
- hypothesis.R: 238 B, Aug 22, 2021, 1:45 PM
- Lesson_3_Data Structures
- Lesson_7_Regression Analysis
- Lesson_8_Classification
- M2_Movie_Metadata.csv: 1.4 MB, Aug 21, 2021, 1:22 PM
- math distribution.jpeg: 17.3 KB, Aug 21, 2021, 12:24 PM
- reading_files.R: 762 B, Aug 19, 2021, 5:19 PM
- regression.R: 1.1 KB, Aug 27, 2021, 11:14 AM
- Text.csv: 4.4 KB, Aug 18, 2021, 10:21 AM
- walmart_project.R: 969 B, Sep 5, 2021, 11:42 AM
- Walmart_Store_sales.csv: 355.2 KB, Sep 5, 2021, 9:32 AM

Console

```
R 4.1.1 ~> fit <- lm(store_sun_q2_2012 ~ store_sun_com + holiday_sales + mtcars_final + q2_2012_sales + q3_2012_sales + sales_by_month + sales_by_month_2010 + sales_by_semester + sales_by_semester_2010 + std_avg_sales + std_avg_sales_sort + store_sun_com + store_sun_q2_2012)
```

Estimate Std. Error t value Pr(>|t|)

(Intercept) -3483483.4 3001974.9 -1.160 0.2479

store_data_sl_no\$sl_no 235.9 1426.9 0.165 0.8690

store_data_sl_no\$CP 19855.8 33547.0 1.466 0.1450

store_data_sl_no\$unemployment 124852.4 59178.7 2.110 0.0367 *

store_data_sl_no\$fuel_price -67463.6 49616.7 -1.360 0.1761

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 151300 on 138 degrees of freedom

Multiple R-squared: 0.08517, Adjusted R-squared: 0.05965

F-statistic: 3.212 on 4 and 138 Df, p-value: 0.01479

> holiday_sales = transform(holiday_sales, check_data = ifelse(weekly_sales > 1041256, "TRUE", "FALSE"))

> view(holiday_sales)

> sales_by_semester = summarise(group_by(walmart_data, Year, Semester), semester_sales = sum(weekly_sales))

'summarise()' has grouped output by 'Year'. You can override using the '.groups' argument.

> view(sales_by_semester)

After data Visualization for Month sales it can be seen December month has the highest sales as compared to other Months and it's a outlier.

Form semester Sales, it can be seen that for year 2010 & 2011 the 2nd semester has better sales than the 1st semester. But for year 2012 the 1st semester has better sales than the 2nd semester.

Q6 -> Statistical Model

For Store 1 – Build prediction models to forecast demand

- Linear Regression – Utilize variables like date and restructure dates as 1 for 5 Feb 2010 (starting from the earliest date in order). Hypothesize if CPI, unemployment, and fuel price have any impact on sales.
- Change dates into days by creating new variable.
- Ans ->
- Hypothesis for Serial number,CPI,Unemployment,Fuel Price

The screenshot shows the RStudio interface with the following components:

- Script Editor:** Contains R code for hypothesis testing and model fitting. The code includes comments for hypotheses regarding serial number, CPI, unemployment, and fuel price, followed by the `lm()` function call.
- Console:** Displays the output of the `lm()` function, including the model summary with coefficients, standard errors, t-values, and p-values. The summary indicates that the model is significant (F-statistic = 0.0479).
- Environment:** Shows the data frame 'walmart_data' with 143 observations and 15 variables. The variables include 'date', 'store_id', 'store_name', 'total_sales', 'cpi', 'unemployment', and 'fuel_price'.
- Files:** Lists the files in the project, including 'walmart_data.csv' and 'walmart_project.R'.

- The Accuracy rate is 93%

The screenshot shows the RStudio interface with the following components:

- Script Editor:** Contains R code for calculating the accuracy rate. The code includes the `predict()` function call and the calculation of the accuracy rate using the `abs()` function.
- Console:** Displays the output of the `predict()` function, showing the accuracy rate is 0.93629134656472, which is approximately 93.6%.
- Environment:** Shows the data frame 'walmart_data' with 143 observations and 15 variables. The variables include 'date', 'store_id', 'store_name', 'total_sales', 'cpi', 'unemployment', and 'fuel_price'.
- Files:** Lists the files in the project, including 'walmart_data.csv' and 'walmart_project.R'.