Weavenow's Case Study Report Duc Anh Bui (Tyler) 2022-09-26 Task We would like to how to predict demand in retail environments using some standard datasets. You are requested to do the following steps: Data processing and modeling factors Common demand prediction methods (including feature selection and regularization) Evaluation and visualization Loading libraries library(tidyverse) library(dplyr) library(tseries) library(forecast) library(randomForest) library(prophet) library(reshape2) Data processing and cleaning data <- read.csv("data_raw.csv")</pre> df <- as.data.frame(data)</pre> week <- as.Date(sapply(strsplit(as.character(df\$week.sku.weekly_sales.feat_main_page.color.price.vendor.functiona</pre> lity), "\\,"), '[',1), "%Y-%m-%d") sku <- as.numeric(sapply(strsplit(as.character(df\$week.sku.weekly_sales.feat_main_page.color.price.vendor.functio nality), "\\,"), '[',2)) weekly_sales <- as.numeric(sapply(strsplit(as.character(df\$week.sku.weekly_sales.feat_main_page.color.price.vendo r.functionality), "\\,"), '[',3)) feat_main_page <- sapply(strsplit(as.character(df\$week.sku.weekly_sales.feat_main_page.color.price.vendor.functio</pre> nality), "\\,"), '[',4) color <- sapply(strsplit(as.character(df\$week.sku.weekly_sales.feat_main_page.color.price.vendor.functionality),</pre> "\\,"), '[',5) price <- as.numeric(sapply(strsplit(as.character(df\$week.sku.weekly_sales.feat_main_page.color.price.vendor.funct ionality), "\\,"), '[',6)) vendor <- as.numeric(sapply(strsplit(as.character(df\$week.sku.weekly_sales.feat_main_page.color.price.vendor.func</pre> tionality), "\\,"), '[',7)) functionality <- sapply(strsplit(as.character(df\$week.sku.weekly_sales.feat_main_page.color.price.vendor.function</pre> ality), "\\,"), '[',8) df <- data.frame(week, sku, weekly_sales, feat_main_page, color, price, vendor, functionality)</pre> write_csv(df, "processed_data.csv") Helper function to monitor the model's performence options(scipen=999) validation_stats_script <- function(y_actual, y_hat) {</pre> r2 <- cor(y_actual,y_hat)^2 mse <- mean((y_actual - y_hat)^2)</pre> mape <- mean(abs((y_actual-y_hat)/y_actual))*100</pre> mae <- mean(abs(y_actual - y_hat))</pre> rmse <- sqrt(mse) res <- list("r2" = signif(r2,2), "mse" = signif(mse,2), "mape" = signif(mape, 2), "mae" = signif(mae, 2), "rmse" = signif(rmse,2)) return(res) Part 1: Weekly Sales of SKU = 11 **ARIMA Model** Although it is widely known that Arima does not perform well with weekly data, it still is one of the basic/ traditional forecasting model for time series More data processing and plotting df_arima <- filter(df, sku == 11)</pre> keep <- c("week", "weekly_sales")</pre> df_arima <- df_arima[keep]</pre> $ggplot(df_arima, aes(x = week, y = weekly_sales)) +$ $geom_line() + ylim(0, 400) +$ scale_x_date(date_labels = "%Y-%m-%d", date_breaks = "2 month") + theme_bw() + theme(legend.title = element_blank(), axis.text.x = element_text(angle=45, vjust=0.5)) 400 · 300 100 Check for stationary of data adf.test(df_arima\$weekly_sales) ## Warning in adf.test(df_arima\$weekly_sales): p-value smaller than printed p-value Augmented Dickey-Fuller Test ## ## data: df_arima\$weekly_sales ## Dickey-Fuller = -4.7482, Lag order = 4, p-value = 0.01## alternative hypothesis: stationary acf(as.ts(df_arima\$weekly_sales), main = "Weekly Sales") **Weekly Sales** 1.0 0.8 9.0 ACF 0.4 0.2 0.0 -0.2 5 15 0 10 20 Lag pacf(as.ts(df_arima\$weekly_sales), main = "Weekly Sales") **Weekly Sales** 0.1 Partial ACF 0.0 -0.1 -0.2 15 5 10 20 Lag Based on the p-val of the ADF test, we reject the null hypothesis (data is non-stationary) when we choose the alpha level of 0.05. We are 95% confidence about this result that the data is stationary. However, the ACF and PACF above say otherwise, there is a correlation at lag 1. There are 3 parameters for the Arima model as Arima(p,d,q) With the result of PACF, it can be seen the p,q are 1 With the result of the ACF, it can be seen the d is 1 So, lets test around the parameters values around (1,3) Fitting the Arima model Since R has an amazing function to fit the best Arima parameters called auto.arima(), let's use it instead arima_mod <- auto.arima(as.ts(df_arima\$weekly_sales))</pre> The best model is Arima(2,1,1) Splitting training and testing datasets In order to test for the validity of the model, we first have to split the dataset in to training and testing datasets. train_index <- 0.7 # set the how much data to put into train data set</pre> n_total <- nrow(df_arima) # number of rows in the dataframe</pre> df_arima_train <- df_arima[1:(train_index*n_total),] # subsetting train dataset</pre> df_arima_test <- df_arima[(train_index*n_total+1):n_total,] # subsetting test dataset</pre> arima_predict <- numeric(n_total-nrow(df_arima_train)) #predetermine the dataset for the predictions from arima m</pre> Fitting the ARIMA model recursively for (i in 1:(n_total-(train_index*n_total))) { df_arima_train1 <- df_arima[1:(train_index*n_total-1+i),]</pre> arima_model <- auto.arima(as.ts(df_arima_train1\$weekly_sales))</pre> pred <- forecast(arima_model,1)</pre> arima_predict[i]<- pred\$mean</pre> Visualize the result df_arima_pred <- tibble("Train" = c(df_arima_train\$weekly_sales, df_arima_test\$weekly_sales),</pre> "Test" = c(df_arima_train\$weekly_sales, df_arima_test\$weekly_sales), "ARIMA" = c(df_arima_train\$weekly_sales, arima_predict), time = df_arima\$week) # create a tibble final1 <- melt(data = df_arima_pred, id.vars = "time") # melt the data into molten data frame</pre> final1 %>% ggplot(aes(x = time, y = value, col = variable, linetype = variable)) +geom_line() + xlab("") + ylab("Sales") + scale_color_manual(values=c("cyan3", "chocolate", "darkgoldenrod1")) + scale_linetype_manual(values=c(1, 4, 2)) + scale_x_date(date_labels = "%Y-%m-%d", date_breaks = "2 month") + theme_bw() + theme(legend.title = element_blank(), axis.text.x = element_text(angle=45, vjust=0.5)) + ggtitle("Weekly SKU 11") + theme(plot.title = element_text(hjust = 0.5)) Weekly SKU 11 300 200 Train · - · Test - - ARIMA 100 · Arima model's validation metric metric <- validation_stats_script(df_arima_test\$weekly_sales, arima_predict)</pre> metric.name <- c("r2", "mse", "mape", "mae", "rmse")</pre> metric.val <- c(metric\$r2, metric\$mse, metric\$mape, metric\$mae, metric\$rmse)</pre> metric.final <- data.frame(metric.name, metric.val)</pre> metric.final metric.name metric.val <chr> <dpl> r2 0.004 6100.000 mse 73.000 mape mae 56.000 78.000 rmse 5 rows As expected, the Arima model does not perform well with weekly data that the r2 only shows 0.004 mean the model only explains 0.4% of the data. With MAPE of 73, the forecast is off by 73% from true value. Similar with MAPE, MAE shows 56 meaning the forecast is off by 56 unit from true value. Part 2: Weekly Total Sales Random Forest Model Random Forest model is an ensemble learning method (use multiple learning algorithm to obtain better predictive performance) of classification, regression and other tasks that operates by constructing a multitude of decision trees at training time. Data processing to find the Weekly Total Sales df_total_sales <- df %>% group_by(week) %>% summarise(total_sales = sum(weekly_sales)) Initial visualization of the Weekly Total Sales by Week df_total_sales %>% ggplot(aes(x= week, y = total_sales)) + geom_line() + scale_x_date(date_labels = "%Y-%m-%d", date_breaks = "2 month") + theme_bw() + theme(legend.title = element_blank(), axis.text.x = element_text(angle=45, vjust=0.5)) 7500 total_s 2500 week Since the data provided is very small, it hard to say if there is any seasonality in the data. Splitting training and testing datasets train_index <- 0.7 n_total <- nrow(df_total_sales)</pre> train <- df_total_sales[1:(train_index*n_total),]</pre> test <- df_total_sales[(train_index*n_total+1):n_total,]</pre> Fitting Random Forest model rd <- randomForest(data = train, total_sales~week) # fit the model</pre> rd_predict <- tail(rd\$predicted, nrow(test)) # forecast the data</pre> Visualize the forecasting reports with Random Forest model df_rd_pred <- tibble("Train" = c(train\$total_sales, test\$total_sales),</pre> "Test" = c(train\$total_sales, test\$total_sales), "Random Forest" = c(train\$total_sales,rd_predict), time = df_total_sales\$week) # create a tibble final2 <- melt(data = df_rd_pred, id.vars = "time") # melt the tibble into molten data frame</pre> final2 %>% ggplot(aes(x = time, y = value, col = variable, linetype = variable)) +geom_line() + xlab("") + ylab("Sales") + scale_color_manual(values=c("cyan3", "chocolate", "darkgoldenrod1")) + scale_linetype_manual(values=c(1, 4, 2)) + scale_x_date(date_labels = "%Y-%m-%d", date_breaks = "2 month") + theme_bw() + theme(legend.title = element_blank(), axis.text.x = element_text(angle=45, vjust=0.5)) + ggtitle("Total Weekly Sales") + theme(plot.title = element_text(hjust = 0.5)) Total Weekly Sales 7500 Train Sales 0000 · – · Test Random Forest Random Forest model's validation metric metric2 <- validation_stats_script(test\$total_sales, rd_predict)</pre> metric.name2 <- c("r2", "mse", "mape", "mae", "rmse")</pre> metric.val2 <- c(metric2\$r2, metric2\$mse, metric2\$mape, metric2\$mae, metric2\$rmse)</pre> metric.final2 <- data.frame(metric.name2, metric.val2)</pre> metric.final2 metric.name2 metric.val2 <chr> <qpl> r2 0.14 8600000.00 mse 60.00 mape mae 2300.00 2900.00 rmse 5 rows r2 only shows 0.15 mean the model only explains 15% of the data. With MAPE of 60, the forecast is off by 60% from true value. Similar with MAPE, MAE shows 2300 meaning the forecast is off by 2300 unit from true value. **TBATS** Another traditional model to look at is the TBATS model. TBATS stands for: T: Trigonometric seasonality B: Box-Cox transformation A: ARIMA errors T: Trend S: Seasonal components Although TBATS performs the best with seasonal data, and is hard to say if there's any seasonality in the data, TBATS is still a great model to learn and use to forecast time series data. Fitting the TBATS model tbats_predict <- numeric(n_total-nrow(train))</pre> for (i in 1:(n_total-(train_index*n_total))) { df_tbats_train1 <- df_total_sales[1:(train_index*n_total-1+i),]</pre> tbats_model <- tbats(as.ts(df_tbats_train1\$total_sales))</pre> pred <- forecast(tbats_model, 1)</pre> tbats_predict[i] <- pred\$mean</pre> Visualize the forecasting reports with TBATS model df_tbats_pred <- tibble("Train" = c(train\$total_sales, test\$total_sales),</pre> "Test" = c(train\$total_sales, test\$total_sales), "TBATS" = c(train\$total_sales, tbats_predict), time = df_total_sales\$week) final3 <- melt(data = df_tbats_pred, id.vars = "time")</pre> final3 %>% ggplot(aes(x = time, y = value, col = variable, linetype = variable)) +geom_line() + xlab("") + ylab("Sales") + scale_color_manual(values=c("cyan3", "chocolate", "darkgoldenrod1")) + scale_linetype_manual(values=c(1, 4, 2)) + scale_x_date(date_labels = "%Y-%m-%d", date_breaks = "2 month") + theme_bw() + theme(legend.title = element_blank(), axis.text.x = element_text(angle=45, vjust=0.5)) + ggtitle("Total Weekly Sales") + theme(plot.title = element_text(hjust = 0.5)) Total Weekly Sales

3900000.00 mse 34.00 mape mae 1400.00 2000.00 rmse 5 rows The r2 only shows 0.22 mean the model only explains 22% of the data. With MAPE of 34, the forecast is off by 34% from true value. Similar with MAPE, MAE shows 1400 meaning the forecast is off by 1400 unit from true value. Prophet

Lastly, the most modern to look at is the Prophet, designed by the Facebook team to help both experienced and non-experienced statistician to

future <- make_future_dataframe(prophet_model, periods = 30, freq = "week") # make future datas for forecasting

prophet_predict <- tail(InvBoxCox(forecast\$yhat, lam),30) # inverse transform the prediction from the prophet mod</pre>

Train

metric.val3

metric.val4

13000000.00

<qpl>

0.29

64.00

2900.00

3600.00

<dbl>

0.22

· - · Test -- TBATS

7500

Sales 5000

2500

TBATS model's validation metric

metric.final3

metric.name3

forecast the time series data.

keep <- c("week", "y") train1 <- train1[keep]</pre>

Visualizing the model against train and test sets

xlab("") + ylab("Sales") +

ggtitle("Total Weekly Sales") +

Prophet model's validation metric

ggtitle("Total Weekly Sales") +

5 rows

underfitness of model.

theme(plot.title = element_text(hjust = 0.5))

Total Weekly Sales

final4 %>%

geom_line() +

train1 <- train

Finding the lambda value for Prophet Model

lam <- BoxCox.lambda(train1\$total_sales, method="loglik") #lambda</pre>

train1\$y <- BoxCox(train1\$total_sales, lam) #transform the data</pre>

Fitting the prophet model and visualizing the performence

prophet_model <- prophet(train1) # fit the model</pre>

forecast <- predict(prophet_model, future) # forecast data</pre>

final4 <- melt(data = df_prophet_pred, id.vars = "time")</pre>

theme_bw() + theme(legend.title = element_blank(),

theme(plot.title = element_text(hjust = 0.5))

scale_linetype_manual(values=c(1, 4, 2)) +

df_prophet_pred <- tibble("Train" = c(train\$total_sales, test\$total_sales),</pre>

time = df_total_sales\$week)

ggplot(aes(x = time, y = value, colour = variable, linetype = variable)) +

scale_color_manual(values=c("cyan3", "chocolate", "darkgoldenrod1")) +

scale_x_date(date_labels = "%Y-%m-%d", date_breaks = "2 month") +

metric4 <- validation_stats_script(test\$total_sales, prophet_predict)</pre>

metric.val4 <- c(metric4\$r2, metric4\$mse, metric4\$mape, metric4\$mae, metric4\$rmse)</pre>

metric.name4 <- c("r2", "mse", "mape", "mae", "rmse")</pre>

metric.final4 <- data.frame(metric.name4, metric.val4)</pre>

Transform the data with the lambda found and change the column names to fit the model's requirment

colnames(train1) <- c("ds", "y") # change the column name because the model required to</pre>

"Test" = c(train\$total_sales, test\$total_sales),

axis.text.x = element_text(angle=45, vjust=0.5)) +

"Prophet" = c(train\$total_sales, tail(prophet_predict,30)),

<chr>

r2

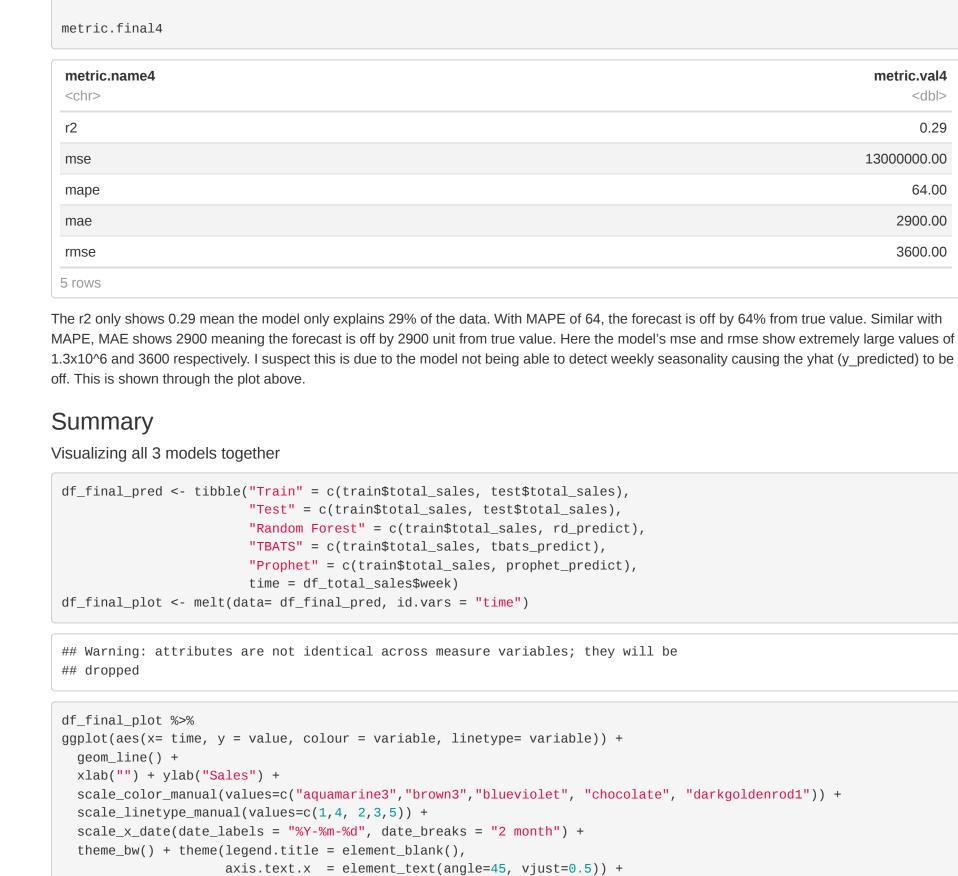
metric3 <- validation_stats_script(test\$total_sales, tbats_predict)</pre>

metric.val3 <- c(metric3\$r2, metric3\$mse, metric3\$mape, metric3\$mae, metric3\$rmse)</pre>

metric.name3 <- c("r2", "mse", "mape", "mae", "rmse")</pre>

metric.final3 <- data.frame(metric.name3, metric.val3)</pre>

Total Weekly Sales 7500 Train Sales ooos · - · Test - - Prophet 2500



7500 Train · – · Test Sales 5000 Random Forest ···· TBATS Prophet

<pre>metric.name_fin <- c("r2", "mse", "mape", "mae", "rmse") rd <- c(metric2\$r2, metric2\$mse, metric2\$mape, metric2\$mae, metric2\$rmse) tbats <- c(metric3\$r2, metric3\$mse, metric3\$mape, metric3\$mae, metric3\$rmse) prophet <- c(metric4\$r2, metric4\$mse, metric4\$mape, metric4\$mae, metric4\$rmse) metric.final_fin <- data.frame(metric.name_fin,rd, tbats, prophet) colnames(metric.final_fin) <- c("Model Metric","Random Forest", "TBATS", "Prophet") metric.final_fin</pre>			
Model Metric	Random Forest	TBATS	Prophet
Model Metric			عالماء
	<qpl></qpl>	<pre><dbl></dbl></pre>	<dpl></dpl>
<chr></chr>	<dbl></dbl>	<dbl> 0.22</dbl>	
<chr><r2< td=""><td></td><td></td><td>0.29</td></r2<></chr>			0.29
<chr> r2 mse</chr>	0.14	0.22	0.29
model Metric <chr> r2 mse mape mae</chr>	0.14 8600000.00	0.22	0.29 13000000.00 64.00 2900.00