# Assignment 04

Forecast Sales

15146314 Yang, Seunghyuck
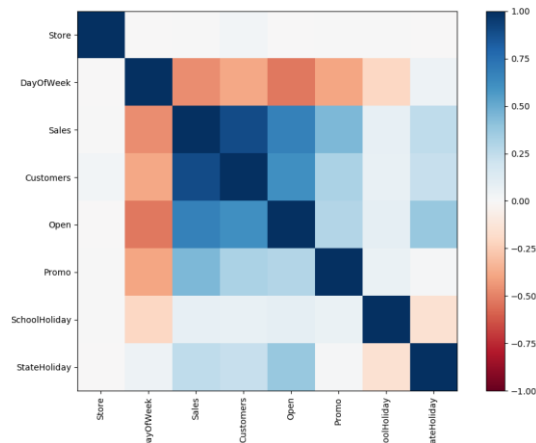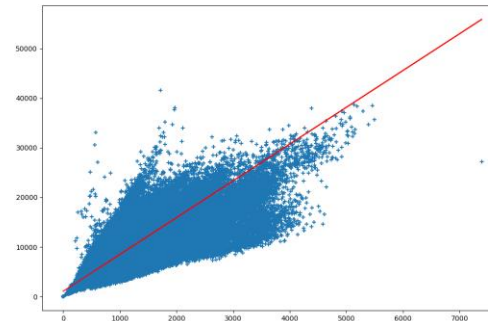
## Correlation Graph



Tier 1: Customers, Open, Promo

Tier 2: SchoolHoliday, StateHoliday, DayOfWeek

## Customer



```
In [46]: train[['Customers', 'Sales']].corr()
Out[46]:
                Customers      Sales
Customers      1.000000   0.894711
Sales          0.894711   1.000000
```
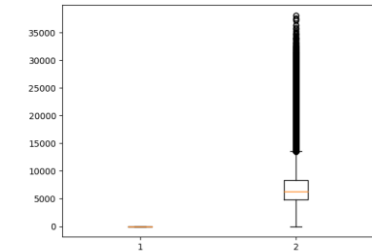
High correlation

Observable data

Select

## Open



```
In [72]: train[(train['Open'] == 1) & (train['Sales'] > 8360)].count()
Out[72]:
Store          211048
DayOfWeek      211048
Date           211048
Sales          211048
Customers      211048
Open           211048
Promo          211048
SchoolHoliday  211048
StateHoliday   211048
dtype: int64
```
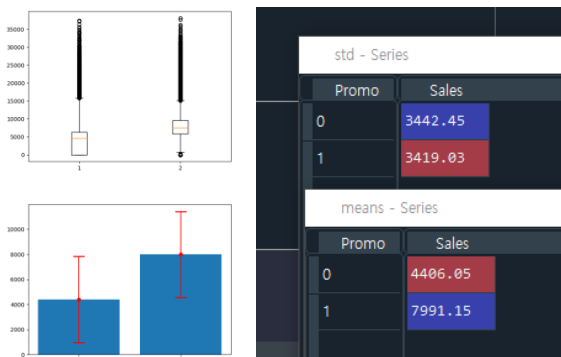
Binary data

Clear relationship

Reasonable correlation

Outliner is too lot

Select without transform

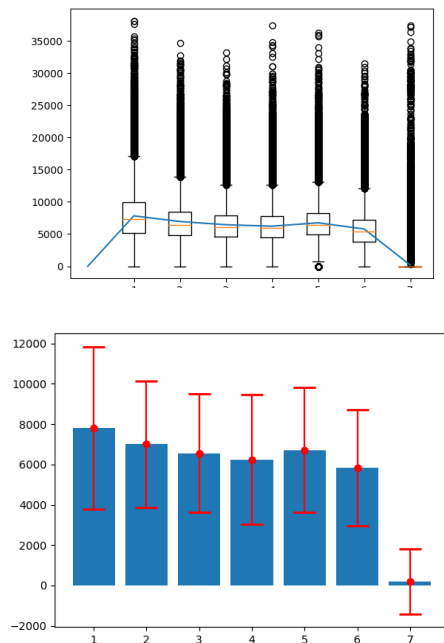# Promo



Binary data

Prominent correlation

Observable gap

Standard deviation is too high
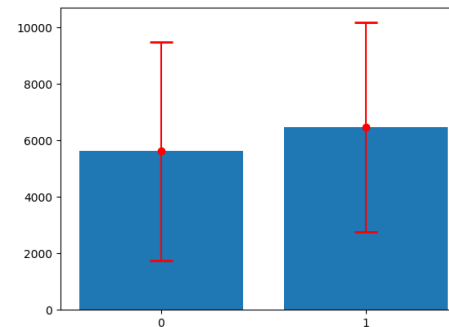
**Discard**

# DayOfWeek



Categorical data

Low correlation

Barely see difference in overall

Clear gap between 7 and others

**Select with transform**
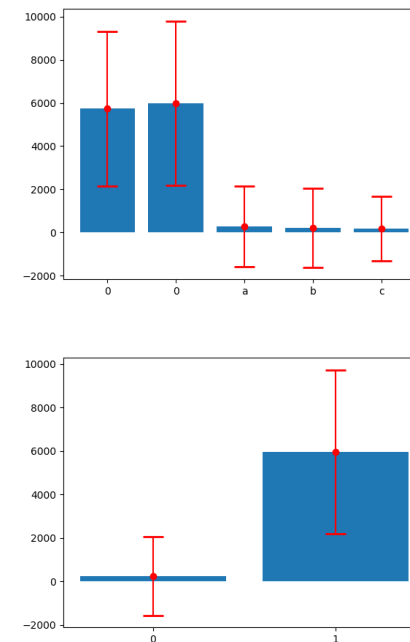
# SchoolHoliday



Binary data

Low correlation

Barely see difference in overall

**Discard**

# StateHoliday



Categorical data

Low correlation

Barely see difference in overall

Low correlation after transform

**Discard**

# Preparing data

# Evaluating variables through Cross Validation

## Logistic Regression

### KFold

### StratifiedKFold



```
discards = ['SchoolHoliday', 'StateHoliday', 'Promo', 'Store']
selects = ['Date', 'Customers', 'Open', 'DayOfWeek']
train = train.drop(discards, axis = 1)

newDay = train['DayOfWeek'] != 7
newDay = newDay.astype(int)
train = train.drop(['DayOfWeek'], axis = 1)
train = pd.concat((train, newDay), axis = 1)

condTrain = (train['Date'] < '2015-01-01')
Xtrain = train[condTrain][selects].drop(['Date'], axis = 1)
ytrain = train[condTrain]['Sales']
Xtest = train[condTrain != True][selects].drop(['Date'], axis = 1)
ytest = train[condTrain != True]['Sales']
```

Transforming data

Split data into 4 pieces

2013 – 2014 / 2015

Xtrain, ytrain, Xtest, ytest

Choose Logistic Regression

KFold, StratifiedKFold

Compare Accuracy

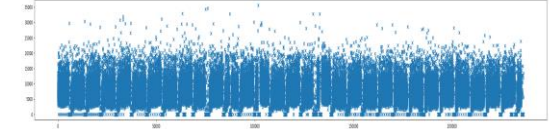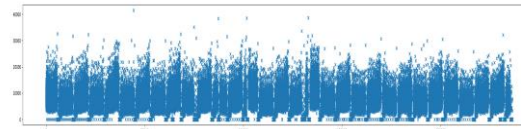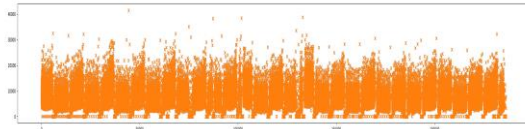Modify variable sets

Find better variable sets
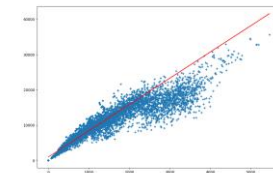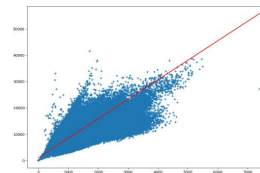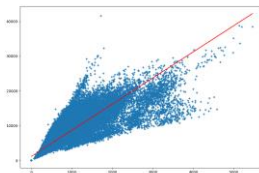
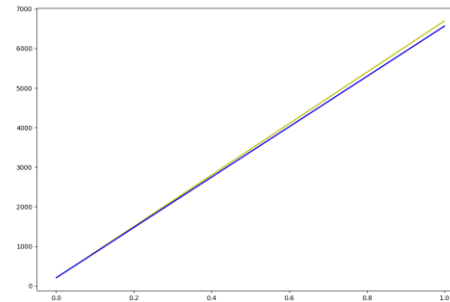# Linear Regression                    Actual data                    KNeighbors Regression



Scatter plot with regression line for each Sales ~ Dates



Scatter plot with regression line for each Sales ~ Customers



Regression lines for each
Sales ~ DayOfWeek (Binary)

R^2: 0.82970197                                                                          R^2: 0.80344976