

# Predictive Analysis - Web Traffic Time Series Forecasting | Kaggle

The goal of this notebook is not to do the best model for each Time Series. It is just a comparison of few models when you have one Time Series. The presentation present a different approaches to forecast a Time Series.

The plan of the notebook is:

- I. Importation & Data Cleaning
- II. Aggregation & Visualisation
- III. Machine Learning Approach
- IV. Basic Model Approach
- V. ARIMA approach (Autoregressive Integrated Moving Average)
- VI. (FB) Prophet Approach
- VII. Keras Starter
- VIII. Comparison & Conclusion

## I. Importation & Data Cleaning

In this first part we will choose the Time Series to work in the others parts. The idea is to find a Time Serie who could be interesting to work with. So in our data we can find 145K Time Series. We will Find a good Time Series to introduce four approaches! So the first step is to import few libraries and the data. The four approaches are Basic Approach / ML Approach / GAM Approach / ARIMA Approach.

```
In [2]: import pandas as pd
import numpy as np
import warnings
import scipy
from datetime import timedelta

# Forecasting with decomposable model
from pylab import rcParams
import statsmodels.api as sm
from statsmodels.tsa.stattools import adfuller

# For machine Learning Approach
from statsmodels.tsa.stattools import lagmat
from sklearn.linear_model import LinearRegression, RidgeCV
from sklearn.ensemble import RandomForestRegressor
from sklearn.metrics import r2_score

# Visualisation
import matplotlib.pyplot as plt
import seaborn as sns

plt.style.use('fivethirtyeight')
warnings.filterwarnings('ignore')

In [4]: # Load the data
train = pd.read_csv("../input/train_1.csv")

In [5]: train_flattened = pd.melt(train[list(train.columns[-50:])+['Page']], id_vars='Page', var_name='date', value_name='Visits')
train_flattened['date'] = train_flattened['date'].astype('datetime64[ns]')
train_flattened['weekend'] = ((train_flattened.date.dt.dayofweek) // 5 == 1).astype(float)

In [6]: # Median by page
df_median = pd.DataFrame(train_flattened.groupby(['Page'])['Visits'].median())
df_median.columns = ['median']

# Average by page
df_mean = pd.DataFrame(train_flattened.groupby(['Page'])['Visits'].mean())
df_mean.columns = ['mean']

# Merging data
train_flattened = train_flattened.set_index('Page').join(df_mean).join(df_median)

In [7]: train_flattened.reset_index(drop=False, inplace=True)

In [8]: train_flattened['weekday'] = train_flattened['date'].apply(lambda x: x.weekday())

In [9]: # Feature engineering with the date
train_flattened['year'] = train_flattened.date.dt.year
train_flattened['month'] = train_flattened.date.dt.month
train_flattened['day'] = train_flattened.date.dt.day

In [10]: train_flattened.head()

This part allowed us to prepare our data. We had created new features that we use in the next steps. Days, Months, Years are interesting to forecast with a Machine Learning Approach or to do an analysis. If you have another idea to improve this first part: Fork this notebook and improve it or share your idea in the comments.
```

## II. Aggregation & Visualisation

```
In [11]: plt.figure(figsize=(50, 8))
mean_group = train_flattened[['Page', 'date', 'Visits']].groupby(['date'])['Visits'].mean()
plt.plot(mean_group)
plt.title('Time Series - Average')
plt.show()

In [12]: plt.figure(figsize=(50, 8))
median_group = train_flattened[['Page', 'date', 'Visits']].groupby(['date'])['Visits'].median()
plt.plot(median_group, color = 'r')
plt.title('Time Series - median')
plt.show()

In [13]: plt.figure(figsize=(50, 8))
std_group = train_flattened[['Page', 'date', 'Visits']].groupby(['date'])['Visits'].std()
plt.plot(std_group, color = 'g')
plt.title('Time Series - std')
plt.show()

In [14]: # For the next graphics
train_flattened['month_num'] = train_flattened['month']
train_flattened['month'].replace('11', '11 - November', inplace=True)
train_flattened['month'].replace('12', '12 - December', inplace=True)

train_flattened['weekday_num'] = train_flattened['weekday']
train_flattened['weekday'].replace(0, '01 - Monday', inplace=True)
train_flattened['weekday'].replace(1, '02 - Tuesday', inplace=True)
train_flattened['weekday'].replace(2, '03 - Wednesday', inplace=True)
train_flattened['weekday'].replace(3, '04 - Thursday', inplace=True)
train_flattened['weekday'].replace(4, '05 - Friday', inplace=True)
train_flattened['weekday'].replace(5, '06 - Saturday', inplace=True)
train_flattened['weekday'].replace(6, '07 - Sunday', inplace=True)

In [15]: train_group = train_flattened.groupby(['month', 'weekday'])['Visits'].mean().reset_index()
train_group = train_group.pivot('weekday', 'month', 'Visits')
train_group.sort_index(inplace=True)

In [16]: sns.set(font_scale=3.5)

# Draw a heatmap with the numeric values in each cell
f, ax = plt.subplots(figsize=(50, 30))
sns.heatmap(train_group, annot=False, ax=ax, fmt="d", linewidths=2)
plt.title('Web Traffic Months cross Weekdays')
plt.show()

This heatmap show us in average the web traffic by weekdays cross the months. In our data we can see there are less activity in Friday and Saturday for December and November. And the biggest traffic is on the period Monday - Wednesday. It is possible to do Statistics Test to check if our intuition is ok. But You have a lot of works!
```

```
In [17]: train_day = train_flattened.groupby(['month', 'day'])['Visits'].mean().reset_index()
train_day = train_day.pivot('day', 'month', 'Visits')
train_day.sort_index(inplace=True)
train_day.dropna(inplace=True)

In [18]: # Draw a heatmap with the numeric values in each cell
f, ax = plt.subplots(figsize=(50, 30))
sns.heatmap(train_day, annot=False, ax=ax, fmt="d", linewidths=2)
plt.title('Web Traffic Months cross days')
plt.show()

With this graph it is possible to see they are two periods with a bigger activity than the rest. The two periods are 25-29 December and 13-14 November. And we can see one period with little activity 15-17 December. They are maybe few outliers during these two periods. You must to investigate more. (coming soon...)
```

## III. ML Approach

The first approach introduces is the Machine Learnin Approach. We will use just a AdaBoostRegressor but you can try with other models if you want to find the best model. I tried with a linear model as like Ridge but ADA model is better. I will be interesting to check if GB or XGB can bit ADA. It is possible to do a Neural Network approach too. But this approach will be done if the kagglers want more !!

```
In [19]: times_series_means = pd.DataFrame(mean_group).reset_index(drop=False)
times_series_means['weekday'] = times_series_means['date'].apply(lambda x: x.weekday())
times_series_means['Date_str'] = times_series_means['date'].apply(lambda x: str(x))
times_series_means[['year', 'month', 'day']] = pd.DataFrame(times_series_means['Date_str'].str.split('-',
2).tolist(), columns = ['year', 'month', 'day'])
date_staging = pd.DataFrame(times_series_means['day'].str.split('-',2).tolist(), columns = ['day', 'othe
r'])
times_series_means['day'] = date_staging['day']+1
times_series_means.drop('Date_str', axis = 1, inplace=True)
times_series_means.head()

The first step for the ML approach is to create the feature that we will predict. In our example we don't predict the number of visits but the difference between two days. The tips to create few features is to take the difference between two days and to do a lag. Here we will take a lag of "diff" seven times. If you have a weekly pattern it is an interesting choice. Here we have few data (2 months so 30 values) and it is a constraint. I done some test and the number 7 is a good choice (weekly pattern?).

In [20]: times_series_means.reset_index(drop=True, inplace=True)

def lag_func(data, lag):
    lag = lag
    X = lagmat(data["diff"], lag)
    lagged = data.copy()
    for c in range(1, lag+1):
        lagged["lag%d" % c] = X[:, c-1]
    return lagged

def diff_creation(data):
    data = np.nan
    data.ix[1:, "diff"] = (data.iloc[1:, 1].as_matrix() - data.iloc[:len(data)-1, 1].as_matrix())
    return data

df_count = diff_creation(times_series_means)

# Creation of 7 features with "diff"
lag = 7
lagged = lag_func(df_count, lag)
last_date = lagged['date'].max()

In [21]: lagged.head()

In [22]: # Train Test split
x0 = ["lag%d" % i for i in range(1, lag+1)] + ['weekday'] + ['day']
split = 0.70
xt = data_lag[(lag+1):][x0]
yt = data_lag[(lag+1):]['diff']
isplit = int(len(xt) * split)
x_train, y_train, x_test, y_test = xt[:isplit], yt[:isplit], xt[isplit:], yt[isplit:]
return x_train, y_train, x_test, y_test, xt, yt

x_train, y_train, x_test, y_test, xt, yt = train_test(lagged)

In [35]: # Linear Model
from sklearn.ensemble import ExtraTreesRegressor, GradientBoostingRegressor, BaggingRegressor, AdaBoostR
egressor
from sklearn.metrics import mean_absolute_error, r2_score

def modelisation(x_tr, y_tr, x_ts, y_ts, xt, yt, model0, model1):
    # Modelisation with all product
    model0.fit(x_tr, y_tr)

    prediction = model0.predict(x_ts)
    r2 = r2_score(y_ts.as_matrix(), model0.predict(x_ts))
    mae = mean_absolute_error(y_ts.as_matrix(), model0.predict(x_ts))
    print ("-----")
    print ("mae with 70% of the data to train:", mae)
    print ("-----")

    # Model with all data
    model1.fit(xt, yt)

    return model1, prediction, model0

model0 = AdaBoostRegressor(n_estimators = 5000, random_state = 42, learning_rate=0.01)
model1 = AdaBoostRegressor(n_estimators = 5000, random_state = 42, learning_rate=0.01)

clr, prediction, clr0 = modelisation(x_train, y_train, x_test, y_test, xt, yt, model0, model1)

In [36]: # Performance 1
plt.style.use('ggplot')
plt.figure(figsize=(50, 12))
line_up = plt.plot(prediction, label='Prediction')
line_down = plt.plot(np.array(y_test), label='Reality')
plt.ylabel('Series')
plt.legend(handles=[line_up, line_down])
plt.title('Performance of predictions - Benchmark Predictions vs Reality')
plt.show()

In [37]: # Prediction
def pred_df(data, number_of_days):
    data_pred = pd.DataFrame(pd.Series(data["date"][data.shape[0]-1 + timedelta(days=1)], columns = ["date"])
    for i in range(number_of_days):
        inter = pd.DataFrame(pd.Series(data["date"][data.shape[0]-1 + timedelta(days=1+2)], columns = ["date"])
        data_pred = pd.concat([data_pred, inter]).reset_index(drop=True)
    return data_pred

data_to_pred = pred_df(df_count, 30)

In [38]: def initialisation(data_lag, data_pred, model, xtrain, ytrain, number_of_days):
    # Initialisation
    model.fit(xtrain, ytrain)

    for i in range(number_of_days-1):
        lag1 = data_lag.tail(1)['diff'].values[0]
        lag2 = data_lag.tail(1)['lag1'].values[0]
        lag3 = data_lag.tail(1)['lag2'].values[0]
        lag4 = data_lag.tail(1)['lag3'].values[0]
        lag5 = data_lag.tail(1)['lag4'].values[0]
        lag6 = data_lag.tail(1)['lag5'].values[0]
        lag7 = data_lag.tail(1)['lag6'].values[0]
        lag8 = data_lag.tail(1)['lag7'].values[0]

        data_pred['weekday'] = data_pred['date'].apply(lambda x:x.weekday())
        weekday = data_pred['weekday'][0]

        row = pd.Series([lag1, lag2, lag3, lag4, lag5, lag6, lag7, lag8, weekday],
        ,['lag1', 'lag2', 'lag3', 'lag4', 'lag5', 'lag6', 'lag7', 'lag8', 'weekday'])
        to_predict = pd.DataFrame(columns = ['lag1', 'lag2', 'lag3', 'lag4', 'lag5', 'lag6', 'lag7', 'lag8',
        'weekday'])
        prediction = pd.DataFrame(columns = ['diff'])
        to_predict = to_predict.append([row])
        prediction = pd.DataFrame(model.predict(to_predict), columns = ['diff'])

    # Loop
    if i == 0:
        last_predict = data_lag["Visits"][data_lag.shape[0]-1 + prediction.values[0][0]]

    if i > 0 :
        last_predict = data_lag["Visits"][data_lag.shape[0]-1 + prediction.values[0][0]]

        data_lag = pd.concat([data_lag, prediction.join(data_pred["date"]).join(to_predict)]).reset_index
        data_lag["Visits"][data_lag.shape[0]-1] = last_predict

    # test
    data_pred = data_pred[data_pred["date"]>data_pred["date"][0]].reset_index(drop=True)
    return data_lag

model_fin = AdaBoostRegressor(n_estimators = 5000, random_state = 42, learning_rate=0.01)

In [39]: lagged = initialisation(lagged, data_to_pred, model_fin, xt, yt, 30)

In [40]: lagged[lagged['diff']<0]
lagged.ix[(lagged.Visits < 0), 'Visits'] = 0

In [41]: df_lagged = lagged[['Visits', 'date']]
df_train = df_lagged[df_lagged['date'] <= last_date]
df_pred = df_lagged[df_lagged['date'] >= last_date]
plt.style.use('ggplot')
plt.figure(figsize=(30, 5))
plt.plot(df_train.date, df_train.Visits)
plt.plot(df_pred.date, df_pred.Visits, color='b')
plt.title('Training time series in red, Prediction on 30 days in blue -- ML Approach')
plt.show()

Finished for the first approach! The ML method requires a lot of work! You need to create the features, the data to collect the prediction, optimisation etc... This method done a good results when there are a weekly pattern identified or a monthly pattern but we need more data.
```

## IV. Basic Approach

For this model We will use a simple model with the average of the activity by weekdays. In general rules the simplest things give good results!

```
In [42]: lagged_basic = lagged[['date', 'Visits', 'weekday']]
lagged_basic_tr = lagged_basic[lagged_basic['date'] < last_date]
lagged_basic_pred = lagged_basic[lagged_basic['date'] >= last_date]
lagged_basic_pred.drop('Visits', inplace=True, axis=1)

In [43]: prediction_by_days = pd.DataFrame(lagged_basic.groupby(['weekday'])['Visits'].mean())
prediction_by_days.reset_index(drop=False, inplace=True)

In [44]: basic_pred = pd.merge(lagged_basic_pred, prediction_by_days, on='weekday')
basic_approach = pd.concat([lagged_basic_tr, basic_pred])

In [45]: plot_basic = np.array(basic_approach[basic_approach['date'] > last_date].sort_values(by='date').Visits)

In [46]: plt.figure(figsize=(30, 5))
plt.plot(plot_basic)
plt.title('Display the predictions with the Basic model')
plt.show()

In [47]: df_lagged = basic_approach[['Visits', 'date']].sort_values(by='date')
df_train = df_lagged[df_lagged['date'] <= last_date]
df_pred = df_lagged[df_lagged['date'] >= last_date]
plt.style.use('ggplot')
plt.figure(figsize=(30, 5))
plt.plot(df_train.date, df_train.Visits)
plt.plot(df_pred.date, df_pred.Visits, color='b')
plt.title('Training time series in red, Prediction on 30 days in blue -- ML Approach')
plt.show()

No optimisation! No choice between linear, Bagging, boosting or others! Just with an average by week days and we have a result! Fast and easily!
```

## V. ARIMA

This part is inspired by: <https://www.analyticsvidhya.com/blog/2016/02/time-series-forecasting-codes-python/> Very good job with the ARIMA models! It is more simple when we have directly a stationary Time series. It is not our case...

We will use the Dickey-Fuller Test. More informations here: [https://en.wikipedia.org/wiki/Dickey%E2%82%9393Fuller\\_test](https://en.wikipedia.org/wiki/Dickey%E2%82%9393Fuller_test)

```
In [48]: # Show Rolling mean, Rolling Std and Test for the stationarity
df_date_index = times_series_means[['date', 'Visits']].set_index('date')

def test_stationarity(tseries):
    plt.figure(figsize=(50, 8))
    #Determining rolling statistics
    roldmean = pd.rolling_mean(tseries, window=7)
    plotstd = pd.rolling_std(tseries, window=7)

    #Plot rolling statistics:
    orig = plt.plot(tseries, color='blue', label='Original')
    mean = plt.plot(roldmean, color='red', label='Rolling Mean')
    std = plt.plot(plotstd, color='black', label = 'Rolling Std')
    plt.legend(loc='best')
    plt.title('Rolling Mean & Standard Deviation')
    plt.show(block=False)

    #Perform Dickey-Fuller test:
    print('Results of Dickey-Fuller Test:')
    dftest = sm.tsa.adfuller(tseries['Visits'], autolag='AIC')
    dfoutput = pd.Series(dftest[0:4], index=['Test Statistic', 'p-value', 'Lags Used', 'Number of Observations Used'])
    for key, value in dfoutput.items():
        dfoutput[['Critical Value (%)', 'key']] = value
    print(dfoutput)

test_stationarity(df_date_index)

Our Time Series is stationary! It is a good news! We can to apply the ARIMA Model without transformations.

Good job! We have a Time Series Stationary! We can apply our ARIMA Model!!!

In [49]: # Naive decomposition of our Time Series as explained above
decomposition = sm.tsa.seasonal_decompose(df_date_index, model='multiplicative', freq = 7)

trend = decomposition.trend
seasonal = decomposition.seasonal
residual = decomposition.resid
rcParams['figure.figsize'] = 30, 20

plt.subplot(411)
plt.title('Observed = Trend + Seasonality + Residuals')
plt.plot(df_date_index, label='Observed')
plt.legend(loc='best')
plt.subplot(412)
plt.plot(trend, label='Trend')
plt.legend(loc='best')
plt.subplot(413)
plt.plot(seasonal, label='Seasonality')
plt.legend(loc='best')
plt.subplot(414)
plt.plot(residual, label='Residuals')
plt.legend(loc='best')
plt.tight_layout()
plt.show()

We expose the naive decomposition of our time series (More sophisticated methods should be preferred). They are several ways to decompose a time series but in our example we take a simple decomposition on three parts. The additive model is Y[t] = T[t] + S[t] + e[t]
The multiplicative model is Y[t] = T[t] x S[t] x e[t] with:

1. T[t]: Trend
2. S[t]: Seasonality
3. e[t]: Residual

An additive model is linear where changes over time are consistently made by the same amount. A linear trend is a straight line. A quadratic or exponential has the same frequency (width of cycles) and amplitude (height of cycles). A multiplicative model is nonlinear, such as quadratic or exponential. Changes increase or decrease over time. A nonlinear trend is a curved line. A non-linear seasonality has an increasing or decreasing frequency and/or amplitude over time. In our example we can see it is not a linear model. So it is the reason why we use a multiplicative model.
```

```
In [ ]: #from statsmodels.tsa.arima_model import ARIMA
#model = ARIMA(df_date_index, order=(7, 1, 0))
#results_AR = model.fit(disp=-1)
#plt.plot(df_date_index, df_pred.Visits, color='blue')
#plt.plot(results_AR.fittedvalues, color='red')
#plt.show()

In [ ]: #forecast = results_AR.forecast(steps = 30)[0]
#plt.figure(figsize=(30, 5))
#plt.plot(pd.DataFrame(np.exp(forecast)))
#plt.title('Display the predictions with the ARIMA model')
#plt.show()

In [ ]: # DataFrame to collect the predictions
df_prediction_arima = df_date_index.copy()

#list_date = []
#for i in range(31):
#    if i > 0:
#        list_date.append(last_date + pd.to_timedelta(i, unit='D'))

#predictions_arima = pd.DataFrame(list_date, columns = ['date'])
#predictions_arima[['Visits']] = 0
#predictions_arima.set_index('date', inplace=True)
#predictions_arima['Visits'] = np.exp(forecast)

#df_prediction_arima = df_prediction_arima.append(predictions_arima)
#df_prediction_arima.reset_index(drop=False, inplace=True)

In [ ]: #df_arima = df_prediction_arima[['Visits', 'index']]
#df_train = df_arima[df_arima['index'] <= last_date]
#df_pred = df_arima[df_prediction_arima['date'] >= last_date]
#plt.style.use('ggplot')
#plt.figure(figsize=(30, 5))
#plt.plot(df_train.index, df_train.Visits)
#plt.plot(df_pred.index, df_pred.Visits, color='b')
#plt.title('Training time series in red, Prediction on 30 days in blue -- ARIMA Model')
#plt.show()

IN PROGRESS...
```

## VI. Prophet

Prophet is a forecasting tool available in python and R. This tool was developed by Facebook. More information on the library here: <https://research.fb.com/prophet-forecasting-at-scale/>

Compared to the two methods this one will be faster. We can forecast a time series with few lines. In our case we will do a forecast and a display the trend of activity on the period and for a week.

```
In [60]: from fbprophet import Prophet
df_date_index = times_series_means[['date', 'Visits']]
df_date_index = df_date_index.set_index('date')
df_prophet = df_date_index.copy()
df_prophet.reset_index(drop=False, inplace=True)
df_prophet.columns = ['ds', 'y']

m = Prophet()
m.fit(df_prophet)
future = m.make_future_dataframe(periods=30, freq='D')
forecast = m.predict(future)
fig = m.plot(forecast)

In [51]: m.plot_components(forecast);

VI. Keras Starter
```

In this part we will use Keras without optimisation to forecast. It is just a very simple code to begin with Keras and a Time Series. For our example we will try just with one layer and 8 Neurons.

```
In [52]: df_dl = times_series_means[['date', 'Visits']]
train_size = int(len(df_dl) * 0.80)
test_size = len(df_dl) - train_size
train, test = df_dl.iloc[0:train_size,:], df_dl.iloc[train_size:len(df_dl),:]
print(len(train), len(test))

In [53]: look_back = 1

def create_dataset(dataset, look_back):
    dataX = []
    dataY = []
    for i in range(len(dataset)-look_back-1):
        a = dataset.iloc[i:(i+look_back), 1].values[0]
        b = dataset.iloc[i+look_back, 1]
        dataX.append(a)
        dataY.append(b)
    return np.array(dataX), np.array(dataY)

trainX, trainY = create_dataset(train, look_back)

In [54]: from keras.layers import Sequential
from keras.layers import Dense

trainX, trainY = create_dataset(train, look_back)
testX, testY = create_dataset(test, look_back)

model = Sequential()
model.add(Dense(8, input_dim=look_back, optimizer='relu'))
model.add(Dense(1))
model.compile(loss='mean_absolute_error', optimizer='adam')
model.fit(trainX, trainY, epochs=150, batch_size=2, verbose=0)

In [32]: trainScore = model.evaluate(trainX, trainY, verbose=0)
print('Train Score: %.2f MSE (%.2f MAE)' % (trainScore, trainScore))
testScore = model.evaluate(testX, testY, verbose=0)
print('Test Score: %.2f MSE (%.2f MAE)' % (testScore, testScore))

In [55]: trainPredict = model.predict(trainX)
testPredict = model.predict(testX)

# shift train predictions for plotting
trainPredictPlot = np.empty_like(df_dl)
trainPredictPlot[:, :] = np.nan
trainPredictPlot[look_back:len(trainPredict)+look_back, :] = trainPredict

# shift test predictions for plotting
testPredictPlot = np.empty_like(df_dl)
testPredictPlot[:, :] = np.nan
testPredictPlot[len(trainPredict)+(look_back*2)+1:len(df_dl)-1, :] = testPredict

# plot baseline and predictions
plt.plot(np.array(df_dl.Visits))
plt.plot(trainPredictPlot)
plt.plot(testPredictPlot)
plt.title('Prediction with Keras')
plt.show()
```

## VII. Comparison & Conclusion

In Progress...