

Linear Models on Time Series

<https://www.kaggle.com/kashnitsky/topic-9-part-1-time-series-analysis-in-python>

I have to build models with fast, good, cheap as my only guiding principle. That means that some of these models will never be considered "production ready" as they demand too much time for data preparation (as in SARIMA) or require frequent re-training on new data (again, SARIMA) or are difficult to tune (good example - SARIMA). Therefore, it's very often much easier to select a few features from the existing time series and build a simple linear regression model or, say, a random forest. It is good and cheap. This approach is not backed by theory and breaks several assumptions (e.g. Gauss-Markov theorem, especially for errors being uncorrelated), but it is very useful in practice and is often used in machine learning competitions.

I. Feature Extraction

The model needs features, and all we have is a 1-dimensional time series. What features can we extract?

- Lags of time series
- Window statistics:
 - Max/min value of series in a window
 - Average/median value in a window
 - Window variance
 - etc.
- Date and time features:
 - Minute of an hour, hour of a day, day of the week, and so on
 - Is this day a holiday? Maybe there is a special event? Represent that as a boolean feature
- Target encoding
- Forecasts from other models (note that we can lose the speed of prediction this way)

Lags of Time Series

Shifting the series n steps back, we get a feature column where the current value of time series is aligned with its value at time $t-n$. If we make a 1 lag shift and train a model on that feature, the model will be able to forecast 1 step ahead from having observed the current state of the series. Increasing the lag, say, up to 6, will allow the model to make predictions 6 steps ahead; however it will use data observed 6 steps back. If something fundamentally changes the series during that unobserved period, the model will not catch these changes and will return forecasts with a large error. Therefore, during the initial lag selection, one has to find a balance between the optimal prediction quality and the length of the forecasting horizon.

```
# Creating a copy of the initial datagrame to make various transformations
data = pd.DataFrame(ads.Ads.copy())
data.columns = ["y"]
```

```
# Adding the lag of the target variable from 6 steps back up to 24
for i in range(6, 25):
    data["lag_{}".format(i)] = data.y.shift(i)
```

hour, day of week, is_weekend (Boolean) features

```
data.index = pd.to_datetime(data.index)
data["hour"] = data.index.hour
data["weekday"] = data.index.weekday
data["is_weekend"] = data.weekday.isin([5,6])*1
data.tail()
```

Target Encoding

another variant for encoding categorical variables: encoding by mean value. If it is undesirable to explode a dataset by using many dummy variables that can lead to the loss of information and if they cannot be used as real values because

of the conflicts like "0 hours < 23 hours", then it's possible to encode a variable with slightly more interpretable values. The natural idea is to encode with the mean value of the target variable. In our example, every day of the week and every hour of the day can be encoded by the corresponding average number of ads watched during that day or hour. It's very important to make sure that the mean value is calculated over the training set only (or over the current cross-validation fold only) so that the model is not aware of the future.

Target encoding might lead to some overfitting... then we can calculate the target encoding not for the whole train set, but for some window instead to solve this issue.

```
def code_mean(data, cat_feature, real_feature):  
    """  
    Returns a dictionary where keys are unique categories of the cat_feature,  
    and values are means over real_feature  
    """  
    return dict(data.groupby(cat_feature)[real_feature].mean())
```

Function that does all three feature extraction

```
def prepareData(series, lag_start, lag_end, test_size, target_encoding=False):  
    """  
    series: pd.DataFrame  
            dataframe with timeseries  
  
    lag_start: int  
            initial step back in time to slice target variable  
            example - lag_start = 1 means that the model  
            will see yesterday's values to predict today  
  
    lag_end: int  
            final step back in time to slice target variable  
            example - lag_end = 4 means that the model  
            will see up to 4 days back in time to predict today  
  
    test_size: float  
            size of the test dataset after train/ test split as percentage of dataset  
  
    target_encoding: boolean  
            if True - add target averages to the dataset  
  
    """  
  
    # copy of the initial dataset  
    data = pd.DataFrame(series.copy())  
    data.columns = ["y"]  
  
    # lags of series  
    for i in range(lag_start, lag_end):  
        data["lag_{}".format(i)] = data.y.shift(i)  
  
    # datetime features  
    data.index = pd.to_datetime(data.index)  
    data["hour"] = data.index.hour  
    data["weekday"] = data.index.weekday  
    data["is_weekend"] = data.weekday.isin([5,6])*1  
  
    if target_encoding:  
        # calculate averages on train set only  
        test_index = int(len(data.dropna()*(1-test_size))  
        data["weekday_average"] = list(map(code_mean(data[:test_index], 'weekday', "y").get, data.weekday))
```

```

data["hour_average"] = list(map(code_mean(data[:test_index], 'hour', "y").get, data.hour))

# drop encoded variables
data.drop(["hour", "weekday"], axis=1, inplace=True)

# train-test split
y = data.dropna().y
X = data.dropna().drop(['y'], axis=1)
X_train, X_test, y_train, y_test = timeseries_train_test_split(X, y, test_size=test_size)

return X_train, X_test, y_train, y_test

```

II. Modelling

Linear Regression

```

from sklearn.linear_model import LinearRegression
from sklearn.model_selection import cross_val_score
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()

# for time-series cross-validation set 5 folds
tscv = TimeSeriesSplit(n_splits=5)

~~Evaluation Metrics~~
from sklearn.metrics import r2_score, median_absolute_error, mean_absolute_error
from sklearn.metrics import median_absolute_error, mean_squared_error, mean_squared_log_error

def mean_absolute_percentage_error(y_true, y_pred):
    return np.mean(np.abs((y_true - y_pred) / y_true)) * 100

```

```

def timeseries_train_test_split(X, y, test_size):
    """
        Perform train-test split with respect to time series structure
    """

    # get the index after which test set starts
    test_index = int(len(X)*(1-test_size))

    X_train = X.iloc[:test_index]
    y_train = y.iloc[:test_index]
    X_test = X.iloc[test_index:]
    y_test = y.iloc[test_index:]

    return X_train, X_test, y_train, y_test

```

```

def plotModelResults(model, X_train=X_train, X_test=X_test, plot_intervals=False, plot_anomalies=False):
    """
        Plots modelled vs fact values, prediction intervals and anomalies
    """

    prediction = model.predict(X_test)

    plt.figure(figsize=(15, 7))
    plt.plot(prediction, "g", label="prediction", linewidth=2.0)
    plt.plot(y_test.values, label="actual", linewidth=2.0)

    if plot_intervals:

```

```

cv = cross_val_score(model, X_train, y_train,
                      cv=tscv,
                      scoring="neg_mean_absolute_error")

mae = cv.mean() * (-1)
deviation = cv.std()

scale = 1.96
lower = prediction - (mae + scale * deviation)
upper = prediction + (mae + scale * deviation)

plt.plot(lower, "r--", label="upper bond / lower bond", alpha=0.5)
plt.plot(upper, "r--", alpha=0.5)

if plot_anomalies:
    anomalies = np.array([np.NaN]*len(y_test))
    anomalies[y_test<lower] = y_test[y_test<lower]
    anomalies[y_test>upper] = y_test[y_test>upper]
    plt.plot(anomalies, "o", markersize=10, label = "Anomalies")

error = mean_absolute_percentage_error(prediction, y_test)
plt.title("Mean absolute percentage error {0:.2f}%".format(error))
plt.legend(loc="best")
plt.tight_layout()
plt.grid(True);

def plotCoefficients(model):
    """
    Plots sorted coefficient values of the model
    """

    coefs = pd.DataFrame(model.coef_, X_train.columns)
    coefs.columns = ["coef"]
    coefs["abs"] = coefs.coef.apply(np.abs)
    coefs = coefs.sort_values(by="abs", ascending=False).drop(["abs"], axis=1)

    plt.figure(figsize=(15, 7))
    coefs.coef.plot(kind='bar')
    plt.grid(True, axis='y')
    plt.hlines(y=0, xmin=0, xmax=len(coefs), linestyle='dashed');

y = data.dropna().y
X = data.dropna().drop(['y'], axis=1)

# reserve 30% of data for testing
X_train, X_test, y_train, y_test = timeseries_train_test_split(X, y, test_size=0.3)

X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)

lr = LinearRegression()
lr.fit(X_train_scaled, y_train)

plotModelResults(lr, X_train=X_train_scaled, X_test=X_test_scaled, plot_intervals=True)
plotCoefficients(lr)

```

OR using prepareData function above....

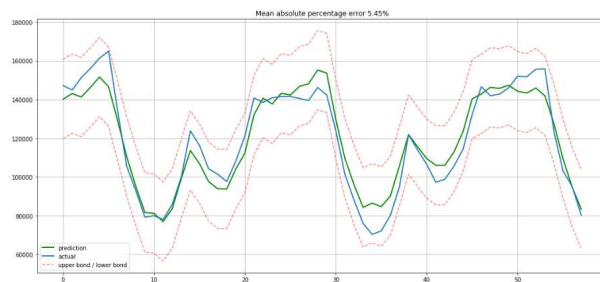
```
X_train, X_test, y_train, y_test = prepareData(ads.Ads, lag_start=6, lag_end=25, test_size=0.3,
target_encoding=True)
```

```
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)
```

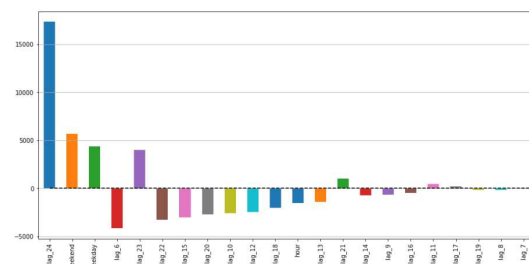
```
lr = LinearRegression()
lr.fit(X_train_scaled, y_train)
```

```
plotModelResults(lr, X_train=X_train_scaled, X_test=X_test_scaled, plot_intervals=True,
plot_anomalies=True)
plotCoefficients(lr)
```

actual v.s. prediction



coefficient bar graphs



Regularized Regression (L1, L2)

~~all the variables and functions below (e.g. X_train etc.) are the continuation from linear regression above~~

```
plt.figure(figsize=(10, 8))
sns.heatmap(X_train.corr());
```

```
from sklearn.linear_model import LassoCV, RidgeCV
```

```
ridge = RidgeCV(cv=tscv)
ridge.fit(X_train_scaled, y_train)
```

```
plotModelResults(ridge,
X_train=X_train_scaled,
X_test=X_test_scaled,
plot_intervals=True, plot_anomalies=True)
plotCoefficients(ridge)
```

```
lasso = LassoCV(cv=tscv)
lasso.fit(X_train_scaled, y_train)
```

```
plotModelResults(lasso,
X_train=X_train_scaled,
X_test=X_test_scaled,
plot_intervals=True, plot_anomalies=True)
plotCoefficients(lasso)
```

XGBoost

Generally, tree-based models handle trends in data poorly when compared with linear models. In that case, you would have to detrend your series first or use some tricks to make the magic happen. Ideally, you can make the series stationary and then use XGBoost. For example, you can forecast trend separately with a linear model and then add predictions from xgboost to get a final forecast.

```
from xgboost import XGBRegressor
```

```
xgb = XGBRegressor()
```

```
xgb.fit(X_train_scaled, y_train)
```

```
plotModelResults(xgb,  
                  X_train=X_train_scaled,  
                  X_test=X_test_scaled,  
                  plot_intervals=True, plot_anomalies=True)
```