

Workflow: preprocessing and modelling of all data

Workflow - all completed using R

Information

Completed
Optional
Figure

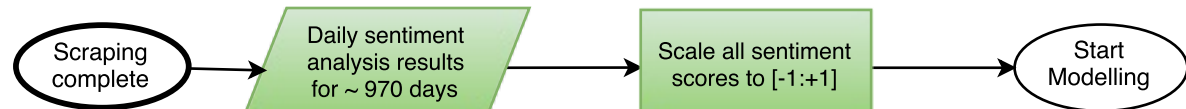


Figure:
Number of tweets

Figure:
Tweet frequency vs. Dow Jones

Collect macro data and combine with sentiment data

Create dummy variables for Mondays and bank holidays

Remove weekends from all data sets

Imputation

Spline
cubic splines

Don't remove
impute later using:
kNN, bagging,
predictive modelling

Last observation carried forward

Create the lagged data sets.
Lag values: 1 to 5

Create data subsets:
Macro, SA (all & avg.) and Mixed

Prune subsets according to pairwise correlation (>90%)

A

Execute model for each subset and each number of lags

Macro data includes a total of 29 covariates. Areas covered include indices (7), commodities (7), energy (2), currency pairs (5), interest rates (6) & market volatility measures (2).

We have sentiment data for weekends, but not for market data. Use of dummy variables may aid the incorporation of the weekend data.

These are the main source of missing data. With weekends removed < 0.1% of data is missing.

All data imputed this way has been created and could still be used. currently left out due to low number and therefore impact of missing data.

Sequential remove predictors with the highest **cumulative** pairwise correlation. The cumulative sums are recalculated after each removal.

Completed so far using the parameter pairs:
{mstop = 5000 & nu = 0.005} - slow descent
{mstop = 2000 & nu = 0.05} - fast(er) descent.

