

Modelling search volumes as a dynamic system responding to external events

Stefan Sabev

Master of Informatics

School of Informatics
University of Edinburgh

2014

Abstract

It is well known that some events might spark people's interest to fly to different destinations. In particular news events or sports events can quite easily make people search for a specific destination - for example the Champions League Quarter final draw increased the number of flight searches from Glasgow to Spain 6 times.

For this project we have collected vast amounts of Twitter data. With this dataset and the flight search dataset provided by Skyscanner it was possible to build a classifier that predicts flight search demand based on what's happening on Twitter. This is a noble approach to predicting flight search volumes utilising the vastness of Social data available.

The potential applications of this are generic prediction of flight search volumes, predicting new events for better marketing and also anomaly detection in traffic caused by events.

Table of Contents

1	Introduction	3
1.1	Synopsis of results	5
1.2	Results for the first iteration of the algorithm	6
1.3	Results for the 2nd iteration	8
1.4	Results for the algorithm with more features	9
1.5	Results on the full dataset	12
1.5.1	With just Twitter counts	12
1.6	Conclusion	12
2	Background and discussion of related work	15
3	Methodology	19
3.1	Data collection	20
3.2	Data processing	21
3.3	Hashtags	21
3.4	Occurrences paired with travel terms	22
3.5	Finding the correlation between those variables	22
3.6	Cleansing the data	24
4	Models	27
4.1	The baseline model	27
4.2	Last 4 Fridays + Twitter Counts	27
4.3	Results	29
4.4	Future improvements	29
5	Future work	31
	Bibliography	33

Chapter 1

Introduction

There aren't that many flight search companies with the global outreach required to collect and aggregate massive volumes of flight search data. And there are even fewer ones that are making some or the whole of their datasets available for research. When you consider that and the fact that online travel is a niche area in itself, one starts to understand why a project of this kind would have been quite hard if not impossible to do up until this date. However as some of those companies grow, they are trying to employ more sophisticated ways of predicting demand and as a result of that their plan capacity or just use that data for revenue forecasting. With the rise of social media, the next logical step would be to start exploring different ways of adding exogenous factors such as twitter data into their models to better their prediction.

In this project you will find the first attempt that someone has made to try using and incorporating social media to predict flight searches. As a consequence of that the effect of social media over online travel can be measured. We will try to show that using this dataset (~480GB at the time of writing, 500 million tweets) we can employ text mining and try to build a classifier that can predict upward and downward shifts in demand for a group of destinations - countries and cities as of the first part. With NLP we will hopefully gain better understanding on how to predict demand and therefore this could be useful not just to flight search companies, but also to airlines and airports.

It will be impossible not to say anything about Skyscanner, since they have kindly provided the flight search data, which I am using in this project.¹ It's an Edinburgh-based company with offices around the whole world, which is slowly but surely positioning itself as a leader in travel search online. It is in the phase of rapid growth and it has been doubling in size and valuation for the past couple of years, which helped secure an investment from Sequioa Capital. [1]

Due to the fact it processes billions of user searches on a monthly basis it is quite important to be able to predict roughly how tomorrow is going to like like - it affects capacity, website performance and costs in general. This has made Skyscanner transform from a flight search company into a data-driven company. With "Big data" comes big responsibility - you need to get the maximum value out of your data rather than

¹**Disclaimer:** I am an employee of the company and have been working there since 2011

just store it on disk. In order to utilise all of it they have developed and deployed a lot of in-house algorithms for anomaly detection, prediction and any other KPI (Key Performance Indicator) that can be measured. However there are still some things that you can't predict from just looking at past data and learning from it.

This particular project was spurred after a discussion that social media could be harnessed to produce some form of aggregation of counts that will allow us to see whether there are going to be any expected spikes to any destinations. Of course that will not be a "one size fits all" approach, since some places such as London, Paris and New York will see steady very high number of search volumes to them. The destinations that are more likely to be better predicted by the classifier with Twitter are the ones that are not constant all year round - Ibiza, Alicante, Dubrovnik and many others. Those destinations have a particular seasonality with spikes around holidays and some events, which the Twitter counts will hopefully capture.

The way I envisioned it in the beginning was to mine the vast quantities of Twitter data in a particular way (all ways described in chapter 3 on page 19) and then take the numbers aggregated on a daily level to produce what we will call in this dissertation the "Twitter counts" for a particular city or overall. The in-house classifier for predicting the searches will hopefully be beaten by the new improved classifier with Twitter factored in. The approach described here could be used to develop a system that monitors the social streams and can be used as a weight in a system of classifiers for predicting the future based on Twitter.

In the first part of the project I have decided to be as pragmatic as possible and explored the most practical of all potential uses - predicting search volumes for every city or country that has been mentioned on Twitter with the the counts and a list of features that are automatically extracted from the public stream. Then based on those numbers a classifier is trained on historical data and then we will also make a prediction for the future.

When the project is extended to be a more real-time system in its second phase, I expect that with the help of some TDT the potential applications of the work here are going to find many applications. In theory, one can easily plug in to a social network public stream and monitor what is happening and what is trending and as soon as something relevant is seen - event in a different country/city or in general something happening in a country one can action in a multitude of ways:

- "X is happening, why not fly there to see it?"
- Develop it further to monitor social media in order to predict spikes in traffic. Appearance on Italian TV caused half an hour outage in 2011!
- If cross-referenced with a capable Customer Relationship Management system you could use the methods described in this project and the results as a data source for a real time marketing solution

The incredible breadth, the opportunity to work with big datasets and the fact that I could work on something novel were the main reasons that made me pick this particular topic for research in my project. It's a very interesting mixture of Machine Learning

and Natural language processing. It has also presented a brilliant opportunity to learn how real life ML and NLP can be used and applied to a big data set and what are the potential benefits of doing so. Of course, it's worth mentioning that there is also the practical aspect of building a solution that could be used to power a solution for predicting search volumes and also anomaly detections with some slight tweaks.

Note - All the aggregated datasets are made available in the GitHub repository. [2] The full Twitter dataset will be made available on request. Of course because of confidentiality I can't publish the dataset which has the Skyscanner searches, since that might break my employee code of conduct. Instead I have anonymised it by normalising as described in the methodology.

1.1 Synopsis of results

As hypothesised, there are two types of destinations:

- One that have a constant search volume to them which doesn't change too much over the course of the year - those are the big cities such as London, New York, Paris and to some extent the searches to some countries.
- Smaller places such as Ibiza, Alicante, Faro and more exotic destinations such as Lusaka, Tobago and many others that have highly seasonal demand.

The difference in the two is graphically represented on the next page in Figure 1.1 and 1.2.

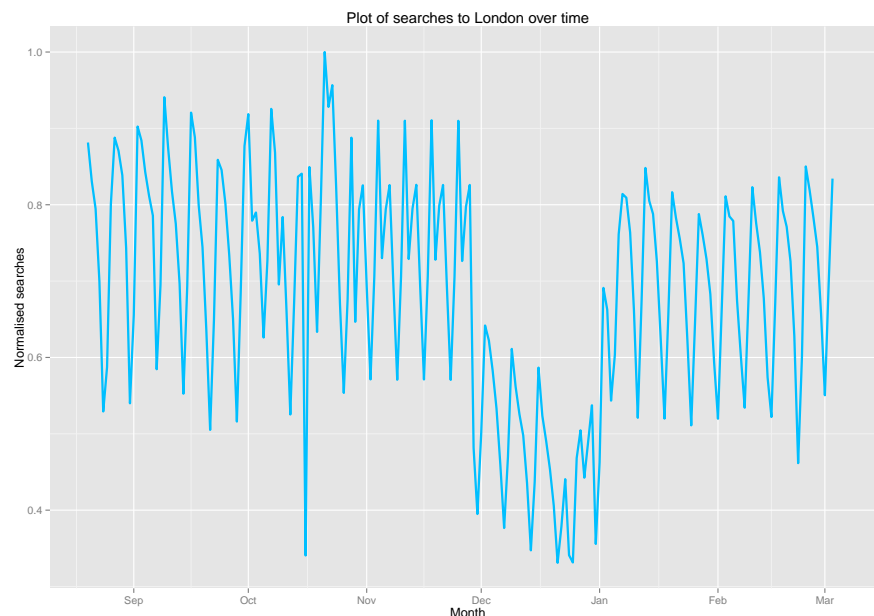


Figure 1.1: Searches to London - as we can see the trundling is going to be almost parallel to the X axis. The slight dip in December is due to seasonality. That is noticed across the industry as a whole

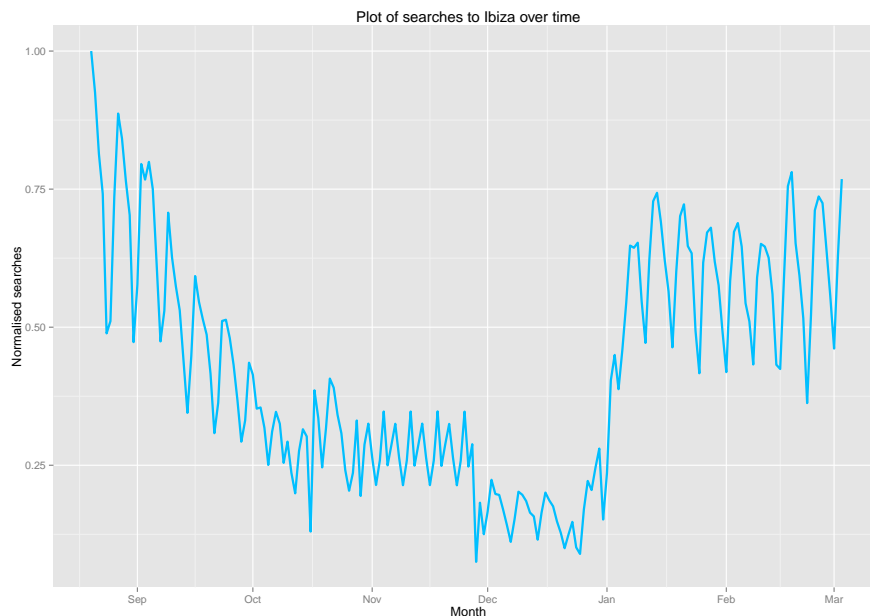


Figure 1.2: Searches to Ibiza - in this plot we can get the idea of seasonality. The dips is observed as soon as enter autumn and then in January the searches to such destinations pick up again as people are starting to plan their summer holidays

Naturally, because of the constant and not so changing nature of the first group my classifier is expected to perform better on the second group. The Root Mean Squared Error from my models on the 2nd group should be smaller or roughly the same as the error generated by the in-house algorithm. It's called "Last 4 Fridays" and essentially it is a slightly tweaked version of an exponentially weighted moving average, which takes into account seasonality. A more detailed discussion as to why this was picked as the baseline can be found in Chapter 4.1 on page 27.

1.2 Results for the first iteration of the algorithm

In order to get started with this problem and generate a first set of results I had to identify a good way of modelling the problem and how to predict the output from some inputs. The first things suggested by my supervisor was to use LASSO [3] in order to determine the weights for all of my inputs. The way the first model works is the following - you still have all the last 4 Fridays, but the Twitter counts are also included. Instead of fixing their weights, I chose to let LASSO determine all the weights for the inputs. This way we get an incremental improvement over the Last 4 Fridays, that is taking into account one of the many exogenous factors that could influence the search volumes either upwards or downwards. I expected that to yield better results for the 2nd group of destinations and possibly for some from group 1. I called the model TwitterDF.

Overall, I was pleasantly surprised by the results. On the full set of 1372 destinations, the first iteration of my model generated better predictions and scored lower RMSE

on 1034 of the destinations. On the other hand the in-house champion Last 4 Fridays performed better on 338 of the destinations. In order to illustrate that difference there is a scatter plot of the RMSEs in Figure 1.3. That means that my model is better on 62% of all the destinations.

In it we can see that there are some destinations where the Twitter model performs much worse than the L4F one (dots far below $y=x$ line). The ones, which are better classified by my model are destinations with a low to medium number of flight searches, since these are the destinations where exogenous factors, such as events can drastically increase or even decrease demand for flights.

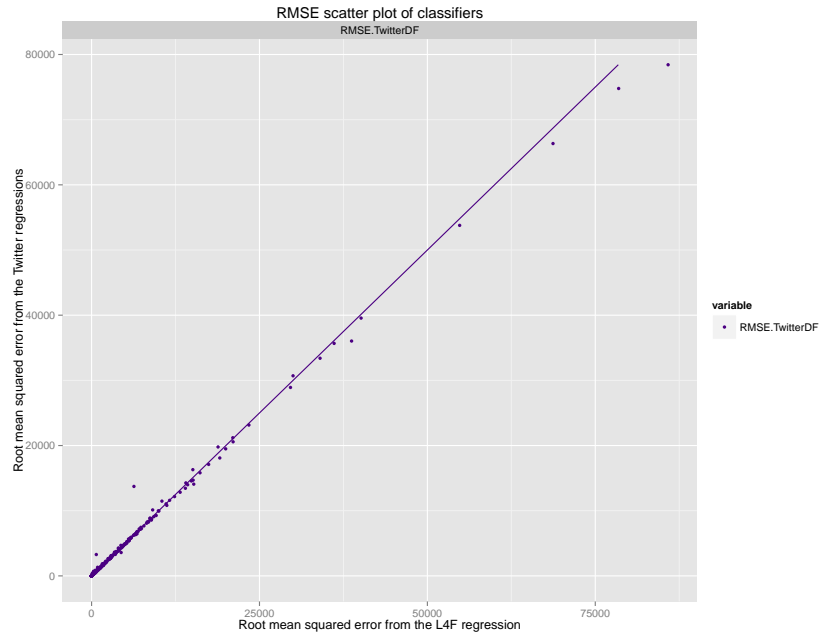


Figure 1.3: Scatter plot of RMSE Twitter Model and RMSE Last 4 Fridays. To put things into perspective, I've also plotted the $y=x$ line to make it easier to visualise the border between the two. Anything above the line is where my classifier performs better and anything below is where the in-house model excels.

In Table 1.1 and 1.2 you can find some picked which illustrate the worst performing classifiers and best performing ones.

If we look at the top 10 destinations in terms of RMSE (and in flight searches) in Table 1.3 we see something quite surprising, but also very nice. The improvement is positive in 9 out of the 10 examples. Even though it's small, the very fact that there is an improvement, means that taking into account exogenous factors is quite important and can contribute positively to the quality of the prediction.

Place	RMSE L4F	RMSE Twitter	Improvement
Parkersburg	3.99	1.82	54.38%
Bismarck	33.41	19.53	41.54%
Marshall Islands	41.63	24.99	39.98%
Cayo Coco	70.44	44.10	37.38%
Yakima	7.44	4.94	33.65%
Jonesboro	3.40	2.34	31.11%
Casper	20.30	14.14	30.34%
Barrow	5.28	3.76	28.81%
Niue	12.50	9.02	27.82%
Cedar City	5.08	3.71	26.91%
Sudbury	8.89	6.56	26.24%

Table 1.1: Results for destinations that yielded the best improvement in percentage terms when classified with the model from this project in comparison to the in-house one. As we can see they are fairly small niche destinations where it seems that people tweeting about it likely to result in a flight search to that place.

Place	RMSE L4F	RMSE Twitter	Improvement
Hayman Island	0.12	3.06	-2511.33%
Airlie Beach	0.65	10.17	-1456.42%
Venezuela	701.98	3277.16	-366.85%
Crooked Creek	0.20	0.93	-362.78%
Pereira	74.94	248.77	-231.97%
Launceston	122.28	404.13	-230.49%
Jodhpur	30.80	74.05	-140.40%
Manaus	295.66	707.50	-139.30%
Del Rio	5.89	13.42	-127.63%
Trondheim	257.36	573.07	-122.68%

Table 1.2: Results for destinations that yielded the biggest negative improvement. Quite an interesting mix destinations. The Venezuela result is attributed to Venezuela trending lately, because of the civil uprisings. Eliminating such "negative" influences is one of my goals for next year.

1.3 Results for the 2nd iteration

The 2nd iteration, was to fix the weights of the last 4 fridays and use the exponential weighting scheme that we use in the L4F classifier. I called that model TwitterCF and it yielded quite interesting results. The scatter comparison for the two models is in Figure 1.4. In Table 1.4 is the overall comparison of which classifier has the lowest RMSE for how many of the total classifiers.

Funnily enough, it's not always the case that the Twitter CF classifier wins over Twitter DF.

Place	RMSE L4F	RMSE TwitterDF	Improvement
Spain	85843.77	78424.46	8.64 %
United States	78479.89	74782.49	4.71%
United Kingdom	68705.88	66334.87	3.45%
Italy	54803.71	53779.50	1.87%
London	40128.75	39553.87	1.43%
Russia	38711.83	36033.21	6.92%
Germany	36113.08	35671.83	1.22%
France	34020.12	33393.17	1.84%
Thailand	29997.31	30700.72	-2.34%
Turkey	29616.37	28912.17	2.38%

Table 1.3: Results for the top destinations by RMSE error of L4F. 9/10 have lower RMSEs with the Twitter model. Using the compound weights has played a drastically big role in decreasing the prediction error.

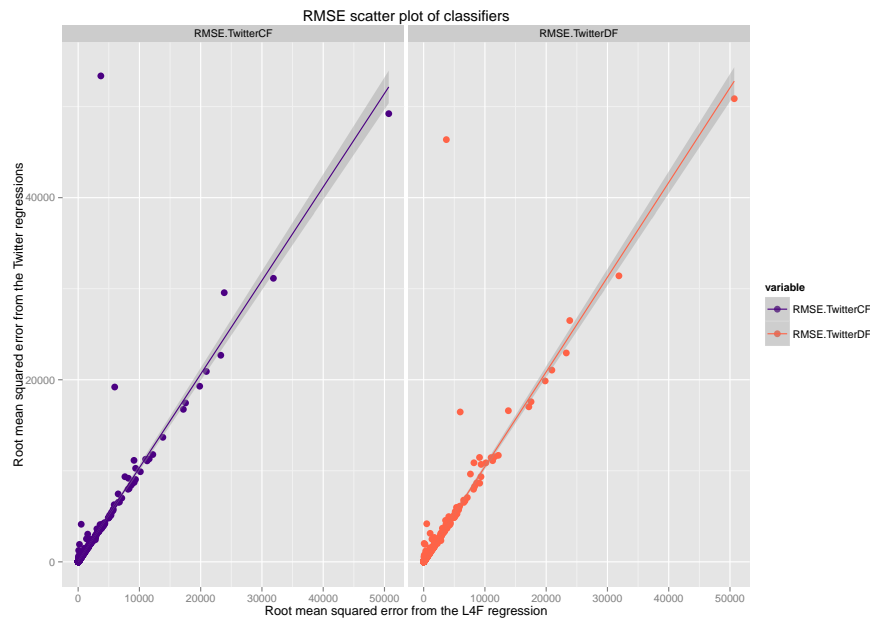


Figure 1.4: The two models compared. TwitterCF stands for compound Friday - taking all of the previous Fridays with predetermined weights and using them as a single weight in the regression. TwitterDF is leaving everything to LASSO.

1.4 Results for the algorithm with more features

In the next iteration of the model I didn't limit it to just twitter counts and searches from previous dates. For this particular model I've used all the words co-occurring with a place name as features. Of course, there is a lot of junk, so I removed all the emoji characters and all the stop words in order to prevent overwhelming the LASSO with a lot of junk features.

In Table 1.5 you can see the improvement in RMSE of the classifier for different values

Method	NumBest
L4F	226
TwitterCF	741
TwitterDF	405
Total	1372

Table 1.4: Comparison of the three methods and who had the best results on the dataset of 1372 points. As we can see, the Twitter classifiers have got a definite majority. 83%+ of the results have the lowest RMSE with the Twitter classifiers

of the penalisation factor α , the number of non-zero weights and also the RMSE for L4F, which is constant, since that is not tweaked. Overall we can see quite a significant improvement in the error as we discard more and more junk features, but it's still not as good as the simple model that is in the first version.

As we can see the more the junk features are penalised the classification improves significantly. However that is still just the beginning. Next year I will apply more NLP on the Twitter dataset trying to capture only the relevant tweets that will result in an improvement of the overall RMSE for the classifier.

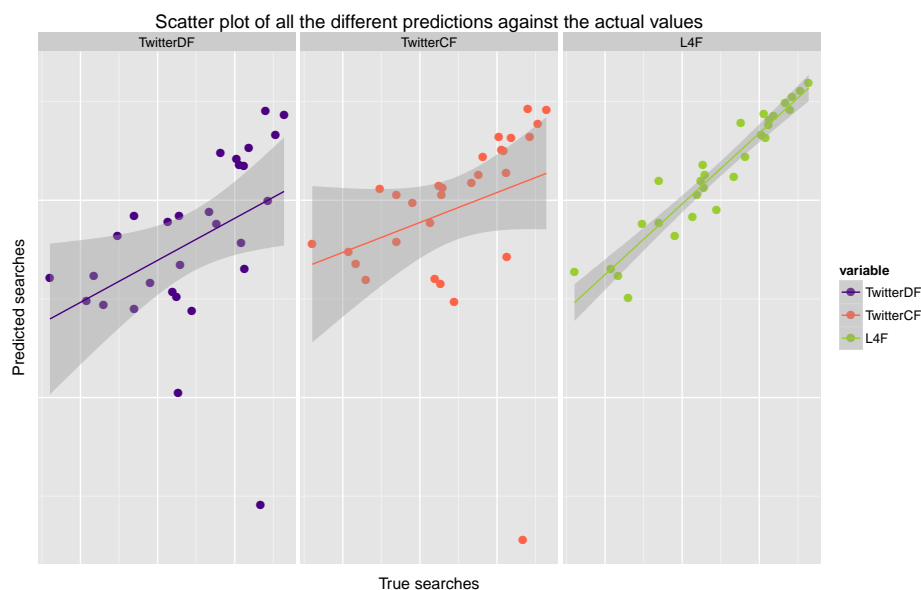


Figure 1.5: Predictions vs actual for a destination from group 1 - London. As discussed previously for this group Twitter models are not as good.

On figures 1.5 and 1.6 one can see the predictions against the actual searches. As we can see the TwitterCF model works

Aplha	> 0 W CF	RMSE TCF	> 0 W DF	RMSE TDF	RMSE L4F
0.5	180	12504.91	192	36323.13	1296.62
1	162	13288.62	160	30995.75	1296.62
2	136	9398.84	142	22176.84	1296.62
5	125	12108.62	121	12817.67	1296.62
10	110	8337.52	102	8254.31	1296.62
20	99	7558.71	87	5367.20	1296.62
50	74	5231.49	69	4232.36	1296.62
125	43	3255.71	41	3681.78	1296.62
250	24	2610.64	32	3230.88	1296.62
500	17	2334.97	15	3513.25	1296.62
1000	8	1742.20	10	3662.42	1296.62
2000	4	1231.04	8	3541.44	1296.62
4000	3	1231.93	6	3446.65	1296.62
8000	2	1237.07	5	3288.45	1296.62
16000	2	1240.17	4	3201.83	1296.62
32000	1	1247.21	4	3198.78	1296.62

Table 1.5: Improvement in RMSE with the expanded feature set for Palma, which is in 2nd group of destinations. As you can see in terms of absolute value, the RMSE for the Twitter classifier with compound fridays - RMSE TCF - has decreased 10.7 time. At $\alpha = 2000$ the compound friday classifier beats Twitter with 4 features (horizontal line). The features with non zee weights decrease from 180 down to 4 to make the better prediction. The RMSE for the classifier with all weights picked by Lasso - RMSE TDF- has decreased 11.3 times, however that is still not as good as the L4F predictor.

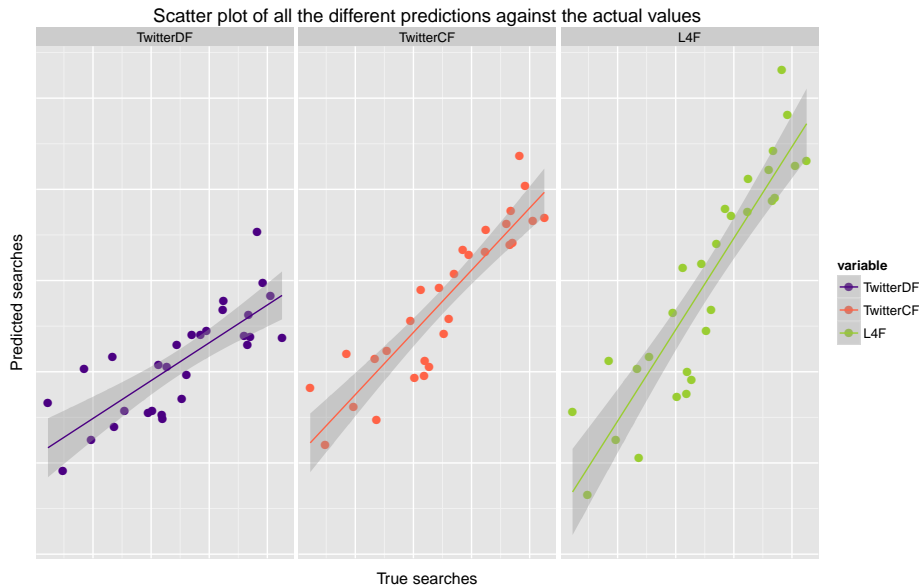


Figure 1.6: Predictions vs actual for a destination from group 2 - Lanzarote. For this group Twitter models make actual gains, albeit small.

1.5 Results on the full dataset

1.5.1 With just Twitter counts

Another interesting benchmark was how it performs on the overall search volumes.

	RMSE TwitterDF	Twitter CF	RMSE L4F	Difference
Overall	1 640 116.72	1 621 390.54	1 674 129.97	3.1 %

Table 1.6: RMSE on the tweet counts for all destinations and searches to all destinations. Difference is the percentage improvement of the best of the twitter classifiers.

As you can see in Figure 1.5 the overall volume of searches is not as volatile as the destinations in group 2, so it would fall into group 1, but surprisingly enough the Twitter classifiers do a marginally better job at predicting.

1.6 Conclusion

The previous results have demonstrated that there is definitely a gain, even though not as great as I expected, in including exogenous factors such as Twitter in the prediction. With the first level of filtering described in chapter 3.4 on page 22 we have achieved better prediction on slightly more than 83% of all the destinations.

During the 2nd phase of the project, I will focus on retrieving only the most relevant tweets from the Twitter dataset and doing better feature extraction that will allow to

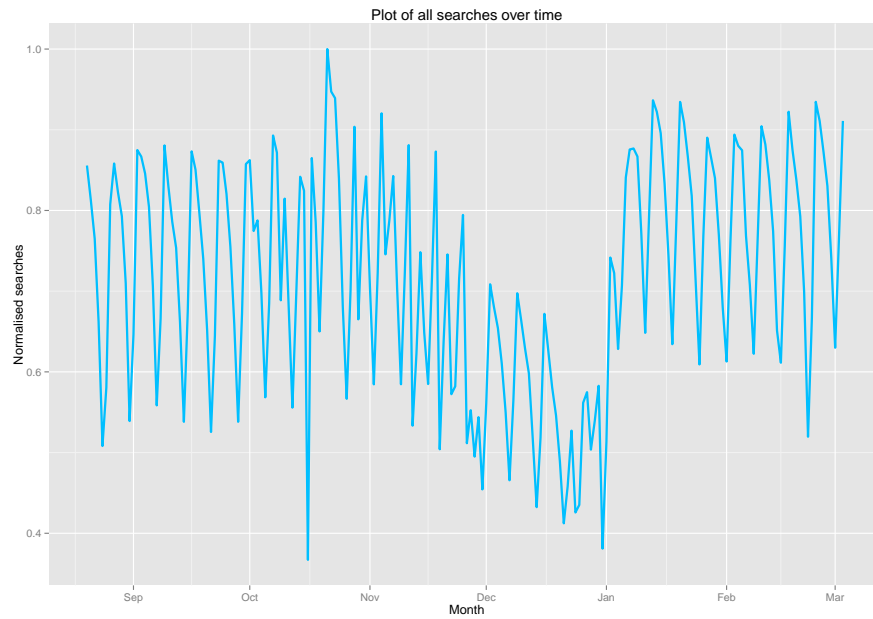


Figure 1.7: Plot of all the searches to all the destinations. This would fall into the first group of steady constant volumes

remove junk features that bring the RMSE up. With that in mind, I am hopeful that I will be able to build the next version of the model that is going to perform just as well on the first group of destinations.

Chapter 2

Background and discussion of related work

With millions and some with billions of users, Social Networks are becoming an increasingly more important in our lives. A big proportion of people use it as their primary way of communicating with the outside world - people will "tweet" about anything, post to Facebook, check in on FourSquare, instagram their food and so on. We have become perfectly fine with externalising our lives and posting every minute detail about our lives online and thusly making it available to everyone else. Of course, there are always exceptions. After all 82% of the world's population are not on Facebook, but those 1.23 billion monthly active users are all using it. Not all of them are avid users and post all of their pictures to it, but a big proportion is. With this amount of data on one's behaviour, life patterns and activities, companies can build a very good understanding of every individual.

As with any other network - be it TV, radio, newspaper - there are always going to be some who are trying their best to commercialise it and benefit from it in some way. Since people are willing to share so much information about themselves online, there are terabytes of data being generated every day on what people did, what things they tweeted about, their latest pictures, etc. As we are entering the "Big Data" age where every company is trying to turn its data into a product or simply establish itself as the leading data supplier for a particular market. Due to the vastness and volumes of the data that these social networks generate they have sprung an entirely new eco-system of its own - companies are now plugging into Twitter, Facebook and all the other networks to figure out everything they can about you. Everyone is now talking about "sentiment analysis" in social media for brands, targeting particular demographics with ads on Facebook, promoting tweets and segmenting the customers into different groups and selling their data to marketers.

A particularly interesting one of that new generation of networks is Twitter. With a base of 200+ million active users it has slowly but surely become one of the most prominent sources of information and news on the web. It's previously beaten traditional news sources on numerous occasions by a few minutes when delivering the latest developments such the Los Angeles earthquake in 2008. [4] They have their data

streams opened up to developers and researchers as well, which is a fantastic opportunity to mine this data set for valuable information. There are plenty of articles on the internet on what people have done with it - demographics research, predicting flu outbreaks, etc. [5] Miles Osborne here at the University of Edinburgh has done quite a lot of research using the Twitter data streams [6] and perhaps the biggest use is in Sasa Petrovic's PhD thesis. [7]

In order to get more familiar with the trending events I have done read a few papers on Topic Detection and Tracking and First News Detection in order to see if I could use and extend it for my case. I'd like to mention Sasa Petrovic's PhD thesis as an excellent paper on the matter. [7] However after careful investigation into the complexity of TDT I decided that this will be implemented in the 2nd part of the project, since I needed to familiarise myself with the dataset, try to see if there is more intelligent filters on the data stream and see how exactly they can be correlated to my 2nd dataset.

As it's been said many times "Predicting is hard, especially about the future". Add a noisy channel as a source of input and one is only making things worse. There have been several attempts to predict something based on Twitter - stock markets [8], flu outbreaks [10], results for the NFL [9] and poll results [11] being the most prominent published papers.

In order to predict the stock markets [8], the authors try to use the cutting edge of the real-time twitter data in order to classify the overall mood and from that to predict whether the DJIA [12] is going to go up or down. In order to achieve it they use OpinionFinder [13], which can measure opinion classifying it into negative or positive. To capture the complexity of human moods, they created their own system called GPOMS, which can capture up to 6 states. In order to make the OF and GPOMS time series comparable, they used their z-scores. They have proven that the GPOMS values are more descriptive of public moods than the OF binary sentiment. However correlation does not imply causation and in order to check whether there is any they use the Granger causality analysis [14]. It rests on the assumption that if X causes Y then changes in X will systematically occur before changes in Y. However in this case they are not testing whether there is actual causation, but rather there is predictive information in the public mood time series that can describe the time deltas in DJIA. Out of all the different mood values Calm has the highest predictive value. After determining that there is definitely a gain in using the Calm values to predict the DJIA delta they use a Self-Organising Fuzzy Neural Network (SOFNN) [15], which is a five-layer hybrid neural network, which has the ability to organise its neurons during the learning process by itself. After tuning and predicting the values on their dataset, they confirmed that Calm has the best predictive value.

However, in their discussion they note that they are not really filtering their users on any particular geographical, language or cultural level. In this project, we have only taken *tweets* from the United Kingdom, which allows us to try to predict where people from the UK want to fly to. As far as their methods are concerned, working with Neural networks is certainly one of the options that I'd like to investigate next year, but because of the fine tuning and additional overhead of using and training and possibly implementing it, I decided to leave this more complicated method for next year.

On the other hand you can do simpler things when you are trying to predict public opinion on a certain topic [11]. The authors make a strong argument that using Twitter you can replace the expensive time inefficient polls. They use public sentiment in order to predict ups and downs in the consumer confidence rating. In order to remove any massive outliers and because of the volatile nature of public sentiment they use something they call "Moving average aggregate sentiment", which smoothers out the daily variations.

On the other hand, the only thing that has been researched in the online travel sector is what is the optimal time to book an airline flight taking into account all the different factors. [16] [17] Both of the referenced papers are using small sets of data which aim to predict what is the optimal time to book. What is written there is the other part of the puzzle - when should you book, but in this particular project we are more interested in when ARE people looking for flights.

Due to the lack of any research in this particular area setting the objectives for this project was very difficult. Planning on how to approach and tackle was in itself a challenge. There is no current proposed method of doing this, so there was no gold standard against which I could benchmark. That made it particularly hard to see whether the classifier is right or wrong and what should I strive to beat.

The Methodology can be found in chapter 3 on page 19 split by all the main subtasks that had to be carried out. The subtasks are described in detail and each one of the sections includes an overview of what was done, how it was done and what difficulties were encountered and how I overcame them.

Chapter 4 on page 27 explains what are the Machine Learning models used to carry out the work for this project. Since there is no model defined anywhere I had to use an in-house algorithm used by Skyscanner to predict the search volumes. It's then described in depth which model worked best and why.

In chapter 5 on page 31, I discuss all the future work that will be carried out for this project both for next year and if someone is to pick it up and start developing on top. In order to do that I have included links to the source code for this project and all the pre-aggregated data sets from Twitter in the repository.

Chapter 3

Methodology

The problem we have at hand is completely new. The fact that there is no previous work and little to none understanding of how it could be tackled left me with a lot of room to manoeuvre. I decided that I'd have to get accustomed to the Twitter dataset and explore and see what can be extracted and how.

That task wasn't as trivial as I expected because of the sheer size of the data. The daily rate at which I am consuming the data is 3.5 GBs (reduced from 6GBs/day), which is not small by any means. So far I have amassed 420 gigabytes of data, which I am using. I have described in detail the attributes I collected in 3.1 on page 20.

The sheer volume of the data gathered imposed a few challenges which I had to tackle:

1. How do I traverse all the data in a clever way that will allow me to keep all the processed data and only process unseen data
2. What data structures to use such that they will hold the data in an efficient way and then output it for later use or for use by other components
3. Building a scalable and performant codebase for the analysis that will allow easy and reproducible experimenting

I have shown what I am using and how I process the data in 3.2 on page 21. My implementation was fast enough, because of the reduction to a smaller working set and it was able to process the whole reduced dataset in about 20 minutes, which is fast enough for re-processing.

After collecting and preparing the data comes the next big question - what do I want to get out of this set? I have tried two ways of filtering it and trying to use it to predict the flight searches:

- Using hashtags which contain a country or city name and taking their counts.
- Taking every tweet that has a city/country name in its text in a conduction with a travel related term from the list.

Both of these methods have obvious advantages and disadvantages - taking only hash-tags will be really fast, but not expressive and representative enough. They are both

explored in 3.3 and 3.4 respectively.

Before building the models there was a final step that had to be done and that was cleansing the data. Because of some teething issues with problems with the Skyscanner data and the fact that the data collections process is not perfect, there are some holes in the data.

For the Skyscanner data set I have used the Last 4 Fridays forecasting method to back-fill and for Twitter I've just used the means to fill in the missing values.

3.1 Data collection

The first and most important part of this project was to start collecting the correct data, which was to be used in building the model later. Twitter offers quite a comprehensive API with a lot of attributes, however in order to reduce the daily volume of data I had to take the most relevant ones for me.

The attributes chosen to collect are:

- Text - the text of the tweet is the most important one, perhaps. Quite a lot of information can be extracted from it alone
- Id - the tweet id. Useful if we want to screen scrape for any additional information or just to provide a tidy small dataset of ids.
- ID Str - String representation of the above.
- Source - What is used to post the tweet. The Twitter website is marked as "web".
- Coordinates - Representation of the geographical location of the Tweet as reported by the utility used to post the tweet.
- Entities - This include hashtags, user mentions and urls included in the tweet. Could be taken from the text, but it's nicer to have them ready.
- Retweet count - the number of times the tweet was retweeted. Useful for any future models.
- Favourited - Indicates whether the tweet was favourited by people - an analogy of this would be a Facebook like.
- Language - The language of the tweet. We are capturing English language tweets at the moment, but put in place for future expansion into the multi-language domain
- Filter level - indicates the level of filtering applied to the stream.
- Place - Shows where the tweet was tweeted from. Gives more detail than coordinates - country, city name, coordinates as well. Not necessarily available for all tweets.

The first stage of the project is to look only at tweets in English coming from the UK. That will allow us to predict and model the flight search volumes in the UK based on the mentions in the Twitter Stream.

The second stage would be to develop the model even further and add multi-language and multi-country support and employ more sophisticated models such as ARIMA or Auto-regressive Vector Models.

By selecting those particular tweet attributes and using the Streaming API, I managed to reduce the daily volume of data from ~6GB of data down to ~3.5GB/day. The amount of data accumulated at the time of writing is ~380 GB. The collector has been running successfully from September 2013, however there are some holes of the data caused by network outages or the script interacting with the Twitter Streaming API crashing.

The data on flight search volumes is kindly provided by Skyscanner. In order to ensure that there are no concerns with confidentiality I have anonymised the data.

3.2 Data processing

3.3 Hashtags

Hashtag is one of the most important constructs by Twitter. Here's an example tweet:

@FunnyChap: Something witty and very well said **#jokeoftheweek**

The structure of the tweets is:

1. The name of the user is denoted with a @ before the username.
2. The text itself.
3. The hashtag (a special Twitter construct) is denoted with a #.

You can filter out and explore twitter content based on hashtags and usually every major event has its own hashtag. The ones which are most mentioned appear in a special section called "Trending".

This option was considered, because if it had worked it would have been what we would call a "quick win". It doesn't require much processing, since what you'd need is just take the ready made list of hashtags and store all the results into a in-memory dictionary, which you'd then split by city/country. That reduced the overall size of the relevant tweets to ~3 GB. What was even greater is that we didn't really require the text information itself, which made the working dataset even smaller.

However the overall counts were quite small and the distribution of the values was very noisy, which means that fitting a model to this would've been quite hard if not impossible.

3.4 Occurrences paired with travel terms

The next slightly more expensive in computational terms option was to look at the actual tweet content and to count the number of times a city/country name appeared in conjunction with a one word form a list of travel-related words:

```
terms =
{ "airport": "",
  "check-in": "",
  "fly": "",
  "land": "",
  "landing": "",
  "plane": "",
  "take off": "",
  "destination": ""
... }
```

The list is modelled as a dictionary, because dictionary lookups are $O(1)$, so iterating over the words in the tweet and checking whether they appear in the city/country dictionary or the travel terms proved to be the most efficient combination.

That gives us about ~10 GB of tweets to work with. The reduced dataset is still much smaller in comparison to the full one. Processing the whole usually takes about a day, because of the extensive lookups in the city/country dictionary and the travel terms one.

3.5 Finding the correlation between those variables

After discovering that the volumes of twitter counts from the second method are not so noisy, I decided to proceed and explore the correlation between the two variables.

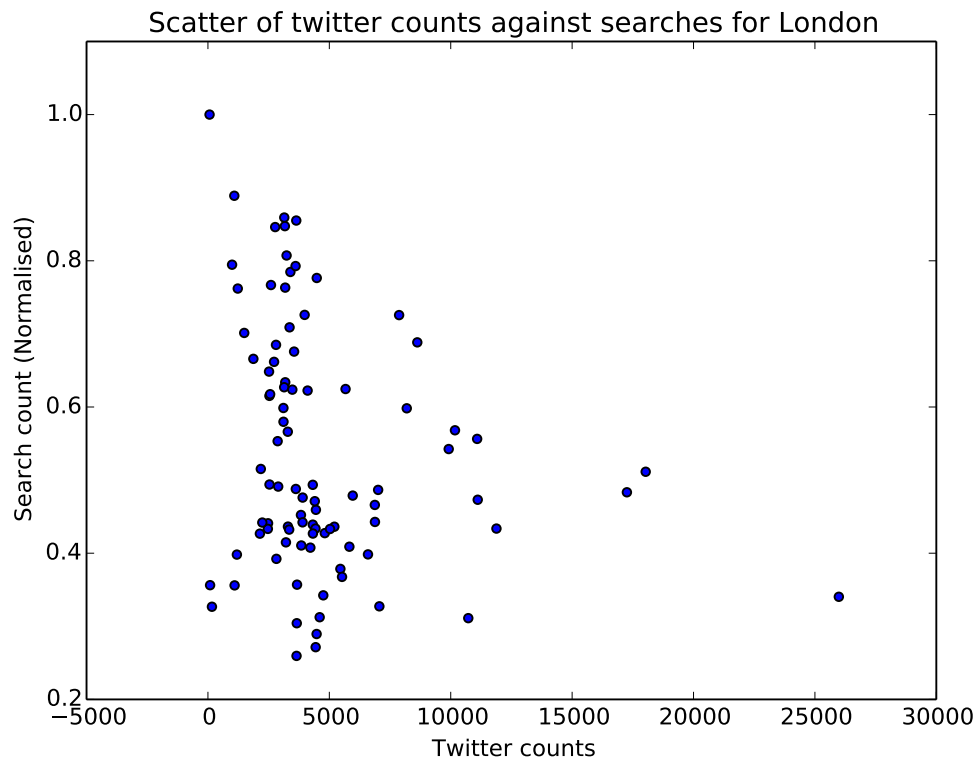
The easiest way to test the statistical correlation between your data would be to carry out the Pearson test.

For London the numbers we got back from that are:

```
In [39]: r_row, p_value
Out [39]: (0.13, 0.28)
```

From this it seems that the situation is truly unrelated. After all a positive correlation of only 13% is something that even a social scientist wouldn't report!

However, when you plot the two things we are trying to correlate something a bit more interesting comes up:



There is definitely a relationship, even though not perfectly linear.

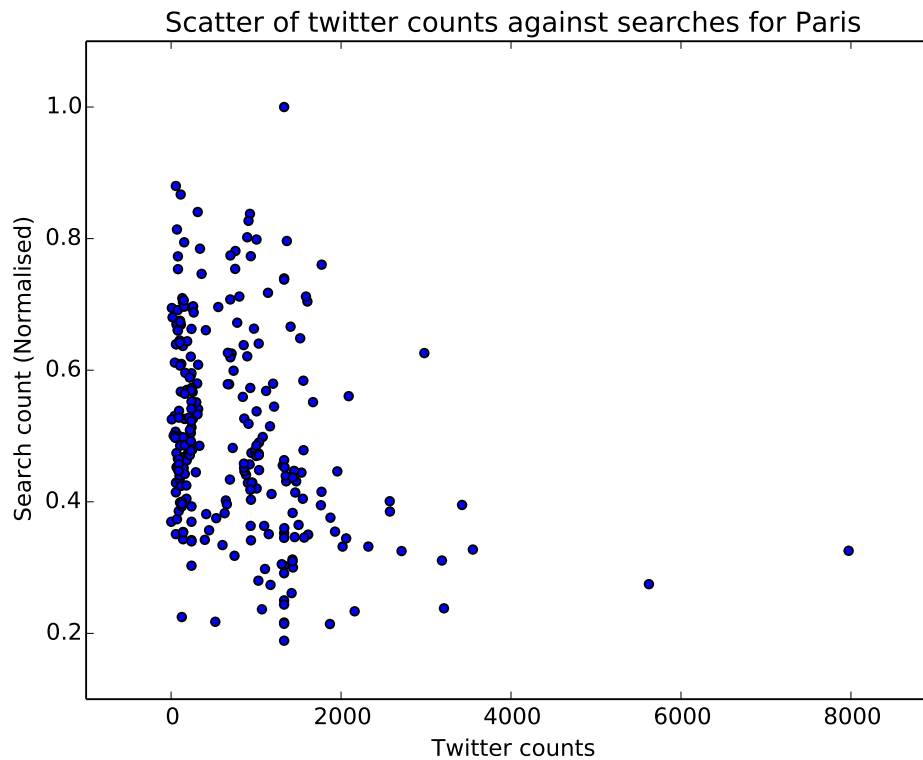
We can observe that there is a well defined cluster where the twitter counts are big and the corresponding searches are relatively big as well. We can't say with great certainty that those two are perfectly correlated, because correlation does not imply causation. So what we needed to do was to delve even further and build a model that would work with those.

Paris was one of the most popular destinations so it was quite interesting to see what the relation is there. The Pearson test yields the following values for Paris:

```
In [24]: r_value, p_value
```

```
Out[24]: (-0.25, 0.0047)
```

And here is the scatter plot:



Here we do see that indeed the correlation is relatively small and negative.

As we can see each of the different cities shows a different behaviour. The implications of that are that the approach we take to building the model must be more versatile - we will aim to build one model for the overall numbers and models on a per city basis.

3.6 Cleansing the data

As mentioned, every time you depend a few external data source there will be some problems with the data, caused by several things.

In this project there were a couple of major sources of problems:

1. The script that collects data from the Streaming API.
2. The searches data coming from Skyscanner being incorrect or partial - not spanning the full date range.

The dates with incomplete data or missing data altogether can seriously impact any regression or statistical test, so it was vital to tidy up the data set by backfilling it. I applied the Last 4 Fridays forecasting algorithm mentioned in the following chapter in order to make the dataset more consistent and easier to work with.

The interesting part was that it completely change the results from the test. For London the values returned by the Pearson test are now:

```
In [23]: r_row, p_value  
Out[23]: (-0.23, 0.02)
```

That is quite interesting and it has a lot of implications on the models we will use afterwards to correlate the two variables. We might have to add a certain "lag" factor to this, which will offset the tweets by correlating the numbers extrapolated from Twitter with the numbers and match them against searches from the following day.

Chapter 4

Models

4.1 The baseline model

The baseline model for predicting the number of redirects Skyscanner has on a daily basis is called Last 4 Fridays. It works in the following way:

1. We want to predict the number of redirects/searches for this Friday.
2. We take the number of redirects/searches Friday from last week, the week before, etc, until we have the counts from the previous 4 weeks for the corresponding day.
3. We then assign weights to those 4 numbers and assign weights - the most recent one will get the highest weight, the one after about a half and so on. The exponential weighting scheme captures short term seasonality.

This model was tested out against more complex classifiers like ARIMA and Autoregressive Vector Methods and it performed really well, but unlike the others it was many times simpler to program and maintain.

4.2 Last 4 Fridays + Twitter Counts

The first model I build was not a radical new approach, but rather an increment on the previous one.

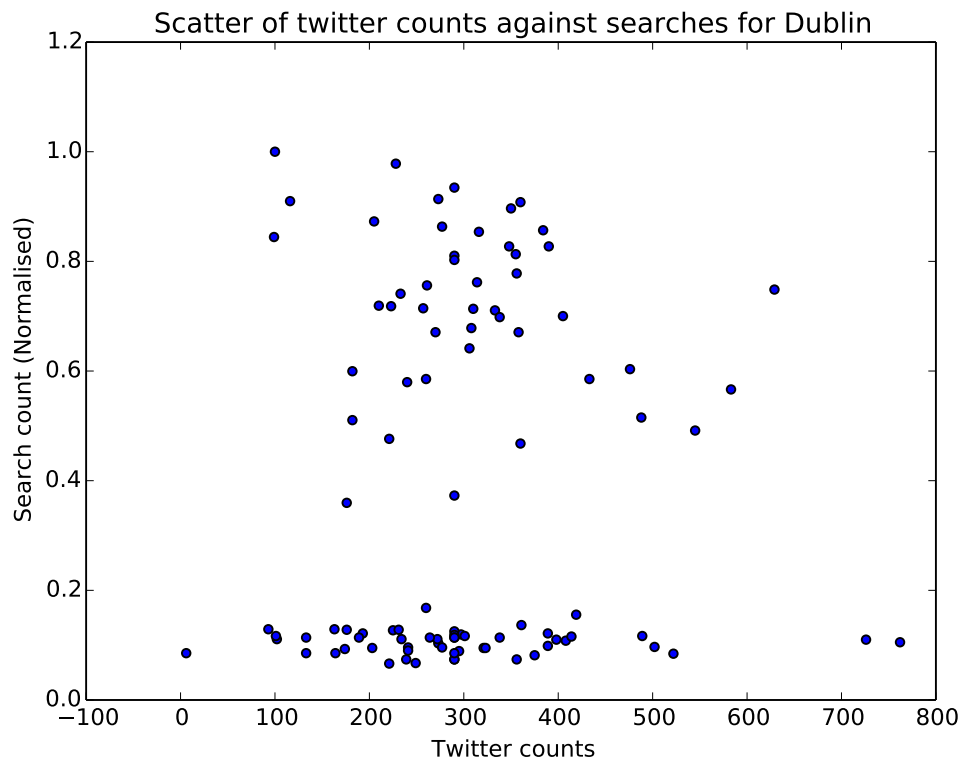
I generated a list of tuples which contained the following information:

- Twitter count for the place - how many times it was mentioned in conjunction with a travel term.
- Same day of the week from last week
- Same day of the week from 2 weeks ago
- Same day of the week from 3 weeks ago

- Same day of the week from 4 weeks ago

What this is is a version of Last 4 Fridays which now has Twitter Counts. Instead of using the exponential weights from the previous algorithm I let LASSO determine their weights. With LASSO I had weights for each of the 272 cities. After measuring the RMSEs for the L4F algorithm and expanded version with Twitter Counts the former performs better on 193 of the 272 cities. In the top 10 L4F+Twitter performs better on 2 out of 10.

Dublin which is in the 4th spot in terms of volumes is quite interesting. I plotted out the counts against searches just to see how it looks:



4.3 Results

And here are the results for the top 10 destinations by volume (they have the highest RMSEs).

City	RMSE L4F+Twitter	RMSE L4F
London	3335.869525	3160.05886
Paris	939.1057276	921.6785638
Barcelona	920.6831066	897.6473097
Milan	760.3436591	760.9286718
Rome	710.7052517	705.8388511
Manchester	574.8431422	572.4116201
Dublin	514.584231	527.884409
Amsterdam	550.2254781	516.2054596
Tenerife	544.6187529	502.0270723
Moscow	499.1803266	495.2139406

This is quite a good start considering the fact that all the weights were automatically determined by the LASSO algorithm and the fact that I have only about 130 data points so far (130 days).

4.4 Future improvements

The next on the list is to expand the feature set with more features from Twitter such as:

- Specific words.
- Trending event or not.
- Investigate whether sentiment will be useful here.

With those I'll be able to better the model and reduce the RMSE across the whole board and hopefully beat the current prediction algorithm for more than 80% of the examples.

Chapter 5

Future work

Bibliography

- [1] BBC, *Sequoia Capital investment values Skyscanner at \$800m*,
<http://www.bbc.co.uk/news/uk-scotland-scotland-business-24380126>
- [2] Stefan Sabev, Code for Inf Project,
<https://github.com/SSabev/HonsProject/>
- [3] Wikipedia, *Lasso (statistics)*,
[http://en.wikipedia.org/wiki/Lasso_\(statistics\)#Lasso_method](http://en.wikipedia.org/wiki/Lasso_(statistics)#Lasso_method)
- [4] Twitter, *Twitter as a news-wire*,
<https://blog.twitter.com/2008/twitter-news-wire>
- [5] Twitter Research, *A collection of publications using Twitter data*,
<https://sites.google.com/site/twitterresearch09/twitter-papers>
- [6] Miles Osborne, *Social media papers*,
<http://homepages.inf.ed.ac.uk/miles/sm-papers.html>
- [7] Sasa Petrovic *Real-time Event Detection in Massive Streams 2012* ,
School of Informatics, University of Edinburgh
<http://homepages.inf.ed.ac.uk/s0894589/petrovic-thesis.pdf>
- [8] Johan Bollen, Huina Mao and Xiao-Jun Zeng,
Twitter mood predicts the stock market,
<http://arxiv.org/pdf/1010.3003v1.pdf>
- [9] Shiladitya Sinha, Chris Dyer , Kevin Gimpel, and Noah A. Smith,
Twitter mood predicts the stock market,
<http://arxiv.org/pdf/1010.3003v1.pdf>
- [10] Michael J. Paul and Mark Dredze,
You Are What You Tweet: Analyzing Twitter for Public Health,
http://www.cs.jhu.edu/~mdredze/publications/twitter_health_icwsm11.pdf
- [11] Brendan O'Connor, Ramnath Balasubramanyan, Bryan R. Routledge, Noah A. Smith,
From Tweets to Polls: Linking Text Sentiment to Public Opinion Time Series,
<http://www.cs.cmu.edu/~nasmith/papers/oconnor+balasubramanyan+routledge+smith.icwsm10.pdf>

- [12] Wikipedia, *Dow Jones Industrial Average*,
http://en.wikipedia.org/wiki/Dow_Jones_Industrial_Average
- [13] OpinionFinder,
<http://mpqa.cs.pitt.edu/opinionfinder/>
- [14] Granger causality,
http://en.wikipedia.org/wiki/Granger_causality
- [15] Gang Leng, Girijesh Prasad, Thomas Martin McGinnity *An on-line algorithm for creating self-organizing fuzzy neural networks* <http://www.sciencedirect.com/science/article/pii/S0893608004001698>
- [16] Etzionni, Tuchinda, Knoblock and Yates, *To Buy or Not to Buy: Mining Airfare Data to Minimize Ticket Purchase Price*, 2003,
<http://knight.cis.temple.edu/~yates/papers/hamlet-kdd03.pdf>
- [17] William Groves and Maria Gini, *Optimal Airline Ticket Purchasing Using Automated User-Guided Feature Selection*
<http://ijcai.org/papers13/Papers/IJCAI13-032.pdf>