Weather and Text for Investment Research Kyle Balkissoon June 2018 Kyle.Balkissoon@ibm.com



Predictive Models are easier to improve by adding meaningful features vs complex statistics (or more quants)



Standard Financial Data

Quant Researchers

- Low Marginal Lift
- Demands free food
 - Too many GARCH discussions

More Data

- Potentially high lift
- No food demands
- Does not talk about GARCH



Weather Stations Globally – Capturing Data





Weather Data

> hea	d(data)																				
ke	y zipcode		class	expir	re_time_gmt	obs_id obs_	name valid_ti	ne_gmt	day_in	d temp	wx_icon	icon_extd	W	x_phrase pi	ressure_tend	pressure_desc	c dewPt h	eat_index rh p	pressure	vis wc	wdir
1 KST	L NA c	bserv	ation		1483260660	KSTL St L	ouis 1483	253460	1	N 27	29	2900	Partly	y Cloudy	NA		21	27 78	29.43	10 27	NA
2 KST	L NA c	bserv	ation		1483264260	KSTL St L	ouis 1483	257060	1	N 28	29	2900	Partly	y Cloudy	NA		21	28 75	29.43	10 28	NA
3 KST	L NA c	bserv	ation		1483267860	KSTL St L	ouis 1483	260660	1	N 27	26	2600		Cloudy	1	Rising	g 20	27 75	29.43	10 27	NA
4 KST	L NA c	bserv	ation		1483271460	KSTL St L	ouis 1483	264260	1	N 28	26	2600		Cloudy	NA		21	28 75	29.43	7 28	NA
5 KST	L NA c	bserv	ation		1483275060	KSTL St L	ouis 1483	267860	1	N 27	26	2600		Cloudy	NA		20	27 75	29.43	8 27	NA
6 KST	L NA c	bserv	ation		1483278660	KSTL St L	ouis 1483	271460	1	N 25	26	2600		Cloudy	1	Rising Rapidly	y 20	25 81	29.44	7 25	70
wdi	r_cardinal	gust v	vspd r	max_te	emp min_tem	p precip_tot	al precip_hrl	y snow_	hrly w	v_desc	feels_li	ke uv_ind	ex qua	lifier qua	lifier_svrty	blunt_phrase t	terse_phr	ase secondary	_swell_he	eight	
1	CALM	NA	0		NA N	Α	A	0	NA	Low		27	0	NA	NA	NA		NA		NA	
2	CALM	NA	0		NA N	Α	A	0	NA	Low		28	0	NA	NA	NA		NA		NA	
3	CALM	NA	0		NA N	A	A	0	NA	Low		27	0	NA	NA	NA		NA		NA	
4	CALM	NA	0		NA N	Α	A	0	NA	Low		28	0	NA	NA	NA		NA		NA	
5	CALM	NA	0		NA N	Α	A	0	NA	Low		27	0	NA	NA	NA		NA		NA	
6	ENE	NA	3		NA N	Α	A	0	NA	Low		25	0	NA	NA	NA		NA		NA	
primary_swell_direction clds primary_swell_height secondary_swell_direction primary_swell_period primary_wave_period water_temp secondary_swell_period primary_wave_height																					
1			NA	CLR		NA			NA			NA		NA	NA		NA		NA		
2			NA	CLR		NA			NA			NA		NA	NA		NA		NA		
3			NA	CLR		NA			NA			NA		NA	NA		NA		NA		
4			NA	CLR		NA			NA			NA		NA	NA		NA		NA		
5			NA	CLR		NA			NA			NA		NA	NA		NA		NA		
6			NA	CLR		NA			NA			NA		NA	NA		NA		NA		
>																					

Weather Data

Calculation of the Synthetic Weather Station – A Representation of weather temporally around a company of interest



Aggregating Geospatial Weather Data to appropriate periods of interest



What formulating the prediction statement looks like:

Gathered weather data for commodity-growing regions



Built dynamic forecasting model using data and analyzed data using Machine Learning

 $TargetVariable_t = StandardPredictors_{t-1} + Precipitation_{t-1} + \dots$

Identified new trading opportunities based on known features and weather resulting in a profitable strategy and a sharpe > 1



News Data: Raw News data in XML Format

<?xml version="1.0" encoding="UTF-8"?>

//itf change.date="June 10, 2005" change.time="19:30" version="-//IPTC//DTD NITF 3.3//EN"> <head> <title>Dan Gurney, Driver and Builder of Racecars, Is Dead at 86</title> <meta content="16gurney" name="slug"></meta> <meta content="Obits" name="dsk"></meta> <meta content="The New York Times" name="cre"></meta> <meta content="NewYork" name="edt"></meta> <meta content="" name="cl"></meta> <meta content="NYT5 Article" name="template"></meta> <meta content="" name="adsensitivity"></meta> <meta content="n" name="pay"></meta>
<meta content="" name="guid"></meta></meta> <meta content="2018-01-16-615692" name="sqn"></meta> <meta content="12" name="print_page_number"></meta>
<meta content="B" name="print_section"></meta> <meta content="Dan Gurney, Dynamic as Both Driver And Builder of Racecars, Is Dead at 86" name="print_headline"></meta>
<meta content="" name="page_title"></meta> <meta content="obituaries" name="section"></meta> <meta content="" name="subsection"></meta> <meta content="Obituaries" name="section display name"></meta> <meta content="" name="subsection_display_name"></meta> <meta content="" name="parent_section_display_name"></meta> <meta content="" name="masthead_url"></meta> <meta content="n" name="comment_flag"></meta> <meta content="2018-01-16 13:39:49" name="updated date"></meta> <meta content="minor" name="update_type"></meta> <meta content="Article" name="asset_type"></meta> <tobject tobject.type="news"> <tobject.property tobject.property.type="current"></tobject.property> </tobject> <docdata management-status="usable"> <date.issue norm="20180116T000000"></date.issue> <date.release norm="20180115T155055"></date.release> <doc.copyright holder="The New York Times Company" year="2018"></doc.copyright> <doc-id id-string="10000003861569" regsrc="NYT"></doc-id> <identified-content> <classifier type="des">Automobile Racing</classifier> <classifier type="des">Indianapolis 500 (Auto Race)</classifier>
<classifier type="des">Deaths (Obituaries)</classifier> <classifier type="tom">Obituary (Obit)</classifier> <classifier type="tom">obituary tobituary </identified-content> <kev-list> <keyword key="Gurney, Dan (1931-2018)"></keyword> <keyword key="Automobile Racing"></keyword> <keyword key="Indianapolis 500 (Auto Race)"></keyword> <keyword key="Formula One"></keyword> <keyword key="National Assn of Stock Car Auto Racing"></keyword> <keyword key="Deaths (Obituaries)"></keyword> </key-list> </docdata> cpubdita date.publication="20180115T000000" ex-ref="https://www.nytimes.com/2018/01/15/obituaries/dan-gurney-driver-and-builder-of-racecars-is-dead-at-86.html?p name="The New York Times" position.section="obituaries" type="web" unit-of-measure="word"></pubdata> </head> <body> <body.head> <hedline> <hli>Dan Gurney, Driver and Builder of Racecars, Is Dead at 86</hli> <hl2 class="col"></hl2> </hedline> <byline>By Frank Litsky</byline> <dateline class="print"></dateline> cabstract>



End-to-end workflow for news data aggregation and analysis

Aggregation



Data is coming from New York Times

Have access to everything published (print and online) historically

Requested articles from 2001 to the end of 2018. All articles include **full text** and **metadata** (section, date published, keywords, etc.)

Cleaning



All tickers are mapped to a **company name, nicknames,** and **subsidiary companies** (Darden to Olive Garden)

Found matches between company names and article text/keywords. Potential for articles to match to multiple companies

Analysis

A surprisingly large fraction of applicatus, even how will matter's degrees and PDD. [In comparison scales and provide the scale of a corry out basis programming tasks, "and cample, Type personally interviewed quarkatus who can't assure "Writes leap that causation from 18, 00's or "What's the number after P [In hexadecimal?" Less trivially, I've Interviewed many candidates who can't asser recursion 18, solve areal problem. These are basis skills, anyone who lades them probably havit? Cam much to programming. Speaking on behalf of software engineers who have 18e iterview prospective meak irrs, I can safely any that we're third of tablish y Garaviditates who can't program their key out of a paper bap. If you can successfully urite a loop that geosfrom 18, 00's eveny language on your resame, can do simple arithmetic without a coloculator, and can use recursion to solve a real problem, you're already about of the packt

10+ features are extracted from article including: sentiment, taxonomy, and word count

R Packages: XML, tm Sentiment: Dictionary Score/Proprietary



Feature creation and engineering for news data

Features



Example features that can be generated on a per taxonomy basis

- 1. Sentiment
- 2. Word Count
- 3. Rolling averages (looking backwards) using N day windows
- 4. Importance Score (combination of sentiment and word count)
- 5. Sentiment Index (3 Buckets)
- 6. Word Count Index (6 Buckets)

There are other features on a **daily basis** such as Sentiment and Word Count Index

Moving Forward



Decay algorithm is used to carry previous observations forward when no new news articles are released. Values decay at **X% rate** each day until a new article is released, resetting the value.

All values decay towards **zero**. For sentiment values that is trending towards **neutral** (1 is positive, -1 is negative). Word count trends towards zero from a positive number to show the **lessening importance/relevance** of an article over time.



Case Study: News & Equity Predictability

Problem Statement: Does news data have additive predictive power on equity returns?

Equity Returns = News Factors + Fama French Carhart Factors

- Target Variable: Daily Returns of 50 Equities
 Pool
- **Predictors**: News data (New York Times), Fama-French Carhart factors
- Period: January 2 2001– December 29 2017
- Takeaways:
 - News factors of average sentiment, word count and importance score are statistically significant
 - Not as significant as FFC factors but still provide some predictive power

	0									
Coefficients:										
	Estimate	Std. Error	t value	Pr(> t)						
(Intercept)	1.275e-04	3.487e-04	0.366	0.7147						
Average_Sentiment	2.149e-02	1.098e-02	1.957	0.0503						
Average Word Count Quintile	1.435e-04	8.261e-05	1.737	0.0824						
Importance Score	-5.814e-03	2.894e-03	-2.009	0.0445	*					
SMB	3.104e-03	9.119e-05	34.033	<2e-16	***					
HML	4.444e-03	8.590e-05	51.732	<2e-16	***					
RMW	-6.942e-03	1.227e-04	-56.569	<2e-16	***					
CMA	-5.752e-03	1.475e-04	-38.998	<2e-16	***					
Signif. codes: 0 '***' 0.00	01 '**' 0.01	l '*' 0.05 '	.' 0.1	'1						
Residual standard error: 0.02144 on 185535 degrees of freedom Multiple R-squared: 0.05766, Adjusted R-squared: 0.05762 F-statistic: 1622 on 7 and 185535 DF, p-value: < 2.2e-16										
	-	_								



