Stonks: Project Description

• What is your theme/topic/goal/issue to be tackled - why is it important to you? My topic is market prediction. This is a fundamentally simple and reasonable goal. However, depending on the data that I harvest from news sources for this goal it could have either a reasonable outcome that is highly helpful or an absurd, chaotic, misrepresentative one. In the first case this project would be more of a market analysis tool and in the other case it would be an absurd experimental data driven art piece detailing the chaotic nature of market reporting and decision making in our very fast paced world.

Markets have a treasure trove of data associated with them. They have existed for hundreds of years. They have millions, billions or trillions of data points created constantly, and they are still being updated at every millisecond. Furthermore, as they continue to be updated, they have more data generated because they expand exponentially. I think this data theme offers a lot of avenues for interesting and revealing outcomes.

This is important to me because analysis and advice regarding markets is usually very misleading, biased or just wrong. This is because markets are one of the most difficult things to predict. I want to attempt to solve this problem or at least show how difficult it is to solve this problem.

• What form will your project evolve into - who is your audience? This project would evolve into a news consolidation platform. It would interest either anyone who is interested in market analysis or anyone who is interested in real time crowdsourcing and data driven art. Its audience is people who would normally scour many different news sources looking for a piece of information that would fit their predetermined belief about the direction of a stock. It would act as a tool that would replace the apophenia associated with looking for patterns in a collection of news data. It would be a scientific tool to screen out this irrationality and lead to possibly a better decision.

Everyone who participates in a market wants to predict them. Specifically, whether this value of its constituents will increase or decrease. While this is a very simplistic goal, since markets are almost completely stochastic, or random, it becomes an almost impossible problem to solve. This is the problem I will attempt to solve for my final project. I think there are quite a lot of people who would be interested in this goal. Markets are huge focal points for many jobs and careers revolve around them. My audience would be anyone who either has an interest or who works in one of these fields. My audience would be anyone who is interested in data consolidation would also be interested in this, but that would be much fewer people, I think...

 Discuss how each of the two readings listed above have inspired/motivated your current choices with regards to the project. Apophenia, a concept from Hito Steyerl's writing, is a major problem when trying to analyze market data. It is so easy to start to see patterns that are just figments of one's imagination even when there are real takeaways to be had. My project would function as a scientific tool to replace apophenia.

Stereyl also foregrounds the concepts of signal and noise. In my project the noise would be irrelevant news regarding market movements and the signal would be news that is relevant and actually operates as predictors. Even in retrospect, it's impossible to know, on an individual basis, which of the millions of headlines predicted a movement in the market, let alone, any one of the market's constituents. However, if given a large enough dataset of daily market movements and daily news collections, I think it might be possible to screen out the noise and find which keywords, topics, or news sources are the most reliable in predicting the market and therefor constitute the signal. Furthermore, I can download archives of the history of news and markets and look at how news predicted markets in the past.

Also, Mimi Onuoha states that, "As we collect more data, we prioritize things that fit patterns of collection." Applied to my dataset, time is prioritized when collecting market data. Every datapoint has an associated a time tamp and sometimes this becomes misleading. This is because time is largely arbitrary when it comes to movements in prices. Usually price changes happen because of events, not times. So, time becomes an almost useless datapoint even though it is the most common. This is also why I am prioritizing collecting data over time. The more data I have, the more accurate my machine prediction is.

What medium(s) do you intend to use and why?
 I intend to use the medium of Data. This data would be displayed on a website so that may be considered the medium instead. I want the final product to be an accessible website that updates market predictions in real time. And/or a feed of consolidated news computationally analyzed.

It's hard to envision this project's front-end result. Despite how complicated the back end could become, the front-end would really be no more than a one-dimensional signal advising the observer to buy or sell the market as a whole.

So, along with this simplistic, capitalist application, I also want to create a frontend for the news aggregation aspect of the project. Maybe an apophenic newspaper. Or machine-readable news. Or rather, a machine digested news. One that aggregates hundreds or even thousands of news sources and spits out the most euphoric/pessimistic stories at the top and most mundane ones at the bottom. Or even a newspaper that combines stories on a given topic into one article resulting in a jumbled, schizophrenic diatribe of sorts.

What is your data: where will you get it; will it be collected - how and why?
 Half of my data will be company valuation data. The other half will be news sources.
 The first half of this data will be a collection of company stock tickers and their associated performances, sourced from company tracking APIs such as NPM google-finance or NPM yahoo-fin. The other half will be data from news sources. Either using

a web crawler and scraper in combination (such as NPM cheerio and NPM Crawler), or a news API (such as 'NewsAPI') that provides news articles in JSON format. The project will be in connecting the two of these datasets in a meaningful way using sentiment analysis.

Another important note is that I plan to analyze news and market data history in order to inform my predictions. If I look at the history of how certain news did or didn't act as a predictor over the course of history, I can use this information to screen for current news that will be a signal for the market. Historical market data is very easily available, news archives are more difficult to find because they may not go as far back in history, but they are still available as well.

 At a very high level: what are the algorithm(s) that will be used and implemented to achieve your intentions?

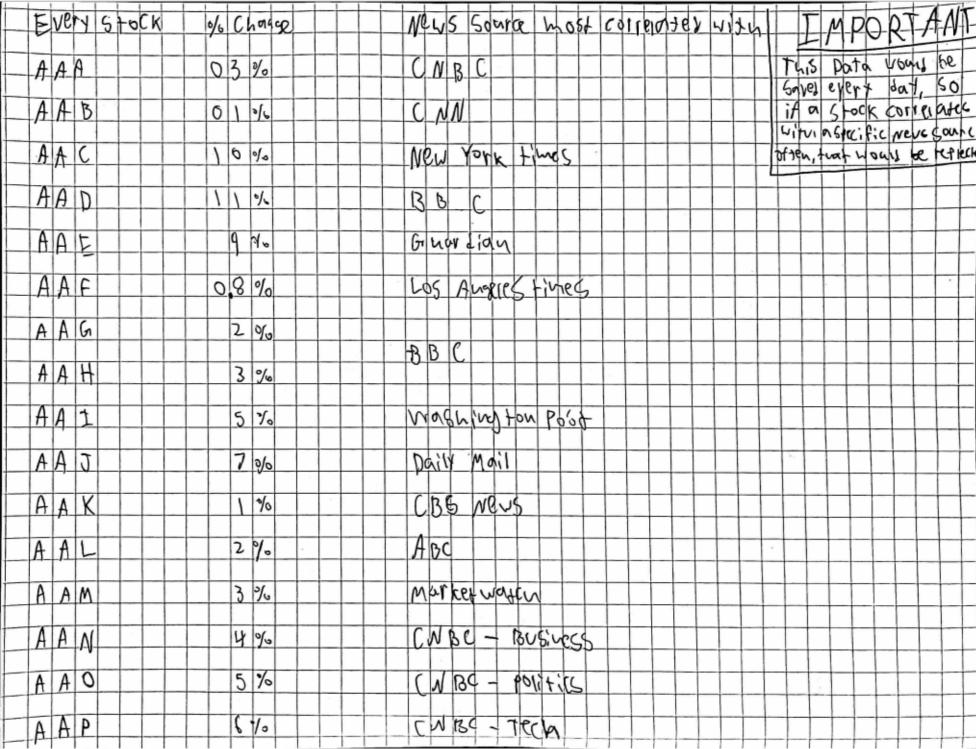
One algorithm that I know I will be using at this point is Sentiment analysis. For this I would use NPM sentiment, the same library that was shown in class. This would allow me to analyze whether a piece of news has a positive or negative sentiment. However, the other datapoint I would need to collect that sentiment analysis would not necessarily provide is which part of the market the news article could be a predictor of - whether it is a general market predictor such as a news article about a pandemic, political upheaval, or a market crash itself. OR whether it is a news article that relates to a sector of the market or even a specific stock. For this I would need to write or find an algorithm that looks for keywords within the headline or possibly even within the content of the article that would classify it as a tech related article, healthcare, automotive, etc. I wonder if there is a way to computationally analyze a news article title to detect whether it falls within a set of predetermined categories. OR whether there are enough articles daily in the world that contain the specific name of a company. Since this could end up being very coding intensive, I will start with creating an algorithm that finds articles that would be predictors of the market in general and not just one part of it.

Two other algorithms that would be fundamental are a news data endpoint and a market data endpoint. The news data endpoint could either come in the form of a web crawler combined with a web scraper (NPM cheerio, NPM crawler for example). These would allow me to access simple strings containing news headlines. The difficult part here is that then I would need to extend my crawler to be able to 'click' on the hotlinks of every news website headline per news source in order to download the content of the article it is associated with. This becomes very coding intensive very quickly. I already tried to get into programming this, and part way through I realized that there are news APIs, such as 'NewsApi,' that offer streams of news data already in JSON format. This API has options to stream from certain regions and only from news that contain certain keywords/phrases.

The Open Data handbook is a project that provides inspiration for my project. In it is described open data as data that is *can be freely used, re-used and redistributed by*

anyone - subject only, at most, to the requirement to attribute and sharealike. I am grateful that open news archives exist. Otherwise, I would have to pay in order to access this history. The open data handbook also specifies that the in order to be open, the data must be available in full. This is very important for my project because, for my algorithms to be as accurate as possible, they must have a full scope of data accessible. This is also why I plan to make my application open as well. So that just as I can build on the shoulders of this open resource others will be able to build on mine.

In the why section, the open data handbook outlines government data as being especially valuable because of its centrality of collection and quantity. Government data is usually open as legislated by law. In my case, since I am building a data application that is intended to analyze corporations, my case is very different. The kind of data that I require is almost as far away from government data as possible. As far away as corporations are from the government I suppose. There are probably various closed data sets that would be far more predictive of the market than news sources. However, since this data is highly valuable, I will not have an affordable route to accessing them.



NEWS SOUNCE (Include	25 hendline news as well as marker centric hers and succo	ntegorics)
CNBC - Business	100 x (Positive) - 6% (Ne of tive)	Wal
- Tech	100 % (Positive) -0 % (Negative)	N2-
- POlitiCS	100x (P05ix 118) - 0x (New ative)	N3
Marketwatch	100% (POSI+ive) - {01,3(N coa+ ive)	M4
CNN		NS
N D C		
CBS News		N B
New York Fines		NA
Daily Mail		NOV
Washington Post		Nu
BBC		N 12-
Guartian		N 13
Los Angeles Times		N IY
New York bost		NIS
ETC.	ч	

1					-			_	_	-	-	-				-							1 1			-	-	-		1
+	Alteran	ativity	_I	18 tco		70	Cu	CCK	ing	(10	1210	M	545		vit	N	_ <	SPEC	4	6	5+01	45,	_	10	Iv	deci	25	or	-
ł	Sedtors	Ceal	6 6	e Ch	AR C	K67	1 1	- 4	- 1		-201		1 1				l I	1 1		40			امما				LUC	160	MICE	2-
+	-Re1044	V100	1 m	Sal	15	HO	A)	h	dh	14	6	N	0.1	100	0	2	十七	e.	1001	(.	hod	60	e Ci	3	c	43	sha	Po	nie	-
+	A160 +	erc o	ve 1	5417	1	a t	waa	t	ul	cli	65:	D	2	70	we.	1	5	8 6	500	3,-	No	12 20	19	108	0,	52	P/7	[5]	C,E	<u> 10.</u>
1	AND 6	ector	s co	all	b		iuke	1	to	N	ew	5	500	ince	2	11/	6-	(0)	e cox	les	0	1 3	FOR	1	Ver	15	50	4	le s	
1	that	are	1)09	607	01	V .	5 pec	if	U	Sec	tot	5	1	ke.	TE	CL	۸.	I	we	かい	49	. 3	OC	11	M	edi	ار کم	fco	It's	ETC.
1	Relate AND 6 Thus The thic	0008	e, W	6 N	our	1 4	oite	orte	1	-he	N	E	WS	0	AT	A	Fiv	50	Teriv	har	11	du	1	he	1	se e	de	1	h -	1
1	whether	u Enj	ices	Cor	121	ate	Wil	tu	1	0	1	W	hid	W.	se	640	15		COLE	1019	9	in	14	W	mic	4	Ne	25	6001	45
_	NEUS GOV	100	Cen	MI SE	the	wey	F	,	TOP	k	149	NOV	16		1	R	VI	hin	her	10	RIC	101 11	V	14	NT	NJO	×			
ļ	CNBC			0					10	70	1,"	1		-		1	-	-	10	-	100	Yo		7		1	1			
ļ	MORK ET W	at ch		0) -	i			ı	Bo	MK	10						_	7	-	100	101				1	Vo	16		
_	CNN			0	-	1			٦	C	11								1	5		1%		Ke	9.2	IND	141	rack	< 01	-
ļ	NBC			(0 -	1			,	Bi	te oi	ih"							7		100	a.c		+	90	ke	ow.	125	con	~wu
ļ	CNN NBC CBG ABC			1	2 -	1				"	69")-	101	0/0		V	vor	25	CO	all	Kon	
	ABC			2) -	1				"50	161									9-	100	10		U	8	WI	1 1	'b 1	wan	MATA
	BBC				0 -	١			- '	AL	nuz	Oh'	,							φ-	13	0/10			Nec	16.	NO	99.0	CITI	u
	New Yor	K Tik	res		5 -	١			_ '	'R	naz o[K	et"								6-	10	· 100		4	100	KS	k	1 4	B ive	a
	NOS Ano	cles Ti	mes	0	-	1			- 1	" Go	916	11:16	2						(0 -	10	360		ìf	+	MOV	8 0	we	Moll	15
	Washing	ton Pt	100		-	1				da	111	1/1								0 -	10	3,/		t	nd4	al	.0	1900	arm	2
	Gnardia			1		1			_ '	, 0	11	'								-6	10	Che		0	N	WO	KhA	ne	WG	7
	NEW YOL	K 602	+	1		1			- 1	K	arda	Gub	101			_	_		,	5 -	16	000		8	acr	(6)	3 a	(11 :	v4 or	200
	HUAFIN	ton	J86		9 -	1					200	4_							7		10	V40		0	NY	W	hetu	ren	the	Y
Ī	HUATURY NO	WS		(-						ell								(5 -	10	036		C	DUI	1	reic	NE	+0	-
Ī	USA TO	4006		1		ı				_ _	eb(<u>t</u>) -	10	090		0	14	Pec	ific	ch	18	
ľ	politica			,	0 -	1				- 1	01	(I								2 -	10	0 96		1	40	50	ock	CU	SAL	
1	Yarbo N	SM3		2	2 -	1					67	<u>d</u>							(5 -		036			145	PC	+0	LANC	in (9
Ť	NPR NO	45		1	5 -	1					et	- (10	0 40				-				
+	Bine Al Di	- 		1	6 -	1					Pd.	-6							(0 -	119	00%								
-	News W	eek		7	0 -	(2	HC								0 -	10	80								
-	ChiCaso	Tribe	the .	1	5 -						6	HC								0 -	10	0 %								
1	Salon	,		1	5 -	1,					e	40								0 -	11	0%	5							
-	News W Chicasos Salon Boston- The Se	com		1	0 ~	1					t	していていい								5 -	10	0 %	,							
	The Se	attle	Tilm	e5 +	-	(e c	1-							1	0 -	- 1	00%	6	-						
۱	1111-0-												_		-		-1-		1-1	_	-	1-1		1		1			1	-

		1	1		_	$\overline{}$			_	_			_		- 1	$\overline{}$					7	1			_	1		_		4
	Alternowilled Kexwilds (a wills took see occur Theirote Significant ?	0 6	d	6	2	H	18		POC	05 P1		1	×	F	7	16)_	L	1	W	01	15	25	U	64	14	be	40		_
	Wild trank are occur	A	CV	05	5		NB	25		0	F	Fel	7	w	-	7.6	/\ A	5	1-1	80	Son	WC.	35	,	T	he	1	WO	nle	_
	Theisofe Simuiricant 9	A	hd	5	il		+1	125	1	10	1			-0-		~	_V1	7	-	-	T	1	Ť	İ	Τ,	\Box				_
					-+-	÷	-11-	~	Ť	j-	-					-		-	-	\vdash	\vdash	+	T	\vdash	\vdash					
	Significant Keyvorus	5	10	k	T	'KR	v/	Ins			_	_	_		C		9	CV	tir	1	T	+	+		\vdash					
		0	-		-11	CICI	+	1000	1	10		۷,۰	_	1	=	-	_	101	1	-	\vdash	+	\top		\vdash	_				
	701	Δ	A				\neg	\neg	+	\vdash	\vdash		_		0	_	T	\vdash		\dagger	†	\dagger	\top	\vdash	\vdash	\vdash			\Box	
		+4	14	1			\dashv	\top	+	+	1				~		Ė	\vdash	\vdash	\dagger	\dagger	+	+	\vdash	\vdash	\vdash				
	Bank	TA	A	F			7		+	+	\vdash				0	_	1	\vdash		†	\dagger	+	+	\vdash	+	+				
		+′-′	1				\neg	+	-	+			_		7		-	\vdash	+	1	t	+	+	+	+	+				
	Cav	1	A	17			\neg	\neg	+	+	\vdash				0	_	1	\vdash	+	\dagger	†	+	+	t	1	+	1			
		1	111	2				\neg	+	1	\vdash	\vdash					Ė	†	†	\dagger	†	+	+	†	T	+			\Box	
	Break	T	60	h			\forall	\dashv	+	\top	\vdash	\vdash			0	_	1	+	+	+	\dagger	\top	+	+	+	+	\vdash		\Box	_
		Τ,	1	1				\neg	\top	\top	\vdash				_		Ť	\top	T	t	T	\top	\top	†	+	+			\Box	
	fez	1	101	UTI	on6	5		\neg		\top	\top	\vdash			0	-	1	T	T	T	†	\top	\top	\top	\top	+				
		\top				Ŭ		\top	\top	\top		T		\vdash	_		ľ		\dagger	\dagger	T		\top	\dagger	\top	\top				
	201	1	480	Hu						1	\top				0	-	1	T	T	T	1	1	T	+	T	\top	\vdash			
		Τ.	T	1										\vdash			Ť	T	T	Ť	Ť		\top	1	1	\top				
	Amorzon	1	8	PC	0<	3						T			0	-	1	\top	T	Ť	\top	\top	\top	\top	\top	\top	\top			
		\top		-							T			\top	Τ	Г	Ť	\top	T	Ť	1	\top		\dagger	\top		\top	\Box	\Box	_
	ROKER	1/	AA	2								T	Г		6	-	1	1	Ť	1	\top		\top		\top	\top	\top			
			7	1					\neg		T	1							\top	Ť	1		\top	\top	1	\top	\top			
	50/148/148	1	10	59	PA	. 10	00							\top	0	-	1	T	1	1	1		\top	\top	\top					
		7	Τ,	T	1								Т	T			Ť					\neg				\top	\top	\top	\Box	
	(arbba	16	AL	T											0	-	1							\top	1	\top	T			Ė
			7'										T		T	T									\top		\top		\Box	Г
	011															T	T	T		1				\top			\top	\top		
														T	T	T						1		1	\top		1	\top	\top	Г
_	Television														\top	T								\top	\top		\top	\top	\top	Γ
_													T			T	T		Ť		7		\top			\top		1	1	T
			T										T		T	1	1			\top	\neg	\neg			+		+		1	
											T		T		T	T	1								1		1		1	1
_					T					Ť			T	\top	1	T	Ť			1		1	\top	7	\top	\top	+	+		T
_			\top	1	\top		1	\Box			_		\neg	_	1		1	\neg	1	-		\rightarrow	-	+	\rightarrow	_	-	-	+	T