

## **Objective**

- 1. Introduce a predictive model of the price of Bitcoin
- 2. Explore significant underlying features of the model
- 3. Provide key insights and takeaways

### **Model Overview**

- 1. Predictive model for the price of Bitcoin
- 2. Standard linear regression
  - As opposed to time series analysis; factors into cross-validation assumptions)
- Three features with high correlation to price of Bitcoin; regularization was evaluated but deemed not necessary

## **Feature Exploration**

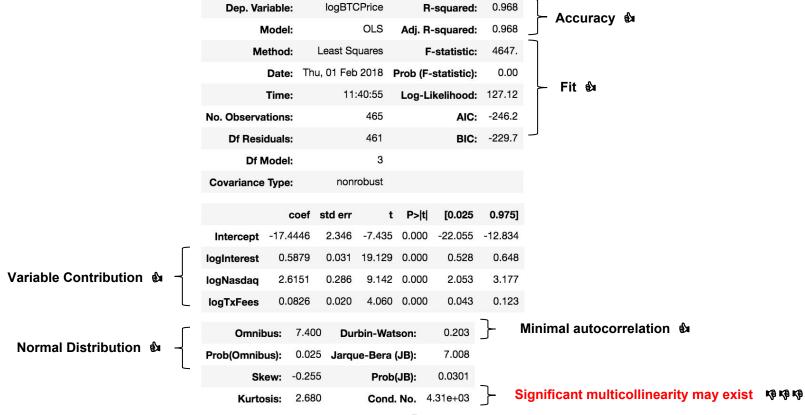
#### Target features not significantly influenced by the price of Bitcoin

BTCPrice	logBTCPrice
----------	-------------

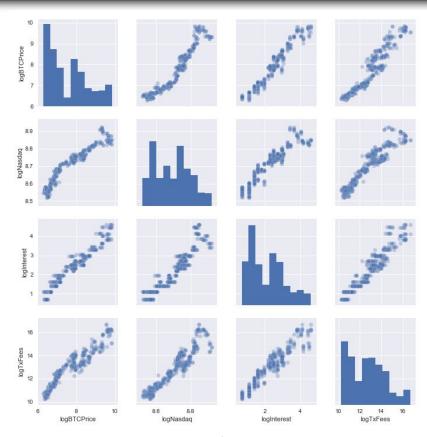
logBTCVol	-0.243729	-0.430323
logGold	-0.169522	-0.295356
logNoTxns	0.448220	0.404641
logAvgBlkSz	0.524849	0.580973
TxFees	0.872554	0.711003
logUniqueAddresses	0.760594	0.768803
Interest	0.935580	0.835610
BTCPrice	1.000000	0.899471
logETHPrice	0.755177	0.937083
logNasdaq	0.793758	0.958180
logTxFees	0.839659	0.958260
Nasdaq	0.813879	0.963822
logHashRate	0.824414	0.967630
logInterest	0.870113	0.977731
logCrypto Market Cap	0.859547	0.987115
logCostperTxn	0.878492	0.987899
logBTCPrice	0.899471	1.000000

	logBTCPrice	logNasdaq	logInterest	logTxFees
logNasdaq	0.958115	1.000000	0.945212	0.948694
logTxFees	0.958548	0.948694	0.957145	1.000000
logInterest	0.977783	0.945212	1.000000	0.957145
logBTCPrice	1.000000	0.958115	0.977783	0.958548

## 3 Features contribute to R<sup>2</sup> of ~97%



## **Correlation (and Multicollinearity)**

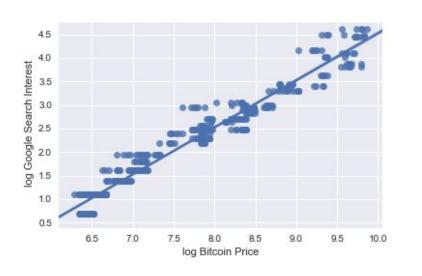


## Google Search Interest at R<sup>2</sup> of 95%

#### **Correlated Relationship**



#### Single Feature Regression



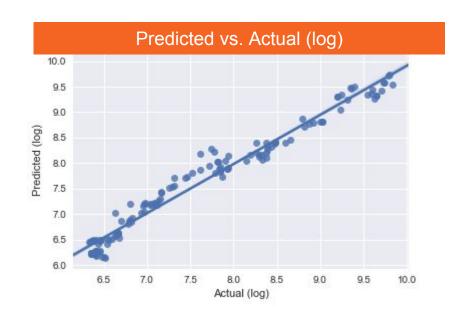
## Nasdaq Index at R<sup>2</sup> of 92%

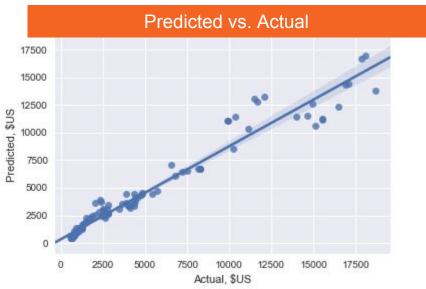


### Single Feature Regression



## **Linear Regression**





## **Key Insights / Takeaways**

### Bitcoin price a function of:

- 1. Google Search Interest
- 2. Nasdaq Composite Index
- 3. Network Transaction Fees

Google Search Interest may be leading indicator

Nasdaq Composite on "bull run" for history of Bitcoin; what happens when the "bear" comes?

## **Next Steps**

- Further explore seasonality in residual data
- Reconfigure model as a time series analysis, complete with price prediction for a certain timeline (requires adjustment of cross-validation)
- Explore social media sentiment (e.g. Twitter) as leading indicator



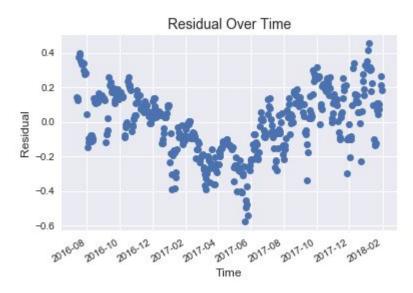
# Appendix

## **Linear Regression**

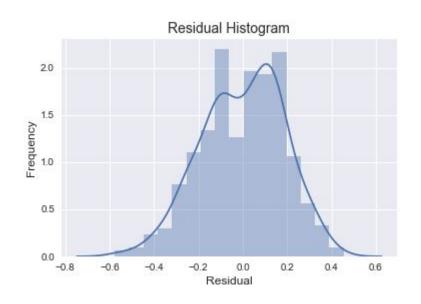
#### With Limits



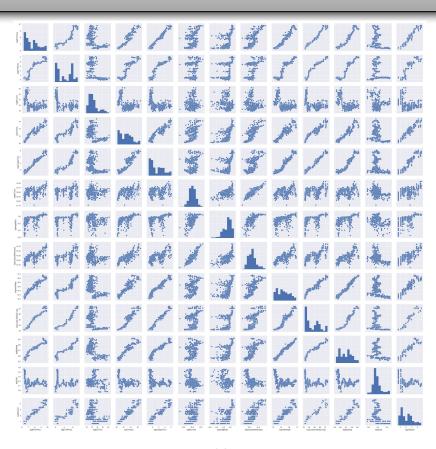
### Residuals



Heteroscedasticity: non-random residual distribution



### **Feature Universe**



## 3 Features contribute to R<sup>2</sup> of ~97%

Dep. Va	riable:		logBTC	Price		R-squared:	0.968
N	/lodel:			OLS	Adj.	Adj. R-squared:	
Me	ethod:		Least Sc	uares		F-statistic:	4647.
	Date:	Thu	ı, 01 Feb	2018	Prob (	F-statistic):	0.00
	Time:		11:	:40:55	Log-	Likelihood:	127.12
No. Observa	itions:			465		AIC:	-246.2
Df Resi	duals:			461		BIC:	-229.7
Df N	/lodel:			3			
Covariance	Туре:		nonr	obust			
		_					
	C	oef	std err	t	P> 1	[0.025	0.975]
Intercept	-17.44	146	2.346	-7.435	0.000	0 -22.055	-12.834
logInterest	0.58	379	0.031	19.129	0.00	0.528	0.648
logNasdaq	2.6	151	0.286	9.142	0.00	2.053	3.177
logTxFees	0.08	326	0.020	4.060	0.00	0.043	0.123
Omnik	ous:	7.400	Dur	bin-Wa	tson:	0.203	
Prob(Omnib	us):	0.02	Jarqu	ıe-Bera	(JB):	7.008	
Sk	ew: -	0.25	5	Prob	(JB):	0.0301	
Kurto	sis:	2.680	)	Cond	. No.	4.31e+03	

	Model Scorecard
Accuracy	R <sup>2</sup> of 96.8%, explaining randomness left in model (SSE) as proportion to variation in data (SST)
Fit	F-statistic P-value very small; data too extreme to fit model by chance alone High Log-Likelihood, low AIC and BIC indicate good fit
Variable contribution	T-test indicates high contribution of variables
Normal Distribution	Omnibus and Jarque-Bera significant at < 0.05, indicating no exact normal distribution of $\epsilon$
Autocorrelation	<b>Durbin-Watson</b> indicates slight positive autocorrelation (common in time series data)
Multicollinearity	Condition Number high, implying matrix may not have a unique, well defined solution

## Google Search Interest at R<sup>2</sup> of 95%

Dep. Va	riable:	logB	TCPrice	F	R-squared	: (	0.954
1	Model:		OLS	Adj. F	R-squared	: (	0.954
Me	ethod:	Least 9	Squares		F-statistic	: 9	9643
	Date:	Tue, 30 J	an 2018	Prob (F	-statistic)	: 3.37e	-31
	Time:	1	7:04:46	Log-l	ikelihood	: 46	6.288
No. Observa	itions:		463		AIC	: -8	38.58
Df Resi	duals:		461		BIC	: -8	30.30
Df N	Model:		1				
Covariance	Туре:	no	nrobust				
	coef	std err		t P> t	[0.025	0.9751	
Intercent	5.5065	0.023	236.676	•	-	5.552	
Intercept	5.5065	0.023	230.070	0.000	5.401	5.552	
logInterest	0.9640	0.010	98.199	0.000	0.945	0.983	
		440 -		*	0.050		
Omnik	ous: 22	.440 <b>[</b>	Durbin-W	atson:	0.259		
Prob(Omnib	<b>us):</b> 0	.000 <b>Ja</b>	rque-Ber	a (JB):	24.987		
Sk	<b>ew:</b> -0	.504	Pro	b(JB):	3.75e-06		
Kurto	sis: 3	.530	Cor	ıd. No.	6.21		

	Model Scorecard
Accuracy	R <sup>2</sup> of 95.4%, explaining randomness left in model (SSE) as proportion to variation in data (SST)
Fit	F-statistic P-value very small; data too extreme to fit model by chance alone High Log-Likelihood, low AIC and BIC indicate good fit
Variable contribution	T-test indicates high contribution of variables
Normal Distribution	Omnibus and Jarque-Bera significant at < 0.000, indicating no exact normal distribution of $\epsilon$
Autocorrelation	<b>Durbin-Watson</b> indicates slight positive autocorrelation (common in time series data)
Multicollinearity	Condition Number low implying unique, well-defined solution

## Nasdaq Index at R<sup>2</sup> of 92%

Dep. Variable	: 1	ogBTCI	Price		R-s	quared:	0.918
Model	:		OLS	A	dj. R-s	quared:	0.918
Method	: Le	ast Squ	uares		F-s	tatistic:	5145.
Date	: Tue,	30 Jan	2018	Pro	b (F-st	tatistic):	3.18e-252
Time	:	17:0	8:24	Lo	g-Lik	elihood:	-90.077
No. Observations	:		463			AIC:	184.2
Df Residuals	:		461			BIC:	192.4
Df Model	:		1				
Covariance Type	:	nonro	bust				
	coef s	td err		t	P> t	[0.025	0.975]
Intercept -75.3	8836	1.156	-65.1	88	0.000	-77.656	-73.111
logNasdaq 9.5	5430	0.133	71.7	30	0.000	9.282	9.804
Omnibus:	21.858	Dur	bin-W	atso	n:	0.060	
Prob(Omnibus):	0.000	Jarqu	ıe-Ber	a (JE	3):	23.621	
Skew:	0.533		Pro	b(JE	<b>3):</b> 7.	42e-06	
Kurtosis:	3.297		Cor	d. N	о.	744.	

	Model Scorecard
Accuracy	R <sup>2</sup> of 91.8%, explaining randomness left in model (SSE) as proportion to variation in data (SST)
Fit	F-statistic P-value very small; data too extreme to fit model by chance alone Low Log-Likelihood / High AIC and BIC may indicate poor fit
Variable contribution	T-test indicates high contribution of variables
Normal Distribution	Omnibus and Jarque-Bera significant at < 0.000, indicating no exact normal distribution of ε
Autocorrelation	<b>Durbin-Watson</b> indicates slight positive autocorrelation (common in time series data)
Multicollinearity	Condition Number high, implying matrix may not have a unique, well defined solution

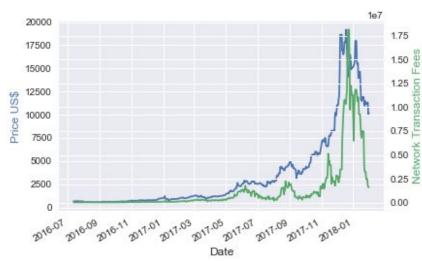
## Transaction Fees at R<sup>2</sup> of 92%

Dep. Va	riable	e:	logB	TCPrice		R-squar	red:	0.91	18
ı	Mode	l:		OLS	Adj. R-squ		red:	0.9	18
М	ethod	d:	Least	Squares		F-statis	tic:	521	3.
	Date	e: T	hu, 01 F	eb 2018	Prob	(F-statis	tic):	1.86e-2	54
	Time	e:	1	9:16:46	Log	-Likeliho	od:	-91.83	37
No. Observa	ations	s:		466		ı	AIC:	187	.7
Df Res	idual	s:		464		E	BIC:	196	0.6
Df I	Mode	l:		1					
Covariance	• Тур	e:	nc	nrobust					
	С	oef	std err	t	P> t	[0.025	0.97	<b>'5</b> ]	
Intercept	-0.06	698	0.107	-0.654	0.513	-0.280	0.1	40	
logTxFees	0.60	011	0.008	72.199	0.000	0.585	0.6	17	
Omni	hue	3.70	05 <b>D</b> u	ırbin-Wa	teon:	0.214			
<u> </u>		0.1				3.662			
Prob(Omnik	jusj.	0.1	Jaro	ue-Bera	(JD):	3.002			
SI	(ew:	0.18	B1	Prob	(JB):	0.160			
Kurto	sis:	2.7	59	Cond	l. No.	101.			

	Model Scorecard
Accuracy	R <sup>2</sup> of 91.8%, explaining randomness left in model (SSE) as proportion to variation in data (SST)
Fit	F-statistic P-value very small; data too extreme to fit model by chance alone Low Log-Likelihood / High AIC and BIC may indicate poor fit
Variable contribution	T-test indicates high contribution of variables
Normal Distribution	Omnibus / JB >0.05 may indicate exact normal distribution of $\epsilon$
Autocorrelation	<b>Durbin-Watson</b> indicates slight positive autocorrelation (common in time series data)
Multicollinearity	Condition Number high, implying matrix may not have a unique, well defined solution

## Transaction Fees at R<sup>2</sup> of 92%

### Correlated Relationship



#### Single Feature Regression



## Methodology

- Explored various features which may have correlation with the price of Bitcoin
- Train-Test split executed using standard Linear Regression methodology rather than as a Time Series, per project requirements
  - The difference stems from the "shuffling", or randomizing, of data in determining split; time series data would not be randomized
  - As such, the price of bitcoin is predicted at a moment in time, rather than predicted on a time series basis

## **Scraping**

- Twitter sentiment analysis
- Coinmarketcap price history
- Bitcoin futures

## **Future Implementations**

- Adjust model to predict price via time series
- Explore valuation of other cryptocurrencies using bitcoin metrics
- Compile sufficient real-time data to execute twitter sentiment analysis over reasonable timeframe (ie 1 week+)

### Other analyses to explore:

- Usage by country
- Bitcoin trading by exchange
- Bitcoin trading by currency
- Use of leverage
- SEASONALITY

### **Process**

- Brainstorm possible interesting correlations, including:
  - Bitcoin-specific metrics (ie transaction fees, volume, block size, hash rate)
  - Bitcoin futures
  - Google search interest
  - Social media sentiment (via Twitter)

### Sources / Resources

- www.coinmarketcap.com
- www.quandl.com
- https://trends.google.com/trends/
- https://marcobonzanini.com/2015/03/09/mining-twitter-data-with-python-part-2/
- http://cs229.stanford.edu/proj2015/029\_report.pdf
- http://text-processing.com/
- https://trends.google.com/trends/explore?q=bitcoin.ethereum