

Big Data in Finance

Sanjiv R. Das
Santa Clara University

UMich - CFLP
Ann Arbor
October 2016

Sanjiv Ranjan Das

- My Ph.D. is in Finance (NYU).
- MS in Computer Science (Berkeley).
- Various areas of work:
 - ① Stochastic processes, continuous time finance.
 - ② Derivatives modeling and risk management.
 - ③ Portfolio theory.
 - ④ Asset pricing.
 - ⑤ Venture capital.
 - ⑥ Credit risk models.
 - ⑦ Machine learning.
 - ⑧ Text Analytics.
 - ⑨ Integrating models of networks and risk with big data.

A Philosophy for Big Data Analytics in Finance

- Using **theory** to develop models to apply big data.
- Question/problems are **primary**, data is secondary.
- **Simplicity, transparency** of models fosters implementability.
- **Dimension reduction** for sufficient statistics.
- Analytics per se is **multidisciplinary**.
- **Disparate** data is the norm.

Three Example Topics

- ① Financial networks.
- ② Zero-revelation prediction of bank malaise.
- ③ Market illiquidity.

Midas Project: Overview

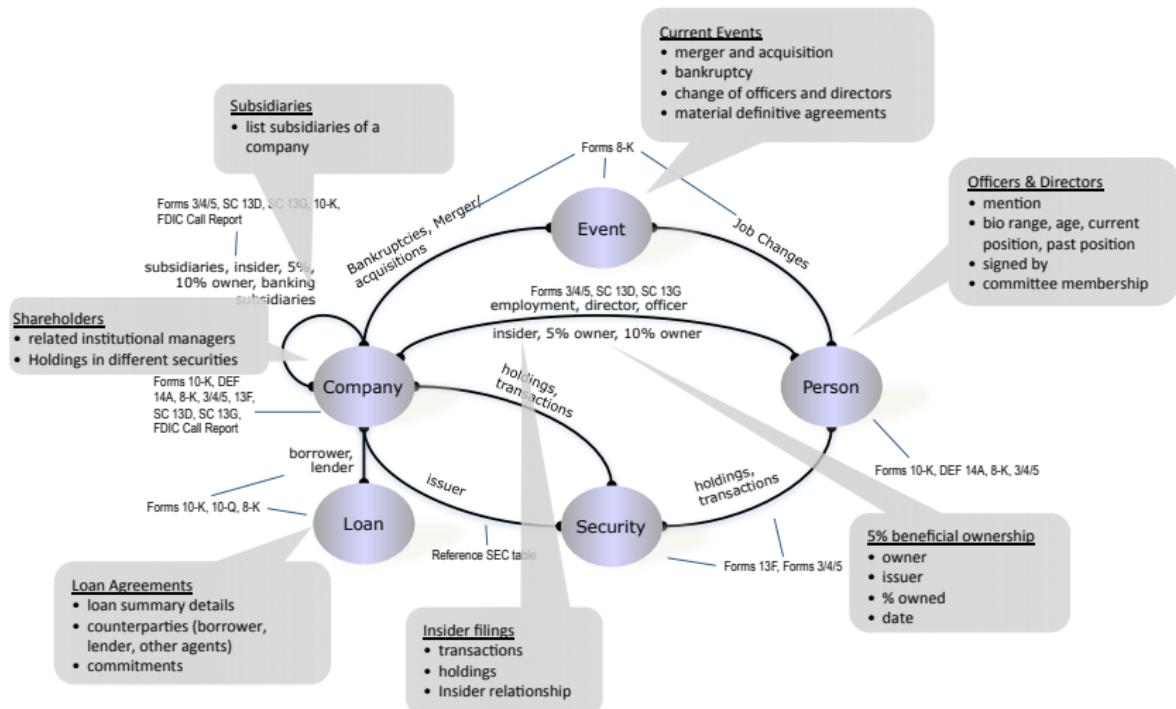
Joint work with IBM Almaden¹

- Focus on financial companies that are the domain for systemic risk (SIFIs).
- Extract information from unstructured text (filings).
- Information can be analyzed at the institutional level or aggregated system-wide.
- Applications: Systemic risk metrics; governance.
- Technology: information extraction (IE), entity resolution, mapping and fusion, scalable Hadoop architecture.

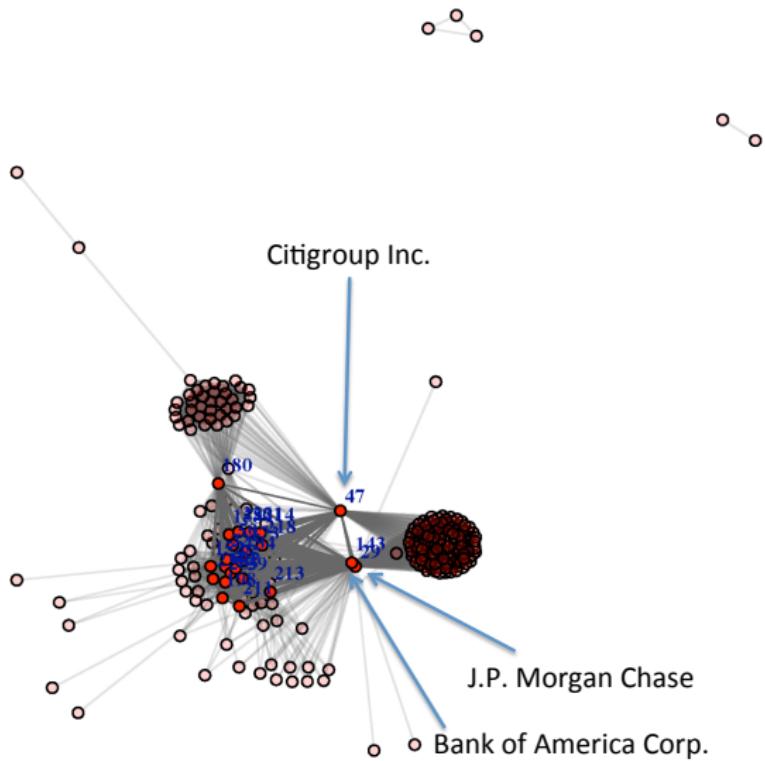
¹ "Extracting, Linking and Integrating Data from Public Sources: A Financial Case Study," (2011), (with Douglas Burdick, Mauricio A. Hernandez, Howard Ho, Georgia Koutrika, Rajasekar Krishnamurthy, Lucian Popa, Ioana Stanoi, Shivakumar Vaithyanathan), *IEEE Data Engineering Bulletin*, 34(3), 60-67. [Proceedings WWW2010, April 26-30, 2010, Raleigh, North Carolina.]

Data

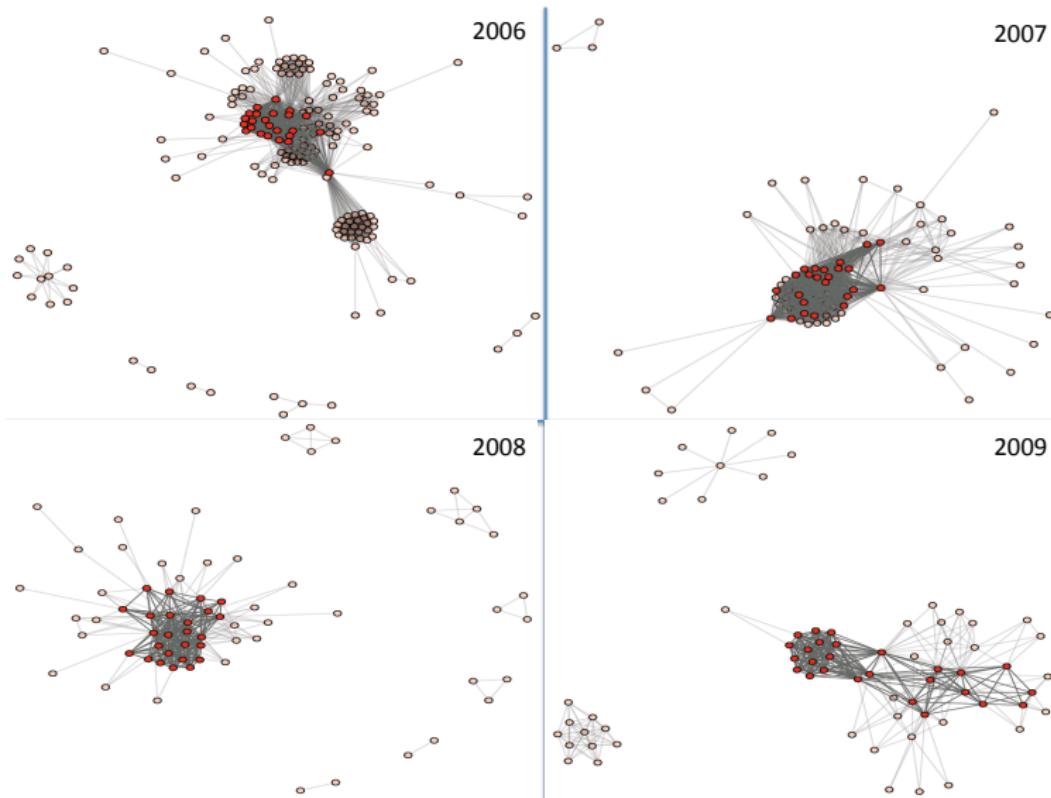
Midas provides Analytical Insights into company relationships by exposing information concepts and relationships within extracted concepts



Loan Network 2005



Loan Network 2006–2009



Systemically Important Financial Institutions (SIFIs)

Year	# Colending banks	# Coloans	Colending pairs	$R = E(d^2)/E(d)$	Diam.
2005	241	75	10997	137.91	5
2006	171	95	4420	172.45	5
2007	85	49	1793	73.62	4
2008	69	84	681	68.14	4
2009	69	42	598	35.35	4

(Year = 2005)		
Node #	Financial Institution	Normalized Centrality
143	J P Morgan Chase & Co.	1.000
29	Bank of America Corp.	0.926
47	Citigroup Inc.	0.639
85	Deutsche Bank Ag New York Branch	0.636
225	Wachovia Bank NA	0.617
235	The Bank of New York	0.573
134	Hsbc Bank USA	0.530
39	Barclays Bank Plc	0.530
152	Keycorp	0.524
241	The Royal Bank of Scotland Plc	0.523
6	Abn Amro Bank N.V.	0.448
173	Merrill Lynch Bank USA	0.374
198	PNC Financial Services Group Inc	0.372
180	Morgan Stanley	0.362
42	Bnp Paribas	0.337
205	Royal Bank of Canada	0.289
236	The Bank of Nova Scotia	0.289
218	U.S. Bank NA	0.284
50	Calyon New York Branch	0.273
158	Lehman Brothers Bank Fsb	0.270
213	Sumitomo Mitsui Banking	0.236
214	Suntrust Banks Inc	0.232
221	UBS Loan Finance Llc	0.221
211	State Street Corp	0.210
228	Wells Fargo Bank NA	0.198

Risk Networks: Definitions and Risk Score

- Assume n nodes, i.e., firms, or “assets.”
- Let $E \in R^{n \times n}$ be a well-defined adjacency matrix. This quantifies the influence of each node on another.
- E may be portrayed as a directed graph, i.e., $E_{ij} \neq E_{ji}$.
 $E_{jj} = 1$; $E_{ij} \in \{0, 1\}$.
- C is a $(n \times 1)$ risk vector that defines the risk score for each asset.
- We define the “risk score” as

$$S = \sqrt{C^\top E C}$$

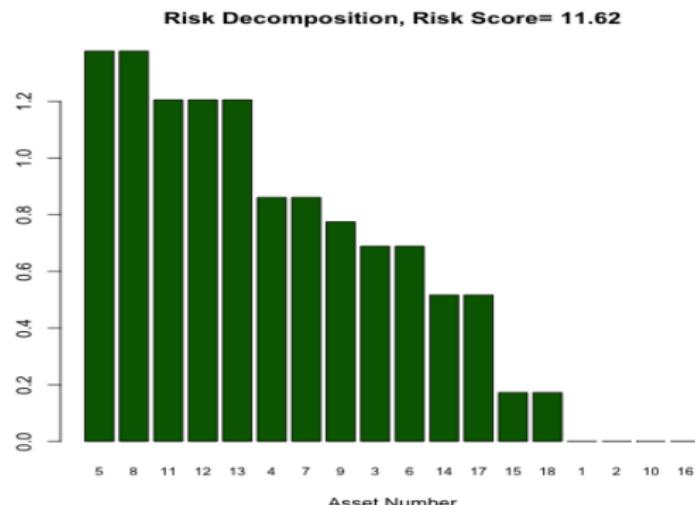
- $S(C, E)$ is linear homogenous in C .

Risk Decomposition

- ① Exploits the homogeneity of degree one property of S .
- ② Risk decomposition (using Euler's formula):

$$S = \frac{\partial S}{\partial C_1} C_1 + \frac{\partial S}{\partial C_2} C_2 + \dots + \frac{\partial S}{\partial C_n} C_n$$

- ③ Plot:



Systemic Risk in Indian Banks

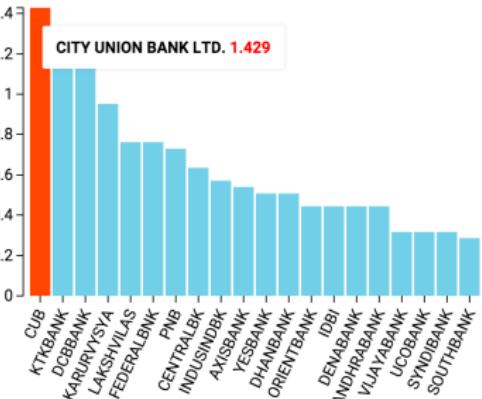
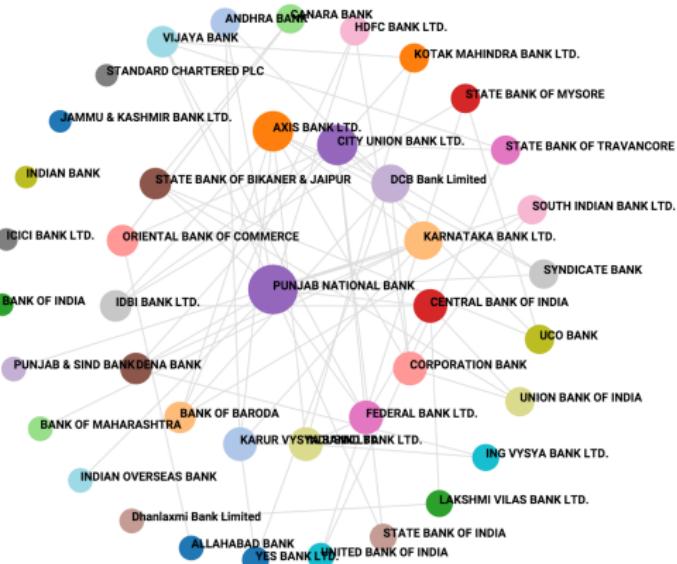
Fragility

2.91

Systemic Risk Score

15.75

Risk Decomposition



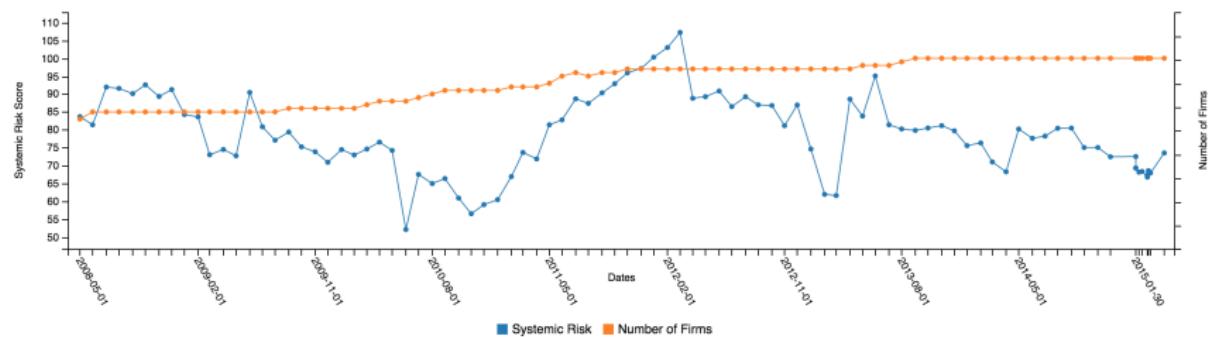
Systemic Risk in India over time

Systemic Risk Dashboard

Segment	Firms	Parameter	Date	<button>Submit</button>
<input type="text"/>	<input type="text"/>	<input type="text"/> Network Plot and	<input type="text"/>	

[SYSTEM CONNECTEDNESS](#)[INDIVIDUAL RISK METRICS](#)[SYSTEMIC RISK TREND](#)[DEFINITIONS](#)

Systemic Risk Trend

[Update](#)

Stochastic Risk Networks in a Structural Framework

- We use the Merton (1974) model to extend the static Das (2016) model to a stochastic network setting (Das, Kim, Ostrov, 2016, wip).
- We extend each node's properties to including size, in addition to the credit score.
- To do this we normalize the S measure.
- This model can be calibrated using the same methods used for the Merton model, or variants such as the Moody's KMV model.

Definitions

Model Data (standard Merton model inputs) for each firm:

- Equity price = $\mathbf{s} = \{s_1, s_2, \dots, s_n\}$
- Equity volatility = $\sigma = \{\sigma_1, \sigma_2, \dots, \sigma_n\}$
- Number of shares = $\mathbf{m} = \{m_1, m_2, \dots, m_n\}$
- Risk free rate = r

Model Variables (all derived from the Merton model):

- n = number of banks in the system
- $\mathbf{a} = n$ -vector with components a_i that represent the assets in bank i (derived from s, σ, m, r).
- $\lambda = n$ -vector with components λ_i that represent the average yearly chance of bank i defaulting (from s, σ, r).
- $\mathbf{E} = n \times n$ matrix with components E_{ij} that represent the probability that if bank j defaults, it will cause bank i to default (from s, σ, r).

Model

- Define \mathbf{c} to be an n -vector with components c_i that represent bank i 's credit risk. More specifically, we define

$$\mathbf{c} = \mathbf{a} \odot \boldsymbol{\lambda},$$

where \odot represents component multiplication; that is, $c_i = a_i \lambda_i$.

- The aggregate systemic risk created by the n banks in our system is

$$R = \frac{\sqrt{\mathbf{c}^\top \mathbf{E} \mathbf{c}}}{\mathbf{1}^\top \mathbf{a}}, \quad (1)$$

where $\mathbf{1}$ is an n -vector of ones, so the denominator $\mathbf{1}^\top \mathbf{a} = \sum_{i=1}^n a_i$ represents the total assets in the n banks.

- r is linear homogenous in $\boldsymbol{\lambda}$.

Network of Top 50 Financial Institutions

Dynamic Risk Networks (2016)

Upload a .csv file with data for banks (in columns) and attributes (in rows) [Click here to see format](#)

File input

no file selected

Systemic Risk Score

5.5667

Please refer to following Paper published for some details [Matrix Metrics: Network-Based Systemic Risk Scoring](#)

Avg Correlation:

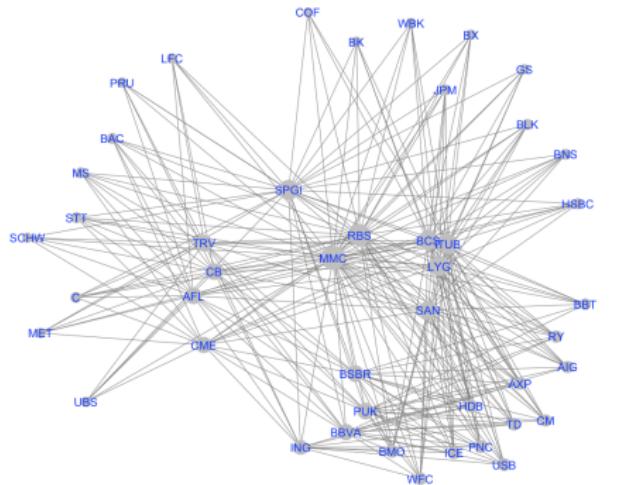


Correlation threshold:



Network Graph

Risk Decomposition



Computational Properties

$$R = \frac{\sqrt{\mathbf{c}^\top \mathbf{E} \mathbf{c}}}{\mathbf{1}^\top \mathbf{a}}, \quad \mathbf{c} = \mathbf{a} \odot \boldsymbol{\lambda}$$

- R is linear homogeneous in $\boldsymbol{\lambda}$: Let α be any scalar constant. If we replace $\boldsymbol{\lambda}$ with $\alpha\boldsymbol{\lambda}$, it immediately follows that \mathbf{c} is replaced by $\alpha\mathbf{c}$, and, by our equation for R , we see that R is replaced by αR .
- Sensitivity of R to changes in $\boldsymbol{\lambda}$: Differentiating our equation for R with respect to $\boldsymbol{\lambda}$

$$\frac{\partial R}{\partial \boldsymbol{\lambda}} = \frac{1}{2} \frac{\mathbf{a} \odot [(\mathbf{E} + \mathbf{E}^T) \mathbf{c}]}{\mathbf{1}^\top \mathbf{a} \sqrt{\mathbf{c}^\top \mathbf{E} \mathbf{c}}}$$

whose components represent the sensitivity of R to changes in each bank's value of λ . This is the basis of **Risk Decomposition**, equal to $(\frac{\partial R}{\partial \lambda} \cdot \boldsymbol{\lambda})$, a vector containing each bank's contribution to R .

Risk Decomposition: 50 Financial Institutions

Dynamic Risk Networks (2016)

Upload a .csv file with data for banks (in columns) and attributes (in rows) [Click here to see format](#)

File input

Choose File no file selected

Compute Scores

Systemic Risk Score

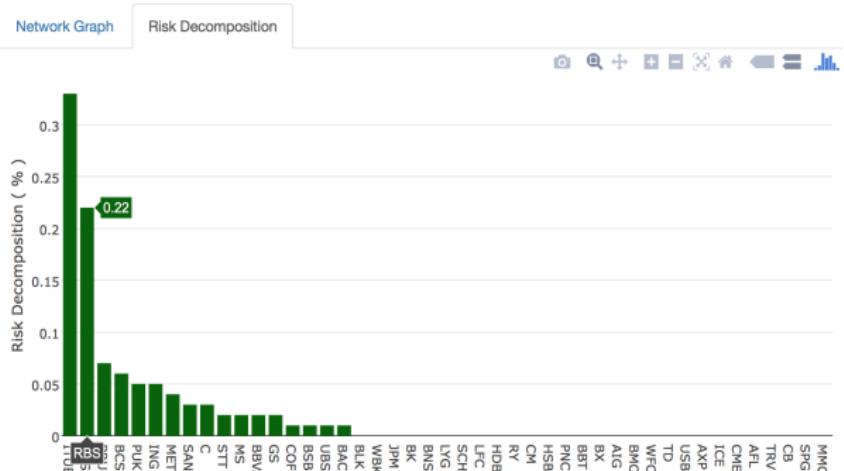
5.5667

Please refer to following Paper published for some details [Matrix Metrics: Network-Based Systemic Risk Scoring](#)

Avg Correlation:



Correlation threshold:

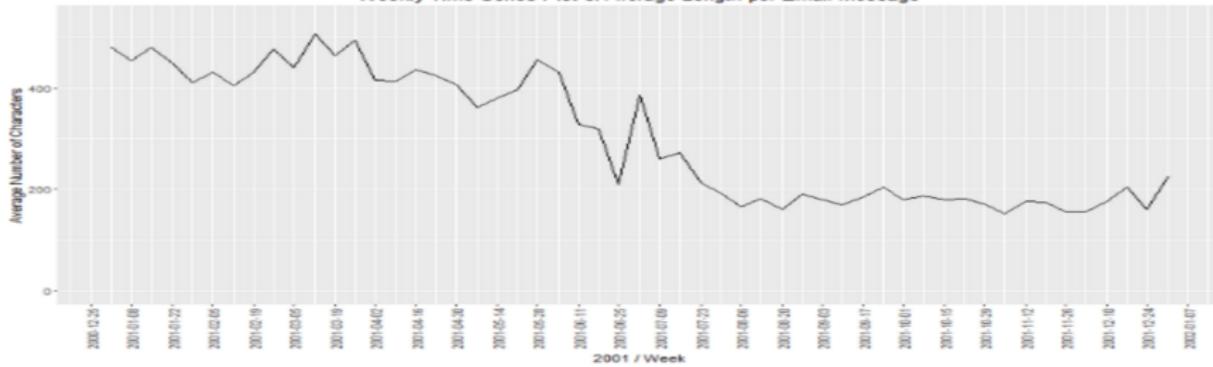
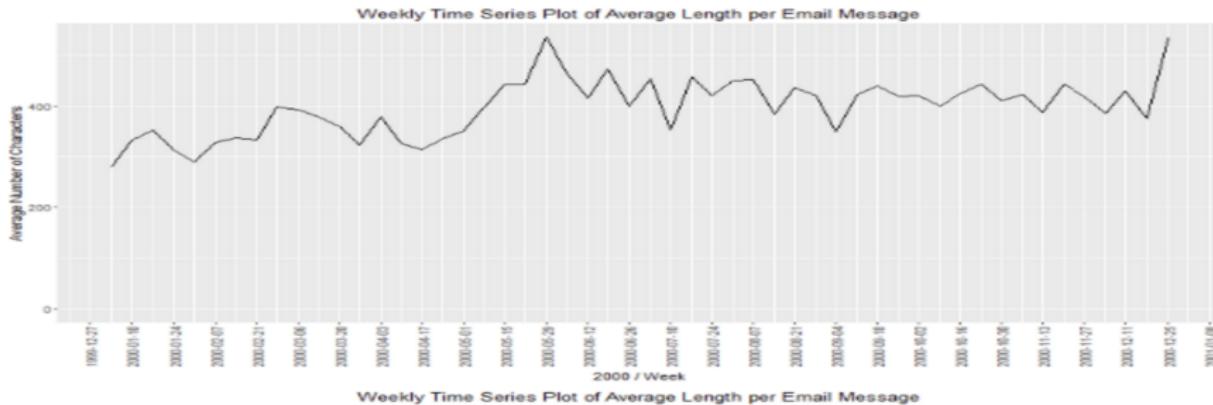


Analyzing Emails for Early Warnings of Failure

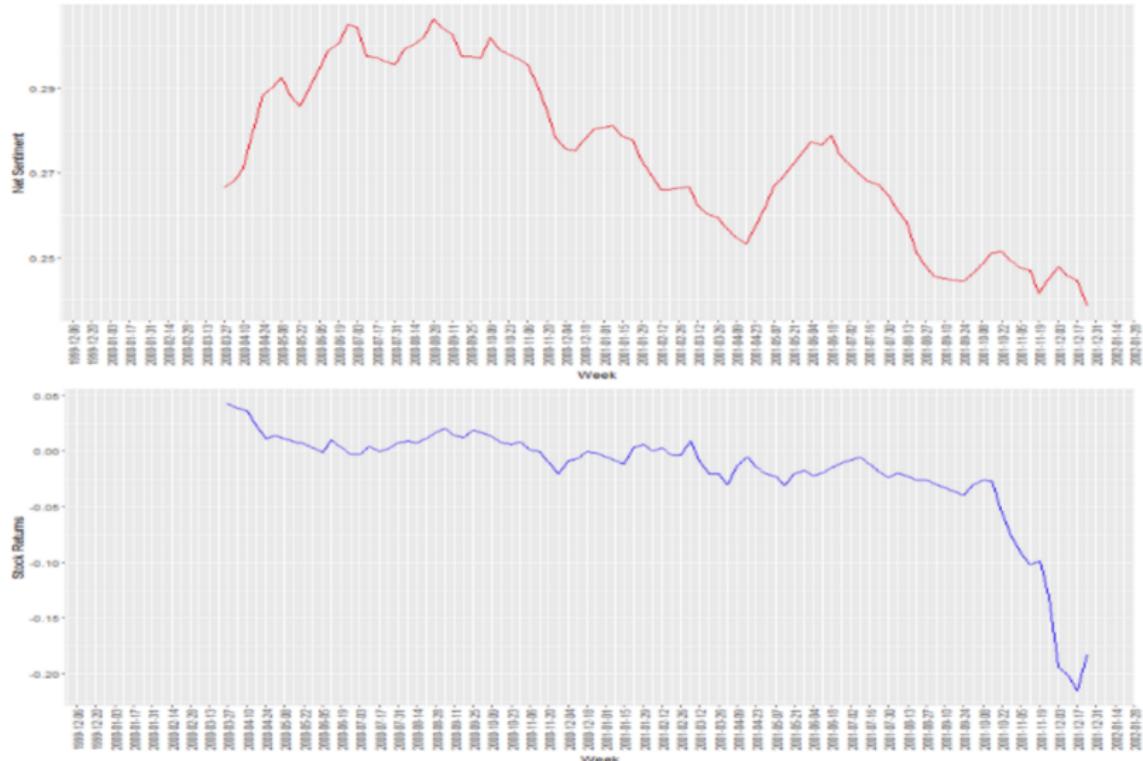
Title: "Zero-Revelation Linguistic Regulation: Detecting Risk Through Corporate Emails and News"

- Financials are often delayed indicators of corporate quality.
- Internal discussion may be used as an early warning system for upcoming corporate malaise.
- Emails have the potential to predict such events.
- Software can analyze vast quantities of textual data not amenable to human processing.
- Corporate senior management may also use these analyses to better predict and manage impending crisis for their firms.
- The approach requires zero revelation of emails.

Enron: Email Length



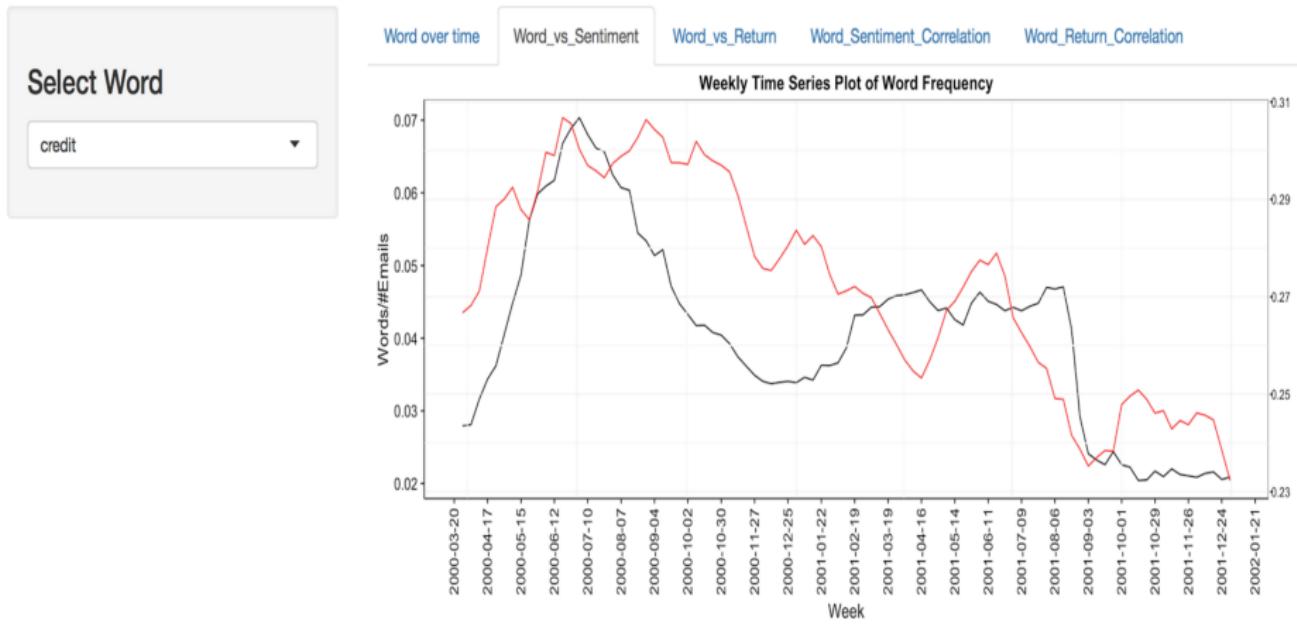
Enron: Sentiment and Returns



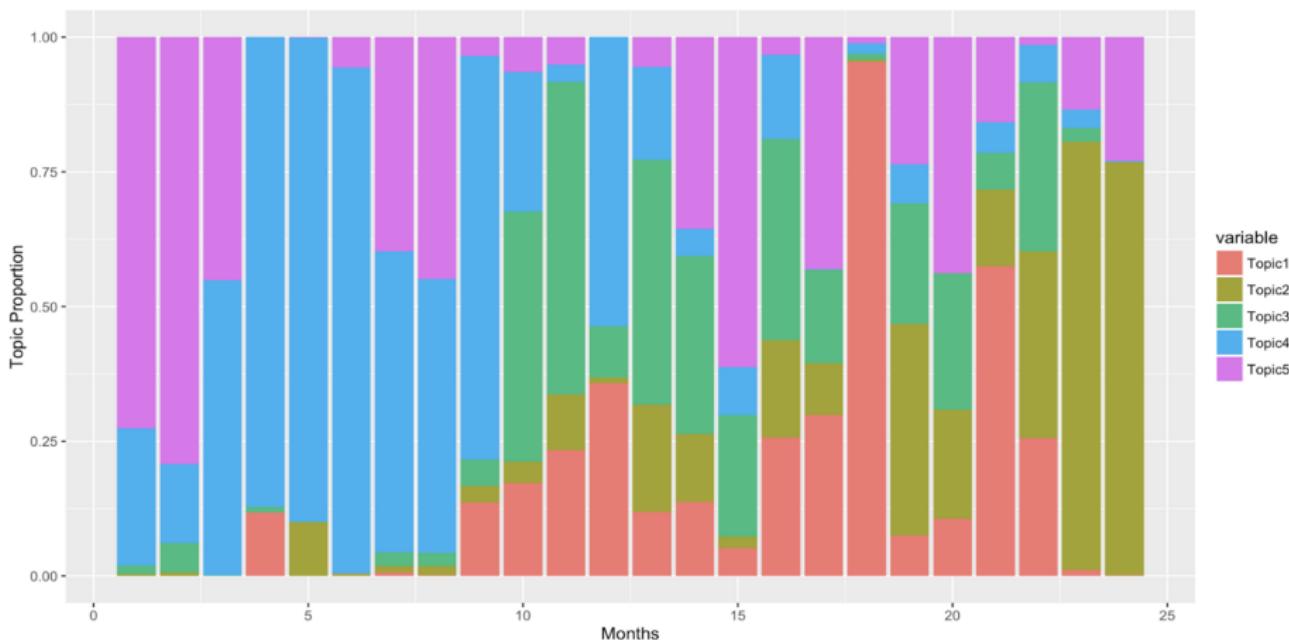
Enron: Returns and Characteristics

Variable	Coefficient Estimate (<i>t</i> -statistic)			
	(1)	(2)	(3)	(4)
<i>MA Net Sentiment</i> t	XXX*** (XXX)	0.575 (0.63)	2.330*** (3.14)	-1.397 (-1.25)
<i>MA Email Length</i> t		0.584*** (2.97)		1.046*** (4.19)
<i>MA Total Emails</i> t			-0.004 (-0.10)	-0.131*** (-2.83)
<i>Intercept</i>		-0.406* (-1.93)	-0.671*** (-3.08)	0.117 (0.43)
Adjusted <i>R</i> -squared	XXX		0.09	0.24
Number of observations	88	88	88	88

Enron: WordPlay



Enron: Topic Analysis



BILLIQ: An Index-Based Measure of Illiquidity

Uses option pricing theory to derive a measure for the cost of immediacy:

$$BILLIQ = -10,000 \times \ln \left[\frac{NAV}{NAV + |ETF - NAV|} \right]$$

~ /Google Drive/Papers/ETFIq/RBILLIQcode - Shiny

<http://127.0.0.1:5808> | [Open in Browser](#) | [Shiny](#)

[Republieh](#) ▾

Index-Based Illiquidity

Input ETF Ticker

Submit

Example of ETF tickers are: LQD, HYG, CSJ, CFT, CIU, AGG, GBF, GVI, MBB, EMB, IIV, BIV, BLV, BND, BSV, etc.

Yield = 3.22

Price = 122.44

NAV = 122.12

BILLIQ = 26.1694621053458 (bps)

The paper that derives this measure of illiquidity is:

George Chacko, Sanjiv Das, Rong Fan (2016), An Index-Based Measure of Liquidity, Journal of Banking and Finance, v68, 162-178.