

Text Selection

Bryan Kelly

Yale University

Asaf Manela

Washington University in St. Louis

Alan Moreira

University of Rochester

October 2018

Motivation

- ▶ Digital text is increasingly available to social scientists
 - ▶ Newspapers, blogs, regulatory filings, congressional records ...
- ▶ Unlike data often used by economists
 - ▶ Text is ultra high-dimensional
 - ▶ Phrase counts are sparse
- ▶ Statistical learning from text requires
 - ▶ Machine learning techniques
 - ▶ Scalable algorithms

This paper

- ▶ Text is often selected by journalists, speechwriters, and others who cater to an audience with limited attention
- ▶ Hurdle Distributed Multiple Regression (HDMR)
 - ▶ Highly scalable approach to inference from big counts data
 - ▶ Includes an economically-motivated **selection equation**
 - ▶ Especially useful when cover/no-cover choice is separate or more interesting than coverage quantity
- ▶ Applications using **newspaper coverage for prediction**
 1. Backcast intermediary capital ratio (He-Kelly-Manela 2017 JFE)
 2. Forecast macroeconomic series (Stock-Watson 2012 JBES)

Related literature

- ▶ We extend machinery developed by Taddy (2012, 2015, 2016) to text selection
 - ▶ Layer economically-motivated [hurdle / selection equation](#) on his Distributed Multinomial Regression (DMR)
 - ▶ Find advantage of HDMR over DMR increases with sparsity
- ▶ Provide new tools to literatures in economics and finance
 - ▶ Finance and media: Antweiler-Frank (2004), Tetlock (2007, 2011), Fang-Peress (2009), Engelberg-Parsons (2011), Dougal et al (2012), Peress (2014), Manela (2014), Fedyk (2018)
 - ▶ Text-based uncertainty: Baker-Bloom-Davis (2016), Manela-Moreira (2017), Hassan et al (2017)
 - ▶ Polarization: Gentzkow-Shapiro (2006), Gentzkow-Shapiro-Taddy
- ▶ Can better control and learn from high-dimensional content

Text data is inherently high-dimensional

Documents

1: Digital text is available.
2: Text is selected!
⋮

⇒

Document-term matrix c

	digital text	text is	is available	is selected	⋮
1:	1	1	1	0	
2:	0	1	0	1	
		⋮			⋮

Text data is inherently high-dimensional

Documents

1: Digital text is available.	⇒	Document-term matrix c				
2: Text is selected!		digital text	text is	is available	is selected	...
⋮		1:	1	1	1	0
		2:	0	1	0	1
			⋮			⋮

Text regression is prone to overfit

- ▶ \mathbf{c}_i vector of counts in d categories for observation i
 - ▶ e.g. c_{ij} is date i newspaper mentions of phrase j (“world war”)
- ▶ \mathbf{v}_i vector of p covariates
 - ▶ e.g. intermediary capital ratio, realized variance on date i
- ▶ Let $v_{iy} \in \mathbf{v}_i$ be a target variable
 - ▶ e.g. intermediary capital ratio
- ▶ Because $d \gg n$, we cannot run an OLS regression

$$v_{iy} = \beta_0 + [\mathbf{c}_i, \mathbf{v}_{i,-y}]' \boldsymbol{\beta} + \varepsilon_i$$

Text inverse regression

- ▶ A text inverse regression approach would instead

1. Regress word counts on covariates

$$\mathbf{c}_i = \lambda (\alpha_j + \mathbf{v}'_i \boldsymbol{\varphi}_j) + \mathbf{v}_i \quad (\text{backward regression})$$

2. Construct low dimensional projection into v_{iy} direction

$$z_{iy} \equiv \sum_j \hat{\varphi}_{jy} c_{ij} \quad (\text{sufficient reduction projection})$$

3. Regress target variable on z_{iy} and other covariates

$$v_{iy} = \beta_0 + [z_{iy}, \mathbf{v}_{i,-y}]' \boldsymbol{\beta} + \varepsilon_i \quad (\text{forward regression})$$

- ▶ $d + p - 1$ dimensional regression reduced to $p + 1$ dimensional!
- ▶ z_{iy} summarizes all textual information relevant for prediction

Text inverse regression

- ▶ A text inverse regression approach would instead

1. Regress word counts on covariates

$$c_i = \lambda (\alpha_j + \mathbf{v}'_i \boldsymbol{\varphi}_j) + \mathbf{v}_i \quad (\text{backward regression})$$

2. Construct low dimensional projection into v_{iy} direction

$$z_{iy} \equiv \sum_j \hat{\varphi}_{jy} c_{ij} \quad (\text{sufficient reduction projection})$$

3. Regress target variable on z_{iy} and other covariates

$$v_{iy} = \beta_0 + [z_{iy}, \mathbf{v}_{i,-y}]' \boldsymbol{\beta} + \varepsilon_i \quad (\text{forward regression})$$

- ▶ $d + p - 1$ dimensional regression reduced to $p + 1$ dimensional!
- ▶ z_{iy} summarizes all textual information relevant for prediction

Text inverse regression

- ▶ A text inverse regression approach would instead

1. Regress word counts on covariates

$$\mathbf{c}_i = \lambda (\alpha_j + \mathbf{v}'_i \boldsymbol{\varphi}_j) + \mathbf{v}_i \quad (\text{backward regression})$$

2. Construct low dimensional projection into v_{iy} direction

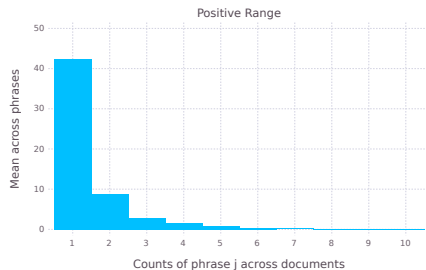
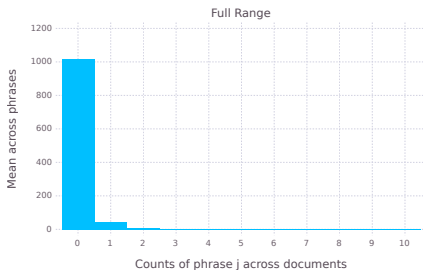
$$z_{iy} \equiv \sum_j \hat{\varphi}_{jy} c_{ij} \quad (\text{sufficient reduction projection})$$

3. Regress target variable on z_{iy} and other covariates

$$v_{iy} = \beta_0 + [z_{iy}, \mathbf{v}_{i,-y}]' \boldsymbol{\beta} + \varepsilon_i \quad (\text{forward regression})$$

- ▶ $d + p - 1$ dimensional regression reduced to $p + 1$ dimensional!
- ▶ z_{iy} summarizes all textual information relevant for prediction

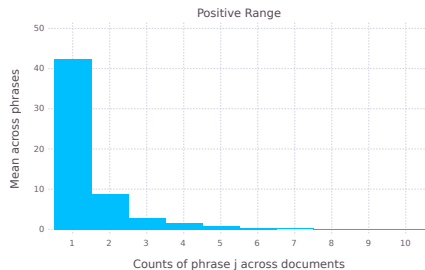
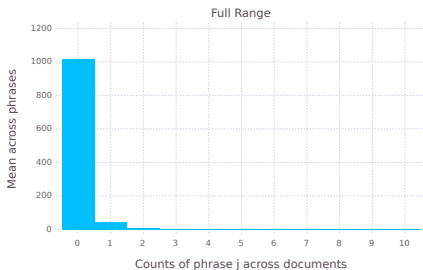
Why would we need a hurdle?



Wall Street Journal, monthly front page text, July 1926 to February 2016

- ▶ Statistics: hurdle better describes text data
 - ▶ Text data often has many more zeros than predicted by Poisson
- ▶ Economics: text is selected
 - ▶ Publishers cater to a boundedly rational reader (Gabaix, 2014)
 - ▶ Politicians select phrases that resonate with voters (Gentzkow-Shapiro-Taddy, 2017)
 - ▶ Censored or socially taboo words (Michel et al, 2011)
 - ▶ Fixed cost of introducing new terms, low marginal cost

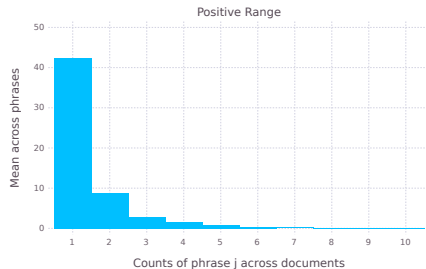
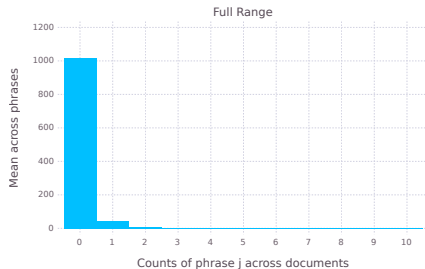
Why would we need a hurdle?



Wall Street Journal, monthly front page text, July 1926 to February 2016

- ▶ **Statistics: hurdle better describes text data**
 - ▶ Text data often has many more zeros than predicted by Poisson
- ▶ **Economics: text is selected**
 - ▶ Publishers cater to a boundedly rational reader (Gabaix, 2014)
 - ▶ Politicians select phrases that resonate with voters (Gentzkow-Shapiro-Taddy, 2017)
 - ▶ Censored or socially taboo words (Michel et al, 2011)
 - ▶ Fixed cost of introducing new terms, low marginal cost

Why would we need a hurdle?



Wall Street Journal, monthly front page text, July 1926 to February 2016

- ▶ Statistics: **hurdle better describes text data**
 - ▶ Text data often has many more zeros than predicted by Poisson
- ▶ Economics: **text is selected**
 - ▶ Publishers cater to a boundedly rational reader (Gabaix, 2014)
 - ▶ Politicians select phrases that resonate with voters (Gentzkow-Shapiro-Taddy, 2017)
 - ▶ Censored or socially taboo words (Michel et al, 2011)
 - ▶ Fixed cost of introducing new terms, low marginal cost

Text selection model

With sparse text, extensive margin may be more informative than intensive margin

- ▶ We suggest a text selection model instead

1. Two part text selection model for counts

$$\mathbf{h}_i^* = f(\kappa_j + \mathbf{w}'_i \boldsymbol{\delta}_j) + \boldsymbol{\omega}_i \quad (\text{Inclusion})$$

$$\mathbf{c}_i^* = \lambda(\alpha_j + \mathbf{v}'_i \boldsymbol{\varphi}_j) + \mathbf{v}_i \quad (\text{Repetition})$$

$$\mathbf{c}_i = \mathbf{c}_i^* \times \mathbf{1}(\mathbf{h}_i^* > 0) = \mathbf{c}_i^* \times \mathbf{h}_i \quad (\text{Observation})$$

2. Construct two low dimensional projections into v_{iy} ($= w_{iy}$)

$$\begin{array}{ll} z_{iy}^0 \equiv \sum_j \hat{\delta}_{jy} h_{ij} & z_{iy}^+ \equiv \sum_j \hat{\varphi}_{jy} c_{ij} \\ \text{Inclusion} & \text{Repetition} \end{array} \quad (\text{SR projections})$$

3. Regress target variable on z_{iy}^+ , z_{iy}^0 and other covariates

$$v_{iy} = \beta_0 + [z_{iy}^0, z_{iy}^+, \mathbf{w}_{i,-y}, \mathbf{v}_{i,-y}]' \boldsymbol{\beta} + \varepsilon_i \quad (\text{forward regression})$$

- ▶ $d + p - 1$ dimensional regression reduced to $p + 2$ dimensional!

Text selection model

With sparse text, extensive margin may be more informative than intensive margin

- ▶ We suggest a text selection model instead

1. Two part text selection model for counts

$$\mathbf{h}_i^* = f(\kappa_j + \mathbf{w}'_i \boldsymbol{\delta}_j) + \boldsymbol{\omega}_i \quad (\text{Inclusion})$$

$$\mathbf{c}_i^* = \lambda(\alpha_j + \mathbf{v}'_i \boldsymbol{\varphi}_j) + \mathbf{v}_i \quad (\text{Repetition})$$

$$\mathbf{c}_i = \mathbf{c}_i^* \times \mathbf{1}(\mathbf{h}_i^* > 0) = \mathbf{c}_i^* \times \mathbf{h}_i \quad (\text{Observation})$$

2. Construct **two** low dimensional **projections** into v_{iy} ($= w_{iy}$)

$$\begin{array}{ll} z_{iy}^0 \equiv \sum_j \hat{\delta}_{jy} h_{ij} & z_{iy}^+ \equiv \sum_j \hat{\varphi}_{jy} c_{ij} \\ \text{Inclusion} & \text{Repetition} \end{array} \quad (\text{SR projections})$$

3. Regress target variable on z_{iy}^+ , z_{iy}^0 and other covariates

$$v_{iy} = \beta_0 + [z_{iy}^0, z_{iy}^+, \mathbf{w}_{i,-y}, \mathbf{v}_{i,-y}]' \boldsymbol{\beta} + \varepsilon_i \quad (\text{forward regression})$$

- ▶ $d + p - 1$ dimensional regression reduced to $p + 2$ dimensional!

Text selection model

With sparse text, extensive margin may be more informative than intensive margin

- ▶ We suggest a text selection model instead

1. Two part text selection model for counts

$$\mathbf{h}_i^* = f(\kappa_j + \mathbf{w}'_i \boldsymbol{\delta}_j) + \boldsymbol{\omega}_i \quad (\text{Inclusion})$$

$$\mathbf{c}_i^* = \lambda(\alpha_j + \mathbf{v}'_i \boldsymbol{\varphi}_j) + \mathbf{v}_i \quad (\text{Repetition})$$

$$\mathbf{c}_i = \mathbf{c}_i^* \times \mathbf{1}(\mathbf{h}_i^* > 0) = \mathbf{c}_i^* \times \mathbf{h}_i \quad (\text{Observation})$$

2. Construct **two** low dimensional **projections** into v_{iy} ($= w_{iy}$)

$$z_{iy}^0 \equiv \sum_j \hat{\delta}_{jy} h_{ij} \quad z_{iy}^+ \equiv \sum_j \hat{\varphi}_{jy} c_{ij} \quad (\text{SR projections})$$

Inclusion Repetition

3. Regress target variable on z_{iy}^+ , z_{iy}^0 and other covariates

$$v_{iy} = \beta_0 + [z_{iy}^0, z_{iy}^+, \mathbf{w}_{i,-y}, \mathbf{v}_{i,-y}]' \boldsymbol{\beta} + \varepsilon_i \quad (\text{forward regression})$$

- ▶ $d + p - 1$ dimensional regression reduced to $p + 2$ dimensional!

Hurdle distributed multiple regression (HDMR)

- ▶ Scale of text data requires convenient functional forms
- ▶ DMR uses independent Poissons to approximate the multinomial, one for each phrase
- ▶ We replace these Poissons with Hurdles (Mullahy, 1986)
- ▶ Hurdle model decomposes into two independent regressions
 1. Inclusion coefs. estimated from coverage indicators h_j and covariates w_i
 2. Repetition coefs. estimated from positive counts c_j and covariates v_i
- ▶ Can be distributed further!
- ▶ Lasso (L_1) regularization for both parts to avoid overfit

Selection bias

- ▶ Coefficients are biased if we use DMR on selected text data
- ▶ Severe bias if omitted variable in w is correlated with v
- ▶ For example, suppose:
 - ▶ FIFA World Cup crowds out financial news (limited attention)
 - ▶ ... and reduces market vol (traders watch it too)
 - ▶ Omitting it would yield biased effect of vol on financial news

Intermediary capital ratio (ICR)

- ▶ Intermediary asset pricing
 - ▶ Theory (Brunnermeier-Pedersen 2009 RFS, He-Krishnamurthy 2013 AER; Brunnermeier-Sannikov, 2014 AER)
 - ▶ Evidence (Adrian-Etula-Muir, 2014 JF; He-Kelly-Manela, 2017 JFE; Muir, 2017 QJE; Haddad-Muir, 2018)
- ▶ He-Kelly-Manela (2017 JFE):
 - ▶ Intermediary capital ratio (ICR) is the aggregate market capital ratio of NY Fed primary dealers
 - ▶ Innovations to the ICR price many asset classes
 - ▶ Suggestive results on predictive ability limited by short time-series starting 1970
- ▶ Can we backcast the ICR using historical newspaper text?
- ▶ Does high ICR predict low future market returns?

Data

Front-page titles and abstracts of the *Wall Street Journal*, 1926-2016

Date	Title	Abstract
20080916	AIG Faces Cash Crisis As Stock Dives 61%	American International Group Inc. was facing a severe cash ...
20080916	AIG, Lehman Shock Hits World Markets ...	The convulsions in the U.S. financial system sent markets ...
20080916	Business and Finance	Central banks around the world pumped cash into money ...
20080916	Keeping Their Powder Dry: Draft Boards ...	The Selective Service System has the awkward task of ...
20080916	Old-School Banks Emerge Atop New ...	Banks are heading "back to basics – to, if you like, the core ...
20080916	World-Wide	Thailand's ruling party chose ousted leader Thaksin's ...
		⋮

HDMR approach to news implied intermediary capital ratio

- ▶ We use HDMR to backcast missing values of ICR with WSJ text
 - + log price dividend ratio (pd_t)
 - + realized variance of financial stocks ($rvfin_t, rvfin_{t-1}$)
- ▶ Heckman selection models are non-parametrically identified
 - ▶ If a continuous variable enters the selection equation but can be excluded from second equation (Gallant-Nychka, 1984)
 - ▶ Proving such a result can be useful, but left for future work
- ▶ We seek an instrument for the inclusion decision
 - ▶ Prior attention to an issue may influence its coverage by the press (Boydston, 2013)
 - ▶ We use prior year realized variance of financial stocks ($rvfin_{t-13 \rightarrow t-1}$)
 - ▶ Assumption: excluded from repetition equation

Predicting ICR with realized variance, pd, and text

- ▶ Estimate **backward regressions**

$$h_{tj}^* = [icr_t, pd_t, rvfin_t, rvfin_{t-1}, rvfin_{t-13 \rightarrow t-1}]' \delta_j + u_{tj} \quad (\text{Inclusion})$$

$$c_{tj}^* = \lambda \left([icr_t, pd_t, rvfin_t, rvfin_{t-1}]' \varphi_j \right) + \varepsilon_{tj} > 0 \quad (\text{Repetition})$$

- ▶ Regress ICR on $z_{ty}^+ \equiv \sum_j \hat{\varphi}_{jy} c_{tj}$, $z_{ty}^0 \equiv \sum_j \hat{\delta}_{jy} h_{tj}$ and covariates

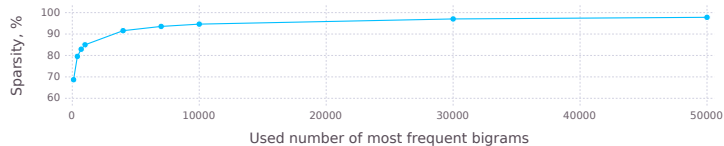
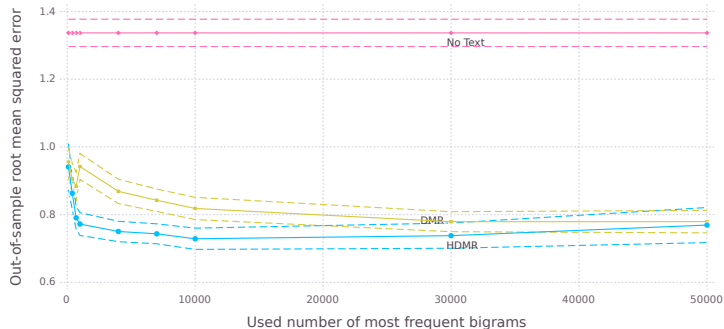
$$icr_t = \left[z_{ty}^+, z_{ty}^0, pd_t, rvfin_t, rvfin_{t-1}, rvfin_{t-13 \rightarrow t-1}, m_t \right]' \beta + v_t$$

(forward regression)

- ▶ Predict **out-of-sample**
 - ▶ Cross-validation with 10 random folds
 - ▶ Pseudo out-of-sample rolling regressions
 - ▶ Report root mean squared error (RMSE)

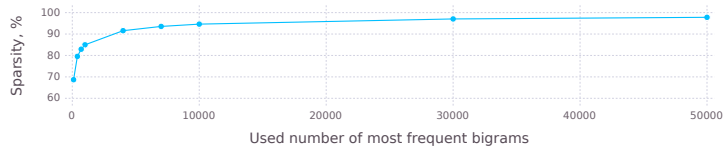
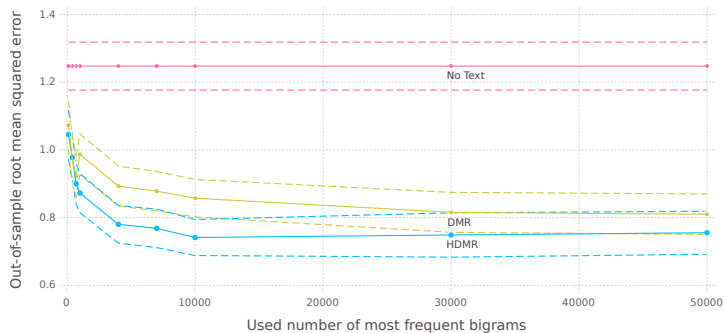
Out-of-sample prediction of ICR with text and covariates

HDMR's out-of-sample fit advantage changes with text sparsity (10-fold cross validation)



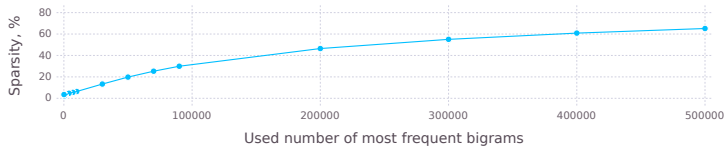
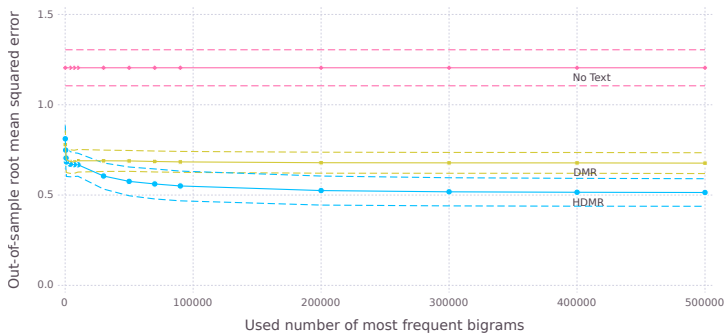
Out-of-sample prediction of ICR with text and covariates

HDMR's out-of-sample fit advantage changes with text sparsity (Pseudo out-of-sample)



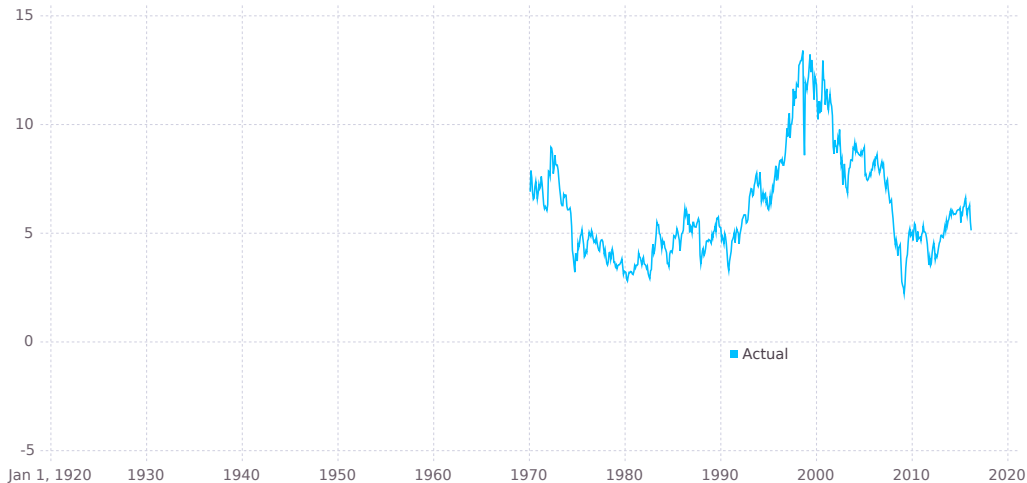
Denser text: HDMR's advantage increases with sparsity

Full WSJ monthly phrase counts, January 1990 to December 2010



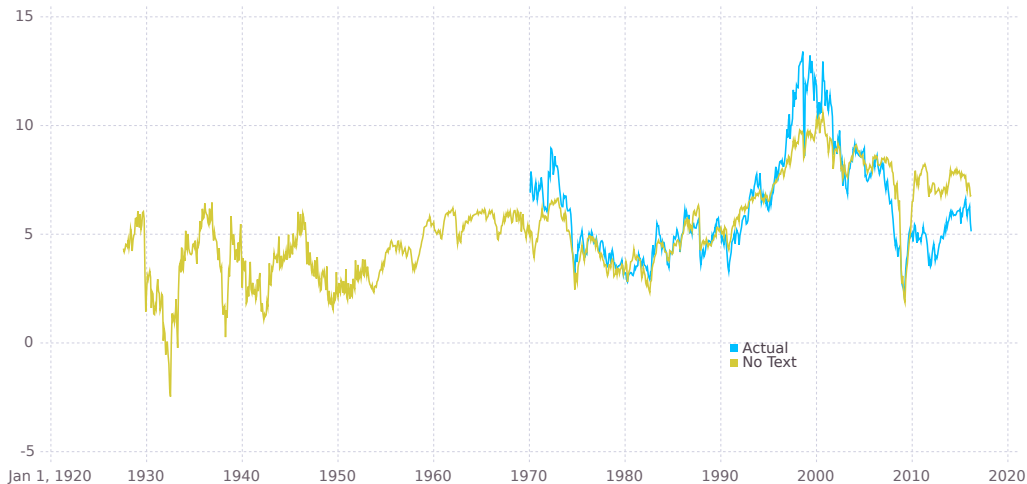
News implied intermediary capital ratio

ICR is available only since 1970 because dealers used to be private



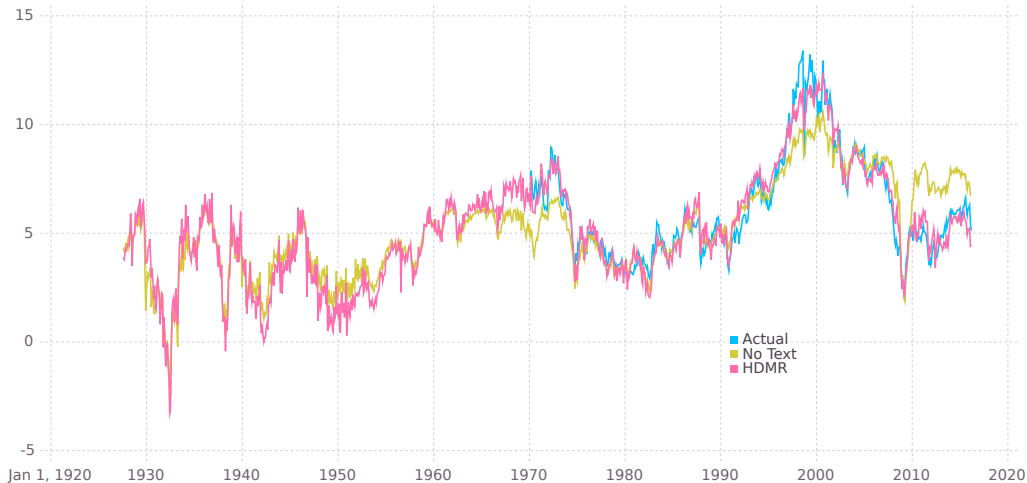
News implied intermediary capital ratio

First stab may be to fit using realized variance and price-dividend ratio without text



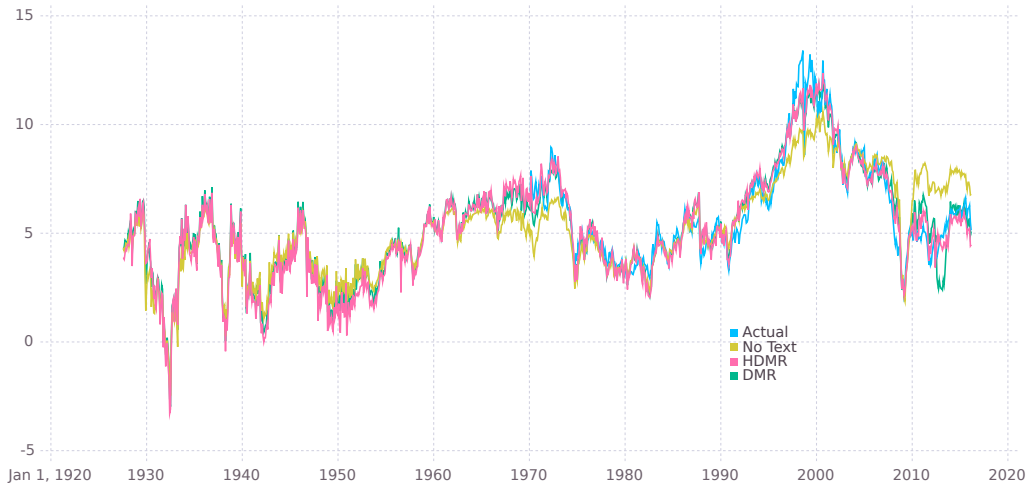
News implied intermediary capital ratio

HDMR gives a different predicted series exploiting text inclusion and repetition



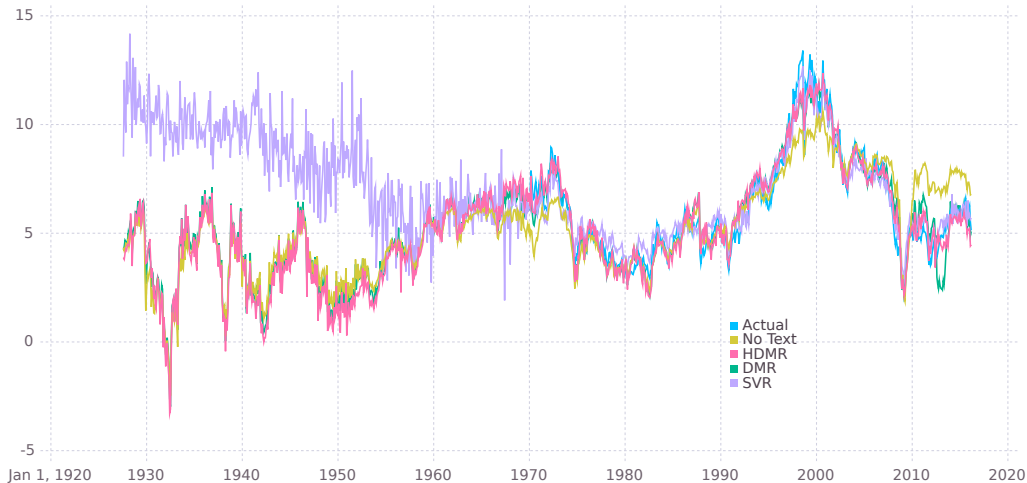
News implied intermediary capital ratio

DMR uses same information as HDMR but does not separate inclusion from repetition



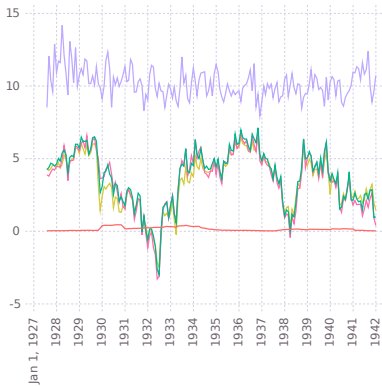
News implied intermediary capital ratio

Support Vector Regression of Manela-Moreira (2017) cannot concentrate on nontext covariates

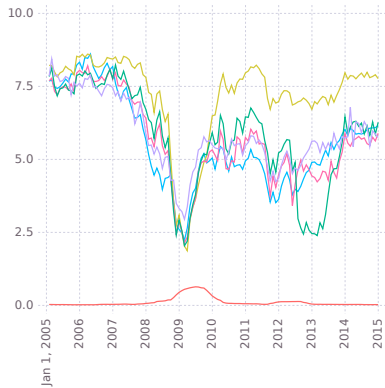


News implied intermediary capital ratio

Great Depression intermediaries were insolvent. Great Recession was almost as bad.



Great Depression



Great Recession

News-implied ICR predicts market returns

Consistent with He-Krishnamurthy (2013), 1σ higher ICR means 4.8pp lower risk premium

	$r_{t \rightarrow t+1}^{em}$			$r_{t \rightarrow t+3}^{em}$			$r_{t \rightarrow t+12}^{em}$		
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
icr_t	-1.29 (0.97)			-1.38 (1.07)			-1.22 (1.03)		
\widehat{icr}_t		-2.05*** (0.74)			-2.11*** (0.77)			-2.19*** (0.80)	
z_t^0			-43.93* (23.82)			-42.68** (21.32)			-43.66** (19.51)
z_t^+			41.20 (71.65)			42.85 (63.02)			38.58 (59.36)
pd_t			-13.92*** (4.95)			-14.68*** (5.14)			-15.19*** (5.19)
$rvfin_{t-1 \rightarrow t}$			-95.23*** (27.52)			-43.59* (23.45)			-8.82 (7.70)
$rvfin_{t-2 \rightarrow t-1}$			43.94 (29.65)			6.23 (21.34)			7.83 (7.22)
$rvfin_{t-13 \rightarrow t-1}$			44.63 (31.76)			42.31* (23.71)			11.47 (20.88)
N	552	841	841	552	841	841	544	833	833
Adjusted R^2	0.14	0.78	1.99	0.84	2.53	3.92	3.01	10.56	12.93

Hodrick (1992) standard errors are in parentheses

News-implied ICR predicts market returns

Similar magnitudes comparing postwar to full sample, but early sample has less text

	$r_{t \rightarrow t+1}^{em}$			$r_{t \rightarrow t+3}^{em}$			$r_{t \rightarrow t+12}^{em}$		
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
icr_t	-1.29 (0.97)			-1.38 (1.07)			-1.22 (1.03)		
\widehat{icr}_t		-1.92** (0.81)			-2.05* (1.23)			-2.14** (0.95)	
z_t^0			-31.56 (28.07)			-34.28 (25.32)			-37.49* (20.82)
z_t^+			-7.95 (84.44)			-2.64 (71.63)			-11.51 (68.51)
pd_t			-13.10** (5.37)			-13.87* (7.64)			-14.51** (5.91)
$rvfin_{t-1 \rightarrow t}$			-27.43 (17.91)			-7.76 (18.51)			-4.73 (8.28)
$rvfin_{t-2 \rightarrow t-1}$			29.31 (19.18)			0.57 (16.82)			-0.93 (8.42)
$rvfin_{t-13 \rightarrow t-1}$			-30.74 (27.38)			-17.60 (40.99)			-21.89 (37.44)
N	552	1,062	1,061	552	1,062	1,061	544	1,054	1,053
Adjusted R^2	0.14	0.43	0.42	0.84	1.50	1.48	3.01	6.64	8.43

Hodrick (1992) standard errors are in parentheses

Explaining the text with ICR-related covariates

WSJ front page monthly, January 1970 to February 2016

Variable	Sparsity	Top positive	Top negative
icr_t^0	0.577	busi bulletin, tax report, labor letter, washington wire, presid clinton	barack obama, presid barack, week output, obama administr, al qaeda
pd_t^0	0.587	barack obama, insid journal, peopl familiar, gold comex, comex troy	labor letter, busi bulletin, tax report, washington wire, ton week
$rvfin_{t-1 \rightarrow t}^0$	0.694	wall street, white hous, steel buy, presid clinton, chief execut	secretari christoph, tutsi rebel, wire wall, secretari perri, & backlog
$rvfin_{t-2 \rightarrow t-1}^0$	0.725	wall street, week mar, rwandan refuge, eastern ukrain, peopl familiar	wire wall, titl front, secretari perri, pictur corp, yitzhak rabin
$rvfin_{t-13 \rightarrow t-1}^0$	0.652	chief execut, presid barack, steel buy, yr trea, trea yld	secretari christoph, dole ks, yitzhak rabin, washington dc, wire clinton
icr_t^+	0.860	ounc dow, lopez obrador, presid clinton, clan leader, bosnia muslim	yr treasuri, treasuri yld, wsj research, c c, bond yr
pd_t^+	0.880	c c, yr treasuri, treasuri yld, avail headlin, bond yr	presid clinton, barrel dow, lopez obrador, wire clinton, bosnia muslim
$rvfin_{t-1 \rightarrow t}^+$	0.883	residenti construct, & unfil, temporari help, earn busi, israel plo	week aug, shimon pere, minist john, haiti militari, unit mine
$rvfin_{t-2 \rightarrow t-1}^+$	0.893	intern telephon, week aug, telegraph corp, clan leader, fix incom	temporari help, radovan karadz, & paperboard, paper &, taleban militia

- ▶ Predicted icr shaped by WSJ front page mentions of
 - ▶ Business and economic news
 - ▶ Government policy
 - ▶ Asset prices
 - ▶ Conditional on contemporaneous valuation ratios and financial volatility
- ▶ Some phrases capture fairly robust features of the data
- ▶ Others are unlikely to be useful for prediction before 2008

Explaining the text with ICR-related covariates

WSJ front page monthly, January 1970 to February 2016

Variable	Sparsity	Top positive	Top negative
icr_t^0	0.577	busi bulletin, tax report, labor letter, washington wire, presid clinton	barack obama, presid barack, week output, obama administr, al qaeda
pd_t^0	0.587	barack obama, insid journal, peopl familiar, gold comex, comex troy	labor letter, busi bulletin, tax report, washington wire, ton week
$rvfin_{t-1 \rightarrow t}^0$	0.694	wall street, white hous, steel buy, presid clinton, chief execut	secretari christoph, tutsi rebel, wire wall, secretari perri, & backlog
$rvfin_{t-2 \rightarrow t-1}^0$	0.725	wall street, week mar, rwandan refuge, eastern ukrain, peopl familiar	wire wall, titl front, secretari perri, pictur corp, yitzhak rabin
$rvfin_{t-13 \rightarrow t-1}^0$	0.652	chief execut, presid barack, steel buy, yr trea, trea yld	secretari christoph, dole ks, yitzhak rabin, washington dc, wire clinton
icr_t^+	0.860	ounc dow, lopez obrador, presid clinton, clan leader, bosnia muslim	yr treasuri, treasuri yld, wsj research, c c, bond yr
pd_t^+	0.880	c c, yr treasuri, treasuri yld, avail headlin, bond yr	presid clinton, barrel dow, lopez obrador, wire clinton, bosnia muslim
$rvfin_{t-1 \rightarrow t}^+$	0.883	residenti construct, & unfil, temporari help, earn busi, israel plo	week aug, shimon pere, minist john, haiti militari, unit mine
$rvfin_{t-2 \rightarrow t-1}^+$	0.893	intern telephon, week aug, telegraph corp, clan leader, fix incom	temporari help, radovan karadz, & paperboard, paper &, taleban militia

- ▶ Predicted icr shaped by WSJ front page mentions of
 - ▶ Business and economic news
 - ▶ Government policy
 - ▶ Asset prices
 - ▶ Conditional on contemporaneous valuation ratios and financial volatility
- ▶ Some phrases capture fairly robust features of the data
- ▶ Others are unlikely to be useful for prediction before 2008

Explaining the text with ICR-related covariates

WSJ front page monthly, January 1970 to February 2016

Variable	Sparsity	Top positive	Top negative
icr_t^0	0.577	busi bulletin, tax report, labor letter, washington wire, presid clinton	barack obama, presid barack, week output, obama administr, al qaeda
pd_t^0	0.587	barack obama, insid journal, peopl familiar, gold comex, comex troy	labor letter, busi bulletin, tax report, washington wire, ton week
$rvfin_{t-1 \rightarrow t}^0$	0.694	wall street, white hous, steel buy, presid clinton, chief execut	secretari christoph, tutsi rebel, wire wall, secretari perri, & backlog
$rvfin_{t-2 \rightarrow t-1}^0$	0.725	wall street, week mar, rwandan refuge, eastern ukrain, peopl familiar	wire wall, titl front, secretari perri, pictur corp, yitzhak rabin
$rvfin_{t-13 \rightarrow t-1}^0$	0.652	chief execut, presid barack, steel buy, yr trea, trea yld	secretari christoph, dole ks, yitzhak rabin, washington dc, wire clinton
icr_t^+	0.860	ounc dow, lopez obrador, presid clinton, clan leader, bosnia muslim	yr treasuri, treasuri yld, wsj research, c c, bond yr
pd_t^+	0.880	c c, yr treasuri, treasuri yld, avail headlin, bond yr	presid clinton, barrel dow, lopez obrador, wire clinton, bosnia muslim
$rvfin_{t-1 \rightarrow t}^+$	0.883	residenti construct, & unfil, temporari help, earn busi, israel plo	week aug, shimon pere, minist john, haiti militari, unit mine
$rvfin_{t-2 \rightarrow t-1}^+$	0.893	intern telephon, week aug, telegraph corp, clan leader, fix incom	temporari help, radovan karadz, & paperboard, paper &, taleban militia

- ▶ Predicted icr shaped by WSJ front page mentions of
 - ▶ Business and economic news
 - ▶ Government policy
 - ▶ Asset prices
 - ▶ Conditional on contemporaneous valuation ratios and financial volatility
- ▶ Some phrases capture fairly robust features of the data
- ▶ Others are unlikely to be useful for prediction before 2008

Explaining the text with ICR-related covariates

WSJ front page monthly, January 1970 to February 2016

Variable	Sparsity	Top positive	Top negative
icr_t^0	0.577	busi bulletin, tax report, labor letter, washington wire, presid clinton	barack obama, presid barack, week output, obama administr, al qaeda
pd_t^0	0.587	barack obama, insid journal, peopl familiar, gold comex, comex troy	labor letter, busi bulletin, tax report, washington wire, ton week
$rvfin_{t-1 \rightarrow t}^0$	0.694	wall street, white hous, steel buy, presid clinton, chief execut	secretari christoph, tutsi rebel, wire wall, secretari perri, & backlog
$rvfin_{t-2 \rightarrow t-1}^0$	0.725	wall street, week mar, rwandan refuge, eastern ukrain, peopl familiar	wire wall, titl front, secretari perri, pictur corp, yitzhak rabin
$rvfin_{t-13 \rightarrow t-1}^0$	0.652	chief execut, presid barack, steel buy, yr trea, trea yld	secretari christoph, dole ks, yitzhak rabin, washington dc, wire clinton
icr_t^+	0.860	ounc dow , lopez obrador, presid clinton, clan leader, bosnia muslim	yr treasuri , treasuri yld, wsj research, c c, bond yr
pd_t^+	0.880	c c, yr treasuri, treasuri yld, avail headlin, bond yr	presid clinton, barrel dow, lopez obrador, wire clinton, bosnia muslim
$rvfin_{t-1 \rightarrow t}^+$	0.883	residenti construct, & unfil, temporari help, earn busi, israel plo	week aug, shimon pere, minist john, haiti militari, unit mine
$rvfin_{t-2 \rightarrow t-1}^+$	0.893	intern telephon, week aug, telegraph corp, clan leader, fix incom	temporari help, radovan karadz, & paperboard, paper &, taleban militia

- ▶ Predicted icr shaped by WSJ front page mentions of
 - ▶ Business and economic news
 - ▶ Government policy
 - ▶ Asset prices
 - ▶ Conditional on contemporaneous valuation ratios and financial volatility
- ▶ Some phrases capture fairly robust features of the data
- ▶ Others are unlikely to be useful for prediction before 2008

Explaining the text with ICR-related covariates

WSJ front page monthly, January 1970 to February 2016

Variable	Sparsity	Top positive	Top negative
icr_t^0	0.577	busi bulletin, tax report, labor letter, washington wire, presid clinton	barack obama, presid barack, week output, obama administr, al qaeda
pd_t^0	0.587	barack obama, insid journal, peopl familiar, gold comex, comex troy	labor letter, busi bulletin, tax report, washington wire, ton week
$rvfin_{t-1 \rightarrow t}^0$	0.694	wall street, white hous, steel buy, presid clinton, chief execut	secretari christoph, tutsi rebel, wire wall, secretari perri, & backlog
$rvfin_{t-2 \rightarrow t-1}^0$	0.725	wall street, week mar, rwandan refuge, eastern ukrain, peopl familiar	wire wall, titl front, secretari perri, pictur corp, yitzhak rabin
$rvfin_{t-13 \rightarrow t-1}^0$	0.652	chief execut, presid barack, steel buy, yr trea, trea yld	secretari christoph, dole ks, yitzhak rabin, washington dc, wire clinton
icr_t^+	0.860	ounc dow, lopez obrador, presid clinton, clan leader, bosnia muslim	yr treasuri, treasuri yld, wsj research, c c, bond yr
pd_t^+	0.880	c c, yr treasuri, treasuri yld, avail headlin, bond yr	presid clinton, barrel dow, lopez obrador, wire clinton, bosnia muslim
$rvfin_{t-1 \rightarrow t}^+$	0.883	residenti construct, & unfil, temporari help, earn busi, israel plo	week aug, shimon pere, minist john, haiti militari, unit mine
$rvfin_{t-2 \rightarrow t-1}^+$	0.893	intern telephon, week aug, telegraph corp, clan leader, fix incom	temporari help, radovan karadz, & paperboard, paper &, taleban militia

- ▶ Predicted icr shaped by WSJ front page mentions of
 - ▶ Business and economic news
 - ▶ Government policy
 - ▶ Asset prices
 - ▶ Conditional on contemporaneous valuation ratios and financial volatility
- ▶ Some phrases capture fairly robust features of the data
- ▶ Others are unlikely to be useful for prediction before 2008

Explaining the text with ICR-related covariates

WSJ front page monthly, January 1970 to February 2016

Variable	Sparsity	Top positive	Top negative
icr_t^0	0.577	busi bulletin, tax report, labor letter, washington wire, presid clinton	barack obama, presid barack, week output , obama administr, al qaeda
pd_t^0	0.587	barack obama, insid journal, peopl familiar, gold comex, comex troy	labor letter, busi bulletin, tax report, washington wire, ton week
$rvfin_{t-1 \rightarrow t}^0$	0.694	wall street, white hous, steel buy, presid clinton, chief execut	secretari christoph, tutsi rebel, wire wall, secretari perri, & backlog
$rvfin_{t-2 \rightarrow t-1}^0$	0.725	wall street, week mar, rwandan refuge, eastern ukrain, peopl familiar	wire wall, titl front, secretari perri, pictur corp, yitzhak rabin
$rvfin_{t-13 \rightarrow t-1}^0$	0.652	chief execut, presid barack, steel buy, yr trea, trea yld	secretari christoph, dole ks, yitzhak rabin, washington dc, wire clinton
icr_t^+	0.860	ounc dow, lopez obrador, presid clinton, clan leader, bosnia muslim	yr treasuri, treasuri yld, wsj research, c c, bond yr
pd_t^+	0.880	c c, yr treasuri, treasuri yld, avail headlin, bond yr	presid clinton, barrel dow, lopez obrador, wire clinton, bosnia muslim
$rvfin_{t-1 \rightarrow t}^+$	0.883	residenti construct, & unfil, temporari help, earn busi, israel plo	week aug, shimon pere, minist john, haiti militari, unit mine
$rvfin_{t-2 \rightarrow t-1}^+$	0.893	intern telephon, week aug, telegraph corp, clan leader, fix incom	temporari help, radovan karadz, & paperboard, paper &, taleban militia

- ▶ Predicted icr shaped by WSJ front page mentions of
 - ▶ Business and economic news
 - ▶ Government policy
 - ▶ Asset prices
 - ▶ Conditional on contemporaneous valuation ratios and financial volatility
- ▶ Some phrases capture fairly robust features of the data
- ▶ Others are unlikely to be useful for prediction before 2008

Explaining the text with ICR-related covariates

WSJ front page monthly, January 1970 to February 2016

Variable	Sparsity	Top positive	Top negative
icr_t^0	0.577	busi bulletin, tax report, labor letter, washington wire, presid clinton	barack obama, presid barack, week output, obama administr, al qaeda
pd_t^0	0.587	barack obama, insid journal, peopl familiar, gold comex, comex troy	labor letter, busi bulletin, tax report, washington wire, ton week
$rvfin_{t-1 \rightarrow t}^0$	0.694	wall street, white hous, steel buy, presid clinton, chief execut	secretari christoph, tutsi rebel, wire wall, secretari perri, & backlog
$rvfin_{t-2 \rightarrow t-1}^0$	0.725	wall street, week mar, rwandan refuge, eastern ukrain, peopl familiar	wire wall, titl front, secretari perri, pictur corp, yitzhak rabin
$rvfin_{t-13 \rightarrow t-1}^0$	0.652	chief execut, presid barack, steel buy, yr trea, trea yld	secretari christoph, dole ks, yitzhak rabin, washington dc, wire clinton
icr_t^+	0.860	ounc dow, lopez obrador, presid clinton, clan leader, bosnia muslim	yr treasuri, treasuri yld, wsj research, c c, bond yr
pd_t^+	0.880	c c, yr treasuri, treasuri yld, avail headlin, bond yr	presid clinton, barrel dow, lopez obrador, wire clinton, bosnia muslim
$rvfin_{t-1 \rightarrow t}^+$	0.883	residenti construct, & unfil, temporari help, earn busi, israel plo	week aug, shimon pere, minist john, haiti militari, unit mine
$rvfin_{t-2 \rightarrow t-1}^+$	0.893	intern telephon, week aug, telegraph corp, clan leader, fix incom	temporari help, radovan karadz, & paperboard, paper &, taleban militia

- ▶ Predicted icr shaped by WSJ front page mentions of
 - ▶ Business and economic news
 - ▶ Government policy
 - ▶ Asset prices
 - ▶ Conditional on contemporaneous valuation ratios and financial volatility
- ▶ Some phrases capture fairly robust features of the data
- ▶ Others are unlikely to be useful for prediction before 2008

Focus on a single phrase for intuition

1σ increase in past year financial vol increases “financial crisis” inclusion odds by 30%

	HDMR		DMR
	Repetition	Inclusion	
Intercept	-9.94	-18.20	-16.66
icr_t	-0.35	-0.59	-0.61
pd_t	1.55	3.76	3.49
$rvfin_{t-1 \rightarrow t}$	1.23	0.85	1.44
$rvfin_{t-2 \rightarrow t-1}$	-0.56	1.07	-0.54
$rvfin_{t-13 \rightarrow t-1}$		2.80	1.26

Backward regressions

⇒

	HDMR	DMR
Repetition	-2.66	-4.73
Inclusion	-4.51	

Forward regressions

Focus on a single phrase for intuition

1σ increase in past year financial vol increases “financial crisis” inclusion odds by 30%

	HDMR		DMR
	Repetition	Inclusion	
Intercept	-9.94	-18.20	-16.66
icr_t	-0.35	-0.59	-0.61
pd_t	1.55	3.76	3.49
$rvfin_{t-1 \rightarrow t}$	1.23	0.85	1.44
$rvfin_{t-2 \rightarrow t-1}$	-0.56	1.07	-0.54
$rvfin_{t-13 \rightarrow t-1}$		2.80	1.26

Backward regressions

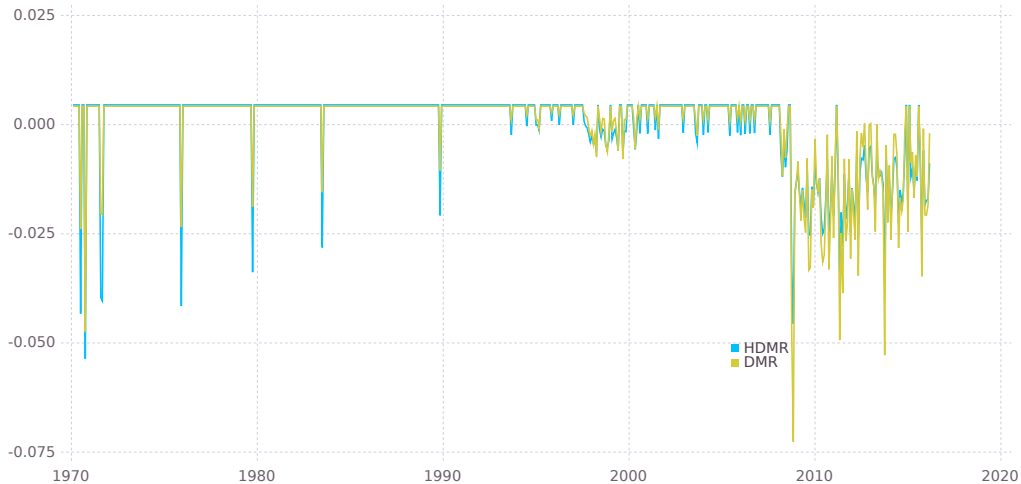
⇒

	HDMR	DMR
Repetition	-2.66	-4.73
Inclusion	-4.51	

Forward regressions

Focus on a single phrase for intuition

“financial crisis” on the front page is bad news for dealers, regardless of repetition



Does newspaper coverage forecast macroeconomic series?

- ▶ Stock-Watson (2012) show that macro forecasts of a simple dynamic factor model are hard to beat

$$Y_{t+h}^h = \beta_0 + [pc_t^1, \dots, pc_t^5]' \beta + \varepsilon_{t+h} \quad (\text{DFM-5})$$

- ▶ We use their data + WSJ text to forecast 1–12 months ahead

$$Y_{t+h}^h = \beta_0 + [z_{tY}^0, z_{tY}^+, pc_t^1, \dots, pc_t^5]' \beta + \varepsilon_{t+h} \quad (\text{HDMR})$$

Main findings

- ▶ Substantial OOS RMSE improvement using text with HDMR relative to DFM-5 for macroeconomic fundamentals
 - ▶ Nonfarm payroll employment forecast is 23–44% better
 - ▶ Housing starts forecast is 45–52% better
- ▶ WSJ text helps predict asset prices directly (stocks, treasuries, currencies) in quarterly/annual horizon but not monthly
- ▶ Advantage of HDMR increases with sparsity of the text
- ▶ Stronger results for nowcasting

Significant improvements in out-of-sample forecasting

HDMR RMSE relative to DFM-5: WSJ full text, 10,000 bigrams

Y_{t+h}^h	$h = 1$	$h = 3$	$h = 12$
IP: total	0.984 (0.372)	0.897** (0.042)	0.775*** (0.000)
Emp: total	0.881** (0.025)	0.828*** (0.001)	0.785*** (0.000)
U: all	0.956 (0.226)	0.825*** (0.003)	0.709*** (0.000)
HStarts: Total	0.716*** (0.000)	0.658*** (0.000)	0.618*** (0.000)
PMI	0.838*** (0.000)	0.852*** (0.000)	0.735*** (0.000)
CPI-ALL	1.110 (1.000)	1.064 (1.000)	1.030 (0.939)
Real AHE: goods	0.985 (0.328)	0.902*** (0.003)	0.629*** (0.000)
FedFunds	0.952* (0.068)	0.842*** (0.000)	0.677*** (0.000)
M1	1.100 (1.000)	1.102 (0.999)	0.990 (0.340)
Ex rate: avg	1.065 (0.999)	0.997 (0.462)	0.881*** (0.001)
S&P 500	1.041 (0.888)	0.941* (0.085)	0.793*** (0.000)
Consumer expect	1.110 (1.000)	1.043 (0.874)	0.980 (0.337)

Diebold-Mariano (1995) p -values are in parentheses

- ▶ Text is informative about future
 - ▶ Short and long run fundamentals
 - ▶ Long run fundamentals and prices
- ▶ Advantage of HDMR increases with text sparsity
- ▶ Text is also useful for nowcasting

Significant improvements in out-of-sample forecasting

HDMR RMSE relative to DFM-5: WSJ full text, 100,000 bigrams

Y_{t+h}^h	$h = 1$	$h = 3$	$h = 12$
IP: total	0.976 (0.313)	0.864*** (0.003)	0.619*** (0.000)
Emp: total	0.769*** (0.000)	0.690*** (0.000)	0.506*** (0.000)
U: all	0.986 (0.412)	0.819*** (0.001)	0.530*** (0.000)
HStarts: Total	0.546*** (0.000)	0.485*** (0.000)	0.519*** (0.000)
PMI	0.829*** (0.000)	0.842*** (0.000)	0.798*** (0.000)
CPI-ALL	1.012 (0.757)	1.051 (0.984)	1.081 (0.998)
Real AHE: goods	1.045 (0.971)	0.985 (0.319)	0.644*** (0.000)
FedFunds	0.996 (0.451)	0.925** (0.028)	0.701*** (0.000)
M1	1.040 (1.000)	1.029 (0.984)	1.058 (0.994)
Ex rate: avg	1.073 (0.999)	1.009 (0.641)	0.738*** (0.000)
S&P 500	1.054 (0.918)	0.881** (0.010)	0.663*** (0.000)
Consumer expect	1.061 (1.000)	1.097 (0.990)	0.852*** (0.002)

Diebold-Mariano (1995) p -values are in parentheses

- ▶ Text is informative about future
 - ▶ Short and long run fundamentals
 - ▶ Long run fundamentals and prices
- ▶ Advantage of HDMR increases with text sparsity
- ▶ Text is also useful for nowcasting

Significant improvements in out-of-sample nowcasting

HDMR RMSE relative to DFM-5: WSJ full text, 100,000 bigrams

Y_{t+h}^h	$h = 1$	$h = 3$	$h = 12$
IP: total	0.975 (0.167)	0.850*** (0.000)	0.626*** (0.000)
Emp: total	0.771*** (0.000)	0.694*** (0.000)	0.524*** (0.000)
U: all	0.982 (0.302)	0.819*** (0.000)	0.542*** (0.000)
HStarts: Total	0.527*** (0.000)	0.509*** (0.000)	0.475*** (0.000)
PMI	0.851*** (0.000)	0.866*** (0.000)	0.847*** (0.000)
CPI-ALL	0.962* (0.082)	1.022 (0.904)	1.098 (0.994)
Real AHE: goods	1.047 (0.910)	1.010 (0.610)	0.622*** (0.000)
FedFunds	0.974 (0.310)	0.888*** (0.009)	0.728*** (0.000)
M1	1.072 (1.000)	1.050 (0.998)	1.064 (0.991)
Ex rate: avg	1.080 (1.000)	1.025 (0.763)	0.743*** (0.000)
S&P 500	1.040 (0.890)	0.902** (0.016)	0.700*** (0.000)
Consumer expect	1.045 (0.998)	1.088 (0.991)	0.851*** (0.000)

Diebold-Mariano (1995) p -values are in parentheses

- ▶ Text is informative about future
 - ▶ Short and long run fundamentals
 - ▶ Long run fundamentals and prices
- ▶ Advantage of HDMR increases with text sparsity
- ▶ Text is also useful for nowcasting

Conclusion

- ▶ Incorporating structural economic restrictions into machine learning methods can improve out-of-sample prediction
- ▶ Hurdle Distributed Multiple Regression (HDMR)
 - ▶ Highly scalable approach to inference from big counts data
 - ▶ Includes an economically-motivated [selection equation](#)
 - ▶ Useful where extensive margin is interesting or more important than intensive margin
- ▶ Applications using [newspaper coverage for prediction](#)
 1. Backcast intermediary capital ratio
 2. Forecast macroeconomic series

Backcasting application

Summary statistics

Variable	Mean	Std	Min	p10	Median	p90	Max	Obs	Available
Phrase counts, c_{tj}	0.086	0.379	0.000	0.000	0.003	0.114	4.576	1075	192607–201602
Phrase indic. h_{tj}	0.054	0.212	0.000	0.000	0.002	0.089	1.000	1075	192607–201602
<i>icr</i>	6.236	2.399	2.230	3.616	5.574	9.578	13.400	557	197001–201605
<i>pd</i>	3.442	0.402	2.213	2.960	3.394	4.017	4.564	1075	192611–201605
<i>rvfin</i> $_{t-1 \rightarrow t}$	0.061	0.144	0.002	0.006	0.022	0.133	2.059	1079	192607–201605
<i>rvfin</i> $_{t-12 \rightarrow t}$	0.061	0.094	0.004	0.010	0.026	0.159	0.636	1068	192706–201605

Macro forecasting

Summary statistics

Variable	Mean	Std	Min	p10	Median	p90	Max	Obs	Available
Phrase counts, c_{tj}	2.971	2.732	0.178	0.590	2.342	6.190	16.909	252	199001–201012
Phrase indic. h_{tj}	0.680	0.409	0.032	0.161	0.860	0.998	1.000	252	199001–201012
IP: total	0.855	3.157	-16.004	-2.697	1.067	4.348	11.810	602	195901–200902
Emp: total	0.606	1.075	-4.177	-0.679	0.769	1.718	5.830	602	195901–200902
U: all	0.004	0.183	-0.700	-0.200	0.000	0.200	0.900	602	195901–200902
HStarts: Total	7.307	0.236	6.192	6.991	7.327	7.602	7.822	602	195901–200902
PMI	52.867	6.927	29.400	44.100	53.500	60.780	72.100	603	195901–200903
CPI-ALL	0.000	1.075	-5.282	-1.162	-0.003	1.091	7.018	602	195902–200903
Real AHE: goods	0.274	1.228	-4.401	-1.113	0.255	1.560	5.883	602	195901–200902
FedFunds	0.002	0.371	-1.560	-0.420	0.010	0.380	1.600	602	195901–200902
M1	0.012	2.184	-10.505	-2.344	-0.000	2.407	7.479	601	195902–200902
Ex rate: avg	-0.170	5.953	-21.103	-8.241	0.038	6.892	21.174	601	195901–200901
S&P 500	1.796	14.410	-91.153	-14.428	2.798	16.943	45.355	603	195901–200903
Consumer expect	-0.053	3.975	-16.500	-4.600	-0.200	4.600	22.500	603	195901–200903

Macro forecasting

WSJ full text, 10,000 bigrams

Months forward:	$h = 1$				$h = 3$				$h = 12$			
Folds:	Random		Rolling		Random		Rolling		Random		Rolling	
Y_{t+h}^h	HDMR	DMR	HDMR	DMR	HDMR	DMR	HDMR	DMR	HDMR	DMR	HDMR	DMR
IP: total	0.984 (0.372)	0.960 (0.118)	0.976 (0.291)	0.991 (0.409)	0.897** (0.042)	0.892** (0.012)	0.899*** (0.009)	0.865*** (0.000)	0.775*** (0.000)	0.793*** (0.000)	0.799*** (0.000)	0.785*** (0.000)
Emp: total	0.881** (0.025)	0.858*** (0.009)	0.926* (0.078)	0.941 (0.153)	0.828*** (0.001)	0.831*** (0.000)	0.874*** (0.007)	0.889** (0.018)	0.785*** (0.000)	0.778*** (0.000)	0.771*** (0.000)	0.826*** (0.004)
U: all	0.956 (0.226)	0.948 (0.189)	0.975 (0.296)	1.007 (0.559)	0.825*** (0.003)	0.843*** (0.010)	0.753*** (0.000)	0.802*** (0.000)	0.709*** (0.000)	0.755*** (0.000)	0.691*** (0.000)	0.753*** (0.000)
HStarts: Total	0.716*** (0.000)	0.653*** (0.000)	0.783*** (0.000)	0.766*** (0.000)	0.658*** (0.000)	0.636*** (0.000)	0.736*** (0.000)	0.773*** (0.000)	0.618*** (0.000)	0.571*** (0.000)	0.815*** (0.000)	0.760*** (0.000)
PMI	0.838*** (0.000)	0.815*** (0.000)	0.967 (0.271)	0.980 (0.377)	0.852*** (0.000)	0.839*** (0.000)	0.880** (0.013)	0.922* (0.084)	0.735*** (0.000)	0.751*** (0.000)	0.888* (0.051)	0.922 (0.132)
CPI-ALL	1.110 (1.000)	1.096 (1.000)	1.116 (1.000)	1.188 (1.000)	1.064 (1.000)	1.091 (1.000)	1.054 (0.992)	1.056 (0.927)	1.030 (0.939)	1.020 (0.819)	1.041 (0.933)	1.003 (0.536)
Real AHE: goods	0.985 (0.328)	0.996 (0.441)	1.030 (0.789)	1.024 (0.793)	0.902*** (0.003)	0.919*** (0.006)	0.989 (0.388)	1.028 (0.744)	0.629*** (0.000)	0.627*** (0.000)	0.884** (0.025)	0.902* (0.056)
FedFunds	0.952* (0.068)	0.918*** (0.005)	0.994 (0.446)	1.056 (0.836)	0.842*** (0.000)	0.835*** (0.001)	0.960 (0.215)	1.011 (0.573)	0.677*** (0.000)	0.665*** (0.000)	0.840*** (0.006)	0.926 (0.146)
M1	1.100 (1.000)	1.082 (1.000)	1.170 (1.000)	1.172 (1.000)	1.102 (0.999)	1.094 (1.000)	1.199 (1.000)	1.134 (1.000)	0.990 (0.340)	1.001 (0.510)	1.057 (0.906)	1.113 (0.986)
Ex rate: avg	1.065 (0.999)	1.059 (1.000)	1.097 (0.966)	1.175 (0.999)	0.997 (0.462)	0.971 (0.123)	1.034 (0.709)	1.080 (0.908)	0.881*** (0.001)	0.835*** (0.000)	0.867** (0.015)	0.863** (0.013)
S&P 500	1.041 (0.888)	1.021 (0.762)	1.089 (0.963)	1.003 (0.541)	0.941* (0.085)	0.940* (0.062)	0.939 (0.123)	0.895*** (0.006)	0.793*** (0.000)	0.785*** (0.000)	0.739*** (0.000)	0.736*** (0.000)
Consumer expect	1.110 (1.000)	1.159 (1.000)	1.180 (1.000)	1.672 (1.000)	1.043 (0.874)	1.022 (0.780)	1.102 (0.990)	1.143 (0.993)	0.980 (0.337)	0.916** (0.036)	0.894** (0.031)	0.850*** (0.003)

Macro forecasting

WSJ full text, 100,000 bigrams

Months forward:	$h = 1$				$h = 3$				$h = 12$			
Folds:	Random		Rolling		Random		Rolling		Random		Rolling	
Y_{t+h}^h	HDMR	DMR	HDMR	DMR	HDMR	DMR	HDMR	DMR	HDMR	DMR	HDMR	DMR
IP: total	0.976 (0.313)	0.961 (0.176)	1.025 (0.718)	1.000 (0.496)	0.864*** (0.003)	0.884** (0.016)	0.900** (0.014)	0.857*** (0.000)	0.619*** (0.000)	0.775*** (0.000)	0.680*** (0.000)	0.769*** (0.000)
Emp: total	0.769*** (0.000)	0.832*** (0.005)	0.933 (0.121)	0.920* (0.073)	0.690*** (0.000)	0.803*** (0.000)	0.811*** (0.000)	0.873*** (0.005)	0.506*** (0.000)	0.737*** (0.000)	0.607*** (0.000)	0.782*** (0.000)
U: all	0.986 (0.412)	0.945 (0.184)	1.023 (0.677)	0.992 (0.441)	0.819*** (0.001)	0.824*** (0.008)	0.921** (0.030)	0.785*** (0.000)	0.530*** (0.000)	0.705*** (0.000)	0.672*** (0.000)	0.726*** (0.000)
HStarts: Total	0.546*** (0.000)	0.647*** (0.000)	0.537*** (0.000)	0.764*** (0.000)	0.485*** (0.000)	0.626*** (0.000)	0.499*** (0.000)	0.762*** (0.000)	0.519*** (0.000)	0.561*** (0.000)	0.529*** (0.000)	0.758*** (0.000)
PMI	0.829*** (0.000)	0.774*** (0.000)	0.990 (0.433)	0.994 (0.466)	0.842*** (0.000)	0.787*** (0.000)	0.982 (0.379)	0.935 (0.153)	0.798*** (0.000)	0.724*** (0.000)	1.001 (0.509)	0.929 (0.169)
CPI-ALL	1.012 (0.757)	1.152 (1.000)	1.065 (1.000)	1.246 (1.000)	1.051 (0.984)	1.130 (1.000)	1.032 (0.936)	1.107 (0.995)	1.081 (0.998)	1.039 (0.926)	1.067 (0.966)	1.034 (0.847)
Real AHE: goods	1.045 (0.971)	1.012 (0.653)	1.158 (1.000)	1.036 (0.845)	0.985 (0.319)	0.905*** (0.001)	1.017 (0.685)	1.000 (0.499)	0.644*** (0.000)	0.595*** (0.000)	0.849*** (0.006)	0.888** (0.040)
FedFunds	0.996 (0.451)	0.906*** (0.002)	1.214 (1.000)	1.036 (0.745)	0.925** (0.028)	0.805*** (0.000)	1.012 (0.605)	0.975 (0.325)	0.701*** (0.000)	0.602*** (0.000)	0.694*** (0.000)	0.846*** (0.008)
M1	1.040 (1.000)	1.169 (1.000)	1.074 (1.000)	1.411 (1.000)	1.029 (0.984)	1.158 (1.000)	1.134 (1.000)	1.327 (1.000)	1.058 (0.994)	1.005 (0.579)	1.118 (0.998)	1.135 (0.993)
Ex rate: avg	1.073 (0.999)	1.080 (1.000)	1.064 (0.922)	1.265 (1.000)	1.009 (0.641)	0.968 (0.122)	1.112 (0.968)	1.130 (0.973)	1.078*** (0.000)	0.805*** (0.000)	0.884** (0.031)	0.851*** (0.008)
S&P 500	1.054 (0.918)	1.031 (0.801)	1.100 (0.980)	1.013 (0.629)	0.881** (0.010)	0.932* (0.056)	0.916** (0.043)	0.876*** (0.003)	0.663*** (0.000)	0.767*** (0.000)	0.632*** (0.000)	0.716*** (0.000)
Consumer expect	1.061 (1.000)	1.206 (1.000)	1.065 (0.997)	1.675 (1.000)	1.097 (0.990)	1.036 (0.862)	1.121 (0.999)	1.166 (0.997)	0.852*** (0.002)	0.909** (0.035)	0.968 (0.305)	0.853*** (0.007)

Macro nowcasting

WSJ full text, 100,000 bigrams

Months forward:	$h = 1$				$h = 3$				$h = 12$			
Folds:	Random		Rolling		Random		Rolling		Random		Rolling	
Y_{t+h}^h	HDMR	DMR	HDMR	DMR	HDMR	DMR	HDMR	DMR	HDMR	DMR	HDMR	DMR
IP: total	0.975 (0.167)	0.942** (0.020)	1.039 (0.824)	0.971 (0.249)	0.850*** (0.000)	0.872*** (0.001)	0.893*** (0.010)	0.855*** (0.000)	0.626*** (0.000)	0.771*** (0.000)	0.701*** (0.000)	0.765*** (0.000)
Emp: total	0.771*** (0.000)	0.820*** (0.000)	0.938 (0.134)	0.873*** (0.009)	0.694*** (0.000)	0.779*** (0.000)	0.830*** (0.000)	0.844*** (0.001)	0.524*** (0.000)	0.726*** (0.000)	0.633*** (0.000)	0.783*** (0.000)
U: all	0.982 (0.302)	0.943 (0.106)	1.035 (0.793)	0.946 (0.122)	0.819*** (0.000)	0.818*** (0.001)	0.949 (0.112)	0.778*** (0.000)	0.542*** (0.000)	0.684*** (0.000)	0.705*** (0.000)	0.714*** (0.000)
HStarts: Total	0.527*** (0.000)	0.641*** (0.000)	0.521*** (0.000)	0.750*** (0.000)	0.509*** (0.000)	0.619*** (0.000)	0.509*** (0.000)	0.752*** (0.000)	0.475*** (0.000)	0.560*** (0.000)	0.507*** (0.000)	0.755*** (0.000)
PMI	0.851*** (0.000)	0.801*** (0.000)	0.981 (0.371)	0.980 (0.380)	0.866*** (0.000)	0.792*** (0.000)	0.976 (0.344)	0.931 (0.140)	0.847*** (0.000)	0.725*** (0.000)	1.023 (0.653)	0.942 (0.202)
CPI-ALL	0.962* (0.082)	1.151 (1.000)	1.044 (0.984)	1.440 (1.000)	1.022 (0.904)	1.135 (1.000)	1.042 (0.965)	1.133 (0.999)	1.098 (0.994)	1.054 (0.975)	1.079 (0.984)	1.030 (0.817)
Real AHE: goods	1.047 (0.910)	1.023 (0.695)	1.210 (1.000)	1.047 (0.869)	1.010 (0.610)	0.931** (0.022)	1.052 (0.927)	0.965 (0.197)	0.622*** (0.000)	0.583*** (0.000)	0.823*** (0.001)	0.889** (0.038)
FedFunds	0.974 (0.310)	0.889*** (0.005)	1.180 (1.000)	1.004 (0.527)	0.888*** (0.009)	0.785*** (0.000)	1.026 (0.727)	0.954 (0.209)	0.728*** (0.000)	0.609*** (0.000)	0.764*** (0.000)	0.879** (0.030)
M1	1.072 (1.000)	1.271 (1.000)	1.091 (1.000)	1.558 (1.000)	1.050 (0.998)	1.138 (1.000)	1.116 (1.000)	2.005 (1.000)	1.064 (0.991)	0.996 (0.459)	1.111 (0.997)	1.110 (0.986)
Ex rate: avg	1.080 (1.000)	1.109 (0.998)	1.042 (0.840)	1.192 (0.999)	1.025 (0.763)	0.973 (0.181)	1.099 (0.948)	1.132 (0.976)	0.743*** (0.000)	0.804*** (0.000)	0.911* (0.079)	0.865** (0.012)
S&P 500	1.040 (0.890)	1.020 (0.725)	1.073 (0.935)	1.024 (0.715)	0.902** (0.016)	0.929** (0.035)	0.930* (0.085)	0.896** (0.011)	0.700*** (0.000)	0.766*** (0.000)	0.671*** (0.000)	0.731*** (0.000)
Consumer expect	1.045 (0.998)	1.234 (1.000)	1.030 (0.929)	1.453 (1.000)	1.088 (0.991)	1.033 (0.866)	1.097 (0.997)	1.220 (1.000)	0.851*** (0.000)	0.909** (0.022)	0.991 (0.445)	0.879** (0.021)