

Text Selection

Bryan Kelly

Asaf Manela

Alan Moreira*

November 2019

Abstract

Text data is ultra-high dimensional, which makes machine learning techniques indispensable for textual analysis. Text is often selected—journalists, speechwriters, and others craft messages to target their audiences’ limited attention. We develop an economically motivated high dimensional selection model that improves learning from text (and from sparse counts data more generally). Our model is especially useful when the choice to include a phrase is more interesting than the choice of how frequently to repeat it. It allows for parallel estimation, making it computationally scalable. A first application revisits the partisanship of US congressional speech. We find that earlier spikes in partisanship manifested in increased repetition of different phrases, whereas the upward trend starting in the 1990s is due to entirely distinct phrase selection. Additional applications show how our model can backcast, nowcast, and forecast macroeconomic indicators using newspaper text, and that it substantially improves out-of-sample fit relative to alternative approaches.

Keywords: Text analysis, machine learning, selection model, high dimension forecast, multinomial regression, hurdle, zero inflation, partisanship, intermediary capital

*Yale University, bryan.kelly@yale.edu; Washington University in St. Louis, amanela@wustl.edu; and University of Rochester, alan.moreira@simon.rochester.edu. We are grateful for helpful comments by Xavier Gabaix, Matthew Gentzkow, Hongyi Liu, Andy Neuhierl (discussant), Jesse Shapiro, Matt Taddy, Paul Tetlock (discussant) and by seminar participants at École Polytechnique Fédérale de Lausanne, Hebrew University, IDC Herzliya, Indiana University, INSEAD, Kansas City Fed, Ohio State, Syracuse University, Tel-aviv University, Rice University, University of Michigan, NBER SI AP, WFA, and the University of Chicago CITE Conference. Computations were performed using the facilities of the Washington University Center for High Performance Computing, which were partially provided through NIH grant S10 OD018091.

1 Introduction

Digital text is increasingly available to social scientists in the form of newspapers, blogs, tweets, regulatory filings, congressional records and more. Two attributes of text differentiate it from other types of data typically used by economists. First, text data is inherently ultra-high dimensional—unique phrases in a corpus (roughly equivalent to the set of variables) often number in the millions—because many words are used to describe similar phenomena. Second, phrase counts are sparse—most phrases have a count of zero in most documents. Due to these attributes, statistical learning from text requires techniques commonly used in machine learning (Gentzkow, Kelly, and Taddy, 2019). This paper proposes a new methodology that fully embraces the sparseness of text data. We show that modeling the extensive margin decision to use a particular word reveals information that leads to large improvements in out-of-sample prediction.

A standard econometric approach for describing counts data is multinomial logistic regression. Unfortunately, this model is computationally intractable in most text-related applications because the number of categories is extremely large. Taddy (2015) shows that one can overcome this dimensionality problem by approximating the multinomial with cleverly shifted independent poisson regressions, one for each word. This *distributed multinomial regression* (DMR) has the great advantage that the poisson regression can be distributed across parallel computing units.

The poisson model has the disadvantage that it provides a poor description of word counts. In particular, there tends to be a much higher proportion of phrases having zero counts than the poisson distribution allows. If one restricts attention to positive counts, however, a (truncated) poisson is a good approximation for the data. This sparsity—the additional probability mass on zero counts—is a feature of many text samples (and is why most text analysis software packages use sparse matrices to store word counts efficiently).¹

We propose a new methodology for text regression. Our *hurdle distributed multinomial regression* (HDMR) model is tailored to the twin challenges of high dimensionality and sparsity of text data. To accommodate the dimensionality of text, we build on Taddy’s DMR insight of independent phrase-level models. However, we replace each phrase-level poisson regression with a hurdle model

¹Zipf (1935) observed that language follows a power law. Rennie, Shih, Teevan, and Karger (2003) suggest a multinomial naive Bayes modification to emulate such power law distributions. Wilbur and Kim (2009) find that in some cases, word “burstiness” is so strong that additional occurrences of a word essentially add no useful information to a classifier.

that has two parts. The first part is a selection equation, which models the text producer's choice of whether or not to include a particular phrase. The second part is a positive counts model, which describes the choice of how many times a word is repeated (conditional on being used at all).²

Our two-part HDMR model thus generalizes DMR by decomposing language choice into an extensive margin (the selection equation) and intensive margin (the positive counts model). Explicitly modeling the extensive margin of phrase choice has two advantages. The statistical advantage is that it introduces a modeling component dedicated to capturing the excess probability mass on zero counts. We use the hurdle approach of [Mullahy \(1986\)](#) to model selection because it allows the extensive and intensive components to be estimated independently, and therefore, can also be distributed at essentially no additional computation cost relative to DMR.³

The economic advantage of our model is that it adapts the selection methodology of [Heckman \(1979\)](#) to a high dimensional setting. HDMR provides a means for estimating models in which sparsity is of first-order economic importance. In text data, HDMR is particularly useful when an author's choice to cover or not cover a topic is more economically interesting than the choice of coverage intensity. For example, newspaper publishers' extensive margin of coverage is informative about their news production technology and the constrained attention of its audience. The selection decision is a key lever that writers use to signal their ideological type to readers (e.g. [Mullainathan and Shleifer, 2005](#); [Gentzkow and Shapiro, 2006](#)). Politicians carefully select phrases that resonate with voters in congressional speech ([Gentzkow, Shapiro, and Taddy, 2019](#)), and fixed costs of using censored or socially taboo words may generate further sparsity ([Michel, Shen, Aiden, Veres, Gray, Pickett, Hoiberg, Clancy, Norvig, Orwant, et al., 2011](#)). Sparsity reflects the suboptimality of writing about a particular topic just a little bit: In order for the signal to be comprehensible there must be a minimum amount of exposition, yet the total amount of text cannot exceed the audience's attention span.⁴ [Gabaix \(2014\)](#) shows that, when agents are boundedly rational, sparsity emerges in equilibrium for a wide variety of economic settings under otherwise general conditions.

In a first application we revisit the partisanship of US congressional speech. [Gentzkow, Shapiro, and Taddy \(2019\)](#) have recently suggested a new measure of partisanship, which relies on DMR to

²We use L1 regularization (lasso) to further manage model dimensionality.

³The hurdle model is widely used, for example, in health economics. [Greene \(2007\)](#) surveys models for counts.

⁴Academic researchers know this well, and therefore tend to use consistent wording to clarify their argument, rather than alternate between synonyms to expound the same contention.

alleviate a finite-sample bias that arises in high-dimensional choice settings. They define partisanship as the ease with which an observer could infer a congressperson's party from a single phrase. They document a fairly low and constant level of partisanship starting in 1873 until the early 1990s when it starts increasing sharply. Our two-part method allows us to refine their findings by decomposing partisanship into two parts. One part measures the ease with which an observer could infer a congressperson's party from a single phrase inclusion, i.e., the fact that the congressperson uses a particular phrase in their vocabulary. The second part focuses on phrase repetition, conditional on inclusion, i.e., how much emphasis they place on a particular issue. We find that repetition partisanship spikes in the late 1890s and late 1920s to levels similar to those seen in the 2000s. But inclusion partisanship remains low for most of the sample, until it starts rising gradually in the 1970s, before spiking up in the 1990s. These findings mean that politicians belonging to different parties have previously emphasized different phrases or issues to a varying degree, while using the same language. What is new about the recent upward trend starting in the 1990s is that Democrats and Republicans choose to communicate with entirely distinct language.

Our model also serves as a basis for exploring the relationships between phrases and numerical covariates such as macroeconomic state variables or financial markets prices, which enter as conditioning variables for the distribution of phrase counts. First, the model doubles as a dimension reduction method for text. Like DMR, HDMR generates *sufficient reductions* of text by projecting phrase counts onto covariates. Sufficient reductions serve as indices that best summarize the text as it relates to each covariate after controlling for other covariates. HDMR produces two such sufficient reductions per covariate: one for word inclusion (the extensive margin) and the other for repetition (the intensive margin). Second, the model can be flipped in order to predict covariates based on text via an inverse regression. This is useful, for example, to backcast key macro variables that have limited histories or missing data, or for “nowcasting” when text is available in a more timely manner than the forecast target.

Our second application illustrates this feature by using HDMR to backcast a measure of equity capitalization in the financial sector using the text of *The Wall Street Journal*. The intermediary capital ratio (*icr*) is the central state variable in the growing literature on intermediary-based asset pricing, and helps explain the behavior of risk premia in a wide array of asset classes (He, Kelly, and Manela, 2017). However, it is only available beginning in 1970. From our long sample of

Wall Street Journal text, we estimate an *icr* series starting in 1927 to investigate the historical interaction between intermediary capital and asset prices.

We find that HDMR gives substantially improved out-of-sample predictions of *icr* compared to DMR, which indicates that modeling selection helps with forecasting in this context. HDMR also outperforms support vector regression (Vapnik, 2000), which Manela and Moreira (2017) use for text-based backcasting of the VIX stock market volatility index. Unlike support vector regression, both DMR and HDMR can concentrate on individual variables that behave differently from word counts (i.e. non-text control variables), but are useful for prediction. We find that the out-of-sample advantage of HDMR over DMR is humped-shaped, increasing with the sparsity of the text up to a point where repetition is rarely observed. On one extreme, if we keep only highly frequent words—a common approach to ad hoc prefiltering of words—the document term matrix becomes dense and the text is similarly described by a selection-free DMR model. At the same time, however, more stringent phrase filters lead to a large deterioration in prediction accuracy from either method. Less filtering (allowing for more phrases) makes for a better prediction model, and also magnifies the benefits of accounting for text selection with a hurdle. However, if so many infrequent phrases are kept, that repetition is rarely observed, the advantage of HDMR over DMR diminishes. At this other extreme, inclusion choice dominates the analysis, and a one-part model like a binomial or a poisson would suffice.

The news-implied intermediary capital ratio series provides for more powerful tests that support central predictions of intermediary asset pricing theory. The results show that times when intermediaries are highly capitalized are “good times” when these marginal investors demand a relatively low premium to hold risky assets. Our findings imply that news text of *The Wall Street Journal* provides a strong signal about stock market risk premia, over and above common financial market covariates such as the price-dividend ratio and stock volatility.

In our third application, we ask whether HDMR can extract information from text of *The Wall Street Journal* that is useful for forecasting US macroeconomic indicators (beyond that of a benchmark principal components method suggested by Stock and Watson, 2012). We find that the information in newspaper text significantly improves over benchmark forecasts for key macroeconomic indicators such as nonfarm payroll employment and housing starts, and that the advantages of text-based information increase as we expand the dimensionality (and along with it, the sparsity)

of the text. In a related analysis, we show that this text is valuable for nowcasting macroeconomic series, which are released at a lower frequency and with a delay relative to news text.

This paper contributes to a rapidly growing literature that incorporates insights from machine learning into econometrics. Recent work applies prediction algorithms in policy analysis (Kleinberg, Ludwig, Mullainathan, and Obermeyer, 2015; Bajari, Nekipelov, Ryan, and Yang, 2015; Athey, 2017; Ludwig, Mullainathan, and Spiess, 2017; Einav, Finkelstein, Mullainathan, and Obermeyer, 2018; Kleinberg, Lakkaraju, Leskovec, Ludwig, and Mullainathan, 2018) and considers parameter uncertainty and causal inference in high dimensional settings (Belloni, Chen, Chernozhukov, and Hansen, 2012; Belloni, Chernozhukov, and Hansen, 2014; Belloni, Chernozhukov, Fernández-Val, and Hansen, 2017; Mullainathan and Spiess, 2017; Athey, Imbens, Pham, and Wager, 2017; Athey, Tibshirani, and Wager, 2019). Our model allows economists interested in analyzing counts data like text, to model selection in the process that generates their data in a flexible, robust and scalable way. We show that adding structure that is well-grounded in economic theory can substantially improve prediction in high dimensional settings.⁵

HDMR allows one to analyze text within a regression framework familiar to economists, and strikes a balance between prediction and interpretation. Dictionary-based approaches are easy to interpret, but are no match to regularized regression in prediction tasks (e.g. Manela and Moreira, 2017; Chebonenko, Gu, and Muravyev, 2018). At the other extreme, neural network models are great nonlinear predictors, but also harder to interpret (Gu, Kelly, and Xiu, 2018, 2019).

Our technology is publicly available via the `HurdleDMR` package for `Julia`, which can be called from many other programming languages like `Python` and `R`. The package allows for computationally efficient distributed estimation of the multiple hurdles over parallel processes, generating sufficient reduction projections, and inverse regressions with selected text. It allows for elastic net type convex combinations of L1 (Lasso) and L2 (Ridge) regularization as in `glmnet` (Friedman, Hastie, and Tibshirani, 2010), and for concave regularization paths as in `gamlr` (Taddy, 2017).

We start Section 2 by describing the DMR model of Taddy (2015), and then extending it to introduce our main contribution, a text selection model, which we refer to as HDMR. We then

⁵See Athey (2018) for a recent survey and Hoberg and Phillips (2016); Hanley and Hoberg (2019); Jiang, Lee, Martin, and Zhou (2019) for recent text analysis applications. We also provide new tools for summarizing high dimensional coverage to a literature studying the media in economics (Gentzkow and Shapiro, 2006; Qin, Strömberg, and Wu, 2018; Durante and Zhuravskaya, 2018) and finance (Antweiler and Frank, 2004; Tetlock, 2007; Fang and Peress, 2009; Engelberg and Parsons, 2010; García, 2013).

illustrate its usefulness with three applications. Section 3 analyses the partisanship of congressional speech. Section 4 backcasts the intermediary capital ratio using newspaper text and analyzes its historical asset pricing properties. Section 5 uses newspaper text to forecast and nowcast key macroeconomic indicators. Section 6 concludes. Appendix A provides further robustness tests and Appendices B and C provides further theoretical detail.

2 A model for text selection

Let \mathbf{c}_i be a vector of counts in d categories for observation i , summing to $m_i = \sum_j c_{ij}$, and let \mathbf{v}_i be a p_v -vector of covariates associated with each observation $i = 1 \dots n$. In a text application, c_{ij} are counts of word or phrase (n-gram) j in document i with attributes \mathbf{v}_i .⁶ An econometrician confronted with modeling counts data may first consider using a multinomial logistic regression:

$$p(\mathbf{c}_i | \mathbf{v}_i, m_i) = MN(\mathbf{c}_i; \mathbf{q}_i, m_i) \text{ for } i = 1 \dots n, \quad (1)$$

$$q_{ij} = \frac{e^{\eta_{ij}}}{\sum_{k=1}^d e^{\eta_{ik}}} \text{ for } j = 1 \dots d, \quad (2)$$

$$\eta_{ij} = \alpha_j + \mathbf{v}_i' \boldsymbol{\varphi}_j. \quad (3)$$

When the number of categories d is very large, as is the case in many natural language processing applications, estimating the parameters of the multinomial, $\boldsymbol{\alpha} = [\alpha_j]$ and $\boldsymbol{\varphi} = [\varphi_{ij}]$, is computationally prohibitive.⁷ Equation (2), which makes sure that word probabilities q_{ij} add up to one, is the main barrier to parallelization across categories because every parameter change must be communicated to all other category estimators.

It is well known that the multinomial can be decomposed into independent poisson conditional on the intensities $e^{\eta_{ij}}$, scaled by a poisson for total word counts m_i ,

$$MN(\mathbf{c}_i; \mathbf{q}_i, m_i) = \frac{\prod_j Po(c_{ij}; e^{\eta_{ij}})}{Po(m_i; \sum_{j=1}^d e^{\eta_{ij}})}. \quad (4)$$

Motivated by this decomposition, Taddy (2015) develops the distributed multinomial regression

⁶For example, in a finance application, v_i could be the market response to an earnings release or return volatility.

⁷For example, our application in Section 3 has a vocabulary d of more than five hundred thousand phrases.

(DMR), a parallel (independent) poisson plug-in approximation to the multinomial,

$$p(\mathbf{c}_i | \mathbf{v}_i, m_i) = MN(\mathbf{c}_i; \mathbf{q}_i, m_i) \approx \prod_j Po(c_{ij}; m_i e^{\eta_{ij}}). \quad (5)$$

The parameters for each category j can then be estimated independently with negative log likelihood

$$l(\alpha_j, \boldsymbol{\varphi}_j | \mathbf{c}_j, \mathbf{v}) = \sum_{i=1}^n \left[m_i e^{\alpha_j + \mathbf{v}_i' \boldsymbol{\varphi}_j} - c_{ij} (\alpha_j + \mathbf{v}_i' \boldsymbol{\varphi}_j) \right]. \quad (6)$$

Intuitively, each independent poisson intensity $\lambda_{ij} = m_i e^{\alpha_j + \mathbf{v}_i' \boldsymbol{\varphi}_j}$ is shifted to account for the fact that all words are more likely to appear in longer (high m_i) documents. Approximation (5) removes the communication bottleneck of recomputing $\sum_{j=1}^d e^{\alpha_j + \mathbf{v}_i' \boldsymbol{\varphi}_j}$ and allows for fast and scalable distributed estimation.

Taddy (2013a, 2015) uses DMR to estimate a low dimensional sufficient reduction projection

$$\mathbf{z}_i = \hat{\boldsymbol{\varphi}}' \mathbf{c}_i$$

and shows that \mathbf{v}_i is independent of \mathbf{c}_i conditional on \mathbf{z}_i . This means that, within this model, \mathbf{z}_i is a sufficient statistic that summarizes all of the content that the text has for predicting the covariates \mathbf{v}_i (or its individual elements). For example, suppose v_{iy} is the first element of \mathbf{v}_i , which is available in a subsample, but needs to be predicted for other subsamples. The first step would be to run a multinomial inverse regression of word counts on the covariates \mathbf{v} in the training sample to estimate $\hat{\boldsymbol{\varphi}}$. Second, estimate a forward regression (linear or higher order)

$$\mathbb{E}[v_{iy}] = \beta_0 + [z_{iy}, \mathbf{v}_{i,-y}, m_i]' \boldsymbol{\beta} \quad (7)$$

where $z_{iy} = \sum_j \hat{\varphi}_{jy} c_{ij}$ is the projection of phrase counts in the direction of v_{iy} . Finally, the forward regression can be used to predict v_{iy} using text and the remaining covariates $\mathbf{v}_{i,-y}$.

2.1 Hurdle distributed multinomial regression

In many cases, and specifically in text applications, the poisson is a poor description of word counts c_{ij} . For example, Figure 1 shows the mean histogram (across documents) for the corpus

we use below, which consists of 10,000 two-word phrases (bigrams) appearing in the title and lead paragraph of front page *Wall Street Journal* articles. The left panel shows a substantial mass point at zero that is hard to reconcile with a poisson. The panel on the right shows that if we restrict attention to positive counts, a (truncated) poisson is a reasonable approximation for the data. In our experience, this sparsity is a feature of many text samples, which is why most text analysis software packages use sparse matrices to store word counts efficiently. As alluded to above, the economics of natural language selection provides many reasons for this sparsity. Furthermore, the economics of content selection suggests that relationships between covariates and text could be better captured by allowing for a separate inclusion choice. For example, the decision to start writing about a topic could be more informative than writing more about a topic.

To model text selection, we replace the independent poissons with a two part hurdle model for counts c_{ij} , which we label the *hurdle distributed multinomial regression (HDMR)*:

$$h_{ij}^* = \gamma_i + \kappa_j + \mathbf{w}_i' \boldsymbol{\delta}_j + v_{ij}, \quad (8)$$

$$h_{ij} = \mathbf{1} \left(h_{ij}^* > 0 \right), \quad (9)$$

$$c_{ij}^* = \lambda \left(\mu_i + \alpha_j + \mathbf{v}_i' \boldsymbol{\varphi}_j \right) + \varepsilon_{ij}, \quad (10)$$

$$c_{ij} = \left(1 + c_{ij}^* \right) h_{ij}. \quad (11)$$

The first two equations describe the choice to include ($h_{ij} = 1$) or exclude ($h_{ij} = 0$) word j in document i , often referred to as the model for zeros or participation. This choice depends on observable covariates $\mathbf{w}_i \in \mathbb{R}^{p_w}$ and an unobservable v_{ij} . Equation (10) is the model for word repetition given inclusion in the document, which can depend on the same or other covariates $\mathbf{v}_i \in \mathbb{R}^{p_v}$ and an unobservable ε_{ij} . The last equation says that we only observe positive counts for included words.⁸

Let Π_0 denote the discrete density for zeros

$$p(h_{ij} = 0 | \mathbf{w}_i) = \Pi_0(\gamma_i + \kappa_j + \mathbf{w}_i' \boldsymbol{\delta}_j).$$

⁸Our two part model is simpler and faster to estimate than Mullahy (1986)'s hurdle, which models positive counts as drawn from a truncated poisson, as opposed to a regular poisson for *counts in excess of one* (repetitions). In a previous draft we used the truncated poisson and found very similar results.

Natural choices for Π_0 are the probit and logit binary choice models. Let P^+ denote the model for word repetition, so that conditional on inclusion,

$$p(c_{ij}^* | \mathbf{v}_i, h_{ij} = 1) = P^+ \left(c_{ij}^*; \lambda \left(\mu_i + \alpha_j + \mathbf{v}_i' \boldsymbol{\varphi}_j \right) \right).$$

Natural choices for P^+ are the poisson and the negative binomial. Combining terms, the joint density is

$$p(c_{ij} | \mathbf{v}_i, \mathbf{w}_i) = [\Pi_0(\gamma_i + \kappa_j + \mathbf{w}_i' \boldsymbol{\delta}_j)]^{1-h_{ij}} \left\{ [1 - \Pi_0(\gamma_i + \kappa_j + \mathbf{w}_i' \boldsymbol{\delta}_j)] P^+(c_{ij} - 1; \lambda(\mu_i + \alpha_j + \mathbf{v}_i' \boldsymbol{\varphi}_j)) \right\}^{h_{ij}} \quad (12)$$

The negative log likelihood takes a convenient form

$$l(\boldsymbol{\mu}, \boldsymbol{\alpha}, \boldsymbol{\varphi}, \boldsymbol{\gamma}, \boldsymbol{\kappa}, \boldsymbol{\delta} | \mathbf{c}, \mathbf{v}, \mathbf{w}) = \sum_{j=1}^d l(\mu_i, \alpha_j, \varphi_j, \gamma_i, \kappa_j, \delta_j | \mathbf{c}_j, \mathbf{v}, \mathbf{w}), \quad (13)$$

$$l(\mu_i, \alpha_j, \varphi_j, \gamma_i, \kappa_j, \delta_j | \mathbf{c}_j, \mathbf{v}, \mathbf{w}) = l^0(\gamma_i, \kappa_j, \boldsymbol{\delta}_j | \mathbf{h}_j, \mathbf{w}) + l^+(\mu_i, \alpha_j, \boldsymbol{\varphi}_j | \mathbf{c}_j, \mathbf{v}), \quad (14)$$

$$l^0(\gamma_i, \kappa_j, \boldsymbol{\delta}_j | \mathbf{h}_j, \mathbf{w}) = - \sum_{i|h_{ij}=0}^n \log \Pi_0(\gamma_i + \kappa_j + \mathbf{w}_i' \boldsymbol{\delta}_j) - \sum_{i|h_{ij}=1}^n \log [1 - \Pi_0(\gamma_i + \kappa_j + \mathbf{w}_i' \boldsymbol{\delta}_j)], \quad (15)$$

$$l^+(\mu_i, \alpha_j, \boldsymbol{\varphi}_j | \mathbf{c}_j, \mathbf{v}) = - \sum_{i|h_{ij}=1}^n \log P^+(c_{ij} - 1; \lambda(\mu_i + \alpha_j + \mathbf{v}_i' \boldsymbol{\varphi}_j)). \quad (16)$$

Note that the coefficients γ_i and μ_i introduce dependence across the log likelihood of different words j . To allow for separation across words, and parallelized estimation, we adapt the argument in Taddy (2015) and use plug in estimators $\hat{\gamma}_i$ and $\hat{\mu}_i$ that approximate the maximum likelihood estimators under certain conditions. We discuss these plug in estimators in Section 2.3 and in Appendix B.

A useful feature of the hurdle is that exclusion ($h_{ij} = 0$) is the only source of zero counts. As a result, it decomposes as in (14) into two parts that can be estimated independently, which facilitates further parallelization.⁹ Specifically, the parameters that govern inclusion ($\kappa_j, \boldsymbol{\delta}_j$) only depend on word j indicators \mathbf{h}_j and on the covariates \mathbf{w} , whereas the parameters that govern repetition ($\alpha_j,$

⁹Zero inflation models are alternative approaches that allow for latent $c_{ij}^* = 0$, in which case zero count observations could result either from exclusion or from inclusion of zero counts. While this distinction is philosophically interesting, the hurdle is more tractable and faster to estimate.

φ_j) only depend on word repetition $c_j - 1 > 0$ and the covariates \mathbf{v} and can be estimated separately in the subsample of word repetitions.

HDMR therefore allows one to estimate text selection in Big Data applications of previously impossible scale, by distributing computation across categories and across the two parts of the selection model.

2.2 Regularization

In many machine learning applications, the feature space (words) is much larger than the number of observations. In our setting, even though each word selection model is “small” with number of parameters growing with the number of covariates \mathbf{w} , our approach requires one model estimate per word. Thus the number of estimated coefficients grows with the number of words and the potential for overfitting is large. In such cases, regularization by penalizing nonzero or large φ and δ coefficients is key to avoid overfit. Our results use L_1 regularization separately for each category j and for each of the two parts of the hurdle

$$\hat{\kappa}_j, \hat{\delta}_j = \arg \min_{\kappa_j, \delta_j} l^0(\hat{\gamma}_i, \kappa_j, \delta_j | \mathbf{h}_j, \mathbf{w}) + n\lambda^0 \sum_{k=1}^{p_w} |\delta_{jk}| \quad \text{where } \lambda^0 \geq 0, \quad (17)$$

$$\hat{\alpha}_j, \hat{\varphi}_j = \arg \min_{\alpha_j, \varphi_j} l^+(\hat{\mu}_i, \alpha_j, \varphi_j | \mathbf{c}_j, \mathbf{v}) + n\lambda^+ \sum_{k=1}^{p_v} |\varphi_{jk}| \quad \text{where } \lambda^+ \geq 0. \quad (18)$$

The penalties λ^0 and λ^+ shrink the loadings toward zero, and because of the Lasso-type L_1 penalties, result in many zero loadings (Tibshirani, 1996).¹⁰ Because the model for positive counts only depends on documents i that include word j , the penalty is normalized by the number of such documents $n^+ \equiv \sum_{i=1}^n h_{ij}$. Fast coordinate descent algorithms for these minimization problems have been proposed by Friedman, Hastie, and Tibshirani (2010), which trace out regularization paths of solutions, one for each of a grid of λ 's, for the class of generalized linear models (GLM, McCullagh and Nelder, 1989). We follow Taddy (2017) in selecting the model that minimizes a corrected AIC, though in relatively modest applications one could use cross validation to select the optimal penalty. To apply the coordinate descent algorithms developed by Friedman, Hastie, and

¹⁰We focus on Lasso penalties here to simplify the exposition. Our `HurdleDMR` package allows for more general elastic net-type regularization as in `glmnet` (Friedman, Hastie, and Tibshirani, 2010), and for concave regularization paths as in `gamlr` (Taddy, 2017).

Tibshirani (2010) to our selection model, we frame its two parts as GLMs: a binomial-logit for the inclusion part and a poisson-log for the repetition part (McCullagh and Nelder, 1989).

2.3 Sufficient reduction projections

For simplicity, in what follows we focus on the case where h_{ij} is binomial (bernoulli) distributed

$$p(h_{ij}|\mathbf{w}_i) = [\Pi_0(\hat{\gamma}_i + \kappa_j + \mathbf{w}'_i \boldsymbol{\delta}_j)]^{1-h_{ij}} [1 - \Pi_0(\hat{\gamma}_i + \kappa_j + \mathbf{w}'_i \boldsymbol{\delta}_j)]^{h_{ij}} \quad (19)$$

with a logit link,

$$\log((1 - \Pi_0)/\Pi_0) = \kappa_j + \mathbf{w}'_i \boldsymbol{\delta}_j + \hat{\gamma}_i. \quad (20)$$

In this case we show in Appendix B that for a large enough number of categories d , the MLE estimator for γ_i converges to $\log\left(\frac{d_i}{d-d_i}\right)$, where $d_i = \sum_{j=1}^d h_{ij}$ is the number of categories used in document i (vocabulary size). We therefore use $\hat{\gamma}_i = \log\left(\frac{d_i}{d-d_i}\right)$ as our plug in estimator.

The distribution of word repetitions c_{ij}^* is assumed to be poisson,

$$P^+(c_{ij}^*; \lambda_{ij}) = Po(c_{ij} - 1; \lambda_{ij}) = \lambda_{ij}^{c_{ij}-1} e^{-\lambda_{ij}}, \quad (21)$$

with log intensity

$$\log \lambda_{ij} = \alpha_j + \mathbf{v}'_i \boldsymbol{\varphi}_j + \hat{\mu}_i. \quad (22)$$

Here we follow Taddy (2015) and use $\hat{\mu}_i = \log(m_i - d_i)$ as our plug in estimator. The only difference from Taddy (2015) is that here we are modeling word repetition so the estimator is $\log(m_i - d_i)$ instead of $\log(m_i)$. See Appendix B for details.

We next show that under these functional forms, the entire empirical content of the text that is useful for predicting a variable in \mathbf{w} or \mathbf{v} , is summarized by two low dimension sufficient statistics.

Result 1. *Assuming a binomial-logit model for inclusion and a poisson-log model for word repetition, the projection $\boldsymbol{\delta} \mathbf{h}_i$ is a sufficient statistic for \mathbf{w}_i and the projection $\boldsymbol{\varphi}(\mathbf{c}_i - \mathbf{h}_i)$ is a sufficient statistic for \mathbf{v}_i , conditional on total counts m_i and vocabulary size d_i . Specifically,*

$$\mathbf{v}_i, \mathbf{w}_i \perp\!\!\!\perp \mathbf{h}_i, \mathbf{c}_i | \boldsymbol{\delta} \mathbf{h}_i, \boldsymbol{\varphi}(\mathbf{c}_i - \mathbf{h}_i), m_i, d_i.$$

Proof. That $\boldsymbol{\varphi}(\mathbf{c}_i - \mathbf{h}_i)$ is a sufficient statistic for \mathbf{v}_i follows immediately from the DMR case discussed in Section 2. To establish sufficiency of $\boldsymbol{\delta h}_i$, note that the likelihood for counts \mathbf{c}_i given observed covariates \mathbf{v}_i and \mathbf{w}_i can be factored into

$$p(\mathbf{c}_i | \mathbf{v}_i, \mathbf{w}_i) = p(\mathbf{h}_i \circ \mathbf{c}_i | \mathbf{v}_i, \mathbf{w}_i) = \phi(\mathbf{c}_i) \psi(\mathbf{h}_i) a(\mathbf{w}_i, d_i) b(\mathbf{v}_i, m_i, d_i) \exp(\mathbf{w}_i' \boldsymbol{\delta h}_i + \mathbf{v}_i' \boldsymbol{\varphi}(\mathbf{c}_i - \mathbf{h}_i)), \quad (23)$$

where

$$\begin{aligned} \psi(\mathbf{h}_i) &= \prod_{j=1}^d e^{(\log(\frac{d_i}{d-d_i}) + \kappa_j) h_{ij}}, \\ \phi(\mathbf{c}_i) &= \prod_{j|h_{ij}=1}^d \exp((c_{ij} - 1)(\log(m_i - d_i) + \alpha_j) - \log((c_{ij} - h_{ij})!)), \\ a(\mathbf{w}_i, d_i) &= \prod_{j=1}^d \frac{1}{1 + e^{\log(\frac{d_i}{d-d_i}) + \kappa_j + \mathbf{w}_i' \boldsymbol{\delta}_j}}, \\ b(\mathbf{v}_i, m_i, d_i) &= \prod_{j|h_{ij}=1}^d \exp(-(m_i - d_i) e^{\alpha_j + \mathbf{v}_i' \boldsymbol{\varphi}_j}), \end{aligned}$$

and we use the fact that the Hadamard product $\mathbf{h}_i \circ \mathbf{c}_i$ is equivalent to \mathbf{c}_i here. Hence, the usual sufficiency factorization (e.g., Schervish, 1995, Theorem 2.21) applies yielding the stated result. See Appendix C for additional details. \square

Result 1 means that once we estimate the HDMR parameters, we can reduce the high-dimension (d) text into low-dimension ($p_v + p_w$) sentiment scores from the text in the direction of the covariates in \mathbf{v} or \mathbf{w} . The projections provide useful summaries of the text, which can be plotted or used as a dimensionality reduction first-step into a more elaborate analysis. As in Taddy (2015), the sufficient statistics $\boldsymbol{\delta h}_i$ and $\boldsymbol{\varphi}(\mathbf{c}_i - \mathbf{h}_i)$ rely on population parameters, but in practice we use plug-in estimators $\hat{\boldsymbol{\delta h}}_i$ and $\hat{\boldsymbol{\varphi}}(\mathbf{c}_i - \mathbf{h}_i)$. Whether these provide useful approximations is a setting-dependent empirical matter.

2.4 Inverse regression for prediction

The end goal of many machine learning or natural language processing applications is out-of-sample prediction. Result 1 provides a guide to supervised learning from text via an inverse regression

of the text on the target variable and other covariates Taddy (2013a). The parameters from the HDMR inverse regression in the training sample are used to form a bivariate sufficient reduction projection of the text. A forward regression (still in the training sample) of the target variable on these projections plus the other covariates is then used to construct the predictor.

More concretely, suppose the target variable $v_{iy} = w_{iy}$ is an element of both \mathbf{v}_i and \mathbf{w}_i . We first construct two univariate sufficient reduction projections $z_{iy}^0 = \boldsymbol{\delta}_y \mathbf{h}_i$ and $z_{iy}^+ = \boldsymbol{\varphi}_y(\mathbf{c}_i - \mathbf{h}_i)$. Because the estimated loadings $\boldsymbol{\delta}_y$ and $\boldsymbol{\varphi}_y$ are partial effects, controlling for the other covariates, $\mathbf{w}_{i,-y}$ and $\mathbf{v}_{i,-y}$, the projections z_{iy}^0 and z_{iy}^+ correspond to partial associations as well. Conditional on the parameters, z_{iy}^0 contains all the information that is useful for predicting v_{iy} from the selection of words used in the text (the extensive margin). Similarly, z_{iy}^+ contains the incremental predictive information in repeating words within document i (the intensive margin). Intuitively, HDMR can learn separately from both the extensive and intensive margins, and use them for more efficient prediction. We would then estimate a forward regression (linear or higher order)

$$\mathbb{E}[v_{iy}] = \beta_0 + [z_{iy}^0, z_{iy}^+, \mathbf{w}_{i,-y}, \mathbf{v}_{i,-y}, m_i, d_i]' \boldsymbol{\beta}, \quad (24)$$

which can be used to predict v_{iy} using text and the remaining covariates $\mathbf{w}_{i,-y}$ and $\mathbf{v}_{i,-y}$. If the target variable is only an element of \mathbf{w} (\mathbf{v}), one would only use z_{iy}^0 (z_{iy}^+) in the forward regression. Note that because the variance of z_{iy}^+ and z_{iy}^0 grow with $m_i - d_i$ and d_i respectively, more extreme forecasts will tend to come from larger documents (high $m_i - d_i$ and/or high d_i documents). As in Taddy (2013b) we adjust for this scale effect by normalizing the sufficient reduction projections by the document size, i.e. $z_{iy}^+ = \boldsymbol{\varphi}_y(\mathbf{c}_i - \mathbf{h}_i)/(m_i - d_i)$ and $z_{iy}^0 = \boldsymbol{\delta}_y \mathbf{h}_i/d_i$.

3 Application I: Partisanship of congressional speech

Gentzkow, Shapiro, and Taddy (2019) have recently suggested a new measure of partisanship, defined as the posterior probability that an observer with a neutral prior expects to identify a speaker's party after hearing the speaker utter a single phrase. In this section we apply our methodology to refine our understanding of partisanship by decomposing speech into the choice of vocabulary (the inclusion decision) and choice of emphasizes (the repetition decision). We use their

data to construct a phrase counts matrix \mathbf{c}_t with rows that represent speakers and columns that represent unique two-word phrases (bigrams). We next provide a brief summary of both methods, and refer the reader to that paper for a more detailed description of the data.¹¹

3.1 Methods

Formally, Gentzkow, Shapiro, and Taddy (2019) model the probability $q_{ij}^{P(i)}$ that a speaker i affiliated with party $P(i) \in \{R, D\}$ in congressional session t utters phrase j , given a vector of speaker characteristics \mathbf{x}_{it} , as

$$q_{jt}^{P(i)}(\mathbf{x}_{it}) = e^{u_{ijt}} / \sum_l e^{u_{ilt}}, \quad (25)$$

$$u_{ijt} = \alpha_{jt} + \mathbf{x}_{it}' \boldsymbol{\gamma}_{jt} + \varphi_{jt} \mathbf{1}_{i \in R_t}, \quad (26)$$

where R_t is the set of Republicans and D_t the set of Democrats who spoke at least one phrase in session t . The u_{ijt} terms can be interpreted as speaker i 's utility from saying phrase j at time t . Their preferred estimator for these probabilities is DMR, which regularizes the group difference coefficients φ_{jt} with a Lasso (L1) penalty. The φ_{jt} 's quantify the divergence in phrase j use between Republicans and Democrats. Given estimates of phrase probabilities, the *partisanship* of speech is defined as

$$\pi_t(\mathbf{x}) = \frac{1}{2} \mathbf{q}_t^R(\mathbf{x})' \boldsymbol{\rho}_t(\mathbf{x}) + \frac{1}{2} \mathbf{q}_t^D(\mathbf{x})' (1 - \boldsymbol{\rho}_t(\mathbf{x})), \quad (27)$$

where

$$\rho_{jt}(\mathbf{x}) = \frac{q_{jt}^R(\mathbf{x})}{q_{jt}^R(\mathbf{x}) + q_{jt}^D(\mathbf{x})}, \quad (28)$$

is the posterior belief that an observer with a neutral prior assigns to a speaker being Republican if she utters phrase j in session t and has characteristics \mathbf{x} .

Average partisanship in session t is the average partisanship across speakers i of $\pi_t(\mathbf{x}_{it})$. It provides a concise summary of how well an observer with knowledge of the true model can guess a speaker's party upon hearing them utter a single phrase. It takes into account both how revealing is a given phrase j about its speaker's identity (ρ_{jt}), and how likely is this phrase to be uttered by a Republican (q_t^R) or a Democrat (q_t^D).

¹¹We thank Matthew Gentzkow's for making the data available in his homepage.

We depart from their approach by assuming that phrase choice is a two part process. The speaker first chooses which phrases to include in a speech, and then separately chooses how many times to repeat each phrase. Both choices depend on the same set of covariates \mathbf{x} as before, and potentially on party affiliation. Specifically, we assume phrase inclusion choice is governed by a product of binomial logits, as in (20), with probabilities

$$q_{jt}^{hP}(\mathbf{x}) = Pr\{\text{inclusion}_j|\mathbf{x}\} = e^{h_{jt}^P(\mathbf{x})} / \sum_l e^{h_{lt}^P(\mathbf{x})}, \quad (29)$$

$$h_{ijt}^{P(i)} = \alpha_{jt}^h + \mathbf{x}_{it}'\boldsymbol{\gamma}_{jt}^h + \varphi_{jt}^h \mathbf{1}_{i \in R_t}, \quad (30)$$

whereas repetition choice is governed by a multinomial, as in (21), with log intensity

$$q_{jt}^{rP}(\mathbf{x}) = Pr\{\text{repetition}_j|\mathbf{x}, \text{inclusion}_j\} = e^{r_{jt}^P(\mathbf{x})} / \sum_l e^{r_{lt}^P(\mathbf{x})}, \quad (31)$$

$$r_{ijt}^{P(i)} = \alpha_{jt}^r + \mathbf{x}_{it}'\boldsymbol{\gamma}_{jt}^r + \varphi_{jt}^r \mathbf{1}_{i \in R_t}. \quad (32)$$

We estimate these conditional probabilities with HDMR, regularizing the group difference coefficients φ_{jt}^h and φ_{jt}^r with Lasso penalties.

We then define two natural variants of partisanship by replacing the phrase probabilities in (27) and (28) with q_{jt}^{hP} or with q_{jt}^{rP} , to form *inclusion partisanship* and *repetition partisanship*. Inclusion partisanship measures what an observer would expect to learn about a speaker's party affiliation from the inclusion of a phrase in a speech. By contrast, repetition partisanship measures what they expect to learn from an additional repetition of a phrase already observed.

3.2 Results

Figure 2 reports the two new HDMR-based measures of congressional speech partisanship. For comparison, we reproduce the Gentzkow, Shapiro, and Taddy (2019) DMR-based measure of partisanship for comparison. It clearly shows that average partisanship based on any utterance is fairly low starting in 1873 until the early 1990s when it starts increasing sharply, as documented by Gentzkow, Shapiro, and Taddy (2019).¹²

¹²The line labeled DMR is quite similar to the one reported in Gentzkow, Shapiro, and Taddy (2019) as their preferred estimator, despite slight differences in the software packages we use to run DMR.

Our decomposition reveals that earlier spikes in partisanship manifested in increased repetition of different phrases, whereas the upward trend starting in the 1990s is due to entirely distinct phrase selection. We find that repetition partisanship spikes to 0.504 in the late 1890s and to 0.505 in the late 1920s. These levels are similar to those seen in the 2000s. These peaks are comparable with the highest level of partisanship based on DMR. But inclusion partisanship remains close to the 0.5 prior for most of the sample, until it starts rising gradually in the 1970s, before spiking up in the 1990s to previously unseen levels.

These findings mean that politicians belonging to the two parties have previously emphasized different phrases or issues to a varying degree, while using similar vocabularies. By contrast, the recent spike in partisanship is quite different as political parties seem to be using distinct vocabularies or speaking about completely different issues.

3.3 Partisan phrases

To better understand these changes in the nature of congressional speech, we next investigate which phrases drive these results. For this purpose, [Gentzkow, Shapiro, and Taddy \(2019\)](#) define *phrase partisanship* as the average effect on the expected posterior that speaker i is Republican of removing a given phrase j ,

$$\zeta_{jt}(\mathbf{x}_{it}) = \frac{1}{2} - \frac{1}{2} \sum_{k \neq j} \left(\frac{q_{kt}^R(\mathbf{x}_{it})}{1 - q_{jt}^R(\mathbf{x}_{it})} + \frac{q_{kt}^D(\mathbf{x}_{it})}{1 - q_{jt}^D(\mathbf{x}_{it})} \right) \rho_{kt}(\mathbf{x}_{it}). \quad (33)$$

Larger positive ζ_{jt} phrases are more likely to be spoken by Republicans and more negative ζ_{jt} phrases by Democrats. Averaging ζ_{jt} across speakers i in a congressional session t provides an intuitive ranking of the most partisan phrases per session.

Table 1 reports the most partisan phrases in several key periods along the history reported in Figure 2. For each session, the first column is based on inclusion partisanship, the second on repetition partisanship, and the third corresponds to any utterance as in [Gentzkow, Shapiro, and Taddy \(2019\)](#). We also report the predicted number of occurrences for each phrase by each party per 100,000 phrases. More partisan phrases, exhibit a sharper imbalance in these predicted occurrences.

The upward spike in inclusion partisanship, visible in Figure 2, is hard to miss in the partisan

phrases of the 110th session of congress. Republicans exclusively use the term “war on terror” to refer to the same war that Democrats call the “war in Iraq.” Republicans say “tax increases” while Democrats say “fiscal response.” Even with their common language, the two parties emphasize different issues in their speeches, as can be seen from the second column. For example, conditional on its inclusion in a speech, Democrats are more likely to repeatedly mention the growing national debt while Republicans emphasize tax payments.

Repetition partisanship attains its highest level at the very end of the sample, during the 114th session of congress, which adjourned in January 2016, at the same time that inclusion partisanship is waning. Republican members of congress at the time are more likely to repeatedly refer to private property and the Internal Revenue Service, while Democrats often mention the Postal Service and Public Service. Despite its decline, inclusion partisanship is still visible as Republicans prefer to discuss tax policy and terrorism, while Democrats are more likely to include gun violence and civil rights in their speeches.

DMR-based partisanship on the rightmost column, is comprised of a mix of both inclusion and repetition choice. For example, during the 114th session, the list of most partisan phrases according to this single metric, combines top inclusion phrases (e.g. “taxpayer dollar,” “gun violence”) and top repetition phrases (e.g. “American people,” “men and women”).

Back in the late 1920s, however, there is little divergence in congressmembers’ choice of phrases to include in their speeches, but there is ample divergence in their choice of phrases to emphasize. For example, during the 70th session, Republicans are much more likely to emphasize the American people and state government, while Democrats are more likely to emphasize Great Britain and the interior department. Finally, an example of the low level of partisanship of the late 1960s is provided by the 90th session of congress—a time when both parties include and repeat similar phrases.

4 Application II: Backcasting the intermediary capital ratio

A rapidly growing body of work finds empirical support for intermediary asset pricing theories (He and Krishnamurthy, 2013; Brunnermeier and Sannikov, 2014; Adrian, Etula, and Muir, 2014; Haddad and Muir, 2018; Baron and Muir, 2018; He and Krishnamurthy, 2018; Koijen and Yogo, 2019). In particular He, Kelly, and Manela (2017) find that a simple two-factor model that in-

cludes the excess stock market return and the aggregate capital ratio of major financial intermediaries—primary dealers—can explain cross-sectional variation in expected returns across a wide array of asset classes. They also present preliminary results on return predictability (time-series regressions), but their conclusions are limited by a relatively short time-series that starts in 1970. Prior to 1970, most primary dealers were private, which precludes a calculation of their capital ratio.

We conjecture that as a publication catering to investors, text that appears on the front page of the *Wall Street Journal* would be informative about the aggregate state of the intermediary sector. Dire language on financial intermediaries’ failure is used to cover unfolding crises like the financial crisis of 2008, the LTCM liquidity crisis following Russia’s default in 1998, and the failure of important dealers like Drexel Burnham Lambert in 1990.

4.1 Data

Our text data includes all titles and lead paragraphs that appear on the front page of the *Wall Street Journal* from July 1926 to February 2016. We include the 10,000 most frequent two-word phrases (bigrams) in separate sentences, after removing case, stopwords, nonletters and Porter stemming. We aggregate the data to the monthly frequency so that \mathbf{c}_t are phrase counts observed during month t .

Figure 1 shows the mean histogram for phrase counts in this sample. The left panel shows that the entire range is highly sparse (has many zeros). The right panel omits zero counts, and shows that a truncated (at zero) poisson distribution is a reasonable approximation for the positive range of counts.

We match this data with the monthly intermediary capital ratio icr_t of He, Kelly, and Manela (2017).¹³ This ratio is our prediction target and is therefore the first element of both covariate vectors \mathbf{v}_t and \mathbf{w}_t . We additionally include in both vectors, two natural covariates that are likely to be correlated with the icr_t : the log price dividend ratio (pd_t) from CRSP, and the realized variance of financial stocks ($rvfin$) over the same month and the prior one. Table 2 reports summary statistics for these variables.

¹³The icr is available at <http://apps.olin.wustl.edu/faculty/manela/data.html>

4.2 Newspaper coverage choice

Our selection model is parametrically identified and therefore technically does not require that different variables be used in the inclusion and repetition equation. However, Heckman (1979) selection models are known to be nonparametrically identified if a continuous variable enters the selection equation but can be excluded from the second equation (Gallant and Nychka, 1987; Heckman and MaCurdy, 1986). Proving such a result in our setting can be useful, but left for future work. Motivated by their insight, we seek an instrument for word inclusion.

Boydston (2013) suggests that prior attention to an issue may influence its coverage by the press. The idea is that once fixed costs such as journalist travel and familiarization with an issue have been incurred, a marginal article is easier to produce. To capture prior attention to the financial sector, we add prior year realized variance of financial stocks ($rvfin_{t-13 \rightarrow t-1}$) as an explanatory variable in the model for inclusion alone, excluding it from the repetition equation. This choice assumes that after conditioning on the intermediary capital ratio, the price-dividend ratio, current and lagged monthly $rvfin$, and on phrase j being included in the *Journal*, the number of times this phrase is repeated does not depend on the prior year's volatility of financial stocks.¹⁴

4.3 Sparsity and out-of-sample fit

A key choice in the preprocessing stage of many text analyses, is to omit words or phrases that rarely appear in the sample. For example, we may keep the X most frequent phrases. From the vantage point of our selection model, this choice is important. If the “cleansed” word counts matrix \mathbf{c} is highly dense because phrases that often do not appear in the text are excluded from the analysis, then the benefit of modeling the extensive margin is likely to be low. Therefore, we assess the improvement in out-of-sample fit as a function of the number of most frequent phrases kept in the sample.

Figure 3 compares the out-of-sample root mean squared error from a 10-fold cross validation exercise. It compares HDMR with DMR, which is provided with the same covariates and text, and with a linear regression of the target on the same covariates but without the text. Both DMR and

¹⁴We have experimented with business news pressure (Manela, 2014), a variant of television news pressure (Eisensee and Stromberg, 2007), that is available for a shorter sample starting 1967, and found similar out-of-sample fit improvements.

HDMR improve considerably over the No Text benchmark and reduce the error by more than 40 percent. We can see that when only a few hundred words are considered in vocabulary, both DMR and HDMR generate similar improvements, but as rarer phrases are included in the vocabulary, selection plays a bigger role, and the benefit from using HDMR increases. The advantage is hump-shaped, and eventually, as rarely-used phrases enter, repetition choice becomes less important (bottom panel), and the out-of-sample fit of HDMR suffers.¹⁵

Because our prediction exercise involves using data in one time period to predict out-of-sample in a different time, cross-validation with random folds may be misleading when both the text and target variable are persistent. For example, if the model relies heavily on the fact that the phrase “subprime mortgage” appears often in the period around the 2008 financial crisis when intermediary capital was low, but not in earlier parts of the sample, then random cross-validation, which would likely include observations around the same period in the test subsamples, may give an overly optimistic measure of out-of-sample fit.

Figure 4, therefore, uses pseudo out-of-sample prediction, starting with the latter half of the sample and rolling back with wider training subsamples, predicting one observation at a time. Even though this validation approach results in somewhat wider confidence intervals, the results are quite similar to random cross-validation. Table 3 focuses on the optimal model by cross-validation, the one that uses the 10,000 most frequent phrases. We find that regardless of the validation method, HDMR provides a significant reduction in out-of-sample root-mean-squared-error.

4.4 News-implied intermediary capital ratio, 1927–2016

Having established that our model produces good out-of-sample fit, we use it to backcast the intermediary capital ratio back to the June 1927, the first month when a full year of financial volatility is available. Figure 5a shows that the intermediary capital ratio predicted by HDMR closely follows the actual one in the period when the latter is available, 1970–2016. Financial variance and the price-dividend ratio, which alone can explain much of the variation in the *icr*, provide a back-bone for the predictor, as can be seen from the No Text benchmark. DMR, HDMR, and SVR, all use these covariates plus the text to improve prediction, but generate somewhat

¹⁵In Section A.2 we show that for richer text, the repetition margin remains important and HDMR’s advantage persists even with half a million phrases.

different time-series. For example, the HDMR predicted values appear lower than those of DMR and feature more negative spikes in the capital ratio. The SVR predictor makes clear that it does not find the variation in *rvfin* and *pd* important.

Figure 5b zooms in on the Great Depression period. With the exception of the SVR predicted series, all predictors indicate that financial intermediaries were substantially undercapitalized between 1929 and 1939. Figure 5c shows a similar pattern, for the more recent Great Recession period, with all predictors showing a sharp fall in the aggregate capital ratio of financial intermediaries, that matches the actual *icr*'s behavior. Interestingly, financial variables without the text overstate the recovery of the *icr* starting in 2010. Compared with HDMR, DMR understates the *icr*'s recovery between 2012 and 2013. To better understand the source of this difference, we next focus on individual phrases and their importance.

4.5 Which phrases are pivotal for out-of-sample fit?

To better understand the improvements in out-of-sample fit that HDMR generates, we report in Table 4 the phrases whose removal from the corpus causes the largest deterioration in out-of-sample root-mean-squared-error.¹⁶ We report results separately for the two validation methods discussed earlier in Section 4.3 because our intuition is that time-series persistence can lead to overstatement of out-of-sample fit using random cross-validation.

A related concern is that language changes over time, so the phrases associated with movements in the *icr* in the fitted sample, are not useful its prediction in the earlier, 1926–1969 period. For example, **subprime mortgage** mentions in the *Wall Street Journal* are unlikely to be related with financial distress prior to the 2007–2009 financial crisis. But the important question for our purposes is whether HDMR can avoid overfitting to such phrases. Manela and Moreira (2017, Section 2.3) analyze this possibility in detail for the same corpus and a different financial variable as target. They find that while language does change over time, the deterioration it induces in out-of-sample fit is quite modest. Examining the list of pivotal phrases can shed light on this issue.

Table 4, Panel (a) shows that some of the predictive ability of the text stems from phrases that capture fairly robust economic fundamentals, such as the positively associated **jobless marri**. Other phrases such as **barack obama**, a US president elected at the peak of the 2008 financial

¹⁶Gentzkow, Shapiro, and Taddy (2019) identify partisan phrases with a similar approach.

crisis, show up as negatively correlated with the *icr*, even though they are unlikely to be useful for its prediction before 2008. Interestingly, his main opponent during the 2012 election, *mitt romney*, is negatively associated too, suggesting that political uncertainty as opposed to specific policies may be the culprit (Pastor and Veronesi, 2012; Baker, Bloom, and Davis, 2016; Manela and Moreira, 2017; Hassan, Hollander, van Lent, and Tahoun, 2017).

Panel (b) uses pseudo out-of-sample rolling validation instead, and shows as conjectured, greater focus on robust fundamentals that are likely to be relevant over the entire sample. Pivotal phrases have to do with government policy (*tax report*, *washington wire*) and economic conditions (*busi bulletin*, *labor letter*). We therefore find this validation approach more likely to approximate true out-of-sample fit.

4.6 Focusing on a single phrase for intuition

For a more intuitive understanding of how inverse regression loadings translate into forward regression prediction, we next focus on a single phrase, *financi crisi*. We expect front page reports of financial crises to be negative signals about the capital ratio of the intermediary sector.

The backward hurdle regression estimates in the first two columns of Table 5a show that the *icr* is indeed negatively correlated with repeated mentions of *financi crisi*, but also that the mere inclusion of this phrase on the front page is a strong negative signal, conditional on the price-dividend ratio and financial volatility. The positive coefficient on $rvfin_{t-13 \rightarrow t-1}$ means that above average prior year financial volatility is followed by higher financial crisis coverage on the front page. Specifically, a one standard deviation in this instrument relative to its mean implies a 22% increase in the log odds ratio for a financial crisis inclusion (from 0.17 to 0.21). The last column shows that a poisson regression (DMR) treats inclusion and repetition as a single object, and assigns a relatively low weight to this instrument compared to contemporaneous financial volatility (both variables are in annual variance units).

These coefficients are used to construct the two sufficient reduction projections, $z_{ty}^0 = \delta_y \mathbf{h}_t / d_t$ and $z_{ty}^+ = \varphi_y(\mathbf{c}_t - \mathbf{h}_t) / (m_t - d_t)$, and plugged into a forward regression of the *icr* on the these and

the remaining covariates, as described in Section 2.4:

$$y_t = b_0 + b_z z_{ty}^+ + b_s z_{ty}^0 + b_v v_{t,-y} + b_m m_t + b_d d_t + \varepsilon_t.$$

The contribution of a single phrase j to the predicted value is therefore

$$\hat{y}_{tj} = b_z \varphi_{jy} (c_{tj} - h_{tj}) / (m_t - d_t) + b_s \delta_{jy} (h_{tj} / d_t).$$

Table 5b reports the forward regression coefficients' products with those of the backward regression, $b_z \varphi_{jy}$ and $b_s \delta_{jy}$ for HDMR, and contrasts it with the corresponding single coefficient product of DMR. We can see that much of the contribution of `financi crisi` to the predicted value in HDMR comes from the extensive margin. A different way to see this is by looking at the time series \hat{y}_{tj} , which appears in Figure 6. A single mention of financial crises is all it takes for HDMR to predict a lower intermediary capital, with repeated mentions having a much lower effect.

4.7 Time-varying risk premia and the intermediary capital ratio

A central prediction of the intermediary asset pricing model (He and Krishnamurthy, 2012, 2013) is that times when intermediaries are highly capitalized are “good times,” when these marginal investors demand a relatively low risk premium to hold investment assets. Preliminary such time-series predictability regression reported in He, Kelly, and Manela (2017) support this prediction, but the short time-series used there limits the power of these tests.

The news-implied intermediary capital ratio allows us to test this prediction in a larger sample that goes back to 1927. Return predictability tests are reported in Table 6 for the monthly, quarterly and annual horizons. Because such regressions use overlapping observations, we use the standard Hodrick (1992) correction to the standard errors. We report results both for the full sample and for the better-studied postwar subsample because our newspaper text is highly sparse in the early sample, so that most of the variation is on the extensive margin. For each horizon, the first column shows that restricting attention to the actual *icr*, which is only available from 1970, yields a negative but statistically insignificant coefficient. By contrast, we find that the news-implied \widehat{icr} , which is available over a much longer period, is significantly negative both over the postwar sample and the

full sample. The point estimates are fairly consistent across horizons and subsamples and imply that a one standard deviation increase in the \widehat{icr} predicts a 4–5 percentage points lower market premium.

To understand better whether the predictive ability comes from covariates that are known predictors like price-dividend ratio or from the text, we regress future stock market excess returns at various horizons on a variant that does not use the text $\widehat{icr}^{\text{No Text}}$, which shows similar results but with lower explanatory power as measured by adjusted R-squared. We also decompose the text-based \widehat{icr} into lagged sufficient reduction projections z_{t-1}^0 , z_{t-1}^+ , and the covariates. From the decomposition we find that the inclusion projection z_{t-1}^0 is a strong predictor of future market returns in the postwar sample, over and above the price-dividend ratio, and that the repetition projection z_{t-1}^+ is statistically significant only in the full sample.

The results imply that there is a set of phrases whose inclusion on the front page of the *Journal* provides a strong signal about stock market risk premia, over and above the valuation ratio (pd). HDMMR provides an efficient way to identify these phrases and their relative weights in a data driven approach while avoiding overfit.

5 Application III: Forecasting key macroeconomic indicators

[Stock and Watson \(2012\)](#) compare forecasts from various forecasting methods designed for a large number of orthogonal predictors with dynamic factor model (DFM) forecasts using a US macroeconomic dataset with 143 quarterly variables spanning 1960–2008. They find that for most series, DFM-5 forecasts, which are based on a simple linear regression of the target on the top 5 principal components of the lagged variables, are superior to shrinkage forecasts. A large literature has explored ways of improving on their results (e.g. [Stock and Watson, 2011, 2016](#); [Kim and Swanson, 2014](#); [Kelly and Pruitt, 2013, 2015](#); [Bitto and Frühwirth-Schnatter, 2018](#); [Chudik, Kapetanios, and Pesaran, 2018](#); [Boot and Nibbering, 2019](#); [Cepni, Güney, and Swanson, 2019](#)).

Our question is different. We ask whether the text that appears in the *Wall Street Journal* contains additional information that is useful for forecasting these macroeconomic indicators beyond that of the DFM-5 benchmark, and whether HDMMR can extract such information from the text.

5.1 Data

Because we are interested in forecasting, we focus on more recent data for which we have a much richer body of text—the full text of all *Wall Street Journal* articles that appear on the front page of the from January 1990 to December 2010. We include two-word phrases (bigrams) in separate sentences, after removing case, stopwords, nonletters and Porter stemming, and aggregate the data to the monthly frequency so that c_t are phrase counts observed during month t .

We match the text with the monthly macroeconomic indicators dataset made available for replication by [Stock and Watson \(2012\)](#). Following their categorization and normalization, while restricting the analysis to monthly series, we use 92 lower-level disaggregated series to compute principal components, and use 12 headline indicators as prediction targets. [Table 7](#) reports summary statistics for these target variables. The transformed series are generally first differences of logarithms (growth rates) for real quantity variables, first differences for nominal interest rates, and second differences of logarithms (changes in rates of inflation) for price series.

5.2 Methods

We compare the out-of-sample root mean squared error for HDMR against the DFM-5 benchmark and against DMR with the same data. We assess out-of-sample fit via (i) random cross-validation and via (ii) pseudo out-of-sample rolling forward predictions (starting with half the available time-series as training sample and forecasting one month ahead). The latter approach is especially appealing given recent evidence of time-varying predictability ([Farmer, Schmidt, and Timmermann, 2019](#)).

Let $Y_{t+\tau}^\tau$ denote the variable to be forecasted in a τ -period ahead forecast $Y_{t+\tau}^\tau$. In each training sample (fold) we use the following methods for fitting.

For DFM-5, we simply form principal components of the entire sample once, and keep the top 5. We expect this to give DFM-5 a slight advantage over the competing methods. We then run an ordinary least squares regression of the target on the PCs:

$$Y_{t+\tau}^\tau = \beta_0 + [pc_t^1, \dots, pc_t^5]' \beta + \varepsilon_{t+\tau} \quad (34)$$

For HDMR, we use these same 5 PCs as well as the target as explanatory variables in inverse

regressions (8) and (10). We then form sufficient reduction projections in the direction of the target. These projections summarize all the information in the text from phrase inclusion (z_{tY}^0) and repetition (z_{tY}^+) that is useful for predicting the target $Y_{t+\tau}^\tau$ after controlling for the PCs. We then run a forward regression (24) of the target on the sufficient reduction projections and the PCs, still using only training sample data:

$$Y_{t+\tau}^\tau = \beta_0 + \left[z_{tY}^0, z_{tY}^+, pc_t^1, \dots, pc_t^5, m_i, d_i \right]' \beta + \varepsilon_{t+\tau} \quad (35)$$

For DMR, we follow essentially the same procedure as for HDMR, with the same variables used to explain the text, and use its single sufficient reduction projection in the forward regression. For each of the three models we use the predicted values from the forward regression model, but applied to out-of-sample validation observations.

5.3 Forecasting results

Table 8 reports out-of-sample RMSE for HDMR, relative to DFM-5, which uses only the PCs without the text, and relative to DMR with the same data. Lower reported ratios indicate larger improvements from using HDMR. We find significant improvements in out-of-sample forecasts for a few variables at the monthly horizon ($\tau = 1$). For example, in Panel (a), total nonfarm payroll employment (“Emp: total”) sees a 17% (18%) improvement from using HDMR relative to DFM-5 and housing starts sees a 32% (31%) improvement using random cross validation (rolling validation).

We find that the text appearing in the *Journal* is more informative about longer horizons, where HDMR generates significant improvements in most forecasts. For example, using the text improves industrial production and employment forecasts by about 20–30% at the annual horizon ($\tau = 12$). Newspaper coverage also considerably improves asset pricing forecasts at this longer horizon, generating large reductions in RMSE for interest rates, exchange rates, and stocks.

We further find that the advantage of HDMR over DMR increases as we increase the dimensionality of the text and use less frequent and more sparse phrases. Because of the computational cost of this exercise, we simply compare a restricted sample with the most frequent 10,000 phrases in Panel (a) with a larger vocabulary of the most frequent 100,000 phrases in Panel (b). We find, for example, that even though the quarterly inflation (CPI-ALL) forecasts of HDMR are comparable

to those of DMR when we use 10,000 phrases, they are substantially better when we use 100,000 phrases. Similar improvements can be seen for M1 and consumer expectations. These findings suggest that HDMR can better learn from high dimensional sparse data like business newspaper coverage, which is selected based on newsworthiness.

5.4 Nowcasting results

Nowcasting is a related but distinct strand of the forecasting literature, which focuses on predicting activity that occurs now but that will only be reported later. A growing body of work starting with [Giannone, Reichlin, and Small \(2008\)](#) and surveyed in [Bańbura, Giannone, Modugno, and Reichlin \(2013\)](#) compares different methods for nowcasting macroeconomic indicators using lagged and contemporaneous macroeconomic series that is available before the actual measurements are released publicly.

We can evaluate whether the text of the *Journal* is informative about the present, by lagging one month both the target and the PCs so that the predictive forward regression becomes

$$Y_t^1 = \beta_0 + \left[z_{tY}^0, z_{tY}^+, pc_{t-1}^1, \dots, pc_{t-1}^5 \right]' \beta + \varepsilon_t. \quad (36)$$

We then follow the same out-of-sample validation procedure of the forecasting exercise. But now, the question we ask is whether the text of the *Journal* that is reported over month t is informative about macroeconomic activity over the same month Y_t^1 .

Table 9 shows that this body of text, and our approach to learning from it in particular, are highly valuable for nowcasting. The improvements in out-of-sample fit are, on average, better than in the forecasting exercise of Table 8. These results suggest that nowcasting with text may offer substantial gains over using macroeconomic series alone. We leave a detailed investigation of nowcasting for future research, as well as whether our methods can improve upon state-of-the-art nowcasting models that use several factor lags, daily data, or mixed frequency VAR, and that deal with the jagged edge of macroeconomic news in realtime.

6 Conclusion

Text data is inherently high-dimensional, which makes machine learning techniques natural tools for its analysis. Text is often selected by journalists, speechwriters, and others who cater to an audience with limited attention.

We develop an economically-motivated high dimensional selection model that can improve machine learning from text in particular and from sparse counts data more generally. Our highly scalable approach to modeling coverage selection is especially useful in cases where inclusion choice is separate or more interesting than repetition choice. It allows one to analyze text within a regression framework familiar to economists, and strikes a balance between prediction and interpretation.

The economics literature has recently gained neat methods from the machine learning literature. Economics can reciprocate. The three applications we analyzed all show considerable gains from injecting a bit of economics into machine learning models.

A Robustness

A.1 Alternative text regressions

Table 3a focuses on the optimal model by cross-validation, the one that uses the 10,000 most frequent phrases and compares HDMR to several benchmarks. For each model we report the measure of fit with and without the text, and the change in the measure of fit.

The first benchmark is DMR, which is provided with the same covariates and text. The improvement from modeling selection with HDMR is a 13 percent reduction in out-of-sample root mean squared error, from 82 to 71 basis points. HDMR provides a 45 percent improvement relative to the No Text benchmark that only uses the other covariates to predict.

The second benchmark model is a “fabricated” variant of HDMR (FHDMR) which adds $h_{ij} = \mathbf{1}(c_{ij} > 0)$ indicators to the text counts matrix \mathbf{c} and then runs DMR as usual with $\tilde{\mathbf{c}} = [\mathbf{c} \ \mathbf{h}]$. If all that HDMR was doing is allow for a nonlinearity of the counts matrix, we would expect FHDMR to do just as well. Instead we find that it generates an 79 basis point RMSE, which is only slightly better than the 82 bp of DMR.

The last benchmark we consider is support vector regression (SVR), which Manela and Moreira (2017) use for a similar backcasting purpose. We follow their approach to calibrating the SVR meta-parameters. Even though we standardize both text and covariates to unit variance, SVR still cannot concentrate on the covariates, which provide first order information on our prediction target. SVR with text improves on an SVR without text, but its 126 basis points error rate is much larger than that of HDMR.

Table 3b reports pseudo-out-sample rolling back validation results, and finds similar improvements in out-of-sample fit from using HDMR.

A.2 Denser text

For a shorter time-series, 1990 to 2010, we can assess HDMR with much denser text—the full *Wall Street Journal*. Figure 7 shows that in this sample, the advantage from using the richer body of text is larger, as it attains lower out-of-sample error rates. These results, however, could also be driven by the different time period. What does seem like a robust conclusion from this comparison is that the advantage of HDMR over DMR increases with the sparsity of the text, which is plotted in the

bottom panel. When only a small number of highly frequent phrases are analyzed, essentially all are included and repeated, so that modeling selection is less important. However, as we add more phrases to the analysis, and the distinction between inclusion and repetition grows in importance, the advantage of HDMR increases. For this denser corpus, even with $d = 500,000$ phrases, about 40 percent of phrases are included and 14 percent are repeated in the average month, and HDMR reduces out-of-sample root mean squared error by 58 percent (121 to 51 basis points) relative to the No Text benchmark, and by 24 percent (68 to 51 bp) relative to DMR.

B What does HDMR approximate?

The method presented in this paper is developed with the goal of computational efficiency. With this goal in mind we model the intensive margin of the word counts using a poisson distribution. Here we show that our HDMR procedure can be thought as maximizing a likelihood that approximates a mixture of d bernoulis and a multinomial. The bernoulis determine the extensive margin, i.e., whether a given word is used in a given document, and the multinomial determines the intensive margin, i.e. how much a word is used in a given document.

Taddy (2015) shows that we can factorize a poisson distribution for the distribution of word counts in a document as a poisson distribution for the total number of words in the document, and a multinomial distribution for the count of each word given the total word count:

$$p(\mathbf{c}_i) = \prod_j Po(c_{ij}, e^{\eta_{ij}}) = MN(\mathbf{c}_i; \mathbf{q}_i, m_i) Po(m_i; \sum_{k=1}^d e^{\eta_{ik}}).$$

Furthermore, he shows that if η_{ij} is of the form $\eta_{ij} = \mu_i + \alpha_j + \mathbf{v}_i' \boldsymbol{\varphi}_j$, and μ_i in the poisson model is estimated at its conditional MLE, $\mu_i^* = \log \left(\frac{m_i}{\sum_j e^{\eta_{ij}}} \right)$, then the coefficients estimates $\boldsymbol{\varphi}_j$ and α_j are the same regardless of whether the distribution is conditioned on m_i . This means that conditional on μ_i^* , estimating the poisson model is equivalent to estimating the multinomial model.

However, choosing $\mu_i^* = \log \left(\frac{m_i}{\sum_j e^{\eta_{ij}}} \right)$ is computationally costly because it makes the d poisson models depend on each other, and therefore makes it impossible to parallelize the model estimation across words, which is essential to make it useful for any Big Data application. Taddy (2015) shows particular cases where μ_i^* degenerates to $\hat{\mu}_i = \log(m_i)$, and shows that in more general cases the

choice $\hat{\mu}_i$ generalizes well in out-of-sample validation.

We now apply this insight above to our setting and show what our hurdle model approximates. We start from the mixture of d binomials for the extensive margin of each word, and a multinomial distribution for the excess word counts of each word. We then show that we obtain the joint distribution associated with (12).

Lets start with the word count \mathbf{c}_i distribution conditional on word count totals m_i and document vocabulary size d_i . Define $\mathbf{h}_i = \mathbf{c}_i \geq 1$, where the inequality is evaluated word by word, so \mathbf{h}_i is d by 1 vector of ones and zeros. Thus \mathbf{h}_i defines the vocabulary of document i as it identifies all words used in the document. Note that $1'\mathbf{h}_i = d_i \leq d$, i.e., each document vocabulary is smaller than the universal vocabulary. A document vocabulary is given by the mixture of d bernoullis with success probabilities given by Π_{ij} for word j in document i . The bernoullis determine whether $h_{ij} = 0$ or $h_{ij} = 1$, i.e., if word j is part of the used vocabulary in document i . This implies that the expected vocabulary size is $E[\sum_{j=1}^d h_{ij}] = \sum_{j=1}^d \Pi_{ij}$. This holds approximately for large enough universal vocabulary d because each draw is independent, i.e., the law of large numbers implies $E[\sum_{j=1}^d h_{ij}] = \sum_{j=1}^d \Pi_{ij} \approx \sum_{j=1}^d h_{ij} = d_i$.

We can then represent the binomial distribution as

$$p(h_i) = \prod_{j=1}^d B(h_{ij}, \Pi_{ij})$$

with $\Pi_{ij} = 1/(1 + e^{-\gamma_i - \epsilon_{ij}})$, this implies the log-likelihood ratio is,

$$\log \left(\frac{\Pi_{ij}}{1 - \Pi_{ij}} \right) = \log \left(\frac{1/(1 + e^{-\gamma_i - \epsilon_{ij}})}{1 - 1/(1 + e^{-\gamma_i - \epsilon_{ij}})} \right) = \log \left(\frac{1/(1 + e^{-\gamma_i - \epsilon_{ij}})}{(e^{-\gamma_i - \epsilon_{ij}})/(1 + e^{-\gamma_i - \epsilon_{ij}})} \right) = \gamma_i + \epsilon_{ij},$$

where γ_i captures document specific vocabulary richness, and ϵ_{ij} capture the relative probability of usage of word j in document i . So for example, a news paper that talks about very broad set of topics will have large γ_i but ϵ_{ij} are all close around zero. While a finance journal will perhaps have a low γ_i but very disperse ϵ_{ij} with some words being significantly more likely to be used and other significantly less likely.

We now show that in similar spirit as Taddy (2015) we can use an approximation to motivate a

plug in estimator for γ_i , which is document specific, and in general would require all the words to be jointly estimated. We start by doing a linear approximation of the probability Π_{ij} as function of ϵ_{ij} ,

$$\Pi_{ij} = 1/(1 + e^{-\gamma_i - \epsilon_{ij}}) \approx 1/(1 + e^{-\gamma_i}) + 1/(1 + e^{-\gamma_i})^2 e^{-\gamma_i} \epsilon_{ij},$$

we then add probabilities within document and across words,

$$\begin{aligned} \sum_{j=1}^d \Pi_{ij} &\approx \sum_{j=1}^d 1/(1 + e^{-\gamma_i}) + \sum_{j=1}^d 1/(1 + e^{-\gamma_i})^2 e^{-\gamma_i} \epsilon_{ij} \\ &\approx d/(1 + e^{-\gamma_i}) + 1/(1 + e^{-\gamma_i})^2 e^{-\gamma_i} \sum_{j=1}^d \epsilon_{ij} \\ &\approx d/(1 + e^{-\gamma_i}), \end{aligned}$$

where we use the the universal vocabulary is large and the the law of large numbers applies to get $\sum_{j=1}^d \epsilon_{ij} \approx 0$. Finally we use our earlier result that $\sum_{j=1}^d \Pi_{ij} \approx d_i$ to get

$$\hat{\gamma}_i = \log\left(\frac{d_i}{d - d_i}\right),$$

which is our plug-in estimator.

We then use a standard logit approximation and obtain the following log-likelihood ratios,

$$\log\left(\frac{\Pi_{ij}}{1 - \Pi_{ij}}\right) = \gamma_i + \epsilon_{ij} = \log\left(\frac{d_i}{d - d_i}\right) + \epsilon_{ij} = \log\left(\frac{d_i}{d - d_i}\right) + \kappa_j + \mathbf{w}'_i \boldsymbol{\delta}_j,$$

where we model the word-document relative word usage ϵ_{ij} as a function of a word specific component κ_j that captures the fact that some words are just more likely to be used unconditionally, and the term $\mathbf{w}'_i \boldsymbol{\delta}_j$ which captures conditional variation in word usage as a function of our covariates w_i .

We described so far the joint distribution for the extensive margin of word usage. Now we describe the the intensive margin, which we model with a multinomial. Specifically the multinomial distribution describes the distribution of word repetitions. Define $c_{ij}^* = c_{ij} - h_{ij}$, the vector of

repetitions, $m_i^* = m_i - d_i$ the total amount of word repetitions in document i and $q_{ij}^* = \frac{h_{ij}q_{ij}}{\sum_{j=1}^d h_{ij}q_{ij}}$ is the share of word counts for word j in document i . Note that $q_{ij}^* = 0$ if $h_{ij} = 0$ since this distribution is conditional on the vocabulary \mathbf{h}_i . Then given the document vocabulary \mathbf{h}_i , the distribution of word counts is

$$p(c_i^* | \mathbf{v}_i, m_i^*, \mathbf{h}_i, d_i) = \frac{m_i^*!}{\prod_{j=1}^d c_{ij}^*!} \prod_{j=1}^d (q_{ij}^*)^{c_{ij}^*} = MN(\mathbf{c}_i^*; \mathbf{q}_i^*, m_i^*)$$

We now follow Taddy (2015) and factorize the joint distribution implied by several independent poisson in terms of a multinomial distribution for the words counts and a poisson distribution for the total word count, and get

$$MN(\mathbf{c}_i^*; \mathbf{q}_i^*, m_i^*) Po(m_i^*; \sum_{k=1}^d e^{\eta_{ik}}) = \prod_{j|h_{ij}=1} Po(c_{ij}^*; e^{\eta_{ij}}), \text{ with } \eta_{ij} = \alpha_j + \mu_i + \mathbf{v}_i' \boldsymbol{\varphi}_j,$$

where μ_i picks up a document specific tendency to repetition when they use any word, α_j picks up a word specific tendency to repetition when used in any document, and $\mathbf{v}_i' \boldsymbol{\varphi}_j$ picks up the conditional association between word repetition and the covariates \mathbf{v}_i .

Using Taddy (2015)'s proposed plug-in estimator for the poisson distribution of total word counts $\hat{\mu}_i = \log(m_i^*) = \log(m_i - d_i)$, we obtain that the multinomial distribution for word repetitions can be represented by the product of independent poisson,

$$MN(\mathbf{c}_i^*; \mathbf{q}_i^*, m_i^*) \approx \prod_{j|h_{ij}=1} Po(c_{ij}^*; m_i^* e^{\alpha_j + \mathbf{v}_i' \boldsymbol{\varphi}_j}).$$

All together this implies that the hurdle model that we bring to the data is

$$\prod_{j=1}^d B(h_{ij}, \Pi_{ij}) MN\left(c_{ij} - h_{ij}; \frac{h_{ij}q_{ij}}{\sum_{j=1}^d h_{ij}q_{ij}}, m_i - d_i\right) Po(m_i - d_i; \sum_{k=1}^d e^{\eta_{ik}}) \approx \prod_{j=1}^d \Pi_{ij}^{1-h_{ij}} \left(\Pi_{ij} Po(c_{ij} - h_{ij}; (m_i - d_i) e^{\alpha_j + \mathbf{v}_i' \boldsymbol{\varphi}_j})\right)^{h_{ij}},$$

$$\text{where } \Pi_{ij} = \frac{1}{1 + \frac{d-d_i}{d_i} e^{-\kappa_j - \mathbf{w}_i' \boldsymbol{\delta}_j}}.$$

C Sufficient Reduction Projection

We now show how these approximations allow us to recover the maximum likelihood estimates for the model parameters even if we separate the big estimation problem in multiple independent problems. Formally, under the conditions of the approximation described above we can factorize the log-likelihood of the joint word distribution as follows:

$$\begin{aligned}
 p(\mathbf{c}_i | m_i, d_i) &= \prod_j \left[1 / (e^{\gamma_i + \epsilon_{ij}} + 1) \right]^{1 - h_{ij}} \left[e^{\gamma_i + \epsilon_{ij}} / (e^{\gamma_i + \epsilon_{ij}} + 1) \text{Po}(c_{ij} - h_{ij}; (m_i - d_i) e^{\eta_{ij}}) \right]^{h_{ij}} \\
 &= \prod_j \frac{1}{(1 + e^{\gamma_i + \epsilon_{ij}})} \prod_j e^{(\gamma_i + \epsilon_{ij}) h_{ij}} \prod_j [\text{Po}(c_{ij} - h_{ij}; (m_i - d_i) e^{\eta_{ij}})]^{h_{ij}} \\
 &= \prod_j \frac{1}{(1 + e^{\log(\frac{d_i}{d-d_i}) + \kappa_j + \mathbf{w}'_i \delta_j})} \prod_j e^{(\log(\frac{d_i}{d-d_i}) + \kappa_j + \mathbf{w}'_i \delta_j) h_{ij}} \prod_j [\text{Po}(c_{ij} - h_{ij}; (m_i - d_i) e^{\eta_{ij}})]^{h_{ij}}.
 \end{aligned}$$

This has two components: the bernoulli component that determines each document vocabulary and the Poisson that determines each document use of this vocabulary. Focusing on the bernoulli component we can write,

$$\begin{aligned}
 \prod_j \frac{1}{(1 + e^{\log(\frac{d_i}{d-d_i}) + \kappa_j + \mathbf{w}'_i \delta_j})} \prod_j e^{(\log(\frac{d_i}{d-d_i}) + \kappa_j + \mathbf{w}'_i \delta_j) h_{ij}} &= \prod_j \frac{1}{(1 + e^{\log(\frac{d_i}{d-d_i}) + \kappa_j + \mathbf{w}'_i \delta_j})} \prod_j e^{(\log(\frac{d_i}{d-d_i}) + \kappa_j) h_{ij}} \prod_j e^{\mathbf{w}'_i \delta_j h_{ij}} \\
 &= a(\mathbf{w}_i, d_i) \psi(\mathbf{h}_i) \prod_{j=1}^d e^{\mathbf{w}'_i \delta_j h_{ij}}
 \end{aligned}$$

where $\psi(\mathbf{h}_i) = \prod_{j=1}^d e^{(\log(\frac{d_i}{d-d_i}) + \kappa_j) h_{ij}}$ and $a(\mathbf{w}_i, d_i) = \prod_{j=1}^d \left(1 + e^{\log(\frac{d_i}{d-d_i}) + \kappa_j + \mathbf{w}'_i \delta_j} \right)^{-1}$. We can similarly factorize the intensive margin component modeled with the poisson distribution as

$$\begin{aligned}
& \prod_j [Po(c_{ij} - h_{ij}; (m_i - d_i)e^{\eta_{ij}})]^{h_{ij}} \\
&= \prod_{j=1|h_{ij}=1}^d \frac{((m_i - d_i)e^{\eta_{ij}})^{c_{ij}-1} e^{-(m_i - d_i)e^{\eta_{ij}}}}{(c_{ij} - 1)!} \\
&= \prod_{j=1|h_{ij}=1}^d \exp((c_{ij} - 1) \log(m_i - d_i) + (c_{ij} - 1)\eta_{ij}) \exp(-(m_i - d_i)e^{\eta_{ij}}) \exp(-\log((c_{ij} - 1)!)) \\
&= \prod_{j=1|h_{ij}=1}^d \exp(-(m_i - d_i)e^{\alpha_j + \mathbf{v}_i' \boldsymbol{\varphi}_j}) \prod_{j=1|h_{ij}=1}^d \exp((c_{ij} - 1)(\log(m_i - d_i) + \alpha_j) - \log((c_{ij} - 1)!)) \prod_{j=1|h_{ij}=1}^d \exp((c_{ij} - 1)(\mathbf{v}_i' \boldsymbol{\varphi}_j)).
\end{aligned}$$

Therefore the joint distribution can be factored in as

$$p(\mathbf{c}_i | \mathbf{v}_i, \mathbf{w}_i) = \phi(\mathbf{c}_i) \psi(\mathbf{h}_i) a(\mathbf{w}_i, d_i) b(\mathbf{v}_i, m_i, d_i) \exp(\mathbf{w}_i' \boldsymbol{\delta} \mathbf{h}_i + \mathbf{v}_i' \boldsymbol{\varphi}(\mathbf{c}_i - \mathbf{h}_i)),$$

where

$$\begin{aligned}
b(\mathbf{v}_i, m_i, d_i) &= \prod_{j=1|h_{ij}=1}^d \exp(-(m_i - d_i)e^{\alpha_j + \mathbf{v}_i' \boldsymbol{\varphi}_j}) \\
\phi(\mathbf{c}_i) &= \prod_{j=1|h_{ij}=1}^d \exp((c_{ij} - 1)(\log(m_i - d_i) + \alpha_j) - \log((c_{ij} - 1)!)).
\end{aligned}$$

This expression is exactly the joint likelihood that our HDMR procedure maximizes.

References

- Adrian, Tobias, Erkki Etula, and Tyler Muir, 2014, Financial intermediaries and the cross-section of asset returns, *Journal of Finance* 69, 2557–2596.
- Antweiler, Werner, and Murray Z. Frank, 2004, Is all that talk just noise? The information content of Internet stock message boards, *Journal of Finance* 59, 1259–1293.
- Athey, Susan, 2017, Beyond prediction: Using big data for policy problems, *Science* 355, 483–485.
- , 2018, The impact of machine learning on economics, in *The Economics of Artificial Intelligence: An Agenda* (University of Chicago Press).
- , Guido Imbens, Thai Pham, and Stefan Wager, 2017, Estimating average treatment effects: Supplementary analyses and remaining challenges, *American Economic Review* 107, 278–81.
- Athey, Susan, Julie Tibshirani, and Stefan Wager, 2019, Generalized random forests, *The Annals of Statistics* 47, 1148–1178.
- Bajari, Patrick, Denis Nekipelov, Stephen P. Ryan, and Miaoyu Yang, 2015, Machine learning methods for demand estimation, *American Economic Review* 105, 481–85.
- Baker, Scott R, Nicholas Bloom, and Steven J Davis, 2016, Measuring economic policy uncertainty, *Quarterly Journal of Economics* 131, 1593–1636.
- Baron, Matthew, and Tyler Muir, 2018, Intermediaries and Asset Prices: Evidence from the U.S., U.K., and Japan, 1870-2016, Working paper.
- Bañbura, Marta, Domenico Giannone, Michele Modugno, and Lucrezia Reichlin, 2013, Now-Casting and the Real-Time Data Flow, in *Handbook of Economic Forecasting*, vol. 2 . pp. 195–237 (Elsevier).
- Belloni, Alexandre, Daniel Chen, Victor Chernozhukov, and Christian Hansen, 2012, Sparse models and methods for optimal instruments with an application to eminent domain, *Econometrica* 80, 2369–2429.
- Belloni, Alexandre, Victor Chernozhukov, Ivan Fernández-Val, and Christian Hansen, 2017, Program evaluation and causal inference with high-dimensional data, *Econometrica* 85, 233–298.
- Belloni, Alexandre, Victor Chernozhukov, and Christian Hansen, 2014, Inference on treatment effects after selection among high-dimensional controls, *The Review of Economic Studies* 81, 608–650.
- Benjamini, Yoav, and Yosef Hochberg, 1995, Controlling the false discovery rate: a practical and powerful approach to multiple testing, *Journal of the Royal statistical society: series B (Methodological)* 57, 289–300.
- Bitto, Angela, and Sylvia Frühwirth-Schnatter, 2018, Achieving shrinkage in a time-varying parameter model framework, *Journal of Econometrics*.
- Boot, Tom, and Didier Nibbering, 2019, Forecasting using random subspace methods, *Journal of Econometrics*.
- Boydston, Amber E., 2013, *Making the News: Politics, the Media, and Agenda Setting* (University of Chicago Press).
- Brunnermeier, Markus K., and Yuliy Sannikov, 2014, A macroeconomic model with a financial sector, *American Economic Review* 104, 379–421.
- Cepni, Oguzhan, I. Ethem Güney, and Norman R. Swanson, 2019, Nowcasting and forecasting GDP in emerging markets using global financial and macroeconomic diffusion indexes, *International Journal of Forecasting* 35, 555–572.

- Chebonenko, Tatiana, Lifeng Gu, and Dmitriy Muravyev, 2018, Text Sentiment's Ability to Capture Information: Evidence from Earnings Calls, SSRN Scholarly Paper ID 2352524 Social Science Research Network Rochester, NY.
- Chudik, Alexander, George Kapetanios, and M. Hashem Pesaran, 2018, A One Covariate at a Time, Multiple Testing Approach to Variable Selection in High-Dimensional Linear Regression Models, *Econometrica* 86, 1479–1512.
- Diebold, Francis X, and Robert S Mariano, 1995, Comparing Predictive Accuracy, *Journal of Business & Economic Statistics* 13, 134–144.
- Durante, Ruben, and Ekaterina Zhuravskaya, 2018, Attack When the World Is Not Watching? US News and the Israeli-Palestinian Conflict, *Journal of Political Economy* 126, 1085–1133.
- Einav, Liran, Amy Finkelstein, Sendhil Mullainathan, and Ziad Obermeyer, 2018, Predictive modeling of US health care spending in late life, *Science* 360, 1462–1465.
- Eisensee, Thomas, and David Stromberg, 2007, News droughts, news floods, and u.s. disaster relief, *Quarterly Journal of Economics* 122, 693–728.
- Engelberg, Joseph, and Christopher A. Parsons, 2010, The Causal Impact of Media in Financial Markets, *Journal of Finance*.
- Fang, Lily, and Joel Peress, 2009, Media Coverage and the Cross-section of Stock Returns, *Journal of Finance* 64, 2023–2052.
- Farmer, Leland, Lawrence Schmidt, and Allan Timmermann, 2019, Pockets of Predictability, Working Paper.
- Friedman, Jerome, Trevor Hastie, and Rob Tibshirani, 2010, Regularization paths for generalized linear models via coordinate descent, *Journal of statistical software* 33, 1.
- Gabaix, Xavier, 2014, A sparsity-based model of bounded rationality, *Quarterly Journal of Economics* 129, 1661–1710.
- Gallant, A Ronald, and Douglas W Nychka, 1987, Semi-nonparametric maximum likelihood estimation, *Econometrica* pp. 363–390.
- García, Diego, 2013, Sentiment during recessions, *Journal of Finance* 68, 1267–1300.
- Gentzkow, Matthew, Bryan T. Kelly, and Matt Taddy, 2019, Text as Data, *Journal of Economic Literature*.
- Gentzkow, Matthew, and Jesse M. Shapiro, 2006, Media bias and reputation, *Journal of Political Economy* 114, pp. 280–316.
- , and Matt Taddy, 2019, Measuring Group Differences in High-Dimensional Choices: Method and Application to Congressional Speech, *Econometrica* 87, 1307–1340.
- Giannone, Domenico, Lucrezia Reichlin, and David Small, 2008, Nowcasting: The real-time informational content of macroeconomic data, *Journal of Monetary Economics* 55, 665–676.
- Greene, William, 2007, Functional form and heterogeneity in models for count data, *Foundations and Trends in Econometrics* 1, 113–218.
- Gu, Shihao, Bryan T. Kelly, and Dacheng Xiu, 2018, Empirical Asset Pricing via Machine Learning, SSRN Scholarly Paper ID 3159577 Social Science Research Network Rochester, NY.
- , 2019, Autoencoder Asset Pricing Models, SSRN Scholarly Paper ID 3335536 Social Science Research Network Rochester, NY.

- Haddad, Valentin, and Tyler Muir, 2018, Do Intermediaries Matter for Aggregate Asset Prices?, Working paper.
- Hanley, Kathleen Weiss, and Gerard Hoberg, 2019, Dynamic Interpretation of Emerging Risks in the Financial Sector, *The Review of Financial Studies*.
- Hassan, Tarek A., Stephan Hollander, Laurence van Lent, and Ahmed Tahoun, 2017, Firm-Level Political Risk: Measurement and Effects, SSRN Scholarly Paper ID 2838644 Social Science Research Network Rochester, NY.
- He, Zhiguo, Bryan Kelly, and Asaf Manela, 2017, Intermediary asset pricing: New evidence from many asset classes, *Journal of Financial Economics* 126, 1–35.
- He, Zhiguo, and Arvind Krishnamurthy, 2012, A model of capital and crises, *The Review of Economic Studies* 79, 735–777.
- , 2013, Intermediary asset pricing, *American Economic Review* 103, 732–770.
- , 2018, Intermediary Asset Pricing and the Financial Crisis, *Annual Review of Financial Economics* 10, 173–197.
- Heckman, James J., 1979, Sample selection bias as a specification error, *Econometrica* 47, 153–161.
- Heckman, James J, and Thomas E MaCurdy, 1986, Labor econometrics, *Handbook of econometrics* 3, 1917–1977.
- Hoberg, Gerard, and Gordon Phillips, 2016, Text-Based Network Industries and Endogenous Product Differentiation, *Journal of Political Economy* 124, 1423–1465.
- Hodrick, Robert J, 1992, Dividend yields and expected stock returns: Alternative procedures for inference and measurement, *Review of Financial Studies* 5, 357–386.
- Jiang, Fuwei, Joshua Lee, Xiumin Martin, and Guofu Zhou, 2019, Manager sentiment and stock returns, *Journal of Financial Economics* 132, 126–149.
- Kelly, Bryan, and Seth Pruitt, 2013, Market expectations in the cross-section of present values, *Journal of Finance* 68, 1721–1756.
- , 2015, The three-pass regression filter: A new approach to forecasting using many predictors, *Journal of Econometrics* 186, 294–316.
- Kim, Hyun Hak, and Norman R. Swanson, 2014, Forecasting financial and macroeconomic variables using data reduction methods: New empirical evidence, *Journal of Econometrics* 178, 352–367.
- Kleinberg, Jon, Himabindu Lakkaraju, Jure Leskovec, Jens Ludwig, and Sendhil Mullainathan, 2018, Human Decisions and Machine Predictions, *The Quarterly Journal of Economics* 133, 237–293.
- Kleinberg, Jon, Jens Ludwig, Sendhil Mullainathan, and Ziad Obermeyer, 2015, Prediction policy problems, *American Economic Review* 105, 491–95.
- Koijen, Ralph S. J., and Motohiro Yogo, 2019, A Demand System Approach to Asset Pricing, *Journal of Political Economy* Status: forthcoming.
- Ludwig, Jens, Sendhil Mullainathan, and Jann Spiess, 2017, Machine learning tests for effects on multiple outcomes, *arXiv preprint arXiv:1707.01473*.
- Manela, Asaf, 2014, The value of diffusing information, *Journal of Financial Economics* 111, 181–199.
- , and Alan Moreira, 2017, News implied volatility and disaster concerns, *Journal of Financial Economics* 123, 137–162.

- McCullagh, Peter, and James A Nelder, 1989, *Generalized Linear Models* (Chapman & Hall).
- Michel, Jean-Baptiste, Yuan Kui Shen, Aviva Presser Aiden, Adrian Veres, Matthew K Gray, Joseph P Pickett, Dale Hoiberg, Dan Clancy, Peter Norvig, Jon Orwant, et al., 2011, Quantitative analysis of culture using millions of digitized books, *Science* 331, 176–182.
- Mullahy, John, 1986, Specification and testing of some modified count data models, *Journal of econometrics* 33, 341–365.
- Mullainathan, Sendhil, and Andrei Shleifer, 2005, The market for news, *American Economic Review* 95, 1031–1053.
- Mullainathan, Sendhil, and Jann Spiess, 2017, Machine learning: an applied econometric approach, *Journal of Economic Perspectives* 31, 87–106.
- Pastor, Lubos, and Pietro Veronesi, 2012, Uncertainty about government policy and stock prices, *Journal of Finance* 67, 1219–1264.
- Qin, Bei, David Strömberg, and Yanhui Wu, 2018, Media Bias in China, *American Economic Review* 108, 2442–2476.
- Rennie, Jason D., Lawrence Shih, Jaime Teevan, and David R. Karger, 2003, Tackling the poor assumptions of naive bayes text classifiers, in *Proceedings of the 20th international conference on machine learning (ICML-03)* pp. 616–623.
- Schervish, Mark J, 1995, *Theory of statistics* (Springer Science & Business Media).
- Stock, James H., and Mark Watson, 2011, Dynamic factor models, *Oxford handbook on economic forecasting*.
- Stock, James H., and Mark W. Watson, 2012, Generalized Shrinkage Methods for Forecasting Using Many Predictors, *Journal of Business & Economic Statistics* 30, 481–493.
- , 2016, Dynamic factor models, factor-augmented vector autoregressions, and structural vector autoregressions in macroeconomics, in *Handbook of macroeconomics*, vol. 2 . pp. 415–525 (Elsevier).
- Taddy, Matt, 2013a, Multinomial inverse regression for text analysis, *Journal of the American Statistical Association* 108, 755–770.
- , 2013b, Multinomial Inverse Regression for Text Analysis, *Journal of the American Statistical Association* 108, 755–770.
- , 2015, Distributed multinomial regression, *Annals of Applied Statistics* 9, 1394–1414.
- , 2017, One-step estimator paths for concave regularization, *Journal of Computational and Graphical Statistics* pp. 1–12.
- Tetlock, Paul C., 2007, Giving content to investor sentiment: The role of media in the stock market, *Journal of Finance* 62, 1139–1168.
- Tibshirani, Robert, 1996, Regression shrinkage and selection via the lasso, *Journal of the Royal Statistical Society. Series B (Methodological)* pp. 267–288.
- Vapnik, N. Vladimir, 2000, *The Nature of Statistical Learning Theory* (Springer-Verlag, New York.).
- Wilbur, W. John, and Won Kim, 2009, The ineffectiveness of within-document term frequency in text classification, *Information Retrieval* 12, 509–525.
- Zipf, George Kingsley, 1935, The psychology of language, *Houghton-Mifflin*.

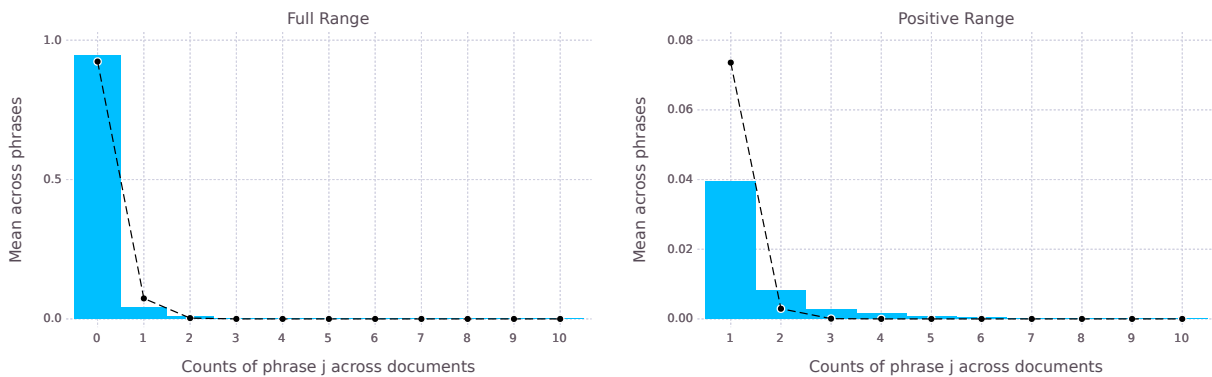


Figure 1: Mean distribution of WSJ front page articles monthly phrase counts

Notes: The figure shows the mean histogram for phrases that appear in the title or lead paragraph of front page *Wall Street Journal* articles, aggregated to form a monthly sample from July 1926 to February 2016. We construct the mean histogram by first calculating a histogram for each phrase across documents, and then averaging over phrases and normalizing to unit scale. The left panel shows that the entire range is highly sparse (has many zeros). The right panel omits zero counts, and shows that a poisson density fitted to the entire range (dashed line) is a poor description of the positive range. The corpus includes the 10,000 most frequent two-word phrases (bigrams) in separate sentences, after removing case, stopwords, nonletters and Porter stemming. Including less frequent phrases makes the corpus sparser and the pattern above more pronounced.

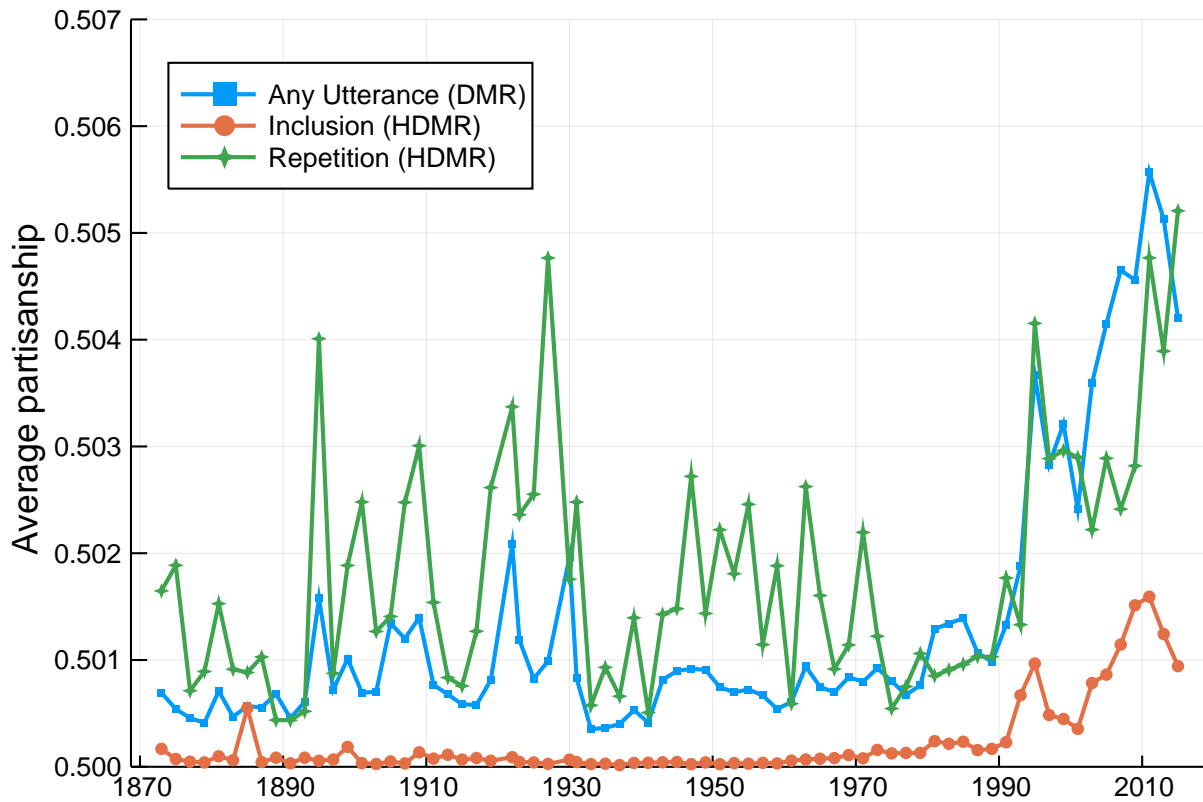


Figure 2: Partisanship of congressional speech

Notes: Average partisanship of US congressional speech is the posterior mass that an observer with a neutral prior expects to place on a speaker's true party after hearing the speaker utter a single phrase averaged over speakers in each congressional session (see Section 3). DMR-based partisanship is the preferred estimator of Gentzkow, Shapiro, and Taddy (2019). HDMR-based estimates decompose divergence between Democrats and Republicans into two: inclusion partisanship measures what an observer would expect to learn about a speaker's party affiliation from the inclusion of a phrase in a speech, while repetition partisanship measures what they expect to learn from an additional repetition of a phrase already observed. The congressional speech data is from Matthew Gentzkow's website.

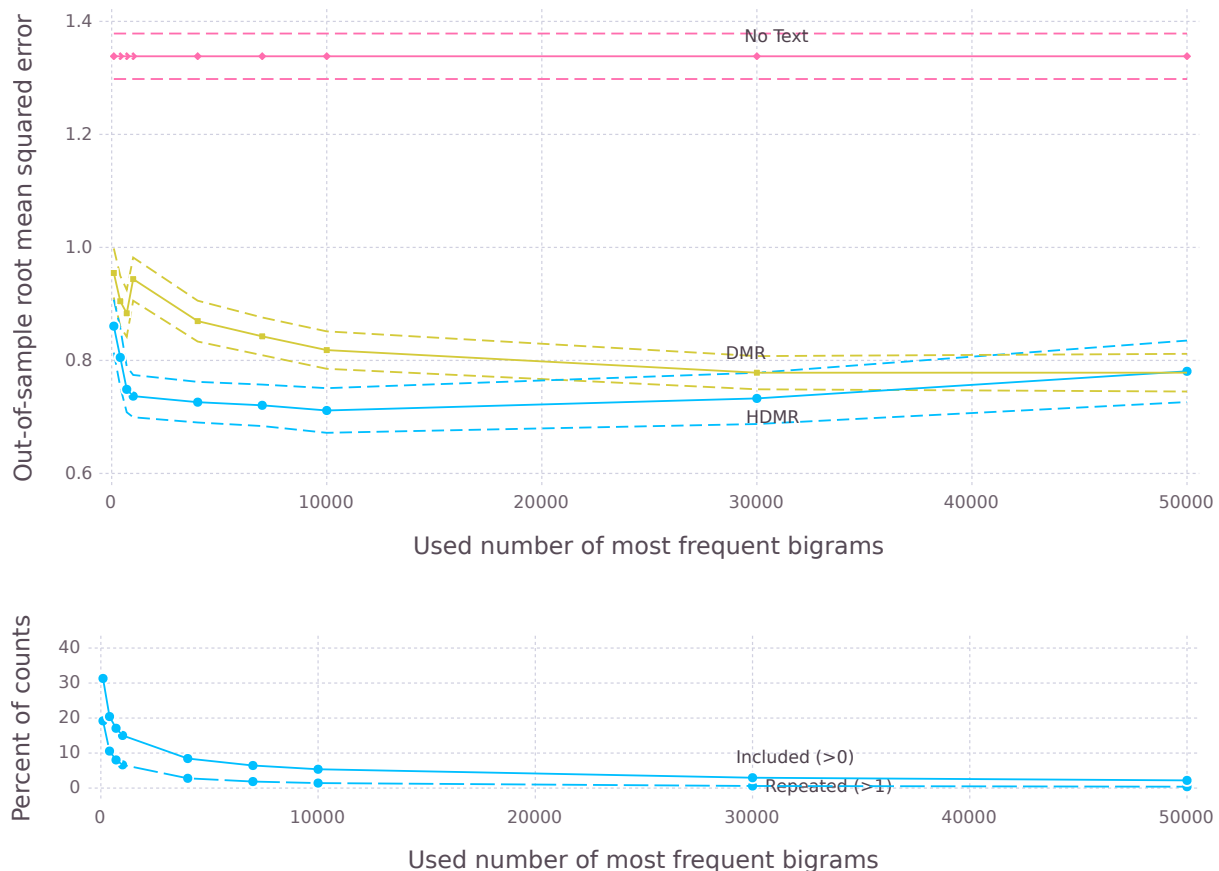


Figure 3: Predicting the intermediary capital ratio with text and covariates: 10-fold cross validation

Notes: The top panel reports out-of-sample root mean squared error from a 10-fold cross validation exercise that tries to predict the intermediary capital ratio (icr_t) using the log price dividend ratio (pd_t), realized variance of financial stocks ($rvfin$) over the same month, over the prior month, and over the prior year, and monthly WSJ front page phrase counts, over the subsample when the capital ratio is available, January 1970 to February 2016. Our proposed model, the hurdle distributed multinomial regression (HDMR) is compared with two benchmarks: (a) The distributed multinomial regression (DMR), which is provided with the same covariates and text, is a state-of-the-art approach to prediction with high-dimensional text, and (b) a linear regression of the target on the same covariates without the text (No Text). The figure shows how the advantage of HDMR in terms of out-of-sample fit changes as a function of the number of most frequent phrases included in the corpus. Dashed lines indicate the 95% confidence interval. The bottom panel shows how sparsity increases with this choice, i.e. it shows that average frequency that a word is used in a document.

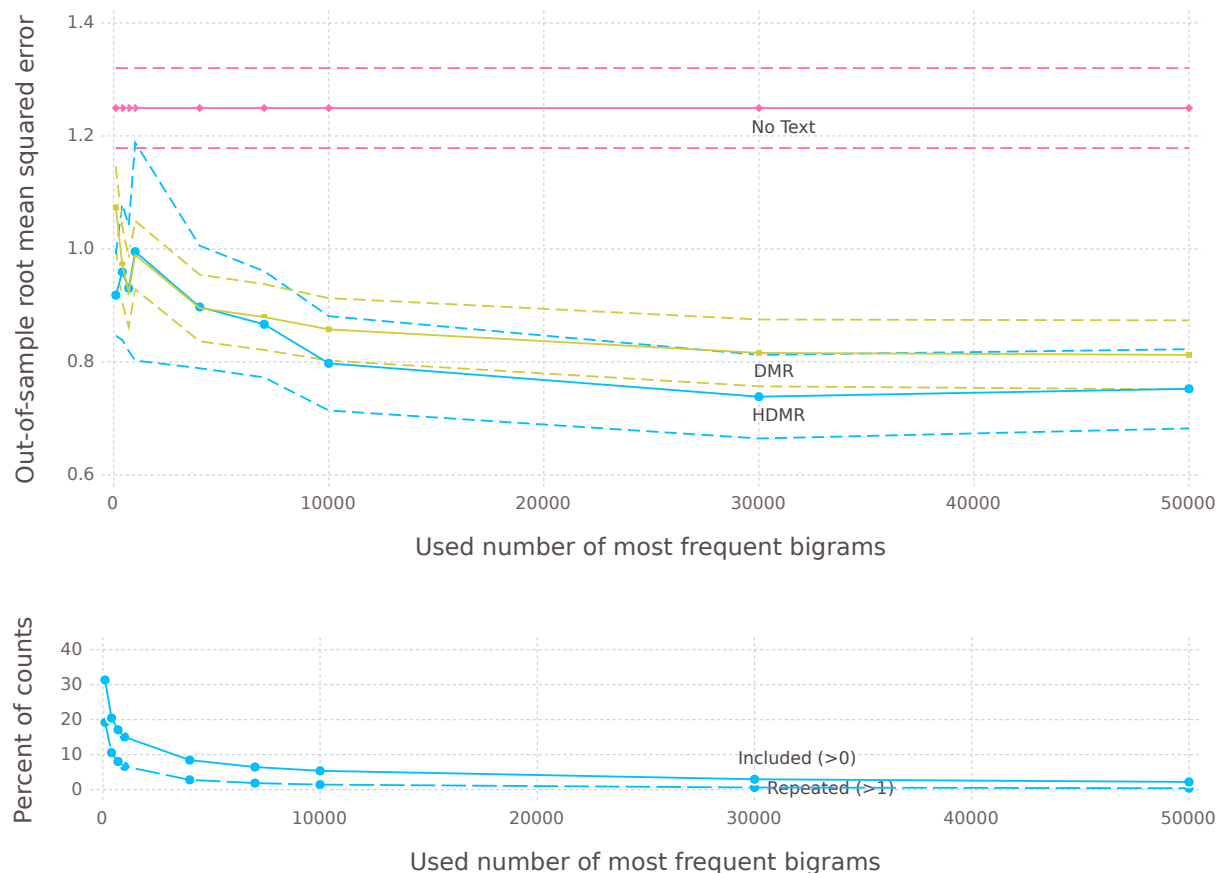
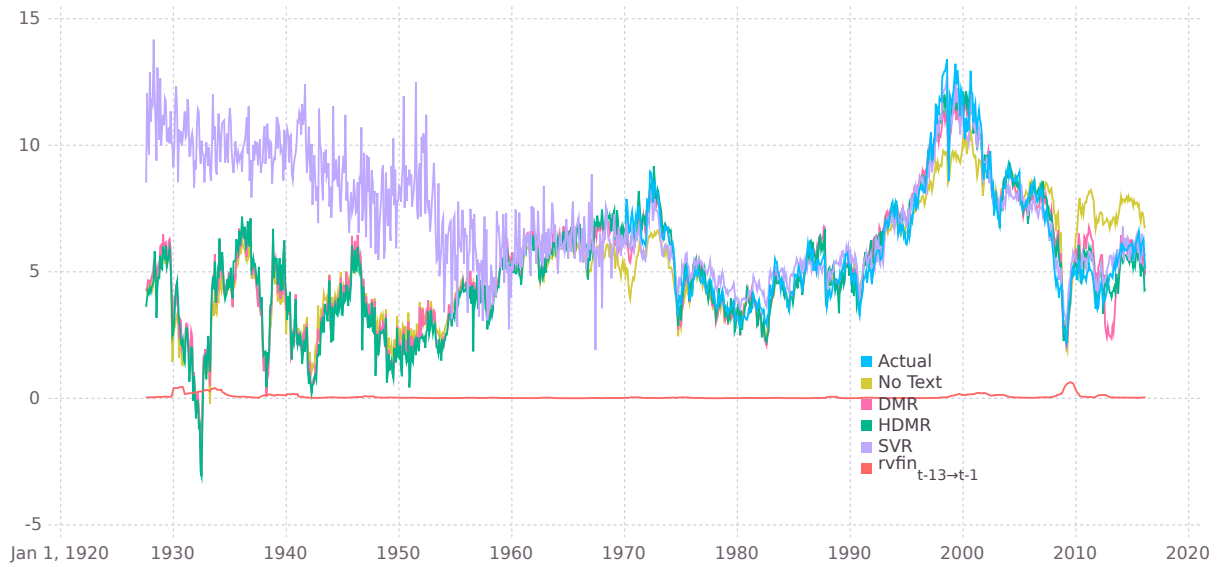
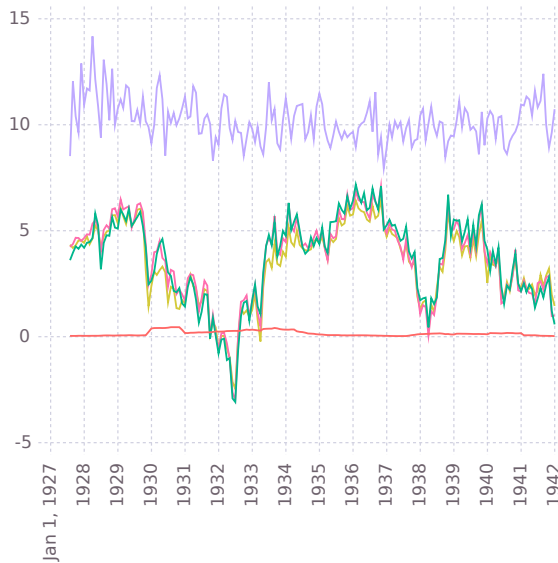


Figure 4: Predicting the intermediary capital ratio with text and covariates: Pseudo out-of-sample

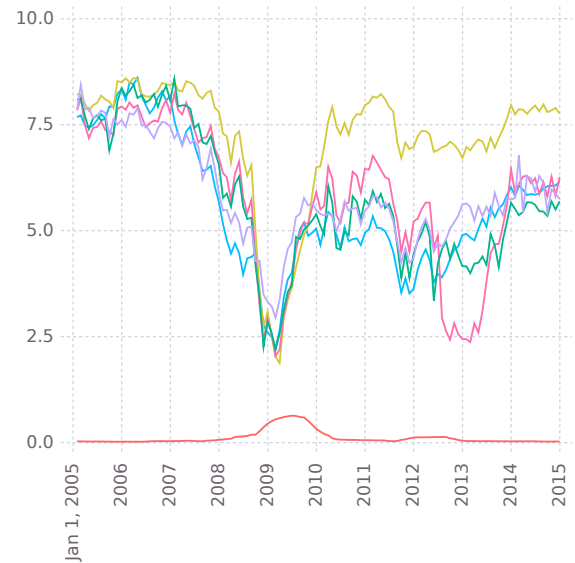
Notes: The top panel reports out-of-sample root mean squared error for predicting the intermediary capital ratio (icr_t) using the log price dividend ratio (pd_t), realized variance of financial stocks ($rvfin$) over the same month, over the prior month, and over the prior year, and monthly WSJ front page phrase counts, over the subsample when the capital ratio is available, January 1970 to February 2016. Unlike the random folds used before for validation, here we assess fit of a pseudo out-of-sample rolling back prediction exercise, starting with the later half of the sample, predicting one observation earlier, then extending the training sample by one earlier observation and rolling backward to assess fit in a backcasting exercise. Our proposed model, the hurdle distributed multinomial regression (HDMR) is compared with two benchmarks: (a) The distributed multinomial regression (DMR), which is provided with the same covariates and text, is a state-of-the-art approach to prediction with high-dimensional text, and (b) a linear regression of the target on the same covariates without the text (No Text). The figure shows how the advantage of HDMR in terms of out-of-sample fit changes as a function of the number of most frequent phrases included in the corpus. Dashed lines indicate the 95% confidence interval. The bottom panel shows how sparsity increases with this choice, i.e. it shows that average frequency that a word is used in a document.



(a) Full sample



(b) Great Depression



(c) Great Recession

Figure 5: Backcasting the intermediary capital ratio with text and covariates

Notes: The figure shows the predicted intermediary capital ratio (\widehat{icr}_t) using the log price dividend ratio (pd_t), realized variance of financial stocks ($rvfin$) over the same month, over the prior month, and over the prior year, and monthly WSJ front page phrase counts, over the extended sample, June 1927 to February 2016. The intermediary capital ratio is only available starting January 1970. Our proposed model, the hurdle distributed multinomial regression (HDMR), which excludes prior year financial stocks variance ($rvfin_{t-13 \rightarrow t-1}$, bottom line) from the repetition equation, is compared with three benchmarks: (a) distributed multinomial regression (DMR, Taddy, 2015), which is provided with the same covariates and text, (b) support vector regression (SVR), and (c) linear regression of the target on the same covariates without the text (No Text).

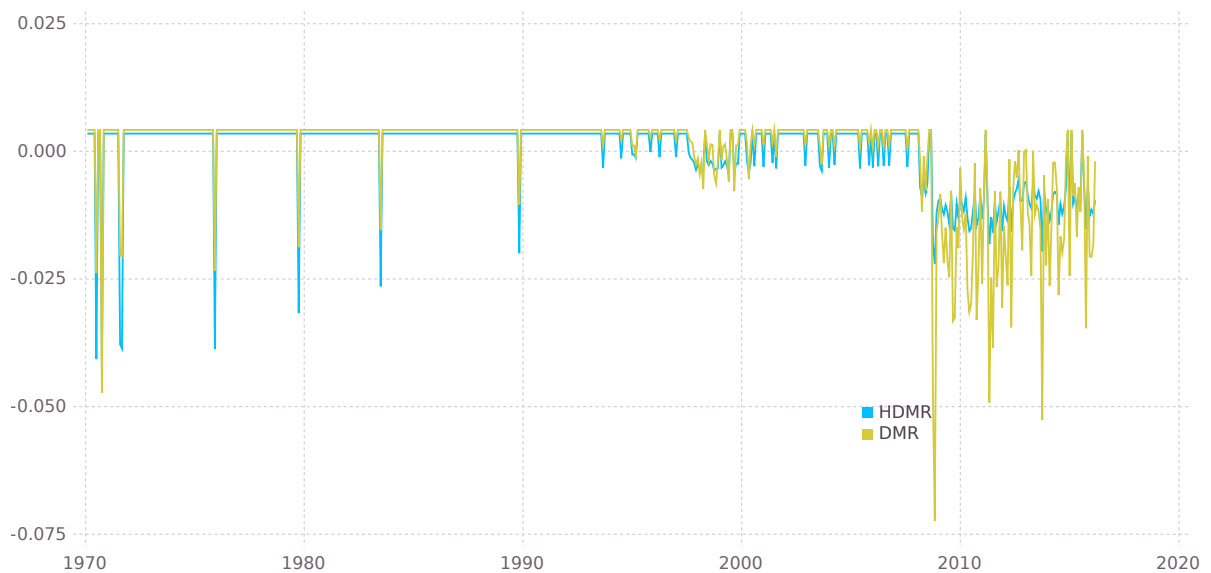


Figure 6: Focusing on the phrase “financial crisis” for intuition

Notes: The figure shows the predicted intermediary capital ratio (\widehat{icr}_t) due only to a single stemmed phrase, “financi crisi.” Our proposed model, the hurdle distributed multinomial regression (HDMR) gives more weight to the mere inclusion of this phrase on the front page of the *Wall Street Journal*, as opposed to its repeated use. Distributed multinomial regression (DMR) estimates, which does not break the variation into inclusion versus repetition, are shown for comparison.

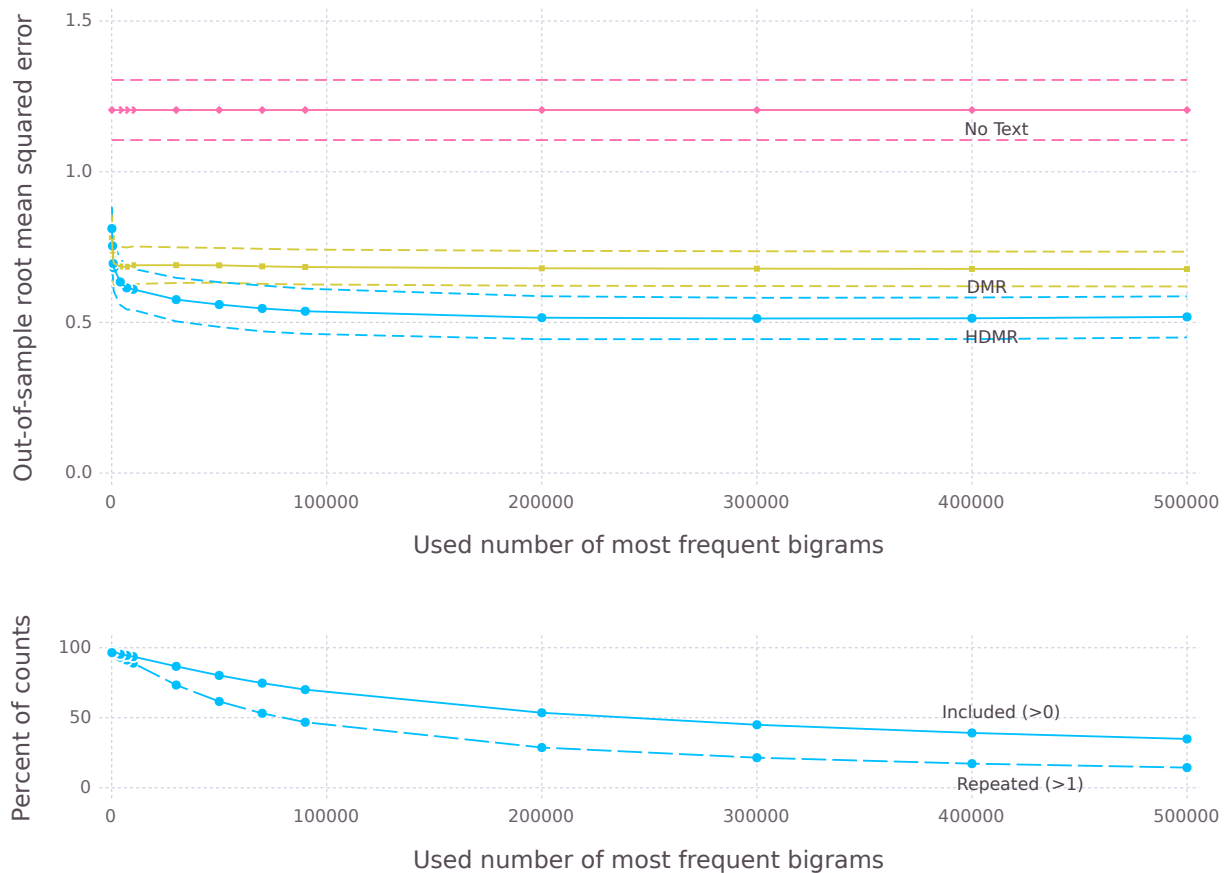


Figure 7: Predicting the intermediary capital ratio with denser text and covariates

Notes: The top panel reports out-of-sample root mean squared error from a 10-fold cross validation exercise that tries to predict the intermediary capital ratio (icr_t) using the log price dividend ratio (pd_t), realized variance of financial stocks ($rvfin$) over the same month, over the prior month, and over the prior year, and all monthly WSJ phrase counts, over the subsample when this text is available, January 1990 to December 2010. Our proposed model, the hurdle distributed multinomial regression (HDMR) is compared with two benchmarks: (a) The distributed multinomial regression (DMR), which is provided with the same covariates and text, is a state-of-the-art approach to prediction with high-dimensional text, and (b) a linear regression of the target on the same covariates without the text (No Text). The figure shows how the advantage of HDMR in terms of out-of-sample fit changes as a function of the number of most frequent phrases included in the corpus. Dashed lines indicate the 95% confidence interval. The bottom panel shows how sparsity increases with this choice, i.e. it shows that average frequency that a word is used in a document.

Table 1: Most partisan phrases in US congressional speech

Inclusion (HDMR)					
Republican	#R	#D	Democratic	#R	#D
subject overflow	1	0	tile farmer	2	3
cove creek	4	3	great kanawha	0	1
reapportion bill	2	2	proper locat	1	2
repres assign	2	1	notion agre	0	0
fertil made	1	0	leav discretionari	0	1
year ago	220	220	lie tri	0	1
assign construct	3	3	might furnish	0	1
fire trap	0	0	fix court	0	0
time manufactur	1	0	suppli san	0	0
auxiliari time	0	0	charact war	0	0

Repetition (HDMR)					
Republican	#R	#D	Democratic	#R	#D
postal servic	600	272	san francisco	679	1947
american peopl	1053	776	great britain	2117	2619
state govern	718	488	interior depart	597	975
pay tax	609	452	secretari interior	1756	1888
intern revenu	553	406	secretari navi	812	929
treasuri depart	998	882	build road	185	295
american citizen	832	717	depart interior	496	586
make differ	365	284	right way	272	293
foreign countri	649	583	bill becom	198	216
oppos bill	390	350	class peopl	172	184

Any Utterance (DMR)					
Republican	#R	#D	Democratic	#R	#D
muscl shoal	93	75	great kanawha	3	35
boulder dam	28	17	arrest drunken	5	35
navi yard	41	31	feder reserv	58	86
roman cathol	15	6	san francisco	23	48
regular armi	34	24	leagu nation	18	33
american peopl	64	55	reserv board	18	32
emerg offic	26	17	great britain	98	112
veteran bureau	37	28	help farmer	12	26
air corp	21	14	control product	6	18
allamerican canal	13	6	ship board	72	83

Inclusion (HDMR)					
Republican	#R	#D	Democratic	#R	#D
distinguish colleagu	53	52	year ago	437	437
high regard	5	5	control crime	2	2
draft card	2	2	sentinel system	0	1
safe street	4	4	member congress	247	248
headstart program	3	3	mani year	240	241
justic forta	1	1	american peopl	175	175
postmast general	6	6	situat must	1	2
nation foundat	1	1	feder bureaucrat	0	0
ambassador goldberg	1	1	charg involv	0	0
join effort	2	2	feder bureaucraci	1	2

Repetition (HDMR)					
Republican	#R	#D	Democratic	#R	#D
treasuri depart	635	494	american peopl	2680	2823
state depart	1198	1062	distinguish friend	384	525
intern revenu	641	542	state legislatur	480	537
support bill	610	516	western state	210	259
depart interior	659	578	san francisco	329	378
postal servic	619	544	bill right	244	285
public servic	799	727	foreign govern	194	230
depart state	413	364	take away	209	240
year age	284	245	secretari interior	773	800
urg passag	193	154	year ago	2977	3002

Any Utterance (DMR)					
Republican	#R	#D	Democratic	#R	#D
educ act	29	22	interest rate	44	54
higher educ	27	22	human right	25	35
vocat educ	15	10	feder reserv	30	40
latin america	30	25	soviet union	82	88
insert record	22	17	american peopl	91	96
job corp	23	19	farm bureau	10	16
distinguish colleagu	32	27	tax increas	37	42
even star	10	7	social secur	93	97
elementari secondari	23	20	public health	20	25
state depart	39	36	war poverti	21	25

Inclusion (HDMR)					
Republican	#R	#D	Democratic	#R	#D
war terror	38	30	war iraq	77	113
tax increas	15	9	civil war	43	58
continent shelf	10	6	dog coalit	1	13
american energi	7	3	african american	14	24
global war	9	5	fiscal respons	23	30
oil shale	5	1	iraq war	18	25
rais tax	11	7	civil right	21	27
nonbind resolut	6	2	support resolut	55	61
general petraeus	11	8	troop iraq	33	38
outer continent	10	6	spend iraq	1	7

Repetition (HDMR)					
Republican	#R	#D	Democratic	#R	#D
pay tax	583	384	american peopl	7253	7770
one thing	1404	1210	pay interest	262	468
year ago	3151	2966	men women	3361	3493
foreign countri	295	163	secretari interior	265	357
san francisco	411	321	nation govern	165	254
privat properti	319	235	interest nation	187	271
get rid	335	258	support bill	1411	1485
circuit court	647	569	postal servic	565	634
american citizen	479	407	peopl countri	543	609
everi year	706	636	step toward	371	429

Any Utterance (DMR)					
Republican	#R	#D	Democratic	#R	#D
tax increas	62	20	health care	237	292
natur gas	48	19	african american	31	63
rais tax	31	10	war iraq	50	79
side aisl	123	106	oil compani	40	64
american energi	18	3	american peopl	257	280
reserv balanc	123	109	nobid contract	2	25
continent shelf	21	8	urg colleagu	155	170
illeg immigr	19	7	paygo rule	10	24
tax rate	15	4	civil war	33	46
outer continent	19	8	children health	31	44

Inclusion (HDMR)					
Republican	#R	#D	Democratic	#R	#D
taxpay dollar	39	27	background check	39	65
american taxpay	20	14	depart homeland	43	60
sponsor terror	12	7	fund bill	21	33
death america	6	1	civil right	19	30
tax dollar	10	6	gun violenc	37	47
radic islam	5	1	african american	13	22
state sponsor	9	4	moment silenc	13	21
job creator	7	3	middl class	22	28
ballist missil	10	6	vote right	7	13
trillion debt	5	1	mass shoot	9	14

Repetition (HDMR)					
Republican	#R	#D	Democratic	#R	#D
american peopl	7100	5451	postal servic	966	1870
men women	2714	2097	year ago	2840	3130
privat properti	475	221	public servic	611	791
intern revenu	452	211	year old	770	925
state depart	1006	776	citizen unit	329	476
support bill	1318	1121	million peopl	910	1039
san francisco	535	384	take away	592	694
author act	609	472	take action	644	741
great state	672	536	honor repres	198	295
bill right	341	210	find way	453	541

Any Utterance (DMR)					
Republican	#R	#D	Democratic	#R	#D
american peopl	272	209	homeland secur	140	203
men women	105	81	climat chang	60	99
colleagu support	108	90	gun violenc	40	75
side aisl	116	98	vote right	26	59
human traffick	42	26	public health	50	80
taxpay dollar	35	19	depart homeland	67	98
radic islam	16	0	african american	41	70
religi freedom	19	4	civil right	32	55
balanc budget	24	10	puerto rico	56	80
nation defens	34	21	afford care	58	81

Notes: We report the ten most Republican and Democratic phrases in select sessions of the US congress. Phrase partisanship is based on congressional speeches and defined in Section 3.3. For each session, we separately sort on inclusion partisanship, on repetition partisanship, and on any utterance partisanship. We also report the predicted number of occurrences per 100,000 phrases for each phrase by Republicans (#R) and Democrats (#D).

Table 2: Summary statistics: Backcasting application

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
icr	6.235	2.399	2.230	3.616	5.548	9.578	13.400	557	197001–201605
pd	3.442	0.402	2.213	2.960	3.394	4.017	4.564	1075	192611–201605
$rvfin_{t-1 \rightarrow t}$	0.061	0.144	0.002	0.006	0.022	0.133	2.059	1079	192607–201605
$rvfin_{t-12 \rightarrow t}$	0.061	0.094	0.004	0.010	0.026	0.159	0.636	1068	192706–201605

Notes: Reported are summary statistics for the *Wall Street Journal* front page articles text and for variables in the monthly sample from July 1926 to May 2016. The corpus includes the 10,000 most frequent two-word phrases (bigrams) in separate sentences, after removing case, stopwords, nonletters and Porter stemming. To summarize the text, we report the mean of per phrase statistics, for counts c_{tj} and for inclusion indicators, $h_{tj} \equiv 1_{\{c_{tj} > 0\}}$. Intermediary capital ratio (icr_t) is the aggregate ratio of market equity to market equity plus book debt of New York Fed primary dealers in percentage terms. The log price over past year dividends (pd_t) is from CRSP. Realized variance ($rvfin$) is the annualized daily variance of financial stock returns over the prior month ($t - 1 \rightarrow t$) or year ($t - 12 \rightarrow t$).

Table 3: Predicting the intermediary capital ratio with text and covariates

(a) Cross-validation with 10 random folds

Model	Out-of-sample			In-sample		
	Text	No Text	Difference	Text	No Text	Difference
HDMR	0.711 (0.020)	1.338 (0.021)	-0.627 (0.019)	0.557 (0.002)	1.322 (0.002)	-0.765 (0.002)
DMR	0.818 (0.017)	1.338 (0.021)	-0.520 (0.022)	0.717 (0.003)	1.322 (0.002)	-0.606 (0.004)
FHDMR	0.789 (0.016)	1.338 (0.021)	-0.549 (0.020)	0.674 (0.002)	1.322 (0.002)	-0.648 (0.003)
SVR	1.256 (0.042)	1.334 (0.020)	-0.078 (0.046)	0.641 (0.005)	1.326 (0.002)	-0.685 (0.005)

(b) Pseudo out-of-sample rolling back

Model	Out-of-sample			In-sample		
	Text	No Text	Difference	Text	No Text	Difference
HDMR	0.798 (0.043)	1.249 (0.036)	-0.452 (0.033)	0.580 (0.001)	1.341 (0.003)	-0.760 (0.002)
DMR	0.858 (0.028)	1.249 (0.036)	-0.392 (0.022)	0.771 (0.003)	1.341 (0.003)	-0.570 (0.002)
FHDMR	0.828 (0.028)	1.249 (0.036)	-0.422 (0.024)	0.733 (0.003)	1.341 (0.003)	-0.608 (0.002)
SVR	1.911 (0.055)	0.956 (0.042)	0.955 (0.065)	0.633 (0.001)	1.358 (0.003)	-0.725 (0.003)

Notes: Reported is in- and out-of-sample root mean squared error (RMSE) for predicting the intermediary capital ratio (icr_t) using the log price dividend ratio (pd_t), realized variance of financial stocks ($rvfin$) over the same month, over the prior month, and over the prior year. Models with text additionally include monthly WSJ front page phrase counts, over the subsample when the capital ratio is available, January 1970 to February 2016. The corpus includes the 10,000 most frequent two-word phrases (bigrams) in separate sentences, after removing case, stopwords, nonletters and Porter stemming. Panel (a) uses 10 random folds for validation while Panel (b) uses pseudo out-of-sample prediction, starting with the latter half of the sample and rolling back one observation at a time. Our proposed model, the hurdle distributed multinomial regression (HDMR), which excludes prior year $rvfin$ from the repetition equation, is compared with three benchmarks: (a) distributed multinomial regression (DMR, Taddy, 2015), which is provided with the same covariates and text, (b) a “fabricated” variant of HDMR which adds $h_{ij} = \mathbf{1}(c_{ij} > 0)$ indicators to the text counts matrix \mathbf{c} and then runs DMR (FHDMR), and (c) support vector regression (SVR). For each model we report RMSE with and without the text, the change in the measure of fit. Standard errors are in parentheses.

Table 4: Pivotal phrases for predicting the intermediary capital ratio out-of-sample

(a) Cross-validation with 10 random folds

Phrase	Δ OOS RMSE	δ	φ	Mean counts	Mean inclusions
busi bulletin	0.006	2.333	0.195	4.863	0.602
barack obama	0.006	-7.054	-0.171	1.224	0.195
tax report	0.006	1.945	0.254	4.929	0.604
euro zone	0.005	-1.170	-1.689	0.508	0.132
washington wire	0.005	1.456	0.184	5.430	0.618
labor letter	0.005	2.090	0.177	3.273	0.535
penn central	0.005	1.375	2.024	0.063	0.042
feder reserv	0.004	-0.443	-0.434	1.629	0.682
jobless marri	0.004	1.583	0.460	0.291	0.260
presid nixon	0.003	1.209	0.528	0.291	0.123
instal credit	0.003	2.514	0.000	0.020	0.020
mitt romney	0.003	-1.477	-1.662	0.221	0.056
iran contra	0.002	0.321	-2.557	0.067	0.042
nixon administr	0.002	1.428	0.788	0.087	0.065
aluminum output	0.002	1.328	0.753	0.206	0.177
credit trend	0.002	3.177	0.000	0.016	0.016
dow industri	0.002	-1.221	-0.385	0.864	0.237
al qaeda	0.002	-1.190	-0.231	1.665	0.264
hedg fund	0.002	-1.303	-0.308	1.186	0.302
lead indic	0.002	1.108	-0.163	0.492	0.391

(b) Pseudo out-of-sample rolling back

Phrase	Δ OOS RMSE	δ	φ	Mean counts	Mean inclusions
tax report	0.094	1.945	0.254	4.929	0.604
busi bulletin	0.084	2.333	0.195	4.863	0.602
labor letter	0.071	2.090	0.177	3.273	0.535
washington wire	0.065	1.456	0.184	5.430	0.618
steel product	0.018	0.923	0.567	0.577	0.512
factori shipment	0.018	0.795	0.036	0.629	0.530
and unfil	0.018	-0.188	0.000	0.219	0.215
hour earn	0.017	0.986	0.000	0.570	0.524
week earn	0.016	0.912	0.000	0.544	0.499
jobless marri	0.015	1.583	0.460	0.291	0.260
lead indic	0.015	1.108	-0.163	0.492	0.391
presid nixon	0.014	1.209	0.528	0.291	0.123
los angel	0.011	0.361	0.255	1.931	0.714
and paperboard	0.011	0.000	0.000	0.174	0.152
paper and	0.011	0.000	0.000	0.174	0.152
and backlog	0.010	-0.776	0.000	0.206	0.206
inventori and	0.009	-0.916	0.000	0.203	0.203
san francisco	0.009	0.364	0.162	1.430	0.638
construct spend	0.008	0.367	-0.598	0.354	0.351
person incom	0.008	0.584	0.294	0.544	0.447

Notes: The table reports the top 20 phrases whose removal from the corpus causes the largest deterioration in out-of-sample root-mean-squared-error (Δ OOS RMSE), when predicting the intermediary capital ratio (icr_t) using text and the log price dividend ratio (pd_t), realized variance of financial stocks ($rvfin$) over the same month, over the prior month, and over the prior year. The text includes monthly WSJ front page phrase counts, over the subsample when the capital ratio is available, January 1970 to February 2016. Panel (a) uses 10 random folds for validation while Panel (b) uses pseudo out-of-sample prediction, starting with the latter half of the sample and rolling back one observation at a time. We also report full sample HDMR coefficients on icr_t for phrase inclusion (δ) and repetition (φ), and mean counts and inclusions across observations. The corpus includes the 10,000 most frequent two-word phrases (bigrams) in separate sentences, after removing case, stopwords, nonletters and Porter stemming.

Table 5: Focusing on the phrase “financial crisis” for intuition

(a) Backward regressions			
	HDMR		DMR
	Repetition	Inclusion	
Intercept	-10.78	-9.62	-16.69
icr_t	-0.43	-0.60	-0.61
pd_t	2.05	4.02	3.50
$rvfin_{t-1 \rightarrow t}$	1.41	0.90	1.44
$rvfin_{t-2 \rightarrow t-1}$	-0.41	1.00	-0.54
$rvfin_{t-13 \rightarrow t-1}$		2.66	1.26
(b) Forward regressions			
	HDMR		DMR
Repetition	-1.07		-4.71
Inclusion	-5.33		

Notes: Panel (a) reports backward HDMR coefficient estimates for the (stemmed) phrase “financial crisis” on the covariates, which excludes prior year financial stocks volatility ($rvfin$) from the repetition equation. Panel (b) reports the forward regression coefficient products with those of the backward regression, $b_z\varphi_{jy}$ and $b_s\delta_{jy}$ for HDMR, and contrasts it with the corresponding single coefficient product of DMR. The hurdle distributed multinomial regression (HDMR) gives more weight to the mere inclusion of this phrase on the front page of the *Wall Street Journal*, as opposed to its repeated use. Distributed multinomial regression (DMR) estimates, which do not break the variation into inclusion versus repetition, are shown for comparison. The corpus includes the 10,000 most frequent two-word phrases (bigrams) in separate sentences, after removing case, stopwords, nonletters and Porter stemming.

Table 6: Time-varying risk premia and the news-implied intermediary capital ratio

(a) Postwar sample, 194601–201602

	$r_{t \rightarrow t+1}^{em}$				$r_{t \rightarrow t+3}^{em}$				$r_{t \rightarrow t+12}^{em}$			
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)	(11)	(12)
icr_t	-3.10 (2.34)				-3.33 (2.57)				-2.94 (2.49)			
\widehat{icr}_t		-4.46 (1.77)				-4.61 (1.83)				-4.96 (1.95)		
$\widehat{icr}_t^{No\ Text}$			-3.99 (1.77)				-4.23 (1.80)				-4.37 (1.87)	
z_t^0				5.86 (2.29)				5.96 (1.71)				5.81 (1.53)
z_t^+				-3.38 (2.05)				-3.22 (2.01)				-3.46 (1.92)
pd_t				-9.11 (2.41)				-9.50 (2.31)				-9.62 (2.23)
$rvfin_{t-1 \rightarrow t}$				-9.87 (2.79)				-4.61 (2.38)				-1.05 (0.78)
$rvfin_{t-2 \rightarrow t-1}$				4.42 (3.01)				0.57 (2.16)				0.73 (0.75)
$rvfin_{t-13 \rightarrow t-1}$				3.79 (2.43)				3.60 (1.80)				1.24 (1.71)
N	552	841	841	841	552	841	841	841	544	833	833	833
Adjusted R^2 , %	0.14	0.63	0.48	2.69	0.84	2.11	1.76	5.97	3.02	9.60	7.36	20.32

(b) Full sample, 192707–201602

	$r_{t \rightarrow t+1}^{em}$				$r_{t \rightarrow t+3}^{em}$				$r_{t \rightarrow t+12}^{em}$			
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)	(11)	(12)
icr_t	-3.10 (2.34)				-3.33 (2.57)				-2.94 (2.49)			
\widehat{icr}_t		-4.25 (1.99)				-4.48 (3.00)				-4.88 (2.43)		
$\widehat{icr}_t^{No\ Text}$			-3.68 (1.99)				-4.04 (3.25)				-4.13 (2.59)	
z_t^0				2.12 (2.36)				1.91 (3.88)				3.03 (2.76)
z_t^+				-3.27 (2.28)				-3.51 (1.90)				-3.97 (1.81)
pd_t				-5.74 (2.21)				-6.03 (3.95)				-6.56 (2.78)
$rvfin_{t-1 \rightarrow t}$				-4.17 (2.60)				-1.31 (2.59)				-0.98 (1.29)
$rvfin_{t-2 \rightarrow t-1}$				4.24 (2.78)				0.08 (2.45)				-0.15 (1.31)
$rvfin_{t-13 \rightarrow t-1}$				-3.45 (2.67)				-2.16 (3.53)				-2.89 (3.51)
N	552	1,062	1,062	1,061	552	1,062	1,062	1,061	544	1,054	1,054	1,053
Adjusted R^2 , %	0.14	0.33	0.23	0.46	0.84	1.16	0.93	1.58	3.02	5.73	4.05	9.81

Notes: Reported are time-series predictability regression estimates of future stock market excess returns ($r_{t \rightarrow t+\tau}^{em}$) at one to twelve month horizons on the intermediary capital ratio (icr_t) in the short sample in which it is available, on the news-implied \widehat{icr}_t that is available for a much longer time-series, or on $\widehat{icr}_t^{No\ Text}$ that conditions only on the non-text covariates. We additionally decompose the \widehat{icr}_t into the sufficient reduction projections z_t^0 , z_t^+ that summarize the text, the log price dividend ratio (pd_t), realized variance of financial stocks ($rvfin$) over the same month, over the prior month, and over the prior year. Explanatory variables are demeaned and scaled to unit variance. The corpus includes the 10,000 most frequent two-word phrases (bigrams) in separate sentences, after removing case, stopwords, nonletters and Porter stemming. Hodrick (1992) standard errors are in parentheses.

Table 7: Summary statistics: Macroeconomic forecasting dataset

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	0.873	3.225	-16.153	-2.760	1.081	4.423	12.172	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	7.306	0.236	6.190	6.990	7.326	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	1.089	-5.401	-1.187	-0.003	1.099	7.175	602	195902–200903
Real AHE: goods	0.295	1.318	-4.700	-1.187	0.272	1.688	6.274	602	195901–200902
FedFunds	0.002	0.371	-1.560	-0.420	0.010	0.380	1.600	602	195901–200902
M1	0.012	2.188	-10.512	-2.347	-0.000	2.414	7.485	601	195902–200902
Ex rate: avg	-0.172	6.017	-21.276	-8.339	0.038	6.961	21.450	601	195901–200901
S&P 500	1.807	14.500	-91.237	-14.479	2.821	17.014	45.411	603	195901–200903
Consumer expect	-0.053	3.975	-16.500	-4.600	-0.200	4.600	22.500	603	195901–200903

Notes: Reported are summary statistics for the *Wall Street Journal* full text and for target variables in the monthly sample of [Stock and Watson \(2012\)](#). The corpus includes the 100,000 most frequent two-word phrases (bigrams) in separate sentences, after removing case, stopwords, nonletters and Porter stemming. To summarize the text, we report the mean of per phrase statistics, for counts c_{tj} and for inclusion indicators, $h_{tj} \equiv 1_{\{c_{tj} > 0\}}$. For each of the 12 categories of variables we use the headline variable as the prediction target. The transformed series are generally first differences of logarithms (growth rates) for real quantity variables, first differences for nominal interest rates, and second differences of logarithms (changes in rates of inflation) for price series.

Table 8: Forecasting macroeconomic series

(a) *Wall Street Journal* full text, 10,000 most frequent phrases

Months forward:	$\tau = 1$				$\tau = 3$				$\tau = 12$			
Folds:	Random		Rolling		Random		Rolling		Random		Rolling	
$Y_{t+\tau}^\tau \setminus$ Benchmark:	DFM-5	DMR	DFM-5	DMR	DFM-5	DMR	DFM-5	DMR	DFM-5	DMR	DFM-5	DMR
IP: total	0.986 (0.680)	1.027 (1.070)	0.956 (0.588)	0.966 (0.208)	0.857 (0.010)	0.961 (0.187)	0.826 (0.000)	0.955 (0.012)	0.697 (0.000)	0.880 (0.000)	0.671 (0.000)	0.855 (0.000)
Emp: total	0.829 (0.001)	0.965 (0.498)	0.825 (0.003)	0.876 (0.000)	0.763 (0.000)	0.917 (0.000)	0.755 (0.000)	0.848 (0.000)	0.682 (0.000)	0.876 (0.013)	0.642 (0.000)	0.782 (0.000)
U: all	0.953 (0.429)	1.005 (1.150)	1.000 (1.490)	0.995 (0.682)	0.805 (0.000)	0.955 (0.164)	0.733 (0.000)	0.910 (0.001)	0.644 (0.000)	0.852 (0.017)	0.619 (0.000)	0.823 (0.000)
HStarts: Total	0.675 (0.000)	1.034 (0.931)	0.690 (0.000)	0.901 (0.000)	0.631 (0.000)	0.993 (0.614)	0.653 (0.000)	0.845 (0.000)	0.622 (0.000)	1.089 (1.089)	0.771 (0.000)	1.015 (0.928)
PMI	0.836 (0.000)	1.026 (0.972)	1.016 (1.459)	1.030 (0.935)	0.854 (0.000)	1.019 (0.849)	1.011 (0.860)	1.093 (1.089)	0.714 (0.000)	0.951 (0.023)	0.820 (0.005)	0.891 (0.000)
CPI-ALL	1.088 (1.333)	0.995 (1.063)	1.139 (1.200)	0.961 (0.230)	1.064 (1.000)	0.977 (0.494)	1.105 (1.091)	1.047 (1.223)	1.028 (0.957)	1.008 (0.922)	1.102 (0.994)	1.104 (0.999)
Real AHE: goods	0.968 (0.463)	0.972 (0.854)	1.021 (1.369)	1.008 (0.796)	0.889 (0.012)	0.964 (0.522)	0.990 (0.828)	0.959 (0.189)	0.614 (0.000)	0.982 (0.414)	0.845 (0.012)	0.939 (0.000)
FedFunds	0.913 (0.001)	0.995 (1.041)	1.042 (1.331)	0.986 (0.551)	0.814 (0.000)	0.976 (0.420)	1.042 (0.991)	1.029 (1.221)	0.647 (0.000)	0.972 (0.371)	0.864 (0.026)	0.938 (0.090)
M1	1.096 (1.091)	1.014 (0.928)	1.182 (1.000)	1.009 (0.786)	1.114 (1.089)	1.019 (0.913)	1.260 (1.000)	1.111 (1.000)	1.008 (0.666)	1.008 (0.972)	1.071 (1.026)	0.962 (0.235)
Ex rate: avg	1.066 (1.200)	1.008 (1.175)	1.094 (1.283)	0.927 (0.023)	1.009 (0.819)	1.038 (0.969)	0.999 (0.852)	0.923 (0.030)	0.889 (0.002)	1.064 (0.999)	0.848 (0.018)	0.981 (0.315)
S&P 500	1.029 (1.218)	1.007 (1.195)	1.069 (1.394)	1.064 (0.989)	0.937 (0.113)	0.995 (0.546)	0.938 (0.287)	1.047 (1.159)	0.750 (0.000)	0.955 (0.027)	0.690 (0.000)	0.937 (0.007)
Consumer expect	1.098 (1.000)	0.946 (0.073)	1.240 (1.091)	0.833 (0.000)	1.032 (1.021)	1.009 (0.839)	1.115 (1.193)	0.975 (0.414)	0.951 (0.152)	1.037 (1.176)	0.865 (0.020)	1.017 (0.858)
Fraction < 1	0.583	0.417	0.333	0.667	0.667	0.667	0.583	0.583	0.833	0.583	0.833	0.750
Fraction < 1 at 95%	0.333	0.000	0.167	0.333	0.583	0.083	0.333	0.417	0.750	0.417	0.833	0.500

Notes: Reported are out-of-sample RMSE for HDMR-based τ -month ahead forecasts that use the lagged text plus the top 5 PCs of the [Stock and Watson \(2012\)](#) monthly macroeconomic series, relative to DFM-5, which uses only these PCs without the text, and relative to DMR with the same data. The transformed macroeconomic series used as prediction targets and for calculation of the PCs are generally first differences of logarithms (growth rates) for real quantity variables, first differences for nominal interest rates, and second differences of logarithms (changes in rates of inflation) for price series. [Diebold and Mariano \(1995\)](#) p -values testing the null hypothesis that the RMSE of HDMR is larger than the benchmark's, and corrected for multiple testing using the [Benjamini and Hochberg \(1995\)](#) procedure, are in parentheses.

Table 8: Forecasting macroeconomic series

(b) *Wall Street Journal* full text, 100,000 most frequent phrases

Months forward:	$\tau = 1$				$\tau = 3$				$\tau = 12$			
Folds:	Random		Rolling		Random		Rolling		Random		Rolling	
$Y_{t+\tau}^\tau \setminus$ Benchmark:	DFM-5	DMR	DFM-5	DMR	DFM-5	DMR	DFM-5	DMR	DFM-5	DMR	DFM-5	DMR
IP: total	0.981 (0.743)	1.020 (0.949)	0.981 (0.916)	0.981 (0.401)	0.833 (0.000)	0.942 (0.044)	0.843 (0.000)	0.985 (0.356)	0.606 (0.000)	0.783 (0.000)	0.638 (0.000)	0.829 (0.000)
Emp: total	0.754 (0.000)	0.905 (0.030)	0.875 (0.019)	0.951 (0.171)	0.670 (0.000)	0.834 (0.000)	0.771 (0.000)	0.881 (0.000)	0.494 (0.000)	0.670 (0.000)	0.560 (0.000)	0.716 (0.000)
U: all	0.941 (0.265)	0.994 (0.646)	0.976 (0.746)	0.982 (0.460)	0.790 (0.000)	0.959 (0.344)	0.866 (0.000)	1.096 (0.989)	0.541 (0.000)	0.766 (0.001)	0.654 (0.000)	0.898 (0.000)
HStarts: Total	0.545 (0.000)	0.842 (0.008)	0.520 (0.000)	0.679 (0.000)	0.510 (0.000)	0.817 (0.008)	0.511 (0.000)	0.671 (0.000)	0.495 (0.000)	0.883 (0.029)	0.515 (0.000)	0.679 (0.000)
PMI	0.804 (0.000)	1.038 (1.006)	0.960 (1.001)	0.962 (0.187)	0.836 (0.000)	1.062 (1.082)	0.931 (0.218)	0.996 (0.667)	1.102 (0.000)	0.797 (1.183)	0.908 (0.053)	0.975 (0.326)
CPI-ALL	1.071 (1.332)	0.931 (0.001)	1.059 (1.199)	0.851 (0.000)	1.069 (1.090)	0.946 (0.006)	1.032 (1.126)	0.934 (0.038)	1.072 (0.998)	1.033 (1.271)	1.053 (1.063)	1.021 (0.842)
Real AHE: goods	1.001 (0.891)	0.989 (0.580)	1.073 (1.656)	1.045 (1.066)	0.977 (0.342)	1.079 (1.184)	1.019 (1.019)	1.019 (0.906)	0.631 (0.000)	1.063 (1.091)	0.808 (0.002)	0.912 (0.000)
FedFunds	0.950 (0.103)	1.049 (1.099)	1.158 (1.091)	1.115 (0.990)	0.873 (0.000)	1.086 (1.311)	0.988 (0.684)	1.014 (0.883)	0.724 (0.000)	1.202 (1.000)	0.707 (0.000)	0.835 (0.000)
M1	1.039 (1.091)	0.887 (0.000)	1.089 (1.000)	0.775 (0.000)	1.026 (1.191)	0.885 (0.000)	1.127 (1.000)	0.848 (0.000)	1.039 (0.991)	1.034 (1.202)	1.074 (0.992)	0.945 (0.116)
Ex rate: avg	1.064 (1.000)	0.986 (0.498)	1.089 (1.476)	0.853 (0.001)	1.076 (1.000)	1.111 (1.000)	1.065 (1.207)	0.938 (0.105)	0.742 (0.000)	0.922 (0.003)	0.830 (0.007)	0.975 (0.215)
S&P 500	1.082 (1.498)	1.048 (0.967)	1.046 (1.743)	1.030 (0.988)	0.892 (0.040)	0.956 (0.046)	0.905 (0.044)	1.033 (0.993)	0.632 (0.000)	0.824 (0.000)	0.612 (0.000)	0.854 (0.000)
Consumer expect	1.029 (1.200)	0.852 (0.000)	1.055 (1.324)	0.712 (0.000)	1.063 (1.316)	1.026 (1.207)	1.057 (1.044)	0.907 (0.028)	0.880 (0.004)	0.968 (0.314)	0.949 (0.212)	1.111 (0.999)
Fraction < 1	0.500	0.667	0.417	0.750	0.667	0.583	0.583	0.667	0.833	0.583	0.833	0.833
Fraction < 1 at 95%	0.250	0.417	0.167	0.417	0.583	0.500	0.417	0.417	0.833	0.500	0.667	0.583

Notes: Reported are out-of-sample RMSE for HDMR-based τ -month ahead forecasts that use the lagged text plus the top 5 PCs of the [Stock and Watson \(2012\)](#) monthly macroeconomic series, relative to DFM-5, which uses only these PCs without the text, and relative to DMR with the same data. The transformed macroeconomic series used as prediction targets and for calculation of the PCs are generally first differences of logarithms (growth rates) for real quantity variables, first differences for nominal interest rates, and second differences of logarithms (changes in rates of inflation) for price series. [Diebold and Mariano \(1995\)](#) p -values testing the null hypothesis that the RMSE of HDMR is larger than the benchmark's, and corrected for multiple testing using the [Benjamini and Hochberg \(1995\)](#) procedure, are in parentheses.

Table 9: Nowcasting macroeconomic series

Folds:	Random		Rolling		Random		Rolling	
$Y_t^1 \setminus$ Benchmark:	DFM-5	DMR	DFM-5	DMR	DFM-5	DMR	DFM-5	DMR
IP: total	0.953 (0.097)	1.011 (1.017)	0.951 (0.243)	0.985 (0.451)	0.948 (0.013)	1.006 (0.732)	0.984 (0.760)	1.011 (0.746)
Emp: total	0.805 (0.000)	0.948 (0.175)	0.771 (0.000)	0.870 (0.000)	0.756 (0.000)	0.921 (0.011)	0.869 (0.012)	0.994 (0.644)
U: all	0.965 (0.437)	1.023 (0.971)	0.927 (0.252)	0.988 (0.514)	0.939 (0.196)	0.995 (0.712)	0.958 (0.719)	1.011 (0.820)
HStarts: Total	0.653 (0.000)	1.010 (1.033)	0.666 (0.000)	0.881 (0.000)	0.530 (0.000)	0.826 (0.000)	0.519 (0.000)	0.693 (0.000)
PMI	0.852 (0.001)	1.015 (0.985)	1.040 (1.846)	1.061 (1.058)	0.818 (0.000)	1.020 (0.778)	0.955 (0.671)	0.974 (0.357)
CPI-ALL	1.028 (1.166)	0.950 (0.094)	1.170 (1.200)	0.935 (0.038)	0.975 (0.175)	0.848 (0.000)	1.061 (1.196)	0.739 (0.000)
Real AHE: goods	1.038 (1.212)	1.021 (0.978)	1.044 (1.439)	1.001 (0.692)	1.025 (1.086)	1.001 (0.683)	1.110 (1.329)	1.059 (1.028)
FedFunds	0.918 (0.095)	1.017 (1.020)	1.076 (1.234)	1.070 (0.981)	0.940 (0.178)	1.058 (0.994)	1.120 (1.090)	1.112 (0.994)
M1	1.123 (1.091)	0.973 (0.694)	1.267 (1.000)	0.842 (0.001)	1.086 (0.999)	0.855 (0.000)	1.122 (1.000)	0.714 (0.000)
Ex rate: avg	1.043 (1.263)	0.980 (0.337)	1.060 (1.303)	0.947 (0.065)	1.037 (1.128)	0.952 (0.067)	1.038 (1.253)	0.868 (0.001)
S&P 500	1.007 (1.022)	0.987 (0.693)	1.033 (1.562)	1.025 (0.979)	1.017 (0.991)	0.995 (0.644)	0.991 (0.820)	0.969 (0.277)
Consumer expect	1.172 (1.000)	0.997 (0.866)	1.271 (1.091)	0.940 (0.045)	1.032 (1.071)	0.835 (0.000)	1.017 (1.371)	0.697 (0.000)
Fraction < 1	0.500	0.500	0.333	0.667	0.583	0.667	0.500	0.667
Fraction < 1 at 95%	0.250	0.000	0.167	0.417	0.333	0.417	0.167	0.417

Notes: Reported are out-of-sample RMSE for HDMR-based nowcasts that use the contemporaneous text plus the top 5 PCs of the [Stock and Watson \(2012\)](#) monthly macroeconomic series, relative to DFM-5, which uses only these PCs without the text, and relative to DMR with the same data. The transformed macroeconomic series used as prediction targets and for calculation of the PCs are generally first differences of logarithms (growth rates) for real quantity variables, first differences for nominal interest rates, and second differences of logarithms (changes in rates of inflation) for price series. [Diebold and Mariano \(1995\)](#) p -values testing the null hypothesis that the RMSE of HDMR is larger than the benchmark's, and corrected for multiple testing using the [Benjamini and Hochberg \(1995\)](#) procedure, are in parentheses.