

# Now- and Backcasting Initial Claims with High-Dimensional Daily Internet Search-Volume Data

Daniel Borup<sup>1</sup>   David Rapach<sup>2,3</sup>   Erik C. Montes Schütte<sup>1</sup>

<sup>1</sup>Aarhus University   <sup>2</sup>Washington University in St Louis   <sup>3</sup>Saint Louis University

2020 Banca d'Italia & Federal Reserve Board  
Joint Conference on Nontraditional Data & Statistical Learning  
November 11, 2020

- ▶ COVID-19 ⇒ economic upheaval in US
  - ▶ March 13, 2020 ⇒ national emergency ⇒ closures of non-essential retail establishments in many parts of US
- ▶ Unemployment insurance initial claims
  - ▶ Spiked in late March ⇒ 6.9M for week ending March 28
    - ▶ Subsequently declined, but remains at elevated levels
  - ▶ Perhaps most closely watched economic variable during COVID-19 crisis
    - ▶ High-frequency (weekly) labor market indicator
    - ▶ US weekly economic indicator ([Lewis, Mertens & Stock 2020](#))

- ▶ We use rich trove of daily internet search-volume data from **Google Trends** to predict US weekly UI
  - ▶ Seek to improve prediction during COVID-19 crisis
- ▶ High-dimensional set of GT terms related to *unemployment*
  - ▶ Daily data for 103 GT terms for most recent 7 days
- ▶ Sequence of now-/backcasts  $\Rightarrow$   $\text{Week}_t$  UI
  - ▶  $\text{Week}_t$  UI  $\Rightarrow$   $\text{Sat}_t$  to  $\text{Sun}_t$
  - ▶ In anticipation of data release by Dept of Labor on  $\text{Thu}_{t+1}$

- ▶ Predictive models
  - ▶  $UI = f(\text{first/second lag of UI, 7 days of 103 GT terms})$
  - ▶ High-dimensional setting  $\Rightarrow$  machine learning
    - ▶  $\# \text{ predictors/inputs} = 7 \times 103 + 1 = 722$
- ▶ Linear specification  $\Rightarrow$  penalized regression
  - ▶ LASSO (Tibshirani 1996)
  - ▶ ENet (Zou & Hastie 2005)
- ▶ Artificial neural nets
  - ▶ Allow for complex, nonlinear predictive relationships

- ▶ Mixed-frequency data  $\Rightarrow$  weekly/daily
  - ▶ U-MIDAS (Forni, Marcellino & Schumacher 2015)
  - ▶ Analyze how information flow affects accuracy of sequence of now-/backcasts
    - ▶ Successive elements include more recent GT data
- ▶ Combine 2 branches of economic forecasting literature
  - ▶ Machine learning (eg, Diebold & Shin 2019; Kotchoni, Leroux & Stevanovic 2019; Medeiros et al forthcoming)
  - ▶ Mixed-frequency data (eg, Clements & Galvao 2008; Forni & Marcellino 2014; Brave, Butters & Justiniano 2019)

- ▶ Emerging literature  $\Rightarrow$  internet search-volume data to predict labor market variables
  - ▶ Choi & Varian (2012)  $\Rightarrow$  small # of GT terms to predict US UI (through mid 2011)
  - ▶ D'Amuri & Marcucci (2017)  $\Rightarrow$  GT terms including *jobs* improve US unemployment rate prediction
  - ▶ Niesert et al (2020)  $\Rightarrow$  array of GT terms useful for predicting unemployment rates in collection of developed countries
  - ▶ Borup & Schütte (forthcoming)  $\Rightarrow$  large # of GT terms & machine-learning tools improve US employment growth prediction

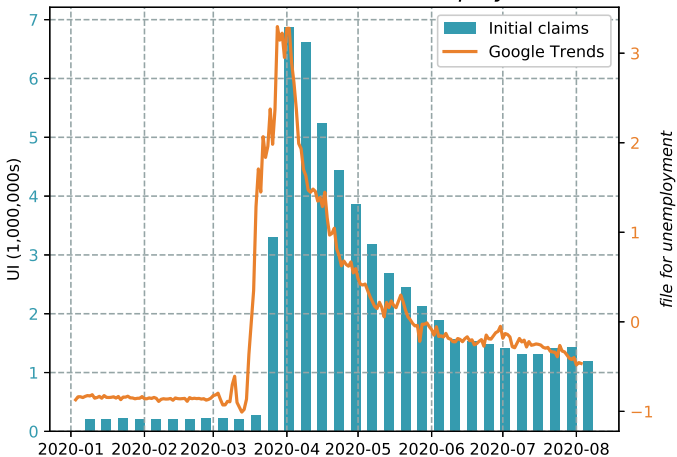
- ▶ Recent studies  $\Rightarrow$  GT terms to predict UI during COVID-19 crisis
  - ▶ Aaronson et al (2020), Goldsmith-Pinkham & Sojourner (2020), Larson & Sinclair (2020)
- ▶ Key elements of our approach
  - ▶ Consider large # of GT terms
  - ▶ Mixed-frequency framework
  - ▶ Machine learning

- ▶ Target  $\Rightarrow$  weekly UI (Sat to Sun)
  - ▶ Thursday 8:30a EST  $\Rightarrow$  Dept of Labor releases UI figure for previous week
- ▶ Seasonally adjusted or non-seasonally adjusted UI?
  - ▶ Subject of debate during COVID-19 crisis ([Rinz 2020](#))
  - ▶ Dept of Labor  $\Rightarrow$  recent  $\Delta$  to additive SA ([Davidson 2020](#))
    - ▶ New procedure only used going forward
  - ▶ We analyze both SA & NSA UI

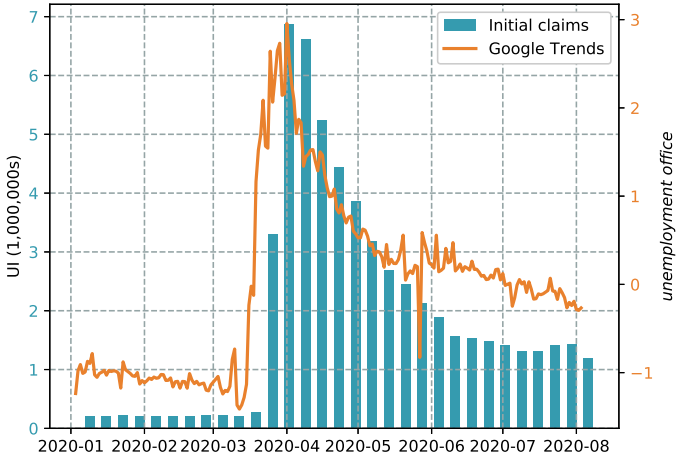


- ▶ Daily search-volume data from GT
  - ▶ Index  $\Rightarrow$  proportion of queries for specific search term within geographical area
  - ▶ Index released with 36-hour delay
    - ▶ Google filters irregular search activity (automated, spam)
- ▶ Start with source term *unemployment*
  - ▶ Google Keyword Planner  $\Rightarrow$  obtain top 15 primitive terms associated with *unemployment*
  - ▶ Expand primitive terms with 25 related terms (*top* category)
  - ▶ Remove low-volume queries  $\Rightarrow$  103 unique terms

### Initial claims and *file for unemployment*



### Initial claims and *unemployment office*



- ▶ Also consider 10 daily macro-financial predictors
  - ▶ INF  $\Rightarrow$  expected Inflation (rate at which T-bill & TIPS achieve same yield)
  - ▶ GLD  $\Rightarrow$  gold return (London afternoon fixing price)
  - ▶ VIX  $\Rightarrow$  CBOE market volatility index
  - ▶ TED  $\Rightarrow$  3-mo LIBOR yield minus 3-mo T-bill yield
  - ▶ TERM  $\Rightarrow$  10-yr T-bond yield minus 3-mo T-bill yield
  - ▶ REITS  $\Rightarrow$  Wilshire US REIT total market index return

- ▶ 10 daily macro-financial predictors (cont'd)
  - ▶ DEF  $\Rightarrow$  Baa-rated corp bond yield minus Aaa-rated corp bond yield
  - ▶ SP500  $\Rightarrow$  S&P 500 index return
  - ▶ EPU  $\Rightarrow$  economic policy uncertainty index (**Baker, Bloom & Davis 2016**)
  - ▶ DIS  $\Rightarrow$  news-based infectious disease equity market volatility tracker (<http://www.policyuncertainty.com/>)

Prediction formation	Backcast/ nowcast	GT used for prediction	Overlap (days)	Available UI release
$\text{Mon}_t$	Nowcast	$\text{Sun}_{t-1}$ to $\text{Sat}_{t-1}$	0	$\text{Week}_{t-2}$
$\text{Tue}_t$	Nowcast	$\text{Mon}_{t-1}$ to $\text{Sun}_t$	1	$\text{Week}_{t-2}$
$\text{Wed}_t$	Nowcast	$\text{Tue}_{t-1}$ to $\text{Mon}_t$	2	$\text{Week}_{t-2}$
$\text{Thu}_t$	Nowcast	$\text{Wed}_{t-1}$ to $\text{Tue}_t$	3	$\text{Week}_{t-1}$
$\text{Fri}_t$	Nowcast	$\text{Thu}_{t-1}$ to $\text{Wed}_t$	4	$\text{Week}_{t-1}$
$\text{Sat}_t$	Nowcast	$\text{Fri}_{t-1}$ to $\text{Thu}_t$	5	$\text{Week}_{t-1}$
$\text{Sun}_{t+1}$	Backcast	$\text{Sat}_{t-1}$ to $\text{Fri}_t$	6	$\text{Week}_{t-1}$
$\text{Mon}_{t+1}$	Backcast	$\text{Sun}_t$ to $\text{Sat}_t$	7	$\text{Week}_{t-1}$
$\text{Tue}_{t+1}$	Backcast	$\text{Mon}_t$ to $\text{Sun}_{t+1}$	6	$\text{Week}_{t-1}$
$\text{Wed}_{t+1}$	Backcast	$\text{Tue}_t$ to $\text{Mon}_{t+1}$	5	$\text{Week}_{t-1}$

$$\blacktriangleright UI_t = f^{(j)} \left( UI_{t-1(2)}, \mathbf{g}_t^{(j)}; \theta^{(j)} \right)$$

$$\blacktriangleright \underbrace{\mathbf{g}_t^{(j)}}_{7K \times 1} = \left[ \mathbf{g}'_{t-j/7} \quad \mathbf{g}'_{t-(j+1)/7} \quad \cdots \quad \mathbf{g}'_{t-(j+6)/7} \right]'$$

$\blacktriangleright \mathbf{g}_{t-i/7} \Rightarrow K$ -vector of GT terms for  $(7-i)$ th day of  $Week_t$

$\blacktriangleright$  Sunday,  $\mathbf{g}_{t-6/7}$ ; Monday,  $\mathbf{g}_{t-5/7}$ ; Tuesday,  $\mathbf{g}_{t-4/7}$ ; Wednesday,  $\mathbf{g}_{t-3/7}$ ; Thursday,  $\mathbf{g}_{t-2/7}$ ; Friday,  $\mathbf{g}_{t-1/7}$ ; Saturday,  $\mathbf{g}_t$

$\blacktriangleright K = 103$

$\blacktriangleright \#$  of predictors  $\Rightarrow 7 \times 103 + 1 = 722$

$\blacktriangleright 7 \times 113 + 1 = 792$  with macro-financial predictors

- ▶ AR benchmark  $\Rightarrow UI_t = \alpha_0 + \alpha_1 UI_{t-1(2)} + \varepsilon_t$
- ▶ Linear  $\Rightarrow UI_t = \beta_0^{(j)} + \beta_{AR}^{(j)} UI_{t-1(2)} + \beta_g^{(j)'} \mathbf{g}_t^{(j)} + \varepsilon_t^{(j)}$ 
  - ▶  $\beta_g^{(j)} \Rightarrow 7K$ -vector of slope coefficients for daily GT terms
  - ▶  $\theta^{(j)} = \left[ \beta_0^{(j)} \quad \beta_{AR}^{(j)} \quad \beta_g^{(j)'} \right]'$
- ▶ U-MIDAS model
  - ▶ Higher-frequency predictors in  $\mathbf{g}_t^{(j)}$  have own coefficients
  - ▶ Machine-learning methods  $\Rightarrow$  flexibly estimate weights (rather than impose lag-polynomial structure) while guarding against overfitting



► LASSO  $\Rightarrow$  penalized regression

$$\text{► } \arg \min_{\boldsymbol{\theta}^{(j)} \in \mathbb{R}^{7K+2}} \frac{1}{2\mathcal{J}} \left\{ \sum_{t=1}^{\mathcal{J}} \left[ \mathbf{U}\mathbf{I}_t - \left( \beta_0^{(j)} + \beta_{\text{AR}}^{(j)} \mathbf{U}\mathbf{I}_{t-1(2)} + \boldsymbol{\beta}_g^{(j)'} \mathbf{g}_t^{(j)} \right) \right] \right\}^2 + \lambda \|\boldsymbol{\beta}^{(j)}\|_1$$

$$\text{► } \boldsymbol{\beta}^{(j)} = \begin{bmatrix} \beta_{\text{AR}}^{(j)} & \boldsymbol{\beta}_g^{(j)'} \end{bmatrix}'$$

►  $\lambda \geq 0 \Rightarrow$  regularization parameter

► Select  $\lambda$  via ERIC (Hui, Warton & Foster 2015)

►  $\ell_1$  penalty  $\Rightarrow$  permits shrinkage to zero (ie, variable selection)

▶ ENet  $\Rightarrow$  refinement of LASSO

$$\text{▶ } \arg \min_{\boldsymbol{\theta}^{(j)} \in \mathbb{R}^{7K+2}} \frac{1}{2T} \left\{ \sum_{t=1}^T \left[ \mathbf{U}_t - \left( \beta_0^{(j)} + \beta_{\text{AR}}^{(j)} \mathbf{U}_{t-1(2)} + \boldsymbol{\beta}_g^{(j)'} \mathbf{g}_t^{(j)} \right) \right] \right\}^2 + \lambda P_\alpha \left( \boldsymbol{\beta}^{(j)} \right)$$

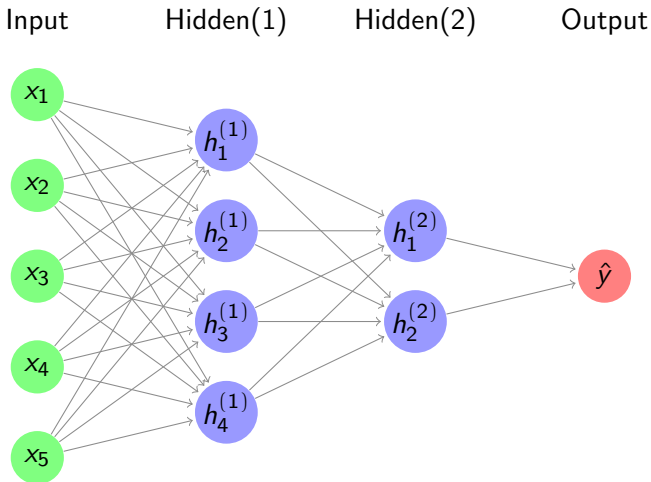
$$\text{▶ } P_\alpha \left( \boldsymbol{\beta}^{(j)} \right) = 0.5(1 - \alpha) \|\boldsymbol{\beta}^{(j)}\|_2^2 + \alpha \|\boldsymbol{\beta}^{(j)}\|_1$$

▶  $0 \leq \alpha \leq 1 \Rightarrow$  blending parameter for  $\ell_1$  &  $\ell_2$  components

▶ Set  $\alpha = 0.5$  (Hastie & Qian 2016)

▶ Compared to LASSO, ENet tends to select highly correlated predictors as group

- ▶ Input layer (set of predictors)  $\Rightarrow x_1, \dots, x_{P_0}$
- ▶ One or more hidden layers  $\Rightarrow$  take signals from neurons in previous layer to generate subsequent signals
  - ▶  $h_m^{(l)} = g\left(w_{m,0}^{(l)} + \sum_{j=1}^{P_{l-1}} w_{mj}^{(l)} h_j^{(l-1)}\right)$ 
    - ▶ Hidden layers  $\Rightarrow l = 1, \dots, L$ ; neurons  $\Rightarrow m = 1, \dots, P_l$
    - ▶ ReLU  $\Rightarrow g(x) = 0$  if  $x < 0$ ,  $g(x) = x$  otherwise
- ▶ Output layer  $\Rightarrow$  translates signals from last hidden layer into prediction
  - ▶  $\hat{y} = w_0^{(L+1)} + \sum_{j=1}^{P_L} w_j^{(L+1)} h_j^{(L)}$



- ▶ Specify 3 ANN architectures
  - ▶ NN1  $\Rightarrow$  1 hidden layer with 104 neurons
  - ▶ NN2  $\Rightarrow$  2 hidden layers with 104 & 10 neurons
  - ▶ NN3  $\Rightarrow$  3 hidden layers with 104, 10 & 3 neurons
- ▶ Training ANN  $\Rightarrow$  estimate (very many) weights
  - ▶ Minimize training-sample MSE augmented with  $\ell_1$  penalty
  - ▶ Adam SGD algorithm (Kingma & Ba 2015)
    - ▶ Python  $\Rightarrow$  keras package

- ▶ Ensemble-Linear
  - ▶ Avg of Linear-LASSO & Linear-ENet
- ▶ Ensemble-ANN
  - ▶ Avg of NN1, NN2 & NN3
- ▶ Ensemble-All
  - ▶ Avg of Linear-LASSO, Linear-ENet, NN1, NN2 & NN3
- ▶ All models fitted using 10-yr rolling window
  - ▶ Full sample  $\Rightarrow$  Jan 2005 to Aug 2020
  - ▶ Out-of-sample period  $\Rightarrow$  Jan 2015 to Aug 2020

## RMSE ratios, seasonally adjusted UI

Prediction	AR RMSE	LASSO	ENet	NN1	NN2	NN3
Mon <sub>t</sub>	642,032	0.8291	<b>0.8111</b>	0.9736	0.9123	0.8130
Tue <sub>t</sub>	642,032	<b>0.7152</b>	0.7181	0.8658	0.9225	0.7477
Wed <sub>t</sub>	642,032	0.6769	<b>0.6698</b>	0.7475	0.7418	0.8340
Thu <sub>t</sub>	445,565	0.6748	<b>0.6598</b>	0.9753	0.6890	0.7610
Fri <sub>t</sub>	445,565	<b>0.6031</b>	0.6167	0.7583	0.6301	0.6435
Sat <sub>t</sub>	445,565	<b>0.4925</b>	0.5130	0.5998	0.5761	0.7204
Sun <sub>t+1</sub>	445,565	<b>0.4651</b>	0.5269	0.5489	0.5591	0.5971
Mon <sub>t+1</sub>	445,565	<b>0.4543</b>	0.5171	0.5355	0.5367	0.6049
Tue <sub>t+1</sub>	445,565	<b>0.4643</b>	0.4743	0.4969	0.5345	0.6017
Wed <sub>t+1</sub>	445,565	0.4601	0.4947	<b>0.4413</b>	0.5384	0.6857

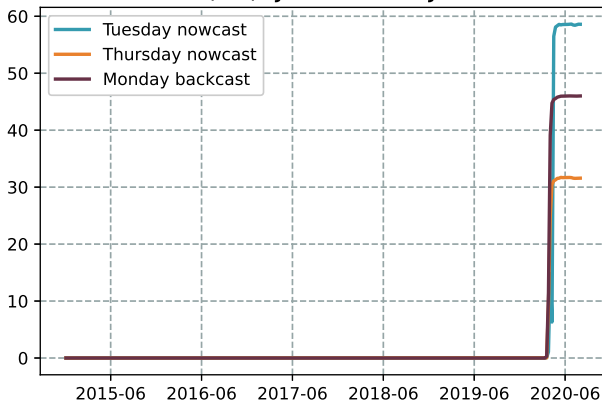
## RMSE ratios, seasonally adjusted UI

Prediction	AR RMSE	Ens-Linear	Ens-ANN	Ens-All
$\text{Mon}_t$	642,032	<b>0.8197</b>	0.8862	0.8422
$\text{Tue}_t$	642,032	<b>0.7145</b>	0.8194	0.7566
$\text{Wed}_t$	642,032	<b>0.6723</b>	0.7345	0.6898
$\text{Thu}_t$	445,565	<b>0.6653</b>	0.7739	0.7124
$\text{Fri}_t$	445,565	<b>0.6059</b>	0.6381	0.6175
$\text{Sat}_t$	445,565	<b>0.5008</b>	0.5892	0.5321
$\text{Sun}_{t+1}$	445,565	<b>0.4913</b>	0.5259	0.5044
$\text{Mon}_{t+1}$	445,565	<b>0.4828</b>	0.5281	0.5018
$\text{Tue}_{t+1}$	445,565	<b>0.4607</b>	0.5167	0.4845
$\text{Wed}_{t+1}$	445,565	<b>0.4743</b>	0.4963	0.4827



## Cumulative differences in squared errors

## LASSO (SA): Jan 2015 to Jul 2020



## RMSE ratios, non-seasonally adjusted UI

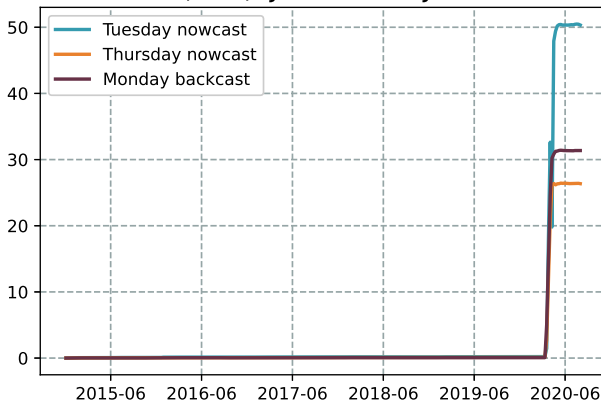
Prediction	AR RMSE	LASSO	ENet	NN1	NN2	NN3
Mon <sub>t</sub>	569,418	0.7448	0.7553	<b>0.6975</b>	0.8412	0.7629
Tue <sub>t</sub>	569,418	0.6821	0.6336	<b>0.6291</b>	0.6830	0.7291
Wed <sub>t</sub>	569,418	0.6262	0.6189	0.6583	0.6285	<b>0.5826</b>
Thu <sub>t</sub>	352,507	0.6478	0.6424	0.7046	<b>0.5230</b>	0.6274
Fri <sub>t</sub>	352,507	0.6014	<b>0.5734</b>	0.6766	0.6071	0.5967
Sat <sub>t</sub>	352,507	0.5525	0.5434	0.6770	<b>0.4330</b>	0.6773
Sun <sub>t+1</sub>	352,507	0.5802	0.5867	<b>0.4357</b>	0.4753	0.4505
Mon <sub>t+1</sub>	352,507	0.4970	0.4805	0.3830	<b>0.3682</b>	0.5540
Tue <sub>t+1</sub>	352,507	0.4958	0.4909	<b>0.3925</b>	0.4632	0.6192
Wed <sub>t+1</sub>	352,507	0.5334	0.5404	<b>0.4473</b>	0.5544	0.6067

## RMSE ratios, non-seasonally adjusted UI

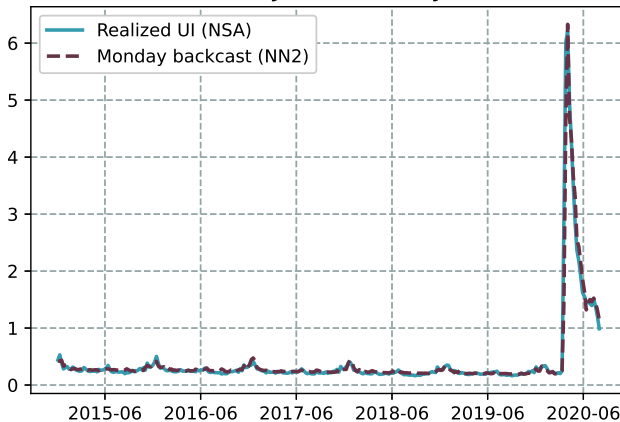
Prediction	AR RMSE	Ens-Linear	Ens-ANN	Ens-All
Mon <sub>t</sub>	642,032	0.7491	0.7583	<b>0.7466</b>
Tue <sub>t</sub>	642,032	0.6815	<b>0.6572</b>	0.6604
Wed <sub>t</sub>	642,032	0.6205	<b>0.6025</b>	0.6030
Thu <sub>t</sub>	445,565	0.6443	<b>0.5472</b>	0.5879
Fri <sub>t</sub>	445,565	0.5865	<b>0.5768</b>	0.5773
Sat <sub>t</sub>	445,565	0.5456	<b>0.4751</b>	0.5043
Sun <sub>t+1</sub>	445,565	0.5831	<b>0.4208</b>	0.4915
Mon <sub>t+1</sub>	445,565	0.4869	<b>0.3887</b>	0.4345
Tue <sub>t+1</sub>	445,565	0.4927	<b>0.4424</b>	0.4631
Wed <sub>t+1</sub>	445,565	0.5355	<b>0.4821</b>	0.5013

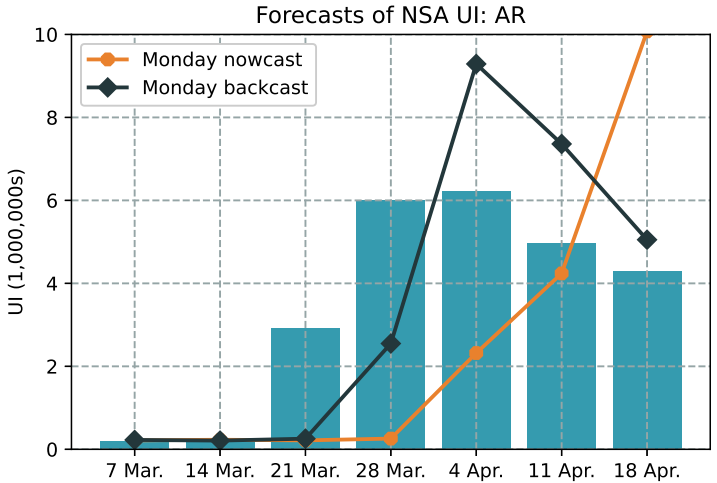
## Cumulative differences in squared errors

### NN2 (NSA): Jan 2015 to Jul 2020

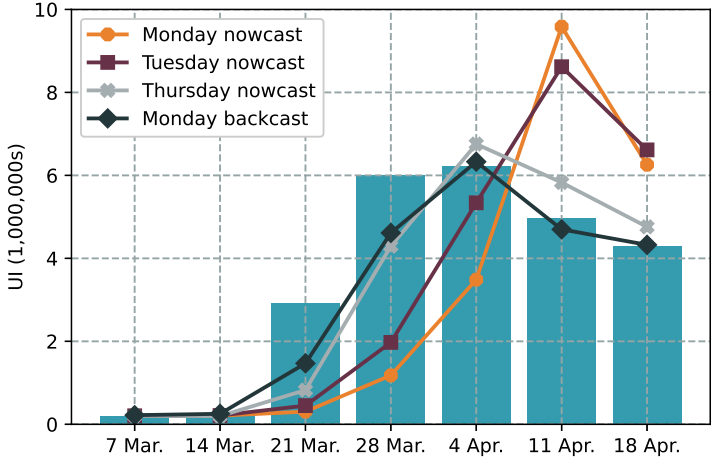


## NN2 (NSA): Jan 2015 to Jul 2020

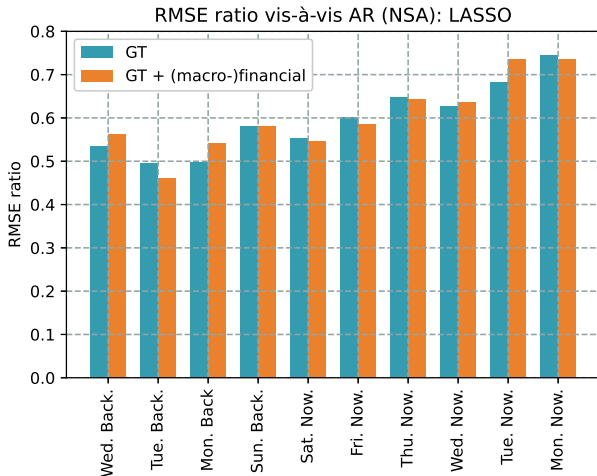




Forecasts of NSA UI: NN2

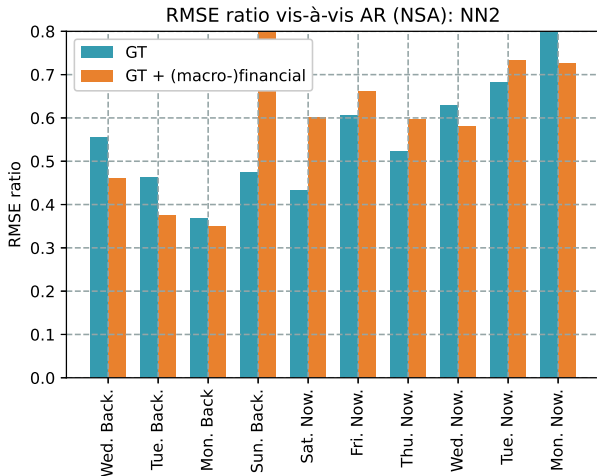


## Daily GT Terms &amp; Macro-Financial Predictors





## Daily GT Terms &amp; Macro-Financial Predictors



- ▶ Partial dependence plot (**Friedman 2001**)
  - ▶ Marginal effect of  $x_s^{(j)}$  on expected value of  $UI_t$  for fitted model

- ▶ Predictors  $\Rightarrow \mathbf{x}_t^{(j)} = \left[ UI_{t-1(2)} \quad \mathbf{g}_t^{(j)} \right]'$

- ▶  $722(792) \times \mathcal{T}$  data matrix  $\Rightarrow \mathbf{X}_{\mathcal{T}}^{(j)} = \left[ \mathbf{x}_1^{(j)} \quad \dots \quad \mathbf{x}_{\mathcal{T}}^{(j)} \right]'$

- ▶  $\mathbf{x}_{C(s)}^{(j)} = \mathbf{x}^{(j)} \setminus x_s^{(j)}$

- ▶  $PD(x_s^{(j)}) = \mathbb{E}_{\mathbf{x}_{C(s)}^{(j)}} \left[ \hat{f}^{(j)}(x_s^{(j)}, \mathbf{x}_{C(s)}^{(j)}) \right]$
- $$= \int_{\mathbf{x}_{C(s)}^{(j)}} \hat{f}^{(j)}(x_s^{(j)}, \mathbf{x}_{C(s)}^{(j)}) p_{C(s)}^{(j)}(\mathbf{x}_{C(s)}^{(j)}) d\mathbf{x}_{C(s)}^{(j)}$$
- ▶  $p_{C(s)}^{(j)}(\mathbf{x}_{C(s)}^{(j)}) = \int_{x_s^{(j)}} p(\mathbf{x}^{(j)}) dx_s^{(j)}$
- ▶  $p(\mathbf{x}^{(j)}) \Rightarrow$  joint probability density for  $\mathbf{x}^{(j)}$
- ▶  $\hat{f}(\mathbf{x}^{(j)}) \Rightarrow$  prediction function for fitted model

- ▶ Estimate via Monte Carlo integration using training sample

- ▶  $\widehat{\text{PD}}\left(x_s^{(j)}\right) = \frac{1}{\mathcal{J}} \sum_{t=1}^{\mathcal{J}} \hat{f}^{(j)}\left(x_s^{(j)}, \mathbf{x}_{C(s),t}^{(j)}\right)$

- ▶ Evaluate at  $x_{s,t}^{(j)}$  for  $t = 1, \dots, \mathcal{J}$  (or set of quantiles)

- ▶ PDP for linear model  $\Rightarrow$  constant slope

- ▶ Horizontal line for predictor not selected by LASSO/ENet

- ▶ Compare PDPs for fitted linear models & ANNs

- ▶ Gauge relative importance of nonlinearities in ANNs

- ▶ Variable importance (Greenwell, Boehmke & McCarthy 2018)

- ▶ 
$$\hat{j}(x_s^{(j)}) = \left\{ \frac{1}{\mathcal{T}-1} \sum_{t=1}^{\mathcal{T}} \left[ \widehat{\text{PD}}(x_{s,t}^{(j)}) - \frac{1}{\mathcal{T}} \sum_{t=1}^{\mathcal{T}} \widehat{\text{PD}}(x_{s,t}^{(j)}) \right]^2 \right\}^{0.5}$$

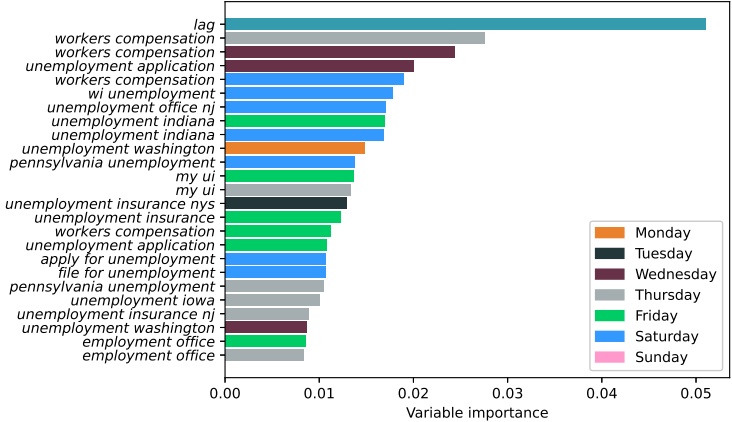
- ▶ Variation in PDP around average (ie, std dev)
  - ▶ Horizontal PDP  $\Rightarrow$  variable importance = 0

- ▶ Scale to facilitate comparison across predictors

- ▶ 
$$\tilde{j}(x_s^{(j)}) = \frac{\hat{j}(x_s^{(j)})}{\sum_{p=1}^P \hat{j}(x_p^{(j)})}$$

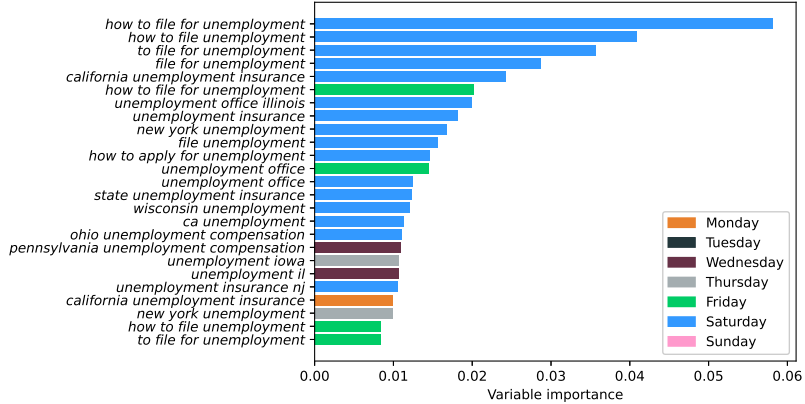
- ▶  $P \Rightarrow$  total # of predictors

Monday backcast: NN2 (NSA) excluding COVID-19 crisis



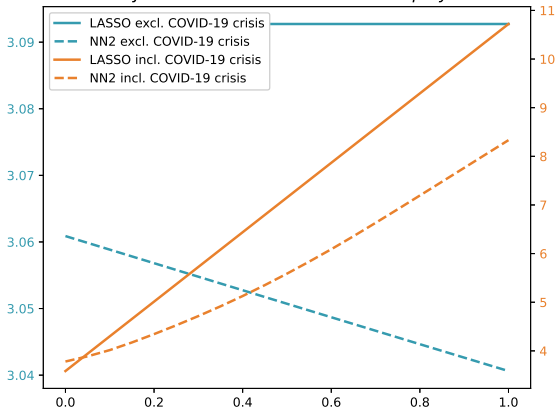
Importance of GT Terms

Monday backcast: NN2 (NSA) including COVID-19 crisis



## Partial dependence plots

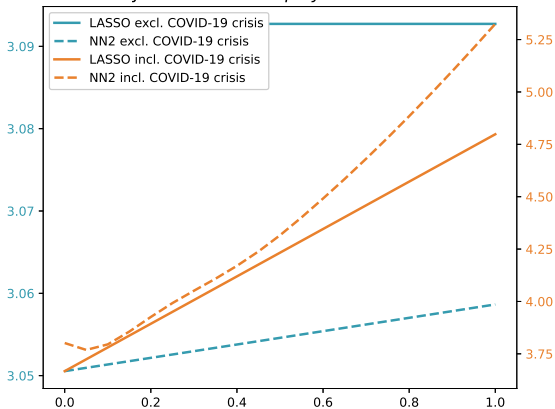
Monday backcast: *how to file for unemployment*





## Partial dependence plots

Monday backcast: *unemployment office illinois*



- ▶ Daily internet search-volume data for 103 GT terms related to *unemployment* improve now-/backcasts of weekly UI
  - ▶ Especially during advent of COVID-19 crisis
- ▶ GT terms become more important during COVID-19 crisis
- ▶ Nonlinearities more relevant for predicting NSA UI
- ▶ Combining mixed-frequency approach (via U-MIDAS) with machine learning appears promising
- ▶ New website provides daily updated now-/backcasts
  - ▶ <https://www.uinowcast.org/>