PREDICTION FEVER

Predicting the Flu via Regional Effects

Shannon Gallagher Statistics Advisors[†]: Ryan Tibshirani & Bill Eddy External Advisor[‡]: Roni Rosenfeld November 6, 2015

Carnegie Mellon University [†]Department of Statistics [‡] MIDAS & Department of Computer Science

How bad is the flu going to be this year?

OUTLINE

- · Background
- · Data
- \cdot ADA Goals
- · Modeling
 - · Empirical Bayes with Regional Effects
 - · Target Model
 - · Posterior Biasing
- · Final Results
- \cdot Future Work

BACKGROUND

The flu is...

- \cdot old
- \cdot deadly
- \cdot costly
- \cdot stochastic

With accurate predictions, the flu is

- · old
- $\cdot \frac{\text{deadly}}{\text{deadly}} \rightarrow \text{manageable}$
 - · Resource allocation
 - $\cdot\,$ Alert health officials
 - Issue warnings
- $\cdot \text{ costly} \rightarrow \text{feasible}$
 - · Fewer sick days
 - · More awareness
- $\cdot \ -\text{stochastic} \rightarrow \text{forecasted}$

DATA

THE CDC COLLECTS DATA VOLUNTARILY FROM PHYSICIANS

Region	Year	Week	wILI
1	2015	25	0.55
2	2015	25	1.45
:	:	:	÷
10	2015	25	0.44
:	:	:	÷
1	2015	45	0.76
2	2015	45	1.57
:	:	:	:
10	2015	45	0.89

Table: Cross Section of Available data.wILI - Weighted Influenza Like Illness

FLU CURVES USUALLY HAVE ONE PROMINENT PEAK



Figure: Examples of wILI curves. From David Farrow's FluV. epicast.org

ADA GOALS

WE WANT TO PREDICT THE WILI FOR THE REMAINING WEEKS OF A SEASON



Figure: Examples of wILI curves. From David Farrow's FluV. epicast.org

for each of the 10 cdc regions



Figure: From cdc.gov

WE HAVE EVIDENCE OF REGIONAL DEPENDENCIES



Figure: 2001 Week 51. This map is generally representative of the other seasons, when the peak flu season begins. Region 6 (Texas +) is hit hard and first. Regions 1 (NE) and 10 (PNW) are barely effected. PA is generally in the middle.

- 1. Allow for regional dependencies
- 2. Predict the wILI values for the remainder of the season
- 3. Predict specific targets
 - · Peak Week
 - · Peak Height/Incidence
- 4. Produce distributional forecasts
- 5. Explain our model to health care professionals

MODELING



Figure: From delphi.midas.cs.cmu.edu.

Empirical Bayes (EB) Framework is based on the assumptions:

1. A flu curve will look like a past curve + modifications + noise

2. The regions are independent from one another



Figure: From delphi.midas.cs.cmu.edu.

Empirical Bayes (EB) Framework is based on the assumptions:

1. A flu curve will look like a past curve + modifications + noise

2. The regions are independent from one another

- 1. Create EB model with regional effects for full curve
- 2. Validate model through cross validation
- 3. Rejoice over great results
- 4. Create another model focusing on targets
- 5. Combine two models (Posterior Biasing)
- 6. Cross validate
- 7. Finish

 $Y_t^{(r,s)} \sim N(\mu_t^{(r,s)}, \sigma^2)$

where

$$\mu_{t}^{(r,s)} = [a_{s} \cdot \alpha_{r}] \cdot f(t - b_{s} - \beta_{r})$$

for week t, region r, season s and priors:

 $\begin{array}{ll} a_{\rm s} \sim {\rm Unif}(2,10) & -{\rm seasonal\ scaling} \\ b_{\rm s} \sim {\rm Unif}\{-6,-5,\ldots,6\} & -{\rm seasonal\ shifting} \\ \alpha_{\rm r} \sim {\rm Unif}(0.25,1.25) & -{\rm regional\ scaling} \\ \beta_{\rm r} \sim {\rm Unif}\{-3,-2,\ldots,3\} & -{\rm regional\ shifting} \\ f \sim {\rm Unif}\{\hat{\rm F}\} & -{\rm smoothed\ observed\ curves} \\ \sigma^2 \sim {\rm Unif}(0.5,2.5) & -{\rm variance} \end{array}$



For given season, given weeks known, the posterior is

$$\mathcal{P}(\text{fut. wILI}|\text{obs. wILI, params}) \propto \frac{1}{\sigma} \exp\left\{\frac{\sum_{\text{regions r}} \sum_{\text{obs. weeks t}} (y - \hat{y})_{r,t}^2}{2\sigma^2}\right\}$$

For given season, given weeks known, the posterior is

$$\mathcal{P}(\text{fut. wILI}|\text{obs. wILI, params}) \propto \frac{1}{\sigma} \exp\left\{\frac{\sum_{\text{regions r}} \sum_{\text{obs. weeks t}} (y - \hat{y})_{r,t}^2}{2\sigma^2}\right\}$$

Recall the priors for the regional variables are

 $\label{eq:arcorrelation} \begin{aligned} & \alpha_{\rm r} \sim {\sf Unif(0.25, 1.25)} & -{\sf regional scaling} \\ & \beta_{\rm r} \sim {\sf Unif\{-3, -2, \dots, 3\}} & -{\sf regional shifting} \end{aligned}$

 \implies If we choose just 5 values for each parameter, then summing over the regions is 25^{10} possible combinations!

For given season, given weeks known, the posterior is

$$\mathcal{P}(\text{fut. wILI}|\text{obs. wILI, params}) \propto \frac{1}{\sigma} \exp\left\{\frac{\sum_{\text{regions r}} \sum_{\text{obs. weeks t}} (y - \hat{y})_{r,t}^2}{2\sigma^2}\right\}$$

Recall the priors for the regional variables are

 $\label{eq:arcorrelation} \begin{aligned} & \alpha_{\rm r} \sim {\sf Unif(0.25, 1.25)} & -{\sf regional scaling} \\ & \beta_{\rm r} \sim {\sf Unif\{-3, -2, \dots, 3\}} & -{\sf regional shifting} \end{aligned}$

 \implies If we choose just 5 values for each parameter, then summing over the regions is <u>25¹⁰</u> possible combinations!

Approximation: Fix the regional variables through min. error.

PREDICTIONS FROM EB WITH REGIONAL EFFECTS LOOK LIKE FLU CURVES



BUT THE APPROXIMATIONS IN IMPLEMENTATION HURT US.



ESTIMATING A POINT IS SIMPLER THAN ESTIMATING A CURVE



Figure: Leave-one-season-out cross validation for EB with regional effects (a) and a targetted regression (b). The x-axis is the weeks from the observed peak and the y-axis is the mean absolute error.

Peak Height_{r,s} =
$$\alpha_0^{(r)} Y_{t,r,s} + \alpha_1^{(r)} Y_{t-1,r,s} + \alpha_2^{(r)} Y_{t-2,r,s} + \alpha_3^{(r)}$$
Week
+ $\alpha_4^{(r)} \cdot Max$ Height $Obs_{t,r,s} + \alpha_5^{(r)} \cdot Max$ Week $Obs_{t,r,s}$
+ $\sum_{i \neq r} \left(\beta_{0,i}^{(r)} Y_{t,i,s} + \beta_{1,i}^{(r)} Y_{t-1,i,s} + \beta_{2,i}^{(r)} Y_{t-2,i,s} \right) + \epsilon,$

- · Model I
 - \cdot Complete independence. Train separate model for each region. (All $\beta_{j,i}=0)$
- Model II
 - · United Model. Same model for all regions. (All $\beta_{j,i} = 0$, $\alpha_i^{(r)} = \alpha_i^{(r')}$ for all r, r')
- · Model III
 - Semi-Dependence. Train separate model for each region and impose sparsity penalty on the $\beta_{j,i}s, j \in \{0, 1, 2\}$.

VARIABLE SELECTION REVEALS REGIONAL DEPENDENCIES



Figure: Arrows colored by size of effect. A line between regions indicates that one of the region's was used in the other region's regression model.

- $\cdot\,$ Spline fit on the lagged values of the region in question
- · Incorporation of regional effects
 - $\cdot\,$ Mixture of Model I and Model III chosen by lowest CV error

ULTIMATELY, WE WANT TO BIAS ESTIMATES/SHRINK POSTERIOR



Visualization of Posterior Biasing

Figure: Image depicting of weighting curves whose peak values are closer to our estimated values. A thicker line represents a larger weight. The blue dot is our estimated value of the peak height and week.

POSTERIOR BIASING YIELDS IMPROVED RESULTS





Figure: CV averaged over the different seasons for the different model types.

FINAL RESULTS

Final Results						
Model	Season CV	CV St. Dev.	Pros	Cons		
EB	179	40	Baseline	Independ.		
EB w/ Reg.	187	48	Regional Effects	Curse of Dim.		
Target	57	7	Simpler	No Curve		
Post. Bias.	91	28	Flexible	Interpret.		

- 1. Allow for regional dependencies Targets use regional data
- 2. Predict the wILI values for the remainder of the season -Competitive Model
- 3. Predict specific targets biased toward targets
 - · Peak Week
 - · Peak Height/Incidence
- 4. Produce distributional forecasts We have entire posterior
- 5. Explain our model to health care professionals will talk to other MIDAS researchers

- $\cdot\,$ Submit predictions to CDC Flu Prediction Contest
- · Quantile Regression
- · More regional data





Calendar Week