

Predicting Alumni Income Based on University-level Data

Niko Miller, Akseli Manninen, Santeri Löppönen



Project Introduction

- Past studies have concluded that a strong connection exists between higher education and income¹
- We further examine this link by regressing median earnings of students 10 years after entry on university-level covariates
- We extend previous studies by considering demographic and geographic factors
- Data on 6681 U.S. universities spans over 9 regions, 59 state post codes, 2430 cities
- Bayesian models:
 - Pooled model for all universities
 - Separate and hierarchical model with region-grouping

¹(Card, 1999)

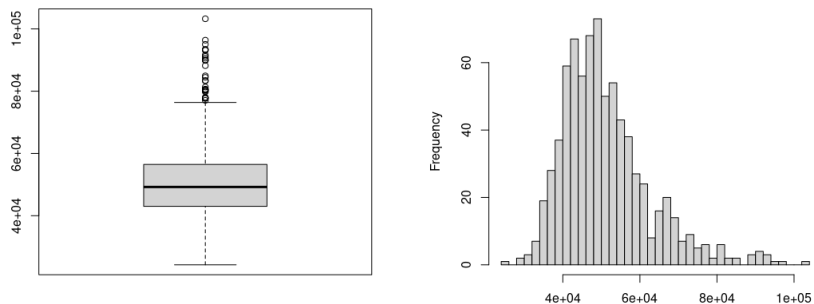
Data Description

- Most recent institutional-level college scorecard data from U.S. Department of Education¹
- 6681 observations/universities and 2989 variables
- Aggregate data for each university
 - Institutional characteristics, enrollment, student aid, costs and student outcomes
- Data quality is quite good, but many variables have lots of NAs
- Why this dataset?

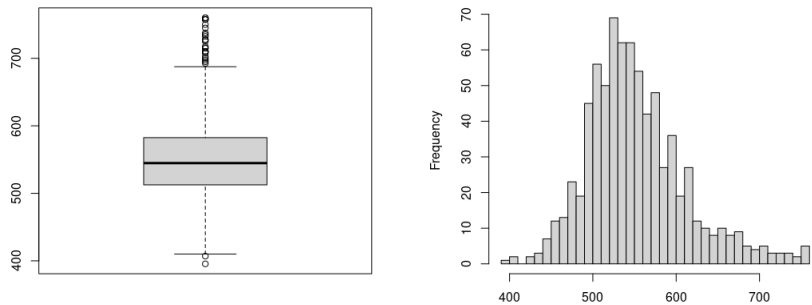
¹Available at: <https://collegescorecard.ed.gov/data>

Data Description

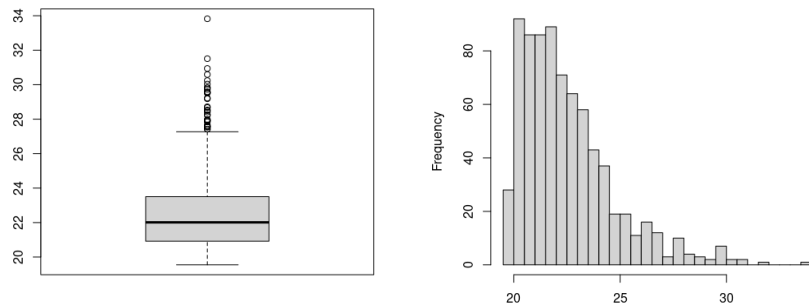
Dependent variable / Median earnings 10 yr. after entry



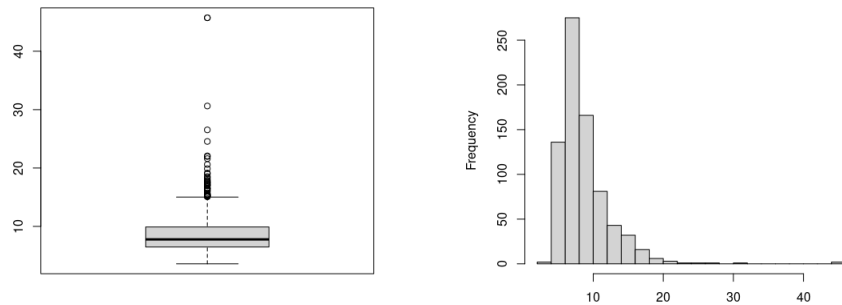
SAT score



Age of entry



Poverty rate



Analysis Problem

- Predict median earnings of alumni with university level covariates
- Multivariate regression setting

- Mathematical model:

- Y is the dependent variable
- X is the data matrix of covariates
- β is the regression coefficient vector
- ε is the residual error
- Hatted variables $\hat{Y}, \hat{\beta}$ indicate LS¹ estimates

$$\mathbf{Y} = \mathbf{X}\beta + \varepsilon$$

$$\begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} 1 & x_{11} & x_{12} & \dots & x_{1q} \\ 1 & x_{21} & x_{22} & \dots & x_{2q} \\ 1 & x_{31} & x_{32} & \dots & x_{3q} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nq} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_q \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \varepsilon_3 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

$$\hat{\mathbf{Y}} = \mathbf{X}\hat{\beta}$$

$$\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$$

Feature Selection

- Raw data had ~3000 variables → feature selection was a focus area
- Feature selection in 3 phases:
 1. Select initial set of features based on common intuition and literature
 2. Assess correlations, linear relationships, and dummy effects of categorical variables
 3. Use stepwise regression to suggest final subset of features

Feature Selection – Phase 1

Initial set of features

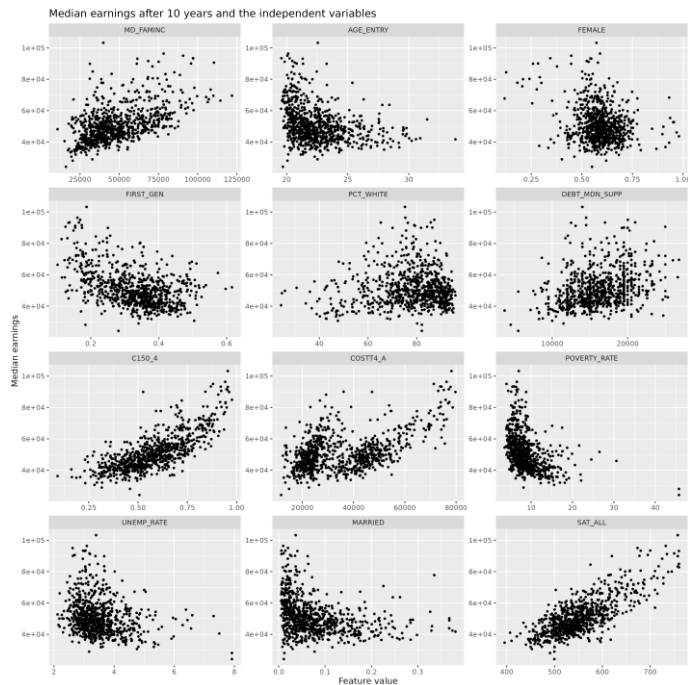
	Name	Data type	Description
1	SATVRMID	integer	Midpoint of SAT scores at the institution (critical reading)
2	SATMTMID	integer	Midpoint of SAT scores at the institution (math)
3	SATWRMID	integer	Midpoint of SAT scores at the institution (writing)
4	MD_FAMINC	double	Median family income
5	AGE_ENTRY	double	Average age of entry
6	FEMALE	double	Share of female students
7	FIRST_GEN	double	Share of first-generation students
8	PCT_WHITE	double	Percent of the population from students' zip codes that is White
9	DEBT_MDN_SUPP	integer	Median debt, suppressed for n=30
10	C150.4	double	Completion rate for first-time, full-time students
11	COSTT4.A	integer	Average cost of attendance (academic year institutions)
12	POVERTY_RATE	double	Poverty rate
13	UNEMP_RATE	double	Unemployment rate
14	MARRIED	double	Share of married students
15	VETERAN	double	Share of veteran students
16	LOCALE	categorical	Locale of institution
17	CCBASIC	categorical	Carnegie Classification – basic
18	CONTROL	categorical	Control of institution

Comments

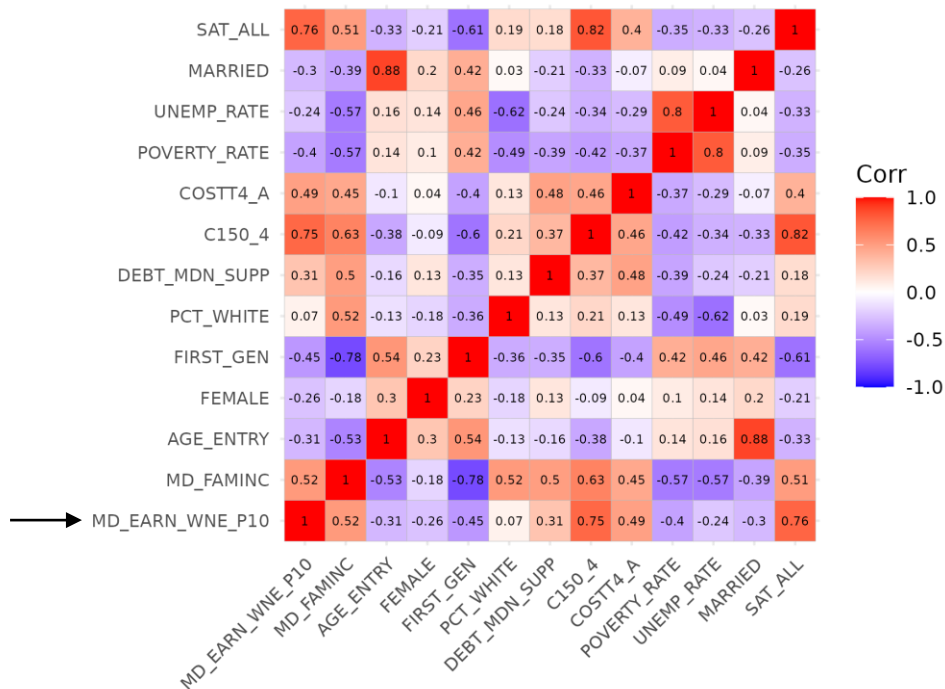
- Ability assumed to have positive effect
- Privileged background assumed to have positive effect
- Age of entry assumed to have adverse effect
- Gender split assumed to have an effect due to empirically observed gender gap
- Marriage and veteran rates assumed to have adverse affect
- Cost and debt assumed to have positive effect
- Assumption that location, control, classification has some effect

Feature Selection – Phase 2

Numerical variables

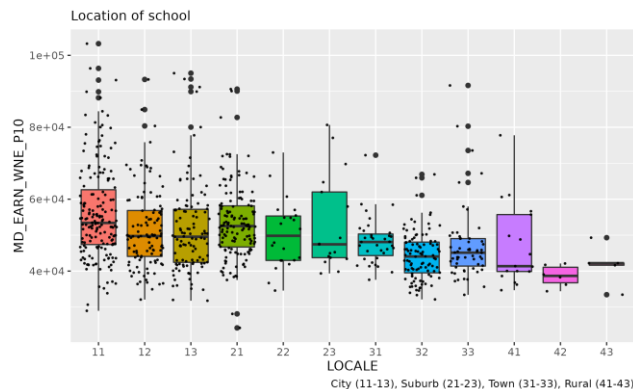


Pearson's correlation



Feature Selection – Phase 2

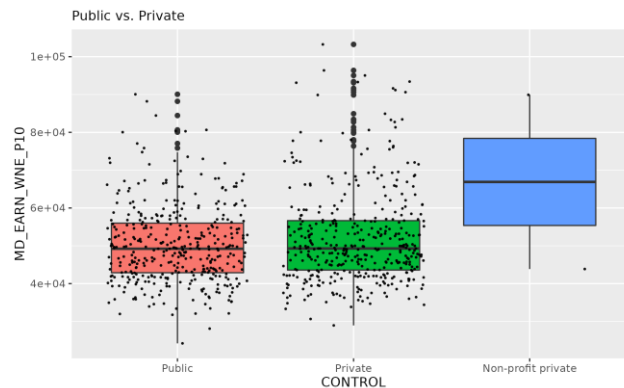
University locale



Comments

- We created a dummy variable URBAN if locale is not rural

University control



Comments

- We created a dummy variable for private institutions (for-profit and non-profit)

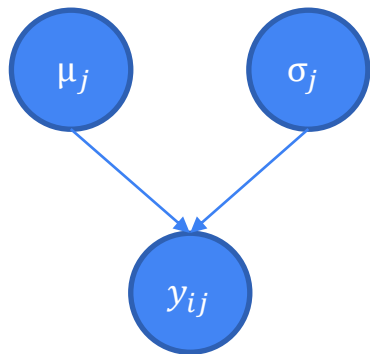
Feature Selection – Phase 3

- Lastly, we used stepwise regression with backward elimination to finetune our model
- Final model suggested by the stepwise regression model:

$$Y = \beta_{SAT_ALL}x_{SAT_ALL} + \beta_{MD_FAMINC}x_{MD_FAMINC} + \beta_{COSTT4_A}x_{COSTT4_A} \\ + \beta_{POVERTY_RATE}x_{POVERTY_RATE} + \beta_{URBAN}x_{URBAN} + \beta_{PRIVATE}x_{PRIVATE} + \varepsilon$$

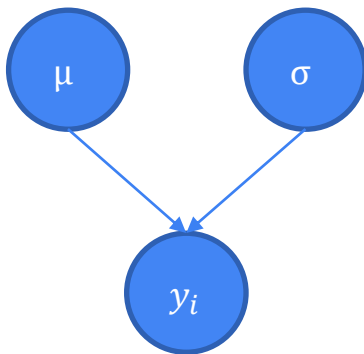
- We proceeded to Stan with this final model

Model Descriptions



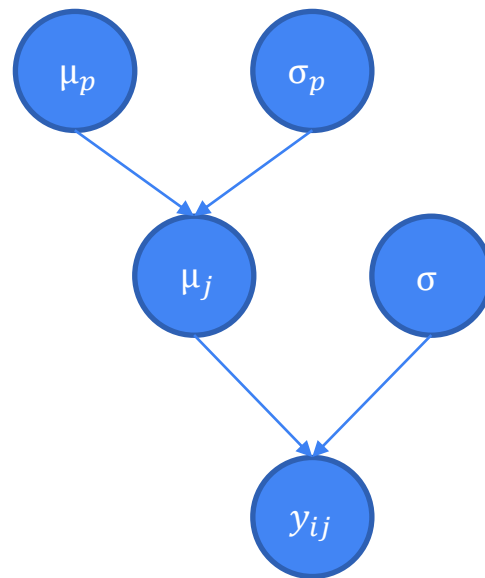
Separate Model

$$\begin{aligned} y_{ij} &| \mu_j, \sigma_j \sim N(\mu_j, \sigma_j^2), \\ \mu_j &= \alpha_j + \mathbf{X}\boldsymbol{\beta}_j \\ \alpha_j, \beta_j, \mu_j, \sigma_j &\sim N \end{aligned}$$



Pooled Model

$$\begin{aligned} y_i &| \mu, \sigma \sim N(\mu, \sigma^2), \\ \mu &= \alpha + \mathbf{X}\boldsymbol{\beta} \\ \alpha, \beta, \mu, \sigma &\sim N \end{aligned}$$



Hierarchical Model

$$\begin{aligned} y_{ij} &| \mu_j, \sigma \sim N(\mu_j, \sigma^2), \\ \mu_j &= \alpha_j + \mathbf{X}\boldsymbol{\beta}_j \\ \alpha_j, \beta_j &| \mu_p, \sigma_p \sim N(\mu_p, \sigma_p^2) \\ \mu_p, \sigma_p &\sim N \end{aligned}$$

Choice of Priors

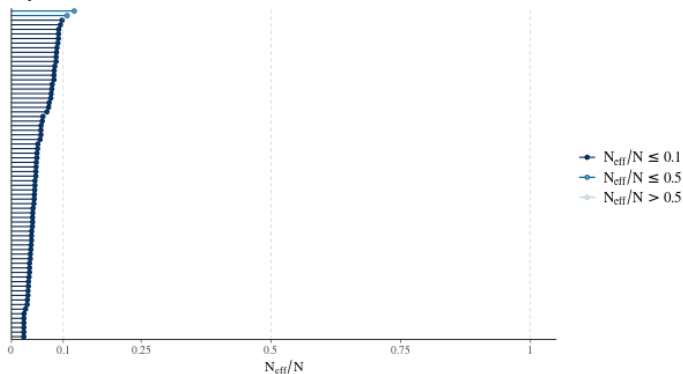
- We tried to conduct background research on the prior choices whenever possible
 - E.g., mean on SAT effect was based on national averages on income and SAT scores
- If no meaningful analysis on effect could be done, we assumed no effect on average
- We chose sufficiently large standard deviations to prevent informativeness in the priors
- We tried to avoid prior-data conflict by keeping the level of informativeness between location and scale parameters constant

Stan Setup

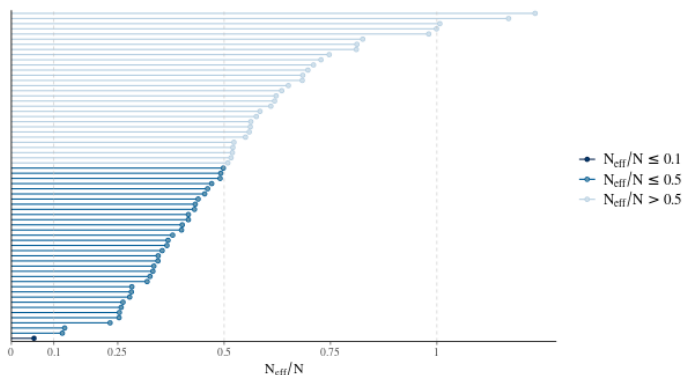
- Computing: JupyterLab
- Stan interface: cmdstanr
- Chains: 4 (default)
- Iterations per chain: 2 000 (default: 1 000 warmup, 1 000 sampling)
- Seed for RNG: 1234

Convergence Diagnostics - ESS

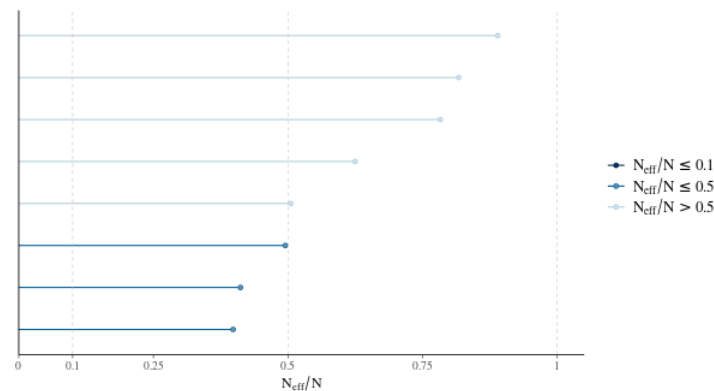
Separate model ESS



Hierarchical model ESS



Pooled model ESS

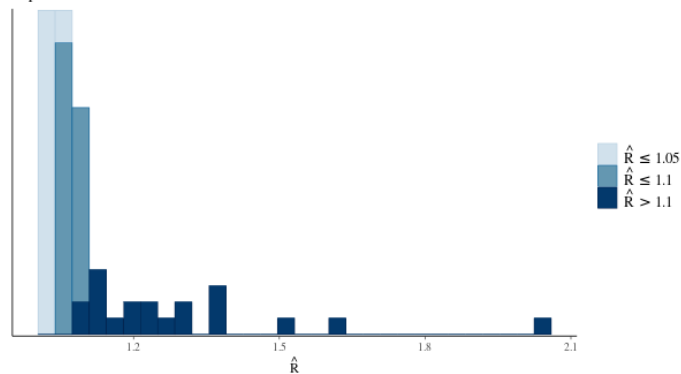


Ratios of effective sample size to total sample size

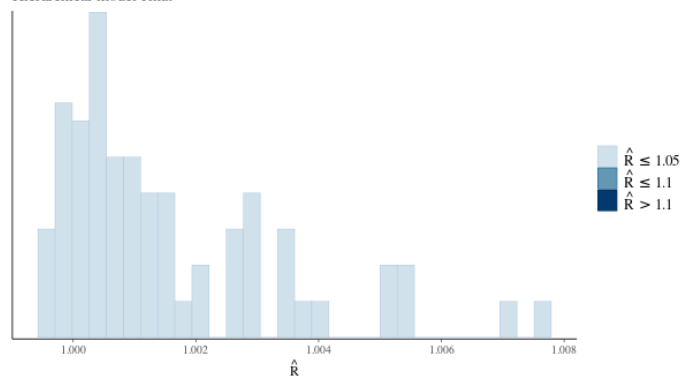
- *light*: between 0.5 and 1 (high)
- *mid*: between 0.1 and 0.5 (good)
- *dark*: below 0.1 (low)

Convergence Diagnostics – Rhat

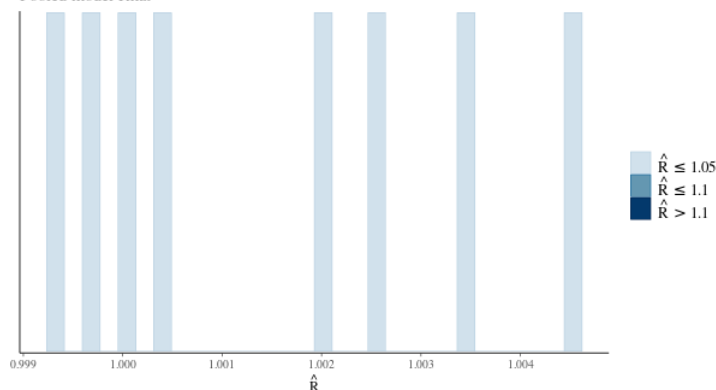
Separate model Rhat



Hierarchical model Rhat



Pooled model Rhat



Rhat values

- *light*: below 1.05 (good)
- *mid*: between 1.05 and 1.1 (ok)
- *dark*: above 1.1 (too high)

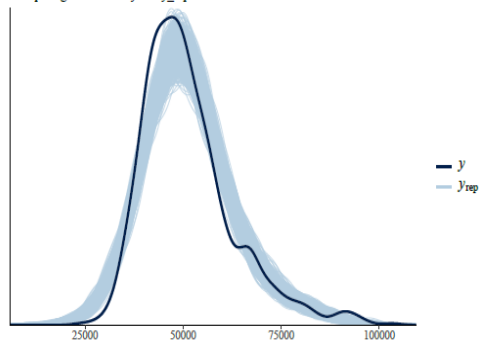
Convergence Diagnostics – HMC specific

- Treedepth:
 - For pooled and hierarchical models satisfactory
 - In separate model all transitions hit the max treedepth
- Divergences:
 - For separate and pooled models no divergent transitions
 - In hierarchical model $\approx 4\%$ divergent transitions

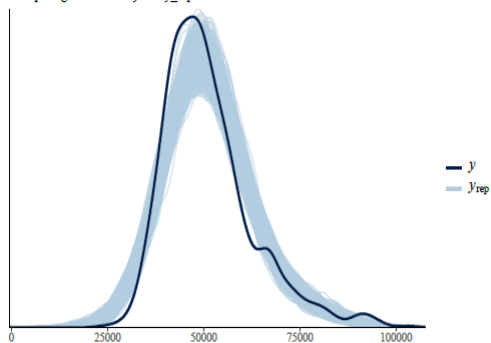
Posterior Predictive Checks

- Observed values vs. posterior predictions

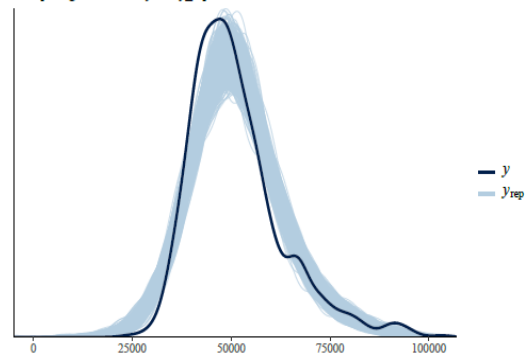
Separate model
Comparing densities of y and y_{rep}



Pooled model
Comparing densities of y and y_{rep}



Hierarchical model
Comparing densities of y and y_{rep}



Model Comparison - LOO

- Model performance order based on log pointwise predictive density (ELPD)

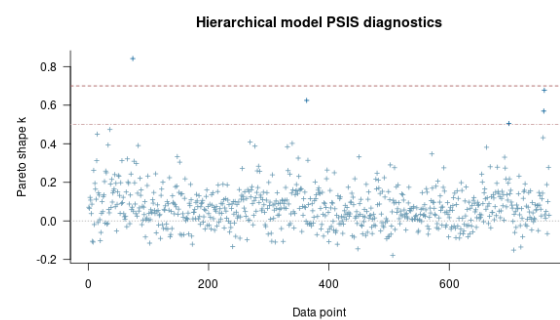
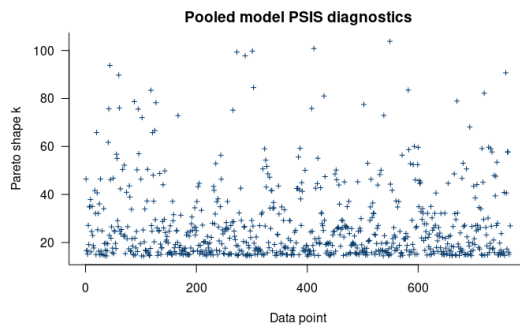
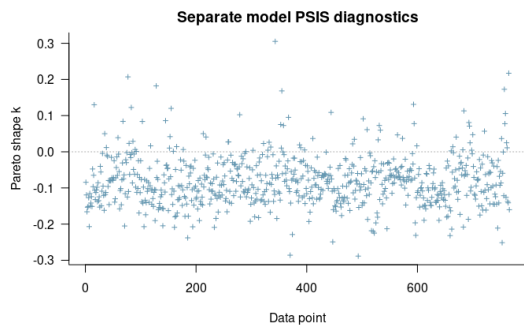
1. **Best:** Separate model (-7700)
2. Hierarchical model
3. Pooled model

	elpd_diff	se_diff
separate	0.0	0.0
hierarchical	-49.7	6.2
pooled	-7450000.0	75700.0

elpd_loo is an estimate of the expected log pointwise predictive density (ELPD). elpd_loo sums individual pointwise log predictive densities.¹

Model Comparison - \hat{k}

- PSIS-LOO estimates of the separate model can be considered reliable since Khat values < 0.5
- PSIS-LOO estimates of the pooled model are likely too optimistic (biased)¹ since Khat values $\gg 0.7$
- PSIS-LOO estimates of the hierarchical model are mostly reliable but few Khat values > 0.7



¹(Gabry, Gelman, Vehtari, 2016.)

Predictive Performance Assessment

Approach

- Train/test split with 763/30 obs.
- Fit all models on train set
- Use fitted model to predict on test set
- Compute Root Mean Squared Error (RMSE)
- Assess model output for RMSE

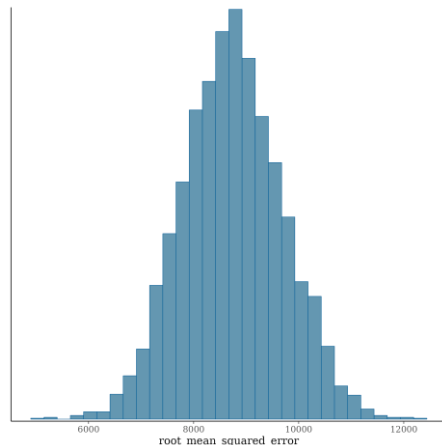
Model output

Description: df [1 x 10]

variable	mean	median	sd	mad	q5	q95	rhat	ess_bulk	ess_tail
<S3: A&id>	<S3: A&id>	<S3: A&id>	<S3: A&id>	<S3: A&id>	<S3: A&id>	<S3: A&id>	<S3: A&id>	<S3: A&id>	<S3: A&id>
root_mean_squared_error	8740.40	8731.02	952.35	948.75	7231.03	10361.85	1.00	3843	4016

1 row

MCMC posterior histogram for RMSE



→ Decent predictions as the sample st.dev. in the complete set (train+test) is \$11.6k for median earnings
1st Separate, 2nd Hierarchical, 3rd Pooled

Prior Sensitivity Analysis

- The width of priors had been a major point of discussion and uncertainty when setting original priors
- Experiment with priors of different widths:
 - Wide: 3x original prior sd
 - Narrow: 0.5x original prior sd
- Narrow priors perform better, Rhat, Khat, RMSE and ESS values improve
- For separate narrow model, all Rhat values reduced to below 1.01 from up to 1.12
- Estimates for coefficients are not affected much
- Wider priors lead to more convergence issues, worse predictive performance and lower ESS

Conclusion: Narrow priors perform better and should be used

Issues and Potential Improvements

- Narrowing down priors had a positive impact → Models could be improved by fine-tuning priors
- Pooled model had extreme \hat{k} values and low elpd → Model could be rebuilt to prevent overfitting
- Separate model took ~1h to sample in Jupyter → Model could be made more efficient computationally
- Test set was small → Could be interesting to see how model fit and predictive performance is affected by altering the size of the test set

Conclusions

Project: Bayesian multivariate linear regression on median alumni earnings 10 years after entry in 793 universities in the U.S.

- Clear link found between education and future earnings; the model has decent predictive power in terms of RMSE

Three models: Pooled, Hierarchical and Separate

- Pooled model is prone to overfitting
- Separate model has issues with convergence and is slow to run, otherwise good performance, preferred by elpd and has lowest RMSE
- Hierarchical model has strong performance overall and doesn't have major issues

Prior sensitivity analysis shows that model performance can be improved by using narrower priors

Contact: Niko Miller: niko.miller@aalto.fi, Akseli Manninen: akseli.manninen@aalto.fi, Santeri Löppönen: santeri.lopponen@aalto.fi

References

- Card, D. (1999). THE CAUSAL EFFECT OF EDUCATION ON EARNINGS. Wolla, S. A., & Sullivan, J. (2017). Education, Income, and Wealth. <https://fred.stlouisfed.org/graph/?g=7yKu>.
- Stan. (n.d) Effective Sample Size [Website]. Retrieved from: https://mc-stan.org/docs/2_19/referencemanual/effective-sample-size-section.html
- Gabry, Gelman, Vehtari. (2016) Practical Bayesian model evaluation using leave-one-out cross-validation and WAIC.