Model assessment, selection and averaging

Part 1: cross-validation Part 2: projection predictive inference

> Aki Vehtari Aalto University, Finland

Slides and extra material at avehtari.github.io/modelselection/

Predicting concrete quality



Predicting cancer recurrence



Model assessment, comparison, selection and averaging

• Modeling complex phenomena with models that are much simpler than the nature (*M*-open)

Model assessment, comparison, selection and averaging

- Modeling complex phenomena with models that are much simpler than the nature (*M*-open)
- Decision theoretical approch in spirit of
 - Lindley, Box, Rubin, Bernardo & Smith, etc.
 - see Vehtari and Ojanen (2012) for more details and references

Stan and 100 package

Computed from 4000 by 20 log-likelihood matrix

Estimate SE elpd_loo -29.5 3.3 p_loo 2.7 1.0

Monte Carlo SE of elpd_loo is 0.1.

Pareto k diagnostic values:

		Count	Pct.	Min .	n_eff
(-lnf, 0.5]	(good)	18	90.0%	899	
(0.5, 0.7]	(ok)	2	10.0%	459	
(0.7, 1]	(bad)	0	0.0%	<na></na>	
(1, Inf)	(very bad)	0	0.0%	<na></na>	

All Pareto k estimates are ok (k < 0.7). See help('pareto-k-diagnostic') for details.

Model comparison: (negative 'elpd_diff' favors 1st model, positive favors 2nd) elpd_diff se -0.2 0.1

Outline

- What is cross-validation
 - Leave-one-out cross-validation (elpd_loo, p_loo)
 - Uncertainty in LOO (SE)
- When is cross-validation applicable?
 - data generating mechanisms and prediction tasks
 - leave-many-out cross-validation
- Fast cross-validation
 - PSIS and diagnostics in loo package (Pareto k, n_eff, Monte Carlo SE)
 - K-fold cross-validation
- Related methods (WAIC, *IC, BF)
- Model comparison and selection (elpd_diff, se)
- Model averaging with Bayesian stacking

Outline

- What is cross-validation
 - Leave-one-out cross-validation (elpd_loo, p_loo)
 - Uncertainty in LOO (SE)
- When is cross-validation applicable?
 - data generating mechanisms and prediction tasks
 - leave-many-out cross-validation
- Fast cross-validation
 - PSIS and diagnostics in loo package (Pareto k, n_eff, Monte Carlo SE)
 - K-fold cross-validation
- Related methods (WAIC, *IC, BF)
- Model comparison and selection (elpd_diff, se)
- Model averaging with Bayesian stacking
- Part 2: Projective Inference in High-dimensional Problems: Prediction and Feature Selection



















 $p(\tilde{y}|\tilde{x} = 18, x, y) = \int p(\tilde{y}|\tilde{x} = 18, \theta) p(\theta|x, y) d\theta$











 $y_{18} - E[p(\tilde{y}|\tilde{x} = 18, x_{-18}, y_{-18})]$



 $y_{18} - E[p(\tilde{y}|\tilde{x} = 18, x_{-18}, y_{-18})]$

Can be use to compute, e.g., RMSE, R², 90% error



 $y_{18} - E[p(\tilde{y}|\tilde{x} = 18, x_{-18}, y_{-18})]$

Can be use to compute, e.g., RMSE, R², 90% error

See LOO-R² at avehtari.github.io/bayes_R2/bayes_R2.html





 $p(\tilde{y}|\tilde{x} = 18, x_{-18}, y_{-18}) = \int p(\tilde{y}|\tilde{x} = 18, \theta) p(\theta|x_{-18}, y_{-18}) d\theta$





 $p(\tilde{y} = y_{18} | \tilde{x} = 18, x, y) \approx 0.07$



 $p(\tilde{y} = y_{18} | \tilde{x} = 18, x, y) \approx 0.07$ $p(\tilde{y} = y_{18} | \tilde{x} = 18, x_{-18}, y_{-18}) \approx 0.03$







 $\sum_{i=1}^{20} \log p(y_i | x_i, x_{-i}, y_{-i}) \approx -29.5$





elpd_loo = $\sum_{i=1}^{20} \log p(y_i | x_i, x_{-i}, y_{-i}) \approx -29.5$ unbiased estimate of log posterior pred. density for new data





 $p_loo = lpd - elpd_loo \approx 2.7$




see Vehtari, Gelman & Gabry (2017a) and Vehtari & Ojanen (2012) for more



LOO is ok for fixed / designed x. SE is uncertainty about y|x.

see Vehtari & Ojanen (2012) and andrewgelman.com/2018/08/03/ loo-cross-validation-approaches-valid/

Distribution for x



LOO is ok for random x. SE is uncertainty about y|x and x.

see Vehtari & Ojanen (2012) and andrewgelman.com/2018/08/03/ loo-cross-validation-approaches-valid/

Distribution for x



LOO is ok for random x. SE is uncertainty about y|x and x.

Covariate shift can be handled with importance weighting or modelling see Vehtari & Ojanen (2012) and andrewgelman.com/2018/08/03/ loo-cross-validation-approaches-valid/

100 package

Computed from 4000 by 20 log-likelihood matrix

	Estimate	SE
elpd_loo	-29.5	3.3
p_loo	2.7	1.0

Monte Carlo SE of elpd_loo is 0.1.

Pareto k diagnostic values:

			Count	Pct.	Min.	n_eff
(-Inf,	0.5]	(good)	18	90.0%	899	
(0.5,	0.7]	(ok)	2	10.0%	459	
(0.7	, 1]	(bad)	0	0.0%	<na></na>	
(1,	lnf)	(very bad)	0	0.0%	<na></na>	

All Pareto k estimates are ok (k < 0.7). See help('pareto-k-diagnostic') for details.









Extrapolation is more difficult



Can LOO or other cross-validation be used with time series?



Leave-one-out cross-validation is ok for assessing conditional model



leave-future-out cross-validation is better for predicting future Bürkner, Gabry and Vehtari (2019)



m-step-ahead cross-validation is better for predicting further future

Bürkner, Gabry and Vehtari (2019)



Can LOO or other cross-validation be used with hierarchical data?











Summary of data generating mechanisms and prediction tasks

- You have to make some assumptions on data generating mechanism
- Use the knowledge of the prediction task if available
- Cross-validation can be used to analyse different parts, even if there is no clear prediction task

see Vehtari & Ojanen (2012) and andrewgelman.com/2018/08/03/ loo-cross-validation-approaches-valid/

Fast cross-validation

- Pareto smoothed importance sampling LOO (PSIS-LOO)
- K-fold cross-validation

see Vehtari, Gelman & Gabry (2017a) and mc-stan.org/loo/





 $\theta^{(s)} \sim p(\theta|x, y)$



 $\theta^{(s)} \sim p(\theta|x, y), \quad p(\tilde{y}|\tilde{x}, x, y) \approx \frac{1}{S} \sum_{s=1}^{S} p(\tilde{y}|\tilde{x}, \theta^{(s)})$



 $\theta^{(s)} \sim p(\theta|x, y), \quad p(\tilde{y}|\tilde{x}, x, y) \approx \frac{1}{S} \sum_{s=1}^{S} p(\tilde{y}|\tilde{x}, \theta^{(s)})$



$$egin{aligned} & heta^{(s)} \sim p(heta|x,y) \ & r_i^{(s)} = p(heta^{(s)}|x_{-i},y_{-i})/p(heta^{(s)}|x,y) \end{aligned}$$



$$\begin{split} \theta^{(s)} &\sim p(\theta|x, y) \\ r_i^{(s)} &= p(\theta^{(s)}|x_{-i}, y_{-i}) / p(\theta^{(s)}|x, y) \propto 1 / p(y_i|x_i, \theta^{(s)}) \end{split}$$



$$\begin{split} \theta^{(s)} &\sim p(\theta|x, y) \\ r_i^{(s)} &= p(\theta^{(s)}|x_{-i}, y_{-i})/p(\theta^{(s)}|x, y) \propto 1/p(y_i|x_i, \theta^{(s)}) \\ \log(1/p(y_i|x_i, \theta^{(s)})) &= -\log_\text{lik}[i] \end{split}$$



$$\begin{split} \theta^{(s)} &\sim p(\theta|x,y) \\ r_i^{(s)} &= p(\theta^{(s)}|x_{-i},y_{-i})/p(\theta^{(s)}|x,y) \propto 1/p(y_i|x_i,\theta^{(s)}) \end{split}$$



$$egin{aligned} & \theta^{(s)} \sim p(\theta|x,y) \ & r_i^{(s)} = p(heta^{(s)}|x_{-i},y_{-i})/p(heta^{(s)}|x,y) \propto 1/p(y_i|x_i, heta^{(s)}) \ & p(y_i|x_i,x_{-i},y_{-i}) pprox \sum_{s=1}^S [w_i^{(s)}p(y_i|x_i, heta^{(s)})] \end{aligned}$$



$$\begin{split} \theta^{(s)} &\sim p(\theta|x, y) \\ r_i^{(s)} &= p(\theta^{(s)}|x_{-i}, y_{-i}) / p(\theta^{(s)}|x, y) \propto 1 / p(y_i|x_i, \theta^{(s)}) \\ p(y_i|x_i, x_{-i}, y_{-i}) &\approx \sum_{s=1}^{S} [w_i^{(s)} p(y_i|x_i, \theta^{(s)})], \text{ where } w \leftarrow \mathsf{PSIS}(r) \end{split}$$



4000 importance weights for leave-18th-out



4000 importance weights for leave-18th-out



4000 importance weights for leave-18th-out



 Pareto k estimates the tail shape which estimates the pre-asymptotic convergence rate of PSIS. Less than 0.7 is ok (Vehtari, Simpson, Gelman, Yao & Gabry (2019))






100 package

Computed from 4000 by 20 log-likelihood matrix

Estimate SE elpd_loo -29.5 3.3 p_loo 2.7 1.0

Monte Carlo SE of elpd_loo is 0.1.

Pareto k diagnostic values:

		Count	Pct.	Min.	n_eff
(-lnf, 0.5]	(good)	18	90.0%	899	
(0.5, 0.7]	(ok)	2	10.0%	459	
(0.7, 1]	(bad)	0	0.0%	<na></na>	
(1, Inf)	(very bad)	0	0.0%	<na></na>	

All Pareto k estimates are ok (k < 0.7). See help('pareto-k-diagnostic') for details.

see more in Vehtari, Gelman & Gabry (2017b) and Vehtari, Simpson, Gelman, Yao & Gabry (2019)

Stan code

$$\log(r_i^{(s)}) = \log(1/\rho(y_i|x_i,\theta^{(s)})) = -\log_\mathsf{lik}[i]$$

Stan code

$$\log(r_i^{(s)}) = \log(1/p(y_i|x_i,\theta^{(s)})) = -\log_\text{lik}[i]$$

```
model {
    alpha ~ normal(pmualpha, psalpha);
    beta ~ normal(pmubeta, psbeta);
    y ~ normal(mu, sigma);
}
generated quantities {
    vector[N] log_lik;
    for (i in 1:N)
        log_lik[i] = normal_lpdf(y[i] | mu[i], sigma);
}
```

Stan code

$$\log(r_i^{(s)}) = \log(1/p(y_i|x_i,\theta^{(s)})) = -\log_\text{lik}[i]$$

```
model {
    alpha ~ normal(pmualpha, psalpha);
    beta ~ normal(pmubeta, psbeta);
    y ~ normal(mu, sigma);
}
generated quantities {
    vector[N] log_lik;
    for (i in 1:N)
        log_lik[i] = normal_lpdf(y[i] | mu[i], sigma);
}
```

RStanARM and BRMS compute log_lik by default

- PSIS-LOO for hierarchical models
 - leave-one-group out is challenging for PSIS-LOO see Merkel, Furr and Rabe-Hesketh (2018) for an approach using quadrature integration

- PSIS-LOO for hierarchical models
 - leave-one-group out is challenging for PSIS-LOO see Merkel, Furr and Rabe-Hesketh (2018) for an approach using quadrature integration
- PSIS-LOO for non-factorizable models
 - Bürkner, Gabry and Vehtari (2018)

- PSIS-LOO for hierarchical models
 - leave-one-group out is challenging for PSIS-LOO see Merkel, Furr and Rabe-Hesketh (2018) for an approach using quadrature integration
- PSIS-LOO for non-factorizable models
 - Bürkner, Gabry and Vehtari (2018)
- PSIS-LOO for time series
 - Approximate leave-future-out cross-validation Bürkner, Gabry and Vehtari (2019)

- PSIS-LOO for hierarchical models
 - leave-one-group out is challenging for PSIS-LOO see Merkel, Furr and Rabe-Hesketh (2018) for an approach using quadrature integration
- PSIS-LOO for non-factorizable models
 - Bürkner, Gabry and Vehtari (2018)
- PSIS-LOO for time series
 - Approximate leave-future-out cross-validation Bürkner, Gabry and Vehtari (2019)
- PSIS-LOO for optimizing, Laplace, ADVI and big data
 - Magnusson, Andersen, Jonasson and Vehtari (2019a) and Magnusson, Andersen, Jonasson and Vehtari (2019b)

- PSIS-LOO for hierarchical models
 - leave-one-group out is challenging for PSIS-LOO see Merkel, Furr and Rabe-Hesketh (2018) for an approach using quadrature integration
- PSIS-LOO for non-factorizable models
 - Bürkner, Gabry and Vehtari (2018)
- PSIS-LOO for time series
 - Approximate leave-future-out cross-validation Bürkner, Gabry and Vehtari (2019)
- PSIS-LOO for optimizing, Laplace, ADVI and big data
 - Magnusson, Andersen, Jonasson and Vehtari (2019a) and Magnusson, Andersen, Jonasson and Vehtari (2019b)
- Pushing the limits of importance sampling through iterative moment matching
 - Paanenen, Piironen, Bürkner, Vehtari (2019)







Bürkner. Gabrv and Vehtari (2019)



Bürkner, Gabry and Vehtari (2019)

K-fold cross-validation

- K-fold cross-validation can approximate LOO
 - all uses for LOO
- K-fold cross-validation can be used for hierarchical models
 - good for leave-one-group-out
- K-fold cross-validation can be used for time series
 - with leave-block-out













kfold_split_stratified()

see Vehtari, Gelman & Gabry (2017a)

WAIC has same assumptions as LOO

see Vehtari, Gelman & Gabry (2017a)

- WAIC has same assumptions as LOO
- PSIS-LOO is more accurate

see Vehtari, Gelman & Gabry (2017a)

- WAIC has same assumptions as LOO
- PSIS-LOO is more accurate
- PSIS-LOO has much better diagnostics

see Vehtari, Gelman & Gabry (2017a)

- WAIC has same assumptions as LOO
- PSIS-LOO is more accurate
- PSIS-LOO has much better diagnostics
- LOO makes the prediction assumption more clear, which helps if K-fold-CV is needed instead

see Vehtari, Gelman & Gabry (2017a)

- WAIC has same assumptions as LOO
- PSIS-LOO is more accurate
- PSIS-LOO has much better diagnostics
- LOO makes the prediction assumption more clear, which helps if K-fold-CV is needed instead
- Multiplying by -2 doesn't give any benefit (Watanabe didn't multiply by -2)

see Vehtari, Gelman & Gabry (2017a)

*IC

- AIC uses maximum likelihood estimate for prediction
- DIC uses posterior mean for prediction
- BIC is an approximation for marginal likelihood
- TIC, NIC, RIC, PIC, BPIC, QIC, AICc, ...

• Like leave-future-out 1-step-ahead corss-validation but starting with 0 observations

• Like leave-future-out 1-step-ahead corss-validation but starting with 0 observations



- Like leave-future-out 1-step-ahead corss-validation but starting with 0 observations
 - which makes it very sensitive to prior



- Like leave-future-out 1-step-ahead corss-validation but starting with 0 observations
 - which makes it very sensitive to prior and
 - unstable in case of misspecified models



- Like leave-future-out 1-step-ahead corss-validation but starting with 0 observations
 - which makes it very sensitive to prior and
 - unstable in case of misspecified models also asymptotically



Cross-validation for model assessment

- CV is good for model assessment when application specific utility/cost functions are used
 - e.g. 90% absolute error
- CV is good for model assessment when application specific utility/cost functions are used
 - e.g. 90% absolute error
- Also useful in model checking in similar way as posterior predictive checking (PPC)
 - model misspecification diagnostics (e.g. Pareto-k and p_loo)
 - checking calibration of leave-one-out predictive posteriors (ppc_loo_pit in bayesplot)

see demos avehtari.github.io/modelselection/

help loo-glossary

If Pareto- $\hat{k} > 0.7$ and

- If p_loo « p then the model is likely to be misspecified. PPC likely to also detect the problem. Try overdispersed model.
- If p_loo N/5, it is likely that the model so flexible that PSIS fails. Try k-fold-CV.
- If p_loo > p, then the model is likely to be badly misspecified. If p«N, then PPC also likely to detect the problem.

• If p_loo > p, then the model is likely to be badly misspecified. If p«N, then PPC also likely to detect the problem.

Roaches with Poisson model

Computed from 4000 by 262 log-likelihood matrix

	Estimate	SE
elpd_loo	-6244.8	727.0
p_loo	292.7	72.9

Monte Carlo SE of elpd_loo is NA.

Pareto k diagnostic values:

		Count	Pct.	Min.	n_eff
(-lnf, 0.5]	(good)	240	91.6%	205	
(0.5, 0.7]	(ok)	7	2.7%	48	
(0.7, 1]	(bad)	8	3.1%	7	
(1, Inf)	(very bad)	7	2.7%	1	

Roaches with Neg-bin model

Computed from 4000 by 262 log-likelihood matrix

	Estimate	SE
elpd_loo	-895.9	37.8
p_loo	6.9	2.7

Monte Carlo SE of elpd_loo is NA.

Pareto k diagnostic values:

		Count	Pct.	Min.	n_eff
(-Inf, 0.5]	(good)	261	99.6%	740	
(0.5, 0.7]	(ok)	0	0.0%	NA	
(0.7, 1]	(bad)	1	0.4%	27	
(1, Inf)	(very bad)	0	0.0%	NA	

If many high k̂, p_loo N/5, it is likely that the model so flexible that PSIS fails. Try k-fold-CV.

Roaches with Poisson and random effect for each observation Computed from 40000 by 262 log-likelihood matrix

	Estimate	SE
elpd_loo	-653.5	24.5
p_loo	189.6	4.7

Monte Carlo SE of elpd_loo is NA.

Pareto k diagnostic values:

		Count	Pct.	Min. n_eff
(-Inf, 0.5]	(good)	0	0.0%	<na></na>
(0.5, 0.7]	(ok)	27	10.3%	762
(0.7, 1]	(bad)	213	81.3%	18
(1, Inf)	(very bad)	22	8.4%	7

Aki.Vehtari@aalto.fi - @avehtari

· Posterior predictive checking is often sufficient



Predicting the yields of mesquite bushes.

Gelman, Hill & Vehtari (2019): Regression and Other Stories, Chapter 11.

Posterior predictive checking is often sufficient



Predicting the yields of mesquite bushes.

Gelman, Hill & Vehtari (2019): Regression and Other Stories, Chapter 11.

- BDA3, Chapter 6
- Gabry, Simpson, Vehtari, Betancourt, Gelman (2018). Visualization in Bayesian workflow. JRSS A, preprint arXiv:1709.01449
- mc-stan.org/bayesplot/articles/graphical-ppcs.html
- betanalpha.github.io/assets/case_studies/principled_bayesian_ workflow.html

Model comparison

- "A popular hypothesis has it that primates with larger brains produce more energetic milk, so that brains can grow quickly" (from Statistical Rethinking)
 - Model 1: formula = kcal.per.g \sim neocortex
 - Model 2: formula = kcal.per.g ~ neocortex + log(mass)

mc-stan.org/loo/articles/loo2-example.html

Aki.Vehtari@aalto.fi - @avehtari



Aki.Vehtari@aalto.fi - @avehtari







• se for elpd_diff is underestimate

- se for elpd_diff is underestimate
- good se estimate may require n>100, or more if model is misspecified

- se for elpd_diff is underestimate
- good se estimate may require n>100, or more if model is misspecified
- normal approximation fails if models are very similar

- se for elpd_diff is underestimate
- good se estimate may require n>100, or more if model is misspecified
- normal approximation fails if models are very similar
- not asymptotically model selection consistent
 - but the selected model has the same predictive performance

Roaches

fit_1 <- stan_glm(y ~ roach100 + treatment + senior, family=neg_binomi
offset=log(exposure2), data=roaches, seed=SEED, refresh=0)</pre>

	Median	MAD_SD
(Intercept)	2.8	0.2
roach100	1.3	0.2
treatment	-0.8	0.2
senior	-0.3	0.3

fit_0 <- stan_glm(y ~ roach100 + senior, family=neg_binomial_2, offset=log(exposure2), data=roaches, seed=SEED, refresh=0)

Posterior predictive checking is often sufficient



Predicting the yields of mesquite bushes.

Gelman, Hill & Vehtari (2019): Regression and Other Stories, Chapter 11.

- BDA3, Chapter 6
- Gabry, Simpson, Vehtari, Betancourt, Gelman (2018). Visualization in Bayesian workflow. JRSS A, preprint arXiv:1709.01449
- mc-stan.org/bayesplot/articles/graphical-ppcs.html
- betanalpha.github.io/assets/case_studies/principled_bayesian_ workflow.html

• For some very simple cases you may assume that true model is included in the list of models considered (*M*-closed)

- For some very simple cases you may assume that true model is included in the list of models considered (*M*-closed)
 - see predictive model selection in *M*-closed case by San Martini and Spezzaferri (1984)

- For some very simple cases you may assume that true model is included in the list of models considered (*M*-closed)
 - see predictive model selection in *M*-closed case by San Martini and Spezzaferri (1984)
 - but you should not force your design of experiment or analysis to stay in the simplified world

- For some very simple cases you may assume that true model is included in the list of models considered (*M*-closed)
 - see predictive model selection in *M*-closed case by San Martini and Spezzaferri (1984)
 - but you should not force your design of experiment or analysis to stay in the simplified world
- In nested case, often easier and more accurate to analyse posterior distribution of more complex model directly avehtari.github.io/modelselection/betablockers.html

- Continuous expansion including all models?
 - and then analyse the posterior distribution directly avehtari.github.io/modelselection/betablockers.html
 - sparse priors like regularized horseshoe prior instead of variable selection

video, refs and demos at avehtari.github.io/modelselection/

- Continuous expansion including all models?
 - and then analyse the posterior distribution directly avehtari.github.io/modelselection/betablockers.html
 - sparse priors like regularized horseshoe prior instead of variable selection video, refs and demos at avehtari.github.io/modelselection/
- Model averaging with BMA or Bayesian stacking?
 - mc-stan.org/loo/articles/loo2-example.html

- Continuous expansion including all models?
 - and then analyse the posterior distribution directly avehtari.github.io/modelselection/betablockers.html
 - sparse priors like regularized horseshoe prior instead of variable selection video, refs and demos at avehtari.github.io/modelselection/
- Model averaging with BMA or Bayesian stacking?
 - mc-stan.org/loo/articles/loo2-example.html
- In a nested case choose simpler if assuming some cost for extra parts?

andrewgelman.com/2018/07/26/

parsimonious-principle-vs-integration-uncertainties/

- Continuous expansion including all models?
 - and then analyse the posterior distribution directly avehtari.github.io/modelselection/betablockers.html
 - sparse priors like regularized horseshoe prior instead of variable selection video, refs and demos at avehtari.github.io/modelselection/
- Model averaging with BMA or Bayesian stacking?
 - mc-stan.org/loo/articles/loo2-example.html
- In a nested case choose simpler if assuming some cost for extra parts? andrewgelman.com/2018/07/26/

parsimonious-principle-vs-integration-uncertainties/

 In a nested case choose more complex if you want to take into account all the uncertainties? andrewgelman.com/2018/07/26/

parsimonious-principle-vs-integration-uncertainties/

Aki.Vehtari@aalto.fi - @avehtari

• Consider the model averaging as a decision problem with aim of maximizing the predictive performance

- Consider the model averaging as a decision problem with aim of maximizing the predictive performance
- Maximize the scoring rule of the predictive distribution for future y

$$\max_{w} S\Big(\sum_{k=1}^{K} w_k p(\tilde{y}|x, y, M_k), p_t(\tilde{y})\Big),$$

- Consider the model averaging as a decision problem with aim of maximizing the predictive performance
- Maximize the scoring rule of the predictive distribution for future y

$$\max_{w} S\Big(\sum_{k=1}^{K} w_k p(\tilde{y}|x, y, M_k), p_t(\tilde{y})\Big),$$

• As we don't know $p_t(\tilde{y})$, we approximate with LOO

- Consider the model averaging as a decision problem with aim of maximizing the predictive performance
- Maximize the scoring rule of the predictive distribution for future y

$$\max_{w} S\Big(\sum_{k=1}^{K} w_k p(\tilde{y}|x, y, M_k), p_t(\tilde{y})\Big),$$

- As we don't know $p_t(\tilde{y})$, we approximate with LOO
- We define the stacking weights as the solution to the following optimization problem:

$$\max_{w} \frac{1}{n} \sum_{i=1}^{n} S\left(\sum_{k=1}^{K} w_{k} \hat{p}(y_{i} | x_{-i}, y_{-i}, M_{k})\right),$$

s.t. $w_{k} \ge 0, \quad \sum_{k=1}^{K} w_{k} = 1.$

• The combined estimation of the predictive density is

$$\hat{p}(\tilde{y}|x,y) = \sum_{k=1}^{K} \hat{w}_k p(\tilde{y}|x,y,M_k)$$

• The combined estimation of the predictive density is

$$\hat{p}(\tilde{y}|x,y) = \sum_{k=1}^{K} \hat{w}_k p(\tilde{y}|x,y,M_k)$$

 When using log-score (corresponding to Kullback-Leibler divergence), we call this stacking of predictive distributions:

$$\max_{w} \frac{1}{n} \sum_{i=1}^{n} \log \sum_{k=1}^{K} w_{k} p(y_{i} | x_{-i}, y_{-i}, M_{k}),$$

s.t. $w_{k} \ge 0, \quad \sum_{k=1}^{K} w_{k} = 1$

• The combined estimation of the predictive density is

$$\hat{p}(\tilde{y}|x,y) = \sum_{k=1}^{K} \hat{w}_k p(\tilde{y}|x,y,M_k)$$

 When using log-score (corresponding to Kullback-Leibler divergence), we call this stacking of predictive distributions:

$$\max_{w} \frac{1}{n} \sum_{i=1}^{n} \log \sum_{k=1}^{K} w_{k} p(y_{i} | x_{-i}, y_{-i}, M_{k}),$$

s.t. $w_{k} \ge 0, \quad \sum_{k=1}^{K} w_{k} = 1$

• We can approximate $p(y_i|x_{-i}, y_{-i}, M_k)$ with PSIS-LOO

• The combined estimation of the predictive density is

$$\hat{p}(\tilde{y}|x,y) = \sum_{k=1}^{K} \hat{w}_k p(\tilde{y}|x,y,M_k)$$

 When using log-score (corresponding to Kullback-Leibler divergence), we call this stacking of predictive distributions:

$$\max_{w} \frac{1}{n} \sum_{i=1}^{n} \log \sum_{k=1}^{K} w_{k} p(y_{i} | x_{-i}, y_{-i}, M_{k}),$$

s.t. $w_{k} \ge 0, \quad \sum_{k=1}^{K} w_{k} = 1$

- We can approximate $p(y_i|x_{-i}, y_{-i}, M_k)$ with PSIS-LOO
- Other cross-validation structures can be used, too
Non-linear model example



Non-linear model example



Bayesian stacking

- In *M*-open case works better than BMA
- In *M*-closed case can have a better small sample performance than BMA

Bayesian stacking

- In *M*-open case works better than BMA
- In *M*-closed case can have a better small sample performance than BMA
- Should be used only for model averaging
 - you may drop models with 0 weights
 - you shouldn't choose the model with largest weight unless it's 1
- Yao, Vehtari, Simpson, & Gelman (2018)

Cross-validation and model selection

- Cross-validation can be used for model selection if
 - small number of models
 - the difference between models is clear

Cross-validation and model selection

- Cross-validation can be used for model selection if
 - small number of models
 - the difference between models is clear
- Do not use cross-validation to choose from a large set of models
 - selection process leads to overfitting

Cross-validation and model selection

- Cross-validation can be used for model selection if
 - small number of models
 - the difference between models is clear
- Do not use cross-validation to choose from a large set of models
 - selection process leads to overfitting
- Overfitting in selection process is not unique for cross-validation

Selection induced bias and overfitting

- Selection induced bias in cross-validation
 - same data is used to assess the performance and make the selection
 - the selected model fits more to the data
 - the CV estimate for the selected model is biased
 - recognised already, e.g., by Stone (1974)

Selection induced bias and overfitting

- Selection induced bias in cross-validation
 - same data is used to assess the performance and make the selection
 - the selected model fits more to the data
 - the CV estimate for the selected model is biased
 - recognised already, e.g., by Stone (1974)
- Performance of the selection process itself can be assessed using two level cross-validation, but it does not help choosing better models

Selection induced bias and overfitting

- Selection induced bias in cross-validation
 - same data is used to assess the performance and make the selection
 - the selected model fits more to the data
 - the CV estimate for the selected model is biased
 - recognised already, e.g., by Stone (1974)
- Performance of the selection process itself can be assessed using two level cross-validation, but it does not help choosing better models
- Bigger problem if there is a large number of models as in covariate selection

Selection induced bias in variable selection



Selection induced bias in variable selection



Aki.Vehtari@aalto.fi - @avehtari

- It's good to think predictions of observables, because observables are the only ones we can observe
- Cross-validation can simulate predicting and observing new data
- Cross-validation is good if you don't trust your model
- Different variants of cross-validation are useful in different scenarios
- Cross-validation has high variance, and **if** you trust your model you can beat cross-validation in accuracy

- It's good to think predictions of observables, because observables are the only ones we can observe
- Cross-validation can simulate predicting and observing new data
- Cross-validation is good if you don't trust your model
- Different variants of cross-validation are useful in different scenarios
- Cross-validation has high variance, and **if** you trust your model you can beat cross-validation in accuracy

- It's good to think predictions of observables, because observables are the only ones we can observe
- Cross-validation can simulate predicting and observing new data
- Cross-validation is good if you don't trust your model
- Different variants of cross-validation are useful in different scenarios
- Cross-validation has high variance, and **if** you trust your model you can beat cross-validation in accuracy

- It's good to think predictions of observables, because observables are the only ones we can observe
- Cross-validation can simulate predicting and observing new data
- Cross-validation is good if you don't trust your model
- Different variants of cross-validation are useful in different scenarios
- Cross-validation has high variance, and **if** you trust your model you can beat cross-validation in accuracy

- It's good to think predictions of observables, because observables are the only ones we can observe
- Cross-validation can simulate predicting and observing new data
- Cross-validation is good if you don't trust your model
- Different variants of cross-validation are useful in different scenarios
- Cross-validation has high variance, and **if** you trust your model you can beat cross-validation in accuracy

- It's good to think predictions of observables, because observables are the only ones we can observe
- Cross-validation can simulate predicting and observing new data
- Cross-validation is good if you don't trust your model
- Different variants of cross-validation are useful in different scenarios
- Cross-validation has high variance, and **if** you trust your model you can beat cross-validation in accuracy

See also Gelman, Hill & Vehtari: Regression and Other Stories

Part 2: Projective Inference in High-dimensional Problems: Prediction and Feature Selection

High dimensional small data

- In the examples *n* = 54...102, *p* = 1536...22283
 - could scale to bigger *n* and bigger *p*

High dimensional small data

- In the examples *n* = 54...102, *p* = 1536...22283
 - could scale to bigger *n* and bigger *p*
- Priors necessary
 - shrinkage priors, hierarchical shrinkage priors
 - dimension reduction with factor models

High dimensional small data

- In the examples *n* = 54...102, *p* = 1536...22283
 - could scale to bigger n and bigger p
- Priors necessary
 - shrinkage priors, hierarchical shrinkage priors
 - dimension reduction with factor models
- The main content of this part: Two stage approach
 - Construct a best predictive model you can ⇒ reference model
 - Feature selection and post-selection inference ⇒ *projection*

Rich model vs feature selection?

- If we care only about the predictive performance
 - Include all available prior information
 - Integrate over all uncertainties
 - No need for feature selection

Rich model vs feature selection?

- If we care only about the predictive performance
 - Include all available prior information
 - Integrate over all uncertainties
 - No need for feature selection
- Variable selection can be useful if
 - need to reduce measurement or computation cost in the future
 - improve explainability

Rich model vs feature selection?

- If we care only about the predictive performance
 - Include all available prior information
 - Integrate over all uncertainties
 - No need for feature selection
- Variable selection can be useful if
 - need to reduce measurement or computation cost in the future
 - improve explainability
- Two options for variable selection
 - Find a minimal subset of features that yield a good predictive model
 - Identify all features that have predictive information

Regularized horseshoe prior

- Horseshoe: can be seen as continuos version of spike-and-slab with *infinite* width slab
 - no shrinkage ($\kappa_j \rightarrow 0$) allows complete separation in logistic model with $n \ll p$



Regularized horseshoe prior

- Horseshoe: can be seen as continuos version of spike-and-slab with *infinite* width slab
 - no shrinkage (κ_j → 0) allows complete separation in logistic model with n ≪ p



- Regularized horseshoe: adds additional *finite* width slab
 - some minimal shrinkage (κ_j > 0) for relevant features, but maintains division to relevant and non-relevant features



Regularized horseshoe

- Piironen and Vehtari (2017). Sparsity information and regularization in the horseshoe and other shrinkage priors. In Electronic Journal of Statistics, 11(2):5018-5051. Online
 - regularized horseshoe
 - how to set the prior based on the sparsity assumption

Why shrinkage priors alone do not solve the variable selection problem

- A common strategy:
 - Fit model with a shrinkage prior
 - Select variables based on marginal posteriors (of the regression coefficients)

Why shrinkage priors alone do not solve the variable selection problem

- A common strategy:
 - Fit model with a shrinkage prior
 - Select variables based on marginal posteriors (of the regression coefficients)
- Problems
 - Marginal posteriors are difficult with correlated features
 - How to do post-selection inference correctly?

Consider data

$$f \sim N(0, 1),$$

 $y \mid f \sim N(f, 1)$
 $x_j \mid f \sim N(\sqrt{\rho}f, 1 - \rho), \quad j = 1, \dots, 25,$
 $x_j \mid f \sim N(0, 1), \quad j = 26, \dots, 50.$

Consider data

$$f \sim N(0, 1),$$

 $\mathbf{y} \mid \mathbf{f} \sim \mathbf{N}(\mathbf{f}, 1)$
 $x_j \mid f \sim N(\sqrt{\rho}f, 1 - \rho), \quad j = 1, \dots, 25,$
 $x_j \mid f \sim N(0, 1), \quad j = 26, \dots, 50.$

• y are noisy observations about latent f

Consider data

$$f \sim N(0, 1),$$

 $y \mid f \sim N(f, 1)$
 $x_j \mid f \sim N(\sqrt{\rho}f, 1 - \rho), \quad j = 1, \dots, 25,$
 $x_j \mid f \sim N(0, 1), \quad j = 26, \dots, 50.$

- y are noisy observations about latent f
- First *p*_{rel} = 25 features are correlated with *ρ* and predictive about *y*

Consider data

$$f \sim N(0, 1),$$

 $y \mid f \sim N(f, 1)$
 $x_j \mid f \sim N(\sqrt{\rho}f, 1 - \rho), \quad j = 1, \dots, 25,$
 $x_j \mid f \sim N(0, 1), \quad j = 26, \dots, 50.$

- y are noisy observations about latent f
- First p_{rel} = 25 features are correlated with ρ and predictive about y
- Remaining 25 features are irrelevant random noise

Consider data

$$f \sim N(0, 1),$$

 $y \mid f \sim N(f, 1)$
 $x_j \mid f \sim N(\sqrt{\rho}f, 1 - \rho), \qquad j = 1, \dots, 25,$
 $x_j \mid f \sim N(0, 1), \qquad \qquad j = 26, \dots, 50.$

- y are noisy observations about latent f
- First p_{rel} = 25 features are correlated with ρ and predictive about y
- Remaining 25 features are irrelevant random noise

Generate one data set $\{x^{(i)}, y^{(i)}\}_{i=1}^{n}$ with n = 50 and $\rho = 0.8$ and assess the feature relevances



A) Gaussian prior, posterior median with 50% and 90% intervals
Example



A) Gaussian prior, posterior median with 50% and 90% intervalsB) Horseshoe prior, same things

Example



- A) Gaussian prior, posterior median with 50% and 90% intervals
- B) Horseshoe prior, same things
- C) Spike-and-slab prior, posterior inclusion probabilities

Example



- A) Gaussian prior, posterior median with 50% and 90% intervalsB) Horseshoe prior, same things
- C) Spike-and-slab prior, posterior inclusion probabilities

Half of the features relevant, but all marginals substantially overlapping with zero

What happens?



Aki.Vehtari@aalto.fi - @avehtari

What happens?



Aki.Vehtari@aalto.fi - @avehtari

What happens?



Aki.Vehtari@aalto.fi - @avehtari

Focus on predictive performance

- Two stage approach
 - Construct a best predictive model you can ⇒ reference model
 - Variable selection and post-selection inference ⇒ projection

Focus on predictive performance

- Two stage approach
 - Construct a best predictive model you can ⇒ reference model
 - Variable selection and post-selection inference ⇒ projection
- Instead of looking at the marginals, find the minimal subset of features which have (almost) the same predictive performance as the reference model

Reference model improves variable selection

Same data generating mechanism, but n = 30, p = 500, $p_{rel} = 150$, $\rho = 0.5$.



Reference model improves variable selection



A) Sample correlation with y vs. sample correlation with f

Reference model improves variable selection



irrelevant x_i , relevant x_j

A) Sample correlation with *y* vs. sample correlation with *f* B) Sample correlation with *y* vs. sample correlation with f_* $f_* =$ linear regression fit with 3 supervised principal components Piironen and Vehtari (2018)

Model simplification technique

- Model simplification technique
- Replace full posterior p(θ | D) with some constrained q(θ) so that the predictive distribution changes as little as possible

- Model simplification technique
- Replace full posterior *p*(θ | *D*) with some constrained *q*(θ) so that the predictive distribution changes as little as possible
- Example constraints
 - $q(\theta)$ can have only point mass at some θ_0
 - \Rightarrow "Optimal point estimates"

- Model simplification technique
- Replace full posterior *p*(θ | *D*) with some constrained *q*(θ) so that the predictive distribution changes as little as possible
- Example constraints
 - q(θ) can have only point mass at some θ₀
 ⇒ "Optimal point estimates"
 - Some features must have exactly zero regression coefficient
 - \Rightarrow "Which features can be discarded"

- Model simplification technique
- Replace full posterior p(θ | D) with some constrained q(θ) so that the predictive distribution changes as little as possible
- Example constraints
 - q(θ) can have only point mass at some θ₀
 ⇒ "Optimal point estimates"
 - Some features must have exactly zero regression coefficient
 ⇒ "Which features can be discarded"
- The decision theoretic idea of conditioning the smaller model inference on the full model can be tracked to Lindley (1968)
 - draw by draw projection introduced by Goutis & Robert (1998), and Dupuis & Robert (2003)
 - see also many related references in a review by Vehtari & Ojanen (2012)



Full posterior for β_1 and β_2 and contours of predicted class probability



Projected point estimates for β_1 and β_2



Projected point estimates, constraint $\beta_1 = 0$



Projected point estimates, constraint $\beta_2 = 0$



Draw-by-draw projection, constraint $\beta_1 = 0$



Draw-by-draw projection, constraint $\beta_2 = 0$

Predictive projection

Replace full posterior *p*(θ | *D*) with some constrained *q*(θ) so that the predictive distribution changes as little as possible

Predictive projection

- Replace full posterior *p*(θ | *D*) with some constrained *q*(θ) so that the predictive distribution changes as little as possible
- As the full posterior $p(\theta \mid D)$ is projected to $q(\theta)$
 - the prior is also projected and there is no need to define priors for submodels separately

Predictive projection

- Replace full posterior *p*(θ | *D*) with some constrained *q*(θ) so that the predictive distribution changes as little as possible
- As the full posterior $p(\theta \mid D)$ is projected to $q(\theta)$
 - the prior is also projected and there is no need to define priors for submodels separately
 - even if we constrain some coefficients to be 0, the predictive inference is conditoned on the information related features contributed to the reference model

• How to select a feature combination?

- How to select a feature combination?
- For a given model size, choose feature combination with minimal projective loss

- How to select a feature combination?
- For a given model size, choose feature combination with minimal projective loss
- Search heuristics, e.g.
 - Monte Carlo search
 - Forward search
 - L1-penalization (as in Lasso)

- How to select a feature combination?
- For a given model size, choose feature combination with minimal projective loss
- Search heuristics, e.g.
 - Monte Carlo search
 - Forward search
 - L₁-penalization (as in Lasso)
- Use cross-validation to select the appropriate model size
 - need to cross-validate over the search paths









Selection induced bias in variable selection



Aki.Vehtari@aalto.fi - @avehtari

Bodyfat: small p example of projection predictive

Predict bodyfat percentage. The reference value is obtained by immersing person in water. n = 251.

Bodyfat: small p example of projection predictive

Predict bodyfat percentage. The reference value is obtained by immersing person in water. n = 251.



Aki.Vehtari@aalto.fi - @avehtari
Marginal posteriors of coefficients



Bivariate marginal of weight and height



The predictive performance of the full and submodels



Aki.Vehtari@aalto.fi - @avehtari

Marginals of projected posterior



Aki.Vehtari@aalto.fi - @avehtari

Projected posterior is not just the conditional of joint



Projection of Gaussian graphical models

 Williams, Piironen, Vehtari, Rast (2018). Bayesian estimation of Gaussian graphical models with projection predictive selection. arXiv:1801.05725



CEU genetic network. BGL: Bayesian glasso; GL: glasso; TIGER: tuning insensitive graph estimation and regression; BMA: Bayesian model averaging; MAP: Maximum a posteriori; Projection: projection predictive

Aki.Vehtari@aalto.fi - @avehtari

More results

- More results projpred vs. Lasso and elastic net: Piironen, Paasiniemi, Vehtari (2018). Projective Inference in High-dimensional Problems: Prediction and Feature Selection. arXiv:1810.02406
- More results projpred vs. marginal posterior probabilities: Piironen and Vehtari (2017). Comparison of Bayesian predictive methods for model selection. Statistics and Computing, 27(3):711-735. doi:10.1007/s11222-016-9649-y.
- projpred for Gaussian graphical models: Williams, Piironen, Vehtari, Rast (2018). Bayesian estimation of Gaussian graphical models with projection predictive selection. arXiv:1801.05725
- More results for Bayes SPC: Piironen and Vehtari (2018). Iterative supervised principal components. 21st AISTATS, PMLR 84:106-114. Online.
- Several case studies for small to moderate dimensional (p = 4...100) small data: Vehtari (2018). Model assessment, selection and inference after selection. https://avehtari.github.io/modelselection/

- Sparse priors do not automate variable selection
 - Don't trust marginal posteriors

- Sparse priors do not automate variable selection
 - Don't trust marginal posteriors
- Reference model + projection can improve feature selection
 - Excellent tradeoff between accuracy and model complexity
 - Useful also for identifying all the relevant features

- Sparse priors do not automate variable selection
 - Don't trust marginal posteriors
- Reference model + projection can improve feature selection
 - Excellent tradeoff between accuracy and model complexity
 - Useful also for identifying all the relevant features
- Well developed for GLMs, but can be used also with other model families

- Sparse priors do not automate variable selection
 - Don't trust marginal posteriors
- Reference model + projection can improve feature selection
 - Excellent tradeoff between accuracy and model complexity
 - Useful also for identifying all the relevant features
- Well developed for GLMs, but can be used also with other model families
- More details and results (+ some theoretical discussion) in the paper
 - Piironen, Paasiniemi, Vehtari (2018). Projective Inference in High-dimensional Problems: Prediction and Feature Selection. arXiv:1810.02406

- Sparse priors do not automate variable selection
 - Don't trust marginal posteriors
- Reference model + projection can improve feature selection
 - Excellent tradeoff between accuracy and model complexity
 - Useful also for identifying all the relevant features
- Well developed for GLMs, but can be used also with other model families
- More details and results (+ some theoretical discussion) in the paper
 - Piironen, Paasiniemi, Vehtari (2018). Projective Inference in High-dimensional Problems: Prediction and Feature Selection. arXiv:1810.02406
- R-package projpred in CRAN and github https://github.com/stan-dev/projpred (easy to use, e.g. with RStan, RStanARM, brms)

References

References and more at avehtari.github.io/masterclass/ and avehtari.github.io/modelselection//

- Model selection tutorial at StanCon 2018 Asilomar
 - more about projection predictive variable selection
- Regularized horseshoe talk at StanCon 2018 Asilomar
- Several case studies
- References with links to open access pdfs