# 7: Evaluating, comparing and expanding models

10/14/19

# Introduction

This chapter focuses mostly on quantifying a model's predictive capabilities for the purposes of model selection and expansion.

# New Notation!

1. $f$ is the true model
2. $y$ is the data we use to estimate our model
3. $\tilde{y}$ is the future (time series) or alternative (not time series) data that we test our predictions on
4. $p_{\text{post}}(\tilde{y}) = p(\tilde{y} \mid y)$
5. $p_{\text{post}}(\theta) = p(\theta \mid y)$
6. $E_{\text{post}}[\cdot]$ is taken with respect to $p(\theta \mid y)$

## Definitions

A **scoring rule/function** $S(p, \tilde{y})$ is a function that takes

1. the distribution you're using to forecast $p$ (ppd, or likelihood with estimated parameters), and
2. a realized value $\tilde{y}$

and then gives you a real-valued number/score/utility. Higher is better, although this convention isn't always followed in the literature.

Keep in mind that the realized value cannot be used to fit the data.

# Examples

Example: $S(p, \tilde{y}) = -(\tilde{y} - E_p[\tilde{y}])^2$

Example: $S(p, \tilde{y}) = \log p(\tilde{y})$

# (Out-of-sample) predictive fit

Future/unseen data is unknown, so we must take the expected score under the true distribution $f$:

$$E_f[S(p, \tilde{y})].$$

A scoring rule is **proper** if the above expectation is maximized when $f = p$.

A scoring rule is **local** if $S(p, \tilde{y})$ only depends on $p(\tilde{y})$ (don't care about events that didn't happen).

Note, when we are dealing with a logarithmic scoring rule, $E[-2 \log p(\tilde{y})]$ is often called an **information criterion.** The book switches back and forth between dealing with expected score, and information criteria.

# Examples

Example: $S(p, \tilde{y}) = -(\tilde{y} - E_p[\tilde{y}])^2$
Most common, perhaps not local or proper for non-Gaussian data.

Example: $S(p, \tilde{y}) = \log p(\tilde{y})$
Obviously local. Proper, too.

# Empirical predictive fit

We are generally not able to evaluate the expectation because we don't know $f$. However, we may be able to wait for new out-of-sample data and use a Monte-Carlo approach:

$$n^{-1} \sum_{i=1}^{n} S(p, \tilde{y}^i) \to E_f[S(p, \tilde{y})]$$

as $n \to \infty$

# Empirical predictive fit

We are generally not able to evaluate the expectation because we don't know $f$. However, we may be able to wait for new out-of-sample data and use a Monte-Carlo approach:

$$n^{-1} \sum_{i=1}^{n} S(p, \tilde{y}^i) \to E_f[S(p, \tilde{y})]$$

as $n \to \infty$

## Definitions

The textbook focuses on $S(p, \tilde{y}) = \log p(\tilde{y})$, and the data are iid (after conditioning on the parameter). They call the following quantity the "elppd:"

### expected log pointwise predictive density

$$E_f[\log p(\tilde{y})] = E_f\left[\log \prod_i p(\tilde{y}_i)\right]$$

$$= \sum_{i=1}^{n} E_f\left[\log p(\tilde{y}_i)\right]$$

- In general, $E_f[\log p(\tilde{y})] \neq \sum_{i=1}^{n} E_f\left[\log p(\tilde{y}_i)\right]$. Sum of "pointwise" ppd not equal to "joint" ppd.

# Problem

For the moment let's use $p(\tilde{y}) = p_{post}(\tilde{y})$

The "elppd" is not obtainable because

1. you don't know $f$
2. you don't have $p_{post}(\tilde{y})$

# Problem

For the moment let's use $p(\tilde{y}) = p_{\text{post}}(\tilde{y})$

The "elppd" is not obtainable because

1. you don't know $f$
2. you don't have $p_{\text{post}}(\tilde{y})$

Using $y$ for $\tilde{y}$, we can come up with a rough elppd estimate called the "lppd"

## log pointwise predictive density

$$\text{lppd} = \log p_{\text{post}}(y) = \sum_{i=1}^{n} \log p_{\text{post}}(y_i)$$

# Problem

There'a also the problem that arises where we cannot evaluate

$$p_{\text{post}}(y) = \int p(y \mid \theta)p(\theta \mid y)\mathrm{d}\theta = E_{\text{post}}[p(y \mid \theta)]$$

The "computed lppd" again uses $y$ for $\tilde{y}$, but it also uses Monte-Carlo to sample from the posterior

## log pointwise predictive density

$$\text{computed lppd} = \log \hat{p}_{\text{post}}(y) = \sum_{i=1}^{n} \log \left( \frac{1}{S} \sum_{j=1}^{S} p(y_i \mid \theta^j) \right)$$

-Biased and probably high variance, though.

# Three problems

Don't know $f$, don't want to wait for $\tilde{y}$...

and unfortunately, plugging the same data that we used for estimation into the predictive distribution might lead us to overfit because this strategy overestimates the average predictive score. What do we do?

# Three problems

Don't know $f$, don't want to wait for $\tilde{y}$...

and unfortunately, plugging the same data that we used for estimation into the predictive distribution might lead us to overfit because this strategy overestimates the average predictive score. What do we do?

However, we can get around this in two ways generally:

1. plug in the already-used $y$ data, but then add an extra penalty term (e.g. AIC, DIC, WAIC, etc.)
2. Cross-Validation: split the data $y$, many different ways, into a train and test set; estimate and evaluate on each split.

## Information Criteria

**AIC** stands for "Akaike's Information Criterion." Let $k$ be the number of parameters:

$$\widehat{\text{elpd}}_{\text{AIC}} = \log p(y \mid \hat{\theta}_{\text{MLE}}) - \overbrace{k}^{\text{penalty}}$$

or

$$\text{AIC} = \underbrace{-2 \log p(y \mid \hat{\theta}_{\text{MLE}})}_{\text{a deviance}} + 2k$$

We estimate $\hat{\theta}_{\text{MLE}}$ using $y$, and we plug $y$ into the log likelihood.

# Information Criteria

**DIC** replaces the point estimate with $\hat{\theta}_{\text{Bayes}} = E[\theta \mid y]$, and replaces the penalty term with $p_{\text{DIC}}$

$$\widehat{\text{elpd}}_{\text{DIC}} = \log p(y \mid \hat{\theta}_{\text{Bayes}}) - p_{\text{DIC}}$$

or

$$\text{DIC} = -2 \log p(y \mid \hat{\theta}_{\text{Bayes}}) + 2p_{\text{DIC}}$$

# Information Criteria

The book gives two ways to estimate $p_{\text{DIC}}$:

1. $p_{\text{DIC}} = 2\left(\log p(y \mid \hat{\theta}_{\text{Bayes}}) - E_{\text{post}}\left[\log p(y \mid \theta)\right]\right)$
2. $p_{\text{DIC alt}} = 2\,\text{Var}_{\text{post}}\left[\log p(y \mid \theta)\right]$

Both of these can be approximated using samples from the posterior.

- Both $p_{\text{DIC}}$ and $p_{\text{DIC alt}}$ are estimated effective number of parameters
- Both reduce to $k$ for linear models with uniform prior distributions

# Information Criteria

$p_{WAIC}$ stands for "Watanabe-Akaike information criterion" or "widely available information criterion"

The book refers to it as the most "fully Bayesian" of the three, probably because it doesn't plug in point estimates into the likelihood instead of integrating.

$$\widehat{\text{elppd}}_{WAIC} = \text{lppd} - p_{WAIC}$$

or

$$\text{WAIC} = -2\text{lppd} + 2p_{WAIC}$$

where lppd is computed by Monte Carlo method as
$\sum_{i=1}^{n} \log\left(\frac{1}{S}\sum_{s=1}^{S} p(y_i \mid \theta^s)\right)$

# Information Criteria

Two ways to estimate

1. $p_{\text{WAIC 1}} = 2 \sum_{i=1}^{n} \left( \log p_{\text{post}}(y_i) - E_{post} \left\{ \log p(y_i \mid \theta) \right\} \right)$
2. $p_{\text{WAIC 2}} = \sum_{i=1}^{n} var_{\text{post}}(\log p(y_i \mid \theta))$

Both of these can be approximated using samples from the posterior.

# Cross-Validation

To assess prediction performance, one may also use **cross-validation**. Here the data is repeatedly partitioned into different training-set-test-set pairs (aka **folds**).

# Cross-Validation

To assess prediction performance, one may also use **cross-validation**.
Here the data is repeatedly partitioned into different training-set-test-set
pairs (aka **folds**).

1. The partitions are nonrandom, test sets are disjoint
2. for each split/estimation/prediction, we never use a data point twice
3. for each split/estimation/prediction, we lose parameter estimation
   accuracy because each training set is smaller than the full set
4. however, we get to average over many prediction scores, which
   reduces variance
5. there is still a bias that we have to estimate (due to small sample size
   but it's usually smaller than AIC/DIC/WAIC/etc.)
6. it can be computationally brutal to calculate for some models

# Cross-Validation

**leave-one-out cross-validation** (loo-cv) is a special case where each test set is of size 1.

This necessarily implies that each training set is of size $n - 1$, and there are $n$ possible splits.

If this ends up being too computationally expensive, it is also possible to do $k$-**fold cross-validation**, which selects $k$ splits/folds. This means the size of each test set is $n/k$, and the size of each training set is $n - n/k$

# Cross-Validation Notation

We only discuss loo-cv...

$p_{\text{post}(-i)}(y_i)$ is the prediction for the $i$th point, using the ppd, which uses the posterior distribution conditioning on all values of the data **except the $i$th**

## Cross-Validation Notation

We only discuss loo-cv...

$p_{\text{post}(-i)}(y_i)$ is the prediction for the $i$th point, using the ppd, which uses the posterior distribution conditioning on all values of the data **except the $i$th**

If this ppd isn't tractable, we can use draws from the posterior as follows:

$$p_{\text{post}(-i)}(y_i) = \frac{1}{S} \sum_{s=1}^{S} p(y_i \mid \theta^s)$$

where $\theta^{is}$ are draws from $p_{\text{post}(-i)}(\theta)$

# Cross-Validation

The Bayesian loo-cv estimate for out-of-sample predictive fit is

$$\text{lppd}_{\text{loo-cv}} = \sum_{i=1}^{n} \log p_{\text{post}(-i)}(y_i)$$

There are also bias-corrected versions as well, that is, $\text{lppd}_{\text{loo-cv}} + b$, where

$$b = lppd - \overline{lppd}_{-i}$$

$$\overline{lppd}_{-i} = \frac{1}{n} \sum_{i=1}^{n} \sum_{j=1}^{n} \log p_{\text{post}(-i)}(y_j)$$

# Comparisons of all comparison criteria

Given score function $S(y, p) = \log p(y)$

- AIC and DIC are based on the **joint likelihood** score function conditioning on a point estimate of $\theta$; WAIC and LOO-CV are based on **individual posterior** score function averaging over the posterior distribution $p(\theta \mid y)$

# Comparisons of all comparison criteria

Given score function $S(y, p) = \log p(y)$

- AIC and DIC are based on the **joint likelihood** score function conditioning on a point estimate of $\theta$; WAIC and LOO-CV are based on **individual posterior** score function averaging over the posterior distribution $p(\theta \mid y)$
- WAIC and LOO-CV require an explicit assumption/requirement that data are independent conditioning on the parameter; not easy to do in some structured-data settings such as time series, spatial and network data

# Comparisons of all comparison criteria

Given score function $S(y, p) = \log p(y)$

- AIC and DIC are based on the **joint likelihood** score function conditioning on a point estimate of $\theta$; WAIC and LOO-CV are based on **individual posterior** score function averaging over the posterior distribution $p(\theta \mid y)$
- WAIC and LOO-CV require an explicit assumption/requirement that data are independent conditioning on the parameter; not easy to do in some structured-data settings such as time series, spatial and network data
- WAIC and LOO-CV are equal asymptotically

# Comparisons of all comparison criteria

Given score function $S(y, p) = \log p(y)$

- AIC and DIC are based on the **joint likelihood** score function conditioning on a point estimate of $\theta$; WAIC and LOO-CV are based on **individual posterior** score function averaging over the posterior distribution $p(\theta \mid y)$
- WAIC and LOO-CV require an explicit assumption/requirement that data are independent conditioning on the parameter; not easy to do in some structured-data settings such as time series, spatial and network data
- WAIC and LOO-CV are equal asymptotically
- AIC and LOO-CV are restrictive when applied to hierarchical models (eight school example)

# Comparisons of all comparison criteria

Given score function $S(y, p) = \log p(y)$

- AIC and DIC are based on the **joint likelihood** score function conditioning on a point estimate of $\theta$; WAIC and LOO-CV are based on **individual posterior** score function averaging over the posterior distribution $p(\theta \mid y)$
- WAIC and LOO-CV require an explicit assumption/requirement that data are independent conditioning on the parameter; not easy to do in some structured-data settings such as time series, spatial and network data
- WAIC and LOO-CV are equal asymptotically
- AIC and LOO-CV are restrictive when applied to hierarchical models (eight school example)

Overall, WAIC is more appealing than the rest of them.