

Statistical Methods - Model evaluation, comparison and selection

Lecture 9

Mattias Villani

Sveriges Riksbank and Stockholm University

May 24, 2010

- Approaches to model comparison and evaluation
- Evaluating prediction performance
- Model selection criteria
- Bayesian model inference

- Given a set of models, how do we compare them?
- Classical testing: likelihood ratio test. Restricted to nested models (some extension to the non-nested case have been proposed). Mostly asymptotic. Pair-wise comparison of models.
- Information criteria. AIC, BIC etc
- Prediction-based comparison. Cross-validation and Bootstrap.
- Bayesian marginal likelihood comparison.

- The best model in a class of models can still be bad ...
- How do we know if a model is 'any good'? First question: good for **what**?
 - Prediction (can we accurately predict a future event? Black box model/algorithm may be OK)
 - Interpretation (can we learn something from the model structure? Is it useful for organizing our thinking? May be OK even if predictions are inaccurate)
 - Policy analysis (what happens with the response when we change the conditioning variables? Both prediction performance and interpretation usually matters here)
- If prediction is our concern, how can we estimate the expected performance of the model in a new data set? How well does the model generalize from in-sample to out-of-sample?
- Weak evaluation of the model: If simulated data from your model looks nothing like actual data (maybe even the same data used to fit the model!), then the model is 'no good'.

■ Components:

- Target variable: Y .
- Covariates/Inputs: X .
- Prediction model $\hat{f}(X)$
- Training sample \mathcal{T} , used to estimate $\hat{f}(X)$
- Loss function: $L[Y, \hat{f}(X)]$ measuring the discrepancy between true Y and the model's prediction $\hat{f}(X)$. Examples:
 - Squared loss $[Y - \hat{f}(X)]^2$
 - Linear loss $|Y - \hat{f}(X)|$
 - Lin-lin loss:

$$L[Y, \hat{f}(X)] = \begin{cases} c_1 \cdot |Y - \hat{f}(X)| & \text{if } Y \leq \hat{f}(x) \\ c_2 \cdot |Y - \hat{f}(X)| & \text{if } Y > \hat{f}(x) \end{cases}$$

- **Generalization error**, or test error, is defined as

$$\text{Err}_{\mathcal{T}} = \text{E}\{L[Y, \hat{f}(X)] | \mathcal{T}\}$$

The expectation is with respect to the joint distribution of (Y, X) . The training set \mathcal{T} is fixed, so $\text{Err}_{\mathcal{T}}$ is the prediction error we can expect on new independent data **given** the particular training sample at hand. In a given application, $\text{Err}_{\mathcal{T}}$ is what we care about.

- The **expected prediction error** is defined as

$$\text{Err} = \text{E}(\text{Err}_{\mathcal{T}}),$$

where the expectation is with respect to the training sample \mathcal{T} . Err measures the average prediction performance of the method/model over all possible training samples. This is usually less relevant when we actually have a particular training sample at hand.

- **Training error** is the average loss over the training sample:

$$e\bar{r} = \frac{1}{N} \sum_{i=1}^N L[y_i, \hat{f}(x_i)].$$

$e\bar{r}$ can be made arbitrarily small by increasing the model complexity and we typically have $e\bar{r} < \text{Err}_{\mathcal{T}}$.

- **In-sample error**

$$\text{Err}_{\text{in}} = \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{Y^0} L[Y^0, \hat{f}(x_i)]$$

where the x_i are the training sample points, but each Y_i^0 is a new response observation drawn from $p(Y^0|X = x_i)$.

- Note that Err_{in} conditions on the values of the covariates observed in the training sample, whereas $\text{Err}_{\mathcal{T}}$ does not.
- Note the distinction between in-sample error and training error, it is easy to confuse them!

- The **optimism** (in the training error) is defined as

$$\text{op} = \text{Err}_{\text{in}} - \bar{\text{err}}$$

and the average optimism is

$$\omega = E_{\mathbf{y}}(\text{op}),$$

where $E_{\mathbf{y}}$ denote the expectation with respect to the outcomes of the y_i , $i \in \mathcal{T}$, but the x 's in the training set are fixed.

- One can show quite generally that

$$\omega = \frac{2}{N} \sum_{i=1}^N \text{Cov}(y_i, \hat{y}_i)$$

Very intuitive: the more effect the observations have on their own fit, the larger the optimism.

- Important relation:

$$E_{\mathbf{y}}(\text{Err}_{\text{in}}) = E_{\mathbf{y}}(e\bar{r}r) + \frac{2}{N} \sum_{i=1}^N \text{Cov}(y_i, \hat{y}_i)$$

which implies that we can estimate the in-sample error by

$$\hat{\text{Err}}_{\text{in}} = e\bar{r}r + \hat{\omega},$$

where $\hat{\omega}$ is an estimate of the average optimism.

- Example: linear smoothers: $\hat{\mathbf{y}} = \mathbf{L}\mathbf{y}$.

$$\begin{aligned}\text{Cov}(\mathbf{y}, \hat{\mathbf{y}}) &= \mathbb{E}[\mathbf{y} - \mathbb{E}(\mathbf{y})][\hat{\mathbf{y}} - \mathbb{E}(\hat{\mathbf{y}})]' \\ &= \mathbb{E}[\mathbf{y} - \boldsymbol{\mu}][\mathbf{L}\mathbf{y} - \mathbf{L}\boldsymbol{\mu}]' \\ &= \text{Cov}(\mathbf{y})\mathbf{L}'\end{aligned}$$

so if $\text{Cov}(\mathbf{y}) = \sigma_\varepsilon^2 \mathbf{I}$, then

$\sum_{i=1}^N \text{Cov}(y_i, \hat{y}_i) = \text{tr} \text{Cov}(\mathbf{y}, \hat{\mathbf{y}}) = \sigma_\varepsilon^2 \text{tr} \mathbf{L} = \sigma_\varepsilon^2 \cdot Df$. Here we have the estimate of the in-sample error

$$\hat{\text{Err}}_{in} = e\bar{r}r + 2 \cdot \frac{Df}{N} \sigma_\varepsilon^2.$$

- This motivates the following definition of Df for general additive-error models $y = f(x) + \varepsilon$ (which may not be linear smoothers)

$$Df(\hat{\mathbf{y}}) = \frac{\text{tr} \text{Cov}(\mathbf{y}, \hat{\mathbf{y}})}{\sigma_\varepsilon^2} = \frac{\sum_{i=1}^N \text{Cov}(y_i, \hat{y}_i)}{\sigma_\varepsilon^2}.$$

- The estimate of the in-sample error is similar to the more generally applicable Akaike Information Criterion (AIC)

$$AIC = -2 \cdot \loglik + 2 \cdot p,$$

where \loglik is the maximum of the likelihood function, and p is the number of effective parameters ($p = Df$ for linear smoothers).

- Note that \hat{Err}_{in} and AIC are in-sample measures and are poor estimators of the generalization error. They can be used to **compare** models, however.
- But AIC is not a consistent model selection criteria ...
- BIC (Bayesian Information Criteria)

$$BIC = -2\loglik + \ln N \cdot p$$

- BIC is consistent and penalizes more complex models more heavily than AIC. BIC can be converted to model probabilities, see next lecture.