

Optimisation of the Training Error

Let us start from the definitions of test time:

$$\text{Err}_\tau = E_{(x_0, y_0)} [L(y_0, \hat{f}(x_0)) | \tau]$$

Test error, where

- x_0, y_0 are new (test) point \rightarrow random
- τ is fixed $\tau = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$

Averaging over all training sets

$$\text{Err} = E_\tau [E_{(x_0, y_0)} [L(y_0, \hat{f}(x_0)) | \tau]]$$

gives the expected error

We also saw the training error

$$\bar{\text{err}} = \frac{1}{N} \sum_{i=1}^N L(y_i, \hat{f}(x_i))$$

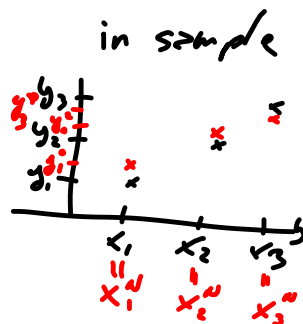
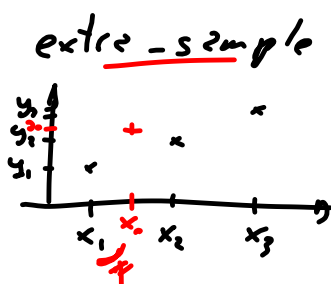
is NOT a good estimate of the test error, because it underestimates it.

Same data used for both training $\hat{f}(x)$ and to test its performance \rightarrow optimistic estimate

The test error can be thought as a extra-sample error (the estimation of the error is computed on new points, $x_0 \neq x_i$)

We are going to evaluate the optimism of $\bar{\text{err}}$ in the in-sample case, i.e. we have new observations in the same points of training set

$$\text{Err}_{\text{in}} = \frac{1}{N} \sum_{i=1}^N E_{y_0} [L(y_i, \hat{f}(x_i)) | \tau]$$



- only y_0 is random, new y_i for x_1, \dots, x_N

Definition: we define optimism the difference between Err_{in} and \bar{err}

$$op := Err_{in} - \bar{err}$$

op is generally positive, as \bar{err} is computed on training data

Definition: average optimism is

$$w := E_y[op]$$

we are computing the expected value over the training set outcome

For a reasonable number of loss functions, including 0-1 loss and squared error, it can be shown that

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

where Cov denotes the covariance \rightarrow Ex 7.4

Note:

- optimism depends on how much y_i affects its own prediction
- the harder we fit the data, the larger the value of $\text{Cov}(\hat{y}_i, y_i)$
 \rightarrow the higher the optimism

As a consequence:

$$E_y[Err_{in}] = E_y[\bar{err}] + \frac{2}{N} \sum_{i=1}^N \text{Cov}(\hat{y}_i, y_i)$$

When \hat{y}_i is obtained by a linear fit in the inputs, $Y = f(x) + \epsilon$

$$\sum_{i=1}^N \text{Cov}(\hat{y}_i, y_i) = d \sigma_\epsilon^2 \quad \leftarrow \text{effective number of parameters}$$

therefore

$$\rightarrow E_y[Err_{in}] = E_y[\bar{err}] + 2 \frac{d}{N} \sigma_\epsilon^2 \quad (*)$$

- optimism increases with the number of inputs;
- " decreases " " sample size.

e.g. in linear regression, the number of covariates

Methods we will see:

- C_p , AIC and BIC estimate the prediction error by estimating the training error and the optimism (work when estimates are linear in their parameters)
- cross-validation and bootstrap-based procedure try to estimate directly the expected error

Note:

- in-sample error is generally NOT of interest (we are mainly interested in new data, including new points in X)
- to select the best model / tuning the complexity parameter, we are more interested in the relative difference in errors rather than the absolute one.

Estimates of Err_{in}

Start from (*)

$$E_y[\bar{Err}_{in}] = E_Y[\bar{err}] + 2 \frac{d}{N} \sigma_\epsilon^2 \quad (*)$$

Write the general form of the in-sample estimate

$$\hat{Err}_{in} = \bar{err} + \hat{w}$$

- when we have linearity and squared errors, from (*)

$$C_p = \bar{err} + 2 \frac{d}{N} \hat{\sigma}_\epsilon^2$$

where:

- \bar{err} is computed through the squared loss;
- d is # of parameters
- $\hat{\sigma}_\epsilon^2$ is an estimate of the noise variance

NB: there are other versions of C_p , different versions may give different numbers, but they all lead to the same model

it is computed using the full model, because it has the smallest bias

Similar idea for AIC (Akaike Information Criterion)

- we start again from (*), but we want to be more general (the error is computed through a likelihood approach)

We are using the asymptotic result ($N \rightarrow \infty$)

$$-2 E[\log P_\theta(Y)] \approx -\frac{2}{N} E\left[\sum_{i=1}^N \log P_\theta(y_i)\right] + 2 \frac{d}{N}$$

- $P_\theta(Y)$ is the family of densities for Y (containing the "true" density)

loglik \hookrightarrow the maximized log-likelihood
 $\ell(\hat{\theta})$
 is the mle maximum likelihood estimate

E.g., in the logistic regression: $AIC = -\frac{2}{N} \loglik + 2 \frac{d}{N}$

Gaussian regression: $AIC \propto C_p$

\hookrightarrow AIC C_p is a special case of AIC

To find the best model, we choose that with the smallest AIC

- straightforward in the simplest cases, we need more attention in more complex situations. In particular, we need to find a reasonable measure for the model complexity

Note: for regularized/penalized regression

$$AIC(\alpha) = \bar{err}(\alpha) + 2 \frac{d(\alpha)}{N} \hat{\sigma}_\epsilon^2$$

- usually minimizing AIC is not the best solution to find the value of the tuning parameter \rightarrow cross-validation is a better approach

The effective number of parameters

- generalized the concept of number of parameters, in order to extend the previous approaches to more complex situation

Suppose

$y = (y_1, \dots, y_N)$ outcome

$\hat{y} = (\hat{y}_1, \dots, \hat{y}_N)$ predictions

Linear methods $\hat{y} = Sy$

where S is a $N \times N$ matrix which:

- depends on X
- independent on Y

Linear regression

$$\hat{y} = \underbrace{X(X^T X)^{-1} X^T}_S y$$

Ridge regression

$$\hat{y} = \underbrace{X(X^T X + \lambda I_p)^{-1} X^T}_S y$$

Effective number of parameters (or effective degrees of freedom)

$$df(S) = \text{trace}(S)$$

value of d

We should replace $\text{trace}(S)$ to d to obtain the correct criterion

$$\text{If } y = f(x) + \varepsilon, \text{Var}(\varepsilon) = \sigma_\varepsilon^2 \rightarrow \sum_{i=1}^N \text{Cov}(\hat{y}_i, y_i) = \text{trace}(S) \sigma_\varepsilon^2$$

$$\text{So } df(\hat{y}) = \frac{\sum_{i=1}^N \text{Cov}(\hat{y}_i, y_i)}{\sigma_\varepsilon^2} \rightarrow \text{Ex 7.5}$$

The Bayesian approach and BIC

The BIC is an alternative criterion to AIC,
↳ Bayesian Information Criterion

$$\frac{1}{N} \text{BIC} = -\frac{2}{N} \log \text{lik} + \log N \cdot \frac{d}{N}$$

$$\text{AIC} = 2$$

$$\text{BIC} = \log(N)$$

Despite they are quite similar, AIC and BIC came from completely different ideas. BIC comes from the Bayesian approach to model selection

$$Pr(M_m | z) \propto Pr(M_m) Pr(z | M_m)$$

$$\propto Pr(M_m) \int Pr(z | M_m, \theta_m) Pr(\theta_m | M_m) d\theta_m$$

To choose between two models, we compare their posterior probabilities

$$\frac{Pr(M_0 | z)}{Pr(M_1 | z)} = \frac{Pr(M_0)}{Pr(M_1)} \cdot \frac{Pr(z | M_0)}{Pr(z | M_1)}$$

Bayes factor

in a large number of cases = 1
(give the same prior probability to the two models)

⇒ the choice between the two models is based on the Bayes factor

Approximating

$$Pr(z | M_m) = \log(z | \hat{\theta}_m, M_m) - \frac{d_m}{2} \log N + O(1)$$

mle

If the loss function is $-2 \log Pr(z | \hat{\theta}, M_m)$, we obtain BIC

• We ^{may} select the model with smallest BIC

correspond to selecting the model with highest posterior probability

Note that

$$\frac{e^{-\frac{1}{2} \text{BIC}_m}}{\sum_{i=1}^M e^{-\frac{1}{2} \text{BIC}_i}}$$

is the posterior probability of selecting the model m

AIC vs BIC

- no clear preference
- BIC leads to a sparser model
- AIC leads to a model with more predictors
- BIC is consistent ($N \rightarrow \infty$, the probability of selecting the true model goes to 1)
- for finite sample size, BIC tends to select a model which is too sparse

Cross-validation

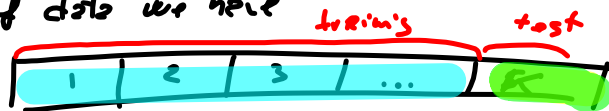
The cross-validation aims to estimate the ~~extre-sample~~ error

$$Err = E[L(Y, \hat{f}(x))]$$

the average test error when $\hat{f}(x)$ is applied to a new sample ^{independent test set}

If we had enough data $\begin{matrix} \swarrow & \searrow \\ \text{training set} & \text{test} \end{matrix}$

↳ in general, we have not, so we mimic this split using the limited amount of data we have



Folds ^{derived using so data but those in fold k}

- divide the observations in k folds
- we use, in turn, $k-1$ folds to train the model (derive $\hat{f}^{-k}(x)$)
- we evaluate the model in the remaining fold

$$CV(\hat{f}) = \sum_{k=1}^K \sum_{i=1}^n L(y_i, \hat{f}^{-k}(x_i))$$

- if $k=2$, two-fold cross-validation
- if $k=N$, leave-one-out cross-validation (LOOCV)
in this case, each observation is a fold

How do we choose k

- bias-variance trade-off
- smaller the k , larger bias, smaller variance
- larger the k , smaller bias, larger variance
(the extreme case is LOOCV, where we use $N-1$ observations for training the model → the training sets are really similar to each other)
- usual choices are $k=5$ or $k=10$

[Fig 7.8]

→ 1. Err vs N

→ the classifier is ok until $\approx \underline{N=100}$ (then is flat)

- if $N=200$, $k=5 \rightarrow$ training $N_t = 160$ \checkmark $160 > \underline{100}$
- if $N=50$, $k=5 \rightarrow$ " $N_t = 40$ \times $40 < \underline{100}$

Notes:

- CV estimates $\mathbb{E} \text{err}$ and not the Err_τ
- if we want to select a tuning parameter via CV

$\hat{f}^{-k}(x, \alpha)$ is the model selected using α and fitted on the observations which do not belong to the k -th fold

$$CV(\hat{f}, \alpha) = \frac{1}{\sum_{k=1}^K n_k} \sum_{k=1}^K \sum_{i=1}^{n_k} L(y_i, \hat{f}(x_i, \alpha))$$

$$\hat{\alpha} = \arg \min_{\alpha} CV(\hat{f}, \alpha)$$

Generalized cross-validation

- convenient approximation to the LOOCV, for squared loss function

LOOCV $\hat{y} = Sy$

$$N \sum_{k=1}^K \frac{1}{n_k} \sum_{i=1}^{n_k} (y_i - \hat{f}^{-i}(x_i))^2 \rightarrow \frac{1}{N} \sum_{i=1}^N (y_i - \hat{f}(x_i))^2 = \frac{1}{N} \sum_{i=1}^N \left(\frac{y_i - \hat{f}(x_i)}{1 - S_{ii}} \right)^2$$

where S_{ii} is the i th term on the diagonal of S

The generalized cross-validation (GCV)

$$GCV(\hat{f}) = \frac{1}{N} \sum_{i=1}^N \left(\frac{y_i - \hat{f}(x_i)}{1 - \text{trace}(S)/N} \right)^2$$

- computational advantages;
- similarities between AIC and GCV \rightarrow Ex 7.7