

## Information Criterion

Ways to quantify error in a given estimator.

One idea

$$\overline{err} = \mathbb{E}_Y \left[ \frac{1}{n} \sum_{i=1}^n L(y_i, \hat{f}(x_i)) \right]$$

In general underestimates error due to training  
and testing on  $\{Y_i\}$ .

Create a new dataset of responses  $y_i^{(new)}$  to get a better estimate  
of the error  $Errin = \mathbb{E}_Y \left[ \frac{1}{N} \sum_{i=1}^N L(\hat{f}(x_i), Y_i) \right]$

View  $y_i^{(new)}$  as R.V. coming from the underlying model.

Def: The difference between the training error and the theoretical error is given by the optimism

$$op = Errin - \overline{err}$$

and the true optimism is given by

$$\mathbb{E}_Y[op] = \frac{2}{N} \sum_{i=1}^n \text{Cov}(\hat{y}_i, y_i)$$

Which yields

$$\mathbb{E}_Y[Errin] = \overline{err} + \frac{2}{N} \sum_{i=1}^n \text{Cov}(y_i, \hat{y}_i)$$

We normally say  $\mathbb{E}_\tau(\text{Err}_n)$  is the CP information criterion.

## VC Dimension Criterion

Rmk: Optimism related to how complex the space  $\mathcal{F}$ .

Suppose  $f(x) = \mathbb{I}_A$  (classification) e.g.  $A = \{x: \beta^T x > 0\}$

Def: The family  $\mathcal{F}$  of indicator functions is said to shatter a collection of points  $C \subset \mathcal{X}$ . If for every subset  $C_1 \subset C$ , there exists

a function  $f_n = \mathbb{I}_A \in \mathcal{F}$  s.t.  $\mathbb{I}_A(x) = \begin{cases} 1 & x \in C_1 \\ 0 & x \in C \setminus C_1 \end{cases}$

Def:  $\mathcal{F}$  has VC dimension if it can shatter any collection of points

Thm:  $\mathbb{E}[L(\hat{f}(X, Y))] = \text{Err}_\tau \leq \frac{\overline{\text{err}}}{(1-\epsilon)_+}$       $\epsilon = \frac{h(\log N/n + 1) - \log m/d}{N}$

$h = \text{VC dimension}$  with probability  $1 - \eta$ .

## Bootstrapping

Define  $z = \begin{pmatrix} x_1^T & y_1 \\ \vdots & \vdots \\ x_n^T & y_n \end{pmatrix}$  with a goal of estimating  $S(z)$ .

Estimate the density  $\hat{p}$  and resample from it to estimate

$$\{\hat{S}(z)\}_{i=1}^B$$