

Bias, Variance, Model Complexity

Quantitative

Loss function: $L(Y, \hat{f}(X)) = \begin{cases} (Y - \hat{f}(X))^2 & \text{squared error} \\ |Y - \hat{f}(X)| & \text{absolute error} \end{cases}$

$$\text{Err}_T = E(L(Y, \hat{f}(X)) | T)$$

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

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

Qualitative

Loss function: $L(G, \hat{G}(X)) = I(G \neq \hat{G}(X))$ (0-1 loss)

$$L(G, \hat{p}(X)) = -2 \sum_{k=1}^K I(G=k) \log \hat{p}_k(X)$$

$$= -2 \log p_G(X) \quad (-2 \times \log\text{-likelihood})$$

Training error: $\bar{\text{err}} = -\frac{2}{N} \sum_{i=1}^N \log \hat{p}_{g_i}(x_i)$

Data-rich situation:



Insufficient data: AIC, BIC, MDL, SRM

Cross-validation and bootstrap

The Bias-Variance Decomposition

· Assumption: $Y = f(X) + \varepsilon$, $E(\varepsilon) = 0$, $\text{Var}(\varepsilon) = \sigma_\varepsilon^2$

· Predicted error of a fit $\hat{f}(X)$ at $X = x_0$:

$$\text{Err}(x_0) = E[(Y - \hat{f}(x_0))^2 \mid X = x_0]$$

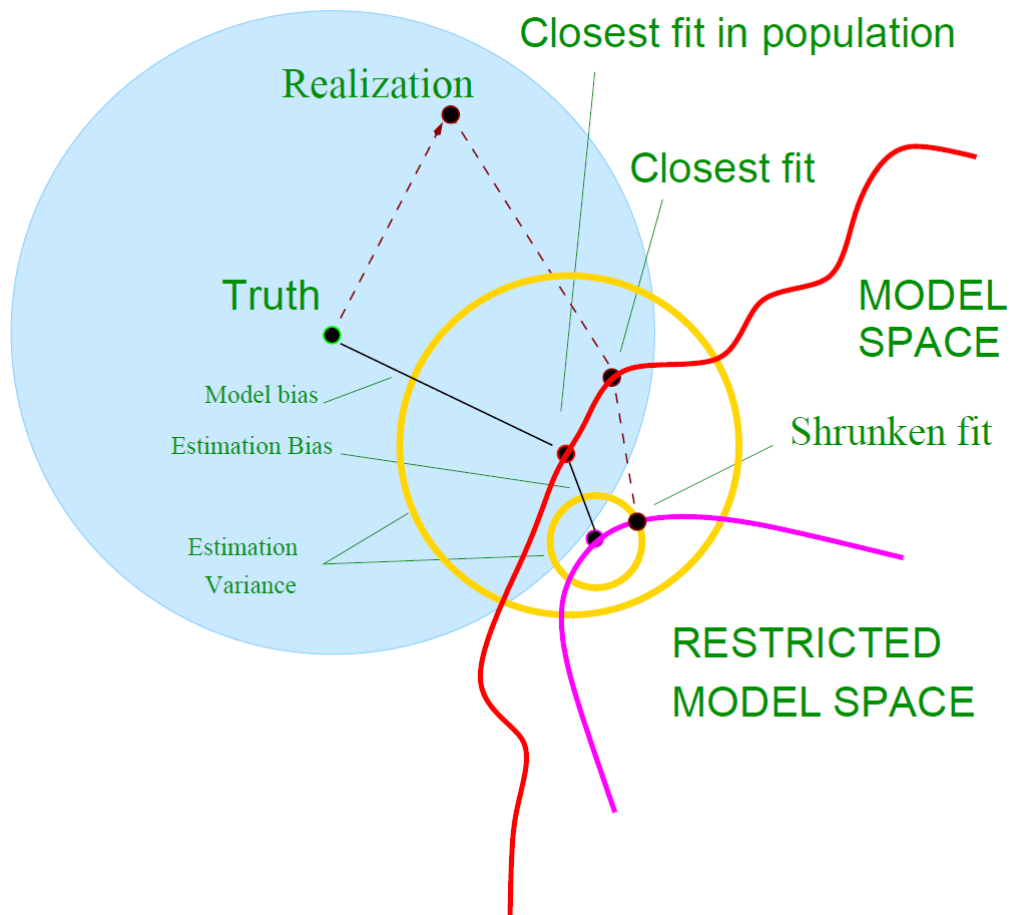
$$= E[(Y - f(x_0) + f(x_0) - \hat{f}(x_0))^2]$$

$$= E[\varepsilon^2] + E[(f(x_0) - \hat{f}(x_0))^2]$$

$$= \sigma_\varepsilon^2 + E[f(x_0) - E(\hat{f}(x_0)) + E(\hat{f}(x_0)) - \hat{f}(x_0)]^2$$

$$= \sigma_\varepsilon^2 + [f(x_0) - E(\hat{f}(x_0))]^2 + E(\hat{f}(x_0) - E[\hat{f}(x_0)])^2$$

$$= \sigma_\varepsilon^2 + \text{Bias}^2 + \text{Var}(\hat{f}(x_0))$$



Optimism of the Training Error Rate

$$\cdot \text{Err}_{\text{in}} = \frac{1}{N} \sum_{i=1}^N E_{Y^0} (Y_i^0 - \hat{f}(x_i))^2$$

$$\bar{\text{err}} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{f}(x_i))^2$$

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

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

\hat{y}_i is affected by y_i

Y_i^0 is not affected by y_i

proof:

$$w = E_y \left(\frac{1}{N} \sum_{i=1}^N E_{Y^0} [L(Y_i^0, \hat{f}(x_i) | T)] \right) - \frac{1}{N} \sum_{i=1}^N L(y_i, \hat{f}(x_i))$$

$$= E_y \left(\frac{1}{N} \sum_{i=1}^N E_{Y^0} [(Y_i^0 - \hat{y}_i)^2] \right) - \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

$$= \frac{1}{N} \sum_{i=1}^N (E_y E_{Y^0} [(Y_i^0)^2] + E_y E_{Y^0} [\hat{y}_i^2] - 2 E_y E_{Y^0} [Y_i^0 \hat{y}_i] \\ - E_y [y_i^2] - E_y [\hat{y}_i^2] + 2 E_y [y_i \hat{y}_i])$$

$$E_y E_{Y^0} [(Y_i^0)^2] = E_y [y_i^2]$$

$$= \frac{2}{N} \sum_{i=1}^N (E(y_i \hat{y}_i) - E(y_i) E(\hat{y}_i))$$

$$E_y E_{Y^0} [Y_i^0 \hat{y}_i] = E(y_i) E(\hat{y}_i)$$

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

$$\cdot \sum_{i=1}^N \text{Cov}(\hat{y}_i, y_i) = d \cdot \sigma_{\epsilon}^2 \quad (d \text{ inputs})$$

Estimates of In-sample Prediction Error

- C_p -Statistics $C_p = \overline{err} + 2 \cdot \frac{d}{N} \hat{\sigma}_\varepsilon^2$

- AIC (Akaike information criterion)

$$-2 \cdot E(\log \Pr_{\hat{\theta}}(Y)) \approx -\frac{2}{N} \cdot E(\log \text{likelihood}) + 2 \cdot \frac{d}{N}$$

$$\Rightarrow AIC = -\frac{2}{N} \sum_{i=1}^N \log P_{\hat{\theta}}(y_i) + 2 \cdot \frac{d}{N}$$

The Effective Number of Parameters

- Linear fitting method: $\hat{y} = Sy$

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

- $\sum_{i=1}^N \text{Cov}(y_i, \hat{y}_i) = \text{tr}[\text{Cov}(y, \hat{y})] = \text{tr}[\text{Cov}(y, Sy)]$

$$= \text{tr}(\sigma_\varepsilon^2 \cdot S) = \sigma_\varepsilon^2 \text{tr}(S)$$

- General definition: $df(\hat{y}) = \frac{\sum_{i=1}^N \text{Cov}(\hat{y}_i, y_i)}{\sigma_\varepsilon^2}$

- With Penalty: $df(\alpha) = \sum_{m=1}^M \frac{\theta_m}{\theta_m + \alpha}$

θ_m is the m -th eigenvalue
of Hessian Matrix

The Bayesian Approach and BIC

- $$\text{BIC} = -2 \log \text{likelihood} + (\log N) \cdot d$$

$$= \frac{N}{\sigma_\varepsilon^2} \left[\overline{\text{err}} + \underbrace{(\log N) \cdot \frac{d}{N}}_{\text{AIC has } k=2} \cdot \sigma_\varepsilon^2 \right]$$

AIC has $k=2$

BIC has $k=\log N$

- A set of candidate model $\{\mathcal{M}_m\}_{m=1}^M$

Known distribution of parameter θ_m $\Pr(\theta_m | \mathcal{M}_m)$

- $$\Pr(\mathcal{M}_m | \mathcal{Z}) \propto \Pr(\mathcal{M}_m) \cdot \Pr(\mathcal{Z} | \mathcal{M}_m)$$

$$\propto \Pr(\mathcal{M}_m) \cdot \int \Pr(\mathcal{Z} | \theta_m, \mathcal{M}_m) \cdot \Pr(\theta_m | \mathcal{M}_m) d\theta_m$$

- Bayes factor:
$$\text{BF}(\mathcal{Z}) = \frac{\Pr(\mathcal{Z} | \mathcal{M}_m)}{\Pr(\mathcal{Z} | \mathcal{M}_l)} = \frac{\Pr(\mathcal{M}_m)}{\Pr(\mathcal{M}_l)} \cdot \frac{\Pr(\mathcal{Z} | \mathcal{M}_m)}{\Pr(\mathcal{Z} | \mathcal{M}_l)}$$

- With Laplace approximation:

$$\log \Pr(\mathcal{Z} | \mathcal{M}_m) = \log \Pr(\mathcal{Z} | \hat{\theta}_m, \mathcal{M}_m) - \frac{d_m}{2} \cdot \log N + O(1)$$

$$\text{BIC} = -2 \log \Pr(\mathcal{Z} | \mathcal{M}_m) = -2 \log \Pr(\mathcal{Z} | \hat{\theta}_m, \mathcal{M}_m) + d_m \cdot \log N$$

- Posterior probability of model \mathcal{M}_m :

$$\Pr(\mathcal{M}_m | \mathcal{Z}) = \frac{e^{-\frac{1}{2} \text{BIC}_m}}{\sum_{l=1}^M e^{-\frac{1}{2} \text{BIC}_l}} = \frac{\Pr(\mathcal{Z} | \mathcal{M}_m) \Pr(\mathcal{M}_m)}{\sum_l \Pr(\mathcal{Z} | \mathcal{M}_l) \Pr(\mathcal{M}_l)}$$

