

Math 342W Lecture 19

Bias-Variance Tradeoff in Regression Modeling ($y \in \mathbb{R}$)

Assumptions

I) S is a realization from Δ , a B.V. mean independent of \vec{x} $E[\Delta | \vec{x}] = 0$

$$\hookrightarrow E[Y | \vec{x}] = E[f(\vec{x}) + \Delta] = E[f(\vec{x})] + E[\Delta | \vec{x}]$$

$$= E[f(\vec{x})] = f(\vec{x}) \rightarrow \text{Conditional Expectation Formula (CEF)}$$

II) Homoskedasticity \approx Constant Variance Constant $\forall \vec{x}$

$$\hookrightarrow \text{Var}[\Delta | \vec{x}] = E[\Delta^2 | \vec{x}] - E[\Delta | \vec{x}]^2 = E[\Delta^2 | \vec{x}] = \sigma^2$$

Lets say we fit a model to \mathbb{D} and obtain g

$$y = g + e = g + (f - g) + S \Rightarrow e = (f - g) + S$$

$$Y = g + (f - g) + \Delta \Rightarrow E = (f - g) + \Delta$$

Realization
to
Random Variable

$$\text{Then we can define Bias}[\vec{x}_*] = E[Y_* - g(\vec{x}_*)] = E[E | \vec{x}_*]$$

$$= E[f - g + \Delta_* | \vec{x}_*] = f - g + E[\Delta_*] = f(\vec{x}_*) - g(\vec{x}_*)$$

$$\text{We can now define } \text{MSE}(\vec{x}_*) = E[(Y_* - g(\vec{x}_*))^2]$$

$$= E[Y_*^2 | \vec{x}_*] - 2g(\vec{x}_*)E[Y_* | \vec{x}_*] + E[g(\vec{x}_*)^2 | \vec{x}_*]$$

$$= f(\vec{x}_*)^2 + 2f(\vec{x}_*)E[\Delta_*] + E[\Delta_*^2] - 2g(\vec{x}_*)[f(\vec{x}_*) + E[\Delta_*]] + g(\vec{x}_*)^2$$

$$= f(\vec{x}_*)^2 + (-2g(\vec{x}_*)f(\vec{x}_*)) + g(\vec{x}_*)^2 + \sigma^2 = [f(\vec{x}_*) - g(\vec{x}_*)]^2 + \sigma^2$$

Now we assume randomness in $\Delta_1, \Delta_2, \dots, \Delta_n, \Delta_*$

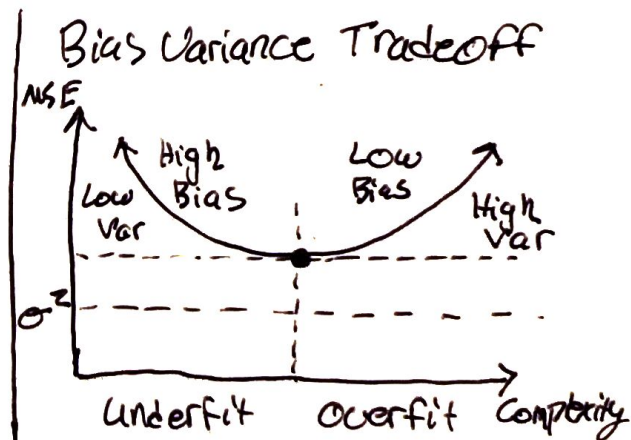
$$\begin{aligned}
 \text{MSE}(\vec{x}_*) &= E_{\Delta_1, \Delta_2, \dots, \Delta_n, \Delta_*} [(Y - G(\vec{x}_*))^2 | \vec{x}_*] \\
 &= E_{\Delta_1, \dots, \Delta_*} [Y^2] - 2 E_{\Delta_1, \dots, \Delta_*} [Y_* G(\vec{x}_*)] + E_{\Delta_1, \dots, \Delta_*} [G(\vec{x}_*)^2] \\
 &= [F(\vec{x}_*) + \sigma^2] - 2 E[Y_* G(\vec{x}_*)] + \text{Var}[G(\vec{x}_*)] + E[G(\vec{x}_*)]^2 \\
 &= [F(\vec{x}_*) - E[G(\vec{x}_*)]]^2 + \text{Var}[G(\vec{x}_*)] + \sigma^2 \\
 &= \sigma^2 + \text{Bias}[G(\vec{x}_*)]^2 + \text{Var}[G(\vec{x}_*)]
 \end{aligned}$$

Now, finally we assume randomness in all previous + X, \vec{x}_*

\Rightarrow This is a random variable model producing $\vec{x}_1, \dots, \vec{x}_n, \vec{x}_*$

$$E_{\vec{x}_1, \dots, \vec{x}_n, \vec{x}_*} [\text{MSE}(\vec{x}_*)]$$

$$\begin{aligned}
 &= E_{X^*s} [\text{Bias}[G(\vec{x}_*)]]^2 \\
 &\quad + E_{X^*s} [\text{Var}[G(\vec{x}_*)]] \\
 &\quad + \sigma^2
 \end{aligned}
 \left. \vphantom{\begin{aligned} &= E_{X^*s} [\text{Bias}[G(\vec{x}_*)]]^2 \\ &\quad + E_{X^*s} [\text{Var}[G(\vec{x}_*)]] \\ &\quad + \sigma^2 \end{aligned}} \right\} \begin{array}{l} \text{Bias} \\ \text{Variance} \\ \text{Decomposition} \end{array}$$



$$\bullet = \sigma^2 + \min \{ \text{bias}^2 + \text{var} \}$$

Overfit : Low Bias & High Variance

Underfit : High Bias & Low Variance

Generalized Linear Models (GLM)

Independently Realized

$$P(\mathcal{D}) = P(Y_1=y_1, Y_2=y_2, \dots, Y_n=y_n \mid \vec{X}_1=\vec{x}_1, \dots, \vec{X}_n=\vec{x}_n)$$

$$= \prod_{i=1}^n P(Y_i=y_i \mid \vec{X}_i=\vec{x}_i) = \prod_{i=1}^n f_{pr}(\vec{x}_i)^{y_i} (1-f_{pr}(\vec{x}_i))^{1-y_i}$$

\mathcal{R} : Maximize $P(\mathcal{D})$

$$v\text{-Bern}(\theta) = \theta^v (1-\theta)^{1-v}$$

Assume $\mathcal{B}_{pr} = \{ \phi(\vec{\omega} \cdot \vec{x}) : \vec{\omega} \in \mathbb{R}^p \}$

Assume $\phi: \mathbb{R} \rightarrow (0, 1)$, called a Link Function

If your algorithm uses the space $\mathcal{H} = \{ \phi(\vec{\omega} \cdot \vec{x}) : \vec{\omega} \in \mathbb{R}^p \}$, then your model is called a generalized linear model

Common Link Functions

① Logistic Link Function $\phi(u) = \frac{e^u}{1+e^u} = \frac{1}{1+e^{-u}}$

② Probit Link Function $\phi(u) = \Phi(u)$, where Φ is CDF of $N(0, 1)$

③ Complementary Log-Log Link Function $\phi(u) = 1 - e^{-e^{-u}}$

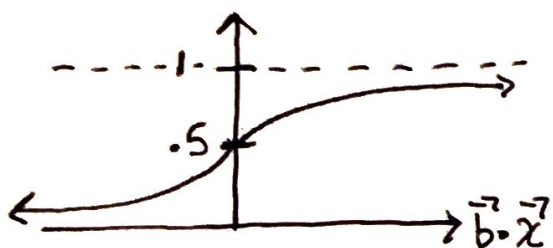
Assume Logistic Link

$$= \prod_{i=1}^n \left(\frac{1}{1+e^{-(\vec{\omega} \cdot \vec{x}_i)}} \right)^{y_i} \left(\frac{1}{1+e^{-(\vec{\omega} \cdot \vec{x}_i)}} \right)^{1-y_i} \Rightarrow$$

$$\mathcal{R}: \vec{b} = \operatorname{argmax}_{\vec{b}} \sum_i \mathcal{L}_i$$

NOTE: No closed form soln' to \vec{b} . Numerical Methods to approximate $\vec{\nabla} P(\mathcal{D}) := \vec{\partial}_{\mathcal{P}1}$

$$g_{pr}(\vec{x}) = \frac{1}{1 + e^{-(\vec{b} \cdot \vec{x})}} \Rightarrow \text{estimates } P(Y_* = 1 | \vec{x}_*)$$



$$\Rightarrow \frac{1}{\hat{p}} = 1 + e^{-(\vec{b} \cdot \vec{x})} \Rightarrow \frac{1}{\hat{p}} - 1 = e^{-(\vec{b} \cdot \vec{x})} \Rightarrow \frac{1 - \hat{p}}{\hat{p}} = e^{-(\vec{b} \cdot \vec{x})}$$

odds against

$$\Rightarrow \frac{\hat{p}}{1 - \hat{p}} = e^{\vec{b} \cdot \vec{x}} \Rightarrow \vec{b} \cdot \vec{x} = \ln\left(\frac{\hat{p}}{1 - \hat{p}}\right)$$

log odds

Now the question is, how do we validate this sort of model?

We have to validate $g_{pr}(\vec{x}_i)$ with y_i , but

$$g_{pr}(\vec{x}_i) = \hat{p} \in (0, 1) \text{ and } y_i \in \{0, 1\}$$

↗ ≠ ↖

Proper Scoring Rules

$S(\hat{p}_i, y_i)$ satisfies $\forall_i f_{pr}(\vec{x}_i) = \arg\max \{S(\hat{p}_i, y_i)\}$

① Brier Score : $S_i = -(y_i - \hat{p}_i)^2 \leq 0$

② Log Score : $S_i = y_i \ln(\hat{p}_i) + (1 - y_i) \ln(1 - \hat{p}_i) \leq 0$

These scoring rules deal with the different spaces for \hat{p}_i and y_i , so we can judge the model performance in a way that makes sense.