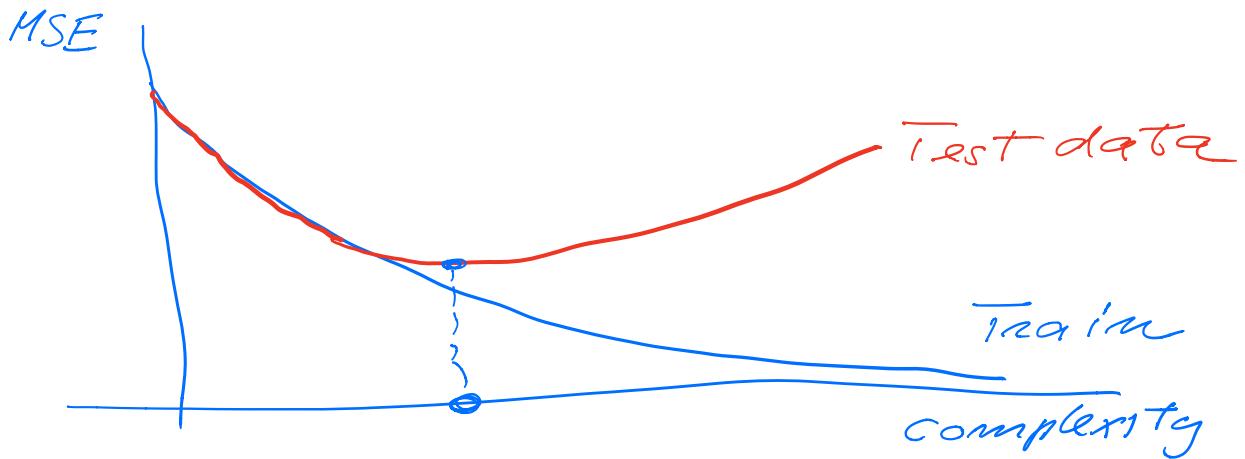


Friday August 28

Review of statistics & Probability Theory



$$MSE = \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \tilde{y}_i)^2$$

$$= E[(y - \tilde{y})^2]$$

$$= \frac{1}{n} \| (y - \tilde{y}) \|_2^2$$

$$\| x \|_2 = \sqrt{\sum_i x_i^2}$$

$$\| x \|_2^2 = \sum_i x_i^2$$

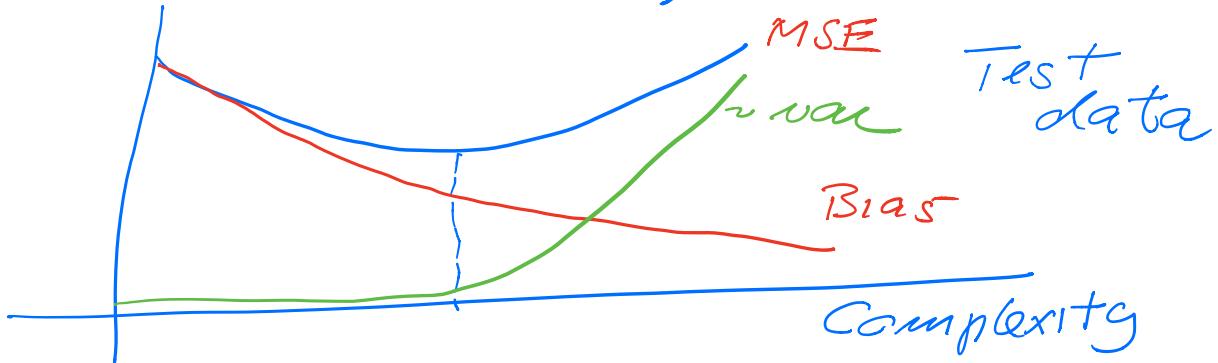
Bias-Variance trade-off

$$MSE = E[(y - E[\tilde{y}])^2] + var(\tilde{y})$$

Bias + σ^2

$$y = f(\cdot) + \varepsilon$$

$$\varepsilon \sim N(0, \sigma^2)$$



Trade-off between complexity and the # of data.

Domain $\mathcal{D} = \{(x_0, y_0), (x_1, y_1) \dots (x_{n-1}, y_n)\}$

$$x \in X \quad y \in Y$$

assume we have a PDF

$$p(x) \text{ and } p(y)$$

$$0 \leq p(x) \leq 1 \quad x \in [a, b]$$

$$\int_a^b p(x) dx = 1$$

Discrete case

$$\sum_{x_i \in X} p(x_i) = 1$$

Expectation value, moments

$$E[x^n] = \int_a^b x^n p(x) dx$$

\bar{x}

$$\sum_{x_i \in X} x_i^n p(x_i)$$

$$M_x = \underline{E[x]} = \langle x \rangle = \mu$$

Variance

$$\text{var}(x) = \sigma_x^2 = \sigma^2 = \int_a^b (x - \mu_x)^2 p(x) dx$$

$$\left(\sum_{x_i \in X} (x_i - M_x)^2 p(x_i) \right)$$

Sample mean

$$\bar{\mu}_x = \frac{1}{n} \sum_{i=0}^{n-1} x_i' \neq \mu_x$$

sample variance

$$\bar{\sigma}_x^2 = \frac{1}{n} \sum_{i=0}^{n-1} (x_i' - \bar{\mu}_x)^2$$

Covariance

$$\text{cov}(x, y) = \iint dx dy p(x, y) \times (x - \mu_x)(y - \mu_y)$$

$$= \frac{\iint dx dy p(x, y) xy}{\mu_x \mu_y}$$

$$= \mathbb{E}[xy] - \mu_x \mu_y$$

i.i.d = independent and identically distributed

$$P(x,y) = \underline{P(x)} \underline{P(y)}$$

\ same distributions

$$\mathbb{E}[xy] = \frac{\int dx p(x) x \int dy p(y) y}{n}$$

$$\mu_x \cdot \mu_y = \mu^2$$

$$\underline{\text{cov}(x,y)} = 0$$

$$\frac{\int \int dx dy x y P(x,y)}{n}$$

$$\mu_x = \int dx p(x)x = \mu_y = \mu$$

$$P(x,y) = \underline{P(x)} \underline{P(y)}$$

if i.i.d. then $\text{cov}(x,y) = 0$

Covariance matrix (2 features)

$$C[x,y] = \begin{bmatrix} \text{cov}(x,x) & \text{cov}(x,y) \\ \text{cov}(y,x) & \text{cov}(y,y) \end{bmatrix}$$

$$\text{cov}(x, x) = \text{var}(x)$$

$$C[x, y] = \begin{bmatrix} \text{var}(x) & \text{cov}(x, y) \\ \text{cov}(y, x) & \text{var}(y) \end{bmatrix}$$

Introduce a scaling function
correlation function

$$K[x, y] = \begin{bmatrix} \frac{\text{cov}(x, x)}{\sqrt{\text{var}(x)\text{var}(x)}} & \frac{\text{cov}(x, y)}{\sqrt{\text{var}(x)\text{var}(y)}} \\ \frac{\text{cov}(y, x)}{\sqrt{\text{var}(x)\text{var}(y)}} & \frac{\text{cov}(y, y)}{\sqrt{\text{var}(y)\text{var}(y)}} \end{bmatrix}$$

\downarrow take values between -1 and 1.

More Probability ---

X or Y

$$P(X \vee Y) = P(X) + P(Y) - \underline{P(X \wedge Y)}$$

Joint Probability

$$P(X, Y) = P(X|Y)P(Y) \\ P(Y|X)P(X)$$

Marginal (sum rule) distribution

$$P(X) = \sum_{y_i \in Y} P(X|Y=y_i) P(Y=y_i)$$

Conditional probability

$$\underbrace{P(X|Y)}_{\text{Posterior likelihood}} = \frac{\underbrace{P(X,Y)}_{\text{Likelihood}}}{\sum_{x_i \in X} P(Y|x=x_i) P(X=x_i)} \quad P(Y) > 0$$

prior

Baye's theorem

Example : mammogram

sensitivity : $X=1$ positive event
 $X=0$ negative

$Y=1$	Breast cancer
$Y=0$	no breast cancer

$\sim 1\% \sim 1\% \rightarrow \sim 0 \sim 2\%$

$$\underline{P(X=1 | Y=1)} = 0.8 \quad \dots$$

$$\underline{P(Y=1 | X=1)} = ?$$

$$P(Y=1) = 0.004 \quad \begin{matrix} P(Y=0) = \\ 1 - P(Y=1) \end{matrix}$$

$$P(X=1 | Y=0) = 0.1$$

Bayes' theorem

$$P(Y=1 | X=1) =$$

$$\frac{P(X=1 | Y=1) P(Y=1)}{P(X=1 | Y=1) P(Y=1) + P(X=1 | Y=0) P(Y=0)}$$

$$= 0.03 \Rightarrow 3\%$$

Maximum Likelihood Estimator

$$P(Y|X) = ?$$

$$\text{Model} \quad P(Y|X, \beta) \rightarrow$$

$$P(D|\beta) \rightarrow \text{iid } (x_i, y_i)$$

$n=1$

$$P(D|\beta) \propto \prod_{i=0}^{n-1} P(y_i|x_i|\beta)$$

β which maximizes $P(D|\beta)$

$$\hat{\beta} = \arg \max_{\beta \in \mathbb{R}^P} P(D|\beta)$$

$$= \arg \max_{\beta \in \mathbb{R}^P} \sum_{i=0}^{n-1} \log P(y_i|x_i|\beta)$$

Instead of maximizing $Z^m g$,
we can minimize

$$C(\beta) = - \sum_{i=0}^{n-1} \log P(y_i|x_i|\beta)$$

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^P} C(\beta)$$

$$P(y_i|x_i|\beta) = ??$$

$$y_i = f(x_i) + \varepsilon_i$$

$$\varepsilon_i \sim N(0, \Sigma^2)$$

$$x_i \in \mathbb{R}^P$$

$$\beta \in \mathbb{R}^P$$

$$E[y_i] = E[x_i|\beta] + E[\varepsilon_i]$$

$$y_i = x_{i*}\beta + \epsilon_i$$

$$\text{var}[y_i] = E[(y_i - E[y_i])^2]$$

$$= E[y_i^2] - (x_{i*}\beta)^2$$

$$= E\left[\frac{(x_{i*}\beta)^2 + 2\sum x_{ij}\beta}{+\epsilon_i^2}\right] = 0$$

$$- (x_{i*}\beta)^2$$

$$\text{var}[y_i] = \sigma^2$$

Assumption is that

$$y_i \sim N(x_{i*}\beta, \sigma^2)$$

$$= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(y_i - x_{i*}\beta)^2}{2\sigma^2}\right]$$

$$= P(y_i | x_i | \beta)$$

$$C(\beta) = \frac{m}{2} \log(2\pi\sigma^2)$$

$$+ \sum_{i=0}^n \left(\frac{(y_i - x_i * \beta)^2}{2\sigma^2} \right)$$

$$\frac{\partial C(\beta)}{\partial \beta} = 0 \Rightarrow$$

$$x^T(y - x\beta) = 0 \Rightarrow$$

$$\hat{\beta} = (x^T x)^{-1} x^T y$$

$$y_i = x_{i*}\beta + \varepsilon_i$$

$$i \rightarrow \begin{bmatrix} x_{i0} & x_{i1} & \dots & x_{ip-1} \\ x_{i0} \\ \vdots \\ x_{i,n-1} \end{bmatrix} \quad \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_{p-1} \end{bmatrix}$$