Comp Sci, Nov 23, 2022

- statistical interpretation
- Resampling techniques

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

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

$f(x)$ is a deterministic function

$$D = \{(x_0, y_0), (x_1, y_1), \ldots (x_{n-1}, y_{n-1})\}$$

Ideally we have PDF $p(x)$, $p(y)$ etc.

$$E[y] = \int_D p(y)\, y\, dy = \mu_y$$

$$\left( \sum_{i \in D} p(y_i) y_i \right)$$

$$var[y] = \int_D p(y)(y - \mu_y)^2$$

$$cov[y, y'] = \int_D p(y, y')(y - \mu_{y})$$
$$\times (y' - \mu_{y'})\, dy\, dy'$$

Sample mean:

$$\bar{\mu}_y = \frac{1}{n} \sum_{i=0}^{n-1} y_i' \neq \mu_y$$

sample variance

$$\text{var}[\bar{y}] = \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \bar{\mu}_y)^2$$
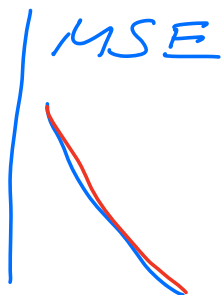
$$\neq \text{var}[\tilde{y}_g^2]$$

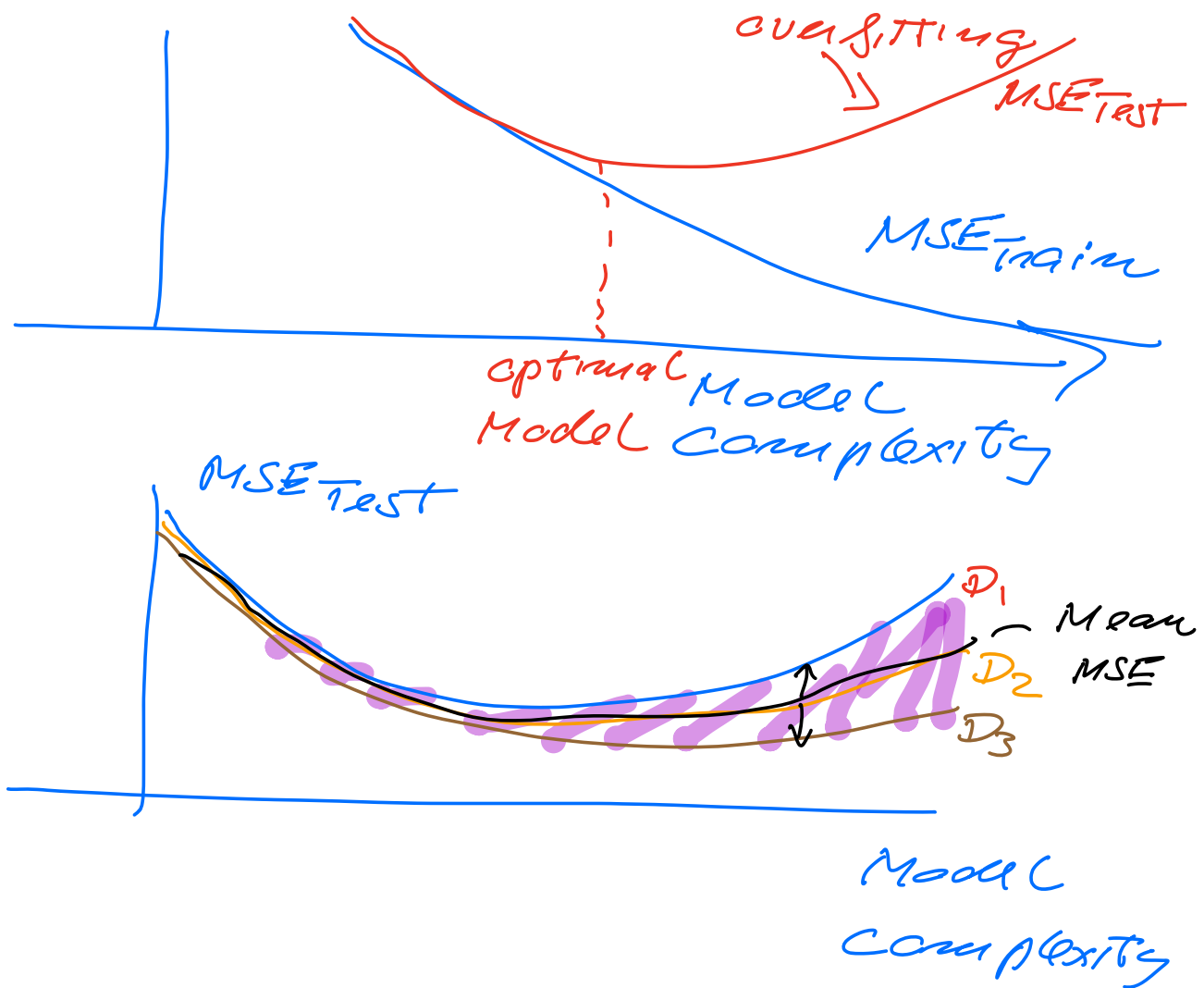$$\mathbb{E}[\bar{y}] = \frac{1}{n} \sum_{i=0}^{n-1} y_i'$$

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

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

Resampling techniques
  — Bootstrap
  — Cross-validation
Look at MSE as example

overfitting

$MSE_{Test}$

$MSE_{Train}$

optimal Model
Model Complexity

$MSE_{Test}$

$D_1$

Mean
MSE

$D_2$

$D_3$

Model
Complexity

# Bootstrap Strategy

Require $D = \{ (x_0, y_0) \ldots (x_n, y_n) \}$

Split in train and test data.

Require $M = $ # bootstrap samples

Have defined $D_{train}$ and

$D_{Test}$

For $i = 1, M$

- Make $D_{Train}(i)$ by randomly selecting with replacement.
$$\begin{cases} D_{Train} = [1, 2, 3, 4, 5] \\ D_{Train}^* = [2, 3, 2, 4, 1] \end{cases}$$

- Train Model $(i)$

- compute $MSE(i)$ (on train) and $MSE_{Test}(i)$

end loop

- compute
$$MSE_{Test} = \frac{1}{M} \sum_{i=0}^{M-1} MSE_{Test}(i)$$

Cross-validation

Define folds = K

Example K = 5

T = Train
T = Test

$D_1$: | T | T | T | T | T |   MSE$_{Test}^{(1)}$

$D_2$: | T | T | T | T | T |   MSE$_{Test}^{(2)}$

$D_3$: | T | T | T | T | T |   ,

$D_4$: | T | T | T | T | T |   ,

$D_5$: | T | T | T | T | T |   MSE$_{Test}^{(5)}$

$$MSE_{Test} = \frac{1}{K} \sum_{i=0}^{K-1} MSE_{Test}^{(i)}$$

———————— ✗ ————————

STATISTICS

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

$$\mathbb{E}[y_i] = \mathbb{E}[f(x_i)] + \mathbb{E}[\varepsilon_i]$$

$$\overset{\shortparallel}{?} \qquad \overset{\shortparallel}{0}$$

$$f(x_i) \simeq \sum_{j=0}^{p-1} x_{ij}\,\beta_j = x_{i*}\beta$$

$$\mathbb{E}[y_i] = x_{i*}\beta \implies \mathbb{E}[y] = X\beta$$

$$\mathrm{var}[y] = \sigma^2 = \mathrm{var}[\varepsilon]$$
( EXERCISE )

$$\mathbb{E}[\beta] = \beta$$

$y_i$ has mean value $X\beta$
and variance $\sigma^2$

$$y \sim N(X\beta, \sigma^2) \implies$$

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

$$= \frac{1}{\sqrt{2\pi\sigma^2}}\, e^{-\frac{(y_i - x_{i*}\beta)^2}{2\sigma^2}}$$

$$= p(y_i \mid \beta)$$

assume $y_i$ are i.i.d.

$$P(D \mid \beta) = \prod_{i=0}^{u-1} P(y_i \mid \beta)$$

$$P(y_i \mid \beta) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - x_{i*}\beta)^2}{2\sigma^2}} \sim -\frac{1}{2} \log$$

$$\hat{\beta} = \underset{\beta \in \mathbb{R}^p}{\arg\max} \; P(D \mid \beta)$$

$$\hat{\beta} = \underset{\beta \in \mathbb{R}^p}{\arg\max} \; \log P(D \mid \beta)$$

on

$$\hat{\beta} = \underset{\beta \in \mathbb{R}^p}{\arg\min} \left( -\log P(D \mid \beta) \right)$$

$$= \sum_{i=0}^{u-1} \left[ \frac{1}{2} \log(2\pi\sigma^2) + \frac{(y_i - x_{i*}\beta^2)}{2\sigma^2} \right]$$

$$= \frac{u}{2} \log(2\pi\sigma^2) + \left\| \frac{(y - X\beta)}{2\sigma^2} \right\|_2^2$$

Taking derivatives wrt $\beta$

$$\Rightarrow \quad \hat{\beta} = \left(X^T X\right)^{-1} X^T \beta$$

$$\left( MSE = \frac{1}{m} \sum_{i=0}^{m-1} \left(y_i - x_{i*}\beta\right)^2 \right.$$

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

# Ridge & Lasso ?

we are going to make
an ansatz (prior) about
the distribution of $\beta$.

Ridge $\quad p(\beta) = \frac{1}{\sqrt{2\pi\gamma^2}} e^{-\beta^2/2\delta^2}$

Lasso $\quad p(\beta) \approx e^{-|\beta|/2\gamma}$
$\qquad$ (Laplace distributi)

Bayes' theorem,

— product rule of probabilities
$\quad p(A,B) = p(A|B)\, p(B)$

$$= P(B|A)\,P(A)$$

- Conditional probability

$$P(A/B) = \frac{P(A,B)}{P(B)}$$

$$\text{if } P(B) > 0$$

$$P(B/A) = \frac{P(A,B)}{P(A)} \quad P(A) > 0$$

if A and B are i.i.d.

$$P(A, B) = P(A)\,P(B)$$

- Marginal distribution

$$P(A) = \sum_{b} P(A/B=b)\,P(B=b)$$

$$P(B) = \sum_{a} P(B/A=a)\,P(A=a)$$

Combining we have
Bayes' theorem

$$P(A|B) = \frac{P(A,B)}{P(B)} =$$

$$\frac{P(B|A)\,P(A)}{\sum_a P(B|A=a)\,P(A=a)}$$

likelihood      prior

posterior ↑

what is optimal $B$
given $D$

$$P(B|D) \propto P(D|B)\,P(B)$$

$$N(X\beta, \sigma^2)$$