

FYS - SIK 4155 September 18

Resampling and statistical analysis

Central Limit theorem

a set of data of i.i.d. events,

$$E[x_i] = \bar{x}_i = \sum p(x_i) x_i$$

simplifying $\bar{x}_i \rightarrow x_i$

$$Z = \frac{x_1 + x_2 + \dots + x_m}{m}$$

x_i they follow $p(x_i)$

what form does $\mathcal{P}(Z)$ take?

$$\mathcal{P}(Z) = \int dx_1 p(x_1) \int dx_2 p(x_2)$$

$$\dots \int dx_m p(x_m)$$

$$\times \delta\left(Z - \frac{x_1 + x_2 + \dots + x_m}{m}\right)$$

$$\delta\left(z - \frac{x_1 + x_2 + \dots + x_m}{m}\right)$$

$$= \frac{1}{2\pi i} \int_{-\infty}^{\infty} dq \exp\left(iq\left(z - \frac{x_1 + \dots + x_m}{m}\right)\right)$$

$$E[x_i] = \mu$$

insert $e^{i\mu q} e^{-i\mu q}$

$$\mathcal{P}(z) = \frac{1}{2\pi i} \int_{-\infty}^{\infty} dq e^{iq(z-\mu)}$$

$$\times \left[\int_{-\infty}^{\infty} dx p(x) e^{iq \frac{(\mu-x)}{m}} \right]^m$$

$$\int_{-\infty}^{\infty} dx p(x) e^{iq(\mu-x)/m}$$

$$= \int_{-\infty}^{\infty} dx p(x) \left[1 + iq \frac{(\mu-x)}{m} - \frac{q^2 (\mu-x)^2}{2m^2} + \dots \right]$$

$$\int dx p(x) x = \mu \quad \wedge \quad \int dx p(x) = 1$$

$$= \left[1 - \frac{\sigma^2 \sigma^2}{2m^2} + \dots \right]$$

$$P(z) = \frac{1}{\sqrt{2\pi \sigma^2/m}} e^{-\frac{(z-\mu)^2}{2\sigma^2/m}}$$

$$\begin{aligned} \text{var}[x] &= \int p(x) (x-\mu)^2 dx \\ &= \sigma^2 \end{aligned}$$

$$\text{var}[z] = \frac{\sigma^2}{m} \Rightarrow$$

$$\text{std} = \sigma/\sqrt{m}$$

Bootstrap Resampling

original sample

$$D = \{z_0, z_1, \dots, z_{n-1}\}$$

(i) Draw a bootstrap sample

$$D_1^* = \{z_0^*, z_1^*, \dots, z_{n-1}^*\}$$

$$\text{compute } \Theta_1^* = \frac{1}{n} \sum_{i=0}^{n-1} z_i^*$$

(ii) Repeat B times
yielding estimators

$$\Theta_1^* \quad \Theta_2^* \quad \dots \quad \Theta_B^*$$

$$\bar{\Theta} = \frac{1}{B} \sum_{j=1}^B \Theta_j^*$$

(iii) can compute
variance

$$S^2 = \frac{1}{B} \sum_{j=1}^B (\Theta_j^* - \bar{\Theta})^2$$

(IV) output S^2 and $\bar{\Theta}$

The bootstrap shows that
it approaches the true
mean and true variance
and other expectation
values.

Bias-variance tradeoff

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

$$\begin{aligned}
&= \text{Bias}[(y_i - \mathbb{E}[\tilde{y}])^2] \\
&\quad + \text{var}[\tilde{y}] + \sigma^2 \\
y &= f(x) + \underset{\substack{\uparrow \\ \text{var}[\varepsilon] = \sigma^2}}{\varepsilon}
\end{aligned}$$

$f(x)$ is deterministic.

$$\begin{aligned}
f(x_i) &\simeq \tilde{y}_i & \mathbb{E}[f] &= f \\
\mathbb{E}[\tilde{y}_i] &= \tilde{y}_i & &= \mathbb{E}[y]
\end{aligned}$$

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

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

if we know p

$$\text{MSE} = \int_D (y - \tilde{y})^2 p(x) dx$$

$$= \int_D (f + \varepsilon - \tilde{y})^2 p(x) dx$$

add and subtract
 $E[\tilde{y}]$

$$\int_D (f + \varepsilon - \tilde{y} + E[\tilde{y}] - E[\tilde{y}])^2 p(x) dx$$

$$E[\varepsilon] = 0 \quad E[\tilde{y}] = \mu_{\tilde{y}}$$

$$= \boxed{\int_D (f(x) - \mu_{\tilde{y}})^2 p(x) dx} \quad \text{Bias}$$

$$+ \boxed{\int_D (\tilde{y}(x) - \mu_{\tilde{y}})^2 p(x) dx} \quad \text{var}[\tilde{y}]$$

$$+ \text{var}[\varepsilon^2]$$

$\sim \sigma^2$

We have a discrete set
and $p(x)$ is unknown
 \Rightarrow sample expectation
values (Reason for
resampling)

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

$$\begin{aligned}
&= \frac{1}{n} \sum_{i=0}^{n-1} \underbrace{(y_i - E[\hat{y}])^2}_{\text{Bias}} \\
&\quad + \frac{1}{n} \sum_{i=0}^{n-1} \underbrace{(\hat{y}_i - E[\hat{y}])^2}_{\text{variance of model}} \\
&\quad + \sigma^2
\end{aligned}$$

Done on test data,

CV

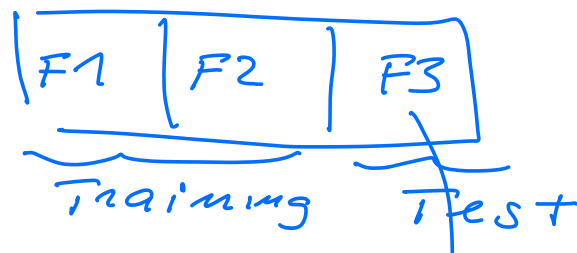
$$X = [1, 2, 3, 4, 5, 6]$$

$K = 3$ number of folds

Fold 1 $[5, 6]$

Fold 2 $[1, 3]$

Fold 3 $[4, 6]$



$$1) \quad \text{Error}_{\text{Test}} = \overset{\downarrow}{\text{Error}_{F_3}}$$

$$2) \quad \underbrace{\left[F_1 \mid F_2 \mid F_3 \right]}_{\text{Test Training}}$$

$$\text{Error}_{\text{Test}} = \text{Error}_{F_1}$$

$$3) \quad \underbrace{\left[F_1 \mid F_2 \mid F_3 \right]}_{\text{Train Test Train}}$$

$$\text{Error}_{\text{Test}} = \text{Error}_{F_2}$$

$$\text{Total error} = \frac{1}{3} \left(\text{Error}_{F_1} + \text{Error}_{F_2} + \text{Error}_{F_3} \right)$$

$$K \sim 5-10$$