

$$\text{Var}(aX + bY) = a^2 \text{Var}(X) + 2ab \text{Cov}(X, Y) + b^2 \text{Var}(Y)$$

chp 5 #1

$$\alpha = \frac{\sigma_y^2 - \sigma_{xy}}{\sigma_x^2 + \sigma_y^2 - 2\sigma_{xy}}, \quad \begin{aligned} \sigma_x^2 &= \text{Var}(X) \\ \sigma_y^2 &= \text{Var}(Y) \\ \sigma_{xy} &= \text{Cov}(X, Y) \end{aligned}$$

show α minimizes

$$\begin{aligned} &\text{Var}(\alpha X + (1-\alpha)Y) \\ &= \alpha^2 \sigma_x^2 + (1-\alpha)^2 \sigma_y^2 + 2\alpha(1-\alpha)\sigma_{xy} \\ \frac{d}{d\alpha} &= [2\alpha \sigma_x^2 + (1-2\alpha + \alpha^2)\sigma_y^2 + (2\alpha - 2\alpha^2)\sigma_{xy}] \\ &= 2\alpha \sigma_x^2 + \sigma_y^2 - 2\alpha \sigma_y^2 + 2\alpha \sigma_{xy} - 2\alpha^2 \sigma_{xy} = 0 \end{aligned}$$

Solving for α

$$\begin{aligned} 2\alpha \sigma_x^2 + \sigma_y^2 - 2\alpha \sigma_y^2 + 2\alpha \sigma_{xy} - 2\alpha^2 \sigma_{xy} &= 0 \\ 2\alpha(\sigma_x^2 + \sigma_{xy} - \sigma_y^2) &= \sigma_y^2 - \sigma_{xy} \\ \alpha &= \frac{\sigma_y^2 - \sigma_{xy}}{\sigma_x^2 + \sigma_{xy} - \sigma_y^2} \quad \checkmark \end{aligned}$$

So this indeed is a minimum

$$\frac{d^2}{d\alpha^2} (\text{Var}(\alpha X + (1-\alpha)Y))$$

$$= 2\sigma_x^2 + 2\sigma_y^2 - 4\sigma_{xy} = 0$$

$$= 2(\sigma_x^2 + \sigma_y^2 - 2\sigma_{xy}) = 0$$

$$\begin{aligned} \text{Var}(X) + \text{Var}(Y) - 2\text{Cov}(X, Y) &\geq 0 \\ 2\text{Var}(X - Y) &\geq 0 \end{aligned}$$

2 a) Probability that the first bootstrap obs
is not the jth obs
 $1 - \frac{1}{n}$

b) $1 - \frac{1}{n}$ for second obs

c) this is because the bootstrap sample
w/ replacement is the product of the observations
so $(1 - \frac{1}{n})_1 (1 - \frac{1}{n})_2 \dots (1 - \frac{1}{n})_n$
 $= (1 - \frac{1}{n})^n$

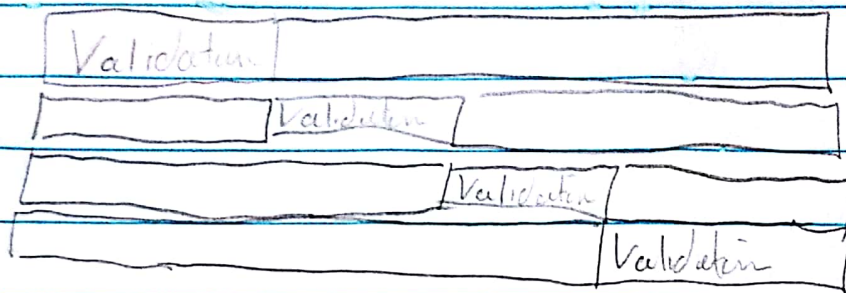
d) when $n=5$ $P(\text{jth obs}) = ?$ in Bootstrap sample
 $\Rightarrow 1 - (1 - \frac{1}{5})^5 = 1 - (.8)^5 = 1 - .32768$
 $= .67232 \approx .672$

e) How about $n=100$?
 $1 - (1 - \frac{1}{100})^{100} = 1 - (.99)^{100} = .634$

f) $n=1000$?
 $1 - (1 - \frac{1}{1000})^{1000} = .6323$

g) R Code

K-fold CV is implemented by randomly dividing the observations into K groups. Each fold has a Validation & Training set that are used to estimate the MSE for that fold. This is repeated for each fold.



$K=4$ CV

b) What are advantages/disadvantages of K-fold CV

i) Compared to the Validation Set approach
In the Validation Set approach there is a lot of variance compared to K-fold CV
So K-fold CV has the advantage

ii) Compared to LOOCV

LOOCV costs more time to calculate compared to K-fold CV, K-fold CV has far less training fit time/iterations. Also LOOCV may give a biased estimate because of the overfitting (high correlation of training sets)