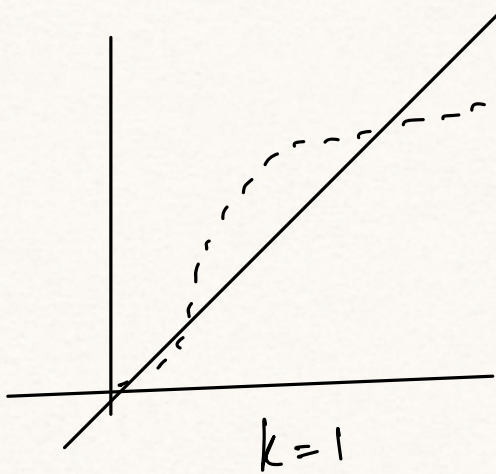
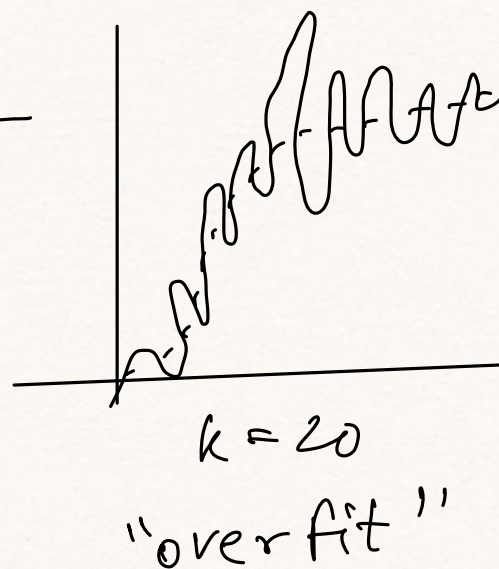
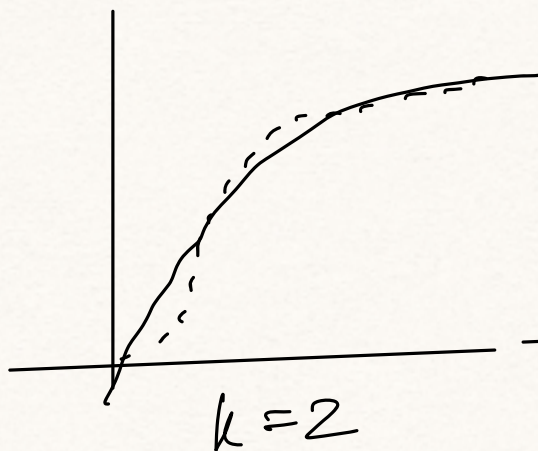


BIAS - VARIANCE

Generalization error (G.E)



"underfit"



- Bias :- Component of G.E due to "expressivity handicap"
- Variance : due to finite sample of training set.

For case of
Squared Error

$$y^{(i)} = f(x^{(i)}) + \epsilon^{(i)}$$

$$E[\epsilon] = 0$$

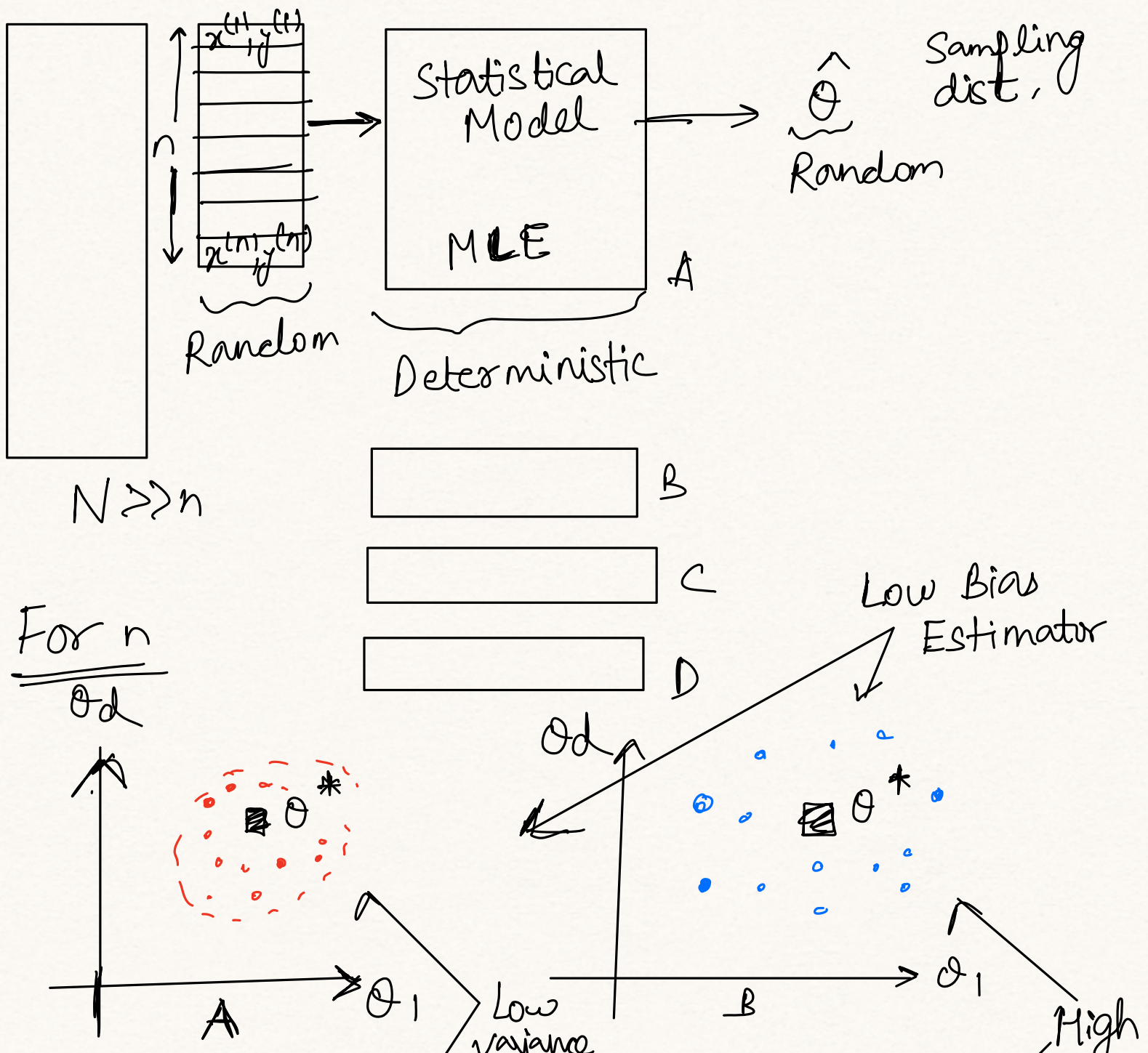
$$\text{Var}[\epsilon] = \sigma^2$$

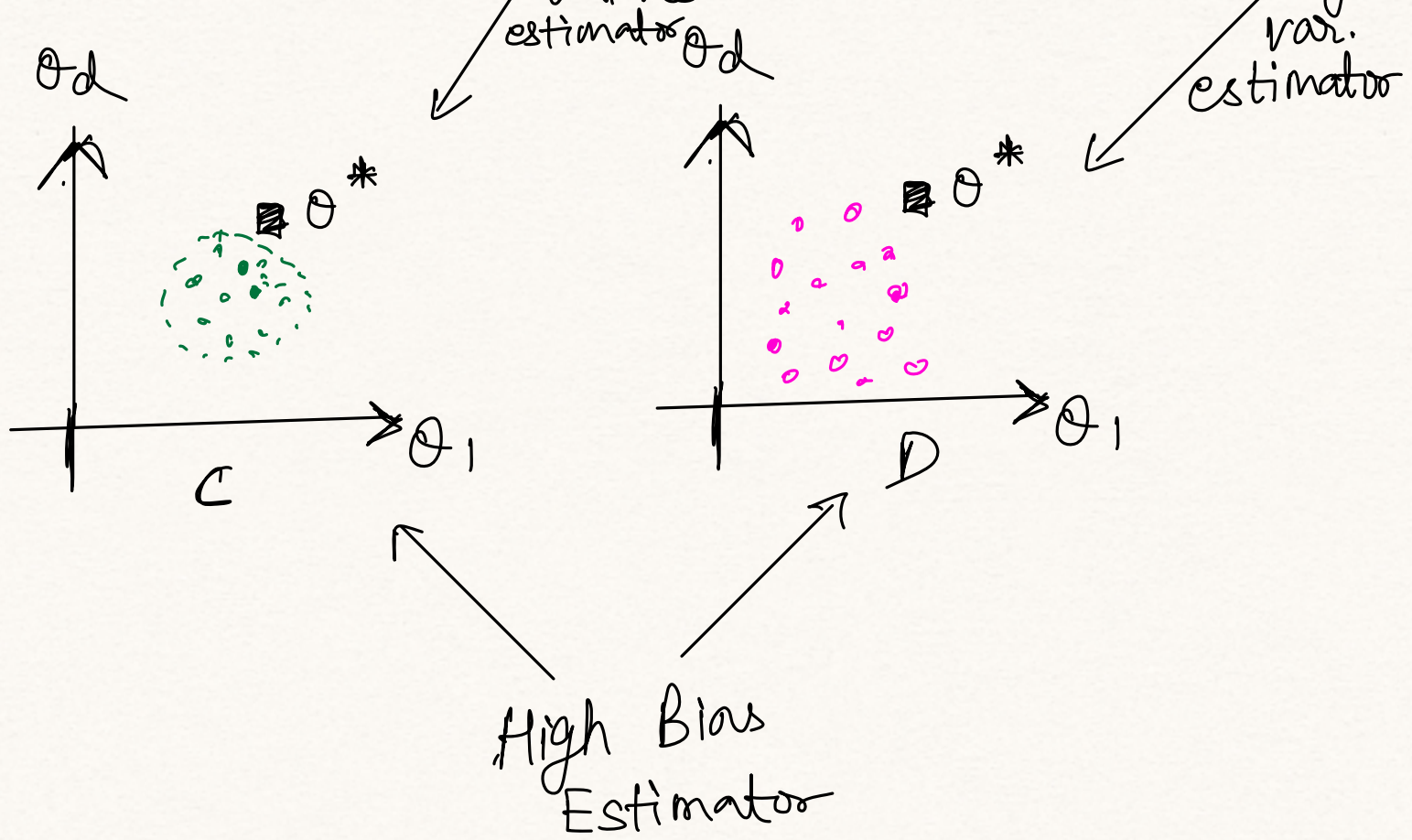
$$\text{Test Error [G.E]} = E[(y - \hat{f}_n(x))^2]$$

[expectation is over all ϵ in
training set & test eq.]

$$= \underbrace{\sigma^2}_{\text{irreducible error}} + \underbrace{\mathbb{E}[f(x) - \hat{f}_n(x)]^2}_{(\text{Bias})^2} + \underbrace{\text{Var}[\hat{f}_n(x)]}_{\text{variance}}$$

$$(x, y) \sim \text{Dist.}(\theta)$$





In case of N ,
 in each case variance will come down
 but bias remains unchanged.
 (comes down, but
 unbiased or biased
 thing remains
 unchanged)

$$E[\hat{\theta} - \theta^*] = \text{Bias}$$

$$\text{Var}[\hat{\theta}] = \text{Variance}$$

If Bias $\rightarrow 0$ as $n \rightarrow \infty$
 \Rightarrow Consistent Estimator

$$\lim_{n \rightarrow \infty} \text{Var}[\hat{\theta}] \xrightarrow{\text{rate}} 0$$

↑ statistical efficiency $\frac{1}{n} / \frac{1}{\sqrt{n}}$

Underfitting \approx High Bias
Overfitting \approx High Variance

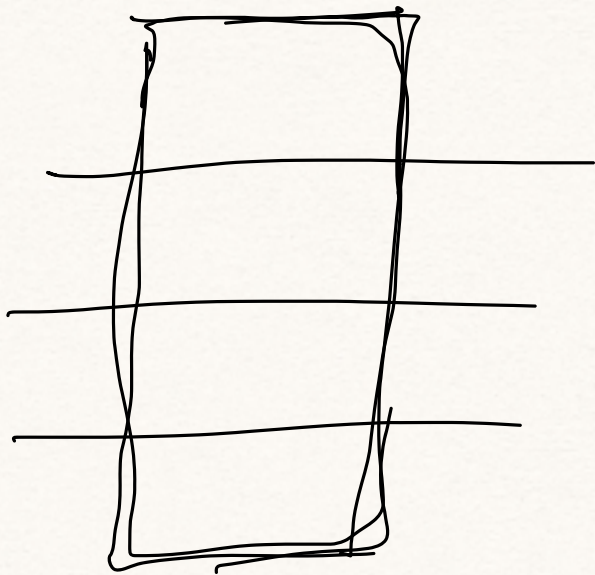
Goal: Do well on Generalization error.

— Cross validation

Split it into $\left\{ \begin{array}{l} \text{Train} \rightarrow \text{Do well on training set} \\ \text{valid/dev} \rightarrow \text{Do well in G.E.} \\ \text{Test} \rightarrow \text{Get an estimate on G.E.} \end{array} \right.$

K-fold Cross validation

Small Data sets



Model-1

fold-1 : valid

fold 2-n : training

Model-2

fold-2 : valid

fold 1,3-n : training.

leave-one-out C.V ($k=n$)

Regularization

Encouraging small $\|\theta\|$

$$J(\theta) = \sum_{i=1}^n (y^{(i)} - h_{\theta}(x^{(i)}))^2 + \lambda \|\theta\|_2^2$$

$$h_{\theta}(x) = \theta^T x$$

$$\sum_{i=1}^d \theta_i^2$$

$$\|\theta\|_1 = \sum_{i=1}^d |\theta_i|$$

L_2 -Regularized Linear Regression

$$J(\theta) = \left(\sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2 \right) + \underbrace{\lambda \|\theta\|_2^2}_{>0}$$

$$\hat{\theta}_n = (X^T X + \lambda I)^{-1} X^T \vec{y}$$

↑
This term is extra

$$(X^T X + \lambda I)^{-1} = U \begin{bmatrix} (\sigma_1^2 + \lambda)^{-1} & & 0 \\ & \ddots & \\ 0 & & (\sigma_d^2 + \lambda)^{-1} \end{bmatrix} U^T$$

$$E[\hat{\theta}_n] = \left[U \begin{bmatrix} \frac{\sigma_1^2}{\sigma_1^2 + \lambda} & & 0 \\ & \ddots & \\ 0 & & \frac{\sigma_d^2}{\sigma_d^2 + \lambda} \end{bmatrix} U^T \right] \theta^*$$

$$\text{cov}(\hat{\theta}_n) = V \begin{bmatrix} \frac{\tau^2 \sigma_1^2}{(\sigma_1^2 + \lambda)^2} & & \\ & \ddots & \\ & & \frac{\tau^2 \sigma_d^2}{(\sigma_d^2 + \lambda)^2} \end{bmatrix} V^T$$

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

$$y = \theta^T x + \varepsilon$$

$$\lambda \uparrow \Rightarrow \text{var} \downarrow, \text{Bias} \uparrow$$

$$\lambda \downarrow \Rightarrow \text{var} \uparrow, \text{Bias} \downarrow$$

$$\text{MSE}[\hat{f}_n] = \underbrace{\tau^2}_{\text{irreducible error}} + \underbrace{\mathbb{E}[\hat{f}_n(x^*) - f(x^*)]^2}_{\text{Bias}^2} + \underbrace{V[\hat{f}_n(x^*)]}_{\text{variance}}$$

$$f(x) = \theta^{*T} x$$

$$\begin{aligned} \text{Bias}(\hat{f}_n(x^*)) &= \mathbb{E}[\hat{f}_n(x^*) - f(x^*)] \\ &= \mathbb{E}[\hat{\theta}_n^T x^* - \theta^{*T} x^*] \\ &= \mathbb{E}[(\hat{\theta}_n - \theta^*)^T] x^* \end{aligned}$$

$$= \text{Bias}(\hat{\theta}_n)^T \boxed{x^*}$$

$$\text{Var}(\hat{f}_n) = x^T \text{COV}(\hat{\theta}_n) x$$

Heuristics for Bias & Variance

Training Error \approx Bias

Cross Validation Error \approx Variance

Training Error

To fight Bias

- * Make model larger
- * Add more features
- * Reduce regularization
- * More complex model

To fight Variance

- * Collect more data
- * Increase regularization
- * Simpler model.

