

# Lec 06: Recap: MAP estimate

$$w_{\text{MAP}}^* \in \arg \min_w \left\{ \frac{1}{2\sigma^2} \|Xw - y\|^2 + \frac{\lambda}{2} \|w\|^2 \right\}$$

$$= \frac{1}{\sigma^2} \left( \frac{1}{2\sigma^2} X^T X + \frac{\lambda}{2} I \right)^{-1} X^T y$$

$A$  is PD if  $\forall x \in \mathbb{R}^d \setminus \{0\}$

$$x^T A x > 0$$

$$v^T A v > 0$$

$$= \frac{1}{2\sigma^2} v^T X^T X v + \frac{\lambda}{2} \|v\|^2$$

$$= \frac{1}{2\sigma^2} \|Xv\|^2 + \frac{\lambda}{2} \|v\|^2 > 0$$

# Bias & Variance

Goal: Estimate  $\hat{y} \rightarrow$  not seen in our training example

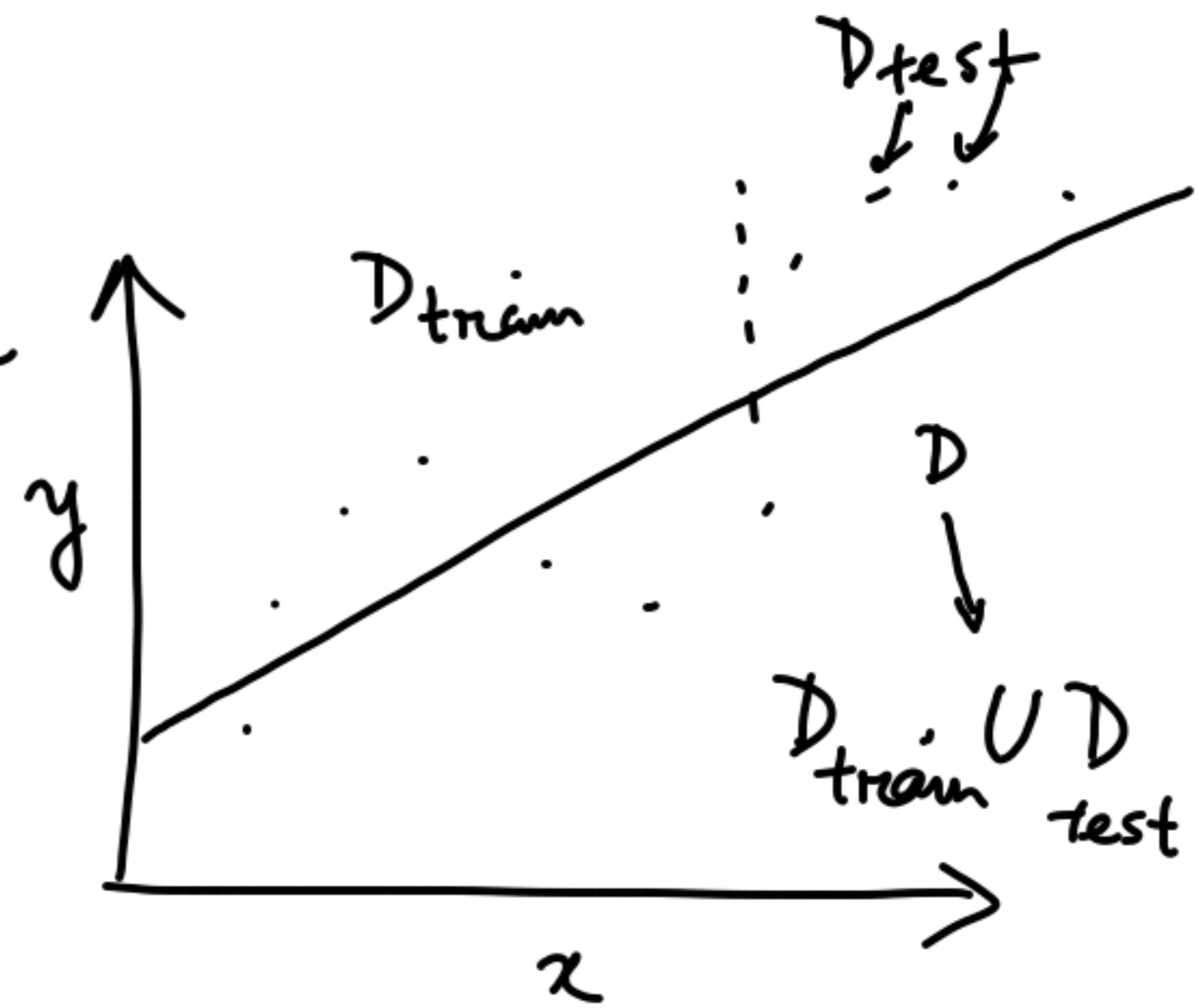
$(\hat{x}, \hat{y}) \rightarrow$  test data point

$(x_i, y_i) \in D_{\text{train}} \rightarrow$  training data points.

$\downarrow$   
 $f_D(\hat{x}) \rightarrow \hat{y}$

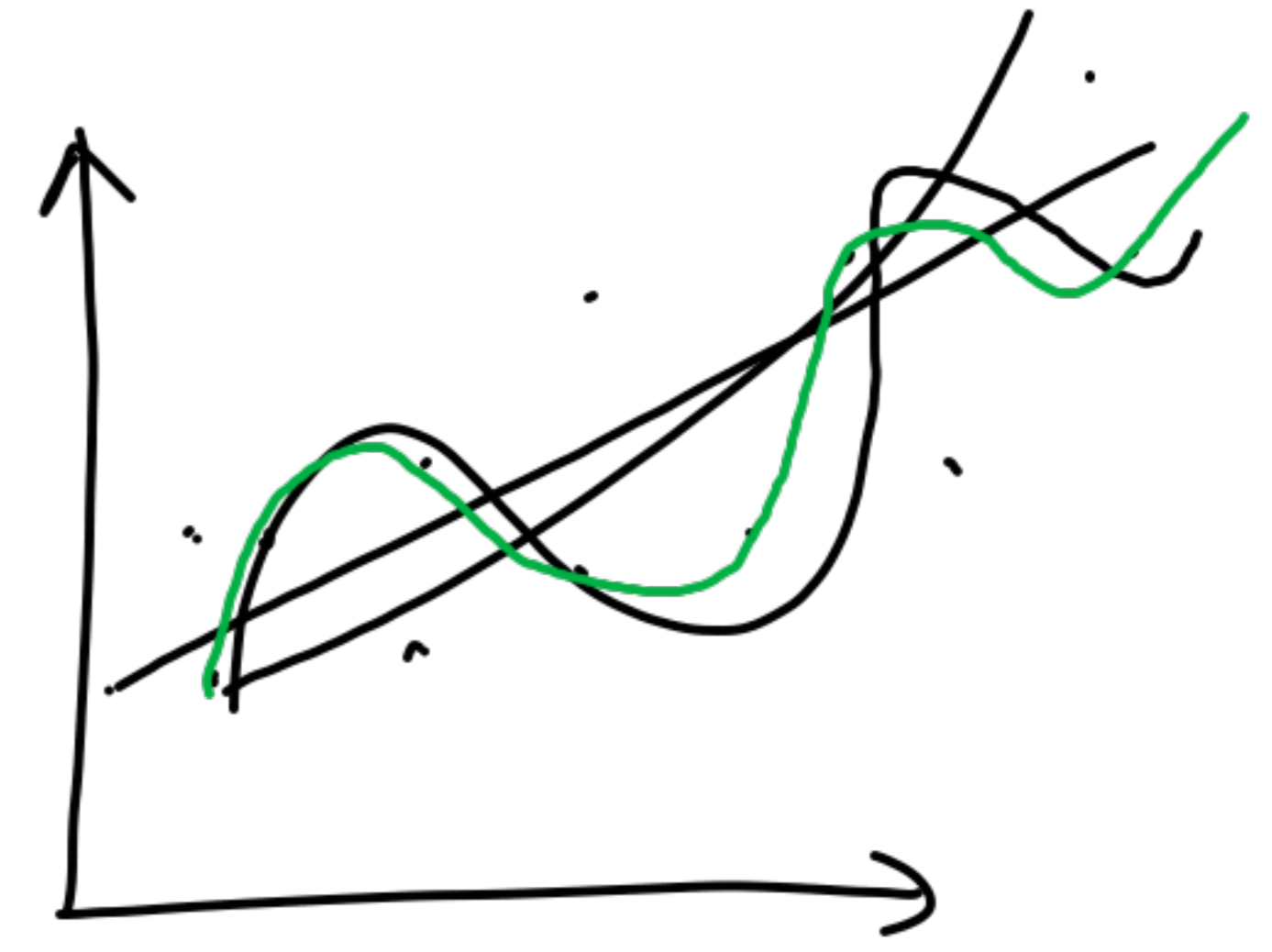
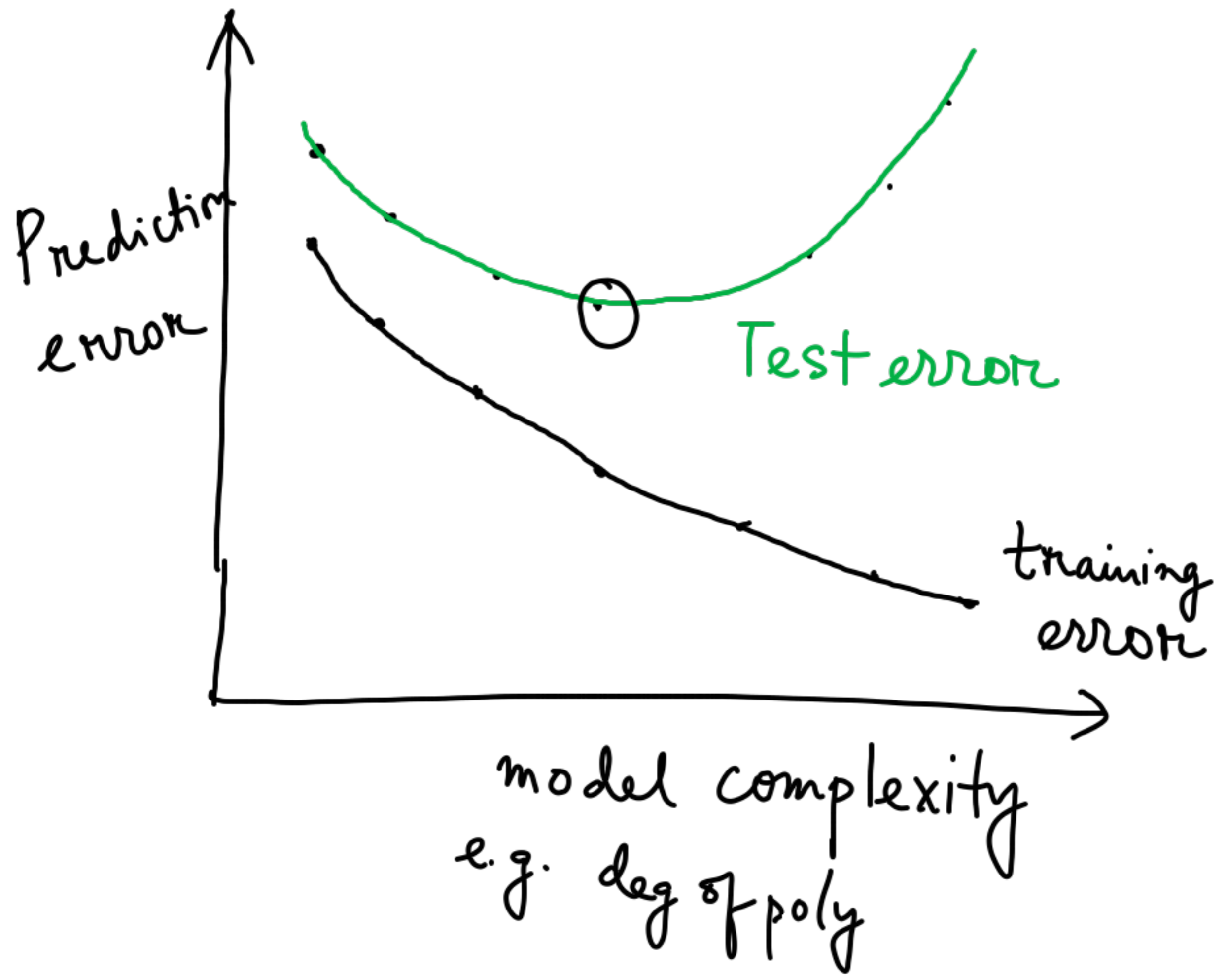
Measure of goodness: ① training error

$$\sum_{i \in D} l(\underbrace{f_D(x_i)}_{\text{prediction}}, \underbrace{y_i}_{\text{actual}})$$



Measure 2: Test error → hold out set  $D_{\text{test}}$

$$\sum_{j \in D_{\text{test}}} l(f_D(x_j), y_j)$$



Test error

- ① Bias
  - ② Variance
  - ③ Noise
- } model dependent

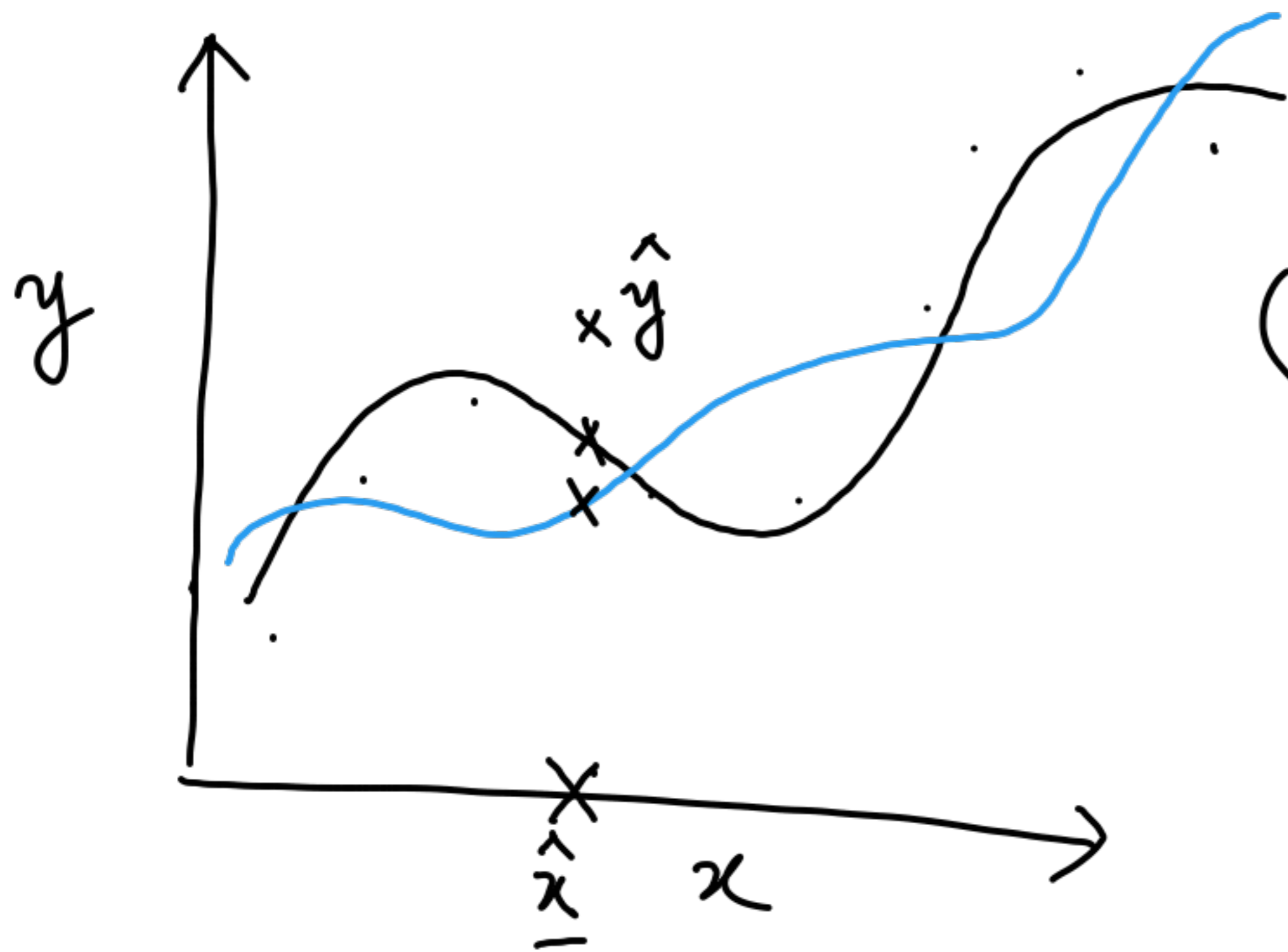
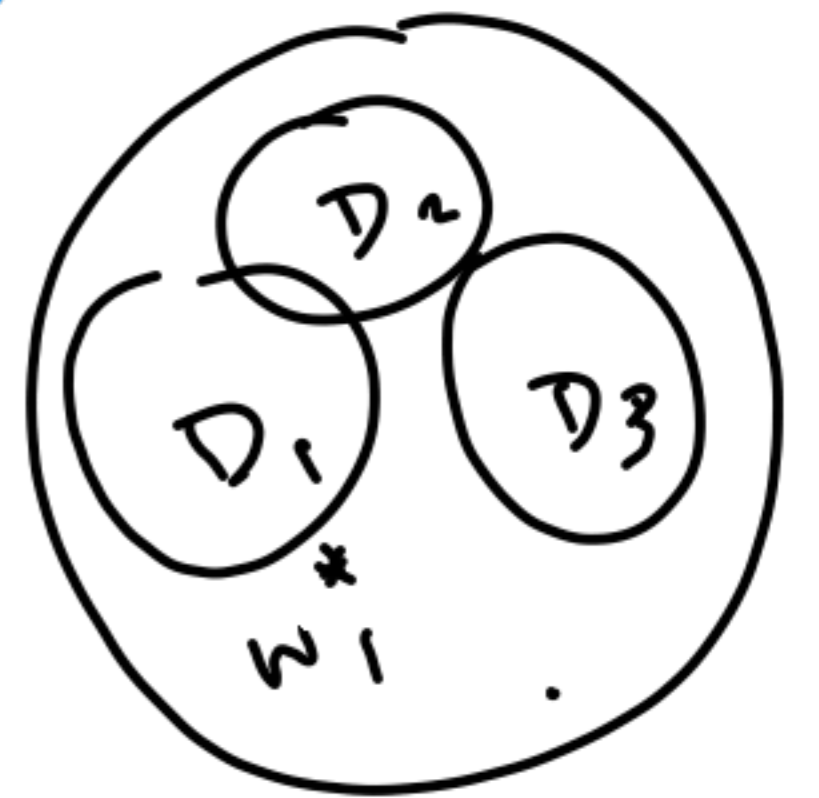
Q: What are Bias and Variance?

$$\frac{1}{3} (w_1^{*T} \hat{x} + w_2^{*T} \hat{x} + \dots)$$

$$\tilde{D} \rightarrow g_D(x)$$

$$(x_i, y_i) \in D$$

$$g_D(x) = \underline{w^T x}$$



$$y = \underline{f(x)} + \epsilon \sim N(0, \sigma^2)$$

Test error:  $\underline{err} = g_D(\hat{x}) - \hat{y}$

Test error

$$\text{err} = g_D(\hat{\lambda}) - \hat{y}$$

$\hat{\lambda}$  fixed

D distribution  
is independent of  
 $\hat{y}$  dist.

$\mathbb{E}$

$$= \left( g_D(\hat{\lambda}) - \mathbb{E}_D \left( g_D(\hat{\lambda}) \right) \right)$$

$$+ \left( \mathbb{E}_D \left( g_D(\hat{\lambda}) \right) - \mathbb{E}(\hat{y}) \right) \rightarrow B = \text{Bias}$$

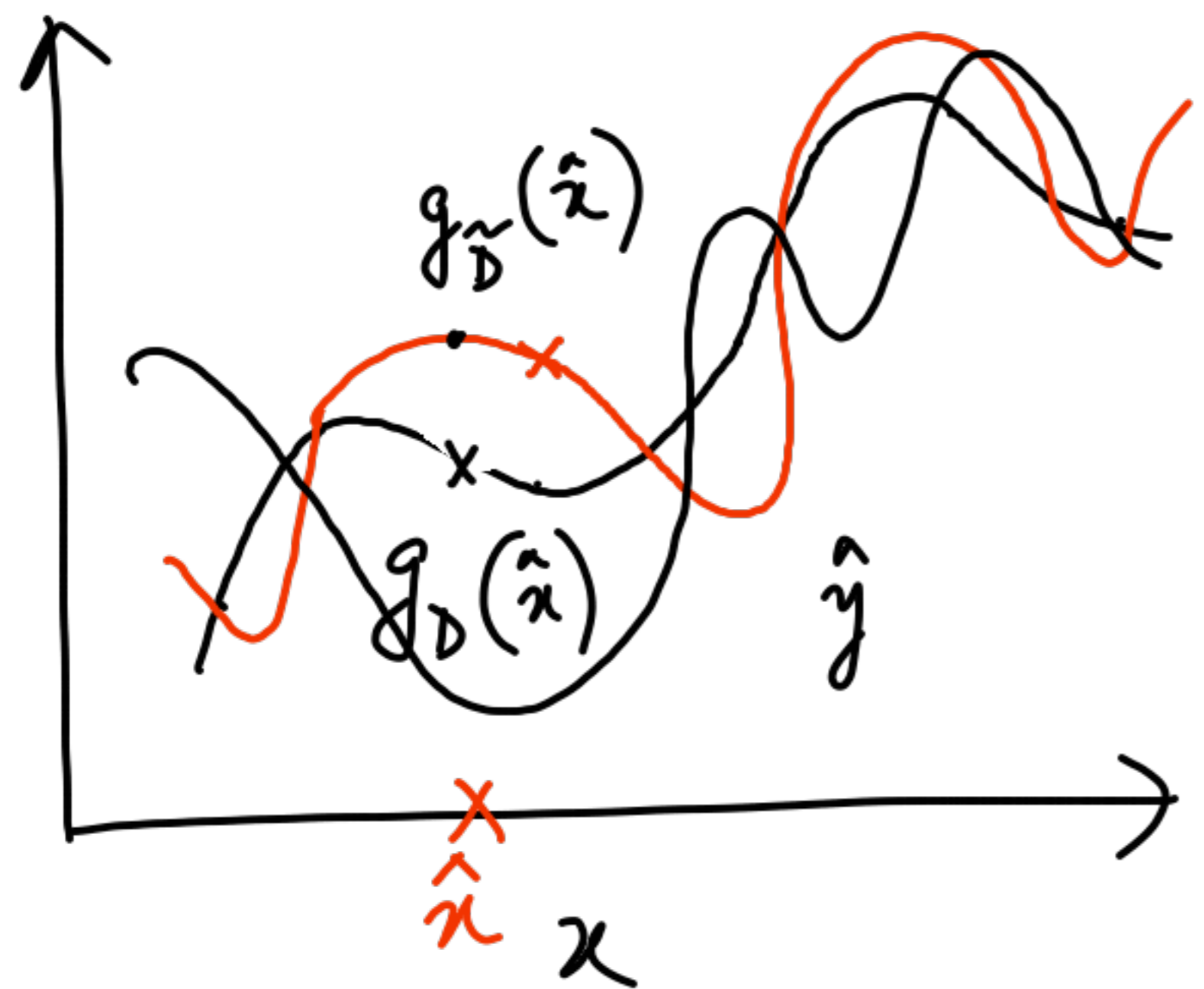
$$+ \left( \mathbb{E}(\hat{y}) - \hat{y} \right) \rightarrow C \text{ noise}$$

$$\mathbb{E}(\text{err}^2) = \underbrace{\mathbb{E}A^2}_{\text{variance}} + \underbrace{\mathbb{E}B^2}_{\text{Bias}} + \mathbb{E}C^2 + 2 \left[ \mathbb{E}(AB) + \mathbb{E}(BC) + \mathbb{E}(CA) \right]$$

$$\mathbb{E}_D \left[ \left( g_D(\hat{x}) - \mathbb{E}_D(g_D(\hat{x})) \right)^2 \right] = \text{Variance of the model}$$

$$\mathbb{E}_D(g_D(\hat{x})) - \mathbb{E}(\hat{y}) = \text{Bias of the model}$$

$$\mathbb{E} \left[ \left( \hat{y} - \mathbb{E} \hat{y} \right)^2 \right] \approx \text{Noise}$$

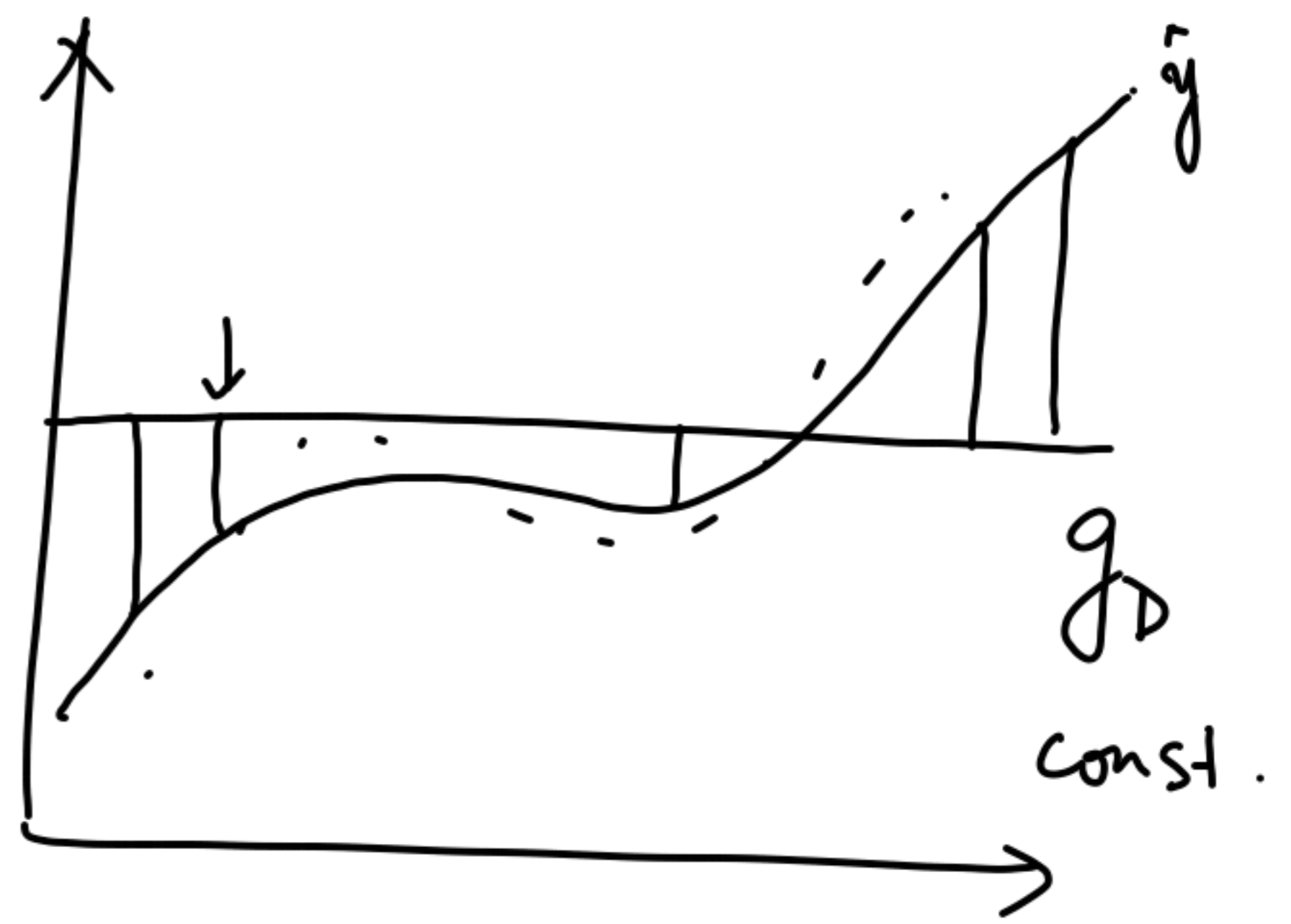


$g_D$   
 $g_D \rightarrow$   
 polynomial  
 of deg 10

Variance high, Bias is low  
 Overfitting

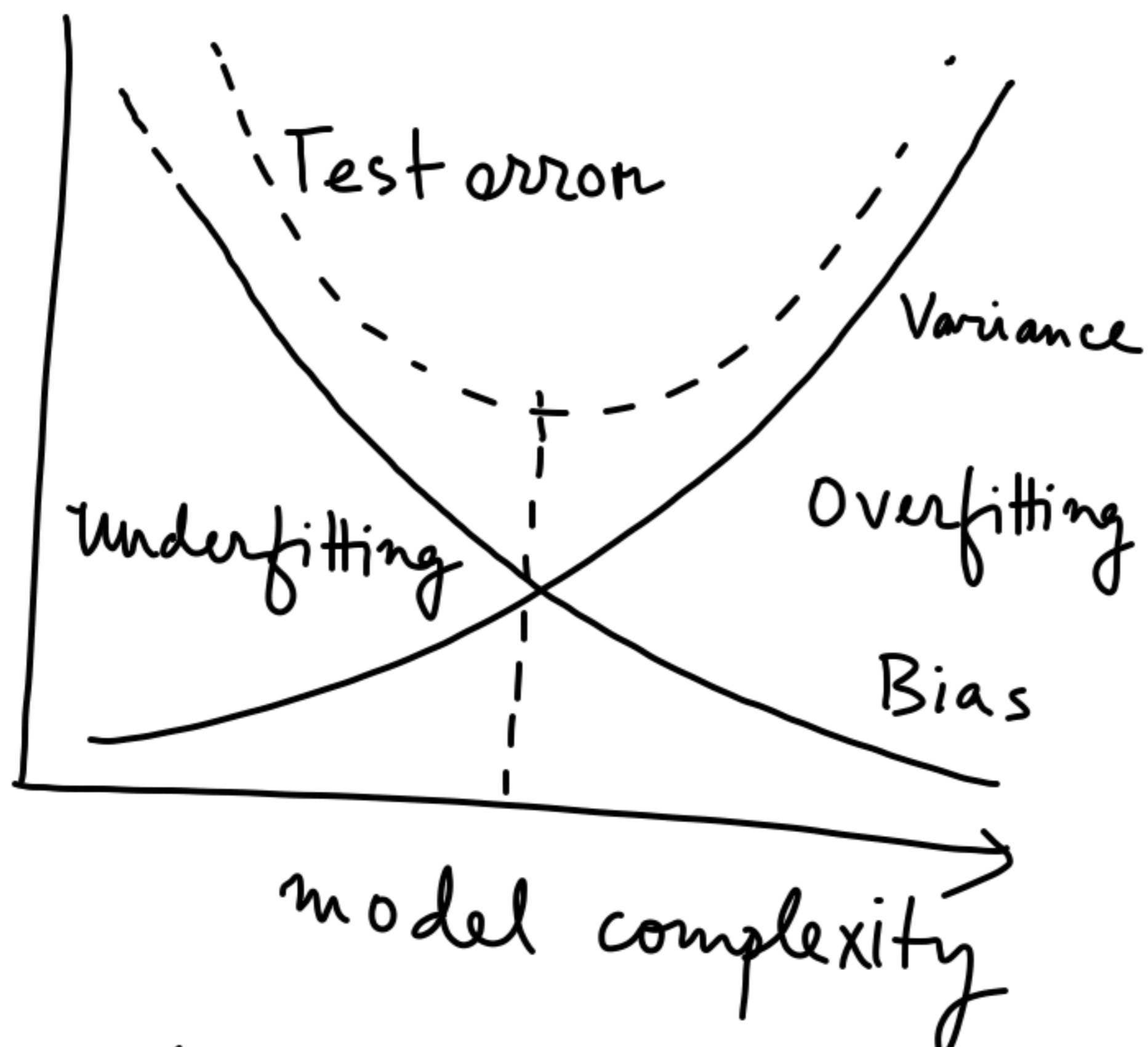
$$y = f(x) + \epsilon$$

$$\mathbb{E} y = f(x)$$



$g_D(x) = w_0$   
 $\text{Var} = \mathbb{E} \left( g_D(\hat{x}) - \mathbb{E}_D (g_D(\hat{x})) \right)^2$   
 Bias is high, variance low  
 Underfitting.





$$E_{\hat{x}, \hat{y}} \ell(g_D(\hat{x}), \hat{y})$$

Regularization for linear regression

$$w_{MLE} \in \operatorname{argmin} \frac{1}{2\sigma^2} \|Xw - y\|^2$$

$$w_{MAP} \in \operatorname{argmin} \left\{ \frac{1}{2\sigma^2} \|Xw - y\|^2 + \frac{\lambda}{2} \|w\|^2 \right\}$$

hyperparameter  $\lambda$

$$\text{Loss}(w) = \frac{1}{n} \text{Loss}_D(w) + \lambda \text{Reg}(w)$$

(Regularized model)

①  $\text{Reg}(w) = \|w\|_2^2 \rightarrow L_2 \text{ norm}$   $\sqrt{\sum_{i=1}^d w_i^2}$

②  $\text{Reg}(w) = \|w\|_1 \rightarrow L_1 \text{ norm}$   $\sum_{i=1}^d |w_i|$

Ridge regression:

$$w^* \in \underset{w}{\text{argmin}} \left\{ \left\| \begin{matrix} \Phi \\ \vdots \\ \Phi_0(x_n) \dots \Phi_m(x_n) \end{matrix} w - y \right\|^2 + \lambda \|w\|_2^2 \right\}$$

LASSO

$$w^* \in \underset{w}{\text{argmin}} \left\{ \left\| \Phi w - y \right\|^2 + \lambda \|w\|_1 \right\}$$

least absolute shrinkage and selection operator  
(LASSO)

equivalent optimization problem

LASSO

$$\min (\Phi^T w - y)^T (\Phi w - y)$$

s.t.

$$\|w\|_1 \leq c_1$$

← const. dependent on  $\lambda$

Ridge

objective is same

$$\text{s.t. } \|w\|_2 \leq c_2$$



