

Recall from last lecture:

$$\begin{aligned} \omega_{MAP}^* &\in \underset{\omega}{\operatorname{argmin}} \left\{ \frac{1}{2\sigma^2} \|x\omega - y\|^2 + \frac{\lambda}{2} \|\omega\|^2 \right\} \\ &= \frac{1}{\sigma^2} \left(\frac{1}{2\sigma^2} X^T X + \frac{\lambda}{2} I \right)^{-1} X^T y \end{aligned}$$

The matrix $\frac{1}{2\sigma^2} X^T X + \frac{\lambda}{2} I$ is invertible

A matrix $A_{d \times d}$ is positive definite if, $\forall x \in \mathbb{R}^d \setminus \{0\}$, $x^T A x > 0$
Any positive definite matrix is invertible

$$\begin{aligned} \text{Let } A &= \frac{1}{2\sigma^2} X^T X + \frac{\lambda}{2} I, \quad v^T A v = \frac{1}{2\sigma^2} v^T X^T X v + \frac{\lambda}{2} \|v\|^2 \\ &= \frac{1}{2\sigma^2} \|Xv\|^2 + \frac{\lambda}{2} \|v\|^2 > 0 \end{aligned}$$

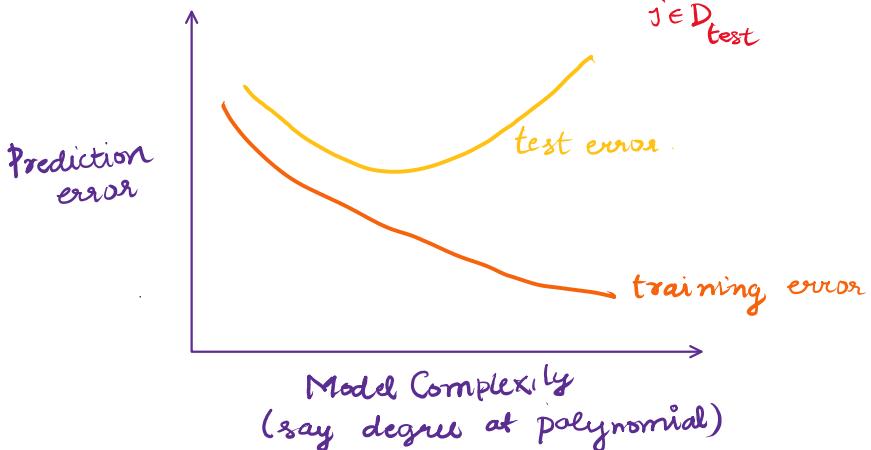
Prediction Errors:

Goal: Estimate \hat{y} , which is not present in training example.
we call (\hat{x}, \hat{y}) as test data point & $(x_i, y_i) \in D_{\text{train}}$ (training data points)
we need our prediction, $g_0(\hat{x})$ to be very close \hat{y} .

Measure of goodness:

i) training error: $\sum_{i \in D} l(g_0(x_i) - y_i)$

ii) Test error: We hold out a set, D_{test} $\sum_{j \in D_{\text{test}}} l(g_0(x_j) - y_j)$



The three primary reasons/causes for test error:

- i) Bias
- ii) Variance
- iii) Noise (persists always / can't get rid of)

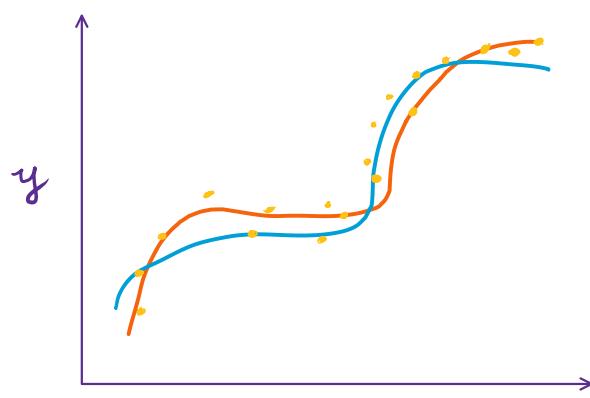
Model dependent

Bias and Variance

Clearly, the predictor function g changes as we change the underlying distribution in which we train it.

Let the actual f^D be

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



$$(x_i, y_i) \in D \Rightarrow g_D(x) = \omega^\top x$$

$$(x_i, y_i) \in \tilde{D} \Rightarrow g_{\tilde{D}}(x)$$

We have the test error, $\text{err} = g_D(\hat{x}) - \hat{y}$, where (\hat{x}, \hat{y}) is a test data point. Assume \hat{x} to be fixed.

$$\begin{aligned} \text{err} &= g_D(\hat{x}) - \hat{y} \\ &= (g_D(\hat{x}) - E_D(g_D(\hat{x}))) \xrightarrow{A} \\ &\quad + (E_D(g_D(\hat{x})) - E(\hat{y})) \xrightarrow{B=\text{bias}} \\ &\quad + (E(\hat{y}) - \hat{y}) \xrightarrow{C} \end{aligned}$$

NB: 'D' distribution is independent of the \hat{y} distribution.

$$E(\text{err}^2) = E[A^2] + E[B^2] + E[C^2] + 2(E(AB) + E(BC) + E(CA))$$

↓ ↓ ↓

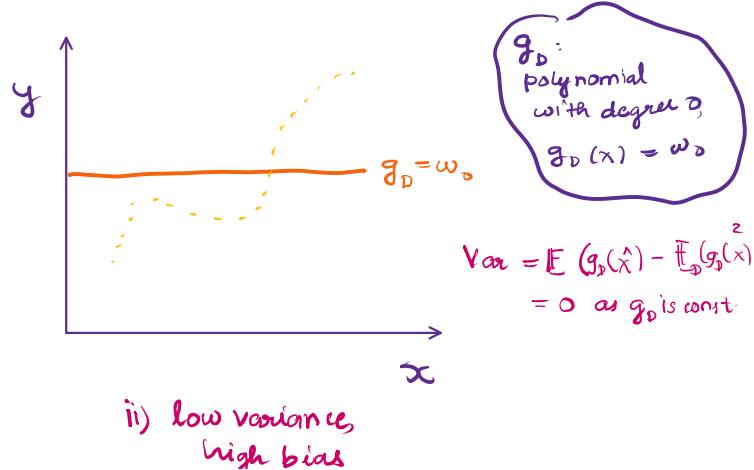
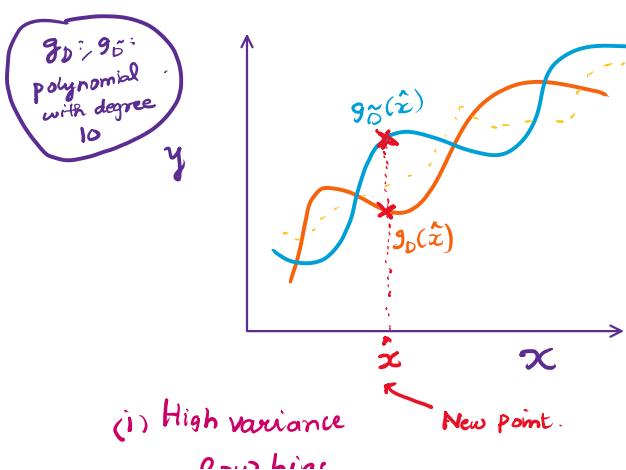
Variance Bias² Noise

Variance of the model: $E_D[(g_D(\hat{x}) - E(g_D(\hat{x})))]^2$

NB: $E(A) = E(C) = 0$
 $E(B)$ is a constant.
 $E(BC) = E(C) E(B) = 0$
 due to independence

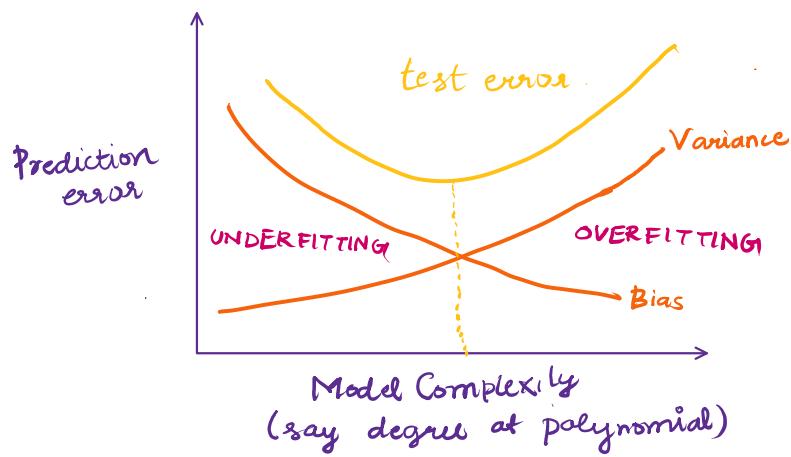
Bias of the model: $E_D(g_D(\hat{x})) - E(\hat{y})$

Noise of the model: $E[(\hat{y} - E(\hat{y}))]^2$



(i) High variance New point.
low bias
(OVERFITTING)

ii) low variance,
high bias
(UNDERFITTING)



Regularization for linear regression

Remember our ultimate procedure is optimization:

$$\begin{aligned} w_{MLE} &\in \arg\min \frac{1}{2\sigma^2} \|Xw - y\|^2 \\ w_{MAP} &\in \arg\min \left\{ \frac{1}{2\sigma^2} \|Xw - y\|^2 + \frac{\lambda}{2} \|w\|^2 \right\} \end{aligned}$$

$$\text{We've: } \text{Loss}(w) = \text{Loss}_D(w) + \lambda \text{Reg}(w)$$

(Regularized model) Regularizer fn
hyperparameter.

Hyperparameter is usually chosen by trial and error, from a grid of possible values.

$$\text{Eq: } \text{loss}(\omega) = E_{\hat{x}, \hat{y}} (l(g_0(\hat{x}), \hat{y}))$$

We've following possibilities at Regularizer for

$$i) \text{ Reg}(w) = \|w\|_2^2 \quad (\text{L}_2 \text{ norm}) = \sqrt{\sum_{i=1}^n w_i^2}$$

$$ii) \text{ Reg}(w) = \|w\|, (\text{L}_1 \text{ norm}) = \sum_{i=1}^d |w_i|$$

$$\text{Ridge regression: } \omega^* \in \underset{\omega}{\operatorname{argmin}} \left\{ \| \Phi \omega - y \|^2 + \lambda \| \omega \|_2^2 \right\}$$

LASSO regression: $\omega^* \in \arg \min_{\omega} \left\{ \|\Phi \omega - y\|^2 + \lambda \|\omega\|_1 \right\}$

↳ least absolute shrinkage and selection operator

$$\Phi = \begin{bmatrix} \phi_0(x_1) & \dots & \phi_m(x_1) \\ \vdots & & \vdots \\ \phi_0(x_n) & \dots & \phi_m(x_n) \end{bmatrix}$$

Equivalent optimization problems:

$$\text{LASSO: } \min_{\omega} (\Phi \omega - y)^T (\Phi \omega - y) \text{ st. } \|\omega\|_1 \leq c_1, \leftarrow \text{const depends on } \lambda$$

$$\text{Ridge: } \min_{\omega} (\Phi \omega - y)^T (\Phi \omega - y) \text{ st. } \|\omega\|_2 \leq c_2$$