



# ISTA 421 + INFO 521

## Introduction to Machine Learning

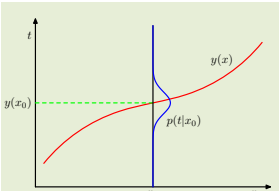
### Lecture 10: Maximum Likelihood Uncertainty 2

**Clay Morrison**  
claytonm@email.arizona.edu  
Harvill 437A  
Phone 621-6609

25 September 2017

 1

## Review

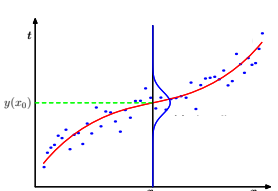


The generating process...

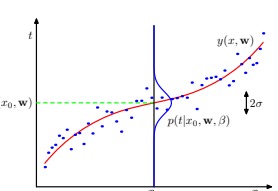
$$\hat{t}_n = \mathbf{w}^\top \mathbf{x}_n + \epsilon_n; \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

$$L = p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2) = \prod_{n=1}^N p(t_n|\mathbf{x}_n, \mathbf{w}, \sigma^2)$$

$$= \prod_{n=1}^N \mathcal{N}(\mathbf{w}^\top \mathbf{x}_n, \sigma^2)$$



... generates data ...



... that we fit a model to

$p(\hat{\mathbf{t}}|\mathbf{X}, \hat{\mathbf{w}}, \hat{\sigma}^2) = \prod_{n=1}^N p(\hat{t}_n|\mathbf{x}_n, \hat{\mathbf{w}}, \hat{\sigma}^2)$   
 $= \prod_{n=1}^N \mathcal{N}(\hat{\mathbf{w}}^\top \mathbf{x}_n, \hat{\sigma}^2)$

prediction

↑

estimated parameters

↑

**Maximum Likelihood Estimates of Params**

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t}$$

$$\hat{\sigma}^2 = \frac{1}{N} (\mathbf{t}^\top \mathbf{t} - \mathbf{t}^\top \mathbf{X} \hat{\mathbf{w}})$$

Predictions:

**The MLE is unique**


$$\frac{\partial^2 \log L}{\partial \mathbf{w} \partial \mathbf{w}^\top} = -\frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X}$$

**Analysis of Uncertainty in Param Estimates via Expectation**

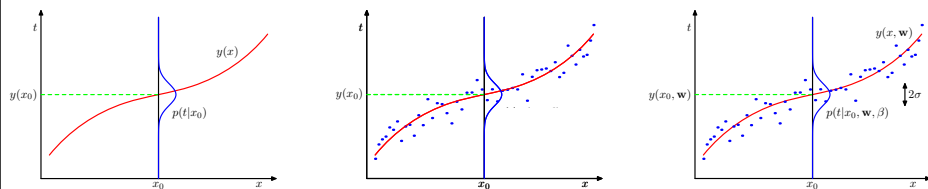
$$\mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{\hat{\mathbf{w}}\} = \mathbf{w}$$

$$\text{cov}\{\hat{\mathbf{w}}\} = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} = -\left(\frac{\partial^2 \log L}{\partial \mathbf{w} \partial \mathbf{w}^\top}\right)^{-1}$$

**The Fisher Information**  $\mathcal{I} = \frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X}$

 2

## Back to: The Generative Picture



The generating process...

$$t_n = \mathbf{w}^\top \mathbf{x}_n + \epsilon_n ; \epsilon \sim \mathcal{N}(0, \sigma^2)$$

... generates data ...

... that we fit a model to

How about modeling the uncertainty in our predictions?

Want model of region of values in which our prediction might fall

estimated parameters

$$p(\hat{\mathbf{t}} | \mathbf{X}, \hat{\mathbf{w}}, \hat{\sigma}^2) = \prod_{n=1}^N p(\hat{t}_n | \mathbf{x}_n, \hat{\mathbf{w}}, \hat{\sigma}^2)$$

$$= \prod_{n=1}^N \mathcal{N}(\hat{\mathbf{w}}^\top \mathbf{x}_n, \hat{\sigma}^2)$$

predictive distribution



3

## Variability in Predictions

- We are predicting 2 values:

$$t_{new} , \sigma_{new}^2$$



4

## Variability in Predictions

- We are predicting 2 values:

$$t_{new}, \sigma_{new}^2$$

$$t_{new} = \hat{\mathbf{w}}^\top \mathbf{x}_{new} \quad \text{Same solution as minimizing mean squared loss}$$



5

## Variability in Predictions

- We are predicting 2 values:

$$t_{new}, \sigma_{new}^2$$

$$t_{new} = \hat{\mathbf{w}}^\top \mathbf{x}_{new} \quad \text{Same solution as minimizing mean squared loss}$$

$$\mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{t_{new}\}$$



6

## Variability in Predictions

- We are predicting 2 values:

$$t_{new}, \sigma_{new}^2$$

$$t_{new} = \hat{\mathbf{w}}^\top \mathbf{x}_{new} \quad \text{Same solution as minimizing mean squared loss}$$

$$\mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{t_{new}\} = \mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{\hat{\mathbf{w}}\}^\top \mathbf{x}_{new}$$



## Variability in Predictions

- We are predicting 2 values:

$$t_{new}, \sigma_{new}^2$$

$$t_{new} = \hat{\mathbf{w}}^\top \mathbf{x}_{new} \quad \text{Same solution as minimizing mean squared loss}$$

$$\begin{aligned} \mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{t_{new}\} &= \mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{\hat{\mathbf{w}}\}^\top \mathbf{x}_{new} \\ &= \mathbf{w}^\top \mathbf{x}_{new} \end{aligned}$$



## Variability in Predictions

- We are predicting 2 values:

$$t_{new}, \sigma_{new}^2$$

$$t_{new} = \hat{\mathbf{w}}^\top \mathbf{x}_{new} \quad \text{Same solution as minimizing mean squared loss}$$

$$\begin{aligned} \mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{t_{new}\} &= \mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{\hat{\mathbf{w}}\}^\top \mathbf{x}_{new} \\ &= \mathbf{w}^\top \mathbf{x}_{new} \end{aligned}$$

The **expected value** of our prediction is the new data attribute multiplied by the **true w**



9

## Predicting the Variance of $t_{new}$

$$\sigma_{new}^2 = \text{var} \{t_{new}\} = \mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{t_{new}^2\} - (\mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{t_{new}\})^2$$



10

## Predicting the Variance of $t_{new}$

$$\sigma_{new}^2 = \text{var} \{t_{new}\} = \mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{t_{new}^2\} - (\mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{t_{new}\})^2$$

Substitute  $t_{new} = \hat{\mathbf{w}}^\top \mathbf{x}_{new}$



## Predicting the Variance of $t_{new}$

$$\sigma_{new}^2 = \text{var} \{t_{new}\} = \mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{t_{new}^2\} - (\mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{t_{new}\})^2$$

Substitute  $t_{new} = \hat{\mathbf{w}}^\top \mathbf{x}_{new}$

$$\begin{aligned} \text{var}\{t_{new}\} &= \mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \left\{ (\hat{\mathbf{w}}^\top \mathbf{x}_{new})^2 \right\} - (\mathbf{w}^\top \mathbf{x}_{new})^2 \\ &= \mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \left\{ \mathbf{x}_{new}^\top \hat{\mathbf{w}} \hat{\mathbf{w}}^\top \mathbf{x}_{new} \right\} - \mathbf{x}_{new}^\top \mathbf{w} \mathbf{w}^\top \mathbf{x}_{new}. \end{aligned}$$



## Predicting the Variance of $t_{new}$

$$\sigma_{new}^2 = \text{var} \{t_{new}\} = \mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{t_{new}^2\} - (\mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{t_{new}\})^2$$

Substitute  $t_{new} = \hat{\mathbf{w}}^\top \mathbf{x}_{new}$

$$\begin{aligned} \text{var}\{t_{new}\} &= \mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \left\{ (\hat{\mathbf{w}}^\top \mathbf{x}_{new})^2 \right\} - (\mathbf{w}^\top \mathbf{x}_{new})^2 \\ &= \mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \left\{ \mathbf{x}_{new}^\top \hat{\mathbf{w}} \hat{\mathbf{w}}^\top \mathbf{x}_{new} \right\} - \mathbf{x}_{new}^\top \mathbf{w} \mathbf{w}^\top \mathbf{x}_{new}. \end{aligned}$$

Substitute  $\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t}$



13

## Predicting the Variance of $t_{new}$

$$\sigma_{new}^2 = \text{var} \{t_{new}\} = \mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{t_{new}^2\} - (\mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{t_{new}\})^2$$

Substitute  $t_{new} = \hat{\mathbf{w}}^\top \mathbf{x}_{new}$

$$\begin{aligned} \text{var}\{t_{new}\} &= \mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \left\{ (\hat{\mathbf{w}}^\top \mathbf{x}_{new})^2 \right\} - (\mathbf{w}^\top \mathbf{x}_{new})^2 \\ &= \mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \left\{ \mathbf{x}_{new}^\top \hat{\mathbf{w}} \hat{\mathbf{w}}^\top \mathbf{x}_{new} \right\} - \mathbf{x}_{new}^\top \mathbf{w} \mathbf{w}^\top \mathbf{x}_{new}. \end{aligned}$$

Substitute  $\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t}$

$$\text{var}\{t_{new}\} = \mathbf{x}_{new}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{ \mathbf{t} \mathbf{t}^\top \} \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{new} - \mathbf{x}_{new}^\top \mathbf{w} \mathbf{w}^\top \mathbf{x}_{new}$$



14

## Predicting the Variance of $t_{new}$

$$\sigma_{new}^2 = \text{var}\{t_{new}\} = \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)}\{t_{new}^2\} - (\mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)}\{t_{new}\})^2$$

Substitute  $t_{new} = \hat{\mathbf{w}}^\top \mathbf{x}_{new}$

$$\begin{aligned}\text{var}\{t_{new}\} &= \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)}\left\{\left(\hat{\mathbf{w}}^\top \mathbf{x}_{new}\right)^2\right\} - \left(\mathbf{w}^\top \mathbf{x}_{new}\right)^2 \\ &= \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)}\left\{\mathbf{x}_{new}^\top \hat{\mathbf{w}} \hat{\mathbf{w}}^\top \mathbf{x}_{new}\right\} - \mathbf{x}_{new}^\top \mathbf{w} \mathbf{w}^\top \mathbf{x}_{new}.\end{aligned}$$

Substitute  $\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t}$

$$\text{var}\{t_{new}\} = \mathbf{x}_{new}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)}\{\mathbf{t} \mathbf{t}^\top\} \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{new} - \mathbf{x}_{new}^\top \mathbf{w} \mathbf{w}^\top \mathbf{x}_{new}$$

On slide 23 of Lec 9, in the derivation of the covariance of  $\hat{\mathbf{w}}$ , we identified  $\mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)}\{\mathbf{t} \mathbf{t}^\top\} = \mathbf{X} \mathbf{w} \mathbf{w}^\top \mathbf{X}^\top + \sigma^2 \mathbf{I}$



15

## Predicting the Variance of $t_{new}$

$$\sigma_{new}^2 = \text{var}\{t_{new}\} = \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)}\{t_{new}^2\} - (\mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)}\{t_{new}\})^2$$

Substitute  $t_{new} = \hat{\mathbf{w}}^\top \mathbf{x}_{new}$

$$\begin{aligned}\text{var}\{t_{new}\} &= \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)}\left\{\left(\hat{\mathbf{w}}^\top \mathbf{x}_{new}\right)^2\right\} - \left(\mathbf{w}^\top \mathbf{x}_{new}\right)^2 \\ &= \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)}\left\{\mathbf{x}_{new}^\top \hat{\mathbf{w}} \hat{\mathbf{w}}^\top \mathbf{x}_{new}\right\} - \mathbf{x}_{new}^\top \mathbf{w} \mathbf{w}^\top \mathbf{x}_{new}.\end{aligned}$$

Substitute  $\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t}$

$$\text{var}\{t_{new}\} = \mathbf{x}_{new}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)}\{\mathbf{t} \mathbf{t}^\top\} \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{new} - \mathbf{x}_{new}^\top \mathbf{w} \mathbf{w}^\top \mathbf{x}_{new}$$

On slide 23 of Lec 9, in the derivation of the covariance of  $\hat{\mathbf{w}}$ , we identified  $\mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)}\{\mathbf{t} \mathbf{t}^\top\} = \mathbf{X} \mathbf{w} \mathbf{w}^\top \mathbf{X}^\top + \sigma^2 \mathbf{I}$

$$\begin{aligned}\text{var}\{t_{new}\} &= \mathbf{x}_{new}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\sigma^2 \mathbf{I} + \mathbf{X} \mathbf{w} \mathbf{w}^\top \mathbf{X}^\top) \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{new} - \mathbf{x}_{new}^\top \mathbf{w} \mathbf{w}^\top \mathbf{x}_{new} \\ &= \sigma^2 \mathbf{x}_{new}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{new} + \mathbf{x}_{new}^\top \mathbf{w} \mathbf{w}^\top \mathbf{x}_{new} - \mathbf{x}_{new}^\top \mathbf{w} \mathbf{w}^\top \mathbf{x}_{new} \\ &= \sigma^2 \mathbf{x}_{new}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{new}.\end{aligned}$$



16



## Predicting the Variance of $t_{new}$

$$\sigma_{new}^2 = \text{var}\{t_{new}\} = \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)}\{t_{new}^2\} - (\mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)}\{t_{new}\})^2$$

Substitute  $t_{new} = \hat{\mathbf{w}}^\top \mathbf{x}_{new}$

$$\begin{aligned}\text{var}\{t_{new}\} &= \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)}\left\{(\hat{\mathbf{w}}^\top \mathbf{x}_{new})^2\right\} - (\mathbf{w}^\top \mathbf{x}_{new})^2 \\ &= \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)}\left\{\mathbf{x}_{new}^\top \hat{\mathbf{w}} \hat{\mathbf{w}}^\top \mathbf{x}_{new}\right\} - \mathbf{x}_{new}^\top \mathbf{w} \mathbf{w}^\top \mathbf{x}_{new}.\end{aligned}$$

Substitute  $\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t}$

$$\text{var}\{t_{new}\} = \mathbf{x}_{new}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)}\{\mathbf{t} \mathbf{t}^\top\} \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{new} - \mathbf{x}_{new}^\top \mathbf{w} \mathbf{w}^\top \mathbf{x}_{new}$$

On slide 23 of Lec 9, in the derivation of the covariance of  $\hat{\mathbf{w}}$ , we identified  $\mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)}\{\mathbf{t} \mathbf{t}^\top\} = \mathbf{X} \mathbf{w} \mathbf{w}^\top \mathbf{X}^\top + \sigma^2 \mathbf{I}$

$$\begin{aligned}\text{var}\{t_{new}\} &= \mathbf{x}_{new}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\sigma^2 \mathbf{I} + \mathbf{X} \mathbf{w} \mathbf{w}^\top \mathbf{X}^\top) \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{new} - \mathbf{x}_{new}^\top \mathbf{w} \mathbf{w}^\top \mathbf{x}_{new} \\ &= \sigma^2 \mathbf{x}_{new}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{new} + \mathbf{x}_{new}^\top \mathbf{w} \mathbf{w}^\top \mathbf{x}_{new} - \mathbf{x}_{new}^\top \mathbf{w} \mathbf{w}^\top \mathbf{x}_{new} \\ &= \sigma^2 \mathbf{x}_{new}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{new}.\end{aligned}$$

Recall:  $\text{cov}\{\hat{\mathbf{w}}\} = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} = -\left(\frac{\partial^2 \log L}{\partial \mathbf{w} \partial \mathbf{w}^\top}\right)^{-1}$



17

## Predicting the Variance of $t_{new}$

$$\sigma_{new}^2 = \text{var}\{t_{new}\} = \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)}\{t_{new}^2\} - (\mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)}\{t_{new}\})^2$$

Substitute  $t_{new} = \hat{\mathbf{w}}^\top \mathbf{x}_{new}$

$$\begin{aligned}\text{var}\{t_{new}\} &= \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)}\left\{(\hat{\mathbf{w}}^\top \mathbf{x}_{new})^2\right\} - (\mathbf{w}^\top \mathbf{x}_{new})^2 \\ &= \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)}\left\{\mathbf{x}_{new}^\top \hat{\mathbf{w}} \hat{\mathbf{w}}^\top \mathbf{x}_{new}\right\} - \mathbf{x}_{new}^\top \mathbf{w} \mathbf{w}^\top \mathbf{x}_{new}.\end{aligned}$$

Substitute  $\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t}$

$$\text{var}\{t_{new}\} = \mathbf{x}_{new}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)}\{\mathbf{t} \mathbf{t}^\top\} \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{new} - \mathbf{x}_{new}^\top \mathbf{w} \mathbf{w}^\top \mathbf{x}_{new}$$

On slide 23 of Lec 9, in the derivation of the covariance of  $\hat{\mathbf{w}}$ , we identified  $\mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)}\{\mathbf{t} \mathbf{t}^\top\} = \mathbf{X} \mathbf{w} \mathbf{w}^\top \mathbf{X}^\top + \sigma^2 \mathbf{I}$

$$\begin{aligned}\text{var}\{t_{new}\} &= \mathbf{x}_{new}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\sigma^2 \mathbf{I} + \mathbf{X} \mathbf{w} \mathbf{w}^\top \mathbf{X}^\top) \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{new} - \mathbf{x}_{new}^\top \mathbf{w} \mathbf{w}^\top \mathbf{x}_{new} \\ &= \sigma^2 \mathbf{x}_{new}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{new} + \mathbf{x}_{new}^\top \mathbf{w} \mathbf{w}^\top \mathbf{x}_{new} - \mathbf{x}_{new}^\top \mathbf{w} \mathbf{w}^\top \mathbf{x}_{new} \\ &= \sigma^2 \mathbf{x}_{new}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{new}.\end{aligned}$$

Recall:  $\text{cov}\{\hat{\mathbf{w}}\} = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} = -\left(\frac{\partial^2 \log L}{\partial \mathbf{w} \partial \mathbf{w}^\top}\right)^{-1}$

So, could be written

$$\sigma_{new}^2 = \mathbf{x}_{new}^\top \text{COV}\{\hat{\mathbf{w}}\} \mathbf{x}_{new}$$



18

## In Summary

$$t_{\text{new}} = \mathbf{x}_{\text{new}}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t} = \mathbf{x}_{\text{new}}^\top \hat{\mathbf{w}}$$

$$\sigma_{\text{new}}^2 = \sigma^2 \mathbf{x}_{\text{new}}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{\text{new}}$$

We estimate this from the data:  $\hat{\sigma}^2$

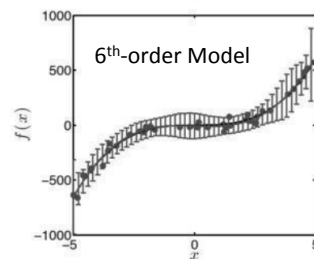
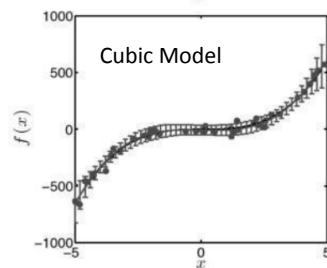
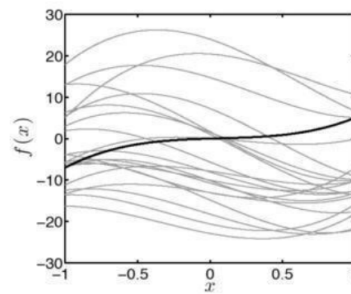
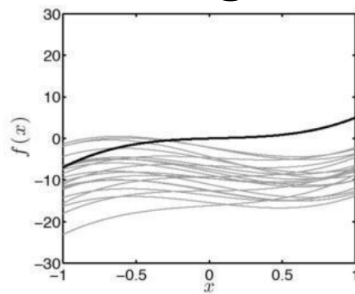
$$\hat{\sigma}^2 = \frac{1}{N} (\mathbf{t}^\top \mathbf{t} - \mathbf{t}^\top \mathbf{X} \hat{\mathbf{w}})$$



19

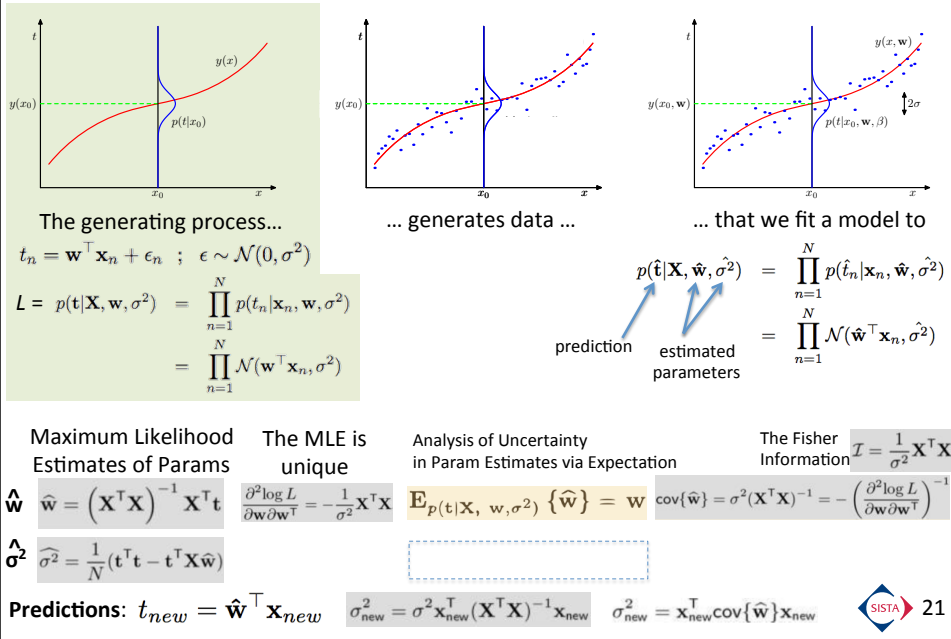
## Plotting Predictive Error Bars

$$t_{\text{new}} \pm \sigma_{\text{new}}^2$$



20

## Review



## Quantifying the Uncertainty in our Estimate of $\hat{\sigma}^2$

$$\mathbf{t} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$\mathbb{E}_{p(\mathbf{t})} \{\mathbf{t}^\top \mathbf{A} \mathbf{t}\} = \text{Tr}(\mathbf{A} \boldsymbol{\Sigma}) + \boldsymbol{\mu}^\top \mathbf{A} \boldsymbol{\mu}$$

$$\mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{\hat{\sigma}^2\} = \frac{1}{N} \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{\mathbf{t}^\top \mathbf{t}\} - \frac{1}{N} \mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{\mathbf{t}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t}\}$$

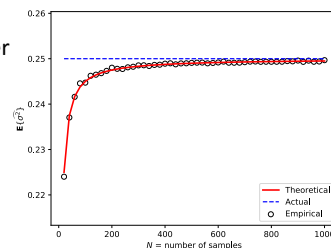
$$\mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{\hat{\sigma}^2\} = \sigma^2 \left(1 - \frac{D}{N}\right)$$

When  $D < N$  (that is, the number of attributes we measure for each data point is *smaller* than the number of data points), then our estimates of the variance will, on average, be lower than the true variance.

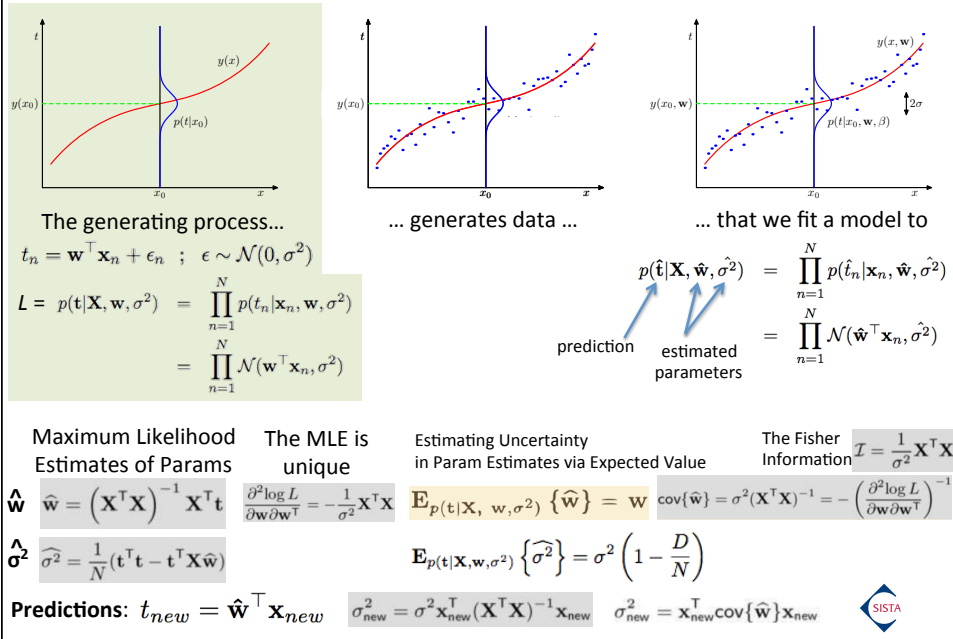
Unlike the estimate for  $\hat{\mathbf{w}}$ , the MLE for  $\hat{\sigma}^2$  is **biased**.

$$D = 2 \text{ and } N = 20$$

$$\mathbb{E}_{p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2)} \{\hat{\sigma}^2\} = \sigma^2 \left(1 - \frac{D}{N}\right) = 0.25 \left(1 - \frac{2}{20}\right) = 0.2250$$



## Review

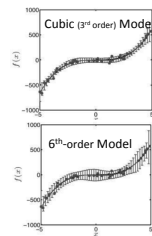


## MLE and Model Selection

$$\log L = -\frac{1}{N} \log 2\pi - N \log \sigma - \frac{1}{2\sigma^2} \sum_{n=1}^N (t_n - \mathbf{w}^\top \mathbf{x}_n)^2$$

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t}$$

$$\hat{\sigma}^2 = \frac{1}{N} \sum_{n=1}^N (t_n - \mathbf{x}_n^\top \hat{\mathbf{w}})^2 = \frac{1}{N} (\mathbf{t}^\top \mathbf{t} - \mathbf{t}^\top \mathbf{X} \hat{\mathbf{w}})$$



## MLE and Model Selection

$$\log L = -\frac{1}{N} \log 2\pi - N \log \sigma - \frac{1}{2\sigma^2} \sum_{n=1}^N (t_n - \mathbf{w}^\top \mathbf{x}_n)^2$$

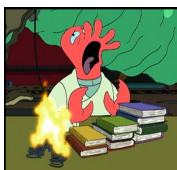
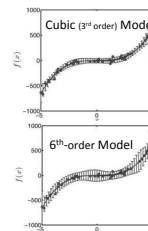
$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t}$$

$$\hat{\sigma}^2 = \frac{1}{N} \sum_{n=1}^N (t_n - \mathbf{x}_n^\top \hat{\mathbf{w}})^2 = \frac{1}{N} (\mathbf{t}^\top \mathbf{t} - \mathbf{t}^\top \mathbf{X} \hat{\mathbf{w}})$$

Plug in  $\hat{\sigma}^2$  to the log likelihood:

$$\begin{aligned} \log L &= -\frac{N}{2} \log 2\pi - \frac{N}{2} \log \hat{\sigma}^2 - \frac{1}{2\hat{\sigma}^2} N \hat{\sigma}^2 \\ &= -\frac{N}{2} (1 + \log 2\pi) - \frac{N}{2} \log \hat{\sigma}^2. \end{aligned}$$

Making  $\hat{\sigma}^2$  smaller makes  $\log L$  larger. Making model more flexible decreases  $\hat{\sigma}^2$ .



## MLE Prefers Complex Models

$$\log L = -\frac{1}{N} \log 2\pi - N \log \sigma - \frac{1}{2\sigma^2} \sum_{n=1}^N (t_n - \mathbf{w}^\top \mathbf{x}_n)^2$$

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t}$$

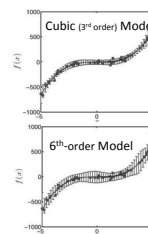
$$\hat{\sigma}^2 = \frac{1}{N} \sum_{n=1}^N (t_n - \mathbf{x}_n^\top \hat{\mathbf{w}})^2 = \frac{1}{N} (\mathbf{t}^\top \mathbf{t} - \mathbf{t}^\top \mathbf{X} \hat{\mathbf{w}})$$

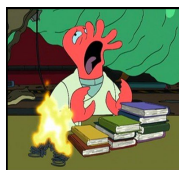
Plug in  $\hat{\sigma}^2$  to the log likelihood:

$$\begin{aligned} \log L &= -\frac{N}{2} \log 2\pi - \frac{N}{2} \log \hat{\sigma}^2 - \frac{1}{2\hat{\sigma}^2} N \hat{\sigma}^2 \\ &= -\frac{N}{2} (1 + \log 2\pi) - \frac{N}{2} \log \hat{\sigma}^2. \end{aligned}$$

Making  $\hat{\sigma}^2$  smaller makes  $\log L$  larger. Making model more flexible decreases  $\hat{\sigma}^2$ .

**Bad news:** Increasing the model complexity will *decrease* the variance!





## MLE Prefers Complex Models

$$\log L = -\frac{1}{N} \log 2\pi - N \log \sigma - \frac{1}{2\sigma^2} \sum_{n=1}^N (t_n - \mathbf{w}^\top \mathbf{x}_n)^2$$

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t}$$

$$\hat{\sigma}^2 = \frac{1}{N} \sum_{n=1}^N (t_n - \mathbf{x}_n^\top \hat{\mathbf{w}})^2 = \frac{1}{N} (\mathbf{t}^\top \mathbf{t} - \mathbf{t}^\top \mathbf{X} \hat{\mathbf{w}})$$

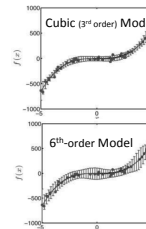
Plug in  $\hat{\sigma}^2$  to the log likelihood:

$$\begin{aligned} \log L &= -\frac{N}{2} \log 2\pi - \frac{N}{2} \log \hat{\sigma}^2 - \frac{1}{2\hat{\sigma}^2} N \hat{\sigma}^2 \\ &= -\frac{N}{2} (1 + \log 2\pi) - \frac{N}{2} \log \hat{\sigma}^2. \end{aligned}$$

Making  $\hat{\sigma}^2$  smaller makes  $\log L$  larger. Making model more flexible decreases  $\hat{\sigma}^2$ .

**Bad news:** Increasing the model complexity will *decrease* the variance!

**Bottom line:** Unfortunately, we can't use MLE to do *model selection*. But, with a *particular* model, MLE will choose the parameters that make the data have the highest overall likelihood under the model.



## The Bias Variance Tradeoff

$$Y = f(X) + \epsilon$$

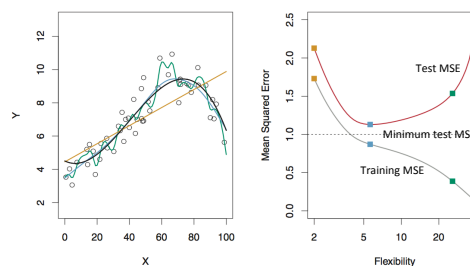


Fig 1: contrasting training with testing error



## The Bias Variance Tradeoff

$$Y = f(X) + \epsilon$$

$$E(Y - \hat{Y})^2 = E[f(X) + \epsilon - \hat{f}(X)]^2$$

$$= \underbrace{[f(X) - \hat{f}(X)]^2}_{\text{Reducible}} + \underbrace{\text{Var}(\epsilon)}_{\text{Irreducible}}$$

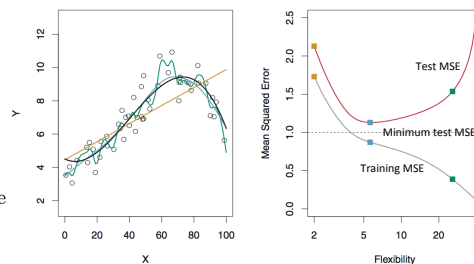


Fig 1: contrasting training with testing error



## The Bias Variance Tradeoff

$$Y = f(X) + \epsilon$$

$$E(Y - \hat{Y})^2 = E[f(X) + \epsilon - \hat{f}(X)]^2$$

$$= \underbrace{[f(X) - \hat{f}(X)]^2}_{\text{Reducible}} + \underbrace{\text{Var}(\epsilon)}_{\text{Irreducible}}$$

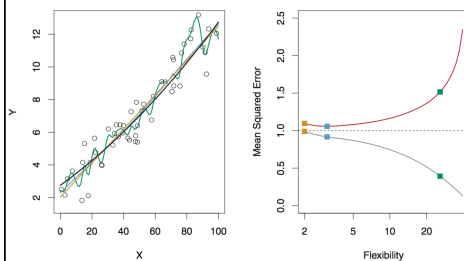
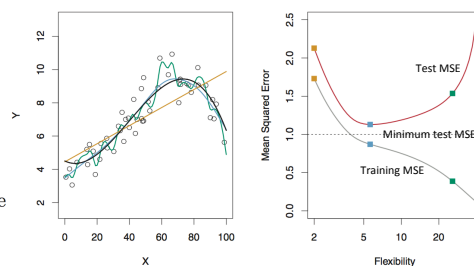


Fig 2: But if data is close to linear, linear fit may do very close to perfect



# The Bias Variance Tradeoff

$$Y = f(X) + \epsilon$$

$$\begin{aligned} E(Y - \hat{Y})^2 &= E[f(X) + \epsilon - \hat{f}(X)]^2 \\ &= \underbrace{[f(X) - \hat{f}(X)]^2}_{\text{Reducible}} + \underbrace{\text{Var}(\epsilon)}_{\text{Irreducible}} \end{aligned}$$

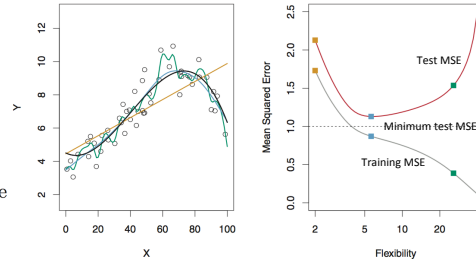
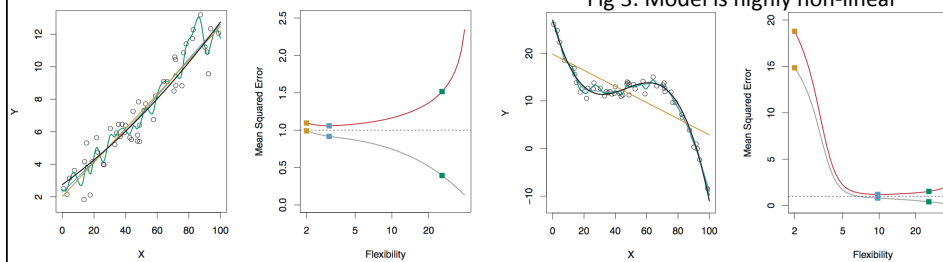


Fig 3: Model is highly non-linear

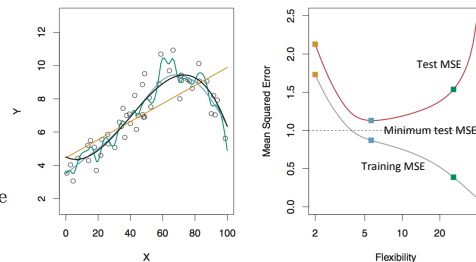


31

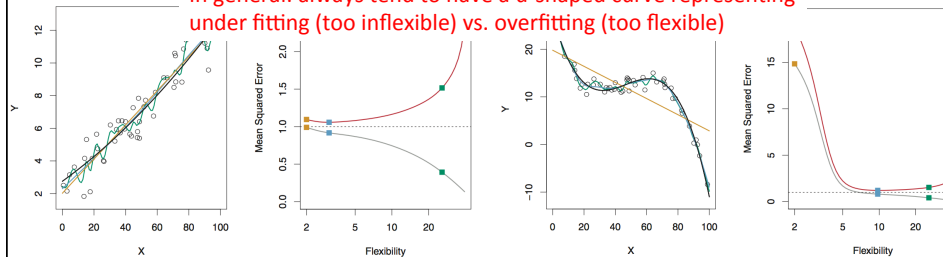
# The Bias Variance Tradeoff

$$Y = f(X) + \epsilon$$

$$\begin{aligned} E(Y - \hat{Y})^2 &= E[f(X) + \epsilon - \hat{f}(X)]^2 \\ &= \underbrace{[f(X) - \hat{f}(X)]^2}_{\text{Reducible}} + \underbrace{\text{Var}(\epsilon)}_{\text{Irreducible}} \end{aligned}$$



In general: always tend to have a u-shaped curve representing under fitting (too inflexible) vs. overfitting (too flexible)



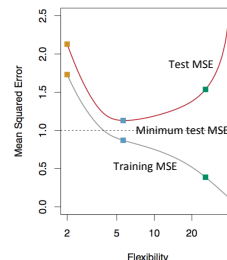
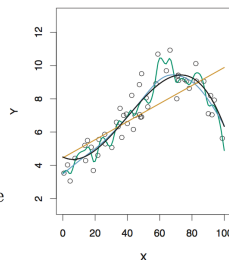
32



# The Bias Variance Tradeoff

$$Y = f(X) + \epsilon$$

$$\begin{aligned} E(Y - \hat{Y})^2 &= E[f(X) + \epsilon - \hat{f}(X)]^2 \\ &= \underbrace{[f(X) - \hat{f}(X)]^2}_{\text{Reducible}} + \underbrace{\text{Var}(\epsilon)}_{\text{Irreducible}} \end{aligned}$$



$$E(y_0 - \hat{f}(x_0))^2 = \text{Var}(\hat{f}(x_0)) + [\text{Bias}(\hat{f}(x_0))]^2 + \text{Var}(\epsilon)$$

Expected Test MSE  
(average test MSE if we repeatedly estimated  $f$  using a large number of training sets)

Model variance  
(how  $f$  changes with different sampled data)

Model Bias  
(error from inflexibility of the model relative to the true  $f$ )

Irreducible variance

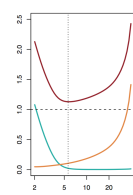
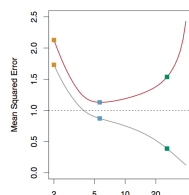
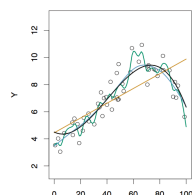
**Bias:** the systematic mismatch between our model and the process that generated the data.

Too simple a model == too high a bias (underfitting)

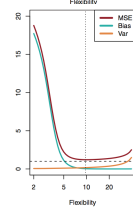
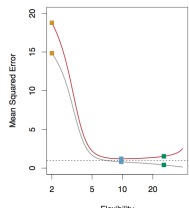
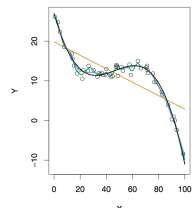
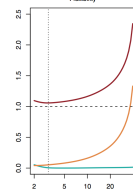
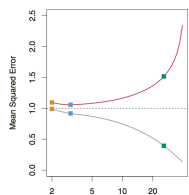
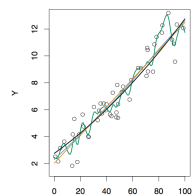
Too complex a model (too many degrees of freedom) == too low a bias (overfitting)

**Variance:** Squared error between model and data

$$E(y_0 - \hat{f}(x_0))^2 = \text{Var}(\hat{f}(x_0)) + [\text{Bias}(\hat{f}(x_0))]^2 + \text{Var}(\epsilon)$$



Blue: Model bias<sup>2</sup>  
Orange: Model variance  
Horizontal dashed:  $\text{Var}(\epsilon)$   
Red: test MSE



**STOP**