

Bias Variance Trade-off

Agenda

- Loss function for regression
- Bias variance trade off
- Debugging variance and bias

Potential Problem – Outliers

How does outlier influence the regression model?

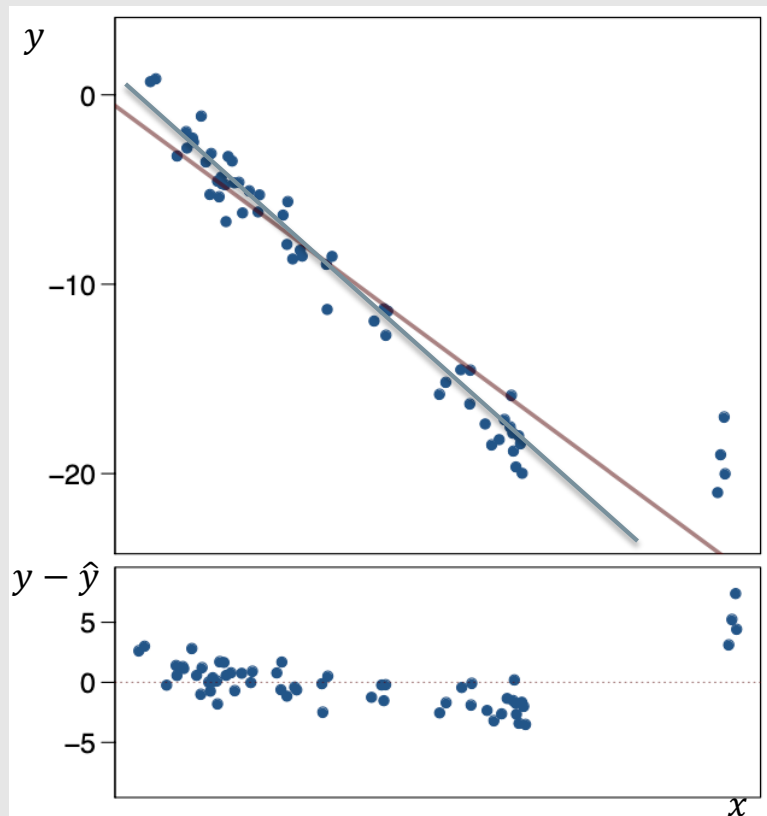
How do we identify outliers?

Residual plot, Data distribution

How do we deal with outliers?

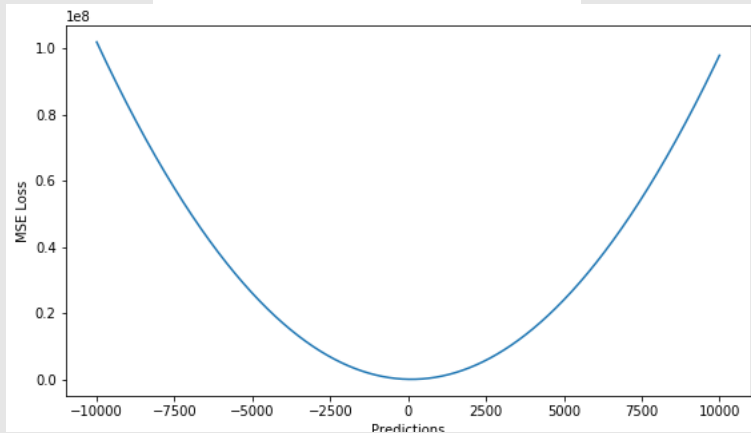
Delete the outlier

More robust models

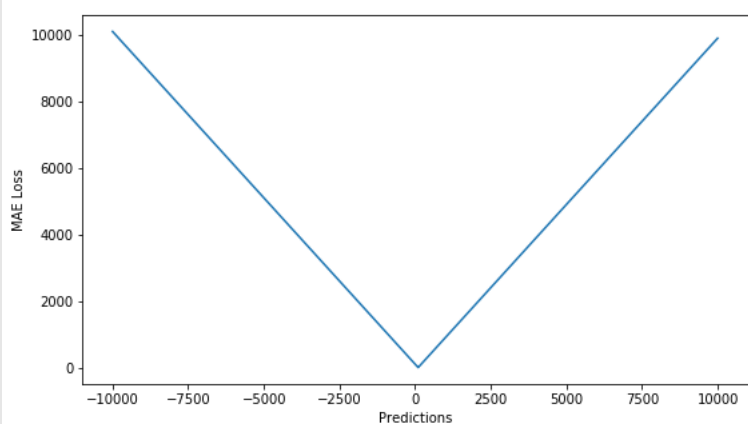


Loss function for regression

$$MSE = \frac{\sum_{i=1}^n (y_i - y_i^p)^2}{n}$$



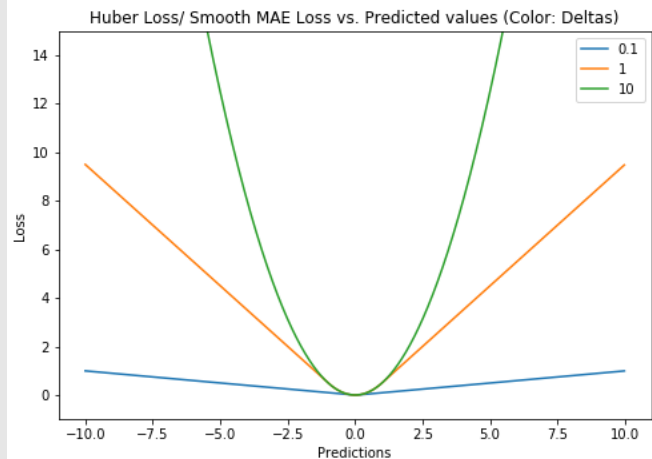
$$MAE = \frac{\sum_{i=1}^n |y_i - y_i^p|}{n}$$



Loss function for regression

Huber Loss (Smooth Mean Absolute Error)

$$L_{\delta}(y, f(x)) = \begin{cases} \frac{1}{2}(y - f(x))^2 & \text{for } |y - f(x)| \leq \delta, \\ \delta |y - f(x)| - \frac{1}{2}\delta^2 & \text{otherwise.} \end{cases}$$



Bias variance trade-off

Linear Regression

Model: $Y = w_0 + w_1X_1 + \dots + w_pX_p + \varepsilon$

Given the training data: $y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$ $X = \begin{pmatrix} \text{---} x_1 \text{---} \\ \vdots \\ \text{---} x_2 \text{---} \end{pmatrix} = \begin{bmatrix} 1 & x_{11} & \dots & x_{1p} \\ 1 & x_{21} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{bmatrix}$

The model can be written as $y = Xw + \varepsilon$

Least square estimation for the parameters:

$$\hat{w} = (X^T X)^{-1} X^T y$$

\hat{w} is a random variable! Even though ground truth w^* is not

Bias and variance in parameter estimation

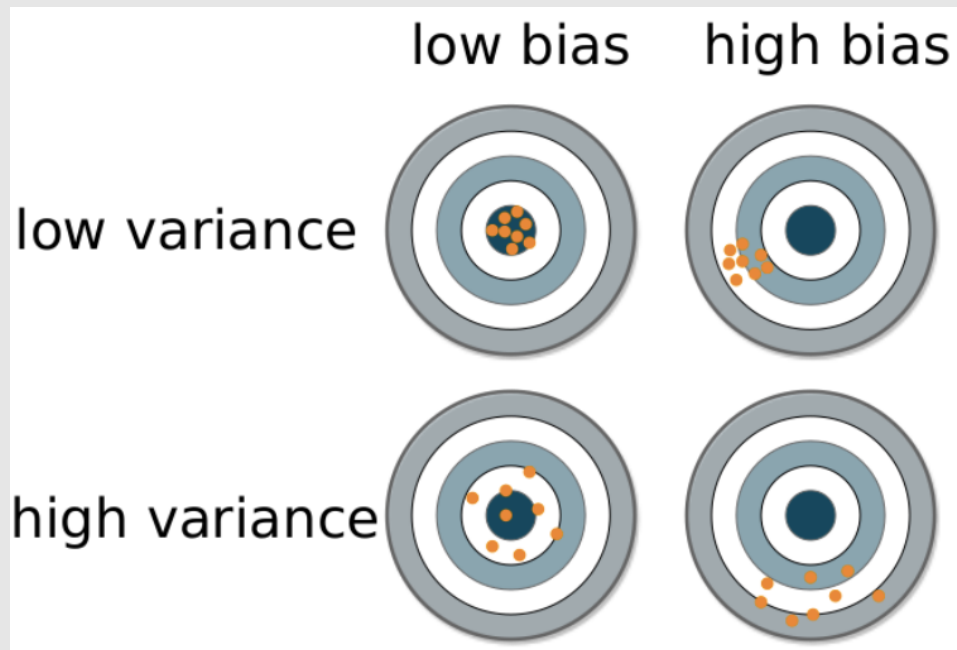
In general given a model with parameter θ , we get an estimator $\hat{\theta}$

- Bias of the estimator: $Bias(\hat{\theta}) = E[\hat{\theta} - \theta^*]$
- Variance of the estimator: $Var(\hat{\theta}) = Cov(\hat{\theta})$

If $Bias(\hat{\theta}) = 0$, then $\hat{\theta}$ is an unbiased estimation for θ

$$\begin{aligned}MSE(\hat{\theta}_n) &= \mathbb{E} \left[\|\hat{\theta}_n - \theta^*\|^2 \right] \\&= \mathbb{E} \left[\|\hat{\theta}_n - \mathbb{E}[\hat{\theta}_n] + \mathbb{E}[\hat{\theta}_n] - \theta^*\|^2 \right] \\&= \mathbb{E} \left[\|\hat{\theta}_n - \mathbb{E}[\hat{\theta}_n]\|^2 + \underbrace{\|\mathbb{E}[\hat{\theta}_n] - \theta^*\|^2}_{\text{Constant}} + 2 \underbrace{(\hat{\theta}_n - \mathbb{E}[\hat{\theta}_n])^T (\mathbb{E}[\hat{\theta}_n] - \theta^*)}_{\text{Zero Mean}} \right] \\&= \mathbb{E} \left[\|\hat{\theta}_n - \mathbb{E}[\hat{\theta}_n]\|^2 \right] + \|\mathbb{E}[\hat{\theta}_n] - \theta^*\|^2 \\&= \mathbb{E} \left[\text{tr} \left[(\hat{\theta}_n - \mathbb{E}[\hat{\theta}_n])(\hat{\theta}_n - \mathbb{E}[\hat{\theta}_n])^T \right] \right] + \|\mathbb{E}[\hat{\theta}_n] - \theta^*\|^2 \\&= \text{tr} \left[\text{Var}(\hat{\theta}_n) \right] + \|\text{Bias}(\hat{\theta}_n)\|^2.\end{aligned}$$

Bias-Variance Decomposition



Bias Variance in Prediction

We have a dataset $D = \{(x^{(i)}, y^{(i)})\}_{i=1}^N$ drawn from a distribution over random variables X and Y , which following the following relation

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

$$E(\varepsilon) = 0$$

$$X \perp \varepsilon$$

Suppose $\hat{f}(x)$ based on the training data D , consider a test point (x_*, y_*) , the squared loss is

$$\begin{aligned} & E_{D, \varepsilon_*} \left[\left(y_* - \hat{f}(x_*) \right)^2 \right] \\ &= E_{D, \varepsilon_*} \left[\left(f(x_*) + \varepsilon_* - \hat{f}(x_*) \right)^2 \right] \\ &= \underbrace{\text{Var}(\varepsilon_*)}_{\text{Irreducible error}} + \underbrace{\left\{ E[\hat{f}(x_*) - f(x_*)] \right\}^2}_{\text{Bias of } \hat{f}} + \underbrace{\text{Var}[\hat{f}(x_*)]}_{\text{Variance of } \hat{f}} \end{aligned}$$

Bias Variance in Prediction

$$E_{D, \varepsilon_*} \left[\left(y_* - \hat{f}(x_*) \right)^2 \right] = \text{Var}(\varepsilon_*) + \{E[\hat{f}(x_*) - f(x_*)]\}^2 + \text{Var}[\hat{f}(x_*)]$$

$$\begin{aligned} & E_{D, \varepsilon_*} \left[y_* - \hat{f}(x_*) \right]^2 \\ &= E_{D, \varepsilon_*} \left[f(x_*) + \varepsilon_* - \hat{f}(x_*) \right]^2 \\ &= E_{D, \varepsilon_*} \left[\underbrace{f(x_*) - \hat{f}(x_*)}_{\text{unrelated with } \varepsilon_*} \right]^2 + E[\varepsilon_*^2] + \underbrace{E[2\varepsilon_* (f(x_*) - \hat{f}(x_*))]}_{\downarrow 0} \\ &= E_D \left[f(x_*) - \hat{f}(x_*) \right]^2 + \text{Var}(\varepsilon_*). \\ &= \left[E(f(x_*) - \hat{f}(x_*)) \right]^2 + \underbrace{\text{Var}(f(x_*) - \hat{f}(x_*))}_{\text{constant}} + \text{Var}(\varepsilon_*). \\ &= \left[E(f(x_*) - \hat{f}(x_*)) \right]^2 + \text{Var}[\hat{f}(x_*)] + \text{Var}(\varepsilon_*). \end{aligned}$$

Bias Variance Trade off

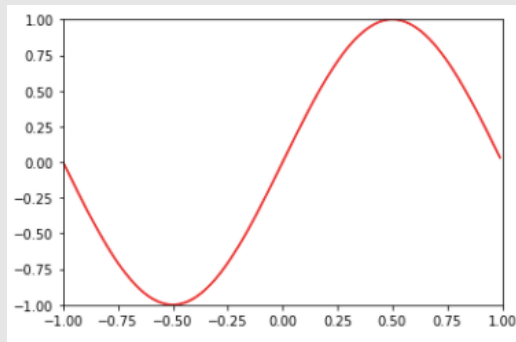
$$E_{D, \varepsilon_*} \left[\left(y_* - \hat{f}(x_*) \right)^2 \right] = \text{Var}(\varepsilon_*) + \left\{ E[\hat{f}(x_*) - f(x_*)] \right\}^2 + \text{Var}[\hat{f}(x_*)]$$

$$E_{X \in D_{test}}(\text{squared loss}) = E_X \left\{ \underbrace{\left[E_D(f(X) - \hat{f}(X)) \right]^2}_{\text{Bias of } \hat{f}} + \underbrace{\text{Var}_D[\hat{f}(X)]}_{\text{Variance of } \hat{f}} + \underbrace{\text{Var}_\varepsilon(\varepsilon)}_{\text{Irreducible error}} \right\}$$

Illustrated Example

Example: Approximate a sine function

True function $f(x) = \sin(\pi x)$, $f: [-1, 1] \rightarrow \mathbb{R}$



You are given two hypotheses of the function to fit:

$$H_0: f(x) = c$$

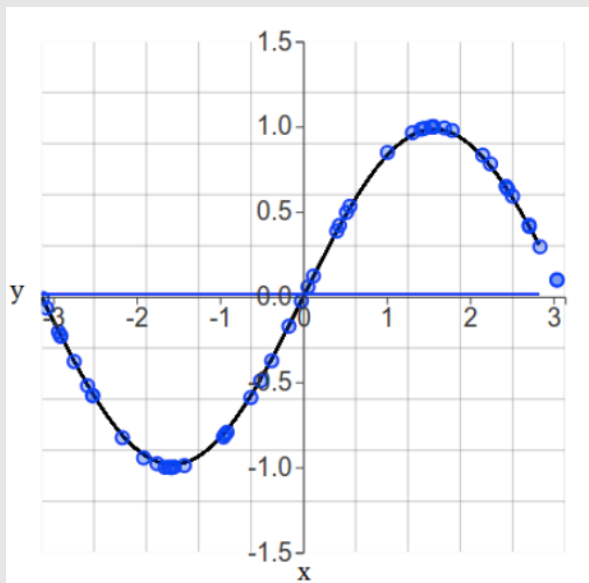
$$H_1: f(x) = w_0 + w_1 x$$

Which leads to better results?

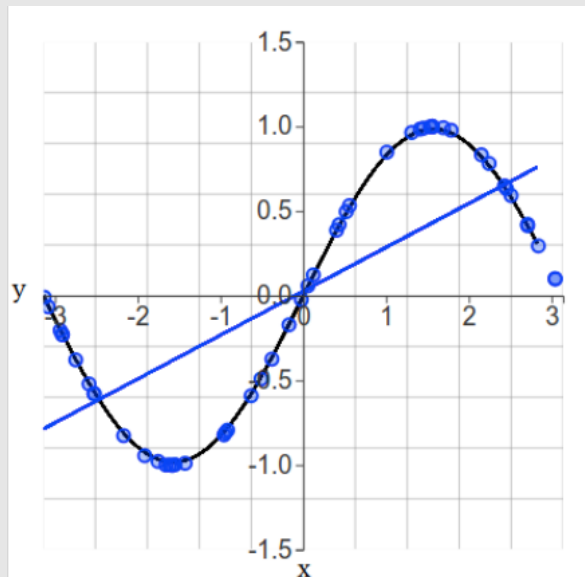
Learning H_0 vs H_0

Data simulation with $M=50$

$$H_0: f(x) = c$$

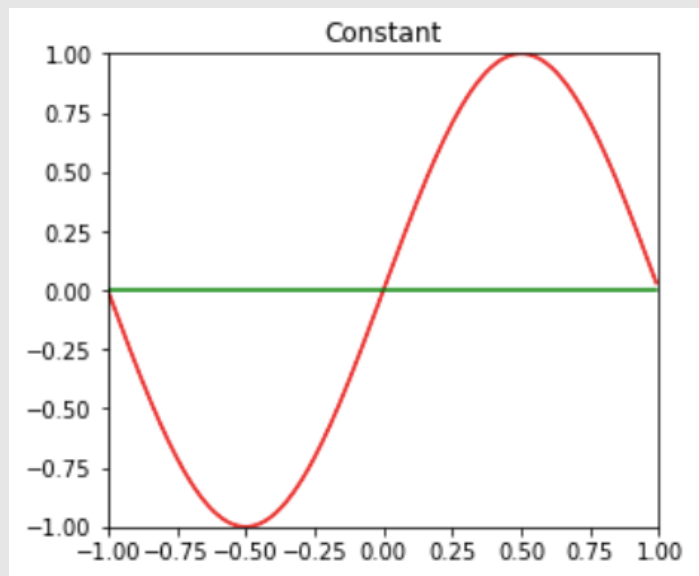


$$H_1: f(x) = w_0 + w_1 x$$

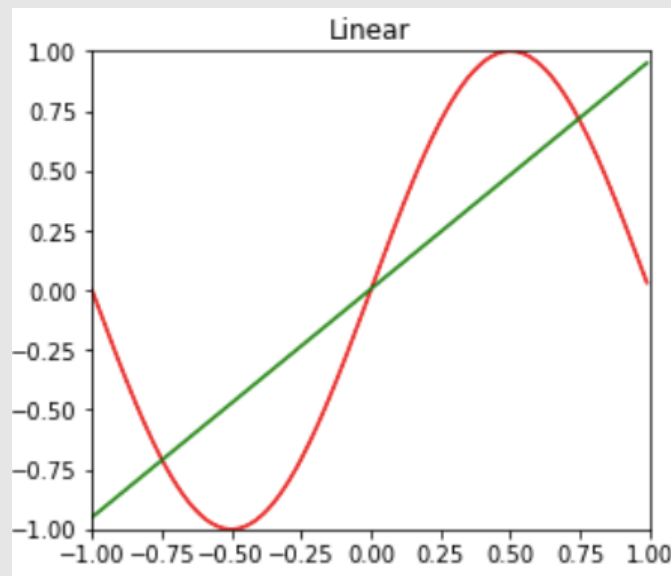


Learning H_0 vs H_0

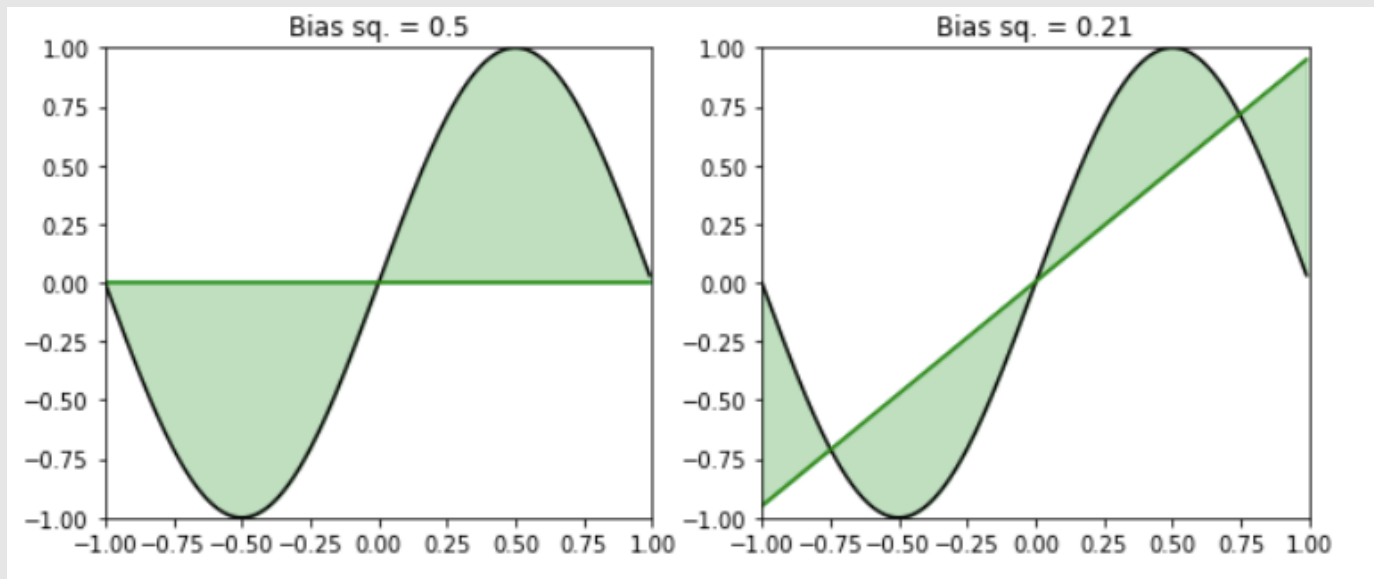
$$H_0: f(x) = c$$



$$H_1: f(x) = w_0 + w_1 x$$



Bias



$$E_X(\text{squared loss}) = E_X \left\{ \underbrace{\left[E_D(f(X) - \hat{f}(X)) \right]^2}_{\text{Bias of } \hat{f}} + \underbrace{\text{Var}_D[\hat{f}(X)]}_{\text{Variance of } \hat{f}} + \underbrace{\text{Var}_\varepsilon(\varepsilon)}_{\text{Irreducible error}} \right\}$$

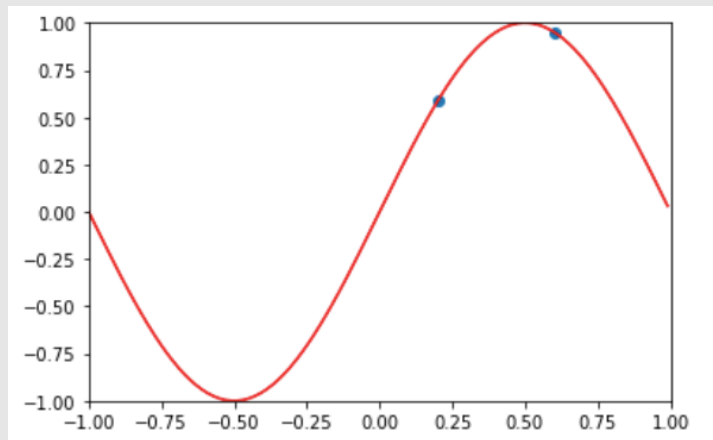
Bias of \hat{f}

Variance of \hat{f}

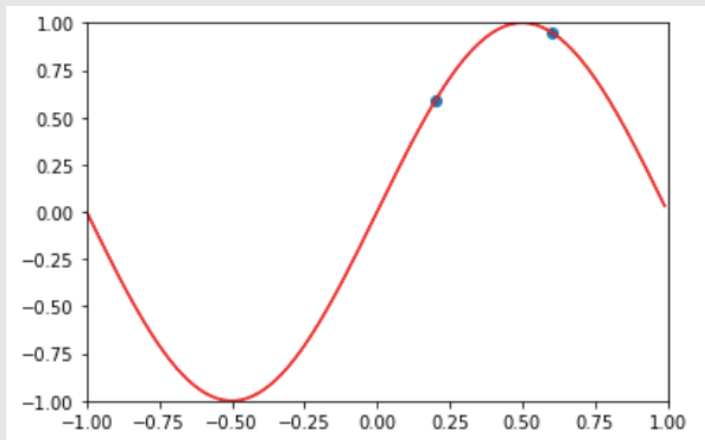
Irreducible error

What if you are only given two points?

$$H_0: f(x) = c$$

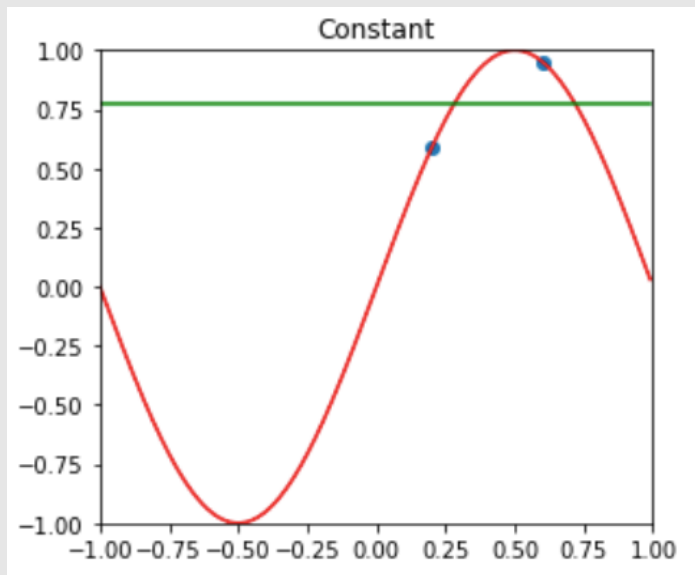


$$H_1: f(x) = w_0 + w_1x$$

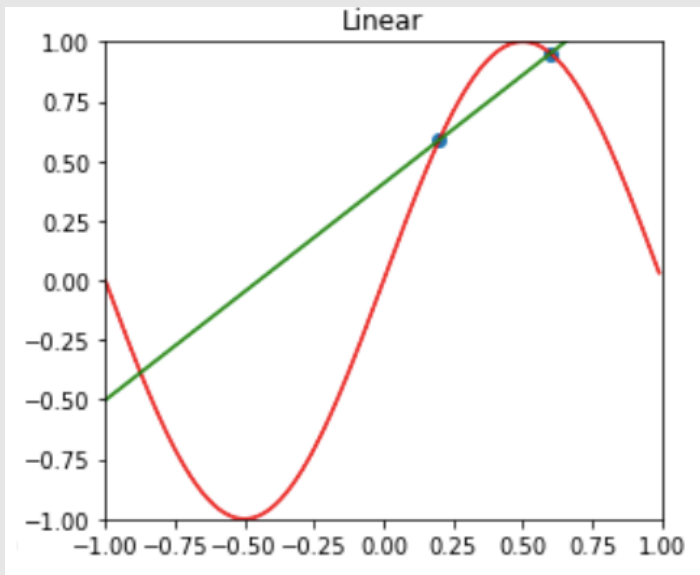


What if you are only given two points?

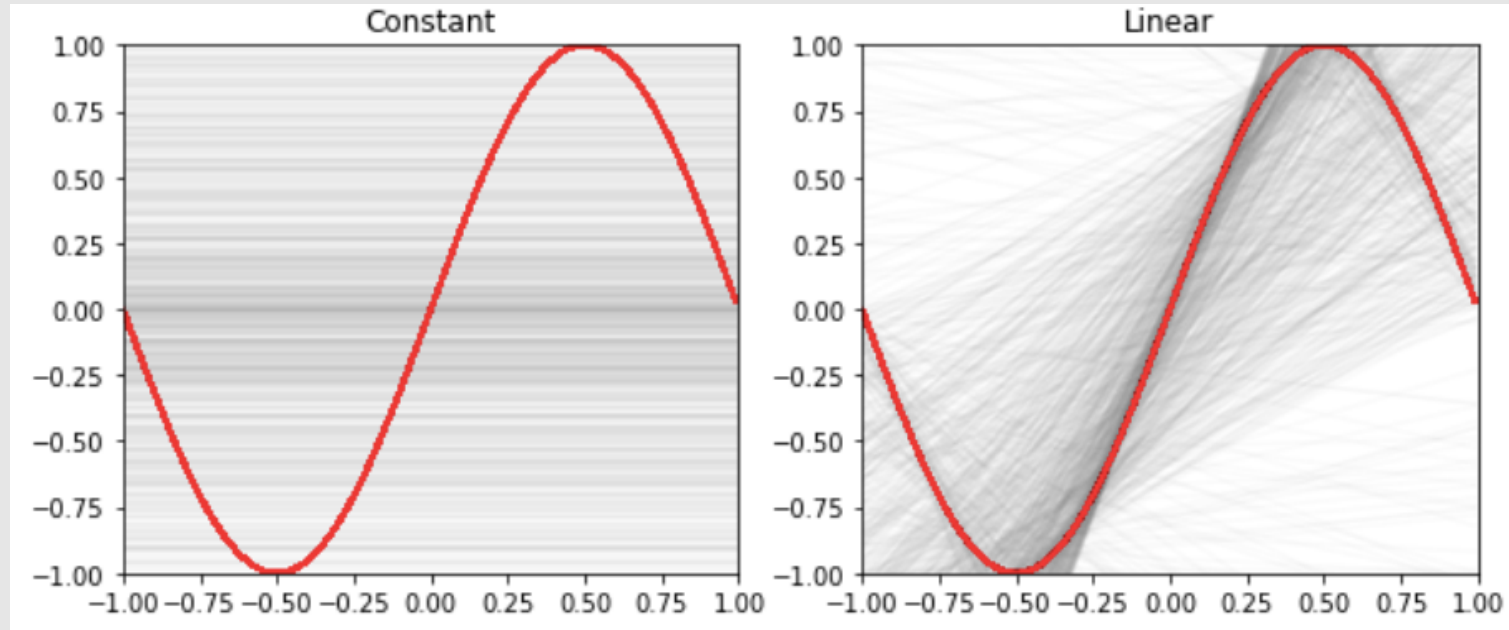
$$H_0: f(x) = c$$



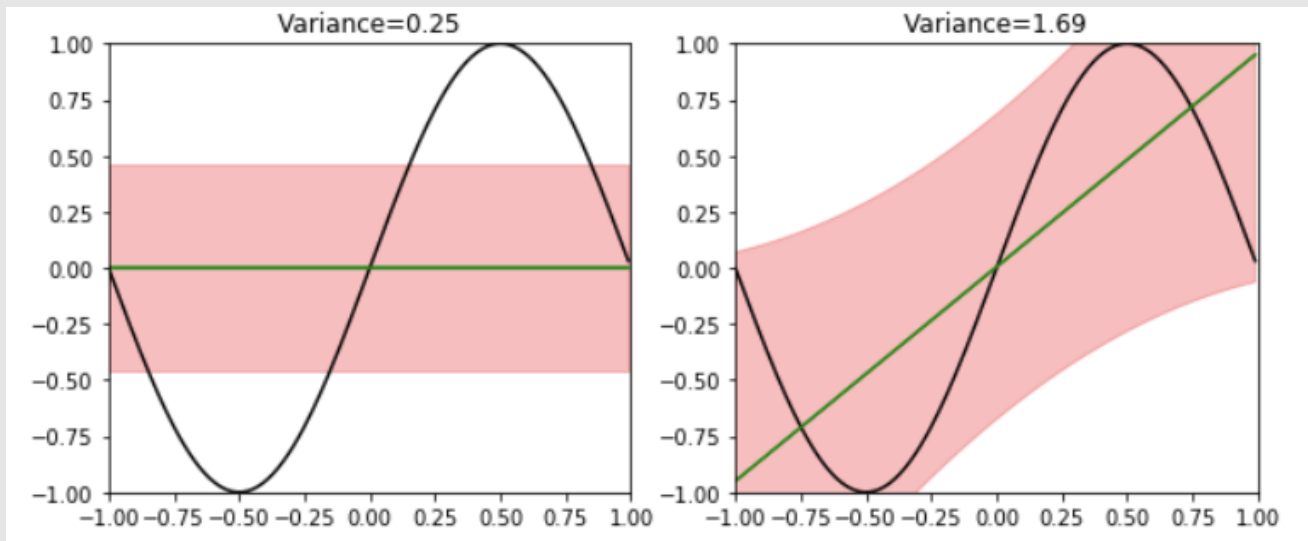
$$H_1: f(x) = w_0 + w_1x$$



Let us repeat the experiment



Variance



$$E_X(\text{squared loss}) = E_X \left\{ \underbrace{\left[E_D(f(X) - \hat{f}(X)) \right]^2}_{\text{Bias of } \hat{f}} + \underbrace{\text{Var}_D[\hat{f}(X)]}_{\text{Variance of } \hat{f}} + \underbrace{\text{Var}_\varepsilon(\varepsilon)}_{\text{Irreducible error}} \right\}$$

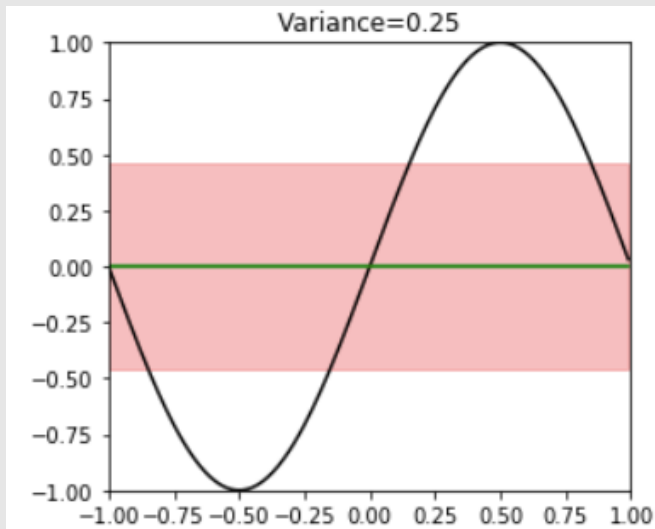
Bias of \hat{f}

Variance of \hat{f}

Irreducible error

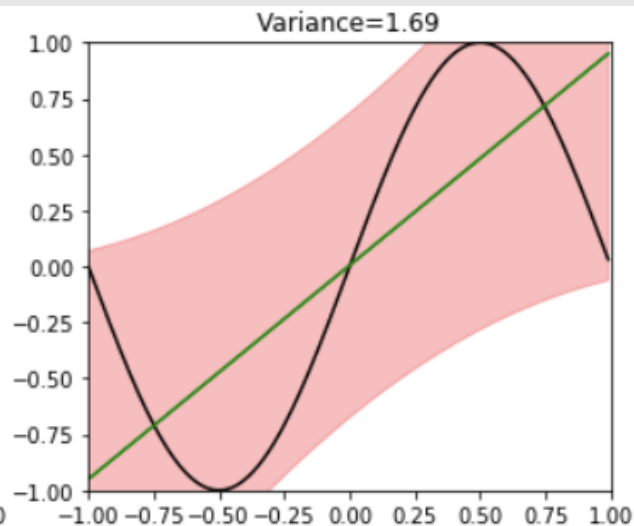
The winner is?

$$H_0: f(x) = c$$



Bias = 0.50
Variance = 0.25

$$H_1: f(x) = w_0 + w_1x$$



Bias = 0.21
Variance = 1.69

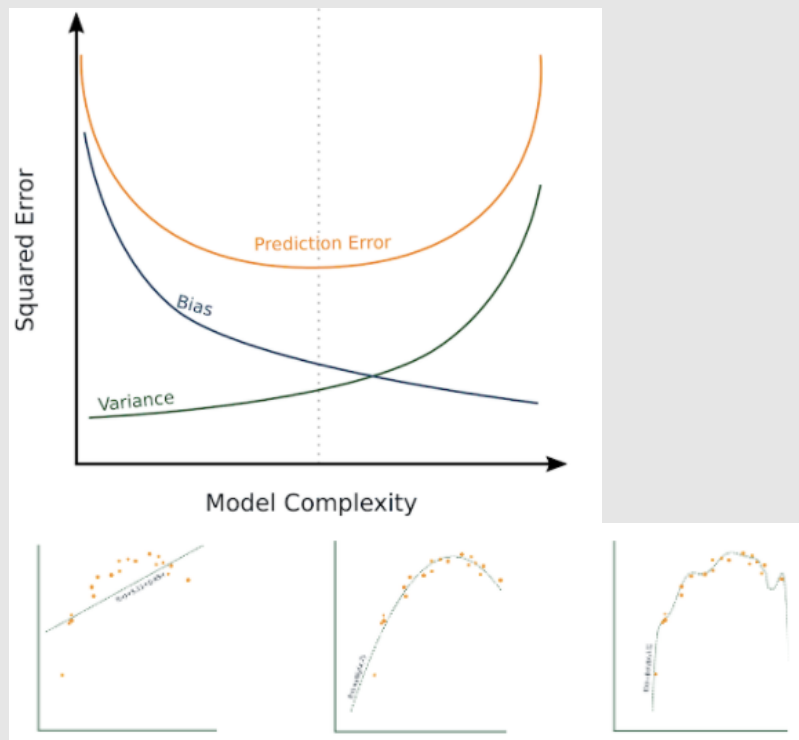
Lesson learned

Match the model complexity to

Data resources not the **response complexity**

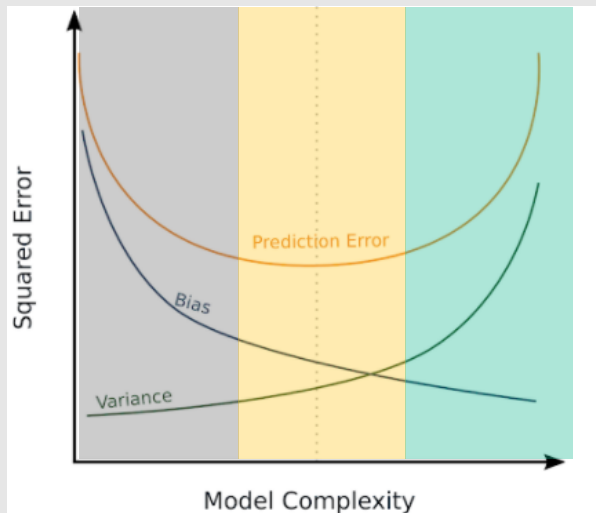
Debugging variance bias

Bias – Variance Trade-off



As we increase model complexity, bias decrease and the variance increase

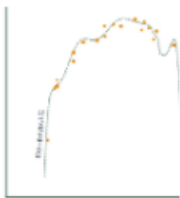
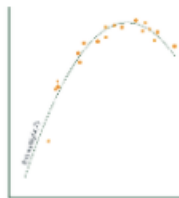
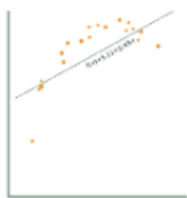
Bias – Variance Trade-off



Regime 1: high bias
Low but consistent
performance
 $\text{Train MSE} \approx \text{Test MSE}$

Regime 2: good trade-off
Acceptable MSE
Consistent MSE

Regime 3: high variance
MSE all over the place
 $\text{Train MSE} \ll \text{Test MSE}$



As we increase model complexity, bias decrease and the variance increase

Model complexity

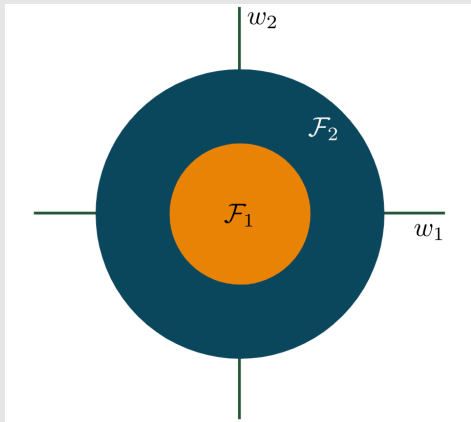
The representational capacity of a set of functions is an indicator of the representational richness within this set of functions.

- It's usually quite easy to compare different models of the same "type". For instance, consider the following two hypothesis sets of functions:

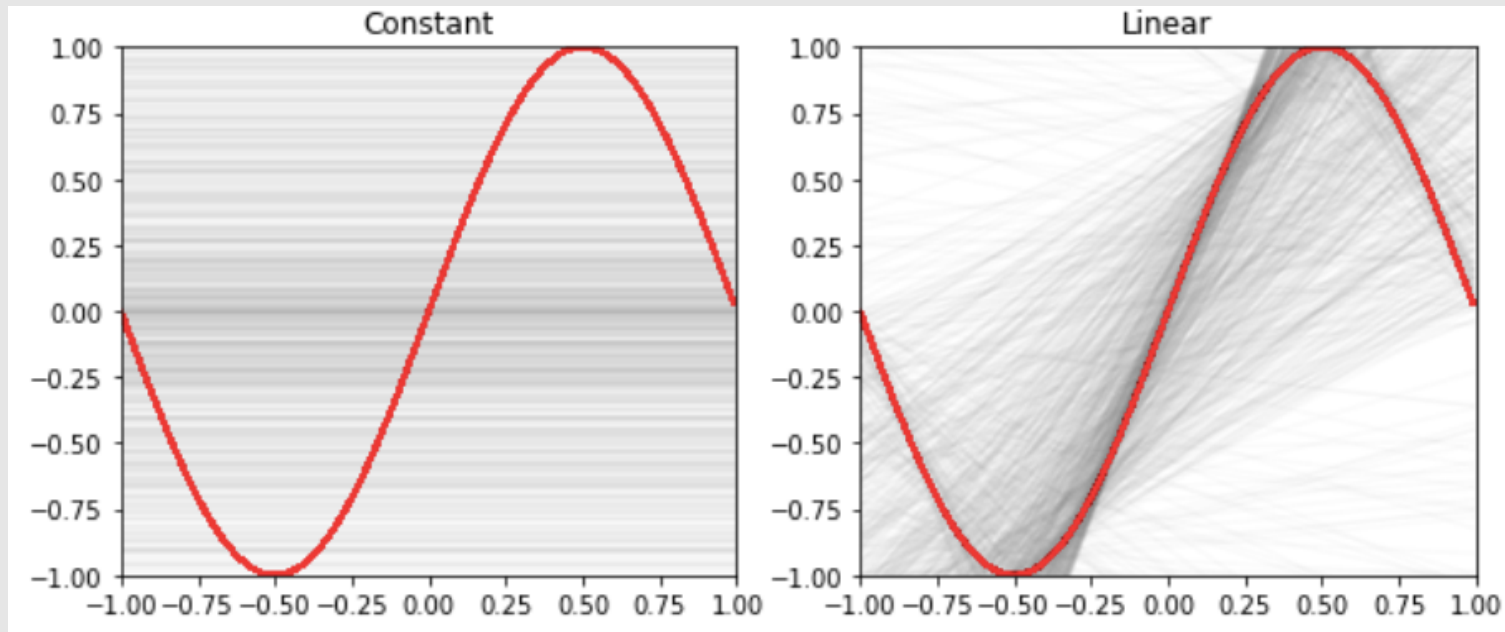
$$\mathcal{F}_1 = \{w \rightarrow w \cdot x \mid \|w\|_2 \leq W\}$$

$$\mathcal{F}_2 = \{w \rightarrow w \cdot x \mid \|w\|_2 \leq 2W\}$$

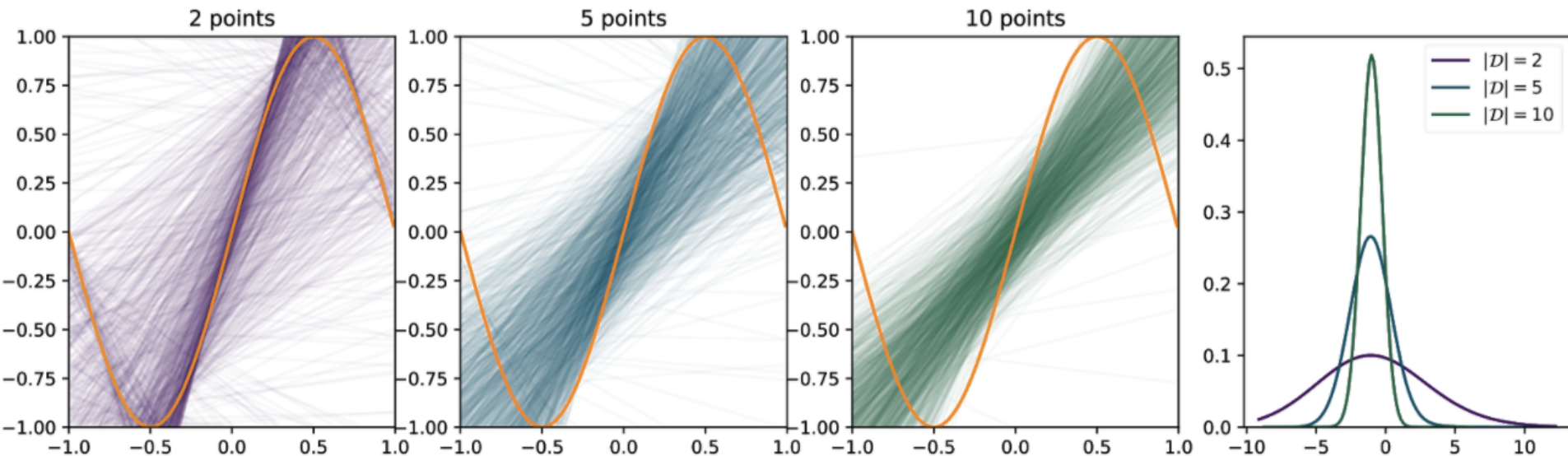
A function family with low capacity, is more likely to underfit.
A function family with high capacity, is more likely to overfit.



Experiment with two samples



Adding more data



As we increase the sample size, the variance decreases!

Theoretical results for linear regression

In linear regression

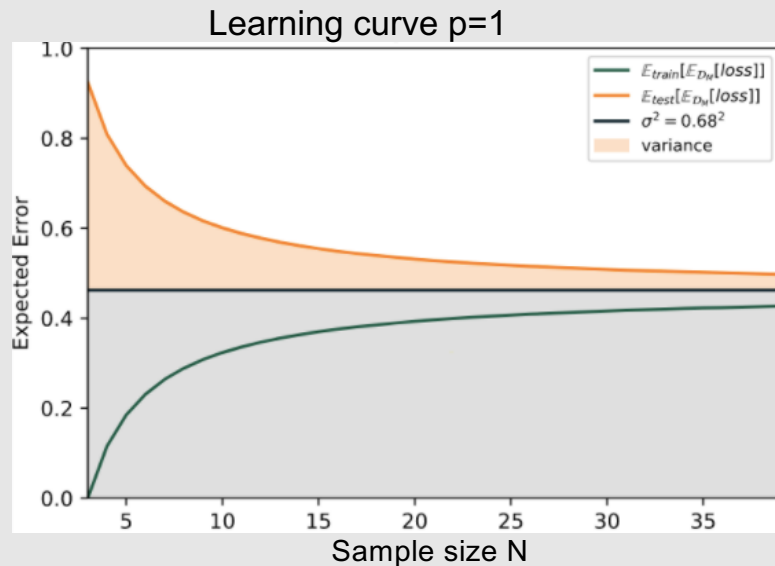
$$y = x^T \beta + \varepsilon$$

$$\hat{\beta} = (X^T X)^{-1} X^T$$

In-sample error: $\sigma^2 \left(1 - \frac{p+1}{N}\right)$

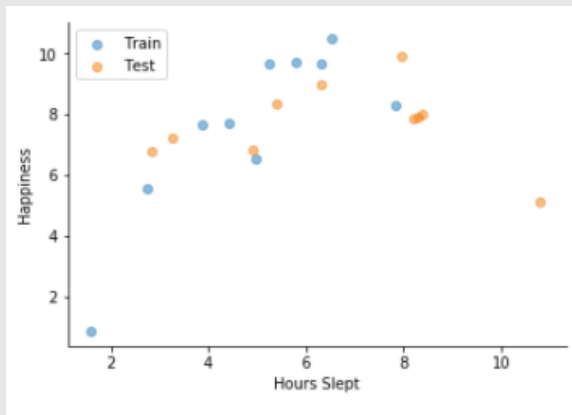
Out-of-sample error: $\sigma^2 \left(1 + \frac{p+1}{N}\right)$

Dataset

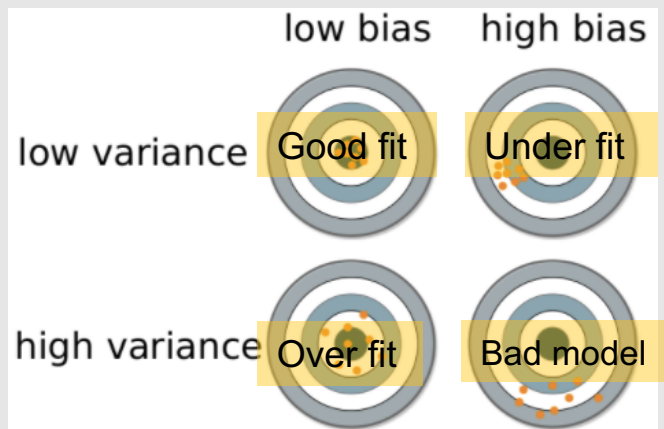


p=1 (with one predictor and the intercept)

Test on a hold-out set



	Underfit	Good fit	Overfit
Training MSE	Bad	Good	Perfect
Validation MSE	Bad	Good	Bad



Solution for high variance

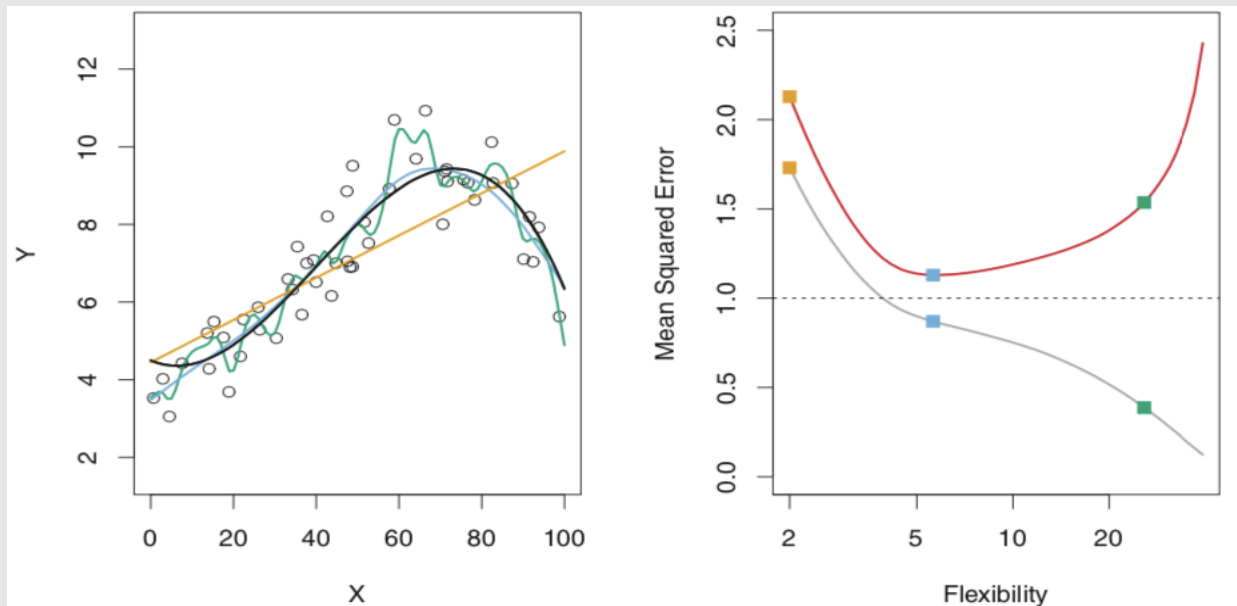
- Add more training data
- Reduce model complexity
- Bagging

Solution for high bias

- Use a more complex model
- Add extra features
- Boosting

Functional View

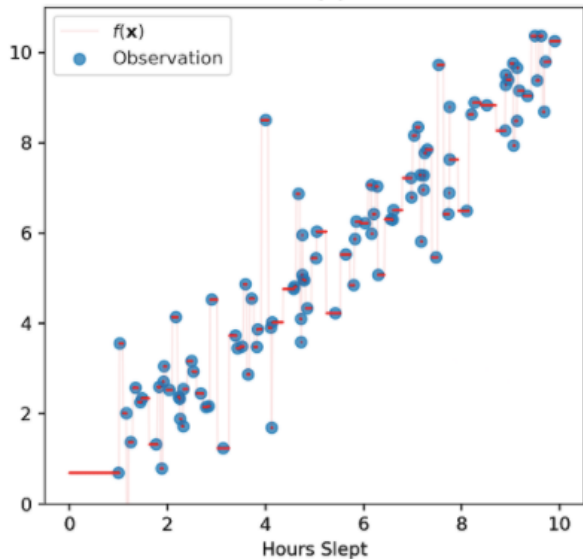
Bias-Variance Trade-off



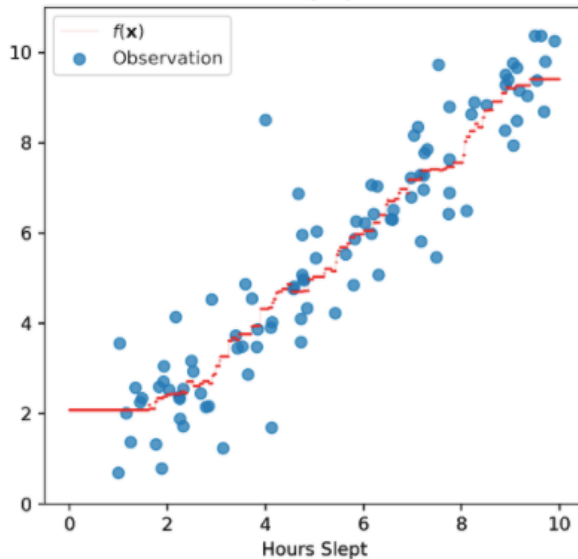
- Black curve is truth. Grey curve on right is training MSE, red curve is testing MSE.
- Orange, blue and green curves/squares correspond to fits different flexibility

Bias-Variance Trade-off

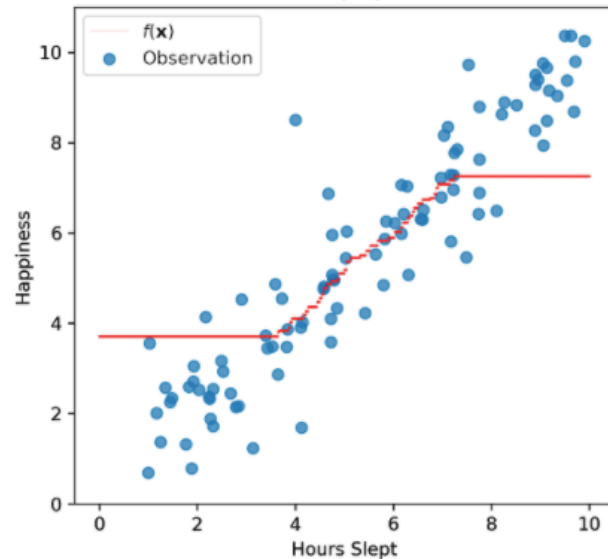
NN(1)



NN(15)



NN(60)



Which is more complex?

Which includes most bias/variance?