# Bias-Variance

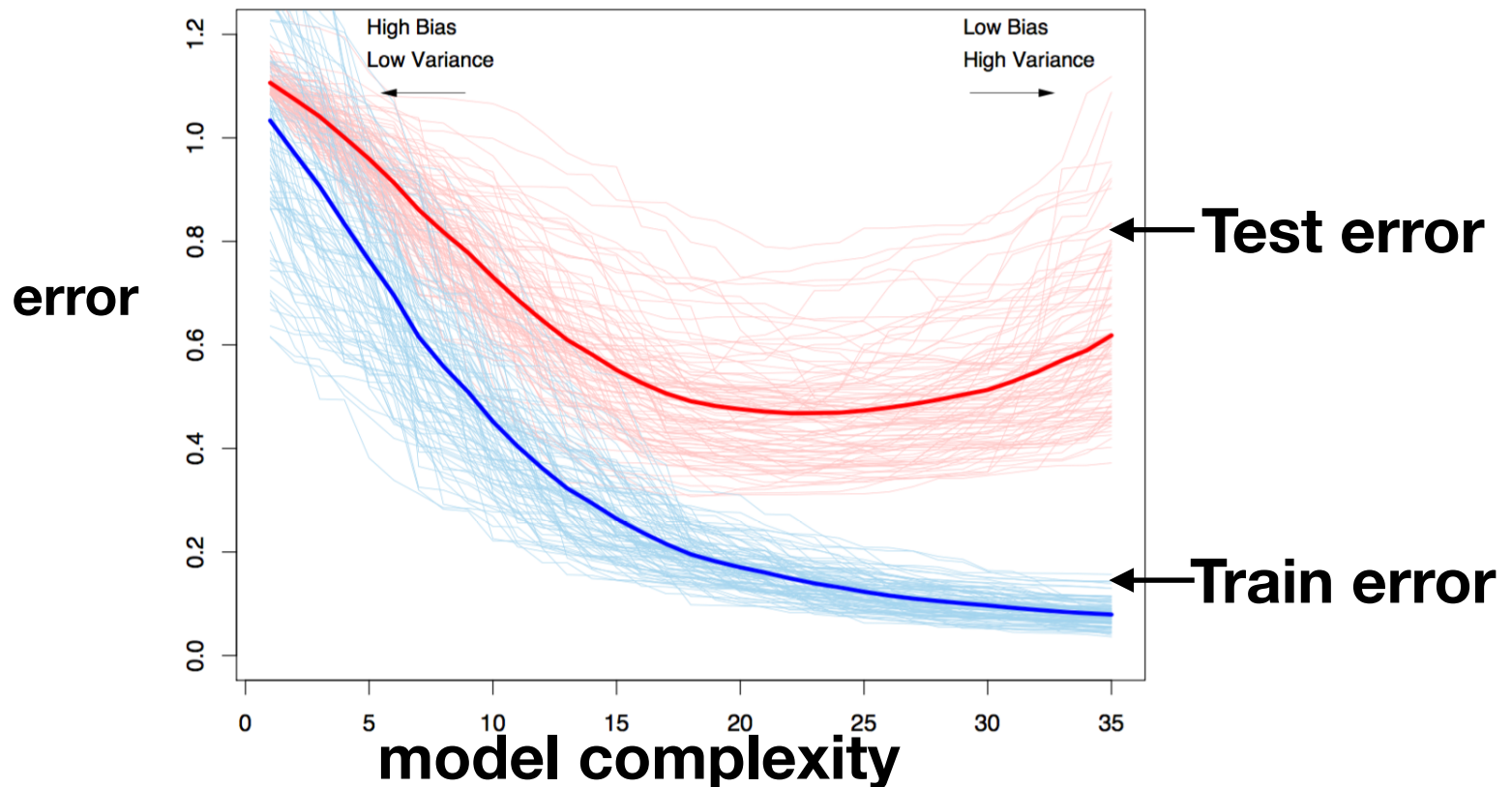| Features | Train MSE | Test MSE |
|:---:|:---:|:---:|
| All | 2640 | 3224 |
| S5 and BMI | 3004 | 3453 |
| S5 | 3869 | 4227 |
| BMI | 3540 | 4277 |
| S4 and S3 | 4251 | 5302 |
| S4 | 4278 | 5409 |
| S3 | 4607 | 5419 |
| None | 5524 | 6352 |

- **test MSE is the primary criteria for model selection**

- Using only 2 features (S5 and BMI), one can get very close to the prediction performance of using all features

- Combining S3 and S4 does not give any performance gain

demo3_diabetes.ipynb

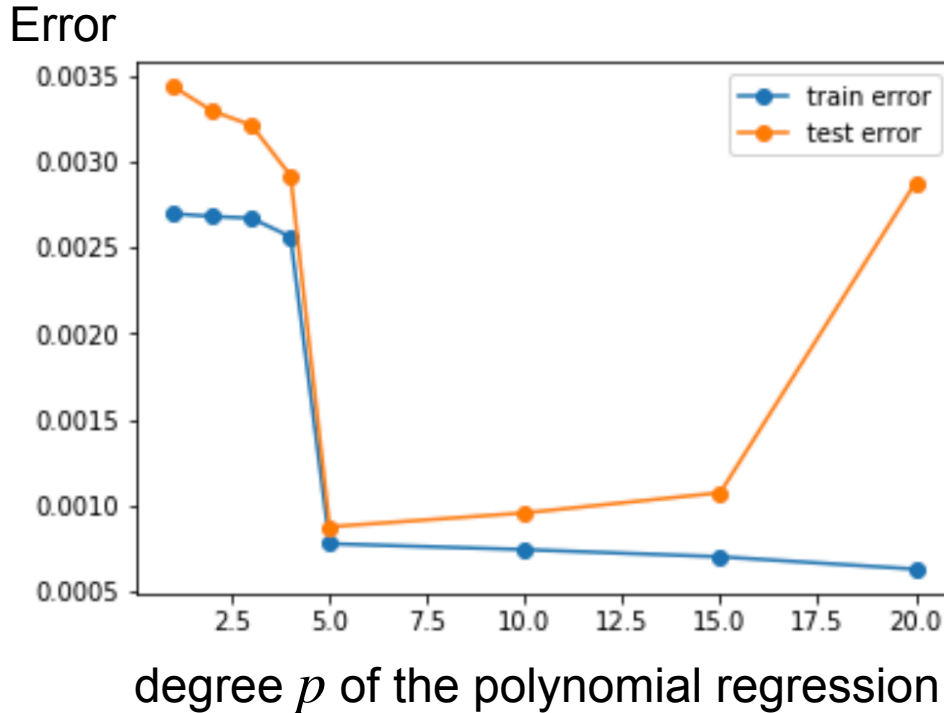# What does the bias-variance theory tell us?

- **Train error** (random variable, randomness from $\mathscr{D}$)
  - Use $\mathscr{D} = \{(x_i, y_i)\}_{i=1}^{n} \sim P_{X,Y}$ to find $\widehat{w}$
  - Train error: $\mathscr{L}_{\text{train}}(\widehat{w}_{\text{LS}}) = \dfrac{1}{|\mathscr{D}|} \sum_{(x_i, y_i) \in \mathscr{D}} (y_i - \widehat{w}^T x_i)^2$

- recall the **test error** is an unbiased estimator of the **true error**

- **True error** (random variable, randomness from $\mathscr{D}$)
  - True error: $\mathscr{L}_{\text{true}}(\widehat{w}) = \mathbb{E}_{(x,y) \sim P_{X,Y}}[(y - \widehat{w}^T x)^2]$

- **Test error** (random variable, randomness from $\mathscr{D}$ and $\mathscr{T}$)
  - Use $\mathscr{T} = \{(x_i, y_i)\}_{i=1}^{m} \sim P_{X,Y}$
  - Test error: $\mathscr{L}_{\text{test}}(\widehat{w}) = \dfrac{1}{|\mathscr{T}|} \sum_{(x_i, y_i) \in \mathscr{T}} (y_i - \widehat{w}^T x_i)^2$

- theory explains **true error**, and hence expected behavior of the (random) **test error**
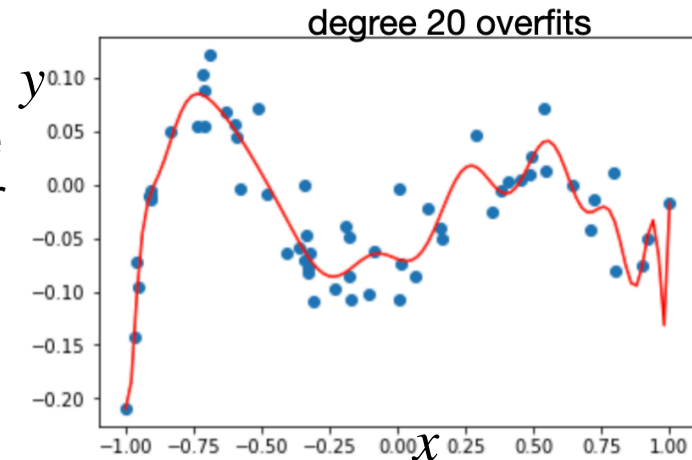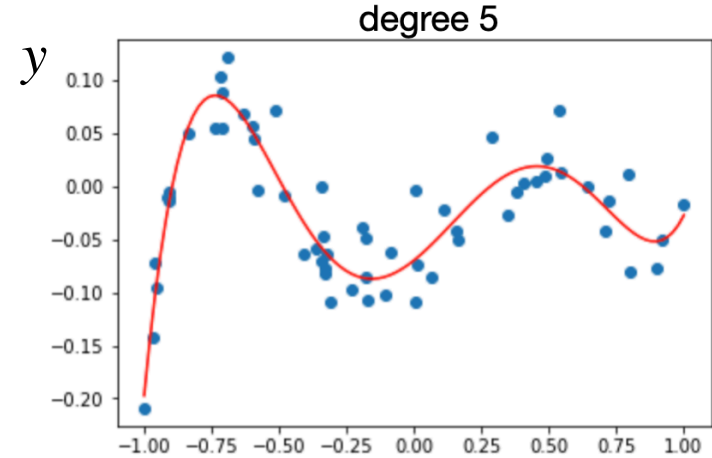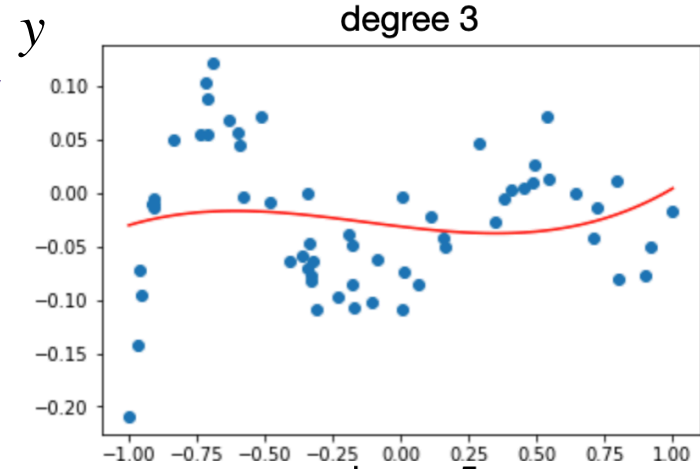
# What does bias-variance theory tell us?

- Train error is optimistically biased (i.e. smaller) because the trained model is minimizing the train error

- Test error is unbiased estimate of the true error, if test data is never used in training a model or selecting the model complexity

- Each line is an i.i.d. instance of $\mathscr{D}$ and $\mathscr{T}$

# Train/test error vs. complexity

Error



degree $p$ of the polynomial regression

- **Model complexity** e.g., degree $p$ of the polynomial model, number of features used in diabetes example
  - Related to the dimension of the model parameter
- **Train error** monotonically decreases with model complexity
- **Test error** has a U shape

# **Statistical learning**

Typical notation:
$X$ denotes a random variable
$x$ denotes a deterministic instance

- Suppose data is generated from a statistical model $(X, Y) \sim P_{X,Y}$

  - and assume we know $P_{X,Y}$ (just for now to explain statistical learning)

- **learning** aims to find a predictor $\eta : \mathbb{R}^d \to \mathbb{R}$ that minimizes

  - expected error $\mathbb{E}_{(X,Y) \sim P_{X,Y}}[(Y - \eta(X))^2]$

  - think of random $(X, Y)$ as a new sample you will encounter when you deployed your learned model, and we care about its average performance

- We assume the function $\eta(x)$ could be anything

  - it can take any value for each $X = x$

- So the optimization can be done separately for each $X = x$

  - $\mathbb{E}_{(X,Y) \sim P_{X,Y}}[(Y - \eta(X))^2] = \mathbb{E}_{X \sim P_X}\Big[\mathbb{E}_{Y \sim P_{Y|X}}[(Y - \eta(x))^2 \,|\, X = x]\Big]$

$$= \int \mathbb{E}_{Y \sim P_{Y|X}}[(Y - \eta(x))^2 \,|\, X = x]\, P_X(x)\, dx$$

Or for discrete $X$, $\qquad = \sum_x P_X(x)\, \mathbb{E}_{Y \sim P_{Y|X}}[(Y - \eta(x))^2 \,|\, X = x]$

Where we used the chain rule: $\mathbb{E}_{X,Y}[f(X, Y)] = \mathbb{E}_X\Big[\mathbb{E}_{Y|X}[f(x, Y) \,|\, X = x]\Big]$
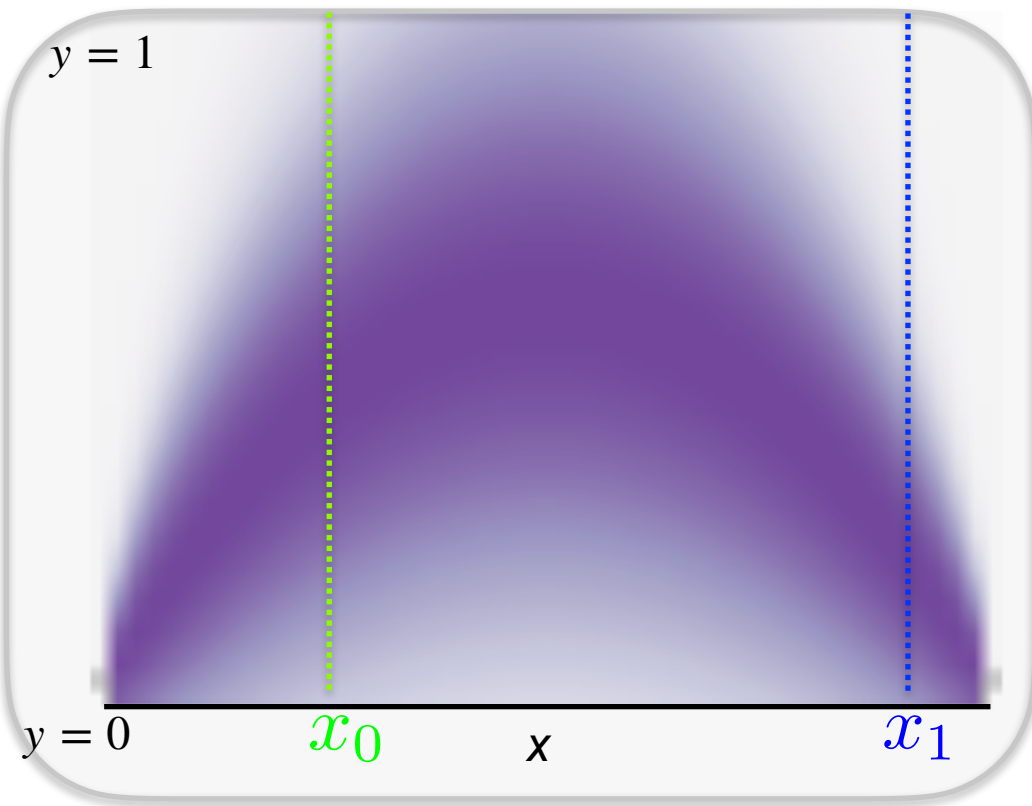
# Statistical learning

- The optimal predictor sets its value for each $X = x$ separately

  - $\eta(x) = \arg\min\limits_{a \in \mathbb{R}} \mathbb{E}_{Y \sim P_{Y|X}}[(Y - a)^2 \,|\, X = x]$

- The optimal solution is $\eta(x) = \mathbb{E}_{Y \sim P_{Y|X}}[Y \,|\, X = x]$,

  which is the best prediction in $\ell_2$-loss/Mean Squared Error

- Claim: $\mathbb{E}_{Y \sim P_{Y|X}}[Y \,|\, X = x] = \arg\min\limits_{a \in \mathbb{R}} \mathbb{E}_{Y \sim P_{Y|X}}[(Y - a)^2 \,|\, X = x]$

- Proof:

- Can't implement optimal statistical estimator $\eta(x) = \mathbb{E}[Y \,|\, X = x]$

  - as we do not know $P_{X,Y}$ in practice

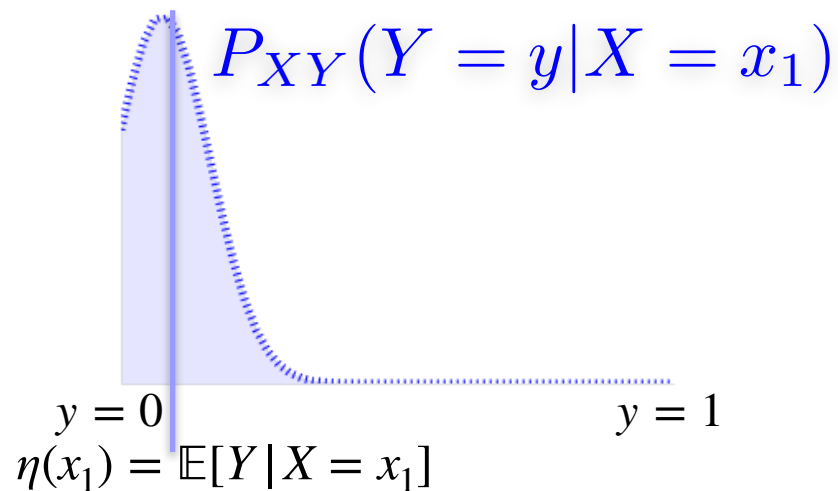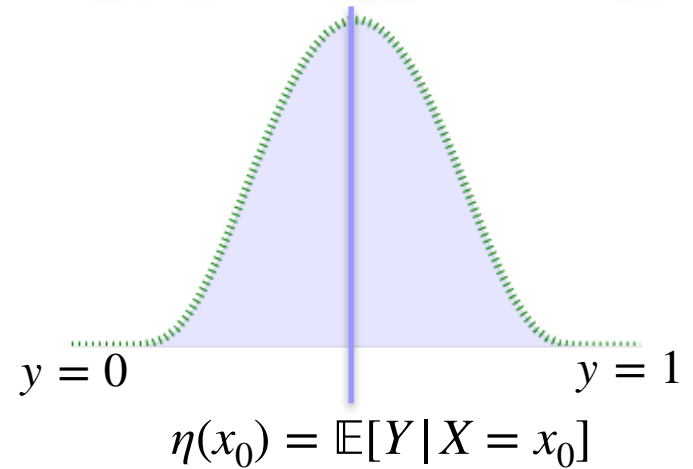- This is only for the purpose of conceptual understanding

# **Statistical Learning**

Ideally, we want to find:

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

$$P_{XY}(X = x, Y = y)$$

$$P_{XY}(Y = y|X = x_0)$$



$y = 1$

$y = 0$       $y = 1$

$$\eta(x_0) = \mathbb{E}[Y|X = x_0]$$

$y = 0$

$x_0$    $x$    $x_1$

$$P_{XY}(Y = y|X = x_1)$$

$y = 0$       $y = 1$

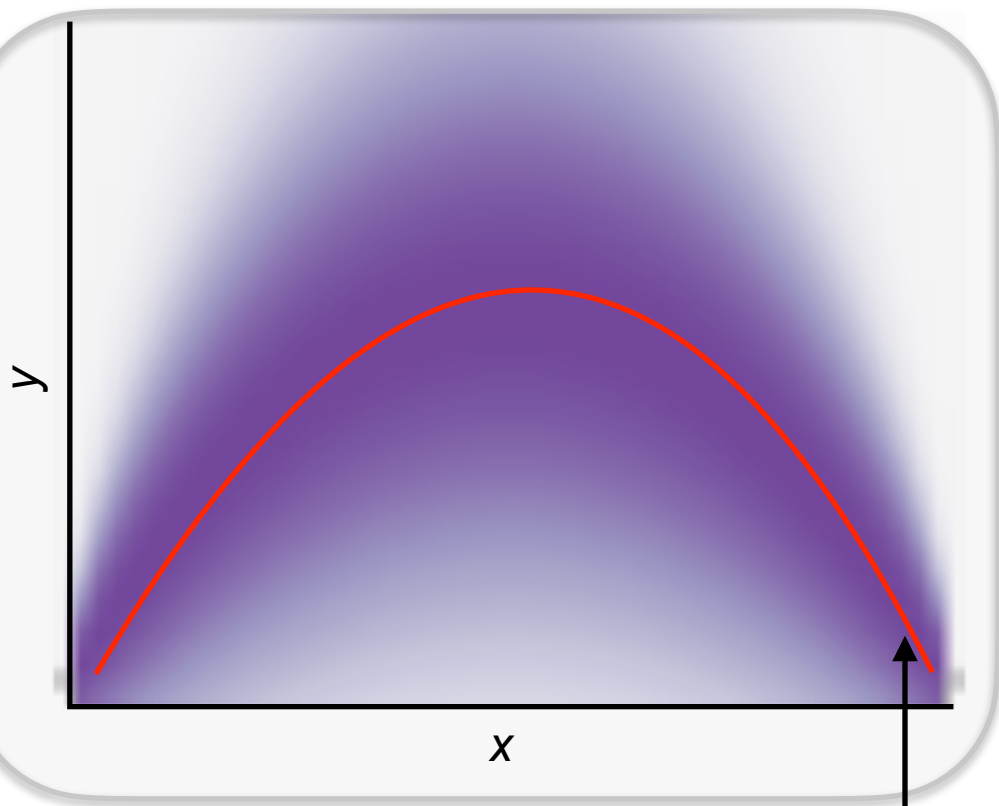$$\eta(x_1) = \mathbb{E}[Y|X = x_1]$$

# Statistical Learning

$$P_{XY}(X = x, Y = y)$$

Ideally, we want to find:

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$
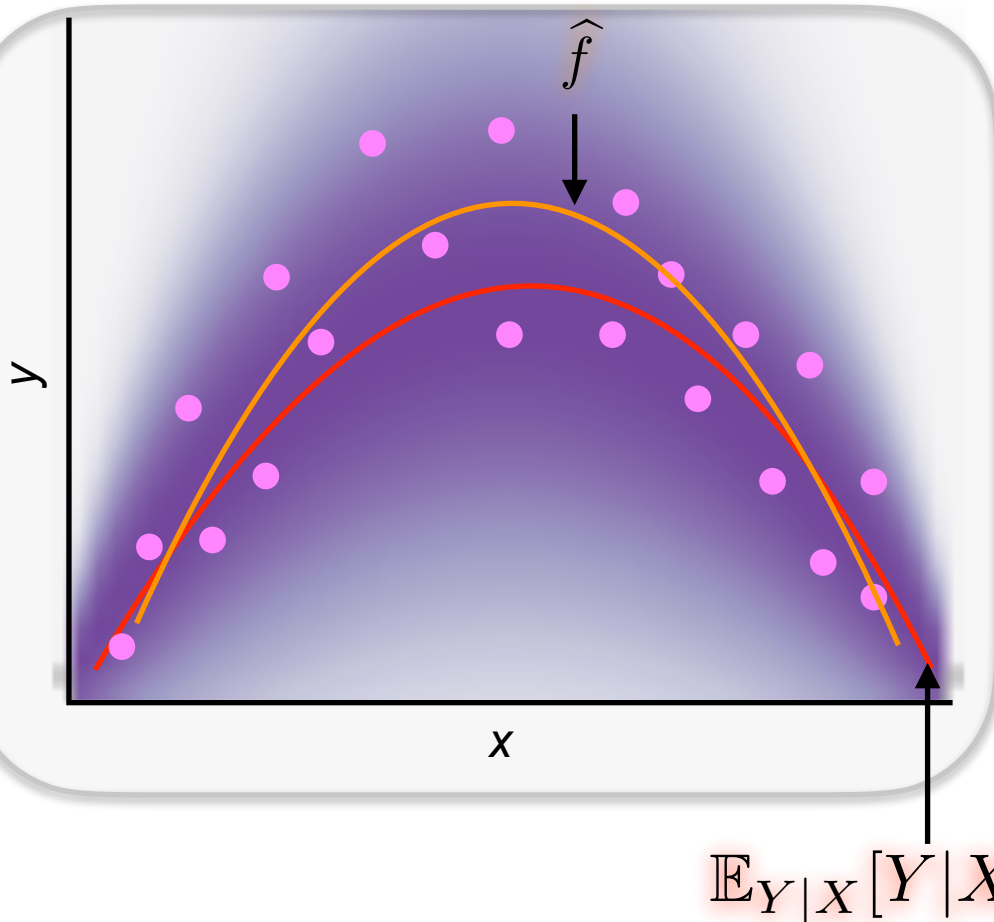
But we do not know $P_{X,Y}$

We only have samples.

$y$

$x$

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

# Statistical Learning

$$P_{XY}(X = x, Y = y)$$

$\widehat{f}$

$$\mathbb{E}_{Y|X}[Y|X = x]$$

Ideally, we want to find:

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

But we only have samples:

$$(x_i, y_i) \overset{i.i.d.}{\sim} P_{XY} \quad \text{for } i = 1, \ldots, n$$

So we need to restrict our predictor to a function class (e.g., linear, degree-$p$ polynomial) to avoid overfitting:
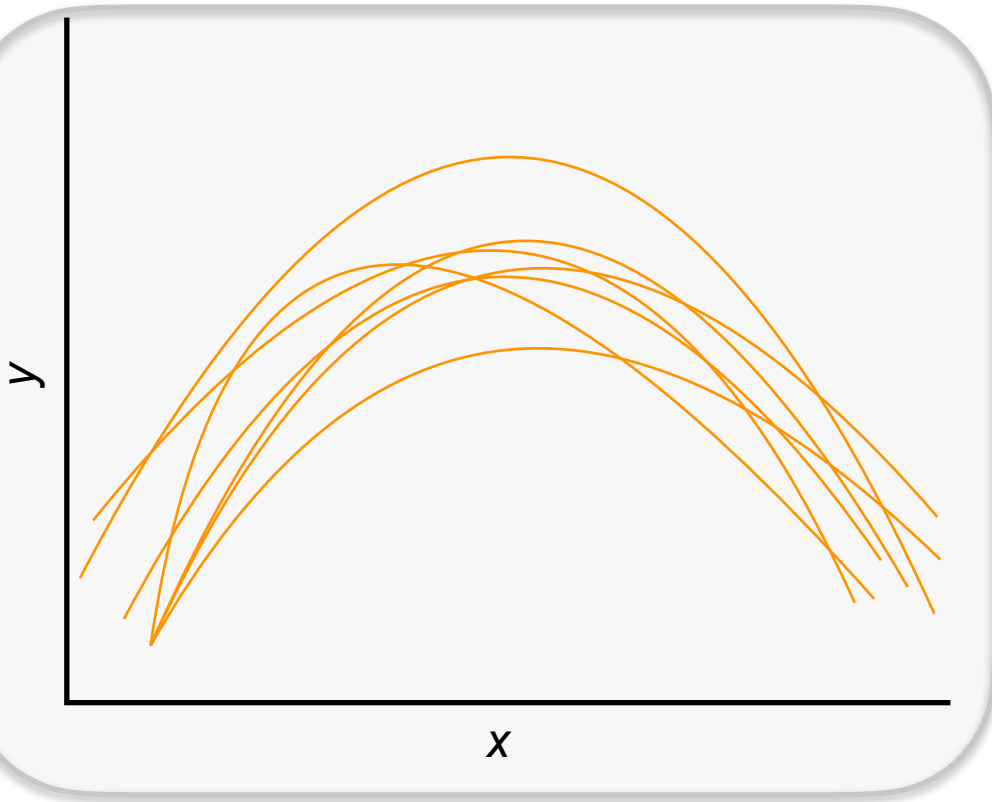
$$\widehat{f} = \arg\min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^{n} (y_i - f(x_i))^2$$

We care about how our predictor performs on future unseen data
True Error of $\hat{f}$: $\mathbb{E}_{X,Y}[(Y - \hat{f}(X))^2]$

**Future prediction error** $\mathbb{E}_{X,Y}[(Y - \hat{f}(X))^2]$ **is random**
**because $\hat{f}$ is random (whose randomness comes from training data $\mathscr{D}$)**

$$P_{XY}(X = x, Y = y)$$



Each draw $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ results in different $\widehat{f}$

# Bias-variance tradeoff

**Ideal predictor**

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

**Learned predictor**

$$\hat{f}_{\mathscr{D}} = \arg\min_{f \in \mathscr{F}} \frac{1}{|\mathscr{D}|} \sum_{(x_i, y_i) \in \mathscr{D}} (y_i - f(x_i))^2$$

- We are interested in the **True Error** of a (random) learned predictor:

$$\mathbb{E}_{X,Y}[(Y - \hat{f}_{\mathscr{D}}(X))^2]$$

- But the analysis can be done for each $X = x$ separately, so we analyze the **conditional true error**:

$$\mathbb{E}_{Y|X}[(Y - \hat{f}_{\mathscr{D}}(x))^2 \,|\, X = x]$$

- And we care about the **average conditional true error**, averaged over training data:

$$\mathbb{E}_{\mathscr{D}}\left[ \mathbb{E}_{Y|X}[(Y - \hat{f}_{\mathscr{D}}(x))^2 \,|\, X = x] \right]$$

written compactly as $\quad = \mathbb{E}[(Y - \hat{f}_{\mathscr{D}}(x))^2]$

# Bias-variance tradeoff

**Ideal predictor**

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

**Learned predictor**

$$\hat{f}_{\mathscr{D}} = \arg\min_{f\in\mathscr{F}} \frac{1}{|\mathscr{D}|} \sum_{(x_i,y_i)\in\mathscr{D}} (y_i - f(x_i))^2$$

- **Average conditional true error**:

$$\mathbb{E}_{\mathscr{D},Y|x}[(Y - \hat{f}_{\mathscr{D}}(x))^2] = \mathbb{E}_{\mathscr{D},Y|x}[(Y - \eta(x) + \eta(x) - \hat{f}_{\mathscr{D}}(x))^2]$$

# Bias-variance tradeoff

**Ideal predictor**
$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

**Learned predictor**
$$\hat{f}_{\mathscr{D}} = \arg\min_{f \in \mathscr{F}} \frac{1}{|\mathscr{D}|} \sum_{(x_i, y_i) \in \mathscr{D}} (y_i - f(x_i))^2$$

- **Average conditional true error**:

$$\mathbb{E}_{\mathscr{D}, Y|x}[(Y - \hat{f}_{\mathscr{D}}(x))^2] = \mathbb{E}_{\mathscr{D}, Y|x}[(Y - \eta(x) + \eta(x) - \hat{f}_{\mathscr{D}}(x))^2]$$

$$= \mathbb{E}_{\mathscr{D}, Y|x}\left[ (Y - \eta(x))^2 + 2(Y - \eta(x))(\eta(x) - \hat{f}_{\mathscr{D}}(x)) + (\eta(x) - \hat{f}_{\mathscr{D}}(x))^2 \right]$$

$$= \mathbb{E}_{Y|x}[(Y - \eta(x))^2] + 2\underbrace{\mathbb{E}_{\mathscr{D}, Y|x}[(Y - \eta(x))(\eta(x) - \hat{f}_{\mathscr{D}}(x))]}_{=0} + \mathbb{E}_{\mathscr{D}}[(\eta(x) - \hat{f}_{\mathscr{D}}(x))^2]$$

(this follows from independence of $\mathscr{D}$ and $(X, Y)$ and
$\mathbb{E}_{Y|x}[Y - \eta(x)] = \mathbb{E}[Y|X = x] - \eta(x) = 0$)

$$= \mathbb{E}_{Y|x}[(Y - \eta(x))^2] \qquad + \qquad \mathbb{E}_{\mathscr{D}}[(\eta(x) - \hat{f}_{\mathscr{D}}(x))^2]$$

**Irreducible error**
(a) Caused by stochastic
label noise in $P_{Y|X=x}$
(b) cannot be reduced

**Average learning error**
Caused by
*(a)* either using too "simple" of a model or
*(b)* not enough data to learn the model accurately

# Bias-variance tradeoff

**Ideal predictor**

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

**Learned predictor**

$$\hat{f}_{\mathscr{D}} = \arg\min_{f \in \mathscr{F}} \frac{1}{|\mathscr{D}|} \sum_{(x_i, y_i) \in \mathscr{D}} (y_i - f(x_i))^2$$

- **Average learning error**:

$$\mathbb{E}_{\mathscr{D}}[(\eta(x) - \hat{f}_{\mathscr{D}}(x))^2] = \mathbb{E}_{\mathscr{D}}\left[\left(\eta(x) - \mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)] + \mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)] - \hat{f}_{\mathscr{D}}(x)\right)^2\right]$$

# Bias-variance tradeoff

**Ideal predictor**

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

**Learned predictor**

$$\hat{f}_{\mathscr{D}} = \arg\min_{f \in \mathscr{F}} \frac{1}{|\mathscr{D}|} \sum_{(x_i, y_i) \in \mathscr{D}} (y_i - f(x_i))^2$$

# Bias-variance tradeoff

### Ideal predictor

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

### Learned predictor

$$\hat{f}_{\mathscr{D}} = \arg\min_{f \in \mathscr{F}} \frac{1}{|\mathscr{D}|} \sum_{(x_i, y_i) \in \mathscr{D}} (y_i - f(x_i))^2$$

- **Average learning error**:

# Bias-variance tradeoff

**Ideal predictor**

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

**Learned predictor**

$$\hat{f}_{\mathscr{D}} = \arg\min_{f \in \mathscr{F}} \frac{1}{|\mathscr{D}|} \sum_{(x_i, y_i) \in \mathscr{D}} (y_i - f(x_i))^2$$

- **Average learning error**:

$$\mathbb{E}_{\mathscr{D}}[(\eta(x) - \hat{f}_{\mathscr{D}}(x))^2] = \mathbb{E}_{\mathscr{D}}\left[ \left( \eta(x) - \mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)] + \mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)] - \hat{f}_{\mathscr{D}}(x) \right)^2 \right]$$

# Bias-variance tradeoff

**Ideal predictor**

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

**Learned predictor**

$$\hat{f}_{\mathscr{D}} = \arg\min_{f \in \mathscr{F}} \frac{1}{|\mathscr{D}|} \sum_{(x_i, y_i) \in \mathscr{D}} (y_i - f(x_i))^2$$

- **Average learning error**:

$$\mathbb{E}_{\mathscr{D}}[(\eta(x) - \hat{f}_{\mathscr{D}}(x))^2] = \mathbb{E}_{\mathscr{D}}\left[ \left( \eta(x) - \mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)] + \mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)] - \hat{f}_{\mathscr{D}}(x) \right)^2 \right]$$

$$= \mathbb{E}_{\mathscr{D}}\left[ \left( \eta(x) - \mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)] \right)^2 + 2(\eta(x) - \mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)])(\mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)] - \hat{f}_{\mathscr{D}}(x)) \right.$$

# Bias-variance tradeoff

**Ideal predictor**

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

**Learned predictor**

$$\hat{f}_{\mathscr{D}} = \arg\min_{f \in \mathscr{F}} \frac{1}{|\mathscr{D}|} \sum_{(x_i, y_i) \in \mathscr{D}} (y_i - f(x_i))^2$$

- **Average learning error**:

$$\mathbb{E}_{\mathscr{D}}[(\eta(x) - \hat{f}_{\mathscr{D}}(x))^2] = \mathbb{E}_{\mathscr{D}}\left[ \left( \eta(x) - \mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)] + \mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)] - \hat{f}_{\mathscr{D}}(x) \right)^2 \right]$$

$$= \mathbb{E}_{\mathscr{D}}\left[ \left( \eta(x) - \mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)] \right)^2 + 2(\eta(x) - \mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)])(\mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)] - \hat{f}_{\mathscr{D}}(x)) \right.$$

$$\left. + (\mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)] - \hat{f}_{\mathscr{D}}(x))^2 \right]$$

# Bias-variance tradeoff

**Ideal predictor**

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

**Learned predictor**

$$\hat{f}_{\mathscr{D}} = \arg\min_{f \in \mathscr{F}} \frac{1}{|\mathscr{D}|} \sum_{(x_i, y_i) \in \mathscr{D}} (y_i - f(x_i))^2$$

- **Average learning error**:

$$\mathbb{E}_{\mathscr{D}}[(\eta(x) - \hat{f}_{\mathscr{D}}(x))^2] = \mathbb{E}_{\mathscr{D}}\left[ \left( \eta(x) - \mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)] + \mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)] - \hat{f}_{\mathscr{D}}(x) \right)^2 \right]$$

$$= \mathbb{E}_{\mathscr{D}}\left[ \left( \eta(x) - \mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)] \right)^2 + 2(\eta(x) - \mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)])(\mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)] - \hat{f}_{\mathscr{D}}(x)) \right.$$

$$\left. + (\mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)] - \hat{f}_{\mathscr{D}}(x))^2 \right]$$

$$= \left( \eta(x) - \mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)] \right)^2 + \mathbb{E}_{\mathscr{D}}\left[ (\mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)] - \hat{f}_{\mathscr{D}}(x))^2 \right]$$

# Bias-variance tradeoff

**Ideal predictor**

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

**Learned predictor**

$$\hat{f}_{\mathcal{D}} = \arg\min_{f \in \mathcal{F}} \frac{1}{|\mathcal{D}|} \sum_{(x_i, y_i) \in \mathcal{D}} (y_i - f(x_i))^2$$

- **Average learning error**:

$$\mathbb{E}_{\mathcal{D}}[(\eta(x) - \hat{f}_{\mathcal{D}}(x))^2] = \mathbb{E}_{\mathcal{D}}\left[ \left( \eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] + \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x) \right)^2 \right]$$

$$= \mathbb{E}_{\mathcal{D}}\left[ \left( \eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] \right)^2 + 2(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])(\mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x)) \right.$$

$$\left. + (\mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x))^2 \right]$$

$$= \underbrace{\left( \eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] \right)^2}_{\textbf{biased squared}} + \underbrace{\mathbb{E}_{\mathcal{D}}\left[ \left( \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x) \right)^2 \right]}_{\textbf{variance}}$$
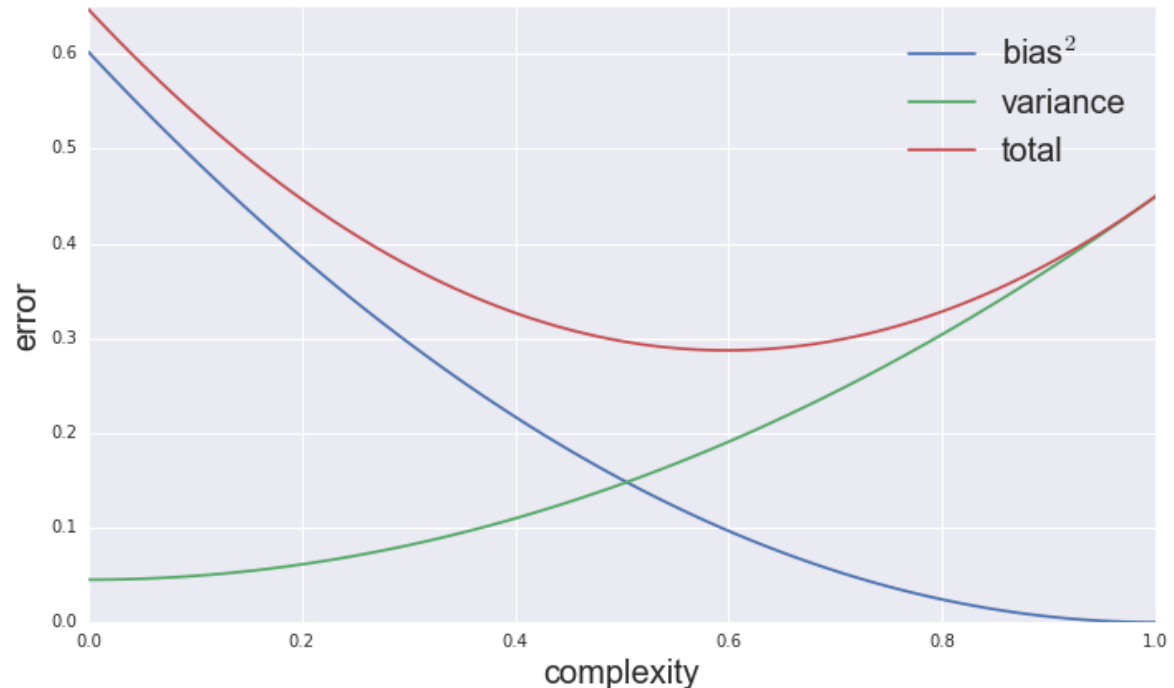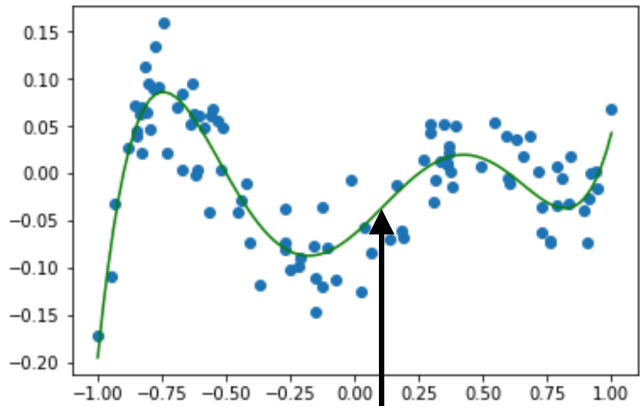
# Bias-variance tradeoff

- **Average conditional true error**:

$$\mathbb{E}_{\mathscr{D},Y|x}[(Y - \hat{f}_{\mathscr{D}}(x))^2] \;=\; \underbrace{\mathbb{E}_{Y|x}\Big[(Y - \eta(x))^2\Big]}_{\textbf{irreducible error}}$$

$$+\;\underbrace{\big(\eta(x) - \mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)]\big)^2}_{\textbf{biased squared}}\;+\;\underbrace{\mathbb{E}_{\mathscr{D}}\Big[\big(\mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)] - \hat{f}_{\mathscr{D}}(x)\big)^2\Big]}_{\textbf{variance}}$$

**Bias squared:**
measures how the predictor is mismatched with the best predictor in expectation

**variance:**
measures how the predictor varies each time with a new training datasets
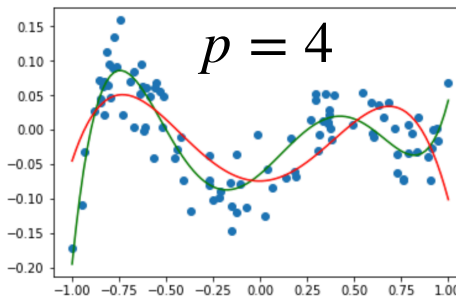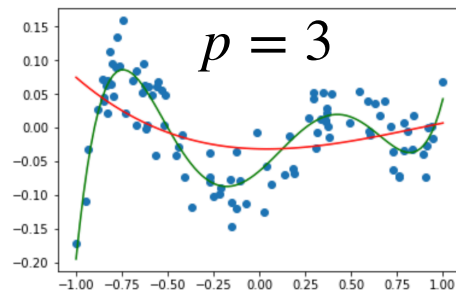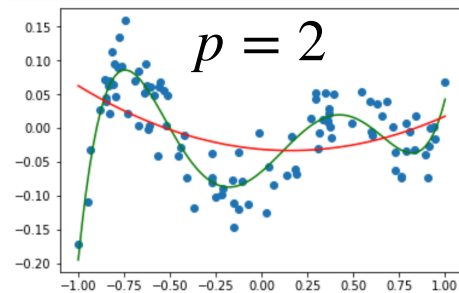
# Questions?

# Test error vs. model complexity

Simple model:
Model complexity is below the complexity of $\eta(x)$

Complex model:
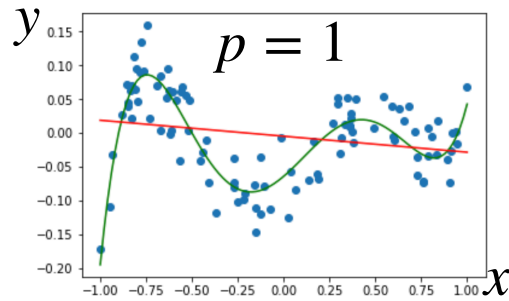


Optimal predictor $\eta(x)$ is degree-5 polynomial

Error

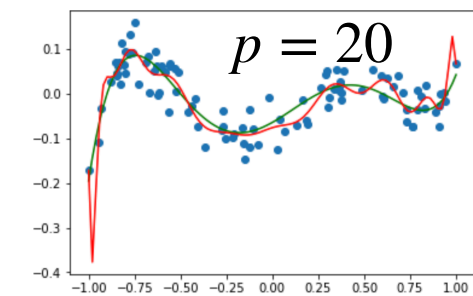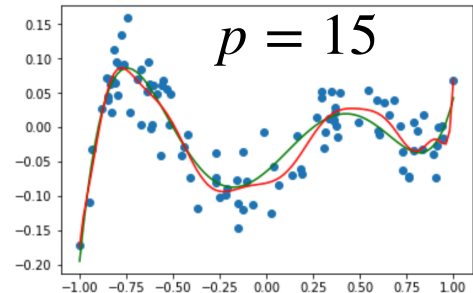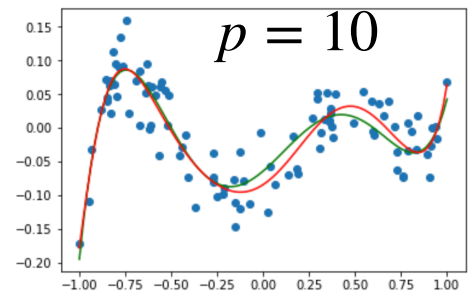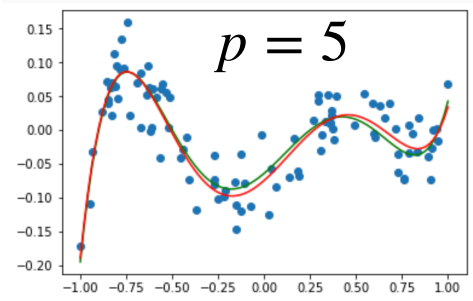Test Error
Train Error

degree $p$ of the polynomial regression

$p = 1$
$p = 2$
$p = 3$
$p = 4$
$p = 5$
$p = 10$
$p = 15$
$p = 20$

demo4_tradeoff.ipynb

# Recap: Bias-variance tradeoff with simple model



(Conceptual) bias variance tradeoff

$\eta(x)$

$p = 4$

Average predictor $\mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)]$

- When model **complexity is low** (lower than the optimal predictor $\eta(x)$)
  - $\text{Bias}^2$ of our predictor, $\left( \eta(x) - \mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)] \right)^2$, is large
  - Variance of our predictor, $\mathbb{E}_{\mathscr{D}}\left[ \left( \mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)] - \hat{f}_{\mathscr{D}}(x) \right)^2 \right]$, is small
  - If we have more samples, then
    - Bias
    - Variance
    - Because Variance is already small, overall test error
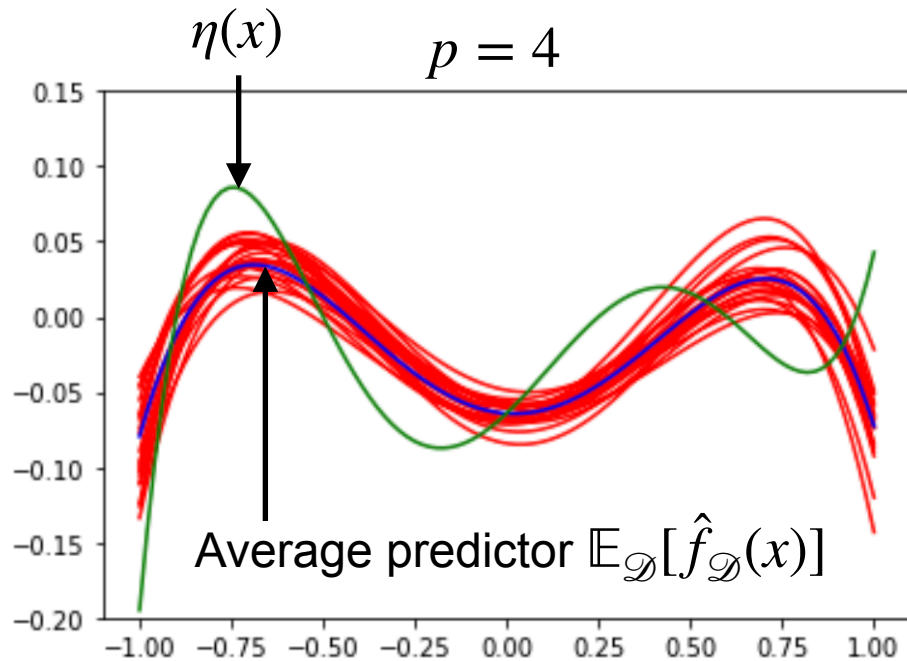
# Recap: Bias-variance tradeoff with simple model

(Conceptual) bias variance tradeoff



$\eta(x)$

$p = 20$

Average predictor $\mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)]$
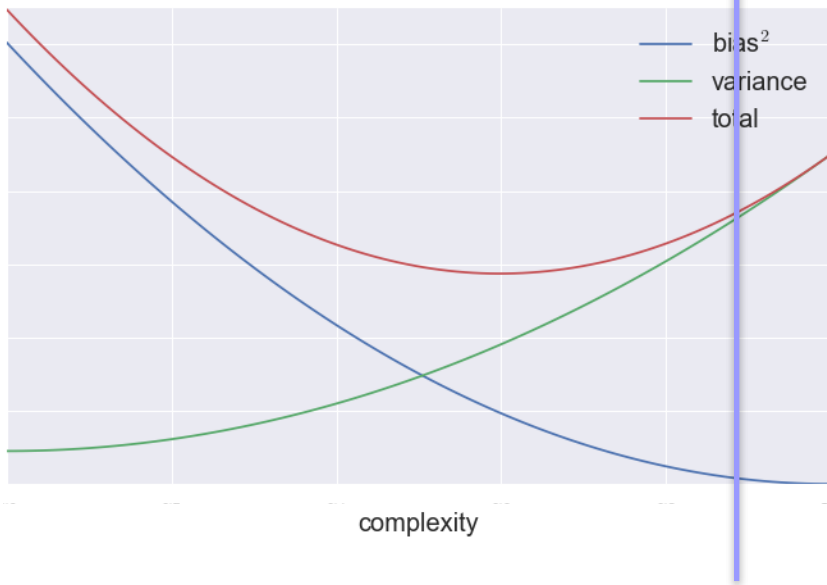
- When model complexity is high (higher than the optimal predictor $\eta(x)$)
  - Bias of our predictor, $\left( \eta(x) - \mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)] \right)^2$, is small
  - Variance of our predictor, $\mathbb{E}_{\mathscr{D}}\left[ \left( \mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)] - \hat{f}_{\mathscr{D}}(x) \right)^2 \right]$, is large
  - If we have more samples, then
    - Bias
    - Variance
    - Because Variance is dominating, overall test error

- let us first fix sample size **N=30**, collect one dataset of size N i.i.d. from a distribution, and fix one training set $S_{\text{train}}$ and test set $S_{\text{test}}$ via 80/20 split

- then we run multiple validations and plot the computed MSEs for all values of **p** that we are interested in

true model complexity



error

**Test error** $\mathscr{L}_{\text{test}}$

**UNDERFIT**          **OVERFIT**

**Training error** $\mathscr{L}_{\text{train}}$

$p_{N=24}^* \simeq 24 - 1$

**Model complexity ( = degree of the polynomial)**

- Given sample size N there is a threshold, $p_N^*$, where training error is zero
- Training error is **always** monotonically non-increasing
- Test error has a trend of going down and then up, but fluctuates

- let us now repeat the process changing the sample size to **N=40** , and see how the curves change



true model complexity

Model complexity ( = degree of the polynomial)

- The threshold, $p_N^*$, moves right
- Training error tends to increase, because more points need to fit
- Test error tends to decrease, because Variance decreases

- let us now fix predictor model complexity **$p=30$**, collect multiple datasets by starting with 3 samples and adding one sample at a time to the training set, but keeping a large enough test set fixed

- then we plot the computed MSEs for all values of train sample size ***Ntrain*** that we are interested in



OVERFIT   UNDERFIT

Test error $\mathscr{L}_{\text{test}}$

Training error $\mathscr{L}_{\text{train}}$

$N_p^* = p + 1 = 31$

train sample size $N_{\text{train}}$

- There is a threshold, $N_p^*$, below which training error is zero (extreme overfit)

- Below this threshold, test error is meaningless, as we are overfitting and there are multiple predictors with zero training error some of which have very large test error

- Test error tends to decrease

- Training error tends to increase

lecture2_polynomialfit.ipynb

# Bias-variance tradeoff for linear models

If $Y_i = X_i^T w^* + \epsilon_i$ and $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$

$$\mathbf{y} = \mathbf{X}w^* + \epsilon$$

$$\widehat{w}_{\text{MLE}} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y} =$$

$$=$$

$$\eta(x) = \mathbb{E}_{Y|X}[Y \mid X = x] =$$

$$\hat{f}_{\mathcal{D}}(x) = x^T\widehat{w}_{\text{MLE}} =$$

# Bias-variance tradeoff for linear models

If $Y_i = X_i^T w^* + \epsilon_i$ and $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$

$$\mathbf{y} = \mathbf{X} w^* + \epsilon$$

$$\widehat{w}_{\text{MLE}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X} w^* + \epsilon)$$

$$= w^* + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon$$

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x] = x^T w^*$$

$$\hat{f}_{\mathcal{D}}(x) = x^T \widehat{w}_{\text{MLE}} = x^T w^* + x^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon$$

- Irreducible error: $\mathbb{E}_{X,Y}[(Y - \eta(x))^2 | X = x] =$
- Bias squared: $\left( \eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] \right)^2 =$
  (is independent of the sample size!)

# Bias-variance tradeoff for linear models

If $Y_i = X_i^T w^* + \epsilon_i$ and $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$

$$\widehat{w}_{\text{MLE}} = w^* + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon$$

$$\eta(x) = x^T w^*$$

$$\hat{f}_{\mathcal{D}}(x) = x^T w^* + x^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon$$

- Variance: $\mathbb{E}_{\mathcal{D}} \left[ \left( \hat{f}_{\mathcal{D}}(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] \right)^2 \right] =$

# Bias-variance tradeoff for linear models

If $Y_i = X_i^T w* + \epsilon_i$ and $\epsilon_i \sim \mathcal{N}(0,\sigma^2)$

$$\widehat{w}_{\text{MLE}} = w* + (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\epsilon$$

$$\eta(x) = x^T w*$$

$$\hat{f}_{\mathscr{D}}(x) = x^T w* + x^T(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\epsilon$$

- Variance: $\mathbb{E}_{\mathscr{D}}\left[\left(\hat{f}_{\mathscr{D}}(x) - \mathbb{E}_{\mathscr{D}}[\hat{f}_{\mathscr{D}}(x)]\right)^2\right] = \mathbb{E}_{\mathscr{D}}[x^T(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\epsilon\epsilon^T\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}x]$

$$= \sigma^2 \mathbb{E}_{\mathscr{D}}[x^T(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}x]$$

$$= \sigma^2 x^T\mathbb{E}_{\mathscr{D}}[(\mathbf{X}^T\mathbf{X})^{-1}]x$$

- To analyze this, let's assume that $X_i \sim \mathcal{N}(0,\mathbf{I})$ and number of samples, $n$, is large enough such that $\mathbf{X}^T\mathbf{X} = n\mathbf{I}$ with high probability and $\mathbb{E}[(\mathbf{X}^T\mathbf{X})^{-1}] \simeq \frac{1}{n}\mathbf{I}$, then

  - Variance is $\dfrac{\sigma^2 x^T x}{n}$, and decreases with increasing sample size $n$

# Regularization

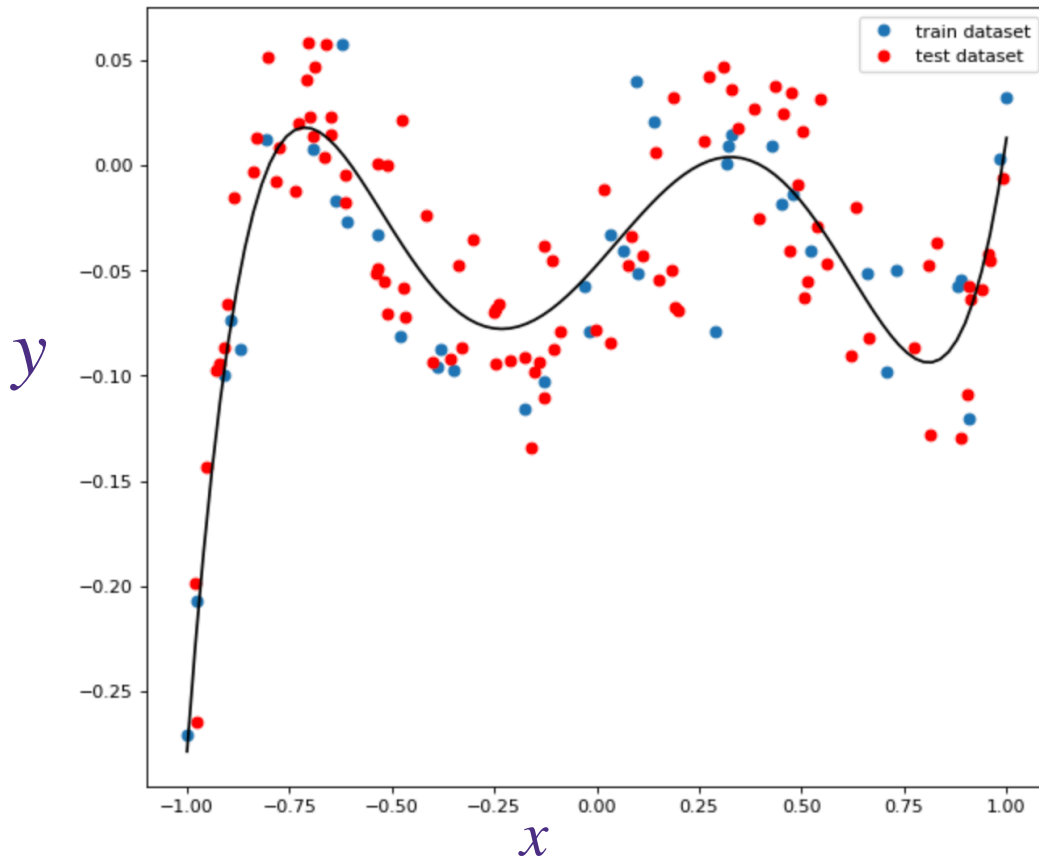# Recap: bias-variance tradeoff

- Consider 100 training examples and 100 test examples
  i.i.d.drawn from degree-5 polynomial features
  $x_i \sim \text{Uniform}[-1,1]$, $y_i \sim f_{w*}(x_i) + \epsilon_i$, $\epsilon_i \sim \mathcal{N}(0,\sigma^2)$

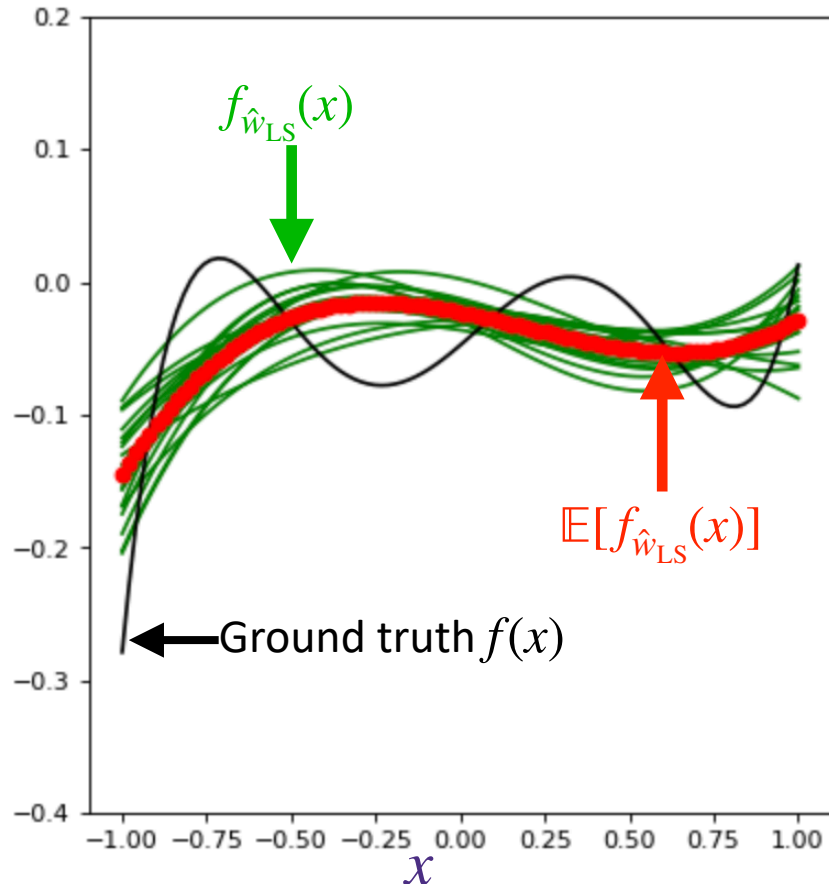$$f_w(x_i) = b* + w_1^*x_i + w_2^*(x_i)^2 + w_3^*(x_i)^3 + w_4^*(x_i)^4 + w_5^*(x_i)^5$$



This is a linear model with features
$$h(x_i) = (x_i, (x_i)^2, (x_i)^3, (x_i)^4, (x_i)^5)$$

# Recap: bias-variance tradeoff
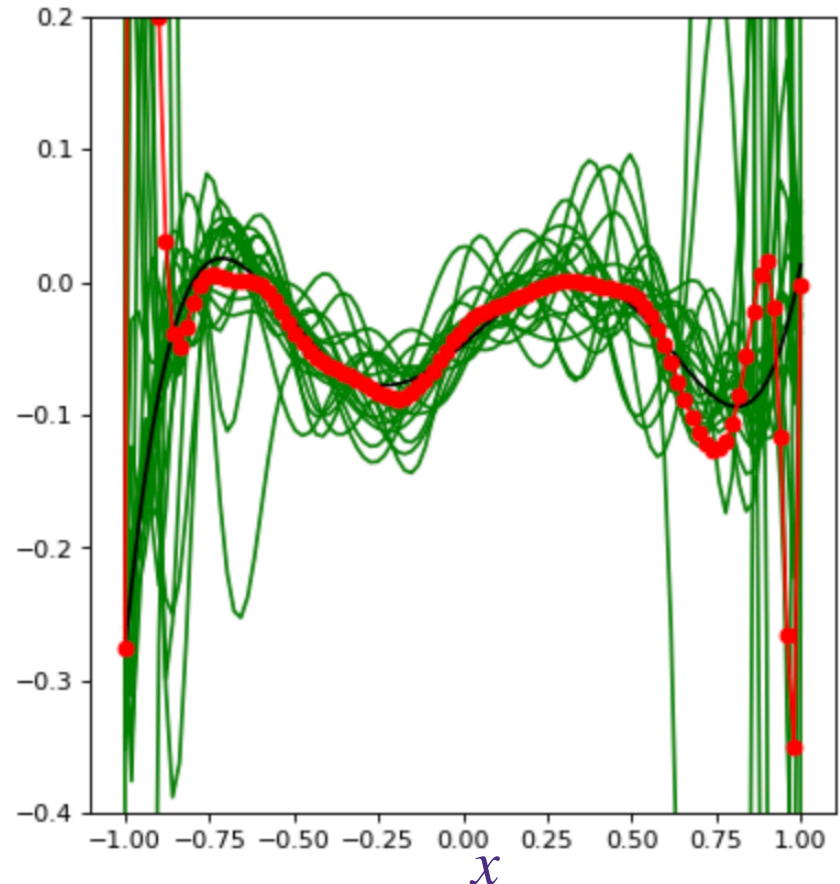
With degree-3 polynomials, we underfit

$\hat{f}_{\hat{w}_{\mathrm{LS}}}(x)$



$f_{\hat{w}_{\mathrm{LS}}}(x)$

$\mathbb{E}[f_{\hat{w}_{\mathrm{LS}}}(x)]$

Ground truth $f(x)$

```
current train error = 0.0036791644380554187
current test error  = 0.0037962529988410953
```

With degree-20 polynomials, we overfit

$\hat{f}_{\hat{w}_{\mathrm{LS}}}(x)$



```
0.0005421686349568773
0.14210029429557927
```
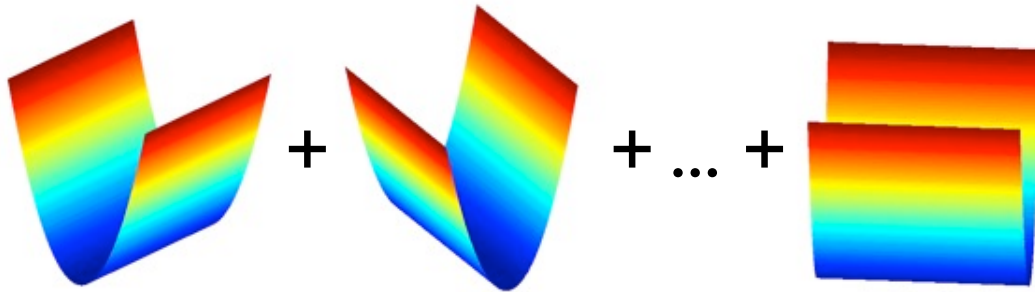
# Sensitivity: how to detect overfitting

- For a linear model,
$$y \simeq b + w_1 x_1 + w_2 x_2 + \cdots + w_d x_d$$
if $|w_j|$ is large then the prediction is sensitive to small changes in $x_j$

- Large sensitivity leads to overfitting and poor generalization, and equivalently models that overfit tend to have large weights

- Note that $b$ is a constant and hence there is no sensitivity for the offset $b$

- In **Ridge Regression,** we use a regularizer $\|w\|_2^2$ to measure and control the sensitivity of the predictor

- And optimize for small loss and small sensitivity, by adding a **regularizer** in the objective (assume no offset for now)

$$\widehat{w}_{ridge} = \arg \min_w \sum_{i=1}^{n} \left( y_i - x_i^T w \right)^2 + \lambda \|w\|_2^2$$
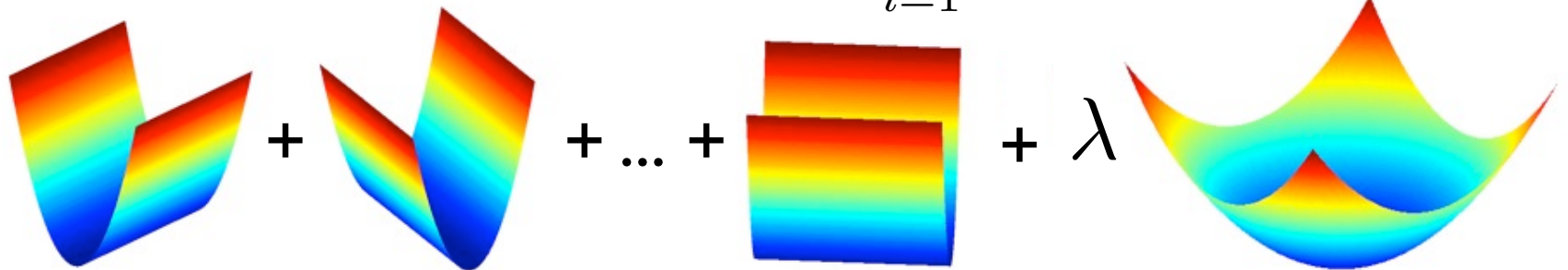
# Ridge Regression

- (Original) Least squares objective:

$$\widehat{w}_{LS} = \arg\min_{w} \sum_{i=1}^{n} \left( y_i - x_i^T w \right)^2$$



$$f_1(\mathbf{w}) \; + \; f_2(\mathbf{w}) + \ldots + \; f_T(\mathbf{w}) = \sum_{t=1}^{T} f_t(\mathbf{w})$$

- Ridge Regression objective:

$$\widehat{w}_{ridge} = \arg\min_{w} \sum_{i=1}^{n} \left( y_i - x_i^T w \right)^2 + \lambda ||w||_2^2$$



$+ \lambda$

$T \quad T$

# Minimizing the Ridge Regression Objective

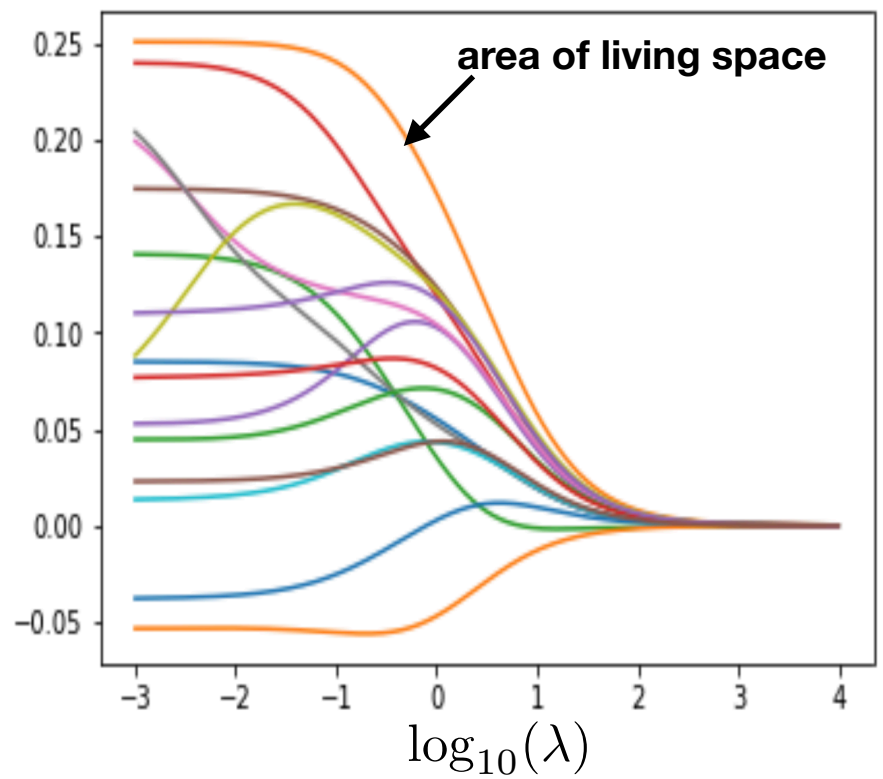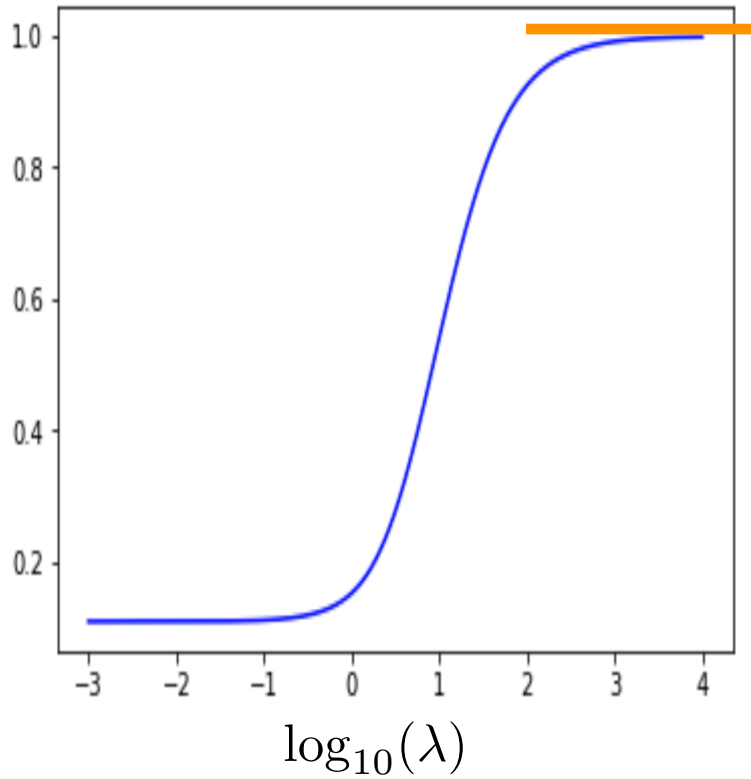$$\widehat{w}_{ridge} = \arg\min_{w} \sum_{i=1}^{n} \left(y_i - x_i^T w\right)^2 + \lambda \|w\|_2^2$$

# Shrinkage Properties

$$\widehat{w}_{ridge} = \arg\min_{w} \sum_{i=1}^{n} \left(y_i - x_i^T w\right)^2 + \lambda \|w\|_2^2$$

$$= (\mathbf{X}^T \mathbf{X} + \lambda I)^{-1} \mathbf{X}^T \mathbf{y}$$

- When $\lambda = 0$, this gives the least squares model
- This defines a family of models hyper-parametrized by $\lambda$
- Large $\lambda$ means more regularization and simpler model
- Small $\lambda$ means less regularization and more complex model

# Ridge regression: minimize $\sum_{i=1}^{n} (w^T x_i - y_i)^2 + \lambda \|w\|_2^2$

training MSE $\quad \dfrac{1}{n} \sum_{i=1}^{n} (y_i - x_i^T \hat{w}_{\text{ridge}}^{(\lambda)})^2$

$w_i$'s



$\log_{10}(\lambda)$

area of living space

$\log_{10}(\lambda)$
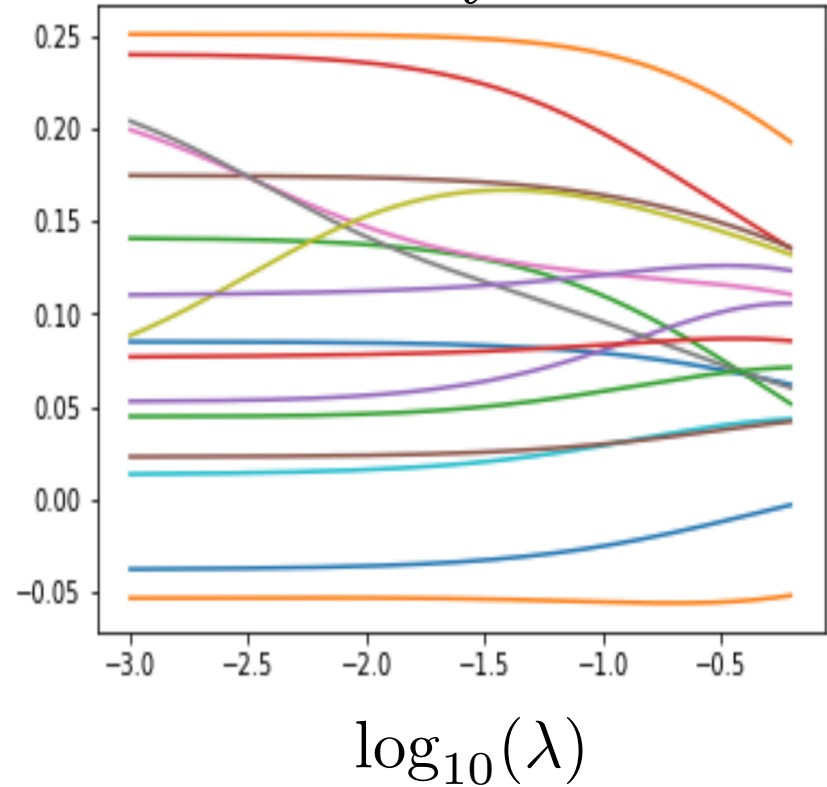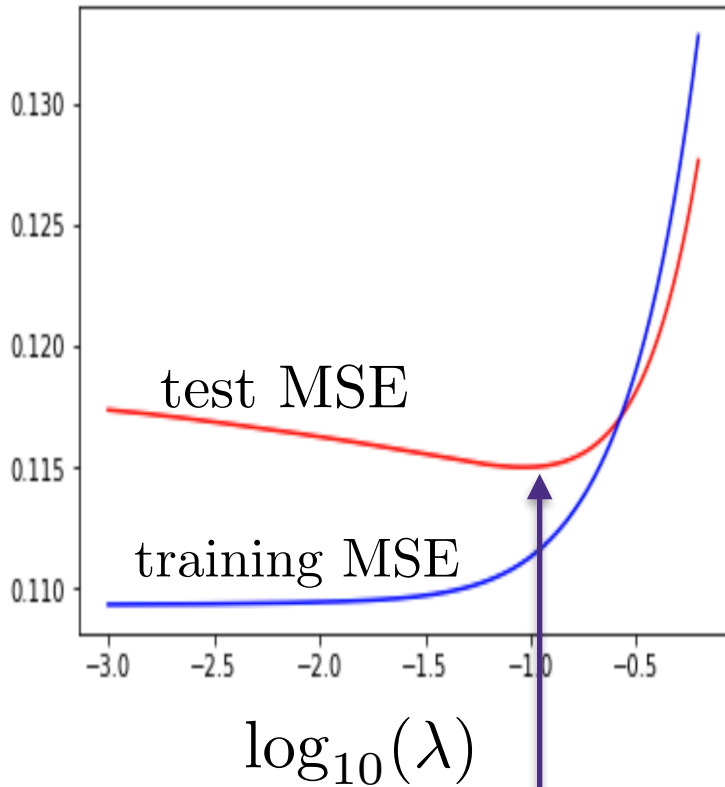
- Left plot: leftmost training error is with no regularization: 0.1093
- Left plot: rightmost training error is variance of the training data: 0.9991
- Right plot: called **regularization path**

# Ridge regression: $\text{minimize} \sum_{i=1}^{n} (w^T x_i - y_i)^2 + \lambda \|w\|_2^2$

$w_i\text{'s}$



$\log_{10}(\lambda)$

test MSE

training MSE

$\log_{10}(\lambda)$

- this gain in test MSE comes from shrinking w's to get a less sensitive predictor
(which in turn reduces the variance)

# Bias-Variance Properties

- Recall: $\hat{w}_{\text{ridge}} = (\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-1}\mathbf{X}^T\mathbf{y}$

- To analyze bias-variance tradeoff, we need to assume probabilistic generative model: $x_i \sim P_X, \ \mathbf{y} = \mathbf{X}w + \epsilon, \ \epsilon \sim \mathcal{N}(0, \sigma^2\mathbf{I})$

- The true error at a sample with feature $x$ is

$$\mathbb{E}_{y, \mathscr{D}_{\text{train}}|x}[(y - x^T\hat{w}_{\text{ridge}})^2 | x]$$

# Bias-Variance Properties

- Recall: $\hat{w}_{\text{ridge}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$

- To analyze bias-variance tradeoff, we need to assume probabilistic generative model: $x_i \sim P_X, \ \mathbf{y} = \mathbf{X}w + \epsilon, \ \epsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$

- The true error at a sample with feature $x$ is

$$\mathbb{E}_{y, \mathcal{D}_{\text{train}}|x}[(y - x^T \hat{w}_{\text{ridge}})^2 \,|\, x]$$

$$= \underbrace{\mathbb{E}_{y|x}[(y - \mathbb{E}[y|x])^2 \,|\, x]}_{\text{Irreducible Error}} + \underbrace{\mathbb{E}_{\mathcal{D}_{\text{train}}}[(\mathbb{E}[y|x] - x^T \hat{w}_{\text{ridge}})^2 \,|\, x]}_{\text{Learning Error}}$$