# 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 $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$

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

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

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

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

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

$$\mathbb{E}\left[ \left( x_i^T \hat{w} + \epsilon - x_i^T w^* \right)^2 \right]$$

$$= \mathbb{E}[\epsilon^2] = \sigma^2$$

- 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!)

$$\mathbb{E}\left[ \left( x^T (\mathbf{X}^T\mathbf{X})^{-1} \mathbf{X}^T \epsilon \right) \right]$$

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

$$a^2 = a\, a^T \text{ when } a \text{ is scalar}$$

$$\left(x^T(X^TX)^{-1}X^T\epsilon\right)^2$$

$$x(X^TX)^{-1}X^T\epsilon\left(x^T(X^TX)^{-1}X^T\epsilon\right)^T$$

$$x^T(X^TX)^{-1}X^T\epsilon\,\epsilon^T X(X^TX)^{-T}x$$

$$(exp\nearrow$$

$$x^T(X^TX)^{-1}X^T\,E(\epsilon\epsilon^T)X(X^TX)^{-T}x$$

$$\sigma^2 \cdot x^T(X^TX)^{-T}\cdot x$$

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

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

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

- 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 $n$
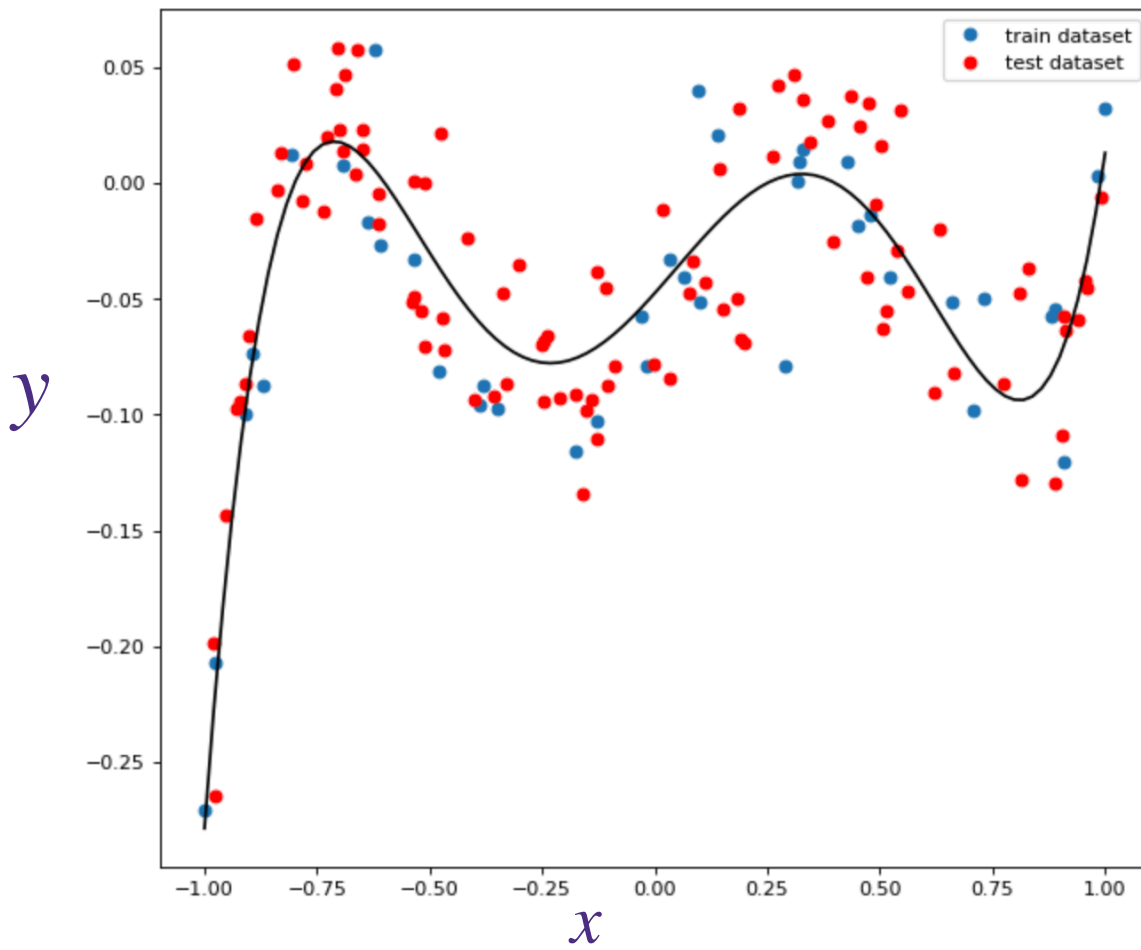
# Regularization

# Recap: bias-variance tradeoff

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

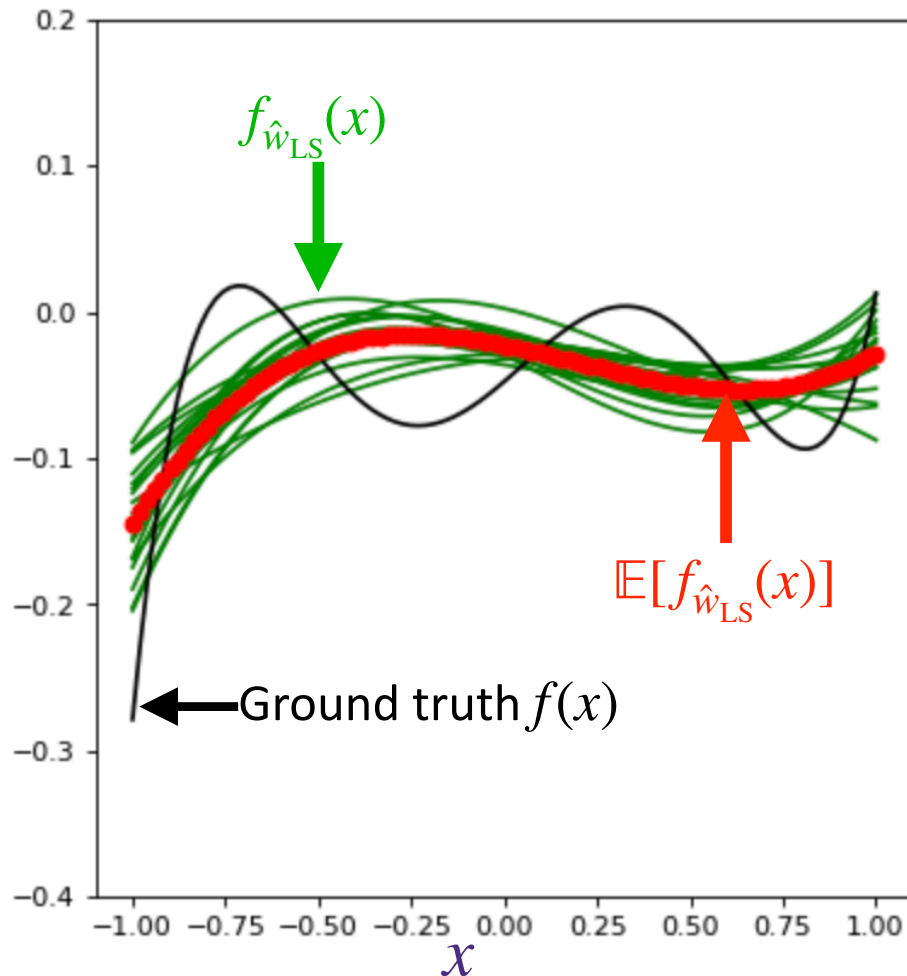$$f_w(x_i) = b* + w_1^* x_i + w_2^*(x_i)^2 + w_3^*(x_i)^3 + w_4^*(x_i)^4 + w_5^*(x_i)^5$$



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
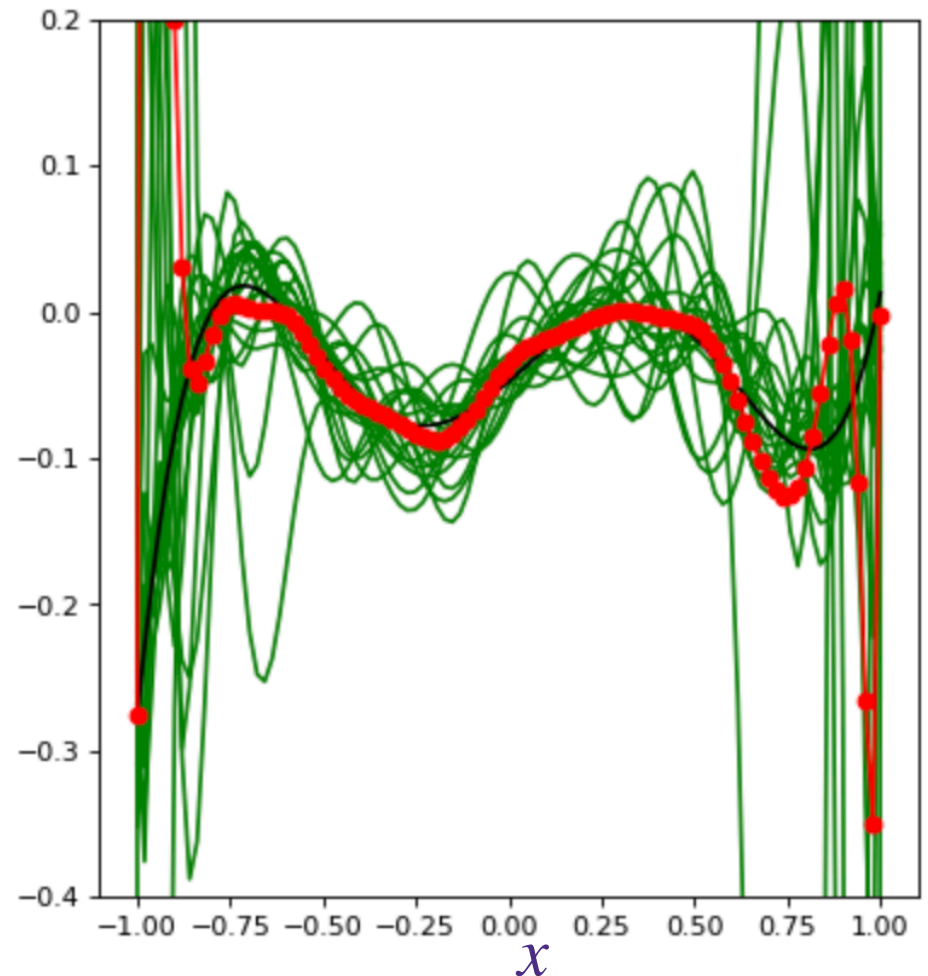
With degree-3 polynomials, we underfit

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



current train error = 0.0036791644380554187
current test error  = 0.0037962529988410953

With degree-20 polynomials, we overfit

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



0.0005421686349568773
0.14210029429557927

$$\mathbb{E}_D \left[ \left( \hat{f}_D(x) - \mathbb{E}_D \left[ \hat{f}_D(x) \right) \right)^2 \right]$$

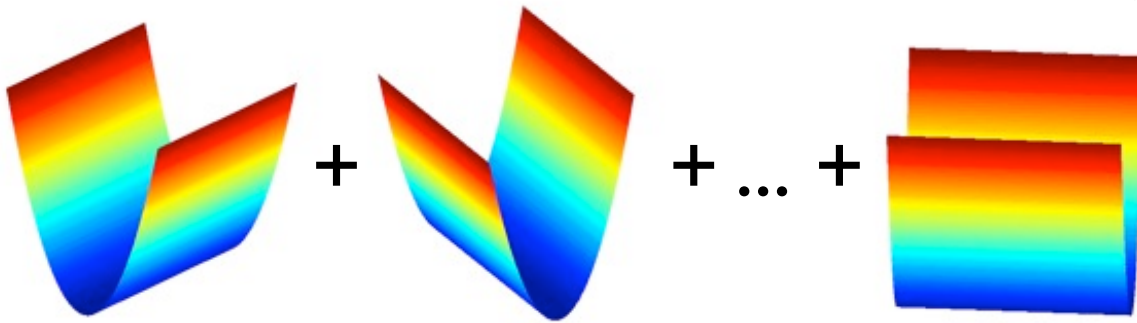# Sensitivity: how to detect overfitting

- 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)

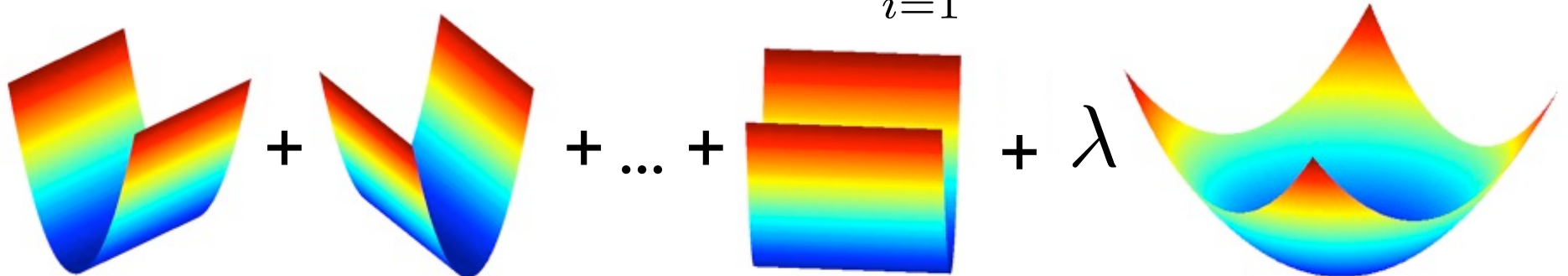$$\widehat{w}_{ridge} = \arg\min_{w} \sum_{i=1}^{n} \left( y_i - x_i^T w \right)^2 + \lambda \|w\|_2^2$$

# 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:

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



$$T \quad T$$

# Minimizing the Ridge Regression Objective

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

Matrix

$$\arg\min_{w} \left(Y - X^T w\right)^T \left(Y - X^T w\right) + \lambda \|w\|_2^2$$

$$\nabla_w \bullet = 2 X^T X w - 2 X^T y + 2\lambda w$$

$$X^T y = X^T X w + \lambda w$$

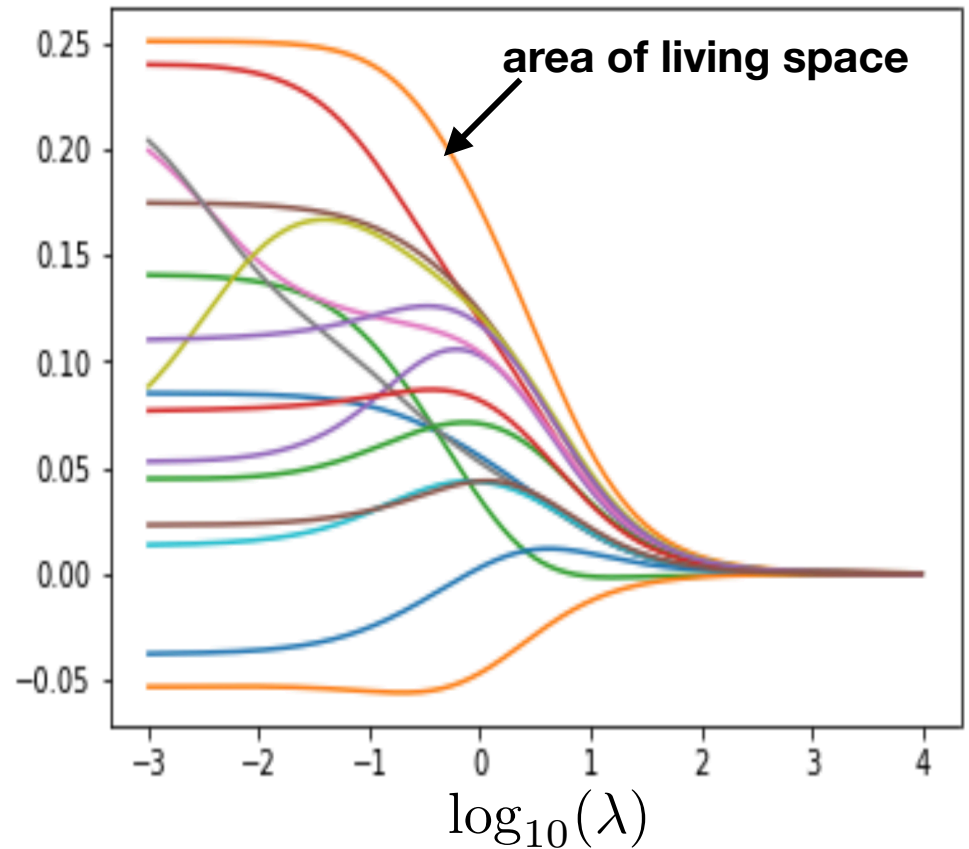$$w = \left(X^T X + \lambda I\right)^{-1} X^T y$$
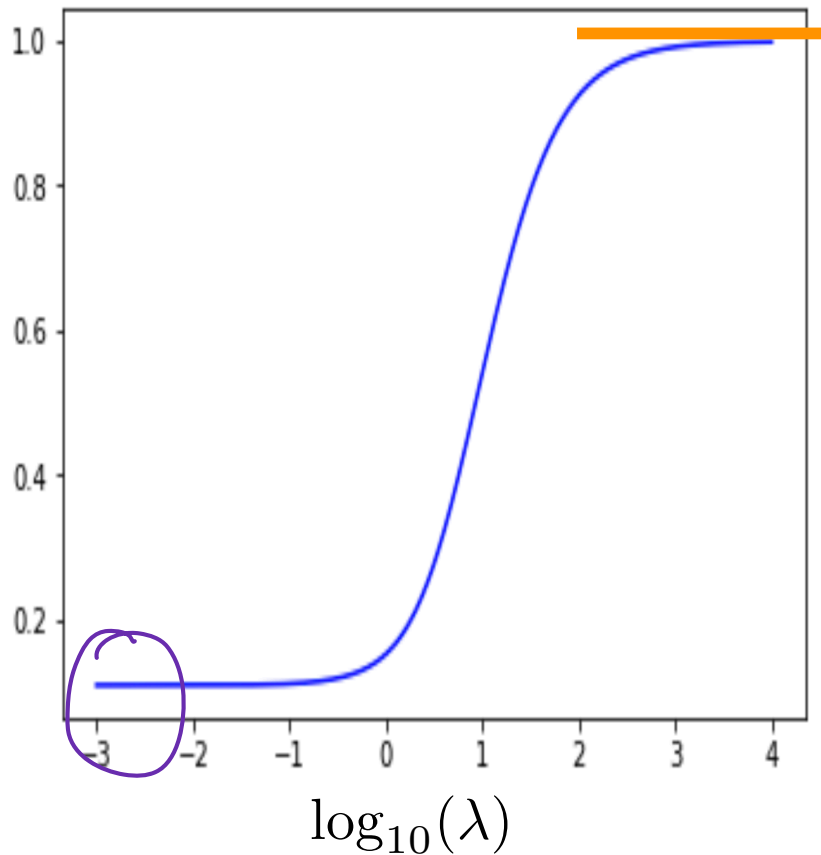
# Shrinkage Properties

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

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

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

training MSE $\quad \dfrac{1}{n} \sum_{i=1}^{n} (y_i - x_i^T \hat{w}_{\text{ridge}}^{(\lambda)})^2$
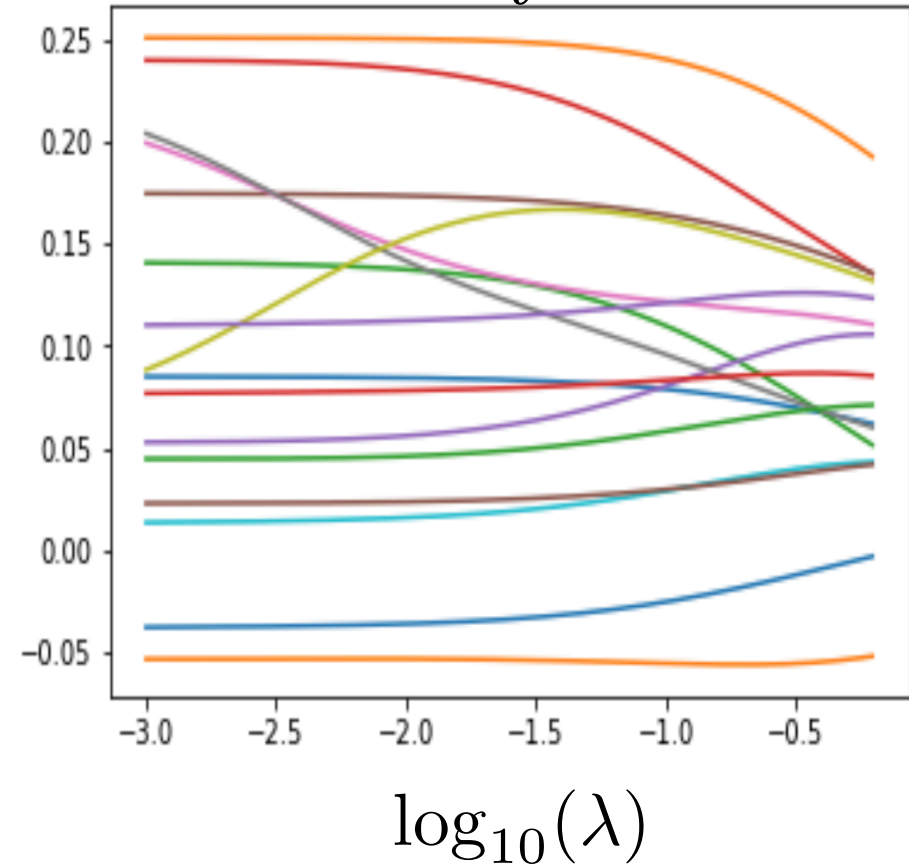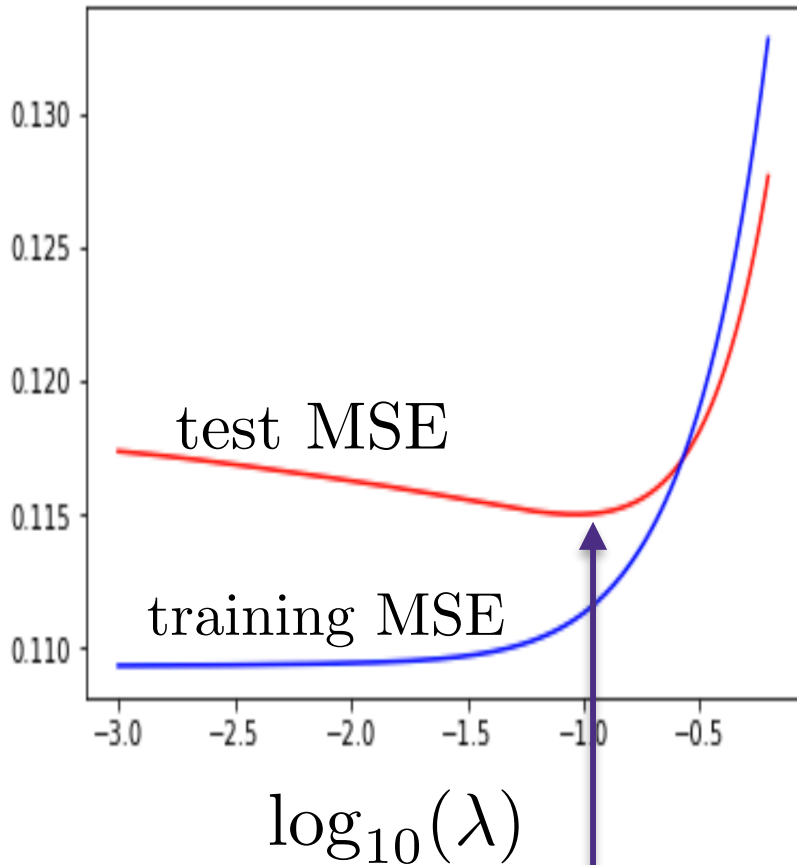
$w_i$'s



area of living space

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

$w_i$'s



- 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

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

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

# 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

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

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

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

# 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

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

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

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

$$= \underline{\sigma^2} + \underline{(x^Tw - \mathbb{E}_{\mathscr{D}_{\text{train}}}[x^T\hat{w}_{\text{ridge}} \,|\, x])^2} + \underline{\mathbb{E}_{\mathscr{D}_{\text{train}}}[(\mathbb{E}_{\tilde{\mathscr{D}}_{\text{train}}}[x^T\hat{w}_{\text{ridge}} \,|\, x] - x^T\hat{w}_{\text{ridge}})^2 \,|\, x]}$$

Irreduc. Error   Bias-squared        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

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

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

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

$$= \underline{\sigma^2} + \underline{(x^Tw - \mathbb{E}_{\mathscr{D}_{\text{train}}}[x^T\hat{w}_{\text{ridge}}|x])^2} + \underline{\mathbb{E}_{\mathscr{D}_{\text{train}}}[(\mathbb{E}_{\tilde{\mathscr{D}}_{\text{train}}}[x^T\hat{w}_{\text{ridge}}|x] - x^T\hat{w}_{\text{ridge}})^2 \,|\, x]}$$

<span style="color:green">Irreduc. Error</span>     <span style="color:blue">Bias-squared</span>          <span style="color:red">Variance</span>

---

Suppose $\mathbf{X}^T\mathbf{X} = n\mathbf{I}$, then $\hat{w}_{\text{ridge}} = (\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-1}\mathbf{X}^T(\mathbf{X}w + \epsilon)$

$$= \frac{n}{n + \lambda}w + \frac{1}{n + \lambda}\mathbf{X}^T\epsilon$$

# 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

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

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

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

$$= \sigma^2 + (x^Tw - \mathbb{E}_{\mathscr{D}_{\text{train}}}[x^T\hat{w}_{\text{ridge}} | x])^2 + \mathbb{E}_{\mathscr{D}_{\text{train}}}[(\mathbb{E}_{\tilde{\mathscr{D}}_{\text{train}}}[x^T\hat{w}_{\text{ridge}} | x] - x^T\hat{w}_{\text{ridge}})^2 | x]$$

(verify at home)

$$= \sigma^2 + \frac{\lambda^2}{(n+\lambda)^2}(w^Tx)^2 + \frac{\sigma^2 n}{(n+\lambda)^2}\|x\|_2^2$$
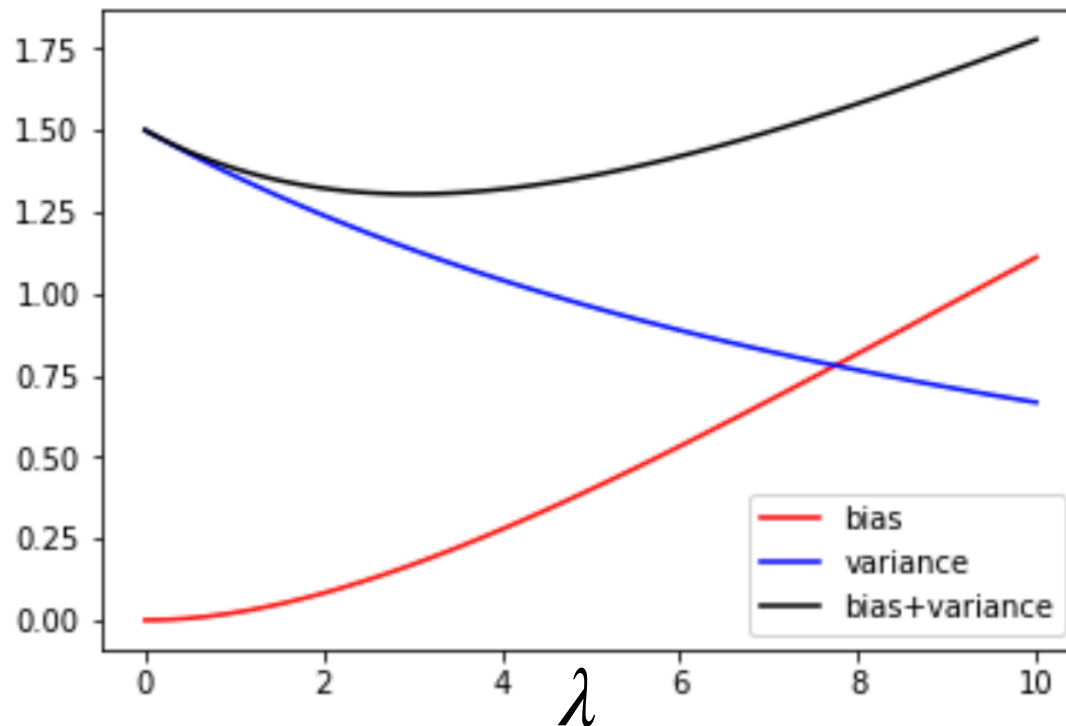
Irreduc. Error      Bias-squared      Variance

# Bias-Variance Properties

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

- True error

$$\mathbb{E}_{y,\mathscr{D}_{\text{train}}|x}[(y - x^T \hat{w}_{\text{ridge}})^2 \,|\, x] = \sigma^2 + \underbrace{\frac{\lambda^2}{(n+\lambda)^2}(w^T x)^2}_{\text{Bias-squared}} + \underbrace{\frac{\sigma^2 n}{(n+\lambda)^2}\|x\|_2^2}_{\text{Variance}}$$



d=10, n=20, $\sigma^2 = 3.0, \|w\|_2^2 = 10$

as $\lambda \to 0$,

$\hat{w}_{\text{ridge}} \to \hat{w}_{\text{LS}}$

as $\lambda \to \infty$

$\hat{w}_{\text{ridge}} \to 0$