

Introduction to Machine Learning

Regularization

Instructor: Pat Virtue

Announcements

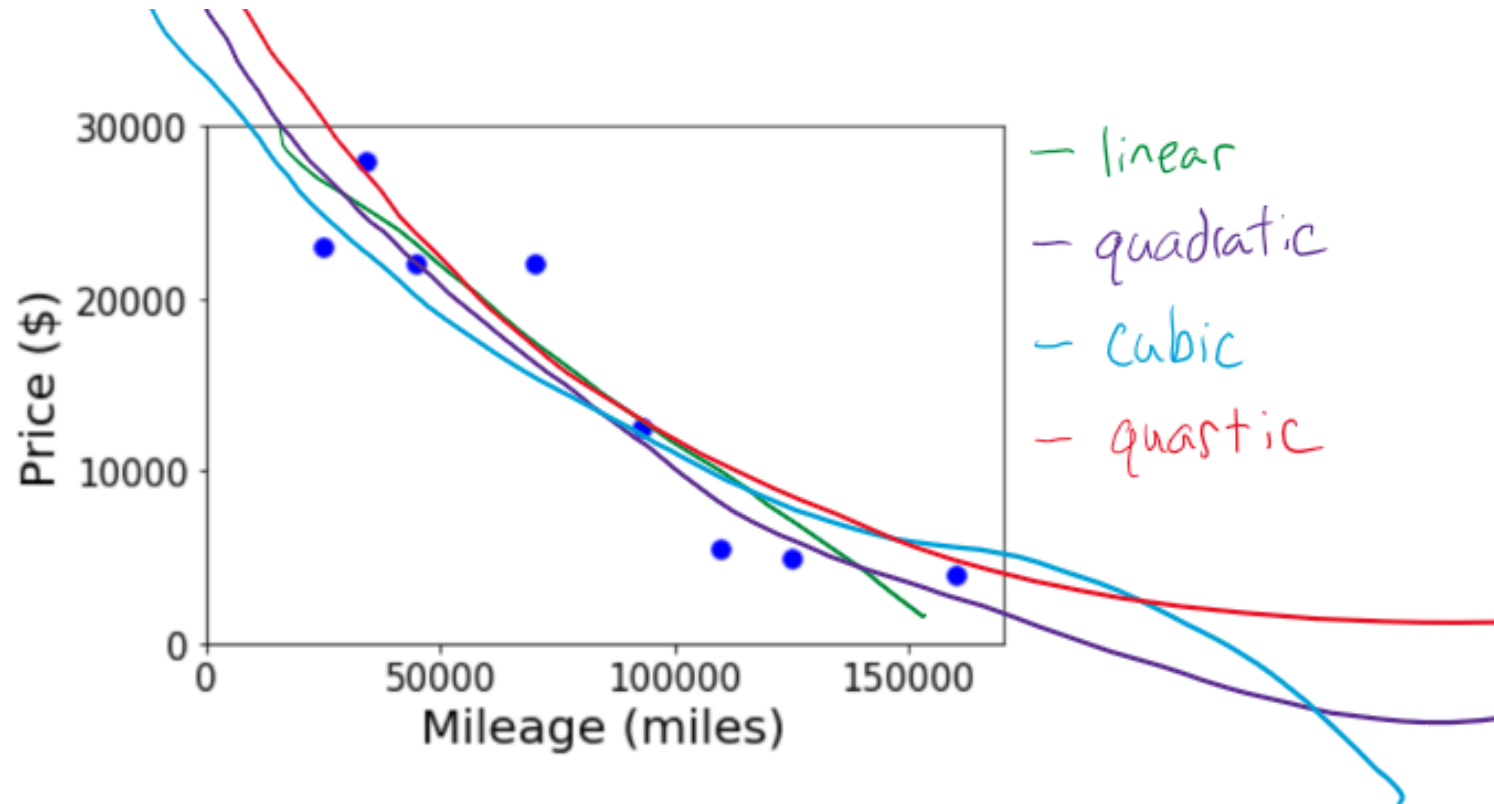
Assignments:

- HW3
 - Planned for release tonight
 - Due Tue, 2/11, 11:59 pm

Overfitting with Polynomial Linear Regression

Better fit training data with higher model complexity

$$w = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ w_3 \end{bmatrix}$$

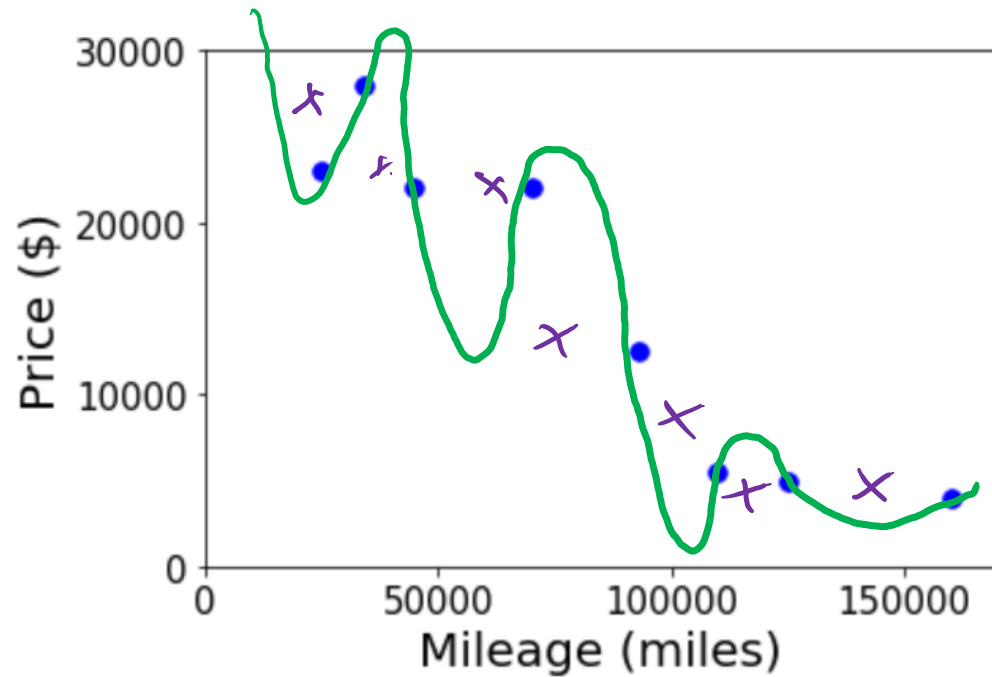


$$X = \begin{bmatrix} 1 & x^{(1)} & x^{(1)2} & x^3 \\ \vdots & x^{(2)} & x^2 & x^3 \\ \vdots & x & \vdots & \vdots \\ \vdots & x^{(n)} & x^2 & x^3 \end{bmatrix}$$

$$y = Xw$$

Overfitting with Polynomial Linear Regression

Better fit training data with higher model complexity



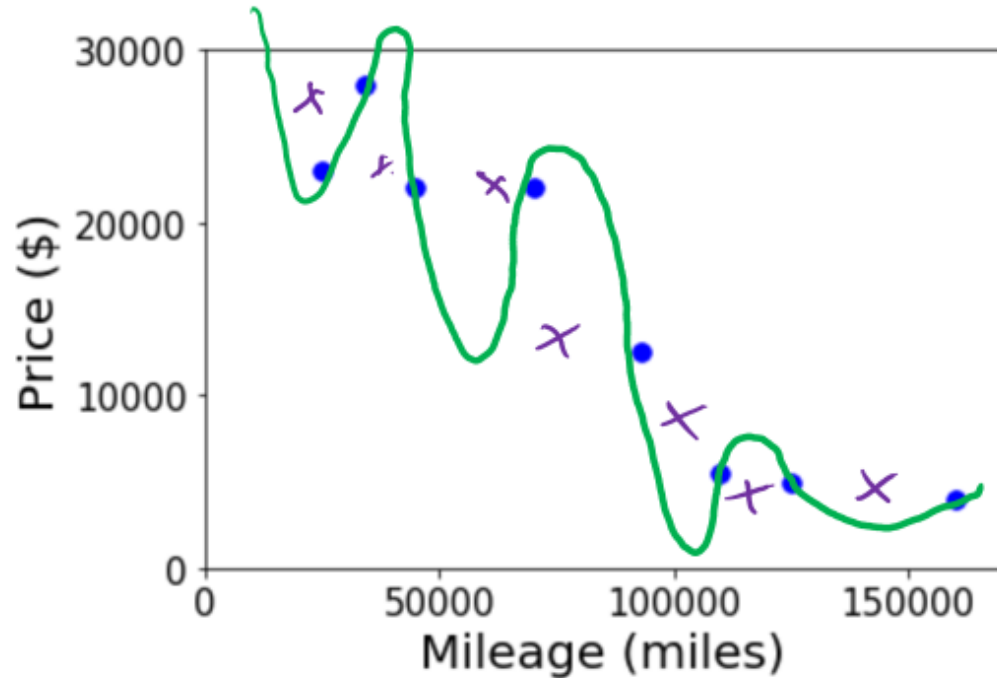
$$y = w_0 + w_1 x + w_2 x^2 + w_3 x^3 + \dots + w_{20} x^{20}$$

How can we deal with overfitting? Use validation. More training data

What are some symptoms of overfitting? Huge weights!

Overfitting with Polynomial Linear Regression

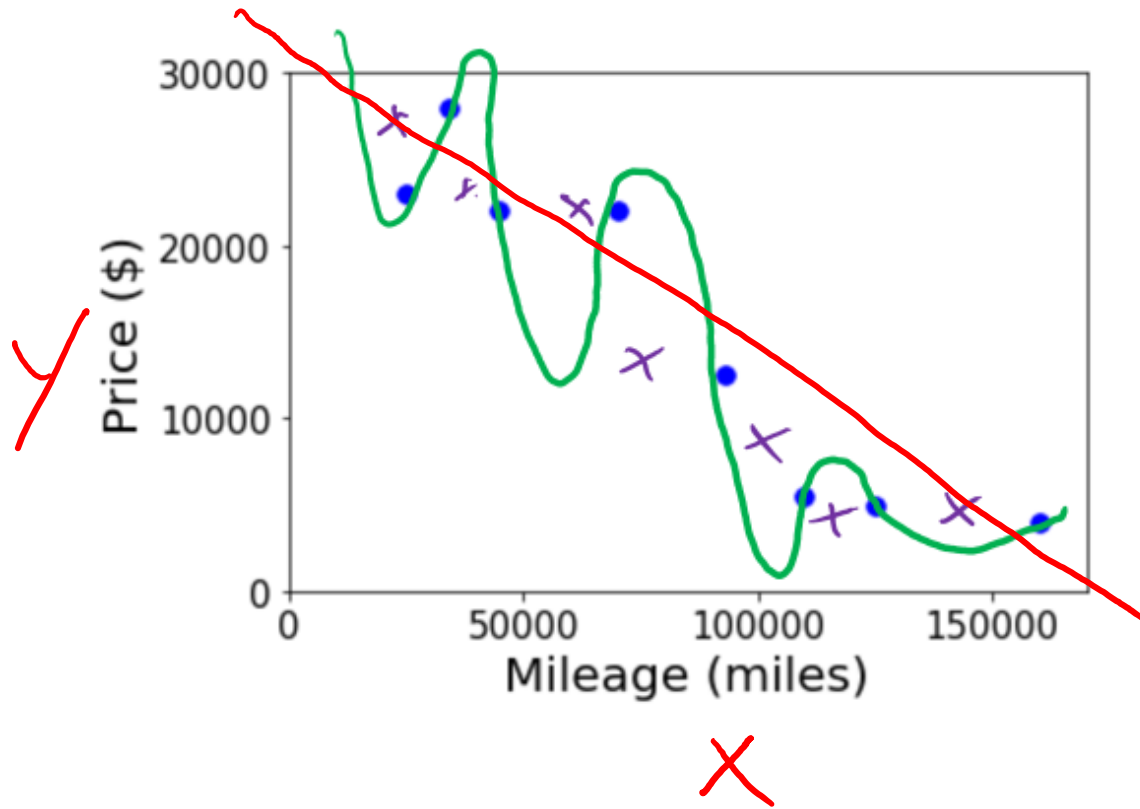
How can we deal with overfitting?



- Use validation set to detect overfitting
- Collect more training data
- Reduce model complexity
 - Lower degree polynomial
 - But then we might underfit ☹️
- Try fitting to many different degrees
 - Use validation data to decide which level of model complexity to use
- Penalize the weights

Overfitting with Polynomial Linear Regression

What are symptoms of overfitting?



- Poor validation score
- HUGE weights!

$$y = mx + b \quad \text{hypo}$$

$$J(m, b) = \|y - (\underline{m}x + \underline{b})\|_2^2$$

Regularization

Combine original objective with penalty on parameters

$$\min_w J(w) + \min_w \text{weight}$$

$$\min_w J(w) + f(w)$$

$$f_1(w) = \sum_m w_m \quad -9000 + 8900$$

$$f_2(w) = \sum_m |w_m|$$

$$f_3(w) = \sum_m w_m^2$$

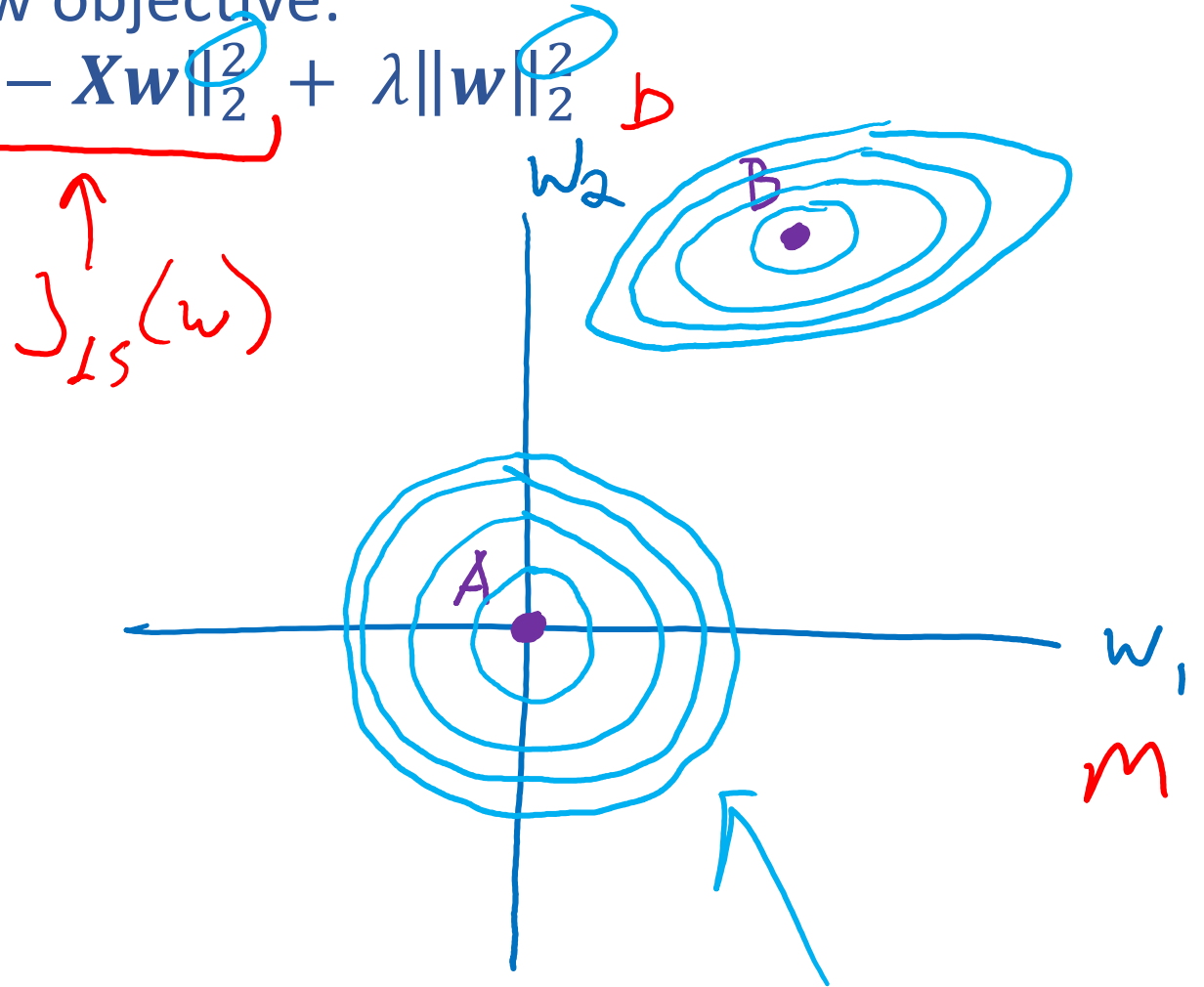
Piazza Poll 1:

Given the optimization of our new objective:

$$\hat{\mathbf{w}} = \min_{\mathbf{w}} \underbrace{\|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2}_{J_{LS}(\mathbf{w})} + \lambda \|\mathbf{w}\|_2^2$$

Select ALL that are true:

- I. As $\lambda \rightarrow 0$, $\hat{\mathbf{w}} \rightarrow$ point A
- II. As $\lambda \rightarrow 0$, $\hat{\mathbf{w}} \rightarrow$ point B
- III. As $\lambda \rightarrow \infty$, $\hat{\mathbf{w}} \rightarrow$ point A
- IV. As $\lambda \rightarrow \infty$, $\hat{\mathbf{w}} \rightarrow$ point B
- V. None of the above
- VI. I have no clue



$$Az + 3z \neq (A+3)z$$

Regularization

$$\|z\|_2^2 = z^T z$$

$$f(z) = z^T A z$$

$$\begin{aligned}\nabla_z f(z) &= 2Az \\ &= 2z^T A\end{aligned}$$

Ridge Regression: Linear regression with ℓ_2 penalty on weights

$$J(\vec{w}) = \left(\|\vec{y} - X\vec{w}\|_2^2 + \lambda \|\vec{w}\|_2^2 \right) \frac{1}{2}$$

$$= \left(\vec{y}^T \vec{y} - 2\vec{w}^T X^T \vec{y} + \vec{w}^T X^T X \vec{w} + \lambda \vec{w}^T \vec{w} \right) \frac{1}{2}$$

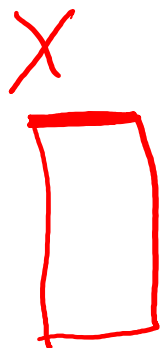
$$\nabla J(w) = -X^T \vec{y} + X^T X \vec{w} + \lambda w = 0$$

$$X^T X w + \lambda w = X^T \vec{y}$$

$$(X^T X + \lambda) w = X^T \vec{y}$$

$$(X^T X + \lambda I) w = X^T \vec{y}$$

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



Ridge Regression

Linear regression with ℓ_2 penalty on weights

$$\begin{aligned} J(\mathbf{w}) &= \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \frac{1}{2} \lambda \|\mathbf{w}\|_2^2 \\ &= \frac{1}{2} [\mathbf{y}^T \mathbf{y} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} + \lambda \mathbf{w}^T \mathbf{w}] \end{aligned}$$

Compute gradient

$$\nabla J(\mathbf{w}) = -\mathbf{X}^T \mathbf{y} + \mathbf{X}^T \mathbf{X} \mathbf{w} + \lambda \mathbf{w}$$

Closed form solution:

$$-\mathbf{X}^T \mathbf{y} + \mathbf{X}^T \mathbf{X} \mathbf{w} + \lambda \mathbf{w} = 0$$

$$\mathbf{X}^T \mathbf{X} \mathbf{w} + \lambda \mathbf{w} = \mathbf{X}^T \mathbf{y}$$

$$(\mathbf{X}^T \mathbf{X} + \lambda) \mathbf{w} = \mathbf{X}^T \mathbf{y} \quad \text{Not quite} \quad (\mathbf{A} + 7)\mathbf{z} \neq \mathbf{A}\mathbf{z} + 7\mathbf{z}$$

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

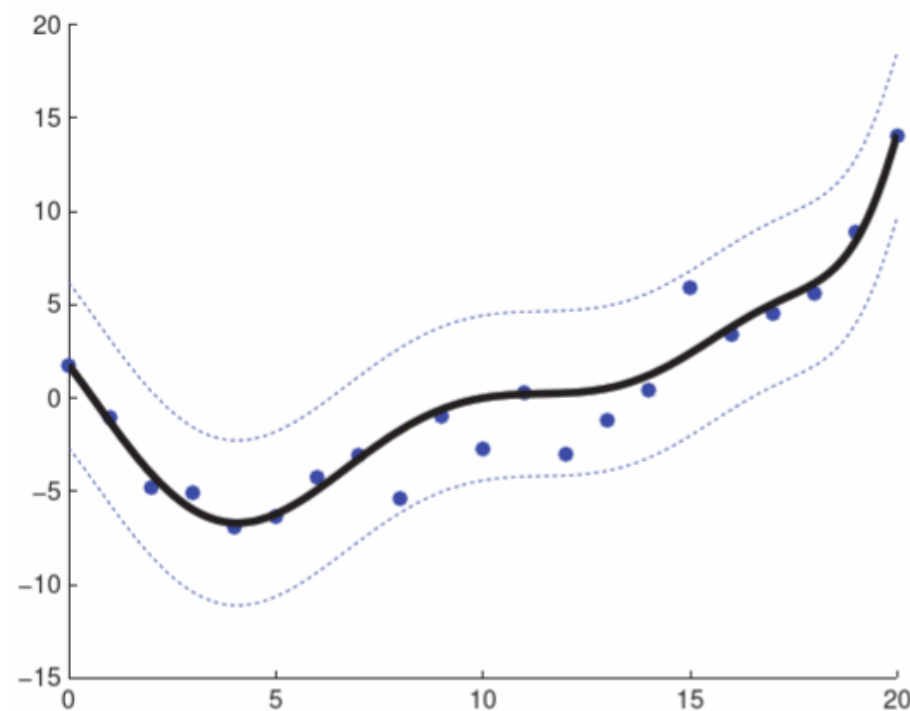
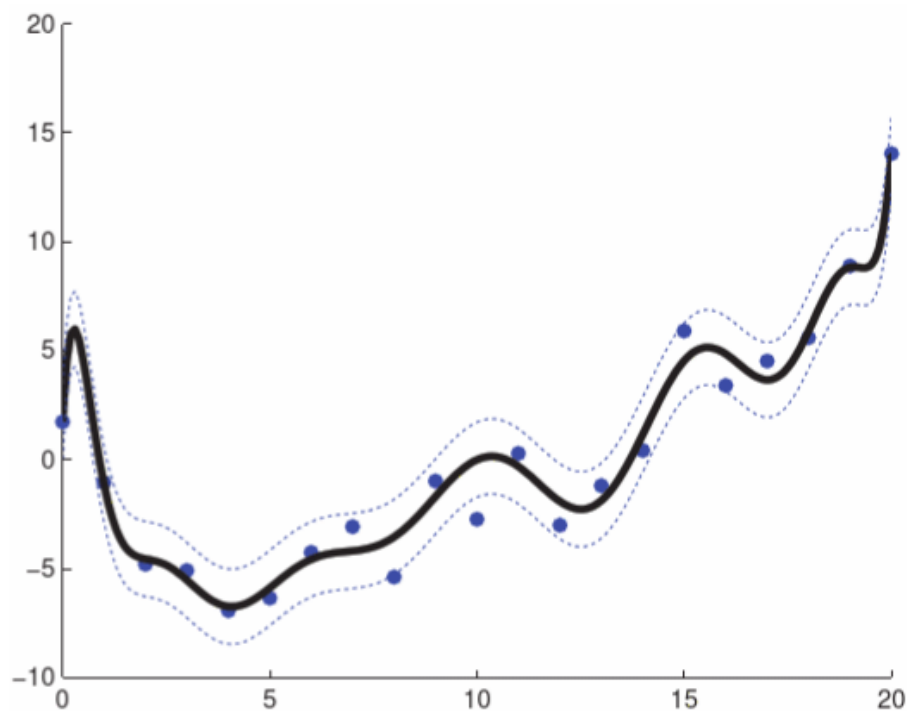
$$\mathbf{w} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

A robust solution to make $\mathbf{X}^T \mathbf{X}$ invertible

Regularization

But how do we choose λ ?

$$\|y - Xw\|_2^2 + \lambda \|w\|_2^2$$

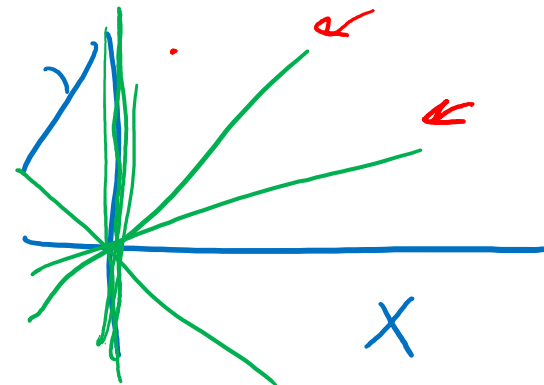


Probabilistic Interpretation

What assumptions are we making about our parameters?

likelihood

$$\rightarrow \underline{f(y | \vec{x}, w, \sigma^2)} = \frac{1}{\sqrt{2\pi\sigma^2}} e^{\left(\frac{-(y - \vec{x}^T w)^2}{2\sigma^2} \right)}$$



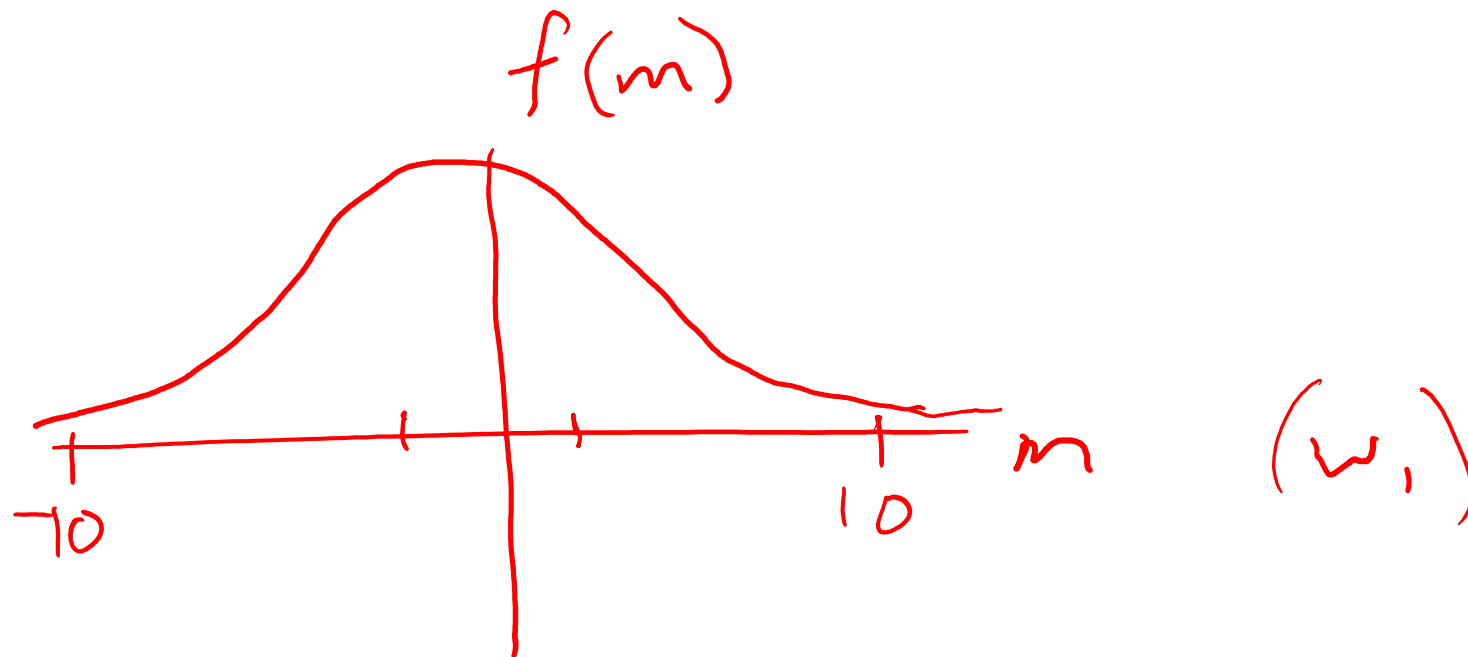
?

$$\rightarrow f(w | \mu_w, \sigma_w^2)$$

$p(\theta)$

$f(y | x, \theta)$

$p(D | \theta)$



MLE and MAP

MLE

$$p(D|\theta)$$

$$\arg \max_{\theta} p(D|\theta)$$

$$p(\theta|D)$$

posterior

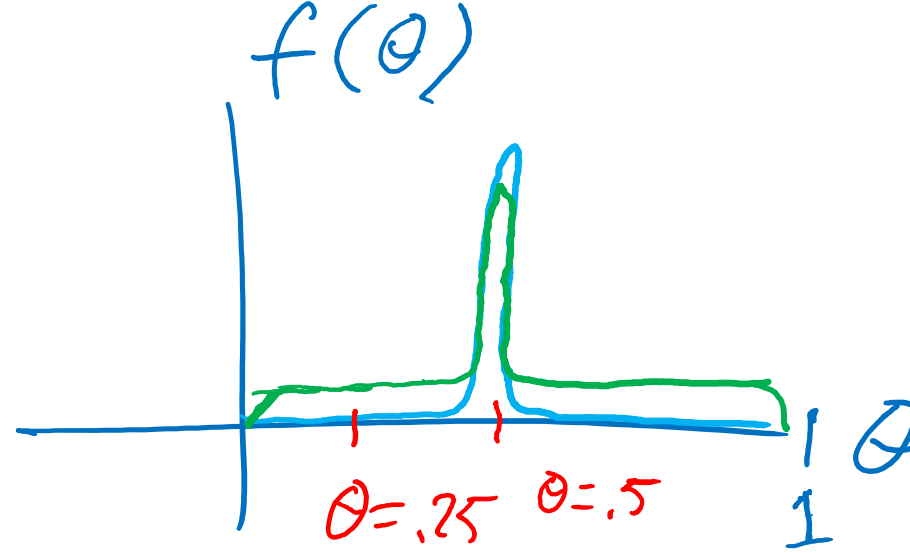
$$\propto p(D|\theta) p(\theta)$$

likelihood prior

Maximum a posteriori

$$\arg \max_{\theta} p(\theta|D) = \arg \max_{\theta} p(D|\theta) p(\theta)$$

Coin Flipping Example



$$f(\theta = 0.25) = 0.1$$

$$f(\theta = 0.5) = 0.8$$

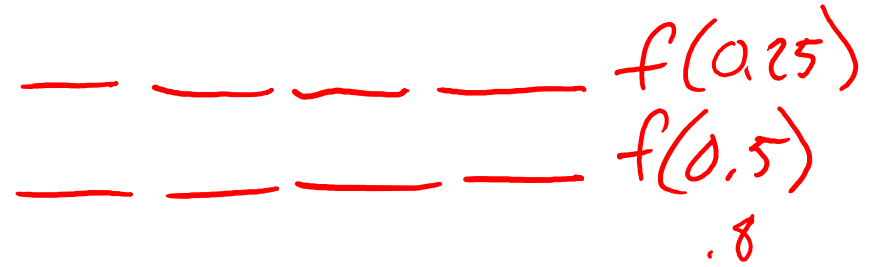
$$P(\theta|D) \propto \prod_{n=1}^N p(D^{(n)}|\theta) p(\theta)$$

Piazza Poll 2:

$$p(\theta | \mathcal{D}) \propto p(\mathcal{D}|\theta)p(\theta)$$

$$\underbrace{p(\theta | \mathcal{D})}_{\text{posterior}} \propto \underbrace{\prod p(\mathcal{D}^{(n)}|\theta)}_{\text{likelihood}} \underbrace{p(\theta)}_{\text{prior}}$$

$$\text{MLE} \quad \max P(\mathcal{D}|\theta)$$
$$\text{MAP} \quad \max P(\theta|\mathcal{D})$$



$\theta=0.25$

$\theta=0.5$

As the number of data points increases, which of the following are true?

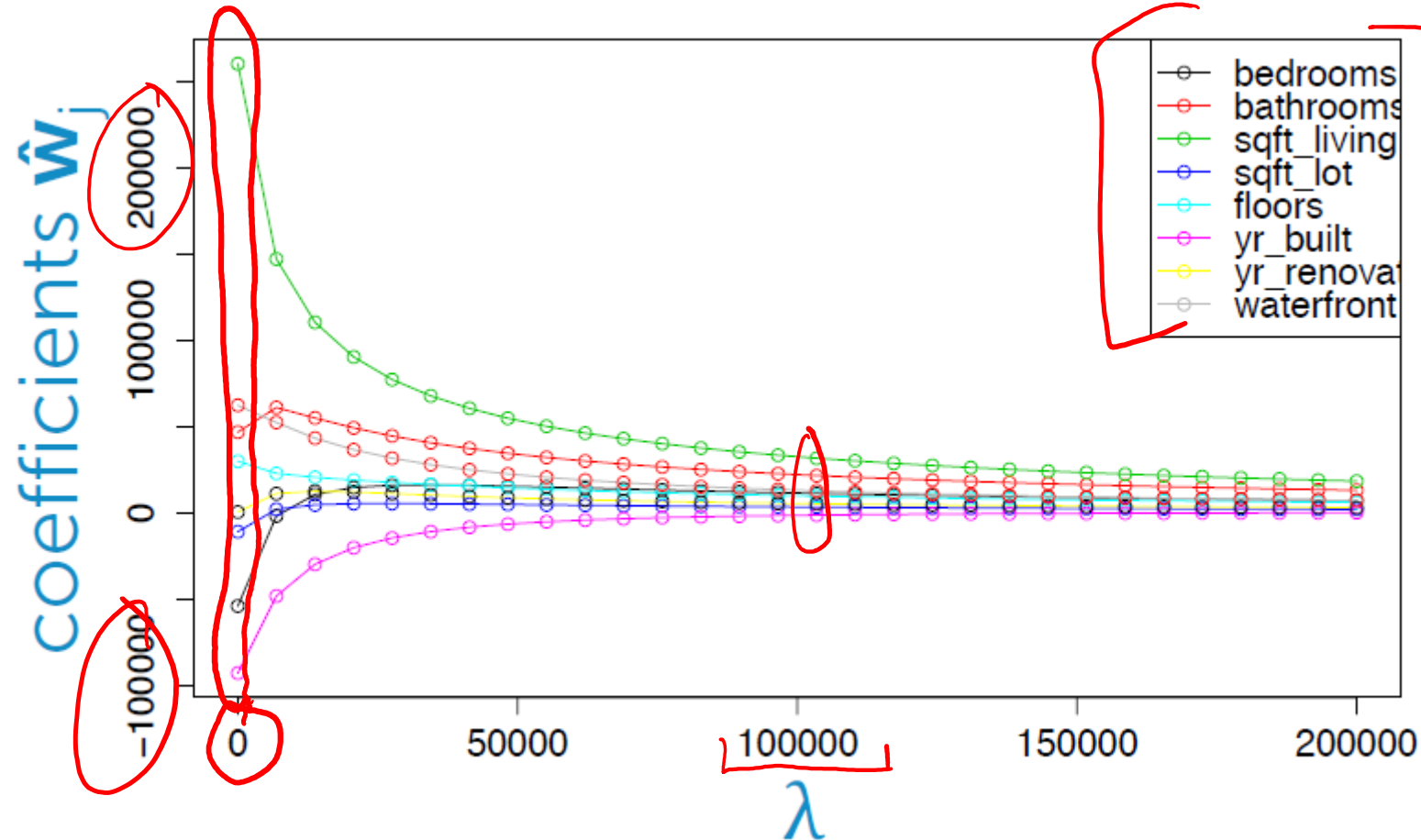
Select ALL that apply

- ☒ A. The MAP estimate approaches the MLE estimate
- ☐ B. The posterior distribution approaches the prior distribution
- ☐ C. The likelihood distribution approaches the prior distribution
- ☒ D. The posterior distribution approaches the likelihood distribution
- ☐ E. The likelihood has a lower impact on the posterior
- ☒ F. The prior has a lower impact on the posterior

Coin Flipping Example

Housing Price Example

Predict housing price from several features



$$w \in \mathbb{R}^m$$

$$x \in \mathbb{R}^{m=8}$$

Housing Price Example

Predict housing price from several features

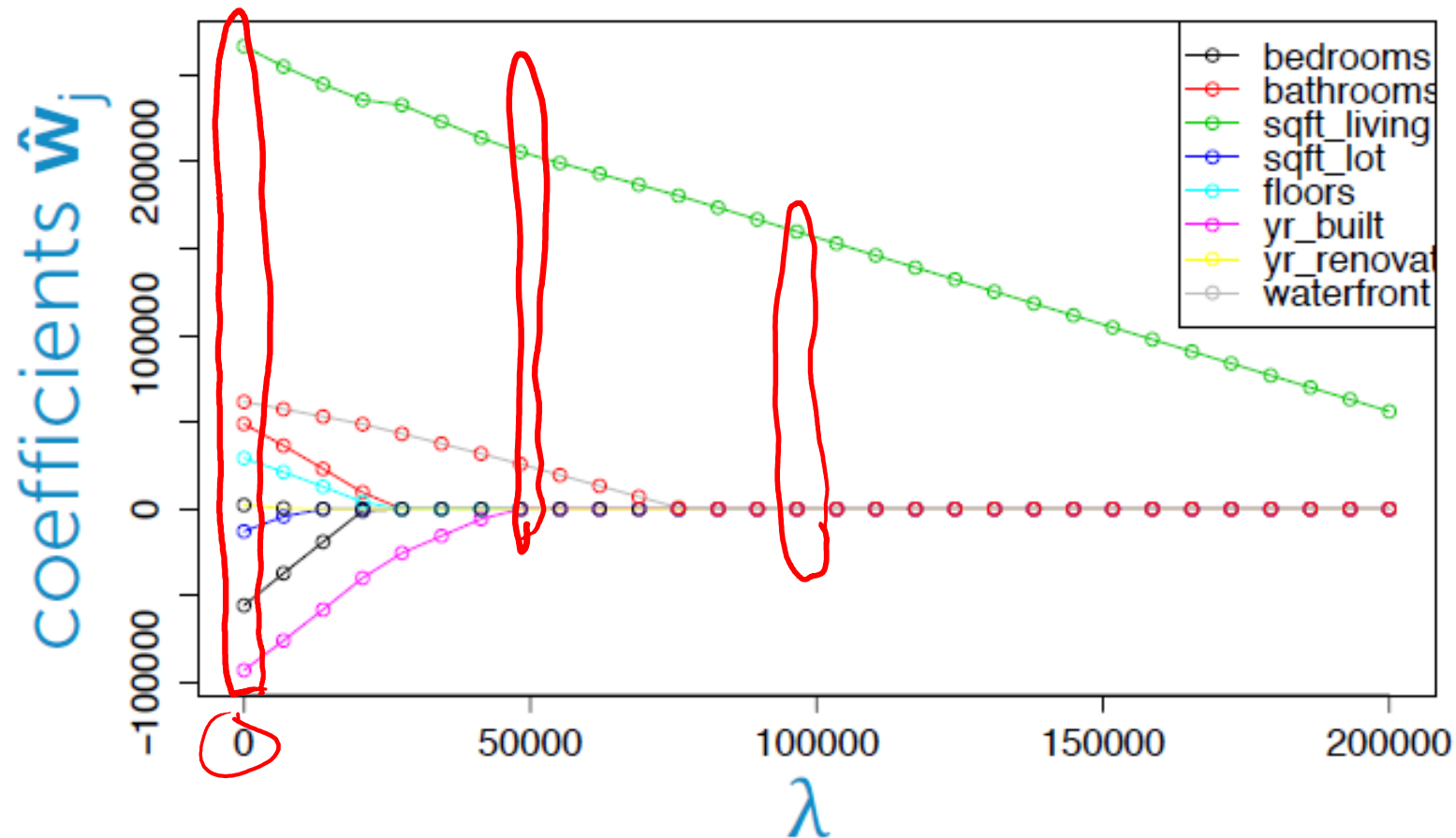


Figure: Emily Fox, University of Washington

Regularization

Combine original objective with penalty on parameters

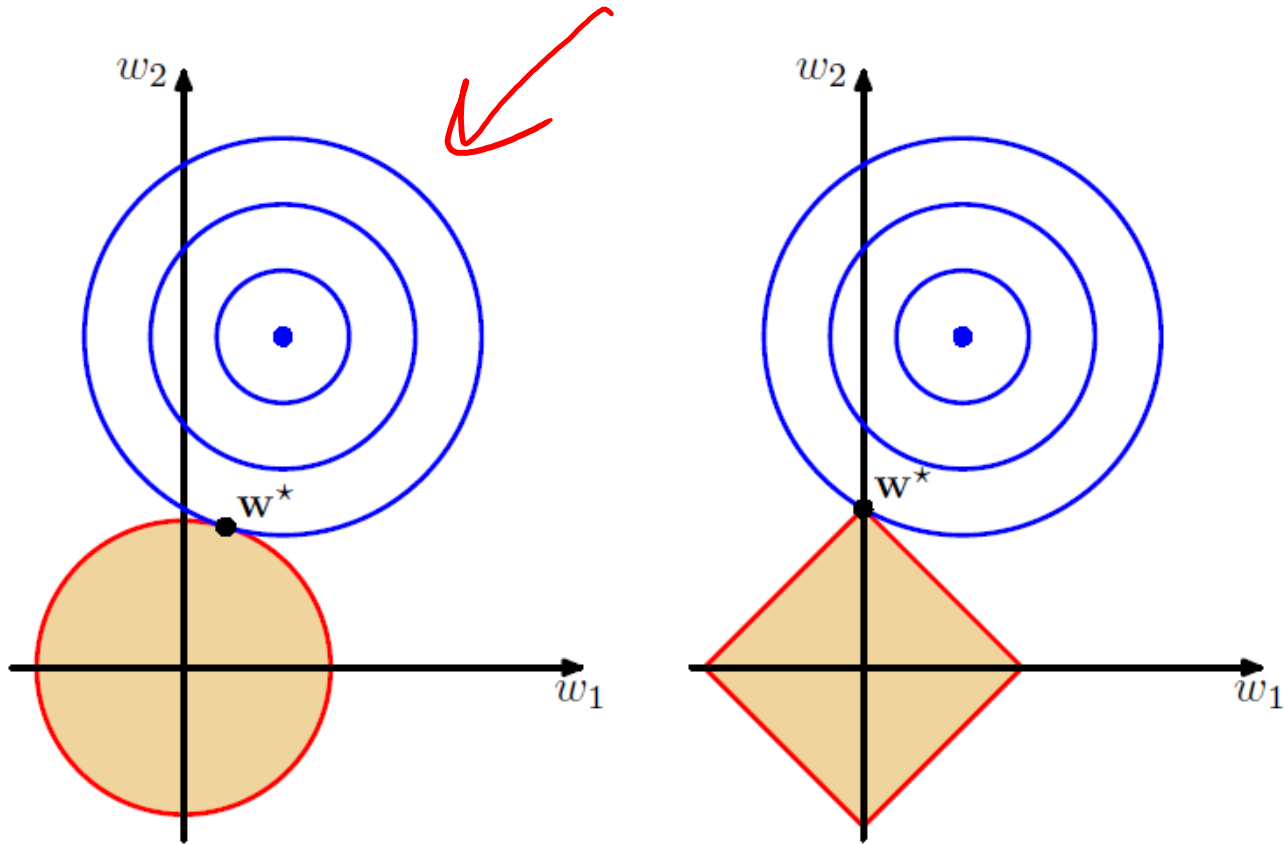
ℓ_2, ℓ_1, ℓ_0 norms

$$\ell_2 \quad \|w\|_2 = \sqrt{\sum w_i^2}$$

$$\ell_1 \quad \|w\|_1 = \sum |w_i|$$

Regularization

Combine original objective with penalty on parameters



LASSO

Linear regression with ℓ_1 penalty on weights

LASSO

Linear regression with ℓ_1 penalty on weights

$$\begin{aligned} J(\mathbf{w}) &= \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \frac{1}{2} \lambda \|\mathbf{w}\|_1 \\ &= \frac{1}{2} [\mathbf{y}^T \mathbf{y} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} + \sum_m |\mathbf{w}_m|] \end{aligned}$$

Probabilistic interpretation

Laplace prior on weights

$$w \sim \text{Laplace}(\mu = 0, b)$$

$$f(w \mid b) = \frac{1}{2b} e^{\left(-\frac{|x|}{b}\right)}$$