

Introduction to Machine Learning

Regularization

Instructor: Pat Virtue

Announcements

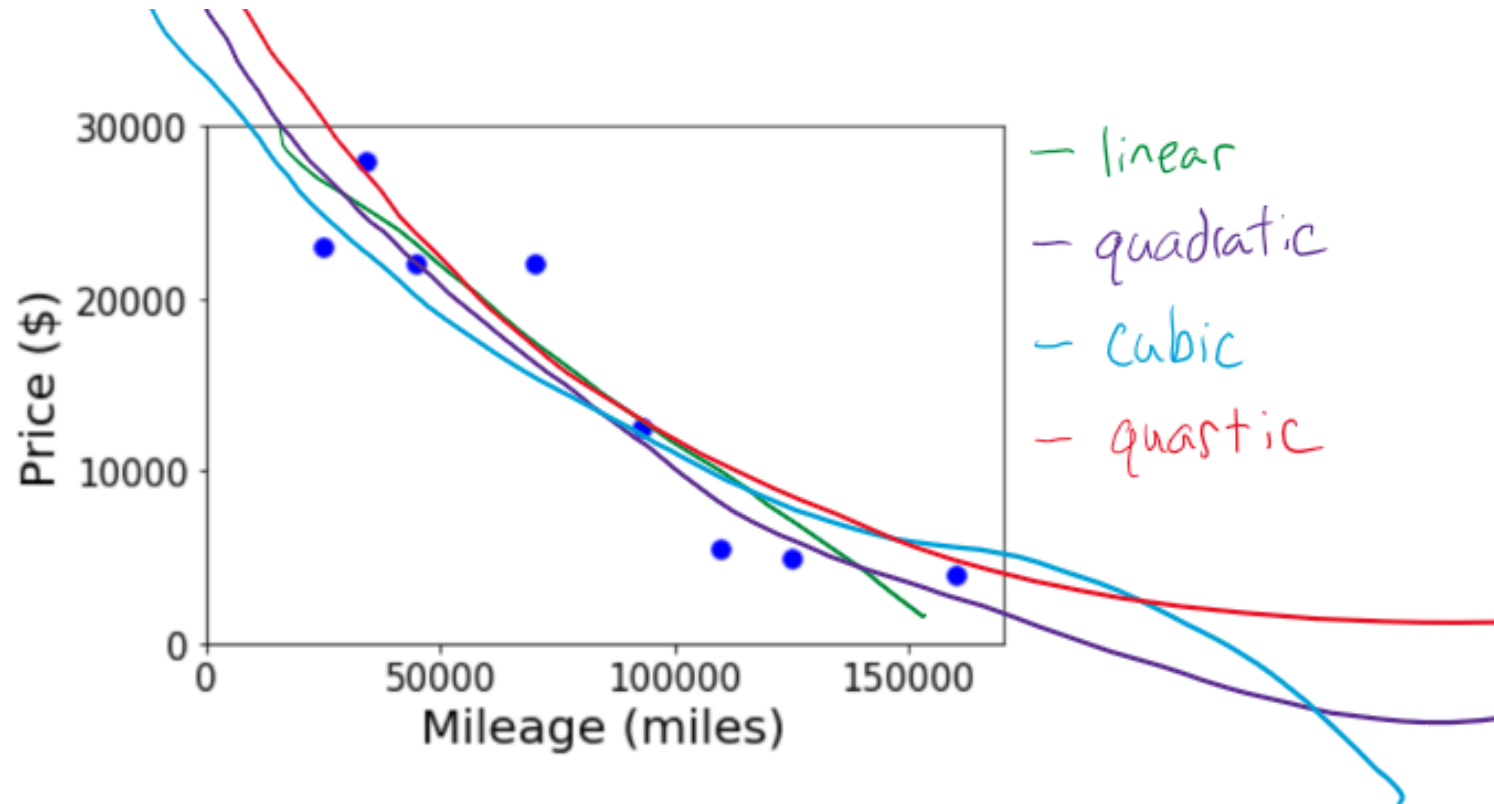
Assignments:

- 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

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

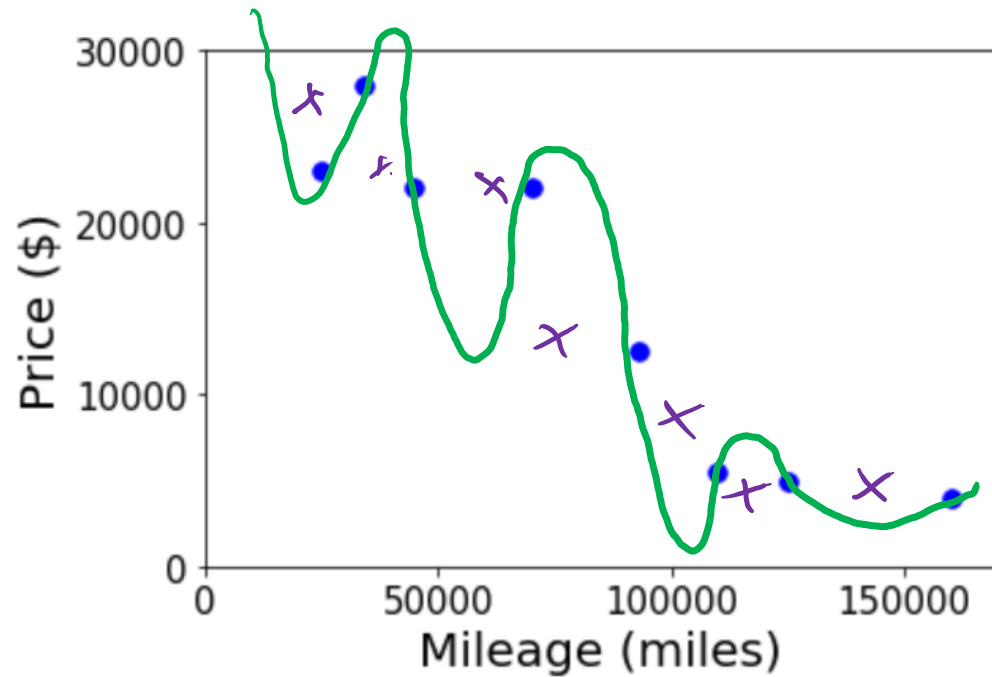


$$X = \begin{bmatrix} 1 & x^{(1)} & x^{(1)2} & x^3 \\ \vdots & x^{(2)} & x^2 & x^3 \\ \vdots & x & \vdots & \vdots \\ \vdots & x^{(n)} & x^2 & x^3 \end{bmatrix}$$

$$y = Xw$$

Overfitting with Polynomial Linear Regression

Better fit training data with higher model complexity



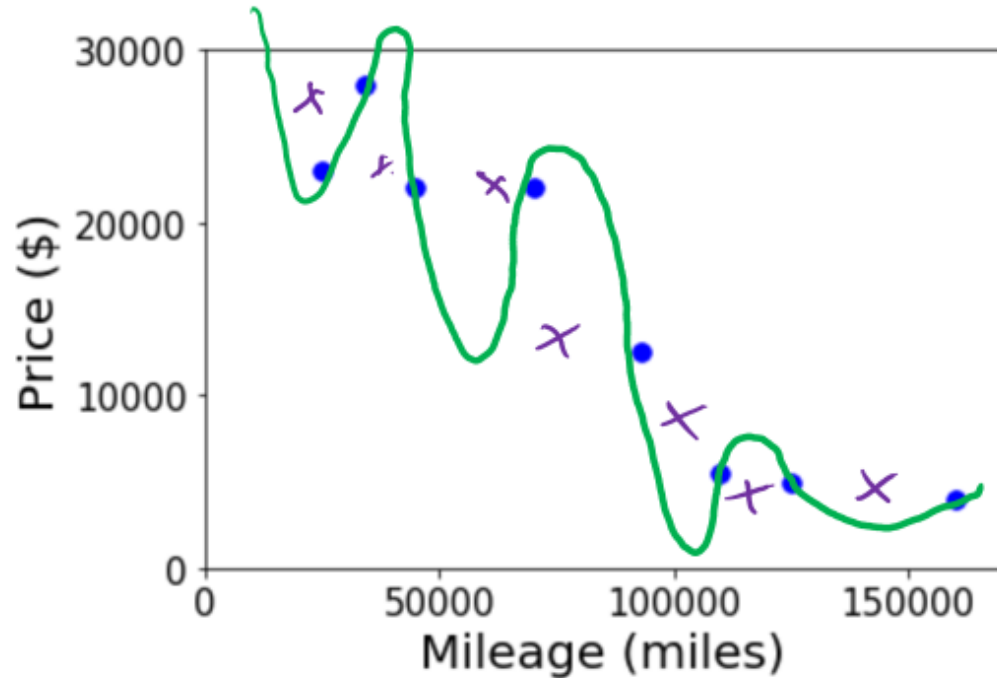
$$y = w_0 + w_1 x + w_2 x^2 + w_3 x^3 + \dots + w_{20} x^{20}$$

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?

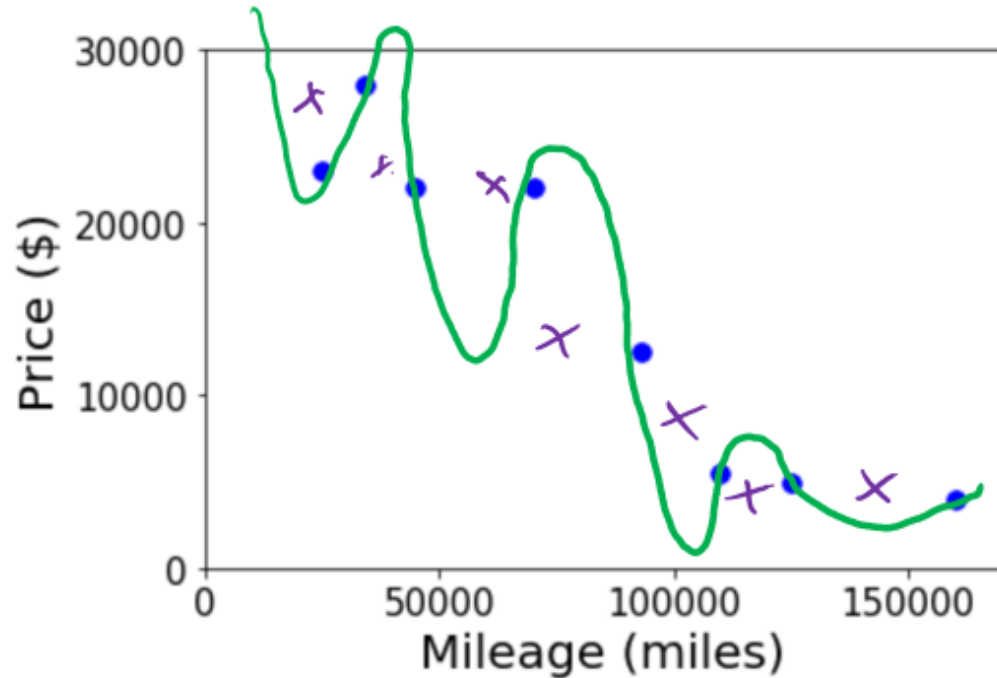


- Use validation set to detect overfitting
- Collect more training data
- Reduce model complexity
 - Lower degree polynomial
 - But then we might underfit ☹️
- Try fitting to many different degrees
 - Use validation data to decide which level of model complexity to use
- Penalize the weights

Overfitting with Polynomial Linear Regression

What are symptoms of overfitting?

- Poor validation score
- HUGE weights!



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:

$$\hat{\mathbf{w}} = \min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$$

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

Regularization

Ridge Regression: Linear regression with ℓ_2 penalty on weights

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:

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

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

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

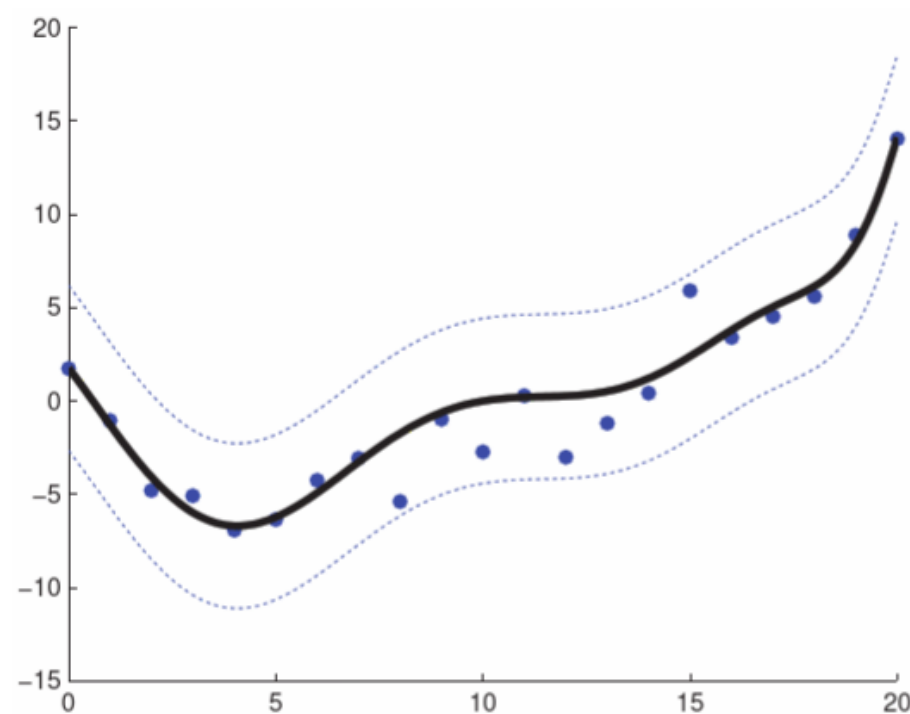
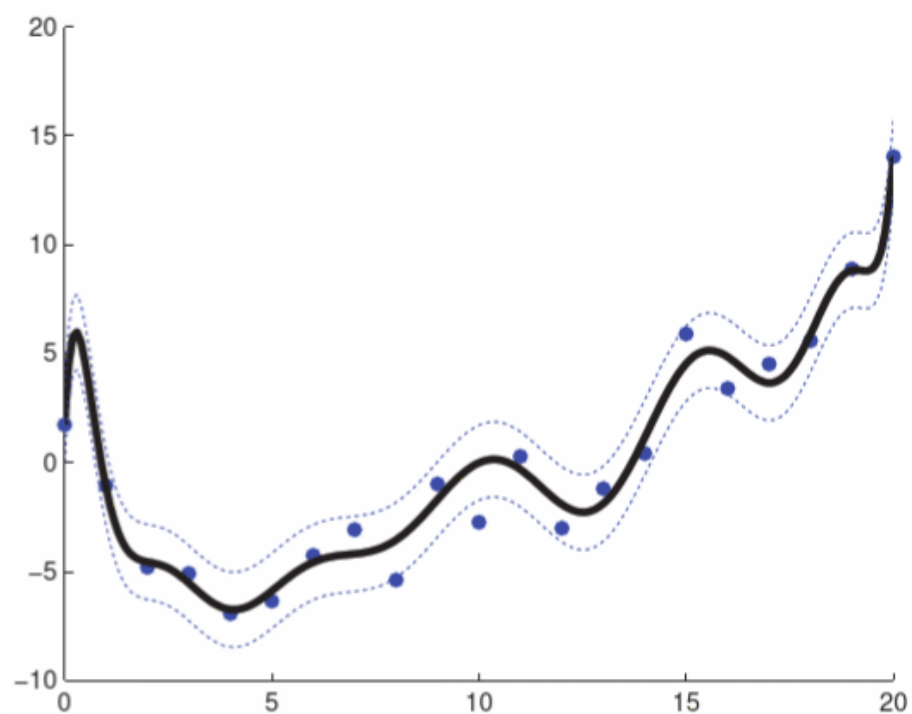
$$(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}) \mathbf{w} = \mathbf{X}^T \mathbf{y}$$

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

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

Regularization

But how do we choose λ ?



Probabilistic Interpretation

What assumptions are we making about our parameters?

MLE and MAP

Coin Flipping Example

Piazza Poll 2:

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

$$p(\theta \mid \mathcal{D}) \propto \prod p(\mathcal{D}^{(n)} \mid \theta) p(\theta)$$

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

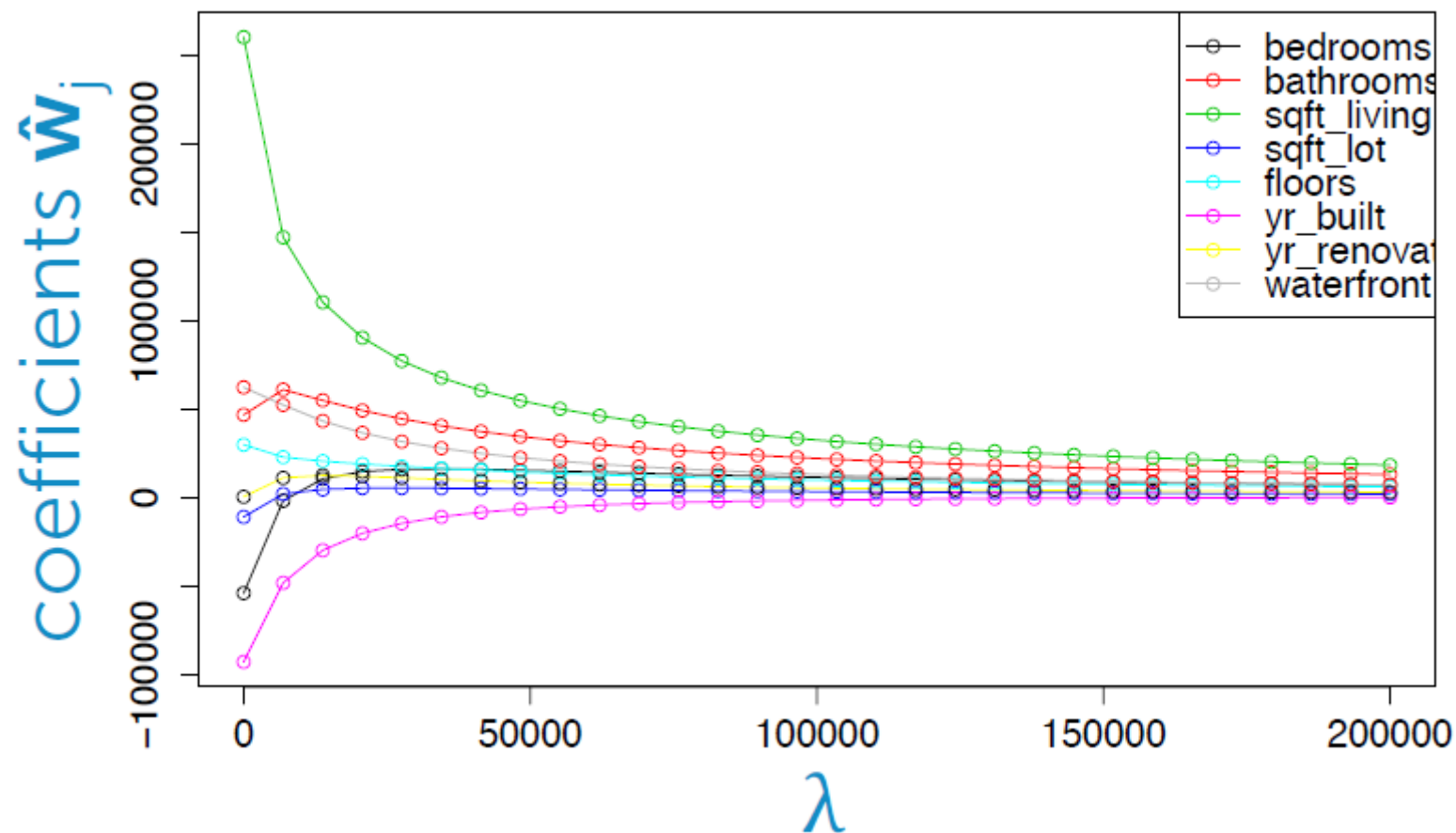


Figure: Emily Fox, University of Washington

Housing Price Example

Predict housing price from several features

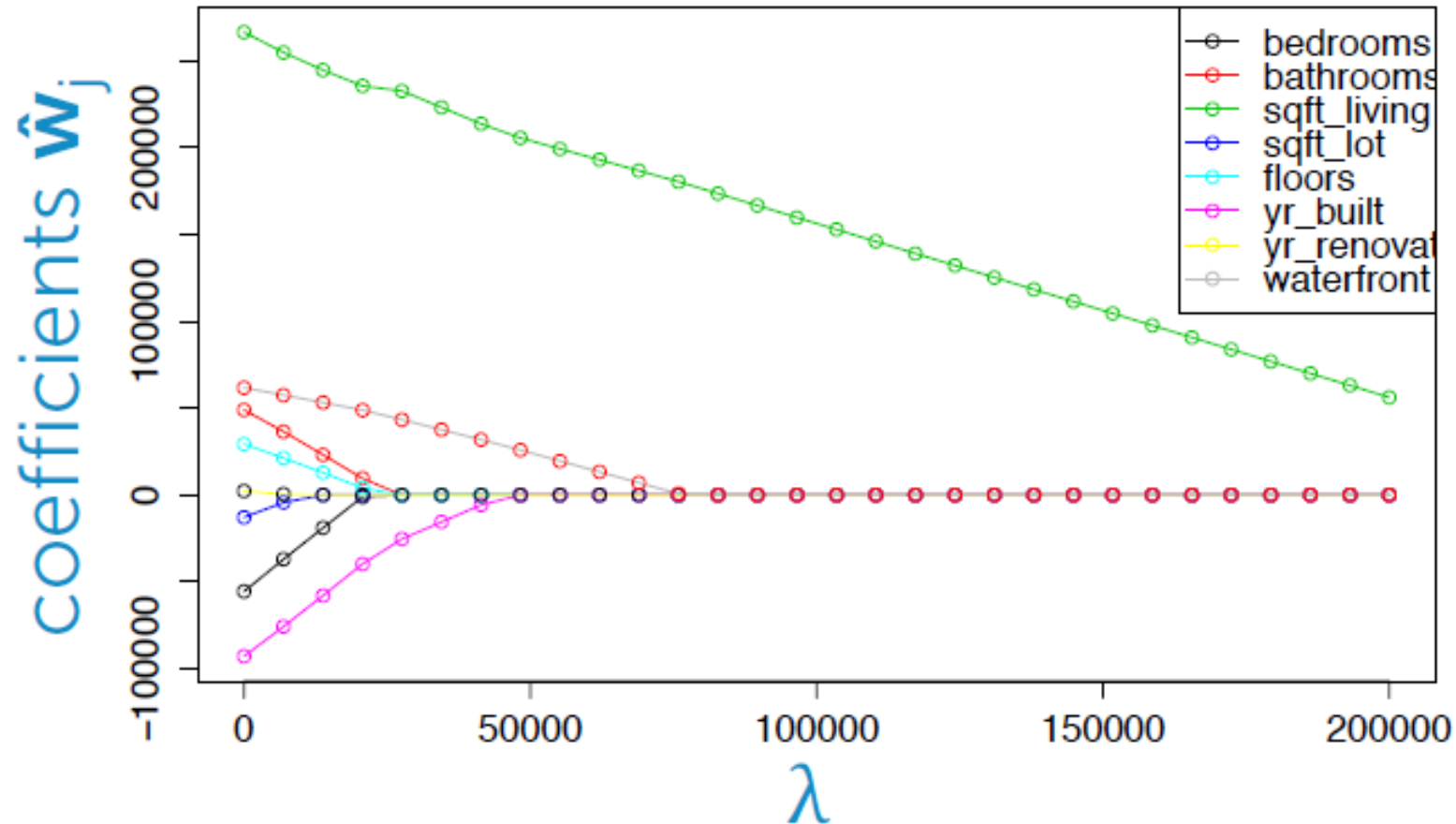


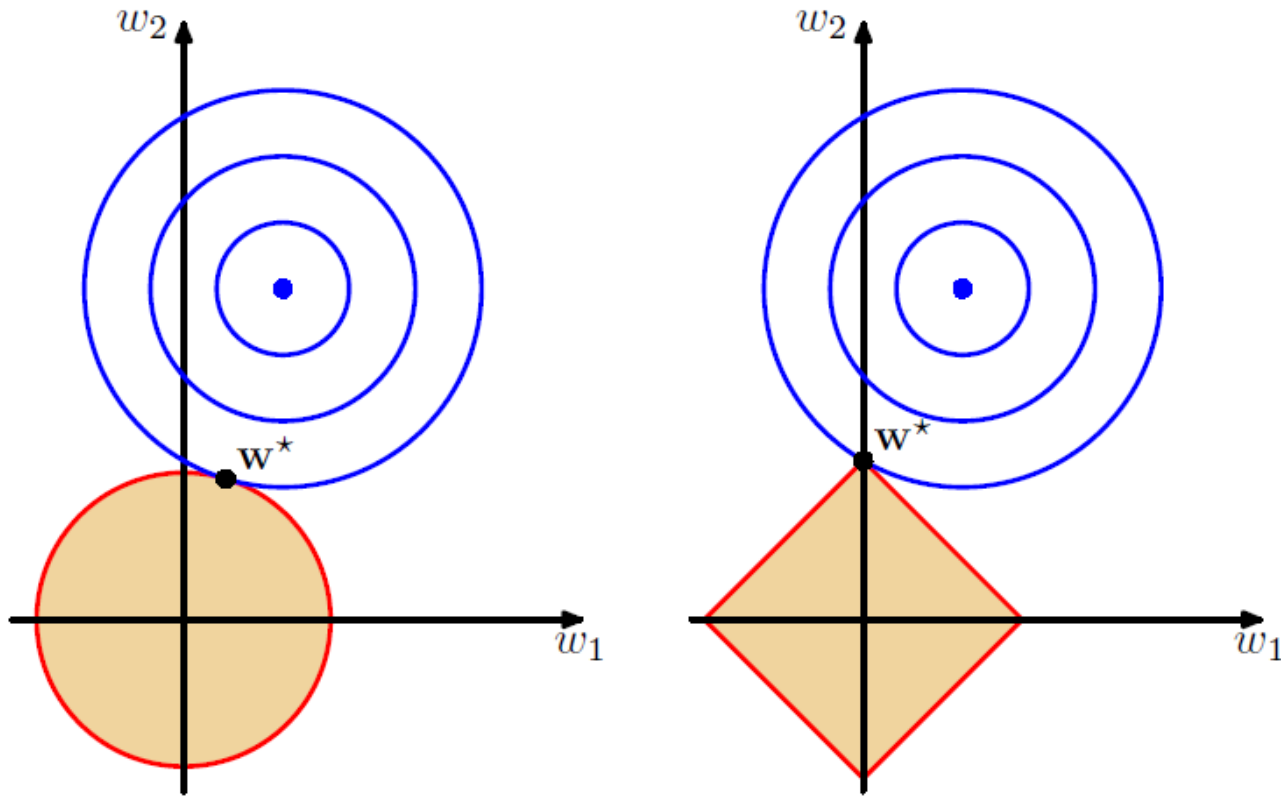
Figure: Emily Fox, University of Washington

Regularization

Combine original objective with penalty on parameters

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

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

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