

Applied Machine Learning

Lecture: 6 Regularization

Ekarat Rattagan, Ph.D.

Slides adapted from Andrew NG, Eric Eaton, Raquel Urtasun, and Patrick Winston

Outline

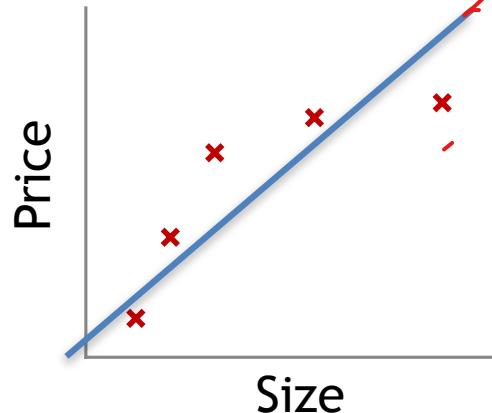
6.1 Overfitting

6.2 Regularization

- Ridge
- Lasso
- Elastic Net

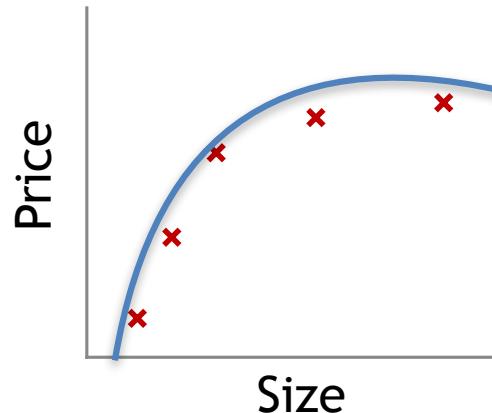
6.1 Overfitting

Example: Linear regression (housing prices)



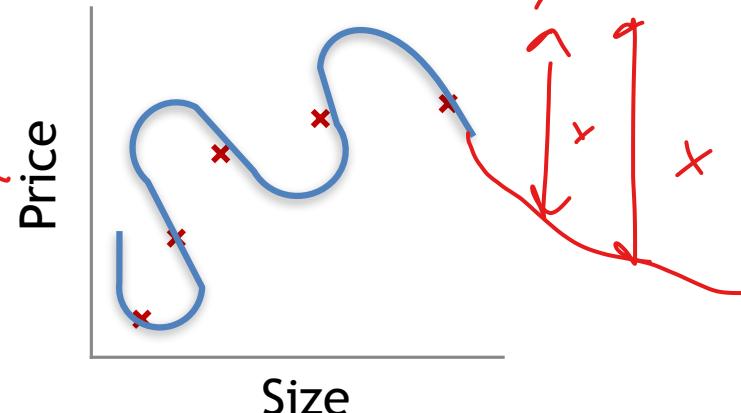
$$\theta_0 + \theta_1 x$$

Underfit
bias



$$\theta_0 + \theta_1 x + \theta_2 x^2$$

Good fit



$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

Overfit
Variance

Overfitting: If we have too many features, the learned hypothesis may fit the training set very well ($J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 \approx 0$), but fail to generalize to new examples.

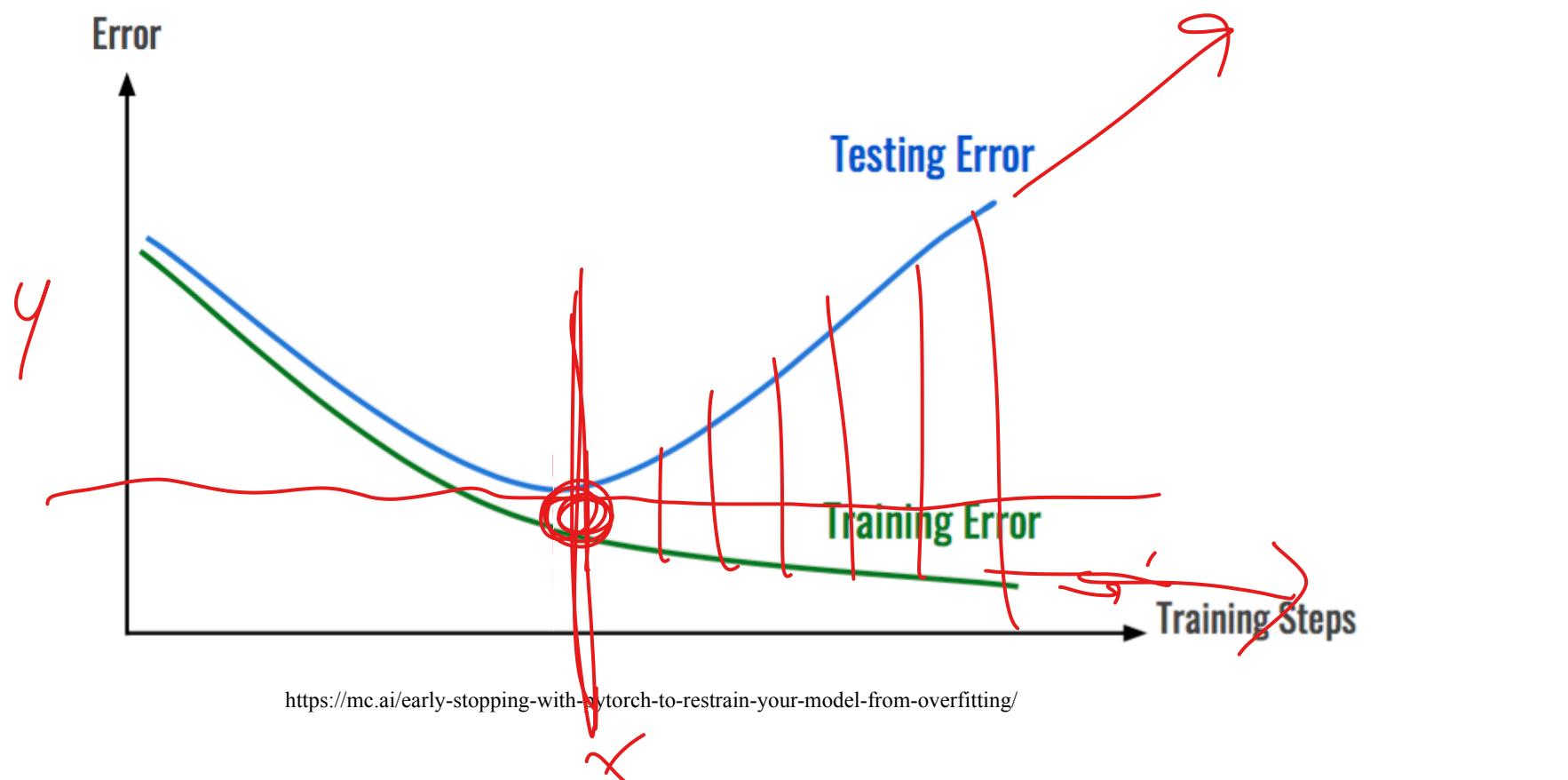
$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

Test

Credit: Andrew NG

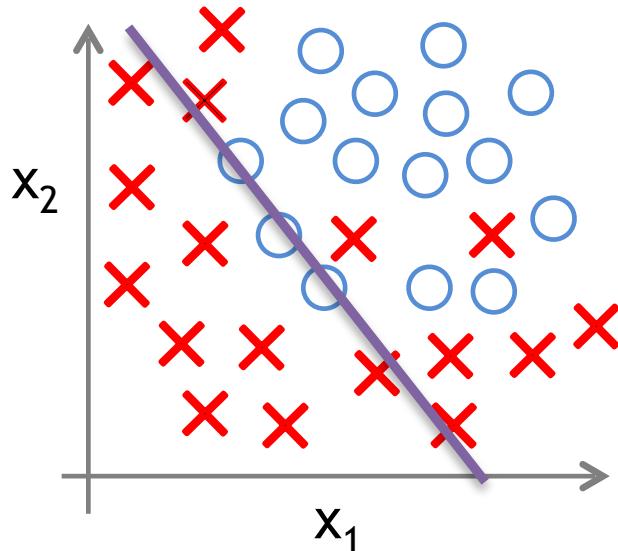
degree

Overfitting example: regression using polynomials



classification

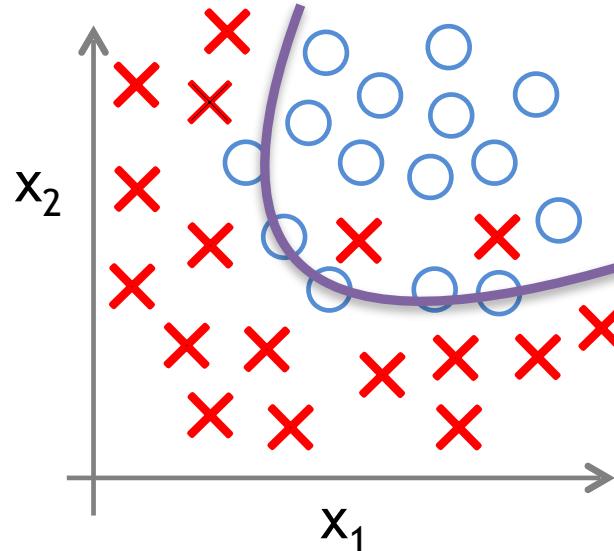
Example: Logistic regression



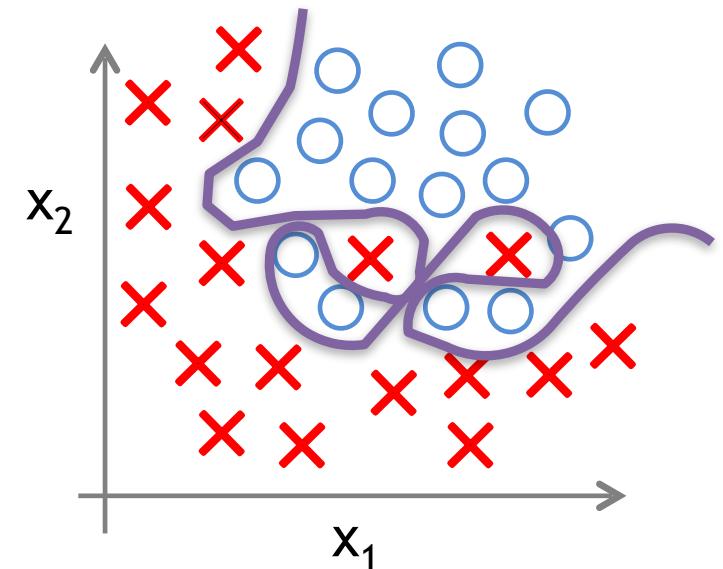
$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

(g = sigmoid function)

6



$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2)$$



$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 + \theta_6 x_1^3 x_2 + \dots)$$

6

Credit: Andrew NG

$x_1, x_2, x_3, \cancel{x_4}$

Addressing overfitting:

Options:

1. Reduce number of features
 - Manually select which features to keep.
 - Model selection algorithm. → Random Forest +
PCA
2. Regularization
 - Keep all the features, but reduce magnitude/values of parameters θ_j .
 - Works well when we have a lot of features, each of which contributes a bit to predicting y .

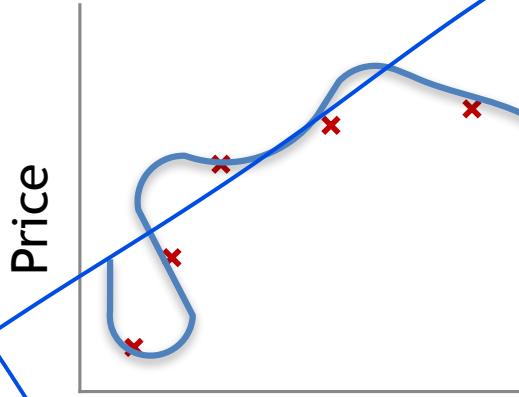
$$\theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

Credit: Andrew NG

$$10x_1 + 100x_2 + 80x_3 \rightarrow 2x_1 + 10x_2 + 8x_3$$

6.2 Regularization

Intuition

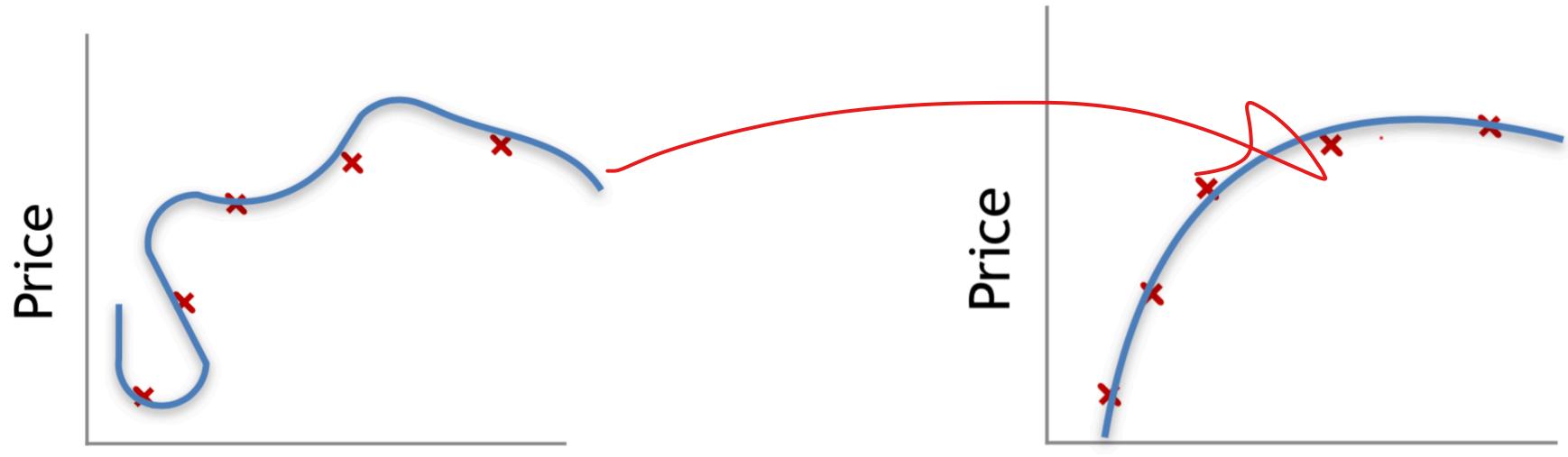


$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

Suppose we penalize and make θ_3, θ_4 really small.
The curve will change.

A good way to reduce overfitting is to regularize the model (i.e., to constrain it)

Intuition



Size of house

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4 \quad \xrightarrow{\hspace{1cm}} \quad \theta_0 + \theta_1 x + \theta_2 x^2$$

Suppose we penalize and make θ_3, θ_4 really small.
The curve will change.

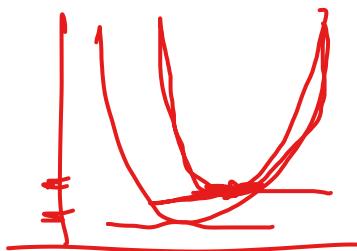
A good way to reduce overfitting is to regularize the model (i.e., to constrain it)

Regularized Linear Models

- Ridge regression
- Lasso regression
- Elastic Net



$$\| \cdot \|_2 = L_2\text{-norm}$$



Ridge regression (squared L_2 -norm)

$$\| \cdot \|_2^2 = \text{square } L_2\text{-norm}$$

- Additional term in the cost function to prevent overfitting
- keep the model weights/coefficients/parameters as small as possible
- L_2 (Euclidean) norm: A square root of the sum of the squared vector values

$$\| \cdot \|_2^2 = \sqrt{x_1^2 + x_2^2 + \dots + x_n^2}$$

Tuning/ridge/penalty parameter controls how much you want to regularize the model

$$+ \lambda \|x\|_2^2$$

new

$$J(\theta) = \frac{1}{2m} \left(\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right)$$

$$J(\theta) = \frac{1}{2m} \sum_{i=0}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^n \theta_j^2$$

$MSE + \lambda \text{ridge}$

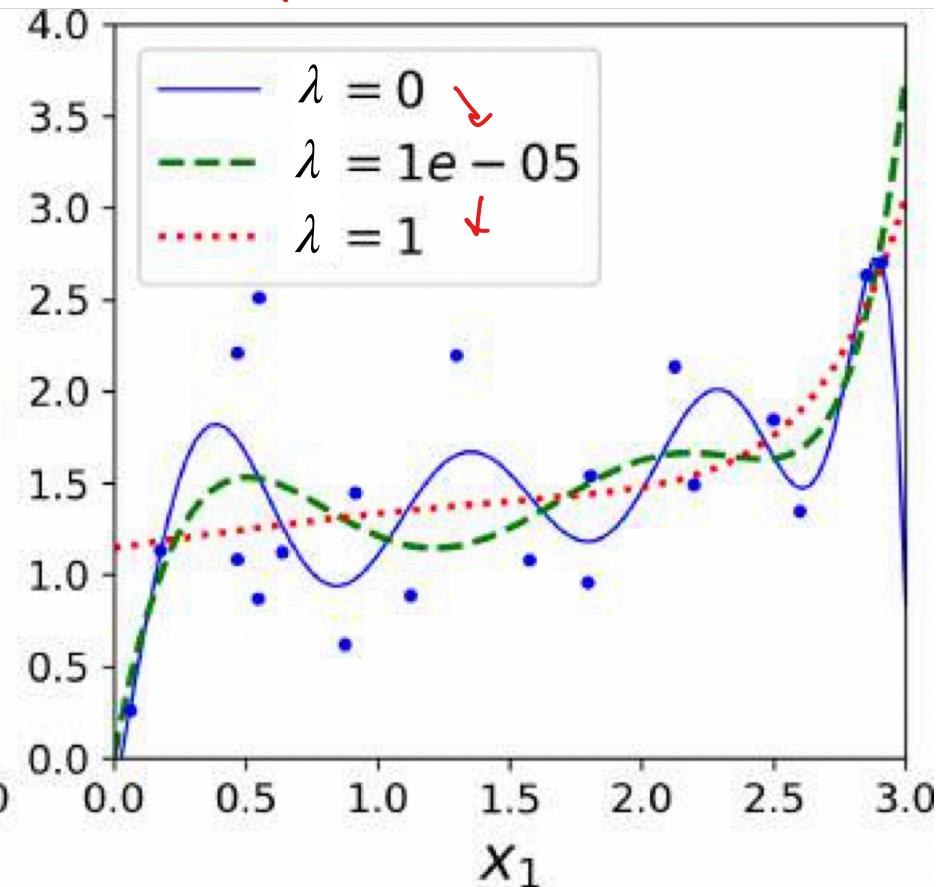
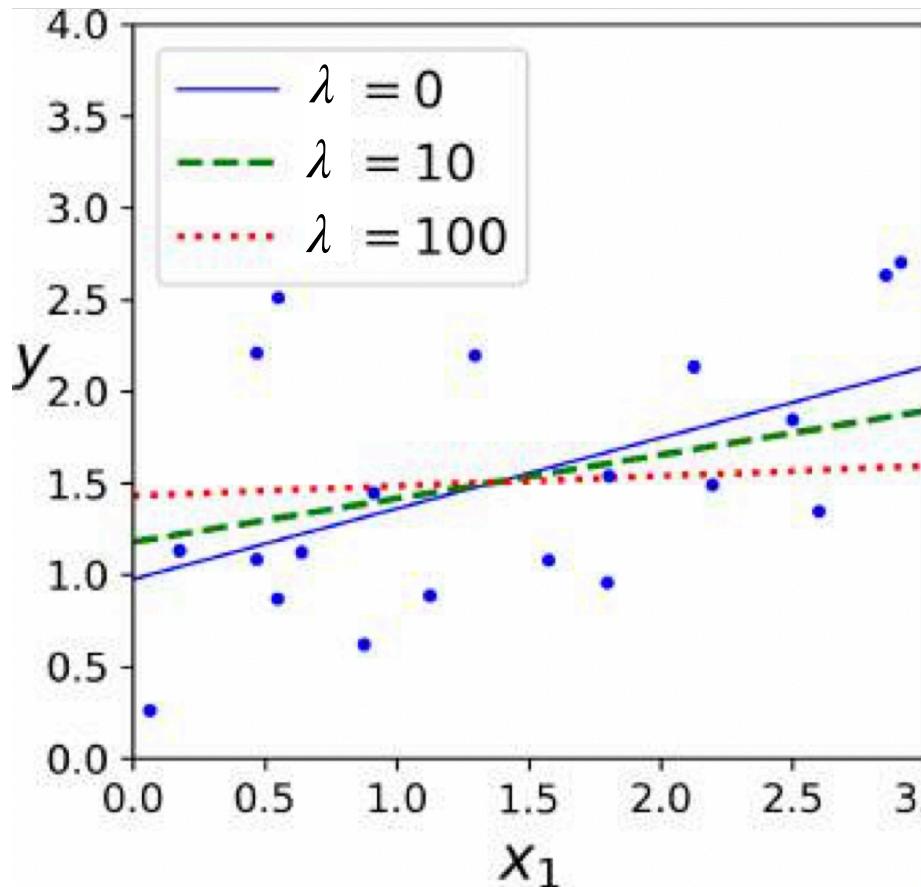
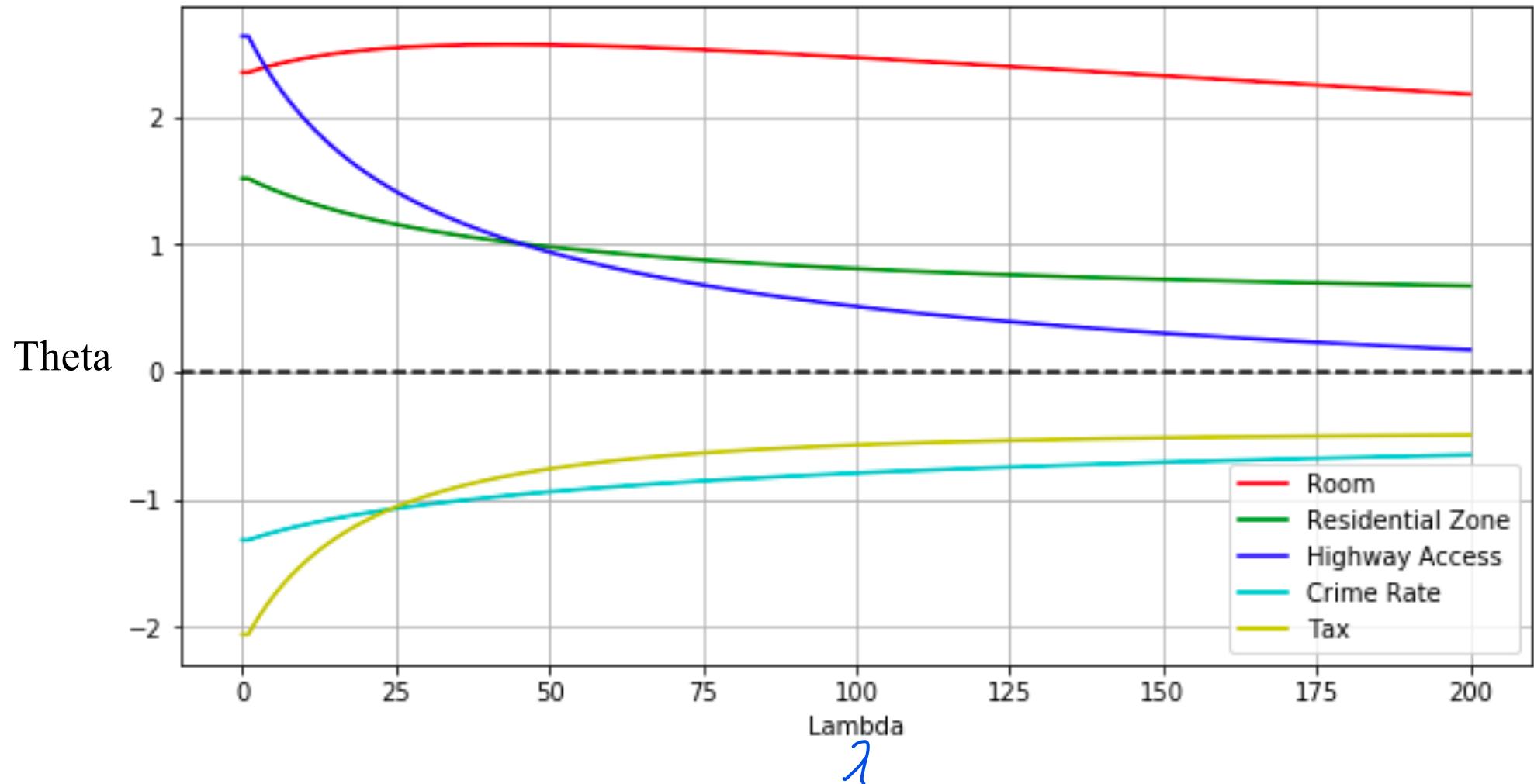


Figure from Hands-on
Machine Learning

Ridge regression

Ridge Regression Trace



Lasso regression (L_1 - norm)

$$\| \cdot \|_1$$

- Least Absolute Shrinkage and Selection Operator (LASSO)
 - Reduce the number of features in a regression model
 - Automatically performs feature selection and outputs a sparse model
 - L_1 (Absolute-value) norm: $\|x\|_1 = |x_1 + x_2 + \dots + x_n|$

new

$$J(\theta) = \frac{1}{2m} \left(\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n |\theta_j| \right)$$

$$J(\theta) = \frac{1}{2m} \sum_{i=0}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n |\theta_j|$$

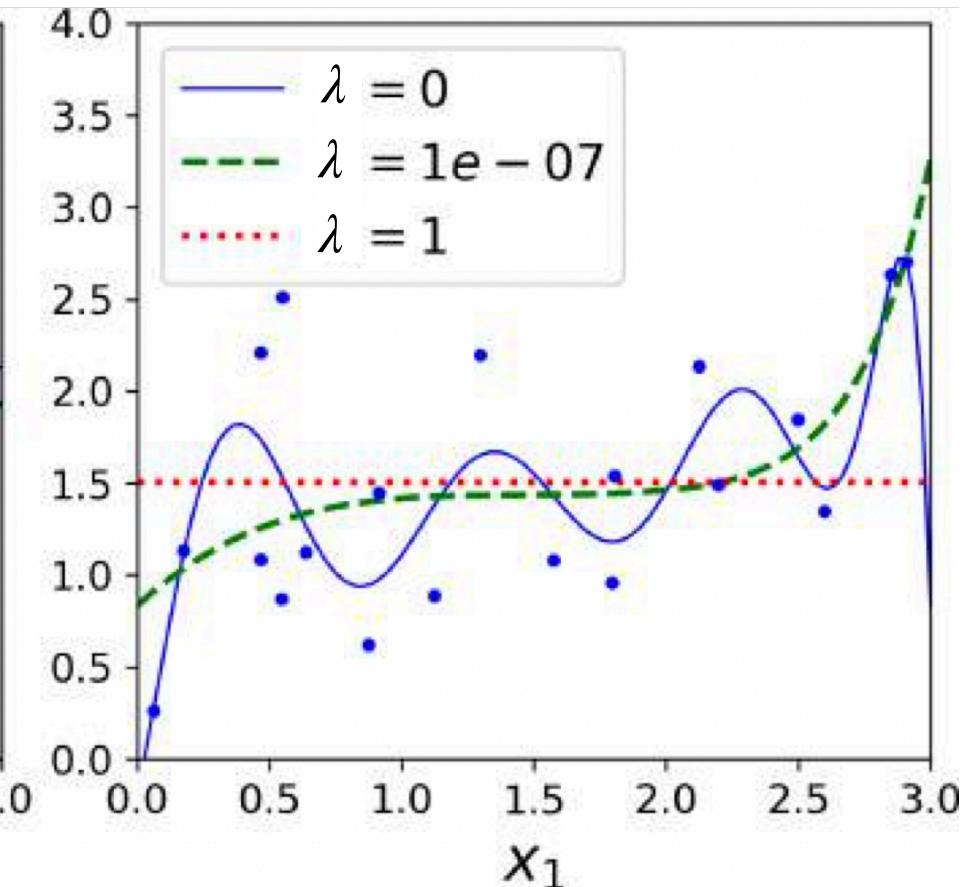
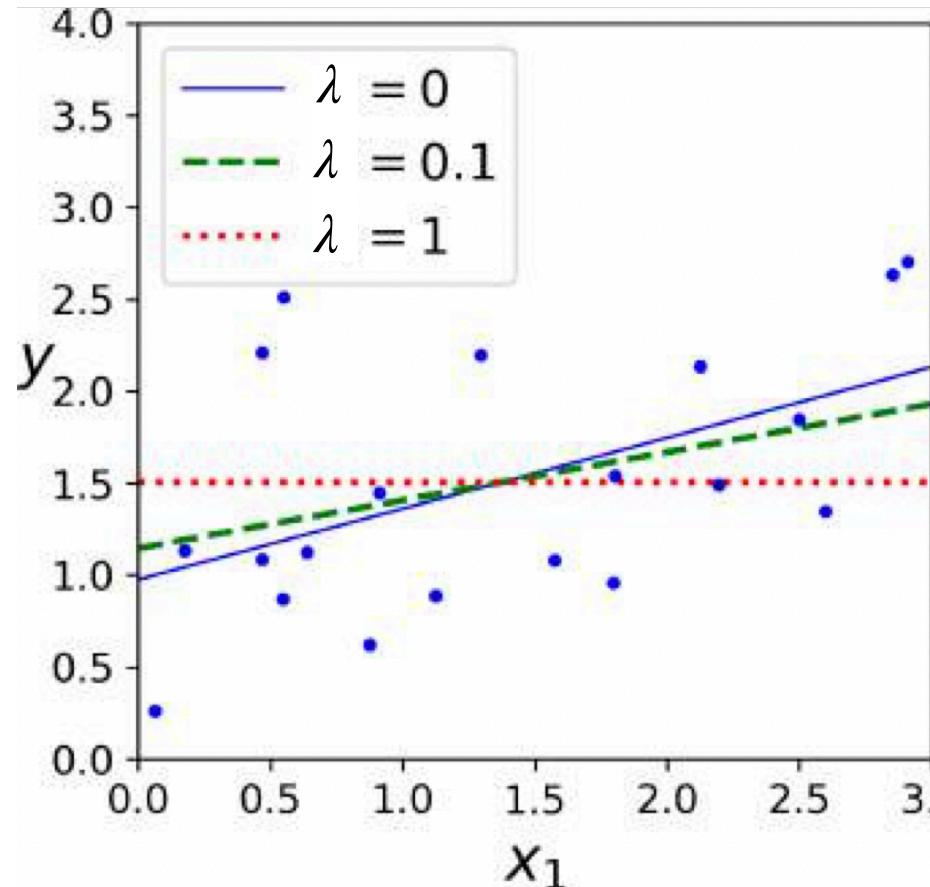
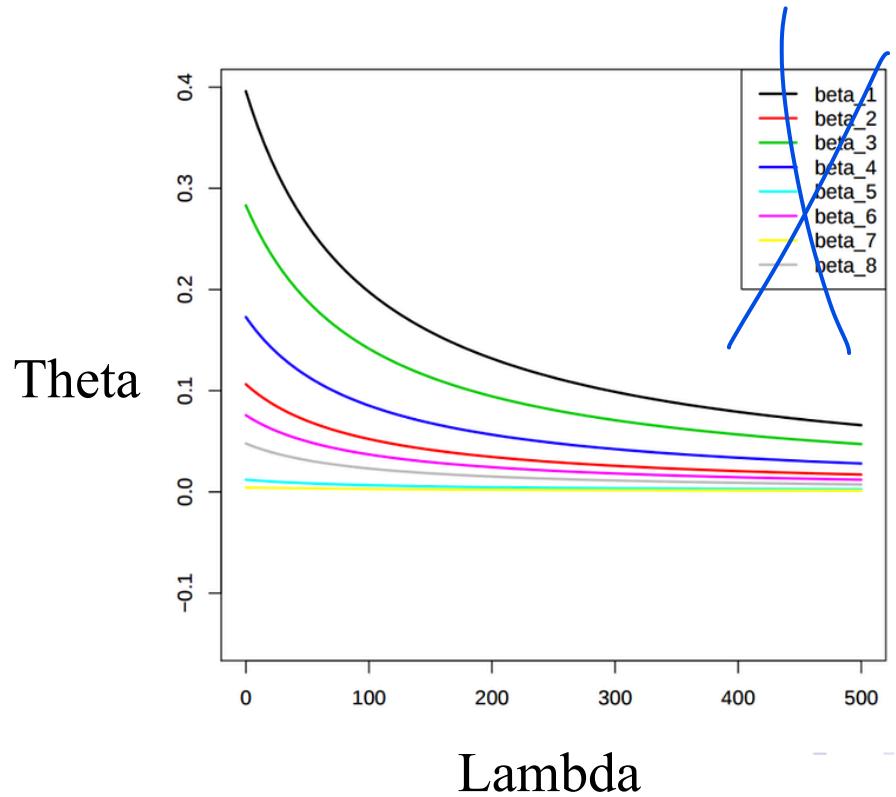


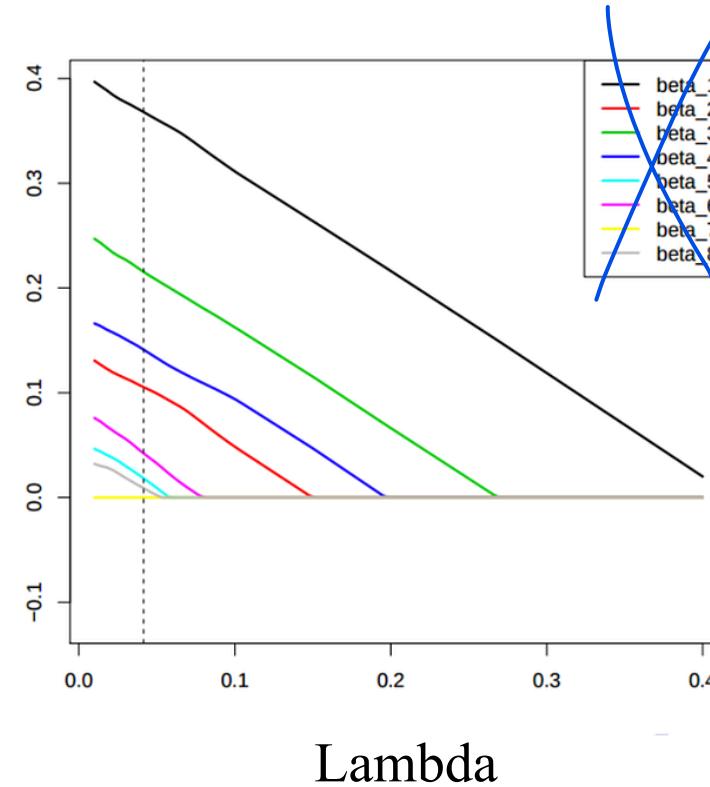
Figure from Hands-on
Machine Learning

Lasso regression

Ridge vs Lasso



Theta



Lambda

When to use what?

Lasso tends to do well if there are a small number of significant parameters and the others are close to zero (hence: when only a few predictors actually influence the response).

Ridge works well if there are many large parameters of about the same value (hence: when most predictors impact the response).

However, in practice, we don't know the true parameter values, so the previous two points are somewhat theoretical. Just run **cross-validation** to select the more suited model for a specific case.

Or... combine the two!

https://shuaili8.github.io/Teaching/VE445/L3_linear%20regression.pdf

Case study



Dmitry Smirnov, Consultant

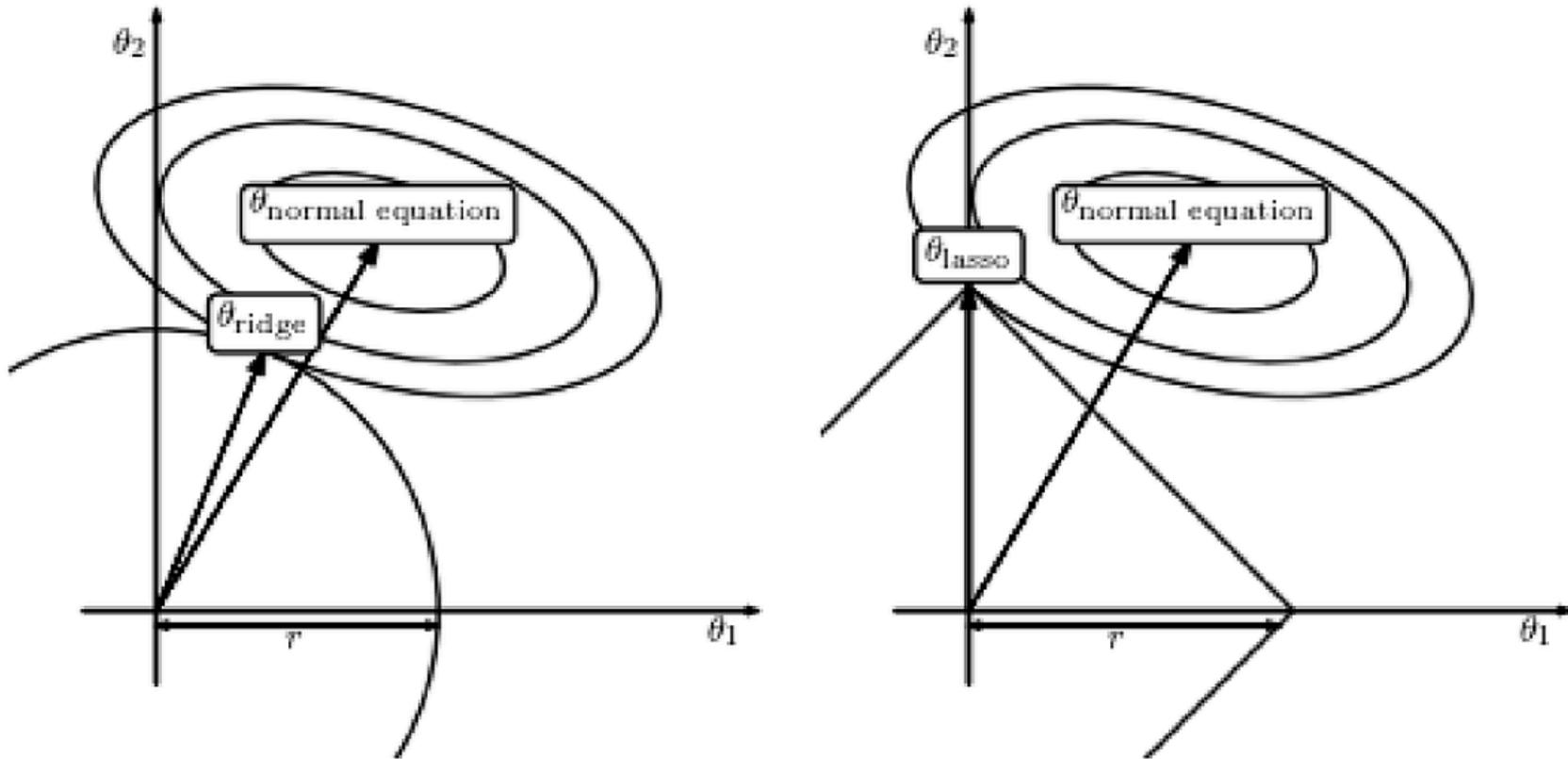
Answered June 19, 2019

Ridge and Lasso are regularized versions of linear regression. In practice, ridge regression leads to lower weights of coefficients, while lasso regressions leads to more coefficients having zero weights.



Which one fits better - depends on your domain, when I was in neuroscience, we often preferred ridge, because when analysing fMRI data you have multiple sources of signal that are located close to each other and having correlated signals, and ridge would provide similar coefficients for these sources, while lasso will drive a lot of these to zero while preserving some that will have higher weights.

Ridge vs Lasso



Elastic Net

- A simple mix of both Ridge and Lasso's regularization terms, and you can control the mix ratio α , where

$$\alpha = \frac{\lambda_2}{\lambda_1 + \lambda_2}$$

λ_2 = Ridge parameter
 λ_1 = Lasso " " $||$

- When $\alpha = 1$, Elastic Net is equivalent to Ridge Regression, and when $\alpha = 0$, it is equivalent to Lasso Regression.



$$J(\theta) = \frac{1}{2m} \left(\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + (1-\alpha) \sum_{j=1}^n |\theta_j| + \alpha \sum_{j=1}^n \theta_j^2 \right)$$

$$J(\theta) = \frac{1}{2m} \sum_{i=0}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + (1-\alpha) \sum_{j=1}^n |\theta_j| + \frac{\alpha}{2} \sum_{j=1}^n \theta_j^2$$

Elastic Net

2-dimensional illustration $\alpha = 0.5$

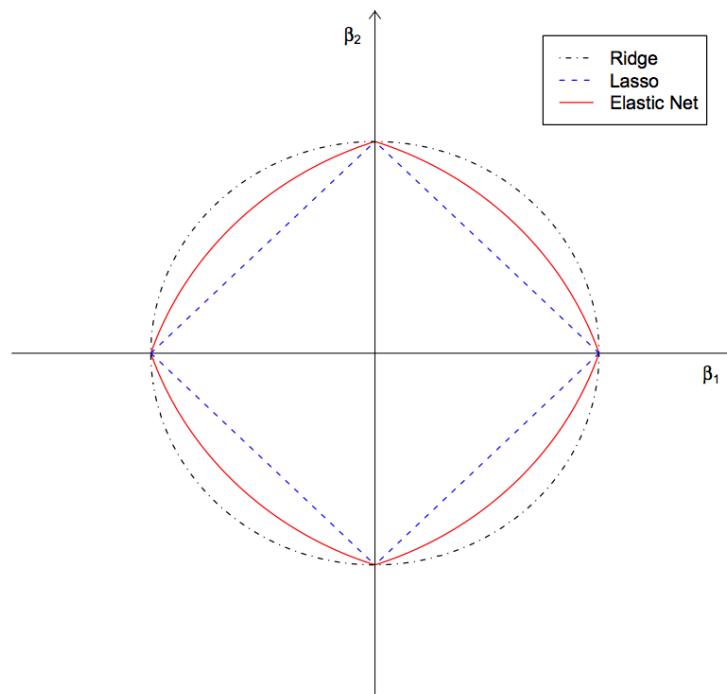
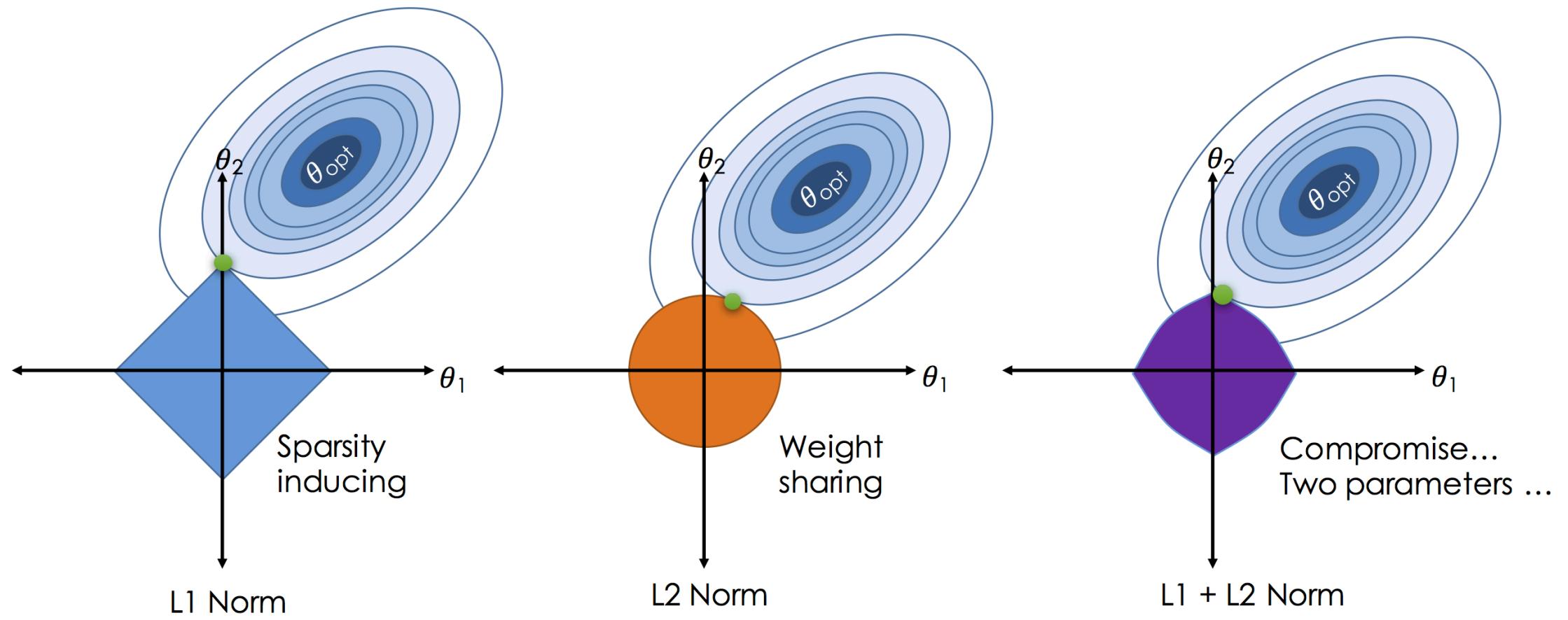


Figure from: Regularization and Variable Selection via the Elastic Net
Hui Zou and Trevor Hastie Department of Statistics Stanford University



When to use what?

So when should you use plain Linear Regression (i.e., without any regularization, Ridge, Lasso, or Elastic Net) ?

It is almost always preferable to have at least a little bit of regularization, Ridge is a good default,

but if you suspect that only a few features are actually useful, you should prefer Lasso or Elastic Net since they tend to reduce the useless features' weights down to zero.

In general, Elastic Net is preferred over Lasso since Lasso may behave erratically when the number of features is greater than the number of training instances or when several features are strongly correlated.

ຈົດປັດທະນາກົດໃຫຍ່

$$\min \frac{\lambda}{2} \sum_{j=1}^n \theta_j^2$$

$$\frac{\lambda}{2} \theta_0^2 + \frac{\lambda}{2} \theta_1^2 + \dots + \frac{\lambda}{2} \theta_n^2$$

ນະນຳ

$$\frac{\partial J(\theta)}{\partial \theta_j} = \text{ກົດທີ່}$$

ເນືອນ
ຕົວຢ່າງ

$$\frac{\partial J(\theta)}{\partial \theta_j} = \text{ກົດທີ່}$$

ທີ່ມີຄົງ

$$\lambda \theta_j =$$

Ridge regression

ridge penalty term
ກົດທີ່ bias ກົດທີ່ θ_0

repeat until convergence {

$$\theta_0 = \theta_0 - \alpha \left((h_\theta(x^{(i)}) - y^{(i)})x_0^{(i)} + \cancel{\lambda \theta_0} \right)$$

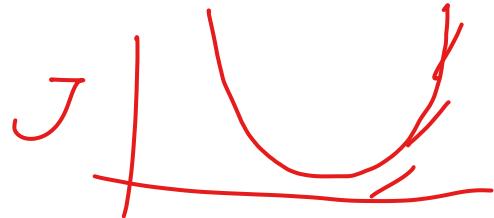
$$\theta_1 = \theta_1 - \alpha \left((h_\theta(x^{(i)}) - y^{(i)})x_1^{(i)} + \lambda \theta_1 \right)$$

.

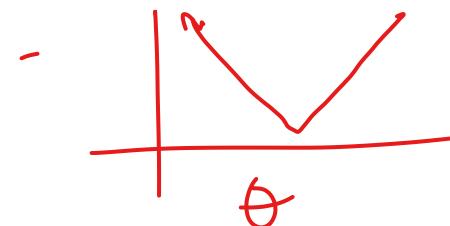
.

$$\theta_n = \theta_n - \alpha \left((h_\theta(x^{(i)}) - y^{(i)})x_n^{(i)} + \lambda \theta_n \right)$$

}



Stochastic Gradient descent



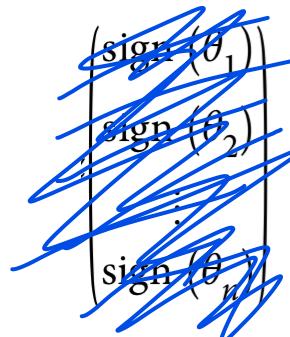
Lasso regression

The Lasso cost function is not differentiable at $\theta_j = 0$, where $j = 0, 1, \dots, n$ but Gradient Descent still works fine if you use a sub-gradient vector instead when any $\theta_j = 0$.

repeat until convergence {

$$\theta_j = \theta_j - \alpha \left[(h_\theta(x^{(i)}) - y^{(i)})x_j^{(i)} + \lambda \underset{j=1}{\overset{n}{\text{sign}(\theta_j)}} \right], \text{ for } j = 0 \dots n$$

}



where $\text{sign}(\theta_j) =$

$$\begin{cases} 1 & \text{if } \theta_j > 0 \\ 0 & \text{if } \theta_j = 0 \\ -1 & \text{if } \theta_j < 0 \end{cases}$$

$$= \begin{cases} \theta_j - \lambda & \text{if } \theta_j > \lambda \\ 0 & \text{if } -\lambda \leq \theta_j \leq \lambda \\ \theta_j + \lambda & \text{if } \theta_j < -\lambda \end{cases}$$

soft-thresholding operator.
[CMU]

25

Stochastic Gradient descent

Elastic Net $\alpha = \frac{\lambda_2}{\lambda_1 + \lambda_2}$

repeat until convergence {

$$\theta_j = \theta_j - \alpha \left[h_{\theta}(x^{(i)}) - y^{(i)} \right] x_j^{(i)} + (1 - \alpha) sign(\theta_j) + \alpha \theta_j$$

}

for j=0 ... n

Reference

- Elastic Net

Zou, Hui, and Trevor Hastie. "Regularization and variable selection via the elastic net." *Journal of the royal statistical society: series B (statistical methodology)* 67, no. 2 (2005): 301-320.

<https://www.youtube.com/watch?v=TJer4ibUZ0I>