

Announcements

Assignments

- HW4
 - Wed, 10/14, 11:59 pm

Survey

- Thanks!
- We'll talk more Wednesday

Plan

Last Time

- Feature engineering
- Regularization with added penalty term

Today

- Wrap-up regularization
- Neural Networks
 - Perceptron
 - Multilayer perceptron
 - Building blocks
 - Objective
 - Optimization

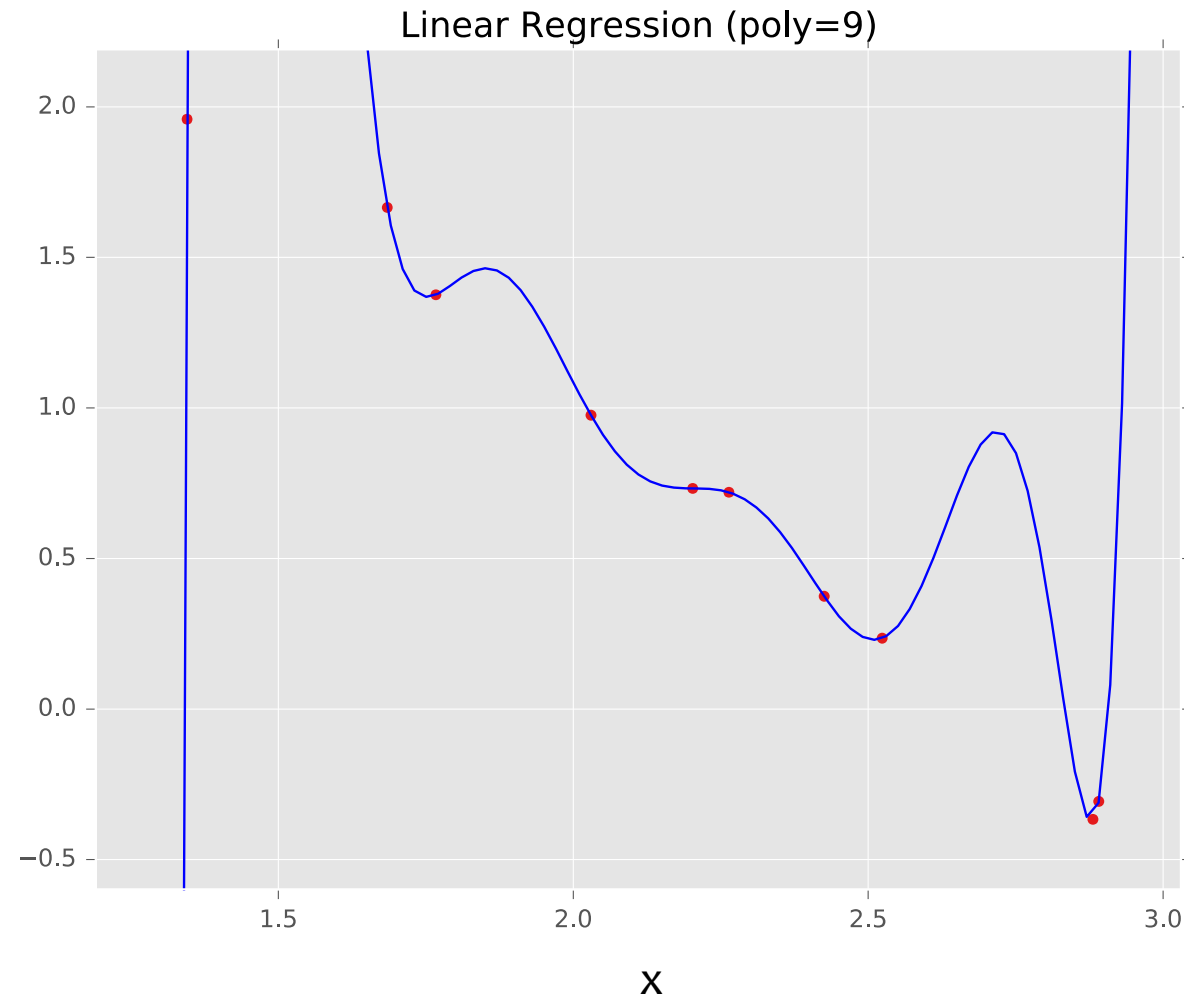
Wrap-up Regularization

Example: Linear Regression

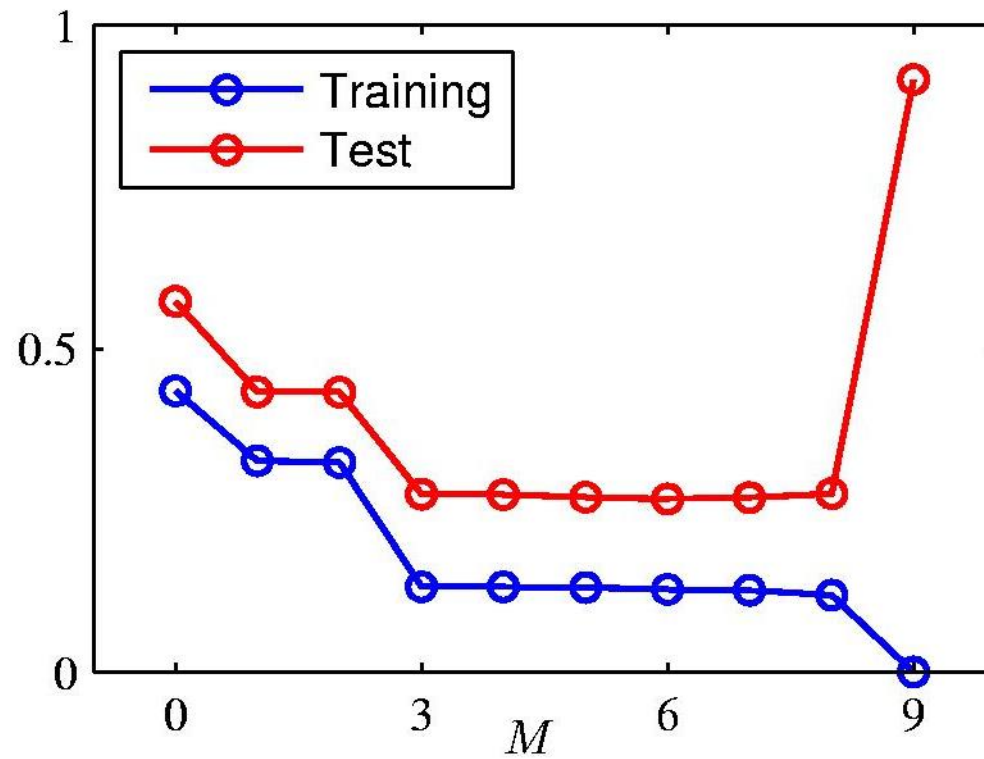
Goal: Learn $y = \mathbf{w}^T \mathbf{f}(\mathbf{x}) + b$
where $\mathbf{f}(\cdot)$ is a polynomial
basis function

y	x	x^2	...	x^9
2.0	1.2	$(1.2)^2$...	$(1.2)^9$
1.3	1.7	$(1.7)^2$...	$(1.7)^9$
0.1	2.7	$(2.7)^2$...	$(2.7)^9$
1.1	1.9	$(1.9)^2$...	$(1.9)^9$

true “unknown”
target function is
linear with
negative slope
and gaussian
noise



Over-fitting



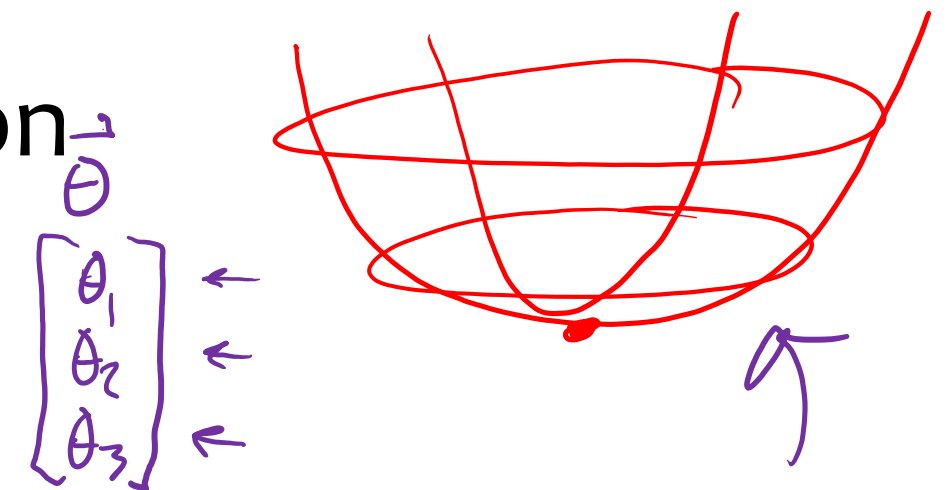
Polynomial Coefficients

	$M = 0$	$M = 1$	$M = 3$	$M = 9$
θ_0	0.19	0.82	0.31	0.35
θ_1		-1.27	7.99	232.37
θ_2			-25.43	-5321.83
θ_3			17.37	48568.31
θ_4				-231639.30
θ_5				640042.26
θ_6				-1061800.52
θ_7				1042400.18
θ_8				-557682.99
θ_9				125201.43

Regularization

Given objective function: $J(\theta)$

Goal is to find: $\hat{\theta} = \operatorname{argmin}_{\theta} \underline{J(\theta)} + \lambda \underline{r(\theta)}$



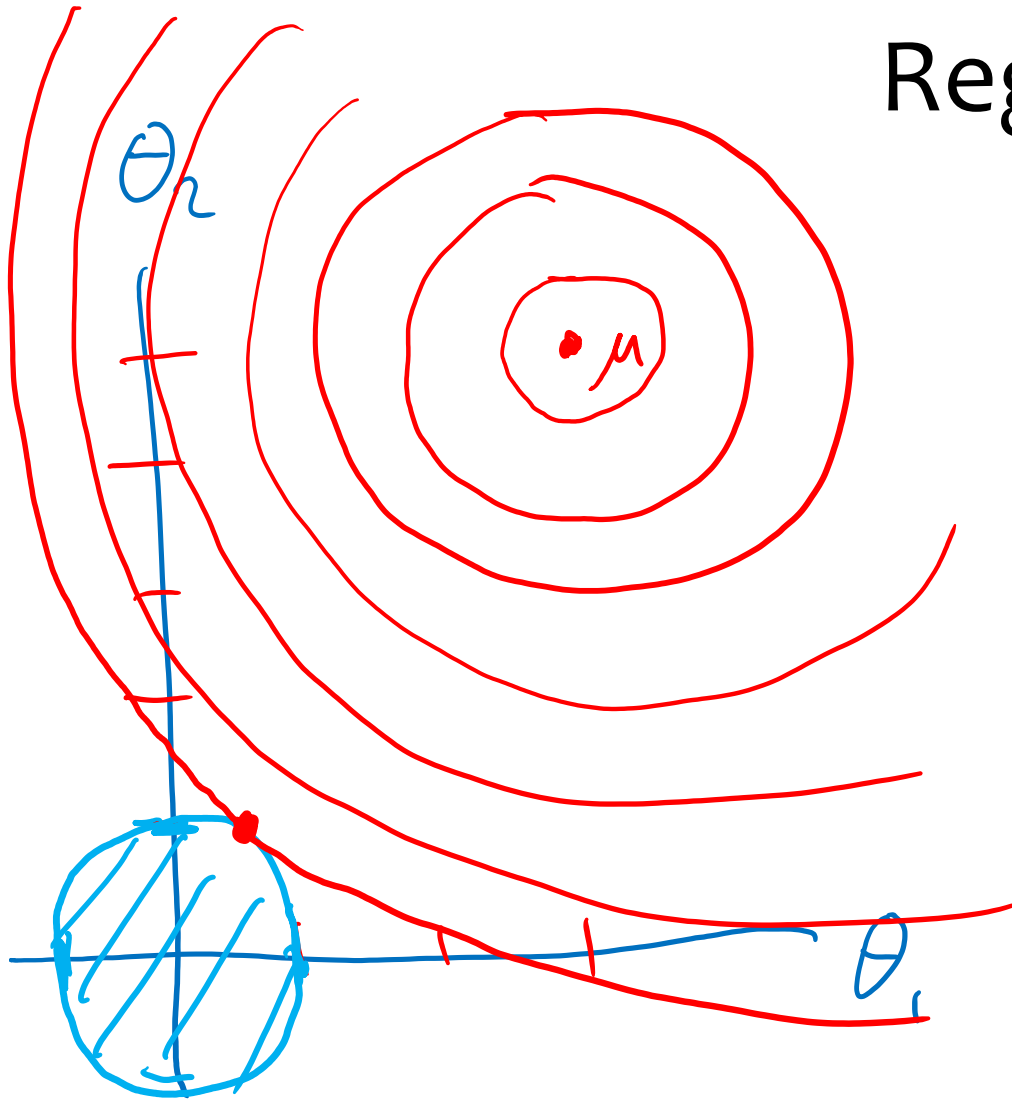
Key idea: Define regularizer $r(\theta)$ s.t. we tradeoff between fitting the data and keeping the model simple

Choose form of $r(\theta)$:

– Example: q-norm (usually p-norm) $r(\theta) = \|\theta\|_q = \left[\sum_{m=1}^M \|\theta_m\|^q \right]^{\left(\frac{1}{q}\right)}$

q	$r(\theta)$	yields parameters that are...	name	optimization notes
L 0	$\ \theta\ _0 = \sum \mathbb{1}(\theta_m \neq 0)$	zero values <i>sparse</i>	L0 reg.	no good computational solutions <i># num nonzeros</i>
1	$\ \theta\ _1 = \sum \theta_m $	zero values <i>sparse</i>	L1 reg.	subdifferentiable
L 2	$(\ \theta\ _2)^2 = \sum \theta_m^2$	small values	L2 reg.	differentiable <i>←</i>

Regularization



$$J(\theta_1, \theta_2) = \|\vec{\theta} - \vec{\mu}\| \quad \mu = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$$

$$\min_{\theta} J(\theta_1, \theta_2)$$

$$\text{s.t.} \quad \|\theta\|_2 \leq 1$$

Previous Piazza Poll

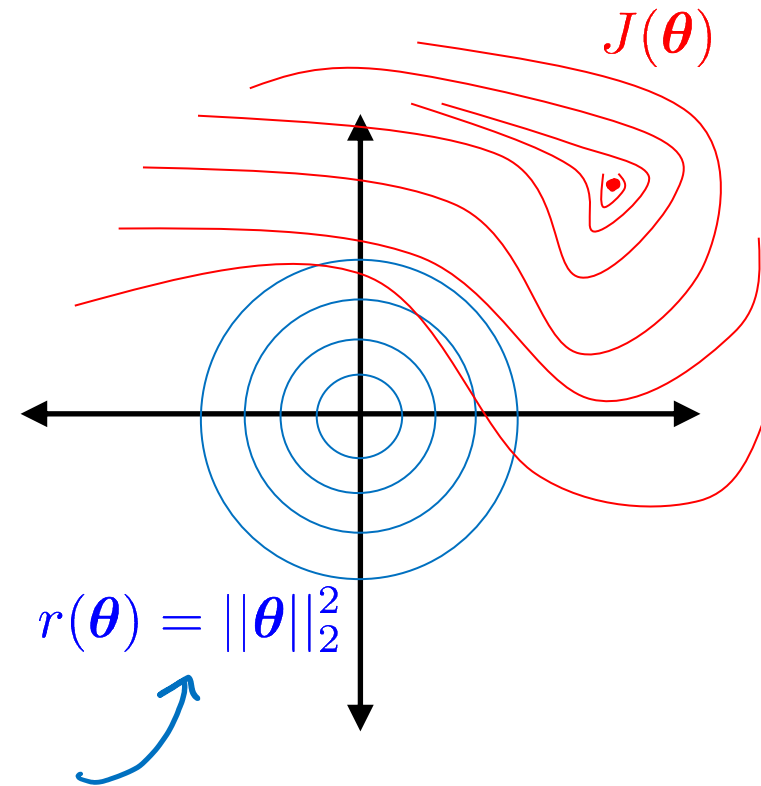
Question:

Suppose we are minimizing $J'(\theta)$ where

$$J'(\theta) = J(\theta) + \lambda r(\theta) \leftarrow$$

As λ increases, the minimum of $J'(\theta)$ will...

- A. ...move towards the midpoint between $J'(\theta)$ and $r(\theta)$
- B. ...move towards the minimum of $J(\theta)$
- ☒ C. ...move towards the minimum of $r(\theta)$
- D. ...move towards a theta vector of positive infinities
- E. ...move towards a theta vector of negative infinities
- F. ...stay the same



Regularization Exercise

In-class Exercise

1. Plot train error vs. regularization weight (cartoon)
2. Plot test error vs. regularization weight (cartoon)



$$\hat{\theta} = \operatorname{argmin}_{\theta} J(\theta) + \lambda r(\theta)$$

Piazza Poll 1

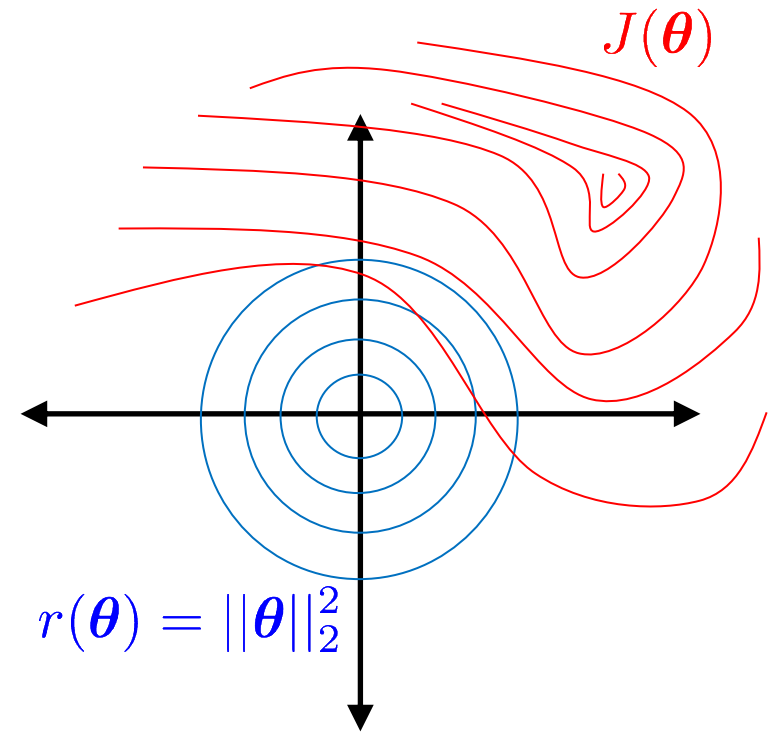
Question:

Suppose we are minimizing $J'(\theta)$ where

$$J'(\theta) = J(\theta) + \lambda r(\theta)$$

As we increase λ from zero, the **validation** error will...

- A. ...increase
- B. ...decrease
- C. ...first increase, then decrease
- D. ...first decrease, then increase
- E. ...stay the same



Regularization

Don't Regularize the Bias (Intercept) Parameter!

- In our models so far, the bias / intercept parameter is usually denoted by θ_0 -- that is, the parameter for which we fixed $x_0 = 1$
- Regularizers always avoid penalizing this bias / intercept parameter
- Why? Because otherwise the learning algorithms wouldn't be invariant to a shift in the y-values

Whitening Data

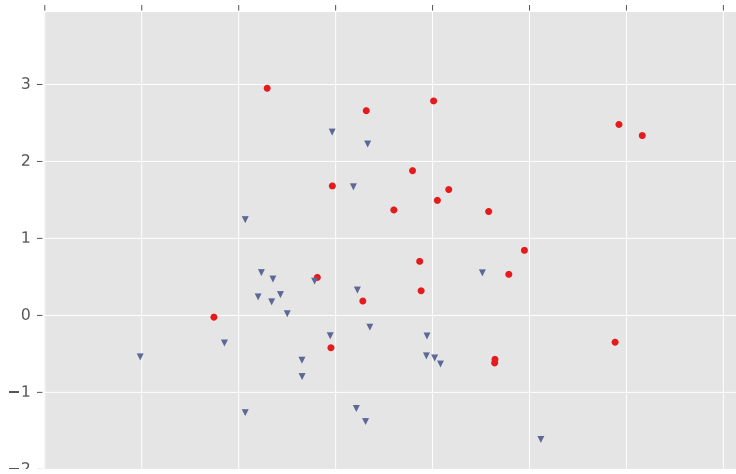
- It's common to *whiten* each feature by subtracting its mean and dividing by its variance
- For regularization, this helps all the features be penalized in the same units
(e.g. convert both centimeters and kilometers to z-scores)

Logistic Regression with Nonlinear Features

Jupyter notebook demo

Example: Logistic Regression

Training
Data

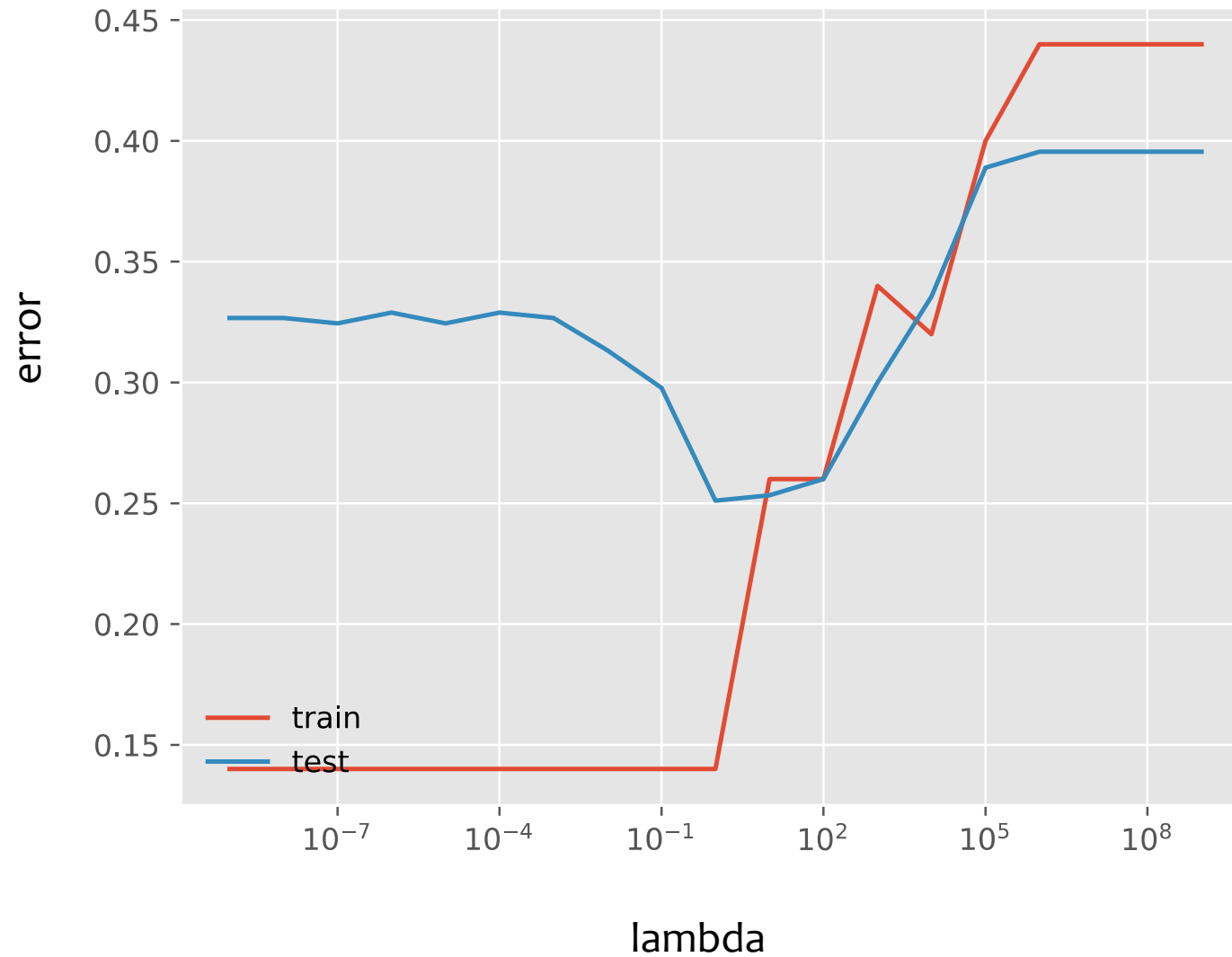


Test
Data

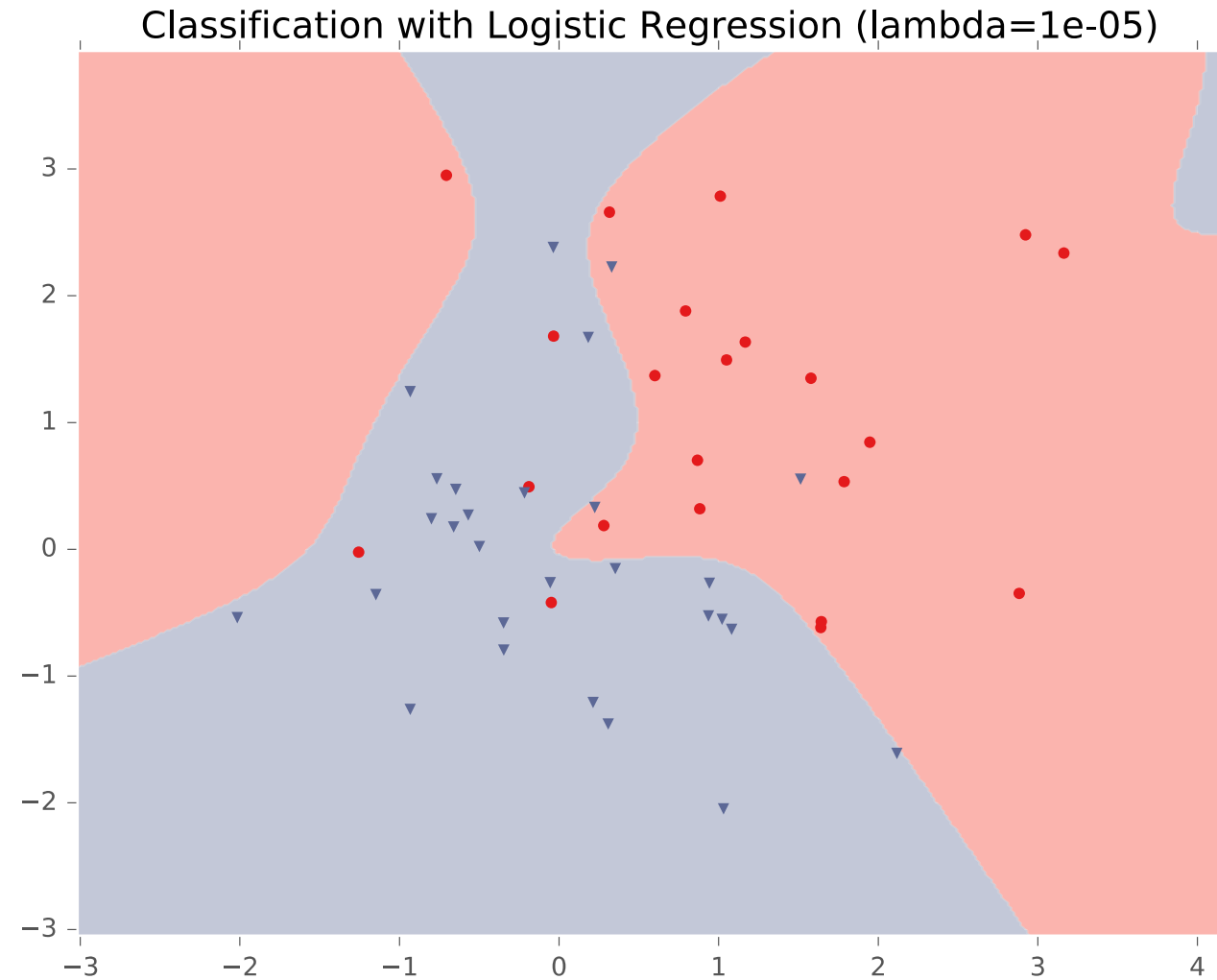


- For this example, we construct **nonlinear features** (i.e. feature engineering)
- Specifically, we add **polynomials up to order 9** of the two original features x_1 and x_2
- Thus our classifier is **linear** in the **high-dimensional feature space**, but the decision boundary is **nonlinear** when visualized in **low-dimensions** (i.e. the original two dimensions)

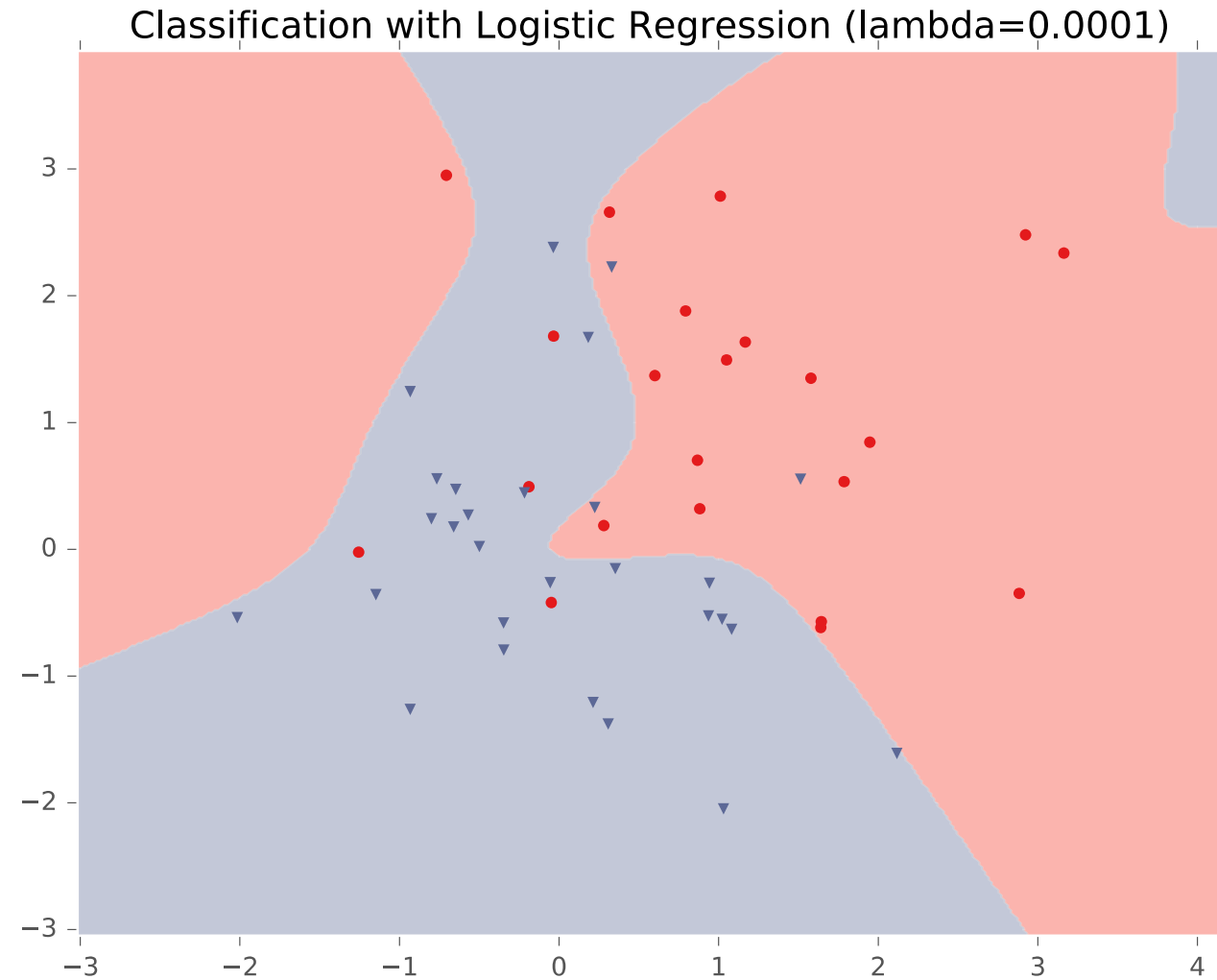
Example: Logistic Regression



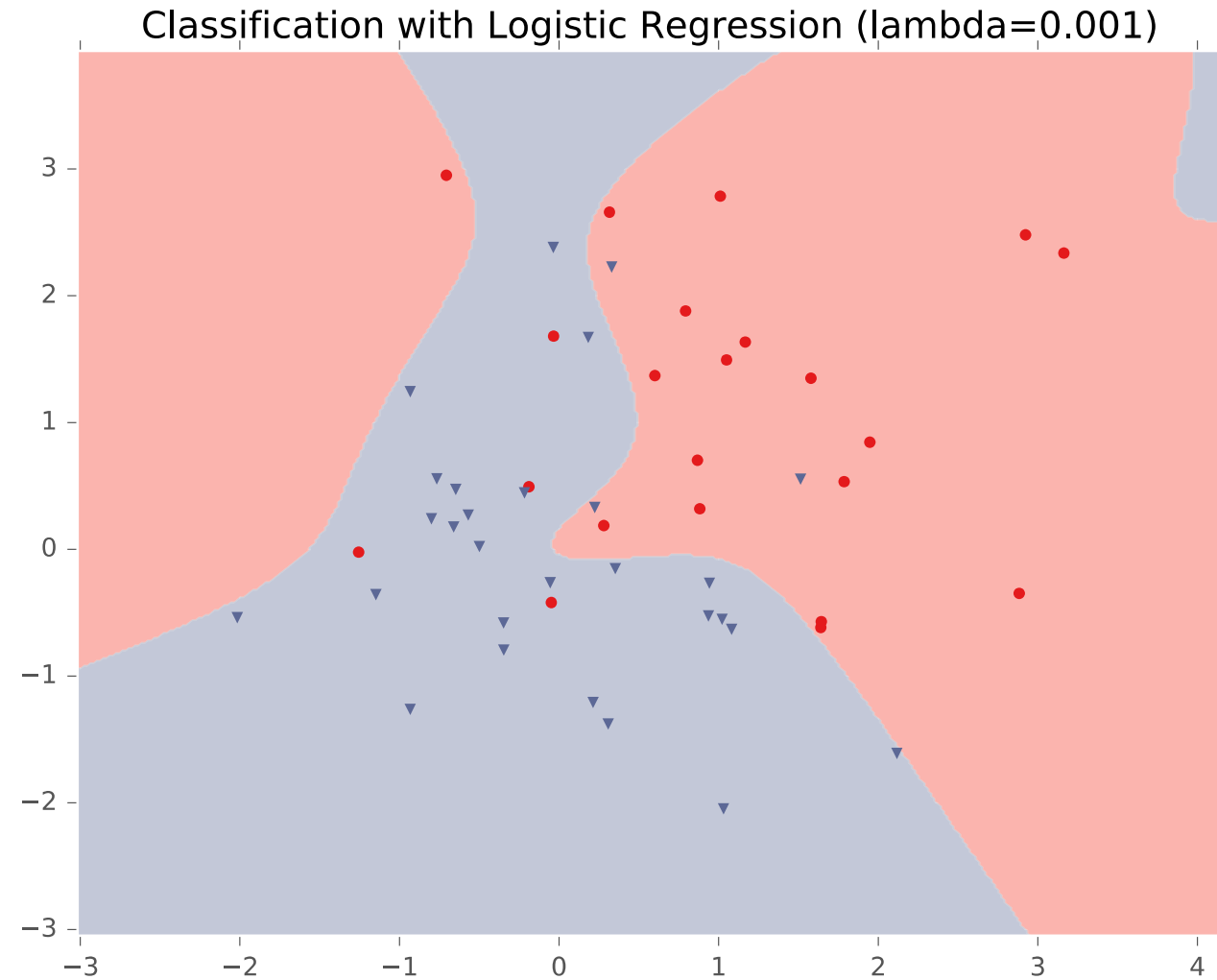
Example: Logistic Regression



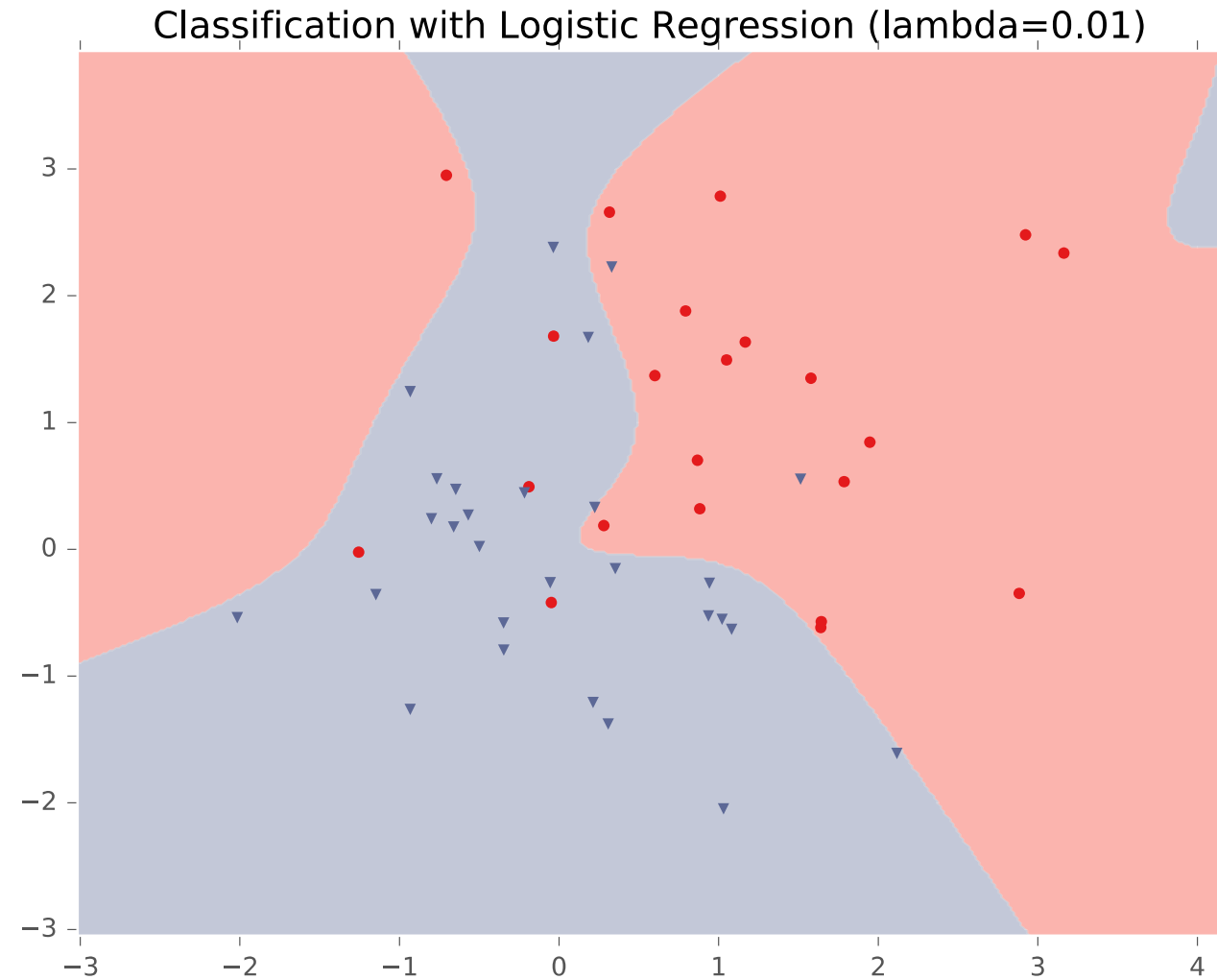
Example: Logistic Regression



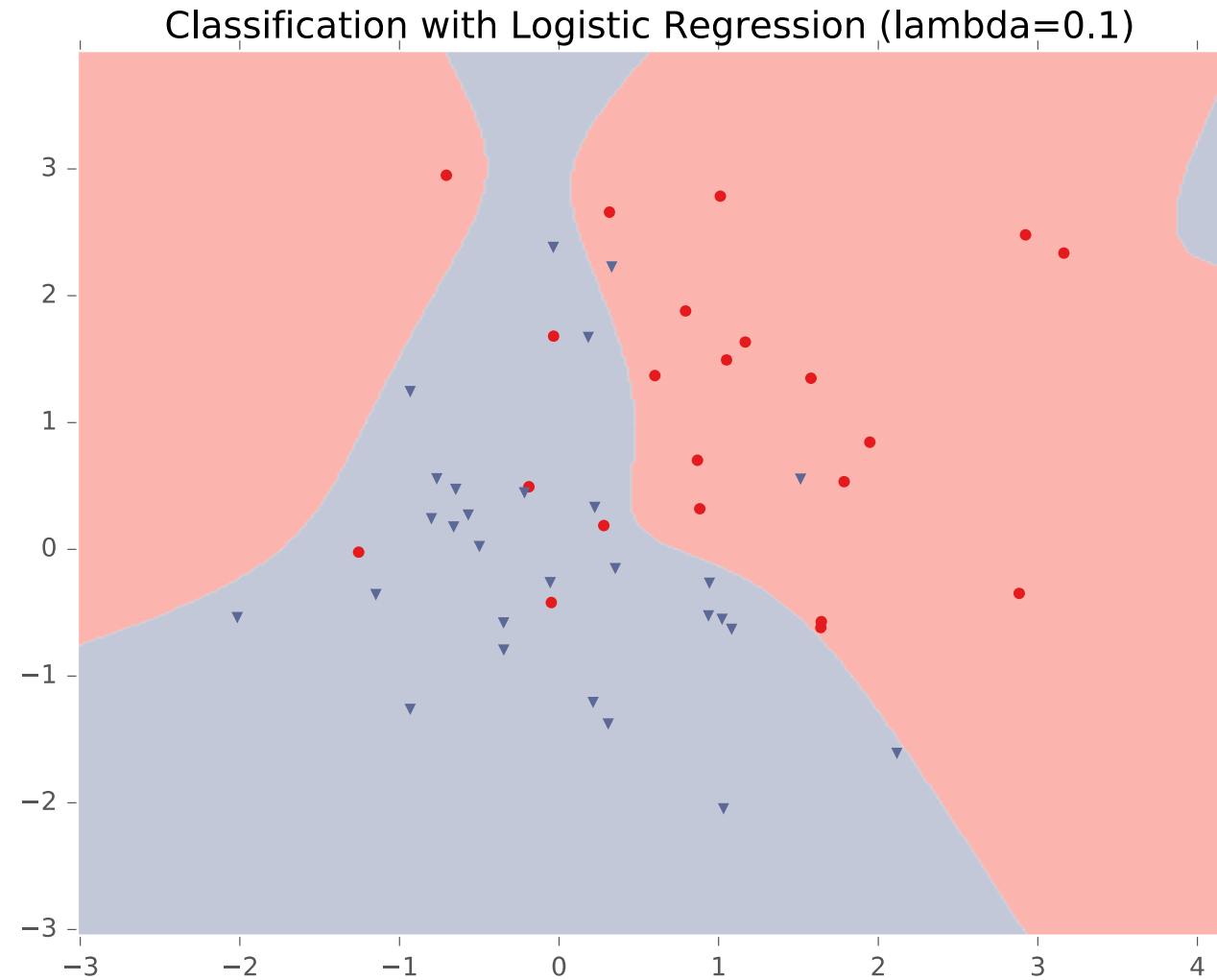
Example: Logistic Regression



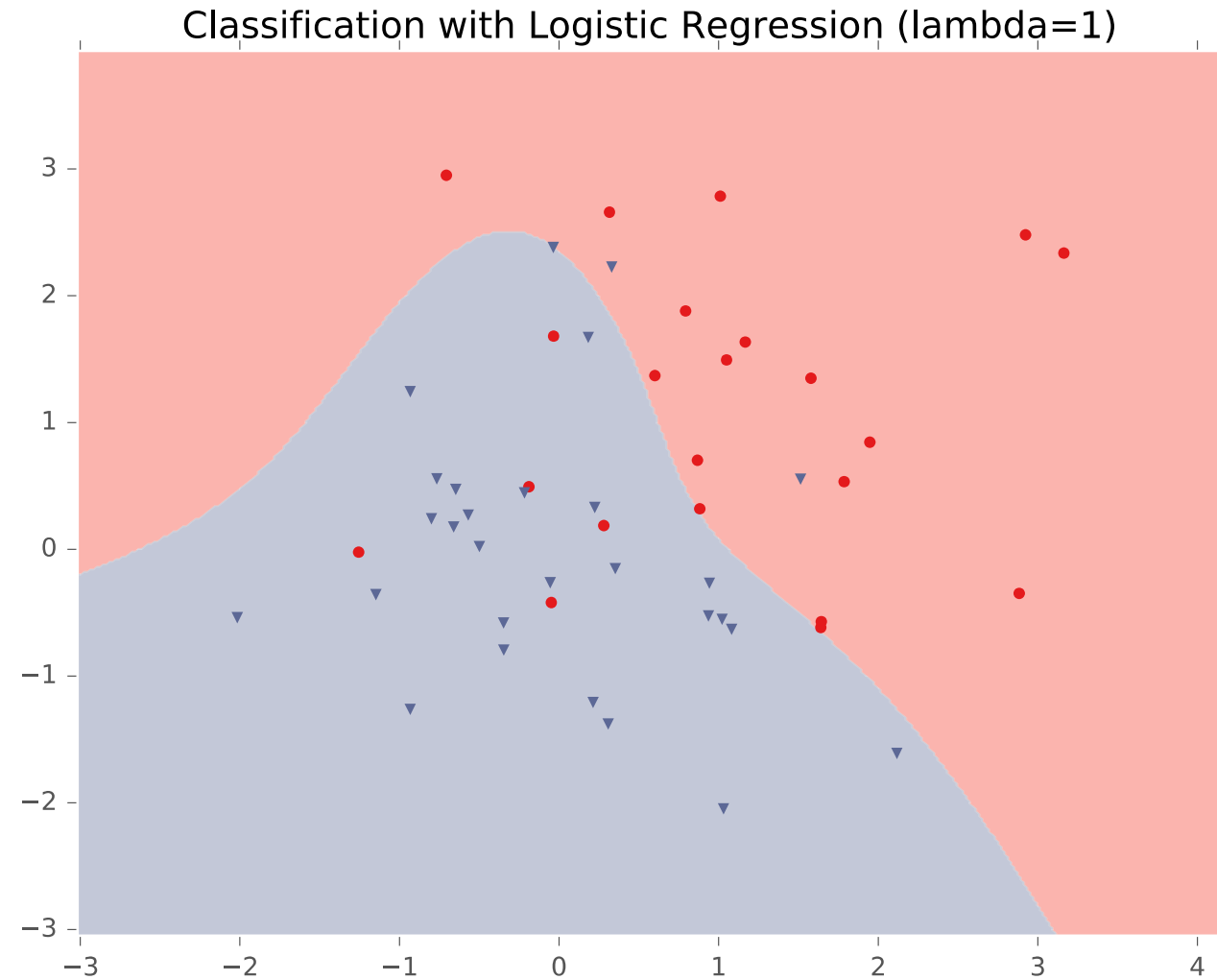
Example: Logistic Regression



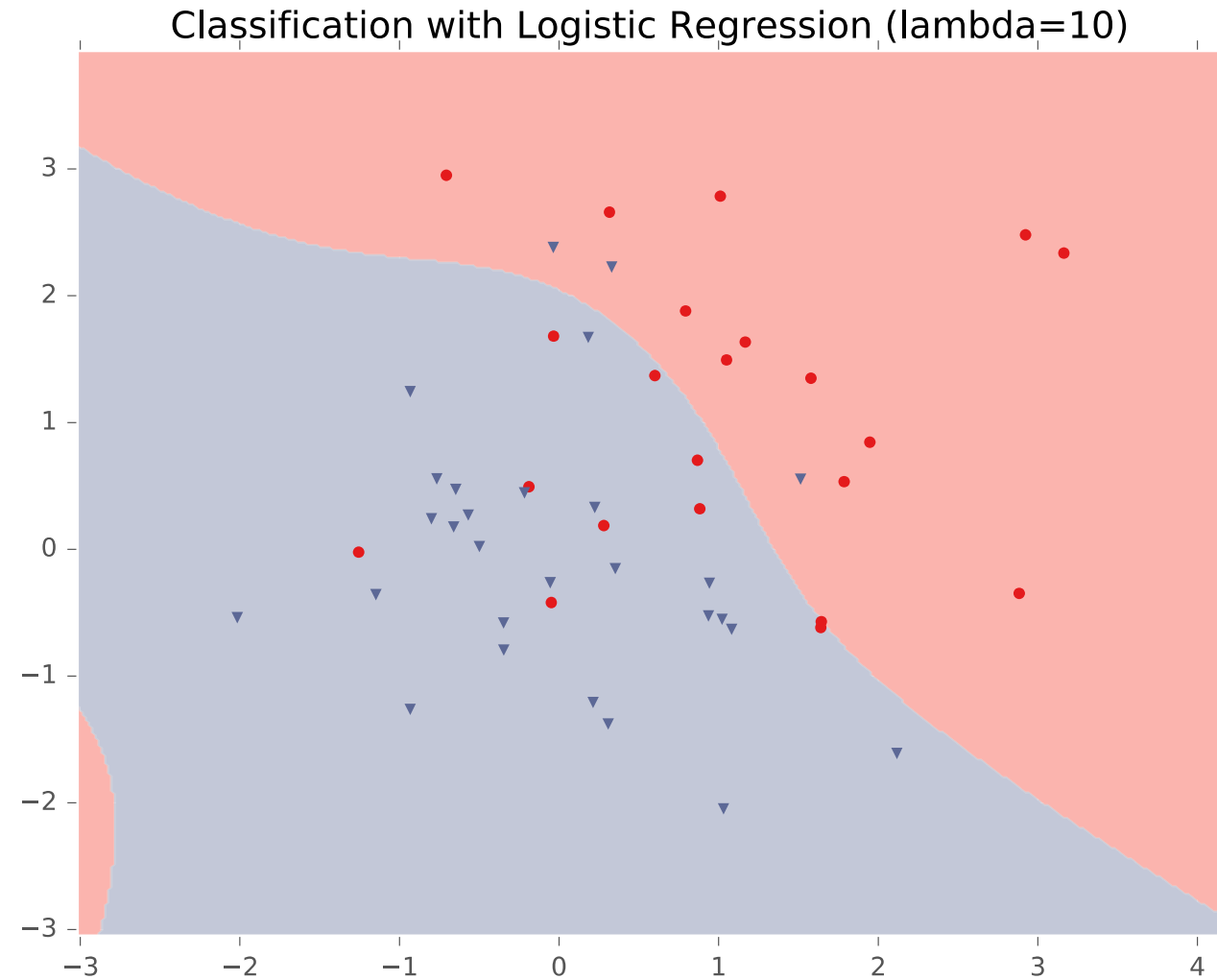
Example: Logistic Regression



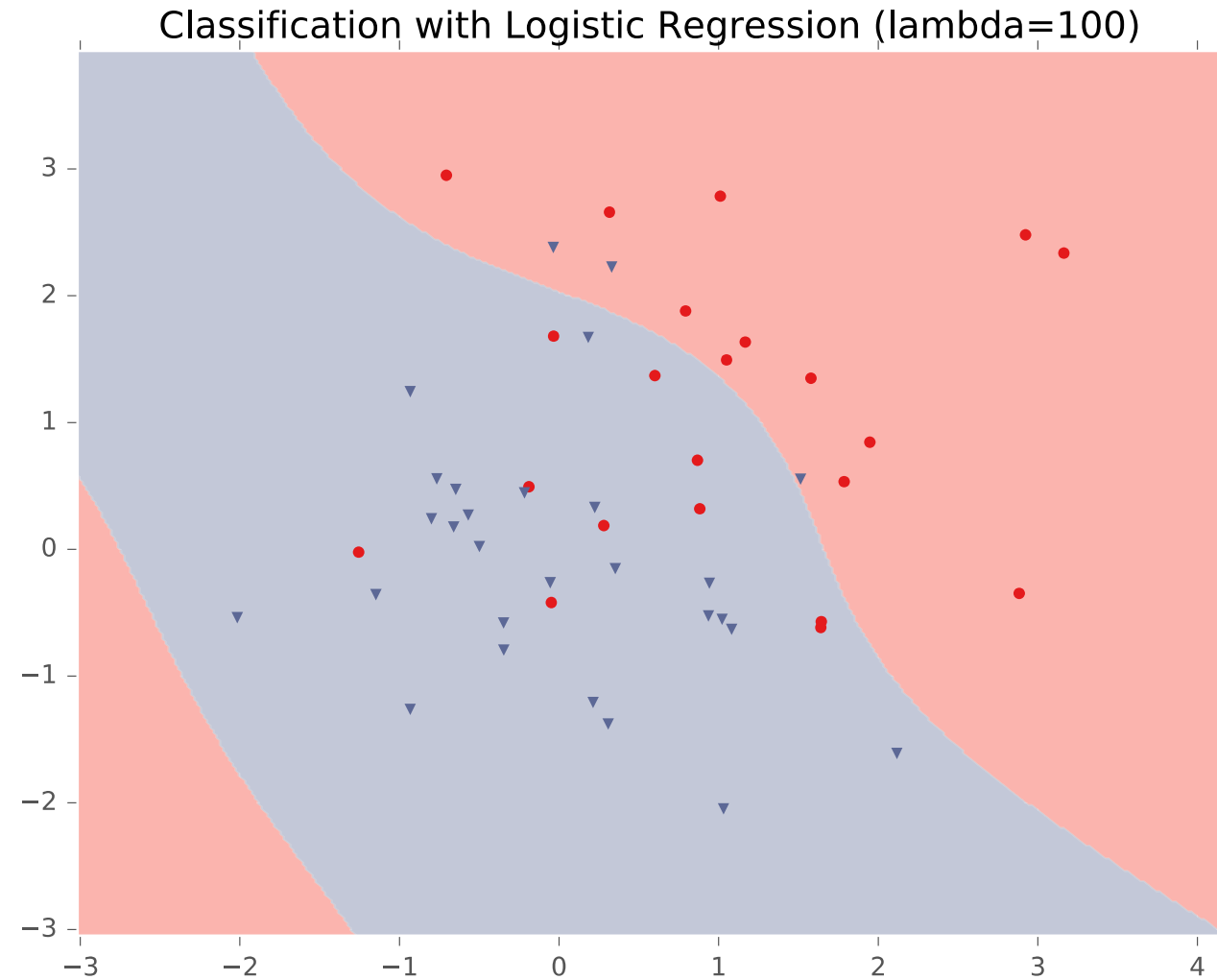
Example: Logistic Regression



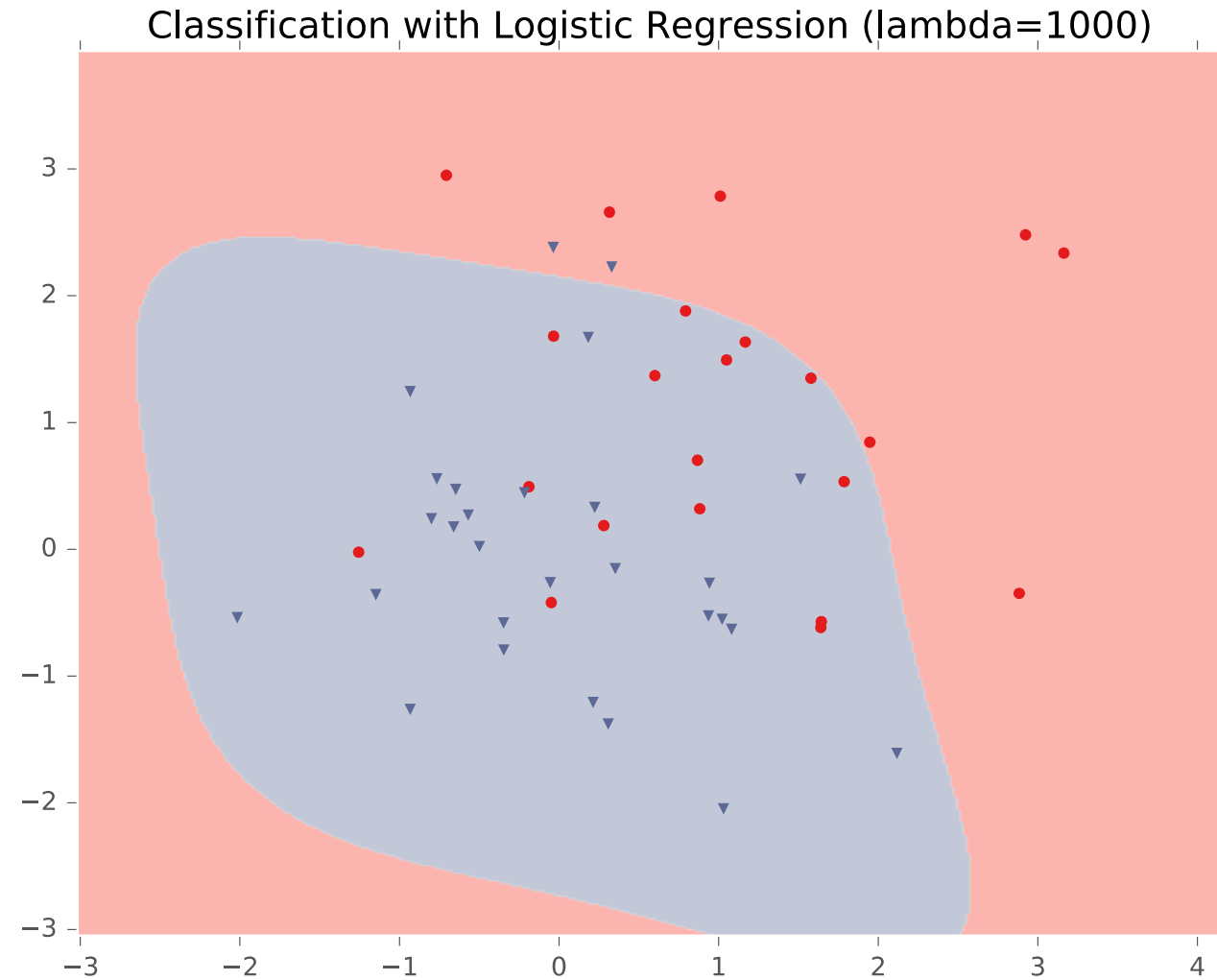
Example: Logistic Regression



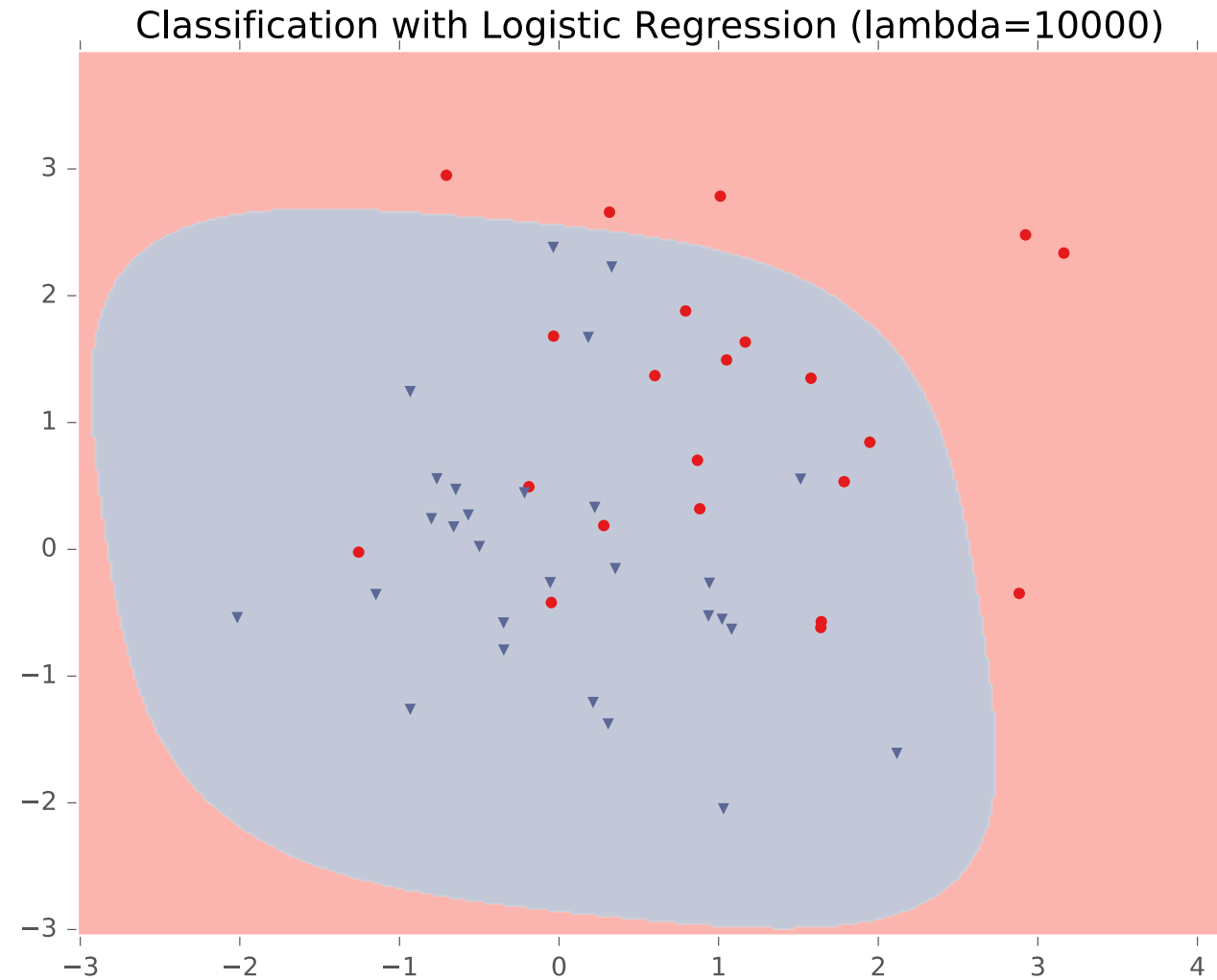
Example: Logistic Regression



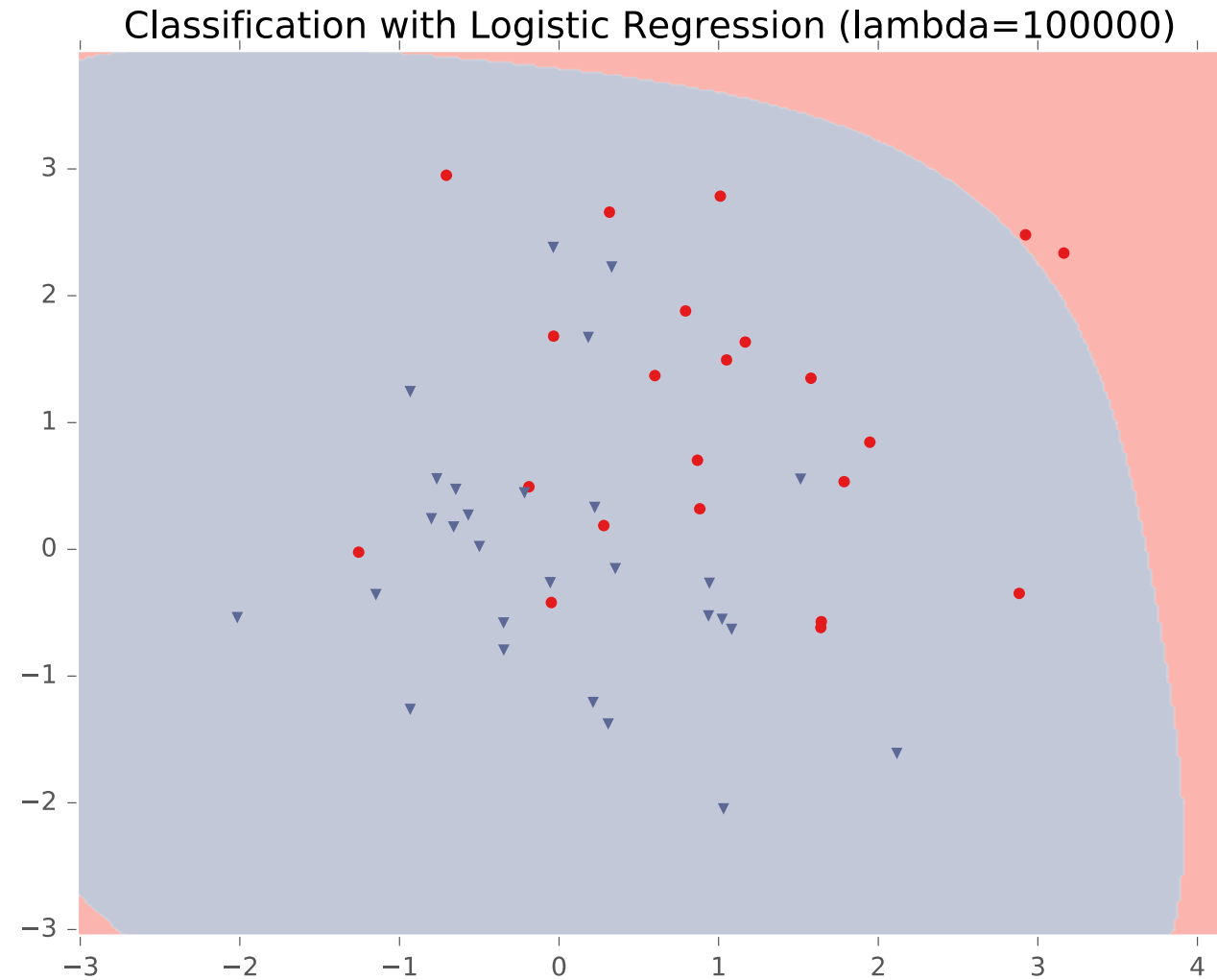
Example: Logistic Regression



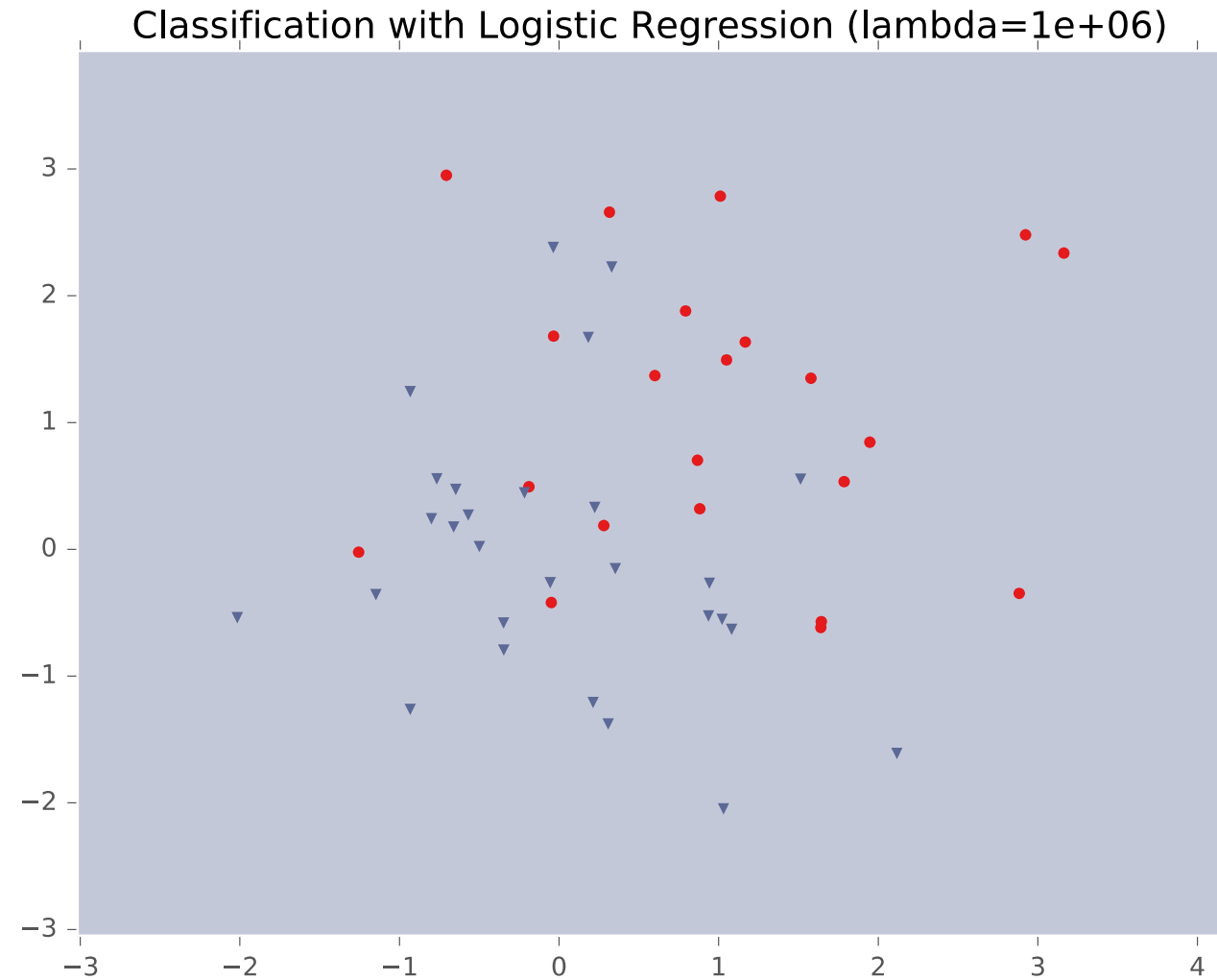
Example: Logistic Regression



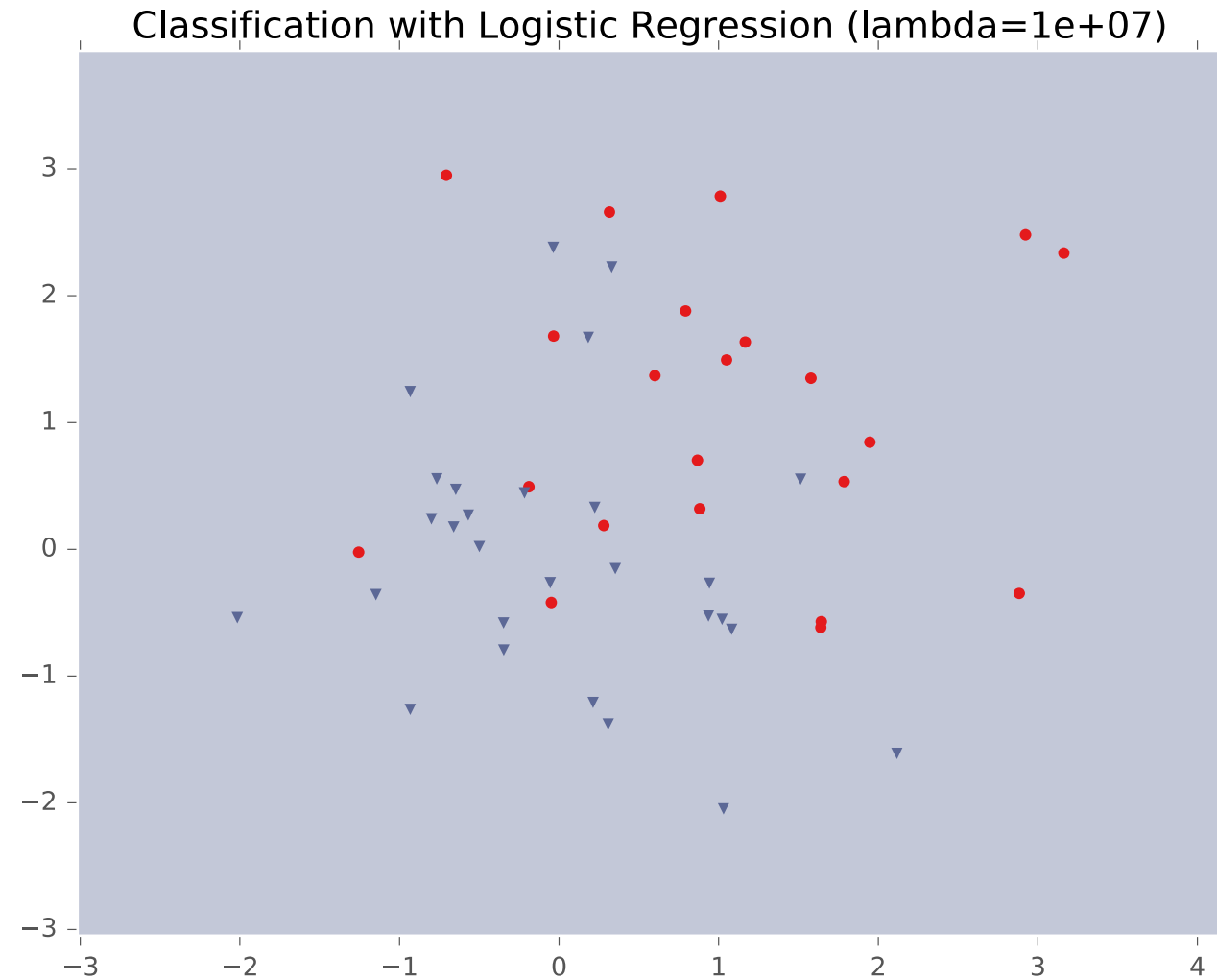
Example: Logistic Regression



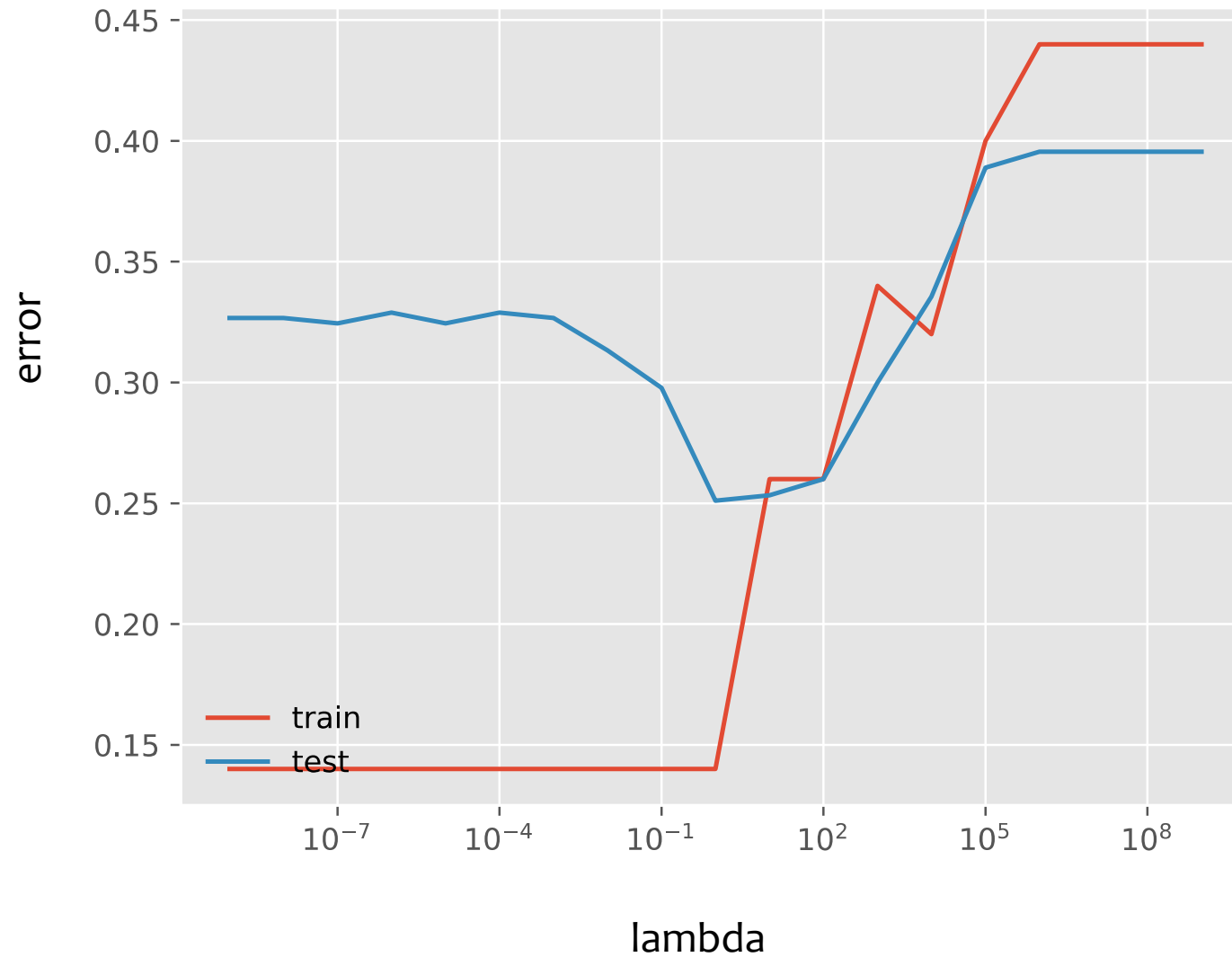
Example: Logistic Regression



Example: Logistic Regression



Example: Logistic Regression



Regularization

Given objective function: $J(\theta)$

Goal is to find: $\hat{\theta} = \underset{\theta}{\operatorname{argmin}} J(\theta) + \lambda r(\theta)$

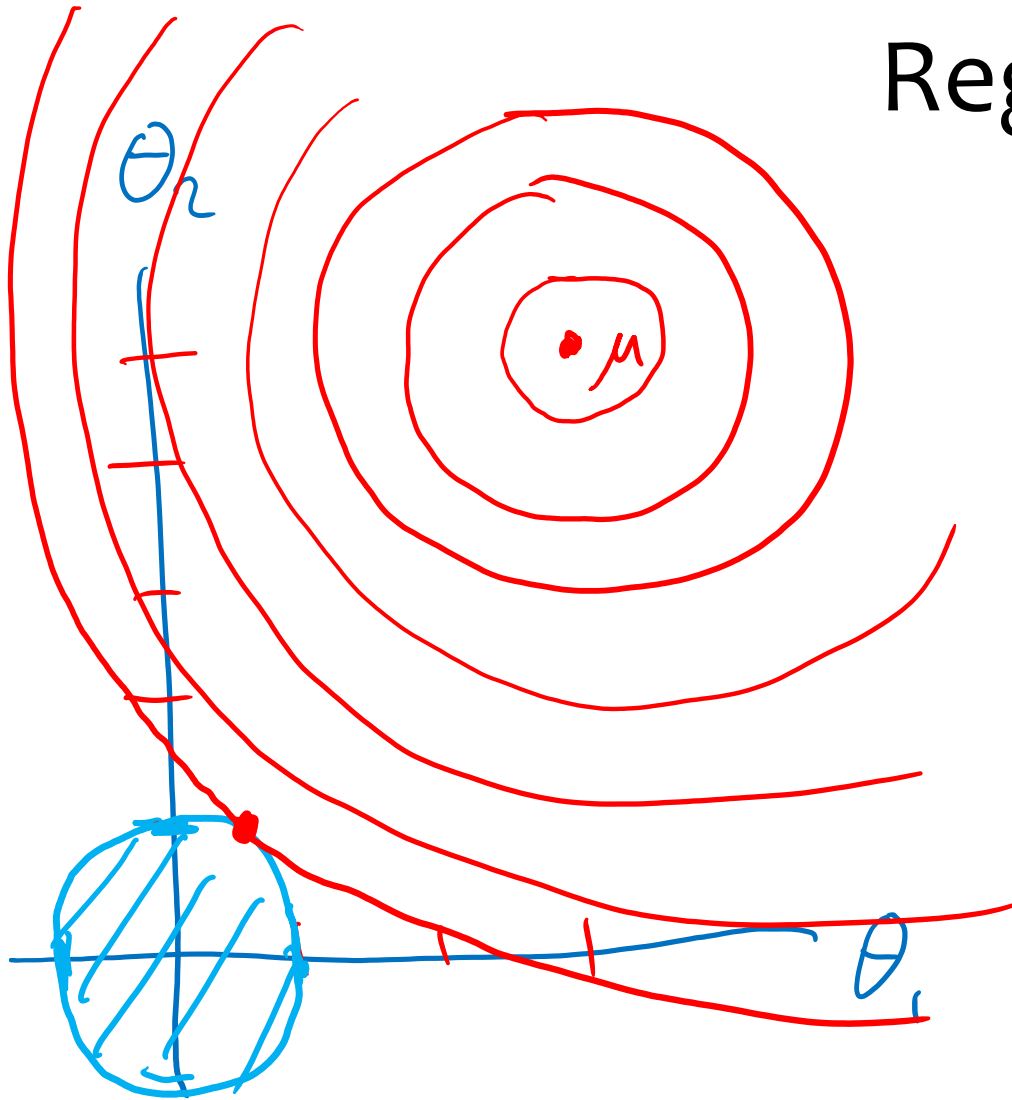
Key idea: Define regularizer $r(\theta)$ s.t. we tradeoff between fitting the data and keeping the model simple

Choose form of $r(\theta)$:

– Example: q-norm (usually p-norm) $r(\theta) = \|\theta\|_q = \left[\sum_{m=1}^M \|\theta_m\|^q \right]^{\left(\frac{1}{q}\right)}$

q	$r(\theta)$	yields parameters that are...	name	optimization notes
0	$\ \theta\ _0 = \sum \mathbb{1}(\theta_m \neq 0)$	zero values	L0 reg.	no good computational solutions
1	$\ \theta\ _1 = \sum \theta_m $	zero values	L1 reg.	subdifferentiable
2	$(\ \theta\ _2)^2 = \sum \theta_m^2$	small values	L2 reg.	differentiable

Regularization



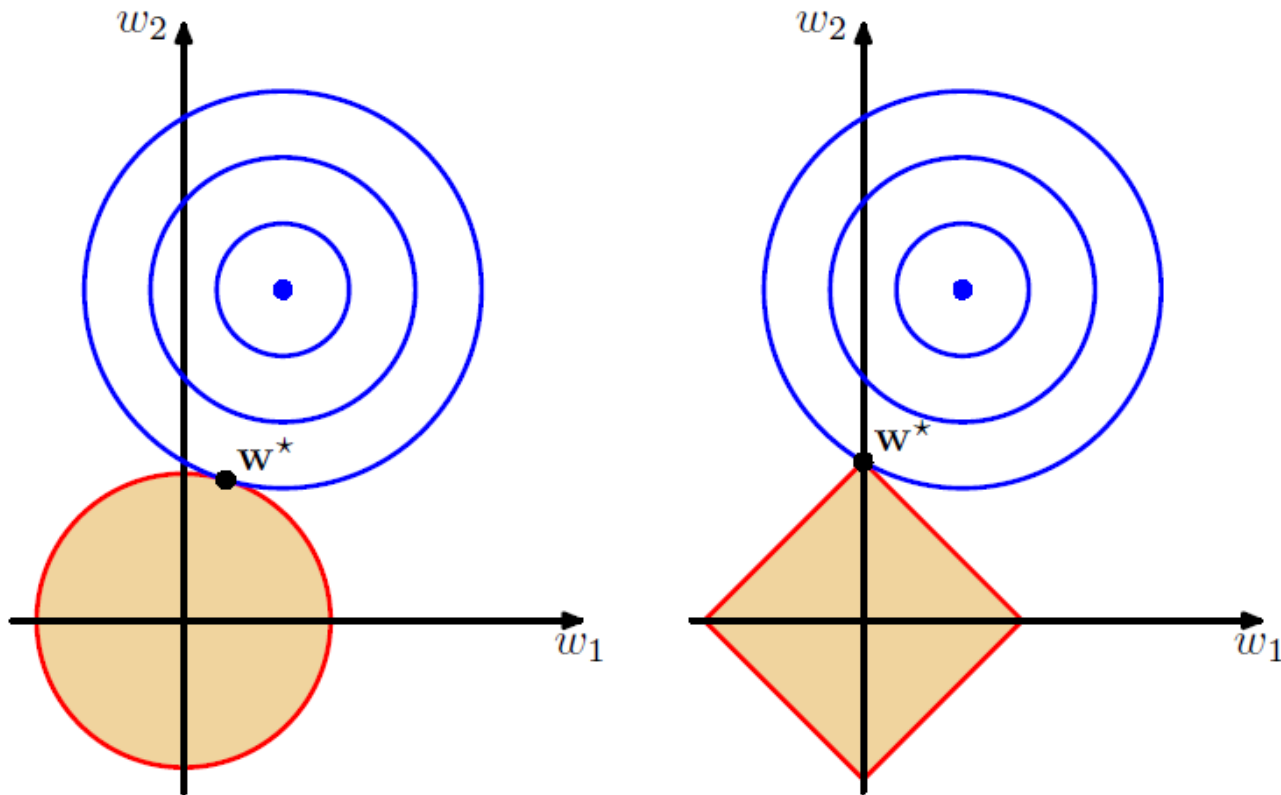
$$J(\theta_1, \theta_2) = \|\bar{\theta} - \vec{\mu}\| \quad \mu = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$$

$$\min_{\theta} J(\theta_1, \theta_2)$$

$$\text{s.t.} \quad \|\theta\|_2 \leq 1$$

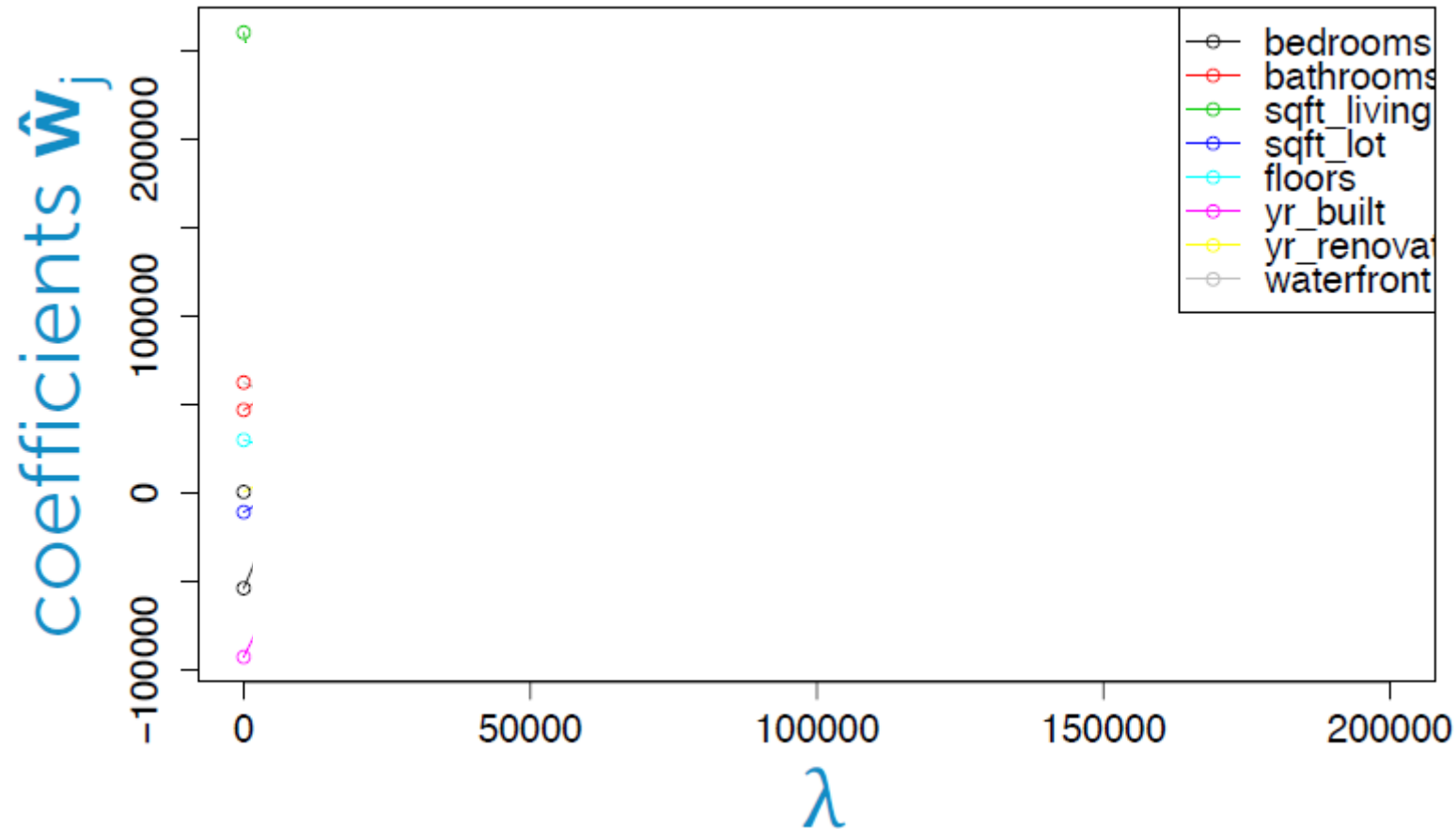
L2 vs L1 Regularization

Combine original objective with penalty on parameters



L2 vs L1: Housing Price Example

Predict housing price from several features



L2 vs L1: Housing Price Example

Predict housing price from several features

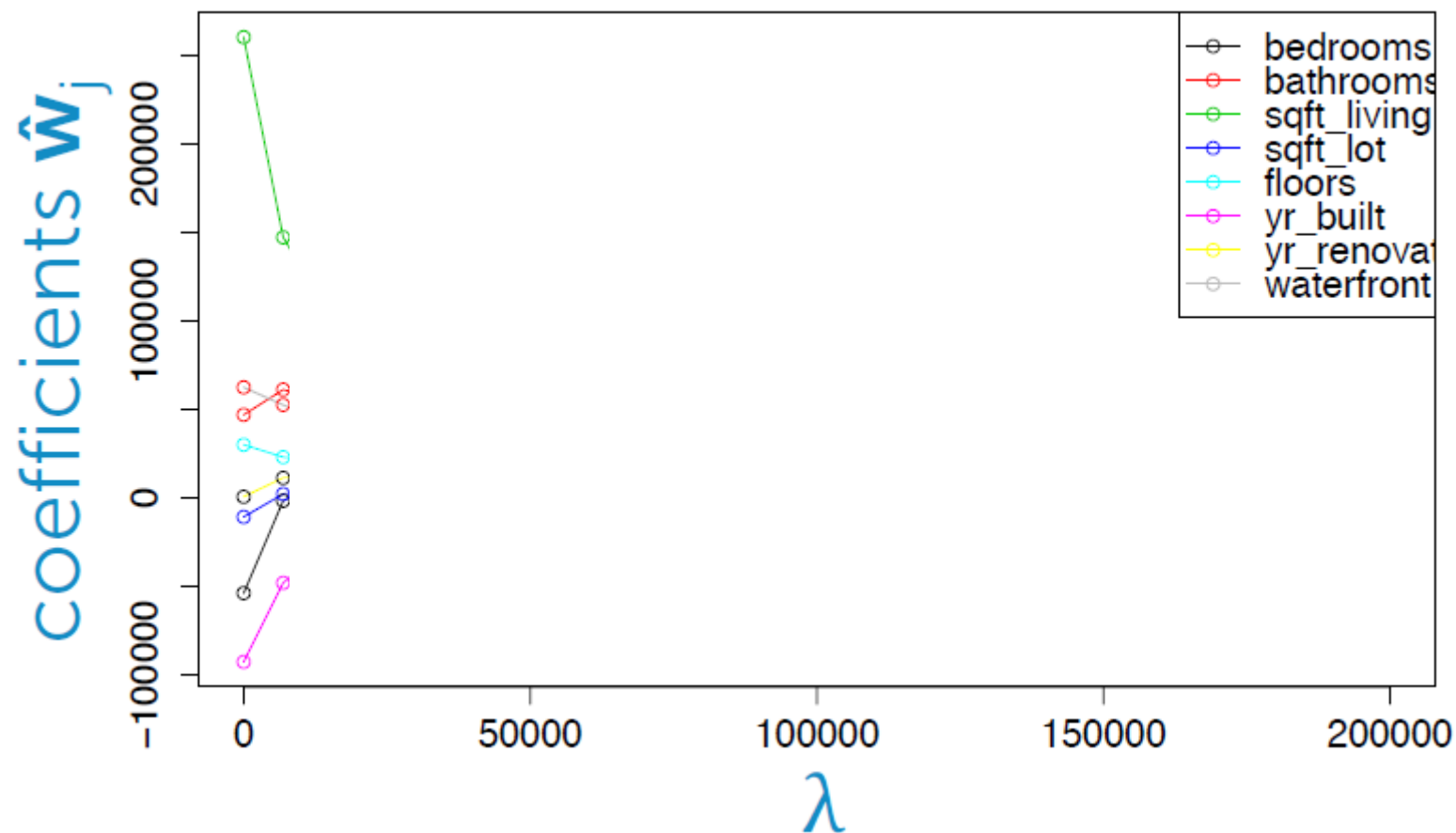


Figure: Emily Fox, University of Washington

L2 vs L1: Housing Price Example

Predict housing price from several features

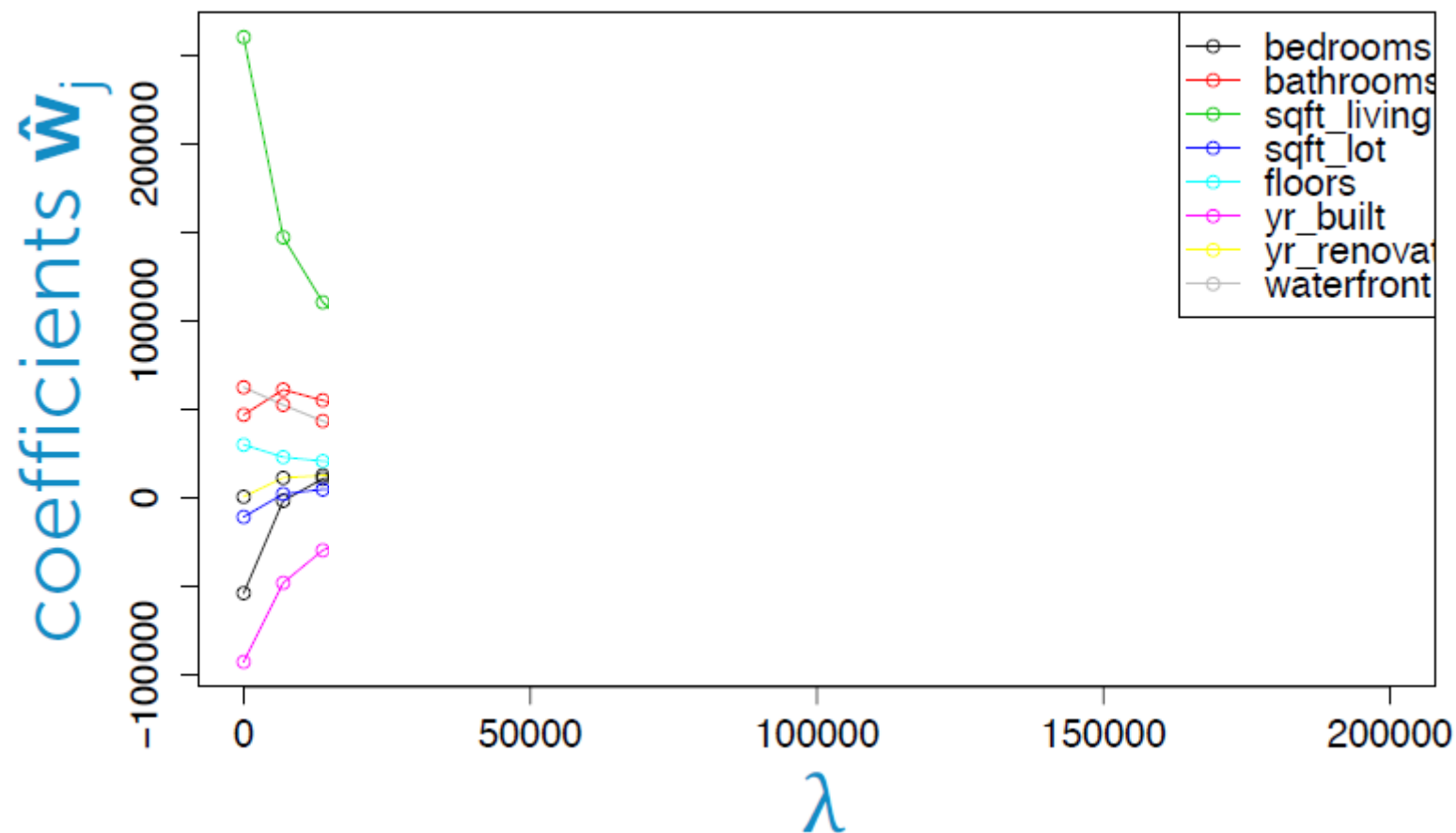


Figure: Emily Fox, University of Washington

L2 vs L1: Housing Price Example

Predict housing price from several features

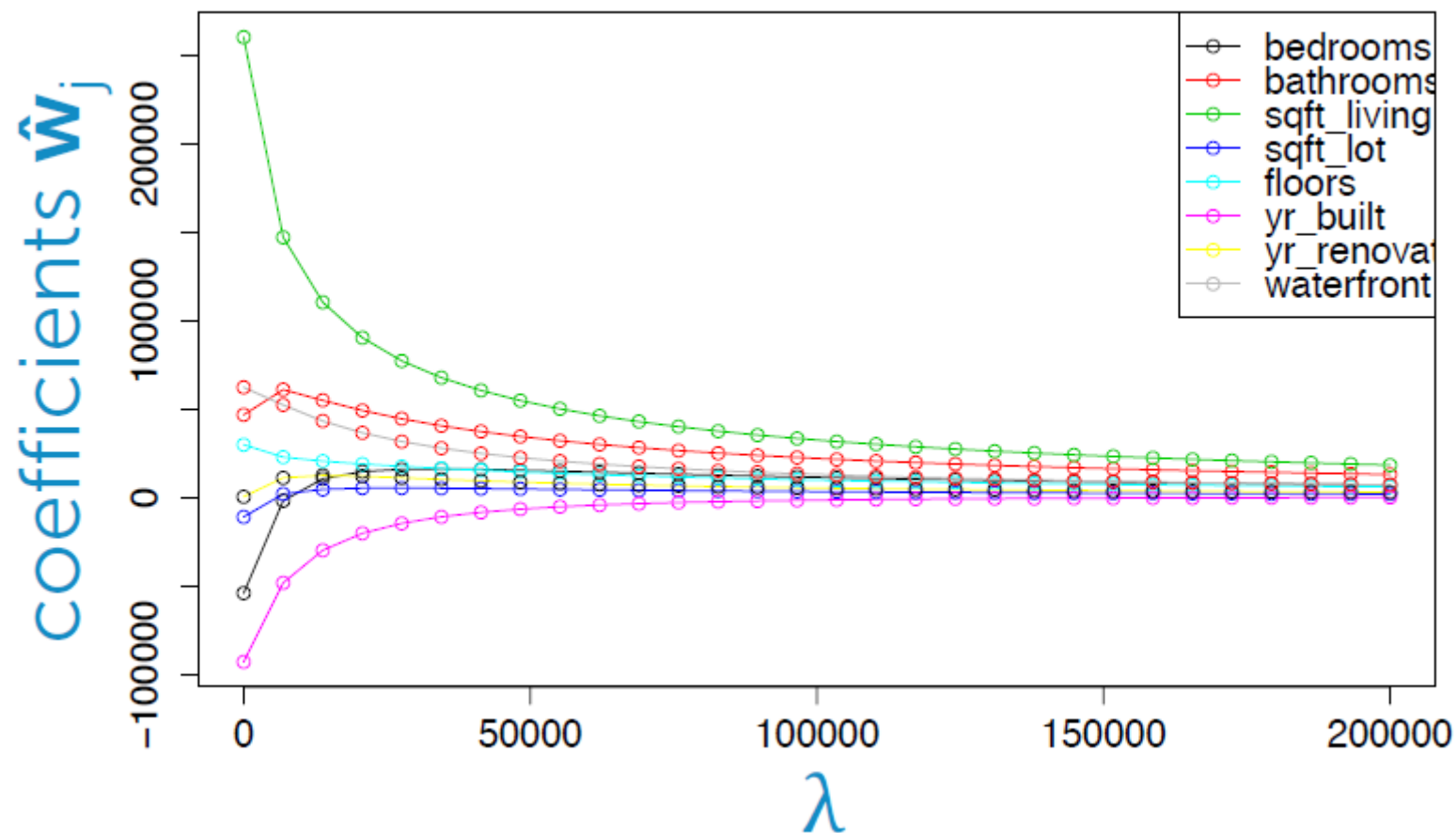
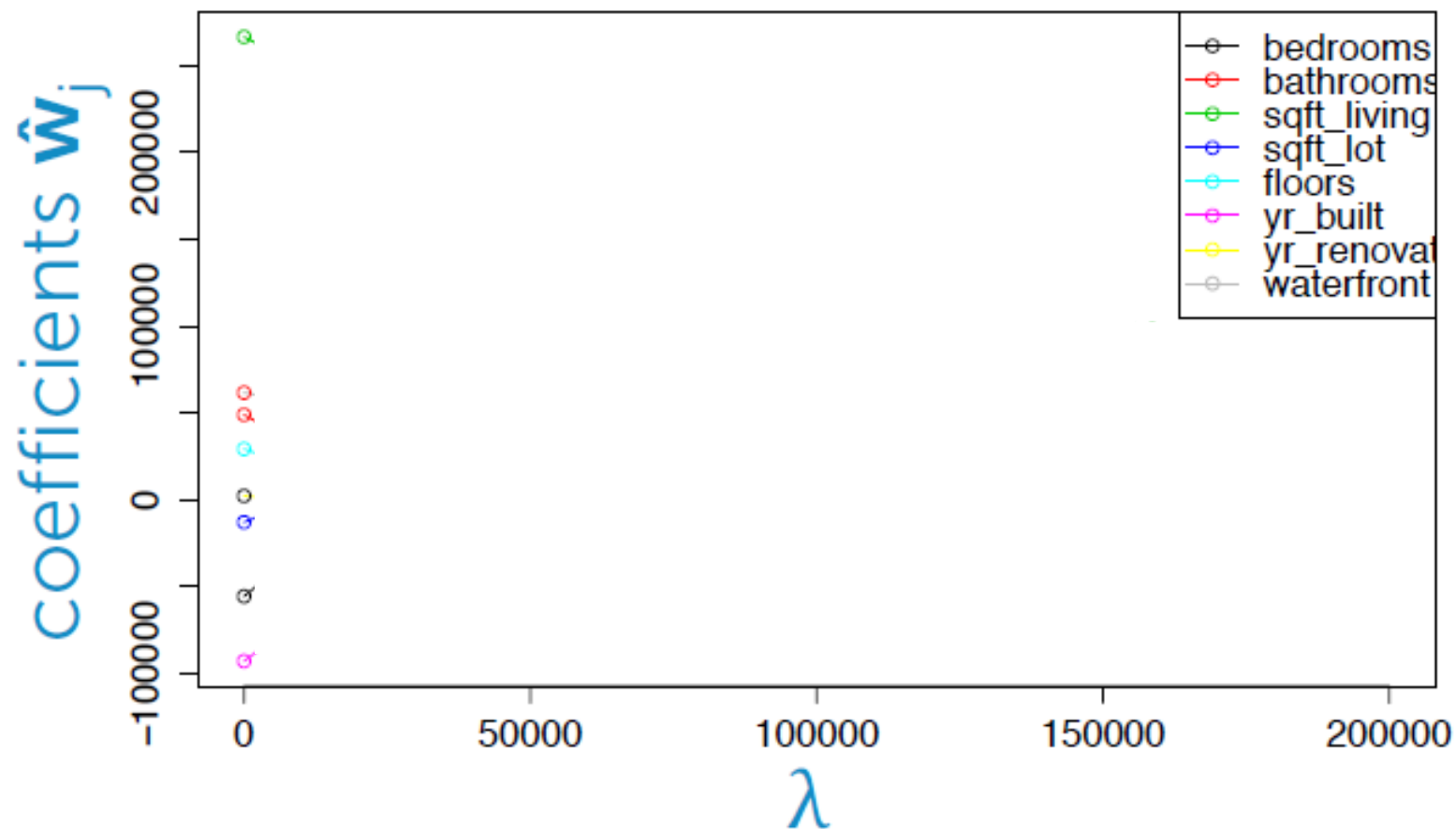


Figure: Emily Fox, University of Washington

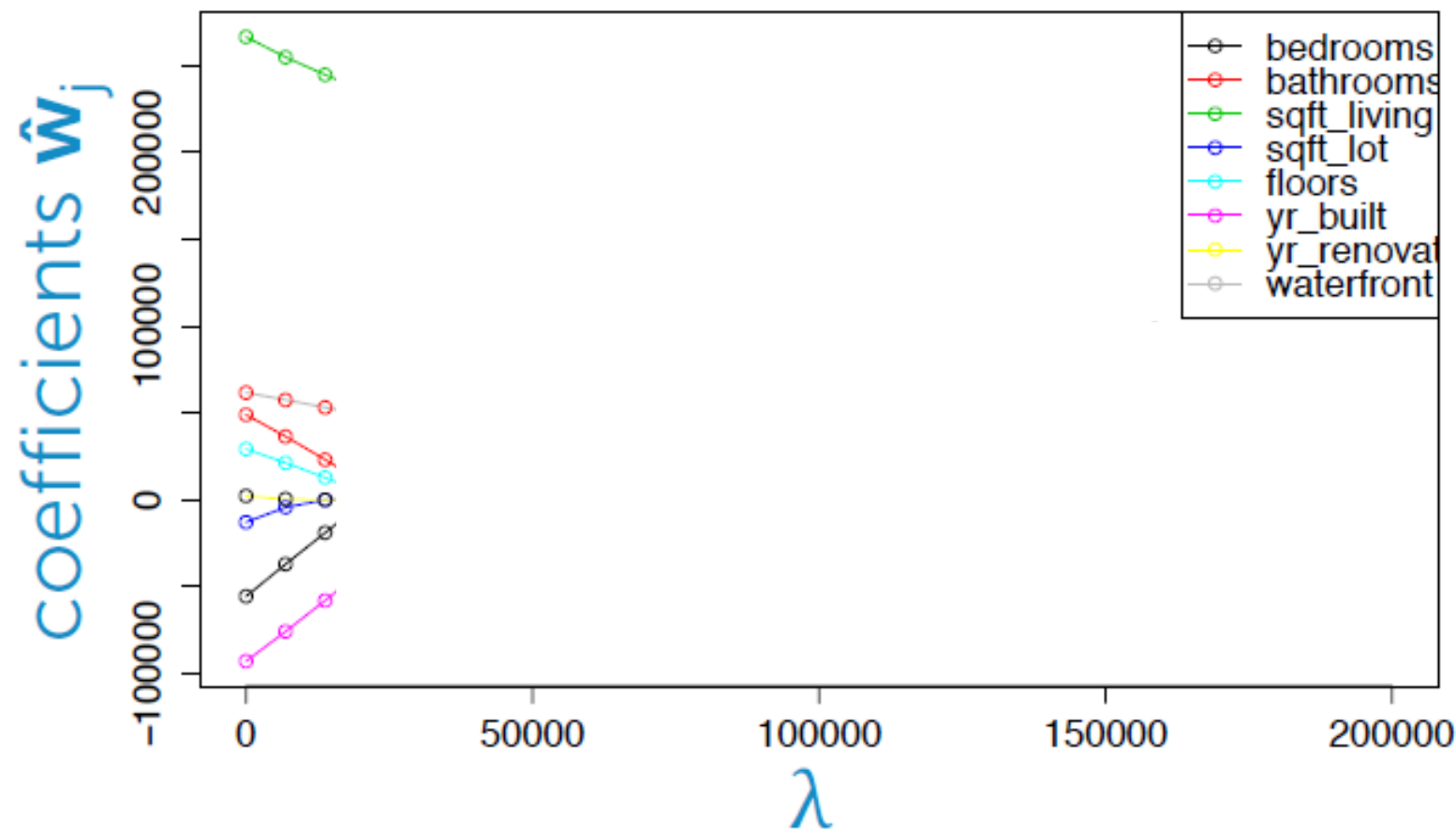
L2 vs L1: Housing Price Example

Predict housing price from several features



L2 vs L1: Housing Price Example

Predict housing price from several features



L2 vs L1: Housing Price Example

Predict housing price from several features

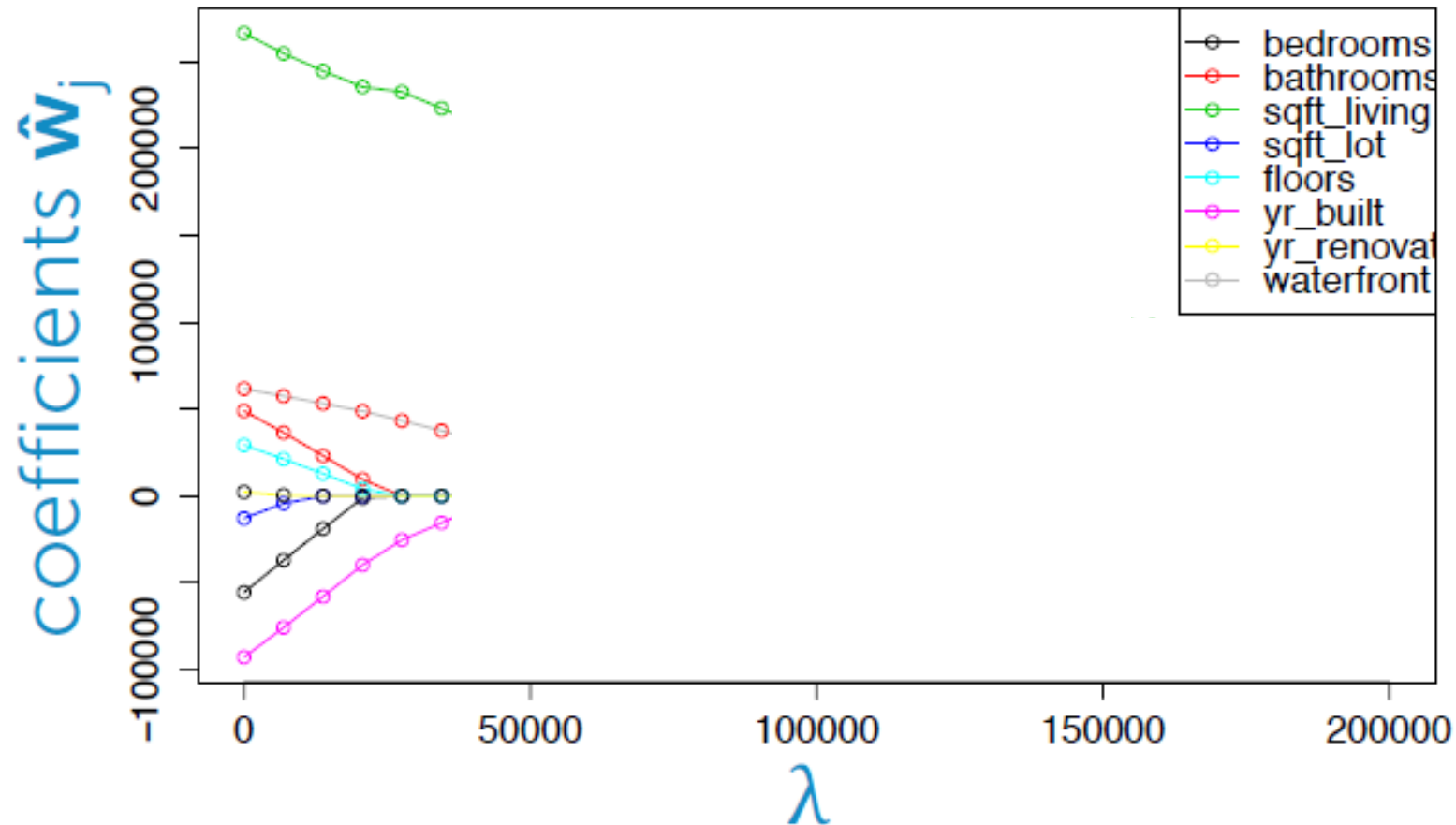
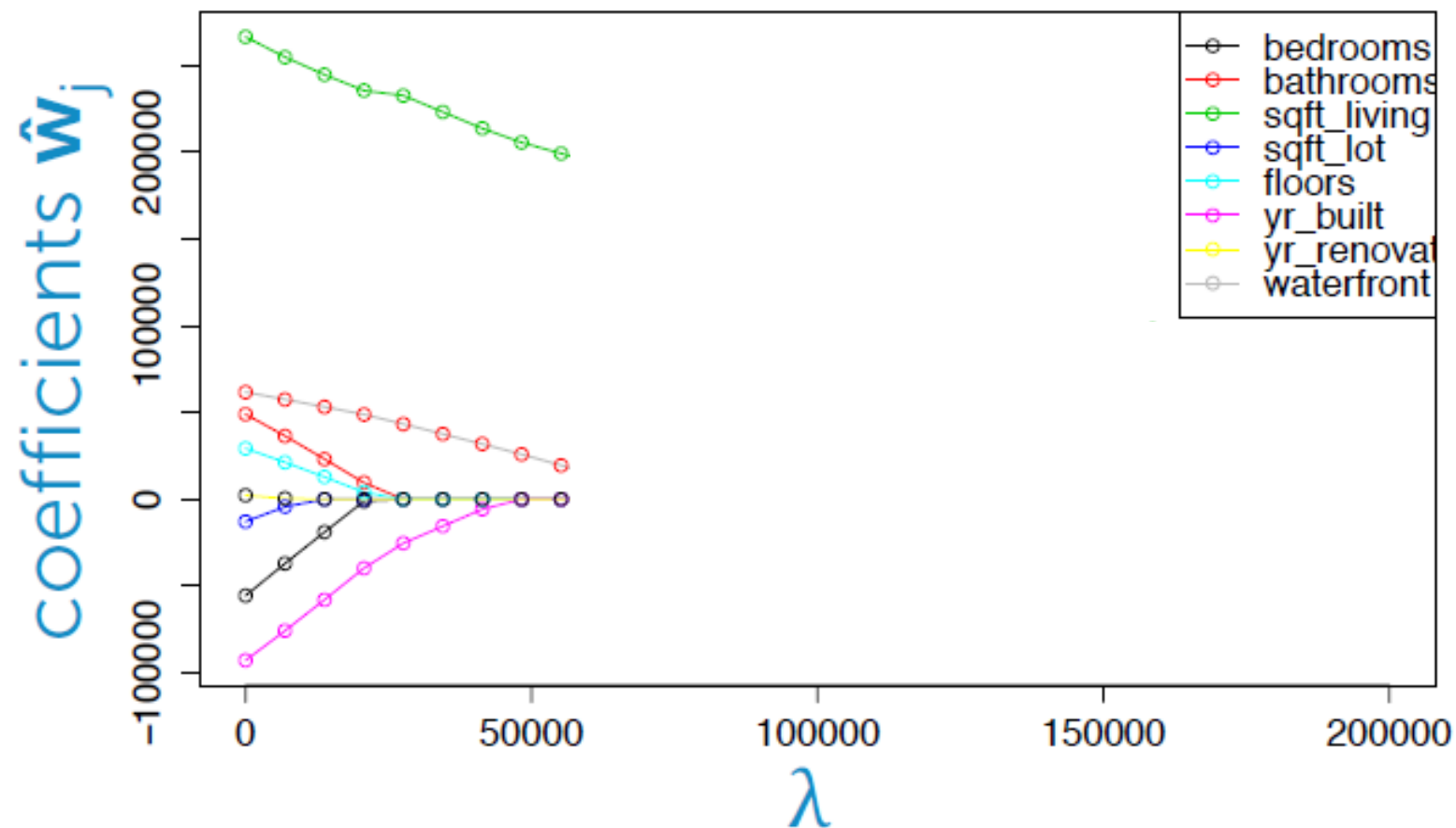


Figure: Emily Fox, University of Washington

L2 vs L1: Housing Price Example

Predict housing price from several features



L2 vs L1: Housing Price Example

Predict housing price from several features

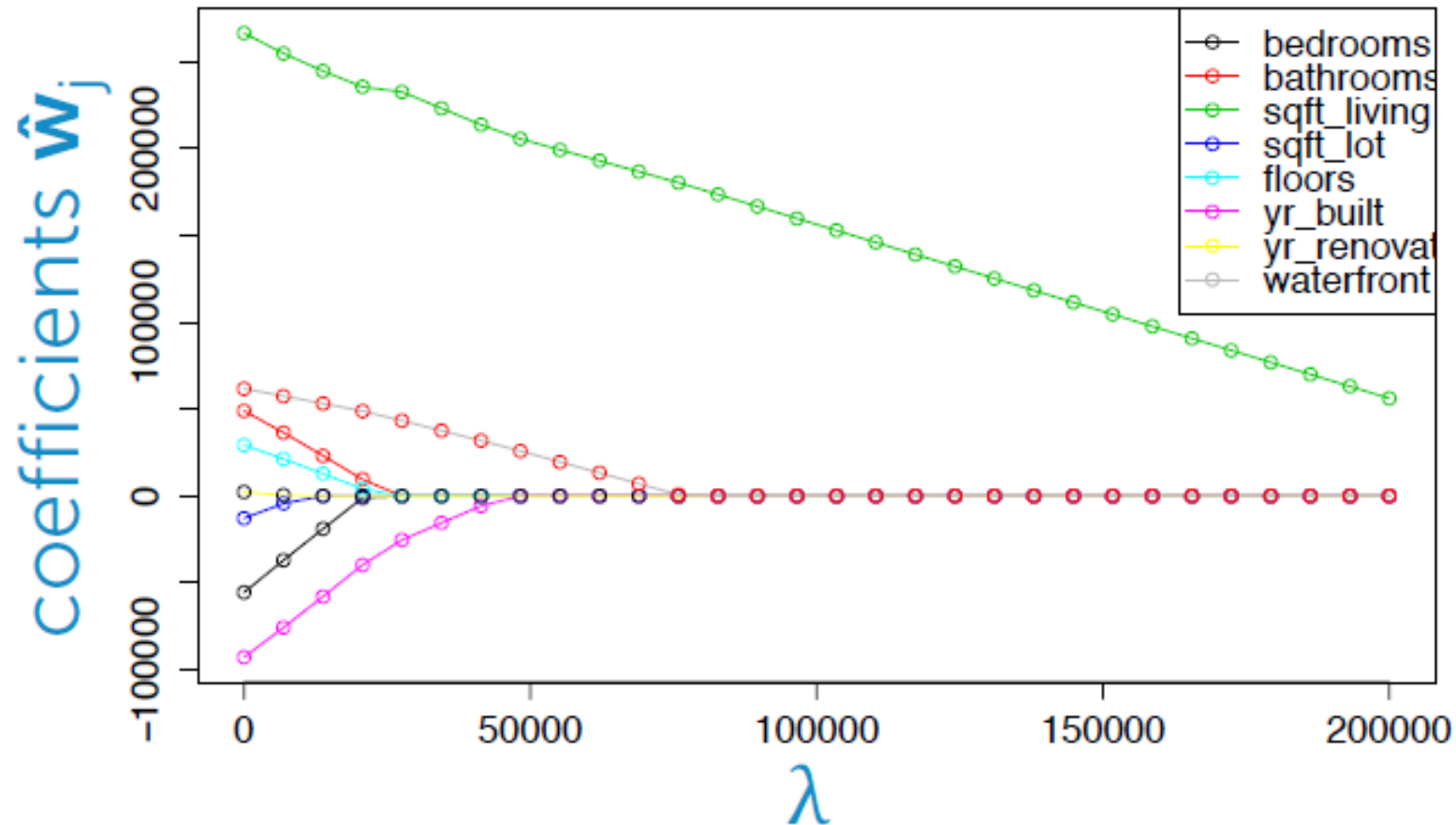


Figure: Emily Fox, University of Washington

Regularization as MAP

- L1 and L2 regularization can be interpreted as **maximum a-posteriori (MAP) estimation** of the parameters
- To be discussed later in the course...

Takeaways

1. **Nonlinear basis functions** allow **linear models** (e.g. Linear Regression, Logistic Regression) to capture **nonlinear** aspects of the original input
2. Nonlinear features **require no changes to the model** (i.e. just preprocessing)
3. **Regularization** helps to avoid **overfitting**
4. (**Regularization** and **MAP estimation** are equivalent for appropriately chosen priors)

Feature Engineering / Regularization Objectives

You should be able to...

- Engineer appropriate features for a new task
- Use feature selection techniques to identify and remove irrelevant features
- Identify when a model is overfitting
- Add a regularizer to an existing objective in order to combat overfitting
- Convert linearly inseparable dataset to a linearly separable dataset in higher dimensions
- Describe feature engineering in common application areas

“Dollar fifty in late charges at the public library”

<https://www.youtube.com/watch?v=e1DnltskkWk>

Breakout rooms: Why is a CMU course better than a library card?

An abstract graphic on the left side of the slide, featuring a sphere-like shape composed of a dense grid of intersecting red, green, and blue lines. The lines are curved and follow the contour of the sphere, creating a complex, woven pattern. The sphere is set against a dark gray background.

Introduction to Machine Learning

Neural Networks

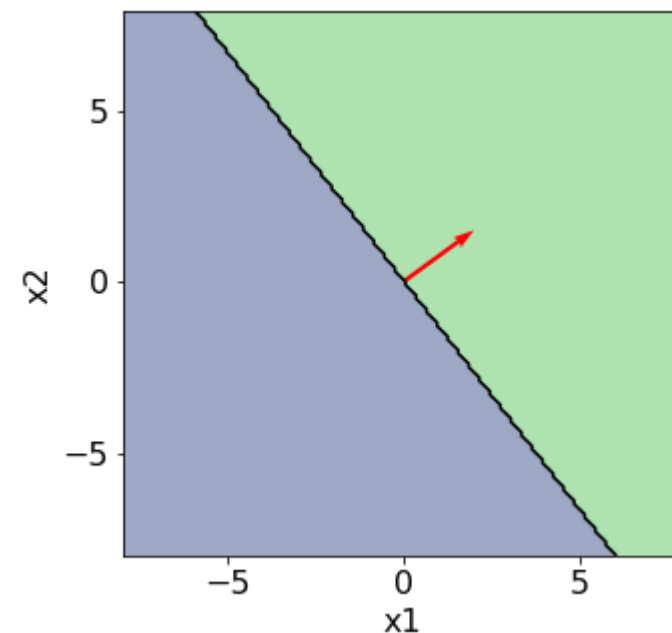
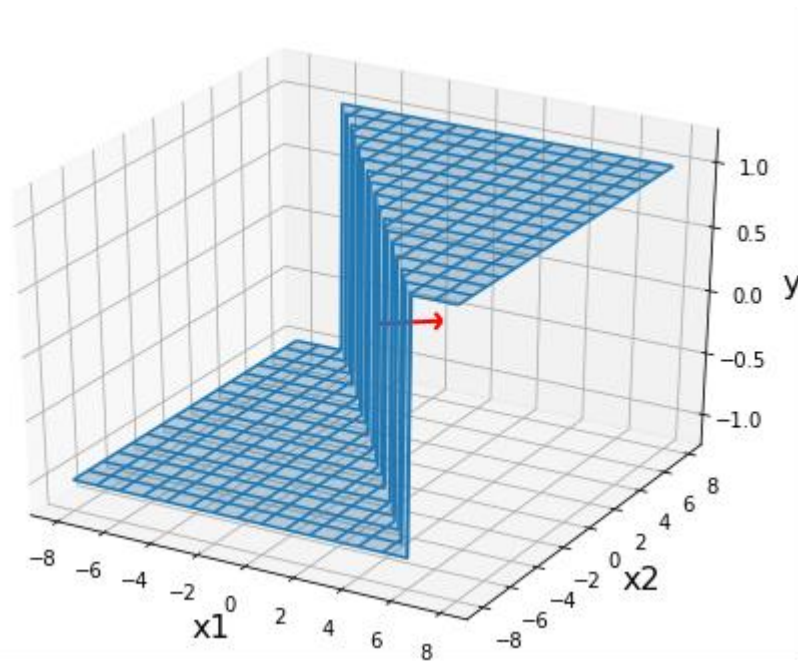
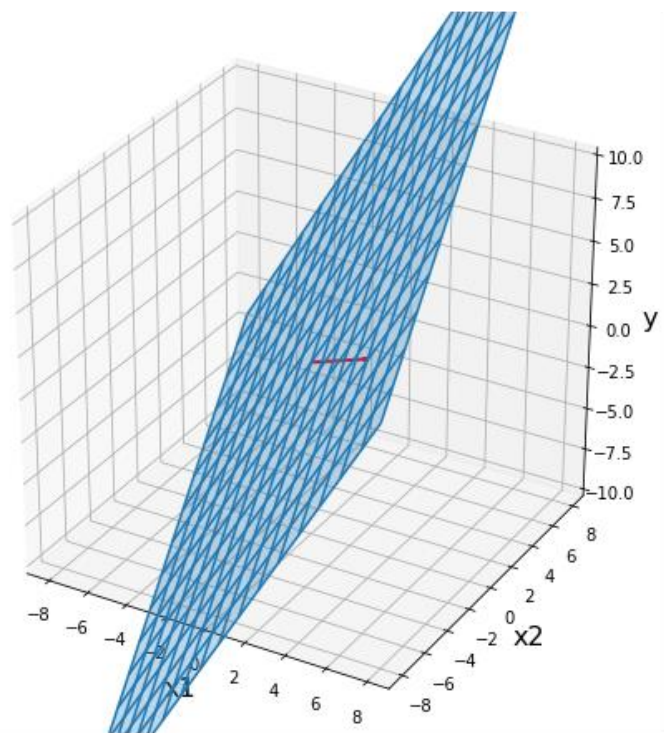
Instructor: Pat Virtue

Perceptron

Classification: Hard threshold on linear model

$$h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} + b)$$

$$\text{sign}(z) = \begin{cases} 1, & \text{if } z \geq 0 \\ -1, & \text{if } z < 0 \end{cases}$$



Piazza Poll 2

Which of the following perceptron parameters will perfectly classify this data?

$$h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} + b)$$

$$\text{sign}(z) = \begin{cases} 1, & \text{if } z \geq 0 \\ -1, & \text{if } z < 0 \end{cases}$$

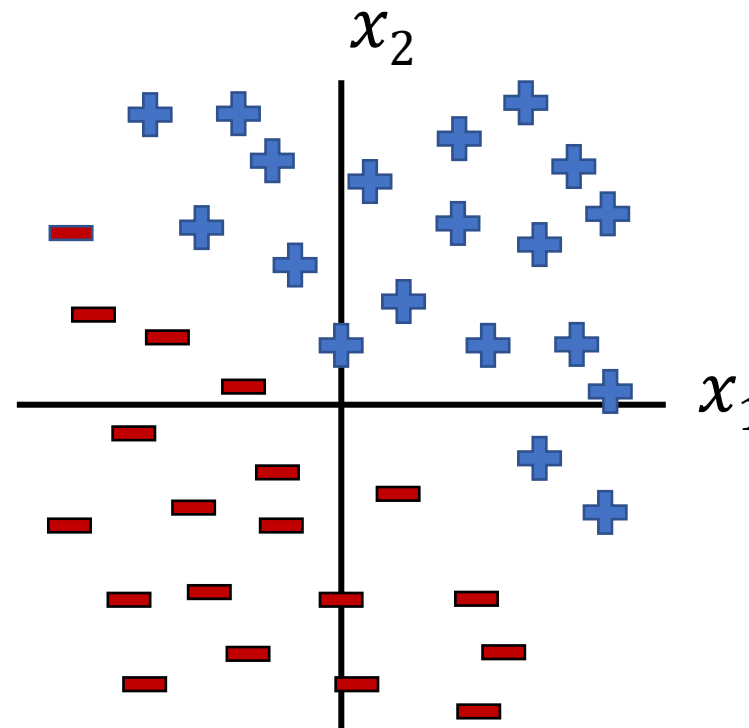
A. $w = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, b = 0$

B. $w = \begin{bmatrix} -1 \\ 1 \end{bmatrix}, b = 0$

C. $w = \begin{bmatrix} 1 \\ -1 \end{bmatrix}, b = 0$

D. $w = \begin{bmatrix} -1 \\ -1 \end{bmatrix}, b = 0$

E. None of the above



Piazza Poll 3

Which of the following perceptron parameters will perfectly classify this data?

$$h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} + b)$$

$$\text{sign}(z) = \begin{cases} 1, & \text{if } z \geq 0 \\ -1, & \text{if } z < 0 \end{cases}$$

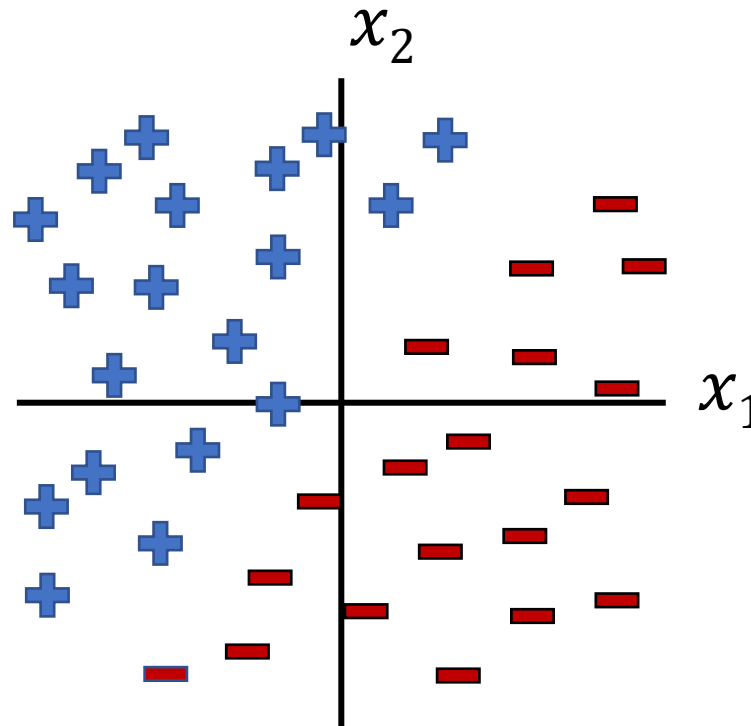
A. $w = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, b = 0$

B. $w = \begin{bmatrix} -1 \\ 1 \end{bmatrix}, b = 0$

C. $w = \begin{bmatrix} 1 \\ -1 \end{bmatrix}, b = 0$

D. $w = \begin{bmatrix} -1 \\ -1 \end{bmatrix}, b = 0$

E. None of the above



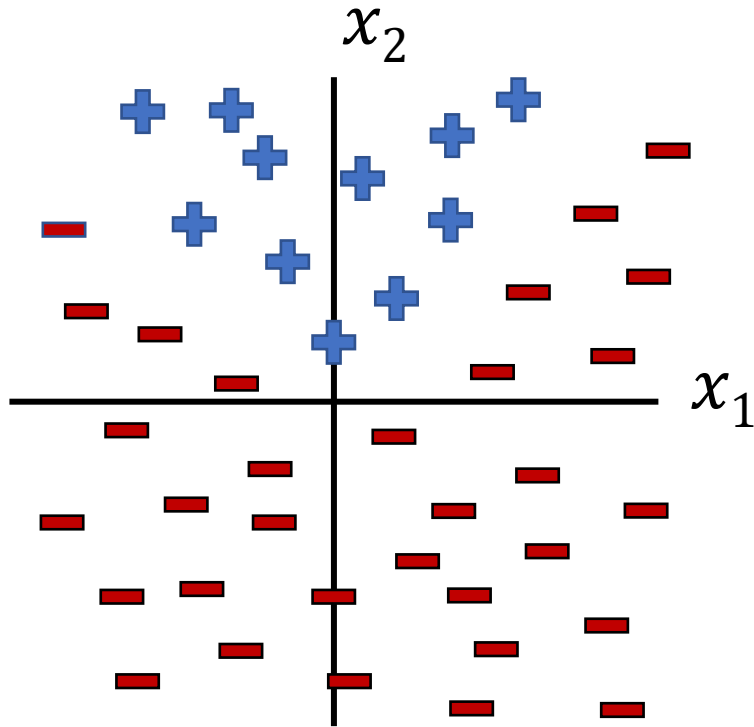
Classification Design Challenge

How could you configure three specific perceptrons to classify this data?

$$h_A(\mathbf{x}) = \text{sign}(\mathbf{w}_A^T \mathbf{x} + b_A)$$

$$h_B(\mathbf{x}) = \text{sign}(\mathbf{w}_B^T \mathbf{x} + b_B)$$

$$h_C(\mathbf{x}) = \text{sign}(\mathbf{w}_C^T \mathbf{x} + b_C)$$



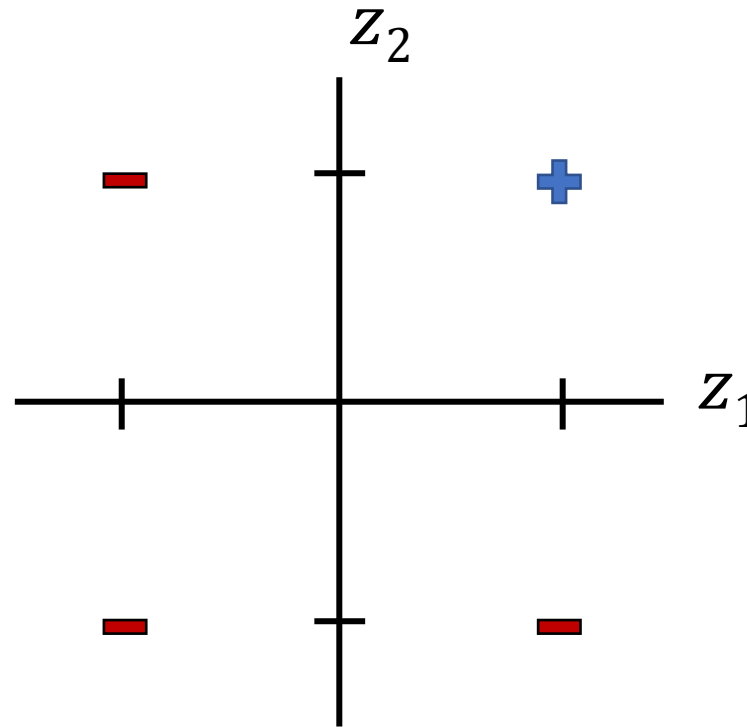
Piazza Poll 4

Which of the following parameters of $h_C(\mathbf{z})$ will perfectly classify this data?

$$h_C(\mathbf{z}) = \text{sign}(\mathbf{w}_C^T \mathbf{z} + b_C)$$

$$\text{sign}(x) = \begin{cases} 1, & \text{if } x \geq 0 \\ -1, & \text{if } x < 0 \end{cases}$$

- A. $w_C = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, b_C = 0$
- B. $w_C = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, b_C = 1$
- C. $w_C = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, b_C = -1$
- D. None of the above



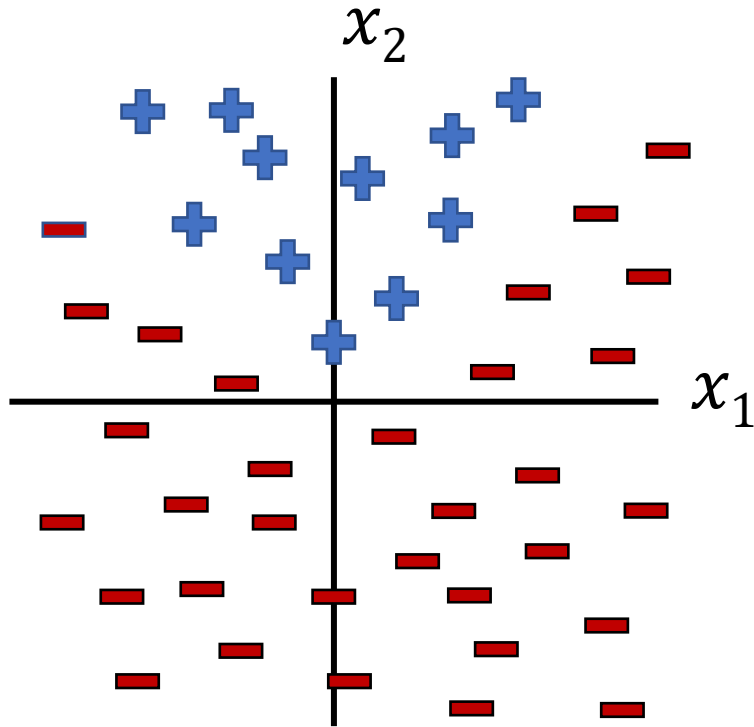
Classification Design Challenge

How could you configure three specific perceptrons to classify this data?

$$h_A(\mathbf{x}) = \text{sign}(\mathbf{w}_A^T \mathbf{x} + b_A)$$

$$h_B(\mathbf{x}) = \text{sign}(\mathbf{w}_B^T \mathbf{x} + b_B)$$

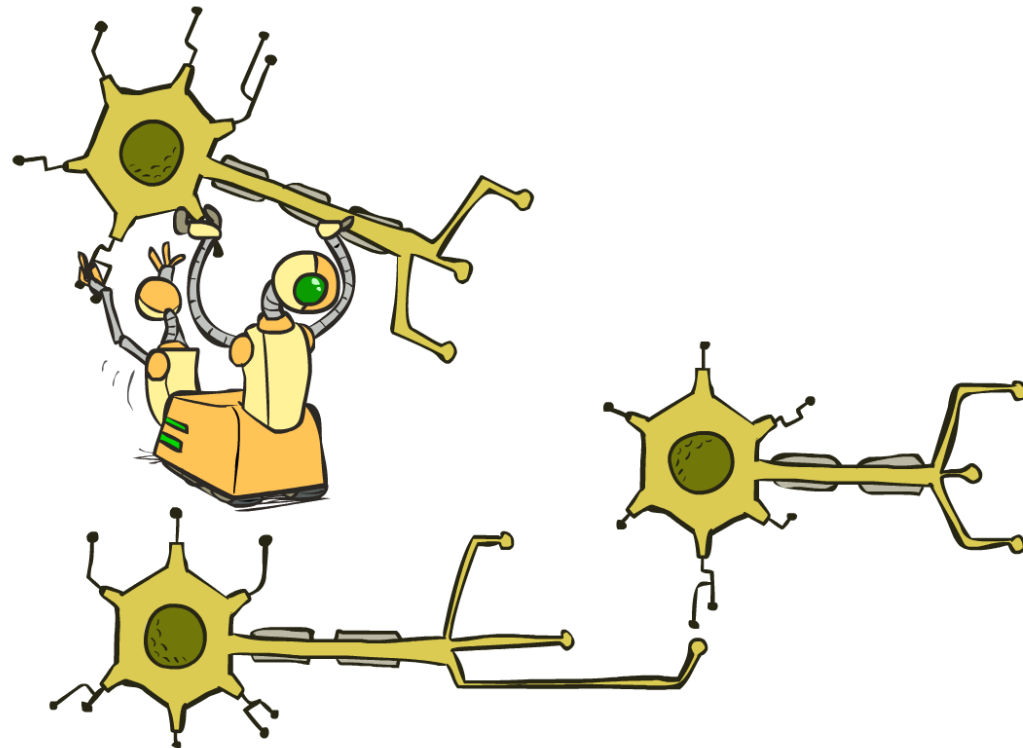
$$h_C(\mathbf{x}) = \text{sign}(\mathbf{w}_C^T \mathbf{x} + b_C)$$



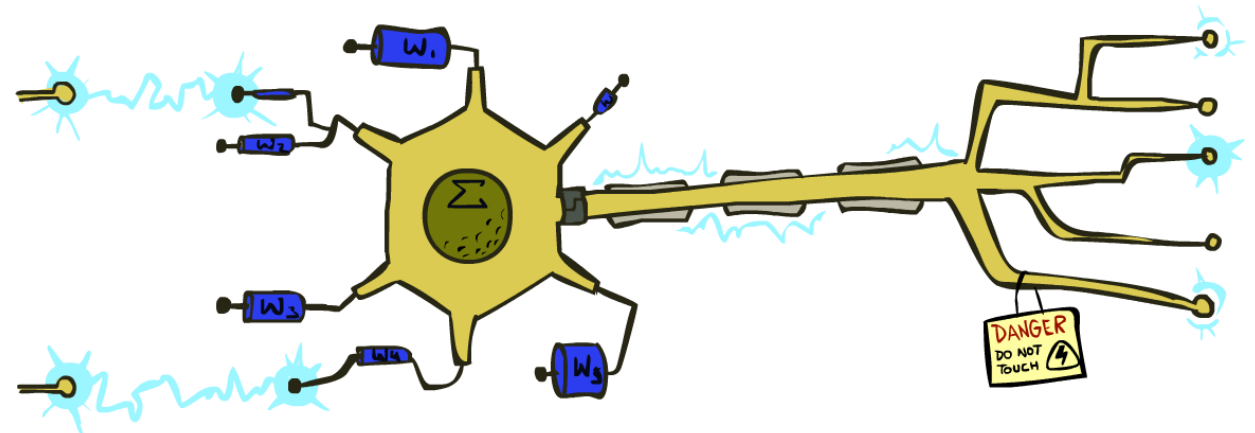
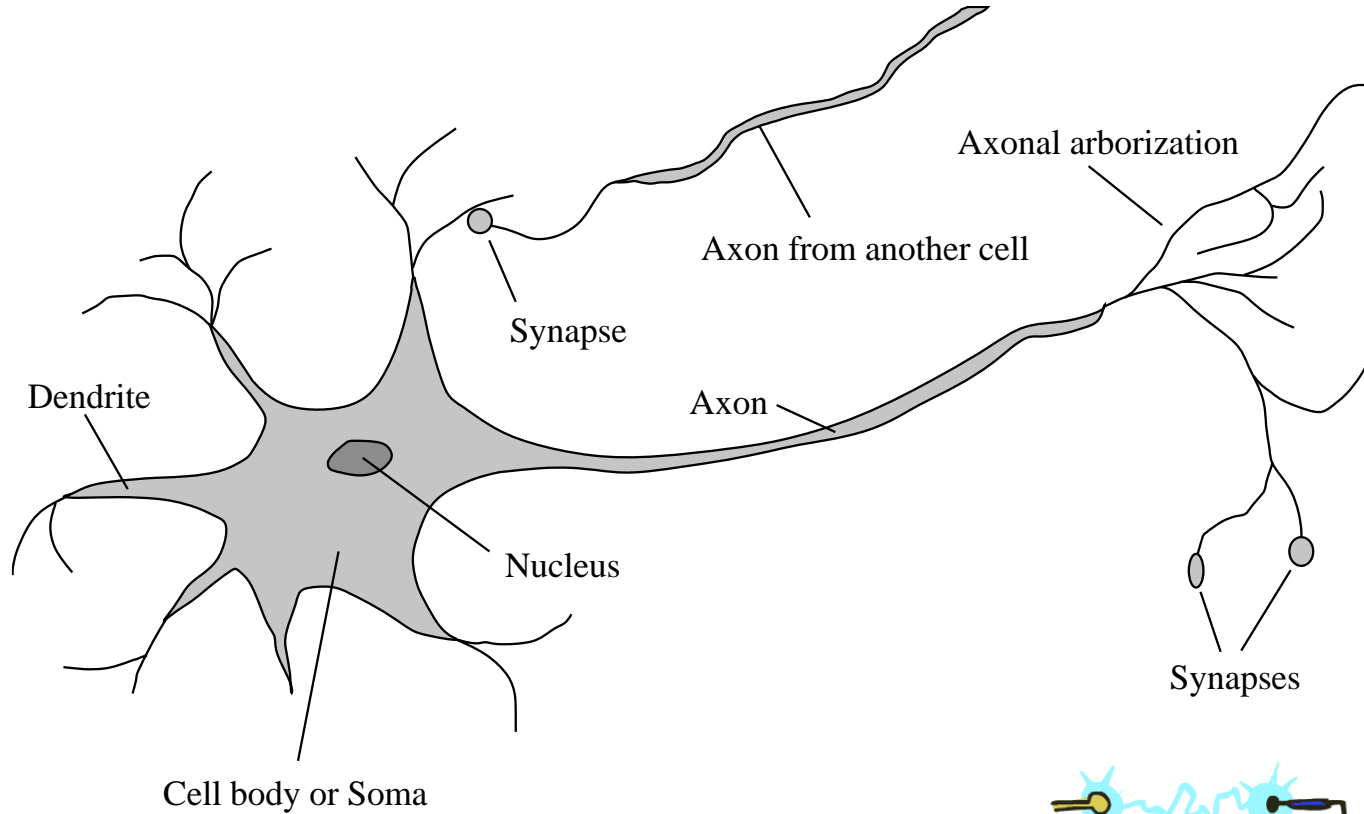
Multilayer Perceptrons

A **multilayer perceptron** is a feedforward neural network with at least one **hidden layer** (nodes that are neither inputs nor outputs)

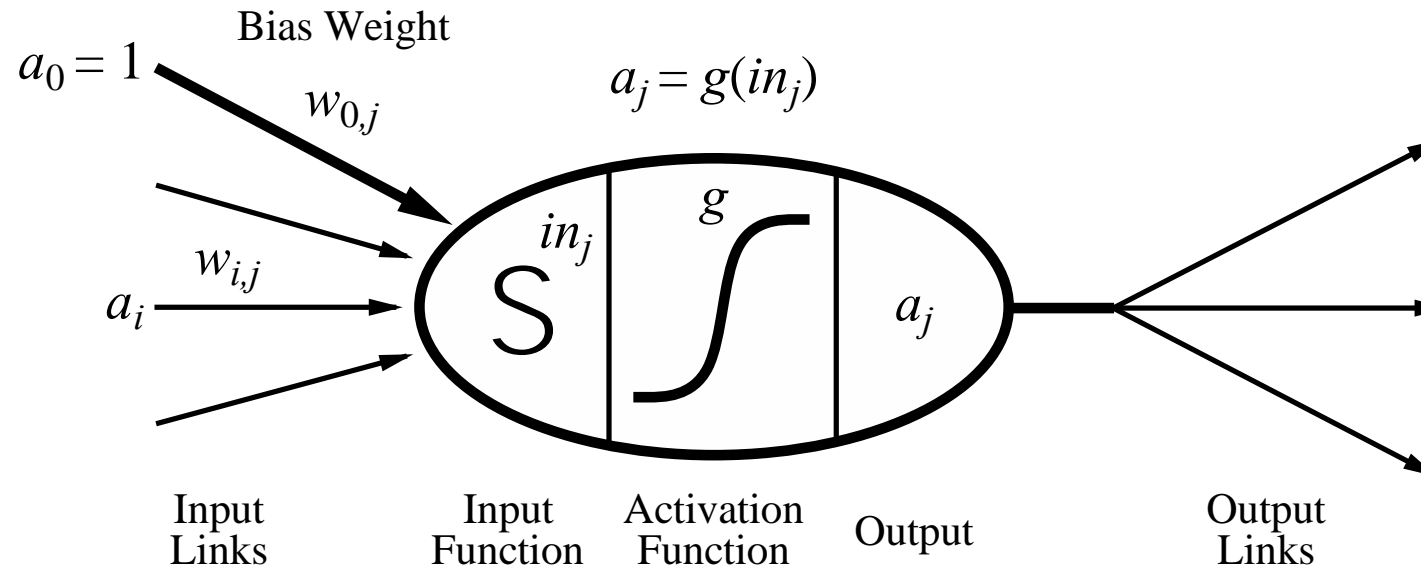
MLPs with enough hidden nodes can represent any function



Very Loose Inspiration: Human Neurons



Simple Model of a Neuron (McCulloch & Pitts, 1943)



Inputs a_i come from the output of node i to this node j (or from “outside”)

Each input link has a **weight** $w_{i,j}$

There is an additional fixed input a_0 with **bias** weight $w_{0,j}$

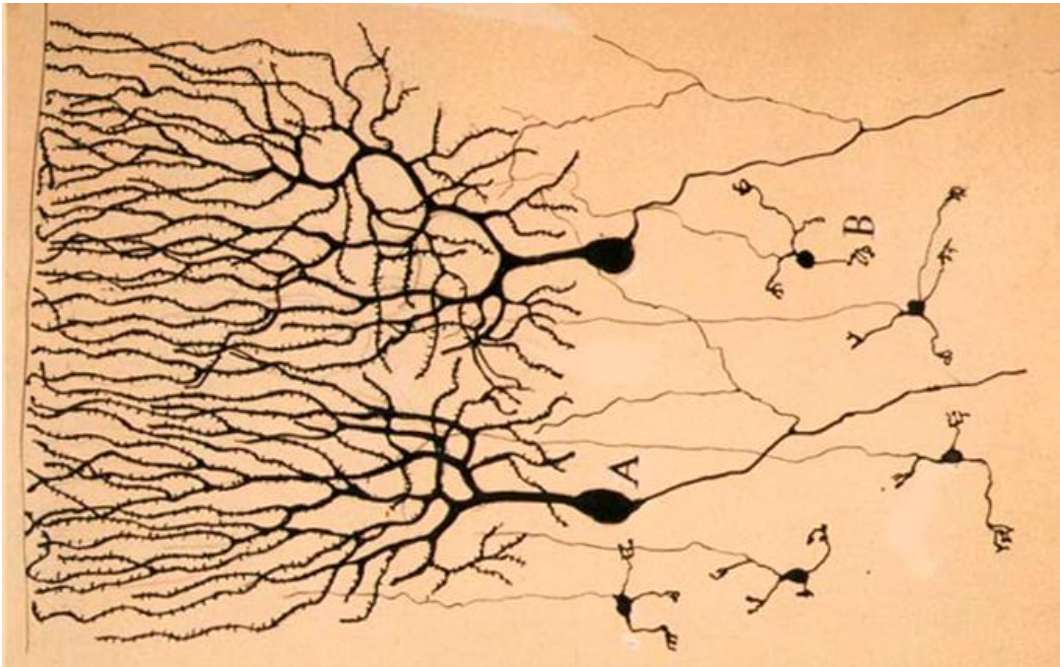
The total input is $in_j = \sum_i w_{i,j} a_i$

The output is $a_j = g(in_j) = g(\sum_i w_{i,j} a_i) = g(\mathbf{w} \cdot \mathbf{a})$

Neural Networks

Inspired by actual human brain

Input
Signal



Output
Signal



DOG



CAT



TREE



CAR



SKY