# ECE ING4
# MACHINE LEARNING

Jeremy Cohen

**Final Week**

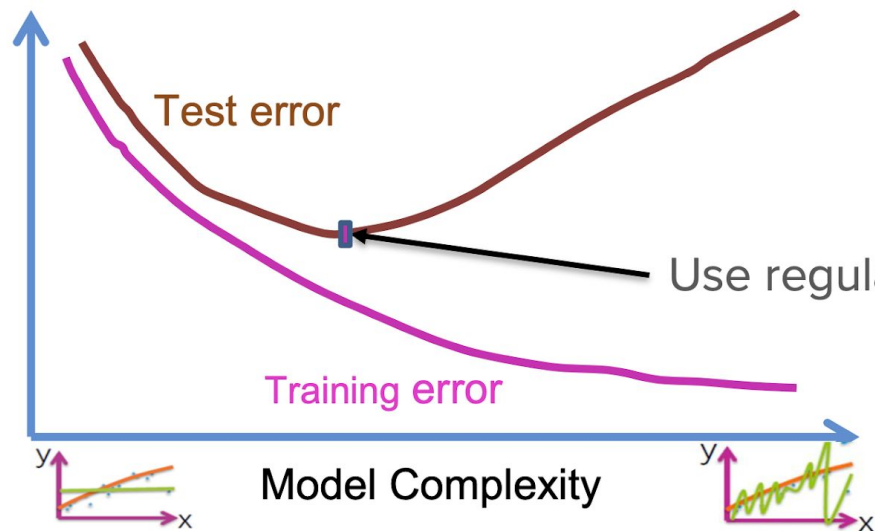- Full Review
- Regularization
- Mini Project

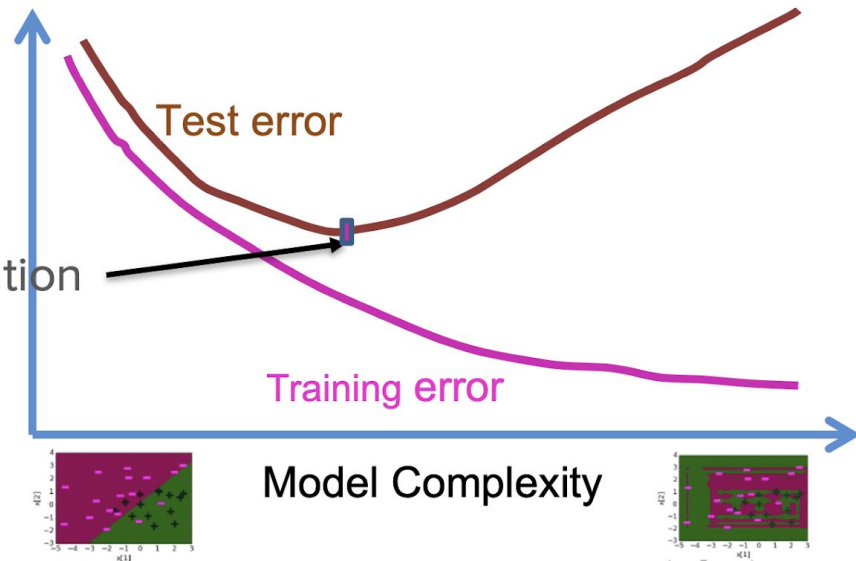# Full Review

# Our Final Course: Regularization
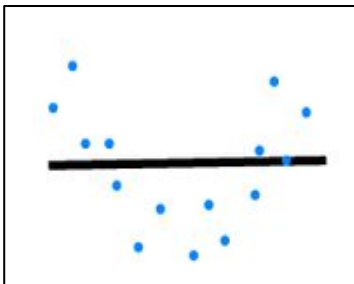
# The Problem

Polynomial Regression

Test error

Use regularization

Training **error**

Model Complexity

Logistic Regression

Test error

Use regularization

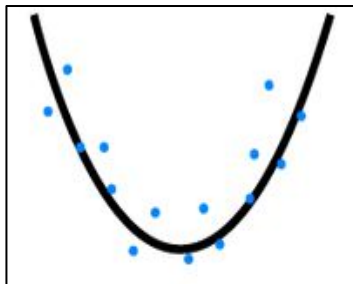Training **error**

Model Complexity

# The Problem



**UNDERFITTING**
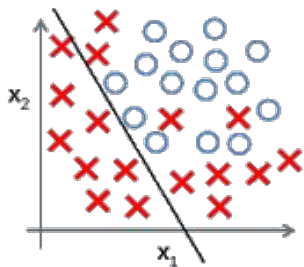
$$\theta_0 + \theta_1 x$$

**JUST RIGHT**

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

**OVERFITTING**

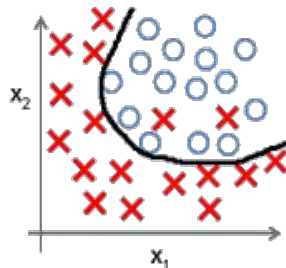$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$
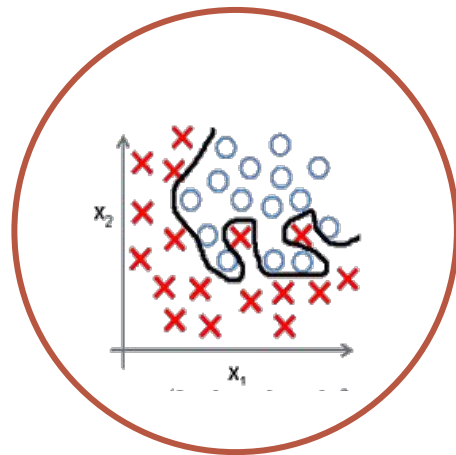
# For Logistic Regression



**UNDERFITTING**

$$h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$
( $g$ = sigmoid function)

**JUST RIGHT**

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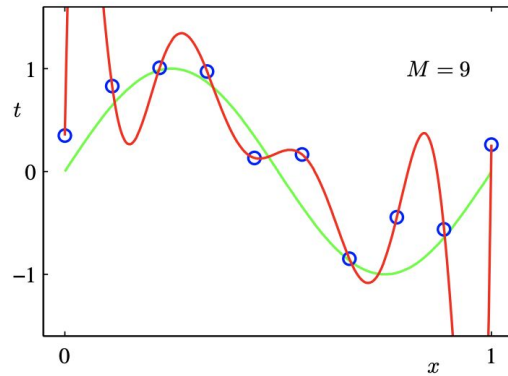$$+\theta_3 x_1^2 + \theta_4 x_2^2$$
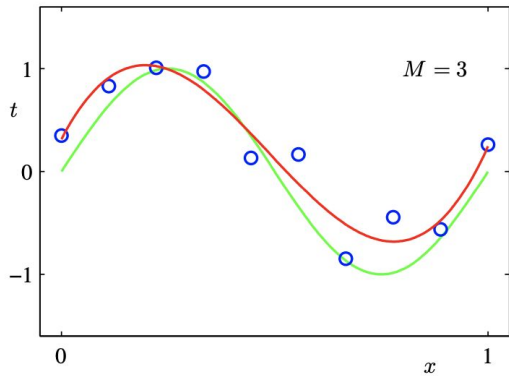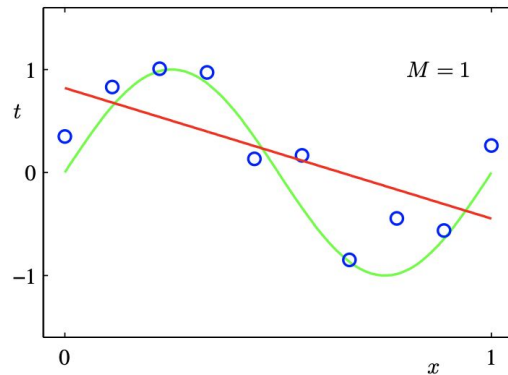$$+\theta_5 x_1 x_2)$$

**OVERFITTING**

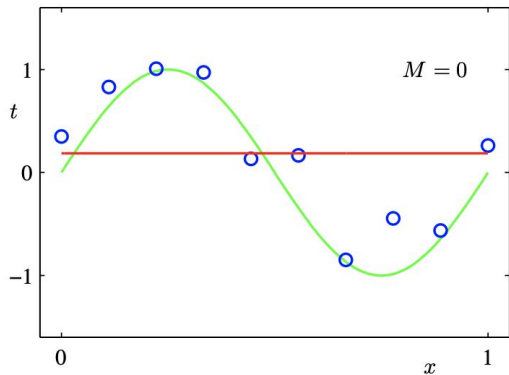$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2$$
$$+\theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2$$
$$+\theta_5 x_1^2 x_2^3 + \theta_6 x_1^3 x_2 + \ldots)$$
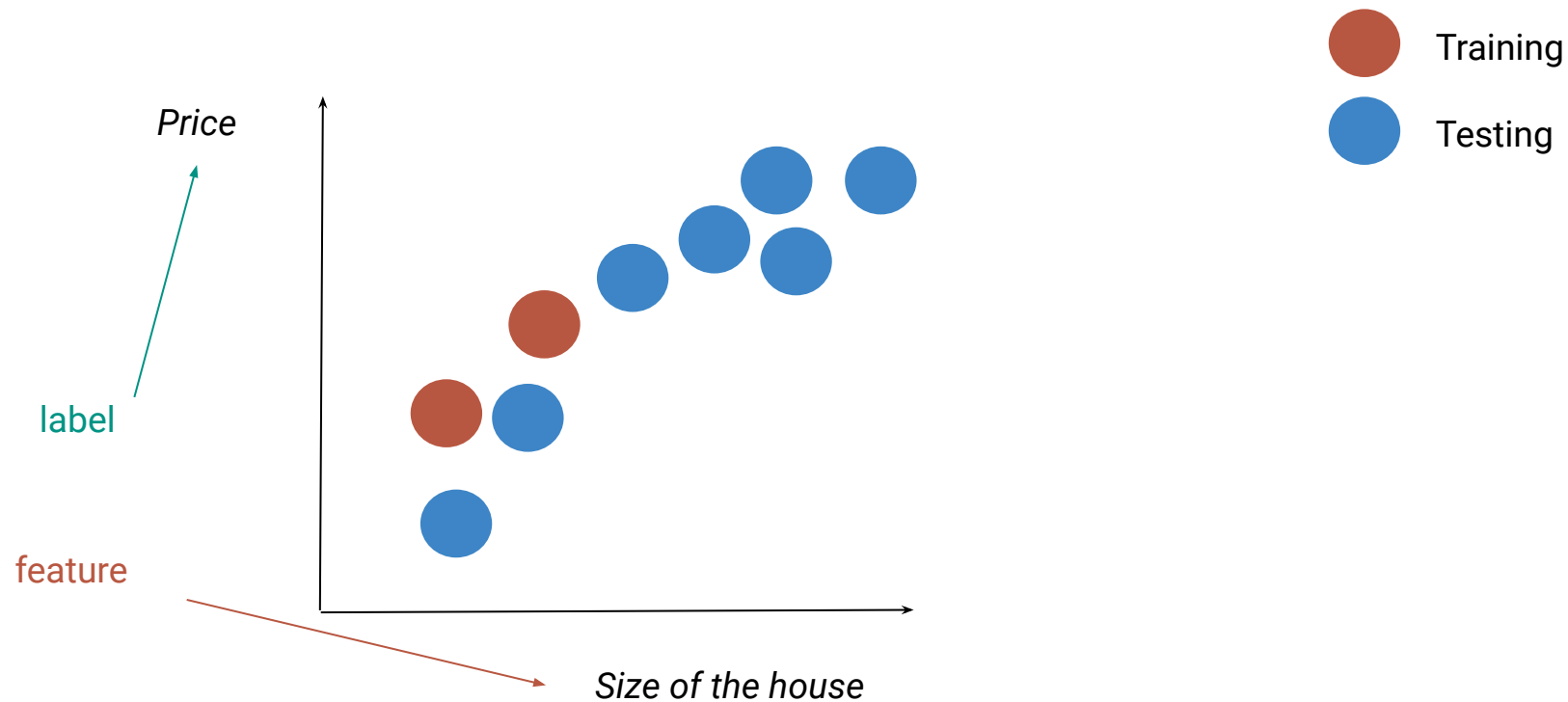
# Overfitting

We fit the training data correctly but fail to generalize

# For Logistic Regression

# Regression

# Regression



$$\text{price} = \theta_0 + \theta_1 * \text{size}$$

Price

Size of the house

Training

Testing

0 bias | some variance

# Regression

price = $\theta_0 + \theta_1$ * size

*Price*

*Size of the house*

Training

Testing

J($\theta$) = Sum of squared errors

We minimize

**the sum of the squared errors**

# Regularization

$$J(\theta) = \text{Sum of squared errors} + \lambda * \text{slope}^2$$



Price

simple regression

ridge regularization

Training

Testing

Size of the house

We minimize

**the sum of the squared errors**

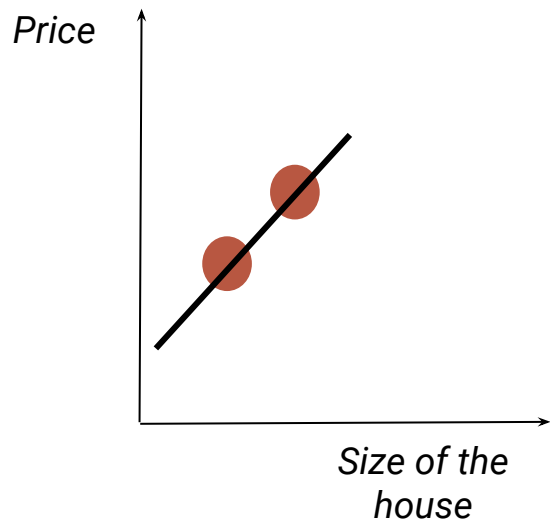**+ the slope ²**

# Regularization

$$J(\theta) = \text{Sum of squared errors} + \lambda * \text{slope}^2$$

**price = $\theta_0 + \theta_1 * \text{size}$**

*Price*

ridge
regularization

*Size of the
house*

adds a penalty

determine how strong the
penalty is

We minimize

**the sum of the squared errors**

**+   the slope $^2$**

# Ridge Regularization

$J(\theta)$ = Sum of squared errors  +  $\lambda$ * slope²

**price = 0.4 + 1.3 * size**

for now: $\lambda$ = 1

*Price*

$J(\theta) = 0 + 1* 1.3²$

$J(\theta) = 1.69$

*Size of the house*

# Ridge Regularization

$$J(\theta) = \text{Sum of squared errors} \quad + \quad \lambda * \text{slope}^2$$

**price = 0.9 + 0.8 * size**

*Price*

for now: $\lambda = 1$

$$J(\theta) = 0.3^2 + 0.1^2 + 1 * 0.8^2$$

$$J(\theta) = 0.74$$

*Size of the house*

# Ridge Regularization

**The price is less sensitive to the size of the house**

# What about λ

**If λ = 0**
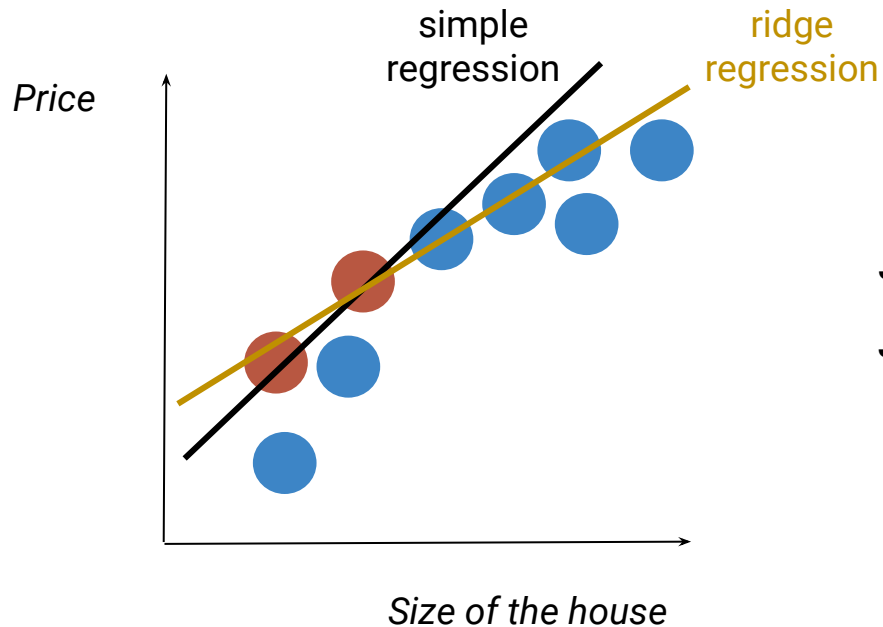


simple regression

ridge regression

*Price*

*Size of the house*

$J(\theta)$ = Sum of squared errors $+ \lambda \text{*slope}^2$

$J(\theta)$ = Sum of squared errors $+ 0$

# What about λ

**If λ = 1**

simple regression      ridge regression

*Price*

*Size of the house*

$J(\theta)$ = Sum of squared errors + λ*slope²

$J(\theta)$ = Sum of squared errors + 1*slope²

# What about λ

**If λ = 2**



$J(\theta)$ = Sum of squared errors + λ*slope²

$J(\theta)$ = Sum of squared errors + 2*slope²

# What about λ

**If λ = 10000**

Price

simple
regression

ridge
regression

Size of the house

$J(\theta)$ = Sum of squared errors $+ \lambda*slope^2$

$J(\theta)$ = Sum of squared errors $+$ $10000*slope^2$

# What about λ

**To estimate λ, we use cross-validation**

# Multiple Features

price = $\theta_0$ + slope1 * size + slope2* rating + slope3 * number of bedrooms+…

We minimize

**the sum of the squared errors**

**+ λ (slope1² + slope2² + slope3²)**

# Multiple features

**2 features - 1 data point**

*Price*

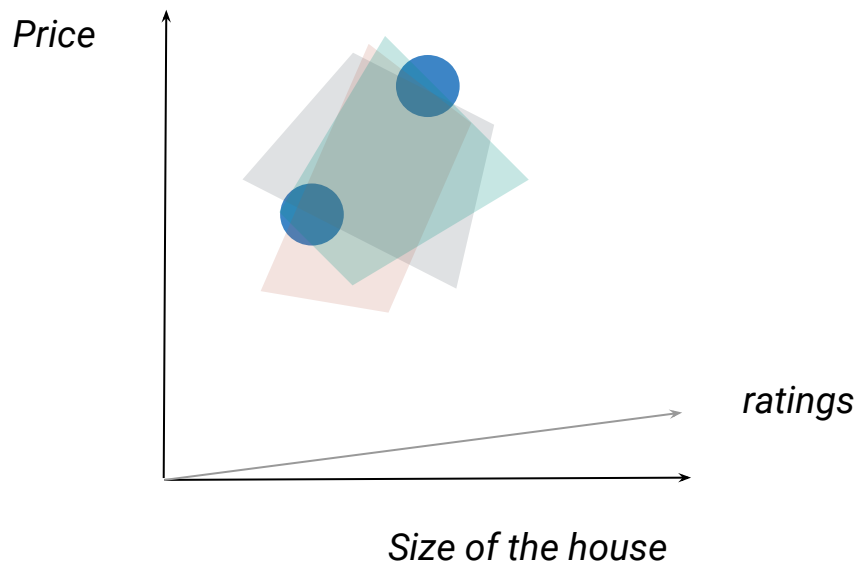*Size of the house*

# Multiple features

**2 features - 2 data points**

# Thousands of features
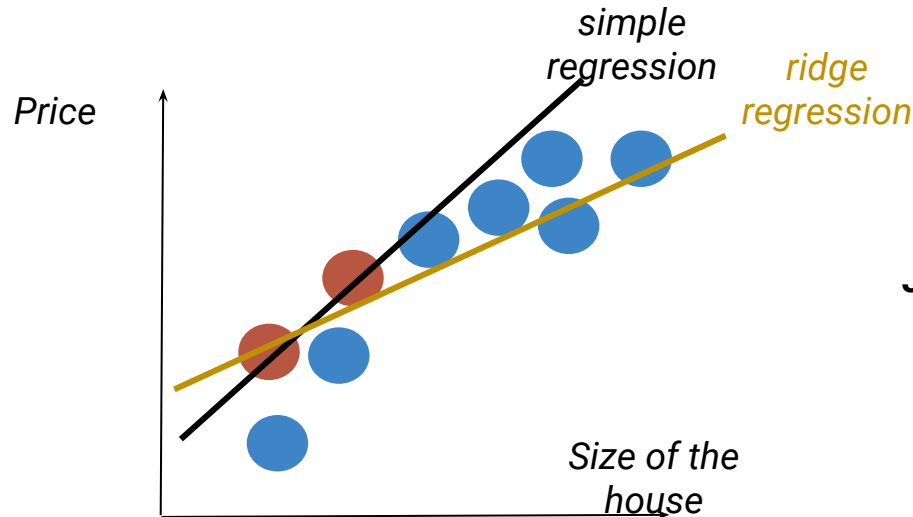
**3 features - 2 data points**

# Thousands of features

**If we have 4 parameters, we need 4 data points**

In case of a dataset with not enough data points compared to the number of features, regularization can help setting some parameters to 0

# Summary

**Ridge Regularization makes the regression less sensitive to the training data (especially when in low number) and helps reduce overfitting by adding a penalty to the cost function.**



$J(\theta) =$ Sum of squared errors $+ \lambda*slope^2$

# Ridge Regularization

**Ridge Regression (L$_2$ regularization)**

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

# Demo

http://madrury.github.io/smoothers/

# Lasso Regularization

**Ridge Regression (L$_2$ regularization)**

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

**Lasso Regression (L$_1$ regularization)**

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

$\lambda$ is the regularization parameter:

- Ridge: Encourages small weights $\theta$ but not exactly 0
- Lasso: "Shrink" some weights $\theta$ exactly to 0

# Gradient Descent

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})x_0^{(i)}$$

$$\theta_j := \theta_j - \alpha \quad \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})x_j^{(i)}$$

With regularization:    $\theta_j := \theta_j \left(1 - \alpha \frac{\lambda}{m}\right) - \alpha \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})x_j^{(i)}$

$\alpha$, $\lambda$ are learning parameters to choose manually

In practice: $(1 - \alpha\lambda/m)$ is between 0.99 and 0.95

# Logistic Regression

**Original Formula**
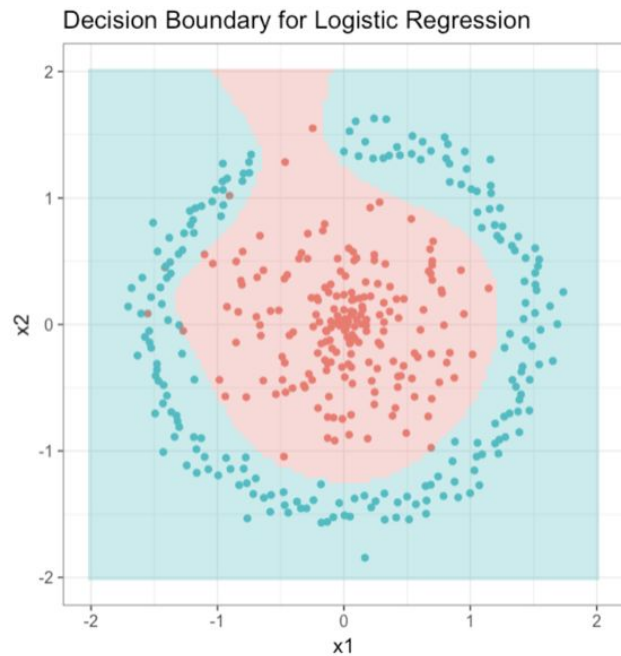
$$J(\theta) = -\frac{1}{m} \sum_{i=1}^{m} \left[ y^{(i)} \log(h_\theta(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)})) \right]$$
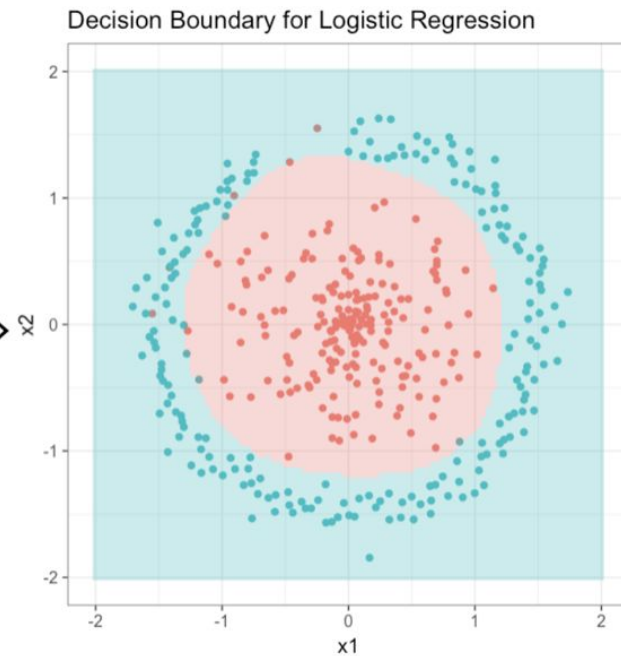
**Updated Formula**

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

*regularization term*

# Logistic Regression

# MINI PROJECT

# What now?

# BOSTON HOUSING PRICE



**Dataset**

## Boston Housing

Concerns housing values in suburbs of Boston

Chad Schirmer • updated 3 years ago (Version 1)

Data     Tasks     Kernels (27)     Discussion     Activity     Metadata

Download (12 KB)     New Notebook

https://www.kaggle.com/schirmerchad/bostonhoustingmlnd

# US ACCIDENTS



**US Accidents (3.0 million records)**

A Countrywide Traffic Accident Dataset (2016 - 2019)

Sobhan Moosavi • updated 13 days ago (Version 3)

Data    Tasks (6)    Kernels (43)    Discussion (6)    Activity    Metadata    Download (1 GB)    **New Notebook**

Dataset    410

https://www.kaggle.com/sobhanmoosavi/us-accidents/kernels
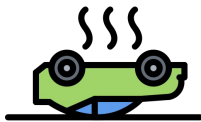
# ENRON

# FREE CHOICE

**Solve a real-world problem of your choice**

# Solve a problem using Machine Learning

**BOSTON HOUSE PRICES**

**US ACCIDENTS**

**ENRON SCANDAL**

**FREE CHOICE**

# Teaming Up is allowed (up to 2)

# FINAL THANK YOU

# ECE STUDENTS GET 50% OFF ALL COURSES!

https://jeremycohen.podia.com

COUPON
**ECE20**

See you soon!