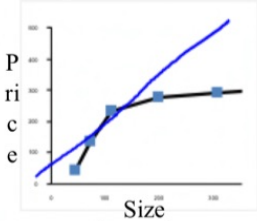
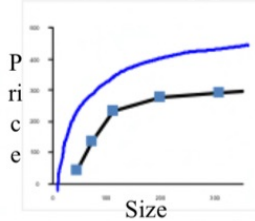


# The problem of over fitting

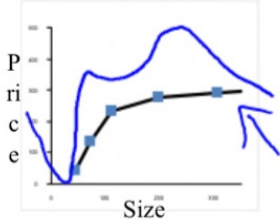
Example: Linear regression (housing prices)



$\rightarrow \theta_0 + \theta_1 x$   
"Underfit" "High bias"



$\rightarrow \theta_0 + \theta_1 x + \theta_2 x^2$   
"Just right"

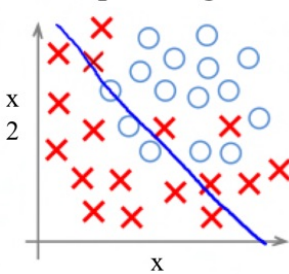


$\rightarrow \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$   
"Overfit" "High variance"

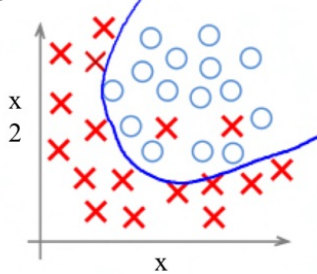
**Overfitting:** If we have too many features, the learned hypothesis may fit the training set very well ( $J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \approx 0$ ), but fail to generalize to new examples (predict prices on new examples).

Andrew

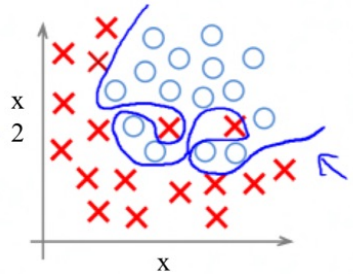
Example: Logistic regression



$\rightarrow h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$   
( $g$  = sigmoid function)  
"Underfit"



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



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

## Question

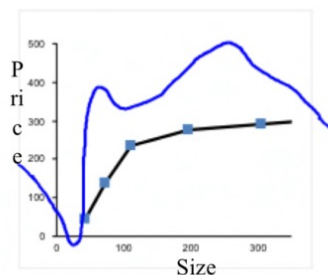
Consider the medical diagnosis problem of classifying tumors as malignant or benign. If a hypothesis  $h_\theta(x)$  has overfit the training set, it means that:

- ☐ It makes accurate predictions for examples in the training set and generalizes well to make accurate predictions on new, previously unseen examples.
- ☐ It does not make accurate predictions for examples in the training set, but it does generalize well to make accurate predictions on new, previously unseen examples.
- ☒ It makes accurate predictions for examples in the training set, but it does not generalize well to make accurate predictions on new, previously unseen examples.
- ☐ It does not make accurate predictions for examples in the training set and does not generalize well to make accurate predictions on new, previously unseen examples.

✓ Correct

### Addressing overfitting:

$x_1$  = size of house  
 $x_2$  = no. of bedrooms  
 $x_3$  = no. of floors  
 $x_4$  = age of house  
 $x_5$  = average income in neighborhood  
 $x_6$  = kitchen size  
:  
:  
 $x_{100}$



more features = harder to plot.

## Addressing overfitting:

Options:

1. Reduce number of features.

- Manually select which features to keep.
- Model selection algorithm (later in course).

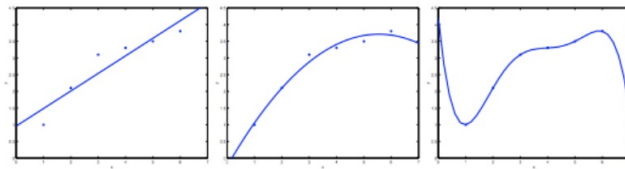
2. Regularization.

- Keep all the features, but reduce magnitude/values of parameters  $\theta_j$

Works well when we have a lot of features, each of which contributes a bit to predicting  $y$ .

## The Problem of Overfitting

Consider the problem of predicting  $y$  from  $x \in \mathbb{R}$ . The leftmost figure below shows the result of fitting a  $y = \theta_0 + \theta_1 x$  to a dataset. We see that the data doesn't really lie on straight line, and so the fit is not very good.



Instead, if we had added an extra feature  $x^2$ , and fit  $y = \theta_0 + \theta_1 x + \theta_2 x^2$ , then we obtain a slightly better fit to the data (See middle figure). Naively, it might seem that the more features we add, the better. However, there is also a danger in adding too many features: The rightmost figure is the result of fitting a 5<sup>th</sup> order polynomial  $y = \sum_{j=0}^5 \theta_j x^j$ . We see that even though the fitted curve passes through the data perfectly, we would not expect this to be a very good predictor of, say, housing prices ( $y$ ) for different living areas ( $x$ ). Without formally defining what these terms mean, we'll say the figure on the left shows an instance of **underfitting**—in which the data clearly shows structure not captured by the model—and the figure on the right is an example of **overfitting**.

Underfitting, or high bias, is when the form of our hypothesis function  $h$  maps poorly to the trend of the data. It is usually caused by a function that is too simple or uses too few features. At the other extreme, overfitting, or high variance, is caused by a hypothesis function that fits the available data but does not generalize well to predict new data. It is usually caused by a complicated function that creates a lot of unnecessary curves and angles unrelated to the data.

This terminology is applied to both linear and logistic regression. There are two main options to address the issue of overfitting:

1) Reduce the number of features:

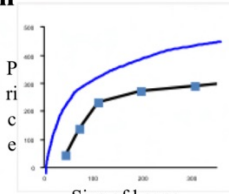
- Manually select which features to keep.
- Use a model selection algorithm (studied later in the course).

2) Regularization

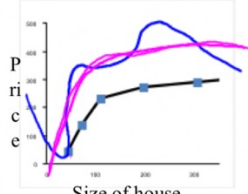
- Keep all the features, but reduce the magnitude of parameters  $\theta_j$ .
- Regularization works well when we have a lot of slightly useful features.

# Cost Function.

## Intuition



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



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

Suppose we penalize and make  $\theta_3, \theta_4$  really small.

$$\rightarrow \min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + 1000 \theta_3^2 + 1000 \theta_4^2$$

$\theta_3 \approx 0$        $\theta_4 \approx 0$

Andrew

Our goal is to minimize the whole function, adding  $\theta_3, \theta_4$  would only make the function become larger.

$\Rightarrow$  Overfitting.

## Regularization.

Small values for parameters  $\theta_0, \theta_1, \dots, \theta_n$

“Simpler” hypothesis

Less prone to overfitting

$$\theta_2, \theta_4 \approx 0$$

Housing:

Features:  $x_1, x_2, \dots, x_{100}$

Parameters:  $\theta_0, \theta_1, \theta_2, \dots, \theta_{100}$

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

~~$\theta_1, \theta_2, \theta_3, \dots, \theta_{100}$~~

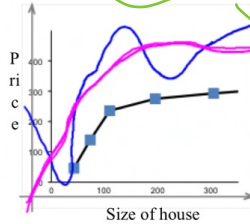
$$\lambda \sum_{j=1}^n \theta_j^2$$

## Regularization.

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

$\min_{\theta} J(\theta)$

regularization parameter



shrink all parameters  
Coz we don't know which  $\theta$  to get.  
regularize parameter.

If  $\lambda$  too large  $\Rightarrow$  fail to fit  
(underfitting)

## Question

In regularized linear regression, we choose  $\theta$  to minimize:

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

What if  $\lambda$  is set to an extremely large value (perhaps too large for our problem, say  $\lambda = 10^{10}$ )?

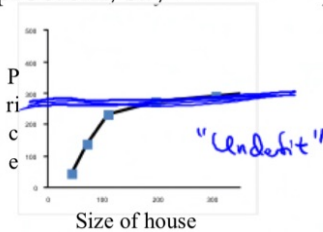
- ☐ Algorithm works fine; setting  $\lambda$  to be very large can't hurt it.
- ☐ Algorithm fails to eliminate overfitting.
- ☒ Algorithm results in underfitting (fails to fit even the training set).
- ☐ Gradient descent will fail to converge.

✓ Correct

In regularized linear regression, we choose  $\theta$  to minimize

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

What if  $\lambda$  is set to an extremely large value (perhaps far too large for our problem, say  $\lambda = 10^{10}$ )?



$\theta_1, \theta_2, \theta_3, \theta_4$   
 $\theta_1 \approx 0, \theta_2 \approx 0$   
 $\theta_3 \approx 0, \theta_4 \approx 0$

$$h_{\theta}(x) = \theta_0$$

$h_{\theta}(x)$

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

if it's too big will over take the all eq.

Andrew

unfitting  
 (fail to fit)

consequently this part will be eq to Zero.

# Cost Function

**Note:** [5:18 - There is a typo. It should be  $\sum_{j=1}^n \theta_j^2$  instead of  $\sum_{i=1}^n \theta_j^2$ ]

If we have overfitting from our hypothesis function, we can reduce the weight that some of the terms in our function carry by increasing their cost.

Say we wanted to make the following function more quadratic:

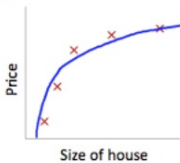
$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

We'll want to eliminate the influence of  $\theta_3 x^3$  and  $\theta_4 x^4$ . Without actually getting rid of these features or changing the form of our hypothesis, we can instead modify our **cost function**:

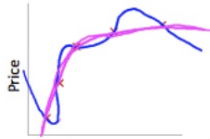
$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + 1000 \cdot \theta_3^2 + 1000 \cdot \theta_4^2$$

We've added two extra terms at the end to inflate the cost of  $\theta_3$  and  $\theta_4$ . Now, in order for the cost function to get close to zero, we will have to reduce the values of  $\theta_3$  and  $\theta_4$  to near zero. This will in turn greatly reduce the values of  $\theta_3 x^3$  and  $\theta_4 x^4$  in our hypothesis function. As a result, we see that the new hypothesis (depicted by the pink curve) looks like a quadratic function but fits the data better due to the extra small terms  $\theta_3 x^3$  and  $\theta_4 x^4$ .

## Intuition



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



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

Suppose we penalize and make  $\theta_3, \theta_4$  really small.

$$\rightarrow \min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + 1000 \theta_3^2 + 1000 \theta_4^2$$

$\theta_3 \approx 0$        $\theta_4 \approx 0$

We could also regularize all of our theta parameters in a single summation as:

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

The  $\lambda$ , or lambda, is the **regularization parameter**. It determines how much the costs of our theta parameters are inflated.

Using the above cost function with the extra summation, we can smooth the output of our hypothesis function to reduce overfitting. If lambda is chosen to be too large, it may smooth out the function too much and cause underfitting. Hence, what would happen if  $\lambda = 0$  or is too small?

Too small  $\Rightarrow$  fail to regularize parameter.



# Regularized Linear Regression

Only penalize  $\theta_1, \dots, \theta_n$ . Not  $\theta_0$   
Do nothing on  $\theta_0$ .

## Gradient descent

Repeat {

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

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

( $j = 1, 2, 3, \dots, n$ )

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

$$1 - \alpha \frac{\lambda}{m} < 1$$

Coz this one will always be positive

$$\therefore 1 - \alpha \frac{\lambda}{m} < 1$$

$\frac{\partial}{\partial \theta_j} J(\theta)$  → regularized

Same as gradient descent pattern.

## Question

Suppose you are doing gradient descent on a training set of  $m > 0$  examples, using a fairly small learning rate  $\alpha > 0$  and some regularization parameter  $\lambda > 0$ . Consider the update rule:

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

Which of the following statements about the term  $(1 - \alpha \frac{\lambda}{m})$  must be true?

☐  $1 - \alpha \frac{\lambda}{m} > 1$

☐  $1 - \alpha \frac{\lambda}{m} = 1$

☒  $1 - \alpha \frac{\lambda}{m} < 1$

☐ None of the above.

✓ Correct



## Normal equation

$$\underline{X} = \begin{bmatrix} (x^{(1)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix} \leftarrow$$

$m \times (n+1)$

$$\underline{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix} \quad \mathbb{R}^m$$

$$\rightarrow \min_{\theta} J(\theta)$$

$$\frac{\partial J(\theta)}{\partial \theta_j} \stackrel{\text{set}}{=} 0$$

$$\Rightarrow \underline{\theta} = \left( \underline{X}^T \underline{X} + \lambda \underbrace{\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}}_{(n+1) \times (n+1)} \right)^{-1} \underline{X}^T \underline{y}$$

e.g.  $n=2$   $\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$

## Non-invertibility (optional/advanced).

Suppose  $m \leq n$ ,  $\leftarrow$   
 (#examples) (#features)

example  $\leq$  features  
 $\Rightarrow$  non-invertible.

$$\underline{\theta} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y}$$

$\underbrace{\hspace{10em}}_{\text{non-invertible / singular}}$

$\underbrace{\hspace{10em}}_{\text{pinv}} \quad \underbrace{\hspace{10em}}_{\text{inv}}$

If  $\lambda > 0$ ,

$$\underline{\theta} = \left( \underline{X}^T \underline{X} + \lambda \begin{bmatrix} 1 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{bmatrix} \right)^{-1} \underline{X}^T \underline{y}$$

invertible.

$\Rightarrow$  overall will be invertible.

# Regularized Linear Regression

**Note:** [8:43 - It is said that  $X$  is non-invertible if  $m \leq n$ . The correct statement should be that  $X$  is non-invertible if  $m < n$ , and may be non-invertible if  $m = n$ .

We can apply regularization to both linear regression and logistic regression. We will approach linear regression first.

## Gradient Descent

We will modify our gradient descent function to separate out  $\theta_0$  from the rest of the parameters because we do not want to penalize  $\theta_0$ .

Repeat {

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

$$\theta_j := \theta_j - \alpha \left[ \left( \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})x_j^{(i)} \right) + \frac{\lambda}{m} \theta_j \right] \quad j \in \{1, 2, \dots, n\}$$

}

The term  $\frac{\lambda}{m} \theta_j$  performs our regularization. With some manipulation our update rule can also be represented as:

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

The first term in the above equation,  $1 - \alpha \frac{\lambda}{m}$  will always be less than 1. Intuitively you can see it as reducing the value of  $\theta_j$  by some amount on every update. Notice that the second term is now exactly the same as it was before.

## Normal Equation

Now let's approach regularization using the alternate method of the non-iterative normal equation.

To add in regularization, the equation is the same as our original, except that we add another term inside the parentheses:

$$\theta = (X^T X + \lambda \cdot L)^{-1} X^T y$$

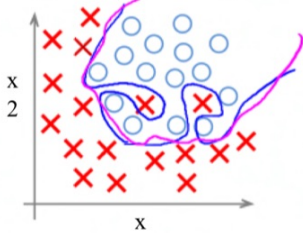
$$\text{where } L = \begin{bmatrix} 0 & & & & \\ & 1 & & & \\ & & 1 & & \\ & & & \ddots & \\ & & & & 1 \end{bmatrix}$$

$L$  is a matrix with 0 at the top left and 1's down the diagonal, with 0's everywhere else. It should have dimension  $(n+1) \times (n+1)$ . Intuitively, this is the identity matrix (though we are not including  $x_0$ ), multiplied with a single real number  $\lambda$ .

Recall that if  $m < n$ , then  $X^T X$  is non-invertible. However, when we add the term  $\lambda \cdot L$ , then  $X^T X + \lambda \cdot L$  becomes invertible.

# Regularized Logistic Regression.

## Regularized logistic regression.



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 + \dots)$$

Cost function:

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

Andrew

as long as add regularization ( $\frac{\lambda}{2m} \sum \theta_j^2$ ), can keep  $\theta$  small.

## Gradient descent

Repeat {

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

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

(j = 1, 2, 3, ..., n)

$\theta_1, \dots, \theta_n$

$\frac{\partial}{\partial \theta_j} J(\theta)$

diff from gradient descent as  $h_{\theta}$  are diff.

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

## Question

When using regularized logistic regression, which of these is the best way to monitor whether gradient descent is working correctly?

- ☐ Plot  $-\left[\frac{1}{m} \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))\right]$  as a function of the number of iterations and make sure it's decreasing.
- ☐ Plot  $-\left[\frac{1}{m} \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))\right] - \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$  as a function of the number of iterations and make sure it's decreasing.
- ☒ Plot  $-\left[\frac{1}{m} \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))\right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$  as a function of the number of iterations and make sure it's decreasing.
- ☐ Plot  $\sum_{j=1}^n \theta_j^2$  as a function of the number of iterations and make sure it's decreasing.

✓ Correct