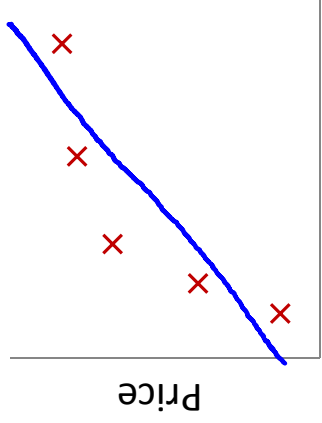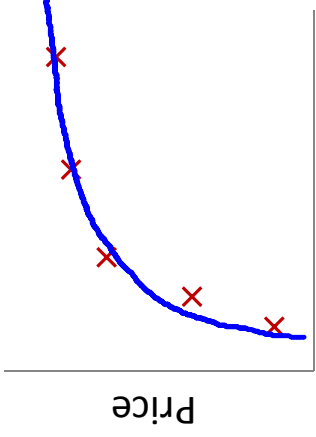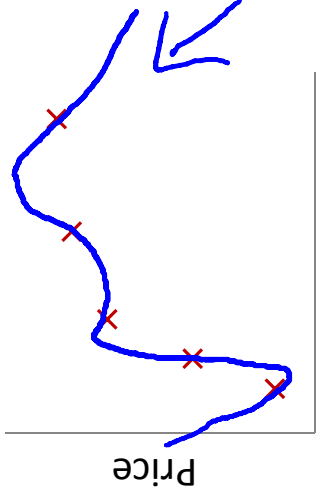# Regularization

# The problem of overfitting

# Example: Linear regression (housing prices)



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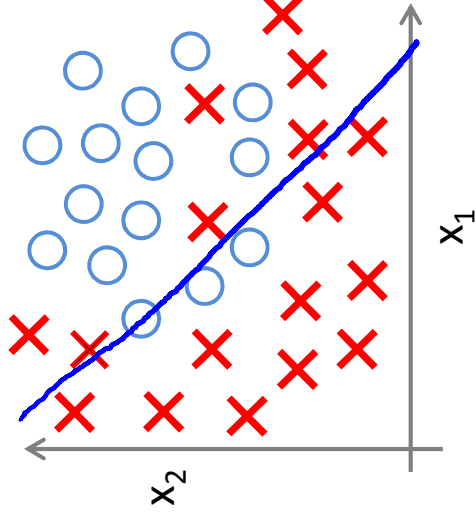"Underfit"  "High bias"

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

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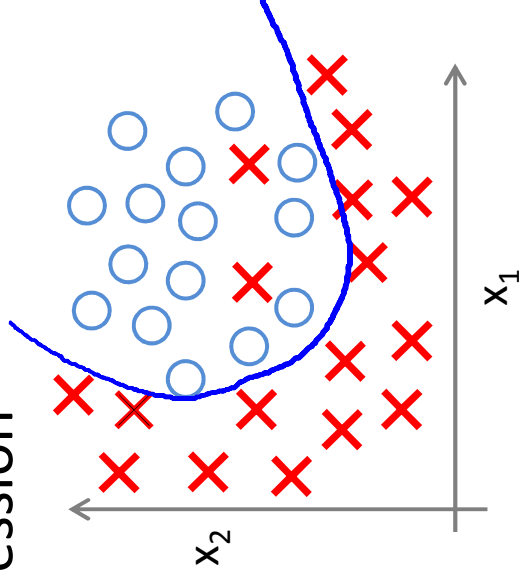"Overfit"  "High variance"

**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).

# Example: Logistic regression
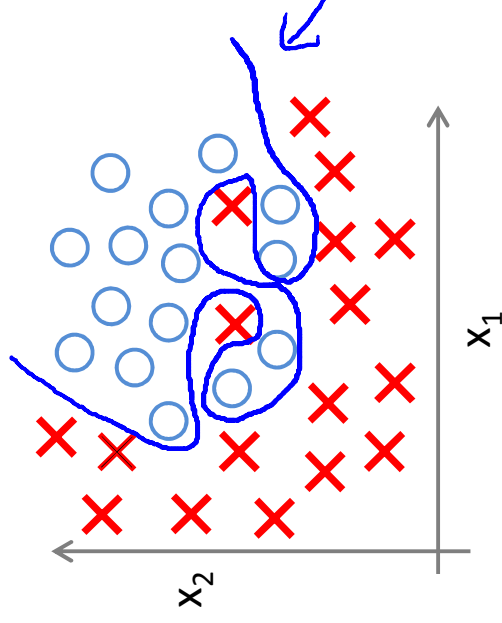


$\rightarrow h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$

( $g$ = sigmoid function)

$\leftarrow$ "Underfit"

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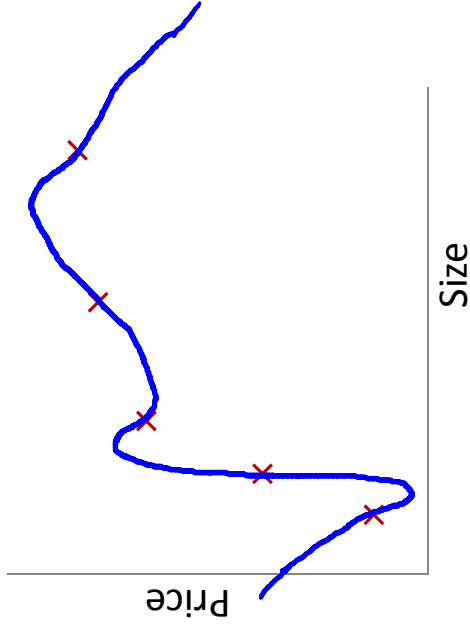$+ \theta_5 x_1 x_2)$

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

"Overfit"

## 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

$\ldots$

$x_{100}$

## Addressing overfitting:
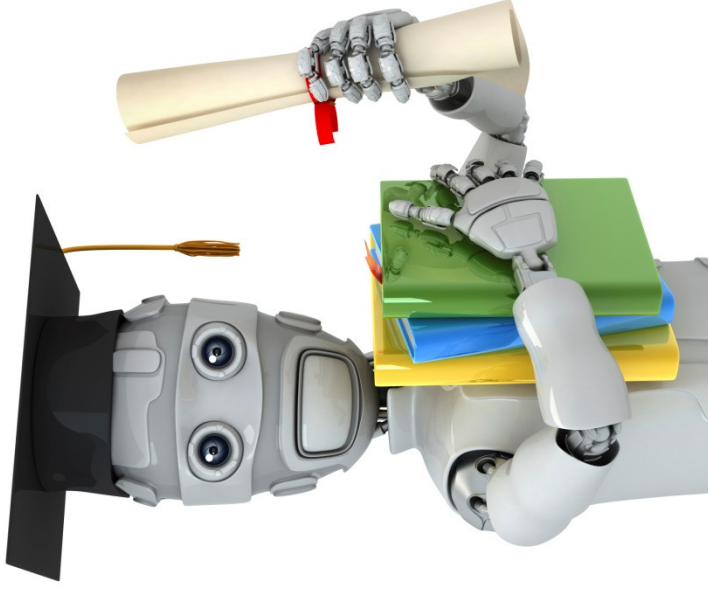
Options:

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 $y$.

# Regularization

# Cost function

Machine Learning

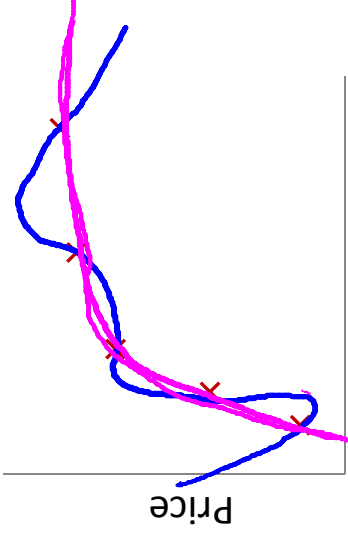# Intuition

Price | Size of house

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

Price | Size of house

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

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

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

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

# Regularization.

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

— "Simpler" hypothesis
— Less prone to overfitting

Housing:

— Features: $x_1, x_2, \ldots, x_{100}$
— Parameters: $\theta_0, \theta_1, \theta_2, \ldots, \theta_{100}$

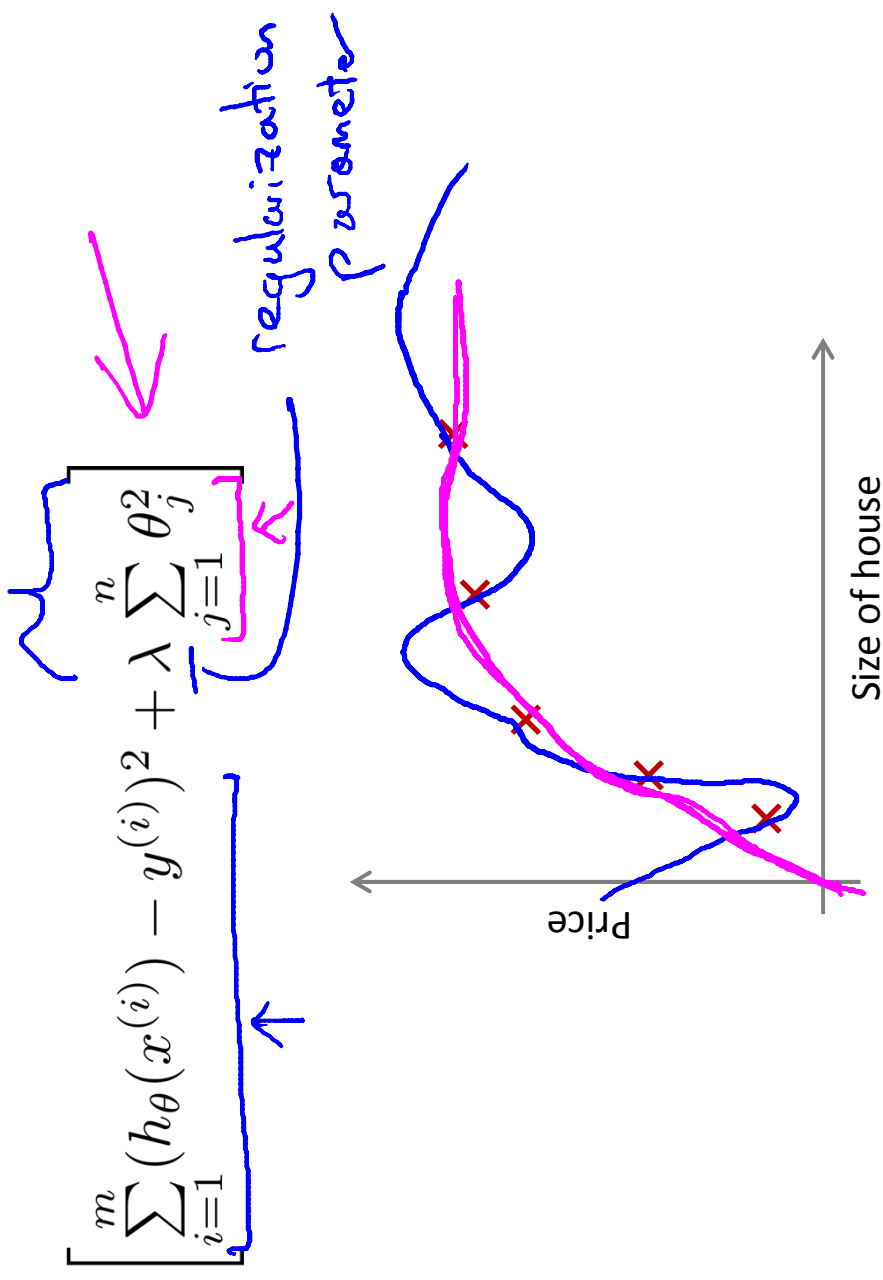$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^{n} \theta_j^2 \right]$$

$\theta_0, \theta_1, \theta_2, \ldots, \theta_{100}$

# Regularization.

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

regularization Parameter

$$\min_\theta J(\theta)$$

Price

Size of house

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

$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^{n} \theta_j^2 \right]$$
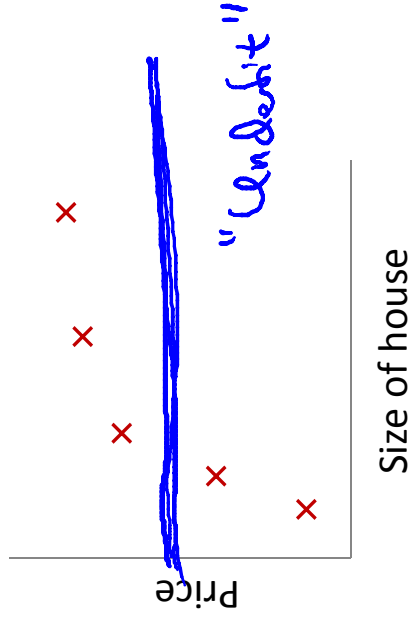
What if $\lambda$ is set to an extremely large value (perhaps for too large for our problem, say $\lambda = 10^{10}$)?

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

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 for too large for our problem, say $\lambda = 10^{10}$)?

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

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

$h_\theta(x) = \theta_0$

$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$

"Underfit"

Price

Size of house

# Regularization

## Regularized linear regression

## Regularized linear regression

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

$$\min_\theta J(\theta)$$

# Gradient descent

Repeat {

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

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

$$(j = \cancel{0}\, 1, 2, 3, \ldots, n)$$

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

}

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

$$0.99$$

$$\theta_j \times 0.99$$

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

$$\rightarrow J(\theta)$$

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

$$\frac{\partial}{\partial \theta_0}$$

$$\theta_1, \theta_2, \ldots, \theta_n$$

**Normal equation**

Andrew Ng

$$X = \begin{bmatrix} (x^{(1)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix} \quad m \times (n+1)$$

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

$$\to \min_\theta J(\theta)$$

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

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

$(n+1) \times (n+1)$

e.g. $n=2$

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

# Non-invertibility (optional/advanced).

Suppose $m \leq n,$

(#examples) (#features)

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

non-invertible / singular

If $\lambda > 0,$

$$\theta = \left( X^T X + \lambda \begin{bmatrix} 0 & & & & \\ & 1 & & & \\ & & 1 & & \\ & & & \ddots & \\ & & & & 1 \end{bmatrix} \right)^{-1} X^T y$$

invertible.

pinv

$\dfrac{1}{K}$

# Regularization

## Regularized logistic regression

Machine Learning

# 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:

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

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

$$\theta_1, \theta_2, \dots, \theta_n$$

# Gradient descent

Repeat {

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

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

$(j = \cancel{0}, 1, 2, 3, \ldots, n)$

$\theta_1, \ldots, \theta_n$

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

}

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

## Advanced optimization

$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$ — theta(1)
— theta(2)
— theta(n+1)

$\to$ function [jVal, gradient] = costFunction(theta);

```
jVal = [code to compute J(θ)];
```

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

$J(\theta)$

$\to$ gradient(1) = [code to compute $\frac{\partial}{\partial \theta_0} J(\theta)$];

$\frac{1}{m} \sum_{i=1}^{m} \left(h_\theta(x^{(i)}) - y^{(i)}\right) x_0^{(i)}$

$\to$ gradient(2) = [code to compute $\frac{\partial}{\partial \theta_1} J(\theta)$];

$\frac{1}{m} \sum_{i=1}^{m} \left(h_\theta(x^{(i)}) - y^{(i)}\right) x_1^{(i)} - \frac{\lambda}{m} \theta_1$

$\to$ gradient(3) = [code to compute $\frac{\partial}{\partial \theta_2} J(\theta)$];

$\left(\frac{1}{m} \sum_{i=1}^{m} \left(h_\theta(x^{(i)}) - y^{(i)}\right) x_2^{(i)}\right) - \frac{\lambda}{m} \theta_2$

$\cdots$

gradient(n+1) = [code to compute $\frac{\partial}{\partial \theta_n} J(\theta)$];

fminunc (@ costFunction")