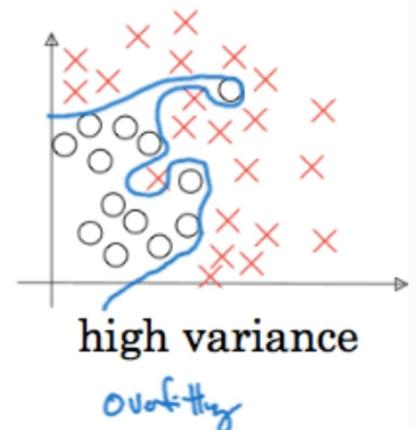
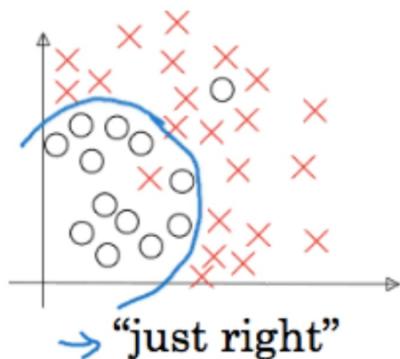
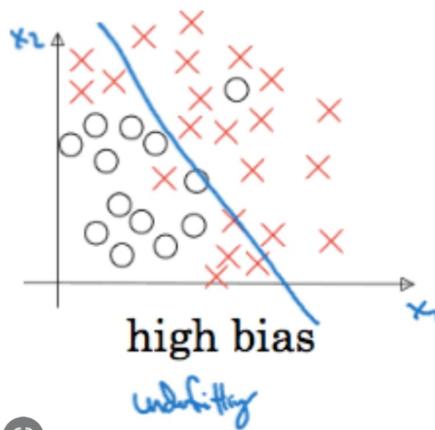


Bias vs Variance

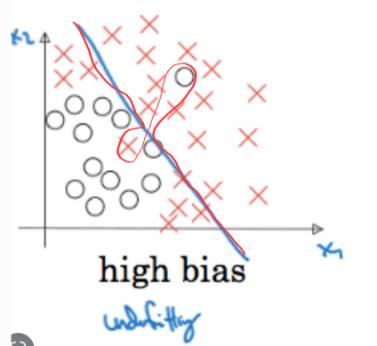


Let's take an example of cat vs dog:-

train error	1%	15%	15%	0.5%
val error	11%	16%	30%	1%
	high variance	high bias	high bias high variance	low bias low variance
Assume human error to be	0.8%	0.8%	0.8%	0.8%

If human error was 1%, then this would have been high variance only & this would have been low bias, low variance.

High bias & High Variance!-

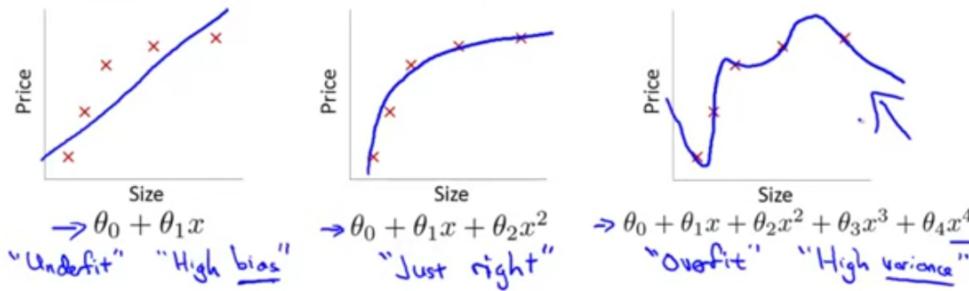


- high bias & high variance

Regularization

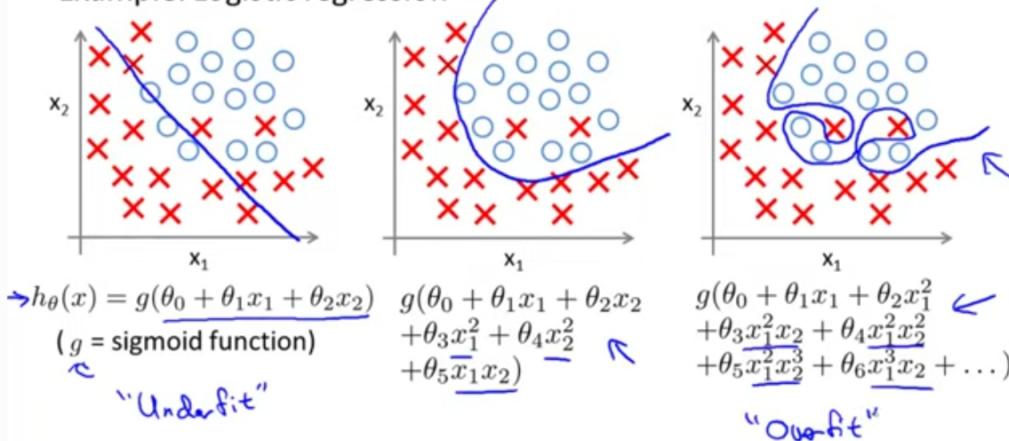
Example

Example: Linear regression (housing prices)



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

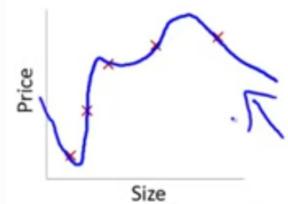


Addressing overfitting

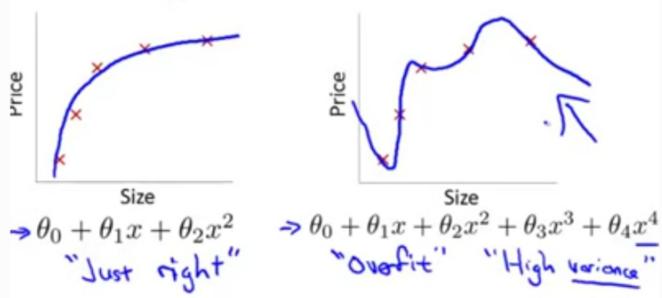
Addressing overfitting:

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



Intuition



so whose we penalize θ_3, θ_4 & make them really small so

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

Now when we min. above we'll end up with $\theta_3, \theta_4 \approx 0$, which will effectively make cost function to $\theta_0 + \theta_1x + \theta_2x^2$, which is good.

Idea

small values for parameters $\theta_0, \dots, \theta_3$

- simpler hypothesis
- less prone to overfitting

Modified cost function:-

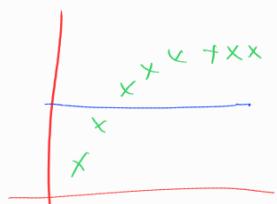
$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (\hat{h}_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

↑ regularization term
↑ regularization parameter

λ provides us a tradeoff b/w fitting the data & keeping the parameters small.

Q What if λ is very large?

then all θ 's will be ≈ 0 & we'll be fitting a line line parallel to x -axis i.e. $\hat{h}_{\theta}(x) = \theta_0$



This is underfitting.

Regularized linear regression

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

let's see gradient descent of linear regression :-

Gradient descent

Repeat {

$$\begin{aligned} & \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)} \quad \frac{\partial}{\partial \theta_0} J(\theta) \\ & \theta_j := \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} - \frac{\lambda}{m} \theta_j \right] \quad (j = \cancel{0}, 1, 2, 3, \dots, n) \\ & \theta_j := \theta_j (1 - \alpha \frac{\lambda}{m}) - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \end{aligned}$$

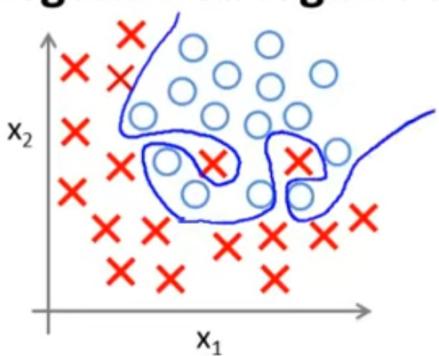
$\approx < 1$

so it becomes $0.9 \theta_j$ or $0.99 \theta_j$

This same as original

Logistic regression Regularized

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:

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

To make it regularize we add $+ \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$

Gradient descent

$$\frac{\partial J(\theta)}{\partial \theta_0}$$

$$\theta_0, \theta_1, \dots, \theta_n$$

Repeat {

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

$$\frac{\partial}{\partial \theta_0} J(\theta)$$

$$\Rightarrow \theta_j := \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} - \frac{\lambda}{m} \theta_j \right] \quad (j = \cancel{0}, 1, 2, 3, \dots, n)$$

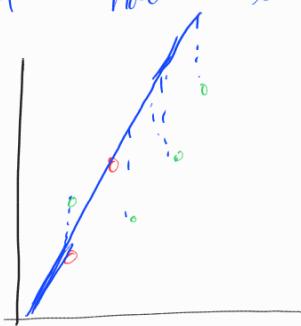
$$h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$

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

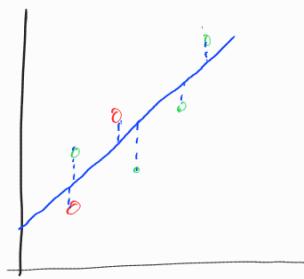
$$= \frac{\partial}{\partial \theta_j} J(\theta)$$

Ridge Regression

with small bias that penalty creates the least square fit will have large variance.

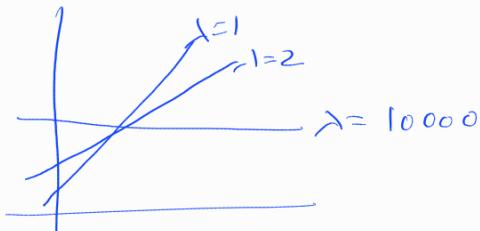


no bias, high variance



little bias, low variance

$$\uparrow \lambda = \downarrow \text{slope}$$



In order to find correct λ we just try a bunch of values for λ & use cross-validation, typically 10-fold cross validation, to determine which one results in lower variance.

Lasso vs Ridge

Regularization
Regularization seeks to control variance by adding a tuning parameter, lambda, or alpha:

LASSO (L1 regularization)

- regularization term penalizes absolute value of the coefficients
- sets irrelevant values to 0
- might remove too many features in your model

Ridge regression (L2 regularization)

- penalizes the size (square of the magnitude) of the regression coefficients
- enforces the B (slope/partial slope) coefficients to be lower, but not 0
- does not remove irrelevant features, but minimizes their impact