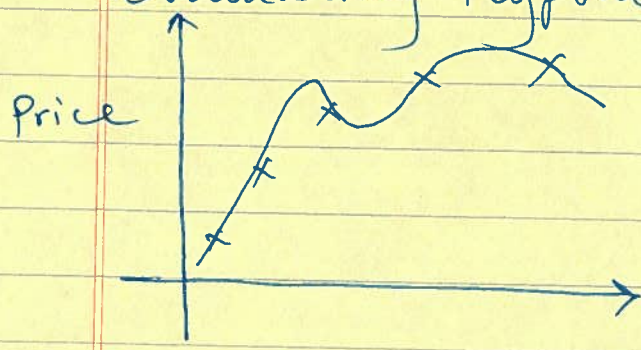## Week 6
## Deciding what to try next.

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

- Get more training examples
- Try smaller set of features
- Try getting additional features
- Try adding polynomial features ($x_1^2, x_2^2, x_1, x_?$
- Try decreasing $\lambda$
- Try increasing $\lambda$

Machine Learning diagnostic:

Evaluating Hypothesis.

Price

Size

Fails to generalize to new examples

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

| 70% | Training Set |
| --- | --- |
| 30% | Test set |

Overfitting training set.
Expect training error $J(\theta)$ to be low & test error $J(\theta)$ high.

→ Learn parameter $\theta$ from training data
(minimizing training error $J(\theta)$)

- Compute test set error:

$$J_{test}(\theta) = \frac{1}{2m_{test}} \sum_{i=1}^{m_{test}} \left( h_\theta(x_{test}^{(i)}) - y_{test}^{(i)} \right)^2$$

Misclassification          (0/1) misclassification error)

$$err(h_\theta(x), y) = \begin{cases} 1 & \text{if } h_\theta(x) \geq 0.5, \quad y=0 \\ & \text{or if } h_\theta(x) < 0.5 \quad y=1 \end{cases} error$$
$$\phantom{err(h_\theta(x), y) = } 0 \quad \text{otherwise}$$

$$\text{Test error} = \frac{1}{m_{test}} \sum_{i=1}^{m_{test}} err\left( h_\theta(x_{test}^{(i)}, y^{(i)}) \right).$$

$d=1$  1.  $h_\theta(x) = \theta_0 + \theta_1 x$  —  $\text{(H)}^{(1)}$     $d = $ degrees of polynomial

$d=2$  2.  $h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \rightarrow \text{(H)}^{(2)}$

$d=3$  3.  $h_\theta(x) = \theta_0 + \theta_1 x + \cdots + \theta_3 x^3 \rightarrow \text{(H)}^{(3)}$

$\vdots$

$d=3$  10.  $h_\theta(x) = \theta_0 + \theta_1 x + \cdots + \theta_{10} x^{10} \rightarrow \text{(H)}^{10}$

\* Training set $\rightarrow$ 80%

Cross Validation set $\rightarrow$ 20%    $M_{cv}$ = # of cv
(cv)                                                        examples
Test set $\longrightarrow$ 20%.    $(x^{(i)}_{cv}, y^{(i)}_{cv})$

$M_{test}$

Training error
$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2$$

cv error
$$J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2$$

Test Error
$$J_{test}(\theta) = \frac{1}{2m_{test}} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2$$
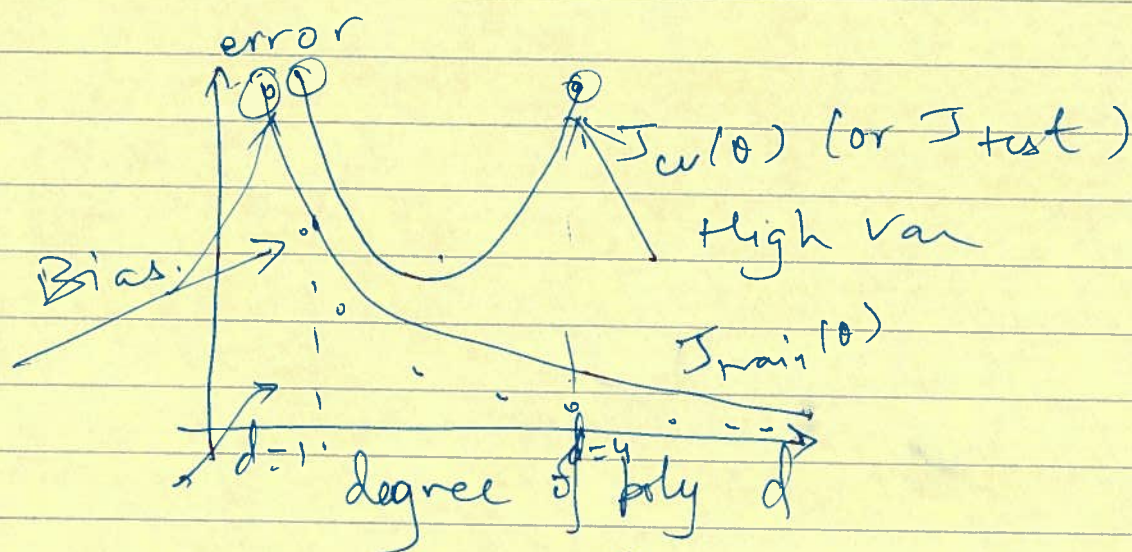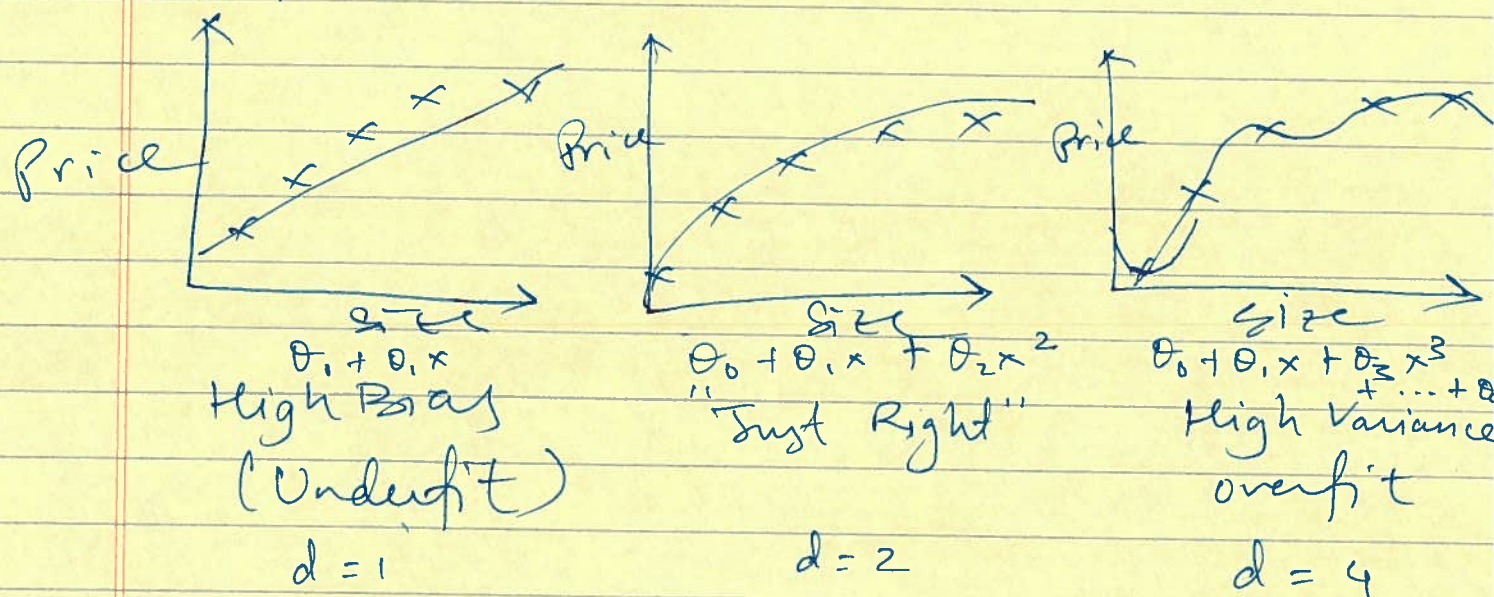
1. #

Test set to measure generalization error.

# Model selection

◁ Training / Validation / Test . Sets.

◁

## Bias Vs Variance or Both.



$\theta_0 + \theta_1 x$
High Bias
(Underfit)

$d = 1$

$\theta_0 + \theta_1 x + \theta_2 x^2$
"Just Right"

$d = 2$

$\theta_0 + \theta_1 x + \theta_3 x^3 + \dots + \theta$
High Variance
overfit

$d = 4$



error

$J_{cv}(\theta)$ (or $J_{test}$)

High Var

$J_{train}(\theta)$

Bias

$d=1$   $d=4$
degree of poly $d$

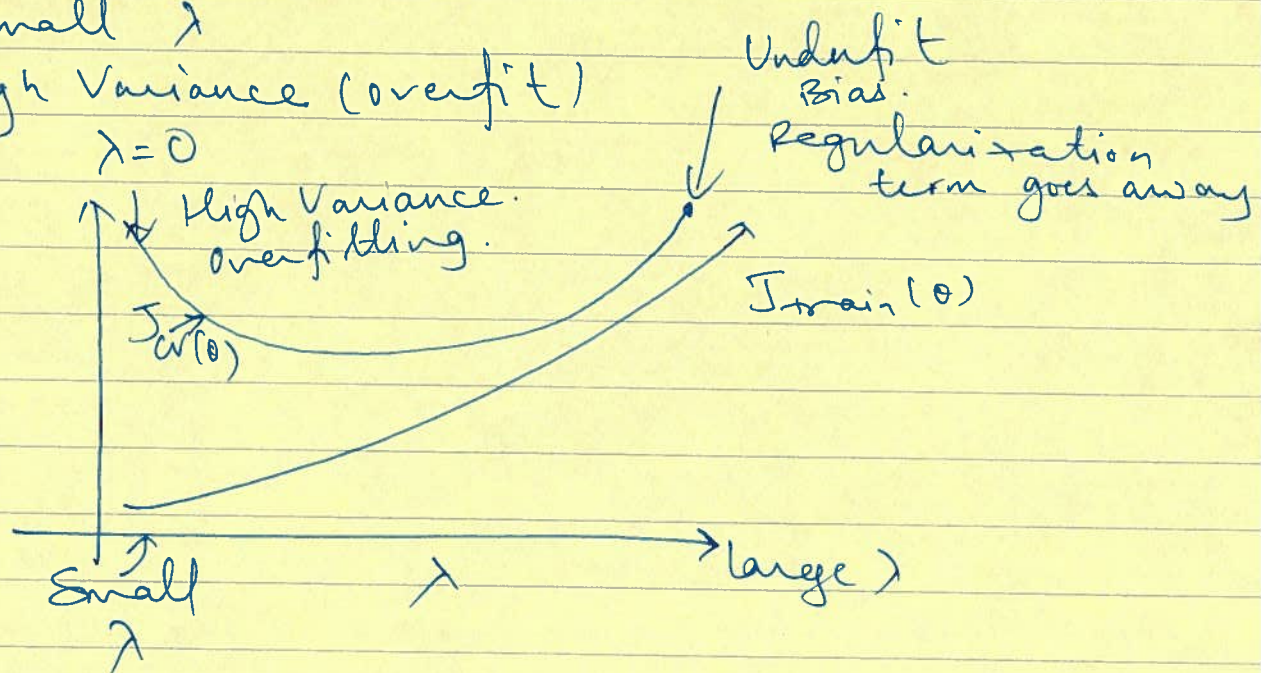| Bias (Underfit) | Variance (Overfit) |
|---|---|
| $J_{train}(\theta)$ } will be high <br> $J_{cv}(\theta)$ } <br> $J_{cv} \approx J_{train}$ | $J_{train}(\theta)$ will be low <br> $J_{cv}(\theta) \gg J_{train}(\theta)$ |

# Regularization Bias/Variance.

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

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2 + \frac{\lambda}{2m} \sum_{j=1}^{m} \theta_j^2$$
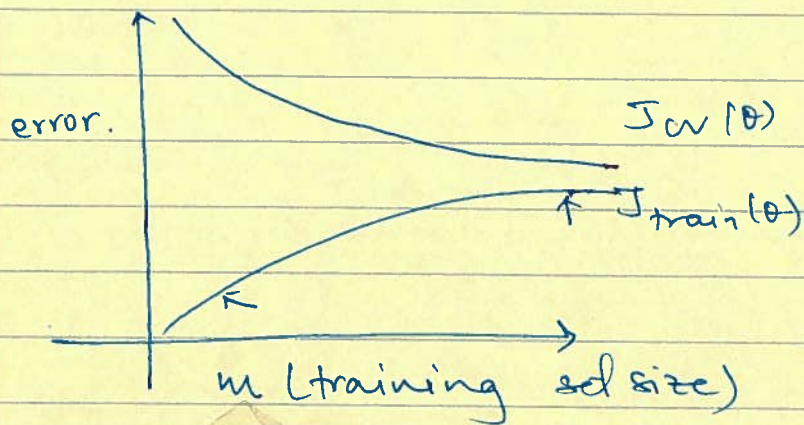
Large $\lambda$
High bias Underfit
$\lambda = 10000$, $\theta_1 \approx 0$, $\theta_2 \approx 0$
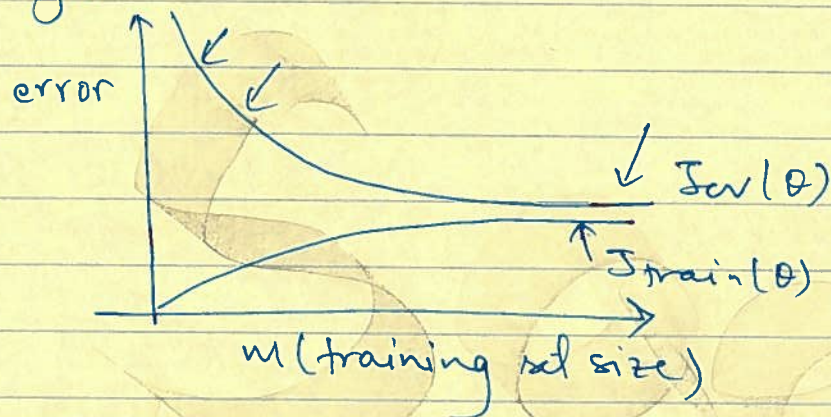$h_\theta(x) \approx \theta_0$

Intermediate $\lambda$
"Just Right"

Small $\lambda$
High Variance (overfit)
$\lambda = 0$

Underfit
Bias.
Regularization
term goes away

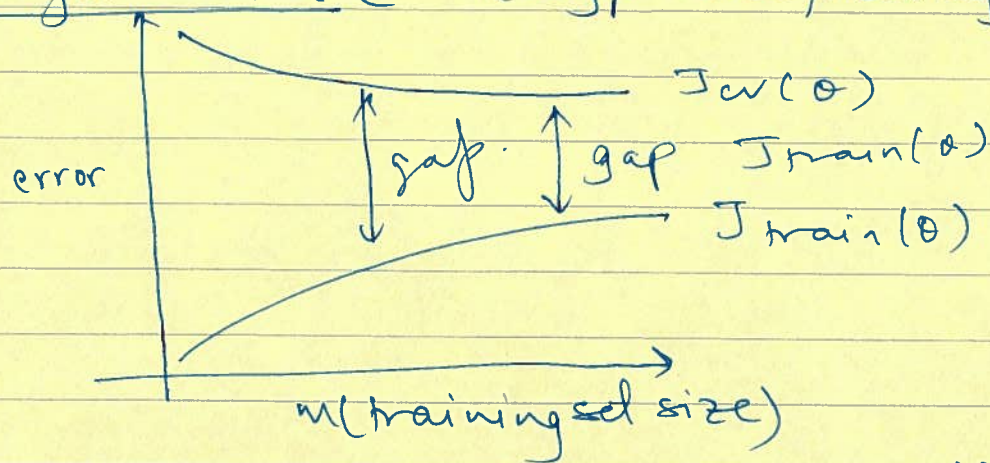# Learning Curves:
## Sanity check



High Bias



if a learning algorithm is suffering from hig bias, getting more training data will not (by itself) help much.
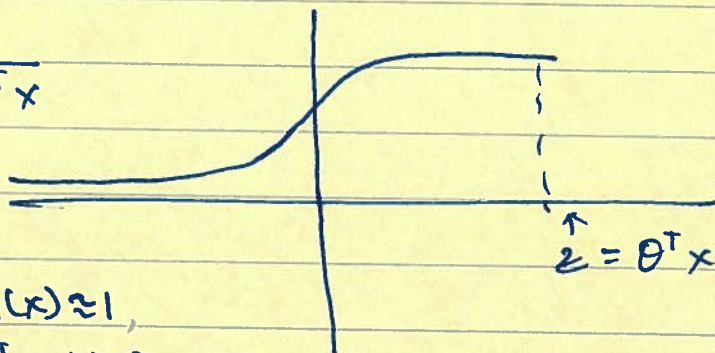
## High Variance    (hyp Overfitting )



$J_{cv}(\theta)$

gap  gap  $J_{train}(\theta)$

$J_{train}(\theta)$

error

m(training set size)

if learning algorithm is suffering from high variance, getting more training data is likely to help.

## Support Vector Machine

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



$z = \theta^T x$

If $y=1$, we want $h_\theta(x) \approx 1$,
$$\theta^T x \gg 0$$

if $y=0$ we want $h_\theta(x) \approx 0$,
$$\theta^T x \ll 0$$

$z = \theta^T x$

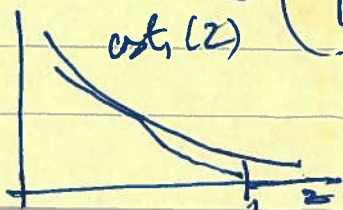Alternative view of logistic regression $(x,y)$

Cost of example:
$$-(y \log h_\theta(x) + (1-y) \log (1-h_\theta(x))) \leftarrow$$

$$= -\left( y \log \frac{1}{1+e^{-\theta^T x}} - (1-y) \log \left(1 - \frac{1}{1+e^{-\theta^T x}}\right) \right)$$

If $y=1$ (want $\theta^T x \gg 0$):

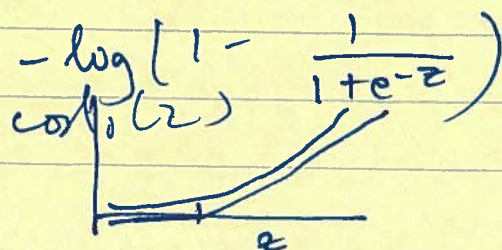If $y=1$ (want $\theta^T x \gg 0$):
$$-\log\left(\frac{1}{1+e^{-z}}\right)$$
$cost_1(z)$

If $y=0$ (want $\theta^T x \ll 0$
$$-\log\left(1 - \frac{1}{1+e^{-z}}\right)$$
$cost_0(z)$

A

## Support Vector machine

$$\min_\theta \frac{1}{m} \sum_{i=1}^{M} y^{(i)} cost_1 (\theta^T x^{(i)}) + (1-y^{(i)}) cost_0 (\theta^T x^{(i)})$$
$$+ \frac{\lambda}{2m} \sum_{j=0}^{n} \theta_j^2$$

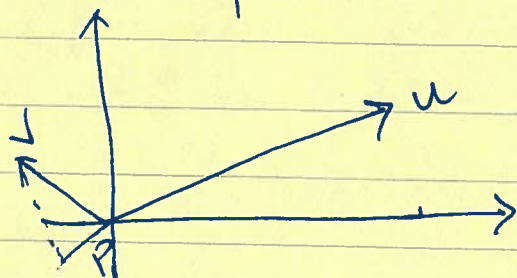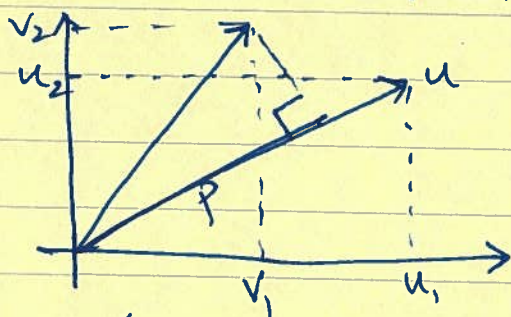$$\min_u (u-s)^2 + 1 \quad \underset{\lambda}{\Rightarrow} u = 5 \qquad \Bigg| \qquad \frac{A + \lambda B}{CA + B} \text{ if } C = \frac{1}{\lambda}$$

$$\min_u 10(u-s)^2 + 10 \Rightarrow u = 5$$

$$\min_\theta C \sum_{i=1}^{M} [ y^{(i)} cost_1 (\theta^T x^{(i)}) + (1-y^{(i)}) cost_0 (\theta^T x^{(i)})]$$
$$+ \frac{1}{2} \sum_{i=1}^{n} \theta_j^2$$

$$h_\theta(x) \begin{cases} 1 \\ 0 \end{cases} \qquad \text{if } \theta^T x \geqslant 0$$

## Math behind SVM
### Vector Inner Prod



$$u = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \qquad v = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

$$u^T v = ?$$

$\|u\| = \underline{length}$ of vec $u$
$$= \sqrt{u_1^2 + u_2^2} \in \mathbb{R}$$

$P =$ length of projection of $v$ onto $u$.

$$u^T v = P \cdot \|u\| = v^T u$$
$$= u_1 v_1 + u_2 v_2 \qquad P \in \mathbb{R}$$
$$P < 0$$

$$\omega = (\sqrt{\omega})^2$$

$$\min_{\theta} \frac{1}{2} \sum_{j=1}^{n} \theta_j^2 = \frac{1}{2}(\theta_1^2 + \theta_2^2) = \frac{1}{2}\left(\sqrt{\theta_1^2 + \theta_2^2}\right)^2 = \frac{1}{2}\|\theta\|^2$$

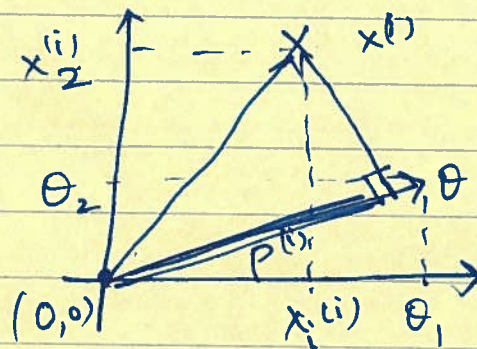$$\text{s.t.} \quad \theta^T x^{(i)} \geq 1 \quad \text{if} \quad y^{(i)} = 1$$
$$\theta^T x^{(i)} \leq -1 \quad \text{if} \quad y^{(i)} = 0$$

$$\begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$$

Simplication $\theta_0 = 0 \qquad n = 2$

$$\theta^T x^{(i)} = ?$$
$$\underset{u^T v}{\uparrow} \quad \updownarrow$$



$$\theta^T x^{(i)} = p^{(i)} \|\theta\|$$
$$= \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)}$$