

SVM

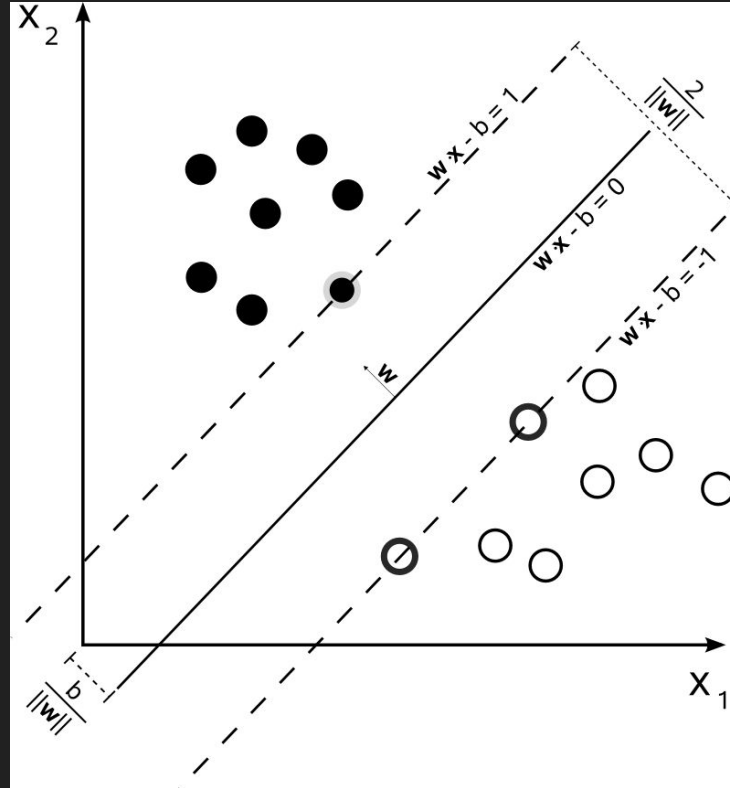
Support Vector Machines

Intuitions and Practical Uses

Support Vector Machines (SVMs) - Basic Ideas

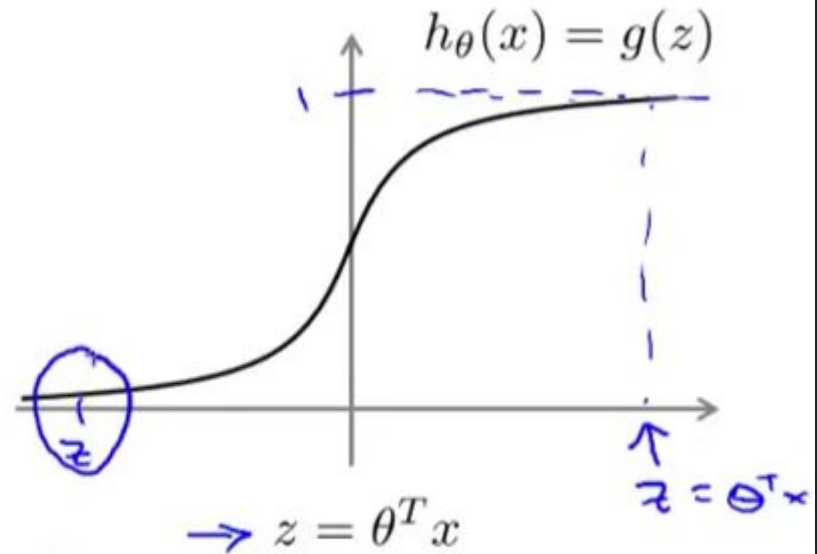
- Supervised learning model
- Used for classification and regression analysis
- We will focus on Classification tonight
- Constructs a hyperplane separating the classes
 - Hyperplane is a subspace of one less dimension than its ambient space
- Non-probabilistic
 - Uses Distances of sample data from class boundary - it is geometric
 - probabilistic classifier is a classifier that is able to predict, given a sample input, **a probability distribution over a set of classes**

Optimization Objectives Video on Coursera



Logistic Regression Function [1:30-2:56]

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



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

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

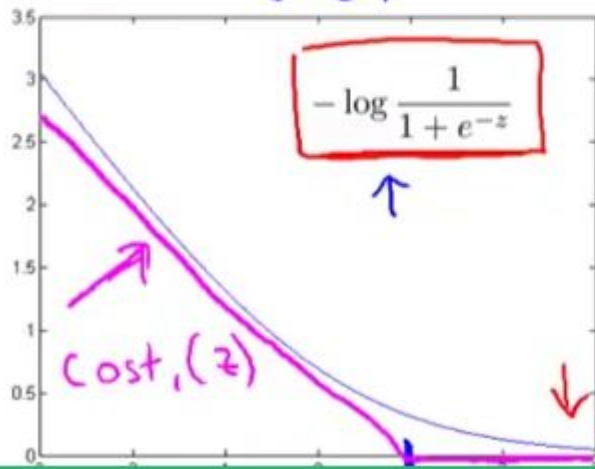
Representation Combines States [3:00-7:50]

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

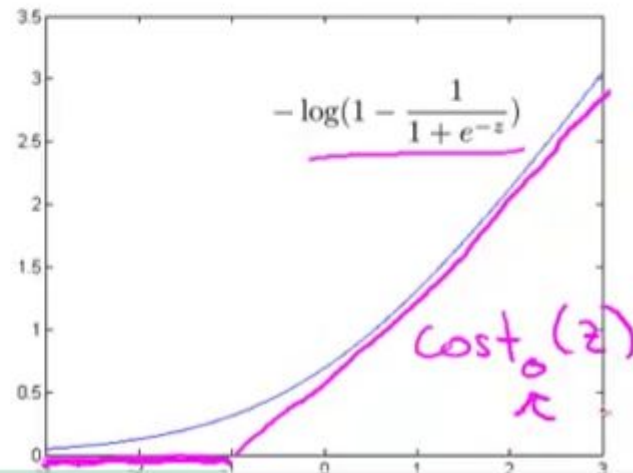
$$= -y \log \frac{1}{1 + e^{-\theta^T x}} - (1 - y) \log \left(1 - \frac{1}{1 + e^{-\theta^T x}}\right) \leftarrow$$

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

$$z = \theta^T x$$




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



SVM vs. Logistic Regression Equations [7:50-13:50]

Support vector machine

Logistic regression:

$$\min_{\theta} \frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \underbrace{\left(-\log h_{\theta}(x^{(i)}) \right)}_{\text{cost}_1(\theta^T x^{(i)})} + (1 - y^{(i)}) \underbrace{\left(-\log(1 - h_{\theta}(x^{(i)})) \right)}_{\text{cost}_0(\theta^T x^{(i)})} \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$


Support vector machine:

$$\min_{\theta} \cancel{\frac{1}{m}} C \sum_{i=1}^m y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

$$\min_u \left(\frac{(u-5)^2}{10} + 1 \right) \rightarrow u=5$$

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

$$\begin{aligned} A + \lambda B &\leftarrow \\ \rightarrow C A + B &\leftarrow \end{aligned}$$

$$C = \frac{1}{\lambda}$$

$$\min_{\theta} C \sum_{i=1}^m \left[y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

SVM Hypothesis [13:50-]

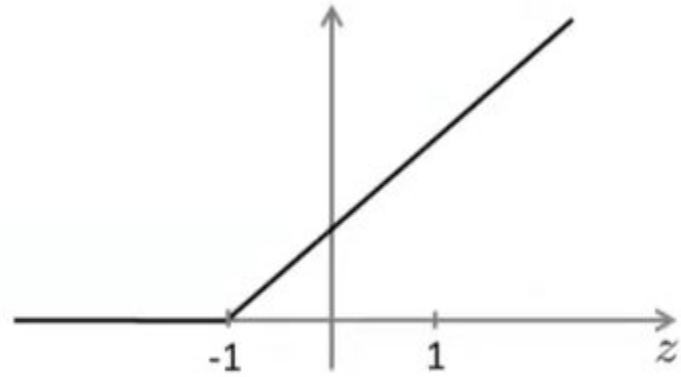
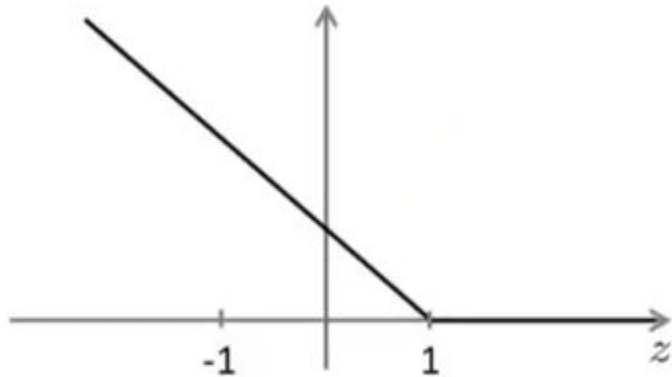
$$\Rightarrow \min_{\theta} C \sum_{i=1}^m \left[y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

Hypothesis:

$$h_{\theta}(x) = \begin{cases} 1 & \text{if } \theta^T x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

SVM Cost Function [0:17-2:31]

$$\min_{\theta} C \sum_{i=1}^m \left[y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

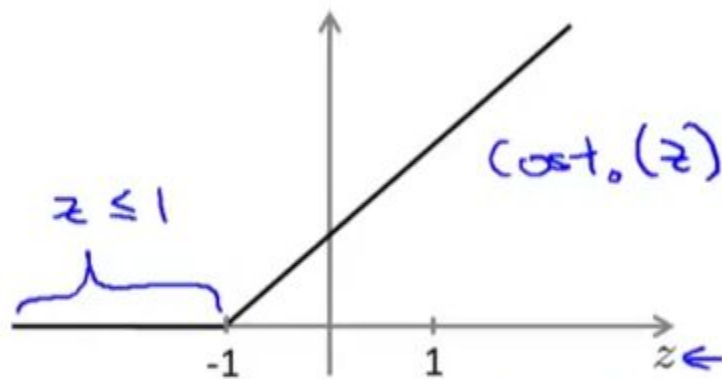
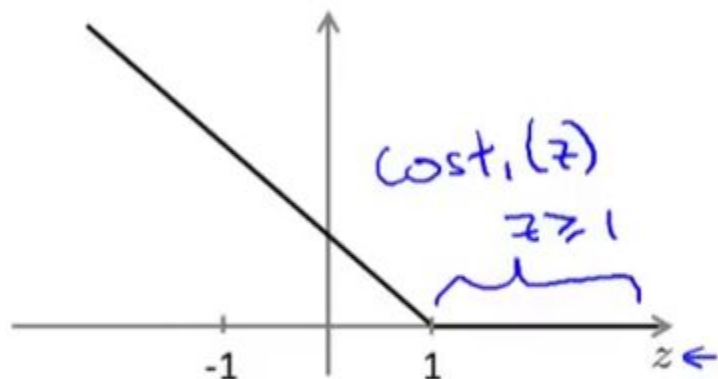


If $y = 1$, we want $\theta^T x \geq 1$ (not just ≥ 0)

If $y = 0$, we want $\theta^T x \leq -1$ (not just < 0)

SVM Cost Function cont

$$\rightarrow \min_{\theta} C \sum_{i=1}^m \left[y^{(i)} \underline{\text{cost}_1(\theta^T x^{(i)})} + (1 - y^{(i)}) \underline{\text{cost}_0(\theta^T x^{(i)})} \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$



➤ If $y = 1$, we want $\theta^T x \geq 1$ (not just ≥ 0)

$$\theta^T x \geq \cancel{0} \quad 1$$

➤ If $y = 0$, we want $\theta^T x \leq -1$ (not just < 0)

$$\theta^T x \leq \cancel{0} \quad -1$$

$$C = 100,000$$

SVM Decision Boundary [2:31-4:39]

$$\min_{\theta} C \left[\sum_{i=1}^m \left[y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) \right] \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

(Note: A blue box highlights the summation term, and a blue arrow points from the \min_{θ} to the box. A blue bracket under the box is labeled $= 0$.)

Whenever $y^{(i)} = 1$:

$$\theta^T x^{(i)} \geq 1$$

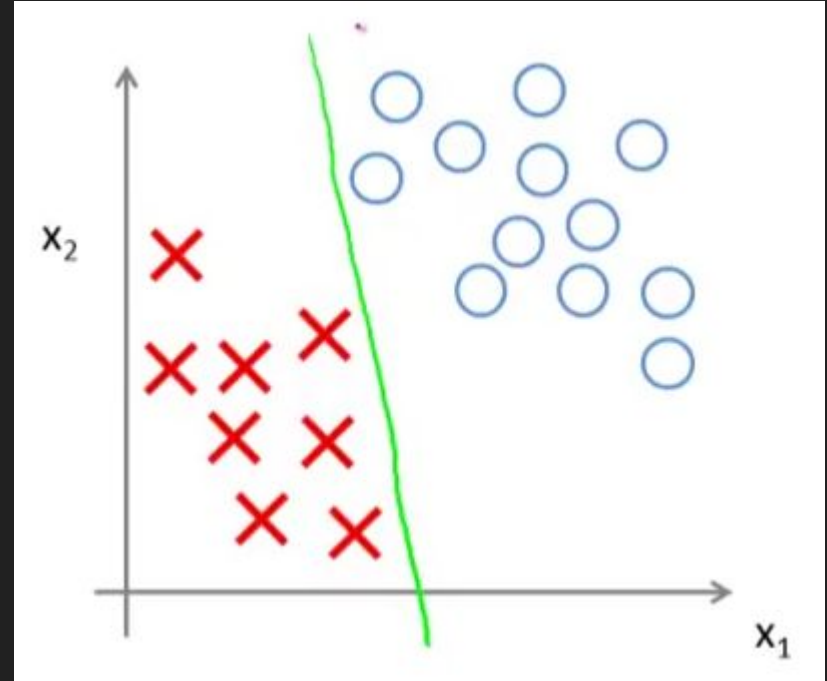
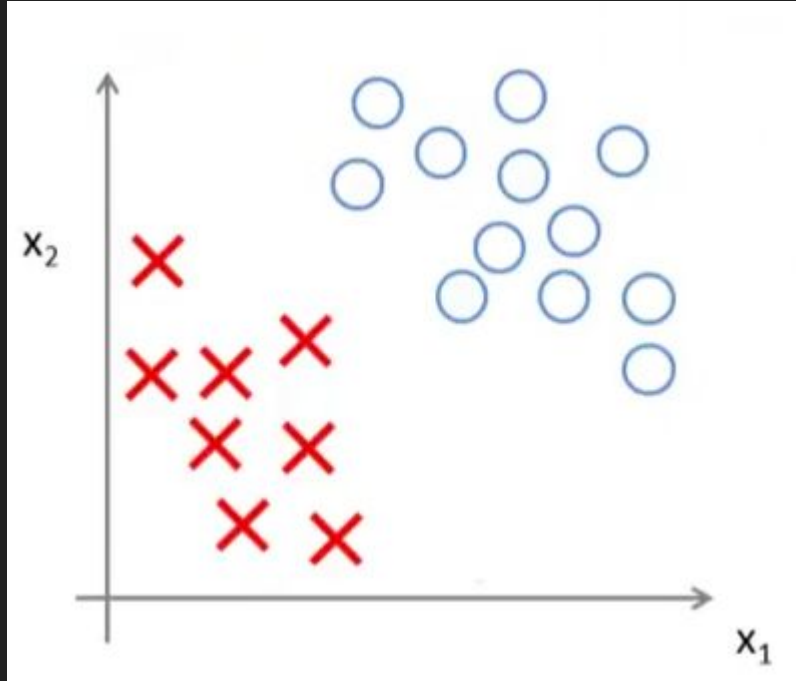
Whenever $y^{(i)} = 0$:

$$\theta^T x^{(i)} \leq -1$$

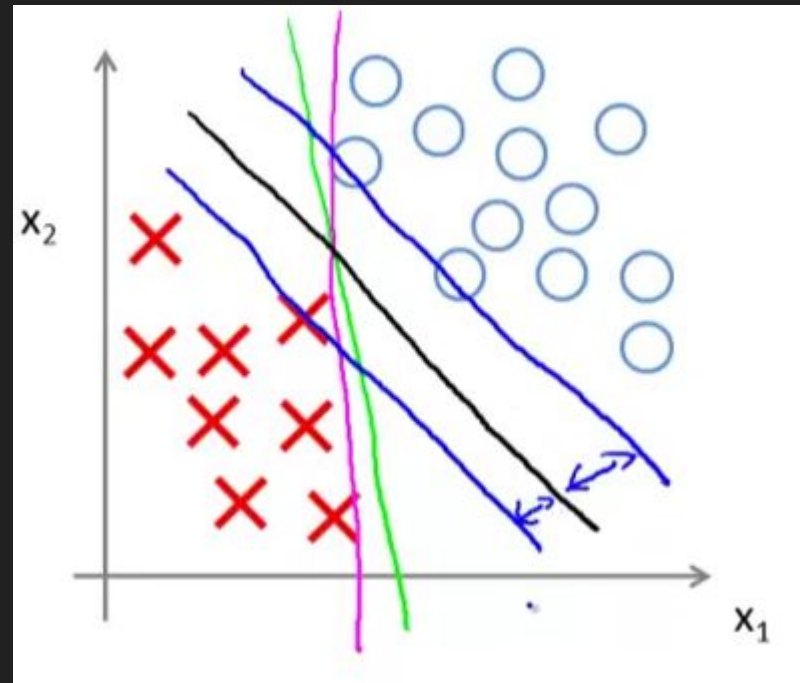
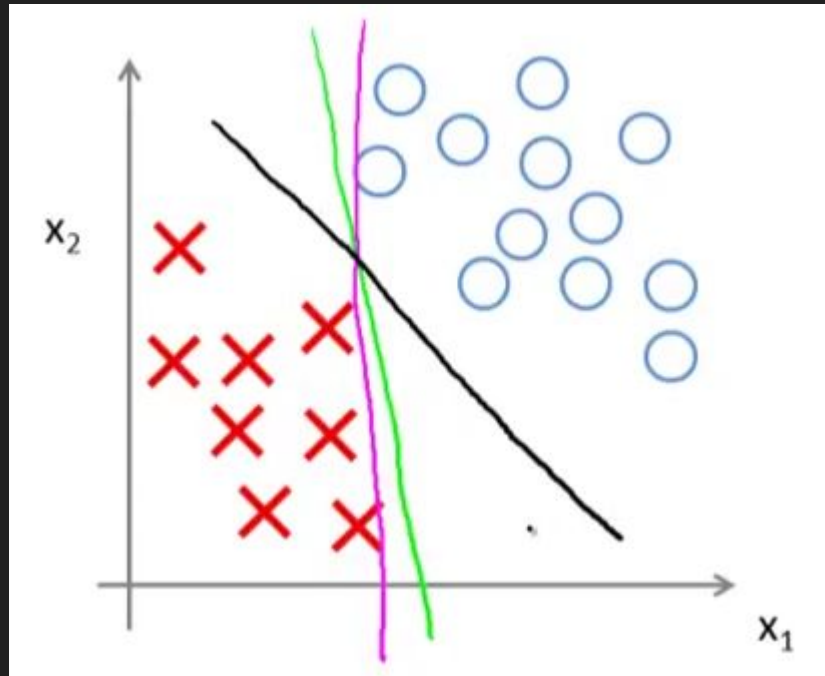
$$\min_{\theta} \cancel{C \times 0} + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

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

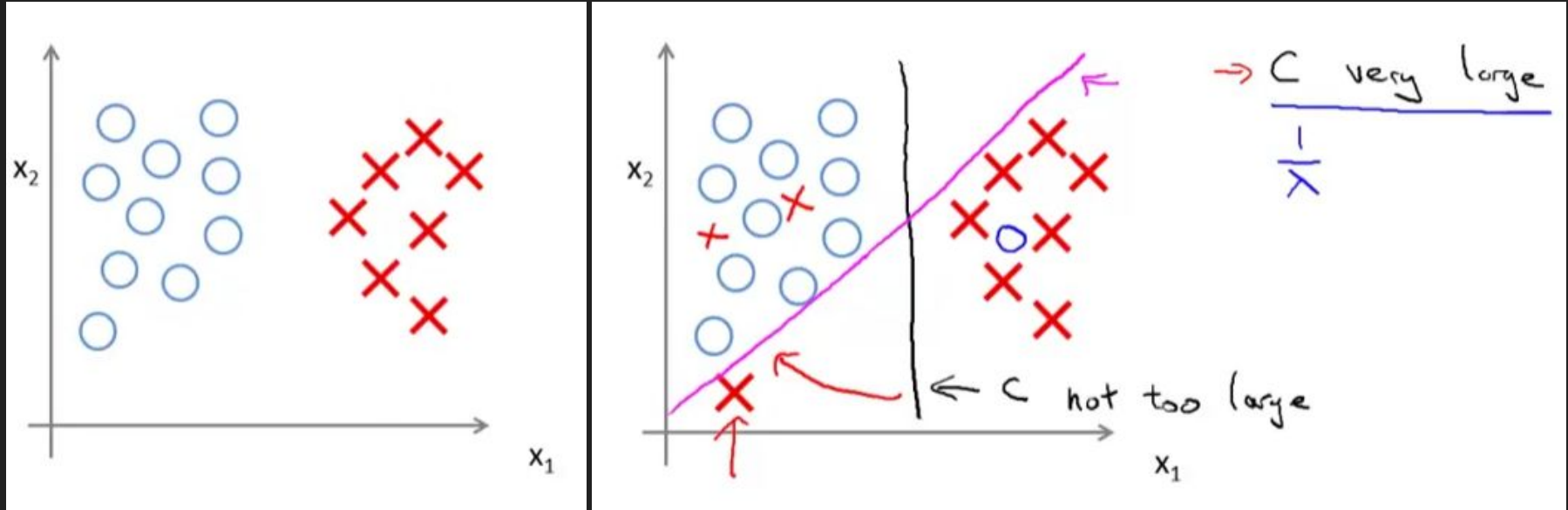
SVM Decision Boundary: Linearly Separable Case



SVM

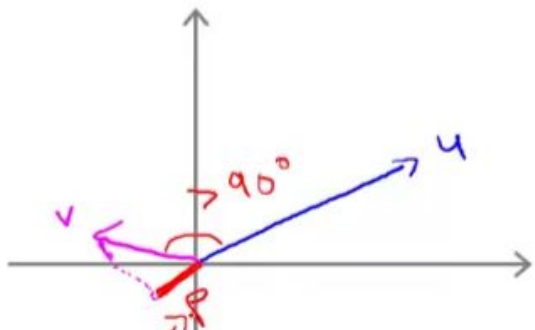
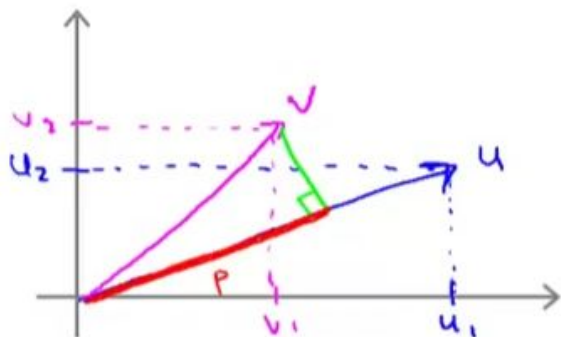


SVM Large Margin Classifier with Outliers [7:19-]



SVM Math - Vector Inner Product [0:00-5:45]

Vector Inner Product



$$\Rightarrow u = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \quad \rightarrow v = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

$$u^T v = ? \quad [u_1 \ u_2] \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

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

p = length of projection of v onto u .

$$\begin{aligned} u^T v &= \underline{p} \cdot \|u\| \leftarrow = v^T u \\ &= u_1 v_1 + u_2 v_2 \leftarrow p \in \mathbb{R} \end{aligned}$$

$$u^T v = p \cdot \|u\|$$

$$p < 0$$

SVM Math - Decision Boundary [5:46-11:11]

SVM Decision Boundary

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

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

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

Simplification: $\theta_0 = 0$. $n=2$

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

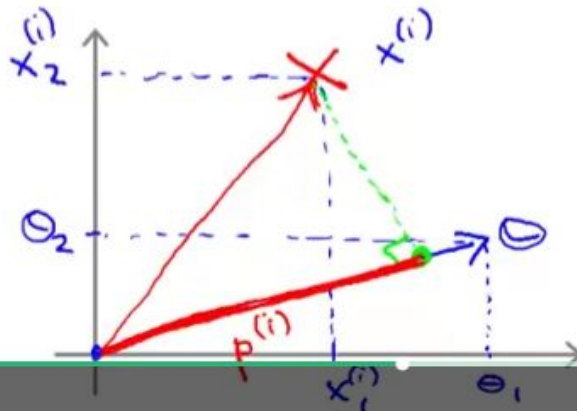
$$= \|\theta\|$$

$$\begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} \quad \theta_0 = 0$$

$$\theta^T x^{(i)} = ?$$

$$\uparrow \quad \uparrow$$

$$u^T v$$



$$\theta^T x^{(i)} = \left[p^{(i)} \cdot \|\theta\| \right]$$

$$= \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)}$$

SVM Math - Decision Boundary

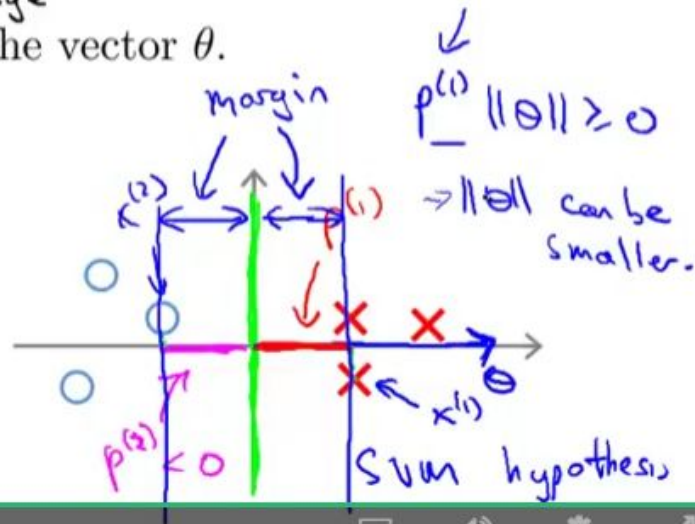
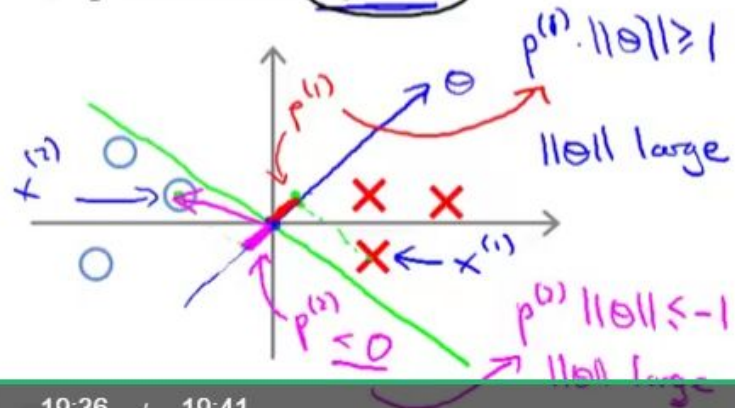
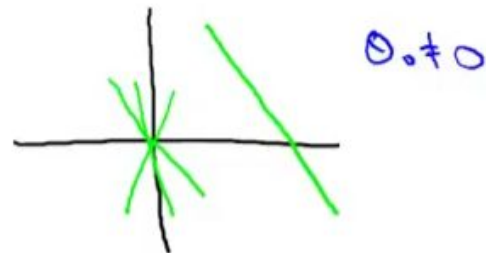
SVM Decision Boundary

$$\min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2 = \frac{1}{2} \|\theta\|^2 \leftarrow$$

$$\text{s.t. } \left. \begin{array}{ll} p^{(i)} \cdot \|\theta\| \geq 1 & \text{if } y^{(i)} = 1 \\ p^{(i)} \cdot \|\theta\| \leq -1 & \text{if } y^{(i)} = -1 \end{array} \right\} \begin{array}{l} C \text{ vary} \\ \text{large} \end{array}$$

where $p^{(i)}$ is the projection of $x^{(i)}$ onto the vector θ .

Simplification: $\theta_0 = 0$



SVM - Using SVM

Use SVM software package (e.g. liblinear, libsvm, ...) to solve for parameters θ .

Need to specify:

→ Choice of parameter C.

Choice of kernel (similarity function):

E.g. No kernel ("linear kernel")

Predict " $y = 1$ " if $\theta^T x \geq 0$

$$\theta_0 + \theta_1 x_1 + \dots + \theta_n x_n \geq 0 \quad \rightarrow \quad \underline{n} \text{ large}, \quad \underline{m} \text{ small} \quad \underline{x \in \mathbb{R}^{n+1}}$$

Gaussian kernel:

$$f_i = \exp\left(-\frac{\|x - l^{(i)}\|^2}{2\sigma^2}\right), \text{ where } l^{(i)} = x^{(i)}.$$

Need to choose σ^2 .

$x \in \mathbb{R}^n$, n small
and/or m large



SVM Kernel (Similarity): Gaussian [4:34-8:25]

Kernel (similarity) functions:

function $f = \text{kernel}(\underline{x1}, \underline{x2})$

$$f = \exp\left(-\frac{\|\underline{x1} - \underline{x2}\|^2}{2\sigma^2}\right)$$

return

$x \rightarrow \begin{matrix} f_1 \\ f_2 \\ \vdots \\ f_m \end{matrix}$

→ Note: Do perform feature scaling before using the Gaussian kernel.

$x \in \mathbb{R}^n$

$$\begin{aligned} \rightarrow \|\underline{x} - \underline{l}\|^2 &= \|\underline{v}\|^2 = v_1^2 + v_2^2 + \dots + v_n^2 \\ &= \underbrace{(x_1 - l_1)^2}_{1000 \text{ feet}^2} + \underbrace{(x_2 - l_2)^2}_{1-5 \text{ bedrooms}} + \dots + (x_n - l_n)^2 \end{aligned}$$

SVM Kernels - Other choices

Other choices of kernel

Note: Not all similarity functions $\text{similarity}(x, l)$ make valid kernels.

- (Need to satisfy technical condition called “Mercer’s Theorem” to make sure SVM packages’ optimizations run correctly, and do not diverge).

Many off-the-shelf kernels available:

- Polynomial kernel:

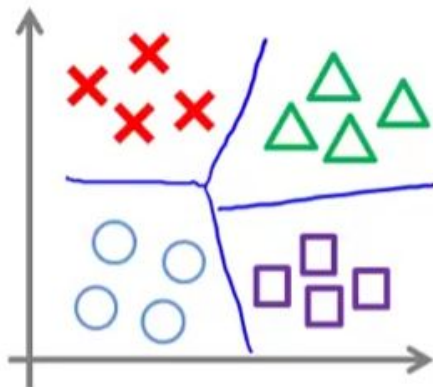
$$k(x, l) = (x^T l)^2, (x^T l)^3, (x^T l + 1)^3, (x^T l + 5)^4, (x^T l + \text{constant})^{\text{degree}}$$

- More esoteric: String kernel, chi-square kernel, histogram intersection kernel, ...

$$\text{sim}(x, l)$$

SVM - Multi-class Classification

Multi-class classification



$$y \in \{1, 2, 3, \dots, K\}$$

↑

Many SVM packages already have built-in multi-class classification functionality.

Otherwise, use one-vs.-all method. (Train K SVMs, one to distinguish $y = i$ from the rest, for $i = 1, 2, \dots, K$), get $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(K)}$
Pick class i with largest $(\theta^{(i)})^T x$

↑ ↑ ↑
 $y=1$ $y=2$ \dots $\theta = K$

Logistic Regression vs. SVMs

Logistic regression vs. SVMs

n = number of features ($x \in \mathbb{R}^{n+1}$), m = number of training examples

→ If n is large (relative to m): (e.g. $n \geq m$, $n = 10,000$, $m = 10 \dots 1,000$)

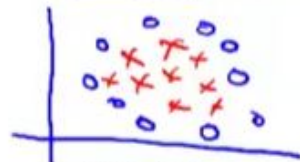
→ Use logistic regression, or SVM without a kernel ("linear kernel")

→ If n is small, m is intermediate: ($n = 1-1,000$, $m = 10-10,000$) ←

→ Use SVM with Gaussian kernel

If n is small, m is large: ($n = 1-1,000$, $m = 50,000+$)

→ Create/add more features, then use logistic regression or SVM without a kernel



→ Neural network likely to work well for most of these settings, but may be slower to train.