

Week 7

- Large Margin Classification:

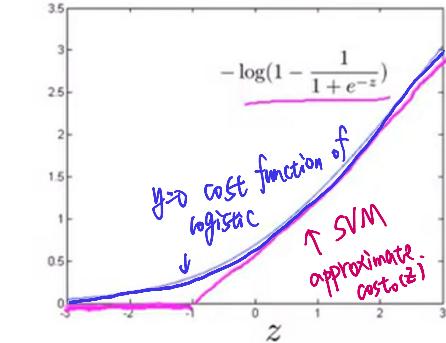
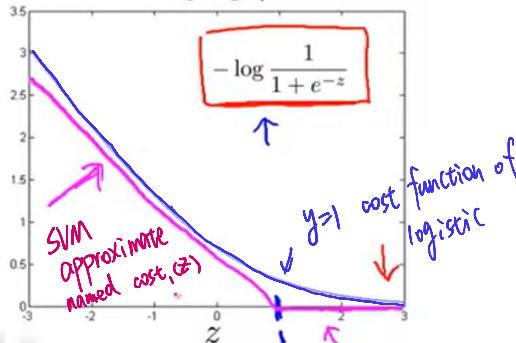
- Optimization Objective:

Alternative view of logistic regression

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

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

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



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:

$$\min_{\theta} \underbrace{\cancel{\frac{1}{m} \sum_{i=1}^m y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)})}}_A + \underbrace{\frac{\lambda}{2} \sum_{j=0}^n \theta_j^2}_B$$

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

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

SVM hypothesis:

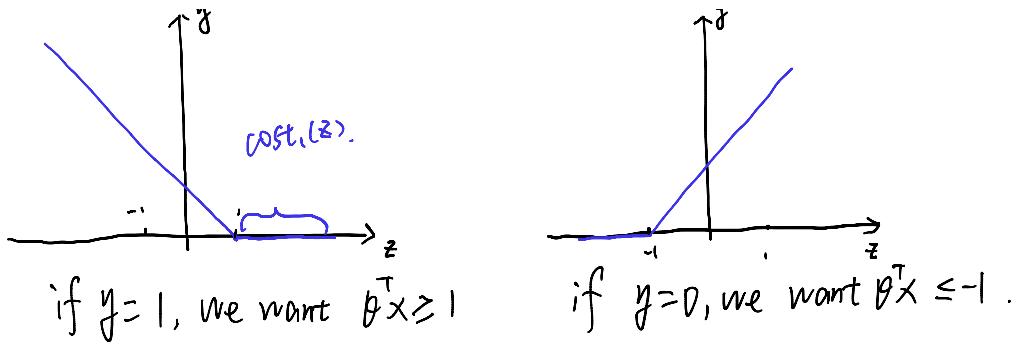
$$\min_{\theta} C \sum_{i=1}^m [y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)})] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

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

A and B just decide how much we care about previous term and last term.

- Large Margin Intuition:



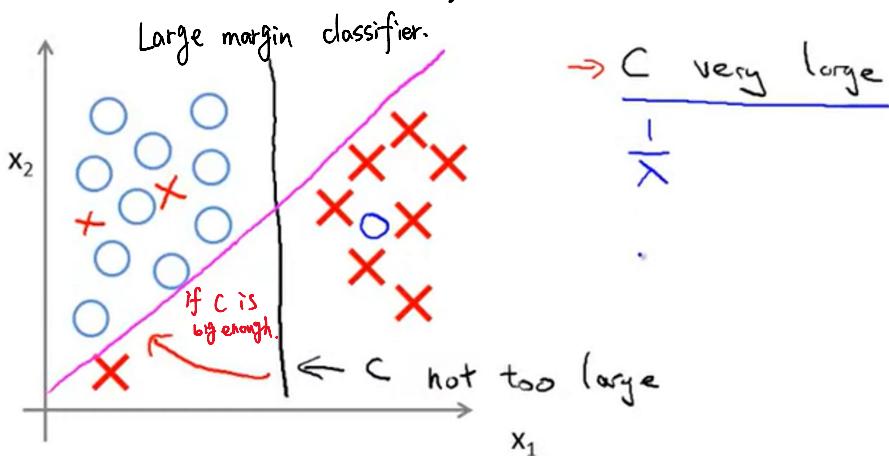
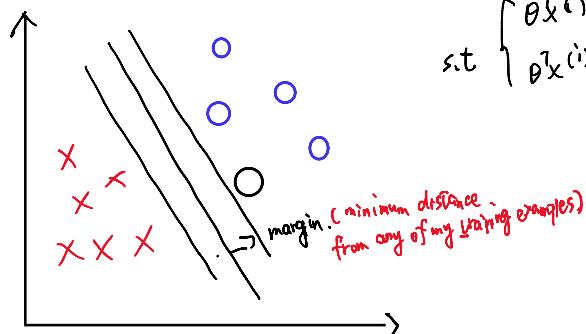


Suppose we set C as a very large value.

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

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

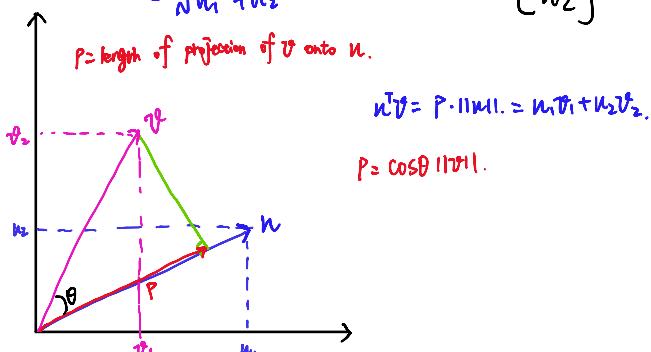
$$\text{s.t. } \begin{cases} \theta^T x^{(i)} \geq 1 & \text{if } y^{(i)} = 1 \\ \theta^T x^{(i)} \leq -1 & \text{if } y^{(i)} = 0 \end{cases}$$



c. Mathematics Behind Large Margin Classification:

$$\|w\| = \text{length of vector} \\ = \sqrt{w_1^2 + w_2^2} \in \mathbb{R}$$

$$w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \quad v = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$



SVM Decision Boundary:

$$\min \frac{1}{2} \sum_{j=1}^n \theta_j^2 \quad (\text{Simplification: } \theta_0 = 0) \quad \text{s.t. } \begin{cases} \theta^T x \geq 1 & \text{if } y^{(i)} = 1 \\ \theta^T x \leq -1 & \text{if } y^{(i)} = 0 \end{cases}$$

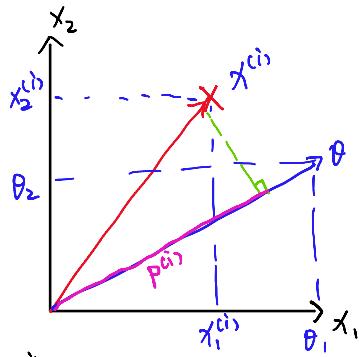
$$= \frac{1}{2} (\theta_1^2 + \theta_2^2) = \frac{1}{2} \underbrace{(\sqrt{\theta_1^2 + \theta_2^2})^2}_{= \|\theta\|^2} = \|\theta\|.$$

what does it do?

$$= \frac{1}{2} \|\theta\|^2.$$

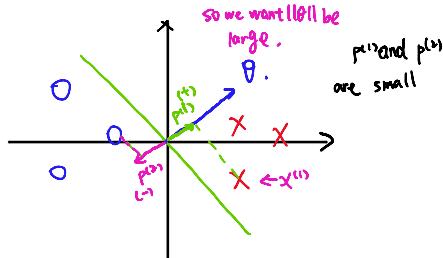
$$\begin{aligned}\theta^T x^{(i)} &= p \cdot \|\theta\| \\ &= \theta_0 x_0^{(i)} + \theta_1 x_1^{(i)}\end{aligned}$$

$$\text{so, } \min \sum_{j=1}^n \theta_j^2 = \frac{1}{2} \|\theta\|^2.$$

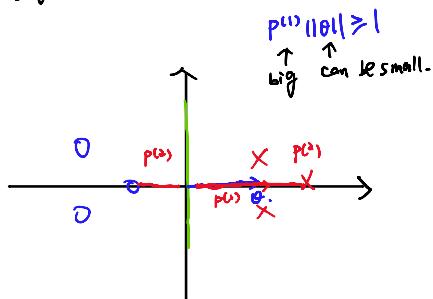


$$\text{s.t. } \begin{cases} p(\theta) \|\theta\| \geq 1 & \text{if } y^{(i)} = 1 \\ p(\theta) \|\theta\| \leq -1 & \text{if } y^{(i)} = 0 \end{cases}$$

① choice I: $p(\theta) \|\theta\| \geq 1$ and $p(\theta)$ is small



② choice II:



• Kernels:

a. Kernels I:

predict $y=1$ if $\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 + \theta_5 x_2^2 + \dots \geq 0$

$$h_\theta(x) = \begin{cases} 1 & \theta_0 + \theta_1 + \dots \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

another form. $\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_4 + \dots$

$$f_1 = x_1, f_2 = x_2, f_3 = x_1 x_2, f_4 = x_1^2, f_5 = x_2^2, \dots$$

Is there a different/better choice of the features f_1, f_2, \dots ?

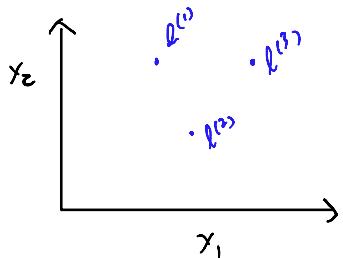
Given x , compute new feature depending on proximity to landmarks $l^{(1)}, l^{(2)}, l^{(3)}$.

$$\text{so, } f_1 = \text{similarity}(x, l^{(1)}) = \exp\left(-\frac{\|x - l^{(1)}\|}{2\sigma^2}\right)$$

$$f_2 = \text{similarity}(x, l^{(2)}) = \exp\left(-\frac{\|x - l^{(2)}\|}{2\sigma^2}\right)$$

$$f_3 = \text{similarity}(x, l^{(3)}) = \exp\left(-\frac{\|x - l^{(3)}\|}{2\sigma^2}\right).$$

Kernel function. (Gaussian Kernel), $k(x, l^{(1)})$.



What these kernels actually do?

Kernels and similarity:

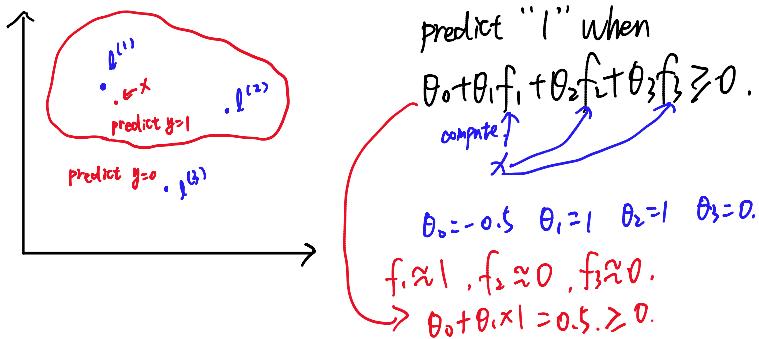
$$f_i = \text{similarity}(x, l^{(i)}) = \exp\left(-\frac{\|x - l^{(i)}\|^2}{2\sigma^2}\right) = \exp\left(-\frac{\sum_{j=1}^n (x_j - l_j^{(i)})^2}{2\sigma^2}\right).$$

If $x \approx l^{(i)}$ $\|x - l^{(i)}\| \approx 0$

$$f_i \approx \exp\left(-\frac{0}{2\sigma^2}\right) \approx 1.$$

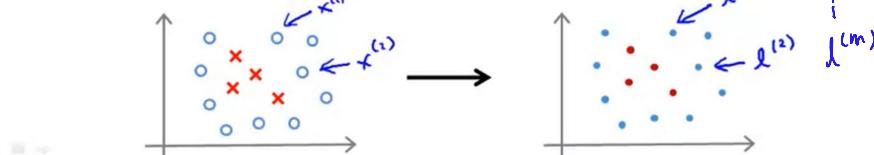
If x is far from $l^{(i)}$

$$f_i = \exp\left(-\frac{(\text{large number})^2}{2\sigma^2}\right) \approx 0$$



b. Kernels II:

Where to get $l^{(1)}, l^{(2)}, l^{(3)}, \dots$?
Where to get $l^{(1)}, l^{(2)}, l^{(3)}, \dots$?



SVM with Kernels

- Given $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})$,
- choose $l^{(1)} = x^{(1)}, l^{(2)} = x^{(2)}, \dots, l^{(m)} = x^{(m)}$.

Given example x :

$$\begin{aligned} \rightarrow f_1 &= \text{similarity}(x, l^{(1)}) \\ \rightarrow f_2 &= \text{similarity}(x, l^{(2)}) \\ \dots \end{aligned}$$

$$f = \begin{bmatrix} f_0 \\ f_1 \\ f_2 \\ \vdots \\ f_m \end{bmatrix} \quad f_0 = 1$$

For training example $(x^{(i)}, y^{(i)})$:

$$\begin{aligned} x^{(i)} \rightarrow f_1^{(i)} &= \text{similarity}(x^{(i)}, l^{(1)}) \\ f_2^{(i)} &= \text{similarity}(x^{(i)}, l^{(2)}) \\ \vdots \\ f_m^{(i)} &= \text{similarity}(x^{(i)}, l^{(m)}) \\ \text{similarity}(x^{(i)}, l^{(i)}) &= \exp\left(-\frac{\|x^{(i)} - l^{(i)}\|^2}{2\sigma^2}\right) = 1 \end{aligned}$$

$$\begin{aligned} x^{(i)} \in \mathbb{R}^{n+1} &\quad (\text{or } \mathbb{R}^n) \\ f^{(i)} = \begin{bmatrix} f_0^{(i)} \\ f_1^{(i)} \\ f_2^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix} & \quad f_0^{(i)} = 1 \end{aligned}$$

Training

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

θ_0 do not regularize θ_0 .

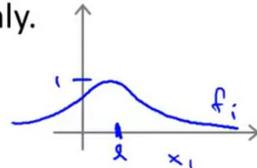
$C = \frac{1}{\lambda}$ Large C : Lower bias, high variance.

Small C : Higher bias, low variance.

σ^2 Large σ^2 : Features f_i vary more smoothly.

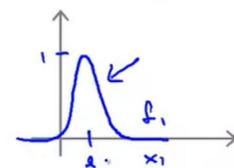
→ Higher bias, lower variance.

$$\exp\left(-\frac{\|x - l^{(i)}\|^2}{2\sigma^2}\right)$$



Small σ^2 : Features f_i vary less smoothly.

Lower bias, higher variance.



• SVMs in Practice:

a. Using An 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 $\underline{\theta^T x} \geq 0$

$$\theta_0 + \theta_1 x_1 + \dots + \theta_n x_n \geq 0$$

n large, m small

$$x \in \mathbb{R}^{n+1}$$

many features, few data, so we want a linear boundary.

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 .

if you have \uparrow features of very different scales.
perform feature scaling.

When to choose Gaussian Kernel.

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



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).

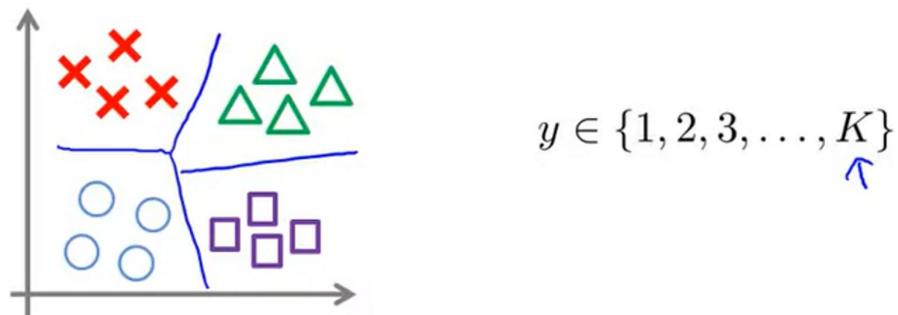
Many off-the-shelf kernels available:

- Polynomial kernel: $k(x, l) = \langle x^T l + \text{constant} \rangle^{\text{degree}}$

$$\langle x^T l \rangle^1, \langle x^T l \rangle^2, \langle x^T l \rangle^3, \langle x^T l \rangle^4, \langle x^T l \rangle^5$$

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

Multi-class classification



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

$$y=1 \quad y=2 \quad \dots \quad \theta=K$$

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. $n \geq m$, $n = 10,000$, $m = 10 \dots 1000$)
→ Use logistic regression, or SVM without a kernel ("linear kernel")
- If n is small, m is intermediate: ($n = 1 \dots 1000$, $m = 10 \dots 10,000$)
→ Use SVM with Gaussian kernel
- If n is small, m is large: ($n = 1 \dots 1000$, $m = 50,000+$)
→ 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.

