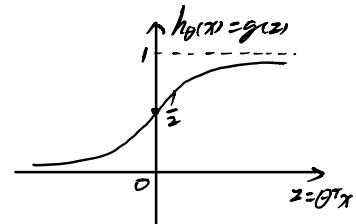


Large Margin classification

alternative view of logistic regression

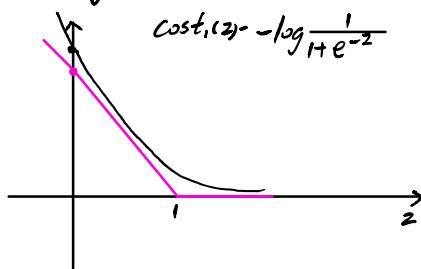
$$h_{\theta}(x) = \frac{1}{1+e^{-\theta^T x}} \quad \begin{array}{l} \text{if } y=1, \text{ we want } h_{\theta}(x) \approx 1, \theta^T x \gg 0 \\ \text{if } y=0, \text{ we want } h_{\theta}(x) \approx 0, \theta^T x \ll 0 \end{array}$$



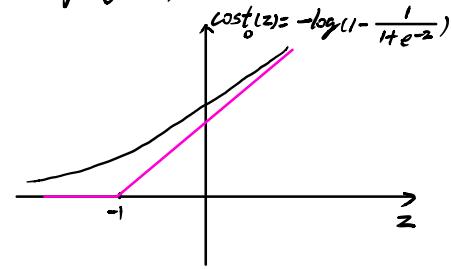
$$\text{single example cost} = -(y \log h_{\theta}(x) + (1-y) \log(1-h_{\theta}(x)))$$

$$= -y \log \frac{1}{1+e^{-z}} - (1-y) \log(1 - \frac{1}{1+e^{-z}})$$

if $y=1$, want $z = \theta^T x \gg 0$



if $y=0$, want $z = \theta^T x \ll 0$



we replace the black line with two magenta line:

$$\text{cost}_0(z) \quad \text{cost}_1(z)$$

logistic regression:

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

Support vector machine:

get rid of m , optima θ will remain unchanged, replace $-\log h_{\theta}(x^{(i)})$ with $\text{cost}_1(\theta^T x^{(i)})$, $\log(1-h_{\theta}(x^{(i)}))$ with $\text{cost}_0(\theta^T x^{(i)})$,

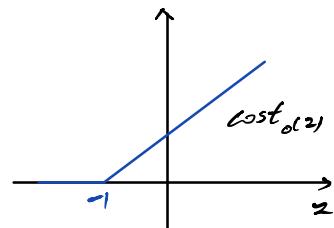
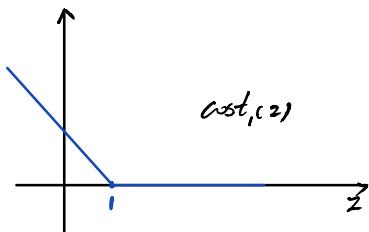
get rid of λ and add C , if $C = \frac{1}{\lambda}$, optima θ remain unchanged.

SVM hypothesis:

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

(this form just for convention)

Large Margin intuition:



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

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

if let C be large, $\sum_{i=1}^m [y^{(i)} \text{cost}_l(\theta^T x^{(i)}) + (1-y^{(i)}) \text{cost}_o(\theta^T x^{(i)})] \approx 0$

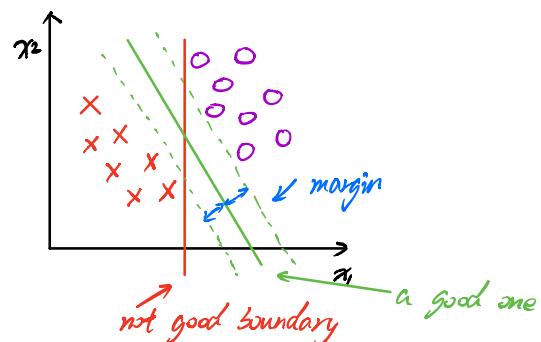
this will be $C \times 0 + \frac{1}{2} \sum_{i=1}^n \theta_i^2$:

$$\min_{\theta} \frac{1}{2} \sum_{i=1}^n \theta_i^2$$

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

SVM decision boundary example:

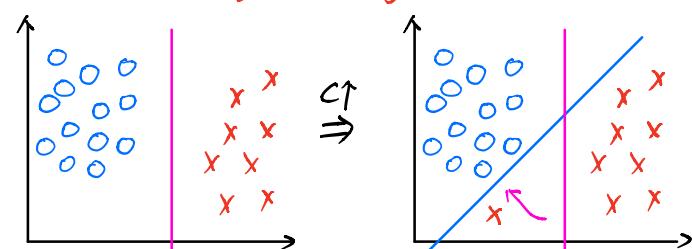
linear separable case:



in presence of outliers

C sensitive

$C \uparrow \nmid \uparrow \lambda \downarrow$ variance

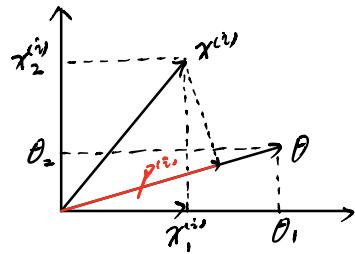


SVM:

$$\min_{\theta} \frac{1}{2} \sum_{i=1}^n \theta_i^2 \quad \text{s.t. } \theta^T x^{(i)} \geq 1 \text{ if } y=1 \\ \theta^T x^{(i)} \leq -1 \text{ if } y=0$$

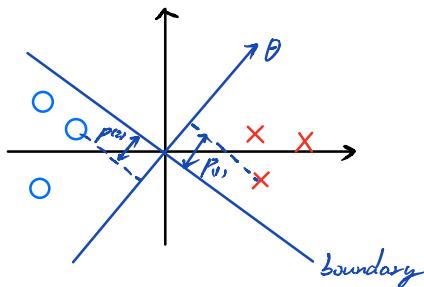
Simplified: $n=2$, and $\theta_0=0$

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



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

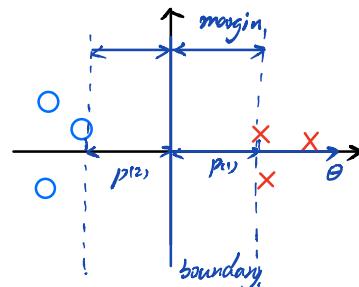
$$\min_{\theta} \frac{1}{2} \|\theta\|^2 \quad \text{s.t. } \begin{cases} p^{(i)} \|\theta\| \geq 1 & \text{if } y=1 \\ p^{(i)} \|\theta\| \leq -1 & \text{if } y=0 \end{cases} \quad (\text{c very large}) \\ \theta_0=0$$



$$p^{(1)} \|\theta\| \geq 1 \quad p^{(2)} \|\theta\| \leq -1$$

$p^{(1)} > 0$, $p^{(2)} < 0$ and $p^{(1)}, |p^{(2)}|$ quite small
so $\|\theta\|$ need to be large, but $\min_{\theta} \frac{1}{2} \|\theta\|^2$
need $\|\theta\|$ to be small.

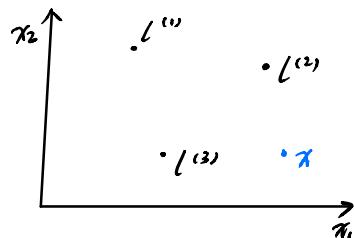
X



$p^{(1)}, |p^{(2)}|$ the biggest, and $\|\theta\|$
can be quite small

✓

Kernel I

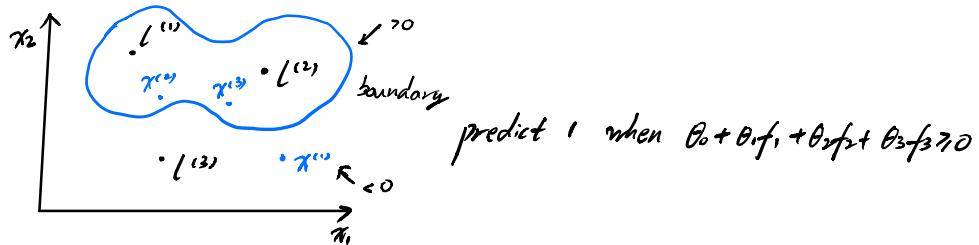
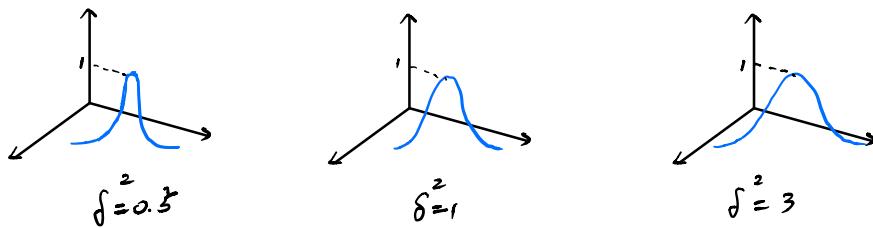


given x , compute new feature depending on proximity to landmarks $l^{(1)}$ $l^{(2)}$ $l^{(3)}$

Gauss kernel:

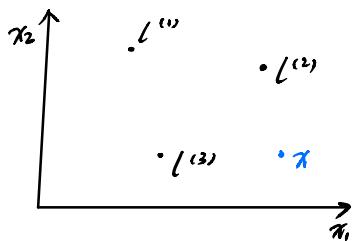
$$f_i = \text{similarity}(x, l^{(i)}) = \exp\left(-\frac{\|x - l^{(i)}\|^2}{2\delta^2}\right)$$

$$\begin{aligned} \text{if } x \approx l^{(i)} & \quad f_i \approx 1 \\ x \text{ far from } l^{(i)}, f_i \approx 0 & \end{aligned} \quad \Rightarrow f_i \in [0, 1]$$



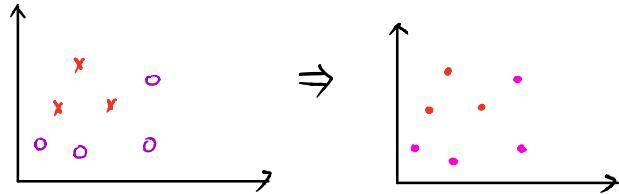
Kernel II:

how to choose landmark? where to get $l^{(1)}, l^{(2)}, l^{(3)}, \dots$?



given set $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

choose $l^{(1)} = x^{(1)}, l^{(2)} = x^{(2)}, l^{(3)} = x^{(3)}, \dots, l^{(m)} = x^{(m)}$



use $f_i = \text{similarity}(x_i | l^{(i)})$: add f_0

$$x^{(i)} \xrightarrow{f} \begin{bmatrix} f_0^{(i)} \\ f_1^{(i)} \\ f_2^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix} \quad x^{(i)} \in \mathbb{R}^{n+1} \quad f_i \in \mathbb{R}^{m+1}$$

predict 1 if $\theta^T f \geq 0$

training: $\min_{\theta} C \sum_{i=1}^m y^{(i)} \text{cost}_+(\theta^T f^{(i)}) + (1-y^{(i)}) \text{cost}_-(\theta^T f^{(i)}) + \frac{1}{2} \sum_{j=1}^m \theta_j^2$ from θ .

$$C \gg 0, \text{ SVM } J(\theta) = \frac{1}{2} \sum_{i=1}^m \theta_i^2$$

$$\sum_{i=1}^m \theta_i^2 = \theta^T \theta = \|\theta\|^2 \quad (\text{no kernel/linear kernel})$$

$$\downarrow$$

$$\theta^T M \theta$$

SVM parameters:

$$C (= \frac{1}{\lambda}) \quad C \uparrow \text{bias} \uparrow \text{variance} \uparrow$$

$$C \uparrow \text{bias} \uparrow \text{variance} \uparrow$$

$$\sigma^2 \quad \sigma^2 \uparrow \text{fi more smoothly, bias} \uparrow \text{variance} \uparrow$$

$$\sigma^2 \downarrow \text{fi less smoothly, bias} \uparrow \text{variance} \uparrow$$

Use SVM:

- choose C , σ^2 , kernel
- do perform feature scaling before using gauss kernel.

multi-class classification

one-vs-all method, train k SVM, distinguish $y=i$ from the others
 $i=1, 2, \dots, k$, get $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(k)}$, pick class i with largest $(\theta^{(i)})^T x$

LR (logistic regression) vs SVM

n feature $x \in \mathbb{R}^{n+1}$, m = training examples

if $n \gg m$: LR, SVM without kernel (linear kernel)

n is small, m is intermediate : SVM with Gaussian kernel

n is small, m is large : create more features, use LR or SVM
without a kernel
neural network is likely to work well for most these setting but slower to train.