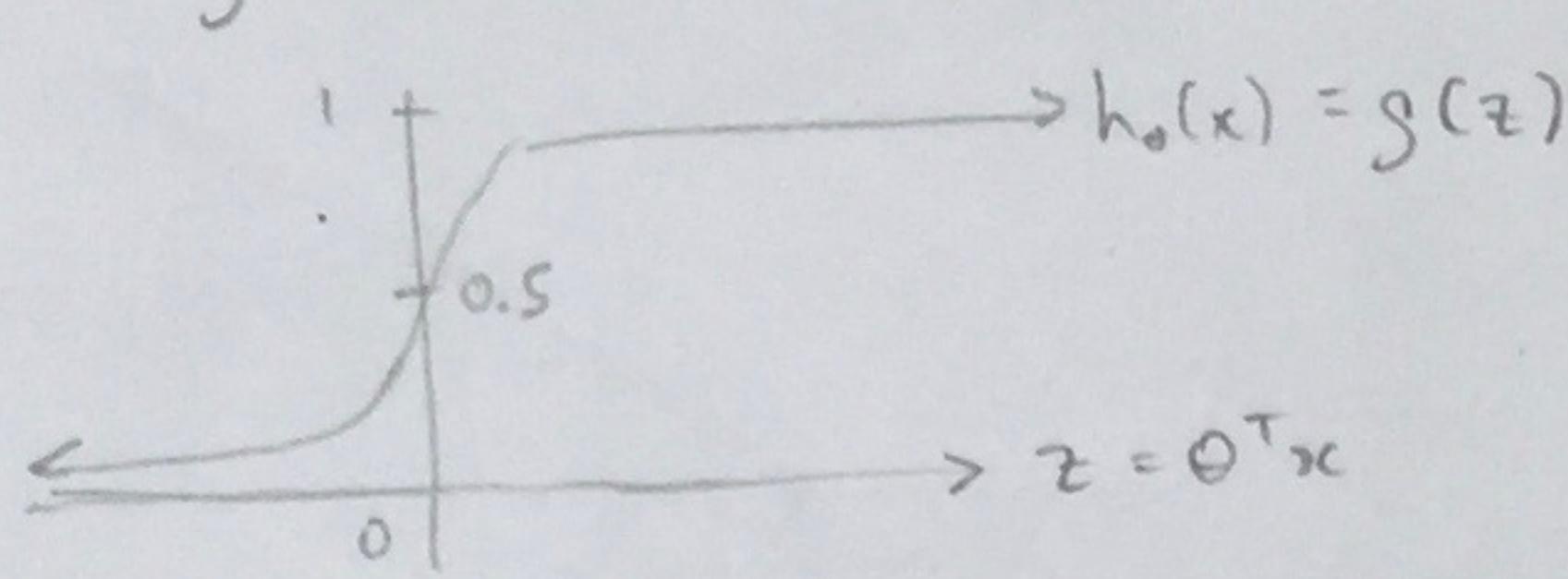


Support Vector Machines: powerful supervised learning algorithm

⇒ recall logistic regression: sigmoid activation function

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

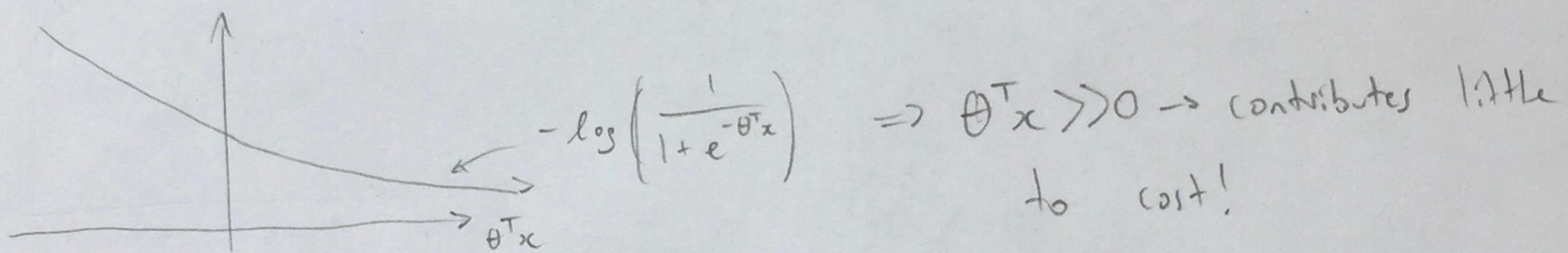


$$\left( g = \frac{1}{1 + e^{-z}} \right)$$

↳ need  $\theta^T x \gg 0$  for  $h_{\theta}(x) \approx 1$  and vice versa for  $h_{\theta}(x) \approx 0$

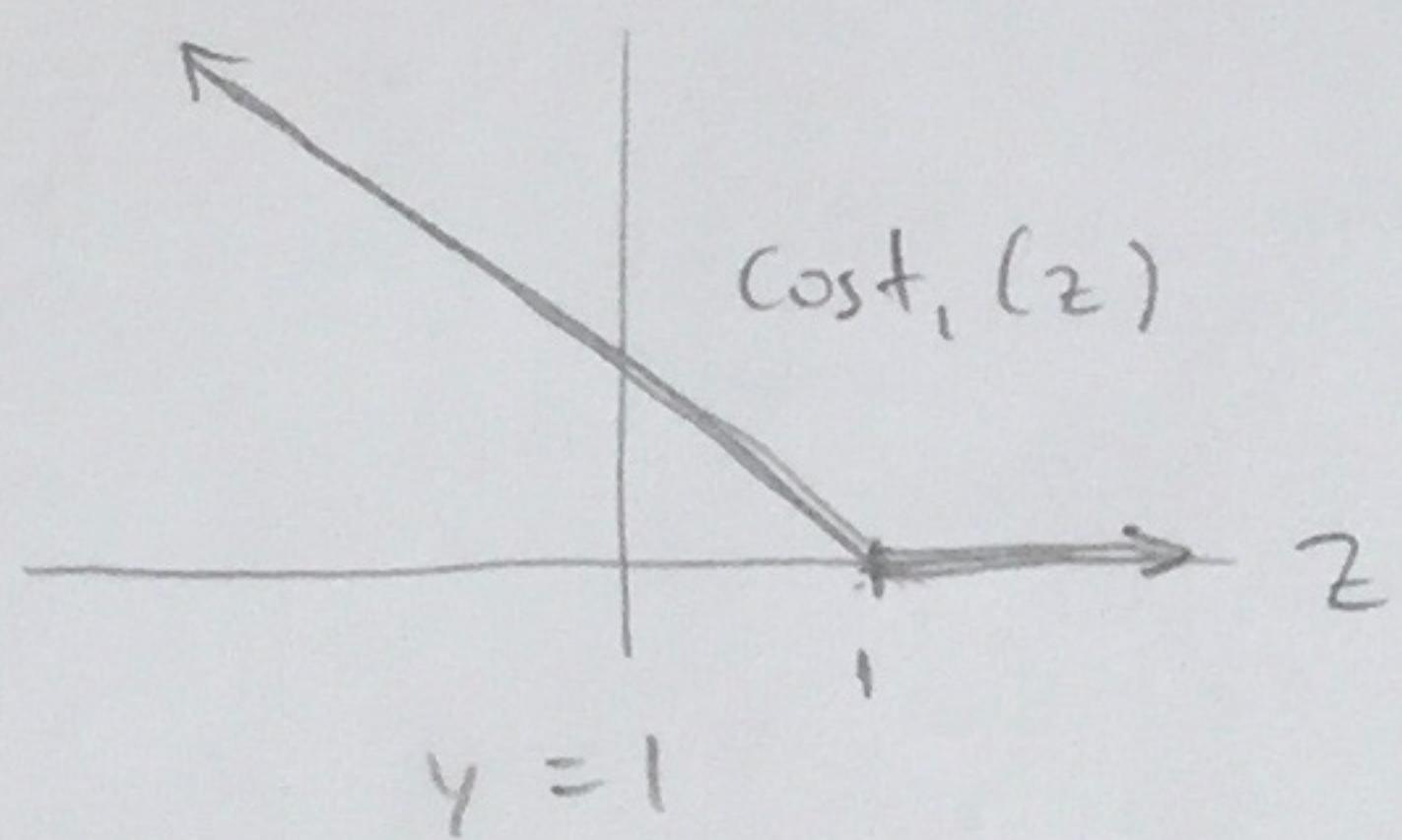
cost of single example  $J_i = - \left[ y^{(i)} \log(h_{\theta}(x)) + (1-y^{(i)}) \log(1-h_{\theta}(x)) \right]$  for  $(x^{(i)}, y^{(i)})$

↳ if  $y=1$ : only term ①:

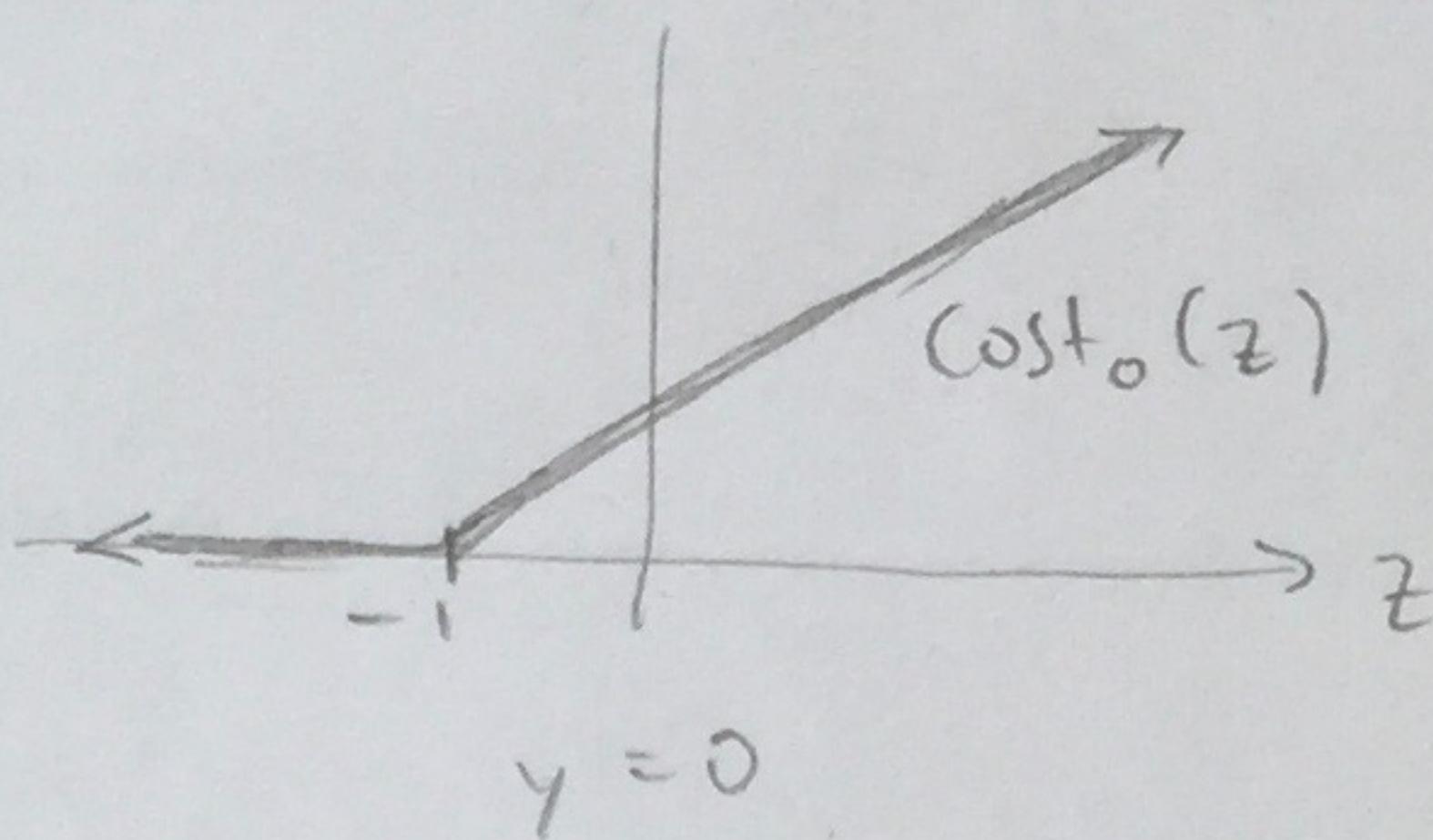


$\Rightarrow$  Note: rather than  $A + \lambda B \rightarrow CA + B$  ( $C \propto \frac{1}{\lambda}$ )

SVM hypothesis  $\begin{cases} 1 & \text{if } \theta^T x \geq 0 \\ 0 & \text{otherwise} \end{cases}$



want  $z = \theta^T x \geq 1$  (not 0)

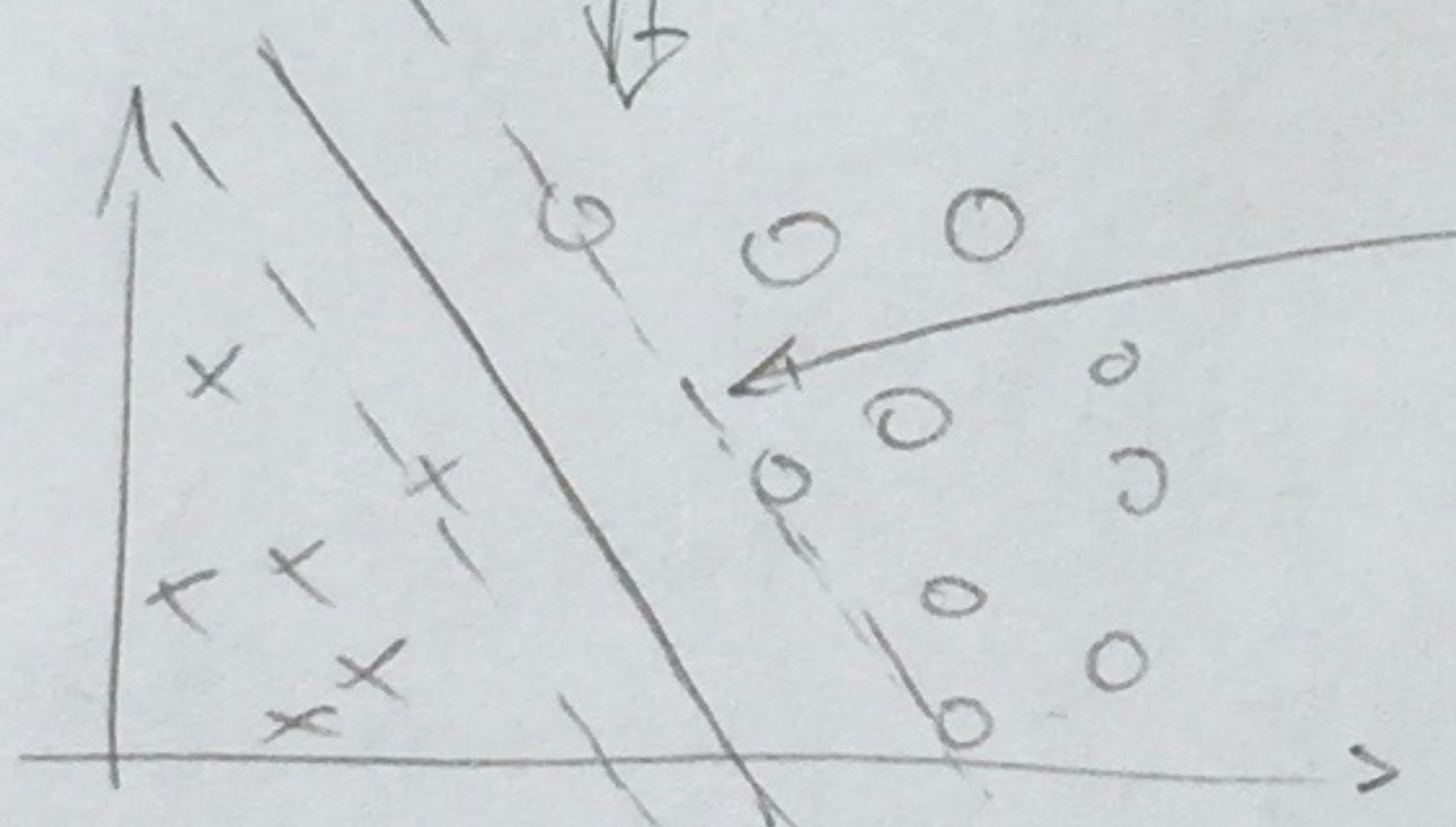


want  $z = \theta^T x \leq -1$  (not 0)

$\hookrightarrow$  SVM's give some "Safety factor" margin ( $|1|$  rather than 0)

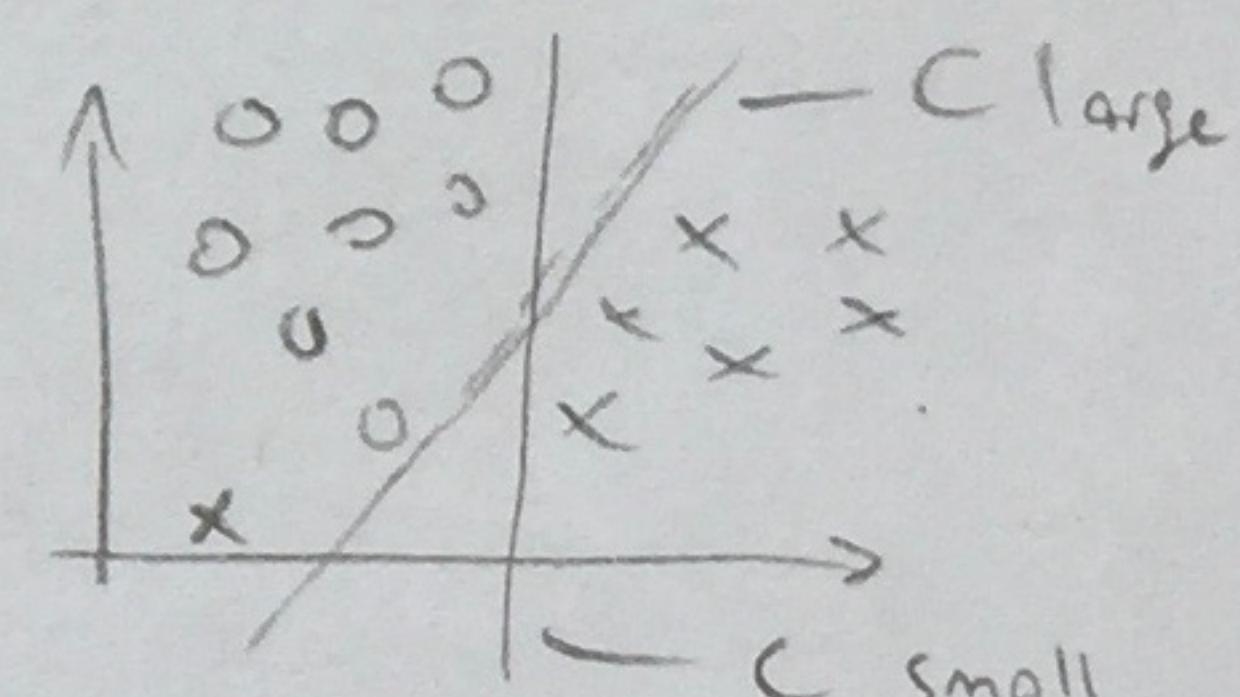
$$\min_{\theta} C \times 0 + \frac{1}{2} \sum_{j=1}^n \theta_j^2 \quad \text{subject to} \quad \begin{aligned} \theta^T x &\geq 1 && \text{if } y^{(i)} = 1 \\ \theta^T x &\leq -1 && \text{if } y^{(i)} = 0 \end{aligned}$$

SVM  
(Linearly Separable)



SVM's ensure maximization of margins (---) "LARGE MARGIN CLASSIFIER" (given  $C$  large)

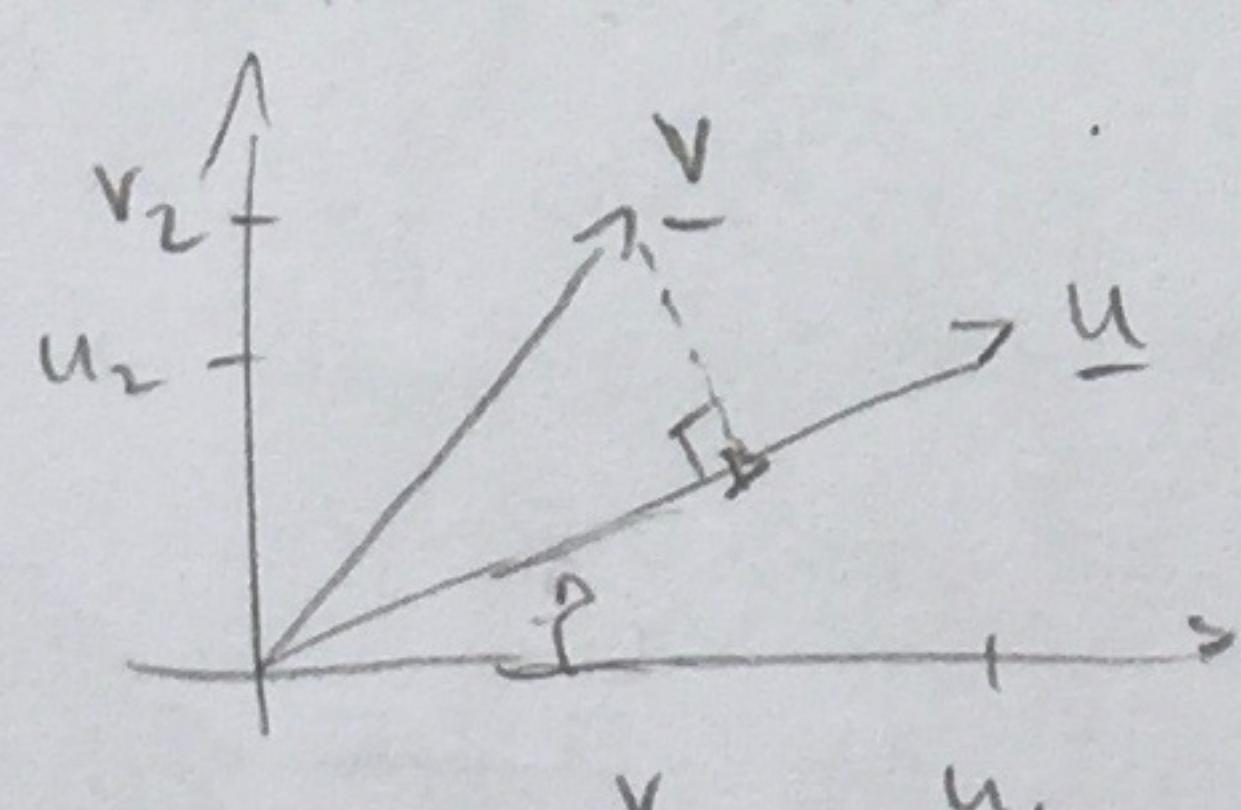
$\Rightarrow$  if  $C$  large, sensitive to outliers  $\rightarrow$



### Mathematics of Large Margin Classifier:

$\Rightarrow$  vector inner product:

$$u^T v \quad (u \cdot v)$$



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

$p$  = length of projection of  $v$  onto  $u$   
(positive or negative)

$$u^T v = p \cdot \|u\| = v^T u = u_1 v_1 + u_2 v_2$$

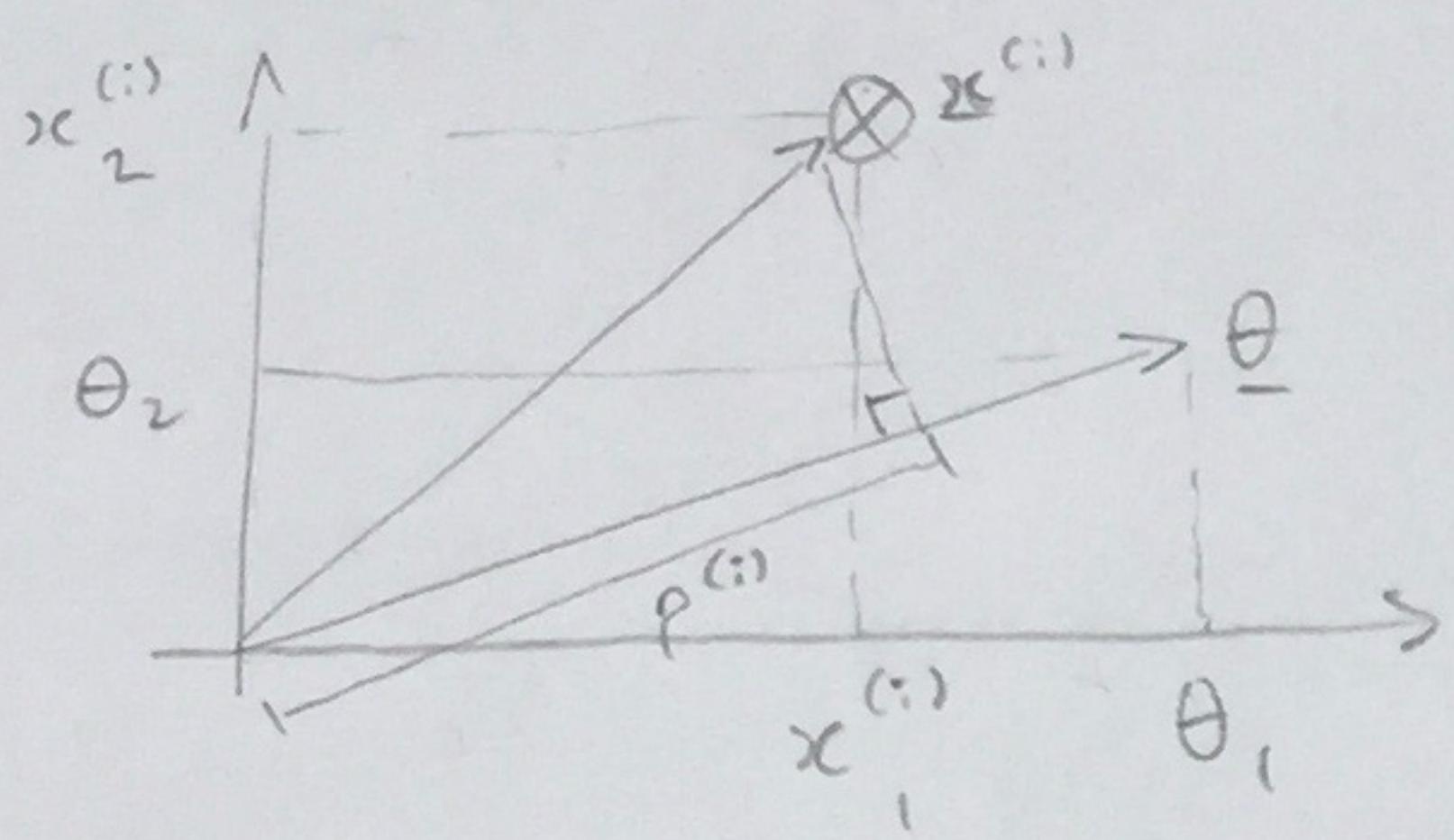
$\Rightarrow$  SVM Decision Boundary  $\rightarrow$  ignore  $\theta_0$ ,  $n=2$

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

such that  $\theta^T x^{(i)} > 1$  for  $y^{(i)} = 1$

$\theta^T x^{(i)} \leq -1$  for  $y^{(i)} = 0$

MINIMIZE  
SQUARED NORM  
OF  $\theta$

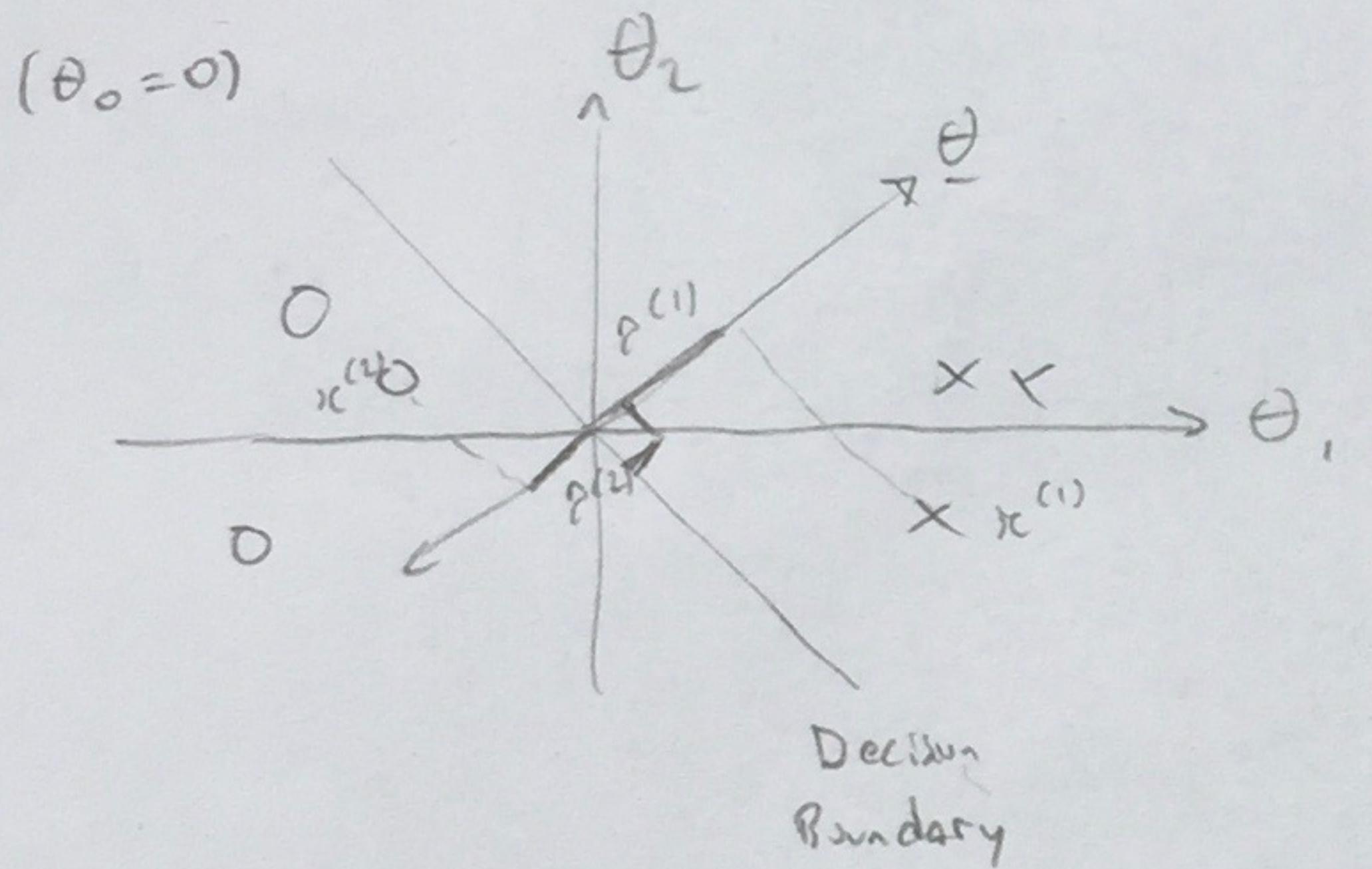


$$\Rightarrow \theta^T x^{(i)} = \rho^{(i)} \|\theta\| = \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)}$$

WHEN C LARGE:

$$\Rightarrow \min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2 = \frac{1}{2} \|\theta\|^2 \text{ such that } \rho^{(i)} \|\theta\| \geq 1 \text{ if } y^{(i)} = 1$$

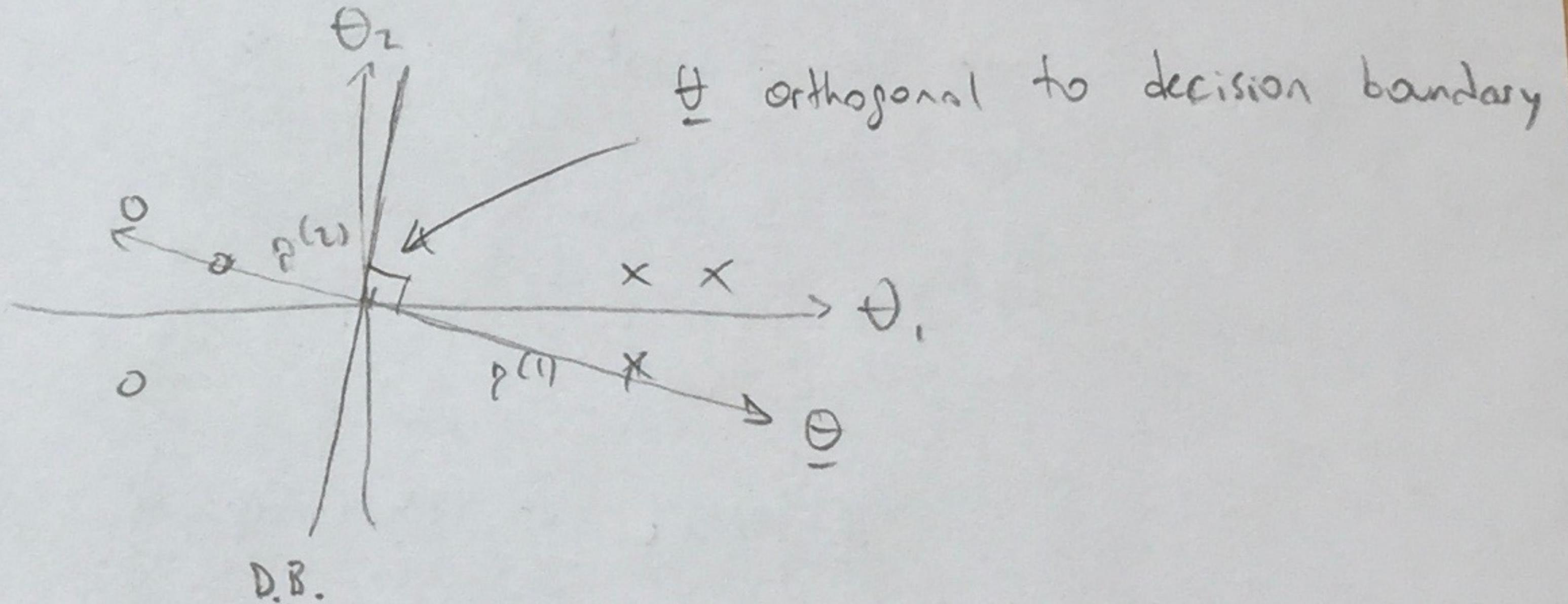
$$\rho^{(i)} \|\theta\| \leq -1 \text{ if } y^{(i)} = 0$$



need  $\rho^{(i)} \|\theta\| \leq -1$ , but  $\rho^{(i)}$  small!

need  $\rho^{(i)} \|\theta\| \geq 1$ , but  $\rho^{(i)}$  small!

BAD DECISION BOUNDARY



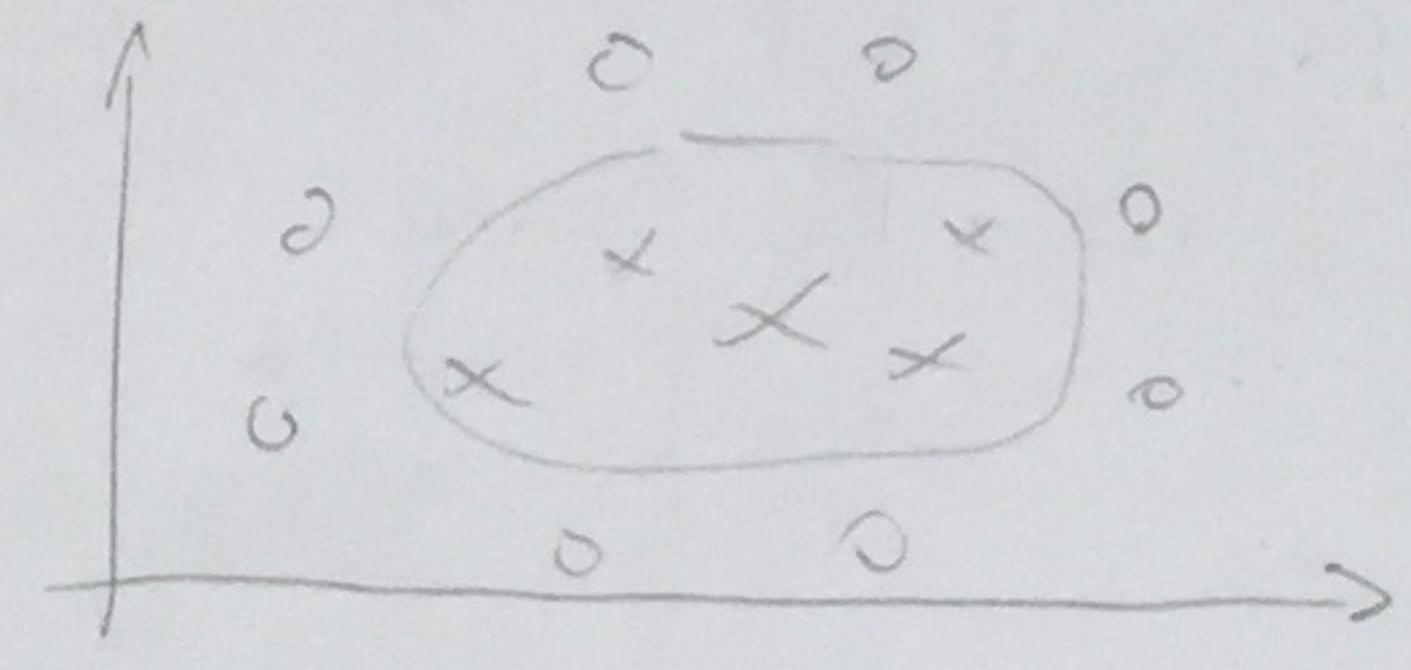
need  $\rho^{(i)} \|\theta\| \geq 1 \rightarrow \rho^{(i)}$  larger, so

$\|\theta\|$  can be smaller! (same for  $\rho^{(i)}$ )

$\Rightarrow$  recall: goal is minimize  $\|\theta\|$

$\hookrightarrow$  BETTER DECISION BOUNDARY

## Non-Linear Decision Boundaries w/ SVM's: KERNELS



predict  $y=1$  if

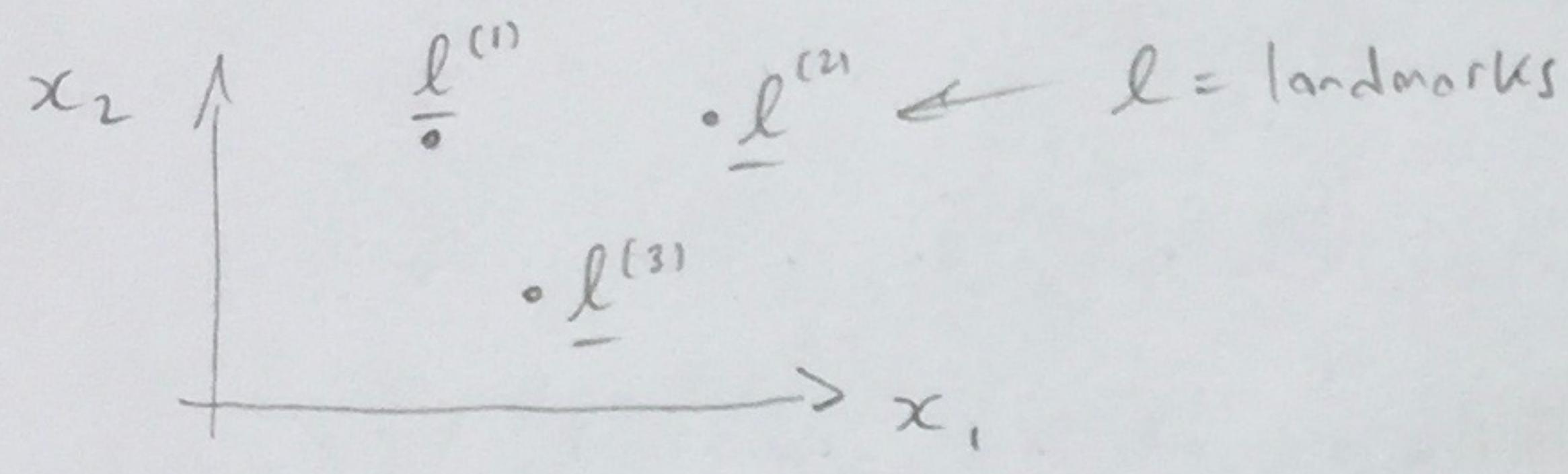
$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \dots + \theta_5 x_2^2 + \dots > 0$$

$f_j = \text{feature}:$

$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 + \dots$$

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

Kernel: given  $\underline{x}$ , compute new feature depending on proximity to landmarks:



euclidean distance  $\underline{w}$

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

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

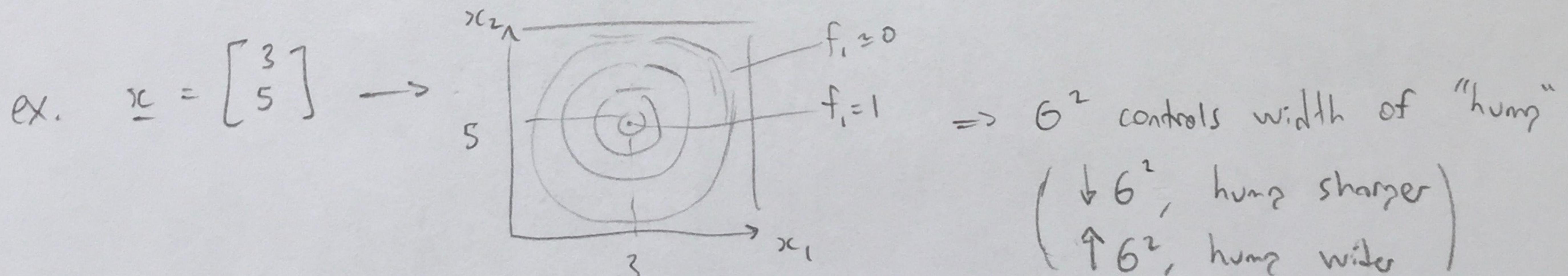
GAUSSIAN  
KERNEL

KERNEL = "SIMILARITY FUNCTION"

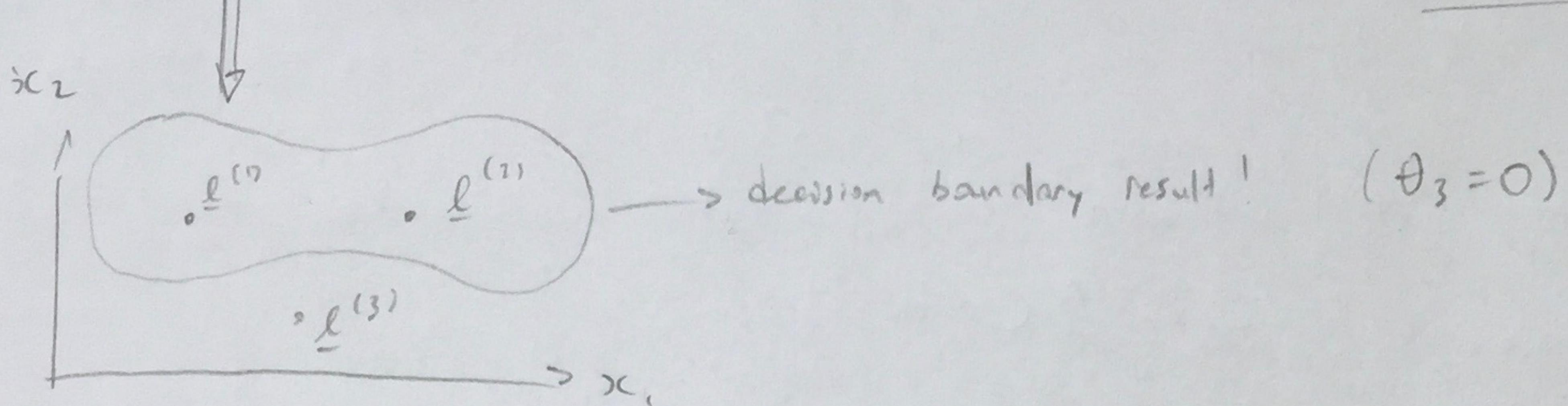
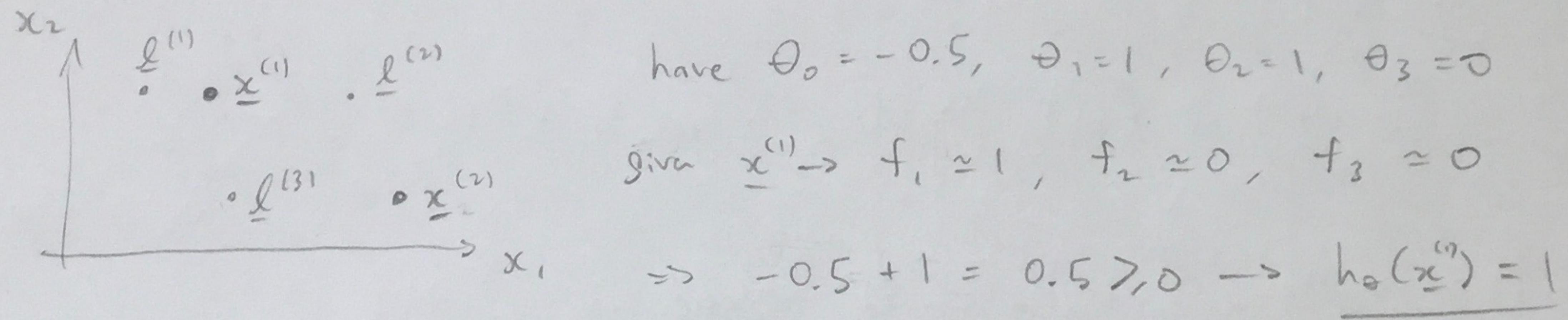
$$* \text{note } \|\underline{x} - \underline{l}^{(1)}\|^2 = \sum_{j=1}^n (x_j - l_j^{(1)})^2$$

$\Rightarrow$  so if  $\underline{x} \approx \underline{l}^{(1)}$   $\rightarrow f_1 \approx 1$ , or if  $\underline{x} \neq \underline{l}^{(1)}$  at all,  $f_1 \approx 0$

$\hookrightarrow$  given one  $\underline{x}$   $\rightarrow$  3 features (one from each landmark)



predict 1 when  $\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 > 0$



$\Rightarrow$  where to get landmarks? How many?

$\hookrightarrow$  put landmarks @ locations of training examples ("m" landmarks)

$\Rightarrow$  measures how close a new point is to data in training set

$\Rightarrow$  given example  $\underline{x}$ :  $f_1 = \text{similarity } (\underline{x}, \underline{l}^{(1)})$  where  $\underline{l}^{(1)} = \underline{x}^{(1)}$   
 $f_2 = \text{similarity } (\underline{x}, \underline{l}^{(2)})$  where  $\underline{l}^{(2)} = \underline{x}^{(2)}$  etc.

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

$[(m+1) \times 1]$

$\Rightarrow$  for  $\underline{x}^{(i)} \rightarrow f_i = \text{sim}(\underline{x}^{(i)}, \underline{l}^{(i)}) \Rightarrow \boxed{\underline{x}^{(i)} \in \mathbb{R}^{n+1} \rightarrow \underline{f}^{(i)} \in \mathbb{R}^{m+1}}$

(one will be  $\text{sim}(\underline{x}^{(i)}, \underline{l}^{(i)}) = 1$ )

$\Rightarrow$  predict  $y=1$  if  $\underline{\theta}^T \underline{f} \geq 0$  where  $\underline{\theta} \in \mathbb{R}^{m+1}$

$\hookrightarrow$  get  $\underline{\theta}$  using SVM!

$\Rightarrow$  SVM for  $\underline{\theta}$ :

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

$n=m!$   
 $\underline{\theta}_0 = 1$  unregularized

$\hookrightarrow$  in reality,  $\sum_{j=1}^n \underline{\theta}_j^2 = \underline{\theta}^T \underline{\theta}$  replaced w/  $\underline{\theta}^T \underline{M} \underline{\theta}$

$\downarrow$   
depends on kernel used,  
runs more efficiently

\* recommends using SVM software for minimization optimization

BUT: must choose  $C!$   $\rightarrow P C \propto \downarrow \lambda$

$\hookrightarrow$  large  $C \rightarrow$  lower bias, high variance

small  $C \rightarrow$  higher bias, low variance

& must choose  $\sigma^2!$   $\rightarrow$  large  $\sigma^2$ ,  $f_i$  varies smoothly

(higher bias, lower variance)

$\hookrightarrow$  small  $\sigma^2$ ,  $f_i$  varies sharply

(lower bias, high variance)

### Practical SVM Use:

$\Rightarrow$  liblinear, libsvm packages to solve for  $\underline{\theta}$

$\hookrightarrow$  specify  $C$ , kernel choice

$\hookrightarrow$  e.g. "no kernel" / linear kernel  $\rightarrow y=1$  if  $\underline{\theta}^T \underline{x} \geq 0, \underline{x} \in \mathbb{R}^{n+1}$

(n large, m small)

$\hookrightarrow$  e.g. gaussian kernel  $\rightarrow$  choose  $\sigma^2$   $f_i = \exp\left(-\frac{\|\underline{x} - \underline{l}^{(i)}\|^2}{2\sigma^2}\right)$

where  $\underline{l}^{(i)} = \underline{x}^{(i)}$

$\hookrightarrow$  must write function for software

(feature scale before using Gaussian!)

$\hookrightarrow$  ensures equal weight for each feature

$\Rightarrow$  Not all kernels work  $\rightarrow$  must satisfy Mercer's theorem (don't diverge)

$\hookrightarrow$  polynomial kernel  $K(x, \underline{e}) = (\underline{x}^T \underline{e})^2, (\underline{x}^T \underline{e})^3, (\underline{x}^T \underline{e} + 1)^3$

(worse than Gaussian)  $\hookrightarrow$  choose degree, constant

$\hookrightarrow$  string, chi-square, histogram intersection, etc.

$\Rightarrow$  Multi-class Classification:

$\Rightarrow$  If not built in to software, use one-vs.-all  $\rightarrow$  get  $\underline{\theta}^{(1)}, \underline{\theta}^{(2)}, \dots, \underline{\theta}^{(k)}$ ,  
choose  $\underline{\theta}_R^{(i)T} \geq$  that's largest class

Logistic Regression vs. SVM's:

$\Rightarrow$  If  $n$  too large (relative to  $m$ ), use log. reg. or SVM w/ linear kernel

$\Rightarrow$  If  $n$  small,  $m$  intermediate  $\rightarrow$  SVM w/ Gaussian  
(1-1000) (10-10,000)

$\Rightarrow$  If  $n$  small,  $m$  large  $\rightarrow$  add more features, use log. reg. or SVM w/ lin. kernel  
(1-100) (50,000+)

$\Rightarrow$  Neural Networks work well for most of these<sup>†</sup>, but slower to train

$\hookrightarrow$  SVM's convex, so global mins always found (unlike N.N.'s)