

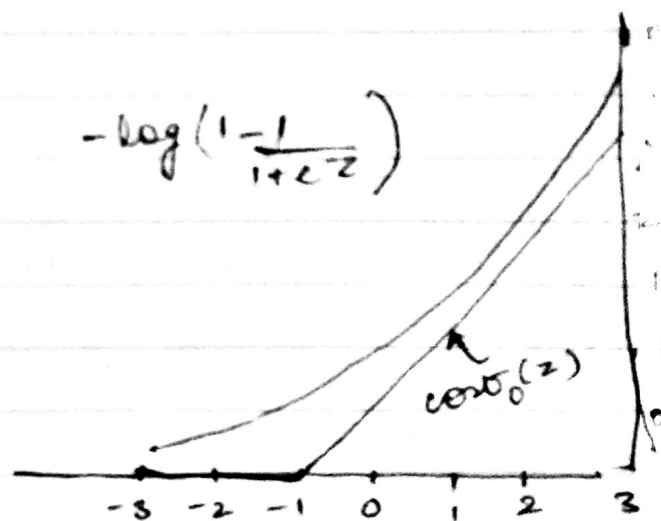
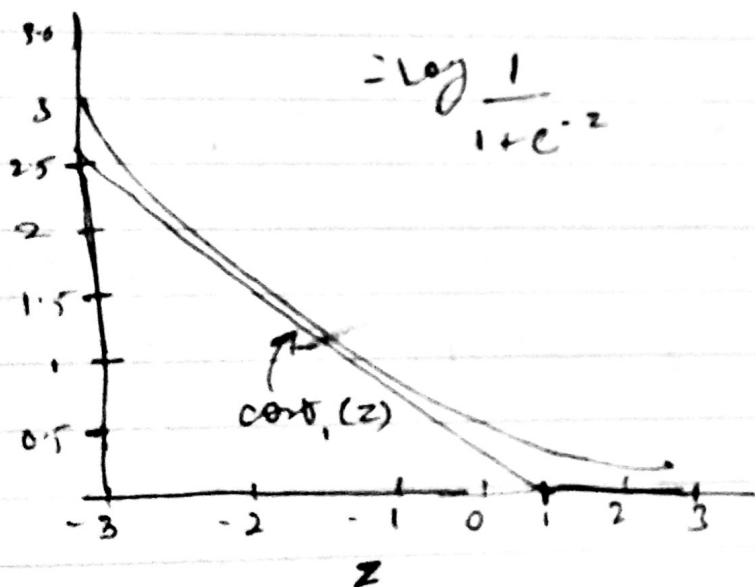
# SUPPORT VECTOR MACHINES

cost function for logistic regression

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

for  $y = 1, \theta^T x \gg 0$

for  $y = 0, \theta^T x \ll 0$



Logistic Regression,

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

$\theta_2 x_2 -$

$$y = mx + c$$
$$\theta_2 x_2 + \theta_1 x_1 + \theta_0 \geq 0$$

$n = 3$

for support vector machines,

$$\min_{\theta} C \sum_{i=1}^m [y^{(i)} \cos \theta(\theta^T x^{(i)}) + (1 - y^{(i)}) \cos \theta(\theta^T x^{(i)})]$$
$$+ \frac{1}{2} \sum_{i=1}^n \theta_i^2$$

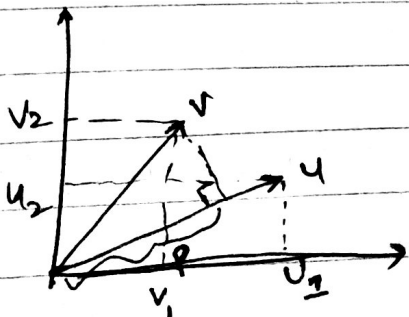
$$h_{\theta}(x) = \begin{cases} 1 & \text{if } \theta^T x \geq 0 \\ 0 & \text{if otherwise.} \end{cases}$$

# SVM is a large margin classifier.

# if  $C$  is very large, the decision boundary will try to incorporate every training example. (gives margin)

Math behind SVM

Vector inner product



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

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

$$x=0$$

$$|0, 2, 0, 2$$

## Kernels

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

for all  $L^{(i)}$ ,  ~~$L^{(i)} = \begin{bmatrix} 0 \\ 2 \\ 0 \\ 2 \end{bmatrix}$~~   $L^{(i)} = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$

$$\therefore f_j = \exp\left(-\sum_{j=1}^n \frac{(x_j - L_j^{(i)})^2}{2\sigma^2}\right)$$

If  $x \approx L^{(i)}$  ← Gaussian kernel

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

$$\approx 1$$

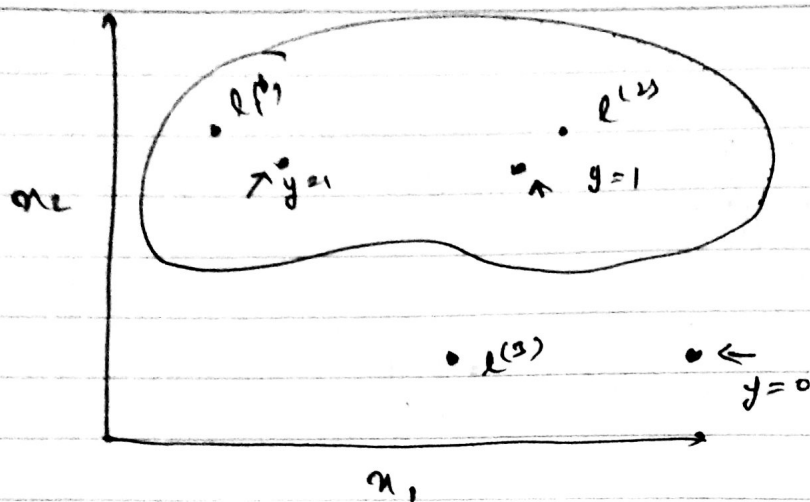
← is a parameter of gaussian kernel = 1 almost all bin

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

$$f_i = \exp\left(-\frac{(\log na)^2}{2\sigma^2}\right) \approx 0$$

$$\approx 0$$

If  $\sigma^2 = 0.5 \rightarrow$  sharp contours  
 $\sigma^2 = 3.0 \rightarrow$  very smooth.



predict 1, when

$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$$

suppose,  $\theta_0 = -0.5$

$$\theta_1 = 1$$

$$\theta_2 = 1$$

$$\theta_3 = 0$$

Examples  $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})$   
 choose  $x^{(1)} = x^{(1)}, x^{(2)} = x^{(2)}, \dots, x^{(m)} = x^{(m)}$

$$f_1 = \sin(x, x^{(1)})$$

$$f_2 = \sin(x, x^{(2)})$$

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

for example  $(x^{(1)}, y^{(1)})$

$$x^{(1)} = \begin{bmatrix} f_1^{(1)} = \sin(x^{(1)}, x^{(1)}) \\ f_2^{(1)} = \sin(x^{(1)}, x^{(2)}) \\ f_3^{(1)} = \sin(x^{(1)}, x^{(3)}) \\ \vdots \\ f_m^{(1)} = \sin(x^{(1)}, x^{(m)}) \end{bmatrix} = 1$$

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

$$f^{(i)} = \begin{bmatrix} f_1^{(i)} \\ f_2^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix}$$

predict  $y=1$  if  $\theta^T f \geq 0$   
 $\nearrow \theta_0 f_0 + \theta_1 f_1 + \theta_2 f_2 \dots$

$$m = m_1$$

Training

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

$\downarrow$   
this could be written  
as  $\theta^T \theta$

$C (= \frac{1}{\lambda})$  ; large  $C$  : High variance  
~~large~~ small  $C$  : Higher bias.

$\sigma^2$  ; large  $\sigma^2$  - vary not smoothly  
Higher bias.

; small  $\sigma^2$  - vary sharply  
Lower bias.

★ Use SVM package (eg., liblinear, libsvm, ...) to solve for  $\theta$ .

Need to specify

→ choice of  $C$

→ choice of kernel.

[if  $n$  is large &  $m$  is small]

Eg: No kernel - predict  $y = 1$  if  $\theta^T x \geq 0$   
(in others, it is  $\theta^T x \geq 1$ )

Gaussian kernel

[ $n$  is small /  $m$  is large]

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

Need to choose  $\sigma^2$ .

# perform feature scaling, if there are too many differences in the features, suppose size of rooms (1000), no of bedrooms (range 1-5).

★ They need to satisfy "Mercer's theorem", to make sure that  $\theta$  is calculated correctly.

### Multi-class classification

# Many SVM has already built-in multi-class classification  
# or use one-vs-all method (like logistic regression).

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

## Logistic regression v/s SVM

( $n = 10,000$ )

( $m = 10, \dots, 1000$ )

# If  $n$  is large (relative to  $m$ ),  
use logistic regression or SVM without a kernel.  
( $n = 1-1000$ ) ( $m = 10, 10,000$ )

# If  $n$  is small,  $m$  is intermediate  
→ use SVM with gaussian kernel.  
( $n = 1, 1000$ ) ( $m = 50,000$ )

# If  $n$  is small,  $m$  is large  
(Create/add more features, then use logistic regression or SVM without a kernel)

\* Neural network will work well, however it is slower to train.

\* In SVM, it will find global minima, not sure about neural network.

---

## UNSUPERVISED LEARNING

# Learning from unlabeled data.

# There's no  $y$ .

Clustering can be used for market analysis.