

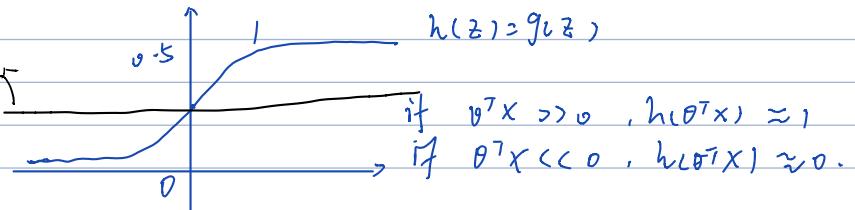
SVM

logistic regression:

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

$$z = \theta^T x$$

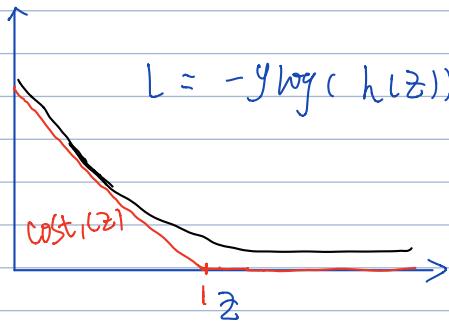
decision boundary
 $g(z) = 0.5$



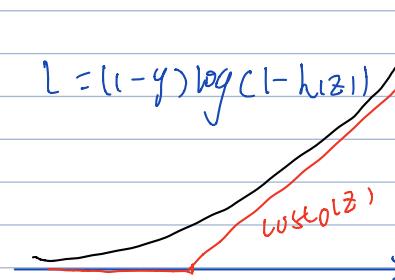
Cost Function:

$$L_i = -y^{(i)} \log(h_\theta(x^{(i)})) - (1-y^{(i)}) \log(1-h_\theta(x^{(i)})) \quad h_\theta(x) = \frac{1}{1+e^{-\theta^T x}}$$

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



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



$$L = \frac{1}{m} \sum_{i=1}^m L_i + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

Logistic Regression Function:

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

Support vector machine:

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

↓

$$\min C \left[\sum y^{(i)} \text{cost}_1(\theta^T x) + (1-y^{(i)}) \text{cost}_0(\theta^T x) \right] + \frac{1}{2} \sum \theta_j^2$$

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

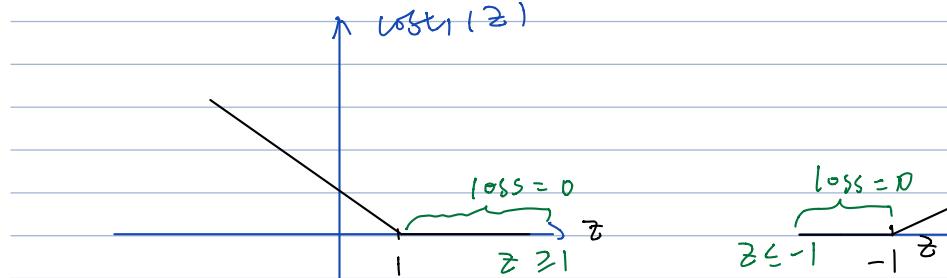
\rightarrow the parameter optimized by SVM.

SVM hypothesis:

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

large margin intuition:

$$\text{loss}_1(z)$$

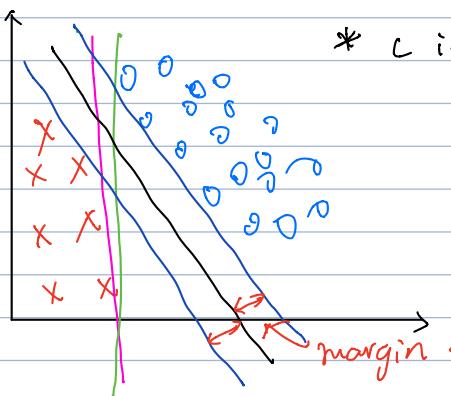


If $y=1$, we want $\theta^T x \geq 1$ (not just $\theta^T x \geq 0$)

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

↓ soft margin / safety factor

linear separable



* C is very large.

当 C 不是很大的时候，SVM 可以忽略一些异常点的影响，得到更好的效果。

甚至当 data 不是线性可分的时候，SVM 也可以给出相容的决策。

what we minimize in the SVM is

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

SVM Decision Boundary:

$$\min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2 = \frac{1}{2} (\theta_1^2 + \theta_2^2) = \frac{1}{2} \underbrace{\|\theta_1^2 + \theta_2^2\|_1^2}_{\|\theta\|} = \frac{1}{2} \|\theta\|^2$$

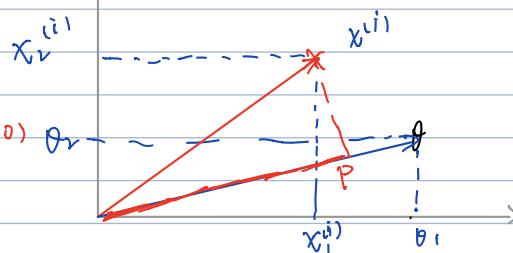
$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} \quad \theta_0 = a$$

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

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

$$\text{if } n=2, \theta_0 = 0$$

the boundary has
to pass origin (0,0)

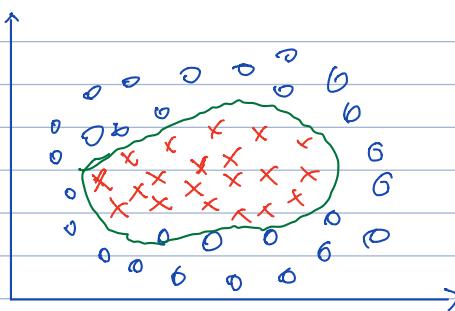


$$\begin{aligned} \theta^T x^{(i)} &= p^{(i)} \cdot \|\theta\| \\ &= \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)} \end{aligned}$$

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

使 $p^{(i)}$ 最大，即样本到 Decision boundary 的距离最大，让样本 vector 投影在 parameters vector 上的值最大。即能找到 θ minimize $\|\theta\|$ ，使 $p^{(i)} \|\theta\|$ 满足约束。

kernels:



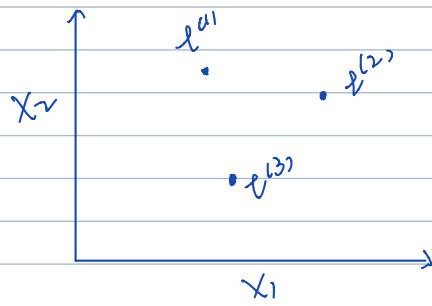
Given the dataset, our model might be
 $\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 + \dots$

If we use a set of new features to replace those: $f_1 = x_1, f_2 = x_2, f_3 = x_1 x_2$. then

$$h(x) = \theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 + \dots$$

How to find better features:

use the similarity between x and a bunch of landmarks $t^{(1)}, t^{(2)}, t^{(3)}$.



Given x :

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

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

$$\|x - t^{(1)}\|^2 = \sum_{j=1}^n (x_j - t_j^{(1)})^2 \rightarrow x \text{ 中所有特征与 } t^{(1)} \text{ 的距离的和.}$$

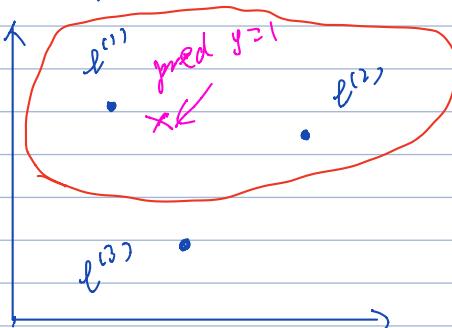
核函数, 此处的为高斯核函数 (Gaussian kernel)

若 x is close to landmarks, $f_1 \approx e^{-0} = 1$

若 x is far from $t^{(1)}$: $f_1 \approx e^{-\infty} = 0$.

注: 只有当 x 与 $t^{(1)}$ 重合时 f_1 才具有最大值, 因改变速率参数 σ^2 的影响.

For example:



predict "1" when

$$h(x) = \theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0.$$

$$\theta_0 = -0.5 \quad \theta_1 = \theta_2 = 1 \quad \theta_3 = 0.$$

$$h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 \geq 0. \Rightarrow y=1$$

在训练中我们采取的特征不是训练实例本身, 而是通过 kernel 计算出来的.

How to choose / set landmarks:

training examples: $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}) \dots, (x^{(m)}, y^{(m)})$

训练集中有 m 个实例, 则我们选择 m 个地标, 并且令 $t^{(1)} = x^{(1)}, t^{(2)} = x^{(2)}, \dots, t^{(m)} = x^{(m)}$
 我们得到的新特征是建立在原有特征与训练集中所有其他特征之间的距离
 基础上得来的, 即:

Given example x .

Given $(x^{(i)}, y^{(i)})$

$$f_1 = \text{similarity}(x, t^{(1)})$$

$$f_2 = \text{similarity}(x, t^{(2)})$$

:

$$f^{(i)} = \begin{cases} f_0^{(i)} = 1 \\ f_1^{(i)} = \text{sim}(x^{(i)}, t^{(1)}) \\ f_2^{(i)} = \text{sim}(x^{(i)}, t^{(2)}) \\ \vdots \\ f_i^{(i)} = \text{sim}(x^{(i)}, t^{(i)}) = 1 \\ \vdots \\ f_m^{(i)} = \text{sim}(x^{(i)}, t^{(m)}) \end{cases}$$

Hypothesis: Given x , compute features $f(-k^{m+1})$

Predict " $y=1$ " if $\theta^T f \geq 0$ $\theta^T f = \theta_0 + \theta_1 f_1 + \theta_2 f_2 + \dots$

Loss Function: $Z = \theta^T f$.

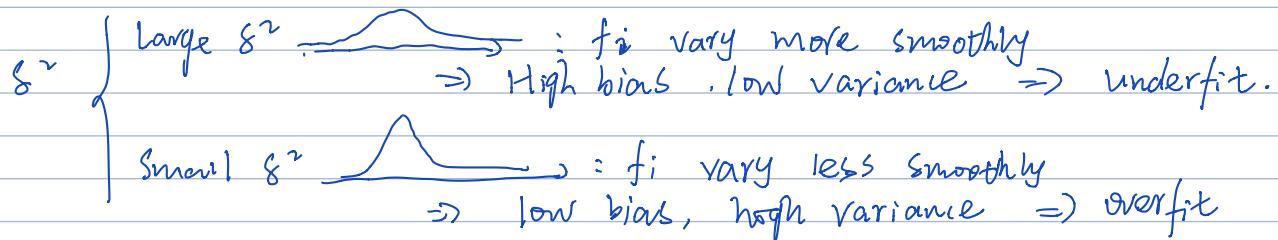
$$\min C \left[\sum y^{(i)} \text{cost}_1(\theta^T f) + (1 - y^{(i)}) \text{cost}_0(\theta^T f) \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

在实际应用中，我们需对正则项进行微调：

$$\sum_{j=1}^n \theta_j^2 = \theta^T \theta \Rightarrow \theta^T \theta : m \text{ 是根据 kernel 算出的一个矩阵.}$$

High bias/variance in SVM:

$$C = \frac{1}{\lambda} \quad \left\{ \begin{array}{l} C \text{ is large} \rightarrow \lambda \text{ is small} \Rightarrow \text{lower bias, high variance} \Rightarrow \text{overfitting.} \\ C \text{ is small} \rightarrow \lambda \text{ is large} \Rightarrow \text{high bias, low variance} \Rightarrow \text{underfit.} \end{array} \right.$$



SVM 的代价函数是凸函数，不存在局部最小值，所以其优化是一种凸优化问题。

A linearly separable dataset can usually be separated by many different lines. Varying the parameter C will cause the SVM's decision boundary to vary among these possibilities. for example, very large C might learn larger values of θ in order to increase the margin on certain examples.