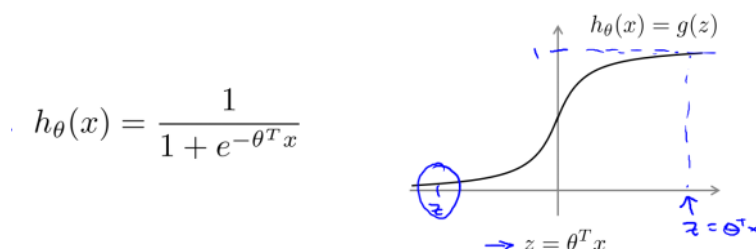


# 12 支持向量机

2022年11月22日 18:14

## 12.1 优化目标

依然从逻辑回归出发，

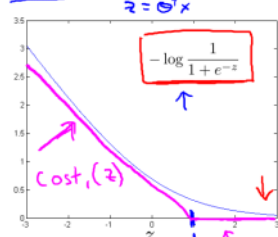


对于一个逻辑回归的模型，要想让假设函数的输出值为1，那么需要 $z$ 远大于0；若想让假设函数的输出为0，那么 $z$ 需远小于0。

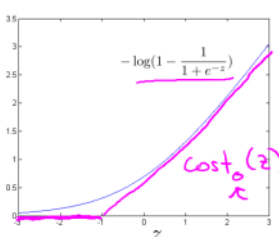
Cost of example:  $-(y \log h_\theta(x) + (1 - y) \log(1 - h_\theta(x)))$  ←

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

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



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



对于一个样本的代价函数，它由两项构成。如果想要 $y$ 为1，那么只有第一项起作用，从而绘制出此时的代价函数曲线；而如果想要 $y$ 为0，那么只有第二项起作用，同样地绘制出曲线。分别在 $z=1$ 和 $-1$ 处取一点，作如图所示的分段折线，可以看出它是能够逼近起始的曲线的。有了这两个函数就可以构建支持向量机 (Support vector machine)。

Logistic regression:

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

Support vector machine:

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

可以看到，和逻辑回归的式子不同的是，将两个代价项替换成了刚刚的所设的函数，同时去掉了相乘的常数 $1/m$ ，显然这个相乘的常数对求取最小值对应的变量是没有影响的；此外，将正则化因子转移到了前面的代价项上，因为本质上正则化因子的设置就是用来平衡第一项和第二项的权重的，所以把它移到前面依然能够行使这样的功能。

最后，和逻辑回归不同的是，支持向量机输出的并不是一个概率，它作出的直接是一个预测，那就是 $y$ 到底是1还是0。

## 12.2 直观上对大间隔的理解

有时候人们会把支持向量机叫做大间隔分类器。

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

$\text{cost}_1(z)$   $z \geq 1$   
 $\text{cost}_0(z)$   $z \leq -1$

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

根据上一节所讲的支持向量机的性质，如果我们想让代价值尽可能地小，那么当  $y=1$  时，只需  $z$  大于等于 1，其对应的代价值就为 0；而当  $y=0$  时，只要  $z$  小于等于 -1，其对应的代价值也就为 0。

我们已知在逻辑回归中，可以通过判断  $\theta^T X$  是否大于零来对其进行分类；而支持向量机则对其提出了更高的要求。

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

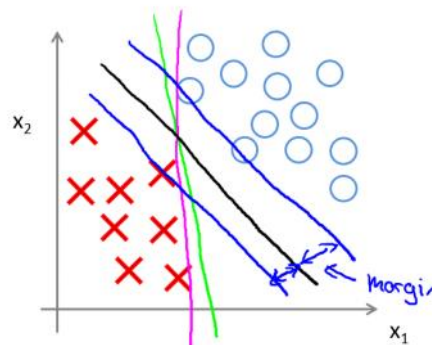
Whenever  $y^{(i)} = 1$ :  $\theta^T x^{(i)} \geq 1$

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

$\min_{\theta} C \times 0 + \frac{1}{2} \sum_{j=1}^n \theta_j^2$   
s.t.  $\theta^T x^{(i)} \geq 1$  if  $y^{(i)} = 1$   
 $\theta^T x^{(i)} \leq -1$  if  $y^{(i)} = 0$

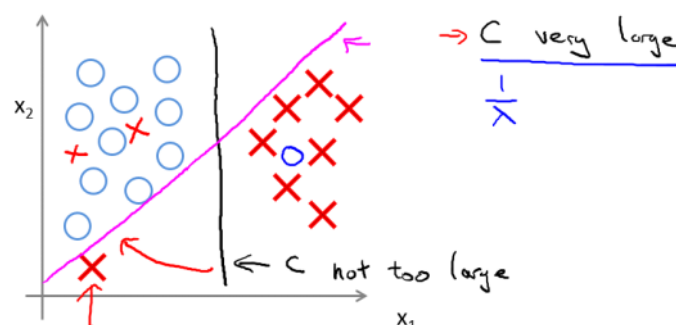
当我们把  $C$  设置得特别大的时候，我们希望找到一个合适的  $\theta$  值使代价项为 0。

### SVM Decision Boundary: Linearly separable case



#### Large margin classifier

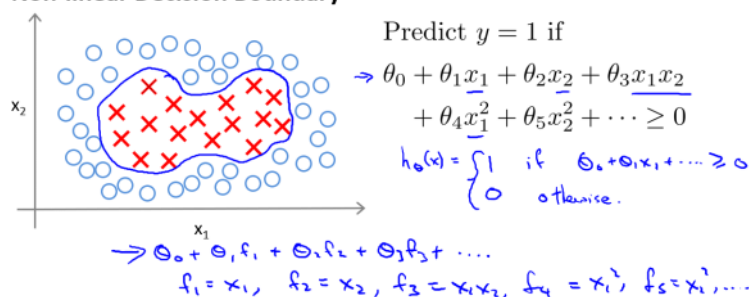
对于决策边界的选择可能有很多种，比如上图粉色和绿色的极端情况，而支持向量机则倾向于选择黑色的那条更稳健的边界。它在数学上的具体表现是它到样本的最小距离要比其他边界更大一些。这个距离叫做支持向量机的间距 (margin)，这使得它在分离数据时具有鲁棒性，因为它会用尽量大的间距来分离。这也是为什么支持向量机被称为最大间距分类器。





对于一个复杂的分类问题，

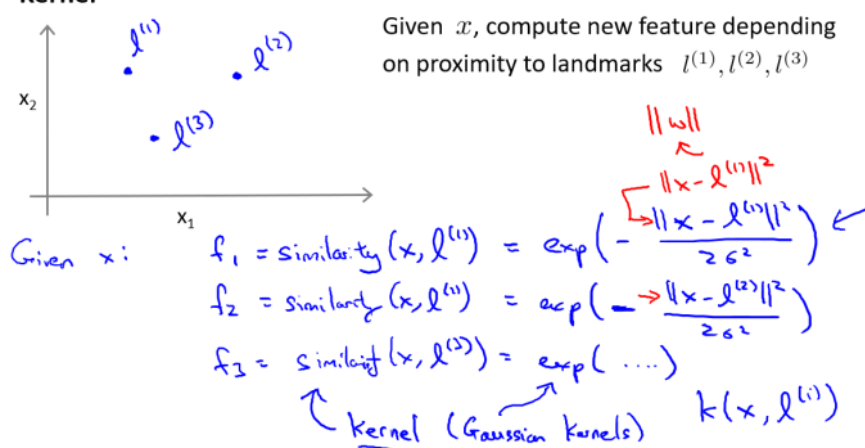
### Non-linear Decision Boundary



Is there a different / better choice of the features  $f_1, f_2, f_3, \dots$ ?

可以通过构造高阶多项式的假设函数来对其进行边界的拟合，此时定义如上蓝色字体所写的函数取代高阶多项式的假设函数。由于我们不知道这些高阶项是否有用以及是否有更好的特征可以嵌入到函数中，所以如何选择是一个问题。

### Kernel



定义 $f$ 为一种 $x$ 和标记(landmark) $l$ 的相似度量。这个相似度量函数就叫做核函数(Kernel)。这里特指高斯核函数。

### Kernels and Similarity

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

If  $x \approx l^{(1)}$ :

$$f_1 \approx \exp\left(-\frac{0^2}{2\sigma^2}\right) \approx 1$$

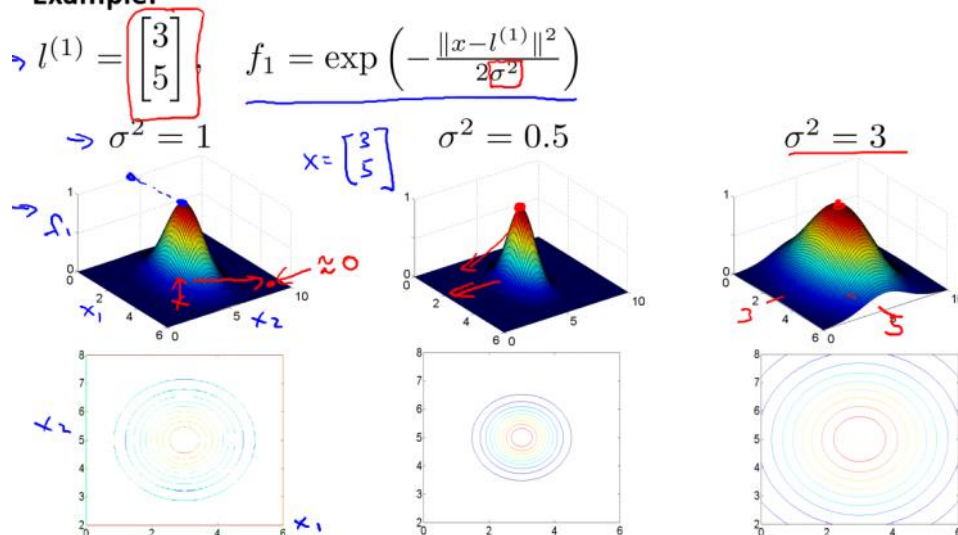
$$\begin{aligned} l^{(1)} &\rightarrow f_1 \\ l^{(2)} &\rightarrow f_2 \\ l^{(3)} &\rightarrow f_3 \end{aligned}$$

If  $x$  is far from  $l^{(1)}$ :

$$f_1 = \exp\left(-\frac{(\text{large number})^2}{2\sigma^2}\right) \approx 0$$

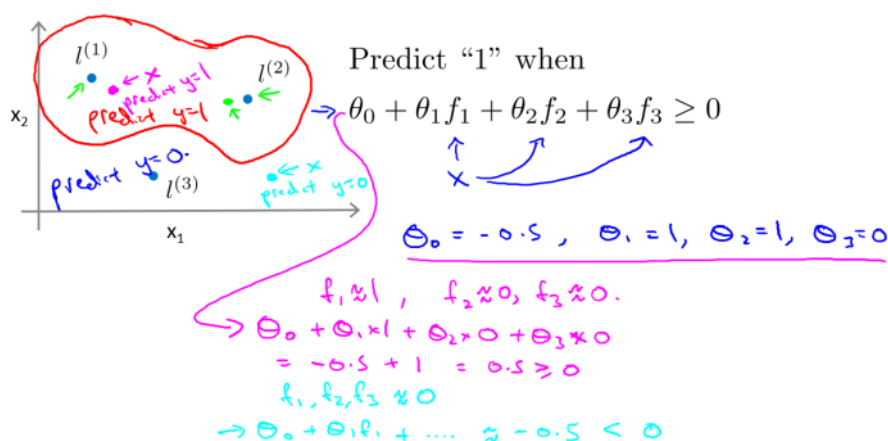
当 $x$ 和 $l$ 很接近时，核函数的值趋于1；当 $x$ 和 $l$ 很远时，核函数的值趋于0。核函数的直观表示如下图：

Example:



可以看出，当 $x$ 与 $l$ 相等时达到最高点；而随着核函数的参数的增加，最高点向周围降低的速度变得越来越缓慢。

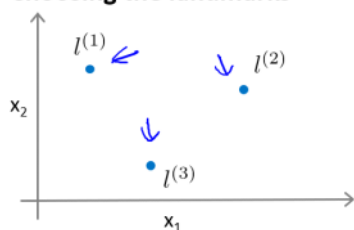
举一个实际预测的例子：



## 12.5核函数2

怎么得到标记点？

Choosing the landmarks

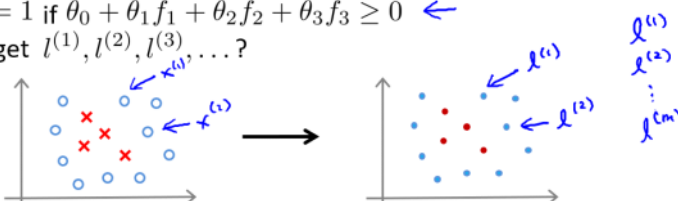


Given  $x$ :

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

Predict  $y = 1$  if  $\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$

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



一般直接将样本作为标记点：

### 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  $x$ :

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

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

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

$$\begin{aligned} f_1^{(i)} &= \sin(x^{(i)}, l^{(1)}) \\ f_2^{(i)} &= \sin(x^{(i)}, l^{(2)}) \\ &\vdots \\ f_m^{(i)} &= \sin(x^{(i)}, l^{(m)}) \end{aligned}$$

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

然后，用  $f$  代替  $x$  进行支持向量机的预测：

Hypothesis: Given  $x$ , compute features  $f \in \mathbb{R}^{m+1}$

→ Predict "y=1" if  $\theta^T f \geq 0$

Training:

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

$$\begin{aligned} - \sum_j \theta_j^2 &= \theta^T \theta \leftarrow \theta = \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_m \end{bmatrix} \quad (\text{ignore } \theta_0) \\ &\quad \theta^T M \theta \quad \|\theta\|^2 \quad m = 10,000 \end{aligned}$$

有研究表明，支持向量机和核函数可以相得益彰以及基于其二者的更高级的优化算法，而逻辑回归和核函数则显得十分缓慢。

老师不推荐自己写最小化代价函数的代码，而应该使用成熟的软件包。

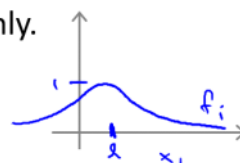
另一个问题是怎么使用支持向量机中的参数。其依据依然是偏差和方差的折中。

### SVM parameters:

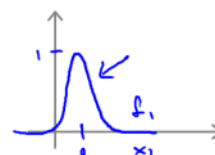
$C (= \frac{1}{\lambda})$ . → Large C: Lower bias, high variance. (small  $\lambda$ )  
→ Small C: Higher bias, low variance. (large  $\lambda$ )

$\sigma^2$  Large  $\sigma^2$ : Features  $f_i$  vary more smoothly.  
→ Higher bias, lower variance.

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



Small  $\sigma^2$ : Features  $f_i$  vary less smoothly.  
Lower bias, higher variance.



同我们在逻辑回归中分析的一样，在逻辑回归中，当  $\lambda$  太小时会得到低偏差和高方差；当  $\lambda$  太大会得到高偏差和低方差。由于  $C$  和  $\lambda$  的反比趋势，对于  $C$  所造成的影响就恰好相反。即大的  $C$  值会倾向于过拟合，小的  $C$  值会倾向于欠拟合。

而对于高斯核函数中  $\sigma^2$  的选取，前面已经分析过，当它更大时核函数会更平滑和缓慢，所得的模型倾向于高偏差和低方差；而当它较小时核函数的变化会剧烈一些、斜率较大，所得模型倾向于低偏差和高方差。

## 12.6 支持向量机的使用

一般建议用已有的软件库实现 SVM 算法，如 `liblinear`, `libsvm`。

特别的，我们需要自己选择  $C$  值；其次，对于有内参的核函数我们需要选择参数值，而



像线性核函数(linear kernel)这种就是无内核参数函数，我们不需要自己选择参数。比如，当特征数量很多而样本数量很少的时候，为了防止过拟合，线性的核函数就是一个好的选择。

核函数的实现：

**Kernel (similarity) functions:**

```

function f = kernel(x1, x2)
    f = exp( - (||x1 - x2||^2) / (2σ^2) )
return
    
```

Handwritten notes:  $x^{(j)} = x^{(j)}$ ,  $x \rightarrow \begin{matrix} f_1 \\ f_2 \\ \vdots \\ f_m \end{matrix}$

Note: Do perform feature scaling before using the Gaussian kernel.

Handwritten derivation for  $\|x - l\|^2$ :

$$\|x - l\|^2 = (x_1 - l_1)^2 + (x_2 - l_2)^2 + \dots + (x_n - l_n)^2$$

Example:  $\Rightarrow 1000 \text{ feet}^2$ ,  $1-5 \text{ bedrooms}$

其中应当注意归一化。

另外，我们选择的核函数应当满足默塞尔定理：

#### 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).

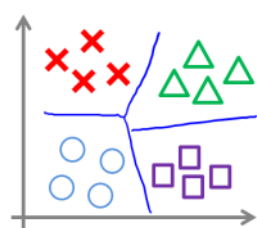
Many off-the-shelf kernels available:

- Polynomial kernel:  $k(x, l) = (x^T l + \text{constant})^{\text{degree}}$   
Handwritten examples:  $(x^T l)^2$ ,  $(x^T l + 1)^2$ ,  $(x^T l + 5)^2$
- More esoteric: String kernel, chi-square kernel, histogram intersection kernel, ...  
Handwritten:  $\text{sim}(x, l)$

这个定理确保所有的SVM算法大类的优化方法并从而迅速得到参数  $\theta$ 。多项式核函数等是其它符合默塞尔定理的核函数。

下面再来讨论多分类的判定边界：

#### Multi-class classification



$$y \in \{1, 2, 3, \dots, K\}$$

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  $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(K)}$

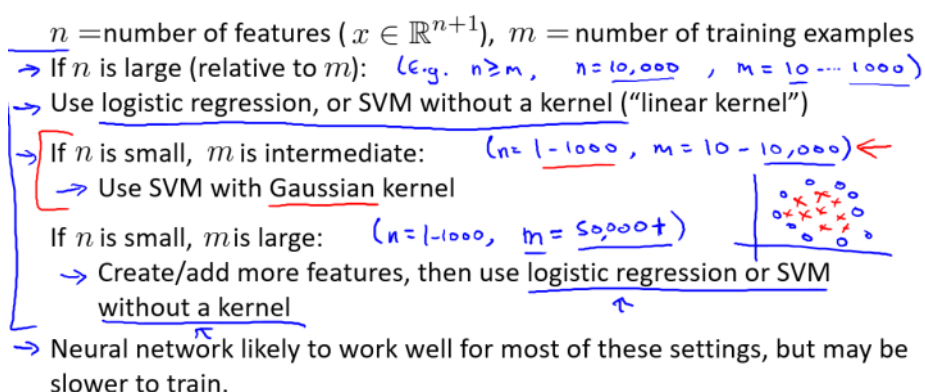
Pick class  $i$  with largest  $(\theta^{(i)})^T x$

Handwritten:  $y=1, y=2, \dots, \theta=K$

一些SVM算法包已经能够很好地做到多分类；而我们也可以采用之前提到的一对多的分类方法，通过k个SVM来进行分类。

最后来比较一下逻辑回归和支持向量机。

### Logistic regression vs. SVMs



- 相较于样本数量，特征数目很多的情况下，我们通常选择逻辑回归或者是使用了线性核的支持向量机。
- 当特征数目很小，而样本数量适中时，用如高斯核函数的有核函数的支持向量机工作得很好。
- 当特征数目很小，但样本数量很大时，由于数据数量太大，用有核函数的支持向量机可能会工作得比较慢。在这种情况下可以尝试手动创建或添加更多特征变量，然后用逻辑回归或者不带核函数的支持向量机。
- 神经网络可能对大多数情况都很有效，但是训练起来会比较慢。