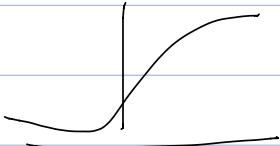


Support Vector Machines

Optimization objective

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

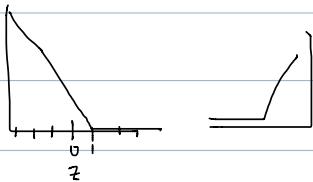


$$z = \theta^T x$$

If $y=1$, we want $h_\theta(x) \approx 1$, $\theta^T x \gg 0$

If $y=0$, we want $h_\theta(x) \approx 0$, $\theta^T x \ll 0$

If $y=1$ ($w_0 + \theta^T y \gg 0$):



SVM: (cost function)

$$\min_{\theta} \frac{1}{m} \sum_{i=1}^m y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) + \frac{\lambda}{2} \sum_{j=0}^n \theta_j^2$$

$$\begin{cases} A + \lambda B & (\text{if } y=1) \\ A + B & (\text{otherwise}) \end{cases}$$

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

Hypothesis:

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

Margin Intuition

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

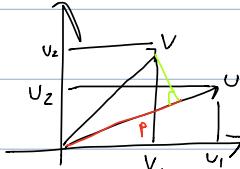
If $y=0$, we want $\theta^T x \leq -1$ ($\text{not just } < 0$)

Mathematics for LM

$$u = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \quad v = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} \quad u^T v = ?$$

$\|u\| = \text{length of vector } u$

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



$p = \text{length of projection of } v \text{ onto } u$

$$\uparrow u^T v = p \cdot \|u\|$$

$$\text{Sign}(t_{+}) = u_1 v_1 + u_2 v_2 \quad p \in \mathbb{R}$$

SVM Decision Boundary

$$w = (\sqrt{w})^2$$

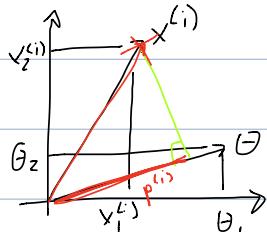
$$\min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2 = \frac{1}{2} (\underbrace{\theta_1^2 + \theta_2^2}_{= \|\theta\|^2}) = \frac{1}{2} \|\theta\|^2$$

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

$$\theta^T x^{(i)} \leq -1 \quad \text{if } y^{(i)} = 0 \quad \boxed{\theta_0 = 0}$$

$$\text{Simplification: } \boxed{\theta_0 = 0} \text{ origin } n=2$$

$$\theta^T x^{(i)} = ?$$



$$\theta^T x^{(i)} = p^{(i)} \cdot \|\theta\| \Leftarrow$$

$$= \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)} \Leftarrow$$

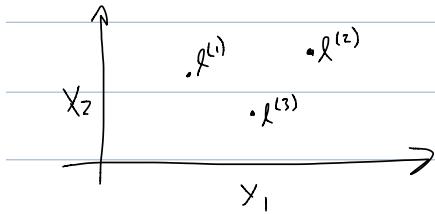
$$\text{s.t. } p^{(i)} \cdot \|\theta\| \geq 1 \quad \text{if } y^{(i)} = 1$$

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

Non-linear Decision Boundary

KERNELS

$$\Theta_0 + \Theta_1 x_1 + \Theta_2 x_2 + \Theta_3 x_1 x_2 + \Theta_4 x_1^2 + \Theta_5 x_2^2 + \dots \geq 0$$



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

$f_2 = \text{similarity}(x, l^{(2)}) = \exp \dots$

$\underbrace{\dots}_{\text{Kernel (gaussian kernels)}} K(x, l^{(i)})$

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

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

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

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

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

Kernels II

Where to get landmarks?

$$\text{Training examples} = \text{landmarks}$$

Given $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(n)}, y^{(n)})$

Choose $l^{(1)} = x^{(1)}, l^{(2)} = x^{(2)}, \dots, l^{(n)} = x^{(n)}$

Given example x :

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

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

$$f_i = \text{sim}(x, f^{(i)})$$

↳ find

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

$$x^{(i)} \rightarrow f_i^{(i)} = \text{sim}(x^{(i)}, f^{(i)})$$

$$f_m^{(i)} = \text{sim}(x^{(i)}, f^{(m)})$$

$$\begin{aligned} & x^{(i)} \in \mathbb{R}^{n+1} \quad f_0^{(i)} = 1 \quad (\text{or } \mathbb{R}^n) \\ & f = \begin{bmatrix} f_0^{(i)} \\ f_1^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix} \end{aligned}$$

Hypothesis: Given x , compute features $f \in \mathbb{R}^{m+1} \quad (\Rightarrow \mathbb{R}^{n+1})$

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

\nearrow meet the line
how?

Training:

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

$\uparrow \theta_0 \leftarrow \text{bias term}$

$$\begin{cases} -\sum_j \theta_j^2 = \theta^T \theta \leftarrow \theta = \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_n \end{bmatrix} \quad [n \times \theta_0] \\ - \frac{\partial J(\theta)}{\partial \theta} \quad \text{Computational efficiency} \end{cases}$$

$m = 10,000$

SVM parameters:

C ($C = \frac{1}{\lambda}$) Large C : Lower bias, high variance (small λ)

Small C : Higher bias, low variance (large λ)

σ^2 Large σ^2 : Features f_i vary more smoothly

Higher bias, lower variance

Small σ^2 : Features f_i vary less smoothly.

Lower bias, higher variance.



Using an SVM

Use SVM software package (liblinear, libsvm, ...) to solve for parameters θ .

Need to specify:

Choice of parameter

Choice of kernel (similarity function):

E.g. No kernel ("linear kernel") if $\log_2 \# \text{features}, \log_2 \# \text{examples}$

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

Gaussian kernel:

if small # features, large # examples

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

Need to choose σ^2 \leftarrow Use feature scaling!

Other kernels:

- Polynomial kernel: $k(x, \ell) = (x^T \ell)^2, (x^T \ell)^3, (x^T \ell + 1)^3, \dots, (x^T \ell + \text{const})^{\text{degree}}$
- String kernel, chi-square kernel, histogram intersection kernel

Multi-class Classification

One vs All

Logistic Regression vs SVMs

$n = \text{number of features} (X \in \mathbb{R}^{n \times 1}), m = \text{number of training examples}$

If n is large (relative to m): $n \approx m$ $n = 10,000$ $m = 10 \dots 10,000$

Use Logistic regression, or SVM without a Kernel (γ_{large})

If n is small, m is intermediate ($n = 1 \dots 1000, m = 10 \dots 10,000$)

Use SVM with Gaussian kernels

If n is small, mislaze: ($n=1-1000, m=50,000$)

(retrn/kd) more features than use logistic regression or SVM without a Kernel

Neural network likely to work well for most of you, may be slow to train