

Logistic Regression

Sid

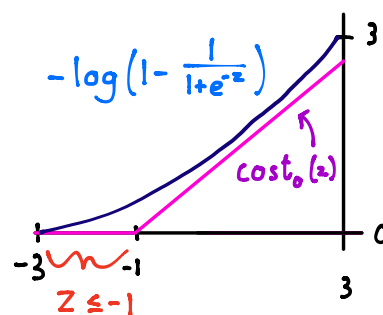
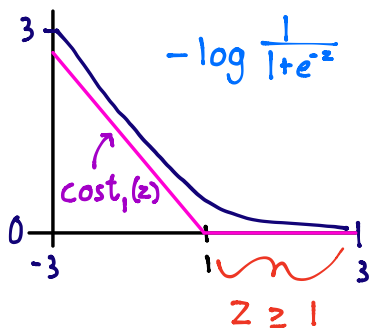
$$h_{\theta}(x) = \frac{1}{1+e^{-\theta^T x}}$$

$$y=1 : h_{\theta}(x) \approx 1, \theta^T x \gg 0$$

$$y=0 : h_{\theta}(x) \approx 0, \theta^T x \ll 0$$

$$\begin{aligned} \text{Cost}(x, y) &: -(y \log h_{\theta}(x) + (1-y) \log (1-h_{\theta}(x))) \\ &= -y \log \frac{1}{1+e^{-\theta^T x}} - (1-y) \log \left(1 - \frac{1}{1+e^{-\theta^T x}}\right) \end{aligned}$$

$$\text{if } y=1 \text{ (want } \theta_{(z)}^T x \geq 1) \quad \text{if } y=0 \text{ (want } \theta^T x \leq -1)$$



Logistic :

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

(convex)

SVM :

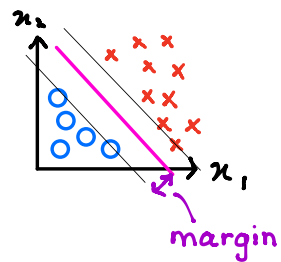
$$\min_{\theta} C \sum_{i=1}^m y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

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

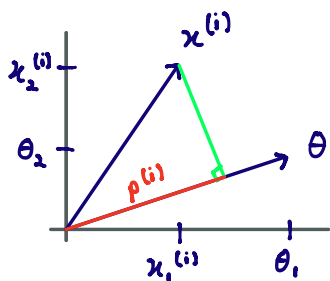
SVM: Large Margin Classifier

$(\theta_0 = 0, n=2)$

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



$$\rho^{(i)} \|\theta\| \leq \begin{cases} \theta^T x^{(i)} \geq 1, & y^{(i)} = 1 \\ \theta^T x^{(i)} \leq -1, & y^{(i)} = 0 \end{cases}$$



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

Kernels : used to make non-linear classifiers using SVM

Given $x \rightarrow$ compute new features (f) depending on proximity to landmarks (l)

$$f^{(i)} = \overset{(k)}{\text{Similarity}}(x, l^{(i)}) = \exp\left(\frac{-\|x - l^{(i)}\|^2}{2\sigma^2}\right)$$

$$\text{kernel (Gaussian kernel)} = \exp\left(-\frac{\sum_{j=1}^n (x_j - l_j^{(i)})^2}{2\sigma^2}\right)$$

$$\bullet x \approx l^{(i)} \rightarrow f_i = \exp\left(-\frac{\approx 0^2}{2\sigma^2}\right) \approx 1$$

$$\bullet x \text{ far from } l^{(i)} \rightarrow f_i = \exp\left(-\frac{\text{large \#}^2}{2\sigma^2}\right) \approx 0$$

$$h_\theta(x) = \theta_1 f_1 + \theta_2 f_2 + \dots$$

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

\downarrow
Choose $l^{(1)} = x^{(1)}, \dots, l^{(m)} = x^{(m)}$

$$x^{(i)} \rightarrow \begin{bmatrix} f_1^{(i)} = \text{sim}(x^{(i)}, l^{(1)}) \\ f_2^{(i)} = \text{sim}(x^{(i)}, l^{(2)}) \\ \vdots \\ f_m^{(i)} = \text{sim}(x^{(i)}, l^{(m)}) \end{bmatrix}$$

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

$$\begin{cases} 1: \theta^T f \geq 0 \\ 0: \theta^T f < 0 \end{cases}$$

$$= \theta^T \theta$$

$$\theta^T M \theta$$

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

$$C (\approx \frac{1}{\lambda})$$

- large C : Low bias, high variance
- small C : High bias, low variance

$$\sigma^2$$

- Large : Features f_i vary smoothly (\uparrow Bias, \downarrow Variance)
- Small : f_i vary less smoothly (\downarrow Bias, \uparrow Variance)

Using a SVM

- use SVM library (liblinear, libsvm, ...)

• Choose C

• Choose kernel (sim. function)

- No kernel ("linear") \rightarrow standard linear classifier
(n is large, m is small) or logistic regression

- Gaussian kernel (n is small, m is large)

• choose σ^2

if too large, add features

$$\text{function } f = \text{kernel}(x_1, x_2)$$

$$f = \exp\left(-\frac{\|x_1 - x_2\|^2}{2\sigma^2}\right)$$

$x^{(i)}$ $x^{(j)} = x^{(j)}$

* feature scaling before using Gaussian kernel

end

Multi-class classification using SVM

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

\rightarrow Train k SVMs to distinguish $y=i \rightarrow \theta^{(1)}, \theta^{(2)}, \dots, \theta^{(k)}$

\rightarrow Pick class i w/ largest $(\theta^{(i)})^T x$