- Alternative view of logistic regression

We'll transform this        in this



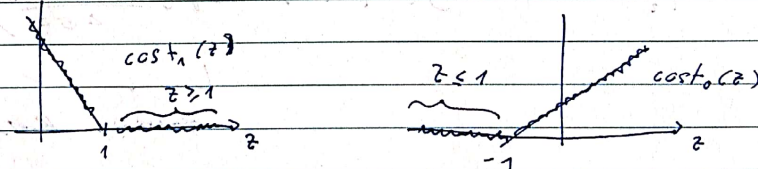$cost_1(z)$        $cost_0(z)$

Logistic reg:

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

SPV

$$\frac{\partial}{\partial \theta} \left[ C \sum_{i=1}^{m} \left[ y^i \, cost_1(\theta^T x^i) + (1 - y^i) \, cost_0(\theta^T x^i) \right] + \frac{1}{2} \sum_{i=1}^{n} \theta_j^2 \right]$$



$cost_1(z)$        $z \leq 1$        $cost_0(z)$

$z \geq 1$

if $y = 1 \Rightarrow \theta^T x \geq 1$  it's what we want, not
if $y = 0 \Rightarrow \theta^T x \leq -1$  just $\theta^T x \geq$ or $\leq 0$, then
                                              we create a proper margin
                                              (decision boundary)

# KERNELS

- Non linear decision boundary

Predict $y=1 \iff \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \ldots \geq 0$

$\Rightarrow h_\theta(x) = \begin{cases} 1 \iff \theta_0 + \ldots \geq 0 \\ 0 \quad else \end{cases}$

Our features are $x_1, x_2, x_3$, but now we'll use others $\{ f_1 = x_1, f_2 = x_2, f_3 = x_1 x_2 \ldots$ for example $\}$

$\rightarrow$ Given $x$, we'll compute new features depending on proximity to landmarks $\ell^1, \ell^2, \ell^3$

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

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

$\underbrace{\qquad\qquad}$ gaussian kernels

- Kernels and similarity
  - if $x \approx \ell^i \Rightarrow f_1 \approx 1$
  - if $x$ far from $\ell^i \Rightarrow f_1 \approx 0$

$$\Rightarrow h_\theta(f) = \begin{cases} 1 \iff \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 + \ldots \geq 0 \\ 0 \quad else \end{cases}$$

- SVM with kernels

$$\underset{\Theta}{\Rightarrow} C\left[\sum_{y=1}^{m} y^i \, cost_1(\Theta^T f^i)\right] + (1-y^i)\, cost_0(\Theta^T f^i) + \frac{1}{2}\sum_{j=1}^{m} \Theta_j^2$$

$C = \frac{1}{\lambda}$ $\begin{cases} C \ggg \Rightarrow \text{ lower bias, high variance} \\ C \lll \Rightarrow \text{ higher bias, low variance} \end{cases}$

$\sigma^2 :$ $\begin{cases} \sigma^2 \gg \Rightarrow f_i \text{ will vary more smoothly} \\ \quad \Rightarrow \text{ higher bias, low variance} \end{cases}$ 

$\sigma^2 \ll \Rightarrow f_i \text{ will vary less smoothly}$ 
$\Rightarrow$ lower bias, high variance

- Using SPV packages
  - Choice of C
  - Choice of kernel

Ex: linear kernel $\Rightarrow y=1 \iff \Theta^T x \geqslant 0$

$\longrightarrow$ - Feature scaling before using Gaussian kernel!!

- SPV multiclass classification
  $\rightarrow$ 1 vs all method

- Logistic regression vs SVM

  $n \equiv n^o$ of features ; $m \equiv n^o$ of training examples

1. - If $n >> m \Rightarrow$ LR or SVM without kernel (linear)

2. - If $n > m \Rightarrow$ SVM with gaussian kernel

3. - If $n < m \Rightarrow$ Create/add more features & to 1

   In python problems this is not typical at all. Training set always way bigger than number of features.
   
   P

\# Neural networks likely to perform nice in every circumstance, but may be slower to train