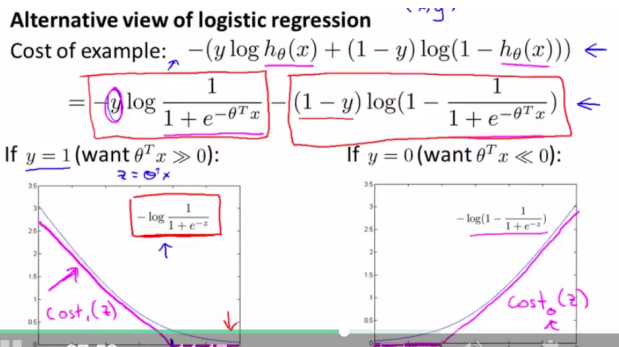- Optimization objective
- Alternative view of logistic regression
  - sigmoid activation
  - if y = 1
    - we want h_theta(x) close to 1
    - Theta'x >> 0
  - if y = 0
    - we want h_theta(x) close to 0
    - Theta'x << 0
- Alternative view of logistic regression
  - new cost function gives support vector machine advantages and easier optimization problem
  - two new cost functions one when y = 1 and one when y = 0



- Support vector machine
  - Logistic regression:



  - replace the cost function with the new cost functions for support vector machines
  - Logistic regression uses lambda to control the relative weighting between the first and second terms
    - large value of lambda gives B a very large weight
  - Support vector machine uses c as an alternative parameter
    - we then minimize CA + B
      - giving C a large value gives A a very large weight
  - Different way of controlling or parameterizing the weights
    - c = 1/lambda

- SVM hypothesis
    - does not output a probability
    - just makes a direct prediction if y = 1 or y = 0

**SVM hypothesis**

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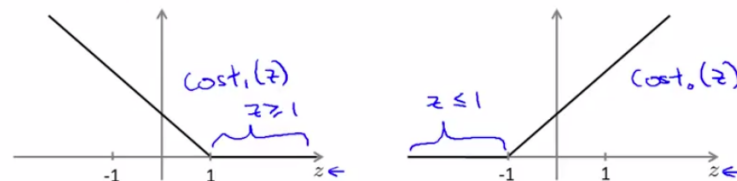
Hypothesis:

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

- Large Margin Intuition
- Support vector machine

**Support Vector Machine**

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

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

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

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

- SVM Decision Boundary
    - set the constant C to be a very large value
    - if C is large highly motivated to choose a term to make sure the first term is equal to zero in the objective function
- SVM Decision Boundary: Linearly separable case
    - margin minimum distance to any of the training examples
    - large margin classifier
- Large Margin classifier in presence of outliers
    - if C is very large with change the decision the boundary for an outlier
    - if C is not very large the SVM will not change the decision boundary for an outlier
- Kernels
    - Given x compute new feature depending on proximity to landmarks l(1), l(2), and l(3)

Given x:

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

$f_2 = similarity(x, l^{(1)}) = \exp\left(-\frac{\|x - l^{(2)}\|^2}{2\sigma^2}\right)$

$f_3 = similarity(x, l^{(3)}) = \exp(\ldots)$

Kernel (Gaussian kernels)

- Kernels and Similarity
  - measure how similar x is to one of the landmarks

**Kernels and Similarity**

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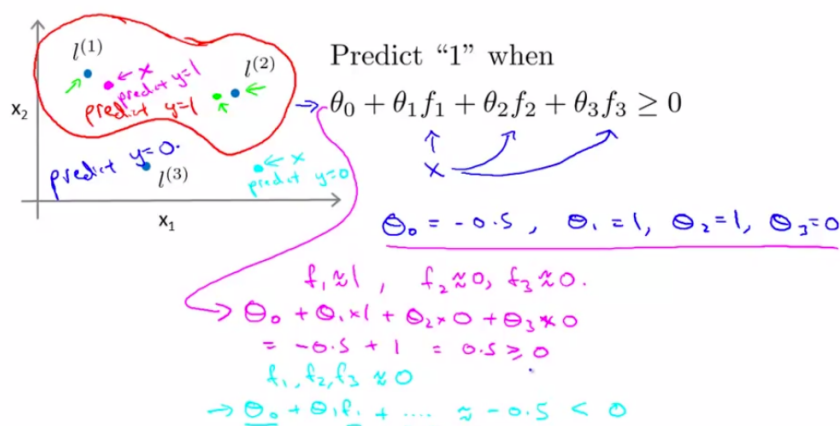
If $x \approx l^{(1)}$ :

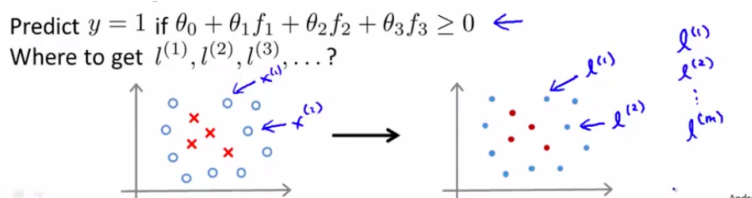$$f_1 \approx \exp\left(-\frac{0^2}{2\sigma^2}\right) \approx 1$$

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

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

- each landmark can now compute new features given the corresponding landmark
- varying sigma
  - .5 the feature falls to zero more rapidly
  - 3 the value of the feature falls away much more slowly



Predict "1" when
$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$$

$\theta_0 = -0.5$ , $\theta_1 = 1$, $\theta_2 = 1$, $\theta_3 = 0$

$f_1 \approx 1$ , $f_2 \approx 0$, $f_3 \approx 0$.

$\to \theta_0 + \theta_1 \times 1 + \theta_2 \times 0 + \theta_3 \times 0$
$= -0.5 + 1 = 0.5 \geq 0$
$f_1, f_2, f_3 \approx 0$
$\to \theta_0 + \theta_1 f_1 + \cdots \approx -0.5 < 0$

- Kernels 2
  - landmarks allowed us to define the similarity function or the kernel
  - take the examples and for every training example we are just going to call it put landmarks as the same locations as the training examples
  - features are going to measure how close something is as seen in training set

Predict $y = 1$ if $\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$ ←
Where to get $l^{(1)}, l^{(2)}, l^{(3)}, \ldots$ ?



- SVM with Kernels
  - Given m training examples
  - choose m landmarks

- Given example x:



Given example $x$:

$f_1 = \text{similarity}(x, l^{(1)})$
$f_2 = \text{similarity}(x, l^{(2)})$

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

- For training example



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

$x^{(i)} \rightarrow \begin{bmatrix} f_1^{(i)} = \text{sim}(x^{(i)}, l^{(1)}) \\ f_2^{(i)} = \text{sim}(x^{(i)}, l^{(2)}) \\ \vdots \\ f_i^{(i)} = \text{sim}(x^{(i)}, l^{(i)}) = \exp(-\frac{0}{2\sigma^2}) = 1 \\ \vdots \\ f_m^{(i)} \quad \text{sim}(x^{(i)}, l^{(m)}) \end{bmatrix}$

$x^{(i)} \in \mathbb{R}^{n+1}$ (or $\mathbb{R}^n$)

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

- SVM with Kernels
  - Hypothesis: Given x, compute features f element of R m+1
  - Training:
    - to obtain the parameters for the support vector machine solve the minimization problem



Training:

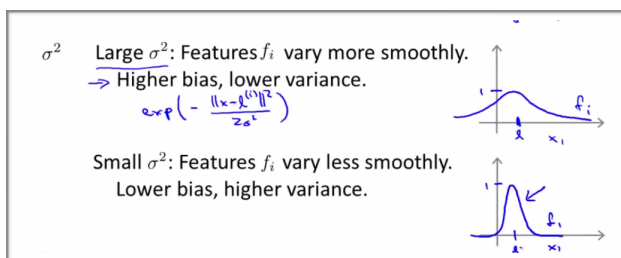$$\min_{\theta} C \sum_{i=1}^{m} y^{(i)} cost_1(\theta^T f^{(i)}) + (1 - y^{(i)}) cost_0(\theta^T f^{(i)}) + \frac{1}{2} \sum_{j=1}^{m} \theta_j^2$$

$n = m$

$\theta^T x^{(i)}$   $\theta^T f^{(i)}$   $\rightarrow \theta_0$

$\sum_j \theta_j^2 = \theta^T \theta \leftarrow \theta = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_m \end{bmatrix}$ (ignore $\theta_0$)

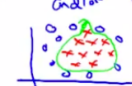$\rightarrow \theta^T M \theta \leftarrow$  $\|\theta\|^2$   $M = 10,000$

- SVM parameters
  - C = 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.
$\rightarrow$ 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.

- Using an SVM
- Support vector machine
  - choice of parameter C
  - choice of kernel
  - e.g. No kernel ("linear kernel")
    - predict "y = 1" if theta_t dot x >= 0

E.g. No kernel ("linear kernel")
Predict "y = 1" if $\theta^T x \geq 0$

$\theta_0 + \theta_1 x_1 + \cdots + \theta_n x_n \geq 0$
$\rightarrow$ n large, m small $\quad x \in \mathbb{R}^{n+1}$

Gaussian kernel:
$$f_i = \exp\left(-\frac{||x - l^{(i)}||^2}{2\sigma^2}\right), \text{ where } l^{(i)} = x^{(i)}.$$
Need to choose $\sigma^2$.

$x \in \mathbb{R}^n$, n small
and/or m large

- Kernel (similarity) functions:

**Kernel (similarity) functions:** $x^{(i)}$ $l^{(j)} = x^{(j)}$
```
function f = kernel(x1,x2)
```
$$f = \exp\left(-\frac{||x1 - x2||^2}{2\sigma^2}\right) \leftarrow$$
```
return
```
$x \rightarrow \begin{matrix} f_1 \\ f_2 \\ \vdots \\ f_m \end{matrix}$

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

  - feature scaling is important
- Other choices of kernel
  - Not all similarity functions 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).
  - polynomial kernel: k(x,l) = (x'l)^2, (x'l)^3, (x'l + 1)^3, (x'l + 5)^4
  - more esoteric: String kernel, chi-square kernel, histogram intersection kernel
- Logistic regression VS SVMS

**Logistic regression vs. SVMs**

$n =$number of features ($x \in \mathbb{R}^{n+1}$), $m =$number of training examples
$\rightarrow$ If $n$ is large (relative to $m$): (e.g. $n \geq m$, $n = 10,000$, $m = 10 \cdots 1000$)
$\rightarrow$ Use logistic regression, or SVM without a kernel ("linear kernel")

$\rightarrow$ If $n$ is small, $m$ is intermediate: ($n = 1 - 1000$, $m = 10 - 10,000$) $\leftarrow$
$\quad \rightarrow$ Use SVM with Gaussian kernel

If $n$ is small, $m$ is large: ($n = 1 - 1000$, $m = 50,000+$)
$\quad \rightarrow$ Create/add more features, then use logistic regression or SVM without a kernel

$\rightarrow$ Neural network likely to work well for most of these settings, but may be slower to train.