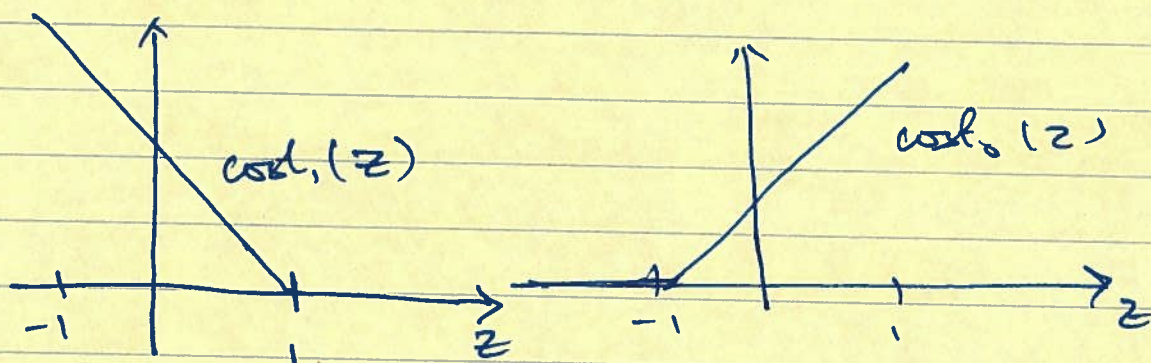Large Margin Intuition:

## SVM

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

$$+ \frac{1}{2} \sum_{j=1}^{n} \theta_j^2$$



$cost_1(z)$

$cost_0(z)$

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

whenever $y^{(i)} = 1$

$\theta^T x^{(i)} \geq 1$

$\min C0 + \frac{1}{2} \sum_{j}^{n} \theta_j^2$

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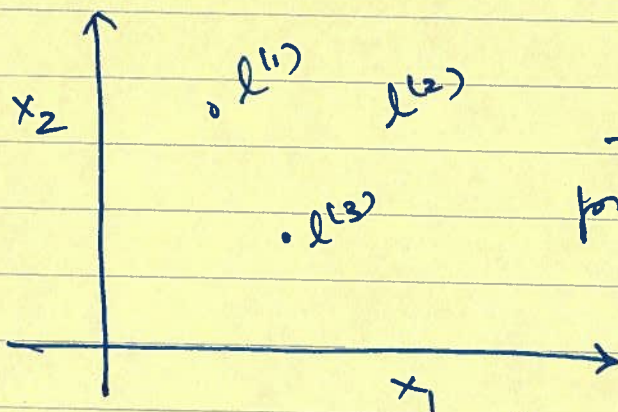
$\theta^T x^{(i)} \leq -1$

# Kernels - I

Non-linear Decision Boundary
Predict $y = 1$ if

$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 + \theta_5 x_2^2 + \cdots \geq$$

$$h_\theta(x) = \begin{cases} 1 & \text{if } \theta_1 + \theta_1 x_1 + \cdots \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 + \cdots$$
$$f_1 = x_1, \quad f_2 = x_2, \quad f_3 = x_1 x_2, \quad f_4 = x_1^2, \quad f_5 = x_2^2$$



Given $x$, compute new features depending on proximity to landmarks $\ell^{(1)}, \ell^{(2)}, \ell^{(3)}$

Given $x$:

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

$K(x, \ell^{(i)})$

$$f_2 = \text{similarity} \left( x, \ell^{(1)} \right)$$

Kernels
(Gaussian kernels) $= \exp \left( - \frac{\| x - \ell^{(1)} \|^2}{2a^2} \right)$

$$f_3 = \text{similarity}(x, \ell^{(3)})$$
$$= \exp \left( - \frac{\| x - \ell^{(3)} \|^2}{2a^2} \right)$$

# Kernels & Similarity

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

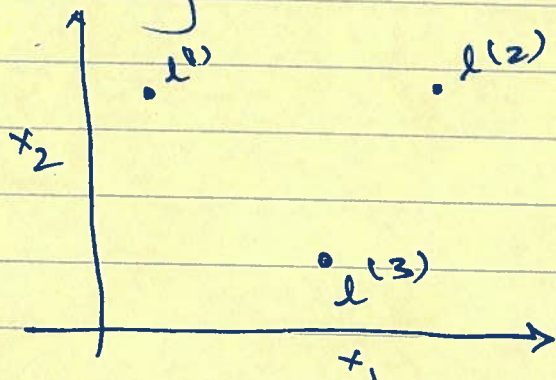$$= \exp\left(-\frac{\sum_{j=1}^{n}(x_j - \ell_j^{(1)})^2}{2\sigma^2}\right)$$

if $x \approx \ell^{(1)}$

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

if $x$ if far from $\ell^{(1)}$:

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

## Choosing the landmark.



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

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

Given $(x^{(1)}, y^{(1)})$, $(x^{(2)}, y^{(2)})$, ..., $(x^{(m)}, y^{(m)})$

choose $l^{(1)} = x^{(1)}$, $l^{(2)} = x^{(2)}$, ..., $l^{(m)} = x^{(m)}$

Given example $x$:

$f_1 = $ similarity $(x, l^{(1)})$  $\qquad f = \begin{bmatrix} f_0 \\ f_1 \\ f_2 \\ \vdots \\ f_m \end{bmatrix} \quad f_0 = 1$

$f_2 = \qquad\qquad\ddots \qquad (x, l^{(2)})$

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

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

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

$f_0^{(i)} = 1$

Hypothesis
Given $x$, compute
features
$f \in \mathbb{R}^{m+1}$
Predict
$y = 1$ if $\theta^T f \geq 0$

Training

$$\sum_j \theta_j = \theta^T \theta \longleftarrow \quad \theta = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_m \end{bmatrix} \quad \text{ignoring } \theta_0.$$

$$\theta^T M \theta \qquad \|\theta\|^2$$

inc c
dec $a^2$
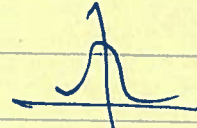
dec c inc $a^2$

$c = \frac{1}{\lambda}$  Large c : Lower bias, high var
Small c : Higher bias, low var.

large c means small $\lambda$
Small c  ..  large $\lambda$.

Large $a^2$ : features $f_i$ vary more
smoothly. High bias, low var.

Small $a^2$, features $f_i$ vary less
smoothly
lower bias, higher var.

If
Overfit:  Dec c, Inc $a^2$

Need to specify parameter $c$:
  Choice of $c$.

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

Gaussian Kernel
  $x \in \mathbb{R}^n$   $n$ small

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

  where $\ell^{(i)} = x^{(i)}$,

Need to choose $\sigma^2$,

$\|x - \ell\|^2$         $v = x - \ell$

$$\|v\|^2 = v_1^2 + v_2^2 + \cdots + v_n^2$$
$$= (x_1 - \ell_1)^2 + (Y_2 - \ell_2)^2 + (x_n - \ell_n)^2.$$

Off Shelf Kernels.
Polynomial Kernels $k(x, \ell) = (x^T \ell)^2$
         $(x^T \ell)^3$  , $(x^T \ell + 1)^3$
         $(x^T \ell + 5)^4$
String kernel, chi-squared Kernel, histogram
  intersection Kernel.

# Multi class Classification

one vs all
  (Train K SVMs)

Logistic Ree vs SVM

If n is large (relative m)

$n$ = num of features
$m$ = num of train ex
  Eg $n \geqslant m$    $n = 10,000$

  $(m = 10 \cdots 1000)$
Use logistic reg or svm without a
                    kernel ("linear kernel"

If n is small, m is intermediate
  Use SVM with Gaussian.
  $n = 1 - 1000$, $m = 10 \sim 10,000$)
if n is small, m is large
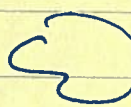  $(n = 1 - 1000$, $m = 50,000$ or $>$)
  create/add more features
  then use logistic reg
    or SVM without a kernel.

Overfit : Dec $c$, inc $a^2$

line $c$
dec $a^2$



dec $c$ inc $a^2$.

$y^{(i)} = 1$, $\theta^T x^{(i)} \geqslant 1$
$y^{(i)} = 0$, $\theta^T x^{(i)} \leqslant -1$

Try using NN with large hid.
Use SVM with Gauss Ker
Create / add new poly fea
Use SVM with a lin ker, without new
$x^{(i)} \in \mathbb{R}^2$, dec bond is st. line.
It is imp to perform feature
norm before using Gau ker
The max value of Gauss ker
$sim(x, l^{(i)})$ is 1