

WEEK-7

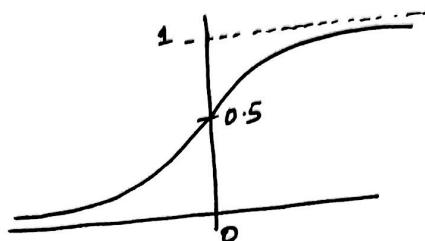
Support Vector Machines

(8)

Large Margin Classification

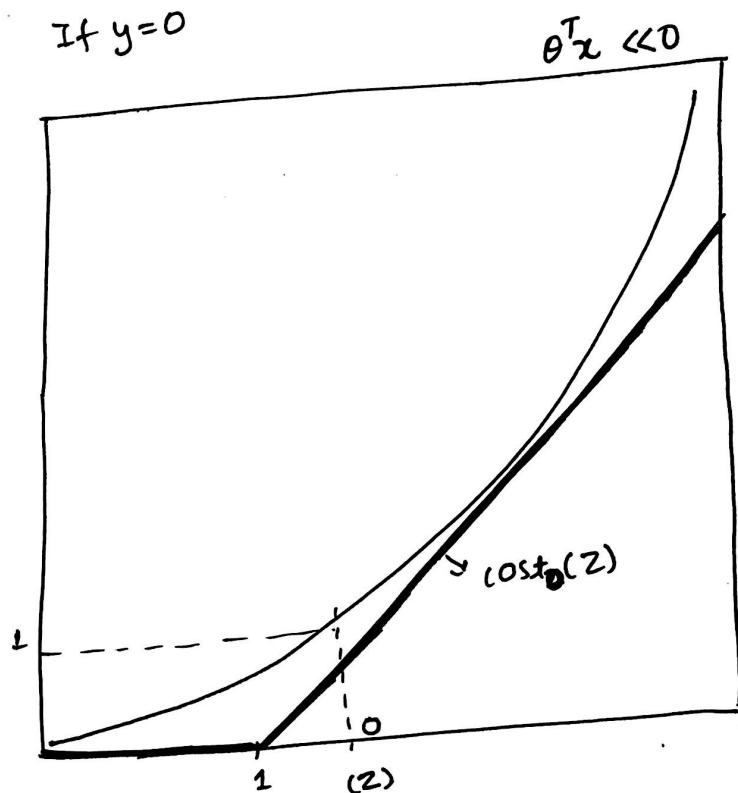
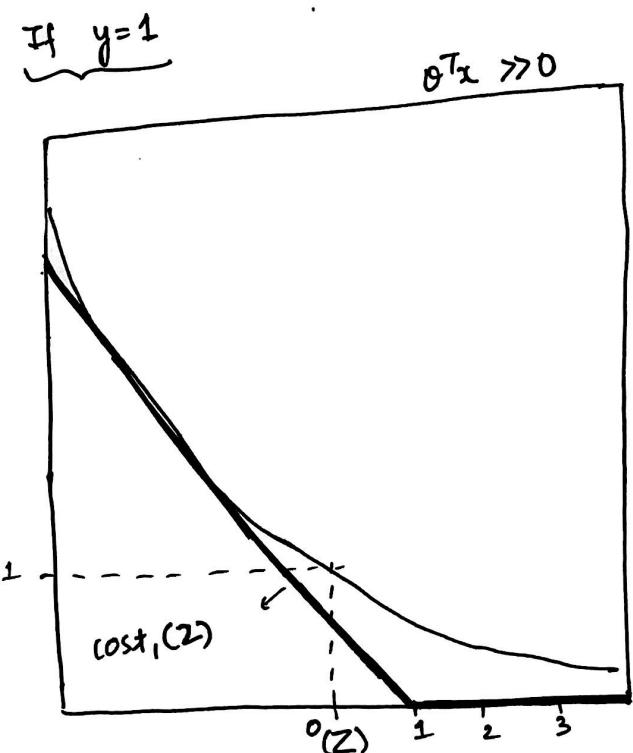
→ Alternate view of logistic regression

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



[if $y=1$, we want $h_{\theta}(x) \approx 1$, $\theta^T x \gg 0$]
 [if $y=0$, we want $h_{\theta}(x) \approx 0$, $\theta^T x \ll 0$]

$$\rightarrow \text{cost} = -y \log\left(\frac{1}{1+e^{-\theta^T x}}\right) - (1-y) \log\left(1 - \frac{1}{1+e^{-\theta^T x}}\right)$$

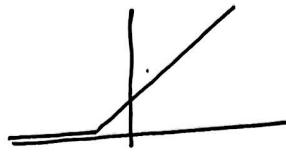
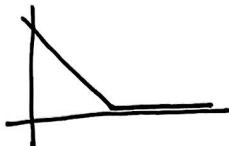


Support Vector Machine (SVM)

Logistic Regression:

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

$\underbrace{cost_1(\theta^T x^i)}$ $\underbrace{cost_0(\theta^T x^i)}$



SVM:

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

(A) (B)

+ $\frac{1}{\lambda}$

controlling the trade off b/w A & B.

$c(A) + B$
controlling the trade off. same as $\lambda \Rightarrow C \approx \frac{1}{\lambda}$

cost fn. of SVM.

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

SVM hypothesis:

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

SVM Intuition :

if $y=1$, we want $\theta^T x \geq 1$ (not just ≥ 0) } extra Safety Margin
 if $y=0$, we want $\theta^T x \leq -1$ (not just < 0) }

SVM Decision Boundary :

Case: C is very large, then $y^i \text{cost}_1(\theta^T x) + (1-y^i)\text{cost}_0(\theta^T x) \approx 0$

whenever $y^i = 1$ and $\theta^T x \geq 1$

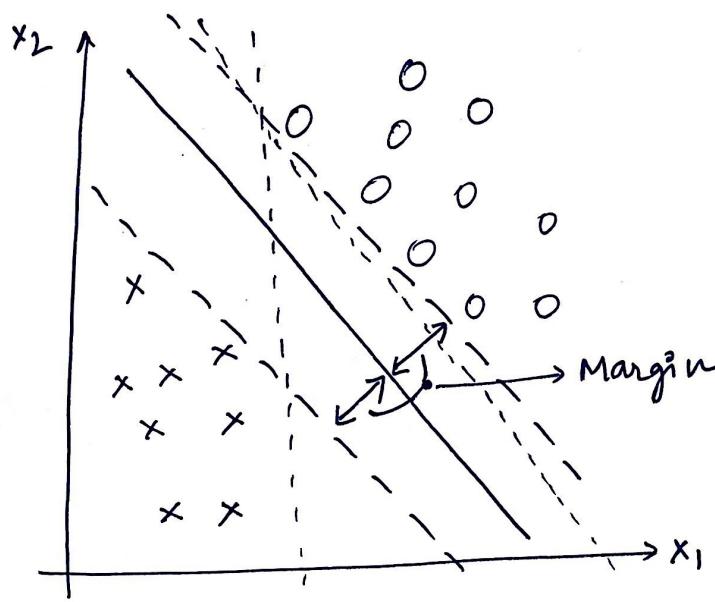
$$\min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

s.t. $\theta^T x \geq 1$

whenever $y^i = 0$ and $\theta^T x \leq -1$

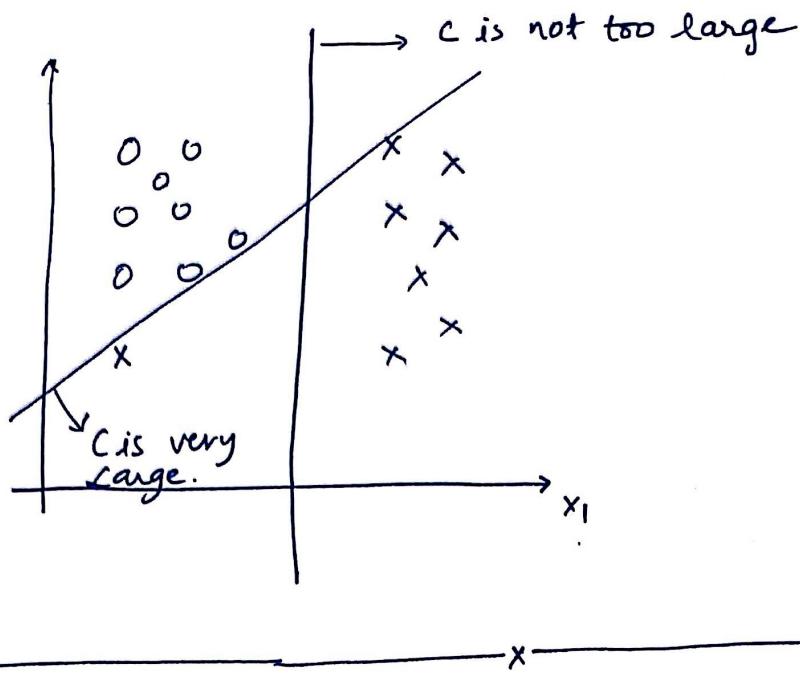
$$\min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

s.t. $\theta^T x \leq -1$

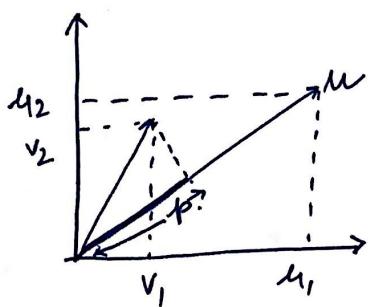


Therefore, SVM's are also called LMC (large margin classifiers)

- Large margin classifier in presence of outliers:



Mathematics behind SVM



$$\begin{aligned} ||u|| &= \text{length of vector } u \\ &= \sqrt{u_1^2 + u_2^2} \end{aligned}$$

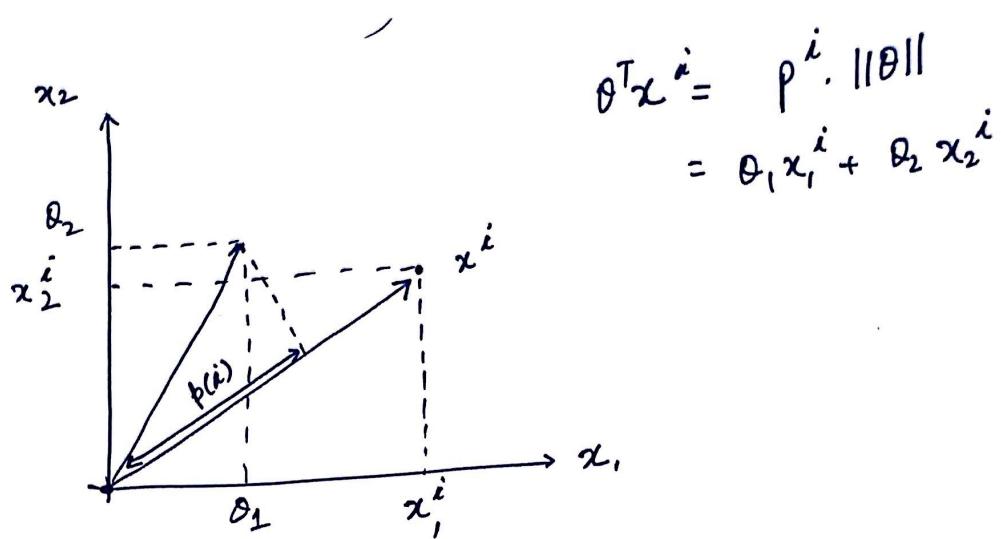
$p = \text{length of projection of } v \text{ onto } u$
signed.

$$\begin{aligned} u^T v &= p \cdot ||u|| \\ &= u_1 v_1 + u_2 v_2 \end{aligned}$$

→ SVM decision Boundary

$$\begin{aligned}
 & \min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2 \quad \left. \right\} \text{simplification, } \theta_0 = 0 \\
 \text{s.t.} \quad & \theta^T x \geq 1 \quad \text{if } y^i = 1 \\
 & \theta^T x \leq -1 \quad \text{if } y^i = 0
 \end{aligned}$$

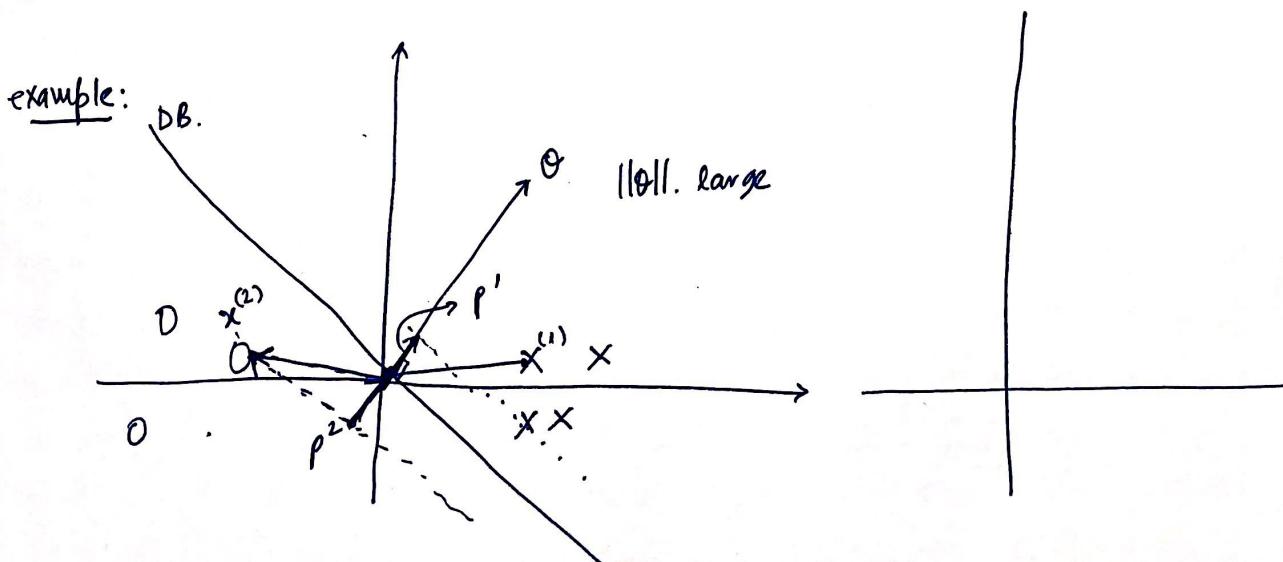
$$\therefore \text{obj. fn} = \frac{1}{2} (\sqrt{\theta_1^2 + \theta_2^2})^2 = \frac{1}{2} \|\theta\|^2$$



∴ ~~s.t.~~ New objective:

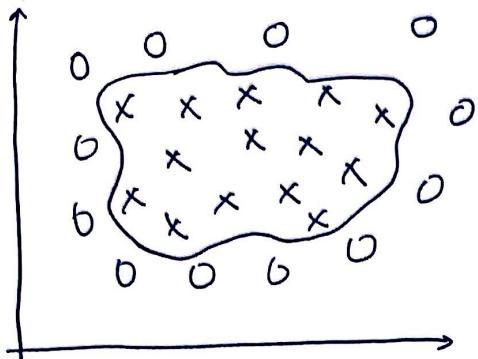
$$\min_{\theta} \frac{1}{2} \cdot \|\theta\|^2$$

$$\begin{aligned}
 \text{s.t.} \quad & p^i \cdot \|\theta\| \geq 1 & \text{if } y^i = 1 \\
 & p^i \cdot \|\theta\| \leq -1 & \text{if } y^i = 0
 \end{aligned}$$



KERNELS

Non-linear Decision Boundary

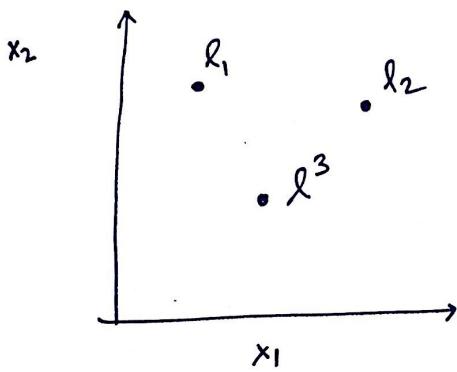


predict $y=1$ if

$$\cancel{\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2} + \dots \theta_{27} x_1^6 x_2^6 \dots \geq 0$$

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

is there a different/better choice of the features?
(higher order polynomials)



Given x , compute new feature depending on proximity to landmarks l^1, l^2, l^3

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

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

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

Kernel (~~Gaussian~~ 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 $x \approx l^1$

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

$$l_1 \rightarrow f_1$$

$$l_2 \rightarrow f_2$$

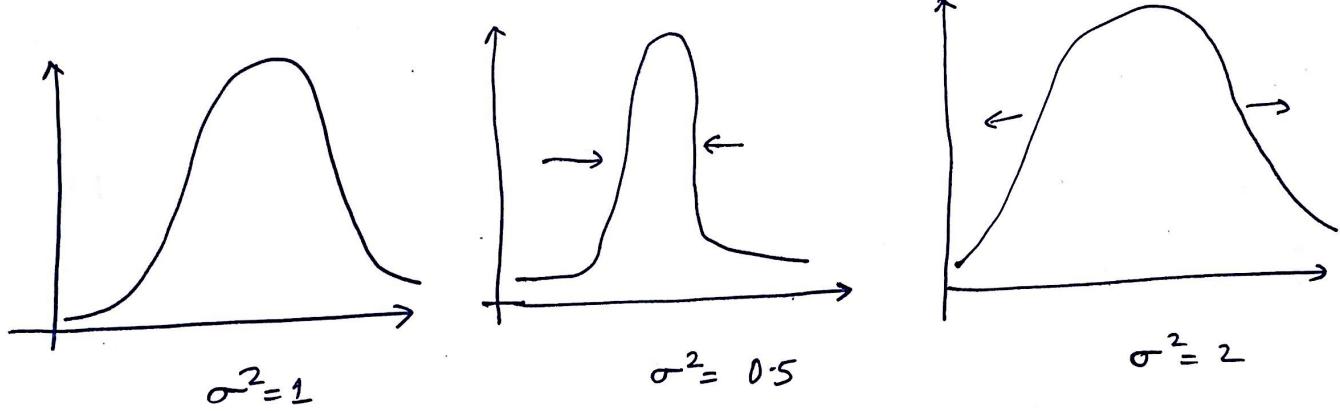
$$l^3 \rightarrow f_3$$

$$\tilde{x}^i$$

If x is far from l^2

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

$\left. \begin{array}{l} \text{new set of features} \\ \text{based on landmarks} \end{array} \right\}$



New Hypothesis :

Predict ① when

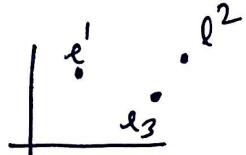
$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0 \quad \Rightarrow \quad \theta_0 = -0.5$$

$$\theta_1 = 1$$

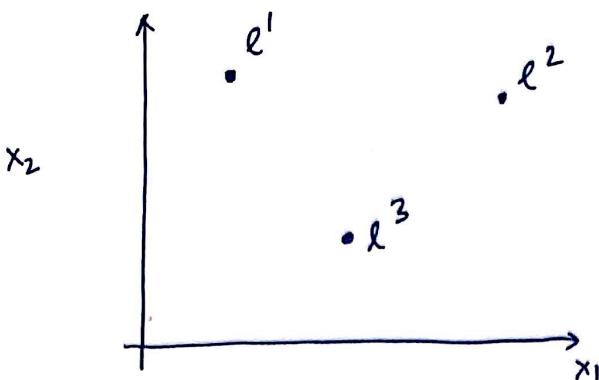
$$\theta_2 = 1$$

$$\theta_3 = 0$$

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



choosing the landmarks:-



$$f_i = \text{similarity} \left(-\frac{\|x - l^i\|^2}{2 \cdot \sigma^2} \right)$$

[for every training example, choose landmark at exactly same position!]

Given: $(x^1, y^1), (x^2, y^2) \dots (x^m, y^m)$

choose: $l^1 = x^1, l^2 = x^2, \dots l^m = x^m$

$$\therefore f = \begin{bmatrix} f_0 \\ f_1 \\ f_2 \\ \vdots \\ f_m \end{bmatrix} \quad f_0 = 1 \quad f \in \mathbb{R}^{m+1}.$$

where

for training example (x^i, y^i)

$$x^i \rightarrow \left\{ \begin{array}{l} f_1^i = \text{sim}(x^i, l^1) \\ f_2^i = \text{sim}(x^i, l^2) \\ \vdots \\ f_m^i = \text{sim}(x^i, l^m) \end{array} \right\} \quad f_i^i = \text{sim}(x^i, l^i) \times 1$$

$x^i \in \mathbb{R}^{n+1} \Rightarrow [f^i \in \mathbb{R}^{m+1}]$ vector
feature vector

$$f_0^i = 1$$

SVM with Kernels (Algorithm)

Hypothesis: Given x , compute features $f \in \mathbb{R}^{m+1}$

predict " $y=1$ " if $\underbrace{\theta^T f}_{} \geq 0$

$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \dots + \theta_m f_m$$

$$\theta \in \mathbb{R}^{m+1}$$

Training:

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

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_m \end{bmatrix}$$

$$\sum_j \theta^T \theta$$

$$\underline{\theta^T M \theta}$$

SVM optimization.
computational
eff.

SVM Parameters:

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

$\left\{ \begin{array}{l} \text{large } C: \text{lower bias, high variance} \\ \text{small } C: \text{higher bias, lower variance} \end{array} \right\}$

σ^2 $\left\{ \begin{array}{l} \text{large } \sigma^2: \text{features } f_i \text{ vary more smoothly} \\ \text{small } \sigma^2: \text{features } f_i \text{ vary less smoothly} \end{array} \right\}$
 Higher bias, lower variance.
 Lower bias, higher variance (overfit)



USING AN SVM

→ Use SVM Software (e.g. liblinear, libsvm...)

→ Need to specify :

- ↳ choice of parameter C
- ↳ choice of kernel.

e.g. No. Kernel. ("linear kernel")

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

(n large, m small)

~~in R^n+1~~
 $x \in \mathbb{R}^{n+1}$

Gaussian Kernel.

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

Need to choose σ^2

$x \in \mathbb{R}^{n+1}$, (n small. and m is large)

* Note: Do perform feature scaling before using the Gaussian Kernel.

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

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

return.

$x_1 \in x^i$

$x_2 \in \ell^i$

- Other choices of Kernel

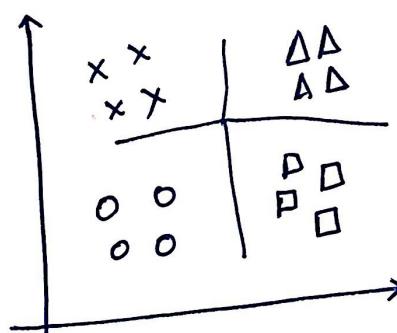
Not all the similarity functions $\text{similarity}(x, l)$ make valid kernels. (Need to satisfy technical condition called "Mercer Theorem" to make sure SVM packages' optimizations run correctly and do not diverge).

Many off-the-shelf kernels available

→ Polynomial Kernel: $k(x, l) = (x^T l)^2$
 ↳ $(x^T l + \text{constant})^{\text{degree}}$.

→ Esoteric: string kernel, chi-square kernel, histogram intersection kernel, ...

MULTI-CLASS CLASSIFICATION



Many SVM packages already have built-in multi-class classification functionality

otherwise, use one-vs-all method

- ① Train K SVM's, one to distinguish $y=i$ from the rest, for $i=1, 2, \dots, K$
- ② get $\theta^1, \theta^2, \dots, \theta^K$
- ③ Pick class i with largest ~~θ^i~~ $(\theta^i)^T \alpha$

Logistic Regression vs. SVMs① if $n \geq m$

$$n = 10,000$$

$$m = 10\dots \text{or } 1000$$

use logistic regression or ~~SVM~~ SVM without a Kernel
(Linear Kernel)

② if n is small, m is intermediate.

$$n = 1-1000$$

$$m = 10-10,000$$

use SVM with Gaussian Kernel.

③ if n is small, m is large

$$n = 1-1000$$

$$m = 50000+$$

Create/add more features, then use Logistic Regression
or SVM without a Kernel.

Similar
Algo's.

{ # Neural network likely to work-well for most of
these settings, but may be slower to train. }

Octave code: for cv test.

```

for C-test = [0.01 0.03 0.1 0.3 1 3 10 30]
for sigma-test = [0.01 0.03 0.1 0.3 1 3 10 30]

    i = i+1
    model = svmTrain(x, y, C-test, @(x1, x2) gaussianKernel
                      (x1, x2, sigma-test));
    % sum package to minimize
    predictions = svmPredict(model, xval)
    prediction-error = mean(double(predictions ~= yval));
    result(i, :) = [C-test sigma-test prediction-error]
endfor
end for.

sorted_result = sortedrows(result, 3)

c = sorted_result(1, 1);
sigma = sorted_result(1, 2);

```