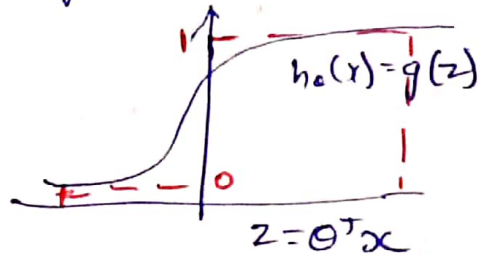


Week 7 - Support Vector Machine

Derivation

Consider logistic regression

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



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

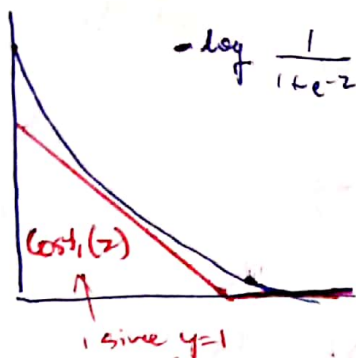
If $y = 0$, we want $h_{\theta}(x) \approx 0$, $\theta^T x \ll 0$

Cost of logistic regression:

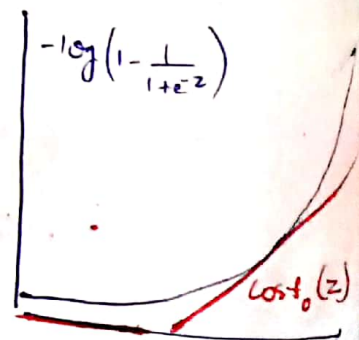
$$-(y \log h_{\theta}(x) + (1-y) \log (1-h_{\theta}(x)))$$

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

If $y = 1$ (want $\theta^T x \gg 0$):



If $y = 0$ (want $\theta^T x \ll 0$):



\therefore Logistic regression:

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

• Can't in terms

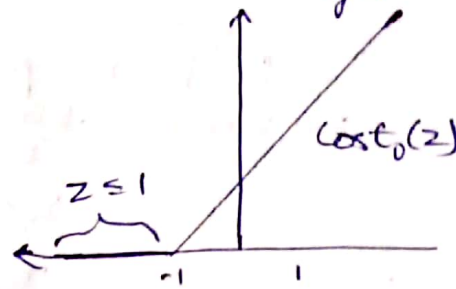
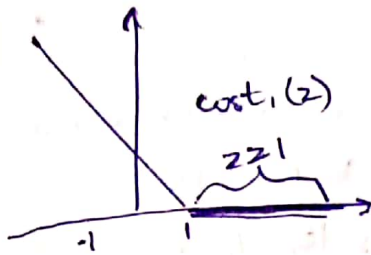
• Before $A + \lambda B \rightarrow$ Now $CA + B$ \leftarrow C is like $\frac{1}{\lambda} + \frac{1}{2m} \sum_{j=1}^n \theta_j^2$

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

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

Support Vector machine

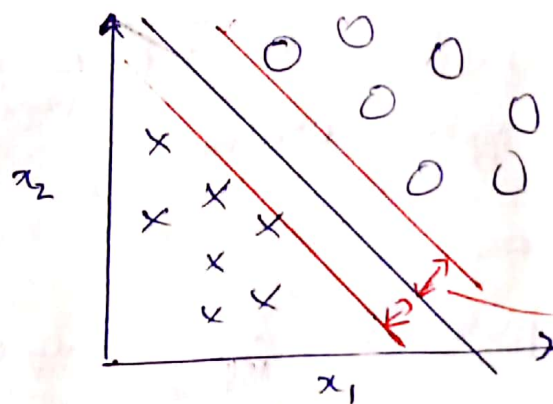
$$\min_{\theta} \left(\sum_{i=1}^m \left[y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2 \right)$$



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

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

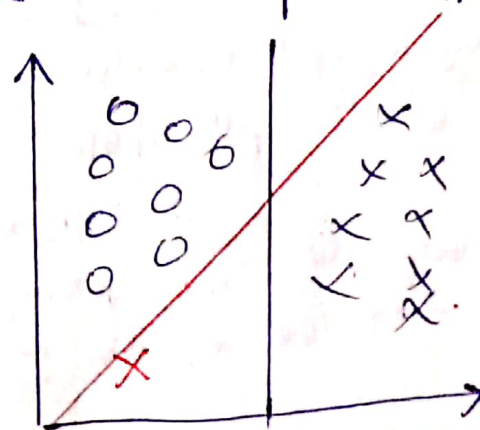
SVM decision boundary:



This margin comes because we are using $\theta^T x \geq 1$ & $\theta^T x \leq -1$ instead of $\theta^T x \geq 0$ & $\theta^T x \leq 0$

∴ It is also called a large margin classifier

How does value of C affect it?



Lets say we introduce a new point x . Then if C is large, the line shifts. If C is small, it remains same.

We may want to keep the line same since we don't want to change it so much just for one data value.

Mathematics behind Large Margin Classification

Vector inner product:

$U^T V$ is called vector inner product

i.e. $[u_1, u_2] \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$

$$U^T V = p \cdot \|u\|$$

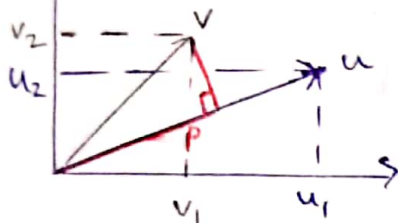
$$= u_1 v_1 + u_2 v_2 = p \cdot \|u\|$$

Consider the 2 vectors V and u .
We project V on u .

Then $\|u\|$ is length of vector u

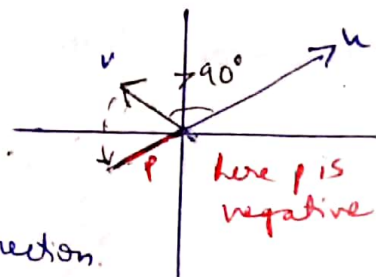
$\|u\| = \sqrt{u_1^2 + u_2^2}$ From Pythagorean theorem

p = length of projection V onto u ,



Here, p can also be negative if the projection angle $> 90^\circ$

So we project V in opposite direction.



SVM decision boundary:

$$\min_{\theta} \frac{1}{2} \sum_{i=1}^n \theta_i^2 \quad \text{Suppose } \theta_0 = 0 \text{ \& } n=2 \text{ (only } \theta_1, \theta_2)$$

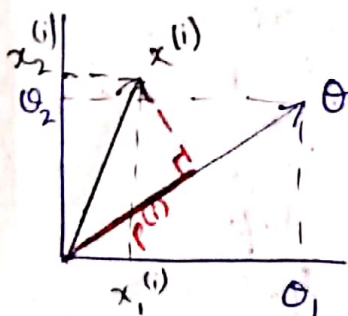
$$\therefore \frac{1}{2} (\theta_1^2 + \theta_2^2) = \frac{1}{2} (\sqrt{\theta_1^2 + \theta_2^2})^2 = \frac{1}{2} \|\theta\|^2$$

This holds true since $\theta_0 = 0$ $\begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix}$

\therefore We are basically trying to minimize length of vector θ

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

Here $p^{(i)}$ is the projection of $x^{(i)}$ onto the vector θ .



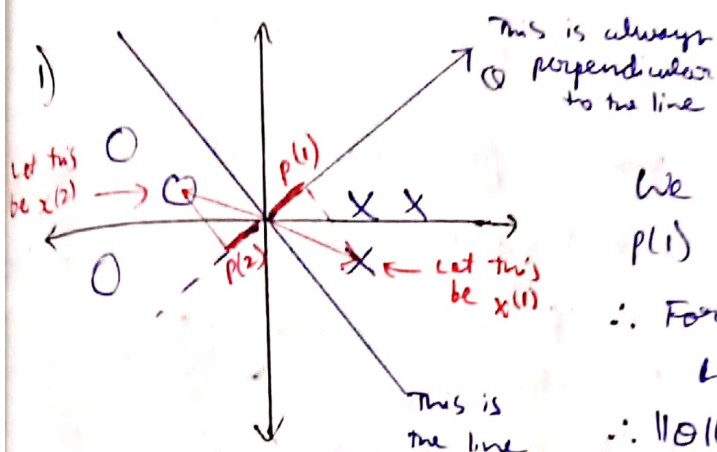
$$\therefore \min_{\theta} \frac{1}{2} \sum_{i=1}^n \theta_i^2 = \frac{1}{2} \|\theta\|^2$$

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

consider 2 cases to see how SVM finds the decision boundary:

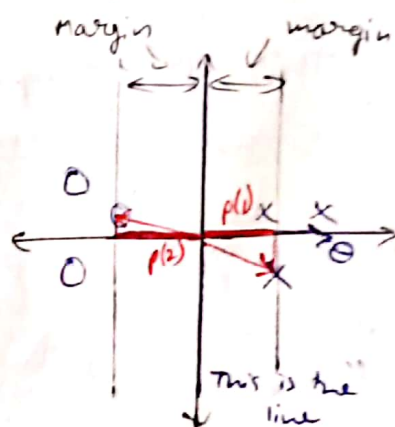
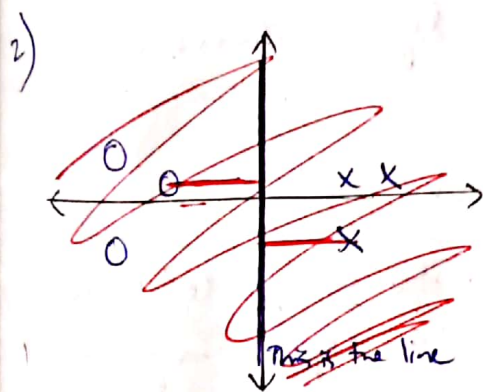
Let $\theta_0 = 0$

This means line needs to pass through origin



We can see here that $p(1)$ and $p(2)$ are small
 \therefore For $p(1) \cdot \|\theta\| \geq 1$
 $\quad \quad \quad \& \quad p(2) \cdot \|\theta\| \leq -1$
 $\therefore \|\theta\|$ needs to be large

But we are minimizing $\|\theta\|$ since, so this line isn't good



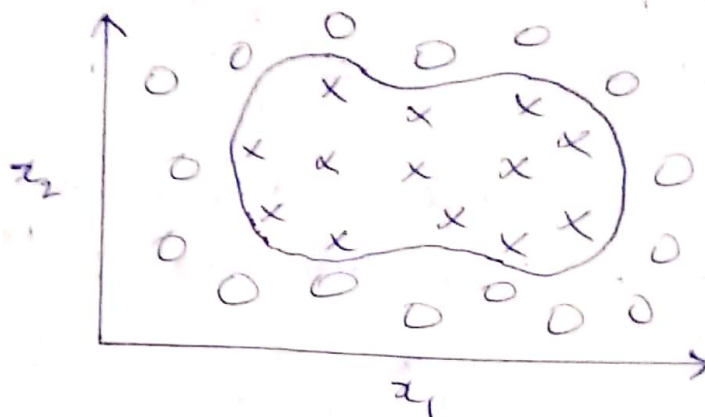
Notice here that $p(1)$ and $p(2)$ are large

\therefore For $p(1) \cdot \|\theta\| \geq 1$ $\|\theta\|$ can be small
 $p(2) \cdot \|\theta\| \leq -1$, \therefore This is what SVM will choose

$p(1), p(2), p(3), \dots$ is basically the margin

Here we have set $\theta_0 = 0$ and confined it to the origin, but it's the same logic even if $\theta_0 \neq 0$ and the line doesn't pass through origin.

Non-linear Boundaries

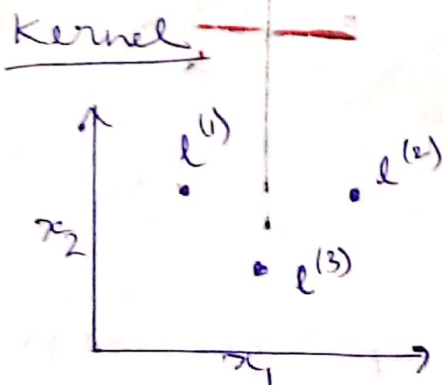


Usually we use polynomial terms
 Predict $y=1$ if,
 $\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 + \theta_5 x_2^2$

$$h_\theta(x) = \begin{cases} 1 & \text{if } \theta_0 + \theta_1 x_1 + \dots \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

We will write this as $\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3$,
 where $f_1 = x_1$, $f_2 = x_2$, $f_3 = x_1 x_2$ and $f_4 = x_1^2$, $f_5 = x_2^2$.

But this is computationally intensive.
 So we use kernels.



Given x , compute new feature depending on proximity to landmarks $l^{(1)}$, $l^{(2)}$, $l^{(3)}$

Given x ,

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

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

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

↑
kernel

↑ Gaussian kernel $k(x, l^{(i)})$

Let's take x_1 ,

$$f_1 = \text{similarity}(x, l^{(1)}) = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right) = \exp\left(-\frac{\sum_{i=1}^n (x_i - l_i^{(1)})^2}{2\sigma^2}\right)$$

If $x \approx l^{(1)}$:

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

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

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

So each landmark ~~defines~~ ^{defines} a new feature f
 $l^{(1)} \rightarrow f_1$, $l^{(2)} \rightarrow f_2$, $l^{(3)} \rightarrow f_3$

SNACK PRINTOUT

Example:

$$\rightarrow l^{(1)} = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$$

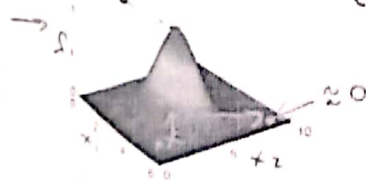
$$\rightarrow \sigma^2 = 1$$

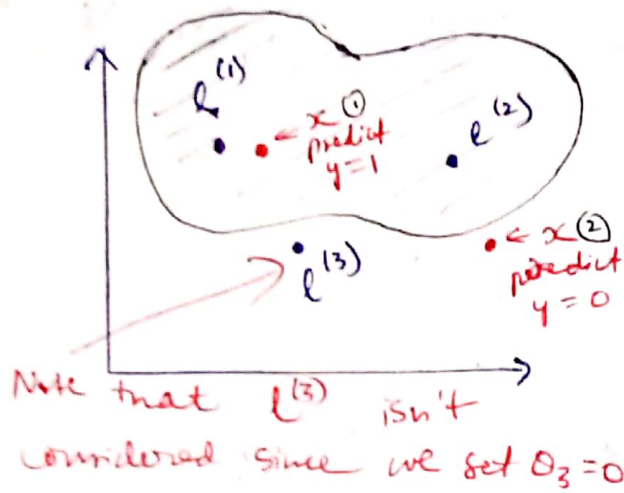
$$f_1 = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right)$$

$$x = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$$

$$\sigma^2 = 0.5$$

$$\sigma^2 = 3$$





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

Let's set
 $\theta_0 = -0.5$
 $\theta_1 = 1$
 $\theta_3 = 0$ as an example

Consider 2 cases:

① Here x is close to $l^{(1)}$

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

$$= \theta_0 + \theta_1 \times 1 + \theta_2 \times 0 + \theta_3 \times 0$$

$$= -0.5 + 1 = 0.5 \geq 0 \quad \therefore y = 1$$

② Here x is far from $l^{(1)}, l^{(2)}$ and $l^{(3)}$

$$\therefore f_1, f_2, f_3 \approx 0$$

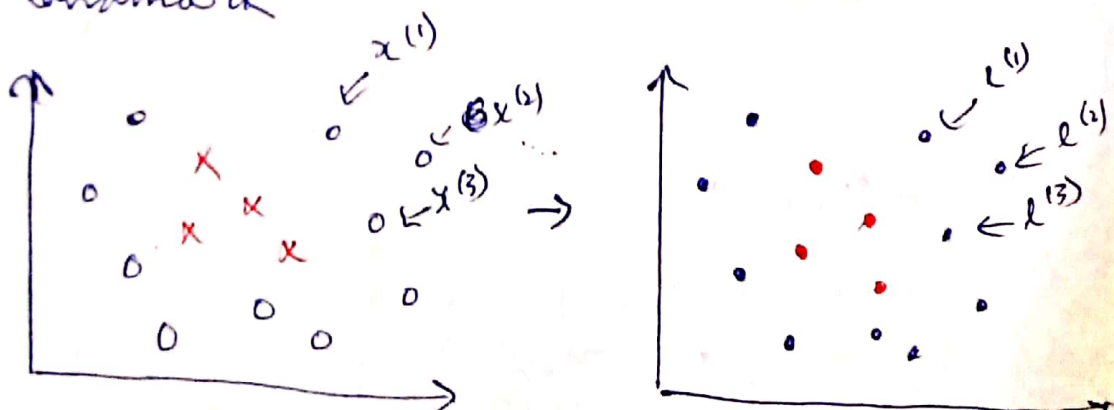
$$= \theta_0 + \theta_1 \times 0 + \theta_2 \times 0 + \theta_3 \times 0$$

$$= -0.5 \leq 0 \quad \therefore y = 0$$

Choosing landmarks:

Before we chose $l^{(1)}, l^{(2)}, l^{(3)}$, but how do we choose it in problem?

We take all the training points as landmark



\therefore We end up with $l^{(1)}, l^{(2)}, l^{(3)}, \dots, l^{(n)}$

Given $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})$
 Choose $l^{(1)} = x^{(1)}, l^{(2)} = x^{(2)}, \dots, l^{(m)} = x^{(m)}$

Given example x ,

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

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

...

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

where $f_0 = 1$
 (we set this like)
 $\theta_0 = 1$

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

$$f_1^{(i)} = \text{sim}(x^{(i)}, l^{(1)})$$

$$f_2^{(i)} = \text{sim}(x^{(i)}, l^{(2)})$$

$$f_i^{(i)} = \text{sim}(x^{(i)}, l^{(i)}) = \exp\left(-\frac{D}{2\sigma^2}\right) = 1$$

somewhere in the middle

$$f_m^{(i)} = \text{sim}(x^{(i)}, l^{(m)})$$

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

Sim with kernels

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

Predict " $y=1$ " if $\theta^T f \geq 0$

$\theta_0 f_0 + \theta_1 f_1 + \dots + \theta_m f_m$

Training:

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

n is equal to m

Notice this, we are using this instead of $\theta^T x^{(i)}$

Mathematically this term is different. They use $\theta^T M \theta$, but that's only to decrease computational intensity.

It has an already built software to minimise included in it.

SVM parameters

- $C (= \frac{1}{\lambda})$ Large C : Lower bias, high variance
Small C : Higher bias, low variance
- σ^2 Large σ^2 : Features f_i vary more smoothly
→ Higher bias, lower variance
Small σ^2 : Features f_i vary less smoothly
→ Lower bias, higher variance

SVM in Practice

- Use SVM software package (e.g. liblinear, libsvm, ...) to solve for parameters C .

Need to specify:

- Choice of parameter C
- Choice of kernel (similarity function):

E.g. ① No kernel (linear kernel):

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

We use this when n (no. of features) is large but m (~~data~~ no. of training examples) is less.

② Gaussian kernel:

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

- Need to choose σ^2

We use this when n is small and/or m is large.

If we choose Gaussian kernel:

The SVM software package will ask you to implement the kernel (similarity) function

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

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

return

Note: Perform feature scaling before using the Gaussian kernel.

The SVM software package will then use this function to generate all the features $f_1, f_2, f_3, \dots, f_n$ and train the SVM.

Some SVM software packages will have inbuilt kernel functions and will allow you to choose

Other choices of kernels:

- All available kernels need to satisfy "Mercer's Theorem" to make sure SVM packages' optimization run correctly

① Polynomial kernel: $K(x, l) = (x^T l + \text{constant})^{\text{degree}}$
 $= (x^T l)^2, (x^T l)^3, (x^T l + 1)^3, \dots$

② More esoteric kernels: string kernel, Chi-square kernel, histogram intersection kernel, ...

Multiclass Classification:

- Many SVM packages already have built-in multi-class classification functionality.
- Otherwise use one-vs-all method (Train K SVMs, one to distinguish $y=i$ from the rest, for $i=1, 2, \dots, K$), get $\hat{y}^{(1)}, \hat{y}^{(2)}, \dots, \hat{y}^{(K)}$. Pick class i with largest $(\hat{y}^{(i)})^T x$.

Logistic regression vs. ~~SVMs~~ SVMs.

n = no. of features, m = no. of training examples

→ If n is large (relative to m):

- Use ~~logistic regression~~ or SVM without a kernel ("linear kernel")

eg. $n=10,000$ $m=10$ or 1000
→ If n is small, m is intermediate:

- Use SVM with Gaussian kernel

eg. $n=1-1000$ $m=10-10,000$
→ If n is small, m is large:

They both perform almost similarly - Then SVM will be slow.

So create/add more features, then use logistic regression or SVM without a kernel

eg. $n=1-1000$ $m=50,000+$

→ Neural networks likely to work well for most of these settings, but may be slower to train.