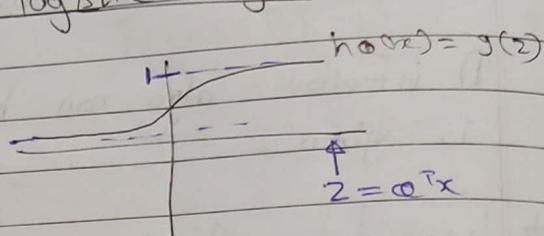


Week 7

C Optimization objective

★ Alternative view of logistic regn.

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



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

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

Cost of example : $-(y \log h_{\theta}(x) + (1-y) \log(1 - h_{\theta}(x)))$

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

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

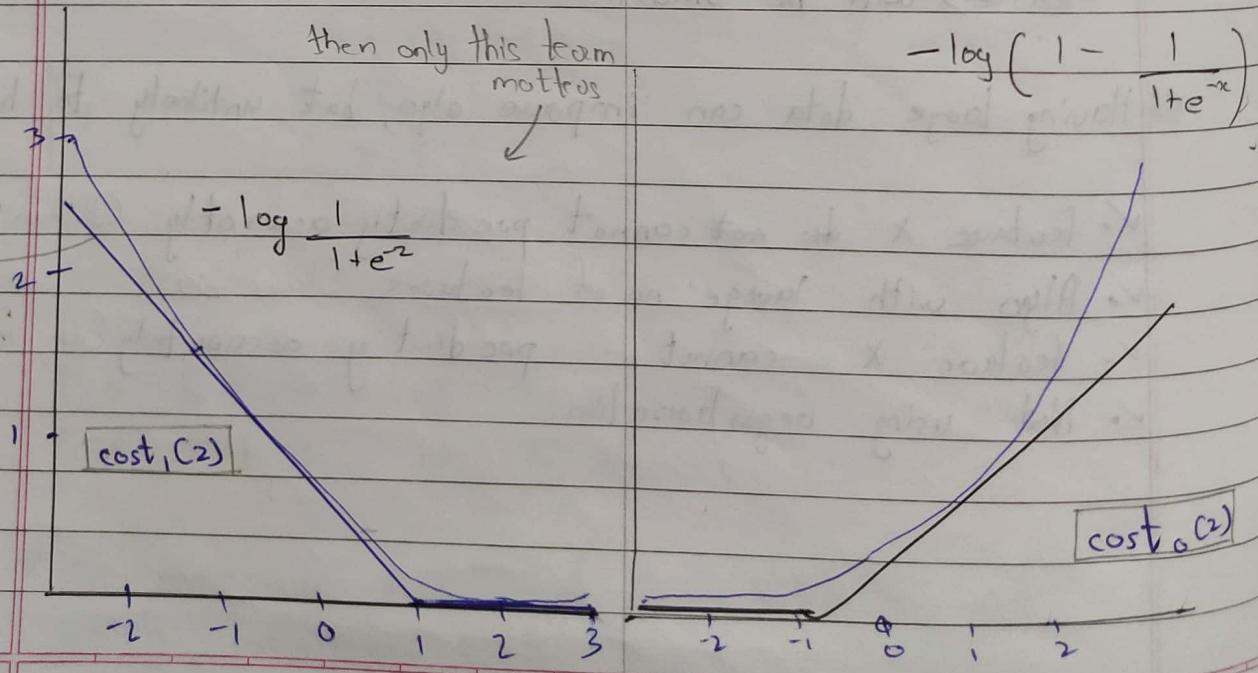
$$z = \theta^T x$$

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

then only this term matters

$$-\log \left(1 - \frac{1}{1 + e^{-z}}\right)$$

slope doesn't matter as off now.



$$\text{cost}_0(z)$$

* logistic regⁿ

$$\min_{\alpha} \frac{1}{m} \left[\sum_{i=1}^m y^{(i)} (-\log h_{\alpha}(x^{(i)})) + (1-y^{(i)}) \log (1-h_{\alpha}(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=0}^n \alpha_j^2$$

* Support vector machine.

$$\min_{\alpha} \frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \text{cost}_c(\alpha^T x^{(i)}) + (1-y^{(i)}) \text{cost}_s(\alpha^T x^{(i)}) \right] + \frac{\lambda}{2m} \sum_{j=0}^n \alpha_j^2$$

Q) $\min_u (u-5)^2 + 10 \leftarrow u = 5$

→ If constant is multiplied, it doesn't matter while finding min of a variable.

So we consider $\frac{1}{m}$

Also, we have crossed out λ , & added C

$A + \lambda B \leftarrow$ (logistic) $\lambda \rightarrow$ contains tradeoff
 $CA + B \leftarrow$ (sum) $C \rightarrow$

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

→ The 2 eqⁿ are not equal, but both will give same value (Optimal # value for the α).

* SVM hypothesis

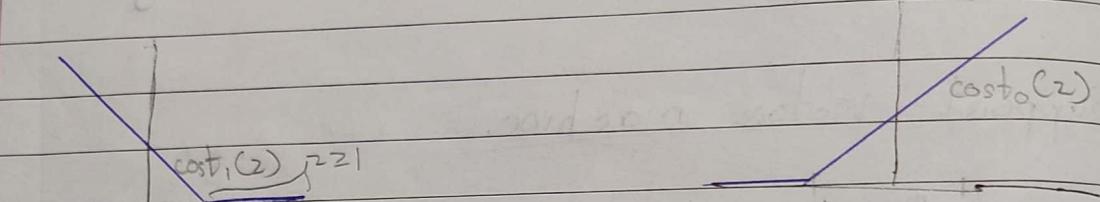
$$\min_{\alpha} C \sum_{i=1}^m [y^{(i)} \text{cost}_c(\alpha^T x^{(i)}) + (1-y^{(i)}) \text{cost}_s(\alpha^T x^{(i)})] + \frac{1}{2} \sum_{j=0}^n \alpha_j^2$$

Hypothesis →

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

-- @ SUM --

$$\min_{\mathbf{Q}} C \sum_{i=1}^m [y_i^{(i)} \text{cost}_1(\mathbf{Q}^T \mathbf{x}^{(i)}) + (1-y_i^{(i)}) \text{cost}_0(\mathbf{Q}^T \mathbf{x}^{(i)})] + \frac{1}{2} \sum_{i=1}^m \alpha_i$$



If $y=1$, we want $\mathbf{Q}^T \mathbf{x} \geq 1$ (not just ≥ 0)

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

What it take to make cost fun ↓.

$\mathbf{Q}^T \mathbf{x} \geq 1$ • For the e.g. i.e. $y=1$, $\text{cost}_1(z)=0$ only when $z \geq 1$

$\mathbf{Q}^T \mathbf{x} \leq -1$ • $y=0$, $\text{cost}_0(z)=0$, $z \leq 1$

→ This builds extra safety factor in SVM.

* SVM decision boundary.

$$\min_{\alpha} C \sum_{i=1}^m [y^{(i)} \text{cost} + \frac{1}{2} \sum_{j=1}^n \alpha_j^2]$$

If we set $C \rightarrow \infty$, we have to make this term 0(min)

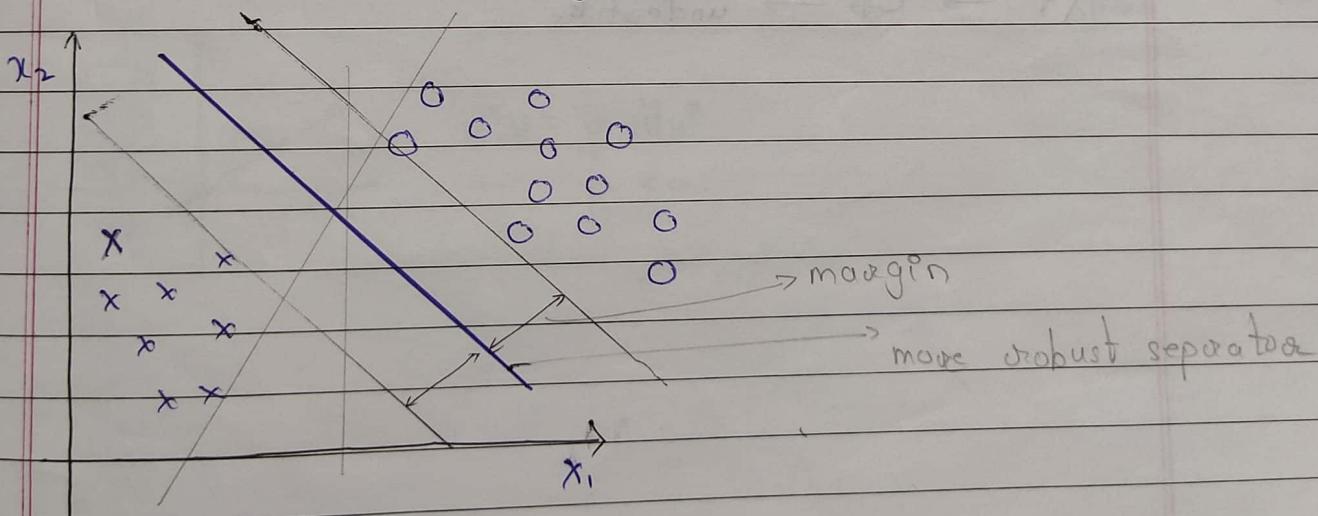
$$\text{when } y^{(i)} = 1: \alpha^T x \geq 1$$

$$\text{when } y^{(i)} = 0: \alpha^T x \leq -1$$

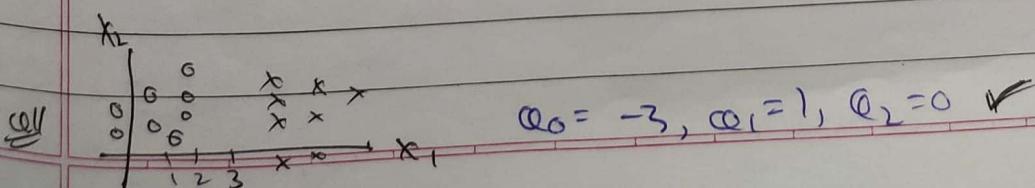
$$\begin{aligned} \min_{\alpha} & C \times 0 + \frac{1}{2} \sum_{j=1}^n \alpha_j^2 \\ \text{s.t. } & \alpha^T x^{(i)} \geq 1 \quad \text{if } y^{(i)} = 1 \\ & \alpha^T x^{(i)} \leq -1 \quad \text{if } y^{(i)} = 0. \end{aligned}$$

optimization

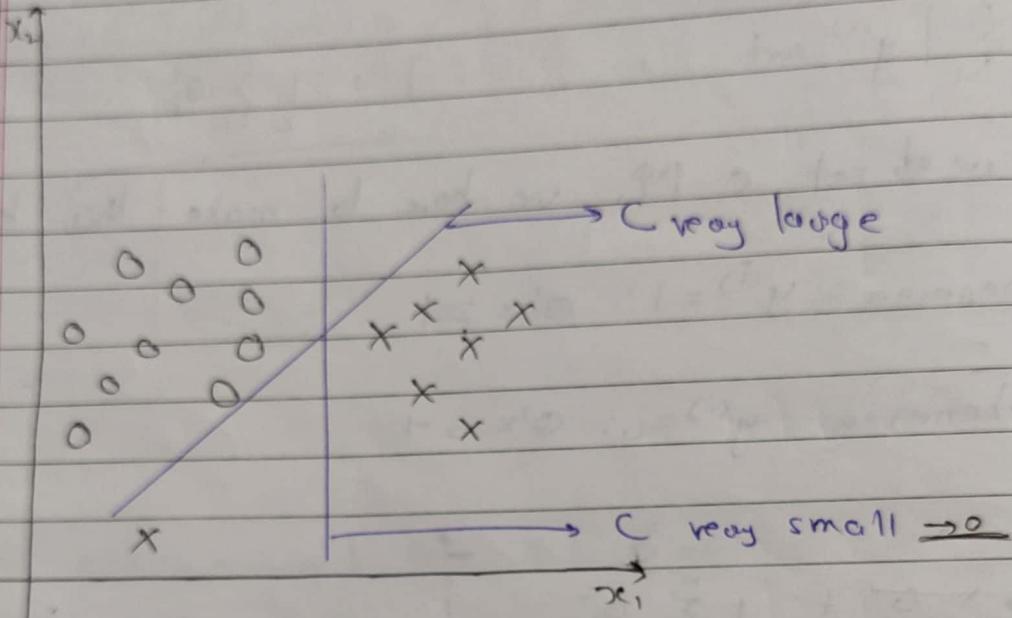
* SVM decision boundary: Linearly separable case.



Large margin classifier:- Consequence of optimization problem
(Good consequence)



* large margin classifier in presence of outliers.



C as if $\frac{1}{\lambda}$

$\lambda \downarrow \rightarrow C^+ \rightarrow$ overfit.

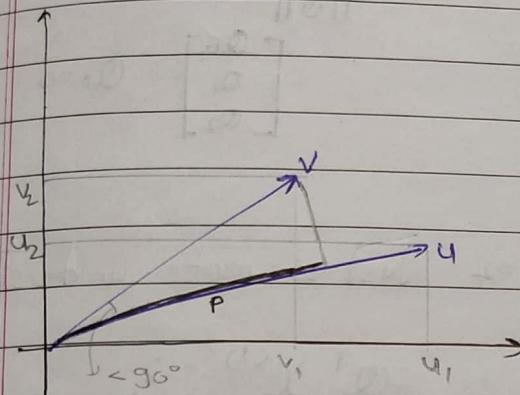
$\lambda^+ \rightarrow C^- \rightarrow$ underfit.

Here we are learning, how doing optimization,
we got margin classifiers.

Sadguru
Page No: 108
Date: 7/8/19

① Mathematics behind large Margin Classification (optimal)

★ Vector inner product.



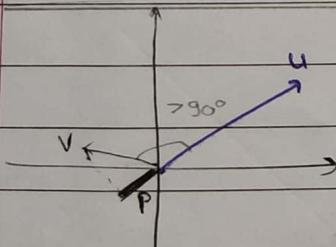
$$u = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \quad v = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

$$\begin{aligned} u^T v &= v^T u = [u_1 \ u_2] \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \\ &= [v_1 \ v_2] \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \end{aligned}$$

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

signed $u^T v = p \|u\| = v^T u = u_1 v_1 + u_2 v_2$

$p > 0, \quad p \in \mathbb{R}$



$$u^T v = p \|u\| \quad p < 0$$

Margin is constant
Margin is constant
Margin is constant

\star SVM Decision Boundary.

$$\min_{\alpha} \frac{1}{2} \sum_{j=1}^n \alpha_j^2 = \frac{1}{2} (\alpha_1^2 + \alpha_2^2) = \frac{1}{2} (\sqrt{\alpha_1^2 + \alpha_2^2})^2 = \frac{1}{2} \|\alpha\|^2$$

 $\|\alpha\|$

s.t. $\alpha^T x^{(i)} \geq 1$ if $y^{(i)} = 1$
 $\alpha^T x^{(i)} \leq -1$ if $y^{(i)} = 0$

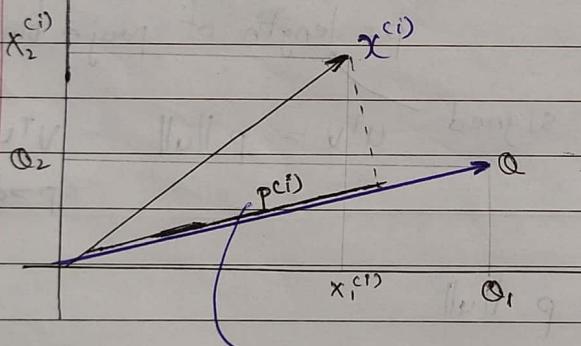
$$\begin{bmatrix} \alpha_0 \\ \alpha_1 \\ \alpha_2 \end{bmatrix}$$

$$\alpha_0 = 0$$

simplification: $\alpha_0 = 0$; $n=2$. Just to make understand better.

$$\alpha^T x^{(i)} = 2.$$

$\uparrow \quad \uparrow$
 $U^T V$



$$\alpha^T x^{(i)} = p^{(i)} \|\alpha\|$$

$$= \alpha_1 x_1^{(i)} + \alpha_2 x_2^{(i)}$$

value of p of i^{th} example.

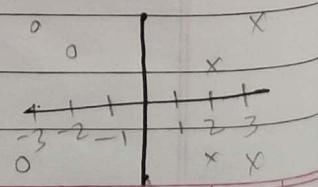
$\min_{\alpha} \frac{1}{2} \sum_{j=1}^n \alpha_j^2$, s.t. $\|\alpha\| p^{(i)} \geq 1$ if $y^{(i)} = 1$

$$\|\alpha\| p^{(i)} \leq -1 \text{ if } y^{(i)} = 0$$

At optimal value of α , value of $\|\alpha\| = ?$.

$\frac{1}{2}$
Ans

$$\frac{1}{2} \times 2 = 1$$



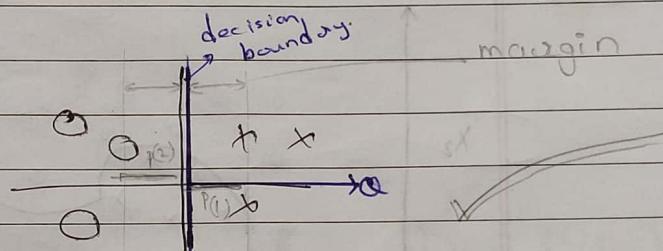
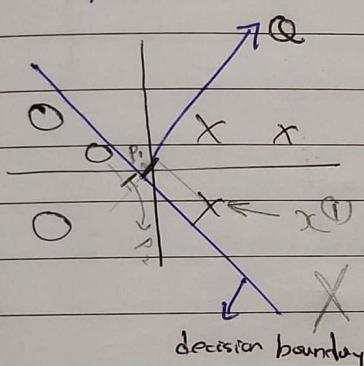
★ SVM Decision Boundary.

$$\min_{\alpha} \frac{1}{2} \sum_{j=1}^n \|\alpha_j\|^2 = \frac{1}{2} \|\alpha\|_2^2$$

$$\text{s.t. } p^{(i)} \cdot \|\alpha\|_2 \geq 1 \quad \text{if } y^{(i)} = 1 \\ p^{(i)} \cdot \|\alpha\|_2 \leq -1 \quad \text{if } y^{(i)} = -1 \quad \left. \right\} \text{C very large}$$

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

simplification: $\alpha_0 = 0$



$p^{(1)}, p^{(2)}$ are bigger, so $\|\alpha\|_2$ can be small now.

$$p^{(i)} \cdot \|\alpha\|_2 \geq 1$$

let say this is decision boundary.

- Not a good decision boundary because examples are near very close to it.
- SVM will not choose such DB. Magnitude of margin is $p^{(1)}, p^{(2)}$. So, in SVM we want large $p^{(1)}$

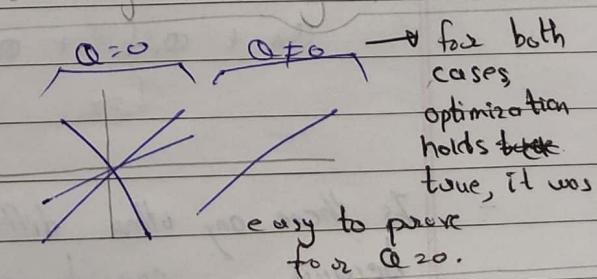
$$p^{(1)} > 0, p^{(2)} < 0$$

For optimization objective we want $p^{(1)} \cdot \|\alpha\|_2 \geq 1$

so $\|\alpha\|_2$ must be very large

$$p^{(2)} \cdot \|\alpha\|_2 \leq -1$$

small must be large

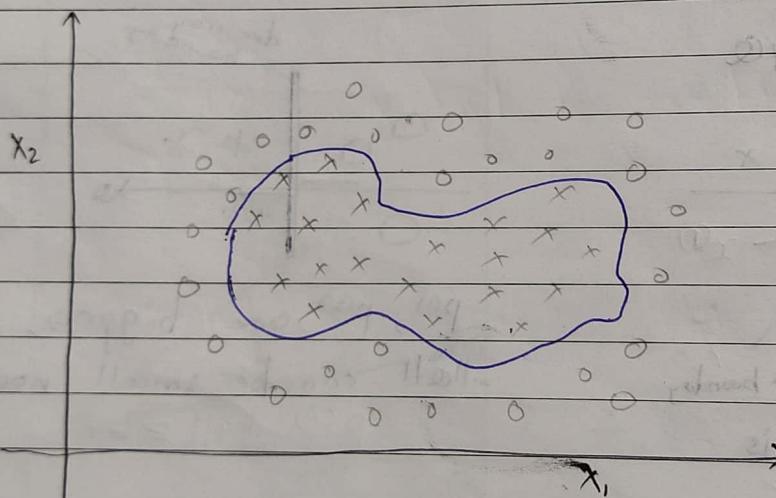


But in optimization objective, we have to make α small. To min we are doing optimization obj.

① Kernels |

★ Hence We are starting at adapting SVM, in order to develop complex non-linear classifiers.
For this we use technique called Kernels.

★ Non-linear Decision Boundary.



Predict $y=1$ if $\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 + \theta_5 x_2^2 \geq 0$

$$h_0(x) = \begin{cases} 1 & \text{if } \theta_0 + \theta_1 x_1 + \dots \geq 0 \\ 0 & \text{otherwise.} \end{cases}$$

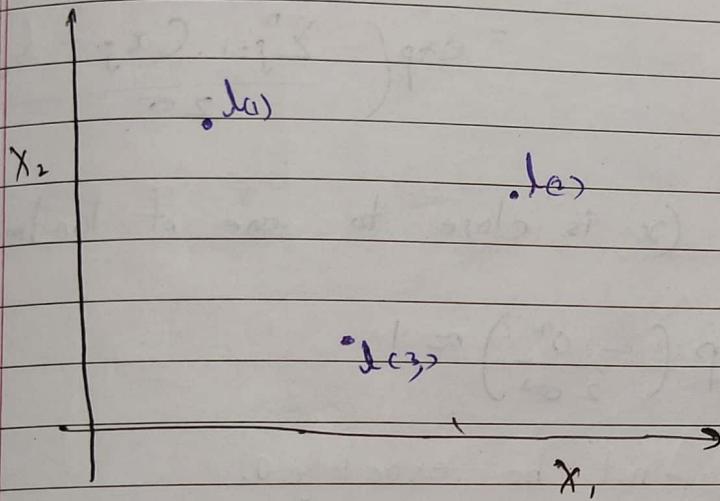
$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 + \dots$$

$$f_1 = x_1 \quad f_2 = x_2 \quad f_3 = x_1 x_2 \quad f_4 = x_1^2 \quad f_5 = x_2^2$$

Is there any other different/better choice of features f_1, f_2, f_3 because, considering higher degree term is expensive.

So, now we are going to find replacement of f_1, f_2, f_3 .

Hence we are going to define only 3 new features, unlike real world problem. Lineup 1, Lineup 2, Lineup 3.



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

We are defining new features as

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

$$= \frac{e^{-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}}}{\sigma}$$

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

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

Gaussian Kernel $\rightarrow k(x, l^{(i)})$

What kernels actually do
PTO

★ Kernels & similarity

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

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

If $x \approx l^{(i)}$: (x is close to one of landmark)

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

distance may not be exactly 0.

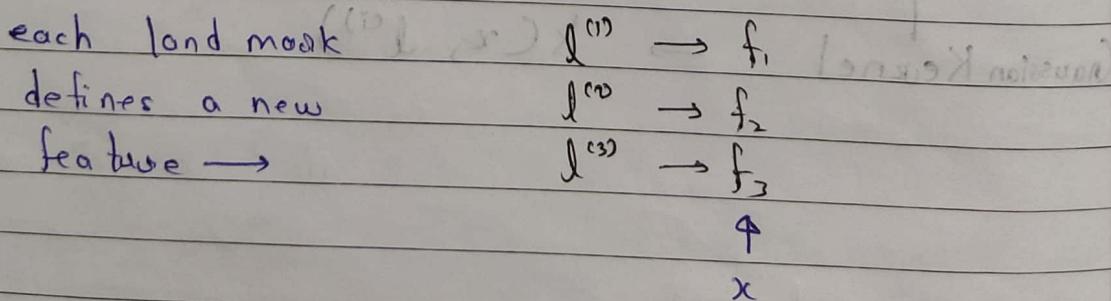
If x is far from $l^{(i)}$:

$$f_i = \exp\left(-\frac{(\text{large no})^2}{2\sigma^2}\right) \approx 0$$

This feature measures how similar x is from one of landmark

$f=0$, x far from landmark

$f=1$, x close to landmark

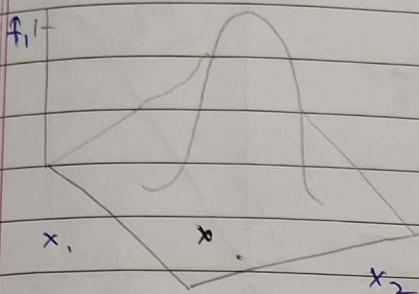


Given one example x , we have to compute 3 features f_1, f_2, f_3 by similarity b/w $l^{(i)}$ & $x^{(i)}$

e.g. $\mathbf{l}^{(1)} = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$, $f_1 = \exp\left(-\frac{\|\mathbf{x} - \mathbf{l}^{(1)}\|^2}{2\sigma^2}\right)$

 $\sigma^2 = 1$

Plotting this feature.



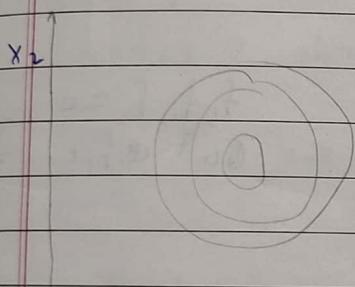
when $\mathbf{x} = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$ exactly,

the $f_1 = 1$ (max)

When \mathbf{x} moves further from $\begin{bmatrix} 3 \\ 5 \end{bmatrix}$

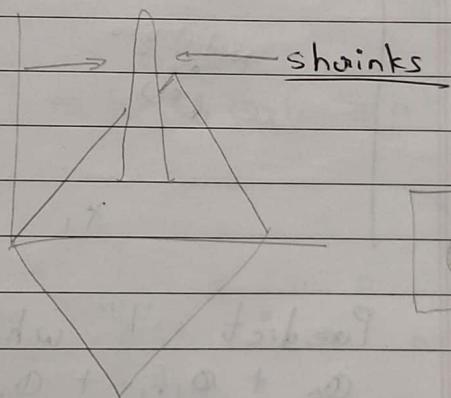
f_1 moves closer to 0.

So, f_1 is measured blw



If $\sigma^2 = 0.5$, everything other is same.

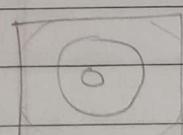
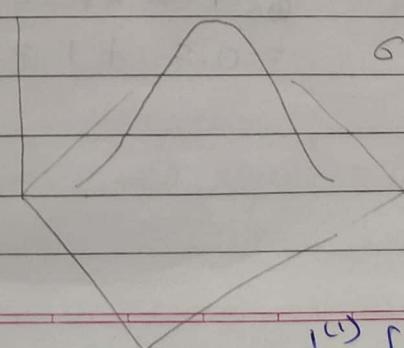
on true plot. $\xrightarrow{x_1}$



Moving from $\mathbf{x} = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$, f_1 falls

close to 0, more rapidly.

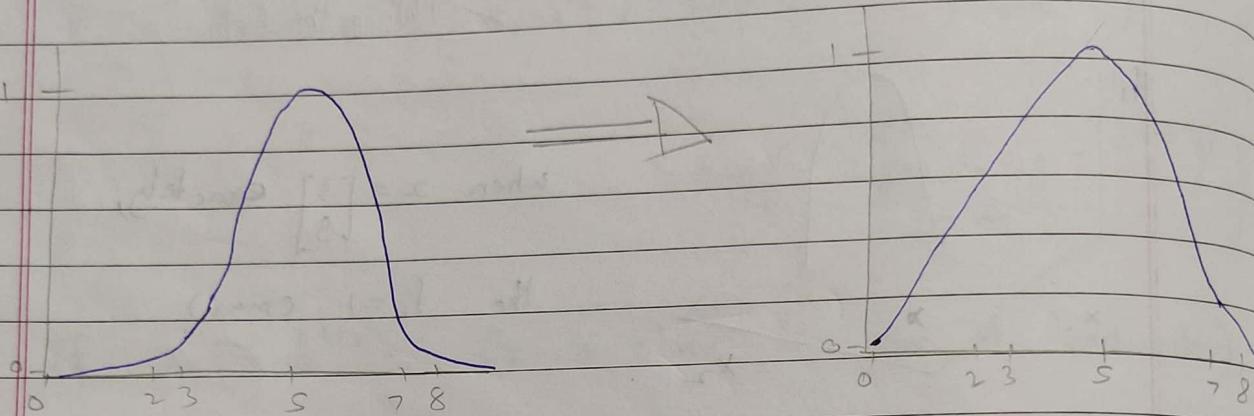
$$\sigma^2 = 3$$



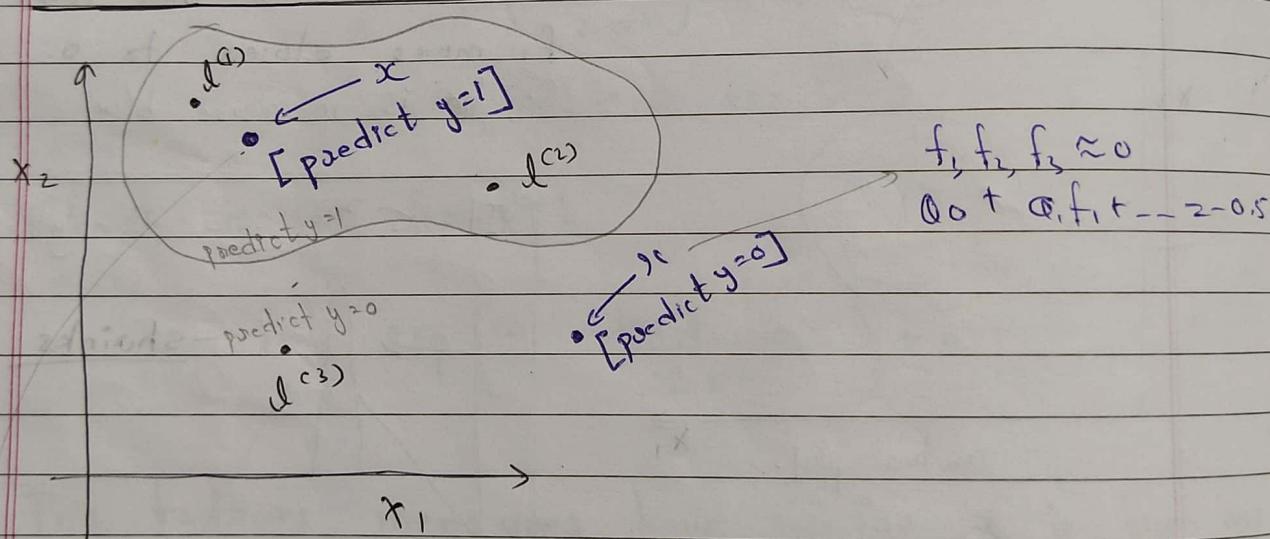
$\mathbf{l}^{(1)}$ falls slowly.

Q1

Consider a 1-D example with one feature x_1 . Suppose $\ell^{(1)} = 5$. To the right is plot of $f_i = \exp(-\gamma(x - \ell^{(i)})^2)$ when $\sigma^2=1$. Suppose we change $\sigma^2=4$.



(2)



Predict "1" when

$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$$

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

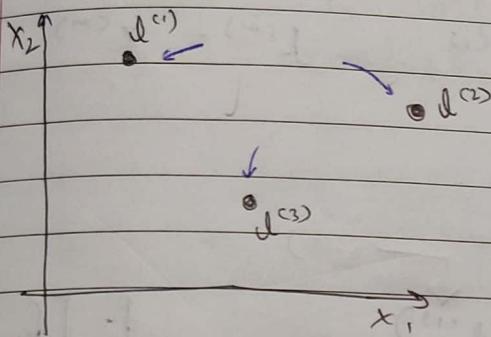
$$\theta_0 + \theta_1 x_1 + \theta_2 x_0 + \theta_3 x_0 \\ = -0.5 + 1 = 0.5 \geq 0 \quad \left| \begin{array}{l} f_1 \approx 1, f_2 \approx 0, f_3 \approx 0 \end{array} \right.$$

closed to $\ell^{(1)}$
far from $\ell^{(2)}, \ell^{(3)}$

Kernels II

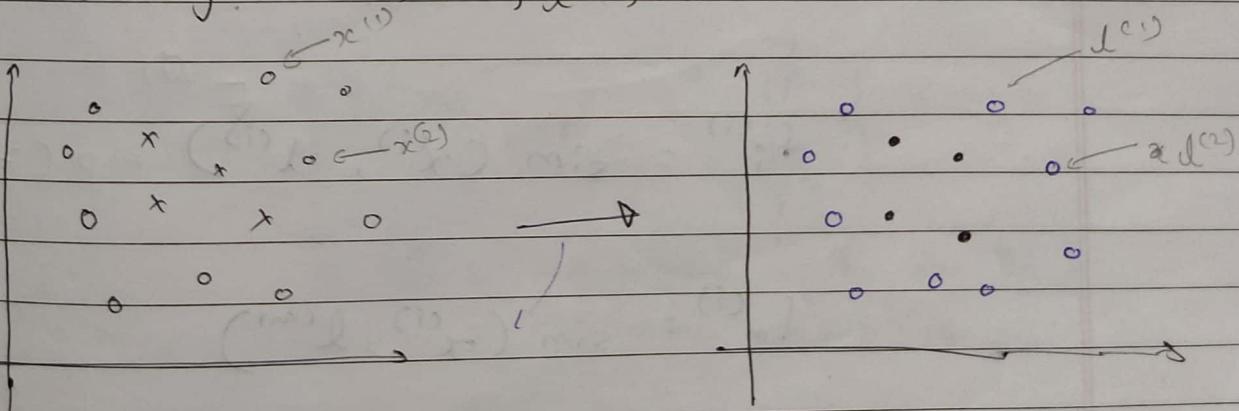
How to choose landmarks directly

Choosing the landmarks

Given x :

$$f_{i(j)} = \text{similarity}(x, l^{(j)})$$

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

Predict $y=1$ if $\alpha_0 + \alpha_1 f_1 + \alpha_2 f_2 + \alpha_3 f_3 \geq 0$ Where to get $l^{(1)}, l^{(2)}, l^{(3)}, \dots$?

We are going to put landmarks at same location as training examples, all.

Color of dots in right figure is not significant

★ SVM with Kernels.

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

choose $\lambda^{(1)} = x^{(1)}, \lambda^{(2)} = x^{(2)}, \dots, \lambda^{(m)} = x^{(m)}$

Given example x^i :

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

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

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

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

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

$$x^{(1)} \rightarrow f_2^{(i)} = \text{sim } (x^{(i)}, \lambda^{(2)})$$

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

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

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

This will represent

my one training example with new

features (f_0, f_1, f_m)

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

* SVM with kernels

. Hypothesis :

Given x , compute features $f \in \mathbb{R}^{M+1}$ $\alpha \in \mathbb{R}^{M+1}$
 predict " $y=1$ " if $\alpha^T f \geq 0$

$$\alpha_0 + \alpha_1 x_1 + \dots + \alpha_m x_m$$

→ This is how we do prediction using α .

How to calculate α .

. Training : (Here we are using new features)

$$\min_{\alpha} C \sum_{i=1}^n y^{(i)} \text{cost}_+(\alpha^T f^{(i)}) + (1-y^{(i)}) \text{cost}_-(\alpha^T f^{(i)}) + \frac{1}{2} \sum_{j=1}^n \alpha_j^2$$

here $n=m$ and $(\alpha_0, \alpha_1, \dots, \alpha_m)$

$$-\sum_j \alpha_j^2 = \alpha^T \alpha \quad \text{if } \alpha = \begin{bmatrix} \alpha_0 \\ \vdots \\ \alpha_m \end{bmatrix} \text{ ignoring } \alpha_0$$

$\alpha^T M \alpha$... sometimes written like this to
optimize computations.

* SVM parameters

$C = \left(\frac{1}{\lambda}\right)$

Prone to overfitting. Large C : lower bias, high variance. (small λ)
Prone to underfitting. Small C : Higher bias, low variance. (large λ)

σ^2

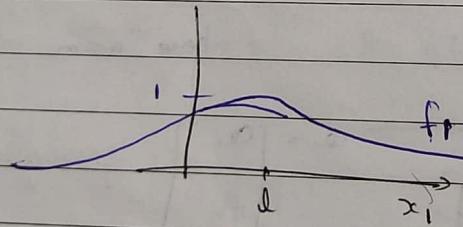
appears in
gaussian
kernel.

• Large σ^2 :

Features fit vary more smoothly.
Higher bias, lower variance.

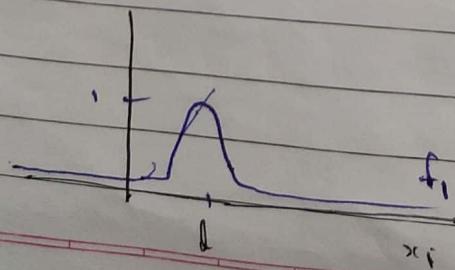
$$\exp\left(-\frac{\|x_i - \bar{x}\|^2}{2\sigma^2}\right) \quad \text{-- If we have only 1 feature}$$

on and landmark there only



• Small σ^2 :

Features fit vary less smoothly.
Lower bias, higher variance.



Q1

If SVM overfits what should we do.

- ✓ Decrease C
- ✓ Increase ϵ^2