

→ Each landmark gives us the features in our hypothesis :

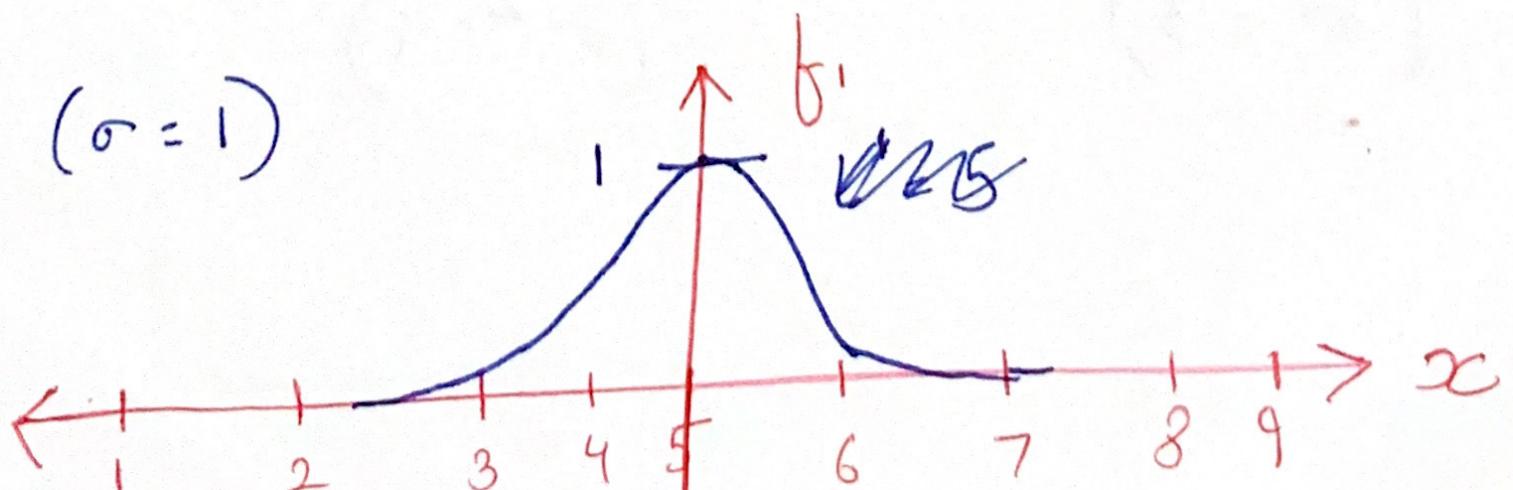
$$\begin{array}{l} L^1 \rightarrow f_1 \\ L^2 \rightarrow f_2 \\ L^3 \rightarrow f_3 \end{array}$$

→ If we plot f_i vs the input vector, we'll observe that as x moves away from L^i , then the feature takes on values close to 0. So, this measures how close x is to this landmark.

Eg: Suppose L_{vec} is a 1-D vector, i.e. $L_{\text{vec}} = [5]$.

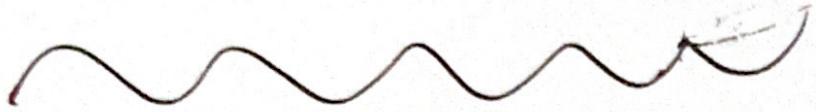
$$\therefore f_1 = \exp \left(-\frac{(x-5)^2}{2\sigma^2} \right)$$

→ Sketch: ($\sigma=1$)



As x moves away from 5, then f_1 takes values closer to 0.

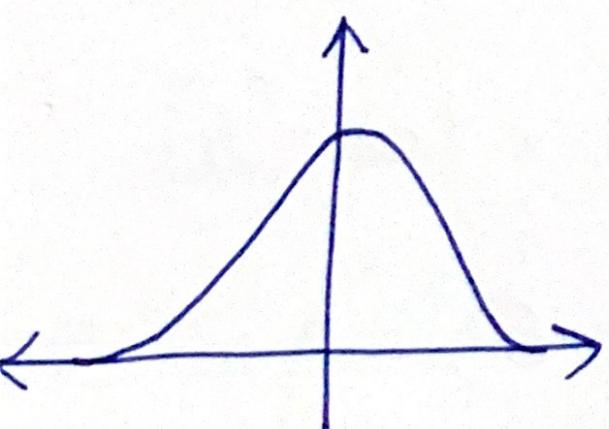
\Rightarrow What does σ^2 do?



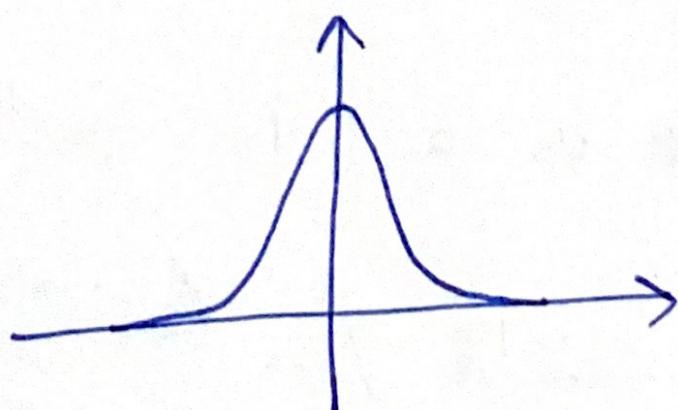
~~$\rightarrow \sigma^2$~~ $\sigma^2 = \text{Variance} \rightarrow \text{Measurement of the spread}$
between numbers in a dataset, i.e. it measures
how far each number in the set
from the mean, & it also under therefore
from every other number in the set.

\rightarrow Variance also defines the steepness of the
rise around the landmark. low variance \Rightarrow more
steepness

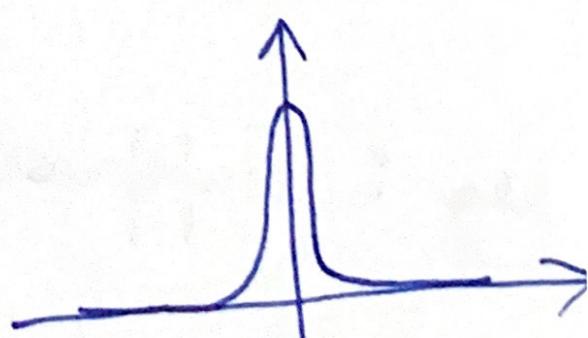
e.g.: for the example on the prev. page,



$$\underline{(\sigma^2 = 3)}$$



$$\underline{(\sigma^2 = 1)}$$



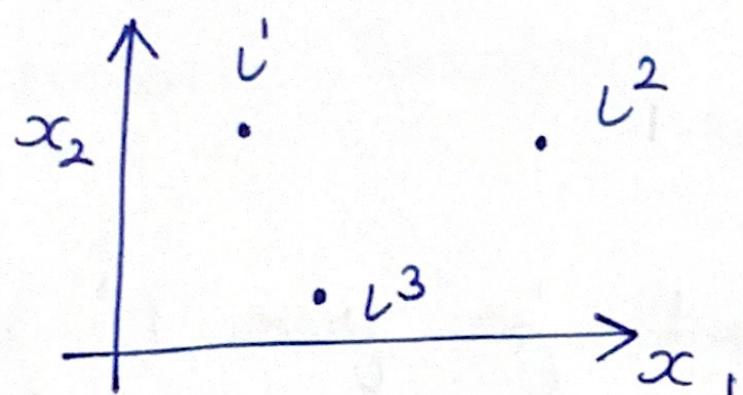
$$\underline{(\sigma^2 = 0.5)}$$

⇒ Given this definition, what kind of hypotheses can we learn?



→ Supp. we have some training examples & we predict $y = 1$ when $\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$

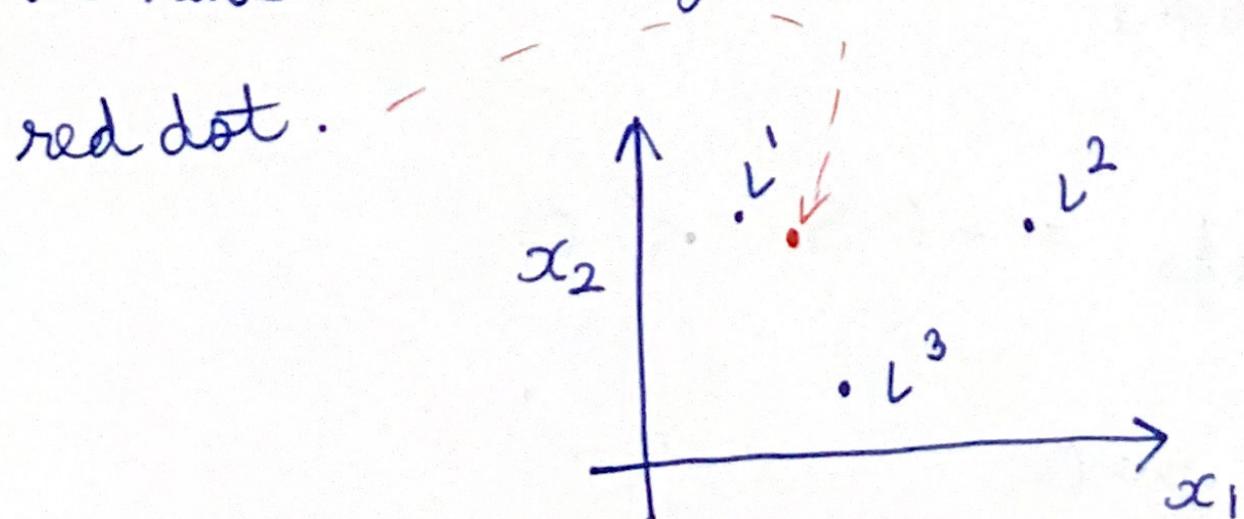
$$\delta \text{supp. } \alpha_0 = 0$$



* Supp.) we've magically obtained the θ vector & for our example, we obtain:

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix} = \begin{bmatrix} -0.5 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$

→ Given our placement of our landmarks, what happens if we have a training example like the one denoted by the red dot.



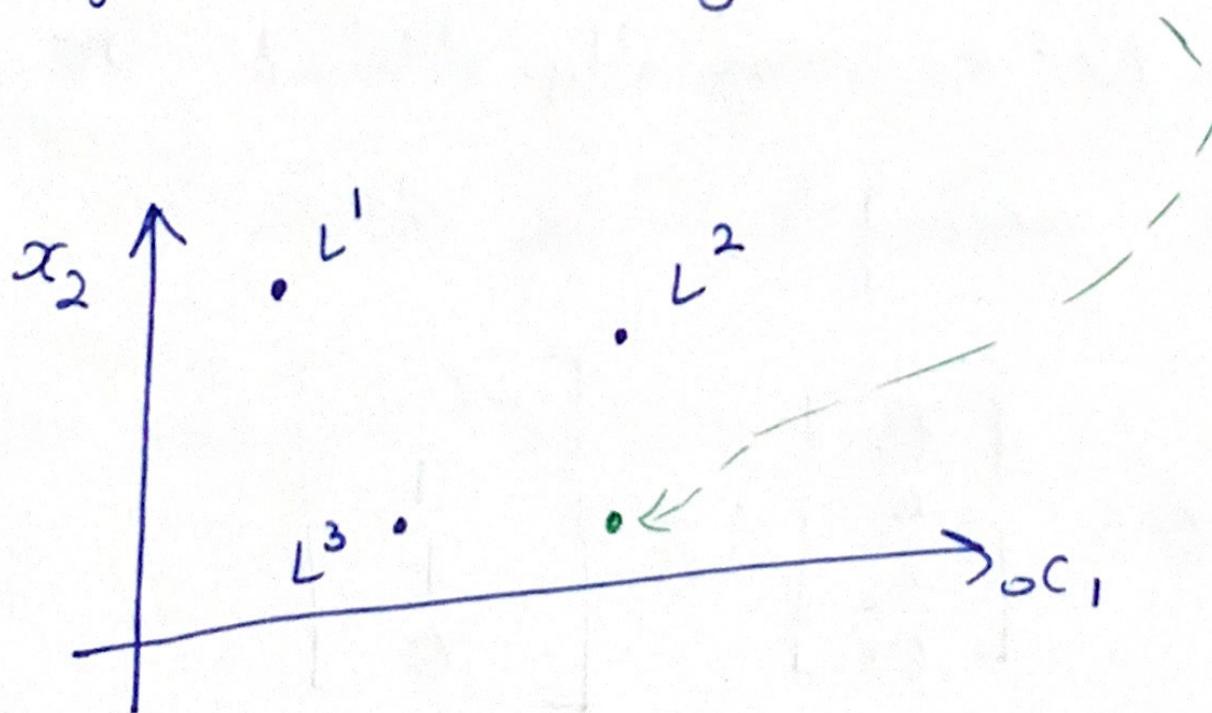
We know $f_1 \approx 1$ as point is near L' , but $f_2 \& f_3$ will be close to 0.

Coming to our hypothesis

$$\theta_0 f_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 = -0.5 + 1 + 0 + 0 = 0.5 \geq 0$$

$$\Rightarrow y = 1$$

→ If we had a point far from L' & L^2 , but near L^3 , we get 0 (denoted by green dot)



Consider Using sin. arguments, $f_1 = f_2 \approx 0$; $f_3 \approx 1$

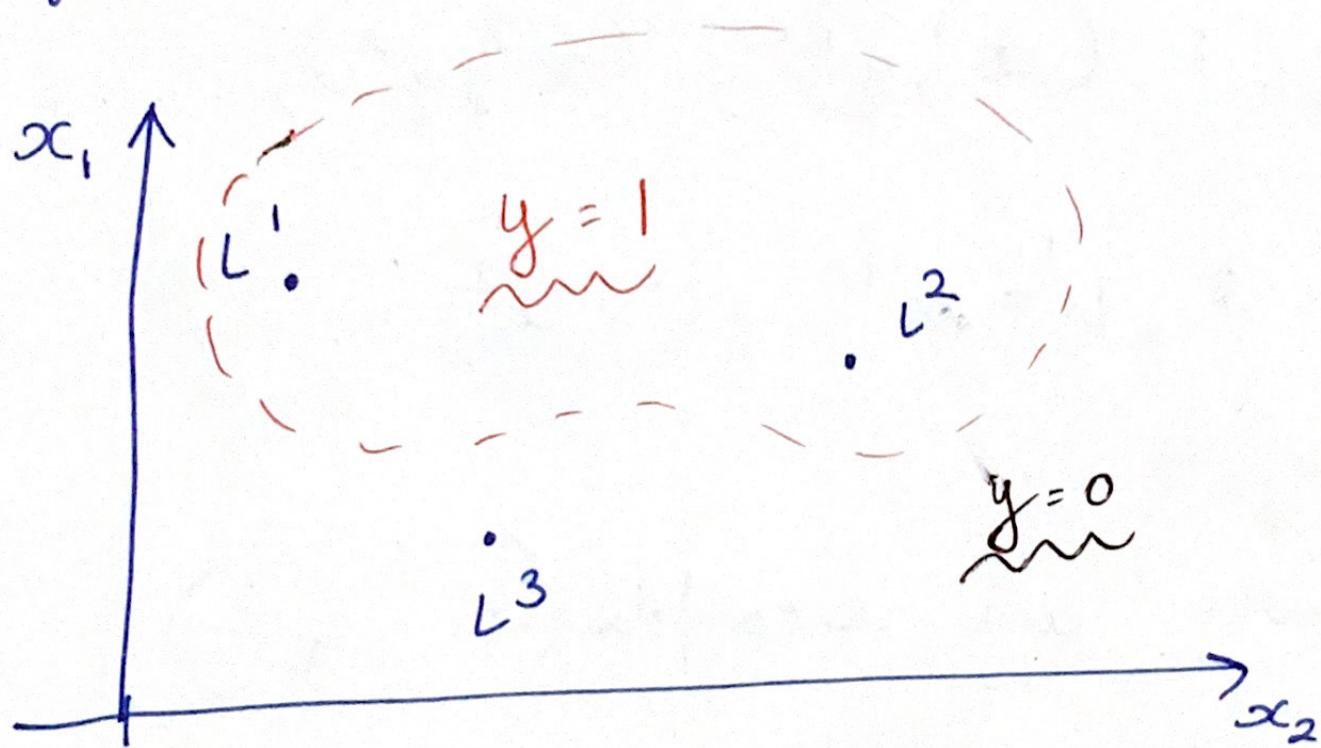
(Computing hypothesis)

$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 = -0.5 < 0$$

$$\Rightarrow y = 0$$

→ Can we see that for points near L^1 ,
our hypothesis predicts $y=1$, & it predicts
 $y=0$ for points near L^3 .

Which means we've created a
non linear decision boundary which looks something
like this :



→ This was cool!! But how do we choose the landmarks ??

\Rightarrow Choosing the landmarks

One way to get the landmarks is to put them in the exact same locations as all the training examples.

This gave

\rightarrow Concretely,

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 an example (x^i, y^i) ,

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

$$f_1^i = k(x^i, l^1)$$

$$f_2^i = k(x^i, l^2)$$

$$f_m^i = k(x^i, l^m) = \exp(0) = 1$$

$$f_m^i = k(x^i, l^m)$$

* Using this feature vector f instead of x
makes using a SVM much more easier.

\Rightarrow SVM with kernels

\rightarrow Hypothesis : Given x_i , compute features $f \in \mathbb{R}^{m+1}$.
Training sample

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

\rightarrow Training :

$$\min_{\alpha} \left\{ C \sum_{y=1}^m y^i \text{Cost}_1(\alpha^T f^i) + (1-y^i) \text{Cost}_0(\alpha^T f^i) + \frac{1}{2} \sum_{j=1}^{m=m} \alpha_j^2 \right\}$$

As ^{we've created} no. of features = no. of training samples

$$\therefore m = m$$

\rightarrow One particular subtlety : If we ignore α_0 then the
foll. is true $\sum_{j=1}^n \alpha_j^2 = \alpha^T \alpha$. What many
implementations rewrite this as $\alpha^T M \alpha$.

* This makes it efficient for large training example.

→ Using kernels to generate f_i is not exclusive to SVMs & may also be applied to algorithms like logistic regression.

However, it is seen that kernels combined with SVMs is much faster than with other algorithms.

⇒ SVM parameters

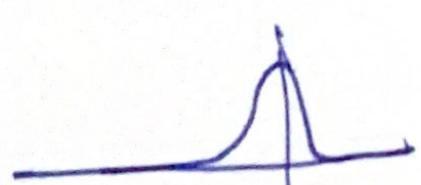
→ $C = \frac{1}{\lambda}$:

* Large $C \Rightarrow$ Small λ : ~~Large~~ Lower bias & higher variance

* Small $C \Rightarrow$ Large λ : Higher bias & lower variance

→ σ^2 : * Large $\sigma^2 \Rightarrow$ features very more smoothly \Rightarrow Higher bias, lower variance

* Small $\sigma^2 \Rightarrow$ feature very less smoothly



\Rightarrow Using a SVM

- Use SVM software package (eg : liblinear, libsvm) to solve for parameters θ . 
- Need to specify:
 - * Choice of parameter C
 - * Choice of kernel (similarity function)
 - No kernel or "Linear" kernel:
 - Predict $y = 1$ if $\theta^T x \geq 0 \rightarrow \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_m x_m \geq 0$
 - Used when n is large, m is small; As ~~if~~ you have a really small & a high dimensional dataset, using a complex model might result in overfitting. So, using a simple linear model makes more sense in this case.
 - Gaussian Kernel : Predict $y = 1$ if $\theta^T f \geq 0$
 - $f_i = \exp\left(-\frac{\|x - l^i\|^2}{2\sigma^2}\right)$, where $l^i = x^i$
 - Need to choose σ^2
 - Used when n is small & m is large.

⇒ Points to remember if using Gaussian Kernel

→ Some SVM packages will expect you to define the kernel function explicitly.

e.g.: For Gaussian kernel,

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

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

return

→ Make sure you perform feature scaling before using a Gaussian kernel.

If you don't, other features with large value will dominate the f value.

e.g.: Let $v = x - l$

$$\text{Then } \|x - l\|^2 = \|v\|^2 = v_1^2 + v_2^2 + \dots + v_n^2$$

$$= (v_1 - 0)^2 + (v_2 - 0)^2 + \dots + (v_n - 0)^2$$

$$= (x_1 - l^1)^2 + (x_2 - l^2)^2 + \dots + (x_n - l^n)^2$$

If $x_i \in [1000, 2000]$ & $x_2 \in [2, 10]$, then $(x_i - l^i)^2$ will result in large values

⇒ Other choices of kernel

→ Note: Not all similarity functions are valid kernels. They must satisfy "Mercer's Theorem" in order to make sure that the SVM packages' optimizations run correctly, & do not diverge.

→ Polynomial Kernel :

$$\text{Similarity } (x, l) = k(x, l) = (x^T l + k)^D$$

; $k \rightarrow \text{constant}$; $D \rightarrow \text{constant degree of polynomial}$

$$\text{eg: } k(x, l) = (x^T l)^2, (x^T l + 1)^2, (x^T l)^3$$

* Uses the idea that if x & l are similar, then their inner product $(x^T l)$ will be large.

* Usually used when x & l are non-negative

* Not used that often & usually performs worse than Gaussian kernel.

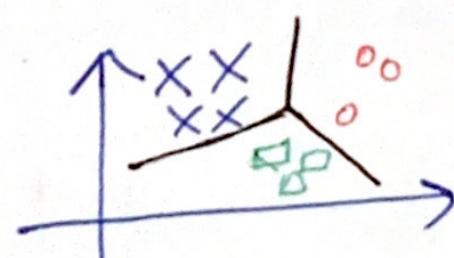
→ More exotic examples:

* String kernel → Used if input data consists of texts or strings.

* chi-square kernel

* histogram intersection kernel

→ Note: While choosing b/w the diff. types of kernel, and choosing parameters such as C, σ^2 etc. choose whatever performs best on the cross validation data.



⇒ Multi-class classification

→ Many SVM libraries have this feature built-in.

→ Otherwise, use one vs all method, discussed in logistic regression.

* Train K SVMs, one to distinguish $y=i$ from the rest $\forall i = [1, 2, \dots, K]$ & get $\alpha^1, \alpha^2, \dots, \alpha^K$

* Pick class i with largest value of $(\alpha^i)^T x$.

\Rightarrow Logistic regression vs SVM

\rightarrow Let m denote no. of features s.t. $x \in \mathbb{R}^{m+1}$ (x_0 included)

$\rightarrow m =$ No. of training examples

\rightarrow If n is large relative to m eg: $n \geq m$, $m=10000$

$$m \in [10, 1000]$$

* Use ~~Kern~~ logistic regression, or SVM without a kernel

(linear kernel).

* Reason: You probably don't have enough ~~data~~ ^{examples} to fit need hypothesis
a very complicated non-linear function.

\rightarrow If n is small, m is intermediate eg: $m \in [1, 1000]$

$$m \in [10, 10000]$$

* Use SVM with Gaussian Kernel, as we have enough examples that we may need a complex non-linear hypothesis

\rightarrow If n is small, m is large: eg: $m \in [1, 1000]$

$$m \in [50000, \infty)$$

* Create/add more features & use logistic regression or SVM without kernel.

\rightarrow A neural net is likely to work well for any of these situations but may be slower to train.