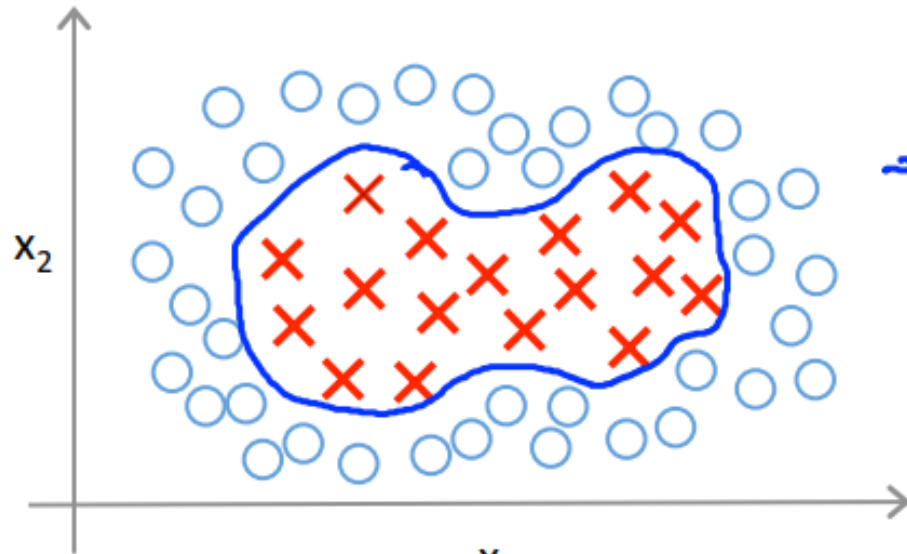


Week 7:

Support Vector Machines

Kernels

Non-linear Decision Boundary



Predict $y = 1$ if

$$\rightarrow \theta_0 + \theta_1 \underline{x_1} + \theta_2 \underline{x_2} + \theta_3 \underline{x_1 x_2} \\ + \theta_4 \underline{x_1^2} + \theta_5 \underline{x_2^2} + \dots \geq 0$$

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

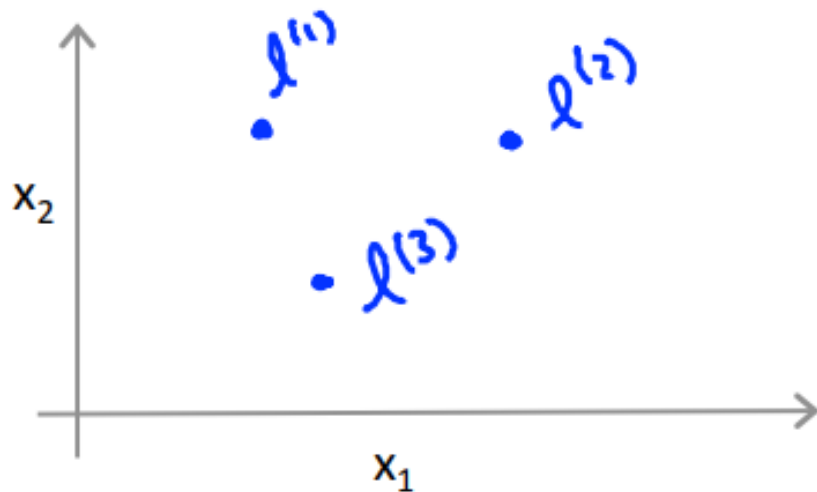
$$\rightarrow \theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 + \dots$$

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

High order polynomial features \rightarrow new features
to reduce computational expense

Measure similarity using Euclidean distance squared

Kernel



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

Given x :

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

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

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

kernel (Gaussian kernels) $k(x, l^{(i)})$

Handwritten annotations: $\|w\|$ with an arrow pointing to the denominator $2\sigma^2$ in the first equation. $\|x - l^{(1)}\|^2$ with an arrow pointing to the numerator of the first equation. $\|x - l^{(i)}\|^2$ with an arrow pointing to the numerator of the second equation.

Kernels and Similarity

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

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

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

$$\begin{array}{lcl} l^{(1)} & \rightarrow & f_1 \\ l^{(2)} & \rightarrow & f_2 \\ l^{(3)} & \rightarrow & f_3 \\ \uparrow & & \uparrow \\ & & \times \end{array}$$

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

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

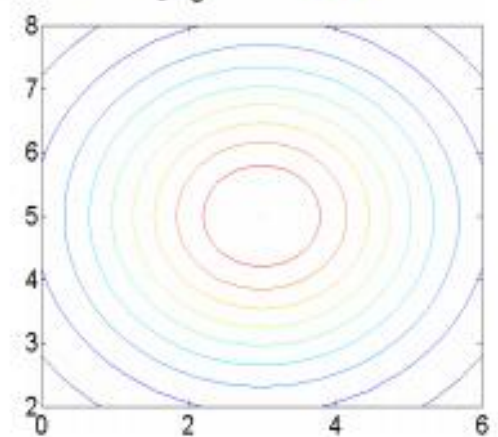
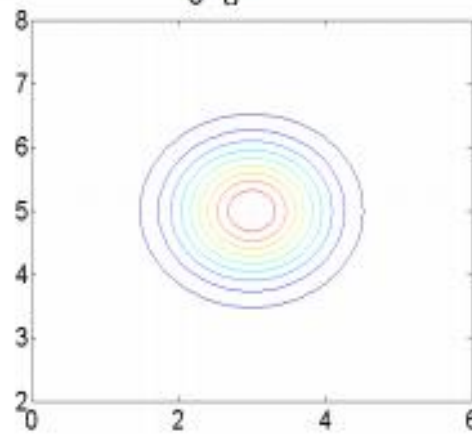
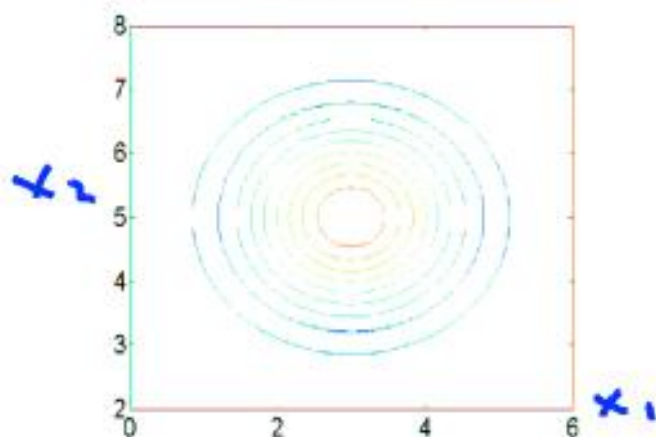
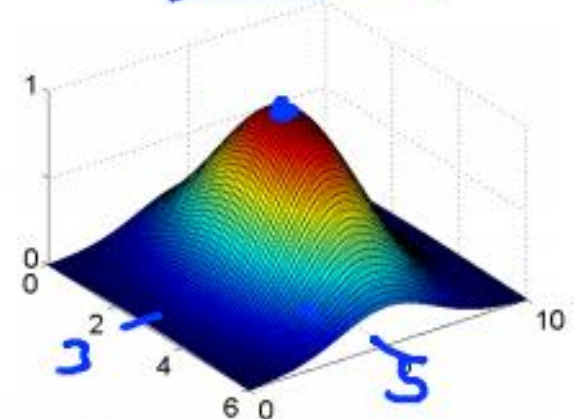
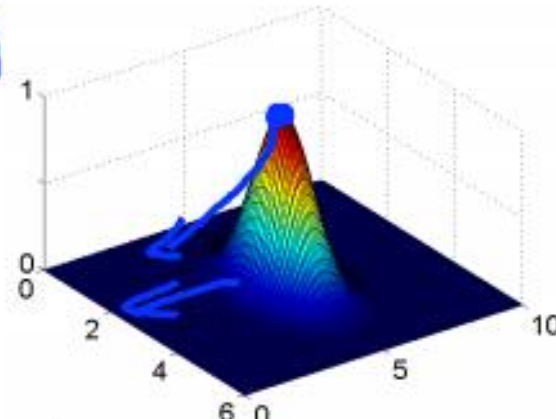
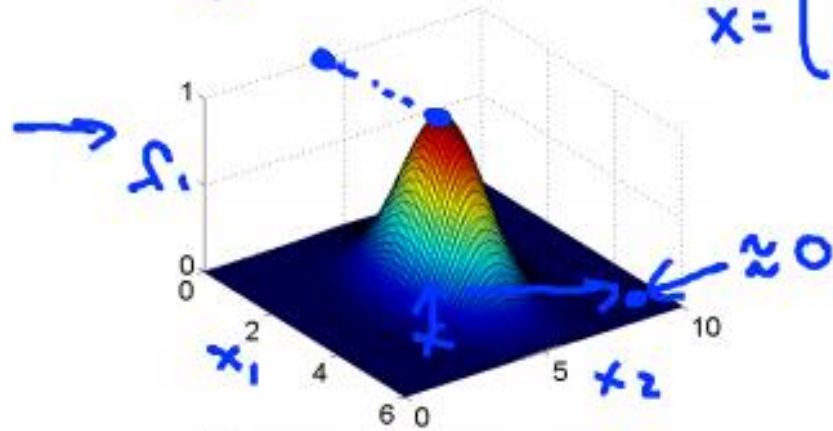
$$\rightarrow l^{(1)} = \begin{bmatrix} 3 \\ 5 \end{bmatrix} \quad f_1 = \exp \left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2} \right)$$

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

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

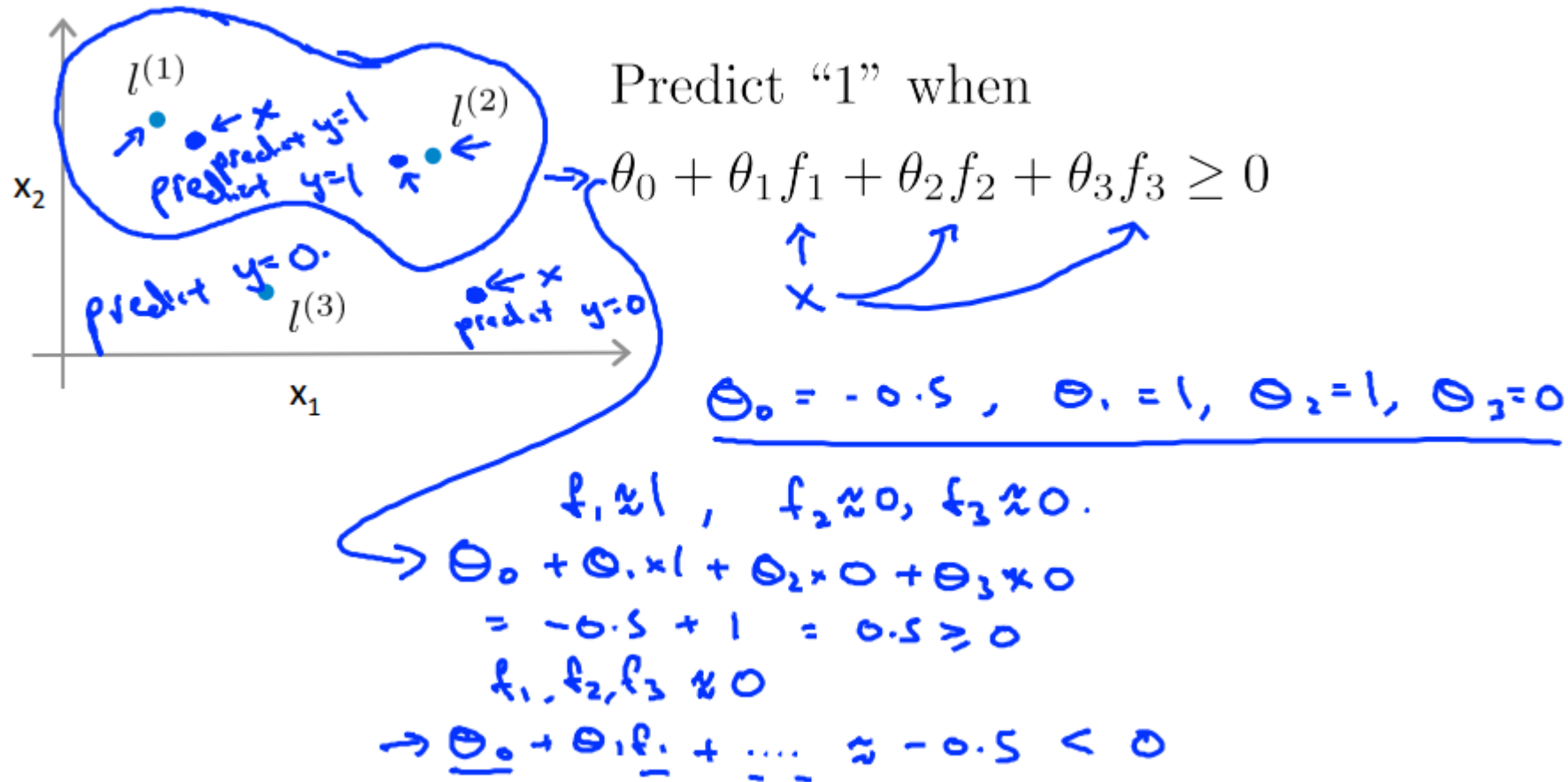
$$\sigma^2 = 0.5$$

$$\sigma^2 = 3$$



Gaussian kernel graph vary by sigma squared

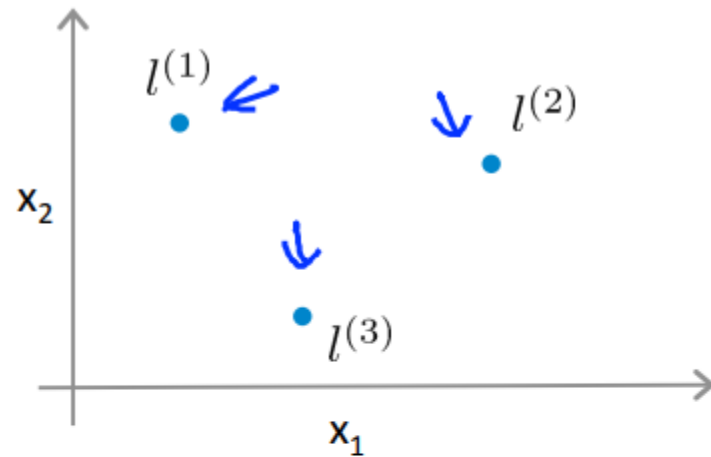
If training example is given



Inner boundary -> predict as 1

Outer boundary -> predict as 0

Choosing landmarks

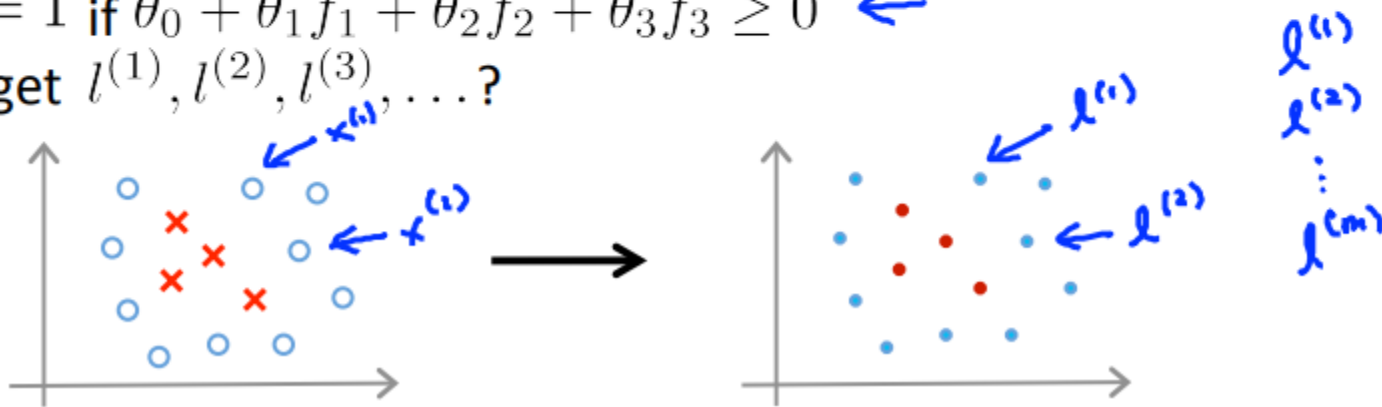


Given x :

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

Predict $y = 1$ if $\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$

Where to get $l^{(1)}, l^{(2)}, l^{(3)}, \dots$?



Set landmarks exact same location as training data

Mapping training data with similarity function

Given example \underline{x} :

$$\begin{aligned} \rightarrow f_1 &= \text{similarity}(x, l^{(1)}) \\ \rightarrow f_2 &= \text{similarity}(x, l^{(2)}) \\ &\vdots \end{aligned}$$

$$f = \begin{bmatrix} f_0 \\ f_1 \\ f_2 \\ \vdots \\ f_m \end{bmatrix} \quad f_0 = 1$$

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

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

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

$$f^{(i)} = \begin{bmatrix} f_0^{(i)} \\ f_1^{(i)} \\ f_2^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix} \quad f_0^{(i)} = 1$$

Andrew N

M-dimensional x vector \rightarrow (M+1)-dimensional feature vector

* $f_0 = 1$

How to get value of theta

SVM with Kernels

Hypothesis: Given \underline{x} , compute features $\underline{f} \in \mathbb{R}^{m+1}$ $\Theta \in \mathbb{R}^{n+1}$
→ Predict "y=1" if $\underline{\theta}^T \underline{f} \geq 0$
 $\Theta_0 f_0 + \Theta_1 f_1 + \dots + \Theta_m f_m$

Training:

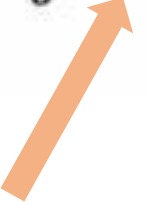
$$\rightarrow \min_{\theta} C \sum_{i=1}^m y^{(i)} \text{cost}_1(\underbrace{\theta^T f^{(i)}}_{\cancel{\Theta^T x^{(i)}}}) + (1 - y^{(i)}) \text{cost}_0(\underbrace{\theta^T f^{(i)}}_{\Theta^T f^{(i)}}) + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$


To get value of theta and minimize it, we can use cost-function

Why not use kernel in other algorithm?

- We 'can' use kernel in other algorithms such as logistic regression
- But computational tricks doesn't generalize to other algorithms
- As a result, computation will be very slow and expensive

Bias variance trade-off

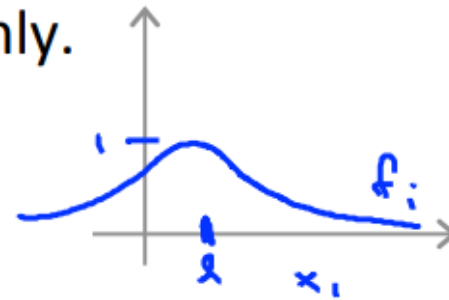
$$\min_{\theta} C \sum_{i=1}^m [y^{(i)} (\text{Cost}_1(\theta^T f^{(i)})) + (1 - y^{(i)}) \text{Cost}_0(\theta^T f^{(i)})] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$


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


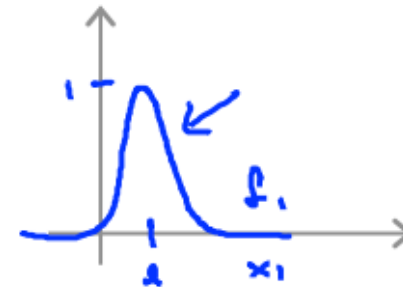
Bias variance trade-off

$C (= \frac{1}{\lambda})$. \rightarrow Large C: Lower bias, high variance. (small λ)
 \rightarrow Small C: Higher bias, low variance. (large λ)

σ^2 Large σ^2 : Features f_i vary more smoothly.
 \rightarrow Higher bias, lower variance.
 $\exp\left(-\frac{\|x - \mu^{(i)}\|^2}{2\sigma^2}\right)$



Small σ^2 : Features f_i vary less smoothly.
Lower bias, higher variance.



Large C = small lambda = tendency to overfit

Small C = large lambda = tendency to underfit

Small sigma = tendency to overfit

Large sigma = tendency to underfit