# Machine Learning Part III
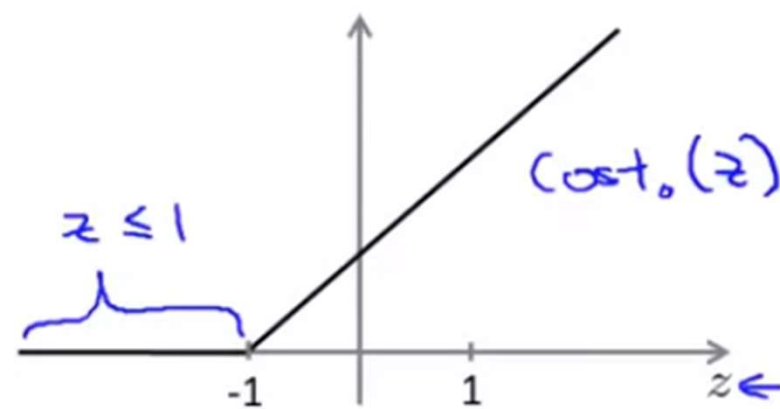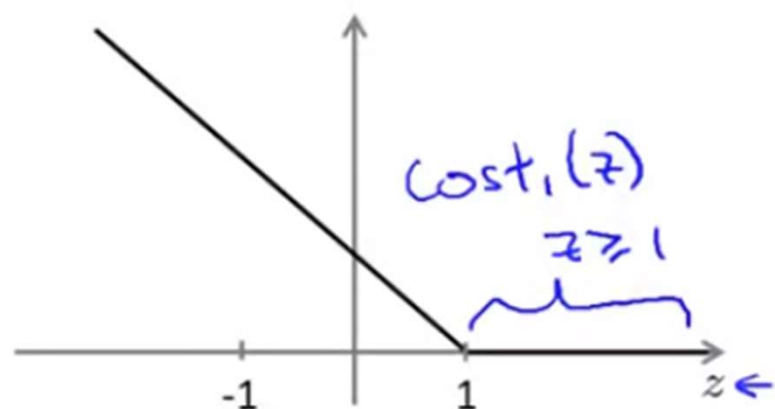
# SVM

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



$cost_1(z)$

$z \geq 1$

$z \leq 1$

$cost_0(z)$

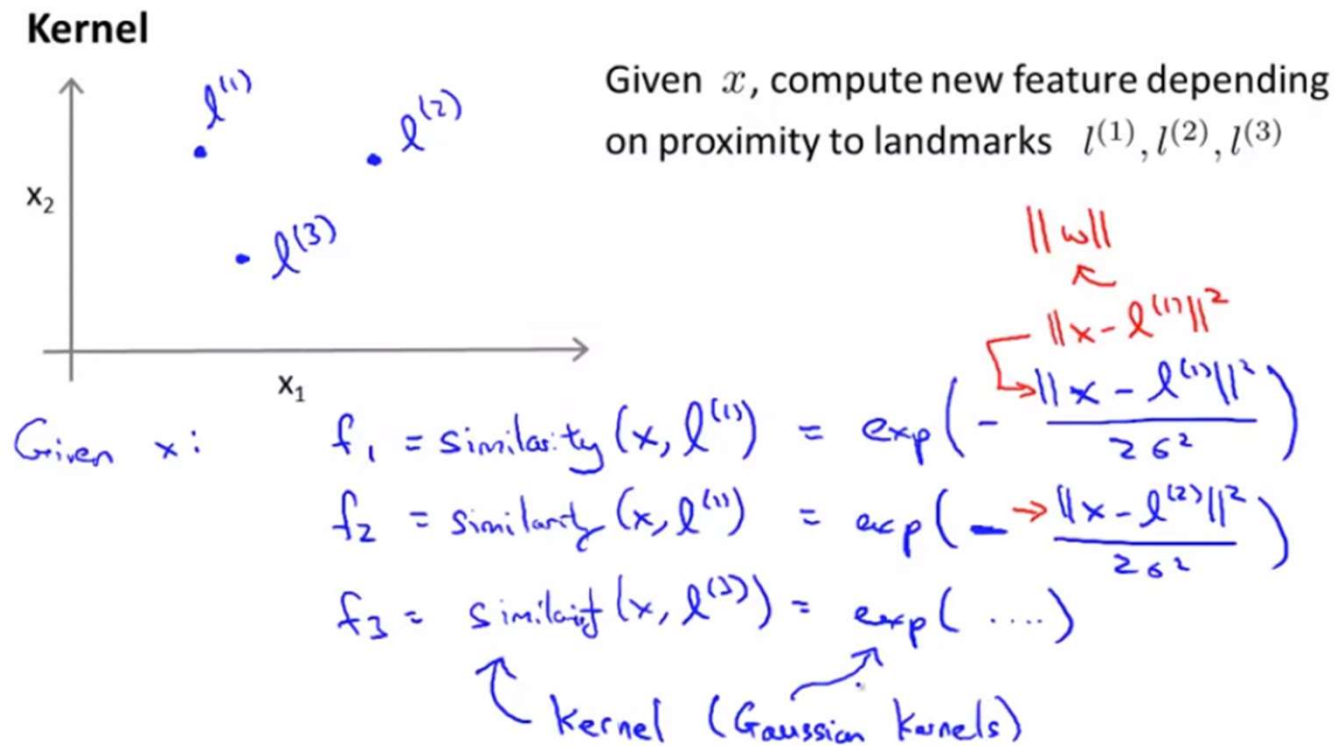$\rightarrow$ If $y = 1$, we want $\theta^T x \geq 1$ (not just $\geq 0$)

$\theta^T x \geq \cancel{0} \; 1$

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

$\theta^T x \leq \cancel{0} -1$

# Kernels (similarity function)

- Use **landmarks** to choose features instead of all combinations of polynomial

**Kernel**

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

$$\|w\|$$

$$- \|x - l^{(1)}\|^2$$

Given $x$:

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

$$f_2 = \text{similarity}(x, l^{(1)}) = \exp\left(- \frac{\|x - l^{(2)}\|^2}{2\sigma^2}\right)$$

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

kernel (Gaussian kernels)

**Kernels and Similarity**

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

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

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

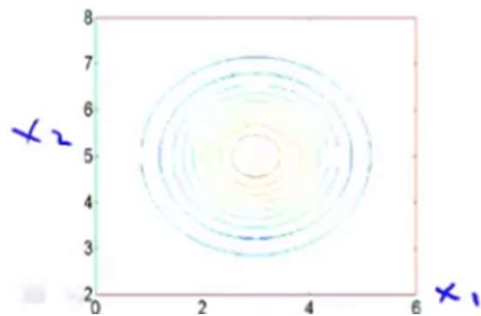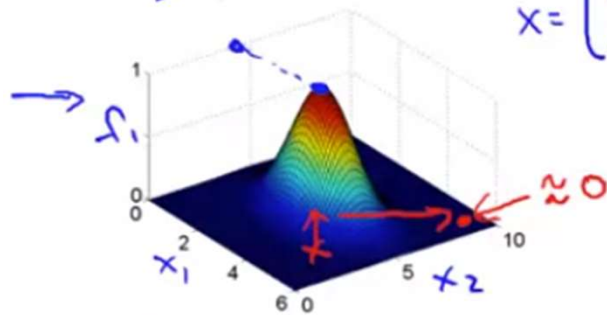If $x$ if far from $l^{(1)}$ :

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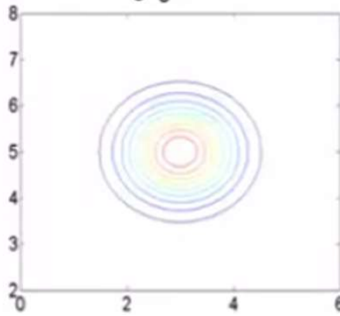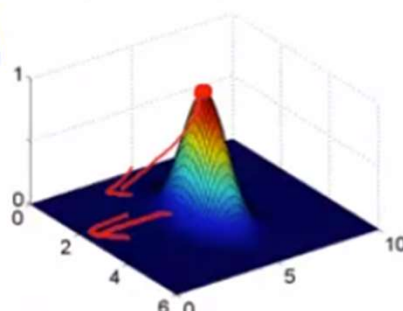
# Effect of sigma

**Example:**

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

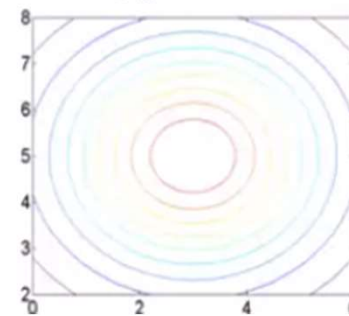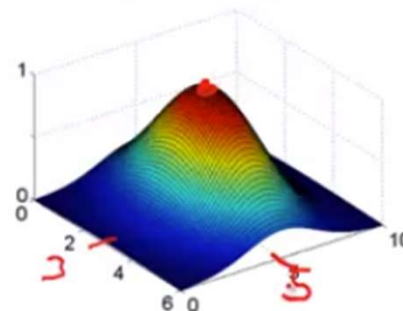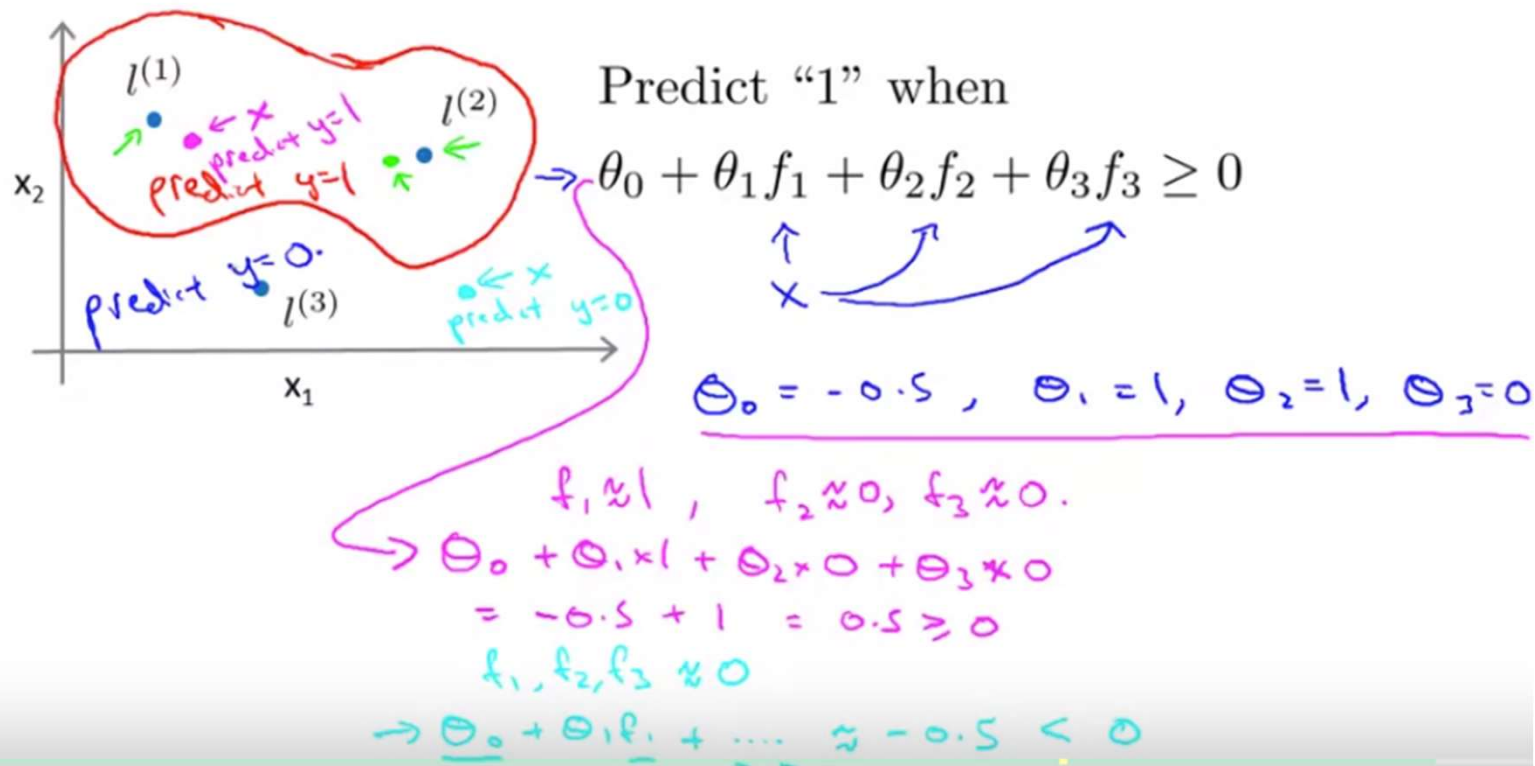$\sigma^2 = 1$  $\quad\quad x = \begin{bmatrix} 3 \\ 5 \end{bmatrix} \quad\quad$ $\sigma^2 = 0.5$  $\quad\quad\quad\quad$ $\sigma^2 = 3$

# Non-linear classifier

$l^{(1)}$

$l^{(2)}$

predict y=1

predict y=1

$x_2$

predict y=0.

$l^{(3)}$

predict y=0

$x_1$

Predict "1" when

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

$x$

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

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

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

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

$f_1, f_2, f_3 \approx 0$

$\theta_0 + \theta_1 f_1 + \ldots \approx -0.5 < 0$

## Choosing the landmarks

$l^{(1)}$ $l^{(2)}$

$x_2$

$l^{(3)}$

$x_1$

Given $x$:

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

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

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

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

$x^{(i)}$

$x^{(1)}$

$l^{(1)}$

$l^{(2)}$

$l^{(1)}$
$l^{(2)}$
$\vdots$
$l^{(m)}$

each data sample is a landmark!
each data sample becomes a feature!

# SVM with Kernels

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

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

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

$n = m$

Training:

$\rightarrow \quad \min_\theta C \sum_{i=1}^{m} y^{(i)} cost_1(\theta^T f^{(i)}) + (1 - y^{(i)}) cost_0(\theta^T f^{(i)}) + \frac{1}{2} \sum_{j=1}^{m} \theta_j^2$

$\theta^T x^{(i)} \quad \theta^T f^{(i)}$

A large *C* parameter tells the SVM to try to classify all the examples correctly
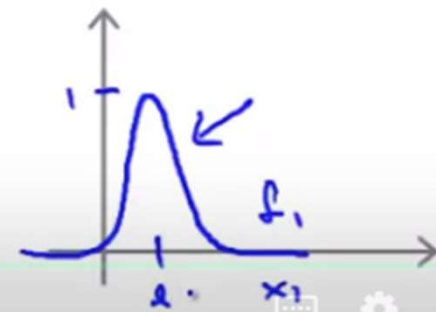
## SVM parameters:

$C\ ( = \frac{1}{\lambda}\ ).$ → Large C: Lower bias, high variance.    (small $\lambda$)

→ Small C: Higher bias, low variance.    (large $\lambda$)

$\sigma^2$    Large $\sigma^2$: Features $f_i$ vary more smoothly.

→ Higher bias, lower variance.

$$\exp\left(-\frac{\|x - \ell^{(i)}\|^2}{2\sigma^2}\right)$$

$f_i$

Small $\sigma^2$: Features $f_i$ vary less smoothly.

Lower bias, higher variance.

$f_i$

# Using SVM

Use SVM software package (e.g. liblinear, libsvm, ...) to solve for parameters $\theta$.

Need to specify:

→ Choice of parameter C.

Choice of kernel (similarity function):

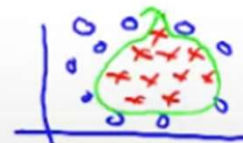E.g. No kernel ("linear kernel")

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

$$\theta_0 + \theta_1 x_1 + \cdots + \theta_n x_n \geq 0$$

→ $n$ large, $m$ small

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

Gaussian kernel:

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

$x \in \mathbb{R}^n$, $n$ small

and/or $m$ large

Need to choose $\sigma^2$.

**Kernel (similarity) functions:**

$x^{(i)}$

$\ell^{(j)} = x^{(j)}$

$f_i$

```
function f = kernel(x1,x2)
```

$$f = \exp\left(-\frac{\|\, x1 - x2 \,\|^2}{2\sigma^2}\right)$$

```
return
```

$x \to$

$f_1$
$f_2$
$\vdots$
$f_m$

→ Note: Do perform feature scaling before using the Gaussian kernel.

$x \in \mathbb{R}^n$

$\to \|x - \ell\|^2$

$v = x - \ell$

$\|v\|^2 = v_1^2 + v_2^2 + \cdots + v_n^2$

$= (x_1 - \ell_1)^2 + (x_2 - \ell_2)^2 + \cdots + (x_n - \ell_n)^2$

1000 feet²

1-5 bedrooms

## Other choices of kernel

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

Many off-the-shelf kernels available:
- Polynomial kernel: $k(x, l) = (x^T l)^2$ to

$$(x^T l)^3, \quad (x^T l + 1)^{\circled{3}}, \quad (x^T l + 5)^{\circled{4}}$$

$$(x^T l + \text{constant})^{\text{degree}}$$

- More esoteric: String kernel, chi-square kernel, histogram intersection kernel, ...

## Logistic regression vs. SVMs

$n =$ number of features ($x \in \mathbb{R}^{n+1}$), $m =$ number of training examples

→ If $n$ is large (relative to $m$): (e.g. $n \geq m$, $n = 10,000$ , $m = 10 \cdots 1000$)
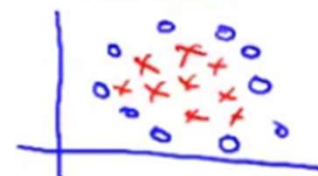
↪ Use logistic regression, or SVM without a kernel ("linear kernel")

→ If $n$ is small, $m$ is intermediate: ($n = 1-1000$ , $m = 10 - 10,000$) ←

↪ Use SVM with Gaussian kernel

If $n$ is small, $m$ is large: ($n = 1-1000$, $m = 50,000+$)

↪ Create/add more features, then use logistic regression or SVM without a kernel

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

## 2.6 Optional (ungraded) exercise: Build your own dataset

- TODO