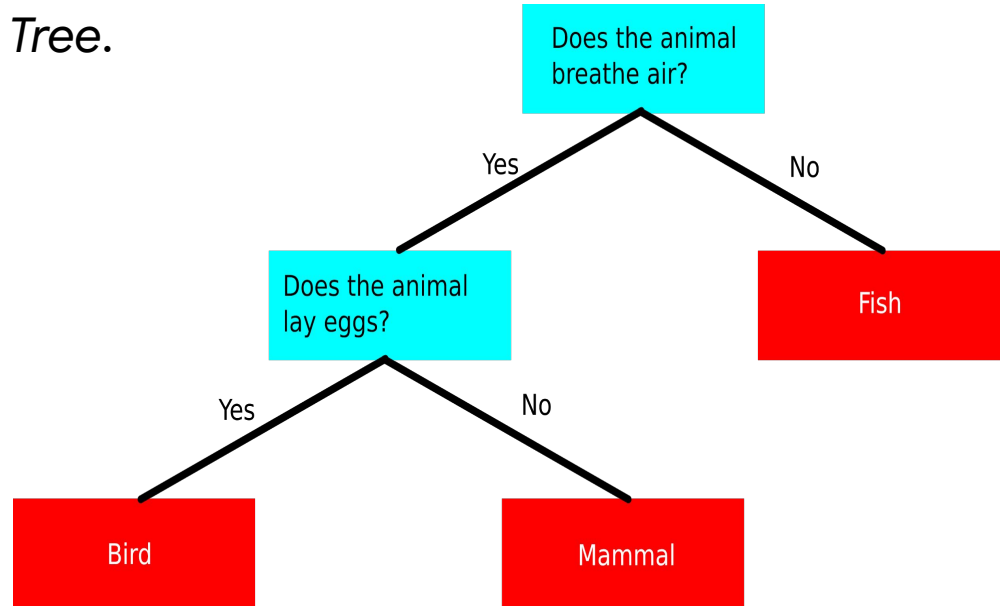What is a Decision Tree?

Predictive Model in tree form that maps items to its target value, starting from root to leaf is known as a *Decision Tree*.

Types of Decision Tree:

- **Classification Tree**
- **Regression Tree**

```
                      Does the animal
                      breathe air?
                   Yes            No
         Does the animal              Fish
         lay eggs?
      Yes          No
   Bird              Mammal
```

# Entropy & Gini

Metrics for above mentioned algorithm can be:

- **Gini Impurity(GI):** Measure of how a randomly chosen element is incorrectly labelled if it was randomly labelled according to its subset distribution.
- **Information Entropy(IE):** Number of bits required in encoding the given data.

$$I_G = 1 - \sum_{j=1}^{c} p_j^2$$

$p_j$: proportion of the samples that belongs to class c for a particular node

**GI**

$$I_H = - \sum_{j=1}^{c} p_j log_2(p_j)$$

$p_j$: proportion of the samples that belongs to class c for a particular node.

**IE**

# Information Gain

- **Information Gain** is used to decide which features to split at each step while building the tree.
- Our objective is to keep the tree as small as possible, thus we choose the split that results in the purest daughter nodes.
- The information value *represents the expected amount of information that would be needed to specify whether a new instance should be classified **yes** or **no**, given that the example reached that node.*

$$\overbrace{IG(T,a)}^{\text{Information Gain}} = \overbrace{H(T)}^{\text{Entropy (parent)}} - \overbrace{H(T|a)}^{\text{Weighted Sum of Entropy (Children)}}$$

$$= -\sum_{i=1}^{J} p_i \log_2 p_i - \sum_{a} p(a) \sum_{i=1}^{J} -\Pr(i|a) \log_2 \Pr(i|a)$$

# How does it work?

*Look on the blackboard*

# Case Study

– Let's consider the example where our objective is to figure out whether we should go out to play **Tennis** or not.

| Day | Outlook | Temperature | Humidity | Wind | PlayTennis |
|-----|---------|-------------|----------|------|------------|
| D1 | Sunny | Hot | High | Weak | No |
| D2 | Sunny | Hot | High | Strong | No |
| D3 | Overcast | Hot | High | Weak | Yes |
| D4 | Rain | Mild | High | Weak | Yes |
| D5 | Rain | Cool | Normal | Weak | Yes |
| D6 | Rain | Cool | Normal | Strong | No |
| D7 | Overcast | Cool | Normal | Strong | Yes |
| D8 | Sunny | Mild | High | Weak | No |
| D9 | Sunny | Cool | Normal | Weak | Yes |
| D10 | Rain | Mild | Normal | Weak | Yes |
| D11 | Sunny | Mild | Normal | Strong | Yes |
| D12 | Overcast | Mild | High | Strong | Yes |
| D13 | Overcast | Hot | Normal | Weak | Yes |
| D14 | Rain | Mild | High | Strong | No |

# Case Study

- First we determine the **information gain** for each candidate attribute(i.e, Outlook, Temperature, Humidity and Wind) and then select the one with **highest** information gain.
- **Attribute Outlook:**

  Outlook_Rain = {2+, 3-}

  Outlook_Sunny = {3+, 2-}

  Outlook_Overcast = {4+, 0-}

Gain(S,Outlook) = Entropy(S) — 5/14 Entropy(S_Outlook_Rain) — 5/14 Entropy(S_Outlook_Sunny) — 4/14 Entropy(S_Outlook_Overcast)

Gain(S,Outlook) = **0.246**

# Case Study

- First we determine the **information gain** for each candidate attribute(i.e, Outlook, Temperature, Humidity and Wind) and then select the one with **highest** information gain.
- **Attribute Temperature:**

  Temperature_Hot = {2+, 2-}

  Temperature_Mild = {4+, 2-}

  Temperature_Cool = {3+, 1-}

Gain(S,Temperature) = Entropy(S) — 4/14 Entropy(S_Temperature_Hot) — 6/14 Entropy(S_Temperature_Mild) — 4/14 Entropy(S_Temperature_Cool)

Gain(S,Temperature) = **0.029**

# Case Study

Similarly, the information gain for all the attributes are as follows :

1. Gain(S,Outlook)          = **0.246**
2. Gain(S,Temperature)   = **0.029**
3. Gain(S,Wind)             = **0.048**
4. Gain(S,Humidity)        = **0.151**

According to the information gain measure, the Outlook attribute provides the best prediction of the target attribute, PlayTennis, over the training examples. Therefore, Outlook is selected as the decision attribute for the root node, and branches are created below the root for each of its possible values (i.e., Sunny, Overcast, and Rain).

# Case Study

The resulting partial decision tree is shown in Figure, along with the training examples sorted to each new descendant node.
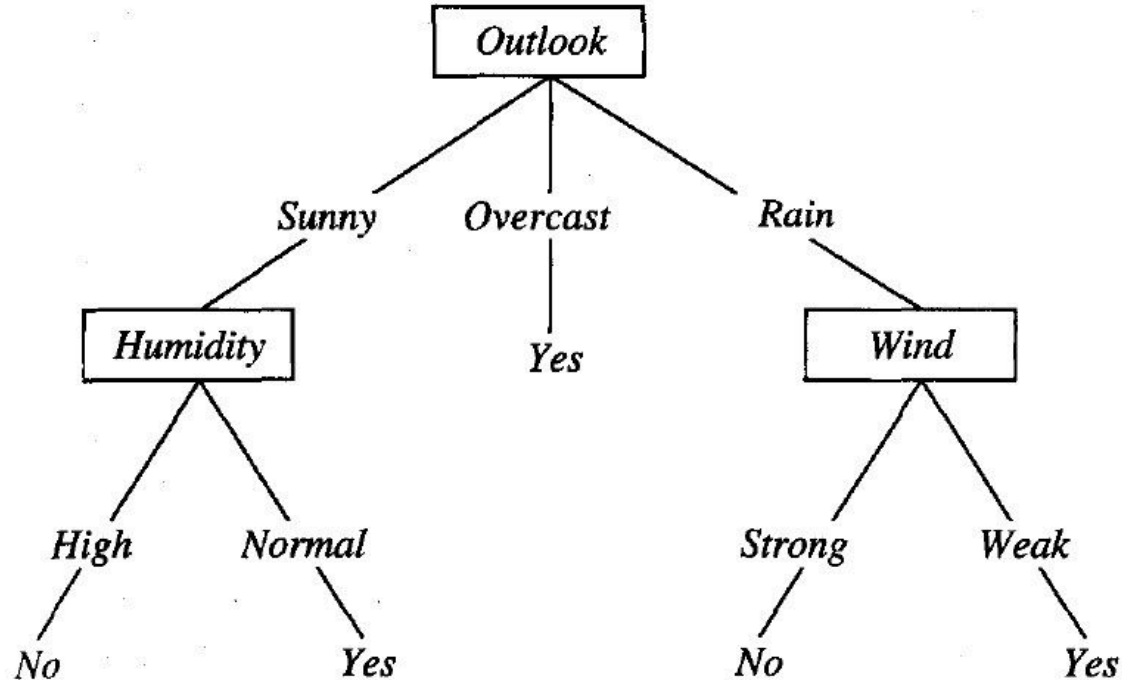


{D1, D2, ..., D14}

[9+,5−]

Outlook

Sunny  Overcast  Rain

{D1,D2,D8,D9,D11}     {D3,D7,D12,D13}     {D4,D5,D6,D10,D14}

[2+,3−]                      [4+,0−]                      [3+,2−]

?          Yes          ?

*Which attribute should be tested here?*

# Case Study

-> Note that every example for which Outlook = Overcast is also a positive example of PlayTennis. Therefore, this node of the tree becomes a leaf node with the classification PlayTennis = Yes. In contrast, the descendants corresponding to Outlook = Sunny and Outlook = Rain still have non-zero entropy, and the decision tree will be further elaborated below these nodes.

-> The process of selecting a new attribute and partitioning the training examples is now repeated for each non terminal descendant node, this time using only the training examples associated with that node.

-> Attributes that have been incorporated higher in the tree are excluded, so that any given attribute can appear at most once along any path through the tree. This process continues for each new leaf node until either of two conditions is met: (1) every attribute has already been included along this path through the tree, or (2) the training examples associated with this leaf node all have the same target attribute value (i.e., their entropy is zero).

# Case Study

– Eventually, our Final decision tree turns out to be like this:
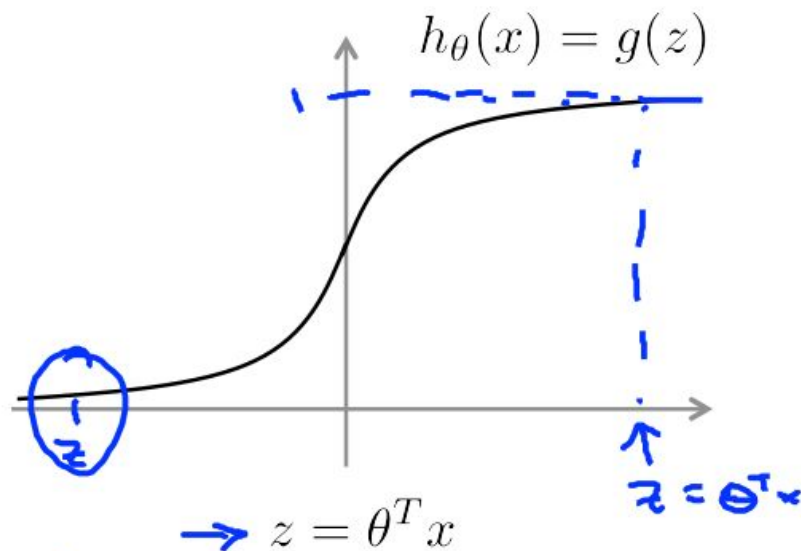
Notebook

# What are Support Vector Machines?

- **Support Vector Machines** construct a hyperplane or set of hyperplanes in a high or infinite-dimensional space, which can be used for classification, regression, or other tasks like outliers detection.
- An **SVM** model represents examples as points in space, mapped in such a way that the examples of the separate categories are divided by a clear gap as wide as possible.
- In addition to performing linear classification, **SVMs** can efficiently perform a non-linear classification using what is called the kernel trick, implicitly mapping their inputs into high-dimensional feature spaces.

# Alternative view of logistic regression

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

$$h_\theta(x) = g(z)$$

$$z = \theta^T x$$

$$z = \theta^T x$$

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

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

# Alternative view of logistic regression

$(x, y)$

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

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

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

$z = \theta^T x$



$-\log \frac{1}{1 + e^{-z}}$

$\text{Cost}_1(z)$

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



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

$\text{Cost}_0(z)$

## Support vector machine

Logistic regression:

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

Support vector machine:

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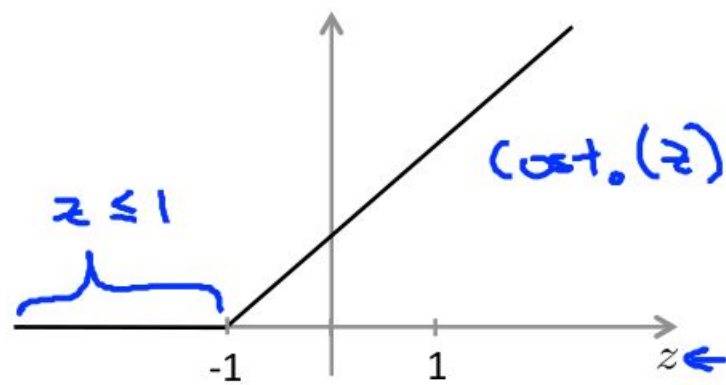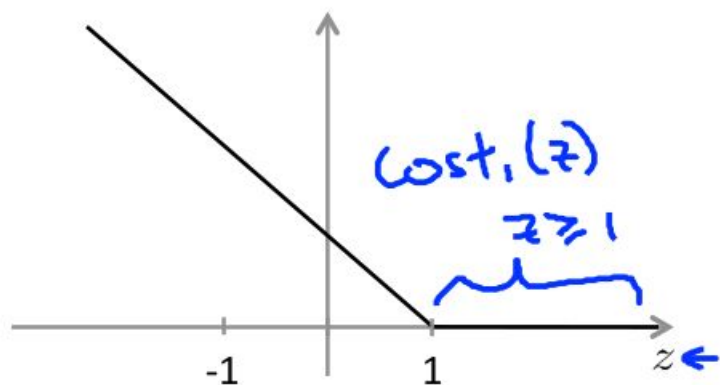
## SVM hypothesis

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

## SVM Decision Boundary

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

$= 0$

Whenever $y^{(i)} = 1$:

$$\theta^T x^{(i)} \geq 1$$

Whenever $y^{(i)} = 0$:

$$\theta^T x^{(i)} \leq -1$$

$$\min_\theta \; \cancel{C \times \theta} + \frac{1}{2}\sum_{i=1}^n \theta_j^2$$

$$s.t. \quad \theta^T x^{(i)} \geq 1 \quad if \quad y^{(i)} = 1$$

$$\theta^T x^{(i)} \leq -1 \quad if \quad y^{(i)} = 0$$

# Hard Margin

– If the training data is linearly separable, we can select two parallel hyperplanes that separate the two classes of data, so that the distance between them is as large as possible. The region bounded by these two hyperplanes is called the "margin", and the maximum-margin hyperplane is the hyperplane that lies halfway between them.

## Support Vector Machine

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



cost$_1$(z)    z ≥ 1

z ≤ 1    cost$_0$(z)

-1    1    z

-1    1    z

→ If $y = 1$, we want $\theta^T x \geq 1$ (not just $\geq 0$)     $\theta^T x \geq \cancel{0} \; 1$

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

C = 100,000

# SVM Decision Boundary: <u>Linearly separable</u> case



Large margin classifier

# Large margin classifier in presence of outliers

# Soft Margin

- To extend SVM to cases in which the data are not linearly separable, we introduce the **hinge loss** function,

$$\max\left(0, 1 - y_i\left(\vec{w} \cdot \vec{x}_i - b\right)\right).$$

- $y_i$ is the i-th target and $\vec{w} \cdot \vec{x}_i - b$ is the current output.
- This function is zero if the constraint in (1) is satisfied, in other words, if vector $x_i$ lies on the correct side of the margin.
- For data on the wrong side of the margin, the function's value is proportional to the distance from the margin.

# Soft Margin

– We then wish to minimise:

$$\left[\frac{1}{n}\sum_{i=1}^{n}\max\left(0, 1 - y_i\left(\vec{w}\cdot\vec{x}_i - b\right)\right)\right] + \lambda\|\vec{w}\|^2,$$

where the parameter λ determines the trade-off between increasing the margin size and ensuring that vector $x_i$ lie on the correct side of the margin. Thus, for sufficiently small values of λ, the second term in the loss function will become negligible, hence, it will behave similar to the hard-margin SVM, if the input data is linearly classifiable.

# Decision Boundary

# Vector Inner Product



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

$$u^T v = ? \qquad [u_1 \quad u_2] \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

$$\|u\| = \text{length of vector } u$$

$$= \sqrt{u_1^2 + u_2^2} \quad \in \mathbb{R}$$

$P = $ length of projection of $v$ onto $u$.

$$u^T v = \underline{P} \cdot \underline{\|u\|} \leftarrow \qquad = v^T u$$

Signed

$$= u_1 v_1 + u_2 v_2 \leftarrow \qquad P \in \mathbb{R}$$

$$u^T v = P \cdot \|u\|$$

$$P < 0$$

# SVM Decision Boundary

$$\omega = \left(\sqrt{\omega'}\right)^2$$

$$\min_{\theta} \frac{1}{2} \sum_{j=1}^{n} \theta_j^2 = \frac{1}{2}\left(\theta_1^2 + \theta_2^2\right) = \frac{1}{2}\left(\sqrt{\theta_1^2 + \theta_2^2}\right)^2 = \frac{1}{2}\|\theta\|^2$$

$$= \|\theta\|$$

s.t. $\quad \theta^T x^{(i)} \geq 1 \qquad$ if $y^{(i)} = 1$

$\quad \longrightarrow \theta^T x^{(i)} \leq -1 \quad$ if $y^{(i)} = 0$

$$\begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} \quad \theta_0 = 0$$

Simplication: $\quad \theta_0 = 0 \qquad \underline{n = 2}$

$\theta^T x^{(i)} = ?$

$\uparrow \qquad \uparrow$

$u^T v$



$$\theta^T x^{(i)} = \boxed{p^{(i)} \cdot \|\theta\|} \leftarrow$$

$$= \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)} \leftarrow$$

# SVM Decision Boundary

$\rightarrow \min\limits_{\theta} \dfrac{1}{2} \sum\limits_{j=1}^{n} \theta_j^2 = \dfrac{1}{2} \|\theta\|^2 \leftarrow$

s.t. $\boxed{p^{(i)} \cdot \|\theta\| \geq 1}$    if $y^{(i)} = 1$

$p^{(i)} \cdot \|\theta\| \leq -1$    if $y^{(i)} = 1$

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

Simplification: $\boxed{\theta_0 = 0}$

$\theta_0 \neq 0$

$\Big\}$ C very large

$p^{(i)} \|\theta\| \geq 0$

$\rightarrow \|\theta\|$ can be smaller.

margin

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

$\|\theta\|$ large

$p^{(i)} \leq 0$

$p^{(i)} \|\theta\| \leq -1$

$\|\theta\|$ large

$p^{(2)} \leq 0$

SVM hypothesis

# Kernels

# Non-linear Decision Boundary



Predict $y = 1$ if

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

$$h_\theta(x) = \begin{cases} 1 & \text{if} \quad \theta_0 + \theta_1 x_1 + \cdots \geq 0 \\ 0 & \text{otherwise.} \end{cases}$$
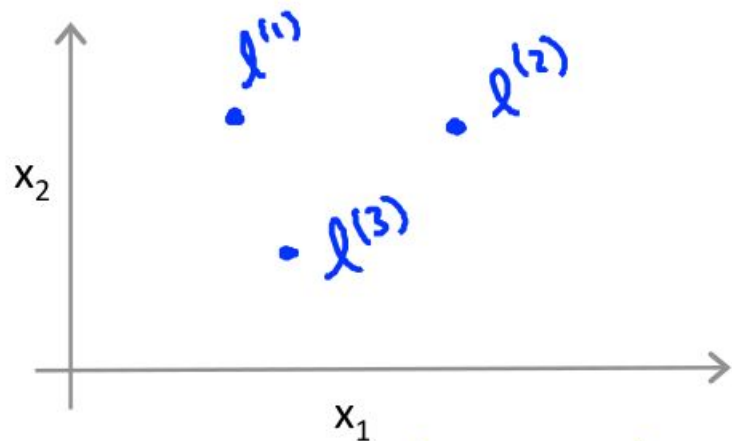
$$\rightarrow \theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 + \cdots$$
$$f_1 = x_1, \quad f_2 = x_2, \quad f_3 = x_1 x_2, \quad f_4 = x_1^2, \quad f_5 = x_2^2, \cdots$$

Is there a different / better choice of the features $f_1, f_2, f_3, \ldots$?

# Kernel



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

$\|w\|$

Given $x$:

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

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

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

Kernel (Gaussian kernels)

$$k(x, l^{(i)})$$

## Kernels and Similarity

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

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

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

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

$$f_1 = \exp\left(-\frac{(\text{large number})^2}{2\sigma^2}\right) \approx 0.$$
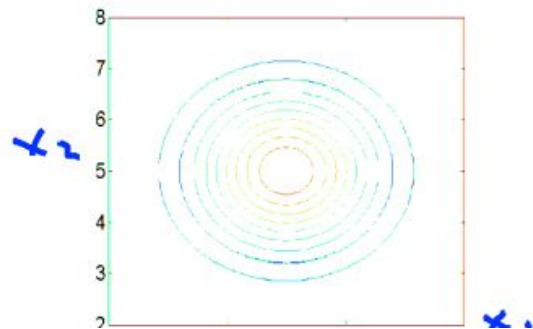
$l^{(1)} \rightarrow f_1$

$l^{(2)} \rightarrow f_2$

$l^{(3)} \rightarrow f_3.$

$x$

**Example:**

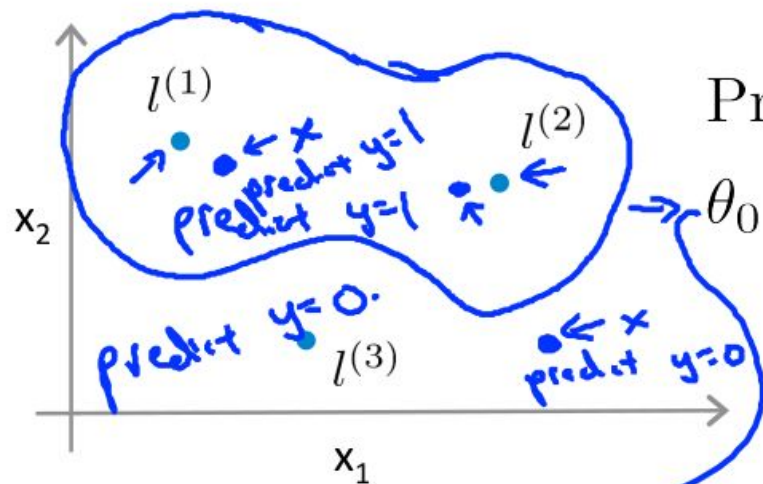$$l^{(1)} = \begin{bmatrix} 3 \\ 5 \end{bmatrix} \qquad f_1 = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right)$$

$\sigma^2 = 1$

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

$\sigma^2 = 0.5 \qquad\qquad \sigma^2 = 3$

Predict "1" when

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

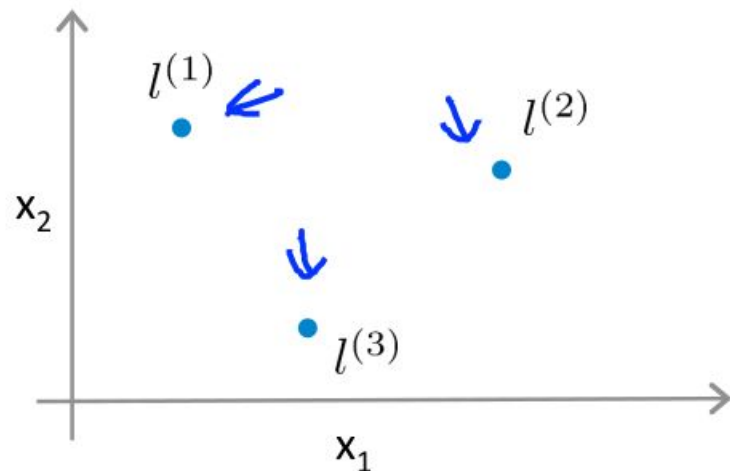$\theta_0 = -0.5$, $\theta_1 = 1$, $\theta_2 = 1$, $\theta_3 = 0$

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

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

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

$f_1, f_2, f_3 \approx 0$

$\rightarrow \theta_0 + \theta_1 f_1 + \cdots \approx -0.5 < 0$

In the figure:

$l^{(1)}$, $l^{(2)}$, $l^{(3)}$

predict y=1

predict y=1

predict y=0.

predict y=0

$x_2$, $x_1$

# Choosing the landmarks
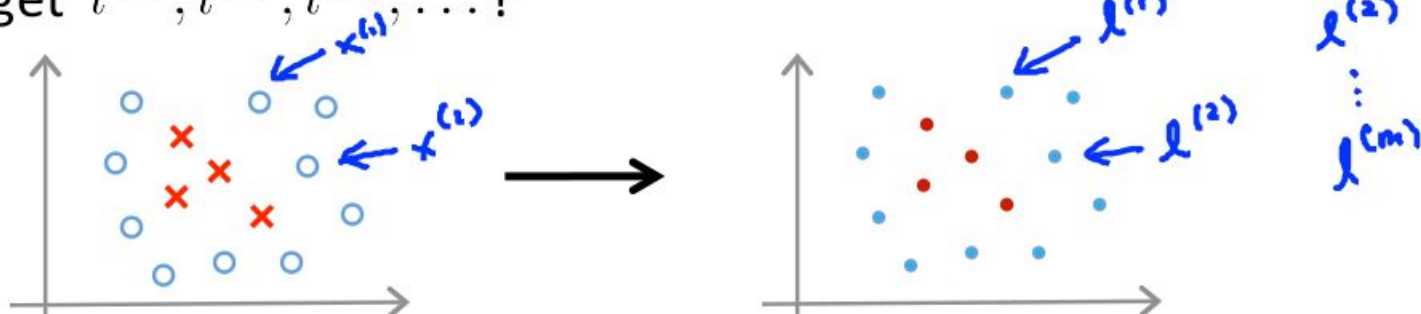


Given $x$:

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

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

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

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

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

## SVM with Kernels

→ Given $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \ldots, (x^{(m)}, y^{(m)})$,

→ choose $l^{(1)} = x^{(1)}, l^{(2)} = x^{(2)}, \ldots, l^{(m)} = x^{(m)}$.

Given example $\underline{x}$:

$$f_1 = \text{similarity}(x, l^{(1)})$$
$$f_2 = \text{similarity}(x, l^{(2)})$$
$$\ldots$$

$\leftarrow x^{(i)}$

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

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

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

$\leftarrow x^{(i)}$

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

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

# 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 $\underline{\theta^T f \geq 0}$

$\Theta_0 f_0 + \Theta_1 f_1 + \cdots + \Theta_m f_m$

Training:

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

$n = m$

$m$

$\Theta^T f^{(i)}$         $\Theta^T f^{(i)}$

$\to \Theta_0$

$$\sum_j \theta_j^2 = \Theta^T \Theta \leftarrow \quad \Theta = \begin{bmatrix} \Theta_1 \\ \vdots \\ \Theta_m \end{bmatrix}$$

(ignore $\Theta_0$]

$\|\Theta\|^2$

$\to \Theta^T M \Theta \leftarrow$

$M = 10,000$

# SVM parameters:

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

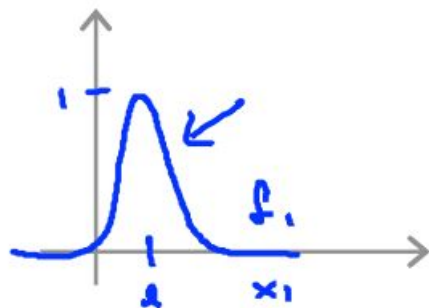→ Small C: Higher bias, low variance.     (large $\lambda$)

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

→ Higher bias, lower variance.

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



Small $\sigma^2$: Features $f_i$ vary less smoothly.
Lower bias, higher variance.

# Logistic Regression vs. SVMs

– Let us consider the following notation:
- **n** = number of features,     **m** = number of training examples
– If **n** is large relative to **m**:
- Use logistic regression, or SVM without a Kernel("linear kernel")
– If **n** is small, **m** is intermediate:
- Use SVM with Gaussian Kernel.
– If **n** is small, **m** is large:
- Create/add more features, then use logistic regression or SVM without a kernel

Presented by:
Shatadru Majumdar (https://github.com/shatadru99)
Purbayan Chowdhury (https://github.com/shivishbrahma)
Soham Biswas (https://github.com/Nibba2018)
Balaka Biswas (https://github.com/BALaka-18)
and
Team Explore ML and Team Innovacion