

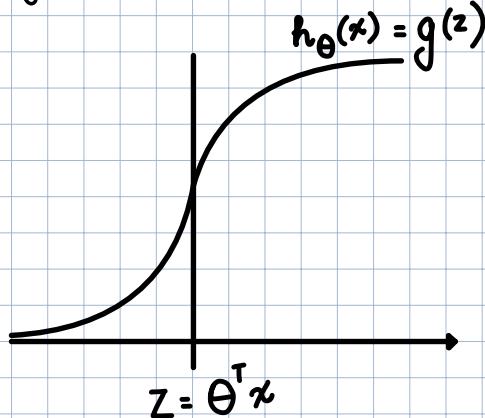
We are going to start from Logistic Regression and modify it to SVM.

From Logistic Regression, we have the hypothesis,

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

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

and the sigmoid graph as :-



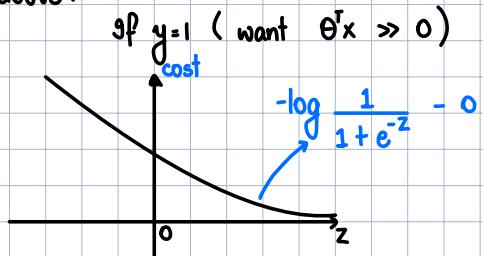
What do we like logistic Regression to do?

- Take an ex. of $y=1$. Then we hope to get $h_{\theta}(x) \approx 1 \Rightarrow \theta^T x \gg 0$.
↓
either in training / test set
for right on the graph
- Similarly for $y=0$, we hope to get $h_{\theta}(x) \approx 0 \Rightarrow \theta^T x \ll 0$.
for left on the graph

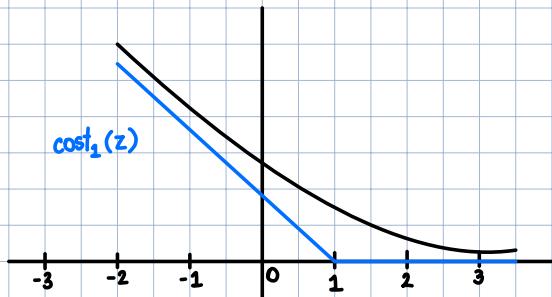
Cost fn. for logistics regression:-

$$\begin{aligned} \text{Cost } (h_{\theta}(x), y) &= -y \log(h_{\theta}(x)) - (1-y) \log(1-h_{\theta}(x)) \quad] \text{ of a single training ex.} \\ &= -y \log\left(\frac{1}{1+e^{\theta^T x}}\right) - (1-y) \log\left(1 - \frac{1}{1+e^{\theta^T x}}\right) \quad \text{i.e. } (x_i, y_i) \end{aligned}$$

Case 1 from above :



Consider the cost fn from above (Case 1) and add a new curve

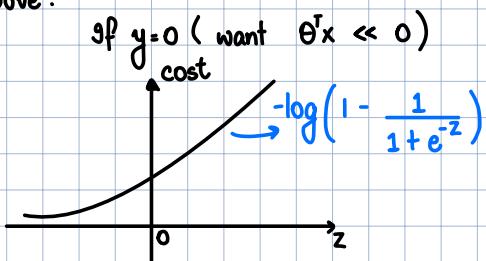


to the plot. We can observe that it is quite closely approximate the curve in black.

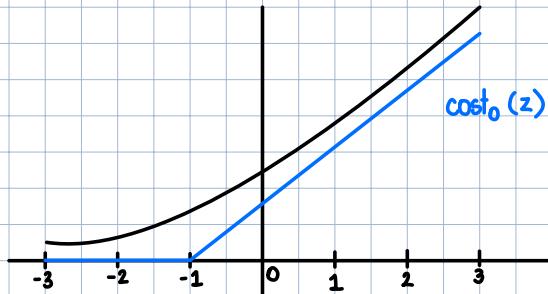
Don't worry about the slope of the blue line, it doesn't matter.

This blue line is the new cost fn. that we are going to use (when $y=1$)

Case 2 from above :



Adding a similar blue curve to this plot, we get,



Minimization Problem for Logistic Regression is :-

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

$\boxed{-y^{(i)} \log(h_{\theta}(x^{(i)}))}$
 $\boxed{(1-y^{(i)}) \log(1-h_{\theta}(x^{(i)}))}$

Minimization Problem for SVM becomes:-

$$\min_{\theta} \frac{1}{m} \sum_{i=1}^m \left[-y^{(i)} \text{Cost}_1(\theta^T x^{(i)}) - (1-y^{(i)}) \text{Cost}_0(\theta^T x^{(i)}) \right] + \frac{\lambda}{2m} \sum_{i=1}^n \theta_i^2$$

We parameterize the minimization problem slightly differently for SVM.

1) Since $\frac{1}{m}$ is constant term, it does not contribute anything to the minimization equation and we can get rid of it.

2) For logistic regression, we have 2 terms to the objective fn. as follows:-

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

Cost₁(θ^Tx⁽ⁱ⁾) Cost₀(θ^Tx⁽ⁱ⁾) 2nd term
1st term ↓ ↓
cost that comes from the training set regularization term
Rewriting ↓
 $\Rightarrow A + \lambda B$

Here, we have a trade off between these 2 terms. What we did was by setting different values for λ (regularization parameter) we could trade-off the relative weight b/w how much we want to fit the training set well (i.e minimizing A) and rest is how much we care about keeping the values of the λ small.

For SVM, just for convention, we are going to use different parameters. We will use,

$$\Rightarrow CA + B, \text{ where } C \text{ is the new convention (or parameter)}$$

Thus, we are going to minimize the above.

The way this would work is :-

- if λ is large, then B would have large weights.
- if C is small, then B would have large weights than A.

Thus, this is just a different way of parameterizing.

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

Finally, unlike logistic regression, SVM doesn't output a probability. Instead, we have the above cost fn. which we minimize to get the parameter θ .

Hypothesis :-

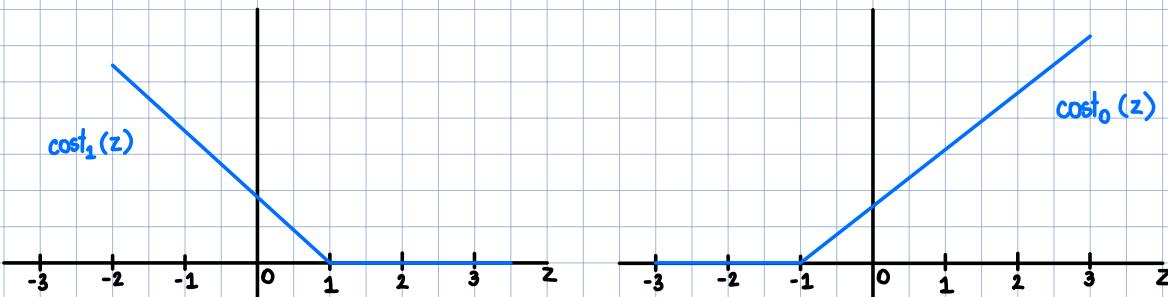
$$h_{\theta}(x) = \begin{cases} 1 & , \text{ if } \theta^T x \geq 0 \\ 0 & , \text{ otherwise} \end{cases}$$

SVM as Large Margin Classifier:-

What we have so far :-

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

Plots of Cost₁, Cost₀



If $y=1$, we want $\theta^T x \geq 1$ (not just ≥ 0)
(like logistic reg)

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

In logistic regression, we said we want $\theta^T x \geq 0$ and $\theta^T x < 0$, but in case of SVM,

we want more q, that is why we choose $\theta^T x \geq 1$ and $\theta^T x \leq -1$.

This gives us more safety or builds a safety margin factor into the SVM.

Case study :-

let's set $C = \text{large value}$. Then we would want to get $A = 0$.

How? We would the above discussed two cases:-

i) whenever $y^{(i)} = 1$, we want $\theta^T x \geq 1$.

ii) whenever $y^{(i)} = 0$, we want $\theta^T x \leq -1$.

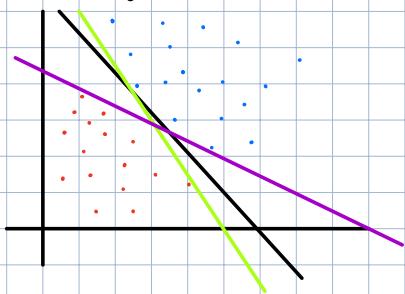
If we meet those conditions, then the minimization problem will

become :- $\min C \cdot 0 + \frac{1}{2} \sum_{i=0}^n \theta_i^2$ subject to (s.t.)

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

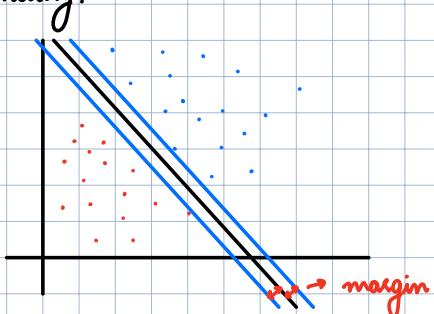
$\theta^T x \leq -1 \quad \text{if } y^{(i)} = 0$

When we try to solve the above min. problem, we get an interesting decision boundary like below:-



This is dataset of +ve and -ve examples. This data is linearly separable, that means \exists a straight line which separate the data perfectly.

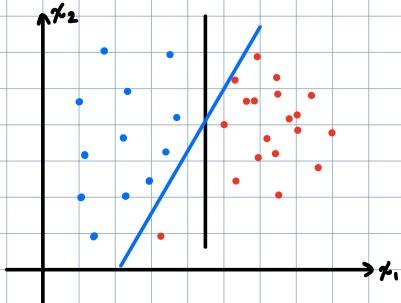
Out of the green, purple and black line, the SVM chooses the black line decision boundary.



Why will SVM choose the black line? because the distance btw the black and blue lines is the max & than green and purple. min distance from any of the examples (pts).

SVM is also called large Margin classifier, which is a consequence the optimization problem.

Large Margin Classifier in Presence of outliers



If C is not too large, then we will get the black line as our decision boundary.

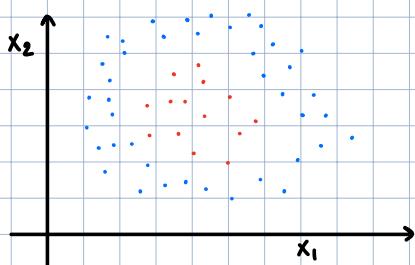
If C was too large, the black line will shift to the blue one.

So SVM will classify points correctly if C is large, otherwise in presence of outliers, the data will not be linearly separable.

Kernels - I

Here, we will start adapting SVM in order to develop complex non-linear classifier.

Non-linear Decision Boundary :-



Suppose we have a training set that looks like the one on left and we want to find a non-linear decision boundary to distinguish the +ve and -ve examples.

We can achieve the above by coming up with a set of complex polynomial features.

For example, predict $y = 1$ if $\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_2^2 + \dots \geq 0$

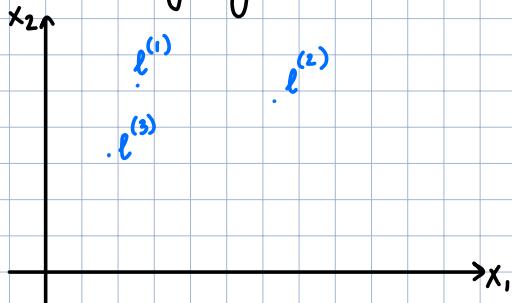
In hypothesis notation,

$$h_{\theta} = \begin{cases} 1 & \text{the above condition} \\ 0 & \text{otherwise} \end{cases}$$

Replace $f_1 = x_1$, $f_2 = x_2$, $f_3 = x_1 x_2$, $f_4 = x_2^2$,

We can see that computing these polynomials can get expensive, so is there a different / better choice of the features f_1, f_2, f_3, \dots ?

In this section, we will only define 3 new features. But for real life problems, we are going to define much larger number.



We manually pick up few points, marked in blue $l^{(i)}$, where $1 \leq i \leq 3$.

Now Given x , let's define $f_i = \text{some measure of similarity between the training example } x \text{ & } l^{(i)}$.

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

$\|\vec{w}\|$ represents length of \vec{w} . i.e. euclidean dist.

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

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

Similarity in the above formula means Kernal, k (mathematically).

The Kernel that we are using here is the Gaussian Kernel.

Kernels and Similarity

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

Now, if $x \approx l^{(i)}$: (i.e if x is close to one of the landmark's, $l^{(i)}$)

- $\|x - l^{(i)}\| \approx 0$

$$\text{So, } f_i \approx \exp\left(-\frac{0}{2\sigma^2}\right) \approx 1$$

Conversely, if x is far from $l^{(i)}$:

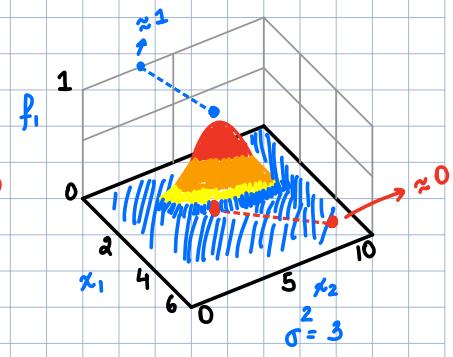
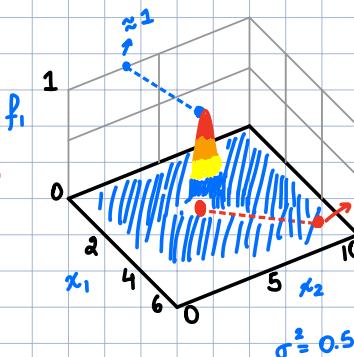
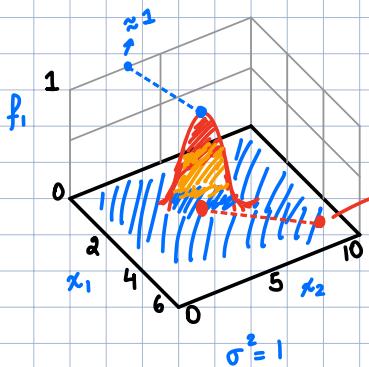
- $f_i = \exp\left(-\frac{(\text{large num})^2}{2\sigma^2}\right) \approx 0$

Now, given x , we can compute three new features f_1, f_2, f_3 given l^1, l^2, l^3 .

Example :-

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

Plot :-



Let definition of the features be as above (Gaussian).

Given the def, let's learn the types of hypothesis we can make

Thus, given the training example x , we are going to compute these features f_1, f_2 and f_3 . Our hypothesis is going to predict

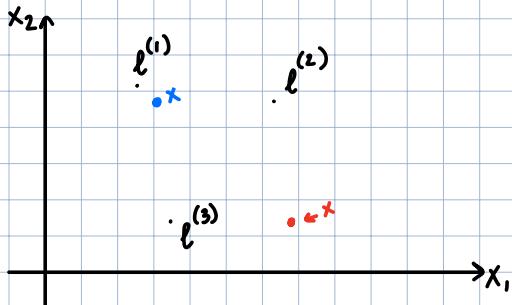
1 when $\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$.

"Predict 1 when

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

Let's say we already ran the algorithm and found that

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



- Added a new training example
- $f_1 \approx 1, f_2 = 0, f_3 = 0$
- Then, $\theta_0 + \theta_1 \cdot 1 + \theta_2 \cdot 0 + \theta_3 \cdot 0$
 $\Rightarrow -0.5 + 1 = 0.5$
- At x , predict $y=1$.

- Added a new training example in red

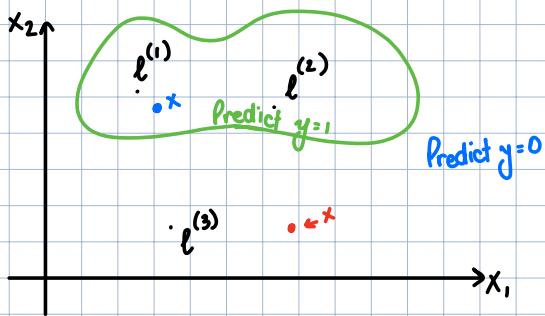
- $f_1 \approx 0, f_2 = 0, f_3 = 0$

- Then, $\theta_0 + \theta_1 \cdot 0 + \theta_2 \cdot 0 + \theta_3 \cdot 0$

$$\Rightarrow -0.5 = -0.5$$

- At x , predict $y=0$.

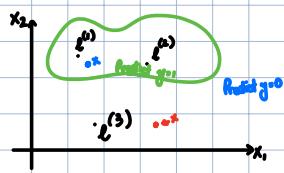
If we keep adding more points, we will see that examples near $l^{(1)}$ and $l^{(2)}$ will predict $y=1$. So the decision boundary becomes:-



Kernels II

So far we have

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

But where to get $l^{(1)}, l^{(2)}, l^{(3)}$ from ?

Here's where .

Given: $(x^{(i)}, y^{(i)})$

Choose:- $l^{(i)} = x^{(i)}$ $1 \leq i \leq m$.

we choose the location of landmarks exactly as the location of M training examples.

Then given : example x :

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

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

$$\text{det } f = \begin{bmatrix} \vdots & \vdots \\ f_1 \\ f_2 \\ \vdots \\ f_m \end{bmatrix}$$

Thus for $(x^{(i)}, y^{(i)})$,

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

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

\vdots

\vdots

\vdots

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

for convention

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

SVM with Kernels

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

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

Training :-

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

Solve the above minimization problem to get the values for our parameters for SVM.

Another way to write the last term above is using linear algebra,

$$\sum_{j=1}^m \theta_j^2 = \theta^T \theta \quad , \text{ where } \theta = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_m \end{bmatrix} \quad \text{ignoring } \theta_0$$

$$= \theta^T M \theta \quad , \text{ where } M \text{ is a matrix which depends upon the Kernel we use.}$$

The reason SVM adds the additional M is that with this modification it allows the SVM to scale to much bigger training sets.

Doesn't really need to know this because SVM models does this trick at the back hand.

SVM parameters :-

Another thing worth knowing about is how do we choose the parameters.

C (similar to $\frac{1}{\Delta}$).

regularization parameter used for logistic regression

Large C : lower bias, high variance \rightarrow prone to overfitting

Small C : Higher bias, lower variance \rightarrow prone to underfitting

σ^2 : (appeared in Gaussian Kernel)

large σ^2 : features f_i vary more smoothly.

Higher bias, lower variance.

Small σ^2 : features f_i vary less smoothly.

lower bias, higher variance.

Logistic regression v.s SVM

When to use them?

Let $n = \# \text{ of features } (x \in \mathbb{R}^{n+1})$

$m = \# \text{ of training examples}$.

- 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 Gaussian kernel

- if n is small, m is large :

Create / add more features, then use logistic regression

Or SVM without a kernel.