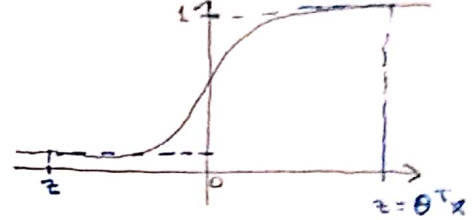


## Support Vector Machine (SVM)

SVM, bazen lojistik regresyon veya sinir ağlarına oranla karmaşık doğrusal olmayan işlevleri öğrenmenin daha temiz ve bazen daha güçlü bir yolunu sunar

logistic regressio:

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

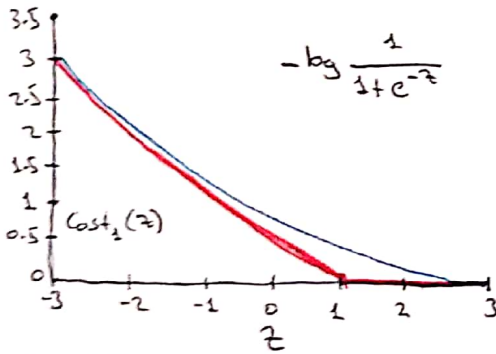


Eğer  $y=1$  ise, biz  $h_{\theta}(x) \approx 1$  olmasını isteriz.  $\theta^T x \gg 0$

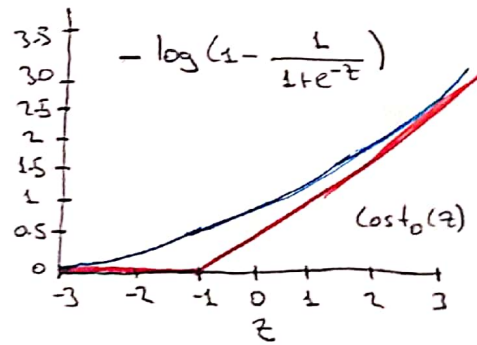
Eğer  $y=0$  ise, biz  $h_{\theta}(x) \approx 0$  olmasını isteriz.  $\theta^T x \ll 0$

$$\text{Cost} = -y \log \frac{1}{1 + e^{-\theta^T x}} - (1-y) \log \left( 1 - \frac{1}{1 + e^{-\theta^T x}} \right)$$

Eğer  $y=1$  ise



Eğer  $y=0$  ise



$z$  büyüdükçe  $\frac{1}{1+e^z}$  değeri 1'e yakınsayacak ve  $\log(1) = 0$  olduğu için cost function değeri 0'a yakın olacaktır.

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

Support vector Machine:

$$\frac{A + \lambda B}{A + B} \quad C = \frac{1}{\lambda}$$

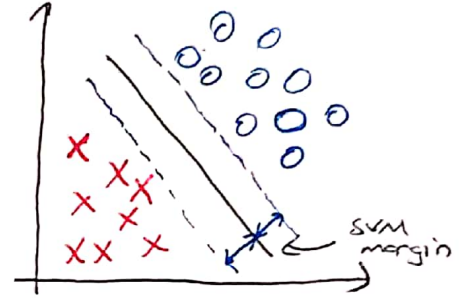
$$\min_{\theta} C \sum_{i=1}^m \left[ y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

Lojistik regresyon gibi belli bir örnek sayısına bölerek optimizasyonunu sağlayabileceğimiz gibi vektör destek makinelerinde olduğu gibi  $C$  gibi bir katsayıyla herperdekte optimize edebiliriz.  $C = \frac{1}{\lambda}$  eşitliğine sahip olarak dizayn edilmiştir.

## SVM Hipotez

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

SVM, bazı large margin classifier olarak adlandırılır. Bunun nedeni sınıflara belli bir uzaklıkta sınıflandırıcı çizgisini çizmesidir. Bu margin optimizasyon sürecisinden kaynaklanmaktadır.



## Large Margin Classification

$$\min_{\theta} C \sum_{i=1}^m [y^{(i)} \cos \theta_1(\theta^T x^{(i)}) + (1-y^{(i)}) \cos \theta_2(\theta^T x^{(i)})] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

$\hookrightarrow 0$

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

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

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

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

Optimizasyon sonucu örneklerin  $y$  değerine göre sonucu yandaki gibi geldiği görülmüştür. Sonuç 0 geldiğinde;

$$\min_{\theta} C \times 0 + \frac{1}{2} \sum_{j=1}^n \theta_j^2 \quad (C \text{ çok büyük seçilmeli})$$

Bu optimizasyon problemini çözmek için  $\theta$ 'lerin küçük seçilmesi gerekir. bunun içinde örneklerden uzak bir boundary çizilmelidir. İşte bu nedenle margin oluşur.

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

$\|\theta\|$

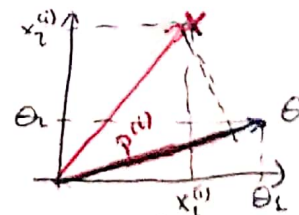
simple example  $\theta_0 = 0$   
 $n=2$

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

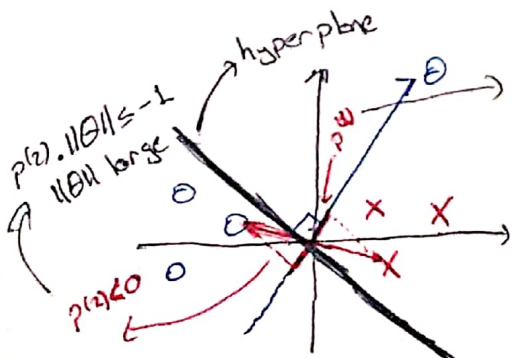
$$\theta^T x^{(i)} \leq -1 \quad \text{if } y^{(i)} = 0$$

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

$$\theta^T x$$



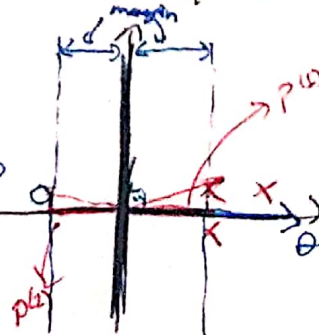
$$\theta^T x^{(i)} = p^{(i)} \cdot \|\theta\| = \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)}$$



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

$\|\theta\|$  large

Örneklerin izdüşümü arttıkça  $\theta$  değerleri artmaktadır. Margin değeri büyümektedir.

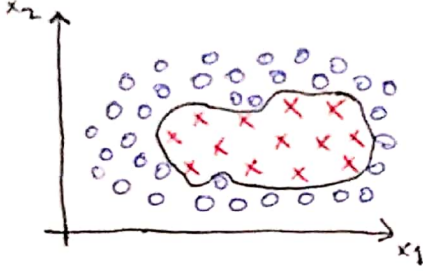


$$p^{(i)} \cdot \|\theta\| \geq 1$$

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

## Kernel 1

Non-linear



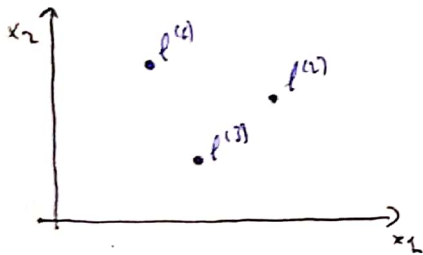
$$h_{\theta}(x) = \begin{cases} 1 & \text{if } \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \dots \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

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

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

Yüksek dereceli polinomları hesaplamak oldukça pahalıdır.

Kernel



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

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

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

Yukarıda kullandığımız benzerlik işlemi bir çekirdek (kernel). Bu çekirdeğin altında Gauss Kernel denir.

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

## Kernel and Similarity

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

Eğer  $x \cong l^{(1)} \Rightarrow f_1 \cong \exp\left(-\frac{0^2}{2\sigma^2}\right) \cong 1$

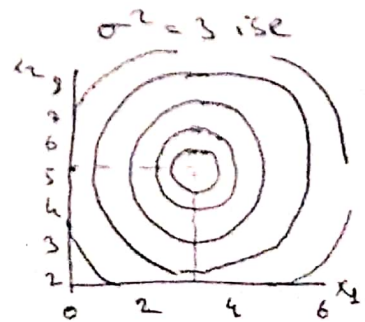
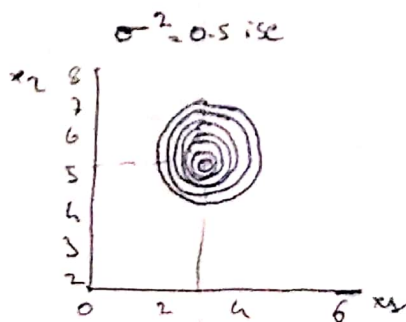
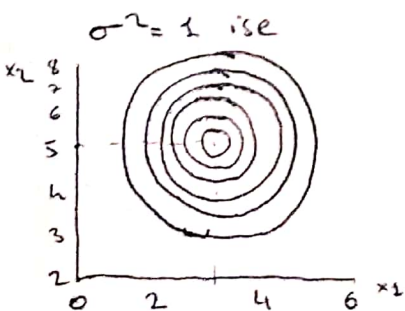
Eğer  $x, l^{(1)}$  den çok uzaksa  $\Rightarrow f_1 = \exp\left(-\frac{(\text{large number})^2}{2\sigma^2}\right) \cong 0$ .

Kernel,  $x$ 'in yer işaretlerinden herhangi birine ne kadar benzediğini ölçmektedir.  $x, l$ 'ye yakınsa 1'e yakınsa, uzaksa sıfıra yakınsa.

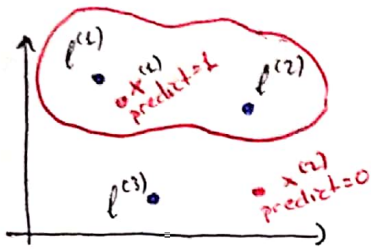
## Example

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

Generally:  $\|x^{(i)} - l^{(i)}\|^2 = \sum_{k=1}^n (x_k^{(i)} - l_k^{(i)})^2$   $f_{(i)} = \exp\left(-\frac{\|x^{(i)} - l^{(i)}\|^2}{2\sigma^2}\right)$



Not  $\sigma$  gauss dağılımıdır. Azaldıkça kontur grafiği daralır ve 0'a daha hızlı iner. Artıkça ise kontur grafiği genişler ve 0'a çok yavaş iner.



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

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

$x^{(1)}$  için

$$f_1 \leq 1, f_2 \leq 0, f_3 \leq 0$$

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

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

$\Rightarrow 1$  predict

$x^{(2)}$  için

$$f_1, f_2, f_3 \leq 0$$

$$= \theta_0 + \theta_1 \times 0 + \dots = -0.5 < 0 \Rightarrow 0 \text{ predict}$$

Yer işaretlelerini nasıl belirleyeceğimiz konusuna gelince eğitim örneklerimize yakın noktalar seçeceğiz. (Yeni cinsi denilebilir)

## SVM with Kernels

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

choose  $l^{(1)} = x^{(1)}, l^{(2)} = x^{(2)}, \dots, l^{(m)} = x^{(m)}$

- $x$  girdisi verildiğinde;

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

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

$\vdots$

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

$$f_0 = 1$$

- Eğitim seti için  $(x^{(i)}, y^{(i)})$ :

$$x^{(i)} \rightarrow \begin{bmatrix} f_1^{(i)} = \text{sim}(x^{(i)}, l^{(1)}) \\ f_2^{(i)} = \text{sim}(x^{(i)}, l^{(2)}) \\ \vdots \\ f_m^{(i)} = \text{sim}(x^{(i)}, l^{(m)}) \end{bmatrix}$$

$$\leftarrow f_i^{(i)} = \text{sim}(x^{(i)}, l^{(i)}) = \exp\left(-\frac{\|x^{(i)} - l^{(i)}\|^2}{2\sigma^2}\right) \approx 1$$

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

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

$$f_0^{(i)} = 1$$

Ortaya çıkan  $f^{(i)}$  vektörü yeni eğitim örneğini temsil edecek özellik vektörüdür.

Hipotez  $z$   $x$  verildiğinde, hesaplanacak features  $f \in \mathbb{R}^{m+1}$

$\rightarrow$  predict "y=1" if  $\theta^T f \geq 0 \Rightarrow \theta_0 f_0 + \theta_1 f_1 + \theta_2 f_2 + \dots$

Ağırlıkları hesaplayabilmek için SVM algoritmasının cost function'ını kullanırız.

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

$$\theta \cdot \theta^T = \|\theta\|^2$$



Not Kernel'li lojistik regresyon gibi algoritmelerde kullanmama nedenimiz oldukça yavaş çalışmasından

## SVM Parameters

$$C (= \frac{1}{\lambda})$$

C Büyükse : lojistik regresyondaki küçük lambda değerine yeni az regularization kullanmaya denk gelir.

- Düşük bias, yüksek variance.

C Küçükse : lojistik regresyondaki büyük lambda kullanmaya denk gelir.

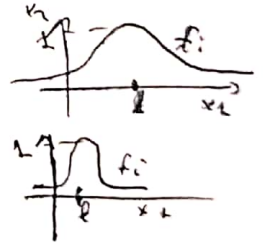
- Yüksek bias, düşük variance.

C çok büyük seçilirse overfit durumuna düşebilir. Küçük seçilirse underfit durumuna düşebilir.

$$\sigma^2$$

$\sigma^2$  Büyükse :  $f_i$  feature'ları o'a yavaşça iner. Yüksek bias, düşük variance.

$\sigma^2$  Küçükse :  $f_i$  feature'ları o'a hızlıca iner. Düşük bias, yüksek variance.



Not SVM algoritması, belirli bir optimizasyon problemi ortaya atar. Bu nedenle parametrelerin  $\Theta$ 'sını kendi yazdığınız yapımla gözölmesi öğrenilmemektedir.

Eğer yazmak zorunda kalırsanız;

- C parametresini seçmelisiniz.
- Kernel seçmek
  - Kernel seçmeyebilirsiniz. Buna "linear kernel" denir.  
predict "yes" if  $\Theta^T x \geq 0$
  - + Linear kernel küçük veri setiniz varsa seçilebilir.
  - Gauss Kernel
  - + Büyük bir veri setiniz varsa seçmek mantıklıdır.

Not Çok farklı ölçeklerde feature'larınız varsa, Gauss Kernel'ini kullanmadan önce feature scaling yapmak önemlidir.

Hangi kernel'i seçerseniz seçin, hepsinin teknik bir koşulu karşılaması gerekir. Buna Mercer Teoremi denir. İhtiyaç duyduğumuz nedeni optimizasyon hilelerine sahip olmasıdır. Yaptığı şey SVM paketlerinin optimizasyon sınıflarını kullanabilmesidir.

## Other Kernel

- Polynomial Kernel

Çok fazla kullanılmaz-  $k(x, l) = (x^T l)^2$

Diğer türleri  $\Rightarrow (x^T l)^3, (x^T l + 1)^3, (x^T l + 5)^4$

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

Genelde kötü performans gösterir.

Negatif olmayan veriler için kullanılır.

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

Not Gökle sınıflandırma için lineer regresyondaki cell in one mantığı kullanılabilir.

## Logistic Regression vs. SVM

$n \Rightarrow$  features  $m \Rightarrow$  training set size

- Eğer özellik sayınız,  $n$  eğitim setinizin boyutundan  $m$  büyükse lojistik regresyon ya da lineer kernel SVM kullanılmalıdır. Bir e-postanın spam olup olmadığını metne göre sınıflandırma örnek verilebilir.

- $n$  küçük ( $1-1000$ ),  $m$  ortalama ( $10-10000$ ) ise SVM Gauss kernel kullanılmalıdır.

- $n$  küçük ( $1-1000$ ),  $m$  büyük ( $10000+$ ) ise SVM Gauss kernel performans süntüsü yosotır. Bunun için daha fazla özellik ekleyerek lojistik regresyon ya da kernel'sız SVM kullanılmalıdır.

Not SVM'in optimizasyon problemi dış bükay optimizasyon problemidir.

Optimize edilmiş bir SVM paketi zehir ağından daha iyi ve hızlı çalışabilir.