

Kernel Methods (SVM)

Here are the methods of Kernel in Machine Learning mentioned below:

1) PCA

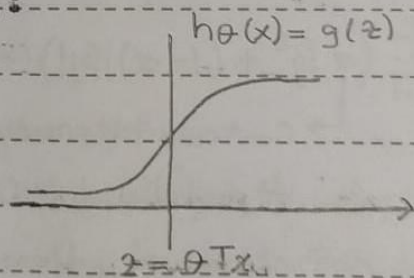
2) SVM

3) Spectral Clustering

SVM (Support Vector Machines)

Logistic Reg. hatirayam:

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



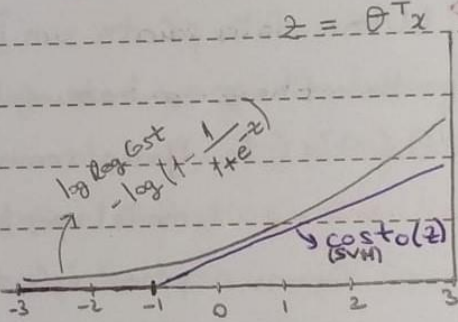
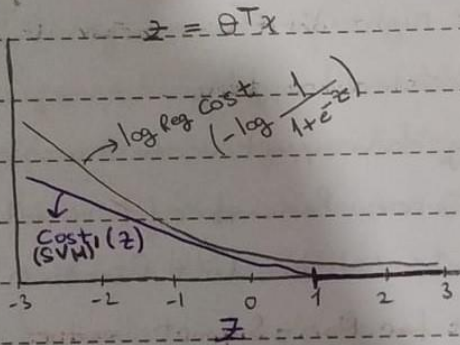
if $y=1$, we want $h_{\theta}(x)=1$, $\theta^T x \gg 0$

if $y=0$, we want $h_{\theta}(x)=0$, $\theta^T x \ll 0$

Cost Func Logistic Reg: $-\left(y \log h_{\theta}(x) + (1-y) \log (1-h_{\theta}(x))\right)$
 $= -y \log \frac{1}{1+e^{-\theta^T x}} - (1-y) \log \left(1 - \frac{1}{1+e^{-\theta^T x}}\right)$

if $y=1$ (we want $\theta^T x \gg 0$)

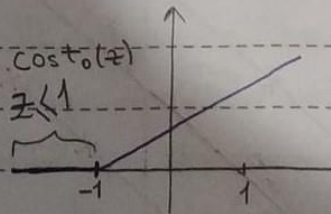
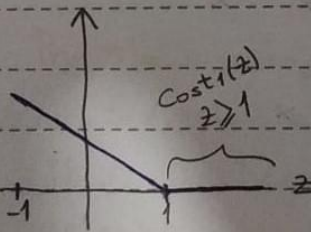
if $y=0$ (we want $\theta^T x \ll 0$)



Cost Func SVM: Hinge loss

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

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



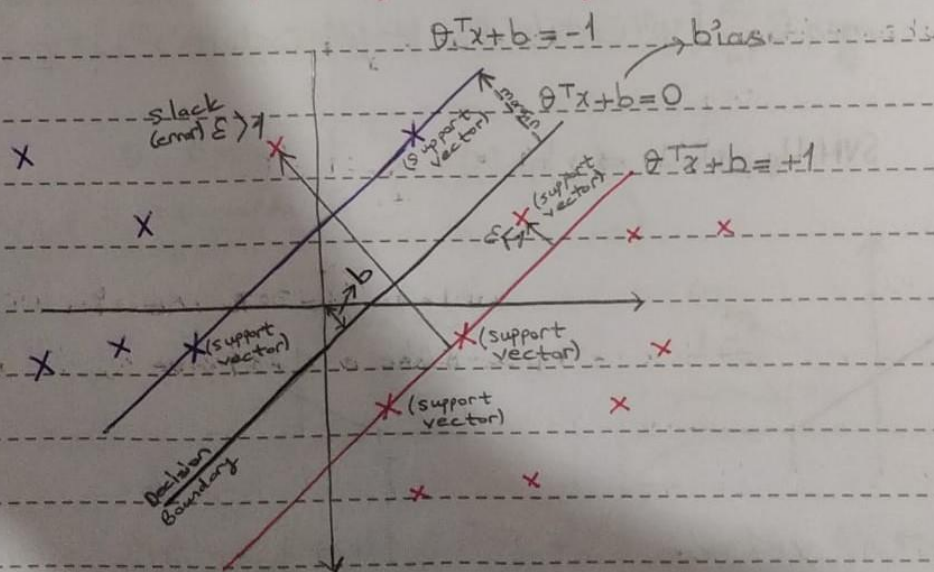
\rightarrow if $y=1$, we want $\theta^T x \gg 1$ (not just ≥ 0)

\rightarrow if $y=0$, we want $\theta^T x \ll -1$ (not just ≤ 0)

\hookrightarrow Eğer bir gözlemi doğru sınıflandıyorsak ve bu gözlem margin'den en az 1 (basically away from decision boundary) uzak ise Cost = 0 olur. Ancak gözlem margin içerisinde veya boundary'inin yakınında ise Cost $\neq 0$ olur.

→ SVM gives some punishment to both incorrect predictions and those close to decision boundary ($0 < \theta^T x < 1$). That is how we call them Support Vectors.

→ When data points are just right on the margin $\theta^T x = 1$; when they are between decision boundary and margin $0 < \theta^T x < 1$. Yani; removing non-support vectors won't affect model performance. Remember! Model fitting process is to minimize the Cost function. Since there is no cost for Non-Support vectors at all, total value of cost function won't be changed by adding or removing them.



→ Observations on the margin, and within the Soft Margin, are called Support Vectors.

Şimdiye dek neler öğrendik?

1) **Maximal Margin classifier**, lineer olarak ayrılabilen sınıfları ayırmak için kullanılır. buay yapar iken "öyle bir boundary bulur ki largest margin olur."

2) Separating hyperplane olmadığı zaman, maximal Margin classifier'ın general bir formu olan **Support vector classifier** kullanılır. **Support Vector Classifier**, **Soft Margin classifier** olarakta bilinir. Bu durumda, bazı gözlemlerin margin'ın veya hyperplane'ın yanlış tarafında olmasına izin veririz.

$\epsilon_1, \epsilon_2, \dots, \epsilon_n$ are **slack variables**, that allow individual observations to be on the wrong side of the margin or hyperplane.

$\epsilon_i \geq 0$;

The **slack variable ϵ_i** tells us where the i th observation is located, relative to hyperplane and relative to the margin.

if $\epsilon_i = 0$ then observation is on the correct side of margin.

if $\epsilon_i > 0$ then " " " " " wrong " " "

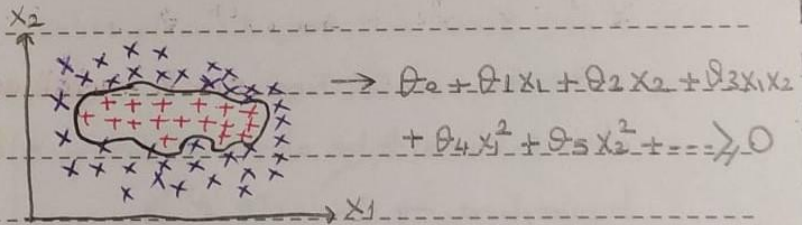
if $\epsilon_i > 1$ then " " " " " " " hyperplane

* **Margin üzerinde olan veya kendi sınıflarına göre margin yanlış tarafında olanlar Support Vector'ler olarak bilinir.**

3) Support Vector Classifier, boundary between two classes linear ise kullanılır. Ama çoğu kez Non-linear class boundary ile karşılaşırız. Yani çoğu zaman Linear Boundary sınıfı ayırmamıza yetmez. Bu durumda **Support Vector Machine (SVM)** kullanılır. SVM, is extension of the Support Vector Classifier that results from enlarging the feature space in a specific way, using Kernels.

↓ Gelin Kernels neymiş bakalım.

Non-Linear Decision Boundary:



$$h_{\theta}(x) = \begin{cases} 1 & \text{if } \theta_0 + \theta_1 x_1 + \dots + \theta_5 x_2^2 \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

Yani burada yeni feature'larımız ne oldu?

$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 + \dots \quad \{ f_i \text{ is new feature} \}$$

$$f_1 = x_1, f_2 = x_2, f_3 = x_1 x_2, \dots$$

→ Bu şekilde yaptığımızda computationally expensive'dir

Daha iyi bir çözüm var mıdır? Bu feature'lara bulmak için.

* Var Kernel (similarity) Functions.

SVM with Kernels:

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

Choose $\ell^{(1)} = x^{(1)}, \ell^{(2)} = x^{(2)}, \dots, \ell^{(m)} = x^{(m)}$

(landmarks'lar direkt training example'lar olarak oluyor.)

Given example x :

$$\left. \begin{aligned} \rightarrow \text{feature 1} = f_1 &= \text{similarity}(x, \ell^{(1)}) \\ \rightarrow \text{feature 2} = f_2 &= \text{similarity}(x, \ell^{(2)}) \\ &\vdots \\ &\vdots \end{aligned} \right\} f = \begin{bmatrix} f_1 \\ f_2 \\ \vdots \\ f_m \end{bmatrix}$$

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

$$f_1^{(i)} = \text{sim}(x^{(i)}, \ell^{(1)})$$

$$x^{(i)} \rightarrow f_2^{(i)} = \text{sim}(x^{(i)}, \ell^{(2)}) \quad \text{Gaussian Kernel}$$

$$\leftarrow f_1^{(i)} = \text{sim}(x^{(i)}, \ell^{(i)}) = \exp\left(-\frac{0}{2\sigma^2}\right) = 1 \quad \left(\begin{array}{l} \text{feature'lardan biri} \\ \text{1 olacak, çünkü kendisi} \\ \text{ile benzerliğine baktık} \end{array} \right)$$

$$f_m^{(i)} = \text{sim}(x^{(i)}, \ell^{(m)})$$

$$f^{(i)} = \begin{bmatrix} f_1^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix} \Rightarrow \text{New feature Vector}$$

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

~~$\theta^T x^{(i)}$~~

* Predict $y=1$ if $\theta^T f \geq 0$; $f \in \mathbb{R}^{m+1}$, $\theta \in \mathbb{R}^{m+1}$

$$\rightarrow \theta_0 + \theta_1 f_1 + \dots + \theta_m f_m$$

* Landmarks'lar training observations ile aynı gözlemler olur.

SVM parameters (Andrew Ng):

C ve σ^2 , sklearn kısmında anlatılacak.

SVM General:

1) P_0 Kernel (Linear Kernel)

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

2) Gaussian Kernel

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

Note:

→ Perform feature scaling before using Gaussian Kernel.

Choose σ^2

3) Other Choices of Kernel

Polynomial $\Rightarrow \text{sim}(x, l) = k(x, l) \Rightarrow (x^T l)^2, (x^T l)^3, (x^T l + 1)^3$

RBF $\Rightarrow k(x, l) = \exp(-\gamma \|x - l^{(i)}\|^2)$

SVM parameters:

1) **C**, default=1 (like log Reg, smaller values specify strong regularization)

C behaves as a Regularization Parameter in SVM.

- A smaller C emphasizes the importance of \mathcal{E} (slack variables).
- A larger C decrease the importance of \mathcal{E} (slack variables).

→ Small values of C result in a wider Margin.

→ Large values of C result in smaller Margin.

$C = \frac{1}{\lambda}$ (similar no regularization)
Large C: Lower Bias, High variance

Small C: Higher Bias, low variance

Logistic Regresyondan farklı olarak, λ regularization term'in önünde kullanılmaz. SVM, C'yi fit term'in önünde kullanır. (Bu bizi değiştirmez aynı şey yooi)

• Logistic Regression:

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

• Support Vector Machine:

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

fit Term (Cost coming from training set) Regularized Term

→ Fit term emphasizes fit the model very well by finding optimal coefficients, and the regularized term controls complexity of the model by constraining the large value of coefficients. There is a trade off between ^①fitting the model well on training dataset and ^②the complexity of the model that may lead to overfitting, which can be adjusted by changing the value of λ or C.

2) kernel, default = "rbf"

Note: rbf kernel is also called as Gaussian kernel

3) degree, default = 3

Degree of the polynomial kernel function

4) gamma, default = "scale"

Note:

$$\left. \begin{aligned} \text{Gaussian RBF} &= \exp \left\{ -\frac{\|x - x^{(i)}\|^2}{2\sigma^2} \right\} \\ \text{Gaussian RBF} &= \exp \left\{ -\gamma \|x - x^{(i)}\|^2 \right\} \end{aligned} \right\} \gamma = \frac{1}{2\sigma^2}$$

* High gamma means more Eğrilik (model heavily impacted by support vectors)
Low " " less Eğrilik (model behave like linear SVM)