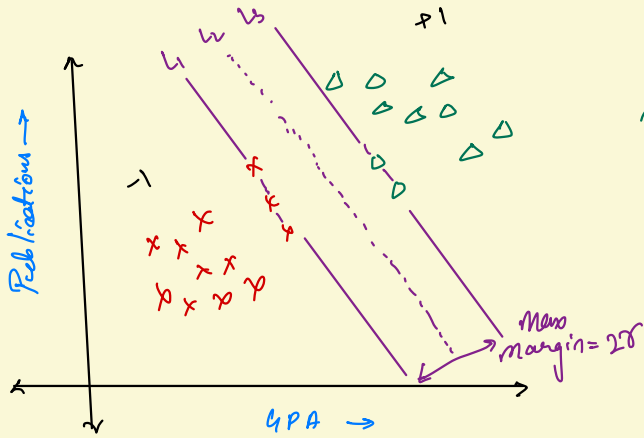


Assignment-1

S. Soundarya

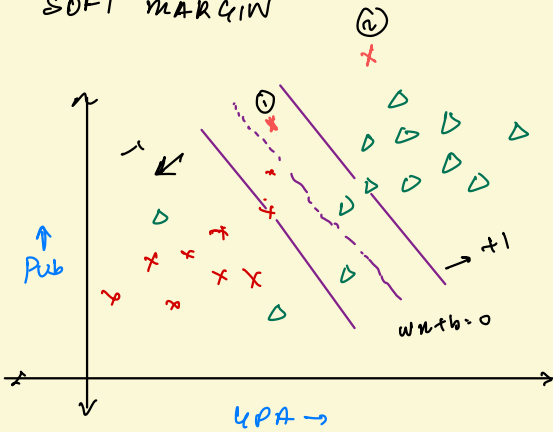
Q8



We try to find a hyperplane that is not dangerously close to both the data to avoid misclassification / we want to max. the margin. (r)
r: shortest distance b/w the margin & the chosen hyperplane.

* Linear SVMs

SOFT MARGIN



For data that are not linearly separable.

Classify data so that type I is on one side of Bound & type II is on the other side but with some misclassification with a penalty for every mistake depending upon the mistake

$$\text{Hinge loss: } \max(0, 1 - y_i^*(w \cdot x_i + b))$$

Incorrectly $y_i(w \cdot x_i + b) < 0$

$$\textcircled{1}: \max(0, 1 - (< 0 & > -1)) \rightarrow \text{lesser}$$

$$\textcircled{2}: \max(0, 1 - (< 0 & < -1)) \rightarrow \text{greater}$$

Soft SVMs:

maximize margin

minimize hinge loss

$$\text{Maximizing } \gamma = \text{minimizing } \frac{1}{\|w\|} \quad (\because \gamma = \frac{1}{\|w\|})$$

L1 SVM

$$\min_{w, b, \xi_i} \frac{1}{2} w^T w + C \sum_{i=1}^N \xi_i$$

$$\xi_i = \max(0, 1 - y_i(w^T x_i + b)) \quad 1 \leq i \leq N \quad \text{PRIMAL}$$

$$\xi_i \geq 0 \quad \xi_i \geq 1 - y_i(w^T x_i + b)$$

DUAL

$$\text{Lagrangian } L, \quad \frac{\partial L}{\partial w} = \frac{\partial L}{\partial \xi_i} = \frac{\partial L}{\partial b} = \frac{\partial L}{\partial \alpha_i} = \frac{\partial L}{\partial \beta_i} = 0$$

DUAL PROBLEM

$$\begin{aligned} \text{minimize}_{\alpha_i} \quad & \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j x_i^T x_j - \sum_{i=1}^N \alpha_i \\ & \sum_{i=1}^N \alpha_i y_i = 0 \quad \text{and} \quad 0 \leq \alpha_i \leq C \end{aligned}$$

VECTORIZED

$$\begin{aligned} \max_{\alpha} \quad & 1^T \alpha - \frac{1}{2} \alpha^T Q \alpha \\ \text{subject to} \quad & y^T \alpha = 0 \\ & \alpha \geq 0 \end{aligned}$$

$$\sum_{i=1}^N \sum_{j=1}^N d_i d_j y_i y_j x_i^T x_j = d^T Q d \quad \text{where} \quad Q = [Q_{ij}]_{n \times n}$$

$$[Q_{ij}]_{n \times n} = [y y^T]_{n \times n} \odot \begin{bmatrix} x_1^T x_1 & \dots & x_1^T x_n \\ & \ddots & \\ x_n^T x_1 & \dots & x_n^T x_n \end{bmatrix}$$

↓
Inner product b/w x_i & x_j

Try out with $N=2$:

$$y y^T = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \begin{bmatrix} y_1 & y_2 \end{bmatrix} = \begin{bmatrix} y_1 y_1 & y_1 y_2 \\ y_2 y_1 & y_2 y_2 \end{bmatrix}$$

$\text{np.dot}(y, y.T)$

$Q = y \otimes x$ // returns element wise

$$\therefore y y^T \odot \begin{bmatrix} x_1^T x_1 & x_1^T x_2 \\ x_2^T x_1 & x_2^T x_2 \end{bmatrix} = \begin{bmatrix} y_1 y_1 x_1^T x_1 & y_1 y_2 x_1^T x_2 \\ y_2 y_1 x_2^T x_1 & y_2 y_2 x_2^T x_2 \end{bmatrix}$$

$$\begin{bmatrix} d_1 & d_2 \end{bmatrix}_{1 \times 2} \begin{bmatrix} y_1 y_1 x_1^T x_1 & y_1 y_2 x_1^T x_2 \\ y_2 y_1 x_2^T x_1 & y_2 y_2 x_2^T x_2 \end{bmatrix}_{2 \times 2} \begin{bmatrix} d_1 \\ d_2 \end{bmatrix}_{2 \times 1} = \begin{bmatrix} d_1 y_1 y_1 x_1^T x_1 + d_2 y_2 y_1 x_2^T x_1 & d_1 y_1 y_2 x_1^T x_2 + d_2 y_2 y_2 x_2^T x_2 \end{bmatrix} \begin{bmatrix} d_1 \\ d_2 \end{bmatrix}$$

$$= d_1^2 y_1 y_1 x_1^T x_1 + d_1 d_2 y_2 y_1 x_2^T x_1 + d_2 d_1 y_1 y_2 x_1^T x_2 + d_2^2 y_2 y_2 x_2^T x_2 = \sum_{i=1}^2 \sum_{j=1}^2 d_i d_j y_i y_j x_i^T x_j$$

• Hence $\sum_{i=1}^2 \sum_{j=1}^2 d_i d_j y_i y_j x_i^T x_j = d^T Q d$

• $d^T d = \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} d_1 \\ d_2 \end{bmatrix} = d_1 + d_2 = \sum_{i=1}^2 d_i$

• $y^T d = \begin{bmatrix} y_1 & y_2 \end{bmatrix} \begin{bmatrix} d_1 \\ d_2 \end{bmatrix} = \sum_{i=1}^2 y_i d_i = 0$

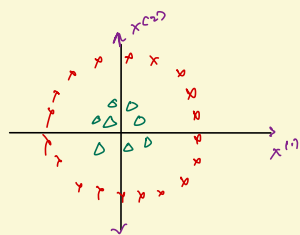
• Plane. $w^* = \sum_{i=1}^N d_i^* y_i^* x_i^*$

$$b^* = y^* - w^{*T} x^*$$

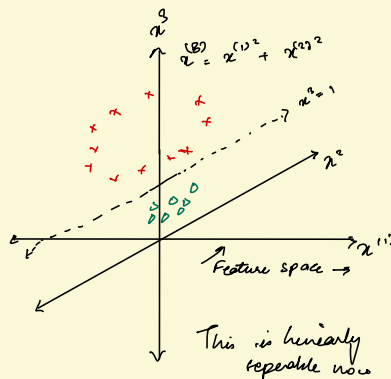
Plane: $w^* x + b^*$

$$= \sum_i \alpha_i y_i \langle \vec{x}_i, \vec{x} \rangle + b$$

Non-linear SVMs



Adding a
new dimension
 $h(x^{(1)}, x^{(2)})$



Idea: Project the data into a higher dim space
so that it's linearly separable in the higher dimension

- > Project data to a higher dim where linearly separable
- > Apply linear SVM to the higher dimension

Aim: we just need the inner products of the transformation.

$$\max_{\alpha} \quad 1^T \alpha - \frac{1}{2} \alpha^T Q \alpha$$

$$\text{Sub. to} \quad y^T \alpha = 0 \\ \alpha \geq 0$$

$$Q = [Q_{ij}]_{n \times n}$$

$$[Q_{ij}]_{n \times n} = [y y^T]_{n \times n} \odot \begin{bmatrix} h(x_1)^T h(x_1) & \dots & h(x_1)^T h(x_n) \\ \vdots & \ddots & \vdots \\ h(x_n)^T h(x_1) & \dots & h(x_n)^T h(x_n) \end{bmatrix}$$

$$\phi: x \rightarrow \mathbb{R}$$

$$\phi(x_i, x_j) = h(x_i)^T h(x_j)$$

Feature space (higher dimensional space)

- First idea was: Learning a non-linear classifier using SVM
- Calculate $h(x)$ for each training example – Find a linear SVM in the feature space.

• Problems:

- Feature space can be high dimensional or even have infinite dimensions, so storing $h(x)$ can be inefficient or even impossible

So we came up with the idea of “Kernels”:

– Kernels are similarity functions that return inner products between the images $h(x_i)$ of data points.

- Choosing kernel (ϕ) is equivalent to choosing $h(x)$ in the feature space.

The kernel trick

- No need to know what $h(x)$ is and what the feature space is.
- No need to explicitly map the data to the feature space. (memory constraint resolved)
- Define a kernel function ϕ and replace the dot product $\langle x, z \rangle$ with a kernel function $\phi(x, z)$ in both training and testing.

Plane: $\omega^* = \sum_{i=1}^N \alpha_i^* y_i h(x_i)$
 $b^* = y_i - \omega^{*T} x_i$
 $\omega^* x + b^*$

$$f(\vec{x}) = \sum_i \alpha_i y_i \Phi(\vec{x}_i, \vec{x}) + b$$

$$\Phi(x_i, x) = h(x_i)^T h(x)$$

Define a function Φ

def phi(x, z):

return np.dot(x, z)

Q = np.dot(y, y.T) +

phi(np.dot(x, x.T))

(a) Linear kernel:

$$\Phi(x, z) = [\alpha x, \alpha z], \alpha \in \mathbb{R}$$

$$\therefore \Phi(x, z) = [h(x), h(z)] = (h(x))^T h(z)$$

$$= [\alpha x_1, \dots, \alpha x_n] \begin{bmatrix} \alpha z_1 \\ \vdots \\ \alpha z_n \end{bmatrix} = \alpha^2 x_1 z_1 + \alpha^2 x_2 z_2 + \dots + \alpha^2 x_n z_n$$

(b) Polynomial kernel:

$$\Phi(x, z) = (x^T z + 1)^P$$

$$\Phi(x, z) = [h(x), h(z)] = (x^T z + 1)^P$$

$$= (x_1 z_1 + x_2 z_2 + \dots + x_n z_n + 1)^P$$

(c) RBF kernel:

$$\Phi(x, z) = \exp(-\gamma \|x - z\|^2) = [h(x), h(z)]$$

$$\Phi(x, z) = e^{-\gamma \left(\sum_{i=1}^N (x_i - z_i)^2 \right)}$$

Q7

logistic Regression for $K=2$:

of samples = m

of features = n

Cost Func

$$J(w, b) = -\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^n y^{[i]} \log(f_{w,b}(x_j^{[i]}) + (1-y^{[i]}) \log(f_{w,b}(x_j^{[i]}))$$

$$\frac{\partial J}{\partial w_j} \text{ \& \& } \frac{\partial J}{\partial b}$$

$$w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix}$$

BINARY
CLASSIFICATION

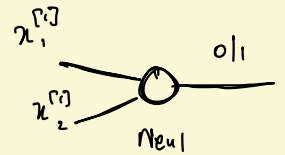
GD Algo:

repeat

$$\{ w_j = w_j - \alpha \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{[i]}) - y^{[i]}) x_j^{[i]}$$

$$b = b - \alpha \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{[i]}) - y^{[i]})$$

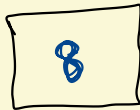
?



K=2

K=10

Eg: for digit recognition:

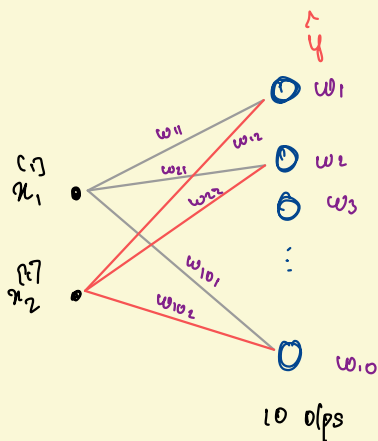


\hat{y}
Pred. o/p.

$$x = \begin{bmatrix} 0.1 \\ 0.0 \\ 0.0 \\ 0.0 \\ 0.0 \\ 0.05 \\ 0.0 \\ 0.8 \\ 0.05 \end{bmatrix}$$

$$y = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

Replace the
Sigmoid function by the softmax function



One Single Sample

$$w = \begin{bmatrix} w_{11} & w_{12} \\ \vdots & \vdots \\ w_{101} & w_{102} \end{bmatrix}_{10 \times 2}$$

$$x = \begin{bmatrix} x_1^{[i]} \\ x_2^{[i]} \end{bmatrix}_{2 \times 1} \quad b = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_{10} \end{bmatrix}$$

$$z^{[i]} = w x + b$$

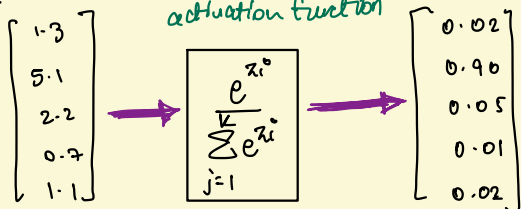
$$z^{[i]} = \begin{bmatrix} z_1^{[i]} \\ z_2^{[i]} \\ \vdots \\ z_{10}^{[i]} \end{bmatrix}$$

$$a_j^{[i]} = \text{softmax}(z^{[i]})$$

$$a_j^{[i]} = \frac{e^{z_j^{[i]}}}{\sum_{j=1}^K e^{z_j^{[i]}}}$$

$$a^{[i]} = \begin{bmatrix} a_1^{[i]} \\ a_2^{[i]} \\ \vdots \\ a_{10}^{[i]} \end{bmatrix}$$

Single sample



$$L = \text{loss}(a_1, \dots, a_N, y) = \begin{cases} -\log a_1, & y=1 \\ -\log a_2, & y=2 \\ \vdots & \vdots \\ -\log a_N, & y=N \end{cases}$$

Loss for a single sample :

$$\sum_{j=1}^K -\log(a_j^{[i]}) y_j^{[i]}$$

Cost function:

$$J(w, b) = -\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^K (-\log(a_j^{[i]})) y_j^{[i]} \quad \text{for } m \text{ samples}$$

Gradient Descent for this \uparrow Refer to the Neural Networks notes for it Derived it then.

MCE: Idea: Is to minimize misclassification

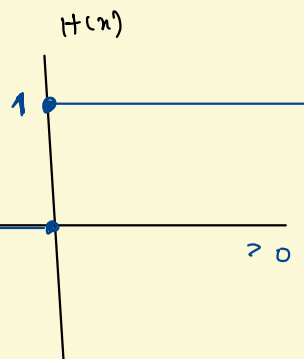
K=2

$$-y_i w^T x_i = \begin{cases} \geq 0, & \text{misclassification} \\ < 0, & \text{correct classification.} \end{cases}$$

$$E_0(w) = \sum_{i=1}^m H(-y_i w^T x_i)$$

\rightarrow counts the no. of misclassification

$y_i = 0 \text{ or } 1$



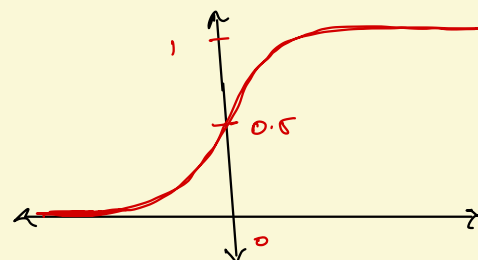
So if we minimize $E_0(w)$ we're done, but $\nabla H = 0$ everywhere except @ origin

Replace heaviside by s.g

$$E_1(w) = \sum_{i=1}^m \ell(-y_i w^T x_i)$$

$$E_1(w) = \sum_{i=1}^m \frac{1}{1 + e^{-y_i w_1 x_1^{[i]} - y_i w_2 x_2^{[i]} \dots - y_i w_n x_n^{[i]}}}$$

\rightarrow Probability of misclassification.



N features.

$$\frac{\partial E_1}{\partial w_j} = \sum_{i=1}^m y_i \ell(-y_i w^T x_i) (1 - \ell(-y_i w^T x_i)) x_j^{[i]}$$

Gradient Descent

repeat

$$\{ w_j = w_j - \sum_{i=1}^m y_i d(y_i w^T x^{(i)}) (1 - d(y_i w^T x^{(i)})) x_j^{(i)}$$

$$b = b - \sum_{i=1}^m y_i d(y_i w^T x^{(i)}) (1 - d(y_i w^T x^{(i)}))$$

}

K=2

$\hat{y} =$
Pred. o/p.

$$\begin{bmatrix} 0.1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0.05 \\ 0 \\ 0.8 \\ 0.05 \end{bmatrix}$$

$y =$

$$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

149

K classes

For each class k the

$\hat{y} = y \begin{cases} \geq 0.5, & \text{correctly classified} \\ < 0.5, & \text{incorrectly classified} \end{cases}$

$$E(w) =$$

Q6 We normalize all the examples (x_i, y_i) so that $\|x_i\| \leq 1 \quad \forall i$

Why Margin in Perceptron Algorithm? We iterate through the data using perceptron algorithm updating the weights not only on mistakes but also on examples x that the current hypothesis gets correct by a margin $< \frac{\gamma}{2}$ (so if an ex is within margin $\frac{\gamma}{2}$ then it's incorrect)

predict +ive $\frac{w_t^T x_i}{\|w_t\|} \geq \frac{\gamma}{2}$, negative if $\frac{w_t^T x_i}{\|w_t\|} < -\frac{\gamma}{2}$

$\frac{y_i w_t^T x_i}{\|w_t\|} < \frac{\gamma}{2}$ then it's a mistake

$$\|w_{t+1}\|^2 \leq \|w_t\|^2 + \|x_i\|^2 + 2y_i w_t^T x_i$$

$$\begin{aligned} \|w_{t+1}\|^2 &\leq \|w_t\|^2 + \|x_i\|^2 + 2y_i w_t^T x_i \\ &\leq \|w_t\|^2 + 2w_t^T x_i + \|x_i\|^2 \end{aligned}$$

$$\leq \|w_t\|^2 \left(1 + \frac{1}{\|w_t\|} \frac{w_t^T x_i}{\|w_t\|} + \frac{\|x_i\|^2}{\|w_t\|^2} \right)$$

$$\|w_{t+1}\| \leq \|w_t\| \left(1 + \underbrace{\frac{2}{\|w_t\|} \frac{\gamma}{2} + \frac{1}{\|w_t\|^2}}_{\geq 0} \right)^{1/2}$$

(for $x > 0 \quad \sqrt{1+x} < 1 + \frac{x}{2}$) \hookrightarrow Verify by MVT

$$\|w_{t+1}\| \leq \|w_t\| \left(1 + \frac{\gamma}{2\|w_t\|} + \frac{1}{2\|w_t\|^2} \right)$$

Hence $\|w_{t+1}\| \leq \|w_t\| + \frac{1}{2\|w_t\|} + \frac{\gamma}{2}$

If the classification is $\frac{y_i^* \omega_t^T x_i^*}{\|\omega_t\|} < \frac{\gamma}{2}$ it's a mistake

$$\frac{\|\omega_t\|}{y_i^* \omega_t^T x_i^*} > \frac{2}{\gamma}$$

If for some i we take $y_i^* \omega_t^T x_i^* = 1$ then, $\|\omega_t\| > \frac{2}{\gamma}$
($\frac{1}{\|\omega_t\|} < \frac{\gamma}{2}$)

$$\|\omega_{t+1}\| \leq \|\omega_t\| + \frac{\gamma}{4} + \frac{\gamma}{2} = \|\omega_t\| + \frac{3\gamma}{4}$$

Repeating it recursively after m mistakes

$$\|\omega_{m+1}\| \leq \|\omega_t\| + \frac{3M\gamma}{4} < \frac{2}{\gamma} + \frac{3M\gamma}{4}$$

Also for a perceptron algo recall

$$M\gamma \leq \|\omega_{m+1}\|$$

$$\text{So, } M\gamma \leq \|\omega_{m+1}\| < \frac{2}{\gamma} + \frac{3M\gamma}{4}$$

$$M\gamma \leq \frac{2}{\gamma} + \frac{3M\gamma}{4}$$

$$\Rightarrow \frac{M\gamma}{4} \leq \frac{2}{\gamma}$$

$$\Rightarrow \boxed{M \leq \frac{8}{\gamma^2}}$$

Q5

$$H(x) = - \int p(x) \ln p(x)$$

$$\int p(x) dx = 1$$

$$\int p(x) x dx = \mu$$

$$\int p(x) (x - \mu) (x - \mu)^T dx = \Sigma$$