

[Piazza] please answer.

(59)

Recap :- Logistic Regression.

• $l(\bar{w}^*, T) \rightarrow$ Log-likelihood of training samples under LR.

$$\nabla_{\bar{w}^*} l = \begin{bmatrix} \frac{\partial}{\partial \bar{w}_0^*} l(\bar{w}^*, T) \\ \vdots \\ \frac{\partial}{\partial \bar{w}_n^*} l(\bar{w}^*, T) \end{bmatrix}$$

$$\frac{\partial l(\bar{w}^*, T)}{\partial \bar{w}_j^*} = \sum_{i=1}^n \bar{x}_j(i) \cdot \left[f(\bar{x}(i)) - \Pr(f(\bar{x}(i)) = 1 \mid \bar{x}(i), \bar{w}^*) \right]$$

- Compare more change for probability of label 'small'.

(-) Discussion session :- (next)

LR vs Perceptron vs Naive Bayes vs SVM

- What are the positives and negatives.

Support Vector Machines.

* references here

- Vladimir Vapnik - Founder(s) of Statistical learning theory.

1992 :- accuracy comparable to NN on 'image' data sets.

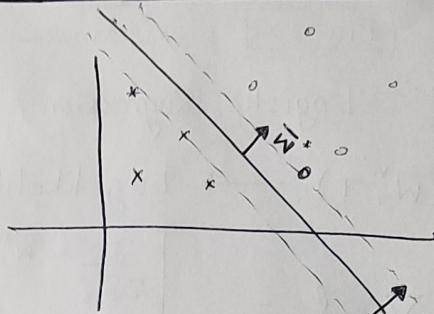
Now :- everywhere ...

Agenda :- 1. Maximum margin classifiers.

2. Non-linear classifiers & Kernel methods.

Maximum-margin classifiers.

Suppose the training data is linearly separable.



Margin of a perfect classifier for the data is the "maximum width" that the boundary could be increased before we hit a data point.

Learn a "large" maximum margin classifier.

- better generalization

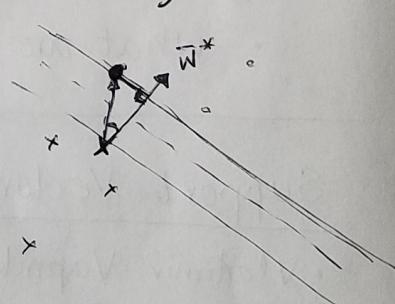
How to find the MMC?

Recall :- decision rule :- $\begin{cases} \bar{x}^* \cdot \bar{w}^* > 0 \rightarrow +1 \\ \bar{x}^* \cdot \bar{w}^* < 0 \rightarrow -1 \end{cases}$

Appropriately normalize

For $f(\bar{x}(i)) = 1, \bar{x}^*(i) \cdot \bar{w}^* \geq 1$

For $f(\bar{x}(i)) = -1, \bar{x}^*(i) \cdot \bar{w}^* \leq -1$.



Suppose, i, i' are such that,

$$\left. \begin{array}{l} \bar{x}^*(i) \cdot \bar{w}^* = 1 \\ \bar{x}^*(i') \cdot \bar{w}^* = -1 \end{array} \right\} \text{distance} = \frac{(\bar{x}^*(i) - \bar{x}^*(i')) \cdot \bar{w}^*}{\|\bar{w}^*\|_2} = \frac{(\bar{x}(i) - \bar{x}(i')) \cdot \bar{w}}{\|\bar{w}\|_2} = \frac{2}{\|\bar{w}\|_2}$$

Formulation :-

$$\boxed{\text{Maximize } \frac{2}{\|\bar{w}\|_2}} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{subject to}$$

(64)

$$f(\bar{x}(i)) \cdot (\bar{x}(i) \cdot \bar{w} - w_0) \geq 1.$$

$$\text{Minimize } \frac{1}{2} \|\bar{w}\|_2^2, \text{ s.t.}$$



Quadratic objective function, and linear constraints.

Lagrangian :-

$$\text{minimize } L_p(\bar{w}, w_0, \alpha_i) = \frac{1}{2} \|\bar{w}\|_2^2 - \sum_{i=1}^n \alpha_i \cdot (f(\bar{x}(i)) \cdot (\bar{x}(i) \cdot \bar{w} - w_0) - 1)$$

$$\text{and } \alpha_i \geq 0$$

$$\frac{\partial L_p}{\partial \bar{w}} = \bar{0} \Rightarrow \bar{w} - \sum_{i=1}^n \alpha_i \cdot f(\bar{x}(i)) \cdot \bar{x}(i) = \bar{0}$$

$$\Rightarrow \bar{w} = \sum_{i=1}^n \alpha_i \cdot (f(\bar{x}(i)) \cdot \bar{x}(i)) - \bar{0}$$

④ KKT conditions: (complementary slackness) [Discussion]

At optimum:- $\alpha_i \cdot f(\bar{x}(i)) \cdot (\bar{x}(i) \cdot \bar{w} - w_0 - 1) = 0$ for all i .

∴ All ~~not~~ samples not at the boundary (support) have $\alpha_i = 0$

SV :- support vectors.

$$\bar{w} = \sum_{i \in SV} \alpha_i f(\bar{x}(i)) \cdot \bar{x}(i)$$

↳ of things.

what if the data is not linearly classifiable?

Add slack variables for misclassification :-

PRIMAL } minimize $\frac{1}{2} \|\vec{w}\|_2^2 + C \sum_{i=1}^n \xi_i$, such that

$$f(\vec{x}(i)) \cdot (\vec{x}^*(i) \cdot \vec{w}^*) \geq 1 - \xi_i, \quad \xi_i \geq 0.$$

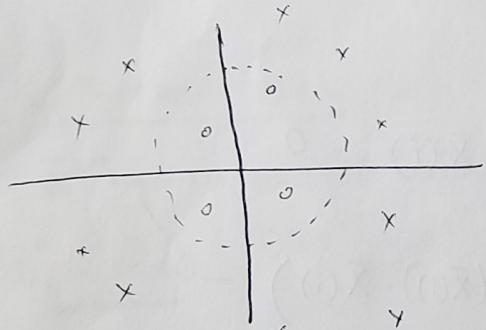
$C \rightarrow$ overfitting parameter

Non-linear SVM :- Kernel trick.

GO TO
DOT PROD

MNIST data-set :-
Linear SVM :- 8.5% error
Polynomial SVM :- 1% error

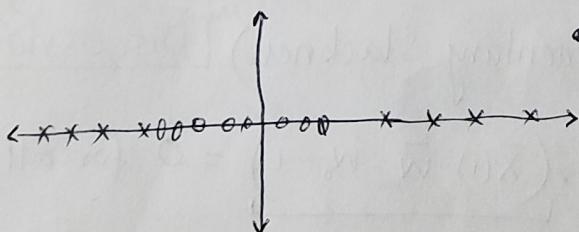
Linear classifiers are not complicated enough sometimes.



- No linear classifier would be great.
- Need non-linear decision regions.

→ obtaining this is hard

- ① ~~to~~ linear SVM's very well understood.



General idea :-

- ① Map the features to a higher dimensional space.
- ② Apply a linear classifier in the new space.

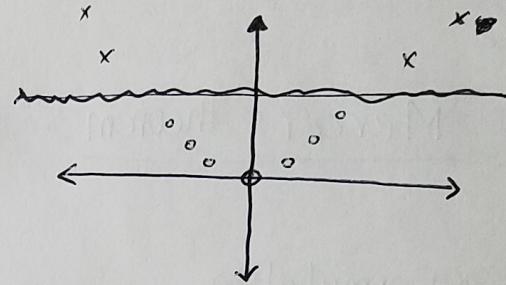
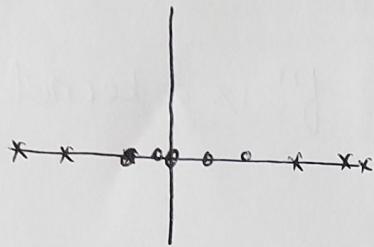
$\phi: \mathbb{R}^k \rightarrow \mathbb{R}^m$, ($m > k$) is the function mapping the features to the new space.

Example :-

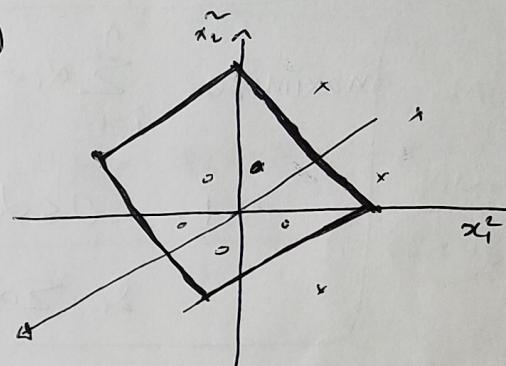
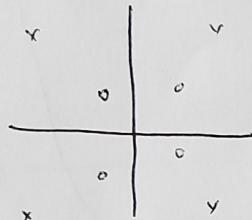
$$x \rightarrow (x, x^2)$$

$$K=1, m=2.$$

(63)



$$(x_1, x_2) \rightarrow (x_1^2, x_2^2, x_1 x_2)$$



The kernel trick :-

\bar{w}^* → weight vector after training the SVM in the 'new' space, (i.e, in m dimensions)
 $= (\bar{w}_0, \bar{w}_1, \dots, \bar{w}_m)$

So, the function for classification is now :-

$$g(\bar{x}) = \bar{w}^* \cdot \phi(\bar{x}) + \bar{w}_0$$

$$= \left[\sum_{i \in SV} \alpha_i \phi(\bar{x}(i)) \right] \cdot \phi(\bar{x}) + \bar{w}_0$$

$$= \sum_{i \in SV} \alpha_i \underbrace{\phi(\bar{x}(i))^T \cdot \phi(\bar{x})}_{\phi(\bar{x}(i))^\top \phi(\bar{x})} + \bar{w}_0$$

We do not need to know the explicit mapping, we only use the dot-product of feature vectors in training & test.

Kernelfunction :- $K(\bar{x}(i), \bar{x}(j)) = \phi(\bar{x}(i))^T \cdot \phi(\bar{x}(j))$

Mercer's Theorem :- Every ^{+semidefinite} f^n is a kernel.

• Formulation

DUAL OF PRIMAL ON P 62

$$\begin{aligned} & \text{maximize}_{\alpha_i} \quad \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum \alpha_i \alpha_j f(\bar{x}(i)) f(\bar{x}(j)) K(\bar{x}(i), \bar{x}(j)) \\ & \text{s.t.} \quad 0 \leq \alpha_i \leq C \\ & \quad \& \sum \alpha_i f(\bar{x}(i)) = 0 \end{aligned}$$

• Linear Kernel :- $K(\bar{x}(i), \bar{x}(j)) = \bar{x}(i)^T \cdot \bar{x}(j)$

Poly " :- $K(\bar{x}(i), \bar{x}(j)) = (1 + \bar{x}(i)^T \bar{x}(j))^P$

Gaussian / RBF :- $K(\bar{x}(i), \bar{x}(j)) = \exp\left(-\frac{\|\bar{x}(i) - \bar{x}(j)\|_2^2}{2\sigma^2}\right)$

DOT PRODUCT :-

Lagrange Dual problem :-

$$\text{maximize} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j f(\bar{x}(i)) f(\bar{x}(j)) \cdot \boxed{\bar{x}(i)^T \cdot \bar{x}(j)}$$

subject to $0 \leq \alpha_i \leq C$

$$\sum_{i=1}^n \alpha_i g_i f(\bar{x}(i)) = 0$$

You work only with the dot product of the vectors, nothing else.