



2 UNIT II

Unit - 2 - II



SPARSE KERNEL MACHINES

Introduction :

Limitation of algorithms based on non linear kernels is that the kernel functions $K(x_n, x_m)$ must be evaluated for all possible pairs x_n and x_m of training points.

Here we will take up kernel based algorithms that have sparse solutions so that prediction for new inputs depend only on the kernel functions evaluated at a subset of the training data points.

Mapping the data onto higher dimensions is helpful in creating hyperplane based classifiers. In such classifiers we can reduce the data sets to support vectors as feature vectors.
Example: $w_1 = 1, -1, 0.5, -0.5, \dots$, $w_2 = -4, 4, -3, 3, -2, 2$

Maximal Margin Classifiers :

Important property of SVM :

Determination of Model parameters correspond to a convex optimisation problem and hence any local solution is also a global optimum. (Connect this to above example).

Note : SVM is a decision machine and hence does not provide posterior probabilities (Alternate Sparse Kernel Machine is Relevance Vector Machine (RVM))

① What is SVM?

SVM is a supervised Learning Algorithm which can be used for classification and Regression problems as Support Vector Classification (SVC) and Support Vector Regression (SVR). Generally it is used for small data sets.



Ideology behind SVM :



SVM is based on the idea of finding a hyperplane that separates the features into different domains.

Terminologies used in SVM .

Support vector points : These are points closest to the hyperplane

Margins : The distance of the vectors from the hyperplane are called margins

Margin maximising hyperplanes :

The hyperplane is called as Margin maximising hyperplane.

Hyper plane : (Decision Surface)

The hyperplane is a function which is used to differentiate between features. In 2D hyperplane is a line & in 3D is a plane and the function which classifies in higher dimension is called a hyper plane.

Hard Margin :

The equations $\pi^- : b + w^T x = 0$, $\pi^+ : b + w^T x = 1$ and $\pi^- : b + w^T x = -1$ are very strict constraints to correctly classify each and every data point. Hence they are called Hard Margins.

- ② Derive the Maximum Margin Classifier expression for binary classification of points in SVM

(OR)

Derive an expression for a Hyperplane that classifies the features of Data for SVM classifier.



Let there be m dimensions

The equation of the hyperplane in M dimension is given by

$$y = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_m x_m$$

$$= w_0 + \sum_{i=1}^m w_i x_i$$

$$= w_0 + w^T x \quad \text{where } w = (w_1, w_2, \dots, w_m)^T$$

$$x = x_1, x_2, \dots, x_m$$

$$= b + w^T x \quad \text{① where } b \text{ is the bias term}$$

If $\Phi(x)$ denotes the feature space transformation, then SVM for 2-class classification problem using linear models takes the form

$$y(x) = w^T \Phi(x) + b \quad \text{②}$$

Let there be N training data x_1, x_2, \dots, x_N

Let the target values be t_1, t_2, \dots, t_N where $t_n \in \{-1, 1\}$.

Assuming the data is linearly separable $\exists w$ and $b \ni$:

$$\begin{cases} y(x_n) > 0 \text{ for points having } t_n = 1 \\ y(x_n) < 0 \quad " \quad " \quad \ni t_n = -1 \end{cases} \quad \text{③}$$

SVM approaches thro' a concept of the MARGIN (defined to be smallest distance between decision boundary and any of the samples) and this margin needs to be maximised.

The L_2 distance of a pt x from the hyperplane

defined by $y(x) = 0$ as in ② is given by $\frac{|y(x)|}{\|w\|}$

Clearly $t_n y(x_n) > 0 \quad \forall n$ from the view of ③

So, Distance from x_n to the decision surface is

$$\frac{t_n y(x_n)}{\|w\|} = \frac{t_n [w^T \Phi(x_n) + b]}{\|w\|}$$

$$\textcircled{4} \quad \frac{w_1 x_1 + w_2 x_2 + w_3}{\sqrt{w_1^2 + w_2^2}}$$

$$w_1 x_1 + w_2 x_2 + w_3 = 0$$

L_2 distance from

$$(x_1, y_1) :$$

$$\frac{w_1 x_1 + w_2 x_2 + w_3}{\sqrt{w_1^2 + w_2^2}}$$

The Margin is given by $\frac{1}{\|w\|}$ distance to closest x_n from data sets.



Claim: Optimise w and b (parameters) to maximise this distance

Hence Maximum Margin is given by

$$\arg \max_{w, b} \left\{ \frac{1}{\|w\|} \min_{x_n} \left[\ln(w^T \phi(x_n) + b) \right] \right\} \quad (5)$$

as w is independent of n , $\frac{1}{\|w\|}$ is taken out of "Min".

If we re-scale $w \rightarrow kw$ and $b \rightarrow kb$ then the distance of DS from x_n $\left[\ln(w^T \phi(x_n) / \|w\|) \right]$ remains unchanged

∴ We can set

$$\ln(w^T \phi(x_n) + b) = 1 \quad (6)$$

In this case all data points satisfy

$$\ln(w^T \phi(x_n) + b) \geq 1 \quad \forall n = 1 \dots N \quad (7)$$

(7) is called Canonical Representation of Hyperplane.

\Rightarrow Maximise $\|w\|^{-1} \Leftrightarrow$ Minimise $\|w\|^2$ and we have an optimization problem

$$\arg \min_{w, b} \frac{1}{2} \|w\|^2 \quad (8)$$

Subjected to (7), a quadratic programming problem

Note :-

① (8) is a quadratic programming problem. in which we need to minimise a quadratic function subject to a set of linear inequality constraints.



Note ②

In order to solve constrained optimization problem, we use Lagrangian multipliers.



$a_n > 0$ with one multiplier for each constraint in ⑦ :

$$L[w, b, a] = \frac{1}{2} \|w\|^2 - \sum_{n=1}^N a_n \{ [\epsilon_n (w^T \phi(x_n) + b) - 1] \}$$
(9)

where $a = (a_1, a_2, \dots, a_N)^T$

$$\frac{\partial L}{\partial w} = 0 \Leftrightarrow w = \sum_{n=1}^N a_n \epsilon_n \phi(x_n)$$
(10)

$$\frac{\partial L}{\partial b} = 0 \Leftrightarrow b = \sum_{n=1}^N a_n \epsilon_n$$
(11)

Note ③ Eliminating w and b from $L[w, b, a]$ using these conditions it gives dual representation of Max Margin problem.

$$L(a) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m \epsilon_n \epsilon_m K(x_n, x_m)$$
(12)

Subjected to constraints

$$a_n \geq 0 \quad \forall n = 1 \text{ to } N$$
(13)

$$\sum_{n=1}^N a_n \epsilon_n = 0$$
(14)

$$\text{Where } K(x, x') = [\phi(x)]^T [\phi(x')]^T$$

Note ④ Solution of quadratic programming problem in M variables has $O(M^3)$ complexity.

Note ⑤ : Formulate the Soft Margin Technique of SVM ?



The given data may not be linearly separable.

⇒ we need an updation so that one function may skip few outliers and classify it.

Hence Introduce a slack variable ξ_i so that

$$y_i (w^T x_i + b) \geq 1 - \xi_i \quad \forall i \quad (15)$$

If $\xi_i = 0$, the points can be considered correctly classified
 $\xi_i > 0$ incorrectly classified

Then ξ_i can be considered as error associated with x_i

$$\text{Average error} = \frac{1}{n} \sum_{i=1}^n \xi_i \quad (16)$$

Thus Mathematical formulation of Soft Margin SVM becomes

$$\underset{w, b}{\text{Minimise}} \quad \frac{1}{2} \|w\|^2 + \sum_{i=1}^n \xi_i \quad (17)$$

$$\text{Subject to } y_i [w^T x_i + b] \geq 1 - \xi_i \quad \forall i = 1, 2, \dots, n$$

Formulation : (Soft Margin Technique)

To find the vector w and the scalar b such that the hyperplane represented by w and b maximises the Margin distance and minimises the loss term subjected to the condition that all points are correctly classified.

NOTE ⑧ Interpret Loss function of SVM?

Loss function Interpretation of SVM

We have $z_i = y_i (w^T x_i + b) \longmapsto ①$

$$z_i > 1 \quad ②$$

12

If $z_i > 1$ then Loss is '0'
 $z_i \leq 1$ then Loss increases



Hence

$$\text{Hinge loss} = \max(0, 1 - z_i)$$



NOTE (9) What is dual form of SVM?

Dual form of SVM :

Consider a data set which is not at all linearly separable.

In logistic regression separation of data is by projecting it on to higher dimension by ADDING RELEVANT FEATURES But in SVM the projection is done on a powerful way which is called 'PRIMAL FORM'

An alternative method is the Dual form of SVM that uses Lagrange's Multiplier Method.

$$\underset{\alpha}{\text{Maximise}} \left[\sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j (x_i^T x_j) \right] \quad \rightarrow (18)$$

Subject to $\alpha_i \geq 0 \quad \forall i = 1 \text{ to } n$

$$\sum_{i=1}^n \alpha_i y_i = 0 \quad \left. \right\} \quad \rightarrow (19)$$

Remark: If $\alpha_i > 0$, then x_i is the support vector

(1) $\alpha_i = 0$, " " is NOT support vector

(2) To solve actual problem we do not require actual data. Instead dot product of all pairs sufficient

(3) Calculation of ' b ' also requires dot product only

(4) Major advantage of Dual SVM is that it depends only on ' α '.



Note (10) What is Kernel Trick?



The Kernel is a way of computing the dot product of two vectors x and y in some (high dimensional) feature space, [Sometimes it is called 'Generalised dot product'].

Note (11) Write a short note on various types of Kernels used in SVM?

| S. NO | Type of Kernel | Expression | Remarks. |
|-------|--|---|---|
| 1. | Linear Kernel | | |
| 2. | Polynomial | $K(x_1, x_2) = (a + x_1^T x_2)^b$ | b - deg of Kernel a = const term. |
| 3. | Radial Basis function. (RBF) / Gaussian Kernel. | $K(x_1, x_2) = \exp(-\gamma \ x_1 - x_2\ ^2)$ | $\ x_1 - x_2\ \rightarrow$ Euclidean distance b/w x_1 & x_2 |

Parameters :

- ① C : Inverse of strength of regularisation
 - Behaviour: As C increases the model gets overfit
 - As C decreases the model gets underfit
- ② γ : Used only for RBF Kernel
 - Behaviour: As γ increases the model gets overfit
 - As γ decreases the model gets underfit



Pros and Cons of SVM:

Pros :

- (1) Really effective in higher dimension
- (2) Effective when the number of features more than training examples
- (3) Best Algorithm when classes are separable
- (4) Hyperplane is affected by support vectors and thus outliers have less impact
- (5) Suited for extreme case binary classification

Cons :

- (1) computationally expensive for larger data sets
- (2) Does not perform well in case of overlapped classes
- (3) Selecting the appropriate kernel function is tricky
- (4) Selection of appropriate hyperparameters requiring sufficient generalisation performance

Preparing data for SVM:

Numerical Conversion : As inputs are numerical instead of categorical, Most commonly used are 'one hot encoding' Label-encoding etc.

Binary Classification :

One Versus the rest method / One versus one method.

Overlapping Class Distributions :



① What is the work of regularisation parameter C in SVM?
 (in box constraints)

It controls the trade-off between minimising training errors and controlling model complexity.

② State the -Kuhn-Tucker conditions of the Lagrangian formulation of SVM and derive the box constraints for overlapping class distributions?

Based on the assumption that training data are linearly separable in feature space $\phi(x)$, the resulting SVM gives exact separation of the training data in original input space x , although corresponding decision boundary is non-linear. But in practice, the class-conditional distributions may overlap. Hence our goal is to maximise the margin while softly penalising points that lie on the wrong side of the margin boundary.

$$\text{Min}(w) = \frac{1}{2} \|w\|^2 + C \sum_{n=1}^N \text{leg}_n \quad \xrightarrow{\textcircled{1}}$$

Where C controls the trade off between the slack variable penalty and the Margin.

The Corresponding Lagrangian formulation is given by

$$L[w, b, a] = \frac{1}{2} \|w\|^2 + C \sum_{n=1}^N \text{leg}_n - \sum_{n=1}^N a_n \{ \text{leg}_n (y_n a_n - 1) + \text{leg}_n \} - \sum_{n=1}^N \mu_n \text{leg}_n \quad \xrightarrow{\textcircled{2}}$$

Where $a_n, \mu_n > 0$ are Lagrangian multipliers.



The KKT conditions are given by



$$a_n > 0 \quad (3)$$

$$b_n y(x_n) - 1 + \xi_n > 0 \quad (4)$$

$$a_n (b_n y(x_n) - 1 + \xi_n) = 0 \quad (5)$$

$$\lambda_n > 0 \quad (6)$$

$$\xi_n > 0 \quad (7)$$

$$\mu_n \xi_n = 0 \quad (8) \quad \forall n = 1 \dots N$$

For optimising,

$$\frac{\partial L}{\partial w} = 0 \Leftrightarrow w = \sum_{n=1}^N a_n x_n + b(x_n) \quad (9)$$

$$\frac{\partial L}{\partial b} = 0 \Leftrightarrow \sum_{n=1}^N a_n x_n = 0 \quad (10)$$

$$\frac{\partial L}{\partial \xi_n} = 0 \Leftrightarrow a_n = C - \mu_n \quad (11)$$

Eliminating w , b and $\{\xi_n\}$ from the Lagrangian,

$$L(a) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m x_n x_m K(x_n, x_m) \quad (12)$$

Now Minimise (12) wrt dual variables $\{a_n\}$, subject to

$$0 \leq a_n \leq C \quad (13)$$

$$\sum_{n=1}^N a_n x_n = 0 \quad (14) \quad \text{where } n = 1, 2, \dots, N$$

(13) are known as 'BOX CONSTRAINTS'.

- (3) Express the 2nd order polynomial Kernel in terms of its components namely vector function and feature space.

SOL:

$$\begin{aligned}
 K(x, z) &= (1 + x^T z)^2 \\
 &= (1 + \gamma_1 z_1 + \gamma_2 z_2)^2 \\
 &= 1 + 2\gamma_1 z_1 + 2\gamma_2 z_2 + \gamma_1^2 z_1^2 + 2\gamma_1 \gamma_2 z_1 z_2 + \gamma_2^2 z_2^2 \\
 &= (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, \sqrt{2}x_1 x_2, x_2^2) (1, \sqrt{2}z_1, \sqrt{2}z_2, z_1^2, \\
 &\quad \sqrt{2}z_1 z_2, z_2^2)^T \\
 &= \phi(x)^T \phi(z)
 \end{aligned}$$

(4) How do you relate SVM with logistic regression?

They can be related by their error functions given by

$$E_{SV}(y_n t_n) = (1 - y_n t_n)_+$$

approximation to the misclassification error in SVM.

In logistic regression, with $p[t = \frac{1}{y}] = \sigma(y)$ and $p[t \neq y] = \sigma(-y)$. For this error function is given by

$$\sum_{n=1}^N E_{LR}(y_n t_n) + \lambda \|w\|^2$$

(5) Write a short note on Multiclass SVM

Commonly used approach focuses constructing K-Separate SVMs in which Kth model $y_k(x)$ is trained using the C_{k-1} class data as positive examples and the data from remaining $K-1$ classes as negative examples. This is one-verses-the-rest approach. Two disadvantages are: 1. Individual classifiers can lead to inconsistent results in which an input is assigned to multiple classes simultaneously. e.g. Due to imbalanced training data the symmetry of original problem could be lost.

Another approach (one-verses-one) is to train $K(K-1)/2$ different 2 classes SVMs on all possible pairs of classes, and then to classify test points according to which class has the highest number of votes. In this approach more amount of significant computation is required.



Regression
problem

⑥ Explain how SVM can be extended to Regression problems? (OR) Write down

the steps involved in converting the general SVM to be applied for a

By preserving the property of Sparseness one can extend SVM to regression problems :-

Step 1 : Introduction of ϵ -intensive error function
having a linear cost associated with errors outside the insensitive region as

$$E_{\epsilon}[y_n - \epsilon] = \begin{cases} 0 & |y(x) - \epsilon| < \epsilon \\ |y(x) - \epsilon| - \epsilon & \text{otherwise} \end{cases} \quad (1)$$

(Refer Fig 2)

for minimising a regularised error function.

$$E_r' = \frac{1}{2} \sum_{n=1}^N (y_n - \epsilon_n)^2 + \frac{\gamma}{2} \|w\|^2 \quad (2)$$

Step ② Reexpressing the optimization function as

$$\text{Min } E_r(w, \epsilon) = \frac{1}{2} \|w\|^2 + C \sum_{n=1}^N E_{\epsilon}(y(x_n) - \epsilon_n) \quad (3)$$

$$\text{Subjected to } \epsilon_n \leq y(x_n) + \epsilon + \hat{\epsilon}_n \quad (4)$$

$$\epsilon_n \geq y(x_n) - \epsilon - \hat{\epsilon}_n \quad (5)$$

where $\hat{\epsilon}_n$ are slack variables, ($\hat{\epsilon}_n$), $\epsilon_n > y(x_n) + \epsilon$ and $\epsilon_n < y(x_n) - \epsilon$ converted to ④ & ⑤ Refer Fig 2.

Step ③ Modification of error function in view of ③,

④ and ⑤ :

$$E_r(w, \epsilon_n, \hat{\epsilon}_n) = \frac{1}{2} \|w\|^2 + C \sum_{n=1}^N (\epsilon_n + \hat{\epsilon}_n) \quad (6)$$

⑥ to be minimised subject to $\epsilon_n, \hat{\epsilon}_n \geq 0$ & ④ and ⑤



Step ④ : Construction and Optimization



of Lagrangian function with Lagrangian multipliers $a_n, \hat{a}_n, \mu_n, \hat{\mu}_n \geq 0$:

$$L[w, b, \epsilon_{gn}, \hat{\epsilon}_{gn}]$$

$$= C \sum_{n=1}^N (\epsilon_{gn} + \hat{\epsilon}_{gn}) + \frac{1}{2} \|w\|^2 - \sum_{n=1}^N (\mu_n \epsilon_{gn} + \hat{\mu}_n \hat{\epsilon}_{gn}) \\ - \sum_{n=1}^N a_n (\epsilon + \epsilon_{gn} + y_n - t_n) - \sum_{n=1}^N \hat{a}_n (\epsilon + \hat{\epsilon}_{gn} + \hat{t}_n - y_n)$$

Step ⑤ Simplification of constraints and setting the derivatives to 0
Substituting for $y(x)$, $w^T \phi(x) + b$, *

$$\frac{\partial L}{\partial w} = 0 \Leftrightarrow w = \sum_{n=1}^N (a_n - \hat{a}_n) \phi(x_n) \quad (8)$$

$$\frac{\partial L}{\partial b} = 0 \Leftrightarrow \sum_{n=1}^N (a_n - \hat{a}_n) = 0 \quad (9)$$

$$\frac{\partial L}{\partial \epsilon_{gn}} = 0 \Leftrightarrow a_n + \mu_n = \epsilon \quad (10)$$

$$\frac{\partial L}{\partial \hat{\epsilon}_{gn}} = 0 \Leftrightarrow \hat{a}_n + \hat{\mu}_n = \epsilon \quad (11)$$

Step ⑥ Construction of Dual : (Elimination of corresponding variables from Lagrangian)

$$\text{Maximise } \bar{L}(a, \hat{a}) = -\frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N (a_n - \hat{a}_n)(a_m - \hat{a}_m) k(x_n, x_m) \\ - \epsilon \sum_{n=1}^N (a_n + \hat{a}_n) + \sum_{n=1}^N (a_n - \hat{a}_n) t_n \quad (12)$$

$$\text{Where } k(x, x') = \phi(x)^T \phi(x')$$

This is a constrained maximization problem. $a_n, \hat{a}_n, \mu_n, \hat{\mu}_n \geq 0$, (10) and (11) yields



Step ⑦ Construction of box constraints:

$$0 \leq a_n \leq c \quad (13)$$

$$0 \leq \hat{a}_n \leq c \quad (14)$$



Step ⑧ Simplification to programmable form of Dual:

Substitution ⑧ in to *

$$y(x) = \sum_{n=1}^N (a_n - \hat{a}_n) k(x, x_n) + b \quad (15)$$

Corresponding Kuhn-Tucker conditions yield

$$a_n [c + l_{g_n} + y_n - \epsilon_n] = 0 \quad (16)$$

$$\hat{a}_n [c + \hat{l}_{g_n} - y_n + \epsilon_n] = 0 \quad (17)$$

$$(c - a_n) l_{g_n} = 0 \quad (18)$$

$$(c - \hat{a}_n) \hat{l}_{g_n} = 0 \quad (19)$$

where $b = \epsilon_n - c - \sum_{m=1}^N (a_m - \hat{a}_m) k(x_n, x_m)$ (20)

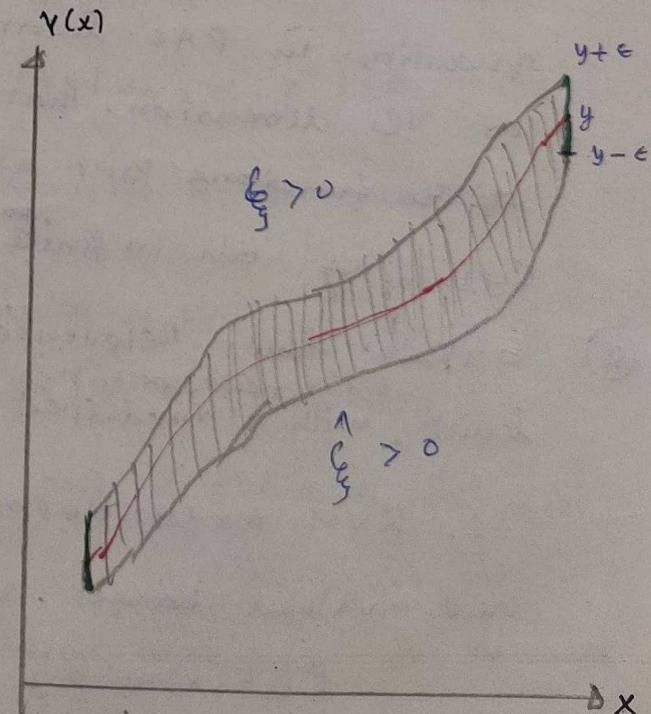
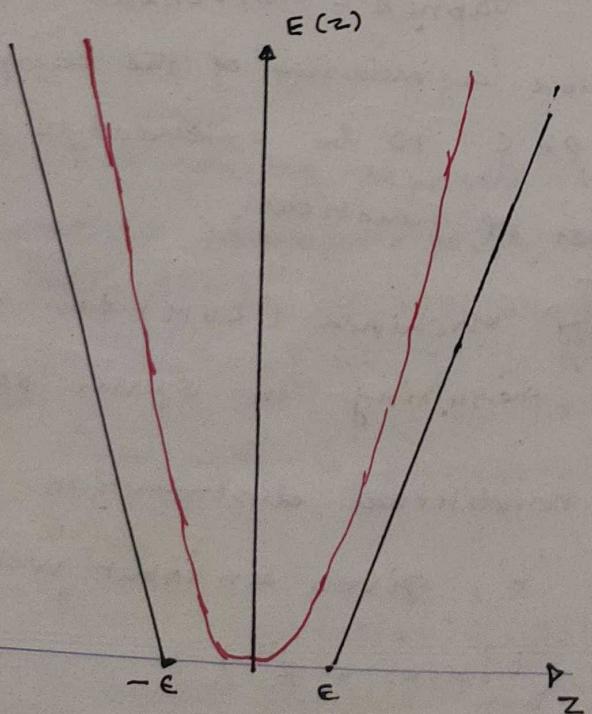


Fig 1: Plot of ϵ -insensitive error fn
in which error increases linearly 21
with distance beyond the insensitive region.



⑦ Write a short note on 'computational Learning Theory' (CLT)



CLT is a framework that provides full motivation to SVM originated from "Probably Approximately Correct" (PAC) whose goal is to understand how large a data set needs to be in order to give 'Good Generalization'. It provides bounds for computational cost of learning.

In PAC learning, we say a function $f(x, D)$ drawn from the space F of such functions on the basis of the training set D has good generalization if its expected error rate is below some threshold ϵ , so that

$$E_{x,r} [I(f(x; D) \neq t)] < \epsilon$$

where $I(\cdot)$ is an indicator function.

PAC learning aims to provide bounds for the minimum size of N of data set needed to meet this criterion. A key quantity in PAC learning is Vapnik-Chervonenkis dimension or VC dimension, that provides a measure of the complexity of the functions that allows PAC to be extended to spaces containing an infinite number of functions.

- ⑧ Explain how Relevance Vector Machines (RVM) can be built ^(for regression) with a modified prior resulting in sparse solutions?

1. RVM model defines a conditional distribution for a real-valued target variable ' t ', given an input vector x :

$$P[t|x, w, \beta] = N(t|y_m, \beta^{-1}) \quad \rightarrow ①$$

Where $\beta = 1/\sigma^2$ being the noise precision and mean given in the linear form:



$$y(x) = \sum_{i=1}^N w_i \phi_i(x) = w^T \phi(x) \xrightarrow{\textcircled{2}}$$



as bias

with non-linear basis functions $\phi_i(x)$ including a constant term of $\textcircled{2}$,

2. General expression takes the SVM-like form

$$y(x) = \sum_{n=1}^N w_n K(x, x_n) + b \xrightarrow{\textcircled{3}} \text{where } b - \text{bias}$$

Here number of parameters is $M = N + 1$

Unlike SVM (a) There is no restriction in positive definite kernel

(b) The basis functions are not tied either in number or location to the training data points

(Distinguish RVM for regression from SVM for regression \uparrow)

Step ② Learning choice of w and b
3. * Let there be N observations of input vector x

* The data matrix x have n th row as x_n^T , $n = 1$ to N .

* Target values be $t = (t_1, t_2, \dots, t_N)^T$

* Then Likelihood function is given by

$$p[t|x, w, \beta] = \prod_{n=1}^N p(t_n|x_n, w, \beta) \xrightarrow{\textcircled{4}}$$

* Introducing separate hyperparameter α_i for each of the weight parameters w_i (instead of single shared hyperparameter as in SVM), the weight prior is given by

$$p[w|\alpha] = \prod_{i=1}^M N(w_i|0, \alpha_i^{-1}) \xrightarrow{\textcircled{5}}$$

Where α_i is precision of w_i and $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_M)^T$

* Using $\xrightarrow{\textcircled{6}}$ for linear regression models, Posterior distribution for the weights is again Gaussian and is given by

$$p[w|t, x, \alpha, \beta] = N(w|m, \Sigma) \xrightarrow{\textcircled{7}}$$

④ Step ② Evaluating mean and covariance of posterior?



Where mean and covariance are given by

$$\mu = \beta \sum \varphi^T t \quad (8)$$

$$\Sigma = (A + \beta \varphi^T \varphi)^{-1} \quad (9)$$



Where φ is $N \times M$ design matrix with $\varphi_{ni} = \varphi_i(x_n)$ and $A = \text{diag}(\alpha_i)$ and K in (3) is the symmetric $(N+1) \times (N+1)$ kernel matrix with elements from $k(x_n, x_m)$

The values of α , and β in (7) are determined using type-2 maximum likelihood, (known as EVIDENCE Approximation)
[(Ans: What is Evidence Approximation in RVM)]

Ans: $\varphi^T t$ is type-2 maximum likelihood that maximises the Marginal likelihood function obtained by integrating weight parameters]

$$\text{Thus } p[t|x, \alpha, \beta] = \int p(t|x, w, \beta) p(w|\alpha) dw \quad (10)$$

As it is convolution of two Gaussians, it can be evaluated through log likelihood in the form

$$\begin{aligned} \log p[t|x, \alpha, \beta] &= \log N[t|0, C] \\ &= -\frac{1}{2} [N \log 2\pi + \log |C| + t^T C^{-1} t] \end{aligned}$$

where $t = (t_1, t_2, \dots, t_N)^T$ and $\longleftarrow (11)$

$$C = \beta^{-1} I + \varphi A^{-1} \varphi^T \quad \longleftarrow (12)$$

For maximising (11), setting derivatives of Marginal likelihood to zero yields

Step ③ of Re-estimating Hyper parameters:

$$\alpha_i^{\text{new}} = \frac{\gamma_i}{m_i^2} \quad \longleftarrow (13)$$

$$(\beta^{\text{new}})^{-1} = \frac{\|t - \varphi \mu\|^2}{N - \sum_i \gamma_i} \quad \longleftarrow (14)$$



Where m_i - i^{th} component of the posterior mean as in (8)

γ_i = a measure corresponding to w_i determined by data defined as $\gamma_i = 1 - \alpha_i \sum_{j \neq i} \gamma_j$ where Z_{ii} is the i^{th} diagonal component of posterior covariance given by (9)

(6) Step ④ Re estimating mean and covariance till suitable convergence criterion is satisfied.

(7) What do you mean 'Relevant Vectors' in RVM?

In 'Relevant Vector Machine', in the general expression

$$y(x) = \sum_{n=1}^N w_n K(x, x_n) + b, \text{ the inputs } x_n \text{ corresponding}$$

to remaining non zero weights are called 'Relevant Vectors'. They are identified through the mechanism of automatic relevant determination.

(8) What is the alternative procedure for improving training speed of RVM?

Ans : Sparsity Mechanism.

(9) Explain the Mathematical Analysis of Mechanism of Sparsity in the context of 'Relevant Vector Machine'

1. Origin of Sparsity in Bayesian Linear Models :

Consider a data set comprising $N = 2$ observations t_1, t_2 together with a model having single basis function $\phi(x)$, with hyper parameter α_i with isotropic noise having precision β .

From marginal likelihood of (11) and (12) of previous derivation of RVM for regression, we can write .



($P(t|a, \beta) = N[t|0, c]$) in which,

Covariance matrix takes the form



$$C = \frac{1}{\beta} I + \frac{1}{\alpha} \Phi \Phi^T \quad \xrightarrow{\text{①}}$$

where Φ denotes the N -dimensional vector $(\phi(x_1), \phi(x_2))^T$
and similarly $t = (t_1, t_2)^T$

This is a zero-mean Gaussian process over t with covariance C . Given a particular observation for t , we need to find a^* and β^* by maximising marginal likelihood on a particular d_i and then determine its stationary points explicitly.

We have

$$\log p[t | X, a, \beta] = -\frac{1}{2} [N \log 2\pi + \log |C| + t^T C^{-1} t] \quad \text{①}$$

$$C = \beta^{-1} I + \Phi A^{-1} \Phi^T \quad \xrightarrow{\text{②}} \quad \text{from ⑫}$$

from ⑪
is previous
bw

$$= \beta^{-1} I + \sum_{j \neq i} \alpha_j^{-1} \Phi_j \Phi_j^T + d_i^{-1} \Psi_i \Psi_i^T$$

$$= C_{-i} + d_i^{-1} \Psi_i \Psi_i^T \quad \xrightarrow{\text{③}}$$

Where Ψ_i denotes the i th column of Φ : $(\phi_i(x_1), \phi_i(x_2), \dots, \phi_i(x_n))$ and C_{-i} represents the matrix C with the contribution from basis function i removed. The determinant and inverse of C can be written as

$$|C| = |C_{-i}| [1 + d_i^{-1} \Psi_i^T C_{-i}^{-1} \Psi_i] \quad \xrightarrow{\text{④}}$$

$$C^{-1} = C_{-i}^{-1} - \frac{C_{-i}^{-1} \Psi_i \Psi_i^T C_{-i}^{-1}}{d_i + \Psi_i^T C_{-i}^{-1} \Psi_i} \quad \xrightarrow{\text{⑤}}$$

With the above equations ① can be written in the form

$$L(a) = \gamma (a_{-i}) + \lambda (a_i) \quad \xrightarrow{\text{⑥}}$$



Where $L(d_i)$ is simply the log marginal likelihood with basis function ψ_i omitted and.

$$\lambda(d_i) = \frac{1}{2} \left[\log d_i - \log(d_i + \beta_i) + \frac{\eta_i^2}{d_i + \beta_i} \right]$$



7

and contains all of the dependence on d_i and

$$\beta_i = \psi_i^T C_{-i}^{-1} \psi_i \quad (8)$$

$$\eta_i = \psi_i^T C_{-i}^{-1} t \quad (9)$$

Here β_i is called sparsity and η_i known as quality of ψ_i

* Large value of β_i relative to $\eta_i \Rightarrow$ Basis function ψ_i is more likely pruned from model.

* 'quality' represents a measure of alignment of ψ_i with error bet. training set and y_i of predictions with ψ_i excluded]

[Ans: Interpret the two terms 'sparsity' and 'quality' in sparsity mechanism.]

Now stationary points are given by

$$\frac{d\lambda(d_i)}{dd_i} = \frac{\eta_i^{-1} \beta_i^2 - (\eta_i^2 - \beta_i)}{2(d_i + \beta_i)^2} = 0 \Rightarrow (10)$$

$$\text{We get } d_i = \frac{\beta_i^2}{\eta_i^2 - \beta_i}$$

(12) Write down Sequential Sparse Bayesian Learning Algorithm:

Sequential Sparse Bayesian Learning Algorithm:

- For solving regression problem Initialise precision parameter β
- Initialise one basis function ψ_i with hyperparameter d_i , set $\eta_i = \beta_i^2 / (\eta_i^2 - \beta_i)$. With the remaining hyperparameters d_j for $j \neq i$ initialised to ∞ so that only ψ_i is included in the model.



3. Evaluate Σ and m along with q_i^2 & s_i for all basis functions.
4. Select a candidate basis function ψ_i
5. If $q_i^2 > s_i$ and α_i finite then update α_i using $\alpha_i = \frac{s_i}{q_i^2 - s_i}$
6. If $q_i^2 > s_i$ and $\alpha_i = \infty$, then add ψ_i to the model and evaluate hyperparameter α_i using $\alpha_i = s_i^2 / (q_i^2 - s_i)$
7. If $q_i^2 \leq s_i$ and α_i finite, then remove ψ_i & set $\alpha_i = \infty$
8. For solving a regression problem update β .
9. check for convergence. If not converged Go to 3

(13) Build a 'Relevance Vector Machine' for classification problem?

Consider a 2 class problem with binary target variable $t \in \{0, 1\}$.

With a linear combination of basis functions transformed by a logistic sigmoid function the model becomes

$$y(x, w) = \sigma [w^T q(x)] \quad \textcircled{1} \quad \text{where } \sigma(\cdot)$$

is the logistic sigmoid function given by $\sigma(z) = \frac{1}{1 + e^{-z}}$.

Initialise a hyperparameter α . Consider Laplace approximation for this fixed α . Mode of the posterior distribution is obtained by maximising :

$$\log p[w|t_{1:n}] = \log \{ p(t|w) p(w|x) \} - \log p(t|x)$$

$$= \sum_{n=1}^N \{ t_n \log y_n + (1-t_n) \log (1-y_n) \} - \frac{1}{2} w^T A w + C \rightarrow \textcircled{2}$$

where $A = \text{diag}(\alpha_i)$. Using iterative reweighted least squares, gradient vector and Hessian matrix of the log posterior distribution



$$(2) \Rightarrow \sigma [\log p(w|x)] = \Phi^T (\mathbf{t} - \mathbf{y}) - \mathbf{w}^T \quad \xrightarrow{\text{H}} (3)$$

$$\nabla \sigma [\log p(w|x)] = -(\Phi^T B \Delta + \mathbf{t}) \quad \xrightarrow{\text{H}} (4)$$

where B is an $N \times N$ diagonal matrix with elements $b_{ii} = y_i(1-y_i)$
the vector $\mathbf{y} = (y_1, \dots, y_N)^T$, and Φ is the design matrix with
elements $\phi_{ni} = \phi_i(x_n)$.