**Neural Networks for Signal Processing-I**

**Prof. Shayan Srinivasa Garani**

**Department of Electronic System Engineering**

**Indian Institute of Science, Bengaluru**

**Lecture – 40**

**ε-Insensitive Loss Function**

Support Vector Machines (SVMs) can indeed be applied to non-linear regression problems. To effectively set up and solve these regression problems using SVMs, it's essential to understand the concept of epsilon-insensitive loss functions. Let's explore how this can be done.

(Refer Slide Time: 01:52)



Our goal is to develop a robust estimator for regression that is insensitive to small fluctuations in the model, essentially a non-parametric approach where specific parameters are not predefined.

One major issue with traditional estimators, such as the least squares estimator, is its sensitivity to outliers. This sensitivity can lead to poor performance, especially with distributions that have long tails. The squared error metric tends to amplify errors, particularly when outliers are present. In contrast, using absolute error can mitigate this problem. Therefore, the objective here is to minimize the maximum degradation in performance, which can be framed as a min-max optimization problem.

(Refer Slide Time: 04:21)



To address this, we construct a loss function L(d, y), where d is the desired response and y is the estimator's output. We define this loss function as |d - y|. From this, we derive the epsilon-insensitive loss function. The epsilon-insensitive loss function operates as follows: if the absolute difference |d - y| is within an epsilon tolerance, the loss is zero; otherwise, the loss is |d - y| - $\epsilon$. This can be visualized as a "bathtub" function. Within the interval from -$\epsilon$ to +$\epsilon$, the loss is zero; outside this interval, the loss increases linearly with |d - y| - $\epsilon$. This defines an epsilon-insensitive region where the estimator is robust to small deviations, effectively disregarding errors within the range of -$\epsilon$ to +$\epsilon$.

Having established the epsilon-insensitive loss function, we can now apply this to SVMs for non-linear regression problems. Consider the non-linear regression model:

$$d = f(x) + v$$

where f is an unknown function, x is the input vector, v represents noise with unknown statistics, and d is the observed response. Given a training set consisting of pairs $(x_i, d_i)$ for i = 1 to n, our goal is to estimate the function f that describes the dependence of d on x.

(Refer Slide Time: 06:34)



The approach involves using the epsilon-insensitive loss function to handle the non-linear regression problem effectively with SVMs, thereby ensuring robustness against small errors and deviations in the model.

In tackling non-linear regression problems with SVMs, we use a linear combination of non-linear basis functions. This approach leverages the idea of incorporating non-linear functions over points in the input space to guide our formulation.

Let's consider y defined as:

$$y = \sum_{j=0}^{m} w_j \phi_j(x)$$

Here, y represents the output, where $w_j$ are the weights and $\varphi_j(x)$ are the non-linear basis functions applied to the input x. This can be seen as the inner product between the weight vector w and the vector of non-linear basis functions $\varphi(x)$. In this context, w is a weight vector including $w_0$ to $w_{m-1}$, with $w_0$ often representing the bias term. The function $\varphi(x)$ comprises $\varphi_0(x)$ to $\varphi_{m-1}(x)$, where $\varphi_0(x) = 1$ to account for the bias.

(Refer Slide Time: 08:42)



Our goal is to minimize the empirical risk given by:

$$\frac{1}{n} \sum_{i=1}^{n} \text{loss}(d_i, y_i)$$

where $\text{loss}(d_i, y_i)$ is the epsilon-insensitive loss function. We aim to minimize this risk subject to the constraint that the norm of the weight vector is less than or equal to a fixed constant $C_0$.

To solve this optimization problem, we introduce two sets of slack variables, $\zeta_i$ and $\zeta_i'$, for i = 1 to n. These slack variables help handle violations of the constraints. The constraints are as follows:

1. $\zeta_i \geq 0$ and $\zeta_i' \geq 0$ (non-negativity of slack variables).

2. Consistent with the epsilon-insensitive loss function, d - y should be greater than $\epsilon$ or d - y - $\epsilon$ should be positive. This translates into constraints that ensure the deviation from the epsilon-insensitive interval is properly accounted for.

(Refer Slide Time: 10:33)



We express these constraints with inequalities:

- For the case where d - y exceeds $\epsilon$, and
- For the case where y - d exceeds $\epsilon$.

Next, we formulate the cost functional. The first term of the cost functional is designed to minimize the number of misclassified points. This idea is similar to our approach for linearly separable patterns, where we formulated an average risk and set up an optimization

problem. The second term of the cost functional is the squared norm of the weight vector, which we seek to minimize to maximize the margin of separation.

(Refer Slide Time: 13:25)



Thus, the cost functional includes both minimizing the empirical risk and the norm of the weight vector. We aim to minimize this cost subject to constraints involving slack variables, Lagrange parameters, and the non-linear transformation of the input vectors. The Lagrangian for this problem incorporates the weight vector w, the slack variables ζ and ζ', and potentially all variables i from 1 to n can be combined into a vector for more efficient computation.

Similarly, we introduce the Lagrange multipliers α and α', as well as γ and γ', which correspond to the constraints on the slack variables. These slack variables must be non-negative, hence the need for Lagrange multipliers to enforce this constraint.

To proceed, we need to set up the Lagrangian for the optimization problem and then optimize with respect to the various variables within it. This involves taking partial

derivatives of the Lagrangian with respect to all variables, setting those derivatives to zero, and solving the resulting equations.

(Refer Slide Time: 14:18)



The cost functional J is quite complex, consisting of six terms. The first two terms pertain to the cost itself, while the next two terms are related to the constraints involving the Lagrange multipliers α and α', which link the weight vector w, the basis functions φ, and the responses d for the two regions in the epsilon-insensitive loss function. The final two terms address the constraints on the slack variables.

In essence, the last four terms of J are constraints related to the epsilon-insensitive loss function, while the first two terms correspond to the cost component.

To optimize this Lagrangian, we take the derivative with respect to the weight vector w and set it to zero. This yields the equation:

$$w = \sum_{i=1}^{n} (\alpha_i - \alpha_i')\phi(x_i)$$

(Refer Slide Time: 16:11)



This result shows that the multipliers $\gamma_i$ and $\gamma_i'$ are related to $\alpha_i$ and $\alpha_i'$ through a constant C, which imposes constraints on $\alpha_i$ and $\alpha_i'$. If the connection isn't immediately clear, it will be explored further in the course homework.

Next, we formulate the dual problem, which involves maximizing a dual cost function Q over the Lagrange multipliers $\alpha_i$ and $\alpha_i'$. This dual function includes the inner product kernel $k(x_i, x_j)$. The formulation of this dual problem is straightforward once we compute the partial derivatives of the Lagrangian from the primal problem with respect to the parameters.

I integrate this into the primal problem and then rearrange it to obtain the dual problem. In the dual problem, we aim to maximize a function, which is essentially the negative of the cost from the primal problem. This kernel function represents the inner product between two functions, specifically the inner product of $\varphi(x_i)$ and $\varphi(x_j)$.

Given this dual problem, it's crucial to understand its significance. We are provided with the training data set, which consists of pairs $(x_i, d_i)$ where i ranges from 1 to n. Our task is

to determine the Lagrange multipliers $\alpha_i$ and $\alpha_i'$ that maximize the dual cost function while satisfying the specified constraint equations.

(Refer Slide Time: 17:02)



Once we have computed the optimal weight vector w, we can then derive the regression function f given by $w^T \varphi(x)$. This function can be compactly expressed as:

$$f(x) = \sum_{i=1}^{n}(\alpha_i - \alpha_i')k(x_i, x)$$

Here, the only free parameters are $\epsilon$ and C. $\epsilon$ is related to the epsilon-insensitive loss function, and C is a scaling constant used for summing the slack variables. These parameters are not optimized within the framework but are set as part of the problem setup.

After formulating and solving the optimization problem using an optimization solver, we obtain the optimal hyperplane. From this hyperplane, we derive the optimal weight vector, which allows us to compute the regression function.

This problem, particularly in the context of non-linear regression, involves multiple constraints, two related to the epsilon-insensitive loss function and two concerning the slack variables. Additionally, the cost function includes terms for minimizing the weight vector norm and the misclassification error. Consequently, solving this problem involves handling a somewhat complex optimization task.

With this module, you should now have a clearer understanding of how to set up and solve these optimization problems in general. Next, we will revisit the XOR problem and explore how SVMs can be applied to solve it through an illustrative example.