

UNIVERSITY OF  
**WATERLOO**



UNIVERSITY OF WATERLOO  
CHERITON SCHOOL OF COMPUTER SCIENCE

## CS 480 - Homework 1

**Prepared by:**

Jianxiang (Jack) Xu [20658861]

12 February 2021

**CS480/680: Introduction to Machine Learning**

## Homework 2

Due: 11:59 pm, February 12, 2021

**Exercise 1:  $k$ -Nearest Neighbours (5 pts)**

The KNN-dataset provides a 4x4 representation of the handwritten digits 5, 6, and 8. Using the notebook provided, train a KNN classifier using 5-fold cross validation to determine the optimal  $k$ -value for this task.

1. (3 pt) Create a graph that shows the average accuracy based on 5-fold cross validation when varying the number of neighbours from 1 to 30.

Ans: [Answer 1.1](#)

2. (1 pt) Report the best hyperparameter found by 5-fold cross-validation and the cross-validation accuracy.

Ans: [Answer 1.2](#)

3. (1 pt) Report the test accuracy based on the best hyperparameter.

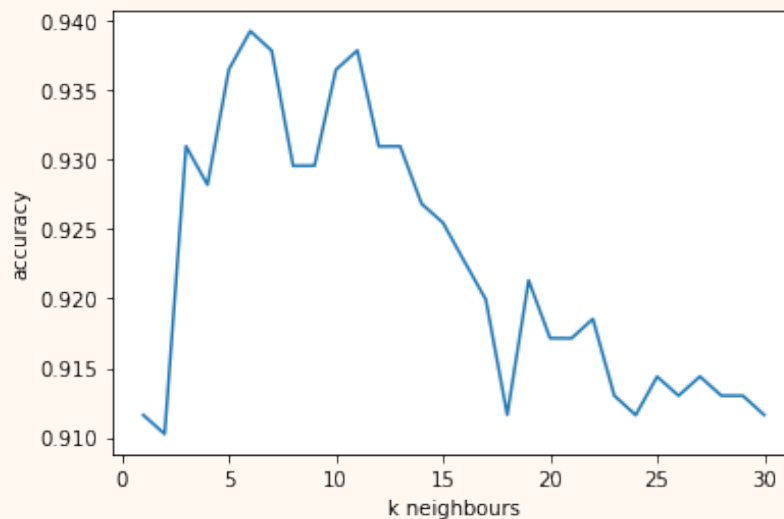
Ans: [Answer 1.3](#)**Answer 1.1: KNN classifier with 5-fold cross validation**

Figure 1-1. KNN Accuracy vs. K-size

See code in the jupyter submission.

**Answer 1.2: Best Hyperparameter**

From the output below, we may find the best hyperparameter is  $k = 5$ , with 93.92% cross-validation accuracy and 92.31% test accuracy.

```
1 best # of neighbours k: 5
2 best cross validation accuracy: 0.9392432950191572
3 test accuracy: 0.9230769230769231
4
```

**Answer 1.3: Test Accuracy**

As reported in [Answer 1.2](#), the final test accuracy is 92.31%.

**Exercise 2: Poisson Regression (4 pts)**

Recall that in logistic regression we assumed the *binary* label  $Y_i \in \{0, 1\}$  follows the Bernoulli distribution:  $\Pr(Y_i = 1|X_i) = p_i$ , where  $p_i$  also happens to be the mean. Under the independence assumption we derived the log-likelihood function:

$$\sum_{i=1}^n (1 - y_i) \log(1 - p_i) + y_i \log(p_i). \quad (1)$$

Then, we parameterized the mean parameter  $p_i$  through the logit transform:

$$\log \frac{p_i}{1 - p_i} = \mathbf{w}^\top \mathbf{x}_i + b, \quad \text{or equivalently} \quad p_i = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x}_i - b)}. \quad (2)$$

Lastly, we found the weight vector  $\mathbf{w}$  and  $b$  by maximizing the log-likelihood function.

In the following we generalize the above idea to the case where  $Y_i \in \mathbb{N}$ , i.e.,  $Y_i$  can take any natural number (for instance, when we are interested in predicting the number of customers or network packages).

1. (1 pt) Naturally, we assume  $Y_i \in \mathbb{N}$  follows the Poisson distribution (with mean  $\mu_i > 0$ ):

$$\Pr(Y_i = k|X_i) = \frac{\mu_i^k}{k!} \exp(-\mu_i), \quad k = 0, 1, 2, \dots \quad (3)$$

Given a dataset  $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^n$ , what is the log-likelihood function (of  $\mu_i$ 's) given  $\mathcal{D}$ ?

Ans: [Answer 2.1](#)

2. (1 pt) Can you give some justification of the parameterization below?

$$\log \mu_i = \mathbf{w}^\top \mathbf{x}_i + b. \quad (4)$$

Ans: [Answer 2.2](#)

3. (1 pt) Based on the above, write down the objective function for Poisson regression. Please specify the optimization variables and whether you are maximizing or minimizing. [Constants can be dropped.]

Ans: [Answer 2.3](#)

4. (1 pt) Compute the gradient of your objective function above and formulate a gradient algorithm for finding the weight vector  $\mathbf{w}$  and  $b$ .

Ans: [Answer 2.4](#)

**Answer 2.1: Log-likelihood Function**

Since  $Y_i = y_i \in \mathbb{N}$  are i.i.d. Poisson random variables, we may obtain the likelihood function ( $L(\mu_i)$ ):

$$L(\mu_i) = \Pr(Y_1 = y_1, \dots, Y_n = y_n | X_1 = x_1, \dots, X_n = x_n) \quad (5)$$

$$= \prod_{i=1}^n \Pr(Y_i = y_i | X_i) \quad (6)$$

$$= \prod_{i=1}^n \frac{\mu_i^{y_i}}{y_i!} \exp(-\mu_i) \quad (7)$$

Since maximizing the likelihood function is equivalent to maximizing its log function, hence we may take log of the entire likelihood function  $L(\mu_i)$  from above to compute the log-likelihood:

$$l(\mu_i) = \log(L(\mu_i)) = \log \left( \prod_{i=1}^n \frac{\mu_i^{y_i}}{y_i!} \exp(-\mu_i) \right) \quad (8)$$

$$= \sum_{i=1}^n (\log(\exp(-\mu_i)) + \log(\mu_i^{y_i}) - \log(y_i!)) \quad (9)$$

$$= \sum_{i=1}^n (-\mu_i + y_i \log(\mu_i) - \log(y_i!)) \quad (10)$$

**Answer 2.2: Justification of the parameterization**

The goal in linear regression is to find some best fit factor  $\mathbf{w}$  and bias  $b$  that will allow us to map  $\mathbf{x}_i$  to some model  $y$ . Similarly, we transfer the space into a distribution model, where the output  $y_i$  is sampled from a normal distribution (Poisson distribution in this case).

Hence,

$$y \sim \mathcal{N}(y; \mu, \sigma^2), \quad \text{where } \mu_i = \mathbf{w}^\top \mathbf{x}_i + b. \quad (11)$$

However, in our case, the distribution is Poisson and  $\mu_i > 0$ . Since  $\mathbf{w}^\top \mathbf{x}_i + b$  has to be mapped to  $(-\infty, \infty)$ ,  $\log(\mu_i)$  is a reasonable mapping strategy to transfer  $\mu_i > 0$  to  $(-\infty, \infty)$ .

Hence,

$$y \sim Po(y; \mu, \sigma^2), \quad \text{where } \log(\mu_i) = \mathbf{w}^\top \mathbf{x}_i + b. \quad (12)$$

**Answer 2.3: Objective Function for Poisson Regression**

From Equation (4) (justified in Answer 2.2), we may now complete the log-likelihood function as:

$$l(\mu_i) = \sum_{i=1}^n (-\mu_i + y_i \log(\mu_i) - \log(y_i!)) \quad (13)$$

$$= \sum_{i=1}^n \left( -e^{(\mathbf{w}^\top \mathbf{x}_i + b)} + y_i(\mathbf{w}^\top \mathbf{x}_i + b) - \log(y_i!) \right) \quad (14)$$

Since  $\log(y_i!)$  term in the equation could be regarded as constant, since it is provided in training data, we may drop it.

$$\therefore \frac{\partial \sum_{i=1}^n -\log(y_i!)}{\partial \mu_i} = 0 \quad (15)$$

$$\therefore l(w, b) = \sum_{i=1}^n \left( -e^{(\mathbf{w}^\top \mathbf{x}_i + b)} + y_i(\mathbf{w}^\top \mathbf{x}_i + b) \right) \quad (16)$$

There is no closed form solution to maximizing the log-likelihood function, but we can convert it to minimization and apply gradient descent.

$$\max_{w, b} l(w, b) = \min_{w, b} -l(w, b) \quad (17)$$

$$= \min_{w, b} \sum_{i=1}^n \left( e^{(\mathbf{w}^\top \mathbf{x}_i + b)} - y_i(\mathbf{w}^\top \mathbf{x}_i + b) \right) \quad (18)$$

**Objective Function**

$$\min_{w, b} \sum_{i=1}^n \left( e^{(\mathbf{w}^\top \mathbf{x}_i + b)} - y_i(\mathbf{w}^\top \mathbf{x}_i + b) \right) \equiv \min_{\mu_i} \sum_{i=1}^n (\mu_i - y_i \log(\mu_i)) \quad (19)$$

**Answer 2.4: Gradient Algorithm**

Let's find the gradient step of the objective function (Equation (19)):

$$\frac{\partial}{\partial \mathbf{w}} \sum_{i=1}^n \left( e^{(\mathbf{w}^\top \mathbf{x}_i + b)} - y_i(\mathbf{w}^\top \mathbf{x}_i + b) \right) = \sum_{i=1}^n x_i e^{w^T x_i + b} - \sum_{i=1}^n y_i x_i \quad (20)$$

$$\frac{\partial}{\partial b} \sum_{i=1}^n \left( e^{(\mathbf{w}^\top \mathbf{x}_i + b)} - y_i(\mathbf{w}^\top \mathbf{x}_i + b) \right) = \sum_{i=1}^n e^{w^T x_i + b} - \sum_{i=1}^n y_i \quad (21)$$

$$(22)$$

Hence, the gradient algorithm for weight vector  $\mathbf{w}$  and  $b$ :

**Gradient Descent Update Function**

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \left( \sum_{i=1}^n x_i e^{w^T x_i + b} - \sum_{i=1}^n y_i x_i \right) \quad (23)$$

$$b \leftarrow b - \eta \left( \sum_{i=1}^n e^{w^T x_i + b} - \sum_{i=1}^n y_i \right) \quad (24)$$

**Exercise 3: Fun with Classification (5 pts)**

For this problem, you are allowed to use `statsmodels` and `sklearn` as directed.

1. (2 pts) Run logistic regression, SVM with  $\ell_2$  regularization with parameter 1 (soft-margin SVM), and SVM with regularization parameter `float('inf')` (hard-margin SVM) on Mystery Dataset A (note that there is only a training dataset and no test dataset). Use `Logit` from `statsmodels` and `SVC` (with linear kernel) from `sklearn`. One of these methods will not work -- mathematically explain why (that is, give an explanation which is more informative than any error message you encounter). How could the associated problems be remedied? Discuss similarities and differences between the solution obtained via these three methods.

Ans: [Answer 3.1](#)

2. (3 pts) Take your solution for the soft-margin SVM from the previous part. For each point in the dataset, take its inner product with the produced coefficient vector, and scale the result by the sign of each point's label (replace 0's with -1's). How many of these values are  $\leq 1$ ? [Be sure you're getting all of them -- there may be numerical precision issues, so if in doubt, err on the side of counting a point.] Based on your answer to these questions, sketch a 2D caricature of what the points and the hyperplane defined by the SVM solution look like. Write the parameter vector solution to the SVM problem as a linear combination of some points in your dataset. How many points did you require? [You may find built-in functions useful for this purpose.]

Compute the solution for the same three methods on Mystery Dataset B. You will again run into issues with one of them -- explain why. [The answer is likely to be simpler than last time.] Find a way to write the parameter vector solution to the SVM problem as a linear combination of some points in your dataset -- do not report the solution itself, but report how many points you used, and how you arrived at this answer. Compare the empirical prediction accuracy (i.e., using 0-1 loss) of the successfully-trained classifiers on the test set.

Ans: [Answer 3.2](#)



**Answer 3.1: Logit and SVM**

It appears 'Logit' from 'statsmodels' terminated with error message of "PerfectSeparationError: Perfect separation detected, results not available", whereas soft-margin and hard-margin SVM are both capable to perform the fitting.

This error happens when a linear combination of the predictors yield a perfect prediction of the response variable or the binary outcome in this case. Simply,  $x_i$  predicts  $y_i$  perfectly. If we perform l2 norm on every single x dataset, we may realize every two alternating points with opposing labels are equal distance from the origin, and we may suspect, there exists such plane through the origin that can perfectly separate the dataset. Hence, the maximum likelihood on  $x_i$  would approach to infinity, therefore, we cannot find a suitable coefficient to maximize the conditional likelihood.

A solution to resolve this is to utilize a form of penalized regression to penalize the likelihood (the whole point of SVM).

The final trained results are shown in Answer 3.1, and we observed that both hard and soft margin SVM provides the exact same hyperplane through the origin as expected. Whereas, the logistic regression failed to compute.

```

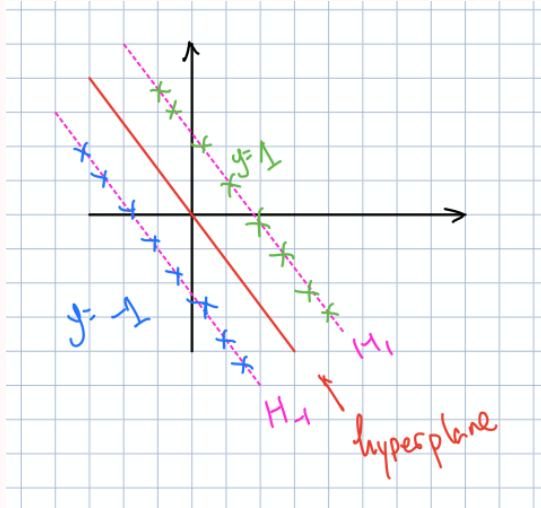
1 w_hard = [[-1.21665426e-01  3.62093715e-04 -8.14547914e-02 -1.30461544e-01
2           1.03679744e-01 -1.52182440e-01  1.04733167e-01  2.37085565e-02
3           1.88904300e-01  4.14259231e-01 -5.74918926e-02  2.53091455e-02
4           -9.57967331e-02  2.78213946e-01  1.26402871e-02  4.24952941e-02
5           -1.08058239e-01 -3.17898534e-02  3.32288614e-03  1.71401564e-01
6           -6.93185355e-02 -2.85727264e-01 -1.11431511e-01 -2.81551300e-02
7           7.14262962e-02  1.98803054e-01 -1.89996218e-01  9.63349060e-02
8           1.20110591e-01 -5.53355480e-02 -4.29452659e-03 -1.10216493e-01
9           -1.57449276e-01  3.64106600e-02 -4.35440021e-02 -1.42084333e-01
10          1.30750877e-01  2.68330435e-02  6.77532548e-02 -3.01716444e-01
11          -1.22263636e-02  2.15382941e-01  1.26537049e-01  2.24429271e-02
12          1.34745230e-01  3.23138170e-02  1.47580798e-01  5.67770933e-02
13          2.91490606e-01 -3.25700173e-02]]
14 b_hard = [-0.]
15 w_soft = [[-1.21665426e-01  3.62093715e-04 -8.14547914e-02 -1.30461544e-01
16            1.03679744e-01 -1.52182440e-01  1.04733167e-01  2.37085565e-02
17            1.88904300e-01  4.14259231e-01 -5.74918926e-02  2.53091455e-02
18            -9.57967331e-02  2.78213946e-01  1.26402871e-02  4.24952941e-02
19            -1.08058239e-01 -3.17898534e-02  3.32288614e-03  1.71401564e-01
20            -6.93185355e-02 -2.85727264e-01 -1.11431511e-01 -2.81551300e-02
21            7.14262962e-02  1.98803054e-01 -1.89996218e-01  9.63349060e-02
22            1.20110591e-01 -5.53355480e-02 -4.29452659e-03 -1.10216493e-01
23            -1.57449276e-01  3.64106600e-02 -4.35440021e-02 -1.42084333e-01
24            1.30750877e-01  2.68330435e-02  6.77532548e-02 -3.01716444e-01
25            -1.22263636e-02  2.15382941e-01  1.26537049e-01  2.24429271e-02
26            1.34745230e-01  3.23138170e-02  1.47580798e-01  5.67770933e-02
27            2.91490606e-01 -3.25700173e-02]]
28 b_soft = [-0.]
29 |w_hard|_2 - |w_soft|_2 = 0.0

```

**Answer 3.2: Conclusion**

**How many of these values are  $\leq 1$ ?**

All the values (2000/2000) are  $\leq 1$

**2D caricature sketch of SVM solution****Parameter vector solution**

2 points needed.

$$y = -0.5x_{1999} + 0.5x_{1998}$$

**Mystery Dataset B**

We ran into problem (infinity loop) with Hard-margin SVM, while Soft-margin SVM and Logistic regression has been successfully fitted. This is caused by the dataset that is inseparable with hard-margin, hence the program ran into infinity loop trying to separate the dataset that is not separable.

192 points needed.

To compare the SVM and logistic regression performance, the empirical metric accuracy is calculated based on test data:  $\text{accuracy} = \frac{\# \text{ right predictions}}{\# \text{ of test data}}$

We find:

Soft-SVM Test Performance (Empirical Accuracy): 97.15%

Logit Test Performance (Empirical Accuracy): 96.95%

Hence, the Soft-SVM has better performance.

**Exercise 4: Support Vector Regression (8 pts)**

Let us consider support vector regression:

$$\min_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n \max\{|y_i - (\mathbf{w}^\top \mathbf{x}_i + b)| - \varepsilon, 0\}, \quad (25)$$

where  $\mathbf{x}_i \in \mathbb{R}^d$ ,  $y_i \in \mathbb{R}$ , and  $\|\mathbf{w}\|_2 := \sqrt{\sum_{j=1}^d w_j^2}$  is the Euclidean norm.

- (2 pts) Derive the Lagrangian dual of the support vector regression loss function (25). Please include intermediate steps so that you can get partial credits.

Ans: [Answer 4.1](#)

In the following you will complete and implement the following gradient algorithm for solving support vector regression in Equation (25):

**Algorithm 1.1: GD for SVR.**

**Input:**  $X \in \mathbb{R}^{n \times d}$ ,  $\mathbf{y} \in \mathbb{R}^n$ ,  $\mathbf{w} = \mathbf{0}_d$ ,  $b = 0$ ,  $\text{max\_pass} \in \mathbb{N}$ , step size  $\eta$

**Output:**  $\mathbf{w}, b$

```

1 for  $t = 1, 2, \dots, \text{max\_pass}$  do
2   for  $i = 1, 2, \dots, n$  do
3     choose step size  $\eta$ 
4     if  $|y_i - (\langle \mathbf{x}_i, \mathbf{w} \rangle + b)| \geq \varepsilon$  then
5        $\mathbf{w} \leftarrow$                                      //  $\mathbf{x}_i$  is the  $i$ -th row of  $X$ 
6        $b \leftarrow$ 
7        $\mathbf{w} \leftarrow$                                      // proximal step
```

Note that this differs a bit from what you've seen so far, in terms of gradient descent. Rather than taking steps based on the entire loss function, we instead take a step based on the unregularized loss, and then perform a projection step based on the regularizer.

- (2 pts) Compute the gradient w.r.t.  $\mathbf{w}$  and  $b$  for each second term in Equation (25). Note that in places where the function is non-differentiable, you might have to compute a sub-gradient.

$$C \sum_{i=1}^n \max\{|y_i - (\mathbf{w}^\top \mathbf{x}_i + b)| - \varepsilon, 0\} \quad (26)$$

Ans: [Answer 4.2](#)

- (1 pt) Find the closed-form solution of the following proximal step:

$$P^\eta(\mathbf{w}) = \underset{\mathbf{z}}{\operatorname{argmin}} \frac{1}{2\eta} \|\mathbf{z} - \mathbf{w}\|_2^2 + \frac{1}{2} \|\mathbf{z}\|_2^2 \quad (27)$$

Ans: [Answer 4.3](#)

- (3 pts) Implement Algorithm 1.1. You should use Ex 1.3 to complete lines 5-6, and Ex 1.4 for line 7. Run it on Mystery Dataset C (this is Mystery Dataset A from Assignment 1, but reused), and report your training error, training loss, and test error. Use  $C = 1$  and  $\varepsilon = 0.1$ .

Ans: [Answer 4.4](#)

**Answer 4.1: Lagrangian Dual Derivation**

Let's substitute the regulation (second) term with a complement slack variables  $(\xi_i, \xi'_i)$  as suggested in lecture to account the delta error outside the  $\varepsilon$  margin.

**Condition:** (28)

$$y_i - \hat{y}_i \leq \xi_i + \varepsilon \quad (29)$$

$$\hat{y}_i - y_i \leq \xi'_i + \varepsilon \quad (30)$$

$$\text{where } \xi_i, \xi'_i > 0, \quad \hat{y}_i = \mathbf{w}^T \mathbf{x}_i + b \quad (31)$$



Hence, we may simplify the formulation Equation (25) into:

$$\min_{\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n (\xi_i + \xi'_i) \quad (32)$$

**where** (33)

$$y_i - \hat{y}_i \leq \xi_i + \varepsilon \quad (34)$$

$$\hat{y}_i - y_i \leq \xi'_i + \varepsilon \quad (35)$$

$$\text{where } \xi_i, \xi'_i > 0, \quad \hat{y}_i = \mathbf{w}^T \mathbf{x}_i + b \quad (36)$$

We may now apply Lagrangian with complement  $(\alpha, \alpha')$  and  $(\beta, \beta')$  multipliers:

$$\begin{aligned} \min_{\mathbf{w}, b, \xi, \xi'} \max_{\alpha, \alpha' \geq 0, \beta, \beta' \leq 0} & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n (\xi_i + \xi'_i) + \sum_{i=1}^n \alpha_i (y_i - \hat{y}_i - \varepsilon - \xi_i) \\ & + \sum_{i=1}^n \alpha'_i (-y_i + \hat{y}_i - \varepsilon - \xi'_i) - \sum_{i=1}^n \beta_i \xi_i - \sum_{i=1}^n \beta'_i \xi'_i \end{aligned} \quad (37)$$

Swap min and max:

$$\begin{aligned} \max_{\alpha, \alpha' \geq 0} \min_{\beta, \beta' \leq 0} & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^n (\xi_i + \xi'_i) + \sum_{i=1}^n \alpha_i (y_i - (\mathbf{w}^T \mathbf{x}_i + b) - \varepsilon - \xi_i) \\ & + \sum_{i=1}^n \alpha'_i (-y_i + (\mathbf{w}^T \mathbf{x}_i + b) - \varepsilon - \xi'_i) - \sum_{i=1}^n \beta_i \xi_i - \sum_{i=1}^n \beta'_i \xi'_i \end{aligned} \quad (38)$$

Take derivative w.r.t.  $w, b, \xi_i, \xi'_i$ :

$$\frac{\partial}{\partial \mathbf{w}} = \mathbf{w} - \sum_{i=1}^n \alpha_i \mathbf{x}_i + \sum_{i=1}^n \alpha'_i \mathbf{x}_i \quad (39)$$

$$\frac{\partial}{\partial b} = -\sum_{i=1}^n \alpha_i + \sum_{i=1}^n \alpha'_i \quad (40)$$

$$\frac{\partial}{\partial \xi_i} = nC - n\alpha_i - n\beta_i \quad (41)$$

$$\frac{\partial}{\partial \xi'_i} = nC - n\alpha'_i - n\beta'_i \quad (42)$$

, and  $\equiv 0$ :

$$\sum_{i=1}^n (\alpha_i - \alpha'_i) \mathbf{x}_i = \mathbf{w} \quad (43)$$

$$\sum_{i=1}^n (\alpha'_i - \alpha_i) = 0 \quad (44)$$

$$\alpha_i + \beta_i = \alpha'_i + \beta'_i = C \quad (45)$$

Substitute back to Lagrangian (Equation (38)):

$$\begin{aligned} \max_{\alpha, \alpha' \geq 0} \min_{\beta, \beta' \leq 0} & \frac{1}{2} \left\| \sum_{i=1}^n (\alpha_i - \alpha'_i) \mathbf{x}_i \right\|_2^2 + C \sum_{i=1}^n (\xi_i + \xi'_i) + \sum_{i=1}^n \alpha_i (y_i - (\mathbf{w}^T \mathbf{x}_i + b) - \varepsilon - \xi_i) \\ & + \sum_{i=1}^n \alpha'_i (-y_i + (\mathbf{w}^T \mathbf{x}_i + b) - \varepsilon - \xi'_i) - \sum_{i=1}^n \beta_i \xi_i - \sum_{i=1}^n \beta'_i \xi'_i \end{aligned} \quad (46)$$

$$\begin{aligned} \max_{\alpha, \alpha' \geq 0} \min_{\beta, \beta' \leq 0} & \frac{1}{2} \left\| \sum_{i=1}^n (\alpha_i - \alpha'_i) \mathbf{x}_i \right\|_2^2 + C \sum_{i=1}^n (\xi_i + \xi'_i) + \sum_{i=1}^n \alpha_i y_i - \sum_{i=1}^n \alpha_i (\mathbf{w}^T \mathbf{x}_i + b) - \sum_{i=1}^n \alpha_i \varepsilon - \sum_{i=1}^n \alpha_i \xi_i \\ & - \sum_{i=1}^n \alpha'_i y_i + \sum_{i=1}^n \alpha'_i (\mathbf{w}^T \mathbf{x}_i + b) - \sum_{i=1}^n \alpha'_i \varepsilon - \sum_{i=1}^n \alpha'_i \xi'_i - \sum_{i=1}^n \beta_i \xi_i - \sum_{i=1}^n \beta'_i \xi'_i \end{aligned} \quad (47)$$

$$\begin{aligned} \max_{\alpha, \alpha' \geq 0} \min_{\beta, \beta' \leq 0} & \frac{1}{2} \left\| \sum_{i=1}^n (\alpha_i - \alpha'_i) \mathbf{x}_i \right\|_2^2 + C \sum_{i=1}^n (\xi_i + \xi'_i) - \sum_{i=1}^n (\alpha_i + \alpha'_i) \varepsilon + \sum_{i=1}^n (\alpha'_i - \alpha_i) y_i \\ & - \sum_{i=1}^n (\alpha_i + \beta_i) \xi_i - \sum_{i=1}^n (\alpha'_i + \beta'_i) \xi'_i \end{aligned} \quad (48)$$

$$\max_{\alpha, \alpha' \geq 0} \frac{1}{2} \left\| \sum_{i=1}^n (\alpha_i - \alpha'_i) \mathbf{x}_i \right\|_2^2 - \sum_{i=1}^n (\alpha_i + \alpha'_i) \varepsilon + \sum_{i=1}^n (\alpha'_i - \alpha_i) y_i \quad (49)$$

Convert to minimization problem:

**The Dual**

$$\min_{C \geq \alpha, \alpha' \geq 0} \quad \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n (\alpha_i - \alpha'_i)(\alpha_j - \alpha'_j) \mathbf{x}_i^T \mathbf{x}_j + \varepsilon \sum_{i=1}^n (\alpha_i + \alpha'_i) - \sum_{i=1}^n (\alpha'_i - \alpha_i) y_i \quad (50)$$

$$\text{s.t.} \quad \sum_{i=1}^n (\alpha'_i - \alpha_i) = 0 \quad (51)$$

**Answer 4.2: Gradient**

$$\frac{\partial}{\partial \mathbf{w}} = \sum_{i=1}^n \begin{cases} 0 & |y_i - \hat{y}_i| \leq \epsilon \\ -Cx_i & y_i > \hat{y}_i, |y_i - \hat{y}_i| > \epsilon \\ Cx_i & \text{otherwise} \end{cases} \quad (52)$$

$$\frac{\partial}{\partial b} = \sum_{i=1}^n \begin{cases} 0 & |y_i - \hat{y}_i| \leq \epsilon \\ -C & y_i > \hat{y}_i, |y_i - \hat{y}_i| > \epsilon \\ C & \text{otherwise} \end{cases} \quad (53)$$

$$\text{where} \quad \hat{y}_i = \mathbf{w}^T \mathbf{x}_i + b \quad (54)$$

**Answer 4.3: Closed Form Solution of the Proximal Step**

$$\frac{\partial}{\partial \mathbf{z}} P^\eta(\mathbf{w}) = \frac{\partial}{\partial \mathbf{z}} \left( \frac{1}{2\eta} \|\mathbf{z} - \mathbf{w}\|_2^2 + \frac{1}{2} \|\mathbf{z}\|_2^2 \right) \quad (55)$$

$$0 = \frac{1}{\eta} (\mathbf{z} - \mathbf{w}) + \mathbf{z} \quad (56)$$

$$\mathbf{z} = \frac{\mathbf{w}}{1 + \eta} \quad (57)$$

**Closed Form Solution**

$$\mathbf{z} = \frac{\mathbf{w}}{1 + \eta} \quad (58)$$



**Answer 4.4: Implementation****Algorithm 2.2: Seudo implementation of GD for SVR.**

**Input:**  $X \in \mathbb{R}^{n \times d}$ ,  $\mathbf{y} \in \mathbb{R}^n$ ,  $\mathbf{w} = \mathbf{0}_d$ ,  $b = 0$ ,  $\text{max\_pass} \in \mathbb{N}$ , step size  $\eta$

**Output:**  $\mathbf{w}, b$

```

8 for  $t = 1, 2, \dots, \text{max\_pass}$  do
9   for  $i = 1, 2, \dots, n$  do
10    choose step size  $\eta$ 
11    if  $|y_i - (\langle \mathbf{x}_i, \mathbf{w} \rangle + b)| \geq \varepsilon$  then
12       $\mathbf{w} \leftarrow \mathbf{w} - \eta C \mathbf{x}_i \times \text{sgn}(\langle \mathbf{w}, \mathbf{x} \rangle + b - y_i)$            //  $\mathbf{x}_i$  is the  $i$ -th row of  $X$ 
13       $b \leftarrow b - \eta C \times \text{sgn}(\langle \mathbf{w}, \mathbf{x} \rangle + b - y_i)$ 
14     $\mathbf{w} \leftarrow \frac{\mathbf{w}}{1+\eta}$                                            // proximal step

```

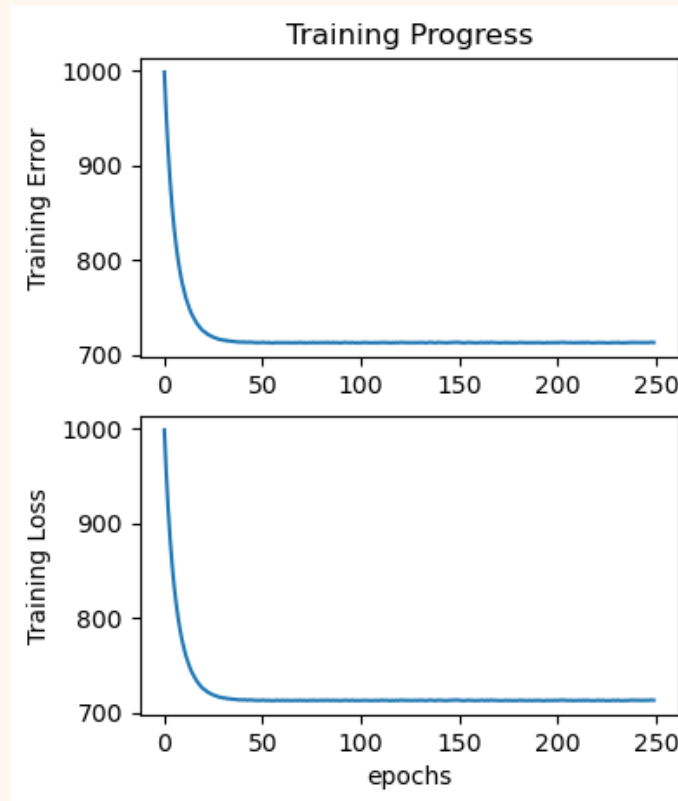


Figure 4-2. Training Progress

```

1 > [result ] T: 249 | Training Error: 712.84065 | Training Loss: 713.37650 | Test Error: 891.05659
2

```