

# Project 1: Kernelized Linear Classification

Stefano Chiesa

## Abstract

This report presents an analysis of kernelized linear classification. I implemented several machine learning algorithms from scratch, including the Perceptron, Support Vector Machines (SVMs) using the Pegasos algorithm, and Regularized Logistic Regression. I analysed the impact of polynomial feature expansion and kernel methods using Gaussian and polynomial kernels. The performance of each model was evaluated based on  $\ell_{0-1}$  loss, variance and runtime.

## 1 Introduction

The goal of this project was to classify labels based on numerical features using kernelized linear classification methods, implementing and testing various machine learning algorithms from scratch to assess their performance and to understand the trade-offs between different approaches.

You can find the code and the Project 1 instructions in the [KernelLinearClassificationML](#) GitHub repository.

## 2 Dataset Exploration and Preprocessing

### 2.1 Dataset Description

The dataset consists of 10,000 instances, each characterized by 10 numerical features ( $x_1$  to  $x_{10}$ ) and a binary target label  $y$ . The target variable  $y$  takes values of 1 or -1, meaning the problem is a binary classification task.

### 2.2 Data Preprocessing

To manage the computational limitations of my machine, I performed a random sampling of the dataset selecting a subset of 1,000 instances using the `sample` method, as shown below:

```
data = data.sample(n=1000, random_state=45)
```

where `data` is a Pandas data frame containing the imported data.

This step was necessary because the original dataset was too large for my PC to process efficiently, and reducing the dataset size allowed the algorithms to run faster while still providing meaningful insights for model training and evaluation. The `random_state` parameter was set to 45 to ensure reproducibility of the results.

Before training and testing each algorithm, I used the function `split_data()` to split the dataset into training and testing sets. The function first introduces a bias term by adding a column of ones to the feature matrix: including this bias term is a common practice in machine learning, as it provides greater flexibility and improves the model's ability to fit the data. The data points are then randomized through shuffling: this shuffling step ensures that both the training and testing sets are representative samples of the overall dataset, mitigating any potential biases. A bias term has been added by introducing a column of ones to the feature matrix. This step improves flexibility, allowing it to not pass from the origin, and fit, allowing the decision boundary to shift up or down. Finally, the dataset is split according to a specified training ratio set to 70% for training and 30% for testing.

## 3 Methodology

To ensure the robustness of my results and to measure the volatility of the model performance, each model was trained and tested 50 times. This approach allowed me to measure the variability in performance and obtain a more reliable estimate of each model's effectiveness. For each model, I followed the following procedure:

- **Initialization:** Each model was initialized with the same hyperparameters;
- **Training and Testing:** The model was trained on the training dataset and evaluated on the testing dataset 50 times. The dataset was split into training and testing sets for each run according to the fixed ratio;
- **Performance Measurement:** After each run, the  $\ell_2$  loss was recorded while the runtime was measured on all 50 runs.

In my implementation of the algorithms, I have included the predictions obtained during the training as an additional output. This feature is not typically provided by default in standard algorithm implementations.

### 3.1 Implemented Algorithms

#### 3.1.1 The Perceptron

##### Pseudocode

The Perceptron algorithm can be summarized with the following pseudocode:

---

**Algorithm 1** Perceptron Learning Algorithm

---

```
1: Input: Feature matrix  $X$ , target labels  $y$ , learning rate  $\eta$ , number of epochs  $T$ 
2: Initialize weights  $w \leftarrow 0$  ▷ Initialize weights to zero
3: for epoch = 1 to  $T$  do
4:   for each sample  $(x_i, y_i)$  in  $(X, y)$  do
5:     Compute prediction:  $\hat{y}_i \leftarrow \text{sign}(w \cdot x_i)$ 
6:     if  $\hat{y}_i \neq y_i$  then
7:       if  $\hat{y}_i = -1$  and  $y_i = 1$  then
8:         Update weights:  $w \leftarrow w + \eta \cdot x_i$ 
9:       else if  $\hat{y}_i = 1$  and  $y_i = -1$  then
10:        Update weights:  $w \leftarrow w - \eta \cdot x_i$ 
11:      end if
12:    end if
13:  end for
14: end for
15: Output: Final weights  $w$ , predictions  $\hat{y} \leftarrow \text{sign}(X \cdot w)$ 
```

---

#### 3.1.2 Support Vector Machines (SVMs) using Pegasos

**Pseudocode** The Pegasos (Primal Estimated sub-Gradient Solver for SVM) algorithm for training a Support Vector Machine (SVM) can be summarized with the following pseudocode:

---

**Algorithm 2** Pegasos Algorithm for SVM

---

```
1: Input: Feature matrix  $X$ , target labels  $y$ , regularization parameter  $\lambda$ , number of epochs  $T$ , batch size  $B$ 
2: Initialize weight vector  $w \leftarrow 0$  ▷ Initialize weights to zero
3: Initialize iteration counter  $t \leftarrow 0$ 
4: Initialize weight sum  $w_{\text{sum}} \leftarrow 0$ 
5: for epoch = 1 to  $T$  do
6:   for each batch in the dataset do
7:     Increment iteration counter:  $t \leftarrow t + 1$ 
8:     Compute learning rate:  $\eta \leftarrow \frac{1}{\lambda t}$ 
9:     Randomly select a batch of size  $B$ 
10:    for each sample  $(x_i, y_i)$  in the batch do
11:      if  $y_i(w \cdot x_i) < 1$  then
12:        Update weights:  $w \leftarrow (1 - \eta\lambda)w + \eta y_i x_i$ 
13:      else
14:        Update weights:  $w \leftarrow (1 - \eta\lambda)w$ 
15:      end if
16:    end for
17:    Update weight sum:  $w_{\text{sum}} \leftarrow w_{\text{sum}} + w$ 
18:  end for
19: end for
20: Compute final weights:  $w \leftarrow \frac{w_{\text{sum}}}{t}$ 
21: Output: Final weights  $w$ , predictions  $\hat{y} \leftarrow \text{sign}(X \cdot w)$ 
```

---

**Mathematical Formulation**

The goal of the SVM is to find a hyperplane that maximizes the margin between two classes while minimizing classification errors. The optimization problem can be expressed as:

$$\min_w \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \max(0, 1 - y_i(w \cdot x_i)) \quad (1)$$

where:

- $w$  is the weight vector to be optimized.
- $\lambda$  is the regularization parameter that controls the trade-off between margin maximization and error minimization.
- $y_i$  is the target label for the  $i$ -th training example.
- $x_i$  is the feature vector for the  $i$ -th training example.

The Pegasos algorithm approximates the solution to this optimization problem using stochastic gradient descent (SGD). At each iteration, a subset (batch) of the data is used to compute an estimate of the gradient, which is then used to update the weight vector. The learning rate  $\eta$  is adjusted over time according to the iteration count  $t$ :

$$\eta_t = \frac{1}{\lambda t} \quad (2)$$

To improve stability and reduce variance, the algorithm keeps track of the cumulative sum of the weight vectors and computes the final weight vector as the average of all weight vectors across iterations.

**3.1.3 Regularized Logistic Regression**

**Pseudocode** The Regularized Logistic Regression algorithm can be summarized with the following pseudocode:

---

**Algorithm 3** Regularized Logistic Regression with Stochastic Gradient Descent

---

```
1: Input: Feature matrix  $X$ , target labels  $y$ , regularization parameter  $\lambda$ , number of epochs  $T$ , batch  
   size  $B$ , initial learning rate  $\eta_0$ , learning rate decay  $\gamma$   
2: Initialize weights  $w \leftarrow$  random small values ▷ Initialize weights randomly  
3: Transform target labels:  $y \leftarrow (y + 1)/2$  ▷ Transform labels from  $\{-1, 1\}$  to  $\{0, 1\}$   
4: for epoch = 1 to  $T$  do  
5:   Shuffle the data  
6:   for each batch  $X_{\text{batch}}, y_{\text{batch}}$  in  $(X, y)$  do  
7:     Update learning rate:  $\eta \leftarrow \eta_0 / (1 + \gamma \cdot t)$  ▷ Decay learning rate  
8:     Compute logits:  $\text{logits} \leftarrow X_{\text{batch}} \cdot w$   
9:     Apply sigmoid function:  $\hat{y} \leftarrow \frac{1}{1 + \exp(-\text{logits})}$   
10:    Compute gradient:  $\text{gradient} \leftarrow \frac{1}{B} \cdot (X_{\text{batch}}^\top \cdot (\hat{y} - y_{\text{batch}})) + \lambda \cdot w$   
11:    Update weights:  $w \leftarrow w - \eta \cdot \text{gradient}$   
12:    Increment iteration counter  $t \leftarrow t + 1$   
13:  end for  
14: end for  
15: Transform predictions:  $\hat{y} \leftarrow 2 \cdot (\hat{y} > 0.5) - 1$  ▷ Transform  $\{0, 1\}$  back to  $\{-1, 1\}$   
16: Output: Final weights  $w$ , predictions  $\hat{y}$ 
```

---

### Mathematical Formulation

The Regularized Logistic Regression algorithm I implemented involves the following steps initialize the weight vector  $w$  with small random values. I did this because the algorithm was giving me performance problems. Initializing to a random small value instead of zero can cause faster convergence.

Then it convert target labels  $y$  from  $\{-1, 1\}$  to  $\{0, 1\}$  to use the logistic regression.

For each batch of data, it computes the logits using:

$$\text{logits} = X \cdot w \quad (3)$$

It applies the sigmoid function to convert logits to probabilities:

$$\hat{y} = \frac{1}{1 + \exp(-\text{logits})} \quad (4)$$

Notice that the exponential argument `logit` has been cropped with a value between  $[-709, +709]$  to avoid overflow and underflow. This is necessary because Python type `float64` can't handle values smaller than  $e^{-709}$  (they are considered 0) and bigger than  $e^{709}$ .

After the Gradient Computation and the Weight Update, the algorithm makes the final predictions by converting probabilities back to class labels  $\{-1, 1\}$ :

$$\hat{y} = 2 \cdot (\hat{y} > 0.5) - 1 \quad (5)$$

## 3.2 Polynomial Feature Expansion

In polynomial feature expansion of degree 2, we transform a dataset with original features  $(x_1, x_2, \dots, x_{10})$  to include additional features that capture non-linear relationships and interactions between the original features. This is accomplished by introducing new features that are:

### 1. Squared Features:

- For each original feature  $x_i$ , we include its square  $x_i^2$ . This allows the model to capture the quadratic effects of the features.

### 2. Pairwise Product Features:

- For each unique pair of features  $(x_i, x_j)$  where  $i \neq j$ , we include the product  $x_i \cdot x_j$ . This enables the model to learn interactions between different features.

The expanded dataset thus has:

- 10 original features
- 10 squared features
- 45 pairwise product features

This results in a total of  $10 + 10 + \frac{10 \times 9}{2} = 65$  features.

### 3.3 Kernel Methods

#### 3.3.1 Kernelized Perceptron

**Pseudocode** The Kernelized Perceptron algorithm can be summarized with the following pseudocode:

---

#### Algorithm 4 Kernel Perceptron

---

```

1: Input: Feature matrix  $x$ , target labels  $y$ , learning rate  $\eta$ , number of epochs  $E$ ,  $\sigma$  (for Gaussian
   kernel), kernel ('gaussian' or 'polynomial'),  $c$  (constant for polynomial kernel),  $d$  (degree for
   polynomial kernel)
2: Initialize weights (alphas)  $\alpha \leftarrow \mathbf{0}$  ▷ Initialize weights to zero
3: Compute Kernel Matrix:
4: if kernel = 'gaussian' then
5:   for  $i = 1$  to  $n$  do
6:     for  $j = i$  to  $n$  do
7:        $K[i, j] \leftarrow \exp\left(-\frac{\|x[i] - x[j]\|^2}{2\sigma^2}\right)$ 
8:     end for
9:      $K \leftarrow K + K^T - \text{diag}(K)$  ▷ Make matrix symmetric
10:  end for
11: else if kernel = 'polynomial' then
12:  for  $i = 1$  to  $n$  do
13:    for  $j = i$  to  $n$  do
14:       $K[i, j] \leftarrow (c + x[i] \cdot x[j])^d$ 
15:    end for
16:     $K \leftarrow K + K^T - \text{diag}(K)$  ▷ Make matrix symmetric
17:  end for
18: end if
19: for epoch = 1 to  $E$  do
20:  for  $i = 1$  to  $n$  do
21:    Compute prediction:  $\hat{y}_i \leftarrow \sum_{j=1}^n \alpha_j y_j K_{ij}$ 
22:    if  $y_i \cdot \hat{y}_i \leq 0$  then
23:      Update weight:  $\alpha_i \leftarrow \alpha_i + \eta$ 
24:    end if
25:  end for
26: end for
27: Compute final predictions:  $\hat{y} \leftarrow \text{sign}\left(\sum_{i=1}^n \alpha_i y_i K\right)$ 
28: Output: Weights (alphas)  $\alpha$ , predictions  $\hat{y}$ 

```

---

### Mathematical Formulation

The Kernel Perceptron algorithm extends the basic perceptron by incorporating a kernel function to map input features into a higher-dimensional space. This allows the perceptron to learn complex decision boundaries.

The kernel matrix  $K$  measures the similarity between pairs of input samples in the feature space transformed by the kernel function.

The Kernel matrix is always symmetric, the matrix  $K$  is adjusted by:

$$K \leftarrow K + K^T - \text{diag}(K) \quad (6)$$

where  $K^T$  is the transpose of  $K$  and  $\text{diag}(K)$  is a diagonal matrix containing the diagonal elements of  $K$ . This adjustment makes the computation faster.

### 3.3.2 Kernelized Pegasos for SVM

**Pseudocode** According to [1] The Kernelized Pegasos algorithm can be summarized with the following pseudocode:

---

**Algorithm 5** Kernel Pegasos Algorithm

---

```
1: Input: Dataset  $S$ , regularization parameter  $\lambda$ , number of iterations  $T$ 
2: Initialize dual variables (alphas)  $\alpha_1 \leftarrow \mathbf{0}$  ▷ Initialize alphas to zero
3: for  $t = 1$  to  $T$  do
4:   Choose  $i_t \in \{1, \dots, |S|\}$  uniformly at random ▷ Randomly select an index
5:   for all  $j = i_t$  do
6:     Set  $\alpha_{t+1}[j] \leftarrow \alpha_t[j]$ 
7:   end for
8:   if  $y_{i_t} \left( \frac{1}{\lambda t} \sum_j \alpha_t[j] y_{i_t} K(x_{i_t}, x_j) \right) < 1$  then
9:     Update  $\alpha$ :
10:      $\alpha_{t+1}[i_t] \leftarrow \alpha_t[i_t] + 1$ 
11:   else
12:     No update:
13:      $\alpha_{t+1}[i_t] \leftarrow \alpha_t[i_t]$ 
14:   end if
15: end for
16: Output: Dual variables  $\alpha_{T+1}$ 
```

---

However, there is a typo (red) in the pseudocode. Here's the corrected version:

---

**Algorithm 6** Kernel Pegasos Algorithm

---

```

1: Input: Dataset  $S$ , regularization parameter  $\lambda$ , number of iterations  $T$ 
2: Initialize dual variables (alphas)  $\alpha_1 \leftarrow \mathbf{0}$  ▷ Initialize alphas to zero
3: for  $t = 1$  to  $T$  do
4:   Choose  $i_t \in \{1, \dots, |S|\}$  uniformly at random ▷ Randomly select an index
5:   for all  $j = 1$  to  $|S|$  do
6:     Set  $\alpha_{t+1}[j] \leftarrow \alpha_t[j]$ 
7:   end for
8:   if  $y_{i_t} \left( \frac{1}{\lambda t} \sum_{j=1}^{|S|} \alpha_t[j] y_j K(x_{i_t}, x_j) \right) < 1$  then
9:     Update  $\alpha$ :
           
$$\alpha_{t+1}[i_t] \leftarrow \alpha_t[i_t] + 1$$

10:   else
11:     No update:
           
$$\alpha_{t+1}[i_t] \leftarrow \alpha_t[i_t]$$

12:   end if
13: end for
14: Output: Dual variables  $\alpha_{T+1}$ 

```

---

To determine whether a row is correctly classified or not, we need to evaluate the decision function using the labels  $y_j$  of all data points and not just  $y_{i_t}$ , the batch labels.

## 4 Hypertuning

The goal is to find the best set of hyperparameters that lead to the best model performance.

### 4.1 Hypertuning Process

The hypertuning process I used involves the following steps:

- **Parameter Grid Definition:** A grid of hyperparameters to explore was defined. This grid includes various values for each hyperparameter that will be evaluated;
- **Cross-Validation Splitting:** the dataset was split into  $k$  folds. For each fold, one part was used as the validation set, and the remaining parts were combined to form the training set;
- **Model Evaluation:** For each combination of hyperparameters, the model was trained using the training folds and evaluated on the validation fold;
- **Average Loss Calculation:** The average loss across all folds for each hyperparameter combination was computed to assess its performance;
- **Best Parameter Selection:** The combination of hyperparameters that results in the lowest average loss was identified;

### 4.2 Results for Each Algorithm

#### 4.2.1 Perceptron

For the Perceptron algorithm, the grid search was performed over learning rates and the number of epochs. The best hyperparameters found were:

- **Learning Rate ( $\eta$ ):** 0.001
- **Number of Epochs:** 500

#### 4.2.2 Pegasos

The grid search for the Pegasos algorithm involved tuning the regularization parameter, number of epochs, and batch size. The optimal hyperparameters were:

- **Regularization Parameter ( $\lambda$ ): 1**
- **Number of Epochs: 500**
- **Batch Size: 100**

#### 4.2.3 Logistic Regression

For Regularized Logistic Regression, the grid search focused on the regularization parameter, number of epochs, and batch size. The best hyperparameters were:

- **Regularization Parameter ( $\lambda$ ): 1**
- **Number of Epochs: 500**
- **Batch Size: 100**

#### 4.2.4 Perceptron (Polynomial Features)

For the Perceptron algorithm, the grid search was performed over learning rates and the number of epochs. The best hyperparameters found were:

- **Learning Rate ( $\eta$ ): 0.001**
- **Number of Epochs: 50**

Notice that it needs fewer epochs

#### 4.2.5 Pegasos (Polynomial Features)

The grid search for the Pegasos algorithm involved tuning the regularization parameter, number of epochs, and batch size. The optimal hyperparameters were:

- **Regularization Parameter ( $\lambda$ ): 1**
- **Number of Epochs: 500**
- **Batch Size: 100**

#### 4.2.6 Logistic Regression (Polynomial Features)

For Regularized Logistic Regression, the grid search focused on the regularization parameter, number of epochs, and batch size. The best hyperparameters were:

- **Regularization Parameter ( $\lambda$ ): 0.01**
- **Number of Epochs: 500**
- **Batch Size: 500**

#### 4.2.7 Kernel Perceptron

The hyperparameter tuning for the Kernel Perceptron algorithm involved optimizing the kernel type, learning rate, and kernel-specific parameters. The best hyperparameters for each kernel type are as follows:



**Gaussian Kernel:**

- **Learning Rate ( $\eta$ ):** 0.001
- **epochs:** 50
- **Width of the Gaussian function ( $\sigma$ ):** 2.0

**Polynomial Kernel:**

- **Learning Rate ( $\eta$ ):** 0.001
- **epochs:** 500
- **Bias ( $c$ ):** 1.0
- **Degree ( $d$ ):** 3

**4.2.8 Kernel Pegasos**

The hyperparameter tuning for the Kernel Pegasos algorithm involved optimizing parameters similar to those for Kernel Perceptron, adding regularization and batch size. The optimal hyperparameters for each kernel type are as follows:

**Gaussian Kernel:**

- **Regularization Parameter ( $\lambda$ ):** 0.001
- **epochs:** 100
- **Width of the Gaussian function ( $\sigma$ ):** 2.0
- **Batch Size:** 50

**Polynomial Kernel:**

- **Regularization Parameter ( $\lambda$ ):** 0.01
- **epochs:** 500
- **Batch Size:** 10
- **Bias ( $c$ ):** 0.5
- **Degree ( $d$ ):** 3

**5 Models Performance**

In this section, I'll show the performance of the models, based on:

- **Error:** measured with the 0-1 loss ( $\ell_{0-1}$ )

$$\ell_{0-1}(y_i, \hat{y}_i) = \begin{cases} 0 & \text{if } y_i = \hat{y}_i \\ 1 & \text{if } y_i \neq \hat{y}_i \end{cases} \quad (7)$$

The overall 0-1 loss  $L_{0-1}$  is then given by:

$$L_{0-1} = \frac{1}{n} \sum_{i=1}^n \ell_{0-1}(y_i, \hat{y}_i) \quad (8)$$

- **Error Variance:**

Let  $L_{0-1} = \{L_{0-1,1}, L_{0-1,2}, \dots, L_{0-1,10}\}$  be a list of 0-1 loss values. The variance of these values can be computed using:

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (L_{0-1,i} - \bar{L}_{0-1})^2 \quad (9)$$

where  $\bar{L}_{0-1}$  is the mean of the 0-1 loss values:

$$\bar{L}_{0-1} = \frac{1}{n} \sum_{i=1}^n L_{0-1,i} \quad (10)$$

- **Runtime:** the time it takes to train and test each algorithm ( $T$ ).

Every algorithm has been run **50 times**, to detect the variance of the Error.

## 5.1 Perceptron, Pegasos and Regularized Logistic Classification

### Perceptron

#### Training

- $L_{0-1}$ : 36.45%
- $\sigma^2$ : 0.0063

$T$ : 0m 0s 596ms

#### Test

- $L_{0-1}$ : 36.69%
- $\sigma^2$ : 0.0060

### Pegasos

#### Training

- $L_{0-1}$ : 45.22%
- $\sigma^2$ : 0.0063

$T$ : 0m 1s 378

#### Test

- $L_{0-1}$ : 44.61%
- $\sigma^2$ : 0.0060

### Regularized Logistic Classification

#### Training

- $L_{0-1}$ : 29.84%
- $\sigma^2$ : 0.0007

$T$ : 0m 0s 72ms

#### Test

- $L_{0-1}$ : 29.52%
- $\sigma^2$ : 0.0013

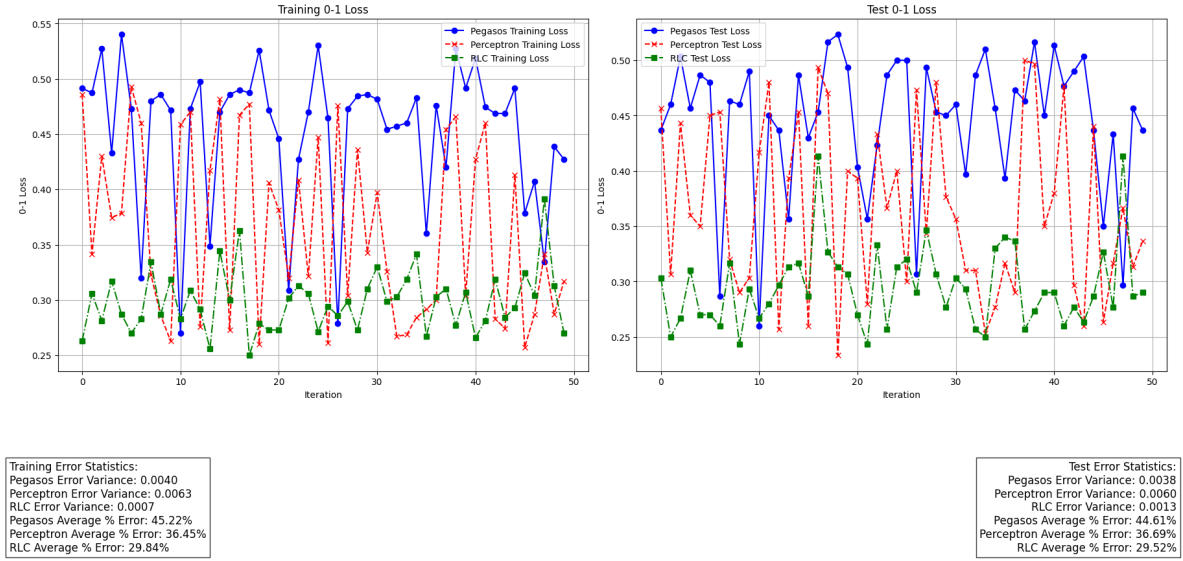


Figure 1: Perceptron, Pegasos and Regularized Logistic Classification Error Distribution

The **Regularized Logistic Classification** is by far the best-performing algorithm among the three. It excels in error rate, variance and runtime, in both training and test.

## 5.2 Perceptron, Pegasos and Regularized Logistic Classification (Polynomial Expansion)

### Perceptron

#### Training

- $L_{0-1}$ : 35.73%
- $\sigma^2$ : 0.0057

$T$ : 0m 0s 569ms

#### Test

- $L_{0-1}$ : 35.57%
- $\sigma^2$ : 0.0044

### Pegasos

#### Training

- $L_{0-1}$ : 44.15%
- $\sigma^2$ : 0.0055

$T$ : 0m 1s 680ms

#### Test

- $L_{0-1}$ : 44.74%
- $\sigma^2$ : 0.0071

### Regularized Logistic Classification

#### Training

- $L_{0-1}$ : 28.75%
- $\sigma^2$ : 0.0007

$T$ : 0m 0s 167ms

#### Test

- $L_{0-1}$ : 29.58%
- $\sigma^2$ : 0.0006

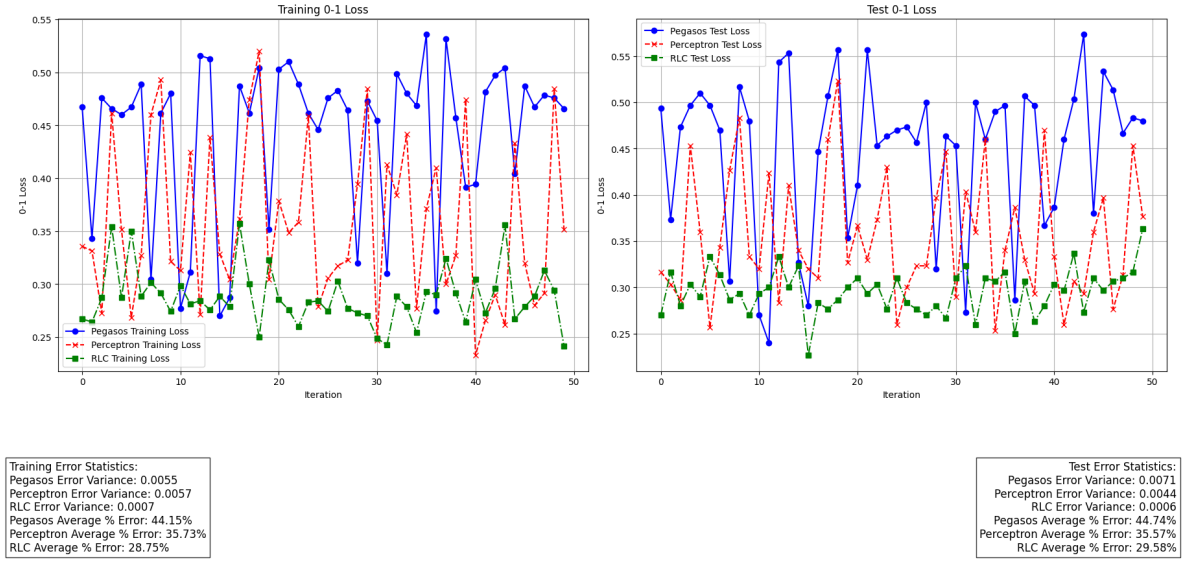


Figure 2: Perceptron, Pegasos and Regularized Logistic Classification Error Distribution after Polynomial Expansion

Perceptron shows a slight improvement in training error after polynomial expansion, but the test error remains similar. Pegasos experiences a slight improvement in training error but a slight increase in test error. Regularized Logistic Classification shows the lowest training error after polynomial the expansion. However, its test error slightly increases.

The Runtime is similar and the differences are probably due to unrelated machine processes.

In summary, **Regularized Logistic Classification** remains the best-performing algorithm in both scenarios, demonstrating better accuracy and consistency.

### 5.3 Kernelized Perceptron

#### Gaussian

##### Training

- $L_{0-1}$ : 0%
- $\sigma^2$ : 0

$T$ : 0m 2s 383ms

#### Polynomial

##### Training

- $L_{0-1}$ : 34.87%
- $\sigma^2$ : 0.0081

$T$ : 0m 3s 127ms

##### Test

- $L_{0-1}$ : 10.77%
- $\sigma^2$ : 0.0002

##### Test

- $L_{0-1}$ : 47.09%
- $\sigma^2$ : 0.0051

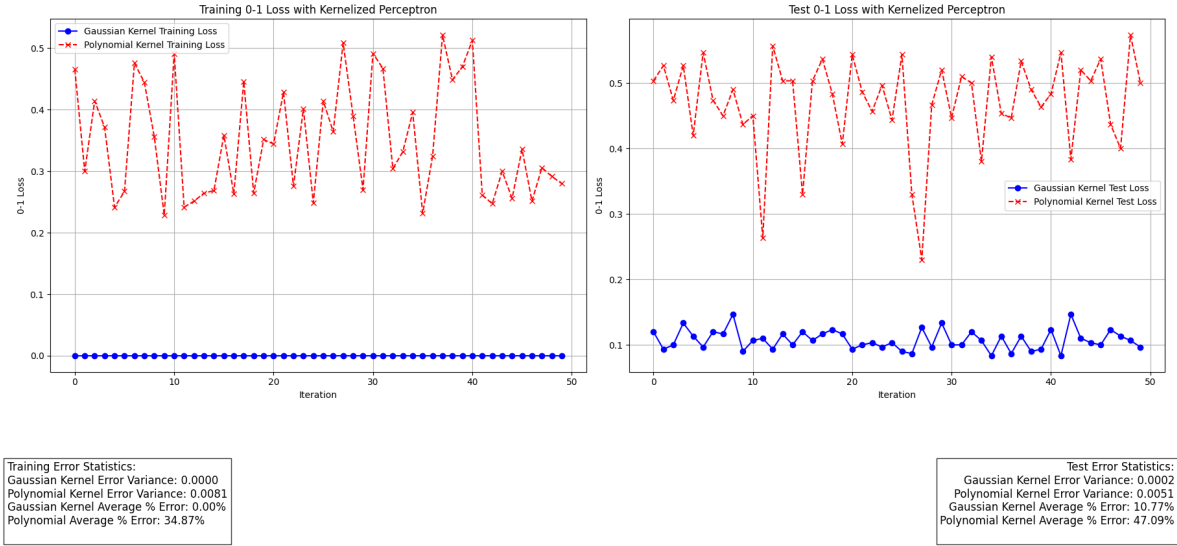


Figure 3: Kernelized Perceptron with Gaussian and Polynomial Kernel

The Kernelized Perceptron performed better with a **Gaussian Kernel** in both training and test indices. It also has a lower runtime. However, the training error of 0% suggests potential overfitting. Besides, the algorithm has the lowest test error even when comparing it to the previous algorithms.

## 5.4 Kernelized Pegasos

### Gaussian

#### Training

- $L_{0-1}$ : 0%
- $\sigma^2$ : 0

$T$ : 0m 2s 711ms

### Polynomial

#### Training

- $L_{0-1}$ : 35.74%
- $\sigma^2$ : 0.0069

$T$ : 0m 3s 926ms

#### Test

- $L_{0-1}$ : 9.65
- $\sigma^2$ : 0.0003

#### Test

- $L_{0-1}$ : 48.87%
- $\sigma^2$ : 0.0024

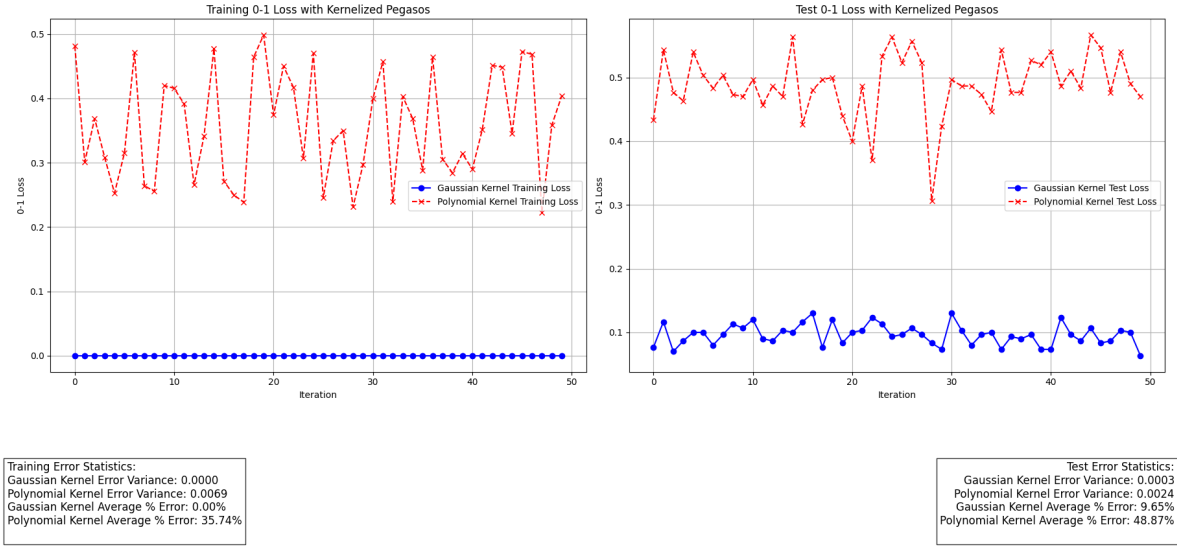


Figure 4: Kernelized Pegasos with Gaussian and Polynomial Kernel

Kernelized Pegasos and Kernelized Perceptron shows similar results: the **Gaussian kernel** achieves a perfect 0% which might suggest potential overfitting, as the model may have learned the training data too well. Nevertheless, the Gaussian kernel performs well on the test set, with a relatively low error rate of 9.65% and the tiniest variance. This means that, despite the risk of overfitting, the model generalizes effectively to unseen data. In contrast, the Polynomial kernel exhibits a much higher training error of 35.74% and a significantly higher test error of 48.87%. The gap between training and test errors combined with a test error close to random classification, indicates that the polynomial kernel doesn't perform well with this dataset.

## 6 Conclusion

In this report, I have evaluated the following classification algorithms: Perceptron, Pegasos, Regularized Logistic Classification, Kernelized Perceptron, and Kernelized Pegasos based on their performance in terms of error, error variance, and runtime.

The Perceptron shows decent performance in terms of error but is outperformed by the Regularized Logistic Classification. In addition, it has a relatively high variance.

Pegasos shows higher error rates compared to the other algorithms, both in the training and testing phases. This indicates that it may not be as effective at minimizing classification errors. In addition, its runtime it's higher than the other similar algorithms in terms of complexity.

**The Regularized Logistic Classification algorithm consistently outperforms the Perceptron and Pegasos** in terms of error and variance. It also shows the fastest runtime. Using a polynomial expanded dataset, Perceptron shows a slight improvement in training error after polynomial expansion, but the test error remains similar. Pegasos experiences a slight improvement in training error but a slight increase in test error. On the other hand, Regularized Logistic Classification shows the lowest training error after polynomial expansion. However, its test error slightly increases. The Runtime is similar for all the algorithms before and after the polynomial expansion; the differences are probably due to unrelated machine processes. **Regularized Logistic Classification still is the best-performing algorithm** in both scenarios.

The Kernelized Perceptron with the Gaussian kernel achieves a perfect training error of 0%, which suggests **overfitting**. Despite this, it performs well on the test set with a low error rate observed (10.77%). The Polynomial kernel, however, shows a suboptimal performance with high training and test errors.

Similar to Kernelized Perceptron, the Kernelized Pegasos with the Gaussian kernel performs excellently in terms of test error and variance, showing the tiniest numbers (9.65% and 0.0003). Nevertheless, it shows a training error of 0% indicative of possible **overfitting**. The Polynomial kernel performs poorly,

with both high training and test errors.

## 6.1 Evaluation

The **Regularized Logistic Classification is the most consistent and efficient algorithm across the board**: it offers the best balance of low error rates, low variance, and efficient run-time. The Kernelized Perceptron and Kernelized Pegasos with Gaussian may suffer from overfitting due to perfect training errors. However, they show the best test performances. On the other hand, the Polynomial kernels in both Kernelized Perceptron and Pegasos, as well as Pegasos in general, are less effective, exhibiting error rates close to the random classification and high variances.

## 7 References

### References

- [1] Shalev-Shwartz, S., Singer, Y., Srebro, N. et al. Pegasos: primal estimated sub-gradient solver for SVM. *Mathematical Programming*, 127, 3–30 (2011). <https://doi.org/10.1007/s10107-010-0420-4>.

## 8 Appendix

I declare that this material, which I now submit for assessment, is entirely my own work and has not been taken from the work of others, save and to the extent that such work has been cited and acknowledged within the text of my work. I understand that plagiarism, collusion, and copying are grave and serious offences in the university and I accept the penalties that would be imposed should I engage in plagiarism, collusion or copying. This assignment, or any part of it, has not been previously submitted by me or any other person for assessment on this or any other course of study.