

CS-534 Machine Learning

Implementation Assignment 1

Author 1 (%) Author 2 (%) Author 3 (%)

Introduction

Given a set of M training pairs $\{\mathbf{x}_i, y_i\}_{i=1}^M$, where \mathbf{x}_i is a set of features $\{x_{ij}\}_{j=1}^N$ and y_i is the corresponding target variable, we wish to find a set of parameters $\mathbf{w} = \{w_1, w_2, \dots, w_n\}$ that minimizes the regularized Sum of Squared Error (SSE) objective

$$J(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^m (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 + \lambda \|\mathbf{w}\|^2.$$

Here, the regularization term $\lambda \|\mathbf{w}\|^2$ is used to counteract overfitting: a common artifact caused by high dimensional features. Overfitting is often manifest in large parameters, thus the regularization term to encourage small weights. λ is a hyper-parameter that controls the influence of regularization on these weights and subsequent objective value and solution. This value is often not known and difficult to determine. The objective of this assignment is to implement the gradient descent algorithm (GDA) on a high dimensional training set and investigate the effect of the regularization parameter via cross-validation and validation using an independent dataset.

Gradient descent algorithm

To implement GDA, we first derive the gradient of the cost function $J(\mathbf{w})$. The gradient vector is of size M , the number of samples in the training set, and each element is equal to the partial derivative of the cost function $J(\mathbf{w})$ with respect to the parameters $\{w_i\}_{i=1}^N$. The chain rule gives

$$\nabla J(\mathbf{w}) = \left[\frac{\partial J}{\partial w_1}, \frac{\partial J}{\partial w_2}, \dots, \frac{\partial J}{\partial w_n} \right]^\top. \quad (1)$$

Thus, the k^{th} element of the gradient vector is

$$\frac{\partial J}{\partial w_k} = \sum_{i=1}^m (\mathbf{w}^\top \mathbf{x}_i - y_i) x_{ik} + 2\lambda w_k. \quad (2)$$

This gives the following regularized gradient of the cost function $J(\mathbf{w})$:

$$\nabla J(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m (\mathbf{w}^\top \mathbf{x}_i - y_i) \mathbf{x}_i + 2\lambda \mathbf{w}. \quad (3)$$

Note that we have scaled the first term by the number of samples in the equation to represent the average gradient per weight. The gradient is now incorporated into the update rule for gradient descent. For a given training pair $\{\mathbf{x}_i, y_i\}$, the batch update rule is

$$w_j =: w_j - \alpha ((\mathbf{w}^\top \mathbf{x}_i - y_i) x_{ij} + 2\lambda w_j). \quad (4)$$

Equations 3 and 4 are used in a batch gradient descent algorithm where the parameters \mathbf{w} are updated simultaneously across features until the norm of $\nabla J(\mathbf{w})$ converges at a small value ϵ . To improve the speed of convergence we will also standardize the features by $s(\mathbf{x}) = \frac{\mathbf{x} - \bar{\mathbf{x}}}{\sigma}$.

Algorithm 1: Batch Gradient Descent

```

1  $\mathbf{w} = \mathbf{w}^0$ ;
2 do
3    $\nabla J(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m (\mathbf{w}^\top \mathbf{x}_i - y_i) \mathbf{x}_i + 2\lambda \mathbf{w}$ ;
4    $\mathbf{w} =: \mathbf{w} - \alpha \nabla J(\mathbf{w})$ ;
5 while  $|\nabla J(\mathbf{w})| \geq \epsilon$ ;
```

In the algorithm, parameters \mathbf{w} are initialized to the zero vector of size M : the number of features in \mathbf{x}_i . These parameters are iteratively updated by the regularized gradient of $J(\mathbf{w})$ scaled by the hyper-parameter α , called the learning rate. An inherent trade-off exists in the selection of this parameter. If the learning rate is too small the algorithm will be slow to converge or may not converge at all. A large learning can result in overstepping the optimum and oscillating out of control to infinity. Prior to tuning the regularization parameter, we will explore different learning rates.

Learning rate selection (10 pts)

Tuning the regularizer (20 pts)

10-fold cross-validation (30 pts)

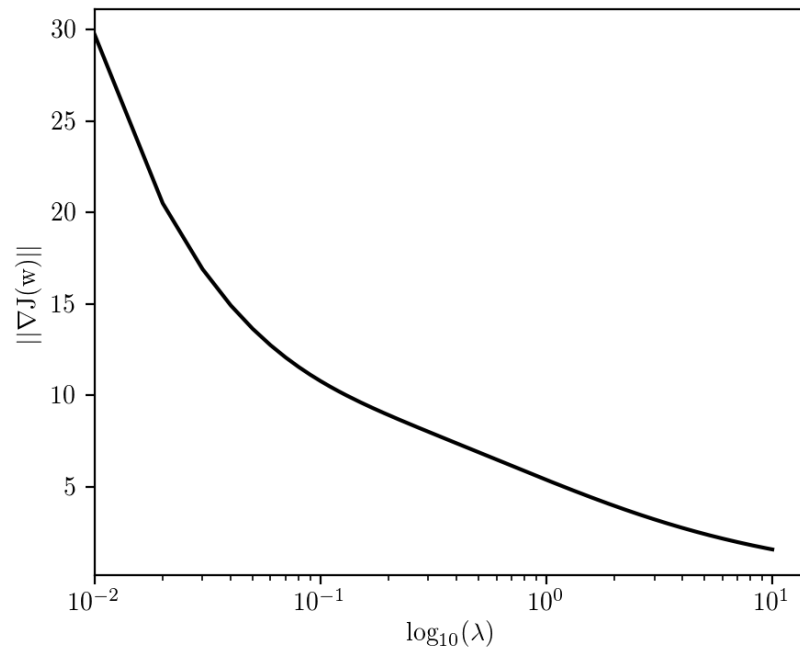


Figure 1: The effect of different λ values on the norm of the gradient.

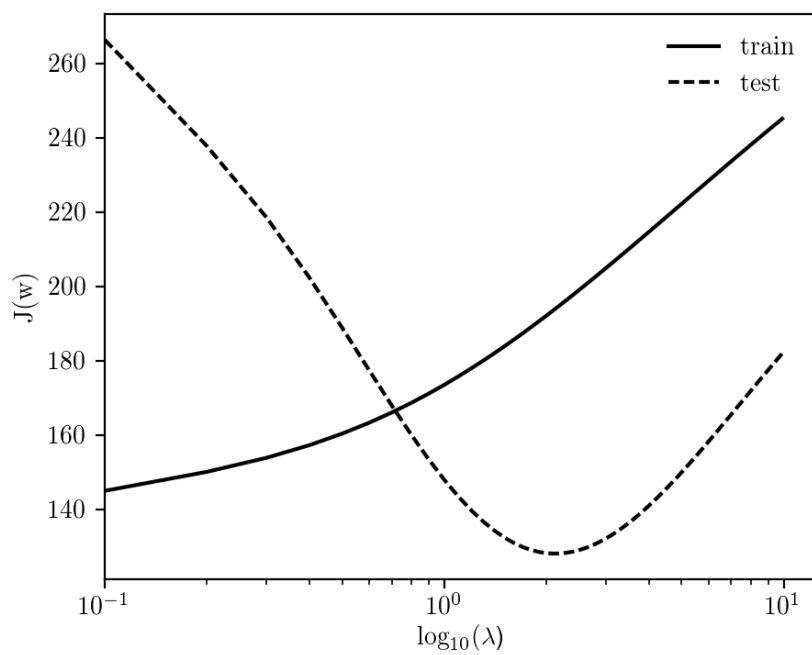


Figure 2: The effect of different λ values on the training SSE (solid) and the test SSE (dashed)