

236606 - Deep Learning - Homework 1

Spring 2018

Submission Instruction

Submission deadline: **24/04/2018 23:59**

TA in charge: Izik Golan

- Submit a **zip** file containing **only** your code and the written answers, in the webcourse site: <http://webcourse.cs.technion.ac.il/236606/>
- Answers must be submitted as a PDF file, edited by a document editor (**not** a scan of handwritten answers)
- Submission is **pairs only**
- There will be **no** late submissions

1 Distance of Point to Hyperplane

In tutorial 1, we calculated the distance of a point from a hyperplane using geometry. The same result can be achieved with constrained optimization.

- A. Given a hyperplane, $H_{\mathbf{w},b} = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{w}^\top \mathbf{x} + b = 0\}$, and a vector $\mathbf{x}_0 \in \mathbb{R}^n$, write the constrained optimization problem whose solution is the Euclidean distance of \mathbf{x}_0 to $H_{\mathbf{w},b}$. Use the following standard form for your optimization problem:

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}, \mathbf{x}_0, \mathbf{w}, b) \\ \text{subject to} \quad & g(\mathbf{x}, \mathbf{x}_0, \mathbf{w}, b) = 0 \end{aligned}$$

- B. Using the method of Lagrange multipliers, solve the optimization problem from 1(A) and show that the solution is identical to the one we showed in Tutorial 1. You can read about the method of Lagrange multipliers at https://en.wikipedia.org/wiki/Lagrange_multiplier.

2 Perceptron

In Lecture 2, we proved that the perceptron algorithm makes at most $1/(\gamma)^2$ mistakes on any sequence of labeled examples that is contained inside the unit hyper-sphere, and linearly-separable by a margin γ . The algorithm does **not** guarantee anything about the margin of the final classifier.

Assume that \mathcal{X} is in the unit hyper-sphere. Let $S \subseteq \mathbb{R}^d \times \{-1, +1\}$ be a sequence of labeled examples. Consider the following variation of the vanilla perceptron algorithm, where the only difference is that we now update \mathbf{w}_{t+1} also for margin error (examples that are closer than $\gamma/2$ to the decision boundary), and not just prediction mistakes:

Margin Perceptron algorithm

```
procedure MARGINPERCEPTRON( $S, \gamma$ )
   $\mathbf{w}_1 \leftarrow \vec{0}$ 
  for  $t = 1, \dots, T$  do
     $\mathbf{x}_t, y_t \leftarrow \text{next}(S)$ 
    if  $\frac{\mathbf{w}_t \cdot \mathbf{x}_t}{\|\mathbf{w}_t\|} \in (-\frac{\gamma}{2}, +\frac{\gamma}{2})$  or  $y_t (\mathbf{w}_t \cdot \mathbf{x}_t) \leq 0$  then
       $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + y_t \mathbf{x}_t$ 
    else
       $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t$ 
    end if
  end for
  return  $\mathbf{w}_{T+1}$ 
end procedure
```

In Lecture 2, we proved that if there exists a unit-vector $\mathbf{w}^* \in \mathbb{R}^d$ such that $\forall (\mathbf{x}, y) \in S, y(\mathbf{w}^* \cdot \mathbf{x}) > 0$, then the number of mistakes M on S made by the online Perceptron algorithm is at most $1/(\gamma)^2$, where

$$\gamma = \min_{\mathbf{x} \in S} |\mathbf{w}^* \cdot \mathbf{x}|$$

- A. Prove that if \mathbf{w}^* is not a unit vector, we can normalize both the sequence of examples, S , and \mathbf{w}^* to be on the unit hyper-sphere, without changing the predictions made by \mathbf{w}^* on each sample.
- B. Prove that if there exists a linear separator with margin γ for S , then the margin perceptron algorithm will make at most $\frac{8}{\gamma^2}$ mistakes on S . **Hint:** prove that $\|\mathbf{w}_{t+1}\| \leq \|\mathbf{w}_t\| + \frac{1}{2\|\mathbf{w}_t\|} + \frac{\gamma}{2}$, and see what happens when $\|\mathbf{w}_t\| \geq \frac{2}{\gamma}$

3 Ridge Regression

Consider the data set $S_m = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$ where $\mathbf{x}_i \in \mathbb{R}^n, y_i \in \mathbb{R}$. Also, assume the samples are centered, i.e., $\forall j = 1, \dots, n : \frac{1}{m} \sum_{i=1}^m (\mathbf{x}_i)_j = 0$. Ridge regression is a variant of linear regression, where a regularization term is added to the loss function. This is also called Tichonov regularization. The minimization objective (loss) for an ERM training strategy is

$$L(\mathbf{w}, b) = \frac{1}{2m} \|\mathbf{X}\mathbf{w} + b\mathbf{1} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_2^2,$$

where \mathbf{X} is the data matrix with row i being \mathbf{x}_i , \mathbf{y} being a vector where $\mathbf{y}_i = y_i$, and $\mathbf{1}$ is a vector of 1s of size m .

- A. Compute $\frac{\partial L}{\partial b}$. Simplify as much as possible. **Hint:** This should not be a function of \mathbf{w} or \mathbf{X}
- B. Compute $\frac{\partial L}{\partial \mathbf{w}}$. Simplify as much as possible. Write your answer in vector form.
- C. Suppose that $\lambda > 0$. Infer the solution obtained by equating the derivatives from 3(A) and 3(B) to 0. Are these solutions a global optimum for L ?
- D. In 3(C), you need to take an inverse of some matrix \mathbf{M} to get the solution for \mathbf{w} . It is well known that a matrix \mathbf{M} is invertible if and only if it is *positive definite*, meaning that $\forall \mathbf{v} \neq \mathbf{0} \in \mathbb{R}^n : \mathbf{v}^\top \mathbf{M} \mathbf{v} > 0$. Prove that if $\lambda > 0$ then \mathbf{M} is invertible. You can read more about positive definite and semi-definite matrices at:
https://en.wikipedia.org/wiki/Positive-definite_matrix.
- E. What advantage of ridge regression over vanilla linear regression does 3(D) highlight?

4 Binary Classification Using Ridge Regression

Use the following script to obtain the MNIST dataset:

```
import pickle, gzip, numpy, urllib.request, json

data_url = "http://deeplearning.net/data/mnist/mnist.pkl.gz"
# Load the dataset
urllib.request.urlretrieve(data_url, "mnist.pkl.gz")
with gzip.open('mnist.pkl.gz', 'rb') as f:
    train_set, valid_set, test_set = pickle.load(f, encoding='latin1')
```

The MNIST dataset contains grayscale images of handwritten digits (0-9). Each image is 28x28 pixels, each of which with a label corresponding to the written digit in it. For the purpose of linear regression, we flatten each image to be a vector of length 784. Since we are dealing with binary classification, consider all even digits to have label +1, and all odd digits to have label -1.

Throughout this exercise, make sure to take advantage of NumPy vectorizations as much as you can.

- A. Pre-process the dataset by subtracting from each sample (including samples from the validation and test sets) the vector which is the mean of all vectors in the training set.
- B. Build a linear regressor by minimizing the objective from question 3 w.r.t the training set. Implement both the analytic solution and a gradient descent solution with 100 gradient steps and learning rate 0.001. Train the regressor on the training data using 8 values of λ taken from a logarithmic scale, $[10^{-5}, 10^2]$ (i.e., take powers of 10).
- C. For the purpose of classification, the classifier which is based on the linear regressor is given by $h_{\mathbf{w},b}(\mathbf{x}) = \text{sign}(\mathbf{x} \cdot \mathbf{w} + b)$. We define the 0-1 loss and squared loss on a set of m labeled examples as:

$$L_{0-1}(\mathbf{w}, b) = \frac{1}{m} \sum_{i=1}^m \mathbf{I}(h_{\mathbf{w},b}(\mathbf{x}_i) \neq y_i)$$
$$L_{\text{squared}}(\mathbf{w}, b) = \frac{1}{m} \|\mathbf{X} \cdot \mathbf{w} + b\mathbf{1} - \mathbf{y}\|_2^2$$

For each value of λ , and each training scheme (gradient descent and analytic), What is the 0-1 loss of the training and validation sets? What is the squared loss?

- D. For each training scheme, write the 0-1 loss and squared loss on the test set, of the model which performed best on the validation set in terms of 0-1 loss.
- E. For the value of λ you found to perform best on the validation set in terms of 0-1 loss using the gradient descent scheme, plot the learning curves consisting of the train and test losses which are defined in 4(C), as a function of the number of gradient steps.
- F. Describe the simplest possible scenario you can think of for which the use of linear regression for classification is a bad idea?

Good Luck!