

Derivation of Newton Algorithm

Exercise 1: Part 1: Newton Algorithm

In Newton algorithm we iteratively update weight vector, using the following equation:

$$\tilde{w}_{t+1} \leftarrow \tilde{w}_t - \eta_t [\nabla^2 f(\tilde{w}_t)]^{-1} \cdot \nabla f(\tilde{w}_t)$$

where $\nabla^2 f(\tilde{w}_t)$ is Hessian and $\nabla f(\tilde{w}_t)$ is gradient.

Given: $f(\tilde{w}) = \frac{1}{h_+} \sum_{i:y_i=1} \log(1 + \exp(-y_i \tilde{w}^T \tilde{x}_i)) + \frac{1}{h_-} \sum_{j:y_j=-1} \log(1 + \exp(-y_j \tilde{w}^T \tilde{x}_j)) + \lambda \|\tilde{w}\|_2^2$

$$P_i = \frac{1}{1 + \exp(-\tilde{w}^T \tilde{x}_i)} \rightarrow \text{Sigmoid function}$$

$$\begin{aligned} \nabla f(\tilde{w}) &= \frac{df(\tilde{w})}{d\tilde{w}} = \frac{1}{h_+} \sum_{i:y_i=1} \frac{-x_i y_i \exp(-y_i \tilde{w}^T \tilde{x}_i)}{1 + \exp(-y_i \tilde{w}^T \tilde{x}_i)} + \frac{1}{h_-} \sum_{j:y_j=-1} \frac{-x_j y_j \exp(-y_j \tilde{w}^T \tilde{x}_j)}{1 + \exp(-y_j \tilde{w}^T \tilde{x}_j)} + 2\lambda \tilde{w} = \\ &= \frac{1}{h_+} \sum_{i:y_i=1} \left(-y_i \tilde{x}_i + \frac{y_i \tilde{x}_i}{1 + \exp(-y_i \tilde{w}^T \tilde{x}_i)} \right) + \frac{1}{h_-} \sum_{j:y_j=-1} \left(-y_j \tilde{x}_j + \frac{y_j \tilde{x}_j}{1 + \exp(-y_j \tilde{w}^T \tilde{x}_j)} \right) + 2\lambda \tilde{w} = \\ &= \frac{1}{h_+} \sum_{i:y_i=1} \left(-\tilde{x}_i + \frac{\tilde{x}_i}{1 + \exp(-\tilde{w}^T \tilde{x}_i)} \right) + \frac{1}{h_-} \sum_{j:y_j=-1} \left(\tilde{x}_j - \frac{\tilde{x}_j}{1 + \exp(-\tilde{w}^T \tilde{x}_j)} \right) + 2\lambda \tilde{w} = \\ &= \frac{1}{h_+} \sum_{i:y_i=1} \left(-\tilde{x}_i + \frac{\tilde{x}_i}{1 + \exp(-\tilde{w}^T \tilde{x}_i)} \right) + \frac{1}{h_-} \sum_{j:y_j=-1} \left(\tilde{x}_j - \tilde{x}_j \left(1 - \frac{1}{1 + \exp(-\tilde{w}^T \tilde{x}_j)} \right) \right) + 2\lambda \tilde{w} = \\ &= \frac{1}{h_+} \sum_{i:y_i=1} \left(\tilde{x}_i \left(-1 + \frac{1}{1 + \exp(-\tilde{w}^T \tilde{x}_i)} \right) \right) + \frac{1}{h_-} \sum_{j:y_j=-1} \left(\tilde{x}_j \frac{1}{1 + \exp(-\tilde{w}^T \tilde{x}_j)} \right) + 2\lambda \tilde{w} = \\ &= \boxed{\frac{1}{h_+} \sum_{i:y_i=1} \tilde{x}_i (p_i - 1) + \frac{1}{h_-} \sum_{j:y_j=-1} \tilde{x}_j p_j + 2\lambda \tilde{w}} \end{aligned}$$

In a vectorized form it can be expressed as:

$$\nabla f(\tilde{w}) = \frac{1}{h_+} X_{\text{pos}}^T (p_{\text{pos}} - 1) + \frac{1}{h_-} X_{\text{neg}}^T (p_{\text{neg}}) + 2\lambda \tilde{w}$$

where

$$X_{\text{pos}} = X[y=1]$$

$$X_{\text{neg}} = X[y=-1]$$

$$p_{\text{pos}} = P[y=1]$$

$$p_{\text{neg}} = P[y=-1]$$

Implementation attached

Alternatively, gradient can be expressed as follows:

$$\nabla f(\underline{w}) = \frac{1}{n_+} \left(\frac{y+1}{2} \right) \times (\underline{p} - \underline{1}) + \frac{1}{n_-} \left(\frac{1-y}{2} \right) \times \underline{p} + 2\lambda \underline{w}$$

Hessian

$$\begin{aligned} \nabla^2 f(\underline{w}) &= \frac{d \nabla f(\underline{w})}{d \underline{w}} = \frac{d}{d \underline{w}} \left(\frac{1}{n_+} \sum_{i:y_i=1} \underline{x}_i (\underline{p}_i - \underline{1}) + \frac{1}{n_-} \sum_{j:y_j=-1} \underline{x}_j \underline{p}_j + 2\lambda \underline{w} \right) = \\ &= \boxed{\frac{1}{n_+} \sum_{i:y_i=1} \underline{x}_i \underline{x}_i^T \underline{p}_i (1-\underline{p}_i) + \frac{1}{n_-} \sum_{j:y_j=-1} \underline{x}_j \underline{x}_j^T \underline{p}_j (1-\underline{p}_j) + 2\lambda I} = \\ &= \boxed{\frac{1}{n_+} \underline{X}_{pos}^T (\underline{p}_{pos} (1-\underline{p}_{pos})) \cdot \underline{X}_p + \frac{1}{n_-} \underline{X}_{neg}^T (\underline{p}_{neg} (1-\underline{p}_{neg})) \underline{X}_{neg} + 2\lambda I} \end{aligned}$$

where $\underline{X}_{pos} = X[y=1] \rightarrow$ rows of matrix X corresponding to $y=1$

$\underline{X}_{neg} = X[y=-1] \rightarrow$ rows of matrix X corresponding to $y=-1$

$$\underline{P} = \begin{bmatrix} \underline{p}_1 \\ \underline{p}_2 \\ \vdots \\ \underline{p}_n \end{bmatrix}; \quad \underline{P}_{pos} = \underline{P}[y=1]; \quad \underline{P}_{neg} = \underline{P}[y=-1]$$

Logistic Regression on CIFAR-10 data set

1. One-vs-All strategy for multiclass classification

The one-vs-all strategy trains 10 classifiers, each time with some c-th class as positive and the rest as negative. For prediction, we run all 10 classifiers and predict with the classifier that has the highest dot product (probability/confidence).

The implemented Newton Algorithm was tested on the CIFAR-10 dataset, using the one-vs-all strategy. Only the first training batch was used for training. The table below summarizes the results of this test.

Regression constant λ	Percent error
0.1	0.7405
1	0.7103
10	0.6815

The results of all tests are significantly better than random prediction, which would give the percent error of 0.9. Out of three values of regularization constant that was tested, $\lambda = 10$ yielded the best performance with 0.6815 percent error, or the accuracy rate of 0.3185.

The implementation of the algorithm is attached to the assignment (in ipython notebook).

2. One-vs-One strategy for multiclass classification

The one-vs-one strategy trains 45 classifiers, each time with some c-th class as positive and some c'-th class as negative. For prediction, we run all 45 classifiers, each of which will vote either for some class c or some class c'. Predict with the class that has most votes (break ties arbitrarily).

The implemented Newton Algorithm was tested on the CIFAR-10 dataset, using the one-vs-one strategy. Only the first training batch was used for training. The table below summarizes the results of this test.

Regression constant λ	Percent error
0.1	0.692
1	0.686
10	0.6668

The results of all tests are significantly better than random prediction, which would give the percent error of 0.9. Similarly to the one-vs-all strategy, out of three

values of regularization constant that was tested, $\lambda = 10$ yielded the best performance with 0.6668 percent error, or the accuracy rate of 0.3332.

The implementation of the algorithm is attached to the assignment (in ipython notebook).

Kernel Logistic Regression on CIFAR-10 data set

The implementation of the kernel binary algorithm is attached in ipython notebook. The tables below summarize the results of the test with three different kernels used.

Regression constant λ	Percent error
0.1	0.438
1	0.438
10	0.438

Table 3.1: Linear kernel

For linear kernel the results of all tests are better than random prediction, which would give the percent error of 0.5. The percent error turned out to be 0.438 for all values of λ that were tested.

Regression constant λ	Percent error
0.1	0.4395
1	0.439
10	0.4365

Table 3.2: Polynomial kernel

For polynomial kernel the results of all tests are better than random prediction, which would give the percent error of 0.5. Out of all tests, the smallest percent error of 0.4365 was obtained when $\lambda=10$.

Regression constant λ	Percent error $\sigma = 1$	Percent error $\sigma = 5$	Percent error $\sigma = 10$
0.1	0.4305	0.4305	0.4305
1	0.453	0.453	0.453
10	0.4625	0.4625	0.4625

Table 3.3: Gaussian kernel

For Gaussian kernel the results of all tests are better than random prediction, which would give the percent error of 0.5. The results were the same for different sigma values. Out of all tests, the smallest percent error of 0.4305 was obtained when $\lambda=0.1$.

Derivation of Kernel Logistic Regression

Part 4: Kernel Logistic Regression

Kernelize: $\tilde{x}_i \rightarrow \phi(\tilde{x}_i)$.

Then the objective function becomes:

$$f(\tilde{w}) = \frac{1}{n_+} \sum_{i:y_i=1} \log(1 + \exp(-y_i \tilde{w}^T \phi(\tilde{x}_i))) + \frac{1}{n_-} \sum_{j:y_j=-1} \log(1 + \exp(y_j \tilde{w}^T \phi(\tilde{x}_j))) + \lambda \|\tilde{w}\|^2$$

$\hookrightarrow \tilde{w} \in \mathbb{R}^n$

$$\tilde{w} = \sum_i \lambda_i \phi(\tilde{x}_i) \Rightarrow \tilde{w}^T \phi(\tilde{x}_i) = \sum_i \lambda_i \underbrace{\phi(\tilde{x}_i)^T \phi(\tilde{x}_i)}_K = \lambda^T K \cdot i \Rightarrow$$

$$\Rightarrow f(\lambda) = \frac{1}{n_+} \sum_{i:y_i=1} \log(1 + \exp(-y_i \lambda^T K \cdot i)) + \frac{1}{n_-} \sum_{j:y_j=-1} \log(1 + \exp(-y_j \lambda^T K \cdot i)) + \lambda^T K \lambda$$

This expression has exactly the same form, as the objective function in part 1. Here we have λ instead of \tilde{w} , $K \cdot i$ instead of \tilde{x}_i and K instead of X .

\Rightarrow Newton algorithm can be applied to solve the minimization problem \Rightarrow

$$\lambda_{t+1} = \lambda_t - [\nabla^2 f(\lambda_t)]^{-1} \cdot \nabla f(\lambda_t)$$

$\nabla^2 f(\lambda_t)$ and $\nabla f(\lambda_t)$ have the same form as $\nabla^2 f(\tilde{w}_t)$ and $\nabla f(\tilde{w}_t)$ in part 1. \Rightarrow following the same derivations, we can show that:

$$\left\{ \begin{array}{l} \nabla f(\lambda_t) = \frac{1}{n_+} K_{pos}^T (\underbrace{p_{pos} - 1}_{\lambda_t}) + \frac{1}{n_-} X_{neg}^T (\underbrace{p_{neg} - 1}_{\lambda_t}) + 2\lambda K \lambda_t \\ \nabla^2 f(\lambda_t) = \frac{1}{n_+} K_{pos}^T (\underbrace{p_{pos}(1-p_{pos})}_{\lambda_t}) K_{pos} + \frac{1}{n_-} K_{neg}^T (\underbrace{p_{neg}(1-p_{neg})}_{\lambda_t}) K_{neg} + 2\lambda K \end{array} \right.$$

$K_{pos} = K[y=1]$ } matrices composed of rows of matrix
 $K_{neg} = K[y=-1]$ } K corresponding to $y=1$ and $y=-1$

Implementation attached

$$\text{New } p_i = \frac{1}{1 + \exp(-\underline{\alpha}_i^T \underline{k}_{::i})} \Rightarrow \underline{p} = \begin{bmatrix} p_1 \\ \vdots \\ p_n \end{bmatrix}$$

$$\underline{p}_{\text{pos}} = \underline{p}[y=1] ; \underline{p}_{\text{neg}} = \underline{p}[y=-1]$$

Time complexity of training is $O(n^2 d)$

Testing

for a new test point \underline{x}' :

$$p = \frac{1}{1 + \exp(-\underline{w}^T \phi(\underline{x}'))} = \boxed{\frac{1}{1 + \exp(-\underline{\alpha}^T \underline{k}(\underline{x}'))}}$$

$$\text{where } k(\underline{x}') = K(\underline{x}, \underline{x}')$$

↳ training samples

Time complexity for testing is $O(nd)$

calculating
kernel vector $k(\underline{x}')$