

Practical task 6.

Optimization: GD vs SGD vs mini-batch SGD Consider the logistic regression method for binary classification. In this part you need to implement gradient descent (GD), stochastic gradient descent (SGD) and mini-batch stochastic gradient descent (mini-batch SGD) methods for log-loss optimization. In these three methods a gradient of a loss function $L(w) = \frac{1}{N} \sum_{i=1}^N l_i(w)$ on each step is obtained differently:

- in GD the whole training dataset is used: $\nabla_w L(w) = \nabla_w \frac{1}{N} \sum_{i=1}^N l_i(w)$,
- in SGD only one random training object is used: $\nabla_w L(w) \approx \nabla_w l_i(w)$ for some random i ,
- in mini-batch SGD a set of M random training objects is used: $\nabla_w L(w) \approx \nabla_w \frac{1}{M} \sum_{i=1}^M l_{k_i}(w)$, where $M < N$ and $\{k_1, \dots, k_M\}$ are some random indexes.

1. Consider the log-loss with L1 and L2 regularizations: $L(w) = \sum_i \ln(1 + \exp(-w^T x_i y_i)) + \gamma \|w\|_1 + \beta \|w\|_2$. Find its derivative and the update rules for GD, SGD and mini-batch SGD.
2. Implement **Python** functions that optimize the loss function with GD, SGD and mini-batch SGD. Prototypes for these functions are the following:

```
w,L = GD(X, y, max_epoch = 1000, alpha = 0.1, gamma = 0, beta = 0, tol = 0.1)
w,L = SGD(X, y, batch_size = 1, max_epoch = 1000, alpha = 0.1, gamma = 0, beta = 0)
```

where **w** is a weight vector, **L** is a vector of loss values after each epoch (starting with its initial value before optimization), **X,y** is a training data, **max_epoch** is a maximum number of epochs, **alpha** is a step size, **gamma** and **beta** are coefficients for regularization and **tol** is a tolerance parameter for GD stopping criterion. Here an epoch is a pass over training data so GD makes one step in one epoch, SGD makes N steps and mini-batch SGD makes N/M steps.

Some hints for you:

- Zeros or small random numbers from $[-\frac{1}{2d}, \frac{1}{2d}]$ should be chosen for weight initialization.
 - The efficient step size for GD is approximately $0.01 - 1$.
 - Step size should be constant for GD and decreasing for SGD, for example, $\alpha/epoch_number$ where α is some constant.
 - In SGD and mini-batch SGD you should randomly pick objects from training data on each step (use `numpy.random.choice` or shuffle training data at the beginning of each epoch and then pick objects successively (use `sklearn.utils.shuffle`)).
 - Stopping criteria: for GD use $L_{old} - L_{new} < tol$, for SGD simply do a particular number of iterations.
 - For code efficiency use numpy vectors to compute gradients.
3. Generate simple 2-dimensional data and check that your functions work properly on it. Don't forget to standardize the data and then add a constant feature to it. Run GD, SGD and mini-batch SGD (with batch size 10 – 100) from the same initial point and compare the results. What can you say about the convergence rate of these two algorithms?

Here you should demonstrate the following plots:

- plot with data points and decision boundary for each method,
- plot of decreasing $L(w)$ for increasing epoch number (for all methods together).

Regularization

1. Load the first dataset. Use `pickle.load`. Fit a logistic regression classifier on the training samples. Use GD with different regularizations (without one, only L1, only L2, L1 and L2), start all the methods from the same initial point (for example, use zero vector as initial weights vector). Don't forget to standardize the data and then add a constant feature to it.
2. Compare the results of the methods on the train and test data, explain the difference.
3. Use the resulting weights vector of GD with L1 regularization to determine two the most important features. Fit the logistic classifier only on these two features (+ the constant one) and visualize the decision boundary. Does L1 regularization help you to chose important features?

Model evaluation In this part of the task you will work with the problem of diabetes diagnostics. Load the diabetes dataset using `pickle.load`. This dataset has the following features:

1. Number of pregnancies
2. Plasma glucose concentration after 2 hours in an oral glucose tolerance test
3. Diastolic blood pressure
4. Triceps skin fold thickness
5. 2-Hour serum insulin
6. Body mass index
7. Diabetes pedigree function
8. Age

Class label is equal to 1 if a person has a diabetes and to -1 otherwise.

1. Train the logistic regression classifier on this dataset. Use mini-batch SGD without regularization. Don't forget to standardize the data and then add a constant feature to it.
2. In diagnostic problems false positive and false negative errors actually have different costs. Let's say, if we make false negative error (don't detect a condition when it is present), then the patient doesn't have a necessary treatment and, if we make false positive error (detect a condition when it isn't present), then the patient simply need to be tested more. Therefore, the cost of false negative error is higher and we care much more about this type of error.

Compute a confusion matrix for fitted classifier. How many errors of each type have you got? Compute a false positive and a false negative rates for this classifier. Why are they so different?

Useful functions: `sklearn.metrics.confusion_matrix`.

3. To change the proportion of errors of different types you can change a threshold a at the prediction rule $y = \sigma(w^T x) > a$, where $a \in [0, 1]$.

Show the ROC-curve of the fitted classifier and a point on it which corresponds to $a = 0.5$ (the one you computed at the previous step). Using ROC-curve choose a so that false negative rate is less than 20% while false positive rate is still small. What accuracy and false positive rate does the final algorithm have?

Useful functions: `sklearn.metrics.roc_curve`.