

Foundations of Data Analysis – SS20

Lab Assignment: Supervised Learning

Due: 9:30 AM on the 29th of April

0 Instructions and the data set

Please follow these submission instructions carefully, so that I can focus on grading and providing helpful feedback. Not following these instructions will result in a -5 point penalty, and I may ask you to resubmit it in the correct format. Upload your submission to Moodle as exactly *two* files, your python code and comments in the file `<last_name>.py` and your write-up containing the requested figures and question answers in `<last_name>.pdf`, replacing `<last_name>` with your last name(s). I suggest you use \LaTeX for the pdf, but you can also hand-write yours answers and include them in the pdf as photos, or use any other markup language/text formatter you would like, as long as the submission is a pdf. Working in groups is allowed (and even encouraged), provided you do *all* of the following:

1. clearly include the names of the people you worked with (inside the write-up, not in the file name)
2. do all of the write-up yourself, in your own words
3. use the group for discussing *ideas* and not just sharing answers

For any questions, email me (Alex) at `alexander.markham@univie.ac.at`, or post in the dedicated lab discussion forum on Moodle.

The data set `lab_iris_data.csv` on Moodle is a modified version the classic Iris flower data set (make sure you use this version and not the version Moritz uploaded during the LDA tutorial). It has two classes that are linearly separable based on three features (whereas the original has three classes total, only two of which are not linearly separable from each other, along with four features).

Your code must be well-documented and readable; see the Style Guide for Python Code for the style conventions. The algorithm is not very complicated to implement and you should not need to use any sources other than the assignment sheet, your previous python tutorial and Pen and Paper assignment, the lectures, and the documentation for numpy and matplotlib. However, if you insist on using other sources, you must thoroughly document and cite them—excessive use will be penalized, and undocumented/uncited use will be considered plagiarism and taken very

seriously. If you're unsure, simply send me an email at `alexander.markham@univie.ac.at` or post your question on the Moodle forum.

The *only* python packages you are allowed to use are `numpy` and `matplotlib` (plus you will need from `mpl_toolkits.mplot3d` `import Axes3D` to use `matplotlib`'s 3-d scatterplot function). I have uploaded to a short video on Moodle to help you get started.

1 Visualize the data (10%)

The first step is to load the data into Python. I suggest using the `loadtxt()` along with the `converters` argument to change the class labels to 0 and 1.

Produce a 3-d scatterplot (see the documentation for help) of the data (each dimension corresponding to a feature), with the data points colored differently according to their class, and put this plot in your write-up.

2 Implement and train a regularized logistic regression model using stochastic gradient descent (30%)

Tasks

Recall that in logistic regression, our hypothesis is $h_{\theta}(\mathbf{x}^{(i)}) = \frac{1}{1 + \exp(-\mathbf{x}^{(i)}\boldsymbol{\theta})}$, where $\boldsymbol{\theta}$ is a column vector of parameters (this replaces \underline{w} and b from the notes¹), returns a probability of the sample $\mathbf{x}^{(i)}$ being from a certain class. So, you are given a matrix $\mathbf{X} \subset \mathbb{R}^{m,n}$ of samples (with m rows of samples and n columns of features) and a vector $\mathbf{y} \subset \{0,1\}^m$ of labels to train the model. Then, using the trained model, you can classify a sample $\mathbf{x}^{(i)}$ in class 1 if $P(y = 1 \mid \mathbf{x}^{(i)}) = h_{\theta}(\mathbf{x}^{(i)}) \geq \frac{1}{2}$ and class 0 otherwise. We define the loss function as

$$\ell_{\text{L.R.}}(\boldsymbol{\theta}, \mathbf{X}, \mathbf{y}) = \sum_{i=1}^m [-y^{(i)} \log(h_{\theta}(\mathbf{x}^{(i)})) - (1 - y^{(i)}) \log(1 - h_{\theta}(\mathbf{x}^{(i)}))] + \frac{\lambda}{2} \sum_{j=1}^n \theta_j^2$$

and we can thus compute the gradient (over the entire training set; but it could analogously be defined for any subset of samples) as follows

$$\frac{\partial}{\partial \boldsymbol{\theta}} \ell_{\text{L.R.}}(\boldsymbol{\theta}, \mathbf{X}, \mathbf{y}) = \sum_{i=1}^m [(h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)}) \mathbf{x}^{(i)\top}] + \lambda \boldsymbol{\theta}.$$

Notice that the above loss function and gradient are slightly different from the Pen and Paper assignment: the model here is *regularized*.

¹do not forget to add a column of 1s to \mathbf{X} for the bias term in $\boldsymbol{\theta}$ so that the decision boundary is not constrained to pass through the origin

Recall the stochastic gradient descent algorithm:

Algorithm 1: Stochastic gradient descent

Data: Training samples $\mathbf{X} = \begin{bmatrix} \mathbf{x}^{(1)} \\ \vdots \\ \mathbf{x}^{(m)} \end{bmatrix}$ and corresponding labels $\mathbf{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix}$

Parameters: T , number of iterations; η_{init} , initial learning rate; and k , batch size

```
1 initialize  $\hat{\theta} = \mathbf{0}$ ;  
2 for  $t$  in  $(1, \dots, T)$  do  
3    $\eta \leftarrow \frac{\eta_{\text{init}}}{\sqrt{t}}$ ;  
4   randomly draw  $k$  samples from  $(\mathbf{X}, \mathbf{y})$  to get  $\mathbf{X}_k$  and  $\mathbf{y}_k$ ;  
5    $\hat{\theta} \leftarrow \hat{\theta} - \eta \cdot \frac{1}{k} \cdot \frac{\partial}{\partial \theta} \ell(\hat{\theta}, \mathbf{X}_k, \mathbf{y}_k)$ ;  
6 end
```

Result: $\hat{\theta}$ that minimizes ℓ

Implement stochastic gradient descent for regularized logistic regression. Make sure you can easily adjust the regularization factor λ , the initial learning rate η , and the number of iterations T .

Train the classifier on the given data set, starting with values $\lambda = 1$, $T = 100$, $k = 20$, and $\eta_{\text{init}} = 0.1$. Answer the following questions in your write-up:

- do these parameter values result in a well-trained model?
- how do you know?
- experiment by trying different values and see if you find a set that performs better or faster.
- explicitly describe how varying each of the four parameters affects training and the final model.

Note: The Gradient Descent algorithm, loss functions, and gradients are usually with respect to the entire training set, but the (Batch) Stochastic Gradient Descent algorithm that you are implementing is a stochastic approximation of the Gradient Descent algorithm, so it approximates the true gradient by computing it and updating θ after each (set of) training sample(s) rather than summing over all samples before updating θ .

3 Plot the separating hyperplane (30%)

Produce another 3-d scatterplot like in Task 1 above, but this time also include the separating hyperplane you found in Task 2, and include this in your write-up.

4 Further questions (30%)

Include your (thoughtful) answers to the following questions in your write-up, with 2–3 sentences to explain your reasoning for each answer:

- Given the above Note in Task 2 about non-stochastic Gradient Descent, why do they think the different variations of the algorithm exist, i.e. when is one more useful than the other?
- Under what conditions do they produce the same trained model?
- How (if at all) would you change the pseudocode in Algorithm 1 to use a different loss function?
- How does the regularized model we used here differ from the result if we had not used a regularization term?