

COMP9417 - Machine Learning

Homework 2: Logistic Regression & Optimization

Introduction In this homework, we will first explore some aspects of Logistic Regression and performing inference for model parameters. We then turn our attention to gradient based optimisation, the workhorse of modern machine learning methods.

Points Allocation There are a total of 25 marks. The available marks are:

- Question 1 a): 2 marks
 - Question 1 b): 3 marks
 - Question 1 c): 2 marks
 - Question 1 d): 3 marks
 - Question 1 e): 1 mark
 - Question 2 a): 2 marks
 - Question 2 b): 3 marks
 - Question 2 c): 1 mark
 - Question 2 d): 0 marks
 - Question 2 e): 3 marks
 - Question 2 f): 3 marks
 - Question 2 g): 2 marks
- Handwritten blue annotations: A bracket groups the first five items with the number "11" next to it. Another bracket groups the last seven items with the number "14" next to it.

What to Submit

- A single PDF file which contains solutions to each question. For each question, provide your solution in the form of text and requested plots. For some questions you will be requested to provide screen shots of code used to generate your answer — only include these when they are explicitly asked for.
- **.py file(s) containing all code you used for the project, which should be provided in a separate .zip file.** This code must match the code provided in the report.
- You may be deducted points for not following these instructions.
- You may be deducted points for poorly presented/formatted work. Please be neat and make your solutions clear. Start each question on a new page if necessary.

- You **cannot** submit a Jupyter notebook; this will receive a mark of zero. This does not stop you from developing your code in a notebook and then copying it into a .py file though, or using a tool such as **nbconvert** or similar.
- We will set up a Moodle forum for questions on this homework. Please read the existing questions before posting new questions. Please do some basic research online before posting questions. Please only post clarification questions. Any questions deemed to be *fishing* for answers will be ignored and/or deleted.
- Please check the Moodle forum for updates to this spec. It is your responsibility to check for announcements about the spec.
- Please complete your homework on your own, do not discuss your solution with other people in the course. General discussion of the problems is fine, but you must write out your own solution and acknowledge if you discussed any of the problems in your submission (including their name and zID).
- As usual, we monitor all online forums such as Chegg, StackExchange, etc. Posting homework questions on these site is equivalent to plagiarism and will result in a case of academic misconduct.

When and Where to Submit

- **Due date: Week 7, Sunday July 18th, 2021 by 11:55pm. Please note that the forum will not be actively monitored on weekends.**
- Late submissions will incur a penalty of 20% per day (from the ceiling, i.e., total marks available for the homework) for the first 5 days. For example, if you submit 2 days late, the maximum possible mark is 60% of the available 25 marks.
- Submission must be done through Moodle, no exceptions.

Question 1. Regularized Logistic Regression & the Bootstrap

In this problem we will consider the dataset provided in `Q1.csv`, with binary response variable Y , and 45 continuous features X_1, \dots, X_{45} . Recall that Regularized Logistic Regression is a regression model used when the response variable is binary valued. Instead of using mean squared error loss as in standard regression problems, **we instead minimize the log-loss**, also referred to as the cross entropy loss. For a parameter vector $\beta = (\beta_1, \dots, \beta_p) \in \mathbb{R}^p$, $y_i \in \{0, 1\}$, $x_i \in \mathbb{R}^p$ for $i = 1, \dots, n$, the log-loss is

$$L(\beta, \beta_0) = \sum_{i=1}^n y_i \ln \left(\frac{1}{s(\beta_0 + \beta^T x_i)} \right) + (1 - y_i) \ln \left(\frac{1}{1 - s(\beta_0 + \beta^T x_i)} \right),$$

where $s(z) = (1 + e^{-z})^{-1}$ is the logistic sigmoid (see Homework 0 for a refresher.) In practice, we will usually add a penalty term, and consider the optimisation:

$$(\hat{\beta}_0, \hat{\beta}) = \arg \min_{\beta_0, \beta} \{CL(\beta_0, \beta) + \text{penalty}(\beta)\} \quad (1)$$

where the penalty is usually not applied to the bias term β_0 , and C is a hyper-parameter. For example, in the ℓ_1 regularisation case, we take $\text{penalty}(\beta) = \|\beta\|_1$ (a LASSO for logistic regression).

- 2 (a) Consider the `sklearn logistic regression implementation` (section 1.1.11), which claims to minimize the following objective:

$$\hat{w}, \hat{c} = \arg \min_{w, c} \left\{ \|w\|_1 + C \sum_{i=1}^n \log(1 + \exp(-\tilde{y}_i(w^T x_i + c))) \right\}. \quad (2)$$

It turns out that this objective is identical to our objective above, but only after re-coding the binary variables to be in $\{-1, 1\}$ instead of binary values $\{0, 1\}$. That is, $\tilde{y}_i \in \{-1, 1\}$, whereas $y_i \in \{0, 1\}$. Argue rigorously that the two objectives (1) and (2) are identical, in that they give us the same solutions ($\hat{\beta}_0 = \hat{c}$ and $\hat{\beta} = \hat{w}$). Further, describe the role of C in the objectives, how does it compare to the standard LASSO parameter λ ? *What to submit: some commentary/your working.*

- 3 (b) Take the **first 500** observations to be your training set, and the rest as the test set. In this part, we will perform cross validation over the choice of C from scratch (**Do not use existing cross validation implementations here, doing so will result in a mark of zero.**)

Create a grid of 100 C values ranging from **$C = 0.0001$ to $C = 0.6$** in equally sized increments, inclusive. For each value of C in your grid, perform 10-fold cross validation (i.e. split the data into 10 folds, fit logistic regression (using the `LogisticRegression` class in `sklearn`) with the choice of C on 9 of those folds, **and record the log-loss on the 10th, repeating the process 10 times.**) For this question, we will take the first fold to be the first 50 rows of the training data, the second fold to be the next 50 rows, etc. Be sure to use ℓ_1 regularisation, and the `liblinear` solver when fitting your models.

To display the results, we will produce a plot: the x-axis should reflect the choice of C values, and for each C , plot a **box-plot** over the 10 CV scores. Report the value of C that gives you the best CV performance. Re-fit the model with this chosen C , and **report both train and test accuracy using this model**. Note that we do **not** need to use the \tilde{y} coding here (the `sklearn` implementation is able to handle different coding schemes automatically) so no transformations are needed before applying logistic regression to the provided data. *What to submit: a single plot, train and test accuracy of your final model, a screen shot of your code for this section, a copy of your python code in `solutions.py`*

- 2 (c) In this part we will compare our results in the previous section to the `sklearn` implementation of gridsearch, namely, the `GridSearchCV` class. My initial code for this section looked like:

```

1 grid_lr = GridSearchCV(estimator=
2                       LogisticRegression(penalty='l1',
3                                           solver='liblinear'),
4                               cv=10,
5                               param_grid=param_grid)
6 grid_lr.fit(Xtrain, Ytrain)

```

However, this gave me a very different answer to the result in (b). Provide two reasons for why this is the case, and then, if it is possible, re-run the code with some changes to give consistent results to those in (b), and if not, explain why. It may help to read through the [documentation](#). What to submit: some commentary, a screen shot of your code for this section, a copy of your python code in `solutions.py`

We next explore the idea of inference. To motivate the difference between prediction and inference, see some of the answers to this [stats.stackexchange](#) post. Needless to say, inference is a much more difficult problem than prediction in general. In the next parts, we will study some ways of quantifying the uncertainty in our estimates of the logistic regression parameters. **Assume for the remainder of this question that $C = 1$, and work only with the training data set ($n = 500$ observations) constructed earlier.**

- 3 (d) In this part, we will consider the nonparametric bootstrap for building confidence intervals for each of the parameters β_1, \dots, β_p . (Do not use existing Bootstrap implementations here, doing so will result in a mark of zero.) To describe this method, let's first focus on the case of $\hat{\beta}_1$. The idea behind the nonparametric bootstrap is as follows:

1. Generate B bootstrap samples from the original dataset. Each bootstrap sample consists of n points sampled **with replacement** from the original dataset, where n is the size of the original dataset.
2. On each of the B bootstrap samples, compute an estimate of β_1 , giving us a total of B estimates which we denote $\tilde{\beta}_1^{(1)}, \dots, \tilde{\beta}_1^{(B)}$.
3. Define the bootstrap mean and standard error respectively:

$$\tilde{\beta}_1 = \frac{1}{B} \sum_{b=1}^B \tilde{\beta}_1^{(b)} \quad \widetilde{\text{s.e.}}(\tilde{\beta}_1) = \sqrt{\frac{1}{B-1} \sum_{b=1}^B (\tilde{\beta}_1^{(b)} - \tilde{\beta}_1)^2}.$$

4. A 90% bootstrap confidence interval for β_1 is then given by the interval:

$$((\tilde{\beta}_1)_L, (\tilde{\beta}_1)_U) = (\text{5th quantile of the bootstrap estimates}, \text{95th quantile of the bootstrap estimates})$$

The idea behind a 90% confidence interval is that it gives us a range of values for which we believe with 90% probability the true parameter lives in that interval. If the computed 95 % interval contains the value of zero, then this provides us evidence that $\beta_1 = 0$, which means that the first feature should not be included in our model.

Take $B = 10000$ and set a random seed of 12 (i.e. `np.random.seed(12)`). Generate a plot where the x -axis represents the different parameters β_1, \dots, β_p , and plot a vertical bar that runs from $(\tilde{\beta}_p)_L$ to $(\tilde{\beta}_p)_U$. For those intervals that contain 0, draw the bar in red, otherwise draw it in blue. Also indicate on each bar the bootstrap mean. Remember to use $C = 1.0$.

What to submit: a single plot, ~~some commentary~~, a screen shot of your code for this section, a copy of your python code in solutions.py

- 1 (e) Comment on your results in the previous section, what do the confidence intervals tell you about the underlying data generating distribution? How does this relate to the choice of C when running regularized logistic regression on this data? Is regularization necessary?

Question 2. Gradient Based Optimization

In this question we will explore some algorithms for **gradient based optimization**. These algorithms have been crucial to the development of machine learning in the last few decades. The most famous example is the backpropagation algorithm used in deep learning, which is in fact just an application of a simple algorithm known as (stochastic) gradient descent. The general framework for a gradient method for finding a minimizer of a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is defined by

$$x^{(k+1)} = x^{(k)} - \alpha_k \nabla f(x_k), \quad k = 0, 1, 2, \dots, \quad (3)$$

where $\alpha_k > 0$ is known as the step size, or learning rate. Consider the following simple example of minimizing $g(x) = 2\sqrt{x^3 + 1}$. We first note that $g'(x) = 3x^2(x^3 + 1)^{-1/2}$. We then need to choose a starting value of x , say $x^{(0)} = 1$. Let's also take the step size to be constant, $\alpha_k = \alpha = 0.1$. Then we have the following iterations:

$$\begin{aligned} x^{(1)} &= x^{(0)} - 0.1 \times 3(x^{(0)})^2((x^{(0)})^3 + 1)^{-1/2} = 0.7878679656440357 \\ x^{(2)} &= x^{(1)} - 0.1 \times 3(x^{(1)})^2((x^{(1)})^3 + 1)^{-1/2} = 0.6352617090300827 \\ x^{(3)} &= 0.5272505146487477 \\ &\vdots \end{aligned}$$

and this continues until we terminate the algorithm (as a quick exercise for your own benefit, **code this up and compare it to the true minimum of the function which is $x_* = -1$**). This idea works for functions that have vector valued inputs, which is often the case in machine learning. For example, when we minimize a loss function we do so with respect to a weight vector, β . **When we take the step-size to be constant at each iteration, this algorithm is called gradient descent.** For the entirety of this question, **do not use any existing implementations of gradient methods, doing so will result in an automatic mark of zero for the entire question.**

- 2 (a) Consider the following optimisation problem:

you need to find the derivative mathematically for this part, no scipy
'pseudo-code' here, so just plug in the alpha/gradients values for a general iteration

$$\min_{x \in \mathbb{R}^n} f(x),$$

where

$$f(x) = \frac{1}{2} \|Ax - b\|_2^2,$$

and where $A \in \mathbb{R}^{m \times n}$, $b \in \mathbb{R}^m$ are defined as

$$A = \begin{bmatrix} 1 & 0 & 1 & -1 \\ -1 & 1 & 0 & 2 \\ 0 & -1 & -2 & 1 \end{bmatrix}, \quad b = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}.$$

Run gradient descent on f using a step size of $\alpha = 0.1$ and starting point of $x^{(0)} = (1, 1, 1, 1)$. You will need to terminate the algorithm when the following condition is met: $\|\nabla f(x^{(k)})\|_2 < 0.001$. In

your answer, clearly write down the version of the gradient steps (3) for this problem. Also, print out the first 5 and last 5 values of $x^{(k)}$, clearly indicating the value of k , in the form:

$$\begin{aligned} k = 0, & \quad x^{(k)} = [1, 1, 1, 1] \\ k = 1, & \quad x^{(k)} = \dots \\ k = 2, & \quad x^{(k)} = \dots \\ & \vdots \end{aligned}$$

What to submit: *an equation outlining the explicit gradient update, a print out of the first 5 and last 5 rows of your iterations, a screen shot of any code used for this section and a copy of your python code in solutions.py.*

- 3 (b) Note that using a constant step-size is sub-optimal. Ideally we would ideally like to take large steps at the beginning (when we are far away from the optimum), then take smaller steps as we move closer towards the minimum. There are many proposals in the literature for how best to choose the step size, here we will explore just one of them called the method of **steepest descent**. This is almost identical to gradient descent, except at each iteration k , we choose

$$\alpha_k = \arg \min_{\alpha \geq 0} f(x^{(k)} - \alpha \nabla f(x^{(k)})).$$

In words, **the step size is chosen to minimize an objective at each iteration of the gradient method, the objective is different at each step since it depends on the current x -value.** In this part, we will run steepest descent to find the minimizer in (a). First, derive an explicit solution for α_k (mathematically, please show your working). Then run steepest descent with the same $x^{(0)}$ as in (a), and $\alpha_0 = 0.1$. Use the same termination condition. Provide the first and last 5 values of $x^{(k)}$, as well as a plot of α_k over all iterations. *What to submit: a derivation of α_k , a print out of the first 5 and last 5 rows of your iterations, a single plot, a screen shot of any code used for this section, a copy of your python code in solutions.py.*

- 1 (c) Comment on the differences you observed, why would we prefer steepest descent over gradient descent? Why would you prefer gradient descent over steepest descent? Finally, explain why this is a reasonable condition to terminate use to terminate the algorithm.

In the next few parts, we will use the gradient methods explored above to solve a real machine learning problem. Consider the data provided in Q2.csv. It contains 414 real estate records, each of which contains the following features:

- transactiondate: date of transaction
- age: age of property
- nearestMRT: distance of property to nearest supermarket
- nConvenience: number of convenience stores in nearby locations
- latitude
- longitude

The target variable is the property price. The goal is to learn to predict property prices as a function of a subset of the above features.

- (d) We need to preprocess the data. First remove any rows with missing values. Then, delete all features except for age, nearestMRT and nConvenience. Then use the `sklearn minmaxscaler` to normalize the features. Finally, create a training set from the first half of the resulting dataset, and a test set from the remaining half. Your end result should look like:

- first row X_train: [0.73059361, 0.00951267, 1.]
- last row X_train: [0.87899543, 0.09926012, 0.3]
- first row X_test: [0.26255708, 0.20677973, 0.1]
- last row X_test: [0.14840183, 0.0103754, 0.9]
- first row Y_train: 37.9
- last row Y_train: 34.2
- first row Y_test: 26.2
- last row Y_test: 63.9

What to submit: a copy of your python code in solutions.py

3 (e) Consider the loss function

$$\mathcal{L}(x, y) = \sqrt{\frac{1}{4}(x - y)^2 + 1} - 1,$$

and consider the linear model

$$\hat{y}_i = w_0 + w_1 x_{i1} + w_2 x_{i2} + w_3 x_{i3}, \quad i = 1, \dots, n.$$

We can write this more succinctly by letting $w = (w_0, w_1, w_2, w_3)^T$ and $x_i = (1, x_{i1}, x_{i2}, x_{i3})^T$, so that $\hat{y}_i = w^T x_i$. The mean loss achieved by our model (w) on a given dataset of n observations is then

$$\mathcal{L}(w) = \frac{1}{n} \sum_{i=1}^n \left[\sqrt{\frac{1}{4}(y_i - w^T x_i)^2 + 1} - 1 \right].$$

We will run gradient descent to compute the optimal weight vector w . The iterations will look like

$$w^{(k+1)} = w^{(k)} - \alpha_k \nabla \mathcal{L}(w).$$

Instead of computing the gradient directly though, we will rely on an automatic differentiation library called JAX. Read the first section of the [documentation](#) to get an idea of the syntax. **Implement gradient descent from scratch and use the JAX library to compute the gradient of the loss function at each step.** You will only need the following import statements:

```
1 # pip install --upgrade pip
2 # pip install jax
3 # you may/may not also need to run: pip install jaxlib
4 # pip install --upgrade jax[cpu]
5
6 import jax.numpy as jnp
7 from jax import grad
8
```

Use $w^{(0)} = [1, 1, 1, 1]^T$, and a step size of 1. Terminate your algorithm when the absolute value of the loss from one iteration to the other is less than 0.0001. Report the number of iterations taken, and the final weight vector. **Further, report the train and test losses achieved by your final model, and produce a plot of the training loss at each step of the algorithm.** What to submit: a single plot, the final weight vector, the train and test **loss** of your final model, a screen shot of your code for this section, a copy of your python code in solutions.py

- 3 (f) Finally, re-do the previous section but with steepest descent instead. In order to compute α_k at each step, you can either use JAX or it might be easier to use the `minimize` function in `scipy` (See lab3). Run the algorithm with the same $w^{(0)}$ as above, and take $\alpha_0 = 1$ as your initial guess when numerically solving for α_k (for each k). Terminate the algorithm when the loss value falls below 2.5. Report the number of iterations it took, as well as the final weight vector, and the train and test losses achieved. Generate a plot of the losses as before and include it. What to submit: a single plot, the final weight vector, the train and test accuracy of your final model, a screen shot of your code for this section, a copy of your python code in `solutions.py`
- 2 (g) In this question we have explored the gradient descent and steepest descent variants of the gradient method. Many other gradient based algorithms exist in the literature. Choose one and describe it. What to submit: some commentary, any plots or code you used in this section (you don't need to write any code, or supply plots, but you may. Also be sure to cite any sources used here.)