

COMP9417 - Machine Learning

Homework 2

Introduction In this homework we first take a closer look at feature maps induced by kernels. We then explore a creative use of the gradient descent method introduced in homework 1. We will show that gradient descent techniques can be used to construct combinations of models from a base set of models such that the combination can outperform any single base model.

Points Allocation There are a total of 28 marks.

- Question 1 a): 2 marks
- Question 1 b): 2 marks
- Question 1 c): 4 marks
- Question 1 d): 2 marks
- Question 1 e): 2 marks
- Question 1 f): 1 mark
- Question 1 g): 1 mark
- Question 2 a): 2 marks
- Question 2 b): 3 marks
- Question 2 c): 3 marks
- Question 3 a): 2 marks
- Question 3 b): 2 marks
- Question 3 c): 2 marks

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 about 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 Moodle announcements 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(s) 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.
- You may **not** use SymPy or any other symbolic programming toolkits to answer the derivation questions. This will result in an automatic grade of zero for the relevant question. You must do the derivations manually.

When and Where to Submit

- Due date: Week 7, Monday **July 10th, 2023 by 5pm**. Please note that the forum will not be actively monitored on weekends.
- Late submissions will incur a penalty of 5% per day **from the maximum achievable grade**. For example, if you achieve a grade of 80/100 but you submitted 3 days late, then your final grade will be $80 - 3 \times 5 = 65$. Submissions that are more than 5 days late will receive a mark of zero.
- Submission must be made on **Moodle**, no exceptions.

Question 1. Gradient Descent for Learning Combinations of Models

In this question, we discuss and implement a gradient descent based algorithm for learning combinations of models, which are generally termed 'ensemble models'. The gradient descent idea is a very powerful one that has been used in a large number of creative ways in machine learning beyond direct minimization of loss functions.

The Gradient-Combination (GC) algorithm can be described as follows: Let \mathcal{F} be a set of base learning algorithms¹. The idea is to combine the base learners in \mathcal{F} in an optimal way to end up with a good learning algorithm. Let $\ell(y, \hat{y})$ be a loss function, where y is the target, and \hat{y} is the predicted value.² Suppose we have data (x_i, y_i) for $i = 1, \dots, n$, which we collect into a single data set D_0 . We then set the number of desired base learners to T and proceed as follows:

(I) Initialize $f_0(x) = 0$ (i.e. f_0 is the zero function.)

(II) For $t = 1, 2, \dots, T$:

(GC1) Compute:

$$r_{t,i} = -\frac{\partial}{\partial f(x_i)} \sum_{j=1}^n \ell(y_j, f(x_j)) \Big|_{f(x_j)=f_{t-1}(x_j), j=1,\dots,n}$$

for $i = 1, \dots, n$. We refer to $r_{t,i}$ as the i -th pseudo-residual at iteration t .

(GC2) Construct a new *pseudo* data set, D_t , consisting of pairs: $(x_i, r_{t,i})$ for $i = 1, \dots, n$.

(GC3) Fit a model to D_t using our base class \mathcal{F} . That is, we solve

$$h_t = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n \ell(r_{t,i}, f(x_i))$$

(GC4) Choose a step-size. This can be done by **either** of the following methods:

(SS1) Pick a fixed step-size $\alpha_t = \alpha$

(SS2) Pick a step-size adaptively according to

$$\alpha_t = \arg \min_{\alpha} \sum_{i=1}^n \ell(y_i, f_{t-1}(x_i) + \alpha h_t(x_i)).$$

(GC5) Take the step

$$f_t(x) = f_{t-1}(x) + \alpha_t h_t(x).$$

(III) return f_T .

We can view this algorithm as performing (functional) gradient descent on the base class \mathcal{F} . Note that in (GC1), the notation means that after taking the derivative with respect to $f(x_i)$, set all occurrences of $f(x_j)$ in the resulting expression with the prediction of the current model $f_{t-1}(x_j)$, for all j . For example:

$$\frac{\partial}{\partial x} \log(x+1) \Big|_{x=23} = \frac{1}{x+1} \Big|_{x=23} = \frac{1}{24}.$$

¹For example, you could take \mathcal{F} to be the set of all regression models with a single feature, or alternatively the set of all regression models with 4 features, or the set of neural networks with 2 layers etc.

²Note that this set-up is general enough to include both regression and classification algorithms.

- (a) Consider the regression setting where we allow the y -values in our data set to be real numbers. Suppose that we use squared error loss $\ell(y, \hat{y}) = \frac{1}{2}(y - \hat{y})^2$. For round t of the algorithm, show that $r_{t,i} = y_i - f_{t-1}(x_i)$. Then, write down an expression for the optimization problem in step (GC3) that is specific to this setting (you don't need to actually solve it).

What to submit: your working out, either typed or handwritten.

- (b) Using the same setting as in the previous part, derive the step-size expression according to the adaptive approach (SS2).

What to submit: your working out, either typed or handwritten.

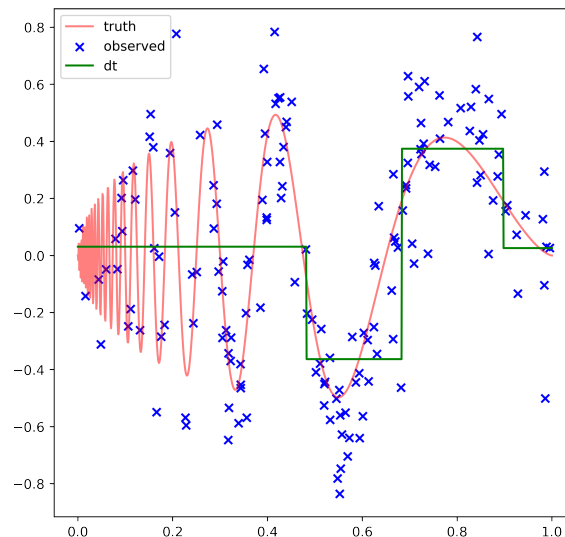
- (c) We will now implement the gradient-combination algorithm on a toy dataset from scratch, and we will use the class of decision stumps (depth 1 decision trees) as our base class (\mathcal{F}), and squared error loss as in the previous parts.³ The following code generates the data and demonstrates plotting the predictions of a fitted decision tree (more details in `q1.py`):

```

1  np.random.seed(123)
2  X, y = f_sampler(f, 160, sigma=0.2)
3  X = X.reshape(-1,1)
4
5  fig = plt.figure(figsize=(7,7))
6  dt = DecisionTreeRegressor(max_depth=2).fit(X,y) # example model
7  xx = np.linspace(0,1,1000)
8  plt.plot(xx, f(xx), alpha=0.5, color='red', label='truth')
9  plt.scatter(X,y, marker='x', color='blue', label='observed')
10 plt.plot(xx, dt.predict(xx.reshape(-1,1)), color='green', label='dt') # plotting
    example model
11 plt.legend()
12 plt.show()
13
```

The figure generated is

³In your implementation, you may make use of `sklearn.tree.DecisionTreeRegressor`, but all other code must be your own. You may use NumPy and matplotlib, but do not use an existing implementation of the algorithm if you happen to find one.



Your task is to generate a 5×2 figure of subplots showing the predictions of your fitted gradient-combination model. There are 10 subplots in total, the first should show the model with 5 base learners, the second subplot should show it with 10 base learners, etc. The last subplot should be the gradient-combination model with 50 base learners. Each subplot should include the scatter of data, as well as a plot of the true model (basically, the same as the plot provided above but with your fitted model in place of dt). Comment on your results, what happens as the number of base learners is increased? You should do this two times (two 5×2 plots), once with the adaptive step size, and the other with the step-size taken to be $\alpha = 0.1$ fixed throughout. There is no need to split into train and test data here. Comment on the differences between your fixed and adaptive step-size implementations. How does your model perform on the different x-ranges of the data?

What to submit: two 5×2 plots, one for adaptive and one for fixed step size, some commentary, and a screen shot of your code and a copy of your code in your .py file.

- (d) Repeat the analysis in the previous question but with depth 2 decision trees as base learners instead. Provide the same plots. What do you notice for the adaptive case? What about the non-adaptive case? *What to submit: two 5×2 plots, one for adaptive and one for fixed step size, some commentary, and a copy of your code in your .py file.*
- (e) Now, consider the classification setting where y is taken to be an element of $\{-1, 1\}$. We consider the following classification loss: $\ell(y, \hat{y}) = \log(1 + e^{-y\hat{y}})$. For round t of the algorithm, what is the expression for $r_{t,i}$? Write down an expression for the optimization problem in step (GC3) that is specific to this setting (you don't need to actually solve it).

What to submit: your working out, either typed or handwritten.

- (f) Using the same setting as in the previous part, write down an expression for α_t using the adaptive approach in (SS2). Can you solve for α_t in closed form? Explain.

What to submit: your working out, either typed or handwritten, and some commentary.

- (g) In practice, if you cannot solve for α_t exactly, explain how you might implement the algorithm. Assume that using a constant step-size is not a valid alternative. Be as specific as possible in your answer. What, if any, are the additional computational costs of your approach relative to using a constant step size?

What to submit: some commentary.

Question 2. Test Set Linear Regression

Recall that the rMSE (root-MSE) of a hypothesis⁴ h on a **test** set $\{(x_i, y_i)\}_{i=1}^n$, where x_i is a p -dimensional vector and y_i is a real number, is defined as

$$\text{rMSE}(h) = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - h(x_i))^2}.$$

For all parts, assume that the x_i 's are known to you, but the y_i 's are not. Suppose however that you are permitted to query $\text{rMSE}(h)$. What this means is that you can query, for any hypothesis h , the rMSE of h on the test set. Suppose further that you know that $\text{rMSE}(z) = c_0$, where z is the hypothesis that returns zero for any input, i.e. $z(x_i) = 0$ for each $i = 1, \dots, n$, and c_0 is some arbitrary positive number.

- (a) Assume you have a set of T hypotheses $\mathcal{H} = \{h_1, h_2, \dots, h_T\}$. You are told that for each i , there is a hypothesis h in \mathcal{H} , such that $h(x_i) = y_i$. In other words, for any point in the test set, there is at least one hypothesis in \mathcal{H} that predicts that point correctly⁵. Suppose that you are permitted to blend predictions of different hypotheses in \mathcal{H} ⁶. Design a naive, brute-force algorithm that constructs a hypothesis g from the elements of \mathcal{H} such that $\text{rMSE}(g) = 0$. How many queries of rMSE do you need to make? How does your algorithm scale (in the worst case) with the test size n ? Describe your algorithm in detail.

What to submit: a description of your algorithm and the number of queries required, either typed or hand-written.

- (b) We now consider a better approach than brute-force. For a given hypothesis h , define

$$h(X) = (h(x_1), h(x_2), \dots, h(x_n))^T$$

$$y = (y_1, y_2, \dots, y_n)^T.$$

We can always compute $h(X)$, but we do not know y . What is the smallest number of queries required to compute $y^T h(X)$? Describe your approach in detail.

What to submit: a description of your approach and the number of queries required, either typed or hand-written.

- (c) Given a set of K hypotheses $\mathcal{H} = \{h_1, \dots, h_K\}$, use your result in the previous part to solve

$$\min_{\alpha_1, \dots, \alpha_K} \text{rMSE} \left(\sum_{k=1}^K \alpha_k h_k \right),$$

and obtain the optimal weights $\alpha_1, \dots, \alpha_K$. Describe your approach in detail, and be sure to detail how many queries are needed and the exact values of the α 's, in terms of X, y and the elements of \mathcal{H} .

⁴The term hypothesis just means a function or model that takes as input x and returns as output a prediction $h(x) = \hat{y}$.

⁵Note that this does not imply that there is some hypothesis in \mathcal{H} that predicts all points correctly.

⁶For example, you could construct a blended hypothesis g that returns the predictions of h_2 on test points 1 to 5, and the predictions of h_5 on points 6 to n .

What to submit: a description of your algorithm and the number of queries required, either typed or hand-written.

Question 3. Newton's Method

Note: throughout this question do not use any existing implementations of any of the algorithms discussed unless explicitly asked to in the question. Using existing implementations can result in a grade of zero for the entire question. In homework 1 we studied gradient descent (GD), which is usually referred to as a first order method. Here, we study an alternative algorithm known as Newton's algorithm, which is generally referred to as a second order method. Roughly speaking, a second order method makes use of both first and second derivatives. Generally, second order methods are much more accurate than first order ones. Given a twice differentiable function $g : \mathbb{R} \rightarrow \mathbb{R}$, Newton's method generates a sequence $\{x^{(k)}\}$ iteratively according to the following update rule:

$$x^{(k+1)} = x^{(k)} - \frac{g'(x^{(k)})}{g''(x^{(k)})}, \quad k = 0, 1, 2, \dots, \quad (1)$$

For example, consider the function $g(x) = \frac{1}{2}x^2 - \sin(x)$ with initial guess $x^{(0)} = 0$. Then

$$g'(x) = x - \cos(x), \quad \text{and} \quad g''(x) = 1 + \sin(x),$$

and so we have the following iterations:

$$\begin{aligned} x^{(1)} &= x^{(0)} - \frac{x^{(0)} - \cos(x^{(0)})}{1 + \sin(x^{(0)})} = 0 - \frac{0 - \cos(0)}{1 + \sin(0)} = 1 \\ x^{(2)} &= x^{(1)} - \frac{x^{(1)} - \cos(x^{(1)})}{1 + \sin(x^{(1)})} = 1 - \frac{1 - \cos(1)}{1 + \sin(1)} = 0.750363867840244 \\ x^{(3)} &= 0.739112890911362 \\ &\vdots \end{aligned}$$

and this continues until we terminate the algorithm (as a quick exercise for your own benefit, code this up, plot the function and each of the iterates). We note here that in practice, we often use a different update called the *dampened* Newton method, defined by:

$$x^{(k+1)} = x^{(k)} - \alpha \frac{g'(x_k)}{g''(x_k)}, \quad k = 0, 1, 2, \dots \quad (2)$$

Here, as in the case of GD, the step size α has the effect of 'dampening' the update.

(a) Consider the twice differentiable function $f : \mathbb{R}^n \rightarrow \mathbb{R}$. The Newton steps in this case are now:

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

where $H(x) = \nabla^2 f(x)$ is the Hessian of f . Explain heuristically (in a couple of sentences) how the above formula is a generalization of equation (1) to functions with vector inputs. *what to submit: Some commentary*

(b) Consider the function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ defined by

$$f(x, y) = 100(y - x^2)^2 + (1 - x)^2.$$

Create a 3D plot of the function using `mplot3d` (see lab0 for example). Further, compute the gradient and Hessian of f . *what to submit: A single plot, the code used to generate the plot, the gradient and Hessian calculated along with all working. Add a copy of the code to solutions.py*

- (c) Using NumPy only, implement the (undampened) Newton algorithm to find the minimizer of the function in the previous part, using an initial guess of $x^{(0)} = (-1.2, 1)^T$. Terminate the algorithm when $\|\nabla f(x^{(k)})\|_2 \leq 10^{-6}$. Report the values of $x^{(k)}$ for $k = 0, 1, \dots, K$ where K is your final iteration. *what to submit: your iterations, and a screen shot of your code. Add a copy of the code to solutions.py*