

Homework 4

36-708, Spring 2023

Due May 5 at 5PM EST

Please attach all code to your homework. In an RMarkdown document for example, this can be done in one line (see [Yihui Xie's website](#) for how to do this).

1 Shapley values

Consider a simple cost allocation problem. Suppose there are n people that wish to share transportation to get from point A to their respective destinations, which are all in succession on the same street. Suppose that the cost of going from point A to the i^{th} person's destination costs c_i , and without loss of generality suppose $c_1 < c_2 < \dots < c_n$. The 'cost', $\nu : 2^{[n]} \rightarrow \mathbb{R}$ of a trip is defined in the following natural way,

$$\begin{aligned}\nu(\emptyset) &= 0 \\ \nu(\{i\}) &= c_i \\ \nu(S) &= c_j \text{ where } j = \max(S)\end{aligned}$$

(a) Define the 'sub-cost' functions:

$$\nu_j(S) = \begin{cases} c_j - c_{j-1} & \text{if } j \leq \max(S) \\ 0 & \text{otherwise,} \end{cases}$$

where $c_0 \equiv 0$ by convention. That is, ν_j is the additional cost for going to the j^{th} stop. It is 0 if nobody needs to go beyond stop j , and otherwise is $c_j - c_{j-1}$. Show that

$$\phi_i(\nu) = \sum_{j=1}^n \phi_i(\nu_j)$$

(b) Show that for any $i < j$, person i is a 'null player' with respect to cost ν_j . Further, show that persons $i_1, i_2 \geq j$ are equivalent with respect to cost ν_j .

(c) Use (b) to show that

$$\phi_i(\nu_j) = \begin{cases} \frac{c_j - c_{j-1}}{n - j + 1} & \text{if } i \geq j \\ 0 & \text{otherwise.} \end{cases}$$

Finally, conclude that

$$\phi_i(\nu) = \sum_{j=1}^i \frac{c_j - c_{j-1}}{n - j + 1}.$$

(d) Interpret the Shapley values. Does such a cost division make sense? If you didn't know about the Shapley values, what would your fair division look like? What are some advantages and disadvantages of your division strategy compared to Shapley's? (If you *would* do the same thing as Shapley, comment on its advantages and disadvantages.)

2 Image segmentation and compression using k -means

In this part, we explore how we can use k -means clustering for image segmentation and compression. The goal of image segmentation is to partition an image into multiple segments, where each segment typically represents an object in the image. Consider the following algorithm for image segmentation.

- Treat each pixel in the image as a point in 3-dimensional space comprising the intensities of the red, blue, and green channels. Treat each pixel in the image as a separate data point.
- Apply k -means clustering and identify the clusters.
- All the pixels belonging to a cluster are treated as a segment in the image.

For any value of k we can re-construct the image by replacing each pixel vector with the red, blue, and green intensity triplet given by the center to which that pixel has been assigned.

- (a) Implement a function that performs k -means++.
- (b) Take any image (you can be creative here!) and show the image re-constructions obtained using various values of k .

3 PCA for face recognition

In this part, we explore how we can use PCA to recognize faces. We play with AT&T database which has 10 near frontal images of 40 individuals under different illuminations per individual. The data can be downloaded from https://36708.github.io/hw4_faces/. We represent each image as a 1024 dimensional vector (the images have 32×32 pixels).

- (a) Implement a function to perform PCA.
- (b) With the mean ($\hat{\mu}_n$) and matrix of eigenvectors (V) learned from the training data, you can project other data points into this eigen space. Let x_{test} be a test point. Then, to project it into the eigenspace, simply subtract the mean vector and multiply by the eigenvector matrix (i.e, compute $V(x - \hat{\mu}_n)$). This will give you a vector of length k . Given any point $z \in \mathbb{R}^k$ in the eigenspace we can reconstruct the point in the original space as $V^T z + \hat{\mu}_n$.

Take a test image and project it into the eigenspace (use $k = 5, 10, 20, 40$). Use the projections to reconstruct the test images. Display the reconstructed images along with the original test image.

4 Neural networks

4.1 Warmup

- (a) Consider a multi-class classification problem. We observe y_i taking on K discrete values alongside covariates $x_i \in \mathbb{R}^p$. One method of modelling this problem is multinomial logistic regression. Here, we learn parameters $\beta_1, \dots, \beta_{K-1}$ following the below structural equation:

$$\mathbb{P}(y_i = k|x_i) = \frac{\exp(\beta_k^T x_i)}{1 + \sum_{l=1}^{K-1} \exp(\beta_l^T x_i)} \text{ for } k \leq K - 1,$$

with $\mathbb{P}(y_i = K|x_i) = 1 - \sum_{k=1}^{K-1} \mathbb{P}(y_i = k|x_i)$. Reformulate this problem in terms of neural networks. Specify the input layers, hidden layers, output layers, and loss function.

4.2 PCA and autoencoders

In this problem, we explore the relation between PCA, kernel PCA and auto-encoder neural networks (trained to output the same vector they receive as input). Consider the following setup.

- Number of data points: n
 - Number of features: d
- (a) Consider an auto-encoder with a single hidden layer of k nodes. Let w_{ij} denote the weight of the edge from the i th input node to the j th hidden node. Similarly, let v_{ij} denote the weight of the edge from the i th hidden node to the j th output node. Show how you can set the activation functions of hidden and output nodes as well as the weights w_{ij} and v_{ij} such that the resulting auto encoder resembles PCA.

4.3 Kernel PCA and autoencoders

Recall that kernel PCA is a non-linear dimensionality reduction technique where a principal vector v_j is computed as a linear combination of training examples in the feature space

$$v_j = \sum_{i=1}^n \alpha_{ij} \phi(x_i)$$

Computing the principal component of a new point x can then be done using kernel evaluations:

$$z_j(x) = \langle v_j, \phi(x) \rangle = \sum_{i=1}^n \alpha_{ij} \langle \phi(x_i), \phi(x) \rangle = \sum_{i=1}^n \alpha_{ij} k(x_i, x)$$

We show that kernel PCA can be represented by a neural network. First we define a kernel node. A kernel node with a vector w_i of incoming weights and an input vector x computes the output $y = k(x, w_i)$.

- (a) Show that, given a data set x_1, \dots, x_n , there exists a network with a single hidden layer and the output of the network is the kernel principal components $z_1(x), \dots, z_k(x)$ for a given input x . Specify the number of nodes in the input, output and hidden layers, the type and activation function of hidden and output nodes, and the weights of the edges in terms of α, x_1, \dots, x_n .
- (b) What is the number of parameters (weights) required to store the network in the previous problem?
- (c) Another way to do non-linear dimensionality reduction is to train an auto encoder with non-linear activation functions (e.g. sigmoid) in the hidden layers instead of using kernels. State one advantage and one disadvantage of that approach compared to kernel PCA.

4.4 Neural network for classification

In this problem, we play with the TensorFlow playground (<http://playground.tensorflow.org>), which is a nice visual tool for training simple Multi Layer Perceptrons (MLPs). Your goal in this problem is to carefully select input features to design the “smallest” MLP classifiers that can achieve low test loss for each of the 4 data sets in TensorFlow playground shows the 4 data sets available on TensorFlow playground). Here “smallest” is defined as having least number of neurons in the network. By low test loss we mean a test loss ≤ 0.1 for the swiss roll data set and a test loss of ≈ 0 for the rest of the data sets. Submit screenshots after your networks achieve the required test loss for each of the following data sets.

- (a) Circles
- (b) Clusters
- (c) Squares
- (d) Swiss Roll

4.5 Neural network for regression

Finally, we use a neural network on a regression problem. This is the most open-ended problem so far and you are essentially free to try anything. You are encouraged to play around with different packages and various settings.

- Dataset: We play with the Ames Housing Dataset. Recall that the raw data is available at <http://jse.amstat.org/v19n3/decock/AmesHousing.txt> and the dataset description can be found at <http://jse.amstat.org/v19n3/decock/DataDocumentation.txt>.
 - You will need to do some preprocessing (including missing values, possible outliers, etc.). Please feel free to make any reasonable preprocessing decisions, so long as you report them.
 - Train and test split: Again use a 3:1 random split.
 - Neural network architecture: You are free to play around with any architecture.
- (a) Report different choices (preprocessing, architecture, optimizers, etc.) you made throughout and the resulting train and test accuracies.
 - (b) Comment on relative advantages and disadvantages of training a neural network versus a stacked regressor you trained in a previous homework.