

## Mini Project 2

Due: May 21, 2023, 11:59PM PT

*Student Name:**Instructor Name: John Lipor*

## Introduction

In this project, you will continue working on the same target problem as mini project 1 (MP1), in which you are sought to predict heat flow residuals and bin them into one of four groups. In MP1, we treated the data as tabular, i.e., as sets of features that may have no natural structure across examples. However, in reality, each feature is an estimated surface across the study area, which is the greater Nevada region. In this project, we will treat one such feature as an image and aim to predict heat flow residuals using a CNN.

The detrended elevation map for the entire study region is shown in Fig. 1, with the wells used to measure heatflow indicated in red. Among geoscientists, there is a belief that changes in the elevation/topography may be indicative of geothermal activity. This is known as a search for “favorable structural settings.” The goal of this project is to see whether a CNN trained to predict heatflow residuals utilizes such favorable structural settings in order to make its predictions. Hence, while part of the aim of this project is to train an accurate predictor, an equally important goal is to examine the resulting saliency maps to see *how* that predictor makes its decisions.

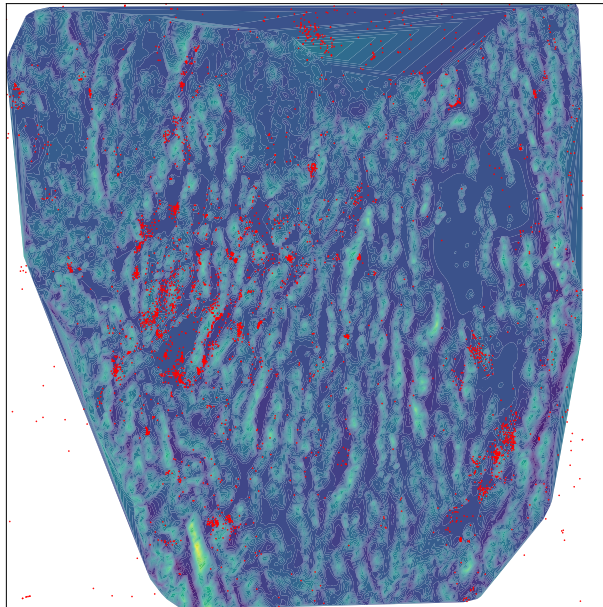


Figure 1: Detrended elevation map for entire study region. Red dots indicate wells used to measure heat flow.

To create the dataset, the detrended elevation map is processed into patches of size  $40 \text{ km} \times 40 \text{ km}$ . For each patch, the corresponding label is taken to be the maximum residual value among all wells in the patch.

Patches containing zero wells are discarded, and the training and test datasets are created in an 80/20 split. For a residual  $r_i$ , we will again utilize the categories

1. low:  $r_i \leq 25$
2. transition:  $r_i \in (25, 50]$
3. high:  $r_i \in (50, 200]$
4. very high:  $r_i > 200$

The resulting labels 1-4 are *ordinal labels*. As with MP1, your final predictor will be evaluated using the following loss

$$\ell(y, \hat{y}) = |y - \hat{y}|, \quad (1)$$

where  $y, \hat{y} \in \{1, \dots, 4\}$ .

Your task is to train an accurate CNN-based predictor of the heatflow residual given an elevation image patch. You are provided with the image and corresponding heat flow residual for each patch in the training set. You will also be provided with the images only of a test dataset. As part of your project, you will submit your predictions on the test set, and I will return your score evaluated using (1).<sup>1</sup>

## Notes

1. Your network must be coded and trained using PyTorch and D2L (not Keras or other similar packages).
2. You may make full use of StackExchange, Medium, or any other online resource you can find.
3. You may copy code from the internet, but you must **understand** what the code is doing and **cite** the source you copied it from.
4. You may **not** copy code from other students in this class. You **may** share MP1 code among yourselves.

## Requirements

You must create a report in either L<sup>A</sup>T<sub>E</sub>X or an approved type setting language (**not** Microsoft Word) that contains the following sections:

1. Problem description
  - Explain what data you have to work with, what algorithms you will use, what your goal is, and why anyone should care.
2. Exploratory data analysis (EDA)
  - **(NEW)** Since we are working with images, EDA will look different and perhaps be less involved than for MP1. You should still examine the ranges of image values (hint: they are not 0-255, which may be important for training... or it may not) and residuals, remove outliers, etc.
  - There are many excellent tutorials on EDA online. Cite whichever ones you found helpful.
  - Be sure to make any significant findings stand out and try to keep this concise.

---

<sup>1</sup>The top scoring project will receive 500 Lipor points.

### 3. Challenges

- What were the challenges you encountered when applying machine learning to this dataset?
- Did these challenges mainly result from the data? From results? Installing libraries to perform preprocessing?

### 4. Approach

- **(NEW)** Upon opening the dataset, you may notice that there are fewer than 300 images to use for training, which will be further reduced when you create training and validation sets. As a result, you **must** describe approaches you took to overcoming the limited dataset size. This need not be through image augmentation,<sup>2</sup> but it could be.
- Provide a thorough but concise description of your approach, including (but not limited to) your approach to wrangling, preprocessing, and improving the performance of your predictor.
- To avoid overfitting, you must use some form of validation (either a single validation set or cross validation) to tune your hyperparameters. Describe how you went about this and why.
- A summary paragraph at the beginning of this section could be helpful.

### 5. Evaluation and summary

- **(NEW)** As stated above, you must provide some insights into which parts of the images your CNN is utilizing for prediction. For example, you may use class activation maps or guided backpropagation to visualize this.
- **(NEW)** You may also wish to display the learned feature maps and filters, but only do so if they are informative and you are able to draw a conclusion from them.
- Make use of any meaningful metrics (e.g., the loss (1), regression MSE, classification error), but only include metrics that tell you something different.
- Consider diving into the data to see if you can determine any trends regarding which points were misclassified.
- Likewise consider if there were any features that were particularly important or unimportant.
- Summarize your solution, describing what worked, what didn't, what the main limitations are (e.g., if nonlinear classifiers would improve results), and any general conclusions about the dataset.

### 6. What I learned

- Describe the main skills/tools that you learned and used for this project and how you learned them.

---

<sup>2</sup>See `torchvision.transforms`.

## Grading

The grading breakdown is below. Each section will be graded according to technical correctness, effort, and creativity. Note that clarity of writing is a major component. You should put yourself in the place of writing for a boss or senior in a workplace. If your writing is terrible, they will assume your work is terrible.

<b>Item</b>	<b>Percentage</b>
Description	5%
EDA	20%
Challenges	10%
Approach	30%
Evaluation	20%
What I learned	5%
Clarity/conciseness of written communication	10%