

Applied Machine Learning HW1: Feature Map and k -NN (20%)

Due Monday April 18 @ 11:59pm on Canvas

(you should finish part 1 by the end of week 1, and part 2 by the end of week 2!)

Instructions:

1. This HW, like all other programming HWs, should be done in Python 3 and numpy only. See the course homepage for a numpy tutorial. If you don't have a Python+numpy installation yourself, you can use the College of Engineering servers by

```
ssh username@access.engr.oregonstate.edu
```

replacing `username` with your ENGR user name (should be identical to your ONID login). See the following for instructions on using SSH from Windows:

<https://it.engineering.oregonstate.edu/accessing-unix-server-using-putty-ssh>

If you don't have an ENGR account, see <https://it.engineering.oregonstate.edu/get-engr-account>.

You're **highly recommended** to set up SSH keys to bypass Duo authentication and password:

<https://it.engineering.oregonstate.edu/ssh-keygen>

The default `python3` on ENGR servers have `numpy`, `matplotlib`, and `pandas` installed.

2. Besides machine learning, this HW also teaches you data (pre-)processing skills and Unix (Linux or Mac OS X) command lines tools. These skills are even more important than machine learning itself for a software engineer or data scientist. As stated in the syllabus, Windows is **not** recommended (no data scientist uses Windows). The instructor and the TAs, like other computer scientists, have very limited Windows expertise, and will not be able to provide much technical assistance for Windows. Use `cmd`, `cygwin`, or `ssh` (see above) if you don't have Linux or Mac OS X on your own computer.

If you prefer to use Jupyter notebook but have difficulty setting it up on your own computer, the TAs have written a tutorial on running jupyter notebook remotely on the server:

https://classes.engr.oregonstate.edu/eecs/spring2022/cs513-400/extra/TAs_jupyter_tutorial.pdf

3. Ask questions on **Slack**. You're encouraged to answer other students' questions as well.
4. Do not use machine learning packages such as `sklearn`, though you can use them to verify your results.
5. Do not use data analysis packages such as `pandas` or `seaborn`. Your code should only depend on standard built-in packages plus `numpy`.
6. Download HW1 data from the course homepage. Unzip it by running `tar -xzf hw1-data.tgz` (here `-x` means "extract" and `z` means "unzip"; you might also be able to open it using WinZip or 7-zip on Windows). It contains:

<code>income.train.txt.5k</code>	training data (5,000 examples, with labels)
<code>income.dev.txt</code>	dev set (1,000 examples, with labels)
<code>income.test.blind</code>	test set (1,000 examples, without labels)
<code>toy.txt & binarize_example.py</code>	toy data and binarization code example
<code>validate.py & random_output.py</code>	tools to validate your test result

The semi-blind test set does not contain the target labels ($>50K$ or $\leq 50K$), which you will need to predict using your best model. Part of your grade is based on your prediction accuracy on test. Please make heavy use of the three python programs supplied, so that you don't need to start from scratch. For example, you should base your Part 1 on `binarize_example.py` and use `validate.py` to verify your test result in Part 4.

If you'd like to download the zip file to your remote ENGR server, use this command on the server:

```
wget https://classes.engr.oregonstate.edu/eecs/spring2022/cs513-400/hw1/hw1-data.tgz
```

7. You should submit a single `.zip` file containing `hw1-report.pdf`, `income.test.predicted`, and all your code. \LaTeX ing is recommended but not required. **Do not forget the debrief section (in each HW).**

1 Data (Pre-)Processing (Feature Map)

This is the most coding-heavy part in HW1, but you don't need to start from scratch. See `binarize_example.py`.

1. Take a look at the data. A training example looks like this:

37, Private, Bachelors, Separated, Other-service, White, Male, 70, England, <=50K

which includes the following 9 input fields plus one output field (y):

age, sector, education, marital-status, occupation, race, gender, hours-per-week, country-of-origin, target

Q: What are the positive % of training data? What about the dev set? Does it make sense given your knowledge of the average per capita income in the US?

2. Q: What are the youngest and oldest ages in the training set? What are the least and most amounts of hours per week do people in this set work? Hint:

```
$ cat income.train.txt.5k | sort -nk1 | head -1
```

Note the `$` is the prompt, not part of the command. `cat` lists all contents of a file (like Windows `type`), `sort -nk1` means sort by the first column numerically, and `head -1` lists the first line. Here `|` means “pipe”, i.e., feeding the output of the previous command as input to the next command.

3. There are two types of fields, *numerical* (*age* and *hours-per-week*), and *categorical* (everything else).¹ The default preprocessing method is to *binarize* all categorical fields, e.g., *race* becomes many *binary* features such as *race=White*, *race=Asian-Pac-Islander*, etc. These resulting features are all binary, meaning their values can only be 0 or 1, and for each example, in each field, there is one and only one positive feature (this is the so-called “one-hot” representation, widely used in ML and NLP).

Q: Why do we need to binarize all categorical fields?

4. Q: If we do not count *age* and *hours*, what's maximum possible Euclidean and Manhattan distances between two training examples? Explain.

5. Why we do **not** want to binarize the two numerical fields, *age* and *hours*? What if we did? How should we define the distances on these two dimensions so that each field has equal weight? (In other words, the distance induced by each field should be bounded by **2** (N.B.: not 1! why?)).

Hint: first, observe that the max distance between two people on a categorical field is **2**. If we simply “normalize” a numerical field by, say, *age* / 100, it might look OK but now what's the max distance between two people on *age*? Are you treating all fields equally?

6. Q: How many features do you have in total (i.e., the dimensionality)? Hint: should be around **90**. **Note:** the target label is not a feature! How many features do you allocate for each of the 9 fields? **Hint:**

```
$ for i in `seq 1 9`; do cat income.train.txt.5k | cut -f $i -d ',' | sort | uniq | wc -l; done
```

Here ``seq 1 9`` returns the sequence `1 2 ... 9`, `cut` is a command to extract specific columns of each line (e.g., `cut -f 2` extracts the second column), and `-d ','` means using comma as the column separator. `uniq` filters out duplicate consecutive rows, and `wc -l` counts the number of lines. So `sort | uniq | wc -l` returns the number of unique rows.

Hint: No need to start from scratch, as I've supplied `binarize_example.py` which is basically a standalone version of the “binarize from scratch” portion of the jupyter notebook I showed in the binarization video. You can run it in a directory that contains `toy.txt` (just `python3 binarize_example.py`) and expand this code to do binarization for all fields.

7. Q: How many features would you have in total if you binarize all fields?

¹In principle, we could also convert *education* to a numerical feature, but we choose **not** to do it to keep it simple.

2 Calculating Manhattan and Euclidean Distances

Hint: you can use the Matlab style “broadcasting” notations in numpy (such as matrix - vector) to calculate many distances in one shot. For example, if A is an $n \times m$ matrix (n rows, m columns, where n is the number of people and m is the number of features), and p is an m -dimensional vector (1 row, m columns) representing the query person, then $A - p$ returns the difference vectors from each person in A to the query person p , from which you can compute the distances:

```
>>> A = np.array([[1,2], [2,3], [4,5]])
>>> p = np.array([3,2])
>>> A - p
array([[ -2,  0],
       [ -1,  1],
       [ 1,  3]])
>>> np.linalg.norm(A-p, axis=1)
array([2.          , 1.41421356, 3.16227766])
```

This is Euclidean distance (what does `axis=1` mean?). You need to figure out Manhattan distance yourself.

To make sure your distance calculations are correct, we provide the following example calculations using the first person in the dev set:

```
$ head -1 income.dev.txt
45, Federal-gov, Bachelors, Married-civ-spouse, Adm-clerical, White, Male, 45, United-States, <=50K
```

The top-3 examples in the training set that are closest to the above person, according to the Manhattan distance, should be the following rows (note that the command `sed -n XXp` prints the XX^{th} line of a file):

```
$ sed -n 4873p income.train.txt.5k
33, Federal-gov, Bachelors, Married-civ-spouse, Adm-clerical, White, Male, 42, United-States, >50K
$ sed -n 4788p income.train.txt.5k
47, Federal-gov, Bachelors, Married-civ-spouse, Adm-clerical, White, Male, 45, Germany, >50K
$ sed -n 2592p income.train.txt.5k
48, Federal-gov, Bachelors, Married-civ-spouse, Prof-specialty, White, Male, 44, United-States, >50K
```

Notice that the first of these three persons matches all categorical fields with the dev person, only differing slightly in the two numerical fields, and the second and third persons match all but one categorical fields. The Manhattan distances of these three people to the dev person are:

```
(45-33) / 50. + (45-42) / 50. = 0.3
(47-45) / 50. + (45-45) / 50. + 1 + 1 = 2.04
(48-45) / 50. + (45-44) / 50. + 1 + 1 = 2.08
```

Coincidentally, these three people are also the top-3 closest according to the Euclidean distances, with the distances being

```
sqrt( ((45-33) / 50.) ** 2 + ((45-42) / 50.) ** 2 ) = 0.24738633753705963
sqrt( ((47-45) / 50.) ** 2 + ((45-45) / 50.) ** 2 + 1 ** 2 + 1 ** 2 ) = 1.4147791347061915
sqrt( ((48-45) / 50.) ** 2 + ((45-44) / 50.) ** 2 + 1 ** 2 + 1 ** 2 ) = 1.4156270695349111
```

Also notice that in both cases, the 3-NN predictions are wrong, as the top-3 closest examples are all >50K.

Finally, remember that you don't really need to sort the distances in order to get the top- k closest examples.

Note: If you couldn't get the same top-3 people as I listed above, it is possible that you included the target label in your feature map.

Questions:

1. Find the five (5) people closest to the last person (in Euclidean distance) in dev, and report their distances:

```
$ tail -1 income.dev.txt
58, Private, HS-grad, Widowed, Adm-clerical, White, Female, 40, United-States, <=50K
```

Note: this obviously means top-5 closest people in the training set. Remember: for any query person, you can only match him/her with people in the training set in order to predict his/her label.

2. Redo the above using Manhattan distance.
3. What are the 5-NN predictions for this person (Euclidean and Manhattan)? Are these predictions correct?

YOU SHOULD FINISH EVERYTHING UP TO HERE BY THE END OF WEEK 2.

3 k -Nearest Neighbor Classification

1. Implement the basic k -NN classifier (with the default Euclidean distance).

Q: Is there any work in training after finishing the feature map?

Q: What's the time complexity of k -NN to test one example (dimensionality d , size of training set $|D|$)?

Q: Do you really need to sort the distances first and then choose the top k ? Hint: there is a faster way to choose top k without sorting.

2. Q: Why the k in k -NN has to be an odd number?
3. Evaluate k -NN on the dev set and report the error rate and predicted positive rate for $k = 1, 3, 5, 7, 9, 99, 999, 9999$, e.g., something like:

```
k=1      dev_err xx.x% (+:xx.x%)
k=3      ...
...
k=9999   ...
```

Q: what's your best error rate on dev, and where did you get it? (Hint: 1-NN dev error should be $\sim 23\%$ and its positive % should be $\sim 27\%$).

4. Now report both training and ~~testing~~ **dev** errors (your code needs to run a lot faster! See Question 4.3 **5.3** for hints. See also week 2 videos for numpy and linear algebra tutorials, in case you're not familiar with the "Matlab"-style of thinking which is inherited by numpy):

```
k=1      train_err xx.x% (+:xx.x%)  dev_err xx.x% (+:xx.x%)
k=3      ...
...
k=9999   ...
```

Q: When $k = 1$, is training error 0%? Why or why not? Look at the training data to confirm your answer.

5. Q: What trends (train and dev error rates and positive ratios, and running speed) do you observe with increasing k ? Do they relate to underfitting and overfitting?

Q: What does $k = \infty$ actually do? Is it extreme overfitting or underfitting? What about $k = 1$?

6. Redo the evaluation using Manhattan distance. Better or worse? Any advantage of Manhattan distance?
7. Redo the evaluation using all-binarized features (with Euclidean). Better or worse? Does it make sense?

4 Deployment

Now try more k 's and take your best model and run it on the semi-blind test data, and produce `income.test.predicted`, which has the same format as the training and dev files.

Q: At which k and with which distance did you achieve the best dev results?

Q: What's your best dev error rates and the corresponding positive ratios?

Q: What's the positive ratio on test?

Part of your grade will depend on the accuracy of `income.test.predicted`.

IMPORTANT: You should use our `validate.py` to verify your `income.test.predicted`; it will catch many common problems such as formatting issues and overly positive or overly negative results:

```
cat income.test.predicted | python3 validate.py
```

We also provided a `random_output.py` which generates random predictions (~50% positive) and it will pass the formatting check, but fail on the positive ratio:

```
$ cat income.test.blind | python3 random_output.py | python3 validate.py
```

which might output:

```
Your file passed the formatting test! :)
```

```
Your positive rate is 49.6%.
```

```
ERROR: Your positive rate seems too high (should be similar to train and dev).
```

```
PLEASE DOUBLE CHECK YOUR BINARIZATION AND K-MEANS CODE.
```

If you test the dev set, it will certainly pass this test (try `cat income.dev.txt | python3 validate.py`).

Our automatic grading system will assume your `incoming.test.predicted` passes this test; if it doesn't, you will receive 0 points for the blind test part.

5 Observations

1. Q: Summarize the major drawbacks of k -NN that you observed by doing this HW. There are a lot!
2. Q: Do you observe in this HW that best-performing models tend to exaggerate the existing bias in the training data? Is it due to overfitting or underfitting? Is this a potentially social issue?
3. Q: What numpy tricks did you use to speed up your program so that it can be fast enough to print the training error? Hint: (a) broadcasting (such as matrix - vector); (b) `np.linalg.norm(..., axis=1)`; (c) `np.argsort()` or `np.argpartition()`; (d) slicing. The main idea is to do as much computation in the vector-matrix format as possible (i.e., the Matlab philosophy), and as little in Python as possible.
4. How many seconds does it take to print the training and dev errors for $k = 99$ on ENGR servers? Hint: use `time python ...` and report the user time instead of the real time. (Mine was about 14 seconds).
5. What is a Voronoi diagram (shown in k -NN slides)? How does it relate to k -NN?

Debriefing (required in your report)

1. Approximately how many hours did you spend on this assignment?
2. Would you rate it as easy, moderate, or difficult?
3. Did you work on it mostly alone, or mostly with other people?
4. How deeply do you feel you understand the material it covers (0%–100%)?
5. Any other comments?