

Fall 2024 CS4641/CS7641 A Homework 1 - Programming Section

Instructor: Dr. Mahdi Roozbahani

Deadline: Friday, September 20th, 11:59 pm EST

- No unapproved extension of the deadline is allowed. Submissions past our 48-hour penalized acceptance period will lead to 0 credit.
- Discussion is encouraged on Ed as part of the Q/A. However, all assignments should be done individually.
- Plagiarism is a serious offense. You are responsible for completing your own work. You are not allowed to copy and paste, or paraphrase, or submit materials created or published by others, as if you created the materials. All materials submitted must be your own.
- All incidents of suspected dishonesty, plagiarism, or violations of the Georgia Tech Honor Code will be subject to the Institute's Academic Integrity procedures. If we observe any (even small) similarities/plagiarisms detected by Gradescope or our TAs, WE WILL DIRECTLY REPORT ALL CASES TO OSI, which may, unfortunately, lead to a very harsh outcome. Consequences can be severe, e.g., academic probation or dismissal, grade penalties, a 0 grade for assignments concerned, and prohibition from withdrawing from the class.

Instructions for the assignment

- This assignment consists of warm-up programming questions designed to get you familiar with our programming homework structure.

Using the autograder

- Grads will typically find three assignments on Gradescope and Undergrads will typically find four assignments:
 - "Assignment X Non-programming": Where you will submit the written portion of the assignment.
 - "Assignment X Programming": Where you will submit any .py files and any program outputs as required by the problem.
 - "Assignment X Programming - Bonus for All": Where you will submit any .py files and any program outputs as required for Bonus for All.
 - "Assignment X Programming - Bonus for Undergrad": Where you will submit any .py files and any program outputs as required for Bonus for Undergrad. (Undergrad Only)
- You will submit your code for the autograder in the Assignment 1 Programming section.
- We provided you .py files and we added libraries in those files please DO NOT remove those lines and add your code after those lines. Note that these are the only allowed libraries that you can use for the homework.
- You are allowed to make as many submissions until the deadline as you like. Additionally, note that the autograder tests each function separately, therefore it can serve as a useful tool to help you debug your code if you are not sure of what part of your implementation might have an issue.

Deliverables and Points Distribution

Q7: Programming Warm-Up [5pts total]

Deliverables:

- `warmup.py`
- `env.pkl`

Parts:

- Setup [2pts] - *programming*
- Num py [3pts] - *programming*
 - Numpy Basics [2pts]
 - Broadcasting [1pts]

7.1 Setup [1pt]

- Deliverable: `env.pkl`

This notebook is tested under `python 3.11`, and the corresponding packages can be downloaded from [miniconda](#). You may also want to get yourself familiar with several packages:

- [jupyter lab](#): provides a web-based IDE with a built-in debugging functionality for jupyter notebooks.
- [numpy](#): a high performance math library backed by C
- [matplotlib](#): a python plotting library

Other packages you may find indispensable in machine learning (and potentially your project) are:

- [scikit-learn](#): provides many classical ML and data analysis algorithms
- [pandas](#): provides many useful tools for organizing and manipulating data
- [seaborn](#): make beautiful plots with less fidgeting in matplotlib
- [plotly](#): another great data visualization package

Please implement the functions that have "raise NotImplementedError", and after you finish the coding, please delete or comment "raise NotImplementedError".

```
#####  
### DO NOT CHANGE THIS CELL ###  
#####  
  
import sys  
  
sys.path.append("./utilities/")  
sys.path.append("warmup.py")  
  
import numpy as np  
  
print("Version Information")  
  
print("python: {}".format(sys.version))  
print("numpy: {}".format(np.__version__))  
  
%load_ext autoreload  
%autoreload 2
```

Version information
python: 3.11.9 (main, Apr 19 2024, 16:48:06) [GCC 11.2.0]
numpy: 1.26.2

7.1.1 Basics, Imports, and Directories

For the following part, you will need to ensure your notebook runtime is started in the correct directory. You can verify this with the following cell.

```
#####  
### DO NOT CHANGE THIS CELL ###  
#####
```

```
# RUN ME #  
import os  
os.getcwd()  
  
'/app/src/teacher_files'
```

In the cell below, import the `PackageUtils` class from `utils.py`.

```
# YOUR CODE HERE #
```

In the cell below, call the `get_packages` method of the `PackageUtils` class to see what packages are installed in this notebook's runtime environment.

```
# YOUR CODE HERE #
```

```
anyio==4.1.0  
argon2-cffi==23.1.0  
argon2-cffi-bindings==21.2.0  
arrow==1.3.0  
astor==0.8.1  
asttokens==2.4.1  
async-lru==2.0.4  
attrs==23.1.0  
autoflake==2.2.1  
autopep8==2.0.4  
Babel==2.14.0  
beautifulsoup4==4.12.2  
black==23.12.0  
bleach==6.1.0  
Bottleneck @ file:///croot/bottleneck_1707864210935/work  
Brotli @ file:///croot/brotli-split_1714483155106/work  
build==1.2.1  
CacheControl==0.14.0  
certifi==2023.11.17  
cffi==1.16.0  
cfgv==3.4.0  
chardet==5.2.0  
charset-normalizer==3.3.2  
cleo==2.1.0  
click==8.1.7  
comm==0.2.0  
contourpy @ file:///croot/contourpy_1700583582875/work  
crashtest==0.4.1  
cryptography==43.0.0
```

cycler @ file:///tmp/build/80754af9/cycler_1637851556182/work
debugpy==1.8.0
decorator==5.1.1
defusedxml==0.7.1
distlib==0.3.8
dulwich==0.21.7
executing==2.0.1
fastjsonschema==2.19.0
filelock==3.13.1
fonttools @ file:///croot/fonttools_1713551344105/work
fqdn==1.5.1
identify==2.5.33
idna==3.6
imagecodecs @ file:///croot/imagecodecs_1695064943445/work
imageio @ file:///croot/imageio_1707247282708/work
importlib-metadata==8.4.0
installer==0.7.0
ipykernel==6.27.1
ipython==8.18.1
ipywidgets==8.1.1
isoduration==20.11.0
isort==5.13.2
jaraco.classes==3.4.0
jedi==0.19.1
jeepney==0.8.0
Jinja2==3.1.2
joblib @ file:///croot/joblib_1718217211762/work
json5==0.9.14
jsonpiter==2.4
jsonschema==4.20.0
jsonschema-specifications==2023.11.2
jupyter==1.0.0
jupyter-console==6.6.3
jupyter-events==0.9.0
jupyter-isp==2.2.1
jupyter_client==8.6.0
jupyter_core==5.5.0
jupyter_server==2.12.1
jupyter_server_terminals==0.5.0
jupyterlab==4.0.9
jupyterlab_widgets==3.0.9
jupyterlab_pygments==0.3.0
jupyterlab_server==2.25.2
keyring==24.3.1
kivsolver @ file:///work/ci_py311/kivsolver_1676827230232/work
lazy_loader @ file:///croot/lazy_loader_1718176737906/work
markdown-it-py==3.0.0
MarkupSafe==2.1.3
matplotlib @ file:///croot/matplotlib-suite_1713336378214/work
matplotlib-inline==0.1.6
mdurl==0.1.2
mistune==3.0.2
mkl-fft @ file:///croot/mkl_fft_1695058164594/work
mkl-random @ file:///croot/mkl_random_1695059800811/work
mkl-service==2.4.0
more-itertools==10.4.0
msgpack==1.0.8
mypy-extensions==1.0.0
nbclient==0.9.0
nbconvert==7.12.0

nbformat==5.9.2
nbqa==1.7.1
nest-asyncio==1.5.8
networkx @ file:///croot/networkx_1720002482208/work
nodeenv==1.8.0
notebook==7.0.6
notebook_shim==0.2.3
numexpr @ file:///croot/numexpr_1696515281613/work
numpy==1.26.2
overrides==7.4.0
packaging==23.2
pandas @ file:///croot/pandas_1718308974269/work/dist/pandas-2.2.2-cp311-cp311-linux_x86_64.whl#sha256=3c7ce50f9f519c785bd4cdd28a0ca71f85a541f3d27b25aa9da770f953e7f2e9
pandocfilters==1.5.0
parso==0.8.3
pastel==0.2.1
pathspec==0.12.1
patsy @ file:///croot/patsy_1718378176128/work
pdf-watermark==2.0.0
pdfkit==1.0.0
pexpect==4.9.0
Pillow==10.1.0
pip @ file:///croot/pip_1723484598856/work
pkginfo==1.11.1
platformdirs==4.1.0
ply==3.11
poethepoet==0.24.4
poetry==1.8.3
poetry-core==1.9.0
poetry-plugin-export==1.8.0
pre-commit==3.6.0
prometheus-client==0.19.0
prompt-toolkit==3.0.43
psutil==5.9.6
ptyprocess==0.7.0
pure-eval==0.2.2
pyclean==2.7.6
pycodestyle==2.11.1
pycparser==2.21
pyflakes==3.1.0
Pygments==2.17.2
pyparsing @ file:///work/ci_py311/pyparsing_1677811559502/work
pypdf==3.17.2
pyproject_hooks==1.1.0
PyQt5==5.15.10
PyQt5-sip @ file:///croot/pyqt-split_1698769088074/work/pyqt_sip
PySocks @ file:///work/ci_py311/pysocks_1676822712504/work
python-dateutil==2.8.2
python-json-logger==2.0.7
pytz @ file:///croot/pytz_1713974312559/work
pyupgrade==3.15.0
PyYAML==6.0.1
pyzmq==25.1.2
qtconsole==5.5.1
QtPy==2.4.1
rapidfuzz==3.9.6
referencing==0.32.0
reportlab==4.0.8
requests==2.31.0
requests-toolbelt==1.0.0

```
rfc3339-validator==0.1.4
rfc3986-validator==0.1.1
rich==13.7.1
rpds-py==0.13.2
scikit-image @ file:///croot/scikit-image_1718285223463/work
scikit-learn @ file:///croot/scikit-learn_1721921875708/work
scipy @ file:///croot/scipy_1717521478074/work/dist/scipy-1.13.1-cp311-cp311-linux_x86_64.whl#sha256:f0a29afe8e78f15653c9afe02349a3462e4e9e7131c5e68b77b70fd68d25e6a2
seaborn @ file:///croot/seaborn_1718302919398/work
SecretStorage==3.3.3
Send2Trash==1.8.2
setuptools==69.0.2
shellingham==1.5.4
sip @ file:///croot/sip_1698675935381/work
spx==1.16.0
sniffio==1.3.0
soupsieve==2.5
stack-data==0.6.3
statsmodels @ file:///croot/statsmodels_1718381181899/work
terminado==0.18.0
threadpoolctl @ file:///croot/threadpoolctl_1719407800858/work
tifffile @ file:///croot/tifffile_1695107451082/work
tinycss2==1.2.1
tokenize-rt==5.2.0
tomli==2.0.1
tomkit==0.13.2
tornado==6.4
tqdm @ file:///croot/tqdm_1724853939799/work
traitlets==5.14.0
trove-classifiers==2024.7.2
tweet-preprocessor==0.6.0
types-python-dateutil==2.8.19.14
typing_extensions @ file:///croot/typing_extensions_1715268824938/work
tzdata @ file:///croot/python-tzdata_1690578112552/work
unicodedata2 @ file:///croot/unicodedata2_1713212950228/work
uri-template==1.3.0
urllib3==2.1.0
virtualenv==20.25.0
vcwidth==0.2.12
webcolors==1.18
webencodings==0.5.1
websocket-client==1.7.0
wheel==0.43.0
widgetsnextension==4.0.9
zipp==3.20.1
```

7.1.2 Local Testing & Debugging

Optional local tests using a small toy dataset are sometimes provided to aid in debugging. The local tests are all stored in `localtests.py`

The autograder is the final arbiter

- There are no points associated with passing or failing the local tests, you must still pass the autograder to get points.
- It is possible to fail the local test and pass the autograder.
 - The autograder may have tolerances to account for minor implementation differences.
 - The reverse is also true, as the autograder may cover a larger number of corner cases.
- You do not need to pass both local and autograder tests to get points, passing the Gradescope autograder is sufficient for credit.

Work smarter, not harder

- Read the stack trace carefully. Often it will tell you exactly what's wrong.
- Understand what the local-test is doing. That way you can develop your own tests.
- Grow beyond the print statement: embrace a debugger. Jupyter-lab has a [built-in debugger](#) which allows you to look at data types, set breakpoints, and examine variables. If using a different IDE, look up your IDE's documentation on how to setup a proper debugger.
- Develop incrementally and test frequently, both locally and on Gradescope. Waiting to complete the whole class before testing can make it hard to isolate errors.

For this problem perform the following in the cell below:

- import `WarmupTests` from the `local tests.py`.
- Run the cell and submit `env.pkl`

```
import unittest

# Import WarmupTests from the local tests.py in the utilities folder.
# END YOUR CODE ABOVE #

#####
### DO NOT CHANGE THIS CELL ###
#####
result = unittest.main(
    argv=["Ignored", "WarmupTests.test_get_packages"], verbosity=1, exit=False)
if not result.result.wasSuccessful():
    sys.exit(1)
```

Passed test_get_packages

Ran 1 test in 0.006s

OK

7.2 Numpy Basics [2pt]

The following exercise will familiarize you with the basics of working with Numpy and navigating the [numpy documentation](#).

In `warmup.py` you will implement several "one-liners" using functions provided by numpy. No points will be awarded on Gradescope for any use of for loops or list comprehensions. Implement the following functions in `warmup.py`:

- `indices_of_k`
- `argmax_1d`
- `mean_rows`
- `sum_squares`

You may test your implementation with the below local tests. These local tests only checks the returned values of your implementation and does not check whether your implementation uses loops. Gradescope will check to make sure your implementation does not use loops (for, while, or list comprehensions).

WARNING: Make sure you match the dimensions of the output given in the comments of required function in `warmup.py`

HINT: Print and see what numpy functions are doing as much as you can!

```
#####
### DO NOT CHANGE THIS CELL ###
#####
import unittest

result = unittest.main(argv=[""], verbosity=1, exit=False)

if not result.result.wasSuccessful():
    sys.exit(1)
```

```
.....
-----
Ran 5 tests in 0.010s
```

```
OK
Correct values for argmax_1d
Passed test_get_packages
Correct values for indices_of_k
Correct values for mean_rows
Correct values for sum_squares
```

7.3 Broadcasting [1pt]

One of the simplest and most common similarity metrics in ML is the Manhattan Distance or [taxicab-distance](#). The function below takes two lists of N points in D dimensional space ($N \times D$ numpy arrays) and computes the Manhattan distance between every possible pair of points. *Hint: you can use this to try creating your own unittests*

Unfortunately such an implementation is too slow for a large dataset. In `fast_manhattan` leverage the broadcasting properties of numpy to create a faster version in a single line.

```
#####
### DO NOT CHANGE THIS CELL ###
#####

import numpy as np

def slow_manhattan(x, y):
    """
    Args:
        x: N x D numpy array
        y: M x D numpy array
    Return:
        dist: N x M numpy array, where dist[i, j] is the Manhattan distance between
        x[i, :] and y[j, :]
    """
    dist = np.empty((x.shape[0], y.shape[0]))
    for i in range(x.shape[0]):
        for j in range(y.shape[0]):
            d = 0
            for k in range(x.shape[1]):
                d += abs(x[i][k] - y[j][k])
            dist[i][j] = d
    return dist
```

Let's test the speed of this naive implementation:


```
%%timeit
```

```
#####  
### DO NOT CHANGE THIS CELL ###  
#####
```

```
x = np.random.rand(100, 3)  
y = np.random.rand(100, 3)  
d = slow_manhattan(x, y)
```

38 ms ± 881 µs per loop (mean ± std. dev. of 7 runs, 10 loops each)

Compare this with the vectorized implementation:

```
#####  
### DO NOT CHANGE THIS CELL ###  
#####
```

```
import varmup
```

```
%%timeit
```

```
#####  
### DO NOT CHANGE THIS CELL ###  
#####
```

```
x = np.random.rand(100, 3)  
y = np.random.rand(100, 3)  
d = varmup.fast_manhattan(x, y)
```

535 µs ± 5.48 µs per loop (mean ± std. dev. of 7 runs, 1,000 loops each)

Finally for `multiple_choice` answer the following by returning the correct integer value:

Which of the following best describes the space and time complexity of `slow_manhattan` compared to `fast_manhattan`:

- return 0 if: `fast_manhattan` has lower space and time complexity.
- return 1 if: `slow_manhattan` has lower space complexity and the same time complexity.
- return 2 if: `fast_manhattan` has higher space complexity and lower time complexity.
- return 3 if: Both are about the same in space and time complexity.

```
#####  
### DO NOT CHANGE THIS CELL ###  
#####
```

```
varmup.multiple_choice()
```