

hw2_coding

March 4, 2025

1 Problem1: The Need For Speed: Vectorization and Numpy

Performing mathematical operations on vectors and matrices is ubiquitous in most machine learning algorithms. Whether it's a simple similarity measure that works by calculating the dot product between two vectors, or deep neural networks, they all involve repeated matrix operations. This makes it imperative that our underlying code design performs matrix operations efficiently.

1.1 1.1 The Perils of Python

While Python is widely the language of choice for machine learning researchers across the globe (thanks to the speed of development and code readability it offers and the support it enjoys from the open-source community), Python as a high-level language on average is much slower than a lower level language like C++. To combat this, libraries like numpy and scipy implement most of the back-end operations they perform in C / C++, while providing wrappers in Python to be able to call underlying C code seamlessly from a Python script.

1.2 1.2 Speed Comparison: Numpy and Python

We highly recommend you to use numpy extensively in this course, it will be difficult to pass the programming portion of Homework 4 without writing most of your matrix operations in numpy. In this section, we'll see why.

Consider you have two vectors $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$. To see how similar they are, as measured by the cosine angle between them, you want to compute their dot product. This translates to the following operation:

$$\mathbf{a} \cdot \mathbf{b} = a_1 b_1 + a_2 b_2 + \dots + a_n b_n$$

When translated to code, notice how the dot product in NumPy is a whopping 100x faster than the native Python!

```
[12]: from timeit import timeit
import numpy as np
import array

VECTOR_SIZE = int(1e8)

# NumPy arrays
a = np.random.rand(VECTOR_SIZE)
```

```

b = np.random.rand(VECTOR_SIZE)

# Python arrays
aArr = array.array('d', a)
bArr = array.array('d', b)

def test_py_arr():
    sum = 0
    for i in range(VECTOR_SIZE):
        # TODO: apply dot product by using aArr and bArr
        sum+=aArr[i]+bArr[i]
        pass

    return sum

def test_np():
    # TODO: apply dot product by using a, b. Using numpy to vectorize to speed_
    ↪up the calculation
    a@b
    pass

# faster than multiprocessing , python lists , or numpy arrays with python loops
# faster than using a range and indexing

def time_dot_product (f):
    return timeit(f, setup=f, number =5) / 5

```

Test to time consumption!

```

[14]: print (f" NumPy = {time_dot_product (test_np) :.2f}s") # about 0.05 s
      print (f" Python on an array = {time_dot_product(test_py_arr) :.2f}s") # much_
      ↪longer time

```

NumPy = 0.03s

Python on an array = 7.31s

2 Problem2: k-Nearest Neighbor (kNN) exercise

Complete and hand in this completed worksheet (including its outputs and any supporting code outside of the worksheet) with your assignment submission.

The kNN classifier consists of two stages:

- During training, the classifier takes the training data and simply remembers it
- During testing, kNN classifies every test image by comparing to all training images and transferring the labels of the k most similar training examples
- The value of k is cross-validated

In this exercise you will implement these steps and understand the basic Image Classification pipeline, cross-validation, and gain proficiency in writing efficient, vectorized code.

```
[16]: # Run some setup code for this notebook.

import random
import numpy as np
from data_utils import load_CIFAR10
import matplotlib.pyplot as plt

# This is a bit of magic to make matplotlib figures appear inline in the
↳notebook
# rather than in a new window.
%matplotlib inline
plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plots
plt.rcParams['image.interpolation'] = 'nearest'
plt.rcParams['image.cmap'] = 'gray'

# Some more magic so that the notebook will reload external python modules;
# see http://stackoverflow.com/questions/1907993/
↳autoreload-of-modules-in-ipython
%load_ext autoreload
%autoreload 2

[17]: # Load the raw CIFAR-10 data.
cifar10_dir = 'datasets/cifar-10-batches-py'

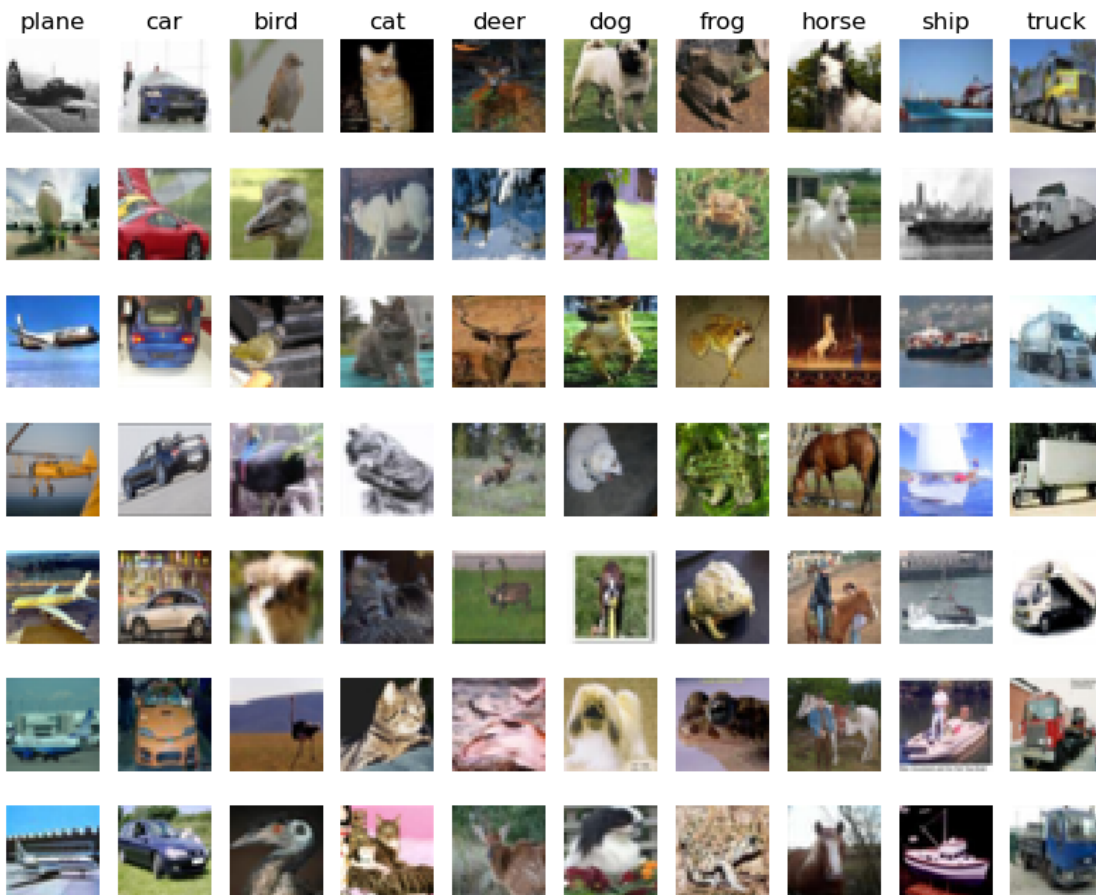
# Cleaning up variables to prevent loading data multiple times (which may cause
↳memory issue)
try:
    del X_train, y_train
    del X_test, y_test
    print('Clear previously loaded data.')
except:
    pass

X_train, y_train, X_test, y_test = load_CIFAR10(cifar10_dir)

# As a sanity check, we print out the size of the training and test data.
print('Training data shape: ', X_train.shape)
print('Training labels shape: ', y_train.shape)
print('Test data shape: ', X_test.shape)
print('Test labels shape: ', y_test.shape)
```

```
Training data shape: (50000, 32, 32, 3)
Training labels shape: (50000,)
Test data shape: (10000, 32, 32, 3)
Test labels shape: (10000,)
```

```
[18]: # Visualize some examples from the dataset.
# We show a few examples of training images from each class.
classes = ['plane', 'car', 'bird', 'cat', 'deer', 'dog', 'frog', 'horse', 'ship', 'truck']
num_classes = len(classes)
samples_per_class = 7
for y, cls in enumerate(classes):
    # print(y, cls)
    idxs = np.flatnonzero(y_train == y) # y_train == y
    idxs = np.random.choice(idxs, samples_per_class, replace=False) # np.random.
    # choice idxs samples_per_class replace=False
    for i, idx in enumerate(idxs):
        plt_idx = i * num_classes + y + 1
        plt.subplot(samples_per_class, num_classes, plt_idx)
        plt.imshow(X_train[idx].astype('uint8'))
        plt.axis('off')
        if i == 0:
            plt.title(cls)
plt.show()
```



```
[19]: # Subsample the data for more efficient code execution in this exercise
num_training = 5000
mask = list(range(num_training))
X_train = X_train[mask]
y_train = y_train[mask]

num_test = 500
mask = list(range(num_test))
X_test = X_test[mask]
y_test = y_test[mask]

print('X_train origin shape : ', X_train.shape)
# Reshape the image data into rows
X_train = np.reshape(X_train, (X_train.shape[0], -1))
X_test = np.reshape(X_test, (X_test.shape[0], -1))
print(X_train.shape, X_test.shape)
```

```
X_train origin shape : (5000, 32, 32, 3)
(5000, 3072) (500, 3072)
```

```
[20]: from k_nearest_neighbor import KNearestNeighbor

# Create a kNN classifier instance.
# Remember that training a kNN classifier is a noop:
# the Classifier simply remembers the data and does no further processing
classifier = KNearestNeighbor()
classifier.train(X_train, y_train)
```

We would now like to classify the test data with the kNN classifier. Recall that we can break down this process into two steps:

1. First we must compute the distances between all test examples and all train examples.
2. Given these distances, for each test example we find the k nearest examples and have them vote for the label

Lets begin with computing the distance matrix between all training and test examples. For example, if there are N_{tr} training examples and N_{te} test examples, this stage should result in a $N_{te} \times N_{tr}$ matrix where each element (i,j) is the distance between the i-th test and j-th train example.

Note: For the three distance computations that we require you to implement in this notebook, you may not use the `np.linalg.norm()` function that numpy provides. Instead, you should just think about the defination of ‘norm’ and implement it by yourself.

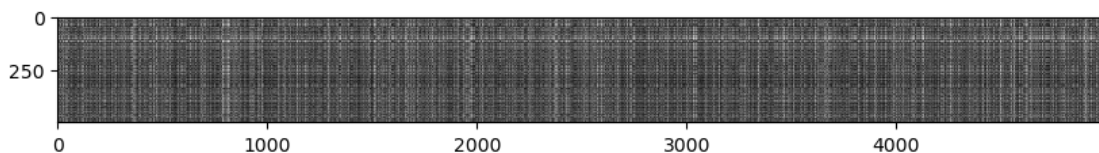
First, open `k_nearest_neighbor.py` and implement the function `compute_distances_two_loops` that uses a (very inefficient) double loop over all pairs of (test, train) examples and computes the distance matrix one element at a time.

```
[22]: # Open hw2_coding/k_nearest_neighbor.py and implement
      # compute_distances_two_loops.
```

```
# Test your implementation:
dists = classifier.compute_distances_two_loops(X_test)
print(dists.shape)
```

(500, 5000)

```
[23]: # We can visualize the distance matrix: each row is a single test example and
      # its distances to training examples
plt.imshow(dists, interpolation='none')
plt.show()
```



Inline Question 1

Notice the structured patterns in the distance matrix, where some rows or columns are visibly brighter. (Note that with the default color scheme black indicates low distances while white indicates high distances.)

- What in the data is the cause behind the distinctly bright rows?
- What causes the columns?

Your Answer : 1. A bright row in the distance matrix means that the corresponding test point is very distant from most or all of the training points. 2. A bright column in the distance matrix means that the corresponding training point is very distant from most or all of the test points.

Now implement the function `predict_labels` in `k_nearest_neighbor.py`, and run the code below:

```
[27]: # We use k = 1 (which is Nearest Neighbor).
y_test_pred = classifier.predict_labels(dists, k=1)

# Compute and print the fraction of correctly predicted examples
num_correct = np.sum(y_test_pred == y_test)
accuracy = float(num_correct) / num_test
print('Got %d / %d correct => accuracy: %f' % (num_correct, num_test, accuracy))
```

Got 137 / 500 correct => accuracy: 0.274000

You should expect to see approximately 27% accuracy. Now lets try out a larger k, say `k = 5`:

```
[29]: y_test_pred = classifier.predict_labels(dists, k=5)
num_correct = np.sum(y_test_pred == y_test)
```

```
accuracy = float(num_correct) / num_test
print('Got %d / %d correct => accuracy: %f' % (num_correct, num_test, accuracy))
```

Got 139 / 500 correct => accuracy: 0.278000

You should expect to see a slightly better performance than with $k = 1$.

In the following parts of this section, please open `k_nearest_neighbor.py`, then implement the functions `compute_distances_one_loop`, and `compute_distances_no_loops`.

```
[32]: # Now lets speed up distance matrix computation by using partial vectorization
# with one loop. Implement the function compute_distances_one_loop and run the
# code below:
dists_one = classifier.compute_distances_one_loop(X_test)
```

```
# To ensure that our vectorized implementation is correct, we make sure that it
# agrees with the naive implementation. There are many ways to decide whether
# two matrices are similar; one of the simplest is the Frobenius norm. In case
# you haven't seen it before, the Frobenius norm of two matrices is the square
# root of the squared sum of differences of all elements; in other words,
```

```
↳ reshape
# the matrices into vectors and compute the Euclidean distance between them.
```

```
difference = np.linalg.norm(dists - dists_one, ord='fro')
```

```
print('One loop difference was: %f' % (difference, ))
```

```
if difference < 0.001:
```

```
    print('Good! The distance matrices are the same')
```

```
else:
```

```
    print('Uh-oh! The distance matrices are different')
```

One loop difference was: 0.000000

Good! The distance matrices are the same

```
[33]: # Now implement the fully vectorized version inside compute_distances_no_loops
# and run the code
dists_two = classifier.compute_distances_no_loops(X_test)
```

```
# check that the distance matrix agrees with the one we computed before:
```

```
difference = np.linalg.norm(dists - dists_two, ord='fro')
```

```
print('No loop difference was: %f' % (difference, ))
```

```
if difference < 0.001:
```

```
    print('Good! The distance matrices are the same')
```

```
else:
```

```
    print('Uh-oh! The distance matrices are different')
```

No loop difference was: 0.000000

Good! The distance matrices are the same

```
[34]: # Let's compare how fast the implementations are
def time_function(f, *args):
```

```

"""
    Call a function f with args and return the time (in seconds) that it took
    to execute.
"""
import time
tic = time.time()
f(*args)
toc = time.time()
return toc - tic

two_loop_time = time_function(classifier.compute_distances_two_loops, X_test)
print('Two loop version took %f seconds' % two_loop_time)

one_loop_time = time_function(classifier.compute_distances_one_loop, X_test)
print('One loop version took %f seconds' % one_loop_time)

no_loop_time = time_function(classifier.compute_distances_no_loops, X_test)
print('No loop version took %f seconds' % no_loop_time)

# You should see significantly faster performance with the fully vectorized
# implementation!

# NOTE: depending on what machine you're using,
# you might not see a speedup when you go from two loops to one loop,
# and might even see a slow-down.

```

```

Two loop version took 18.545219 seconds
One loop version took 32.122436 seconds
No loop version took 0.145172 seconds

```

2.0.1 Cross-validation

We have implemented the k-Nearest Neighbor classifier but we set the value $k = 5$ arbitrarily. We will now determine the best value of this hyperparameter with cross-validation.

```

[51]: num_folds = 5
k_choices = [1, 3, 5, 8, 10, 12, 15, 20, 50, 100] # Stored all possible choices
        of k

X_train_folds = []
y_train_folds = []
#####
# TODO:
# Split up the training data into folds. After splitting, X_train_folds and
# y_train_folds should each be lists of length num_folds, where
# y_train_folds[i] is the label vector for the points in X_train_folds[i].
# Hint: Look up the numpy array_split function.
#####

```



```

# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
X_train_folds=np.array_split(X_train,num_folds)
y_train_folds=np.array_split(y_train,num_folds)
pass

# *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

# A dictionary holding the accuracies for different values of k that we find
# when running cross-validation. After running cross-validation,
# k_to_accuracies[k] should be a list of length num_folds giving the different
# accuracy values that we found when using that value of k.
k_to_accuracies = {}

#####
# TODO:
# Perform k-fold cross validation to find the best value of k. For each
# possible value of k, run the k-nearest-neighbor algorithm num_folds times,
# where in each case you use all but one of the folds as training data and the
# last fold as a validation set. Store the accuracies for all fold and all
# values of k in the k_to_accuracies dictionary.
#####
# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
for k in (k_choices):
    k_to_accuracies[k]=[]
    for i in range(num_folds):
        X_train_fold = np.concatenate([X_train_folds[j] for j in
↪range(num_folds) if j != i])
        y_train_fold = np.concatenate([y_train_folds[j] for j in
↪range(num_folds) if j != i])
        X_test_fold=X_train_folds[i]
        y_test_fold=y_train_folds[i]

        classifier.train(X_train_fold, y_train_fold)
        dists = classifier.compute_distances_no_loops(X_test_fold)

        y1_test_pred = classifier.predict_labels(dists, k)

        num_correct = np.sum(y1_test_pred == y_test_fold)
        accuracy = float(num_correct) / num_test
        k_to_accuracies[k].append(accuracy)

pass

```

```
# *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
```

```
# Print out the computed accuracies
```

```
for k in sorted(k_to_accuracies):
```

```
    for accuracy in k_to_accuracies[k]:
```

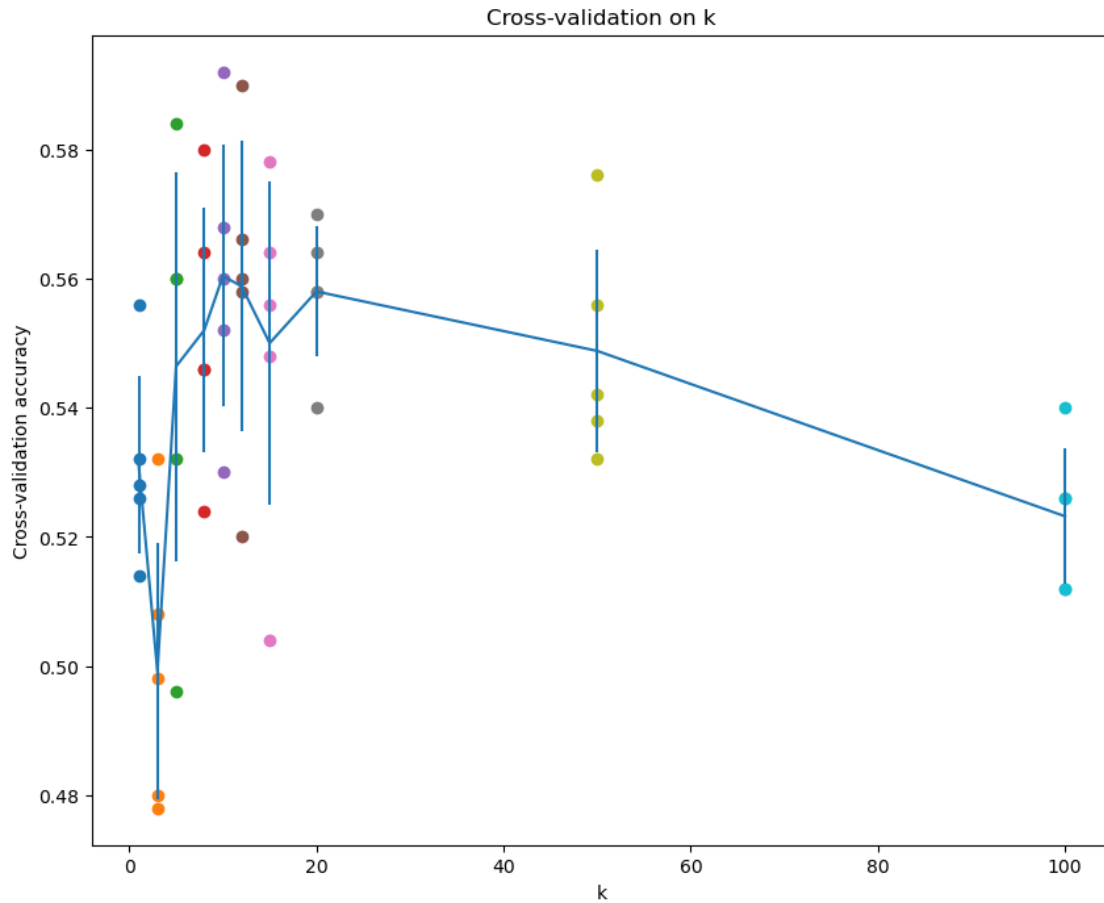
```
        print('k = %d, accuracy = %f' % (k, accuracy))
```

```
k = 1, accuracy = 0.526000
k = 1, accuracy = 0.514000
k = 1, accuracy = 0.528000
k = 1, accuracy = 0.556000
k = 1, accuracy = 0.532000
k = 3, accuracy = 0.478000
k = 3, accuracy = 0.498000
k = 3, accuracy = 0.480000
k = 3, accuracy = 0.532000
k = 3, accuracy = 0.508000
k = 5, accuracy = 0.496000
k = 5, accuracy = 0.532000
k = 5, accuracy = 0.560000
k = 5, accuracy = 0.584000
k = 5, accuracy = 0.560000
k = 8, accuracy = 0.524000
k = 8, accuracy = 0.564000
k = 8, accuracy = 0.546000
k = 8, accuracy = 0.580000
k = 8, accuracy = 0.546000
k = 10, accuracy = 0.530000
k = 10, accuracy = 0.592000
k = 10, accuracy = 0.552000
k = 10, accuracy = 0.568000
k = 10, accuracy = 0.560000
k = 12, accuracy = 0.520000
k = 12, accuracy = 0.590000
k = 12, accuracy = 0.558000
k = 12, accuracy = 0.566000
k = 12, accuracy = 0.560000
k = 15, accuracy = 0.504000
k = 15, accuracy = 0.578000
k = 15, accuracy = 0.556000
k = 15, accuracy = 0.564000
k = 15, accuracy = 0.548000
k = 20, accuracy = 0.540000
k = 20, accuracy = 0.558000
k = 20, accuracy = 0.558000
k = 20, accuracy = 0.564000
k = 20, accuracy = 0.570000
k = 50, accuracy = 0.542000
```

```
k = 50, accuracy = 0.576000
k = 50, accuracy = 0.556000
k = 50, accuracy = 0.538000
k = 50, accuracy = 0.532000
k = 100, accuracy = 0.512000
k = 100, accuracy = 0.540000
k = 100, accuracy = 0.526000
k = 100, accuracy = 0.512000
k = 100, accuracy = 0.526000
```

```
[53]: # plot the raw observations
for k in k_choices:
    accuracies = k_to_accuracies[k]
    plt.scatter([k] * len(accuracies), accuracies)

# plot the trend line with error bars that correspond to standard deviation
accuracies_mean = np.array([np.mean(v) for k,v in sorted(k_to_accuracies.
    ↪items())])
accuracies_std = np.array([np.std(v) for k,v in sorted(k_to_accuracies.
    ↪items())])
plt.errorbar(k_choices, accuracies_mean, yerr=accuracies_std)
plt.title('Cross-validation on k')
plt.xlabel('k')
plt.ylabel('Cross-validation accuracy')
plt.show()
```



```
[ ]: # Based on the cross-validation results above, choose the best value for k,
# retrain the classifier using all the training data, and test it on the test
# data. You should be able to get above 28% accuracy on the test data.

best_k = 10
# from the graph above, we could see that when k = 10, it has the best result_
↳ through cross-validation

classifier = KNearestNeighbor()
classifier.train(X_train, y_train)
y_test_pred = classifier.predict(X_test, k=best_k)

# Compute and display the accuracy
num_correct = np.sum(y_test_pred == y_test)
accuracy = float(num_correct) / num_test
print('Got %d / %d correct => accuracy: %f' % (num_correct, num_test, accuracy))
```

Inline Question 2

Which of the following statements about k -Nearest Neighbor (k -NN) are true in a classification setting, and for all k ? Select all that apply. 1. The decision boundary of the k -NN classifier is linear. 2. The training error of a 1-NN will always be lower than or equal to that of 5-NN. 3. The test error of a 1-NN will always be lower than that of a 5-NN. 4. The time needed to classify a test example with the k -NN classifier grows with the size of the training set. 5. None of the above.

Your Answer :

2,4

Your Explanation : 1. the decision boundary of k -NN is nonlinear. 2. The training error of 1-NN must be 0. 3. We can find counter-example in the result above. 4. It is true because the algorithm needs to compute the distance between the test example and every training example to find the nearest neighbors.

[]: