

1. **Directed graphical models and probability inference** I have 4 tops: a red sweater, a blue T-shirt, a green hoodie, and a white tank top. I need your help to decide what to wear.

(a) (1 pts) I decide what to wear based on four factors:

- if it's raining
- if I want to take a walk outside
- if I feel sick
- the day of the week it is

Using the Naive Bayes assumption, draw a graphical model that indicates how I will make a decision of what to wear each day.

(b) (1 pts) To be more specific,

- I only wear the green hoodie when I walk outside, regardless of all other factors.
- If I feel sick, I will walk outside 10% of the time. If I feel well, I will walk outside 60% of the time.
- When it rains, I feel sick 70% of the time; otherwise, I feel sick 15% of the time.

Draw a corresponding graphical model for determining whether I wear a green hoodie. Given that it is raining, infer the probability that I am wearing a green hoodie.

(c) (1 pts) The probability that I wear a tank top, independently of all the other clothes, is 75% if it's raining and 25% if it's not raining.

- Today is Monday and it is raining.
- The probability that it will rain, given that the previous day rained, is 70%. The probability that it will rain, given that the previous day did not rain, is 10%.

Draw a graphical model predicting whether I will wear a tank top on Wednesday, and calculate this probability.

## 2. Clustering(3 pts)

- Open the python notebook `mnist_dimred_clustering.ipynb`, and load the MNIST data by running the first cell. Don't change the way I've formatted it, or the checksums won't work.
- Fill in the function for determining the Euclidean distance between any sample point and the entire training data. From the print function, you should get a value of `160239119987.1912`.
- Using this function, code up a K-Means method, using an initialization of

```
mu = X[:K,:]
```

where  $K = 10$  is the number of clusters. Run this for 10 iterations.

- Define the *purity* of a set  $\mathcal{S}$  as the following fraction:

$$\text{class\_purity}(\mathcal{S}) = \max_i \frac{|\{y_j = i, j \in \mathcal{S}\}|}{|\mathcal{S}|}$$

that is, it is the size of the largest same-label subset of  $\mathcal{S}$  divided by the total size of  $\mathcal{S}$ . Report the purity of your clustering method after running k-means *on only the first 25 datapoints*, up to 3 digits after the decimal. Plot also the clustering result.

- Now, put the above section in a nice function, as we will use it again and again to evaluate future embeddings.
- At this point, we should also de-mean the data, as it will improve the embedding performances:

```
X = X - np.outer(np.ones(X.shape[0]), np.mean(X, axis=0)).
```

- Using `np.linalg.svd`, implement PCA, and reduce the data dimension to  $d = 10, 100$ , and  $500$ . Re-run K-Means to get new class memberships. (You may want to enact the option `full_matrices = False` when computing the SVD.)

- Use random hashing (as promoted by the JL lemma) and reduce the feature dimension to 10,100,500 dimensions, and cluster with MNIST.
- Use Isomap, LLE, or spectral embeddings to reduce dimension to 10,100,500, and cluster with MNIST. You only need to do one of the three, but you do need to code it up from scratch. When constructing the nearest neighbor graph, use Euclidean distance, and pick a threshold that seems reasonable to you. (Can be distance thresholding or nearest number of neighbors.) For isomap, you may find `scipy.sparse.csgraph.shortest_path` useful, if you can create the appropriate CSR formatted adjacency matrix. For all other hyperparameters, make a choice that seems reasonable to you.
- Use sklearn's t-SNE to reduce dimension to 1,2,3, and cluster with MNIST. In addition, plot something that shows the clustering effect, either in 1,2, or 3-D space. Keep most default hyperparameter settings to make your answers comparable.
- For these three situations, return the purity for each method and dimension size.
- Comment on the effect of the different dimensionality reduction schemes in the MNIST clustering task. What are the tradeoffs, in terms of performance and computational complexity?

3. **Hidden Markov Model spellchecker(4 pts)** In this exercise we will make a spell-checker using a HMM. To do this, download `alice_nlp_release.ipynb` and follow the instructions.

- Read through the first two blocks to get an idea of what the task is. The idea is to go through the corrupted corpus, identify words which have probably been corrupted, and correct them probabilistically.
- In the 4th box, fill in the functions to construct the word probabilities (weighted frequencies in uncorrupted corpus) and transition matrix (which gives  $\Pr(\text{word} \mid \text{prev word})$ ). If done correctly, the lines printed out should read

```
prob. of "alice" 0.014548615047424706
prob. of "queen" 0.002569625514869818
prob. of "chapter" 0.0009069266523069947
```

with smoothing

```
prob. of "the alice" 0.00025406504065040653
prob. of "the queen" 0.016514227642276422
prob. of "the chapter" 0.012957317073170731
```

no smoothing

```
prob. of "the alice" 0.0
prob. of "the queen" 0.03968253968253968
prob. of "the chapter" 0.0
prob. of "the hatter" 0.031135531135531136
```

- In the 5th box, fill in the function for computing the emission probability. The first 10 words closest to Alice should be
- ```
['abide', 'alice', 'above', 'voice', 'alive', 'twice', 'thick', 'dance', 'stick', 'prize']
```
- Construct and run your Hidden Markov Model spell checker using the functions computed for the prior probabilities, emission probabilities, and transition probabilities. List some words whose spelling was corrected correctly, and some examples where the spell-correcter did not work as expected. Report the recovery rate of the “fixed” corpus.

## Challenge!

### 1. Correlated mixture of sequential experts.

I want to buy a yacht, but I'm not sure if it's a good idea given the economy. So, I decide to question  $m$  consultants. Each consultant has more-or-less the same qualifications, and they come in one at a time.

The first consultant comes in my office. I ask, "Should I buy a yacht?" She says yes with probability  $p$ .

On her way out the building, she meets the second consultant. They chat briefly, and she leaves, he comes, and the process repeats. Each time, I ask the consultant if I should buy a yacht, and receive "yes" with probability  $p$ ; each time, the consultant chats briefly with the next consultant. *However*, the answers now are *not* i.i.d., but rather each expert's answer is correlated with the answer of experts he/she chatted with in the lobby.

At the end of the day, I have met with  $m$  consultants. I will make a decision whether to buy a yacht based on majority rule. We will now calculate the probability that I will buy a yacht.

- (a) We model the answer of each consultant as  $Y_i = 1$  if the  $i$ th consultant recommended "yes", and  $Y_i = -1$  otherwise. Show that if distribution

$$\Pr(Y_1 = 1) = p, \quad \Pr(Y_i = 1|Y_{i-1} = 1) = c + p - cp, \quad \Pr(Y_i = -1|Y_{i-1} = -1) = cp - p + 1$$

then  $\Pr(Y_i = 1) = p$  for all  $i$ .

- (b) The Pearson correlation coefficient between a random variable  $U$  and  $V$  can be expressed as

$$\text{corr}(U, V) = \frac{\mathbb{E}[UV] - \mathbb{E}[U]\mathbb{E}[V]}{\sqrt{\text{var}(U)\text{var}(V)}}.$$

Show that the correlations between each pair of sequential experts  $\text{corr}(Y_i, Y_{i-1}) = c$ ,

Hint: Go ahead and use a symbolic calculator, like WolframAlpha, to simplify messy expressions.

- (c) For  $m = 3$ , what is the probability that I will buy a yacht, in terms of  $P_{11} = \Pr(Y_i = 1|Y_{i-1} = 1), P_{00} = \Pr(Y_i = -1|Y_{i-1} = -1)$  and  $p$ ?
- (d) **Code.** Compute exactly the probability that I will buy a yacht, for  $m = 10$ , in Python or MATLAB. Generate a plot that shows the probability that I will buy a yacht, sweeping  $c \in [-1, 1]$ . (Note that there are cases where  $p$  and  $c$  are an infeasible pair, and can be detected when probabilities are not in the range (0,1)—these cases should not be plotted.) Do this for several interesting values of  $p$ . Comment on how the decision changes as a function of  $c$ ,  $p$ , and  $m$ .
- (e) Simulate the sequential advisers, and give a numerical estimate of what my decision will be if  $m = 25$  and  $m = 100$ . Generate a similar plot, using these numerical estimates of  $\Pr(\text{I will buy a yacht})$ . Use your own discretion to decide how many trials you need, and what values of  $c$  and  $p$  are useful.