

### Exercise 1

i) Lagrangian:  $g(\lambda, \omega) = L(\omega) - \lambda \left( \sum_{i=1}^n \omega_i - 1 \right)$

Now we want to solve for the lagrangian to find the argmax  $L(\omega)$  given the constraints, we can do so by solving the following:

$$\nabla g(\lambda, \omega) = 0 \Leftrightarrow \begin{cases} \nabla L(\omega) - \lambda \nabla g(\omega) = 0 \\ g(\omega) = 0 \end{cases} \quad (\text{I})$$

so now:

$$g(\lambda, \omega) = \sum_{i=1}^n c_i \log \omega_i - \lambda \left( \sum_{i=1}^n \omega_i - 1 \right)$$

taking the derivative in terms of  $\omega$ :

$$\begin{aligned} \frac{\partial g(\lambda, \omega)}{\partial \omega_a} &= \sum_{i=1}^n \frac{c_i}{\omega_i} - \lambda \sum_{i=1}^n 1 \quad \text{cw-} \\ &= \left[ \frac{c_1}{\omega_1}, \dots, \frac{c_n}{\omega_n} \right]^T - \lambda [1, \dots, 1]^T = 0 \end{aligned}$$

Now rearranging this gives us:

$$\left[ \frac{c_1}{\omega_1}, \dots, \frac{c_n}{\omega_n} \right]^T = \lambda [1, \dots, 1]^T$$

$$\nabla L(\omega) - \lambda \nabla g(\omega) = 0$$

where  $\lambda = \left[ \frac{c_1}{\omega_1}, \dots, \frac{c_n}{\omega_n} \right]^T$  and  $\omega \in \mathbb{F}$  giving us extrema value which satisfies our constraint and thus satisfies the first part of (I), and furthermore we need  $g(\omega) = \sum_{i=1}^n \omega_i - 1 = 0$ , so then we can see the constraint set  $\mathcal{G} = \{\omega \in \mathbb{R}^n \mid \sum_{i=1}^n \omega_i = 1\}$  is satisfied giving us the desired result. Now

$$\nabla L(\omega) = \left[ \frac{c_1}{\omega_1}, \dots, \frac{c_m}{\omega_n} \right]^T = 0$$

states that  $\omega \in \mathcal{F}$  is an extrema, lets use the second derivative test to show that  $\omega$  is a maximum

$$\begin{aligned} \nabla_{\omega\omega} L(\omega) &= \sum_{i=1}^m \frac{c_{iab}}{\omega_{iab}} \\ (\text{II}) \quad &= \sum_{i=1}^m -\frac{(c_{iab})_i}{(\omega_{iab})_i^2} \end{aligned}$$

we know that  $c_1, \dots, c_m \geq 0$

denominator will always be positive since it is squared

And we can see that  $\nabla_{\omega\omega} L(\omega) < 0$ , by the second derivative test since our second derivative is negative at our only extrema  $\omega \in \mathcal{F}$ , we see that our extrema is a global maximum over  $\mathcal{F}$ .

ii) Note that in part i), we did not explicitly mention that  $\omega_i > 0$ , but it must be assumed implicitly as  $\log(\omega_i)$  must take in positive  $\omega_i$  values as logarithms cannot have non-positive arguments. So adding this fact to part i) allows us to explicitly state that

$$\hat{\omega} = \left[ \frac{c_1}{\sum_{i=1}^m c_i}, \dots, \frac{c_2}{\sum_{i=1}^m c_i} \right]$$

is a global maximum, as (II) shows us that the second derivative will give us a negative value where  $\omega \in \mathcal{B}$

## Exercise 2

i) we have that  $N=14$

Now let:  $\hat{q}_0 = 0 = \text{'does not buy computer'}$   
 $\hat{q}_1 = 1 = \text{'buys computer'}$

where  $c_0 := \sum_{s=1}^{14} \mathbb{1}(y_s = '0') = 4$

$$c_1 := \sum_{s=1}^{14} \mathbb{1}(y_s = '1') = 10$$

now normalizing the above quantity gives us  
a prior PMF of:

$$[\hat{q}_0, \hat{q}_1] = [0.29, 0.71]$$

does not  
buy computer

buys computer

ii)

Class conditional PMF for not buying computer:

age	income	student	credit rating
$\leq 30$ [31, 40]	low med high	no yes	fair excel
$> 40$			

$$[\hat{q}_{0,1}, \dots, \hat{q}_{0,10}] = \frac{1}{16} [2, 0, 2, 1, 1, 2, 3, 1, 1, 3]$$

Class conditional PMF for buying computer:

age	income	student	credit rating
$\leq 30$ [31, 40]	low med high	no yes	fair excel
$> 40$			

$$[\hat{q}_{1,1}, \dots, \hat{q}_{1,10}] = \frac{1}{40} [3, 4, 3, 3, 5, 2, 5, 5, 7, 3]$$

iii) Compute the predictive probabilities for the following:

$x_1 = \text{"age} \leq 30, \text{medium income, student, fair credit rating"}$



$$\phi(x_1) = [1, 0, 0, 0, 1, 0, 0, 1, 1, 0]$$

$$P(Y=0 | \phi(x_1)) \propto \left[ \frac{2 \cdot 0 \cdot 2 \cdot 1 \cdot 1 \cdot 2 \cdot 3 \cdot 1 \cdot 1 \cdot 3}{16^4} (0.29) \right] = 7.93 \cdot 10^{-6}$$

$$P(Y=1 | \phi(x_1)) \propto \left[ \frac{3 \cdot 4 \cdot 3 \cdot 3 \cdot 5 \cdot 2 \cdot 5 \cdot 5 \cdot 7 \cdot 3}{40^4} (0.71) \right] = 0.00014$$

Now normalizing gives us:

$$P(Y=0 | \phi(x_1)) = \frac{7.93 \cdot 10^{-6}}{0.00014793} = 0.05$$

$$P(Y=1 | \phi(x_1)) = \frac{0.00014}{0.00014793} = 0.95$$

I predict that test sample  $x_1$  will choose to buy a computer as our model gave us a 95% chance that this will happen. The factor that most likely contributed to this is the fair credit rating.

$x_2 = \text{"age} \in [31,40], \text{low income, not student, fair credit rating"}$



$$\phi(x_2) = [0, 1, 0, 1, 0, 0, 1, 0, 1, 0]$$

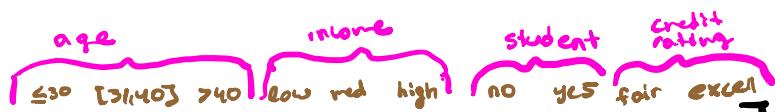
$$P(Y=0 | \phi(x_2)) \propto \left[ \frac{2 \cdot 0 \cdot 2 \cdot 1 \cdot 1 \cdot 2 \cdot 3 \cdot 1 \cdot 1 \cdot 3}{16^4} (0.29) \right] = 0.00001$$

$$P(Y=1 | \phi(x_2)) \propto \left[ \frac{3^{\circ} \cdot 4^{\circ} \cdot 3^{\circ} \cdot 5^{\circ} \cdot 2^{\circ} \cdot 5^{\circ} \cdot 5^{\circ} \cdot 7^{\circ} \cdot 3^{\circ}}{40^4} (0.71) \right] = 0.00011$$

$$P(Y=0 | \phi(x_2)) = \frac{0.00001}{0.00012} = 0.08$$

$$P(Y=1 | \phi(x_2)) = \frac{0.00011}{0.00012} = 0.92$$

I predict  $x_2$  will end up buying the computer. Once again the factor that affected this prediction the most is the 'fair credit rating'.



$$\phi(x_3) = [0, 0, 1, 0, 0, 1, 1, 0, 0, 1, 1]$$

$$P(Y=0 | \phi(x_3)) \propto \left[ \frac{2^{\circ} \cdot 0^{\circ} \cdot 2^{\circ} \cdot 1^{\circ} \cdot 1^{\circ} \cdot 2^{\circ} \cdot 3^{\circ} \cdot 1^{\circ} \cdot 1^{\circ} \cdot 3^{\circ}}{16^4} (0.29) \right] = 0.00015$$

$$P(Y=1 | \phi(x_3)) \propto \left[ \frac{3^{\circ} \cdot 4^{\circ} \cdot 3^{\circ} \cdot 5^{\circ} \cdot 2^{\circ} \cdot 5^{\circ} \cdot 5^{\circ} \cdot 7^{\circ} \cdot 3^{\circ}}{40^4} (0.71) \right] = 0.00002$$

$$P(Y=0 | \phi(x_3)) = \frac{0.00015}{0.00017} = 0.88$$

$$P(Y=1 | \phi(x_3)) = \frac{0.00002}{0.00017} = 0.12$$

I predict that  $x_3$  will NOT buy a computer, as our model gave a 88% of this happening. The attributes that contributed to this the most is the fact that they are NOT students and their excellent credit rating.

iv) I went ahead and used the dataset given to us by the professor. To process the data I used the CountVectorizer that returns a matrix w/ words as rows and has 05 columns.

```
#vectorizer = TfidfVectorizer(stop_words=stopwords_list)
vectorizer = CountVectorizer(stop_words=stopwords_list)
vectors = vectorizer.fit_transform(data_cleaned).transpose() # words x docs # in the form of sparse matrix
idx_to_word = np.array(vectorizer.get_feature_names()) # list of words that corresponds to feature coordinates

print('">>>> vectors.shape', vectors.shape)
i = 3
print('newsgroups_labels[i] (category)', newsgroups_labels[i])
#print('>>> data_cleaned[i] (article txt)', data_cleaned[i])
print('>>> vectors[:,i] \n', vectors[:,i])
offset = 5545
C0 = 0
C1 = 0
q0 = [0] * 45534
q1 = [0] * 45534
q0_total = 0
q1_total = 0
y_test = [0] * 20
q0_test = [0] * 45534
C_test = 0

#gets our train data to compute
for i in range(len(newsgroups_labels)-offset): #range ensures we get about 120 train cases
    if newsgroups_labels[i] == 0: #only looking for entries with label 1 or 0
        C0 += 1
        for j in range(45534):
            q0[j] += vectors[j,i]
            q0_total += vectors[j,i]
    elif newsgroups_labels[i] == 1: #see above ^
        C1 += 1
        for j in range(45534):
            q1[j] += vectors[j, i]
            q1_total += vectors[j,i]

q0_prior_pmf = (C0/(C0+C1))
q1_prior_pmf = (C1/(C0+C1))
q0_inflated = [x+1 for x in q0] #Added one to all entires in order to avoid multiplication by 0 in denominator
q1_inflated = [x+1 for x in q1]
#print('q0_total: ', q0_total)
#print('q1_total: ', q1_total)
#gets our test data
test_size = 0
for i in range(len(newsgroups_labels)-offset, len(newsgroups_labels)-(offset-40)): # -- range ensures we get about 30 test cases
    if newsgroups_labels[i] == 0 or newsgroups_labels[i] == 1:
        print('True label: ', newsgroups_labels[i], end=' ')
        C_test += 1
        y_word_count = 0
        numerator_0 = 1
        numerator_1 = 1
        for j in range(45534): # -- goes through the row of words, computes numerator
            if(vectors[j,i] != 0):
                y_word_count += vectors[j,i]
                exp = vectors[j,i]
                #if(q0[j] != 0 and q1[j] != 0):
                #    numerator_0 *= (q0[j]**exp)
                #    numerator_1 *= (q1[j]**exp)
                numerator_0 *= (q0_inflated[j]**exp)
                numerator_1 *= (q1_inflated[j]**exp)
                #print(' This is final numerator 0: ',numerator_0)
                #print(' This is final numerator 1: ',numerator_1)
        predict_prob_0 = ((numerator_0*q0_prior_pmf)/(q0_total ** y_word_count)) # -- computed predicted probabilities
        predict_prob_1 = ((numerator_1*q1_prior_pmf)/(q1_total ** y_word_count))
        pp0_norm = (predict_prob_0/(predict_prob_0 + predict_prob_1))
        pp1_norm = (predict_prob_1/(predict_prob_0 + predict_prob_1))
        print('Predicted probability of 0: ', pp0_norm, end=' ')
        print('Predicted probability of 1: ', pp1_norm)
```

I did although experience a handful amount of overflow issues in the final computation, which would sometimes return negative probabilities, but for the most cases of our test cases the model did a relatively well job predicting the newspaper articles.