RECITATION 7 HIDDEN MARKOV MODELS

10-601: Introduction to Machine Learning $\frac{11}{11}/2022$

Version: 1.2

1 HMMs

You are given the following training data:

win_C league_C Liverpool_D

win_C Liverpool_D league_C

Liverpool_D win_C

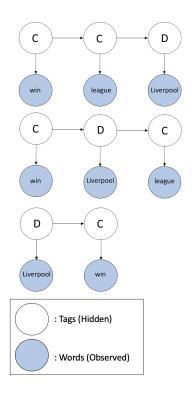


Figure 1: Visualization of Sequences

You are also given the following observed (validation) data: Liverpool win league

1.1 Initial, Emission, and Transition Matrices

Let each observed state $x_t \in \{1, 2, 3\}$, where 1 corresponds to win, 2 corresponds to league, and 3 corresponds to Liverpool. Let each hidden state $Y_t \in \{C, D\}$, where $s_1 = C$ and $s_2 = D$.

First, we need to estimate the HMM parameters - the initial probabilities: π , the transition probability matrix: **B**, and the emission probability matrix: **A**. Remember that we use MLE estimation to do so:

•
$$\hat{C}_k = \frac{N(Y_1^{(i)} = s_k)}{N} \ \forall \ i, k$$

•
$$\hat{B}_{jk} = \frac{N(Y_t^{(i)} = s_k, Y_{t-1}^{(i)} = s_j)}{N(Y_{t-1}^{(i)} = s_j)} \ \forall \ i, t > 1, j, k$$

•
$$\hat{A}_{jk} = \frac{N(X_t^{(i)} = k, Y_t^{(i)} = s_j)}{N(Y_t^{(i)} = s_j)} \ \forall \ i, t, j, k$$

Note: When learning an HMM, we add 1 to each count to make a pseudocount. This improves performance when evaluating unseen cases in the validation set or test set.

- 1. Find the initial matrix π . Recall that $\pi_j = P(Y_1 = s_j)$.
 - Find count matrix and pseudocount matrix:

$$\begin{array}{ccc}
Count & & Count \\
C & & \xrightarrow{\text{Pseudocount}} & C
\end{array}$$

$$D & D$$

• Normalize:

$$\pi = \begin{pmatrix} C \\ D \end{pmatrix}$$

• Find count matrix and pseudocount matrix:

$$\begin{array}{ccc}
Count & & Count \\
C & \begin{bmatrix} 2 \\ 1 \end{bmatrix} & \xrightarrow{\text{Pseudocount}} & C & \begin{bmatrix} 3 \\ 2 \end{bmatrix}
\end{array}$$

• Normalize:

$$\pi = \begin{array}{cc} C & \begin{bmatrix} 3/5 \\ D \end{bmatrix}$$

- 2. Find the transition matrix **B**. Recall that $B_{jk} = P(Y_t = s_k \mid Y_{t-1} = s_j)$.
 - Find count matrix and pseudocount matrix:

• Normalize:

$$\mathbf{B} = \begin{array}{cc} & C & D \\ \\ C & \\ D & \end{array}$$

• Find count matrix and pseudocount matrix:

$$\begin{array}{ccc}
 & C & D \\
C & \begin{bmatrix} 1 & 2 \\ 2 & 0 \end{bmatrix} & \xrightarrow{\text{Pseudocount}} & C & \begin{bmatrix} 2 & 3 \\ 3 & 1 \end{bmatrix}
\end{array}$$

• Normalize:

$$B = \begin{array}{c} C & D \\ C & \begin{bmatrix} 2/5 & 3/5 \\ 3/4 & 1/4 \end{bmatrix}$$

- 3. Find the emission matrix **A**. Recall that $A_{jk} = P(X_t = k \mid Y_t = s_j)$.
 - Find count matrix and pseudocount matrix:

• Normalize:

$$\mathbf{A} = \begin{array}{c} \text{win league Liverpool} \\ C \end{array}$$

• Find count matrix

• Normalize:

$$A = \begin{array}{ccc} & & \text{win league Liverpool} \\ C & & \begin{bmatrix} 1/2 & 3/8 & 1/8 \\ 1/6 & 1/6 & 2/3 \end{bmatrix} \end{array}$$

1.2 The Forward Algorithm

One type of inference problem that can be answered by an HMM is **Evaluation** - computing the probability of a sequence of observations. We calculate the likelihood of observing the validation sequence:

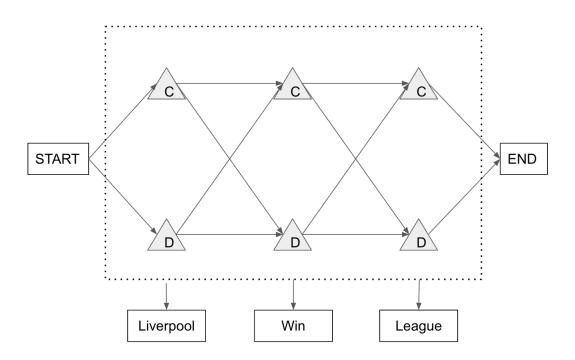
Liverpool win league

To do so, we calculate the forward probability matrix α . Recall that

$$\alpha_t(s_k) = P(x_{1:t}, Y_t = s_k)$$

We have the following bottom-up dynamic programming algorithm to calculate the forward probabilities:

```
for t = 1,...,T: for j = 1,...,J: if t == 1: \alpha_1(s_j) = \pi_j * A_{j,x_1} else: \alpha_t(s_j) = A_{j,x_t} * \sum_k \alpha_{t-1}(s_k) * B_{k,j}
```



First, use the algorithm as defined above to calculate $\alpha_1(C)$ and $\alpha_1(D)$.

•
$$\alpha_1 = \begin{bmatrix} \alpha_1(C) \\ \alpha_1(D) \end{bmatrix} = \begin{bmatrix} \pi_C * A_{C,x_1} \\ \pi_D * A_{C,x_1} \end{bmatrix} = \begin{bmatrix} \pi_C * A_{C,Liverpool} \\ \pi_D * A_{D,Liverpool} \end{bmatrix} = \begin{bmatrix} 3/5 \cdot 1/8 \\ 2/5 \cdot 2/3 \end{bmatrix} = \begin{bmatrix} 0.075 \\ 0.2667 \end{bmatrix}$$

Observe that this can be vectorized as $\pi \odot A_{x_1}$.

Indeed, the way **B** and **A** are constructed allows us to also vectorize the computation of the forward probabilities for $1 < t \le T$:

$$\mathbf{A}_{,x_t} \odot (\mathbf{B}^T \boldsymbol{lpha}_{t-1})$$

•
$$\alpha_2 = \begin{bmatrix} \alpha_2(C) \\ \alpha_2(D) \end{bmatrix} = \mathbf{A}_{,x_2} \odot (\mathbf{B}^T \boldsymbol{\alpha}_1) = \begin{bmatrix} 1/2 \\ 1/6 \end{bmatrix} \odot \begin{pmatrix} \begin{bmatrix} 2/5 & 3/4 \\ 3/5 & 1/4 \end{bmatrix} \begin{bmatrix} 0.075 \\ 0.2667 \end{bmatrix} \end{pmatrix} = \begin{bmatrix} 0.1150 \\ 0.0186 \end{bmatrix}$$

•
$$\alpha_3 = \mathbf{A}_{,x_3} \odot (\mathbf{B}^T \alpha_2) = \begin{bmatrix} 3/8 \\ 1/6 \end{bmatrix} \odot \begin{pmatrix} \begin{bmatrix} 2/5 & 3/4 \\ 3/5 & 1/4 \end{bmatrix} \begin{bmatrix} 0.1150 \\ 0.0186 \end{bmatrix} \end{pmatrix} = \begin{bmatrix} 0.0225 \\ 0.0123 \end{bmatrix}$$

To find the likelihood of observing the validation sequence, all we need are the final forward probabilities:

$$\begin{split} &P(X_1 = \texttt{Liverpool}, X_2 = \texttt{win}, X_3 = \texttt{league}) \\ &= \sum_{y_3 \in \{C,D\}} P(x_1 = \texttt{Liverpool}, x_2 = \texttt{win}, x_3 = \texttt{league}, Y_3 = y_t) \\ &= \sum_{y_t \in \{C,D\}} \alpha_3(y_t) \\ &= 0.0225 + 0.0123 \\ &= 0.0348 \end{split}$$

1.3 The Backward Algorithm

Another type of inference problem that can be answered by an HMM is computing **Marginals** - computing the marginal probability distribution for a hidden state, given a sequence of observations. Recall that

$$P(Y_t = s_k | \vec{x}) = \frac{\alpha_t(s_k)\beta_t(s_k)}{P(\vec{x})}$$

Therefore, along with the forward probability matrix α , we need to find the backward probability matrix β , where

$$\beta_t(s_k) = P(x_{t+1:T}|Y_t = s_k)$$

We have a similar bottom-up dynamic programming algorithm to calculate the backward probabilities:

```
for t = T,...,0:

for j = 1,...,k:

    if t == T:

    \beta_T(s_j) = 1

else:

    \beta_t(s_j) = \sum_{k=1}^J A_{k,x_{t+1}} \beta_{t+1}(s_k) B_{j,k}
```

Conveniently, there is also a matrix expression for the vector of backward probabilities for a given time step t < T:

$$\mathbf{B}(\mathbf{A}_{.x_{t+1}}\odot\boldsymbol{\beta}_{t+1})$$

•
$$\beta_3 = \begin{bmatrix} \beta_3(C) \\ \beta_3(D) \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

•
$$\boldsymbol{\beta}_2 = \begin{bmatrix} \beta_2(C) \\ \beta_2(D) \end{bmatrix} = \mathbf{B}(\mathbf{A}_{,x_3} \odot \boldsymbol{\beta}_3) = \begin{bmatrix} 2/5 & 3/5 \\ 3/4 & 1/4 \end{bmatrix} \begin{pmatrix} \begin{bmatrix} 3/8 \\ 1/6 \end{bmatrix} \odot \begin{bmatrix} 1 \\ 1 \end{bmatrix} \end{pmatrix} = \begin{bmatrix} 0.25 \\ 0.3229 \end{bmatrix}$$

•
$$\beta_1 = \mathbf{B}(\mathbf{A}_{,x_2} \odot \beta_2) = \begin{bmatrix} 2/5 & 3/5 \\ 3/4 & 1/4 \end{bmatrix} \left(\begin{bmatrix} 1/2 \\ 1/6 \end{bmatrix} \odot \begin{bmatrix} 0.25 \\ 0.3329 \end{bmatrix} \right) = \begin{bmatrix} 0.0823 \\ 0.1072 \end{bmatrix}$$

Now, we have our α and β matrices:

$$\alpha = \begin{bmatrix} C & D \\ 1 & \begin{bmatrix} 0.0750 & 0.2667 \\ 0.1150 & 0.0186 \\ 0.0225 & 0.0123 \end{bmatrix}$$

$$\boldsymbol{\beta} = \begin{bmatrix} C & D \\ 1 & \begin{bmatrix} 0.0823 & 0.1072 \\ 0.2500 & 0.3229 \\ 1.0000 & 1.0000 \end{bmatrix}$$

1. What is $P(Y_2 = C|\vec{x})$?

$$P(Y_2 = C|\vec{x}) = \frac{\alpha_2(C)\beta_2(C)}{P(\vec{x})} = \frac{0.1150 \times 0.2500}{0.0348} = 0.8261$$

2. What is $P(Y_2 = D|\vec{x})$?

$$P(Y_2 = D|\vec{x}) = \frac{\alpha_2(D)\beta_2(D)}{P(\vec{x})} = \frac{0.0186 \times 0.3229}{0.0348} = 0.1726$$

3. What is $P(Y_3 = C | \vec{x})$?

$$P(Y_2 = C|\vec{x}) = \frac{\alpha_2(C)\beta_2(C)}{P(\vec{x})} = \frac{0.0225 \times 1}{0.0348} = 0.6466$$

4. What is the minimum Bayes risk (MBR) decoder prediction for Y_2 ?

This is the argmax over all marginal probabilities for the 2nd hidden state; here $P(Y_2 = C|\vec{x}) > P(Y_2 = D|\vec{x})$ so the prediction is hidden state C.

1.4 The Viterbi Algorithm

Instead of finding the most likely hidden state at some time t, we may instead want to find the most likely sequence of hidden states. This is known as **Viterbi Decoding** - computing the most probable assignment of hidden states, given a sequence of observations.

The sequence of words you observe is again the same: Liverpool win league

However, you are only given the tag of the last word: league_C

1. Recall that:

$$\omega_t(s_k) = \max_{y_{1:t-1}} P(x_{1:t}, y_{1:t-1}, y_t = s_k)$$

$$b_t(s_k) = \arg\max_{y_{1:t-1}} P(x_{1:t}, y_{1:t-1}, y_t = s_k)$$

Using the formulae above and the first order Markov assumption, derive a recursive definition for $\omega_t(s_k)$ and $b_t(s_k)$ that will let you employ bottom-up dynamic programming.

First, we decompose the joint probability:

$$P(x_{1:t}, y_{1:t-1}, y_t) = P(x_t, x_{1:t-1}, y_t, y_{t-1}, y_{1:t-1})$$

= $P(x_t|y_t)P(y_t|y_{t-1})P(x_{1:t-1}, y_{1:t-2}, y_{t-1})$

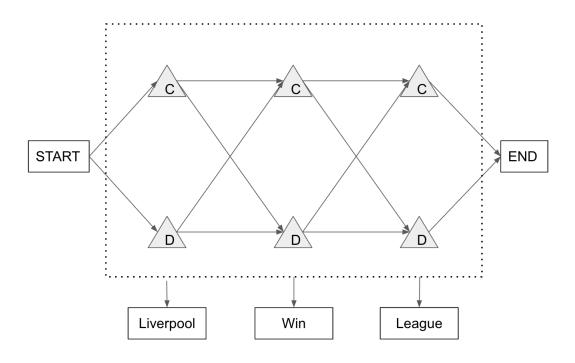
(First order Markov Assumption)

Now, for maximizing:

$$\begin{aligned} &\omega_t(s_k) = \max_{s_j} (P(x_t | y_t = s_k) P(y_t = s_k | y_{t-1} = s_j) \max_{y_{1:t-2}} (P(x_{1:t-1}, y_{1:t-2}, y_{t-1} = s_j))) \\ \Longrightarrow \\ &\omega_t(s_k) = \max_{s_j} A_{s_k, x_t} \cdot B_{s_j, s_k} \cdot \omega_{t-1}(s_j) \end{aligned}$$

$$b_t(s_k) = \underset{s_j}{\arg\max} A_{s_k, x_t} \cdot B_{s_j, s_k} \cdot \omega_{t-1}(s_j)$$

Below is the trellis corresponding to the given data:



- 2. Annotate the trellis at the nodes that correspond to:
 - (a) $\omega_1(C)$
 - (b) $\omega_1(D)$
 - (c) $\omega_2(C)$
 - (d) $\omega_2(D)$
 - (e) $\omega_3(C)$
 - (f) $\omega_3(D)$

- 3. Find the most likely sequence of tags given the observed data:
 - (a) Set up the matrices ω and b

$$\omega = \begin{bmatrix} \omega_0 & 0 & 1 \\ \omega_2 & \omega_3 & -\omega_2 \end{bmatrix}$$

and

Initialize $w_0(START) = 1$

(b) Solve for matrix entries using Dynamic Programming:

$$\begin{split} \omega_1(C) &= \max_{s_j \in \mathtt{C,D,START}} P(x_1 = \mathtt{Liverpool}|Y_1 = \mathtt{C}) \omega_0(s_j) P(Y_1 = \mathtt{C}) \\ &= \frac{1}{8} \cdot 1 \cdot \frac{3}{5} \end{split}$$

$$b_1(\mathbf{C}) = \mathtt{START}$$

$$\begin{split} \omega_1(D) &= \max_{s_j \in \mathtt{C},\mathtt{D},\mathtt{START}} P(x_1 = \mathtt{Liverpool}|Y_1 = \mathtt{D}) w_0(s_j) P(Y_1 = \mathtt{D}) \\ &= \frac{2}{3} \cdot 1 \cdot \frac{2}{5} \\ &= \frac{4}{15} \end{split}$$

$$b_1(D) = \texttt{START}$$

$$\begin{split} \omega_2(C) &= \max_{s_j \in \mathtt{C},\mathtt{D}} P(x_2 = \mathtt{win} | Y_2 = \mathtt{C}) \omega_1(s_j) P(Y_2 = \mathtt{C} | Y_1 = s_j) \\ &= \max \left(\frac{1}{2} \cdot \frac{3}{40} \cdot \frac{2}{5}, \frac{1}{2} \cdot \frac{4}{15} \cdot \frac{3}{4} \right. \\ &= \frac{1}{10} \end{split}$$

$$b_2(\mathbf{C}) = \mathbf{D}$$

$$\begin{split} \omega_2(D) &= \max_{s_j \in \mathtt{C},\mathtt{D}} P(x_2 = \mathtt{win} | Y_2 = \mathtt{D}) \omega_1(s_j) P(Y_2 = \mathtt{D} | Y_1 = s_j) \\ &= \max \left(\frac{1}{6} \cdot \frac{3}{40} \cdot \frac{3}{5}, \frac{1}{6} \cdot \frac{4}{15} \cdot \frac{1}{4} \right. \\ &= \frac{1}{90} \end{split}$$

$$b_2(D) = D$$

$$\begin{split} \omega_3(C) &= \max_{s_j \in \mathtt{C},\mathtt{D}} P(x_3 = \mathtt{league}|Y_3 = \mathtt{C}) \omega_2(s_j) P(Y_3 = \mathtt{C}|Y_2 = s_j) \\ &= \max\left(\frac{3}{8} \cdot \frac{1}{10} \cdot \frac{2}{5}, \frac{3}{8} \cdot \frac{1}{90} \cdot \frac{3}{4} \right. \\ &= \frac{3}{200} \end{split} \bigg)$$

$$b_3(\mathbf{C}) = \mathbf{C}$$

$$\begin{split} \omega_3(D) &= \max_{s_j \in \mathtt{C},\mathtt{D}} P(x_3 = \mathtt{league}|Y_3 = \mathtt{D}) \omega_2(s_j) P(Y_3 = \mathtt{D}|Y_2 = s_j) \\ &= \max\left(\frac{1}{6} \cdot \frac{1}{10} \cdot \frac{3}{5}, \frac{1}{6} \cdot \frac{1}{90} \cdot \frac{1}{4}\right) \\ &= \frac{1}{100} \end{split} \right)$$

$$b_3(D) = C$$

Now, to figure out the order, we set $\hat{y}_t = b_{t+1}(\hat{y}_{t+1})$

$$\hat{y}_{T+1} = \mathtt{END}$$

$$\begin{array}{l} \hat{y}_3 = b_4(\mathtt{END}) \\ = \mathbf{C} \end{array}$$

$$\hat{y}_2 = b_3(\mathbf{C}) \\
= \mathbf{C}$$

$$\begin{array}{l} \hat{y}_1 = b_2(\mathbf{C} \\ = \mathbf{D} \end{array}$$

$$\hat{y}_0 = b_1(D)$$
 $= START$

So, the most likely sequence is: START-D-C-C-END

2 Working in Log-space

2.1 Motivation

Some of the probabilities we work with in Homework 7 about are tiny and some of them are much larger. We tend to work with the tiny ones in log-space and only get back probabilities if we really need them for some other purpose. Throughout hw7 you will keep your probabilities in log-space. In this section we will motivate why we use log-space for small values.

Given the following series of probability values:

$P(x_1=1)$	$P(x_2 = 1 \mid x_1 = 1)$	$P(x_3 = 1 \mid x_2 = 1, x_1 = 1)$
0.002	0.004	0.003

We want to find $P(x_1 = 1, x_2 = 1, x_3 = 1)$. Suppose we have a calculator which only has 4 decimal places of precision, so it can only store values of format X.XXXX

1. What is the correct value of $P(x_1 = 1, x_2 = 1, x_3 = 1)$ without any precision limits?

$$P(x_1 = 1, x_2 = 1, x_3 = 1) = P(x_3 = 1 \mid x_2 = 1, x_1 = 1) * P(x_2 = 1 \mid x_1 = 1) * P(x_1 = 1) = 0.003 * 0.004 * 0.002 = 0.000000024$$

2. What is the value of $P(x_1 = 1, x_2 = 1, x_3 = 1)$ using our faulty calculator?

$$P(x_1 = 1, x_2 = 1)$$
= $P(x_2 = 1 \mid x_1 = 1)P(x_1 = 1) = 0.004 * 0.002 = 0.0000$
 \Rightarrow Truncated!

 $P(x_1 = 1, x_2 = 1, x_3 = 1) = 0.0000 * 0.003 = 0.0000$
 \Rightarrow Truncated again!

3. How do the values of $P(x_1 = 1, x_2 = 1, x_3 = 1)$ from part (1) and (2) compare?

No precision limits: $P(x_1 = 1, x_2 = 1, x_3 = 1) = 0.000000024$

Faulty calculator: $P(x_1 = 1, x_2 = 1, x_3 = 1) = 0.0000$

4. What is the value of $P(x_1 = 1, x_2 = 1, x_3 = 1)$ if we perform the same computation but in log space?

$$\log (P(x_1 = 1, x_2 = 1, x_3 = 1))$$

$$= \log(x_1 = 1) + \log(P(x_2 = 1 \mid x_1 = 1)) + \log(P(x_3 = 1 \mid x_2 = 1, x_1 = 1))$$

$$= \log(0.002) + \log(0.004) + \log(0.003)$$

$$= -6.2146 - 5.8091 - 5.5215$$

$$= -17.5452$$

If we were to recover our value of
$$P(x_1 = 1, x_2 = 1, x_3 = 1) = e^{\log(P(x_1 = 1, x_2 = 1, x_3 = 1))} = e^{-17.5452} = 0.000000024$$

This is good! But we can use the log sum exp trick to extend its use to even smaller scales.

2.2 Forward and Backward Algorithm in Log Space

In the forward algorithm, recall that the entries in α can be computed using the bottom-up dynamic programming algorithm:

- $\bullet \ \alpha_1(j) = \pi_j A_{jx_1}$
- For t > 1, $\alpha_t(j) = A_{jx_t} \sum_{k=1}^{J} \alpha_{t-1}(k) B_{kj}$
- 1. Derive $\log (\alpha_1(j))$ in terms of $\log(\pi_j)$ and $\log(A_{jx_1})$ $\log (\alpha_1(j)) = \log (\pi_j A_{jx_1}) = \log(\pi_j) + \log(A_{jx_1})$
- 2. Derive $\log (\alpha_t(j))$ in terms of $\log (\alpha_{t-1}(k))$ and $\log A_{kj}$ $\log (\alpha_t(j))$ $= \log (A_{jx_t} \sum_{k=1}^{J} \alpha_{t-1}(k) B_{kj})$ $= \log(A_{jx_t}) + \log (\sum_{k=1}^{J} \alpha_{t-1}(k) B_{kj})$

$$= \log(A_{jx_t}) + \log\left(\sum_{k=1}^{J} e^{\log\left(\alpha_{t-1}(k)B_{kj}\right)}\right)$$
$$= \log(A_{jx_t}) + \log\left(\sum_{k=1}^{J} e^{\log\left(\alpha_{t-1}(k)\right) + \log\left(B_{kj}\right)}\right)$$

In the backward algorithm, we also have a similar bottom-up dynamic programming algorithm:

- $\bullet \ \beta_T(j) = 1$
- For $1 \le t \le T 1$, $\beta_t(j) = \sum_{k=1}^J A_{kx_{t+1}} \beta_{t+1}(k) B_{jk}$
- 1. Derive $\log (\beta_T(j))$ $\log (\beta_T(j)) = \log(1) = 0$
- 2. Derive $\log (\beta_t(j))$ in terms of $\log(A_{kx_{t+1}})$, $\log (\beta_{t+1}(k))$, and $\log(B_{jk})$ $\log (\beta_t(j))$ $= \log (\sum_{k=1}^J A_{kx_{t+1}} \beta_{t+1}(k) B_{jk})$ $= \log \left(\sum_{k=1}^J e^{\log (A_{kx_{t+1}} \beta_{t+1}(k) B_{jk})}\right)$

$$= \log \left(\sum_{k=1}^{J} e^{\log(A_{kx_{t+1}}) + \log \left(\beta_{t+1}(k)\right) + \log(B_{jk})} \right)$$

After transforming the equations into log form, you may discover calculations of the following type:

$$\log \sum_{i} \exp\left(v_{i}\right)$$

This may be programmed as is, but $\exp(v_i)$ may cause underflow when v_i is large and negative. One way to avoid this is to use the log-sum-exp trick.

The log-sum-exp trick simply adds the maximum value in the vector to the log probabilities as follows:

$$\log \sum_{i} \exp(v_i - m)) + \max_{i}(v_i)$$

3 Dynamic Programming

DP Notebook

To access this Colab notebook you will need to be logged into Google Drive with your Andrew email.