
CSE 250a. Assignment 4

Out: Tue Oct 23

Due: Tue Oct 30 (in class)

4.1 Maximum likelihood estimation of a multinomial distribution

A $2N$ -sided die is tossed many times, and the results of each toss are recorded as data. Suppose that in the course of the experiment, the n^{th} side of the die is observed C_n times. For this problem, you should assume that the tosses are identically, independent distributed (i.i.d.) according to the probabilities of the die.

(a) **Log-likelihood**

Let $X \in \{1, 2, 3, \dots, 2N\}$ denote the outcome of a toss, and let $p_n = P(X = n)$ denote the probabilities of the die. Express the log-likelihood $\mathcal{L} = \log P(\text{data})$ of the observed results in terms of the probabilities p_n and the counts C_n .

(b) **Maximum likelihood estimate**

Derive the maximum likelihood estimates of the die's probabilities p_n . Specifically, maximize your expression for the log-likelihood \mathcal{L} in part (a) subject to the constraints

$$\sum_{n=1}^{2N} p_n = 1, \quad p_n \geq 0.$$

You should use a Lagrange multiplier to enforce the linear equality constraint, but it is sufficient to observe that the resulting solution is nonnegative.

(c) **Even versus odd**

Compute the probability $P(X \in \{2, 4, 6, \dots, 2N\})$ that the roll of a die is *even* and also the probability $P(X \in \{1, 3, 5, \dots, 2N - 1\})$ that the roll of a die is *odd*. Show that these two probabilities are equal when

$$\sum_{n=1}^{2N} (-1)^n p_n = 0.$$

(d) **Maximum likelihood estimate**

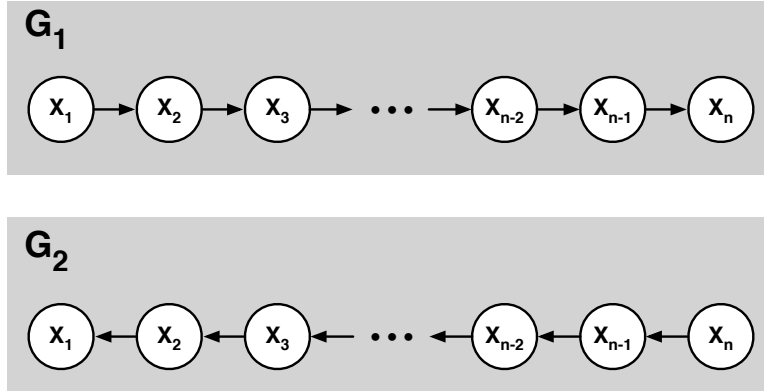
Suppose it is known a priori that the probability of an *even* toss is equal to that of an *odd* toss. Derive the maximum likelihood estimates of the die's probabilities p_n subject to this constraint. Specifically, maximize your expression for the log-likelihood \mathcal{L} in part (a) subject to the constraints

$$\sum_{n=1}^{2N} p_n = 1, \quad \sum_{n=1}^{2N} (-1)^n p_n = 0, \quad p_n \geq 0.$$

Hint: introduce two Lagrange multipliers, one for each linear equality constraint.

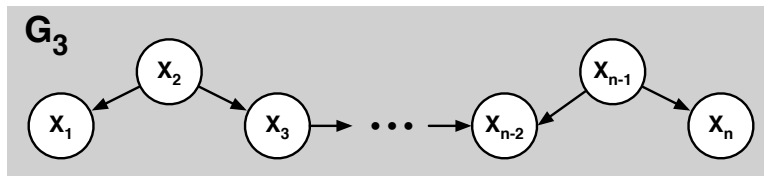
4.2 Maximum likelihood estimation in belief networks

Consider the two DAGs shown below, G_1 and G_2 , over the same nodes $\{X_1, X_2, \dots, X_n\}$, that differ only in the direction of their edges.



Suppose that we have a (fully observed) data set $\{x_1^{(t)}, x_2^{(t)}, \dots, x_n^{(t)}\}_{t=1}^T$ in which each example provides a complete instantiation of the nodes in these DAGs. Let $\text{COUNT}_n(x)$ denote the number of examples in which $X_n = x$, and let $\text{COUNT}_n(x, x')$ denote the number of examples in which $X_n = x$ and $X_{n+1} = x'$.

- Express the maximum likelihood estimates for the CPTs in G_1 in terms of these counts.
- Express the maximum likelihood estimates for the CPTs in G_2 in terms of these counts.
- Using your answers from parts (a) and (b), show that the maximum likelihood CPTs for G_1 and G_2 from this data set give rise to the same joint distribution over the nodes $\{X_1, X_2, \dots, X_n\}$.
- Suppose that some but not all of the edges in these DAGs were reversed, as in the graph G_3 shown below. Would the maximum likelihood CPTs for G_3 also give rise to the same joint distribution? (*Hint*: does G_3 imply all the same statements of conditional independence as G_1 and G_2 ? Look in particular at node X_{n-2} .)



4.3 Statistical language modeling

In this problem, you will explore some simple statistical models of English text. Download and examine the data files on Piazza for this assignment. (Start with the `readme.txt` file.) These files contain unigram and bigram counts for 500 frequently occurring tokens in English text. These tokens include actual words as well as punctuation symbols and other textual markers. In addition, an “unknown” token is used to represent all words that occur outside this basic vocabulary. For this problem, as usual, you may program in the language of your choice.

- (a) Compute the maximum likelihood estimate of the unigram distribution $P_u(w)$ over words w . Print out a table of all the tokens (i.e., words) that start with the letter “A”, along with their numerical unigram *probabilities* (not counts). (You do not need to print out the unigram probabilities for all 500 tokens.)
- (b) Compute the maximum likelihood estimate of the bigram distribution $P_b(w'|w)$. Print out a table of the five **most likely words to follow the** word “THE”, along with their numerical bigram probabilities.
- (c) Consider the sentence “**Last week the stock market fell by one hundred points.**” Ignoring punctuation, compute and compare the log-likelihoods of this sentence under the unigram and bigram models:

$$\mathcal{L}_u = \log \left[P_u(\text{last}) P_u(\text{week}) P_u(\text{the}) \dots P_u(\text{one}) P_u(\text{hundred}) P_u(\text{points}) \right]$$

$$\mathcal{L}_b = \log \left[P_b(\text{last}|\langle s \rangle) P_b(\text{week}|\text{last}) P_b(\text{the}|\text{week}) \dots P_b(\text{hundred}|\text{one}) P_b(\text{points}|\text{hundred}) \right]$$

In the equation for the bigram log-likelihood, the token $\langle s \rangle$ is used to mark the beginning of a sentence. Which model yields the highest log-likelihood?

- (d) Consider the sentence “**The nineteen officials sold fire insurance.**” Ignoring punctuation, compute and compare the log-likelihoods of this sentence under the unigram and bigram models:

$$\mathcal{L}_u = \log \left[P_u(\text{the}) P_u(\text{nineteen}) P_u(\text{officials}) \dots P_u(\text{sold}) P_u(\text{fire}) P_u(\text{insurance}) \right]$$

$$\mathcal{L}_b = \log \left[P_b(\text{the}|\langle s \rangle) P_b(\text{nineteen}|\text{the}) P_b(\text{officials}|\text{nineteen}) \dots P_b(\text{fire}|\text{sold}) P_b(\text{insurance}|\text{fire}) \right]$$

Which pairs of adjacent words in this sentence are not observed in the training corpus? What effect does this have on the log-likelihood from the bigram model?

- (e) Consider the so-called *mixture* model that predicts words from a weighted interpolation of the unigram and bigram models:

$$P_m(w'|w) = (1 - \lambda)P_u(w') + \lambda P_b(w'|w),$$

where $\lambda \in [0, 1]$ determines how much weight is attached to each prediction. Under this mixture model, the log-likelihood of the sentence from part (d) is given by:

$$\mathcal{L}_m = \log \left[P_m(\text{the}|\langle s \rangle) P_m(\text{nineteen}|\text{the}) P_m(\text{officials}|\text{nineteen}) \dots P_m(\text{fire}|\text{sold}) P_m(\text{insurance}|\text{fire}) \right].$$

Compute and plot the value of this log-likelihood \mathcal{L}_m as a function of the parameter $\lambda \in [0, 1]$. From your results, deduce the optimal value of λ to two significant digits.

- (f) Submit a hard copy of your source code for the previous parts of this problem.

4.4 Markov modeling

In this problem, you will construct and compare unigram and bigram models defined over the four-letter alphabet $\mathcal{A} = \{a, b, c, d\}$. Consider the following 16-token sequence \mathcal{S} :

$$\mathcal{S} = \text{"a a d d c c b b b b c c d d a a"}$$

(a) Unigram model

Let τ_ℓ denote the ℓ th token of this sequence, and let $L = 16$ denote the total sequence length. The overall likelihood of this sequence under a unigram model is given by:

$$P_U(\mathcal{S}) = \prod_{\ell=1}^L P_1(\tau_\ell),$$

where $P_1(\tau)$ is the unigram probability for the token $\tau \in \mathcal{A}$. Compute the maximum likelihood estimates of these unigram probabilities on the training sequence \mathcal{S} . Complete the table with your answers.

τ	a	b	c	d
$P_1(\tau)$				

(b) Bigram model

Suppose that the overall likelihood of the sequence \mathcal{S} under a bigram model is computed by:

$$P_B(\mathcal{S}) = P_1(\tau_1) \prod_{\ell=2}^L P_2(\tau_\ell | \tau_{\ell-1}),$$

where $P_2(\tau' | \tau)$ is the bigram probability that token $\tau \in \mathcal{A}$ is followed by token $\tau' \in \mathcal{A}$. Compute the maximum likelihood estimates of these bigram probabilities on the training sequence \mathcal{S} . Complete the table with your answers.

τ'

τ

$P_2(\tau' \tau)$	a	b	c	d
a	$\frac{2}{3}$	0	0	$\frac{1}{3}$
b				
c				
d				

(c) **Likelihoods**

Consider again the training sequence \mathcal{S} , as well as three test sequences \mathcal{T}_1 , \mathcal{T}_2 , and \mathcal{T}_3 of the same length, shown below. Note that \mathcal{T}_2 and \mathcal{T}_3 contain bigrams (underlined) that are not in the training sequence \mathcal{S} .

$$\begin{aligned}\mathcal{S} &= \text{"a a d d c c b b b b c c d d a a"} \\ \mathcal{T}_1 &= \text{"b c b c b c b c b c b c b c"} \\ \mathcal{T}_2 &= \text{"a a a a b b b b c c c c d d d d"} \\ \mathcal{T}_3 &= \text{"b a b a b a b a b a b a b a b a"}\end{aligned}$$

Consider the probabilities of these sequences under the unigram and bigram models from parts (a) and (b) of this problem (i.e., the models that you estimated from the training sequence \mathcal{S}). For each of the following, indicate whether the probability on the left is equal ($=$), greater ($>$), or less ($<$) than the probability on the right.

Note: you can (and should) answer these questions without explicitly computing the numerical values of the expressions on the left and right hand sides.

$P_U(\mathcal{S})$ $P_U(\mathcal{T}_1)$

$P_U(\mathcal{S})$ $P_U(\mathcal{T}_2)$

$P_U(\mathcal{S})$ $P_U(\mathcal{T}_3)$

$P_B(\mathcal{S})$ $P_B(\mathcal{T}_1)$

$P_B(\mathcal{S})$ $P_B(\mathcal{T}_2)$

$P_B(\mathcal{T}_2)$ $P_B(\mathcal{T}_3)$

$P_B(\mathcal{S})$ $P_U(\mathcal{S})$

$P_B(\mathcal{T}_1)$ $P_U(\mathcal{T}_1)$

$P_B(\mathcal{T}_2)$ $P_U(\mathcal{T}_2)$

$P_B(\mathcal{T}_3)$ $P_U(\mathcal{T}_3)$

(d) **Likelihoods**

Consider the model obtained by linear interpolation (or mixing) of the unigram and bigram models estimated in part (a) of this problem:

$$P_M(\tau'|\tau) = (1 - \lambda)P_1(\tau') + \lambda P_2(\tau'|\tau),$$

with mixing coefficient $\lambda \in [0, 1]$. For a sequence of tokens of length L , the mixture model computes the log-likelihood as:

$$\mathcal{L} = \log P_1(\tau_1) + \sum_{\ell=2}^L \log P_M(\tau_\ell|\tau_{\ell-1}).$$

Naturally, this value varies as a function of the coefficient λ . For λ near zero, it is close to the log-likelihood of the unigram model; for λ near one, it is close to that of the bigram model. This last part of this problem asks you to consider, for each of the sequences below, the *qualitative* behavior of the mixture model's log-likelihood as a function of $\lambda \in [0, 1]$. (For instance, is this function constant, or if not, where do its maximum and minimum occur?)

The plots below illustrate four possible behaviors of the mixture model's log-likelihood as a function of $\lambda \in [0, 1]$. For each sequence below, indicate the one plot (either A, B, C, or D) that sketches the correct qualitative behavior.

$\mathcal{S} = \text{"a a d d c c b b b b c c d d a a"}$

☐

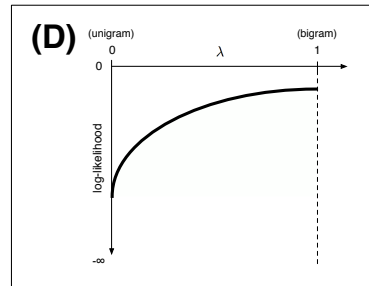
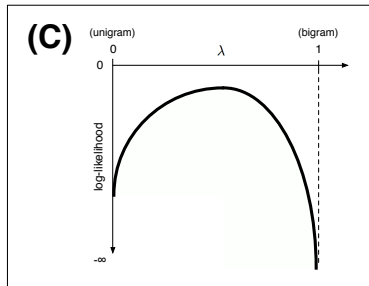
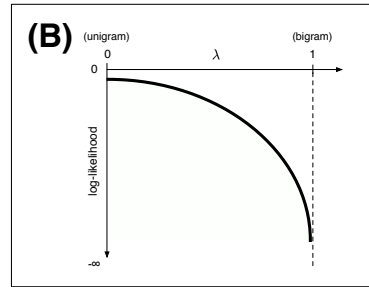
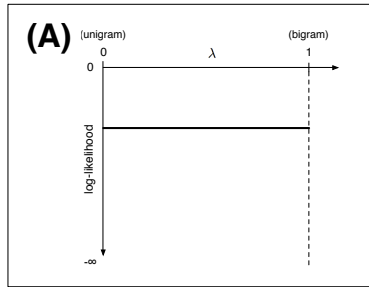
$\mathcal{T}_1 = \text{"b c b c b c b c b c b c b c b c"}$

☐

$\mathcal{T}_2 = \text{"a a a a b b b b c c c c d d d d"}$

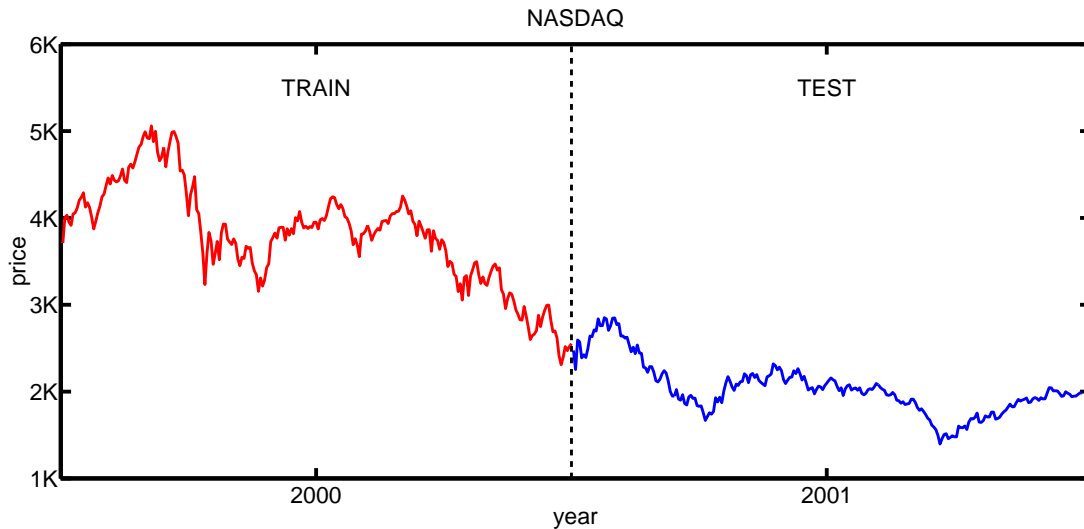
☐

$\mathcal{T}_3 = \text{"b a b a b a b a b a b a b a b a"}$

☐


4.5 Stock market prediction

In this problem, you will apply a simple linear model to predicting the stock market. From the course web site, download the files `nasdaq00.txt` and `nasdaq01.txt`, which contain the NASDAQ indices at the close of business days in 2000 and 2001.



(a) Linear coefficients

How accurately can the index on one day be predicted by a linear combination of the three preceding indices? Using only data from the year 2000, compute the linear coefficients (a_1, a_2, a_3) that maximize the log-likelihood $\mathcal{L} = \sum_t \log P(x_t | x_{t-1}, x_{t-2}, x_{t-3})$, where:

$$P(x_t | x_{t-1}, x_{t-2}, x_{t-3}) = \frac{1}{\sqrt{2\pi}} \exp \left[-\frac{1}{2} \left(x_t - a_1 x_{t-1} - a_2 x_{t-2} - a_3 x_{t-3} \right)^2 \right],$$

and the sum is over business days in the year 2000 (starting from the fourth day).

(b) Mean squared prediction error

For the coefficients estimated in part (a), compare the model's performance (in terms of mean squared error) on the data from the years 2000 and 2001. Would you recommend this linear model for stock market prediction?

(c) Source code

Turn in a print-out of your source code. You may program in the language of your choice, and you may solve the required system of linear equations either by hand or by using built-in routines (e.g., in Matlab, NumPy).
