

Assignment 8

1 Structured Prediction [25 points]

Consider a part-of-speech (PoS) tagging task, in which instances \mathbf{x} are sentences and labels \mathbf{y} are sequences of PoS tags. For example, a labeled example for this problem might be

$\mathbf{y} :$ N V D N
 $\mathbf{x} :$ John threw the ball

Given several such labeled examples as training data, the goal is to learn a prediction model which given a new sentence, can accurately predict the sequence of PoS tags for the sentence.

Suppose we have a vocabulary of M words and a set of K possible PoS tags. Given a sentence-label pair (\mathbf{x}, \mathbf{y}) , let's say that for each possible word u and each possible PoS tag k , we define a joint feature $\phi_{u,k}(\mathbf{x}, \mathbf{y})$ that counts how many times u appears in \mathbf{x} together with k in the corresponding position in \mathbf{y} . Thus if \mathbf{x} and \mathbf{y} are of length T , then denoting by x_t the t -th word in \mathbf{x} and by y_t the t -th PoS tag in \mathbf{y} , we have

$$\phi_{u,k}(\mathbf{x}, \mathbf{y}) = \sum_{t=1}^T \mathbf{1}(x_t = u, y_t = k).$$

Similarly, for each possible pair of words (u, u') and each possible pair of PoS tags (k, k') , we define a joint feature $\phi_{u,u',k,k'}(\mathbf{x}, \mathbf{y})$ that counts how many times (u, u') appears in consecutive positions in \mathbf{x} together with (k, k') in the corresponding positions in \mathbf{y} . Thus

$$\phi_{u,u',k,k'}(\mathbf{x}, \mathbf{y}) = \sum_{t=1}^{T-1} \mathbf{1}(x_t = u, x_{t+1} = u', y_t = k, y_{t+1} = k').$$

Thus there are a total of $d = MK + M^2K^2$ features (note that different sentence-label pairs in the training set can be of different lengths T ; any sentence-label pair is represented by the same number of features). For any given sentence-label pair (\mathbf{x}, \mathbf{y}) , only a small number of these features are non-zero. We denote by $\phi(\mathbf{x}, \mathbf{y})$ the (sparse) d -dimensional joint feature vector derived from (\mathbf{x}, \mathbf{y}) as above.

Recall that, given a joint feature representation as above, both conditional random fields (CRFs) and StructSVM learn a weight vector $\mathbf{w} \in \mathbb{R}^d$, and given a new sentence \mathbf{x} of length T , predict a label $\hat{\mathbf{y}}$ (a PoS tag sequence of length T) as follows:

$$\hat{\mathbf{y}} \in \arg \max_{\mathbf{y}} \text{score}_{\mathbf{w}}(\mathbf{x}, \mathbf{y})$$

where

$$\begin{aligned} \text{score}_{\mathbf{w}}(\mathbf{x}, \mathbf{y}) &= \mathbf{w}^\top \phi(\mathbf{x}, \mathbf{y}) \\ &= \sum_{u,k} w_{u,k} \cdot \phi_{u,k}(\mathbf{x}, \mathbf{y}) + \sum_{u,u',k,k'} w_{u,u',k,k'} \cdot \phi_{u,u',k,k'}(\mathbf{x}, \mathbf{y}). \end{aligned}$$

1. (10 points) Suppose you have learned a weight vector \mathbf{w} as above (using either a CRF or StructSVM learning algorithm). You are given a new sentence $\mathbf{x} = \text{"The dog chased the cat"}$, and are considering two possible labelings for it:

$$\begin{array}{ccccc} \mathbf{y}^1 : & \text{N} & \text{N} & \text{V} & \text{D} & \text{N} \\ \mathbf{x} : & \text{The} & \text{dog} & \text{chased} & \text{the} & \text{cat} \end{array} \quad \begin{array}{ccccc} \mathbf{y}^2 : & \text{D} & \text{N} & \text{V} & \text{D} & \text{N} \\ \mathbf{x} : & \text{The} & \text{dog} & \text{chased} & \text{the} & \text{cat} \end{array}$$

Thus the feature representations $\phi(\mathbf{x}, \mathbf{y}^1)$ and $\phi(\mathbf{x}, \mathbf{y}^2)$ have the following non-zero components (verify that you understand this):

$$\begin{array}{ll} \phi_{\text{the}, \text{N}}(\mathbf{x}, \mathbf{y}^1) = 1 & \phi_{\text{the}, \text{dog}, \text{N}, \text{N}}(\mathbf{x}, \mathbf{y}^1) = 1 \\ \phi_{\text{dog}, \text{N}}(\mathbf{x}, \mathbf{y}^1) = 1 & \phi_{\text{dog}, \text{chased}, \text{N}, \text{V}}(\mathbf{x}, \mathbf{y}^1) = 1 \\ \phi_{\text{chased}, \text{V}}(\mathbf{x}, \mathbf{y}^1) = 1 & \phi_{\text{chased}, \text{the}, \text{V}, \text{D}}(\mathbf{x}, \mathbf{y}^1) = 1 \\ \phi_{\text{the}, \text{D}}(\mathbf{x}, \mathbf{y}^1) = 1 & \phi_{\text{the}, \text{cat}, \text{D}, \text{N}}(\mathbf{x}, \mathbf{y}^1) = 1 \\ \phi_{\text{cat}, \text{N}}(\mathbf{x}, \mathbf{y}^1) = 1 & \\ \\ \phi_{\text{the}, \text{D}}(\mathbf{x}, \mathbf{y}^2) = 2 & \phi_{\text{the}, \text{dog}, \text{D}, \text{N}}(\mathbf{x}, \mathbf{y}^2) = 1 \\ \phi_{\text{dog}, \text{N}}(\mathbf{x}, \mathbf{y}^2) = 1 & \phi_{\text{dog}, \text{chased}, \text{N}, \text{V}}(\mathbf{x}, \mathbf{y}^2) = 1 \\ \phi_{\text{chased}, \text{V}}(\mathbf{x}, \mathbf{y}^2) = 1 & \phi_{\text{chased}, \text{the}, \text{V}, \text{D}}(\mathbf{x}, \mathbf{y}^2) = 1 \\ \phi_{\text{cat}, \text{N}}(\mathbf{x}, \mathbf{y}^2) = 1 & \phi_{\text{the}, \text{cat}, \text{D}, \text{N}}(\mathbf{x}, \mathbf{y}^2) = 1 \end{array}$$

Suppose the relevant components of the weight vector \mathbf{w} you have learned are as follows:

$$\begin{array}{ll} w_{\text{the}, \text{N}} = -6 & w_{\text{the}, \text{dog}, \text{N}, \text{N}} = -2 \\ w_{\text{the}, \text{D}} = 7 & w_{\text{the}, \text{dog}, \text{D}, \text{N}} = 4 \\ w_{\text{dog}, \text{N}} = 3 & w_{\text{dog}, \text{chased}, \text{N}, \text{V}} = 0 \\ w_{\text{chased}, \text{V}} = 2 & w_{\text{chased}, \text{the}, \text{V}, \text{D}} = 3 \\ w_{\text{cat}, \text{N}} = 4 & w_{\text{the}, \text{cat}, \text{D}, \text{N}} = 4 \end{array}$$

Find $\text{score}_{\mathbf{w}}(\mathbf{x}, \mathbf{y}^1)$ and $\text{score}_{\mathbf{w}}(\mathbf{x}, \mathbf{y}^2)$. Show your calculations. Which of the two labelings receives a higher score under \mathbf{w} ?

2. (15 points) Again, suppose you have learned a weight vector \mathbf{w} as above (using either a CRF or StructSVM learning algorithm). In general, given a new sentence \mathbf{x} of length T , there are K^T possible labelings \mathbf{y} , and a brute-force search over all of them to find a labeling $\hat{\mathbf{y}}$ with maximal score under \mathbf{w} is expensive. Give an efficient algorithm to find a highest-scoring label sequence $\hat{\mathbf{y}}$ under \mathbf{w} . Your algorithm should need only $O(K^2T)$ computations.

(Hint: Use a dynamic programming approach similar to the Viterbi algorithm for finding the most likely hidden state sequence in an HMM. The inputs to your problem are a sentence \mathbf{x} of length T and a weight vector $\mathbf{w} \in \mathbb{R}^d$ (with components indexed by $w_{u,k}$ and $w_{u,u',k,k'}$ as described above); the output is a predicted PoS tag sequence $\hat{\mathbf{y}}$, also of length T , with maximal score under \mathbf{w} .)

2 Semi-Supervised Learning [35 points]

In this problem you will consider a semi-supervised extension of Naïve Bayes that incorporates unlabeled data via EM. In particular, you will simulate one step of the EM algorithm on a small data set.

For simplicity, consider a binary classification task with 2-dimensional Boolean instances $\mathbf{x} \in \{0,1\}^2$ and labels $y \in \{\pm 1\}$. Recall that the Naïve Bayes classifier assumes a generative probabilistic model of the form

$$p(\mathbf{x}, y) = p(y) \prod_{j=1}^2 p(x_j | y).$$

In our setting, there are 5 parameters, which we will collectively denote as θ :

$$\begin{aligned}\theta_{+1} &= \mathbf{P}(Y = +1); \\ \theta_{j|k} &= \mathbf{P}(X_j = 1 | Y = k), \quad \text{for each feature } j \in \{1, 2\} \text{ and each label } k \in \{\pm 1\}.\end{aligned}$$

Thus the probability of a labeled example $\mathbf{x} = (1, 0), y = -1$ under the above model would be

$$\mathbf{P}(Y = -1) \mathbf{P}(X_1 = 1 | Y = -1) \mathbf{P}(X_2 = 0 | Y = -1) = (1 - \theta_{+1})\theta_{1|-1}(1 - \theta_{2|-1}).$$

Standard (Supervised) Naïve Bayes. In the supervised setting, given labeled training data $S_L = ((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_{m_L}, y_{m_L}))$, one computes maximum likelihood estimates as follows:

$$\begin{aligned}\hat{\theta}_{+1} &= \frac{1}{m_L} \sum_{i=1}^{m_L} \mathbf{1}(y_i = +1) \\ \hat{\theta}_{j|k} &= \frac{\sum_{i=1}^{m_L} \mathbf{1}(y_i = k, x_{ij} = 1)}{\sum_{i=1}^{m_L} \mathbf{1}(y_i = k)}\end{aligned}$$

Given a new instance \mathbf{x} , one then uses the estimated parameters together with Bayes rule to compute the probability $P(Y = +1 | \mathbf{x})$ under the learned model, and classifies the instance as $+1$ if this probability is greater than $\frac{1}{2}$ and as -1 otherwise.

Semi-Supervised Naïve Bayes. In the semi-supervised setting, given labeled training data S_L as above and additional unlabeled data $S_U = (\mathbf{x}_{m_L+1}, \dots, \mathbf{x}_{m_L+m_U})$, one treats the missing labels in S_U as unobserved variables and uses EM to find maximum likelihood parameter estimates. In particular, one starts with initial parameter estimates $\hat{\theta}^0$ learned as above from the labeled data S_L alone, and then iteratively estimates posterior distributions on the missing labels of the unlabeled data points (E-step) and updates the parameter estimates via a weighted maximum likelihood estimation (M-step). Specifically, on each round t , in the E-step, for each unlabeled example \mathbf{x}_i one computes the posterior distribution over labels, $q^t(k | \mathbf{x}_i) = P(Y_i = k | \mathbf{x}_i; \hat{\theta}^t)$, under the current parameter estimates $\hat{\theta}^t$; in the M-step, one then updates the parameter estimates as follows:

$$\begin{aligned}\hat{\theta}_{+1}^{t+1} &= \frac{1}{m_L + m_U} \left(\sum_{i=1}^{m_L} \mathbf{1}(y_i = +1) + \sum_{i=m_L+1}^{m_L+m_U} q^t(+1 | \mathbf{x}_i) \right) \\ \hat{\theta}_{j|k}^{t+1} &= \frac{\sum_{i=1}^{m_L} \mathbf{1}(y_i = k, x_{ij} = 1) + \sum_{i=m_L+1}^{m_L+m_U} q^t(k | \mathbf{x}_i) \cdot \mathbf{1}(x_{ij} = 1)}{\sum_{i=1}^{m_L} \mathbf{1}(y_i = k) + \sum_{i=m_L+1}^{m_L+m_U} q^t(k | \mathbf{x}_i)}\end{aligned}$$

On convergence, one uses the final parameter estimates to make predictions on new instances in the same manner as before.

Problem. Suppose you have a training sample of 8 labeled examples and 4 unlabeled examples, distributed as follows (note that since there are only 4 possible instances in our simple setup and 2 possible labels, some examples are repeated in the training sample below; this would be unlikely in real data, but the key ideas present in this example would carry over to real data as well):

Labeled data S_L	Unlabeled data S_U
$\mathbf{x} = (x_1, x_2), y$	$\mathbf{x} = (x_1, x_2)$
(1, 1), +1	(1, 1)
(1, 1), +1	(1, 1)
(1, 0), +1	(0, 0)
(0, 0), +1	(0, 0)
(1, 0), -1	
(0, 1), -1	
(0, 0), -1	
(0, 0), -1	
8	4

- (5 points) Calculate the initial maximum likelihood parameter estimates based on the labeled data only: $\hat{\boldsymbol{\theta}}^0 = (\hat{\theta}_{+1}^0, \hat{\theta}_{1|+1}^0, \hat{\theta}_{2|+1}^0, \hat{\theta}_{1|-1}^0, \hat{\theta}_{2|-1}^0)$.
- (8 points) For each instance $\mathbf{x} = (x_1, x_2)$ that appears in the unlabeled data, compute the posterior distribution over the label under the parameter estimates computed in the first part above. In particular, compute each of the following:

$$\begin{aligned}
q^0(+1 | \mathbf{x} = (1, 1)) &= P(Y = +1 | \mathbf{x} = (1, 1); \hat{\boldsymbol{\theta}}^0) \\
q^0(+1 | \mathbf{x} = (0, 0)) &= P(Y = +1 | \mathbf{x} = (0, 0); \hat{\boldsymbol{\theta}}^0)
\end{aligned}$$

Show your calculations.

(Hint: Use Bayes' rule.)

- (10 points) Using the results of the first two parts above, find the updated parameter estimates after one step of EM: $\hat{\boldsymbol{\theta}}^1 = (\hat{\theta}_{+1}^1, \hat{\theta}_{1|+1}^1, \hat{\theta}_{2|+1}^1, \hat{\theta}_{1|-1}^1, \hat{\theta}_{2|-1}^1)$. Show your calculations.
- (6 points) The log-likelihood of the labeled and unlabeled data $S = (S_L, S_U)$ under parameter estimates $\hat{\boldsymbol{\theta}}^t$ is given by

$$\ln p(S; \hat{\boldsymbol{\theta}}^t) = \sum_{i=1}^{m_L} \ln p(\mathbf{x}_i, y_i; \hat{\boldsymbol{\theta}}^t) + \sum_{i=m_L+1}^{m_L+m_U} \ln \left(\underbrace{\sum_{y_i \in \{\pm 1\}} p(\mathbf{x}_i, y_i; \hat{\boldsymbol{\theta}}^t)}_{\ln p(\mathbf{x}_i; \hat{\boldsymbol{\theta}}^t)} \right).$$

Write an expression for the log-likelihood of the given training data as a function of the 5 parameters $\hat{\theta}_{+1}^t, \hat{\theta}_{1|+1}^t, \hat{\theta}_{2|+1}^t, \hat{\theta}_{1|-1}^t, \hat{\theta}_{2|-1}^t$ for any t .

- (6 points) Using the expression derived in the fourth part above, calculate the log-likelihood of the given training data under both the initial parameter estimates $\hat{\boldsymbol{\theta}}^0$ computed in the first part above and the updated parameter estimates $\hat{\boldsymbol{\theta}}^1$ computed in the third part above. After one step of EM, has the log-likelihood of the data increased or decreased?

(Hint: You might like to write a small script which given the 5 parameters as input, returns the log-likelihood value, according to the expression you derived in the fourth part above, as output. You can then plug in the values of the parameter estimates obtained above for $t = 0$ and $t = 1$ into your script to obtain the corresponding log-likelihoods.)