

HW4: Word Segmentation

Virgile Audi
vaudi@g.harvard.edu
Nicolas Drizard
nicolasdrizard@g.harvard.edu

April 16, 2016

1 Introduction

The goal of this assignment is to tackle the NLP task of identifying and labeling contiguous segments of text. We will use sequence models and a dynamic programming method to find the best scoring sequence.

2 Problem Description

The idea is here to label continuous sequence of words with BIO tagging of different entities. The entities are the following:

1. PER: a person
2. LOC: a location
3. ORG: an organization
4. MISC:

Furthermore, this tagging method identifies the continuous group of words belonging to the same entity: the prefix B stop the current tag and begins a new one whereas the prefix I continues adding to the previous tag. However, in our solution we just cared about predicting the entity tag and then we were grouping the contiguous predictions into the same entity because the training text does not contain any B-tag.

3 Model and Algorithms

We used three different methods to solve this problem. The first two are the equivalent of first the Naive Bayes and second the logistic regression from text classification tasks. The last one introduces a customized way to train a neural architecture for this task.

3.1 Hidden Markov Model

We implement here a standard first order hidden Markov Model. The hidden states are the tags and the observed states are the features we built (word counts, capitalization...). The model can be represented with the following graphical model and requires two distributions: emission and transition.

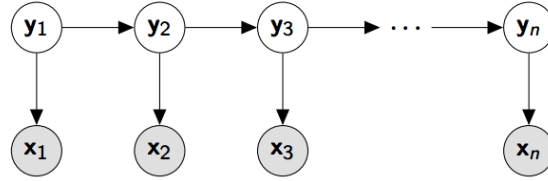


Figure 1: Graphical model of 1st order HMM with one feature

We represent the two distributions with multinomial as they model feature counts. As a result, we can infer them simply with the maximum likelihood estimator:

$$p(x_i = \delta(f) | y_i = \delta(c)) = \frac{F_{f,c}}{F_{.,c}}$$

$$p(y_i = \delta(c_i) | y_{i-1} = \delta(c_{i-1})) = \frac{T_{c_{i-1},c_i}}{T_{c_{i-1},.}}$$

with T_{c_{i-1},c_i} the counts of class c_{i-1} preceding class c_i and $F_{f,c}$ the counts of emission f with class c .

If we consider multiple features, then we still assume that the features are independent with each other (it's the main assumption in the Naive Bayes approach also). Only the emission distribution is changed and we can combine the probability together:

$$p(x_i = (\delta(f_1), \delta(f_2)) | y_i = \delta(c)) = p(x_i = \delta(f_1) | y_i = \delta(c)) p(x_i = \delta(f_2) | y_i = \delta(c)) = \frac{F_{f_1,c}}{F_{.,c}} \frac{F_{f_2,c}}{F_{.,c}}$$

3.2 Maximum-Entropy Markov Model

Next, we implemented a Maximum-Entropy Markov Model. The objective of the MEMM is to evaluate at each time step a distribution over the possible tags using features of the current word, denoted as $feat(x_i)$ and the tag of the previous word, c_{i-1} , using multi-class logistic regression, i.e.

$$p(y_i | y_{i-1}, feat(x_i)) = \text{softmax}([feat(x_i), c_{i-1}] \mathbf{W} + \mathbf{b})$$

3.3 Structured Perceptron

3.4 Viterbi algorithm

The search algorithm that we implemented is the dynamic programming algorithm named after Andrew Viterbi. Its main difference with a greedy approach is that it evaluates at every step and for every previous state, the best possible next step. This guarantees a solution closer to the true optimal solution. The pseudo-code of the algorithm is given by:

```
procedure VITERBIWITHBP
   $\pi \in \mathbb{R}^{n+1 \times \mathcal{C}}$  initialized to  $-\infty$ 
   $bp \in \mathcal{C}^{n \times \mathcal{C}}$  initialized to  $\epsilon$ 
   $\pi[0, \langle s \rangle] = 0$ 
  for  $i = 1$  to  $n$  do
    for  $c_{i-1} \in \mathcal{C}$  do
      compute  $\hat{y}(c_{i-1})$ 
      for  $c_i \in \mathcal{C}$  do
         $score = \pi[i-1, c_{i-1}] + \log \hat{y}(c_{i-1})_{c_i}$ 
        if  $score > \pi[i, c_i]$  then
           $\pi[i, c_i] = score$ 
           $bp[i, c_i] = c_{i-1}$ 
  return sequence from  $bp$ 
```

4 Experiments

4.1 Feature Engineering

The original paper suggests several features to use. We focus on the word counts and a capitalization feature. We defined our capitalization feature as follow:

- 1 : word in low caps;
- 2 : whole word in caps;
- 3 : first letter in cap;
- 4 : one cap in the word;
- 5 : other

We then produced an embedding of the word counts using a pre-trained version.

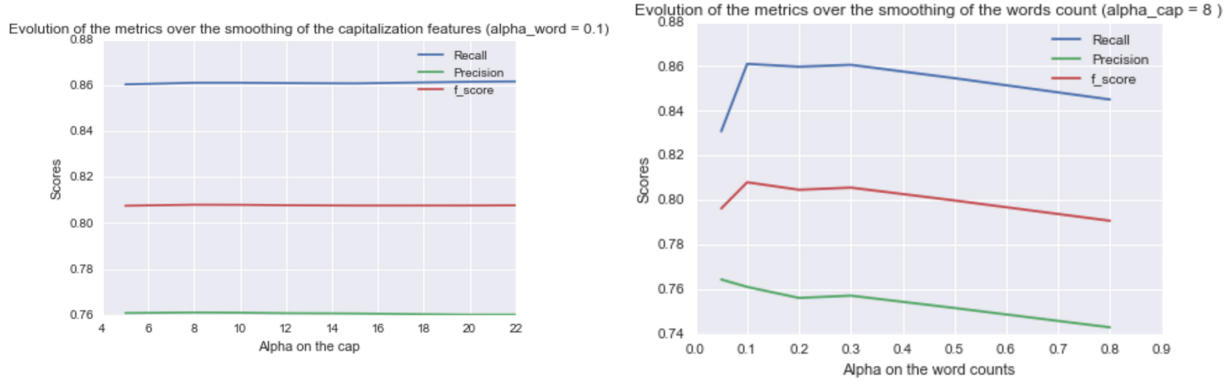
4.2 Model Evaluation

As used in the Kaggle competition, we used the f-score with the precision and recall measure to evaluate our model while tuning the hyperparameters. A positive prediction stands for a label (in the notation of the task, everything which is not the **O** tag):

1. recall: ratio of the true positive predictions among the positives tags in the correct sequence
2. precision: ratio of the true positive predictions among the positive predictions,
3. f-score (with $\beta = 1$): harmonic mean of the precision and the recall, i.e. $f_1 = \frac{2pr}{p+r}$

4.3 Hidden Markov Model

There is only the smoothing parameter α and eventually feature selection here to tune here. We evaluate the impact of adding more features and run experiments with different alpha values to tune them. One important details is to make sure to use a specific smoothing parameter for each distribution, i.e a smoothing parameter may be applied to the transition matrix but also to the emission matrix of each different feature. Each of this distribution has a different tail and need a different smoothing. For instance, the transition matrix need a very small α (around 0.1) because we are pretty confident in it but the capitalizations feature need one much bigger (around 20) because the counts are already high.



We notice that the model is less sensitive to the changes of the smoothing parameter on the capitalization feature as on the word counts. This is pretty reasonable as the feature counts are much higher in the capitalization feature than in the word counts. Tuning this parameter provides a model with a f-score of **0.808**. Using only the word counts features provide a best f-score of **0.764**.

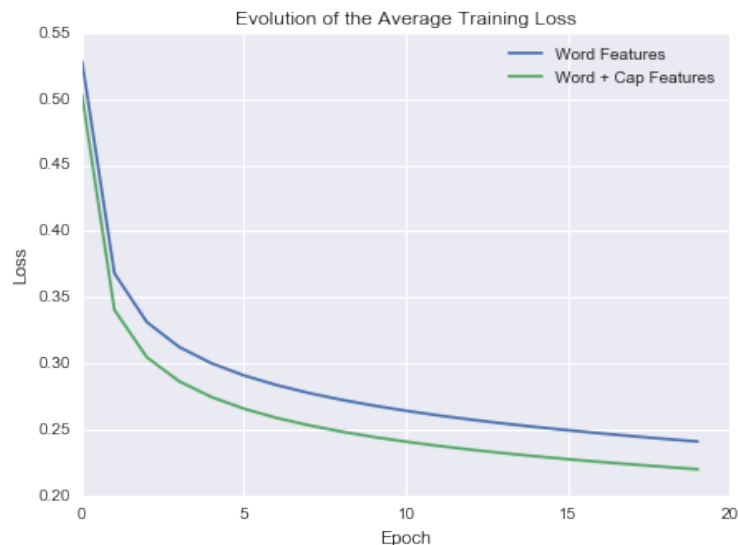
We obtained a Kaggle score on the test set of :

$$K_{HMM} = 0.48365$$

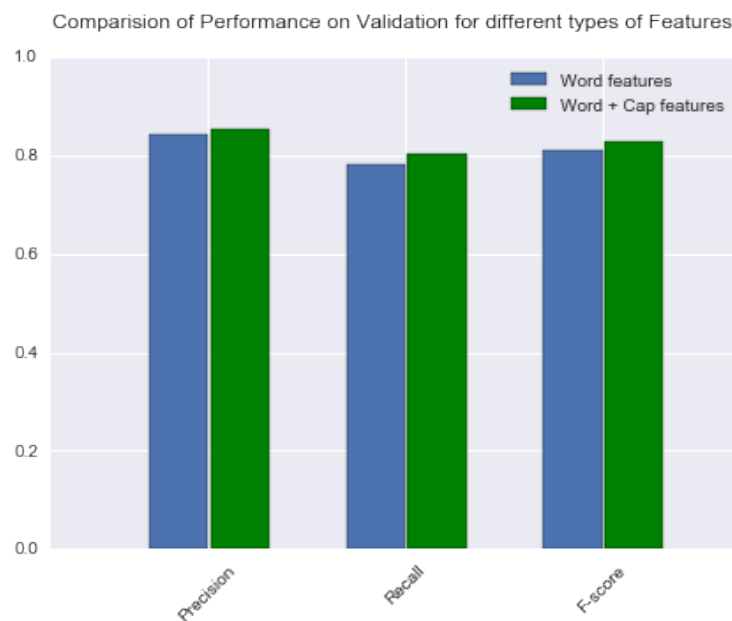
4.4 Maximum-Entropy Markov Model

We coded the MEMM using the nn module and trained using stochastic gradient descent. We also used the Glove embeddings using a lookup table. As for the HMM, we used two different sets of features, i.e. the words and the words and capitalisation of the words. We observed that the training algorithm converges quite rapidly, and that if adding caps to the features helped decrease the loss, the impact was not as strong as expected. Nevertheless, we trained the model

on 20 epochs in order to learn the embeddings for the $\langle s \rangle$ and $\langle \backslash s \rangle$ "words" added during pre-processing.



We evaluated the performance of these two models using the f-score presented above:



Adding caps yielded better results on both Precision and Recall and therefore on the f-score. But as we expected from the small differences in loss, we did not observe an important increase on the f-score. These results were later confirmed on the test set, as the kaggle score obtained for these two models were:

$$K_{nocaps} = 0.52057 \quad \text{and} \quad K_{caps} = 0.55482$$

which are both slightly better than the results of the HMM.

4.5 Structured Perceptron

5 Conclusion

This segmentation task gave us the opportunity to implement different recurrent neural network architectures but also to compare them with more traditional method. Whereas the count based and even the simple neural network models are pretty fast to train they still provide interesting results. The results provided by the three variants of RNN were interesting to illustrate the influence of gates and memory in such networks. The gated recurrent network ended as the best model on this task. One future work could be to stack more layers to our recurrent architecture or to implement a network with a dynamic memory part to give more flexibility in how the model uses the information it already processed.

Appendices

Preprocessing:

```
1  #!/usr/bin/env python
2
3  """NER Preprocessing
4  """
5
6  import numpy as np
7  import h5py
8  import argparse
9  import sys
10 import re
11 import codecs
12
13 # Your preprocessing, features construction, and word2vec code.
14
15
16 FILE_PATHS = {"CONLL": ("data/train.num.txt",
17                          "data/dev.num.txt",
18                          "data/test.num.txt",
19                          "data/tags.txt")}
20 args = {}
21
22
23 def main(arguments):
24     global args
25     parser = argparse.ArgumentParser()
```

```

26         description=__doc__ ,
27         formatter_class=argparse.RawDescriptionHelpFormatter)
28     parser.add_argument('dataset', help="Data set",
29                         type=str)
30     args = parser.parse_args(arguments)
31     dataset = args.dataset
32     train, valid, test, tag_dict = FILE_PATHS[dataset]
33
34     filename = args.dataset + '.hdf5'
35     with h5py.File(filename, "w") as f:
36         f['train_input'] = train_input
37         f['train_output'] = train_output
38         if valid:
39             f['valid_input'] = valid_input
40             f['valid_output'] = valid_output
41         if test:
42             f['test_input'] = test_input
43         f['nfeatures'] = np.array([V], dtype=np.int32)
44         f['nclasses'] = np.array([C], dtype=np.int32)
45
46
47 if __name__ == '__main__':
48     sys.exit(main(sys.argv[1:]))

```

Hidden Markov Model:

```

1  — Documentation:
2  — ——— How to call it from the command line?
3  — For example:
4  — $ th count_based.lua -N 5
5  — Other argument possible (see below)
6  —
7  — ——— Is there an Output?
8  — By default, the predictions on the test set are saved in hdf5 format
   as classifier .. opt.f
9
10 — Only requirements allowed
11 require("hdf5")
12 require 'helper.lua';
13
14 cmd = torch.CmdLine()
15
16 — Cmd Args
17 cmd:option('-datafile', 'data/words_feature.hdf5',
18            'Datafile with features in hdf5 format')
19 cmd:option('-alpha-t', 0.1, 'Smoothing parameter alpha in the

```

```

    transition counts')
20 cmd:option('-alpha_w', 2, 'Smoothing parameter alpha in the word counts
    ')
21 cmd:option('-alpha_c', 20, 'Smoothing parameter alpha in the caps
    counts')
22 cmd:option('-test', 0, 'Boolean (as int) to ask for a prediction on
    test, will be saved in submission in hdf5 format')
23 cmd:option('-datafile_test', 'submission/v_seq_hmm', 'Smoothing
    parameter alpha in the word counts')
24 cmd:option('-nfeatures', 2, 'Number of type of features to use')
25 cmd:option('-cv', 0, 'Boolean (as int) to run a cross validation
    pipeline')
26
27
28
29 — Formating as log-probability and smoothing the input
30 function format_matrix(matrix, alpha)
31     local formatted_matrix = matrix:clone():type('torch.DoubleTensor')
32     formatted_matrix:add(alpha)
33     — Normalize
34     local norm_mat = torch.expandAs(formatted_matrix:sum(1),
        formatted_matrix)
35     formatted_matrix:cdiv(norm_mat)
36     return formatted_matrix:log()
37 end
38
39 — log-scores of transition and emission
40 — corresponds to the vector y in the lecture notes
41 — i: timestep for the computed score
42 function score_hmm(observations, i, emissions, transition, C, nfeatures
    )
43     local observation_emission = torch.zeros(C)
44     for k=1,nfeatures do
45         observation_emission:add(emissions[k][observations[{i,k}]]))
46     end
47     observation_emission = observation_emission:view(C, 1):expand(C, C)
48     — NOTE: allocates a new Tensor
49     return observation_emission + transition
50 end
51
52 — Viterbi algorithm.
53 — observations: a sequence of observations, represented as integers
54 — logscore: the edge scoring function over classes and observations in
    a history-based model
55 function viterbi(observations, logscore, emissions, transition,

```



```

nfeatures)
56   local y
57   — Formating tensors
58   local initial = torch.zeros(transition:size(2), 1)
59   — initial started with a start of sentence: <t>
60   initial[{8,1}] = 1
61   initial:log()
62
63   — number of classes
64   C = initial:size(1)
65   local n = observations:size(1)
66   local max_table = torch.Tensor(n, C)
67   local backpointer_table = torch.Tensor(n, C)
68
69   — first timestep
70   — the initial most likely paths are the initial state distribution
71   — NOTE: another unnecessary Tensor allocation here
72   local init_pred = initial:clone()
73   for i=1,nfeatures do
74       init_pred:add(emissions[i][observations[{1,i}]]))
75   end
76   local maxes, backpointers = init_pred:max(2)
77   max_table[1] = maxes
78
79   — remaining timesteps ("forwarding" the maxes)
80   for i=2,n do
81       — precompute edge scores
82       y = logscore(observations, i, emissions, transition, C,
83                   nfeatures)
84       scores = y + maxes:view(1, C):expand(C, C)
85       — compute new maxes (NOTE: another unnecessary Tensor
86       allocation here)
87       maxes, backpointers = scores:max(2)
88       — record
89       max_table[i] = maxes
90       backpointer_table[i] = backpointers
91   end
92   — follow backpointers to recover max path
93   local classes = torch.Tensor(n)
94   maxes, classes[n] = maxes:max(1)
95   for i=n,2,-1 do
96       classes[i-1] = backpointer_table[{i, classes[i]}]
97   end

```

```

98
99     return classes
100 end
101
102 — Prediction pipeline
103 function predict(observations, emissions, transition, alphas, nfeatures
104 )
105     — Formating model parameters (log and alpha smoothing)
106     — Alphas is a tensor : {alpha_t, alpha_w, alpha_c}
107     emissions_cleaned = {}
108     for i=1,nfeatures do
109         emissions_cleaned[i] = format_matrix(emissions[i], alphas[i+1])
110     end
111     local transition_cleaned = format_matrix(transition, alphas[1])
112     return viterbi(observations, score_hmm, emissions_cleaned,
113         transition_cleaned, nfeatures)
114 end
115 — Cross validation pipeline
116 function cross_validation(observations, emissions, transitions,
117     true_classes,
118     alphas_table, alpha_t)
119     — alphas_table is a table of tensor with the range of parameters
120     to use
121     — Current implementation for 2 features only
122     — alphas_table = {alpha_w_tensor, alpha_c_tensor}
123     — Return a tensor with first columns the alpha value and last the
124     score for each
125     local nfeatures = #alphas_table
126     local v_seq_dev, precision, recall, f
127     local alphas = torch.DoubleTensor(3)
128     local size1 = alphas_table[1]:size(1)
129     local size2 = alphas_table[2]:size(1)
130     local num_evaluations = size1*size2
131
132     — Columns for 2 features are (alphas_w_value, alphas_c_value,
133     f_score, precision, recall)
134     local scores = torch.DoubleTensor(num_evaluations, nfeatures+3)
135
136     for i=1,size1 do
137         alpha_w = alphas_table[1][i]
138         for k=1,size2 do
139             alpha_c = alphas_table[2][k]

```

```

137         alphas:copy(torch.Tensor({alpha_t, alpha_w, alpha_c}))
138         v_seq_dev = predict(observations, emissions, transition,
139                             alphas, nfeatures)
140         precision, recall = compute_score(v_seq_dev, true_classes)
141         f = f_score(precision, recall)
142
143         — Filling the scores tensor
144         scores[{(i-1)*size2+k, 1}] = alpha_w
145         scores[{(i-1)*size2+k, 2}] = alpha_c
146         scores[{(i-1)*size2+k, 3}] = f
147         scores[{(i-1)*size2+k, 4}] = precision
148         scores[{(i-1)*size2+k, 5}] = recall
149     end
150 end
151 return scores
152 end
153
154 function main()
155     — Parse input params
156     opt = cmd:parse(arg)
157
158     — Reading file from pre-processing
159     myFile = hdf5.open(opt.datafile, 'r')
160     data = myFile:all()
161     emission_w = data['emission_w']
162     emission_c = data['emission_c']
163     — Table of emission tensor (one tensor per feature)
164     emissions = {emission_w, emission_c}
165     — Assertion on number of features
166     nfeatures = opt.nfeatures
167     if nfeatures > #emissions then
168         error('Too many features specified')
169     end
170     print('Number of features used: '..nfeatures)
171     transition = data['transition']
172     input_matrix_train = data['input_matrix_train']
173     input_matrix_dev = data['input_matrix_dev']
174     input_matrix_test = data['input_matrix_test']
175     myFile:close()
176
177     — Parameters:
178     true_classes = input_matrix_dev:narrow(2,5,1):clone():view(
179         input_matrix_dev:size(1))

```

```

180 — contain in each column feature observation
181 — (same order as the feature emission tensor in the emissoins
    table)
182 observations = input_matrix_dev:narrow(2,3,nfeatures):clone()
183 — Alpha parameter
184 alphas = torch.Tensor({opt.alpha_t, opt.alpha_w, opt.alpha_c})
185
186 — Prediction on dev
187 v_seq_dev = predict(observations, emissions, transition, alphas,
    nfeatures)
188 precision, recall = compute_score(v_seq_dev, true_classes)
189 f = f_score(precision, recall)
190
191 print('Prediction on dev')
192 print('Precision is : '..precision)
193 print('Recall is : '..recall)
194 print('F score (beta = 1) is : '..f)
195
196 — Cross validation
197 if (opt.cv == 1) then
198     alphas_table = {}
199     — alpha_w
200     alphas_table[1] = torch.Tensor({0.05, 0.1, 0.2, 0.3, 0.5, 0.8})
201     — alpha_c
202     alphas_table[2] = torch.Tensor({5, 8, 10, 12, 15, 20, 22})
203
204     scores = cross_validation(observations, emissions, transitions,
        true_classes,
205                               alphas_table, opt.alpha_t)
206     print(scores)
207
208     — Saving the score
209     myFile = hdf5.open('plot_scores.hdf5', 'w')
210     myFile:write('scores', scores)
211     myFile:close()
212     print('CV on dev saved in '..'plot_scores.hdf5')
213 end
214
215 — Prediction on test
216 if (opt.test == 1) then
217     print('Prediction on test')
218     observations_test = input_matrix_test:narrow(2,3,nfeatures):
        clone()
219     v_seq_test = predict(observations_test, emissions, transition,
        alphas, nfeatures)

```

```

220      — Saving predicted sequence on test
221      myFile = hdf5.open(opt.datafile_test , 'w')
222      myFile:write('v_seq_test', v_seq_test)
223      myFile:write('v_seq_dev', v_seq_dev)
224      myFile:close()
225      print('Sequence predicted on test saved in '..opt.datafile_test
226          )
227
228  end
229
230  main()

```

Helper:

```

1  — function to evaluate the predicted sequence
2  — need to compute precision and recall (class 1 stands for negative
   class)
3  function compute_score(predicted_classes , true_classes)
4      local n = predicted_classes:size(1)
5      local right_pred = 0
6      local positive_true = 0
7      local positive_pred = 0
8      for i=1,n do
9          if predicted_classes[i] > 1 then
10             positive_pred = positive_pred + 1
11         end
12         if true_classes[i] > 1 then
13             positive_true = positive_true + 1
14         end
15         if (true_classes[i] == predicted_classes[i]) and true_classes[i]
16             > 1 then
17             right_pred = right_pred + 1
18         end
19         — Verbose
20         — print('positive_true: '.. positive_true)
21         — print('positive_pred: '.. positive_pred)
22         — print('right_pred: '..right_pred)
23         local precision = right_pred/positive_pred
24         local recall = right_pred/positive_true
25         return precision , recall
26     end
27
28     function f_score(precision , recall)
29         return 2*precision*recall/(precision+recall)

```

30 end