# Homework 2, Intro to NLP, 2017

**This is due at 11pm on Tuesday, October 10. Please see detailed submission instructions below. 100 points total.**

### *How to do this problem set:*

- What version of Python should I use? 2.7
- Most of these questions require writing Python code and computing results, and the rest of them have textual answers. To generate the answers, you will have to fill out supporting files, `vit_starter.py`,`classperc.py` and `structperc.py`.
- Write all the answers in this ipython notebook. Once you are finished (1) Generate a PDF via (File -> Download As -> PDF) and upload to Gradescope (2)Turn in `vit_starter.py`, `classperc.py`, `structperc.py` and `hw_2.ipynb` on Moodle.
- **Important:** Check your PDF before you turn it in to gradescope to make sure it exported correctly. If ipython notebook gets confused about your syntax it will sometimes terminate the PDF creation routine early. You are responsible for checking for these errors. If your whole PDF does not print, try running `$jupyter nbconvert --to pdf hw_1.ipynb` to identify and fix any syntax errors that might be causing problems.
- **Important:** When creating your final version of the PDF to hand in, please do a fresh restart and execute every cell in order. Then you'll be sure it's actually right. One convenient way to do this is by clicking `Cell -> Run All` in the notebook menu.

```
In [1]:  "%autoreload"

Out[1]:  '%autoreload'
```

### *Academic honesty*

- We will audit the Moodle code from a few dozen students, chosen at random. The audits will check that the code you wrote and turned on Moodle generates the answers you turn in on your Gradescope PDF. If you turn in correct answers on your PDF without code that actually generates those answers, we will consider this a potential case of cheating. See the course page for honesty policies.
- We will also run automatic checks of code on Moodle for plagiarism. Copying code from others is considered a serious case of cheating.

# 1. HMM (15 points)

Answer the following questions using the transition matrix $T$ and emission probabilities $E$ below. Below, $\Delta$ and $\square$ are two output variables, $A$ and $B$ are two hidden states; $s_n$ refers to the $n^{th}$ hidden state in the sequence and $o_n$ refers to the $n^{th}$ observation.

For all the questions in this section, write answer and show your work.

**Question 1.1 (2 points)**

Does $P(o_2 = \Delta | s_1 = B) = P(o_2 = \Delta | o_1 = \square)$?

False

$P(o_2 = \Delta | s_1 = B)$ = 0.4*0.5* + *0.4*0.3 = 0.32

$P(o_2 = \Delta | o_1 = \square)$ = 0.224

$P(o_2 = \Delta | s_1 = B) \neq P(o_2 = \Delta | o_1 = \square)$

**Question 1.2 (2 points)**

Does $P(s_2 = B | s_1 = A) = P(s_2 = B | s_1 = A, o_1 = \Delta)$?

True

$P(s_2 = B | s_1 = A)$ = 0.4

$P(s_2 = B | s_1 = A, o_1 = \Delta) =$

$P(s_2 = B, s_1 = A, o_1 = \Delta) / P(s_1 = A, o_1 = \Delta)$

Solving the above equation we get $P(s_2 = B | s_1 = A, o_1 = \Delta)$ = 0.4

That means: -

$P(s_2 = B | s_1 = A) = P(s_2 = B | s_1 = A, o_1 = \Delta)$

**Question 1.3 (3 points)**

Does $P(o_2 = \Delta | s_1 = A) = P(o_2 = \square | s_1 = A, s_3 = A)$

False

$P(o_2 = \Delta | s_1 = A) = 0.2 * 0.5 + 0.3 * 0.3 = 0.19$

$P(o_2 = \Box | s_1 = A, s_3 = A) = P(o_2 = \Box, s_3 = A | s_1 = A)/P(s_3 = A | s_1 = A)$

$= (P(o_2 = \Box, s_2 = A, s_3 = A | s_1 = A) + P(o_2 = \Box, s_2 = B, s_3 = A | s_1 = A))/P(s_3 = A | s_1 = A)$

$(P(o_2 = \Box, s_2 = A, s_3 = A | s_1 = A) = 0.02$

$P(o_2 = \Box, s_2 = B, s_3 = A | s_1 = A) = 0.084$

$P(o_2 = \Box | s_1 = A, s_3 = A) = 0.65$

$P(o_2 = \Delta | s_1 = A)! = P(o_2 = \Box | s_1 = A, s_3 = A)$

**Question 1.4 (3 points)**

Compute the probability of observing $\Box$ as the first emission of a sequence generated by an HMM with transition matrix $T$ and emission probabilities $E$.

1st emmision refer to the emmision by the s1.

we need to calculate $P(o_1 = \Box)$

$P(o_1 = \Box) = P(o_1 = \Box, s_1 = A) + P(o_1 = \Box, s_1 = B)$

$P(o_1 = \Box) = 0.6$

**Question 1.5 (5 points)**

Compute the probability of the first state being $A$ given that the last token in an observed sequence of length 2 was the token $\Delta$.

we need to compute $p(s_1 = A | o_2 = \Delta)$

$p(s_1 = A | o_2 = \Delta) = p(s_1 = A, o_2 = \Delta)/P(o_2 = \Delta)$

$p(s_1 = A | o_2 = \Delta) = 0.3725$

# 2. Viterbi (log-additive form) (20 points)

One HMM chain is shown on the left. The corresponding **factor graph** version is shown on the right. This simply shows the structure of the $A$ and $B_t$ log-prob tables and which variables they express preferences over. $A$ is the **transition factor** that has preferences for the two neighboring variables; for example, $A(y_1, y_2)$ shows how happy the model is with the transition from $y_1$ to $y_2$. The same transition preference function is used at all positions $(t-1, t)$ for each $t = 2..T$. $B_t$ is the **emission factor** that has preferences for the variable $y_t$. As a goodness function it is e.g. $B_1(y_1)$, $B_2(y_2)$, etc.

Let $\vec{y} = (y_1, y_2, \ldots, y_T)$, a proposed tag sequence for a $T$ length sentence. The total goodness function for a solution $\vec{y}$ is

$$G(\vec{y}) = \sum_{t=1}^{T} B_t(y_t) + \sum_{t=2}^{T} A(y_{t-1}, y_t)$$

**Question 2.1 (2 points)**

Define $A$ and $B_t$ in terms of the HMM model, such that $G$ is the same thing as $\log p(\vec{y}, \vec{w})$ under the HMM.

$A(x) = log(p(w_t/x))$ for $x \geq 2 \; x \in I$

$B_t(x) = log(p(w_t/x)) \; x \geq 1 \; x \in I$

**Question 2.2 (18 points)**

Implement additive log-space Viterbi by completing the **viterbi()** function. It takes in tables that represent the $A$ and $B$ functions as input. We give you an implementation of $G()$ in **vit_starter**, you can check to make sure you understand the data structures, and also the exhaustive decoding algorithm too. Feel free to add debugging print statements as needed. The main code runs the exercise example by default.

When debugging, you should make new A and B examples that are very simple. This will test different code paths. Also you can try the **randomized_test()** from the starter code.

Look out for negative indexes as a bug. In python, if you use an index that's too high to be in the list, it throws an error. But it will silently accept a negative index ... it interprets that as indexing from the right.

In [2]: 
```python
# Implement the viterbi() function in vit_starter.py and then run this cell to
 show your output

from vit_starter import *

if __name__=='__main__':
    A = {(0,0):3, (0,1):0, (1,0):0, (1,1):3}
    Bs= [ [0,1], [0,1], [30,0] ]
    # that's equivalent to: [ {0:0,1:1}, {0:0,1:1}, {0:30,1:0} ]

    y = exhaustive(A, Bs, set([0,1]))
    print "Exhaustive decoding:", y
    print "score:", goodness_score(y, A, Bs)
    y = viterbi(A, Bs, set([0,1]))
    print "Viterbi    decoding:", y
```

```
Exhaustive decoding: [0, 0, 0]
score: 36
Viterbi    decoding: [1, 1, 0]
```

**Copy and paste the viterbi function that you implemented in `vit_starter.py`.**

```
In [3]: def viterbi(A_factor, B_factors, output_vocab):
            """
            A_factor: a dict of key:value pairs of the form
                {(curtag,nexttag): score}
            with keys for all K^2 possible neighboring combinations,
            and scores are numbers.  We assume they should be used ADDITIVELY, i.e. in
        log space.
            higher scores mean MORE PREFERRED by the model.

            B_factors: a list where each entry is a dict {tag:score}, so like
            [ {Noun:-1.2, Adj:-3.4}, {Noun:-0.2, Adj:-7.1}, .... ]
            each entry in the list corresponds to each position in the input.

            output_vocab: a set of strings, which is the vocabulary of possible output
            symbols.

            RETURNS:
            the tag sequence yvec with the highest goodness score
            """

            N = len(B_factors)    # length of input sentence

            # viterbi log-prob tables
            V = [{tag:None for tag in output_vocab} for t in range(N)]
            # backpointer tables
            # back[0] could be left empty. it will never be used.
            back = [{tag:None for tag in output_vocab} for t in range(N)]
            V[0]={tag:B_factors[0][tag] for tag in output_vocab}
            for t in range(1,N):
                for tag in output_vocab:
                    max_vit=-1*float("inf")
                    for prev_tag in output_vocab:
                        vit_score= V[t-1][prev_tag] + A_factor[prev_tag,tag] + B_facto
        rs[t][tag]

                        if vit_score > max_vit:
                            max_vit=vit_score
                            V[t][tag] = vit_score
                            back[t][tag]= prev_tag
            best_tag=dict_argmax(V[N-1])
            bestseq=[best_tag]
            for t in range(N-1,0,-1):
                best_tag = back[t][best_tag]
                bestseq.append(best_tag)
            bestseq.reverse()
            return bestseq
```

# 3. Averaged Perceptron (5 points)

We will be using the following definition of the perceptron, which is the multiclass or structured version of the perceptron. The training set is a bunch of input-output pairs $(x_i, y_i)$. (For classification, $y_i$ is a label, but for tagging, $y_i$ is a sequence). The training algorithm is as follows:

For T iterations, iterate through each $(x_i, y_i)$ pair in the dataset, and for each,

1. Predict $y^* := \arg\max_{y'} \theta^T f(x_i, y')$
2. If $y_i \neq y^*$: then update $\theta := \theta^{(old)} + rg$

where $r$ is a fixed step size (e.g. $r = 1$) and $g$ is the *gradient vector*, meaning a vector that will get added into $\theta$ for the update, specifically

$$g = \underbrace{f(x_i, y_i)}_{\text{feats of true output}} - \underbrace{f(x_i, y^*)}_{\text{feats of predicted output}}$$

Both in theory and in practice, the predictive accuracy of a model trained by the structured perceptron will be better if we use the average value of $\theta$ over the course of training, rather than the final value of $\theta$. This is because $\theta$ wanders around and doesn't converge (typically), because it overfits to whatever data it saw most recently. After seeing $t$ training examples, define the *averaged parameter vector* as

where $\theta_{t'}$ is the weight vector after $t'$ updates. (We are counting $t$ by the number of training examples, not passes through the data. So if you had 1000 examples and made 10 passes through the data in order, the final time you see the final example is $t = 10000$.) For training, you still use the current $\theta$ parameter for predictions. But at the very end, you return the $\bar{\theta}$, not $\theta$, as your final model parameters to use on test data.

Directly implementing equation (1) would be really slow. So here's a better algorithm. This is the same as in Hal Daume's CIML chapter on perceptrons, but adapted for the structured case (as opposed to Daume's algorithm, which assumes binary output). Define $g_t$ to be the update vector $g$ as described earlier. The perceptron update can be written

$$\theta_t = \theta_{t-1} + rg_t$$

Thus the averaged perceptron algorithm is, using a new 'weightsums' vector $S$,

1. Initialize $t = 1, \theta_0 = \vec{0}, S_0 = \vec{0}$
2. For each example $i$ (iterating multiples times through dataset),

   - Predict $y^* = \arg\max_{y'} \theta^T f(x_i, y')$
   - Let $g_t = f(x_i, y_i) - f(x_i, y^*)$
   - Update $\theta_t = \theta_{t-1} + rg_t$
   - Update $S_t = S_{t-1} + (t-1)rg_t$
   - $t := t + 1$
3. Calculate $\bar{\theta}$ based on $S$

In an actual implementation, you don't keep old versions of $S$ or $\theta$ around ... above we're using the $t$ subscripts above just to make the mathematical analysis clearer.

Our proposed algorithm computes $\bar{\theta}_t$ as

For the following problems, feel free to set $r = 1$ just to simplify them.

For following questions write only math answers, no code required.

**Question 3.1** (1 point)

What is the computational advantage of computing $\bar{\theta}$ using Equation (2) instead of directly implementing Equation (1)?

Ans :- So we dont need to store any intermediate $\theta_t$ in order to calculate $\bar{\theta}_t$, we will use final $s_t$ and $\theta_t$ to calculate $\bar{\theta}_t$

Now we'll show this works, at least for early iterations.

**Question 3.2** (1 point)

What are $\bar{\theta}_1$, $\bar{\theta}_2$, $\bar{\theta}_3$, and $\bar{\theta}_4$? Please derive them from the Equation (1) definition, and state them in terms of $g_1$, $g_2$, $g_3$, and/or $g_4$.

$$\bar{\theta}_1 = g_1$$

$$\bar{\theta}_2 = (2 * g_1 + g_2)/2$$

$$\bar{\theta}_3 = (3 * g_1 + 2 * g_2 + g_3)/3$$

$$\bar{\theta}_4 = (4 * g_1 + 3 * g_2 + 2 * g_3 + g_4)/4$$

**Question 3.3** (1 point)

What are $S_1$, $S_2$, $S_3$, and $S_4$? Please state them in terms of $g_1$, $g_2$, $g_3$, and/or $g_4$.

$$s_1 = 0$$

$$s_2 = g_2$$

$$s_3 = g_2 + 2 * g_3$$

$$s_4 = g_2 + 2 * g_3 + 3 * g_4$$

$$s_T = \sum_{t=1}^{T}(t-1) * g_t$$

**Question 3.4** (2 points)

Show that Equation (2) correctly computes $\bar{\theta}_3$ and $\bar{\theta}_4$.

$$\bar{\theta}_3 = (3 * g_1 + 2 * g_2 + g_1)/3$$

$$\theta_3 = g_1 + g_2 + g_1$$

$$\bar{\theta}_3 = \theta_3 - (g_2 + 3 * g_3)/3$$

$$= \bar{\theta}_3 - s_3/3$$

The same is true for the $\bar{\theta}_4$ using the formula

$$\bar{\theta}_T = \theta_T - s_T/T$$

**Question 3.5** (2 Extra Credit points)

Use proof by induction to show that this algorithm correctly computes $\bar{\theta}_t$ for any $t$.

$$\bar{\theta}_1 = \theta_1 - s_1 = g_1$$

$$\bar{\theta}_2 = \theta_2 - s_2/2 = g_1 + g_2/2$$

We can assume now

$$\bar{\theta}_T = \theta_T - s_T/T$$

$$\bar{\theta}_{T+1} = \theta_{T+1} - s_{T+1}/(T+1)$$

$$RHS$$

$$\bar{\theta}_{T+1} = \sum_{t=1}^{T} \theta_t/(t+1) = \bar{\theta}_t * t/(t+1) + \theta_{t+1}/(t+1)$$

$$LHS$$

$$\theta_{T+1} - s_{T+1}/(T+1) = \theta_t + g_{t+1} - (s_t + g_{t+1})/(t+1)$$

$$\theta_t - s_t/(t+1) + g_{t+1}/(t+1)$$

putting $s_t = t(\theta_t - \bar{\theta}_t)$

$$\theta_t - s_t/(t+1) + g_{t+1}/(t+1) = t * \theta_t/(t+1) + (\theta_t + g_{t+1})/(t+1) = \bar{\theta}_{t+1}$$

So, by induction we can say,

$$\bar{\theta}_{T+1} = \theta_{T+1} - s_{T+1}/(T+1)$$

# 4. Classifier Perceptron (20 points)

Implement the averaged perceptron for document classification, using the same sentiment analysis dataset as you used for HW1. On the first two questions, we're asking you to develop using only a subset of the data, since that makes debugging easier. On the third question, you'll run on the full dataset, and you should be able to achieve a higher accuracy compared to your previous Naive Bayes implementation. Starter code is provided in `classperc.py`.

**Question 4.1** (8 points)

Implement the simple, non-averaged perceptron. Run your code on **the first 1000 training instances** for 10 passes through the training data. For each pass, report **the training and test set accuracies**.

```
In [8]:  "autoreload"
         from classperc import *
         training_set = construct_dataset(train=True)
         test_set = construct_dataset(train=False)

         vanilla_dict = train(training_set[0:1000], stepsize=10, numpasses=10, do_avera
         ging=False, devdata=test_set[0:1000])
         # Please find the accuracy in below
```

```
[constructing dataset...]
        reading data from large_movie_review_dataset/train\pos
        reading data from large_movie_review_dataset/train\neg
[dataset constructed.]
[constructing dataset...]
        reading data from large_movie_review_dataset/test\pos
        reading data from large_movie_review_dataset/test\neg
[dataset constructed.]
[training...]
        Training iteration 0
TR RAW EVAL: 1000/1000 = 1.0000 accuracy
DEV RAW EVAL: 1000/1000 = 1.0000 accuracy
        Training iteration 1
TR RAW EVAL: 1000/1000 = 1.0000 accuracy
DEV RAW EVAL: 1000/1000 = 1.0000 accuracy
        Training iteration 2
TR RAW EVAL: 1000/1000 = 1.0000 accuracy
DEV RAW EVAL: 1000/1000 = 1.0000 accuracy
        Training iteration 3
TR RAW EVAL: 1000/1000 = 1.0000 accuracy
DEV RAW EVAL: 1000/1000 = 1.0000 accuracy
        Training iteration 4
TR RAW EVAL: 1000/1000 = 1.0000 accuracy
DEV RAW EVAL: 1000/1000 = 1.0000 accuracy
        Training iteration 5
TR RAW EVAL: 1000/1000 = 1.0000 accuracy
DEV RAW EVAL: 1000/1000 = 1.0000 accuracy
        Training iteration 6
TR RAW EVAL: 1000/1000 = 1.0000 accuracy
DEV RAW EVAL: 1000/1000 = 1.0000 accuracy
        Training iteration 7
TR RAW EVAL: 1000/1000 = 1.0000 accuracy
DEV RAW EVAL: 1000/1000 = 1.0000 accuracy
        Training iteration 8
TR RAW EVAL: 1000/1000 = 1.0000 accuracy
DEV RAW EVAL: 1000/1000 = 1.0000 accuracy
        Training iteration 9
TR RAW EVAL: 1000/1000 = 1.0000 accuracy
DEV RAW EVAL: 1000/1000 = 1.0000 accuracy
[learned weights for 0 features from 1000 examples.]
```

**Question 4.2** (8 points)

Implement the averaged perceptron. Run your code on **the first 1000 training instances** for 10 passes through the training data. For each pass, compute the $\bar{\theta}$ so far, and report its **test set accuracy**.

In [9]:
```
average_dict = train(training_set[0:1000], stepsize=10, numpasses=10, do_avera
ging=True, devdata=test_set[0:1000])
#please find the accuracy below
```

```
[training...]
        Training iteration 0
TR RAW EVAL: 1000/1000 = 1.0000 accuracy
DEV RAW EVAL: 1000/1000 = 1.0000 accuracy
DEV AVG EVAL: 1000/1000 = 1.0000 accuracy
        Training iteration 1
TR RAW EVAL: 1000/1000 = 1.0000 accuracy
DEV RAW EVAL: 1000/1000 = 1.0000 accuracy
DEV AVG EVAL: 1000/1000 = 1.0000 accuracy
        Training iteration 2
TR RAW EVAL: 1000/1000 = 1.0000 accuracy
DEV RAW EVAL: 1000/1000 = 1.0000 accuracy
DEV AVG EVAL: 1000/1000 = 1.0000 accuracy
        Training iteration 3
TR RAW EVAL: 1000/1000 = 1.0000 accuracy
DEV RAW EVAL: 1000/1000 = 1.0000 accuracy
DEV AVG EVAL: 1000/1000 = 1.0000 accuracy
        Training iteration 4
TR RAW EVAL: 1000/1000 = 1.0000 accuracy
DEV RAW EVAL: 1000/1000 = 1.0000 accuracy
DEV AVG EVAL: 1000/1000 = 1.0000 accuracy
        Training iteration 5
TR RAW EVAL: 1000/1000 = 1.0000 accuracy
DEV RAW EVAL: 1000/1000 = 1.0000 accuracy
DEV AVG EVAL: 1000/1000 = 1.0000 accuracy
        Training iteration 6
TR RAW EVAL: 1000/1000 = 1.0000 accuracy
DEV RAW EVAL: 1000/1000 = 1.0000 accuracy
DEV AVG EVAL: 1000/1000 = 1.0000 accuracy
        Training iteration 7
TR RAW EVAL: 1000/1000 = 1.0000 accuracy
DEV RAW EVAL: 1000/1000 = 1.0000 accuracy
DEV AVG EVAL: 1000/1000 = 1.0000 accuracy
        Training iteration 8
TR RAW EVAL: 1000/1000 = 1.0000 accuracy
DEV RAW EVAL: 1000/1000 = 1.0000 accuracy
DEV AVG EVAL: 1000/1000 = 1.0000 accuracy
        Training iteration 9
TR RAW EVAL: 1000/1000 = 1.0000 accuracy
DEV RAW EVAL: 1000/1000 = 1.0000 accuracy
DEV AVG EVAL: 1000/1000 = 1.0000 accuracy
[learned weights for 0 features from 1000 examples.]
```

**Question 4.3** (4 points)

Graph four curves on the same plot, using the **full dataset**:

- accuracy of the vanilla perceptron on the training set
- accuracy of the vanilla perceptron on the test set
- accuracy of the averaged perceptron on the test set
- accuracy of your Naive Bayes classifier from HW1 (you don't need to re-run it; just take the best accuracy from your previous results).

The x-axis of the plot should show the number of iterations through the training set and the y-axis should show the accuracy of the classifier. For this part of the HW run your code on **the entire dataset** (all instances). Since Naive Bayes doesn't require multiple passes through the data just produce a single horizontal line showing its overall accuracy. Make sure your plot has a title, a label on the x-axis, a label on the y-axis and a legend showing which line is which. Explain verbally what's happening in this plot.

In [2]:
```python
from classperc import *
training_set = construct_dataset(train=True)
test_set = construct_dataset(train=False)
vanilla_dict = train(training_set, stepsize=10, numpasses=1000, do_averaging=False, devdata=test_set)
average_dict = train(training_set, stepsize=10, numpasses=1000, do_averaging=True, devdata=test_set)
"%autoreload"
%matplotlib inline
from matplotlib import pyplot as plt
#from classperc import plot_accuracy_vs_iteration
#plot_accuracy_vs_iteration(vanilla_dict, vanilla_dict, average_dict, [0.843])
plt.plot(vanilla_dict['train_acc'], [x for x in range(1, 1001)], label = "vanilla perceptron on the training set")
plt.plot(vanilla_dict['test_acc'], [x for x in range(1, 1001)], label = "vanilla perceptron on the test set")
plt.plot(average_dict['test_acc'],[x for x in range(1, 1001)], label = "averaged perceptron on the test set")
plt.plot([0.843 for x in range(1, 1001)],[x for x in range(1, 1001)] , label = "Naive Bayes classifier")
plt.xlabel("Number of iterations")
plt.ylabel("Accuracy of the classifier")
plt.title("Accuracy of classifier for number of iterations")
plt.legend()

## The average perceptron is performing better as compared to the vanilla perceptrons. The vanilla perceptron on the training set,
## is clearly overfitting because its performance significantly degrades in the test set.
```

```
[constructing dataset...]
        reading data from large_movie_review_dataset/train\pos
        reading data from large_movie_review_dataset/train\neg
[dataset constructed.]
[constructing dataset...]
        reading data from large_movie_review_dataset/test\pos
        reading data from large_movie_review_dataset/test\neg
[dataset constructed.]
[training...]
        Training iteration 0
        Training iteration 100
        Training iteration 200
        Training iteration 300
        Training iteration 400
        Training iteration 500
        Training iteration 600
        Training iteration 700
        Training iteration 800
        Training iteration 900
[learned weights for 78380 features from 2000 examples.]
[training...]
        Training iteration 0
        Training iteration 100
        Training iteration 200
        Training iteration 300
        Training iteration 400
        Training iteration 500
        Training iteration 600
        Training iteration 700
        Training iteration 800
        Training iteration 900
[learned weights for 78888 features from 2000 examples.]
```

Out[2]:  <matplotlib.legend.Legend at 0xe9ca4a8>



# 5. Structured Perceptron with Viterbi (40 points)

In this problem, you will implement a part-of-speech tagger for Twitter, using the structured perceptron algorithm. Your system will be not too far off from state of the art performance, coding it all up yourself from scratch!

The dataset comes from http://www.ark.cs.cmu.edu/TweetNLP/ (http://www.ark.cs.cmu.edu/TweetNLP/) and is described in the papers listed there (Gimpel et al.~2011 and Owoputi et al.~2013). The Gimpel article describes the tagset; the annotation guidelines on that webpage describe it futher.

Your structured perceptron will use your Viterbi implementation from 2.2 as a subroutine. If that's buggy, this will cause many problems here---your perceptron will have really weird behavior. (This happened to us when designing your assignment!) If you have problems, try using the greedy decoding algorithm, which we provide in the starter code. Make sure to note which decoding algorithm you're using in your writeup.

The starter code is `structperc.py` and it assumes the two data files `oct27.train` and `oct27.dev` are in the same directory. (For simplicity we're just going to use this `dev` set as our test set.)

**Question 5.1** (2 points)

First let's do a little data analysis to establish the **most common tag** baseline accuracy. Using a small script, load up the dev dataset (oct27.dev) using the function `structperc.read_tagging_file` (from `import structperc`). Calculate the following: What is the most common tag, and what would your accuracy be if you predicted it for all tags?

```
In [19]:  ##Write your code here and show the output
          import structperc
          from collections import defaultdict
          ret= structperc.read_tagging_file("oct27.dev")
          tag_count=defaultdict(int)
          for tokens,tags in ret:
              for tag in tags:
                  tag_count[tag] += 1

          most_common_tag,frequency = max(tag_count.iteritems(),key=lambda k:k[1])
          print "Most common tag : ", most_common_tag
          accuracy = 1.0*frequency/(sum(tag_count.values())) * 100
          print 'Accuracy if most common tag is the output tag : {0}%'.format(accuracy)
```

```
Most common tag :  V
Accuracy if most common tag is the output tag : 15.5712212316%
```

The structured perceptron algorithm works very similarly as the classification version you did in the previous question, except the prediction function uses Viterbi as a subroutine, which has to call feature extraction functions for local emissions and transition factors. There also has to be a large overall feature extraction function for an entire structure at once. The following parts will build up these pieces. First, we will focus on inference, not learning.

**Question 5.2** (2 points)

We provide a barebones version of `local_emission_features`, which calculates the local features for a particular tag at a token position. You can run this function all by itself. Make up an example sentence, and call this function with it, giving it a particular index and candidate tag. Show the code for the function call you made and the function's return value, and explain what the features mean (just a sentence or two).

```
In [20]: ##Show the code for function call with output
         from structperc import local_emission_features
         def test_local_emission():
             example_sentence = "This is test"
             tag = 'A'
             tokens = example_sentence.split(" ")
             print(local_emission_features(0, tag, tokens))
         test_local_emission()
```
```
{'tag=A_curword=This': 1, 'tag=A_biasterm': 1}
```

# This means that we have created a feature f(TAG,CURWORD,null) in the form of a dictionary with

$key = "tag = \boxed{\_}curword="$.

**Question 5.3** (2 points)

Implement `features_for_seq()`, which extracts the full feature vector $f(x, y)$, where $x$ is a sentence and $y$ is an entire tagging sequence for that sentence. This will add up the feature vectors from each local emissions features for every position, as well as transition features for every position (there are $N - 1$ of them, of course). Show the output on a very short example sentence and example proposed tagging, that's only 2 or 3 words long.

To define $f(x, y)$ a little more precisely: If $f^{(B)}(t, x, y)$ means the local emissions feature vector at position $t$ (i.e. the `local_emission_features` function), and $f^{(A)}(y_{t-1}, y_t, y)$ is the transition feature function for positions $(t - 1, t)$ (which just returns a feature vector where everything is zero, except a single element is 1), then the full sequence feature vector will be the vector-sum of all those feature vectors:

$$f(x, y) = \sum_t^T f^{(B)}(t, x, y) + \sum_{t=2}^T f^{(A)}(y_{t-1}, y_t)$$

You implemented $f^{(B)}$ above. You probably don't need to bother implementing $f^{(A)}$ as a standalone function. You will have to decide on a particular convention to encode the name of a transition feature. For example, one way to do it is with string concatenation like this, `"trans_%s_%s" % (prevtag, curtag)`, where prevtag and curtag are strings. Or you could use a python tuple of strings, which works since tuples have the ability to be keys in a python dictionary.

In other words: the emissions and transition features will all be in the same vector, just as keys in the dictionary that represents the feature vector. The transition features are going to be the count of how many times a particular transition (tag bigram) happened. The emissions features are going to be the vector-sum of all the local emission features, as calculated from `local_emission_features`.

```
In [21]:  ##Show the call to your function and output
          example_sentence = "Test this Code"
          tags=["V","A","N"]
          tokens=example_sentence.split()
          structperc.features_for_seq(tokens,tags)
```

```
Out[21]:  defaultdict(int,
                      {'tag=A_biasterm': 1,
                       'tag=A_curword=this': 1,
                       'tag=N_biasterm': 1,
                       'tag=N_curword=Code': 1,
                       'tag=V_biasterm': 1,
                       'tag=V_curword=Test': 1,
                       ('A', 'N'): 1,
                       ('V', 'A'): 1})
```

**Question 5.4** (4 points)

Look at the starter code for `calc_factor_scores`, which calculates the A and B score functions that are going to be passed in to your Viterbi implementation from problem 2, in order to do a prediction. The only function it will need to call is `local_emission_features`. It should NOT call `features_for_seq`. Why not?

Ans :- We want to get the total set of features, feature_for_seq function assumes that the given tokens in the first argument are labelled by the labelseq in the second argument and thats not our objective here in the calc_factor_scores

**Question 5.5** (6 points)

Implement `calc_factor_scores`. Make up a simple example (2 or 3 words long), with a simple model with at least some nonzero features (you might want to use a `defaultdict(float)`, so you don't have to fill up a dict with dummy values for all possible transitions), and show your call to this function and the output.

```
In [12]:  import structperc
          example_sentence = "Test this code"
          weights= { ('V','A'): 5, "tag=V_curword=Test":20}
          tokens = example_sentence.split()
          print structperc.calc_factor_scores(tokens,weights)
```

```
(defaultdict(<type 'float'>, {('V', 'A'): 5}), [defaultdict(<type 'float'>,
 {'!': 0, '#': 0, '$': 0, '&': 0, ',': 0, 'A': 0, '@': 0, 'E': 0, 'D': 0,
  'G': 0, 'M': 0, 'L': 0, 'O': 0, 'N': 0, 'P': 0, 'S': 0, 'R': 0, 'U': 0, 'T':
0, 'V': 20, 'Y': 0, 'X': 0, 'Z': 0, '^': 0}), defaultdict(<type 'float'>,
 {'!': 0, '#': 0, '$': 0, '&': 0, ',': 0, 'A': 0, '@': 0, 'E': 0, 'D': 0,
  'G': 0, 'M': 0, 'L': 0, 'O': 0, 'N': 0, 'P': 0, 'S': 0, 'R': 0, 'U': 0, 'T':
0, 'V': 0, 'Y': 0, 'X': 0, 'Z': 0, '^': 0}), defaultdict(<type 'float'>,
 {'!': 0, '#': 0, '$': 0, '&': 0, ',': 0, 'A': 0, '@': 0, 'E': 0, 'D': 0,
  'G': 0, 'M': 0, 'L': 0, 'O': 0, 'N': 0, 'P': 0, 'S': 0, 'R': 0, 'U': 0, 'T':
0, 'V': 0, 'Y': 0, 'X': 0, 'Z': 0, '^': 0})])
```

**Question 5.6** (4 points)

Implement `predict_seq()`, which predicts the tags for an input sentence, given a model. It will have to calculate the factor scores, then call Viterbi as a subroutine, then return the best sequence prediction. If your Viterbi implementation does not seem to be working, use the implementation of the greedy decoding algorithm that we provide (it uses the same inputs as `vit_starter.viterbi()`).

```
In [13]:  ## Copy and paste predict_seq() function here
          def predict_seq(tokens, weights):
              """
              IMPLEMENT ME!
              takes tokens and weights, calls viterbi and returns the most likely
              sequence of tags
              """
              # once you have Ascores and Bscores, could decode with
              # predlabels = greedy_decode(Ascores, Bscores, OUTPUT_VOCAB)
              (Ascores, Bscores) = calc_factor_scores(tokens, weights)
              return viterbi(Ascores, Bscores, OUTPUT_VOCAB)
```

OK, you're done with the inference part. Time to put it all together into the parameter learning algorithm and see it go.

**Question 5.7** (14 points)

Implement `train()`, which does structured perceptron training with the averaged perceptron algorithm. You should train on oct27.train, and evaluate on oct27.dev. You will want to first get it working without averaging, then add averaging to it. Run it for 10 iterations, and print the devset accuracy at each training iteration. Note that we provide evaluation code, which assumes `predict_seq()` and everything it depends on is working properly.

For us, here's the performance we get at the first and last iterations, using the features in the starter code (just the bias term and the current word feature, without case normalization).

```
Training iteration 0
DEV RAW EVAL: 2556/4823 = 0.5300 accuracy
DEV AVG EVAL: 2986/4823 = 0.6191 accuracy
...
Training iteration 9
DEV RAW EVAL: 3232/4823 = 0.6701 accuracy
DEV AVG EVAL: 3341/4823 = 0.6927 accuracy
Learned weights for 24361 features from 1000 examples
```

```
In [24]: "%autoload"
         import structperc
         structperc.train_model()
```

```
Training iteration 0
DEV RAW EVAL: 2695/4823 = 0.5588 accuracy
DEV AVG EVAL: 2784/4823 = 0.5772 accuracy
Training iteration 1
DEV RAW EVAL: 2441/4823 = 0.5061 accuracy
DEV AVG EVAL: 2946/4823 = 0.6108 accuracy
Training iteration 2
DEV RAW EVAL: 2976/4823 = 0.6170 accuracy
DEV AVG EVAL: 3057/4823 = 0.6338 accuracy
Training iteration 3
DEV RAW EVAL: 2740/4823 = 0.5681 accuracy
DEV AVG EVAL: 3086/4823 = 0.6399 accuracy
Training iteration 4
DEV RAW EVAL: 2793/4823 = 0.5791 accuracy
DEV AVG EVAL: 3088/4823 = 0.6403 accuracy
Training iteration 5
DEV RAW EVAL: 2739/4823 = 0.5679 accuracy
DEV AVG EVAL: 3099/4823 = 0.6425 accuracy
Training iteration 6
DEV RAW EVAL: 3168/4823 = 0.6569 accuracy
DEV AVG EVAL: 3114/4823 = 0.6457 accuracy
Training iteration 7
DEV RAW EVAL: 3050/4823 = 0.6324 accuracy
DEV AVG EVAL: 3121/4823 = 0.6471 accuracy
Training iteration 8
DEV RAW EVAL: 2966/4823 = 0.6150 accuracy
DEV AVG EVAL: 3128/4823 = 0.6486 accuracy
Training iteration 9
DEV RAW EVAL: 2932/4823 = 0.6079 accuracy
DEV AVG EVAL: 3129/4823 = 0.6488 accuracy
Learned weights for 21457 features from 1000 examples
gold , acc 0.9840 (492/500)
gold V acc 0.9587 (720/751)
gold & acc 0.9560 (87/91)
gold P acc 0.9500 (418/440)
gold D acc 0.9295 (290/312)
gold O acc 0.9189 (306/333)
gold T acc 0.8056 (29/36)
gold L acc 0.8000 (52/65)
gold R acc 0.7081 (148/209)
gold E acc 0.5385 (28/52)
gold ! acc 0.5051 (50/99)
gold A acc 0.4979 (119/239)
gold N acc 0.4288 (283/660)
gold $ acc 0.3488 (30/86)
gold G acc 0.2462 (16/65)
gold ^ acc 0.1576 (49/311)
gold # acc 0.0962 (5/52)
gold @ acc 0.0247 (6/243)
gold U acc 0.0110 (1/91)
gold S acc 0.0000 (0/5)
gold X acc 0.0000 (0/4)
gold Z acc 0.0000 (0/9)
gold ~ acc 0.0000 (0/170)
word                gold pred
----                ---- ----
@ciaranyree          @    V     *** Error
```

```
          it                  O    O
          was                 V    V
          on                  P    P
          football            N    N
          wives               N    V        *** Error
          ,                   ,    ,
          one                 $    $
          of                  P    P
          the                 D    D
          players             N    V        *** Error
          and                 &    &
          his                 D    D
          wife                N    N
          own                 V    N        *** Error
          smash               ^    V        *** Error
          burger              ^    V        *** Error
          word                gold pred
          ----                ---- ----
          RT                  ~    Y        *** Error
          @TheRealQuailman    @    V        *** Error
          :                   ~    ,        *** Error
          Currently           R    V        *** Error
          laughing            V    V
          at                  P    P
          Laker               ^    ^
          haters              N    V        *** Error
          .                   ,    ,
```

**Question 5.8** (6 points)

Print out a report of the accuracy rate for each tag in the development set. We provided a function to do this
`fancy_eval`. Look at the two sentences in the dev data, and in your writeup show and compare the gold-
standard tags versus your model's predictions for them. Consult the tagset description to understand what's
going on. What types of things does your tagger get right and wrong?

To look at the examples, you may find it convenient to use `show_predictions` (or write up the equivalent
manually). For example, after 1 iteration of training, we get this output from the first sentence in the devset. (After
investigating TV shows that were popular in 2011 when the tweet was authored, we actually think some of the
gold-standard tags in this example might be wrong.)

```
word                gold pred

----                ---- ----
@ciaranyree         @    @
it                  O    O
was                 V    V
on                  P    P
football            N    ^       *** Error
wives               N    N

,                   ,    ,
one                 $    $
of                  P    P
the                 D    D
players             N    N
and                 &    &
his                 D    D
wife                N    N
own                 V    V
smash               ^    D       *** Error
burger              ^    N       *** Error
```

To do this part, you may find it useful to save your model's weights with pickle.dumps (or json.dumps) and have a
short analysis script that loads the model and devdata to do the reports. If you have to re-train each time you
tweak your analysis code, it can be annoying.

Ans :- Please refer the above cell for above for the solution

**Question 5.9** (OPTIONAL: 4 Extra Credit points)

Improve the features of your tagger to improve accuracy on the development set. This will only require changes to `local_emission_features`. Implement at least 4 new types of features. Report your tagger's accuracy with these improvements. Please make a table that reports accuracy from adding different features. The first row should be the basic system, and the last row should be the fanciest system. Rows in between should report different combinations of features. One simple way to do this is, if you have 4 different feature types, to run 4 experiments where in each one, you add only one feature type to the basic system. For example:

Hint: if you make features about the first character of a word, that helps a lot for the # (hashtag) and @ (at-mention) tags. The URL tag is easy to get too with a similar form of character affix analysis. Character affixes help lots of other tags too. Also, if you have a feature that looks at the word at position $t$, you can make new versions of it that look to the left or right of the $t^{th}$ position in question: for example, 'word_to_left=the'.

In [ ]: