## CPTS 580 Structured Prediction: Algorithms and Applications, Spring 2017
## Homework #1
Due Date: Tue, Feb 14

NOTE 1: Please use a word processing software (e.g., Microsoft word or Latex) to write your answers and submit a printed copy to me. The rationale is that it is sometimes hard to read and understand the hand-written answers. Thanks for your understanding.

NOTE 2: Please ensure that all the graphs are appropriately labeled (x-axis, y-axis, and each curve). The caption or heading of each graph should be informative and self-contained.

1. (**15 points**) Tom wants to apply the LaSO framework to the problem of co-reference resolution that arises in natural language processing. Coreference resolution is a structured prediction problem where the set of mentions $m_1, m_2, \cdots, m_D$ extracted from a document $D$ correponds to a structured input $x$ and the structured output $y$ corresponds to a partition of the mentions into a set of clusters $C_1, C_2, \cdots, C_k$. Each mention $m_i$ belongs to exactly one of the clusters $C_j$. Tom needs your help in designing a concrete search space. Can you give some concrete search space definitions, and discuss their pros and cons in terms of learning an accurate heuristic function to guide the search process?

2. (**85 points**) Please implement the online structured perceptron training algorithm for sequence labeling problems and experiment with two sequence labeling datasets: handwriting recognition, and text-to-speech mapping.

   In a sequence labeling problem, the structured input $x = (x_1, x_2, \cdots, x_T)$ is a sequence of input tokens, where each input token $x_i$ is represented as a $m$-dimensional feature vector; and the structured output $y = (y_1, y_2, \cdots, y_T)$ is a sequence of output labels, where each output label $y_i$ comes from a label set $\{1, 2, \cdots, k\}$. You were provided with a set of training examples $\mathcal{D} = \{(x, y^*)\}$, where $y^*$ is the correct structured output for the structured input $x$.

   You need to learn a scoring function $S(x, y) = w \cdot \phi(x, y)$, where $\phi(x, y) \in \Re^d$ is a joint feature representation over a structured input $x$ and candidate structured output $y \in Y(x)$ and $w \in \Re^d$ corresponds to the weights (or parameters) of the cost function. Essentially, you need to learn the weights $w \in \Re^d$ from the given training data.

---
**Algorithm 1** Randomized Greedy Search (RGS) Inference
---
**Input**: $x$ = structured input, $\phi$ = joint feature function, $w$ = weights of scoring function, $R$ = number of restarts
**Output**: $\hat{y}$, best scoring structured output
  1: Initialize the best scoring output $\hat{y}$ randomly
  2: Initialize the best score $S_{best}$ as $S(x, \hat{y}) = w \cdot \phi(x, \hat{y})$
  3: **for** $R$ iterations **do**
  4:    Pick a random structured output $y_{start}$ as starting point
  5:    **while** NOT reached local optima **do**
  6:      Perform a greedy search step by considering all one-label changes.
  7:      Score each candidate structured output (corresponding to one-label change) $y$:
           $S(x, y) = w \cdot \phi(x, y)$
  8:      Pick the best scoring candidate output and continue search
  9:      Update the best scoring output $\hat{y}$ and the corresponding best score $S_{best}$
 10:    **end while**
 11: **end for**
 12: **return** the best scoring output $\hat{y}$
---

**Algorithm 2** Online Structured Perceptron Training

---

**Input**: $\mathcal{D}$ = Training examples, $\phi$ = joint feature function, $R$ = number of restarts, $\eta$ = learning rate, $MAX$ = maximum training iterations

**Output**: $w$, weights of the scoring function

1: Initialize the weights of the scoring function $w = 0$
2: **for** MAX iterations or until convergence **do**
3:   **for** each training example $(x, y^*) \in \mathcal{D}$ **do**
4:     Make prediction: $\hat{y}$ = RGS-Inference$(x, \phi, w, R)$
5:     Compare $\hat{y}$ and $y^*$ to check for error
6:     If error, perform weight update:
       $w = w + \eta(\phi(x, y^*) - \phi(x, \hat{y}))$
7:   **end for**
8: **end for**
9: **return** weights $w$

---

(a) Implement the structured perceptron training algorithm and the randomized local search inference algorithm as shown in the above pseudo-code.

(b) Plot the Hamming accuracy over the training and testing set as a function of the number of online learning iterations (say MAX=100, $\eta = 0.01$, $R = 20$); for both handwriting and text-to-speech mapping problem.

(c) Repeat (b) with different feature representations $\phi$.
First-order: unary + pairwise features ($d = m \cdot k + k^2$)
Second-order: unary + pairwise + triple features ($d = m \cdot k + k^2 + k^3$)
Third-order: unary + pairwise + triple + quadruple features ($d = m \cdot k + k^2 + k^3 + k^4$)

(d) Fix the number of online learning iterations (MAX=100) and learning rate ($\eta = 0.01$) to reasonable values (based on (b)); and plot the Hamming accuracy for first-order representation as a function of the number of restarts $R$=10, 25, 50, 100, 200.

(e) How will you diagnose the performance of learning algorithm? Please list your ideas and explain your intuition and rationale. Implement your ideas to perform diagnosis and test your hypotheses.

(f) Please feel free to try other ways of randomizing the local search and list your observations in comparison to RGS.