

## Inference and Representation, Fall 2017

Structured Prediction for Part-of-Speech Tagging, MLE and Max-Ent, Max-Sum BP

Zhuoru Lin  
zlin@nyu.edu

---

*Disclaimer: I adhered to NYU honor code in this assignment.*

### 1. Part-of-speech tagging using SSVM

## **2. Max-Product**

### 3. Exponential families

The log-likelihood of the data is:

$$\log[p(x, \theta)] = \sum_{n=1}^L (\langle \theta, f(x^{(n)}) \rangle - \log[Z(\theta)]) \quad (1)$$

The maximum likelihood estimation  $\theta_{ML}$  must satisfied:

$$\frac{\partial p}{\partial \theta}(\theta_{ML}) = 0 \quad (2)$$

$$\sum_{n=1}^L f(x^{(n)}) - L \sum_x \frac{\exp[\langle \theta_{ML}, f(x^n) \rangle]}{Z(\theta_{ML})} f(x^n) = 0 \quad (3)$$

$$\frac{1}{L} \sum_{n=1}^L f(x^{(n)}) = \sum_x p(x|\theta_{ML}) f(x^n) \quad (4)$$

Technically, we also need to show that this is in fact an maximum by showing the second-derivative. This is ignored here.

4. **Maximum entropy distribution** The maximum entropy distribution can be formulated as an functional optimization problem as:

$$\begin{aligned} & \arg \max_p \int_x \log(p(x))p(x)dx \\ & \text{with } \int_x p(x)dx = 1 \\ & \text{and } \int p(x)f_k(x) = \alpha_k \text{ for all } k \end{aligned}$$

The Lagrangian is:

$$\mathcal{L} = \int_x \log(p(x))p(x)dx - \sum_k (\theta_k \int p(x)f_k(x) - \alpha_k) - (\lambda \int p(x)dx - 1). \quad (5)$$

The Lagrange multiplier states that the optimum achieves zero derivatives of Lagrangian. Taking the functional derivative, we have:

$$\log(p(x)) + 1 - \lambda - \sum_k \theta_k f_k(x) = 0 \quad (6)$$

Therefore we must have:

$$p(x) = \exp(\sum_k \theta_k f_k(x) + \lambda - 1) \quad (7)$$

$$= \frac{1}{Z} \exp(\langle \theta, f(x) \rangle) \quad (8)$$

with  $Z = \exp(1 - \lambda)$