

CSC 446 Notes: Lecture 16

1 Logistic Regression a.k.a. Maximum Entropy

$$P(y|x) = \frac{1}{Z_x} e^{\sum_i \lambda_i f_i(x,y)}$$

$$Z_x = \sum_y e^{\sum_i \lambda_i f_i(x,y)}$$

e.g., In NLP, we may use

$$f_{100}(x, y) = \begin{cases} 1 & \text{if } x=\text{word that ends with 'tion' and } y=\text{Noun,} \\ 0 & \text{if otherwise.} \end{cases}$$

The above is more general binary classification where we are only deciding whether an example belongs to a class or NOT. Here, we can have features contributing to multiple classes according to their weights. For binary classification, the decision boundary is linear, as with perceptron or SVM. A major difference from SVMs is that, during training, every example contributes to the objective function, whereas in SVMs only the examples close to the decision boundary matter.

If we plot this function, we get a sigmoid-like graph. We can draw analogy between maximum entropy and neural network, and consider features as the input nodes in the neural network.

If we take the log of equation (1), we get a linear equation

$$\log P = \sum_i \lambda_i f_i + c$$

What should λ be?

$$\begin{aligned} \max_{\lambda} \log \left(\prod_{n=1}^N P(y_n|x_n) \right)^{\frac{1}{N}} &= \max_{\lambda} \sum_n \frac{1}{N} \log \left(\frac{1}{Z_x} e^{\sum_i \lambda_i f_i} \right) \\ &= \max_{\lambda} \frac{1}{N} \sum_n \left(\sum_i \lambda_i f_i - \log Z_x \right) \end{aligned}$$

In the above equation, maximizing λ_i is easy, but maximizing Z_x is not. To maximize for λ_i , we turn it into a concave form and find the point where the derivative w.r.t. λ is zero (hill climbing)

$$L = \frac{1}{N} \sum_n \sum_i \lambda_i f_i - \log \sum_y e^{\sum_i \lambda_i f_i(x_n, y)} \quad (1)$$

$$\frac{\partial L}{\partial \lambda_j} = \frac{1}{N} \sum_n f_j - \frac{1}{Z_x} \frac{\partial}{\partial \lambda_j} \left(\sum_y e^{\sum_i \lambda_i f_i} \right) \quad (2)$$

$$= \frac{1}{N} \sum_n \left[f_j - \frac{1}{Z_x} \sum_y f_j e^{\sum_i \lambda_i f_i(y, x_n)} \right] \quad (3)$$

$$= \frac{1}{N} \sum_n f_j - \sum_y f_j P(y|x_n) \quad (4)$$

$$= \frac{1}{N} \sum_n f_j(x_n, y_n) - \sum_y f_j P(y|x_n) \quad (5)$$

We can morph eq. 5 into expectation form by defining joint probability as follows:

$$P(y, x) = P(y|x) \tilde{P}(x) \quad (6)$$

$$\tilde{P}(x) = \frac{1}{N} \sum_n I(x_n = x) \quad (7)$$

$$\tilde{P}(y|x) = \frac{c(x_n = x, y_n = y)}{c(x_n = x)} \quad (8)$$

Rewriting eq. 5 in expectation form, we get:

$$\frac{\partial L}{\partial \lambda_j} = E_{\tilde{P}}[f_j] - E_P[f_j] \quad (9)$$

where the first term (before the minus) is a constant, and the complexity of calculating the second term depends on the number of classes in the problem. Now we have:

$$\lambda \leftarrow \lambda + \eta \frac{\partial L}{\partial \lambda} \quad (10)$$

We will justify why we chose log linear form instead of something else. Assume we want to find the maximum entropy subject to constraints on the feature expectations:

$$\max_{P(y|x)} H(y|x) \quad (11)$$

$$\min_{P(y|x)} -H(y|x) \quad (12)$$

$$\text{s.t. } E_{\tilde{P}}[f_i] = E_P[f_i] \quad \forall i \quad (13)$$

$$\sum_y P(y|x) = 1 \quad \forall x \quad (14)$$

To put this into words, we want to build a model such that for each feature, our model should match the training data. We have

$$H(y|x) = \sum_{x,y} \tilde{P}(x) P(y|x) \log \frac{1}{P(y|x)} \quad (15)$$

Find the maximum entropy of the above equation as follows

$$L(P, \lambda, \mu) = f_0 + \sum_j \lambda_j f_j \quad (16)$$

$$= \sum_{x,y} \tilde{P}(x) P(y|x) \log P(y|x) \quad (17)$$

$$+ \sum_i \lambda_i \left(\sum_{x,y} \tilde{P}(x) \tilde{P}(y|x) f_i - \tilde{P}(x) P(y|x) f_i \right) \quad (18)$$

$$+ \sum_x \tilde{P}(x) \mu_x \left(\sum_y P(y|x) - 1 \right) \quad (19)$$

$$\frac{\partial L}{\partial P(y|x)} = \tilde{P}(x) (\log P(y|x) + 1) - \sum_i \lambda_i \tilde{P}(x) f_i + \tilde{P}(x) \mu_x = 0 \quad (20)$$

$$\log P(y|x) = -1 + \sum_i \lambda_i f_i - \mu_x \quad (21)$$

$$P(y|x) = e^{-1-\mu_x} e^{\sum_i \lambda_i f_i} \quad (22)$$

$$= \frac{1}{Z_x} e^{\sum_i \lambda_i f_i} \quad (23)$$

The above result shows that maximum entropy has log-linear form. If we solve the dual of the problem

$$g(\lambda, \mu) = E_P \left[\sum_i \lambda_i f_i - 1 - \mu_x \right] + E_{\tilde{P}} \left[\sum_i \lambda_i f_i \right] - E_P \left[\sum_i \lambda_i f_i \right] + E_P [\mu_x] - \sum_x \tilde{P}(x) \mu_x \quad (24)$$

$$= E_P [-1 - \mu_x] + E_{\tilde{P}} \left[\sum_i \lambda_i f_i \right] + E_P [\mu_x] - \sum_x \tilde{P}(x) \mu_x \quad (25)$$

$$= E_P [-1] + E_{\tilde{P}} \left[\sum_i \lambda_i f_i \right] - \sum_x \tilde{P}(x) \mu_x \quad (26)$$

$$= - \sum_{x,y} \tilde{P}(x) e^{-1-\mu_x} e^{\sum_i \lambda_i f_i} + E_{\tilde{P}} \left[\sum_i \lambda_i f_i \right] - \sum_x \tilde{P}(x) \mu_x \quad (27)$$

$$(28)$$

Solving analytically for μ_x that maximizes g :

$$0 = \frac{\partial g}{\partial \mu_x} = \tilde{P}(x) \left(\sum_y -e^{-1-\mu_x+\sum_i \lambda_i f_i} - 1 \right) \quad (29)$$

$$\mu_x = \log \sum_y e^{\sum_i \lambda_i f_i} - 1 \quad (30)$$

Substituting μ_x into g :

$$g(\lambda, \mu) = - \sum_{x,y} \tilde{P}(x) \frac{1}{\sum_y e^{\sum_i \lambda_i f_i}} e^{\sum_i \lambda_i f_i} + E_{\tilde{P}} \left[\sum_i \lambda_i f_i \right] - \sum_x \tilde{P}(x) \left(\log \sum_y e^{\sum_i \lambda_i f_i} - 1 \right) \quad (31)$$

$$= -1 + E_{\tilde{P}} \left[\sum_i \lambda_i f_i \right] - \sum_x \tilde{P}(x) \log \left(\sum_y e^{\sum_i \lambda_i f_i} \right) + 1 \quad (32)$$

$$= E_{\tilde{P}} \left[\sum_i \lambda_i f_i \right] - \sum_x \tilde{P}(x) \log \left(\sum_y e^{\sum_i \lambda_i f_i} \right) \quad (33)$$

$$= E_{\tilde{P}} \left[\sum_i \lambda_i f_i - \log \left(\sum_y e^{\sum_i \lambda_i f_i} \right) \right] \quad (34)$$

$$= E_{\tilde{P}} \left[\sum_i \lambda_i f_i - \log Z_x \right] \quad (35)$$

$$= E_{\tilde{P}} [\log P(y|x)] \quad (36)$$

$$= L \quad (37)$$

Thus solving the dual of the entropy maximization problem consists of maximizing the likelihood of the training data with a log-linear functional form for $P(y|x)$.

Phyo Thiha 4/12; DG 4/13