

CS460/626 : Natural Language Processing/Speech, NLP and the Web

Lecture 34: A very useful Maximum Likelihood Function

Pushpak Bhattacharyya
CSE Dept.,
IIT Bombay

11th Nov, 2012

Observed variable

- #Observation = N

$$X : \langle X_1, X_2, X_3 \dots X_N \rangle$$

Each x_i is a categorical distribution of k outcomes

$$X_i = \langle X_{i1}, X_{i2}, X_{i3}, \dots, X_{ij}, \dots X_{iK} \rangle$$

$$x_{ij} \in \{0,1\}, \sum_{j=1}^M x_{ij} = 1$$

Hidden variable

- Each observation from M ‘sources’, giving rise to M values form a ‘categorised distribution’ for unobserved variables.

$$Z_i = \langle z_{i1}, z_{i2}, z_{i3}, \dots, z_{ij}, \dots z_{im} \rangle$$

$$z_{ij} \in \{0,1\}, \sum_{j=1}^M z_{ij} = 1$$

Variables: The complete picture

- Observed

$X : X_1, X_2, X_3, \dots, X_i, \dots, X_N$

- Unobserved

$Z : Z_1, Z_2, Z_3, \dots, Z_i, \dots, Z_N$

- Complete data:

– (X, Z)

Parameters

- π_i = prob. of choosing the i th ‘source’

$$\sum_{i=1}^M \pi_i = 1$$

- $p_{i1}, p_{i2}, p_{i3}, \dots, p_{ij}, \dots, p_{iK}$ are probabilities of the ‘outcomes’ from the ‘ j th’ source

$$\sum_{l=1}^K p_{jl} = 1$$

Example

- Sources are 10 dice
- Each dice has 6 outcomes
- $M = 10$
- $K = 6$
- Suppose #observations = 20
 - then $N = 20$

Likelihood formulation

- Maximum likelihood of complete data

$$l(\theta) = P(X, Z : \theta)$$

for the i^{th} item

$$P(X_i, Z_i : \theta) = \prod_{j=1}^M \left[\pi_j \left(\prod_{k=1}^K P_{jk}^{x_{ik}} \right) \right]^{z_{ij}}$$

$$p_{j1}, p_{j2}, p_{j3}, \dots, p_{jl}, \dots, p_{jK}$$

$$X_i = \langle x_{i1}, x_{i2}, x_{i3}, \dots, x_{ij}, \dots, x_{iK} \rangle$$

$$Z_i = \langle z_{i1}, z_{i2}, z_{i3}, \dots, z_{ij}, \dots, z_{iK} \rangle$$

Bernoulli like trial

$$P(x_i, z_i : \theta) = \prod_{j=1}^M \pi_j \left(p_{j1}^{x_{i1}} p_{j2}^{x_{i2}} \cdots p_{jk}^{x_{ik}} \right)^{z_{ij}}$$

MLE

$$l(\theta) = P(X, Z : \theta) = \prod_{i=1}^N \prod_{j=1}^M \pi_j \left(\prod_{k=1}^K p_{jk}^{x_{ik}} \right)^{z_{ij}}$$

Variable definition

- i goes over observation
- j goes over source/observation
- k goes over outcome/source/observation

Log likelihood of MLE

$$ll(\theta) = \sum_{i=1}^N \sum_{j=1}^M \log \pi_j + z_{ij} \sum_{k=1}^K x_{ik} \log p_{jk}$$

Already proved:-

We have to maximize expectation wrt z of LL(θ)

[Consequence marginalization of $P(X: \theta)$ wrt z]

Expectation of log likelihood

$$E_z(LL(\theta))$$

Now

$$E_z(f(z)) = f(E(z)) \text{ if } f \text{ is linear in } z$$

$$E_z(f(z)) \stackrel{\Delta}{=} \sum_z P(z)f(z) = f\left[\sum_z P(z)z\right]$$

Optimization of expectation wrt constraints

$$E_z(LL(\theta)) = \sum_{i=1}^N \left[\sum_{j=1}^M E(z_{ij}) \left(\log \pi_j + \sum_{k=1}^K x_{ik} \log P_{jk} \right) \right]$$

Maximize $E_z(LL(\theta))$ subject to the constraints

$$\sum_{j=1}^M \pi_j = 1 \quad (1)$$

$$\sum_{k=1}^K \log P_{ik} = 1 \quad (2)$$

$$\sum_{j=1}^M z_{ij} = 1 \quad (3)$$

Introduction of Lagrange multiplier

$$\lambda_1 \left(\sum_{j=1}^M \pi_j - 1 \right) \quad (\text{A})$$

$$\sum_{i=1}^N \lambda_{2i} \left(\sum_{j=1}^M z_{ij} - 1 \right) \quad (\text{B})$$

$$\sum_{i=1}^N \lambda_{3i} \left(\sum_{k=1}^K p_{ik} - 1 \right) \quad (\text{C})$$

Maximization and expectation step

M - step

Maximize

$E_z(LL(\theta)) + (A) + (B) + (C)$ with respect to π_j, p_{jk}

$j = 1, 2, \dots, M$

$k = 1, 2, \dots, K$

E - step

$$E(z_{ij}) = \frac{\pi_j \left(\prod_{k=1}^K P_{jk}^{x_{ik}} \right)^{z_{ij}}}{\sum_{j=1}^M \pi_j \left(\prod_{k=1}^K P_{jk}^{x_{ik}} \right)^{z_{ij}}}$$

Two coin example

$$P(z = z_i \mid x = x_i)$$

$$E(z_i) = \frac{P \times P_1^{x_i} (1 - P_1)^{(1-x_i)}}{P \times P_1^{x_i} (1 - P_1)^{(1-x_i)} + (1 - P) \times P_2^{x_i} (1 - P_2)^{(1-x_i)}}$$