

Maximum Entropy and Maximum Likelihood

Sargur Srihari

srihari@cedar.buffalo.edu

Topics

- Maximum Entropy Formulation of Markov Network Parameter Estimation
- Motive for alternative formulation
 - Maximum Entropy subject to Constraints

Maximum Entropy Formulation of Markov Network Parameter Estimation

- Basic formulation is: maximum likelihood parameter estimation
- Alternative formulation provides insight
- Insight is that

“Distribution that maximizes entropy subject to constraints (of expectation) is the Gibbs distribution”

Maximizing entropy subject to constraints is dual of maximizing likelihood given structural constraints

Entropy Definition

- Entropy is a measure of information
 - It is the *Expected value* (average value) of information generated by the distribution
 - Where Information is (log) reciprocal of probability
 - For a random variable X with distribution $P(X)$

$$H_p(X) = E_p \left[\log \frac{1}{P(X)} \right] = \sum_x P(x) \log \frac{1}{P(X)}$$

- where we treat $0 \bullet \log 1/0 = 0$
- Entropy is the Lower bound on the average number of bits needed to encode values of X
 - Central Theorem of Information Theory by Shannon
 - i.e., Expected code length cannot be less than $H_p(X)$

Motivation of Max Entropy with Constraints

- Suppose we have some summary statistics of an empirical distribution
 - Ex 1: published in a census report
 - Marginal distributions of single variables or certain pairs, or other events of interest
 - Ex 2: average final grade of students in class (70), correlation of final grades with homework scores
 - (correlation coefficient, $r = 0.6$, on scale of -1 to +1)
- But no access to full data set
- Find a typical distribution that satisfies these constraints and use it to answer other queries⁵

Maximum Entropy Formulation

- Goal is to select a distribution that:
 1. Satisfies given constraints
 2. Has no additional structure or information
- We formulate the following problem:
 - Find $Q(\mathcal{X})$
 - Maximizing $H_Q(\mathcal{X})$
 - Subject to $E_Q[f_i] = E_D[f_i], i=1, \dots, k$
- These are called *expectation constraints*

Solution to Max Entropy Estimation

- Somewhat surprisingly
- The distribution Q^* , the maximum entropy distribution satisfying the expectation constraints, $E_Q[f_i] = E_D[f_i]$, is a Gibbs distribution $Q^* = P_{\hat{\theta}}$ where

$$P_{\hat{\theta}}(\chi) = \frac{1}{Z(\hat{\theta})} \exp \left\{ \sum_i \hat{\theta}_i f_i(\chi) \right\}$$

where $\hat{\theta}$ is the m.l.e. parameterization relative to \mathcal{D}

Duality of Max Ent and Max Likelihood

- The two problems
 1. *Maximizing entropy* subject to expectation constraints
 2. *Maximizing likelihood* given structural constraints on the distribution– are *convex duals of each other*
- For the maximum likelihood parameters $\hat{\theta}$

$$H_{P_{\hat{\theta}}}(\chi) = -\frac{1}{M} \ell(\hat{\theta} : D)$$

Convex Duality

- Optimization problems are viewed as
- Primal Problem
 - E.g., Lagrangian form with specified constraints
- Dual Problem
 - Provides lower bounds to the primal problem
- In general optimal values are not the same
- For convex optimization the duality gap is zero