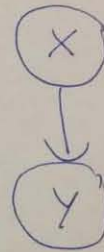


6.438: Recitation - 12

EM-algorithm:

I. Problem Setting: We have some hidden variables which we call X (can be a vector) and some observed variables which we call Y .

G:



We have a known graph structure for variables X and Y , such that the joint distribution of X and Y factorizes according to the graph. Thus the joint probability of X and Y can be written as a function of the (unknown) graph parameters θ .

$$\textcircled{1} P(X, Y) \equiv P(X, Y; \theta)$$

For simplicity, we look at the special case where $\theta = (\theta_x; \theta_{Y|X})$, so that the joint probability factorizes as:

$$\textcircled{2} P_{X,Y}(x, y; \theta) = P_X(x; \theta_x) \cdot P_{Y|X}(y|x; \theta_{Y|X})$$

②

In other words, the variables X and Y satisfy the graph structure of G . Note that many many popular models, such as Naive Bayes and HMMs, satisfy this assumption.

II. Objective: We are given a large number of samples of the Y variable, which we represent as $\{y^{(i)}\}$, $i \in \{1, 2, \dots, S\}$.

These samples provide us with a ~~chara~~ characterization of the marginal distribution of Y , i.e. $P_Y(\cdot)$.

Our aim is to find parameters $\hat{\theta} = (\hat{\theta}_x, \hat{\theta}_{Y|X})$ which produce a marginal distribution for Y close to what we obtain from samples.

More concretely, we look at the following ML objective: find the value $\hat{\theta}$ of θ which maximizes the probability of generating the observed samples. That is,

$$\begin{aligned} \textcircled{3} \quad \hat{\theta}_{ML} &\triangleq \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^S P_Y(y^{(i)}; \theta) \\ &= \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^S \left(\sum_x P_{X,Y}(x, y^{(i)}; \theta) \right) \end{aligned}$$

$$\textcircled{4} \therefore \hat{\Theta}_{ML} = \underset{\Theta}{\operatorname{argmax}} \prod_{i=1}^S \left(\sum_x P_x(x; \theta_x) \cdot P_{y|x}(y^{(i)} | x; \theta_{y|x}) \right)$$

It should be noted that $\textcircled{4}$ is a difficult problem to solve in general. We look at one way of solving it approximately viz. EM-algorithm

III. EM Algorithm:

The EM algorithm is a meta-algorithm for estimating graphical model parameters from partial observations viz solving problem $\textcircled{4}$. Since the equations are already derived in the lecture notes, we will go over them quickly.

Consider problem $\textcircled{4}$, which we ~~can~~ write in a slightly shorter form.

$$\hat{\Theta}_{ML} = \underset{\Theta}{\operatorname{argmax}} \prod_{i=1}^S \left(\sum_x P_{x,y}(x, y_i^{(i)}; \theta_{x,y}) \right)$$

$$= \underset{\Theta}{\operatorname{argmax}} \sum_{i=1}^S \log \left(\sum_x P_{x,y}(x, y_i^{(i)}; \theta_{x,y}) \right)$$

(Taking logarithms, since $\log(\cdot)$ is a monotone function)

(4)

Now consider any set of S distributions over X , which we denote as $q_{x|y}(\cdot|y^{(i)})$
 $\forall i \in \{1, 2, \dots, S\}$

Note that this distribution is defined as a function of the sample $y^{(i)}$. Then, we can write our objective as:

$$\begin{aligned} \hat{\theta}_{ML} &= \operatorname{argmax}_{\theta} \sum_{i=1}^S \log \left(\sum_x q_{x|y}(x|y^{(i)}) \cdot \frac{P_{x,y}(x, y^{(i)}; \theta_{x,y})}{q_{x|y}(x|y^{(i)})} \right) \\ &= \operatorname{argmax}_{\theta} \sum_{i=1}^S \log \left(E_{q_{x|y}(\cdot|y^{(i)})} \left[\frac{P_{x,y}(x, y^{(i)}; \theta_{x,y})}{q_{x|y}(x|y^{(i)})} \right] \right) \\ &= \operatorname{argmax}_{\theta} f(\theta), \end{aligned}$$

$$\text{where } f(\theta) = \sum_{i=1}^S \log \left(E_{q_{x|y}(\cdot|y^{(i)})} \left[\frac{P_{x,y}(x, y^{(i)}; \theta_{x,y})}{q_{x|y}(x|y^{(i)})} \right] \right)$$

By Jensen's inequality, since $\log(\cdot)$ is concave,

$$\textcircled{5} \quad f(\theta) \geq \sum_{i=1}^S E_{q_{x|y}(\cdot|y^{(i)})} \left[\frac{P_{x,y}(x, y^{(i)}; \theta_{x,y})}{q_{x|y}(x|y^{(i)})} \right]$$

with equality iff $\frac{P_{x,y}(x, y^{(i)}; \theta_{x,y})}{q_{x|y}(x|y^{(i)})} = \text{constant}$,
 [Not a function of x]

$$\textcircled{6} \quad \text{i.e. iff } q_{x|y}(x|y^{(i)}) \propto P_{x,y}(x, y^{(i)}; \theta_{x,y})$$

(5)

[The ~~equality~~ last statement follows from the observation that $q_{x|y}^{(i)}(x|y^{(i)}) \propto P_{X,Y}(x, y^{(i)})$, and $q_{x|y}(\cdot)$ is a probability distribution]

We are now ready to give the EM algorithm. Let $g(\theta, q)$ be the RHS in equation (5)

$$\text{i.e. } g(\theta, q) = \sum_{i=1}^S \mathbb{E}_{q_{x|y}(\cdot|y^{(i)})} \log \left[\frac{P_{X,Y}(x, y^{(i)}, \theta_{X,Y})}{q_{x|y}(x|y^{(i)})} \right]$$

As we know, $g(\theta, q)$ is a lower bound on the true objective function viz. $f(\theta)$. We also know that $\max_{q(\cdot|y)} g(\theta, q) = f(\theta)$, achieved when $q(\cdot|y^{(i)}) = P_{X,Y}(\cdot|y^{(i)})$. EM algorithm has 2 steps, one where we maximize over $q(\cdot|y^{(i)})$, and the other where we maximize over θ . It is an ~~EM~~ iterative algorithm to maximize the ~~EM~~ step: objective function $g(\theta, q)$.

[For those interested, this procedure of iteratively maximizing a function with 2 sets of parameters, by alternately holding one set of parameters fixed and maximizing over the other, is called 'alternating maximization!']

(6)

EM in a nutshell:

Objective function: $g(\theta, q_{x|y}(\cdot|y^{(i)}))$

$$= \sum_{i=1}^S E_{q_{x|y}(\cdot|y^{(i)})} \log \left[\frac{P_{x,y}(x, y^{(i)}; \theta_{x,y})}{q_{x|y}(x|y^{(i)})} \right]$$

Initialization: Assign values to parameters $\theta_{x,y}$, call it $\hat{\theta}_{x,y}^{(0)}$.

E-step: Given $\hat{\theta}_{x,y}^{(t)}$, set

$$q_{x|y}^{(t+1)}(\cdot|y^{(i)}) = \operatorname{argmax}_{q_{x|y}(\cdot|y^{(i)})} g(\hat{\theta}_{x,y}^{(t)}, q_{x|y}(\cdot|y^{(i)})) \quad \forall i \in \{1, \dots, S\}$$

From Equation (6), we already know this optimal value viz,

$$q_{x|y}^{(t+1)}(\cdot|y^{(i)}) = P_{x|y}(\cdot|y^{(i)}; \hat{\theta}_{x,y}^{(t)}) \quad \forall i \in \{1, 2, \dots, S\}$$

M-step: Given $q_{x|y}^{(t+1)}(\cdot|y^{(i)})$, set

$$\hat{\theta}_{x,y}^{(t+1)} = \operatorname{argmax}_{\theta} g(\theta, q_{x|y}^{(t+1)}(\cdot|y^{(i)}))$$

We can get a slightly simplified expression by leaving out the term in the denominator of the $\log(\cdot)$ in $g(\theta, q)$, since it does not depend on θ .

$$\hat{\theta}_{x,y}^{(t+1)} = \operatorname{argmax}_{\theta_{x,y}} \sum_{i=1}^S E_{q_{x|y}^{(t+1)}(\cdot|y^{(i)})} \log(P_{x,y}(x, y^{(i)}; \theta_{x,y}))$$

IV. Solution to EM via Sampling ^⑦

EM is, strictly speaking, a meta-algorithm. It converts our original hard problem into 2 simpler problems viz. E-step and M-step, but does not tell us how to solve these steps.

One fairly general method (but not necessarily efficient) of solving EM subproblems is via Sampling. Recall the M-step of EM:

$$\textcircled{7} \quad \hat{\theta}_{x,y}^{(t+1)} = \operatorname{argmax}_{\theta_{x,y}} \sum_{i=1}^S E_{\substack{\gamma^{(t+1)} \\ \gamma_{x|y}(\cdot|y^{(i)})}} \log (P_{x,y}(x,y^{(i)}; \theta_{x,y}))$$

Note that the E-step simply sets $\gamma_{x|y}^{(t+1)}(\cdot|y^{(i)})$ to $P_{x|y}(x|y^{(i)}; \hat{\theta}_{x,y}^{(t)})$, so we do not write it explicitly. Thus, our aim is to solve M-step where q is described as above.

We will do this by generating a large number of samples of x , for each observed sample $y^{(i)}$ $\forall i \in \{1, 2, \dots, S\}$. These samples of x , when juxtaposed with the original samples of y , will give us a large number of samples of (x, y) , which are representative of the original distribution $P_{x,y}(x, y; \hat{\theta}_{x,y}^{(t)})$. Then, our task

(8)

of determining $\hat{\theta}^{(t+1)}$ simply boils down to finding Maximum Likelihood estimates in the fully observed setting, a problem which we have already solved.

* Note: It seems difficult to give a formal definition of the quantity $P_{x,y}(x,y; \hat{\theta}^{(t)}, D)$. Thus, it is for intuition only.

Let us do the above formally. For each observation $y^{(i)} \forall i \in \{1, 2, \dots, S\}$, we generate N_i samples of x from the distribution $q_{x|y}^{(t+1)}(\cdot | y^{(i)})$. Call these samples $\{x^{(i,1)}, x^{(i,2)}, \dots, x^{(i,N_i)}\}$.

If N_i is sufficiently large, then we can write

$$E_{q_{x|y}^{(t+1)}(\cdot | y^{(i)})} (\log(P_{x,y}(x, y^{(i)}; \theta_{x,y}))) \\ \approx \frac{1}{N_i} \sum_{j=1}^{N_i} \log(P_{x,y}(x^{(i,j)}, y^{(i)}; \theta))$$

The above statement merely says that the empirical expectation approximately equals the true expectation, when N_i is large. This is basically the Law of Large Numbers (LLN).

(9)

Substituting the above in equation (7), we get

$$\hat{\theta}^{(t+1)} = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^S \sum_{j=1}^{N_i} \log(P_{x,y}(x^{(i,j)}, y^{(i)}, \theta))$$

Now let $N_i = N \forall i$, i.e. generate the same number of x samples for each $y^{(i)}$. Then, the above expression reduces to:

$$\begin{aligned} \hat{\theta}^{(t+1)} &= \underset{\theta}{\operatorname{argmax}} \frac{1}{N} \sum_{i=1}^S \sum_{j=1}^N \log(P_{x,y}(x^{(i,j)}, y^{(i)}, \theta)) \\ &= \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^S \sum_{j=1}^N \log(P_{x,y}(x^{(i,j)}, y^{(i)}, \theta)) \end{aligned}$$

$$(8) \quad \hat{\theta}^{(t+1)} = \underset{\theta}{\operatorname{argmax}} \sum_{(x,y) \in \text{"Samples"}} \log(P_{x,y}(x, y; \theta))$$

Where "Samples" is the set of all (x, y) -samples we have generated so far, suitably put together.

That is,

$$\text{Samples} = \left\{ \begin{array}{l} \{x^{(1,1)}, y^{(1)}\}, \{x^{(1,2)}, y^{(1)}\}, \dots, \{x^{(1,N)}, y^{(1)}\} \\ \vdots \\ \{x^{(S,1)}, y^{(S)}\}, \{x^{(S,2)}, y^{(S)}\}, \dots, \{x^{(S,N)}, y^{(S)}\} \end{array} \right\}$$

Now we make the striking observation - problem (8) is simply the problem of finding ML estimate of θ , when the true data given to us is "Samples"!

(10)

For the particular case of directed graphical models, the ML estimates are especially easy to compute. The parameters (θ) are simply conditional probabilities, and their ML estimates are empirical conditional probabilities obtained from the data. (in this case, data being the "Samples" set)

Only 1 important question remains viz
How to sample?

Sampling from $q_{x|y}(y^{(i)})$:

The first thing to note is that given any pair of values (x, y) of the random variables (X, Y) , we can calculate the probability of this ~~value~~ pair of values, given the parameters θ i.e. $P_{x,y}(x, y) = P_{x,y}(x, y; \theta_{x,y})$

[The above assumes a directed G.M.; if we had an undirected G.M., we could only calculate probabilities upto a normalization constant. But even that will be good enough]

(9) Since $\hat{q}_{x|y}^{(t+1)}(\cdot | y^{(i)}) = p_{x|y}(x | y^{(i)}; \hat{\theta}_{x,y}^{(t)})$ (11)

$$\propto p_{x,y}(x, y^{(i)}; \hat{\theta}_{x,y}^{(t)}),$$

we know $q(\cdot | y^{(i)})$ upto a normalization constant. That is, given any ^{possible} value x of r.v. X , we can calculate $\hat{q}^*(x) = \frac{\hat{q}(x | y)}{Z}$, where Z is some normalization constant.

Q. How do we sample from a distribution that we can calculate upto a normalization constant?

There are 2 methods we have learnt to handle this situation:

1. MCMC.
2. Importance Sampling: This does not actually generate samples from $q(\cdot)$, but allows us to compute expectation of any function w.r.t. $q(\cdot)$, which is almost always our final aim.

For simplicity, we focus on (1) method 1 i.e. MCMC. In the metropolis-hastings rule, all we needed to create a M.C. with stationary distribution $q(\cdot)$, (Markov Chain)

(12)

was the ratio of probabilities of 2 states
i.e. $\frac{\hat{q}(x_b)}{\hat{q}(x_a)}$, where x_b, x_a are any 2 possible
values of X .

But this is simply equal to $\frac{\hat{q}^*(x_b)}{\hat{q}^*(x_a)}$, which we
can calculate using \hat{q} equation (9).

This shows that we can easily use MCMC
to generate our samples of X corresponding
to each $y^{(i)}$, which we needed in our
algorithm. This concludes our discussion
of Sampling.

Exercise: Can you use Gibbs Sampling for this
task?

V. Solution to EM via Sum-Product (for HMMs)

If you have understood the previous section,
this should be easy to follow. The key
idea here is ~~to~~ that for certain well-
structured graphical models, we can compute
all the empirical probabilities directly,
which we would obtain from sampling.

MIT OpenCourseWare
<http://ocw.mit.edu>

6.438 Algorithms for Inference
Fall 2014

For information about citing these materials or our Terms of Use, visit: <http://ocw.mit.edu/terms>.