

# EM Algorithm for Labeled IRL

James MacGlashan

## 1 Introduction

In this document, I explored the EM algorithm approach to solving our labeled IRL problem. As a review, I first include the description of our model that we've previously defined and then provide the EM algorithm approach.

## 2 Probability Model

The motivation for this model is that we consider the the final label ( $L$ ) that a user gives a trajectory of size  $N$  to be some random function of what they thought about each of the action selections ( $A$ ) exhibited in the trajectory. However these step-wise evaluations ( $X$ ) are unobserved in the data.

We model the probability that an action is evaluated as good or not as proportional to its selection probability according a softmax policy computed for the reward function with parameters  $\theta$ . Specifically:

$$\Pr(x_i = +1|s, a, \theta) = \pi(s, a|\theta) \quad (1)$$

$$\Pr(x_i = -1|s, a, \theta) = 1 - \pi(s, a|\theta), \quad (2)$$

where  $\pi(s, a|\theta)$  is the softmax policy over Q-values computed for the reward function parameterized by  $\theta$ :

$$\pi(s, a|\theta) = \frac{e^{\beta Q(s, a|\theta)}}{\sum_{a'} e^{\beta Q(s, a'|\theta)}}, \quad (3)$$

$\beta$  is a selectable parameter, and  $Q(s, a|\theta)$  is the Q-value computed for the reward function parameterized by  $\theta$ .

For the probability distribution of  $L$ , given the sequence of  $N$  step-wise labels, we would like a distribution that has the property that as more step-wise labels are positive, the probability of a positive trajectory label increases (and vice versa). Although there are many possible distributions that satisfy this property, for concreteness, we choose the sigmoid function. That is,

$$\Pr(L = +1|X_1, \dots, X_n) = \frac{1}{1 + e^{-\phi \sum_i^N X_i}} \quad (4)$$

$$\Pr(L = -1|X_1, \dots, X_n) = 1 - \Pr(L = +1|X_1, \dots, X_n), \quad (5)$$

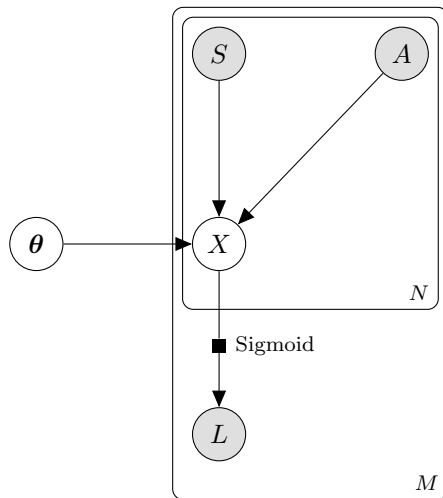


Figure 1: Plate Diagram of our Probability Model

where  $\phi$  is a selectable parameter that tunes how quickly of a majority of step-wise labels increases/decreases the probability of trajectory assignment. For example, when  $\phi = 0$ , trajectory labels are assigned uniformly randomly independently of step-wise labels. As  $\phi \rightarrow \infty$ , the sigmoid converges to a step function in which a trajectory containing even one more positive step-wise label than negative step-wise labels will deterministically cause a positive trajectory label (and vice versa).

### 3 EM Algorithm

The problem with estimating the  $\theta$  parameters of our reward function is that we have a latent variable vector  $X$  (or rather, some of the elements of the  $X$  vector are latent, and some may be observed), which prevents us from easily computing the likelihood of the model and maximizing parameters for it. The EM approach to solving this problem is to first choose values for  $\theta$ ; then choose a new  $\theta$  that maximizes the expected value of the log likelihood function where the distribution of the expectation is the probability distribution of latent variables (the  $X$ s in our case) given the observations available and previous  $\theta$  choice; and then repeating this process. The maximization process can be performed using gradient ascent.

To formalize this process for our problem, first note that the likelihood our parameters (and state-action sequence) given an  $\mathbf{x}$  vector and label  $l$  is

$$\mathcal{L}(\mathbf{s}, \mathbf{a}, \theta | l, \mathbf{x}) = \Pr(l | \mathbf{x}) \prod_i \Pr(x_i | s_i, a_i, \theta) \quad (6)$$

and the log likelihood is

$$\log \mathcal{L}(\mathbf{s}, \mathbf{a}, \boldsymbol{\theta} | l, \mathbf{x}) = \log \Pr(l | \mathbf{x}) + \sum_i \log \Pr(x_i | s_i, a_i, \boldsymbol{\theta}). \quad (7)$$

Additionally, the gradient of the log likelihood is,

$$\nabla_{\boldsymbol{\theta}} \log \mathcal{L}(\mathbf{s}, \mathbf{a}, \boldsymbol{\theta} | l, \mathbf{x}) = \sum_i \frac{\nabla_{\boldsymbol{\theta}} \Pr(x_i | s_i, a_i, \boldsymbol{\theta})}{\Pr(x_i | s_i, a_i, \boldsymbol{\theta})}. \quad (8)$$

To simplify the EM algorithm description, I will introduce the notation  $\mathbf{x}_k$  to indicate the subset of observed elements in an  $\mathbf{x}$  vector, and  $\mathbf{x}_u$  to represent a possible assignment to the subset of the unobserved values of an  $\mathbf{x}$  vector. Using this notation, the expected value of the log likelihood under some candidate parameter  $\boldsymbol{\theta}'$  for missing  $X$  elements distributed according to  $\boldsymbol{\theta}$  is

$$\begin{aligned} E_{\mathbf{x}_u \sim \Pr(\mathbf{x}_u | l, \mathbf{x}_k, \mathbf{s}, \mathbf{a}, \boldsymbol{\theta})} [\log \mathcal{L}(\mathbf{s}, \mathbf{a}, \boldsymbol{\theta}' | l, \mathbf{x})] \\ = \sum_{\mathbf{x}_u} \Pr(\mathbf{x}_u | l, \mathbf{x}_k, \mathbf{s}, \mathbf{a}, \boldsymbol{\theta}) \log \mathcal{L}(\mathbf{s}, \mathbf{a}, \boldsymbol{\theta}' | l, \mathbf{x}) \end{aligned}$$

Given that, the corresponding EM algorithm operating on a single trajectory is as follows (it is generalized to many trajectories by simply summing over each trajectory).

---

**Algorithm 1** Labeled-IRL EM Algorithm

---

**Require:** initial  $\boldsymbol{\theta}_0$ , and data  $\mathbf{s}, \mathbf{a}, \mathbf{x}_k, l$

**for**  $t = 0$  to  $K$  **do**

$\boldsymbol{\theta}_{t+1} \leftarrow \arg \max_{\boldsymbol{\theta}'} \sum_{\mathbf{x}_u} \Pr(\mathbf{x}_u | l, \mathbf{x}_k, \mathbf{s}, \mathbf{a}, \boldsymbol{\theta}_t) \log \mathcal{L}(\boldsymbol{\theta}', \mathbf{s}, \mathbf{a} | l, \mathbf{x}_k, \mathbf{x}_u)$

**end for**

---

To compute the expected value, we need to enumerate each of the possible assignments to the unknown  $\mathbf{x}$  elements and compute the probability of them given the observed data and model parameters  $\boldsymbol{\theta}$ . This probability is computed as

$$\begin{aligned} \Pr(\mathbf{x}_u | l, \mathbf{x}_k, \mathbf{s}, \mathbf{a}, \boldsymbol{\theta}) &= \frac{\Pr(l | \mathbf{x}_k, \mathbf{x}_u) \Pr(\mathbf{x}_u | \mathbf{s}, \mathbf{a}, \boldsymbol{\theta}) \Pr(\mathbf{x}_k | \mathbf{s}, \mathbf{a}, \boldsymbol{\theta}) \Pr(\mathbf{s}, \mathbf{a}, \boldsymbol{\theta})}{\Pr(l | \mathbf{x}_k, \mathbf{s}, \mathbf{a}, \boldsymbol{\theta}) \Pr(\mathbf{x}_k | \mathbf{s}, \mathbf{a}, \boldsymbol{\theta}) \Pr(\mathbf{s}, \mathbf{a}, \boldsymbol{\theta})} \\ &= \frac{\Pr(l | \mathbf{x}_k, \mathbf{x}_u) \Pr(\mathbf{x}_u | \mathbf{s}, \mathbf{a}, \boldsymbol{\theta})}{\Pr(l | \mathbf{x}_k, \mathbf{s}, \mathbf{a}, \boldsymbol{\theta})} \\ &= \frac{\Pr(l | \mathbf{x}_k, \mathbf{x}_u) \prod_i \Pr(x_{u,i} | s_i, a_i, \boldsymbol{\theta})}{\Pr(l | \mathbf{x}_k, \mathbf{s}, \mathbf{a}, \boldsymbol{\theta})}. \end{aligned}$$

A straight forward computation of  $\Pr(l | \mathbf{x}_k, \mathbf{s}, \mathbf{a}, \boldsymbol{\theta})$  requires marginalizing over all possible assignment to the unknown  $X$  elements; however, it can be computed efficiently using dynamic programming because the probability of a label is a sigmoid function that operates on the sum of the  $x$  elements. I will not cover the dynamic programming computation of it here.

Unfortunately, even with an efficient means to compute  $\Pr(l|\mathbf{x}_k, \mathbf{s}, \mathbf{a}, \boldsymbol{\theta})$ , when the number of unknown  $X$  variables is large, the number of  $\mathbf{x}_u$  assignments enumerated in the expectation's outer sum grows exponentially, and the product series over each of unknown element probabilities in the above equation ( $\prod_i \Pr(\mathbf{x}_{u,i}|s_i, a_i, \boldsymbol{\theta})$ ) can have underflow issues. A resolution to this problem is to estimate the expectation with sampling. Monte Carlo sampling is unfortunately intractable because it is not easy to sample from  $\Pr(\mathbf{x}_u|l, \mathbf{x}_k, \mathbf{s}, \mathbf{a}, \boldsymbol{\theta})$ ; moreover, it would not address the underflow issue in the product series. However, it is easy to sample from  $\Pr(\mathbf{x}_u|\mathbf{s}, \mathbf{a}, \boldsymbol{\theta})$  (removing the conditioning on the label), which we can use in importance sampling. With importance sampling, we can replace the expectation computation with the sample-estimate

$$\frac{1}{C} \sum_j^C \frac{\Pr(\mathbf{x}_u^j|l, \mathbf{x}_k, \mathbf{s}, \mathbf{a}, \boldsymbol{\theta})}{\Pr(\mathbf{x}_u^j|\mathbf{s}, \mathbf{a}, \boldsymbol{\theta})} \log \mathcal{L}(l, \mathbf{x}_k, \mathbf{x}_u^j|\mathbf{s}, \mathbf{a}, \boldsymbol{\theta}). \quad (9)$$

where  $\mathbf{x}_u^j$  is sample from the distribution  $\Pr(\mathbf{x}_u|\mathbf{s}, \mathbf{a}, \boldsymbol{\theta})$ .

Note that this simplifies a bit further too:

$$\begin{aligned} \frac{\Pr(\mathbf{x}_u^j|l, \mathbf{x}_k, \mathbf{s}, \mathbf{a}, \boldsymbol{\theta})}{\Pr(\mathbf{x}_u^j|\mathbf{s}, \mathbf{a}, \boldsymbol{\theta})} &= \frac{\Pr(l|\mathbf{x}_k, \mathbf{x}_u) \Pr(\mathbf{x}_u^j|\mathbf{s}, \mathbf{a}, \boldsymbol{\theta})}{\Pr(l|\mathbf{x}_k, \mathbf{s}, \mathbf{a}, \boldsymbol{\theta})} \frac{1}{\Pr(\mathbf{x}_u^j|\mathbf{s}, \mathbf{a}, \boldsymbol{\theta})} \\ &= \frac{\Pr(l|\mathbf{x}_k, \mathbf{x}_u)}{\Pr(l|\mathbf{x}_k, \mathbf{s}, \mathbf{a}, \boldsymbol{\theta})} \end{aligned}$$

Consequently, we have removed the product series from the expectation weight, thereby avoiding underflow issues. Also, as noted previously, the  $\Pr(l|\mathbf{x}_k, \mathbf{s}, \mathbf{a}, \boldsymbol{\theta})$  term can be computed efficiently with dynamic programming.

Now we can write a tractable EM algorithm where we can compute the maximization using gradient ascent.

---

**Algorithm 2** Labeled-IRL Approximate EM Gradient Ascent Algorithm

---

**Require:** initial  $\boldsymbol{\theta}_0$ ; data  $\mathbf{s}, \mathbf{a}, \mathbf{x}_k, l$ ; and learning rate  $\alpha$

```

for  $t = 0$  to  $K$  do
  draw  $j = 1$  to  $C$  samples of  $\mathbf{x}_u^j \sim \Pr(\mathbf{x}_u|\mathbf{s}, \mathbf{a}, \boldsymbol{\theta}_t)$ 
  for  $j = 1$  to  $C$  do
     $w_j \leftarrow \frac{\Pr(l|\mathbf{x}_k, \mathbf{x}_u^j)}{\Pr(l|\mathbf{x}_k, \mathbf{s}, \mathbf{a}, \boldsymbol{\theta}_t)}$  ▷ Expectation step
  end for
   $\boldsymbol{\theta}' \leftarrow \boldsymbol{\theta}_t$ 
  for  $1$  to  $M$  do ▷ Gradient ascent maximization loop
     $\boldsymbol{\theta}' \leftarrow \boldsymbol{\theta}' + \alpha \frac{1}{C} \sum_j^C w_j \sum_{x_i \in \mathbf{x}_k \cup \mathbf{x}_u^j} \frac{\nabla_{\boldsymbol{\theta}'} \Pr(x_i|s_i, a_i, \boldsymbol{\theta}')}{\Pr(x_i|s_i, a_i, \boldsymbol{\theta}')}$ 
  end for
   $\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}'$ 
end for
```

---

## 4 DP Algorithm For Label Probability

One of the terms that must be computed for the EM algorithm weight is  $\Pr(l \mid \mathbf{x}_k, \mathbf{s}, \mathbf{a}, \boldsymbol{\theta})$ . A straightforward computation of this value would involve marginalizing over  $\mathbf{x}_u$ , which grows exponentially large as the number of unknown feedbacks increases. However, by exploiting the fact that  $\Pr(l \mid \mathbf{x}_k, \mathbf{x}_u, \mathbf{s}, \mathbf{a}, \boldsymbol{\theta})$  is a function of the sum of the feedback values, the marginalization can be reduced to a summation that iterates over a number of values that is a linear function of the number of unobserved feedbacks. To demonstrate, first note that  $\Pr(l \mid \mathbf{x}, \mathbf{s}, \mathbf{a}, \boldsymbol{\theta})$ , where all feedback values are known, is defined as

$$\Pr(l \mid \mathbf{x}, \mathbf{s}, \mathbf{a}, \boldsymbol{\theta}) = S \left( \sum_{x \in \mathbf{x}} x \right), \quad (10)$$

where  $S$  is the sigmoid function. Therefore,  $\Pr(l \mid \mathbf{x}_k, \mathbf{s}, \mathbf{a}, \boldsymbol{\theta})$ , where some of the feedbacks are unobserved, can be expressed by marginalizing over the possible *sums* of the unknown feedbacks, rather than marginalizing over all possible feedback assignments:

$$\Pr(l \mid \mathbf{x}_k, \mathbf{s}, \mathbf{a}, \boldsymbol{\theta}) = \sum_{\tau = -|\mathbf{x}_u|}^{|\mathbf{x}_u|} S \left( \tau + \sum_{x \in \mathbf{x}_k} x \right) \Pr(\tau \mid \mathbf{s}_u, \mathbf{a}_u, \boldsymbol{\theta}), \quad (11)$$

where  $\Pr(\tau \mid \mathbf{s}_u, \mathbf{a}_u, \boldsymbol{\theta})$  is the probability that the unknown feedbacks will sum to  $\tau$ , given the states and actions for the unknown feedbacks, and the model parameters  $\boldsymbol{\theta}$ .

Computing this simpler marginal distribution requires computing  $\Pr(\tau \mid \mathbf{s}_u, \mathbf{a}_u, \boldsymbol{\theta})$ : the probability distribution of the sum of possible feedback values. Fortunately, this probability can be expressed recursively and computed efficiently with a dynamic programming algorithm.

To describe this recursive relationship, note that if we knew the probability distribution of the sum of feedback values for all unknown feedbacks except the last of them, then we could compute the probability that *all* of the feedbacks would sum to some value  $\tau$  as the probability that the remaining feedbacks sum to  $\tau - 1$  and the last feedback is a  $+1$  or the remaining feedbacks sum to  $\tau + 1$  and the last feedback is a  $-1$ . That is,

$$\begin{aligned} \Pr(\tau_n = \tau \mid \mathbf{s}_n, \mathbf{a}_n, \boldsymbol{\theta}) &= \Pr(x_n = +1 \mid \mathbf{s}_n, \mathbf{a}_n) \Pr(\tau_{n-1} = \tau_n - 1 \mid \mathbf{s}_{n-1}, \mathbf{a}_{n-1}, \boldsymbol{\theta}) \\ &\quad + \Pr(x_n = -1 \mid \mathbf{s}_n, \mathbf{a}_n) \Pr(\tau_{n-1} = \tau_n + 1 \mid \mathbf{s}_{n-1}, \mathbf{a}_{n-1}, \boldsymbol{\theta}), \end{aligned} \quad (12)$$

where  $n$  is the number unobserved feedbacks,  $\tau_i$  is the random variable specifying the sum of the first  $i$  unobserved feedbacks;  $s_i$  and  $a_i$  represents the  $i$ th state and action that are associated with unobserved feedback  $x_i$ ;  $\mathbf{s}_i, \mathbf{a}_i$  represents the set of the first  $i$  state and actions associated with the unobserved feedbacks; and finally, where the probability of the first unobserved feedback summing to

$\tau$  is defined as the probability that the first feedback takes that value:

$$\Pr(\tau_1 = \tau \mid \mathbf{s}_n, \mathbf{a}_n, \boldsymbol{\theta}) = \begin{cases} \Pr(x_n = +1 \mid \mathbf{s}_n, \mathbf{a}_n) & \text{if } \tau_1 = 1 \\ \Pr(x_n = -1 \mid \mathbf{s}_n, \mathbf{a}_n) & \text{if } \tau_1 = -1 \\ 0 & \text{otherwise} \end{cases} \quad (13)$$

We can compute this recursive probability distribution using dynamic programming in which we build a  $2n + 1 \times n$  matrix, where cell  $i, j$  specifies the value for  $\Pr(\tau_j = i \mid \mathbf{s}_j, \mathbf{a}_j, \boldsymbol{\theta})$ . The values for the matrix are then filled out column by column, from  $j = 1$  to  $n$ . Consequently, using this DP algorithm, we can compute the  $\Pr(l \mid \mathbf{x}_k, \mathbf{s}, \mathbf{a}, \boldsymbol{\theta})$  term required for our EM weight in only  $O(|\mathbf{x}_u|^2)$  time, rather than  $O(2^{|\mathbf{x}_u|})$ .

## 5 Other Importance Sampling Distributions

One of the potential problems with the importance sampling distributions being drawn from the policy under the theta values that is as the agent optimizes for theta, it will produce more deterministic action selection, leading to progressively more biased samples that may not adequately cover the space. Indeed, initial experiments with the current distribution show that in some cases, the log likelihood starts getting worse, whereas when using the full marginalization, it consistently improves.

To combat this issue, we should explore other sampling distributions that better span the space. One potential issue is that with a different sampling distribution, we must include the product series of the individual feedback probabilities, which could produce underflow. However, if we simply choose a slightly noisier sampling feedback distribution than the feedback distribution under the current parameters, then we might avoid underflow problems.

That is, let a vector for unknown feedbacks  $\mathbf{x}_u$  be drawn from some sampling distribution  $f$  that independently draws feedbacks as a function of each state-action pair and the current parameters  $\boldsymbol{\theta}$ :

$$f(\mathbf{x}_u \mid \mathbf{s}, \mathbf{a}, \boldsymbol{\theta}) = \prod_i f_i(x_{u,i} \mid s_i, a_i, \boldsymbol{\theta}). \quad (14)$$

Then using  $f$  as the importance sampling distribution, we have the following importance sampling weights for each sample  $\mathbf{x}_u^j$ :

$$\begin{aligned} w_j &= \frac{\Pr(l \mid \mathbf{x}_k, \mathbf{x}_u^j) \Pr(\mathbf{x}_u^j \mid \mathbf{s}, \mathbf{a}, \boldsymbol{\theta})}{\Pr(l \mid \mathbf{x}_k, \mathbf{s}, \mathbf{a}, \boldsymbol{\theta})} \frac{1}{f(\mathbf{x}_u^j \mid \mathbf{s}, \mathbf{a}, \boldsymbol{\theta})} \\ &= \frac{\Pr(l \mid \mathbf{x}_k, \mathbf{x}_u^j) \prod_i \Pr(x_{u,i}^j \mid s_i, a_i, \boldsymbol{\theta})}{\Pr(l \mid \mathbf{x}_k, \mathbf{s}, \mathbf{a}, \boldsymbol{\theta})} \frac{1}{\prod_i f_i(x_{u,i}^j \mid s_i, a_i, \boldsymbol{\theta})} \\ &= \frac{\Pr(l \mid \mathbf{x}_k, \mathbf{x}_u^j)}{\Pr(l \mid \mathbf{x}_k, \mathbf{s}, \mathbf{a}, \boldsymbol{\theta})} \prod_i \frac{\Pr(x_{u,i}^j \mid s_i, a_i, \boldsymbol{\theta})}{f_i(x_{u,i}^j \mid s_i, a_i, \boldsymbol{\theta})} \end{aligned} \quad (15)$$

Therefore, as long as  $\frac{\Pr(x_{u,i}^j | s_i, a_i, \theta)}{f_i(x_{u,i}^j | s_i, a_i, \theta)}$  is reasonably near 1, there shouldn't be underflow issues. This condition may be held if the single feedback distribution  $f_i$ , is also based on a softmax policy distribution of the current weights theta, but is more uniform (e.g., a small  $\beta$  value.) That is,

$$f_i(x_i | s_i, a_i, \theta) = \begin{cases} \frac{e^{\hat{\beta}Q(s,a|\theta)}}{\sum_{a'} e^{\hat{\beta}Q(s,a'|\theta)}} & \text{if } x_i = 1 \\ 1 - \frac{e^{\hat{\beta}Q(s,a|\theta)}}{\sum_{a'} e^{\hat{\beta}Q(s,a'|\theta)}} & \text{if } x_i = -1 \end{cases}, \quad (16)$$

where  $\hat{\beta} < \beta$  and  $\beta$  is the normal softmax policy parameter used by our probability model (Equation 3), such that  $f_i$  is a more noisy, uniform response.