

# Learning Markov Logic Networks

**Log-linear models** Suppose we are given factors  $\phi_1, \phi_2, \dots$ , each of which maps a state  $\mathbf{x}_i$  of the random variables  $\mathbf{X}$  to a non-negative number. Note that each factor may take a different subset of the variables as input. This is a factor graph (or equivalently a Markov network, depending on how you draw the graph) whose semantics define a joint probability distribution as follows.

$$\Pr(\mathbf{x}) = \frac{\prod_i \phi_i(\mathbf{x}_i)}{Z} \text{ where } Z = \sum_{\mathbf{x}} \prod_i \phi_i(\mathbf{x}_i). \quad (1)$$

The normalizing constant  $Z$  is called the partition function. We will find it convenient to represent this factor graph as a *log-linear model*. In such a model, the factors  $\phi_i$  are represented as exponentiated weighted features, with real-valued weights denoted  $w_i$  and features denoted  $f_i(\mathbf{x})$ :

$$\phi_i(\mathbf{x}_i) = \exp(w_i f_i(\mathbf{x}_i)) \quad (2)$$

The joint distribution can then be written as

$$\Pr(\mathbf{x}) = \frac{\prod_i \exp(w_i f_i(\mathbf{x}_i))}{Z} \quad (3)$$

$$\Pr(\mathbf{x}) = \frac{\exp(\sum_i w_i f_i(\mathbf{x}_i))}{Z} \quad (4)$$

where now

$$Z = \sum_{\mathbf{x}} \exp\left(\sum_i w_i f_i(\mathbf{x}_i)\right). \quad (5)$$

One benefit of this transformation is that we can write  $\sum_i w_i f_i(\mathbf{x}_i)$  instead as the dot product  $\mathbf{w} \cdot \mathbf{f}$  of a weight vector  $\mathbf{w}$  and a feature vector  $\mathbf{f}$  (a linear model), which is a common notation in machine learning.

**Markov logic networks** Syntactically, a Markov logic network (MLN) is a set of weighted first-order logic formulas  $\{w_i \alpha_i\}$  where  $\alpha_i$  is an arbitrary formula. One way to define MLN semantics is to think of them as templates for a log-linear model. Specifically, each grounding of the free variables in  $\alpha_i$  defines a Boolean feature. For example

$$\alpha = [\text{Smokes}(x) \wedge \text{Friends}(x, y) \Rightarrow \text{Smokes}(y)] \quad (6)$$

has many groundings for  $x$  and  $y$ . For any particular choice, say

$$\theta = \{x \setminus \text{Alice}, y \setminus \text{Bob}\} \quad (7)$$

the sentence

$$\alpha\theta = [\textit{Smokes}(\textit{Alice}) \wedge \textit{Friends}(\textit{Alice}, \textit{Bob}) \Rightarrow \textit{Smokes}(\textit{Bob})] \quad (8)$$

is either true or false in a world. Therefore, we can use it as a feature:

$$f(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{x} \models \alpha\theta \\ 0 & \text{otherwise} \end{cases} \quad (9)$$

If  $\alpha$  contains quantifiers, the semantics change:

$$\alpha = [\exists y. \textit{Smokes}(x) \wedge \textit{Friends}(x, y) \Rightarrow \textit{Smokes}(y)] \quad (10)$$

$$\alpha\{x \setminus \textit{Alice}\} = [\exists y. \textit{Smokes}(\textit{Alice}) \wedge \textit{Friends}(\textit{Alice}, y) \Rightarrow \textit{Smokes}(y)] \quad (11)$$

Still, the sentences obtained after grounding the free variables are either true or false in a world, and are therefore Boolean features.

Let  $g_{ij}$  denote the feature induced by the  $j$ th grounding of formula  $\alpha_i$ . Then, if the MLN semantics are simply defining a log-linear model with weights  $w_i$  for each of these features obtained by grounding the formula  $\alpha$ , we get the following.

$$\Pr(\mathbf{x}) = \frac{\exp(\sum_k w_k f_k(\mathbf{x}_k))}{Z} \quad (12)$$

$$\Pr(\mathbf{x}) = \frac{\exp(\sum_i \sum_j w_i g_{ij}(\mathbf{x}_{ij}))}{Z} \quad (13)$$

$$\Pr(\mathbf{x}) = \frac{\exp(\sum_i w_i \sum_j g_{ij}(\mathbf{x}_{ij}))}{Z} \quad (14)$$

$$\Pr(\mathbf{x}) = \frac{\exp(\sum_i w_i n_i(\mathbf{x}_i))}{Z} \quad \text{where } Z = \sum_{\mathbf{x}} \exp\left(\sum_i w_i n_i(\mathbf{x}_i)\right) \quad (15)$$

and where

$$n_i(\mathbf{x}) = \sum_j g_{ij}(\mathbf{x}_{ij}) \quad (16)$$

Essentially,  $n_i$  is a new feature, which is counting how many groundings of  $\alpha_i$  are true in world  $\mathbf{x}$ . That new feature has the same weight  $w_i$  as all the grounding features it summarizes.

**Maximum-Likelihood Weight Learning** Suppose we want to learn maximum-likelihood weights given the dataset  $\mathcal{D} = \{\mathbf{x}^1, \mathbf{x}^2, \dots\}$ , that is, find

$$w = \arg \max_w \Pr(\mathcal{D}; w) \quad (17)$$

$$= \arg \max_w \log \Pr(\mathcal{D}; w) \quad (18)$$

$$= \arg \max_w \log \prod_k \Pr(\mathbf{x}^k; w) \quad (19)$$

$$= \arg \max_w \sum_k \log \Pr(\mathbf{x}^k; w) \quad (20)$$

$$= \arg \max_w \sum_k \log \left( \frac{\exp(\sum_i w_i n_i(\mathbf{x}_i^k))}{Z} \right) \quad (21)$$

$$= \arg \max_w \sum_k \log \left( \exp \left( \sum_i w_i n_i(\mathbf{x}_i^k) \right) \right) - \log(Z) \quad (22)$$

$$= \arg \max_w \sum_k \sum_i w_i n_i(\mathbf{x}_i^k) - \log(Z) \quad (23)$$

Intuitively, the maximum likelihood weights trade off two terms: increasing the weight of the data  $\sum_i w_i n_i(\mathbf{x}_i)$  vs. how much this increase also increases the partition function  $Z$ .

This optimization problem has no easy closed-form solution, hence we will do gradient-based optimization. In order to do so, we need to compute all partial derivatives of the form

$$\frac{\partial \log \Pr(\mathbf{x}; w)}{\partial w_i} = \frac{\partial \sum_i w_i n_i(\mathbf{x}_i) - \log(Z)}{\partial w_i} \quad (24)$$

$$= \frac{\partial \sum_i w_i n_i(\mathbf{x}_i)}{\partial w_i} - \frac{\partial \log(Z)}{\partial w_i} \quad (25)$$

$$= n_i(\mathbf{x}_i) - \frac{\partial \log(Z)}{\partial w_i} \quad (26)$$

That is, the partial derivative of the likelihood given a single example w.r.t. weight  $w_i$  is the difference between  $n_i(\mathbf{x}_i)$ , the number of true groundings of the formula in the database, and the partial derivative of the log-partition function. We can further simplify the latter:

$$\frac{\partial \log(Z)}{\partial w_i} = \frac{\partial \log(\sum_{\mathbf{x}} \exp(\sum_i w_i n_i(\mathbf{x}_i)))}{\partial w_i} \quad (27)$$

$$= \frac{\left( \frac{\partial \sum_{\mathbf{x}} \exp(\sum_i w_i n_i(\mathbf{x}_i))}{\partial w_i} \right)}{\sum_{\mathbf{x}} \exp(\sum_i w_i n_i(\mathbf{x}_i))} \quad (28)$$

$$= \frac{\left( \frac{\partial \sum_{\mathbf{x}} \exp(\sum_i w_i n_i(\mathbf{x}_i))}{\partial w_i} \right)}{Z} \quad (29)$$

$$= \frac{\sum_{\mathbf{x}} \frac{\partial \exp(\sum_i w_i n_i(\mathbf{x}_i))}{\partial w_i}}{Z} \quad (30)$$

$$= \frac{\sum_{\mathbf{x}} n_i(\mathbf{x}_i) \exp(\sum_i w_i n_i(\mathbf{x}_i))}{Z} \quad (31)$$

$$= \sum_{\mathbf{x}} n_i(\mathbf{x}_i) \frac{\exp(\sum_i w_i n_i(\mathbf{x}_i))}{Z} \quad (32)$$

$$= \sum_{\mathbf{x}} n_i(\mathbf{x}_i) \Pr(\mathbf{x}) \quad (33)$$

$$= \mathbb{E}[n_i(\mathbf{x}_i)] \quad (34)$$

The partial derivative of the log-partition function is simply the expectation of the corresponding feature in the distribution parameterized by  $w$ . Hence, the partial derivative of the likelihood for a single data point  $\mathbf{x}$  is

$$\frac{\partial \log \Pr(\mathbf{x}; w)}{\partial w_i} = n_i(\mathbf{x}_i) - \mathbb{E}_w[n_i(\mathbf{x}_i)]. \quad (35)$$

It is the difference between the count in the data  $n_i(\mathbf{x}_i)$  and the expected count  $\mathbb{E}_w[n_i(\mathbf{x}_i)]$  in the distribution parametrized by the current weights  $w$ . It is clear that at the maximum-likelihood weights  $w^*$ , these partial derivatives are all zero, and therefore

$$n_i(\mathbf{x}_i) = \mathbb{E}_{w^*}[n_i(\mathbf{x}_i)]. \quad (36)$$

Therefore, algorithms that maximize this equation are also called *moment matching* algorithms: they ensure that the “moment”  $n_i(\mathbf{x}_i)$  in the data distribution and the learned distribution are equal.

To do gradient optimization in practice, we need the ability to estimate the expectation  $\mathbb{E}[n_i(\mathbf{x}_i)]$ . Recall that

$$n_i(\mathbf{x}) = \sum_j g_{ij}(\mathbf{x}_{ij}) \quad (37)$$

and that therefore

$$\mathbb{E}[n_i(\mathbf{x}_i)] = \sum_{\mathbf{x}} n_i(\mathbf{x}_i) \Pr(\mathbf{x}) \quad (38)$$

$$= \sum_{\mathbf{x}} \sum_j g_{ij}(\mathbf{x}_{ij}) \Pr(\mathbf{x}) \quad (39)$$

$$= \sum_j \sum_{\mathbf{x}} g_{ij}(\mathbf{x}_{ij}) \Pr(\mathbf{x}) \quad (40)$$

$$= \sum_j \mathbb{E}[g_{ij}(\mathbf{x}_{ij})] \quad (41)$$

Finally, because  $g_{ij}$  is a logical sentence (and thus it is a Boolean random variable), we have that

$$\mathbb{E}[n_i(\mathbf{x}_i)] = \sum_j \Pr(g_{ij}(\mathbf{x}_{ij})). \quad (42)$$

Thus, the expectation of the count feature can be computed as the sum of the probabilities of all the groundings of the corresponding formula.

**Bayesian Weight Learning** One issue with the approach above is with the scenario where the count  $n_i(\mathbf{x}_i) = 0$  for some sentence  $\alpha_i$ . There is no way to make  $\mathbb{E}_{w^*}[n_i(\mathbf{x}_i)] = 0$  in a log linear-model, and gradient-based learning will simply keep reducing  $w_i$ , pushing it towards negative infinity, and never converge. Moreover, we do not even want to make  $\mathbb{E}_{w^*}[n_i(\mathbf{x}_i)] = 0$ : it means we are clearly overfitting the training data.

Bayesian learning assumes a prior distribution on the parameters:  $\Pr(w)$ . The likelihood model is still of the form above, now denoted  $\Pr(\mathbf{x}|w)$ . The joint distribution over data and parameters

is therefore  $\Pr(\mathbf{x}, w)$ . Learning in this context is simply obtaining the posterior distribution on parameters  $\Pr(w|\mathbf{x})$ . The parameters we seek to learn are now either the expected value of  $w$  in this posterior (Bayesian estimates), or the most likely value of  $w$  (maximum-a-posteriori estimates):

$$w = \arg \max_w \Pr(w|\mathbf{x}) \quad (43)$$

$$= \arg \max_w \log \Pr(w|\mathbf{x}) \quad (44)$$

$$= \arg \max_w \log \frac{\Pr(\mathbf{x}|w) \Pr(w)}{\Pr(\mathbf{x})} \quad (45)$$

$$= \arg \max_w \log \Pr(\mathbf{x}|w) + \log \Pr(w) - \log \Pr(\mathbf{x}) \quad (46)$$

Since  $\Pr(\mathbf{x})$  does not depend on  $w$ , this is also.

$$w = \arg \max_w \log \Pr(\mathbf{x}|w) + \log \Pr(w) \quad (47)$$

The first term in this objective function is again the log-likelihood, as in maximum-likelihood learning, which we already know how to optimize. The second term aims to learn weights  $w$  that are more likely a priori. Depending on the assumptions on the distribution of this prior (typically, that  $\Pr(w)$  is a Gaussian), this second term can be written as a norm of the weight vector (typically L2 norm). Hence, Bayesian learning amounts to simply adding a regularization term to the weight learning loss function objective.