

# Notes

Andrew Kontaxis, Chris White

November 26, 2017

## **1 Background**

## 1.1 Sum Product Algorithm

Suppose we have a graphical model which is a connected undirected tree; in that case we can choose an arbitrary ordering for the nodes and write:

$$p(x_1, x_2, \dots, x_n) = \prod_i \psi_i(x_i) \prod_{i,j \in E} \psi_{i,j}(x_i, x_j)$$

Suspend disbelief about probabilistic interpretations for a moment and suppose we simply want to compute the quantity

$$p(x_s) := \sum_{i \neq s} p(x_1, x_2, \dots, x_s, \dots, x_n)$$

where we interpret the sum as being over the *state space* of the corresponding variables. We can write

$$\begin{aligned} p(x_1, x_2, \dots, x_s, \dots, x_n) &= \psi_s(x_s) \prod_{i \neq s} \psi_i(x_i) \prod_{i,j \in E} \psi_{i,j}(x_i, x_j) \\ &= \psi_s(x_s) \prod_{i \in \mathcal{N}(s)} \psi_i(x_i) \psi_{i,s}(x_i, x_s) \omega(T_i) \end{aligned}$$

where  $\mathcal{N}(s)$  denotes the set of *neighbors* of node  $s$  and  $\omega(T_i)$  is a *weighting* of the subtree containing node  $i$  formed by removing node  $s$ ; this weighting is a function of all variables in  $T_i$ . Let us focus our attention on a single  $i \in \mathcal{N}(s)$  for a moment, and imagine marginalizing out only the nodes in  $T_i$  first:

$$\begin{aligned} \sum_{T_i} p(x_1, x_2, \dots, x_s, \dots, x_n) &= \sum_{T_i} \psi_s(x_s) \prod_{k \in \mathcal{N}(s)} \psi_k(x_k) \psi_{k,s}(x_k, x_s) \omega(T_k) \\ &= \kappa(x_s, x_{V \setminus T_i}) \sum_{T_i} \psi_i(x_i) \psi_{i,s}(x_i, x_s) \omega(T_i) \\ &= \kappa(x_s, x_{V \setminus T_i}) \sum_{x_i} \psi_i(x_i) \psi_{i,s}(x_i, x_s) \sum_{T_i \setminus x_i} \omega(x_i, T_i) \end{aligned}$$

We should now see that in order to complete this computation, the weighting need only be given to us as a function of  $x_i$  alone. I.e., we can write

$$\kappa(x_s, x_{V \setminus T_i}) \sum_{x_i} \psi_i(x_i) \psi_{i,s}(x_i, x_s) \omega(x_i)$$

where  $\omega(x_i)$  is given by

$$\omega(x_i) := \sum_{T_i \setminus i} \prod_{k \in \mathcal{N}(i) \setminus s} \psi_k(x_k) \psi_{k,i}(x_k, x_i) \omega(T_k)$$

and now we begin to see the recursive nature of our task. We can then proceed to marginalize out the variables in each of the other subtrees resulting in an expression of the form

$$\psi_s(x_s) \prod_{i \in \mathcal{N}(s)} \left( \sum_{x_i} \psi_i(x_i) \psi_{i,s}(x_i, x_s) \omega(x_i) \right)$$

Consequently, to compute the marginal  $p(x_s)$ , each neighbor  $i$  of  $s$  needs to pass a “message” to node  $s$  which is a function purely of  $x_s$ , specifying the “weighting” of the subtree  $T_i$  *conditional on* the value  $x_s$ :

$$\mu_{i \rightarrow s}(x_s) := \sum_{x_i} \psi_i(x_i) \psi_{i,s}(x_i, x_s) \prod_{k \in \mathcal{N}(i) \setminus s} \mu_{k \rightarrow i}(x_i)$$

Note that the messages  $\mu$  are proxies for the “true” weightings  $\omega$  and that whenever  $\mu \equiv \omega$  we have a fixed point and can compute

$$p(x_s) = \psi_s(x_s) \prod_{i \in \mathcal{N}(s)} \mu_{i \rightarrow s}(x_s).$$

This distinction is particularly important in the case when we want to apply this algorithm to graphs with cycles; otherwise we can iteratively update the messages by beginning at the leaves with the messages

$$\mu_{\ell \rightarrow q}(x_q) = \sum_{x_\ell} \psi_\ell(x_\ell) \psi_{\ell,q}(x_\ell, x_q)$$

### 1.1.1 Sum Product Algorithm on Factor Graphs

For the general case in which potential functions are defined on more than just nodes and edges, it is convenient to think in terms of *factor graphs*. In this case, messages take the form

$$\begin{aligned} \mu_{\psi_\alpha \rightarrow s}(x_s) &:= \sum_{x_M} \psi_\alpha(x_s, x_M) \prod_{m \in M} \mu_{m \rightarrow \psi_\alpha}(x_m) \\ \mu_{s \rightarrow \psi_\alpha}(x_s) &:= \prod_{\psi \in \mathcal{N}(x_s) \setminus \psi_\alpha} \mu_{\psi \rightarrow s}(x_s). \end{aligned}$$

Note that we recover the above messages by observing that

$$\begin{aligned} \mu_{\psi_{i,s} \rightarrow x_s}(x_s) &= \sum_{x_i} \psi_{i,s}(x_i, x_s) \mu_{i \rightarrow \psi_{i,s}}(x_i) \\ &= \sum_{x_i} \psi_{i,s}(x_i, x_s) \psi_i(x_i) \prod_{k \in \mathcal{N}(i) \setminus s} \mu_{\psi_{i,k} \rightarrow i}(x_i) \end{aligned}$$

and then changing subscript notation.

### 1.1.2 An Example: Independent Sets on Graphs

An *independent set* of vertices on a graph is a subset of vertices such that no two are adjacent. This example will be particularly fruitful in future discussions; consider a binary tree whose nodes are assigned weights. We imagine this as an undirected graphical model by putting a binary random variable at each node, with potential functions given by:

$$\begin{aligned} \psi_i(x_i) &:= \exp(w_i x_i) \\ \psi_{ij}(x_i, x_j) &:= \chi(x_i + x_j \leq 1) \end{aligned}$$

Here  $\chi(\cdot)$  is the indicator function for whether the condition holds. Let us see what the sum-product algorithm looks like for this model; first imagine a parent of two leaf nodes. Applying the above formulas, each leaf (here denoted  $\ell$  and  $r$ ) sends the message

$$\begin{aligned}\mu_{\ell \rightarrow p}(x_p) &= \exp(w_\ell)\chi(x_p = 0) + 1 \\ \mu_{r \rightarrow p}(x_p) &= \exp(w_r)\chi(x_p = 0) + 1\end{aligned}$$

and consequently the message the  $x_p$  sends to its parent  $x_P$  is given by

$$\begin{aligned}\mu_{p \rightarrow P}(x_P) &= \sum_{x_p} (\exp(w_\ell)\chi(x_p = 0) + 1) (\exp(w_r)\chi(x_p = 0) + 1) \exp(w_p x_p) \chi(x_p + x_P \leq 1) \\ &= (\exp(w_\ell) + 1) (\exp(w_r) + 1) + \exp(w_p) \chi(x_P = 0)\end{aligned}$$

If we pause here and suppose  $x_p$  is the only node connected to  $x_P$ , we see that the marginal for  $x_P$  would be proportional to

$$\exp(w_P x_P) \cdot [(\exp(w_\ell) + 1) (\exp(w_r) + 1) + \exp(w_p) \chi(x_P = 0)]$$

Observe that each term of this summation corresponds to an allowable configuration for an independent set. For example, when  $x_P = 1$  the leading term  $\exp(w_P + w_\ell + w_r)$  corresponds to the independent set  $\{P, \ell, r\}$ .  $\exp(w_\ell)$  occurs in precisely 2 terms because regardless of whether  $P$  is chosen,  $\ell$  can be chosen in two possible independent sets. Lastly we note that the constant term corresponds to the null set.

## 1.2 Max Product Algorithm

note that sum-product derivation carries over to maxes because the only property we relied upon was operator of product is product of operator; moreover, if we specialize to binary integer programming problems suffices to track message differences in the binary case

### 1.2.1 Extended Example: Affinity Propagation

Consider the non-convex optimization problem

$$\begin{aligned} \max_{x_i} \quad & \sum_{i \neq j} s_{ij} x_{ij} + \lambda \sum_j x_{jj} \\ \text{subject to} \quad & \sum_j x_{ij} = 1 \\ & x_{ij} \leq x_{jj} \\ & x_{ij} \in \{0, 1\} \end{aligned}$$

If we imagine  $s_{ij}$  as a similarity measure between two data points  $i$  and  $j$ , then this problem can be thought of as a clustering problem in which every cluster contains a representative data point (those points for which  $x_{jj} = 1$ ).

## 2 Exponential Families

## 2.1 Exponential Families as Max Entropy Distributions

Suppose we are given a collection of *sufficient statistics* or *potentials*  $\phi_\alpha : \mathcal{X} \rightarrow \mathbb{R}$  and a collection of *mean parameters* for each statistic,  $\mu_\alpha$ . There are possibly multiple probability distributions  $p$  over  $\mathcal{X}$  satisfying

$$\mathbb{E}_p[\phi_\alpha(X)] = \mu_\alpha \quad \forall \alpha \in \mathcal{I}$$

or there are none. For example, if  $\phi_1(X) = X$  and  $\phi_2(X) = X^2$  then by Jensen's inequality we must have

$$\mu_2 \geq \mu_1^2.$$

Thus not every combination of  $\mu$  and  $\phi$  will admit an admissible  $p$ . Moreover, if we assume  $p$  is *absolutely continuous* with respect to some base measure over  $\mathcal{X}$ , then this further restricts our admissible  $p$ . Continuing with the above example, if  $\mu_1^2 = \mu_2$  and  $\mathcal{X} = \mathbb{R}$  with Lebesgue measure then we have

$$\int x^2 dp(x) = \left( \int x dp(x) \right)^2$$

which implies the random variable  $X$  is constant almost surely<sup>1</sup>, i.e.  $p$  is a dirac delta function which is *not* absolutely continuous with respect to Lebesgue measure.

Supposing there exists at least one such distribution with mean parameters  $\{\mu_\alpha\}_{\alpha \in \mathcal{I}}$ , it is natural to ask for the distribution with the *maximal* amount of uncertainty satisfying the above mean parameter constraint (we want to encode the least amount of additional information above and beyond the constraints). Using *Shannon entropy* as our measure of “uncertainty” we are led to the optimization problem:

$$p^* := \arg \max_p -\mathbb{E}_p[\log(p)] \quad \text{s.t.} \quad \mathbb{E}[\phi_\alpha(X)] = \mu_\alpha \quad \forall \alpha \in \mathcal{I}$$

### 2.1.1 Discrete Case

First suppose that  $\mathcal{X}$  is finite and discrete. In this case  $p$  is simply a non-negative vector which sums to 1 and we form the Lagrangian

$$\mathcal{L}(p, \lambda, \tau) := -\sum_{i=1}^n \log(p_i) p_i + \sum_{\alpha \in \mathcal{I}} \lambda_\alpha (\mathbb{E}[\phi_\alpha(X)] - \mu_\alpha) + \tau (\sum_{i=1}^n p_i - 1) + \sum_{i=1}^n \gamma_i p_i.$$

We have the necessary stationarity conditions

$$\frac{\partial \mathcal{L}}{\partial p_i} = -\log(p_i) - 1 + \sum_{\alpha \in \mathcal{I}} \lambda_\alpha \phi_\alpha(X_i) + \tau + \gamma_i = 0$$

which implies that  $p^*$  is of the form

$$p^* = \exp \left( \sum_{\alpha \in \mathcal{I}} \lambda_\alpha \phi_\alpha + \tau + \sum_{i=1}^n \gamma_i \chi_i - 1 \right) > 0,$$

i.e., it is in the exponential family with *canonical* or *exponential* parameters  $\lambda_\alpha$  and  $\tau$  is chosen to enforce the appropriate normalization. Additionally, note that because  $p^* > 0$  we must have that  $\gamma_i = 0$  for all  $i$  by complementary slackness.

---

<sup>1</sup>For strictly convex functions  $\phi$ , we have that  $\phi(x) \geq \phi(y) + \phi'(y)(x - y)$  with equality if and only if  $y = x$ . Using  $y = \mathbb{E}[X]$  and integrating both sides yields the claim.

### 2.1.2 Continuous Case

The continuous case is more subtle; for example, with  $\phi_1$  and  $\phi_2$  as above, if  $\mu_2 = \mu_1^2$ , there does not exist a feasible density. Instead, we will prove a weaker claim: if there exists a feasible density of the form  $p^*(x) = \exp(\theta^T \phi - A(\theta))$  then it is necessarily the unique maximizer. Let  $p(x)$  be any other feasible density. Then we have

$$\begin{aligned}
 - \int p \log p \, dx &= - \int p \log \frac{p}{p^*} \, dx - \int p \log p^* \, dx \\
 &= - \int p \log \frac{p}{p^*} \, dx - \int p (\theta^T \phi - A(\theta)) \, dx \\
 &= - \int p \log \frac{p}{p^*} \, dx - \theta^T \bar{\alpha} + A(\theta) \\
 &= - \int p \log \frac{p}{p^*} \, dx - \int p^* \log p^* \, dx
 \end{aligned}$$

and as  $-\int p \log \frac{p}{p^*} \, dx < 0$  for any  $p \neq p^*$  the claim is proven.