

ECE 175B: Probabilistic Reasoning and Graphical Models

Lecture 15: Inference on Discrete Variable Serial Chains III

Florian Meyer & Ken Kreutz-Delgado

*Electrical and Computer Engineering Department
University of California San Diego*

Querying a Distribution Encoded in a Graph

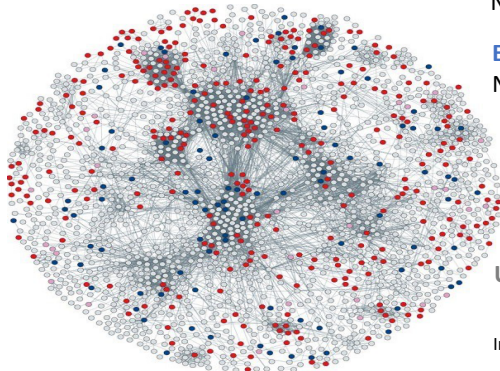
$$P_{\mathbf{X}}(x) = P_{\mathbf{W}, \mathbf{E}}(w, e) = P_{\mathbf{V}, \mathbf{Y}}(v, y) = P_{\mathbf{Y}, \mathbf{E}, \mathbf{Z}}(y, e, z)$$

Query nodes = \mathbf{Y} = "Outcome"

Non-Query nodes = $\mathbf{V} = \mathbf{X} - \mathbf{Y}$

Evidence nodes = \mathbf{E}

Non-Evidence nodes = $\mathbf{W} = \mathbf{X} - \mathbf{E}$



$$\mathbf{Y} \cap \mathbf{E} = \emptyset$$

$$\mathbf{Z} = \mathbf{X} - (\mathbf{Y} \cup \mathbf{E})$$

Unobserved nodes = \mathbf{Z}

In the standard manner we set $P_{\mathbf{X}}(x) = P(x)$, etc.

Two Basic Queries

- Probability of Outcome $Y = y$ Given Evidence $E = e$

Compute the value of $P(y|e)$

- Most Probable Value of Outcome $Y = y$ Given Evidence $E = e$

Compute the most probable outcome

$$\hat{y}_{\text{map}} = \arg \max_y P(y|e)$$

How do we do this in an efficient, computationally tractable manner?

Queries and Serial-Chain Factorization (SCF)

For distributions that have a **Serial-Chain Factorization** (SCF)

$$P(\mathbf{x}) = P(\mathbf{w}, \mathbf{e}) = \frac{1}{Z} \prod_{k=0}^{N-1} \psi_{k,k+1}(w_k, w_{k+1}) \iff P(\mathbf{w} | \mathbf{e}) = \frac{1}{Z(\mathbf{e})} \prod_{k=0}^{N-1} \psi_{k,k+1}(w_k, w_{k+1})$$

we can computationally efficiently compute $P(w_k | \mathbf{e})$ via the very similar Message Passing Algorithm (MPA) and Belief Propagation (BP) “forward–backward” algorithms.

The computation of $P(w_k | \mathbf{e})$ directly answers the query, “What is the probability that W_k takes the value w_k given the evidence?”, which then enables the answer to the query, “What is the most probable value of W_k given the evidence?”:

$$\hat{w}_k^{(\text{map})} = \arg \max_{w_k} P(w_k | \mathbf{e}).$$

The MAP estimate $\hat{w}_k^{(\text{map})}$ gives the Bayes optimal, *Minimum Probability of Error* (MPE), estimate of the state (“symbol”) W_k at the *single* node (“time”) k given *all* of the evidence $\mathbf{e} = \{e_1, \dots, e_k, \dots, e_N\}$. This is a very important result in digital cellular communication where it is known as the minimum symbol error rate solution.

An even more frequently used estimate of W_k in cellular communications, because it is less computationally expensive while given performance nearly as good as the MAP estimate, is provided by the **Viterbi Algorithm** (VA) described in a latter slide.

Variable Elimination – The Sum-Product Algorithm

When $P(\mathbf{w}|\mathbf{e})$ is a SCF distribution, the marginal distribution $P(w_k|\mathbf{e})$ has the **Sum-Product** form,

$$\begin{aligned} P(w_k|\mathbf{e}) &\propto \sum_{\mathbf{w} \setminus w_k} \prod_{k=0}^{N-1} \psi(w_k, w_{k+1}) \\ &= \sum_{w_1} \cdots \sum_{w_{k-1}} \sum_{w_{k+1}} \cdots \sum_{w_N} \psi(w_0, w_1) \cdots \psi(w_{k-1}, w_k) \psi(w_k, w_{k+1}) \cdots \psi(w_{N-1}, w_N) \\ &= \underbrace{\left(\sum_{w_1} \cdots \sum_{w_{k-1}} \psi(w_0, w_1) \cdots \psi(w_{k-1}, w_k) \right)}_{\triangleq m^-(w_k) = \text{message to } w_k \text{ from the past}} \underbrace{\left(\sum_{w_{k+1}} \cdots \sum_{w_N} \psi(w_k, w_{k+1}) \cdots \psi(w_{N-1}, w_N) \right)}_{\triangleq m^+(w_k) = \text{message to } w_k \text{ from the future}} \end{aligned}$$

which results in the simple MPA “sum-product” recursions

$$m^-(w_k) = \sum_{w_{k-1}} \psi(w_{k-1}, w_k) m^-(w_{k-1}) \quad \text{and} \quad m^+(w_k) = \sum_{w_{k+1}} \psi(w_k, w_{k+1}) m^+(w_{k+1}).$$

The variables in $\mathbf{w} \setminus w_k$ are one-by-one “eliminated” (efficiently dealt with) via a sum over a product.

There are other such **Variable Elimination structures**, e.g., **Max-Product**, **Min-Product**, **Max-Sum**, and **Min-Sum**, that each lead to corresponding efficient algorithms. As we shall see, they are closely related. The Max-Product algorithm is discussed in Barber Section 5.2.1 and the Max-Sum algorithm in Bishop Section 8.4.5. Variable Elimination is discussed in Koller Chapter 9.

The Max-Product Algorithm – I

Given $P(\mathbf{w}|\mathbf{e})$ we are often interested in efficiently computing

$$\check{\mathbf{w}}^{(\text{map})} = \arg \max_{\mathbf{w}} P(\mathbf{w}|\mathbf{e}) = \arg \max_{\mathbf{w}} P(\mathbf{w}, \mathbf{e})$$

Note that $\check{\mathbf{w}}^{(\text{map})}$ is an estimate of the *entire sequence of states* (“symbols”) viewed as a single entity $\mathbf{w} = \{w_1, \dots, w_N\}$. I.e., we are asking for an *optimal sequence estimate*.

In general it can happen that the optimal sequence estimate does not produce the optimal single-node k (“time k ”) estimate:

$$\check{w}_k^{(\text{map})} \neq \hat{w}_k^{(\text{map})}.$$

Consistent with cellular comm we can call the optimal **sequence estimate** $\check{\mathbf{w}}^{(\text{map})}$ the Viterbi solution.

When $P(\mathbf{w}, \mathbf{e})$ is a SCF distribution, the optimization $\max_{\mathbf{w}} P(\mathbf{w}, \mathbf{e})$ has the **Max-Product** form,

$$\begin{aligned} \max_{\mathbf{w}} P(\mathbf{w}, \mathbf{e}) &= \max_{w_1, \dots, w_N} \prod_{k=0}^{N-1} \psi(w_k, w_{k+1}) \\ &= \max_{w_1} \cdots \max_{w_{k-1}} \max_{w_{k+1}} \cdots \max_{w_N} \psi(w_0, w_1) \cdots \psi(w_{k-1}, w_k) \psi(w_k, w_{k+1}) \cdots \psi(w_{N-1}, w_N) \\ &= \max_{w_1} \underbrace{\psi(w_0, w_1)}_{\psi(w_1)} \left[\cdots \left[\max_{w_k} \psi(w_{k-1}, w_k) \underbrace{\left[\max_{w_{k+1}} \psi(w_k, w_{k+1}) \left[\cdots \left[\max_{w_N} \psi(w_{N-1}, w_N) \left[1 \right] \right] \right] \right]}_{c_k(w_k)} \right] \cdots \right] \end{aligned}$$

$\underbrace{\hspace{15em}}_{c_{k-1}(w_{k-1})}$

The Max-Product Algorithm – II

The function $c(w_k)$ is the (multiplicative) optimal “cost-to-go” and is computed “backwards” recursively via

$$c(w_{k-1}) = \max_{w_k} \psi_{k-1,k}(w_{k-1}, w_k; \mathbf{e}) c(w_k)$$

for $k = N - 1, \dots, 1$ with boundary condition $c(w_N) = 1$. This shows that a sequence of Max-Product operations “eliminates” w_k from consideration in a step-by-step manner. Note that under our assumption that w_k is categorical, **at each step k we can (and do) readily compute the function $c_k(w_k)$** as a look up table of values for each value that w_k can take. Thus at the end of the recursion we will have computed a sequence of cost-to-go functions $c_1(w_1), \dots, c_k(w_k), \dots, c_{N-1}(w_{N-1})$.

Having determined the optimal cost-to-go functions $c_k(w_k)$ we use them to sequentially produce $\check{w}_k = \check{w}_k^{(\text{map})}$ via a process called *backtracking* or *backtracing*:

$$\text{Viterbi solution } \check{\mathbf{w}} \text{ has components} = \left\{ \begin{array}{l} \check{w}_1 = \arg \max_{w_1} \psi_{0,1}(w_0, w_1) c_1(w_1) \\ \check{w}_2 = \arg \max_{w_2} \psi(\check{w}_1, w_2) c_2(w_2) \\ \vdots \\ \check{w}_k = \arg \max_{w_k} \psi_{k-1,k}(\check{w}_{k-1}, w_k) c_k(w_k) \\ \vdots \\ \check{w}_N = \arg \max_{w_N} \psi_{N-1,N}(\check{w}_{N-1}, w_N) c_N(w_N) \end{array} \right.$$

with $\psi_{0,1}(w_0, w_1) = \psi_1(w_1)$ and $c_N(w_N) = 1$.

Variable Elimination Variants of Max-Product Algorithm

There are different, but equivalent, ways to formulate the Viterbi solution optimization problem:

$$\tilde{\mathbf{w}}^{(\text{map})} = \arg \max_{\mathbf{w}} P(\mathbf{w}, \mathbf{e}) = \arg \max_{\mathbf{w}} \ln P(\mathbf{w}, \mathbf{e}) = \arg \min_{\mathbf{w}} \left(-\ln P(\mathbf{w}, \mathbf{e}) \right)$$

The first variant leads to the **Max-Sum** formulation:

$$\max_{\mathbf{w}} P(\mathbf{w}, \mathbf{e}) = \max_{w_1, \dots, w_N} \sum_{k=0}^{N-1} \ln \psi(w_k, w_{k+1})$$

The second variant leads to the **Min-Sum** formulation:

$$\max_{\mathbf{w}} P(\mathbf{w}, \mathbf{e}) = \min_{w_1, \dots, w_N} \sum_{k=0}^{N-1} \left(-\ln \psi(w_k, w_{k+1}) \right)$$

These formulations are well-defined because of our running assumption that all distributions are positive which ensures that the (evidence-dependent) potentials

$\psi(w_k, w_{k+1}) = \psi_{k,k+1}(w_k, w_{k+1}; \mathbf{e})$ are always strictly positive. They lead to efficient algorithms which all produce the Viterbi solutions but which can have different computational properties.

We now look in more detail at the Min-Sum formulation. The resulting *Min-Sum Algorithm* is famously known as the **Viterbi Algorithm** (VA), after its inventor Andrew Viterbi. Viterbi is a world-famous former UCSD professor who co-founded Qualcomm with another former UCSD professor Irwin Jacobs.

The Viterbi Algorithm - I

Define

$$\ell(w_k, w_{k+1}) \stackrel{\text{def}}{=} -\ln \psi(w_k, w_{k+1})$$

Then the Min-Sum formulation yields

$$\begin{aligned}\max_{\mathbf{w}} P(\mathbf{w}, \mathbf{e}) &= \min_{w_1, \dots, w_N} \sum_{k=0}^{N-1} \ell(w_k, w_{k+1}) \\&= \min_{w_1} \cdots \min_{w_{k-1}} \min_{w_{k+1}} \cdots \min_{w_N} \left[\ell(w_0, w_1) + \cdots + \ell(w_{k-1}, w_k) + \ell(w_k, w_{k+1}) + \cdots + \ell(w_{N-1}, w_N) \right] \\&= \min_{w_1} \underbrace{\ell(w_0, w_1)}_{\ell(w_1)} + \left[\min_{w_2} \ell(w_1, w_2) + \cdots \right. \\&\quad \left. + \underbrace{\left[\min_{w_k} \ell(w_{k-1}, w_k) + \underbrace{\left[\min_{w_{k+1}} \ell(w_k, w_{k+1}) + \left[\cdots \left[\min_{w_N} \ell(w_{N-1}, w_N) + [0] \right] \cdots \right] \right]}_{O_k(w_k)} \right]}_{O_{k-1}(w_{k-1})} \right] \cdots \end{aligned}$$

Here the function $O_k(w_k)$ is the (additive) optimal cost to go and is computed backwards recursively as

$$O_{k-1}(w_{k-1}) = \min_{w_k} \left(\ell_{k-1,k}(w_{k-1}, w_k; \mathbf{e}) + O_k(w_k) \right)$$

for $k = N-1, \dots, 0$ with boundary condition $O(N) = 0$.

The Viterbi Algorithm - II

Note that the recursion for the additive cost to go $O_k(w)$, which applies in the negative-log domain, looks like the negative logarithm of the recursion for the multiplicative cost to go $c_k(w_k)$ which applies in the original domain.

As before, now that we've determined the optimal cost-to-go functions $O_k(w_k)$ we use backtracking to sequentially produce the Viterbi estimates $\check{w}_k = \check{w}_k^{(\text{map})}$:

$$\text{Viterbi solution } \check{\mathbf{w}} \text{ has components} = \begin{cases} \check{w}_1 = \arg \min_{w_1} \left(\ell_{0,1}(w_0, w_1) + O_1(w_1) \right) \\ \check{w}_2 = \arg \min_{w_2} \left(\ell_{1,2}(\check{w}_1, w_2) + O_2(w_2) \right) \\ \vdots \\ \check{w}_k = \arg \min_{w_k} \left(\ell_{k-1,k}(\check{w}_{k-1}, w_k) + O_k(w_k) \right) \\ \vdots \\ \check{w}_N = \arg \min_{w_N} \left(\ell_{N-1,N}(\check{w}_{N-1}, w_N) + O_N(w_N) \right) \end{cases}$$

with $\ell_{0,1}(w_0, w_1) = \ell_1(w_1)$ and $O_N(w_N) = 0$.