# Viterbi Decoding

In this section we prove that for label sequence learning problems, decoding the top scoring output sequence

$$\hat{y} = \operatorname*{argmax}_{\bar{y} \in \mathcal{Y}(x)} f(x, \bar{y})$$

can be performed by a Viterbi algorithm (Forney, 1973; Schwarz and Chow, 1990) in $\mathcal{O}(T|\Sigma|^2)$, where $f$ is a generalized linear model

$$f(x, y) = \langle \mathbf{w}, \Phi(x, y) \rangle.$$

For notational convenience we will use the notation for $\Phi(x, y)$ introduced by Altun et al. (2003b). Given two labels $\sigma, \tau \in \Sigma$, we define

$$\phi^A_{\sigma,\tau}(y_{t-1}, y_t) = [[y_{t-1} = \sigma \wedge y_t = \tau]]$$
$$\phi^B_{\sigma,j}(x_t, y_t) = [[y_t = \sigma]]\psi_j(x_t)$$

As described in Section 3.2.2, $\psi_j(x)$ extracts characteristics of $x$ (e.g., feature $\psi_{234}(x) = 1$ if token $x_t$ starts with a capital letter and 0 otherwise). We will refer to the vector $\psi(x) = (\ldots, \psi_j(x), \ldots)^\mathsf{T}$ and denote the inner product by means of $k(x, x') = \langle \psi(x), \psi(x') \rangle$.

Analogously to Equation 3.23, Altun et al. (2003b) define the joint feature representation $\Phi(x_i, y_i)$ of the $i$-th sequence as the sum of all feature vectors $\Phi(x_i, y_i|t) = (\ldots, \phi^A_{\sigma,\tau}(y_{i,t-1}, y_{i,t}), \ldots, \phi^B_{\sigma,j}(x_{i,t}, y_{i,t}), \ldots)^\mathsf{T}$ extracted at time $t$

$$\Phi(x_i, y_i) = \sum_{t=1}^{T_i} \Phi(x_i, y_i|t).$$

The feature map in Equation B.5 gives rise to the following inner product in input-output space that decomposes into a *label-label* and a *label-observation*

part,

$$\langle \Phi(x, y), \Phi(x', y') \rangle$$

$$= \sum_{s,t} \sum_{\sigma,\tau} [[y_{s-1} = \sigma \wedge y_s = \tau]] \cdot [[y'_{t-1} = \sigma \wedge y'_t = \tau]]$$

$$+ \sum_{s,t} \sum_{\sigma} [[y_s = \sigma]] \psi(x_s) \cdot [[y'_t = \sigma]] \psi(x'_t)$$

$$= \sum_{s,t} [[y_{s-1} = y'_{t-1} \wedge y_s = y'_t]] + \sum_{s,t} [[y_s = y'_t]] k(x_s, x'_t).$$

The Viterbi algorithm uses dynamic programming to maintain a trellis in which nodes correspond to hidden states versus times. Each entry stores the score of the most probable path leading to that node at a certain time, see Figure B.1. Once the computation reaches the end of the sequence, it backtracks the most likely path through the trellis and returns the highest scoring label sequence that generates the input sequence. The following Proposition B.1 shows that the argument of the maximum can be computed by a Viterbi algorithm.

**Proposition** *Given $n$ input-output pairs of sequences of length $T_i$ for $1 \leq i \leq n$, let $\Sigma$ denote the output alphabet with $|\Sigma| < \infty$. Let $f$ be defined as*

$$f(x, y) = \sum_{i=1}^{n} \sum_{\substack{\bar{y} \in \mathcal{Y}(x_i) \\ \bar{y} \neq y_i}} \alpha_i(\bar{y}) \left( \langle \Phi(x_i, y_i), \Phi(x, y) \rangle - \langle \Phi(x_i, \bar{y}), \Phi(x, y) \rangle \right),$$

*with the joint feature map $\Phi(x, y)$ as in Equation B.5. Then for all $\alpha_i(\bar{y}) \geq 0$ and any $x \in \mathcal{X}$,*

$$\hat{y} = \underset{\bar{y} \in \mathcal{Y}(x)}{\mathrm{argmax}} \, f(x, \bar{y})$$

*can be computed with a Viterbi algorithm in time $\mathcal{O}(T|\Sigma|^2)$.*

**Proof** The model $f$ has the form

$$f(x, y) = \sum_{i=1}^{n} \sum_{\substack{\bar{y} \in \mathcal{Y}(x_i) \\ \bar{y} \neq y_i}} \alpha_i(\bar{y}) \left( \langle \Phi(x_i, y_i), \Phi(x, y) \rangle - \langle \Phi(x_i, \bar{y}), \Phi(x, y) \rangle \right)$$

$$= \sum_{i=1}^{n} \sum_{\substack{\bar{y} \in \mathcal{Y}(x_i) \\ \bar{y} \neq y_i}} \alpha_i(\bar{y}) \left( \sum_{s,t} \left( [[y_{i,s} = y_t]] - [[\bar{y}_s = y_t]] \right) k(x_{i,s}, x_t) \right.$$

$$\left. + \sum_{s,t} [[y_{i,s-1} = y_{t-1} \wedge y_{i,s} = y_t]] - [[\bar{y}_{s-1} = y_{t-1} \wedge \bar{y}_s = y_t]] \right).$$

We make the dependency on labels $\sigma, \tau \in \Sigma$ explicit by summing over all transitions and observation states

$$f(x, y) = \sum_{\sigma, \tau \in \Sigma} \sum_{i} \sum_{\substack{\bar{y} \in \mathcal{Y}(x_i) \\ \bar{y} \neq y_i}} \alpha_i(\bar{y}) \left( \sum_{s,t} ([[y_{i,s} = \sigma]] - [[\bar{y}_s = \sigma]]) [[y_t = \tau]] k(x_{i,s}, x_t) \right.$$

$$+ \sum_{s,t} \left( [[y_{i,s-1} = \sigma \wedge y_{i,s} = \tau]] - [[\bar{y}_{s-1} = \sigma \wedge \bar{y}_s = \tau]] \right) [[y_{t-1} = \sigma \wedge y_t = \tau]] \right).$$

The transition scores from label $\sigma$ to label $\tau$ are now given by

$$a(\sigma, \tau) = \sum_{i=1}^{n} \sum_{\substack{\bar{y} \in \mathcal{Y}(x_i) \\ \bar{y} \neq y_i}} \alpha_i(\bar{y}) \left( \sum_{t=1}^{T_i} [[y_{i,t-1} = \sigma \wedge y_{i,t} = \tau]] - [[\bar{y}_{t-1} = \sigma \wedge \bar{y}_t = \tau]] \right)$$

and observation scores for label $y_s = \sigma$ and observation $x_s$ by

$$b(\sigma, x) = \sum_{i=1}^{n} \sum_{t=1}^{T_i} \sum_{\substack{\bar{y} \in \mathcal{Y}(x_i) \\ \bar{y} \neq y_i}} \alpha_i(\bar{y}) ([[y_{i,t} = \sigma]] - [[\bar{y}_t = \sigma]]) k(x_{i,t}, x).$$

The hypothesis $f(x, y)$ can be rewritten in terms of transition scores $a(\sigma, \tau)$ and observation scores $b(\sigma, x)$

$$f(x, y) = \underbrace{\sum_{\sigma, \tau \in \Sigma} a(\sigma, \tau) \sum_{s=1}^{T} [[y_{s-1} = \sigma \wedge y_s = \tau]]}_{=: f_a(x,y)} + \underbrace{\sum_{s=1}^{T} \sum_{\sigma \in \Sigma} [[y_s = \sigma]] b(\sigma, x_s)}_{=: f_b(x,y)}.$$

where $f_a$ weights the occurences of neighboring labels in $y$ by corresponding scores of the model and $f_b$ determines how well observations $x_s$ fit to their labels $y_s$ given the model. To decode the top scoring sequence we define

$$\delta_t(\sigma) = \max_{y_1, \ldots, y_{t-1}} f(x, y_1, \ldots, y_{t-1}, y_t = \sigma),$$

that is, $\delta_t(\sigma)$ denotes the top scoring partial sequence up to position $t - 1$ where $y_t = \sigma$. We first show by induction that

$$\delta_{t+1}(\sigma) = \max_{\tau \in \Sigma} [\delta_t(\tau) + a(\tau, \sigma)] + b(\sigma, x_{t+1})$$

holds. The initialization is simply given by

$$\delta_0(\sigma) = 0, \quad \forall \sigma \in \Sigma$$
$$\delta_1(\sigma) = \max_{\tau \in \Sigma} [\delta_t(\tau) + a(\tau, \sigma)] + b(\sigma, x_{t+1})$$
$$= a(\epsilon, \sigma) + b(\sigma, x_1).$$

The recursion step is given for $2 \leq t \leq T$ by

$$\delta_t(\sigma) = \max_{y_1, \ldots, y_{t-1}} f(x, y_1, \ldots, y_{t-1}, y_t = \sigma)$$

$$= \max_{y_1, \ldots, y_{t-1}} \sum_{\tau, \bar{\tau} \in \mathcal{Y}} a(\tau, \bar{\tau}) \sum_{s=2}^{t-1} [[y_{s-1} = \tau \wedge y_s = \bar{\tau}]]$$

$$+ \sum_{\tau \in \Sigma} a(\tau, \sigma)[[y_{t-1} = \tau \wedge y_t = \sigma]]$$

$$+ \sum_{s=1}^{t-1} \sum_{\tau \in \Sigma} [[y_s = \tau]] b(\tau, x_s) + [[y_t = \sigma]] b(\sigma, x_t)$$

$$= \max_{\sigma^\star} \max_{y_1, \ldots, y_{t-2}} \sum_{\tau, \bar{\tau} \in \mathcal{Y}} a(\tau, \bar{\tau}) \sum_{s=2}^{t-2} [[y_{s-1} = \tau \wedge y_s = \bar{\tau}]]$$

$$+ \sum_{\tau \in \Sigma} a(\tau, \sigma^\star)[[y_{t-2} = \tau \wedge y_{t-1} = \sigma^\star]]$$

$$+ a(\sigma^\star, \sigma)[[y^{t-1} = \sigma^\star \wedge y^t = \sigma]]$$

$$+ \sum_{s=1}^{t-2} \sum_{\tau \in \Sigma} [[y_s = \tau]] b(\tau, x_s) + b(\sigma^\star, x_{t-1}) + b(\sigma, x_t)$$

$$= \max_{\sigma^\star} \left[ \max_{y_1, \ldots, y_{t-2}} f(x, y_1, \ldots, y_{t-2}, y_{t-1} = \sigma^\star) + a(\sigma^\star, \sigma) \right] + b(\sigma, x_t)$$

$$= \max_{\sigma^\star} \left[ \delta_{t-1}(\sigma^\star) + a(\sigma^\star, \sigma) \right] + b(\sigma, x_t).$$

Thus, the top scoring sequence has the score

$$\max f(x, y) = \max_{\sigma \in \Sigma} \delta_T(\sigma).$$

We only sketch the extension to the argument of the maximum since it is analoguous to the regular Viterbi algorithm. We introduce path variables $\varphi_t(\sigma)$ that are initialized by $\varphi_1(\sigma) = \epsilon$ for all $\sigma \in \Sigma$. The sequence $\varphi_t(\sigma)$ is then defined recursively for $2 \leq t \leq T$ by

$$\varphi_t(\sigma) = \operatorname*{argmax}_{\sigma^\star \in \Sigma} \left[ \delta_{t-1}(\sigma^\star) + a(\sigma^\star, \sigma) \right].$$

Once the $\delta_t(\sigma)$ of Proposition      are fixed, the optimal label sequence can be found by backtracking

$$y_T^\star = \operatorname*{argmax}_{\sigma \in \Sigma} \delta_T(\sigma)$$

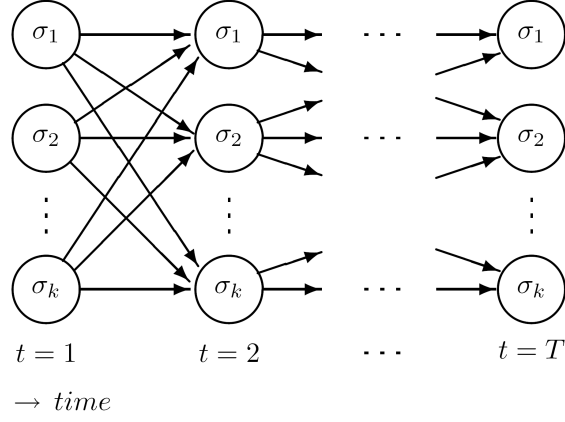$$y_t^\star = \varphi_{t+1}(y_{t+1}^\star) \quad \text{for} \quad t = T - 1, \ldots, 1.$$

**Figure B.1:** *Visualization of a trellis over the alphabet $\Sigma = \{\sigma_1, \ldots, \sigma_k\}$.*

Given the transition matrix $[\mathbf{A}]_{\sigma,\tau} = a(\sigma, \tau)$ and the observation matrix $[\mathbf{B}_{\chi}]_{\sigma,t} = b(\sigma, x_t)$ for input $\chi$, the computation of $\delta$ and $\varphi$ for a fixed $t$ and $\sigma \in \Sigma$ involves visiting $|\Sigma|$ predecessors; thus, for a sequence of length $T$ the time needed is in $\mathcal{O}(T|\Sigma|^2)$. This concludes the proof. $\square$