

Probabilistic Time Series Analysis: Lab 5

Tim Kunisky

September 26, 2018

The Discrete Hidden Markov Model

Same setting as with LDS, only now all variables are discrete:

$$\begin{array}{ccccccc} \cdots & \rightarrow & z_{t-2} & \rightarrow & z_{t-1} & \rightarrow & z_t & \rightarrow & \cdots \\ & & \downarrow & & \downarrow & & \downarrow & & \\ \cdots & & x_{t-2} & & x_{t-1} & & x_t & & \cdots \end{array}$$

States in $\{1, \dots, K\}$; observations in $\{1, \dots, H\}$.

$$\mathbb{P}(z_t = \ell \mid z_{t-1} = k) = A_{k\ell} \quad (\mathbf{A} \in \mathbb{R}^{K \times K} \text{ is the } \textit{transition matrix})$$

$$\mathbb{P}(z_1 = k) = \pi_k \quad (\boldsymbol{\pi} \text{ is the } \textit{initial distribution})$$

$$\mathbb{P}(x_t = i \mid z_t = k) = B_{ki} \quad (\mathbf{B} \in \mathbb{R}^{K \times H} \text{ is the } \textit{emission matrix})$$

The model parameters are $\boldsymbol{\theta} = \{\mathbf{A}, \mathbf{B}, \boldsymbol{\pi}\}$.

Exercise: Let's write down the joint distribution $\mathbb{P}(\mathbf{z}, \mathbf{x} \mid \boldsymbol{\theta})$.

Computation 1: Likelihood Function

The likelihood is

$$\mathbb{P}(\mathbf{x} \mid \boldsymbol{\theta}) = \sum_{k_1=1}^K \cdots \sum_{k_T=1}^K \mathbb{P}(z_1 = k_1, \dots, z_T = k_T, \mathbf{x} \mid \boldsymbol{\theta}).$$

How do we compute it without summing over K^T terms? We use the Markov property of the z_t , which allows a recursion:

$$\begin{aligned}\alpha_t(k) &= \mathbb{P}(x_1, \dots, x_t, z_t = k). \\ &= \sum_{\ell=1}^K \mathbb{P}(x_{[1:t-1]}, z_{t-1} = \ell) \mathbb{P}(z_t = k \mid z_{t-1} = \ell) \mathbb{P}(x_t \mid z_t = k) \\ &= \sum_{\ell=1}^K \alpha_{t-1}(\ell) A_{\ell,k} B_{k,x_t}.\end{aligned}$$

At the end we can take $\mathbb{P}(\mathbf{x} \mid \boldsymbol{\theta}) = \sum_{k=1}^K \alpha_T(k)$.

Likelihood Function with Linear Algebra

The recursion for the likelihood is:

$$\alpha_t(k) = \sum_{\ell=1}^K \alpha_{t-1}(\ell) A_{\ell,k} B_{k,x_t}.$$

This is simple to express in terms of matrix operations: let's say $\alpha_t \in \mathbb{R}^K$, then

$$\begin{aligned} \alpha_t^\top &= \alpha_{t-1}^\top A \text{diag}(B_{\bullet, x_t}) \\ &= \alpha_1^\top A \text{diag}(B_{\bullet, x_t}) \text{diag}(B_{\bullet, x_{t-1}}) \cdots \text{diag}(B_{\bullet, x_2}) \end{aligned}$$

Every entry of each $\text{diag}(B_{\bullet, x_t})$ is at most 1 (it is a probability!) and $\|A\| \leq 1$, so these quantities will shrink quickly—hence the normalization step from lecture.

Computation 2: MAP Estimation / Decoding

Now, we see \mathbf{x} and want to compute

$$\hat{\mathbf{z}} = \text{optimizer in } \max_{k_1} \cdots \max_{k_T} \mathbb{P}(z_1 = k_1, \dots, z_T = k_T, \mathbf{x} \mid \boldsymbol{\theta}).$$

Formally, this is just like the last problem, but with “sum” \rightarrow “max.” So, we can use the same recursion:

$$\begin{aligned} v_t(\ell) &= \max_{k_1, \dots, k_{t-1}} \mathbb{P}(\mathbf{x}_{[1:t]}, \mathbf{z}_{[1:t-1]} = k_{[1:t-1]}, z_t = \ell) \\ &= \max_{k_{t-1}} \max_{k_1, \dots, k_{t-2}} \mathbb{P}(\mathbf{x}_{[1:t-1]}, \mathbf{z}_{[1:t-2]} = k_{[1:t-2]}, z_{t-1} = k_{t-1}) \\ &\quad \mathbb{P}(z_t = \ell \mid z_{t-1} = k_{t-1}) \mathbb{P}(x_t \mid z_t = \ell) \\ &= \max_{k_{t-1}} v_{t-1}(k_{t-1}) A_{k_{t-1}, \ell} B_{\ell, x_t}. \end{aligned}$$

However: we want to actually learn $\hat{\mathbf{z}}$, so we also keep track of the optimizers k_i as we go along.

Coding Assignment

- Find `labs/lab5/lab5-student.ipynb` on the course Github. (It will be easier if you clone the whole repository.)
- We're going to do very real data science! The task is **part of speech (POS) tagging**.
- This tags English words by their grammatical function in a sentence (nouns, verbs, adjectives, adverbs, conjunctions, etc.), which is not always easy: "I will **book** the hotel" vs. "I will read the **book**."
- You will implement HMM for this on the standard Wall Street Journal dataset, containing 39,815 sentences for training and 1,700 sentences for testing.