

Probabilistic Time Series Analysis: Lab 3

Tim Kunisky

September 19, 2018

The Linear Gaussian Hidden Markov Model

The Linear Gaussian Hidden Markov Model

- High level: the Kalman filter is an *algorithm for computing certain conditional distributions* for certain time series models. It is not a model! It is not a way to do inference!

The Linear Gaussian Hidden Markov Model

- High level: the Kalman filter is an *algorithm for computing certain conditional distributions* for certain time series models. It is not a model! It is not a way to do inference!
- Let's review the model that we covered Kalman filtering for:

$$\begin{array}{ccccccc} \cdots & \rightarrow & \mathbf{z}_{t-2} & \rightarrow & \mathbf{z}_{t-1} & \rightarrow & \mathbf{z}_t & \rightarrow & \cdots \\ & & \downarrow & & \downarrow & & \downarrow & & \\ \cdots & & \mathbf{x}_{t-2} & & \mathbf{x}_{t-1} & & \mathbf{x}_t & & \cdots \end{array}$$

$$\mathbf{z}_t = \mathbf{A}\mathbf{z}_{t-1} + \mathbf{w}_t$$

$$\mathbf{x}_t = \mathbf{C}\mathbf{z}_t + \mathbf{v}_t$$

$$\mathbf{w}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{Q})$$

$$\mathbf{v}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{R})$$

$$\mathbf{z}_0 \sim \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma})$$

Inference vs. Filtering vs. Smoothing

Inference vs. Filtering vs. Smoothing

- **Question:** Describe the goal of the inference task for this model.

Inference vs. Filtering vs. Smoothing

- **Question:** Describe the goal of the inference task for this model.
 - Kalman filtering will not do this for you! This is what we will use the EM algorithm for.

Inference vs. Filtering vs. Smoothing

- **Question:** Describe the goal of the inference task for this model.
 - Kalman filtering will not do this for you! This is what we will use the EM algorithm for.
- Once we know the parameters $A, C, Q, R, \mu_0, \Sigma$, we still need the distributions of the unobserved latent states z_t .

Inference vs. Filtering vs. Smoothing

- **Question:** Describe the goal of the inference task for this model.
 - Kalman filtering will not do this for you! This is what we will use the EM algorithm for.
- Once we know the parameters $A, C, Q, R, \mu_0, \Sigma$, we still need the distributions of the unobserved latent states z_t .
- There are two basic settings:

Inference vs. Filtering vs. Smoothing

- **Question:** Describe the goal of the inference task for this model.
 - Kalman filtering will not do this for you! This is what we will use the EM algorithm for.
- Once we know the parameters $A, C, Q, R, \mu_0, \Sigma$, we still need the distributions of the unobserved latent states \mathbf{z}_t .
- There are two basic settings:
 - **Filtering:** You have observations of $\mathbf{x}_1, \dots, \mathbf{x}_t$, and want to compute the distribution of \mathbf{z}_t , e.g. if you are trying to estimate something about a physical system “live” or “online.”

Inference vs. Filtering vs. Smoothing

- **Question:** Describe the goal of the inference task for this model.
 - Kalman filtering will not do this for you! This is what we will use the EM algorithm for.
- Once we know the parameters $A, C, Q, R, \mu_0, \Sigma$, we still need the distributions of the unobserved latent states \mathbf{z}_t .
- There are two basic settings:
 - **Filtering:** You have observations of $\mathbf{x}_1, \dots, \mathbf{x}_t$, and want to compute the distribution of \mathbf{z}_t , e.g. if you are trying to estimate something about a physical system “live” or “online.”
 - **Smoothing:** You have observations of $\mathbf{x}_1, \dots, \mathbf{x}_T$, and want to compute the distribution of \mathbf{z}_t with $t < T$, e.g. if you have gathered an entire time series of observations and want to estimate the latent states afterwards.

Filtering

Filtering

- **Goal:** $\mathbb{P}[\mathbf{z}_n \mid \mathbf{x}_1, \dots, \mathbf{x}_n] \sim \mathcal{N}(\boldsymbol{\mu}_{n|n}, \boldsymbol{\Sigma}_{n|n})$, compute RHS parameters for all n .

Filtering

- **Goal:** $\mathbb{P}[\mathbf{z}_n \mid \mathbf{x}_1, \dots, \mathbf{x}_n] \sim \mathcal{N}(\boldsymbol{\mu}_{n|n}, \boldsymbol{\Sigma}_{n|n})$, compute RHS parameters for all n .
- **Prediction Step:** $\mathbb{P}[\mathbf{z}_n \mid \mathbf{x}_1, \dots, \mathbf{x}_{n-1}] \sim \mathcal{N}(\boldsymbol{\mu}_{n|n-1}, \boldsymbol{\Sigma}_{n|n-1})$.

$$\boldsymbol{\mu}_{n|n-1} = \mathbf{A}\boldsymbol{\mu}_{n-1|n-1}$$

$$\boldsymbol{\Sigma}_{n|n-1} = \mathbf{A}\boldsymbol{\Sigma}_{n-1|n-1}\mathbf{A}^\top + \mathbf{Q}$$

Filtering

- **Goal:** $\mathbb{P}[\mathbf{z}_n \mid \mathbf{x}_1, \dots, \mathbf{x}_n] \sim \mathcal{N}(\boldsymbol{\mu}_{n|n}, \boldsymbol{\Sigma}_{n|n})$, compute RHS parameters for all n .
- **Prediction Step:** $\mathbb{P}[\mathbf{z}_n \mid \mathbf{x}_1, \dots, \mathbf{x}_{n-1}] \sim \mathcal{N}(\boldsymbol{\mu}_{n|n-1}, \boldsymbol{\Sigma}_{n|n-1})$.

$$\boldsymbol{\mu}_{n|n-1} = \mathbf{A}\boldsymbol{\mu}_{n-1|n-1}$$

$$\boldsymbol{\Sigma}_{n|n-1} = \mathbf{A}\boldsymbol{\Sigma}_{n-1|n-1}\mathbf{A}^\top + \mathbf{Q}$$

- **Observation Step:** Add conditioning on \mathbf{x}_n .

$$\tilde{\mathbf{x}}_n = \mathbf{x}_n - \mathbf{C}\boldsymbol{\mu}_{n|n-1}$$

$$\tilde{\mathbf{R}}_n = \mathbf{C}\boldsymbol{\Sigma}_{n|n-1}\mathbf{C}^\top + \mathbf{R}$$

$$\mathbf{K}_n = \boldsymbol{\Sigma}_{n|n-1}\mathbf{C}^\top \tilde{\mathbf{R}}_n^{-1}$$

$$\boldsymbol{\mu}_{n|n} = \boldsymbol{\mu}_{n|n-1} + \mathbf{K}_n \tilde{\mathbf{x}}_n$$

$$\boldsymbol{\Sigma}_{n|n} = (\mathbf{I} - \mathbf{K}_n\mathbf{C})\boldsymbol{\Sigma}_{n|n-1}$$

Smoothing

Smoothing

- **Goal:** $\mathbb{P}[\mathbf{z}_n \mid \mathbf{x}_1, \dots, \mathbf{x}_N] \sim \mathcal{N}(\hat{\boldsymbol{\mu}}_n, \hat{\boldsymbol{\Sigma}}_n)$, compute RHS parameters, assuming you know them for all smaller n .

Smoothing

- **Goal:** $\mathbb{P}[\mathbf{z}_n \mid \mathbf{x}_1, \dots, \mathbf{x}_N] \sim \mathcal{N}(\hat{\boldsymbol{\mu}}_n, \hat{\boldsymbol{\Sigma}}_n)$, compute RHS parameters, assuming you know them for all smaller n .
- **Forward Path:** Run filtering on the entire sequence, from $n = 0$ to $n = N$.

Smoothing

- **Goal:** $\mathbb{P}[\mathbf{z}_n \mid \mathbf{x}_1, \dots, \mathbf{x}_N] \sim \mathcal{N}(\hat{\boldsymbol{\mu}}_n, \hat{\boldsymbol{\Sigma}}_n)$, compute RHS parameters, assuming you know them for all smaller n .
- **Forward Path:** Run filtering on the entire sequence, from $n = 0$ to $n = N$.
- **Backward Path:** Step backwards from N to 1, making corrections to the filtering predictions iteratively. At N , the filtering and smoothing predictions match, and we propagate the “knowledge of the full series” backwards through the predictions.

$$\hat{\boldsymbol{\mu}}_N = \boldsymbol{\mu}_{N|N}$$

$$\hat{\boldsymbol{\Sigma}}_N = \boldsymbol{\Sigma}_{N|N}$$

$$\mathbf{F}_n = \boldsymbol{\Sigma}_{n|n} \mathbf{A}^\top \boldsymbol{\Sigma}_{n+1|n}^{-1}$$

$$\hat{\boldsymbol{\mu}}_n = \boldsymbol{\mu}_{n|n} + \mathbf{F}_n (\hat{\boldsymbol{\mu}}_{n+1} - \mathbf{A} \boldsymbol{\mu}_{n|n})$$

$$\hat{\boldsymbol{\Sigma}}_n = \boldsymbol{\Sigma}_{n|n} + \mathbf{F}_n (\hat{\boldsymbol{\Sigma}}_{n+1} - \boldsymbol{\Sigma}_{n+1|n}) \mathbf{F}_n^\top$$

Inference: Review of the EM Algorithm

Inference: Review of the EM Algorithm

- This is a way of doing *maximum likelihood*, maximizing

$$L(\boldsymbol{\theta}) = \log \mathbb{P}[\boldsymbol{x} \mid \boldsymbol{\theta}] = \log \int_{\boldsymbol{z}} \mathbb{P}[\boldsymbol{x}, \boldsymbol{z} \mid \boldsymbol{\theta}] d\boldsymbol{z}.$$

Inference: Review of the EM Algorithm

- This is a way of doing *maximum likelihood*, maximizing

$$L(\boldsymbol{\theta}) = \log \mathbb{P}[\boldsymbol{x} \mid \boldsymbol{\theta}] = \log \int_{\boldsymbol{z}} \mathbb{P}[\boldsymbol{x}, \boldsymbol{z} \mid \boldsymbol{\theta}] d\boldsymbol{z}.$$

- The idea is to introduce a distribution over \boldsymbol{z} , called \mathbb{Q} , and take an inequality to make the optimization easier:

Inference: Review of the EM Algorithm

- This is a way of doing *maximum likelihood*, maximizing

$$L(\boldsymbol{\theta}) = \log \mathbb{P}[\mathbf{x} \mid \boldsymbol{\theta}] = \log \int_{\mathbf{z}} \mathbb{P}[\mathbf{x}, \mathbf{z} \mid \boldsymbol{\theta}] d\mathbf{z}.$$

- The idea is to introduce a distribution over \mathbf{z} , called Q , and take an inequality to make the optimization easier:

$$\begin{aligned} L(\boldsymbol{\theta}) &= \log \int_{\mathbf{z}} Q[\mathbf{z}] \frac{\mathbb{P}[\mathbf{x}, \mathbf{z} \mid \boldsymbol{\theta}]}{Q[\mathbf{z}]} d\mathbf{z} \\ &= \log \mathbb{E}_Q \left[\frac{\mathbb{P}[\mathbf{x}, \mathbf{z} \mid \boldsymbol{\theta}]}{Q[\mathbf{z}]} \right] \\ &\geq \mathbb{E}_Q \left[\log \frac{\mathbb{P}[\mathbf{x}, \mathbf{z} \mid \boldsymbol{\theta}]}{Q[\mathbf{z}]} \right] \\ &= \mathbb{E}_Q [\log \mathbb{P}[\mathbf{x}, \mathbf{z} \mid \boldsymbol{\theta}]] - \mathbb{E}_Q [\log Q[\mathbf{z}]] . \end{aligned}$$

Inference: Review of the EM Algorithm

- Why would this be a good idea??

Inference: Review of the EM Algorithm

- Why would this be a good idea??
- When is there approximate equality in Jensen's inequality:

$$\log \mathbb{E}_{\mathbf{Q}}[f(\mathbf{z})] \approx \mathbb{E}_{\mathbf{Q}}[\log f(\mathbf{z})]?$$

Inference: Review of the EM Algorithm

- Why would this be a good idea??
- When is there approximate equality in Jensen's inequality:

$$\log \mathbb{E}_{\mathbf{Q}}[f(\mathbf{z})] \approx \mathbb{E}_{\mathbf{Q}}[\log f(\mathbf{z})]?$$

- **When $f(\mathbf{z}) \approx \text{constant}$.**

Inference: Review of the EM Algorithm

- Why would this be a good idea??
- When is there approximate equality in Jensen's inequality:

$$\log \mathbb{E}_{\mathbb{Q}}[f(\mathbf{z})] \approx \mathbb{E}_{\mathbb{Q}}[\log f(\mathbf{z})]?$$

- **When $f(\mathbf{z}) \approx \text{constant}$.**
- What does this mean for us? We want

$$\frac{\mathbb{P}[\mathbf{x}, \mathbf{z} \mid \boldsymbol{\theta}]}{\mathbb{Q}[\mathbf{z}]} \approx c(\mathbf{x}, \boldsymbol{\theta})$$

Inference: Review of the EM Algorithm

- Why would this be a good idea??
- When is there approximate equality in Jensen's inequality:

$$\log \mathbb{E}_{\mathbb{Q}}[f(\mathbf{z})] \approx \mathbb{E}_{\mathbb{Q}}[\log f(\mathbf{z})]?$$

- **When $f(\mathbf{z}) \approx \text{constant}$.**
- What does this mean for us? We want

$$\frac{\mathbb{P}[\mathbf{x}, \mathbf{z} \mid \boldsymbol{\theta}]}{\mathbb{Q}[\mathbf{z}]} \approx c(\mathbf{x}, \boldsymbol{\theta})$$

or...

$$\mathbb{P}[\mathbf{x}, \mathbf{z} \mid \boldsymbol{\theta}] \approx \mathbb{Q}[\mathbf{z}]c(\mathbf{x}, \boldsymbol{\theta}),$$

so the likelihood should approximately factorize!

Inference: Review of the EM Algorithm

- Why would this be a good idea??
- When is there approximate equality in Jensen's inequality:

$$\log \mathbb{E}_{\mathbb{Q}}[f(\mathbf{z})] \approx \mathbb{E}_{\mathbb{Q}}[\log f(\mathbf{z})]?$$

- **When $f(\mathbf{z}) \approx \text{constant}$.**
- What does this mean for us? We want

$$\frac{\mathbb{P}[\mathbf{x}, \mathbf{z} \mid \boldsymbol{\theta}]}{\mathbb{Q}[\mathbf{z}]} \approx c(\mathbf{x}, \boldsymbol{\theta})$$

or...

$$\mathbb{P}[\mathbf{x}, \mathbf{z} \mid \boldsymbol{\theta}] \approx \mathbb{Q}[\mathbf{z}]c(\mathbf{x}, \boldsymbol{\theta}),$$

so the likelihood should approximately factorize!

- **General principle:** The EM algorithm encodes an approximation of the likelihood as factorizing, and optimizes over this approximation.

Coding Assignment

- Find `labs/lab3/lab3-student.ipynb` on the course Github. (It will be easier if you clone the whole repository.)
- Try running everything, but you will only be graded on what you fill in for the missing bits of code marked with a `TODO`. This time, these are all methods of one class for doing sampling, filtering, and smoothing for the model we talked about.