

Homework 3: Inference & Learning for Temporal State Space Models

Brown University CS242: Probabilistic Graphical Models, Fall 2014

Homework due at 11:59pm on October 30, 2014

In this problem set, we explore inference and learning algorithms for the *switching state space model*, which is summarized by the graphical model of Fig. 1. Our observations are a sequence of real-valued vectors, $y_t \in \mathbb{R}^p$, which depend on corresponding latent state vectors $x_t \in \mathbb{R}^d$. The temporal evolution of the state and observation vectors is influenced by discrete “switch” variables $z_t \in \{1, 2, \dots, K\}$ as follows:

$$\begin{aligned} p(y_t \mid x_t, z_t = k) &= \text{Norm}(y_t \mid C_k x_t, R_k) \\ p(x_t \mid x_{t-1}, z_t = k) &= \text{Norm}(x_t \mid A_k x_{t-1}, Q_k) \end{aligned}$$

If $K = 1$ so that $z_t = 1$ for all times, we recover a standard linear Gaussian state space model. In this case, we omit the subscripts above and denote the model parameters by $\{A, C, Q, R\}$. More generally, the discrete switch variables evolve according to a Markov chain:

$$p(z_t \mid z_{t-1} = k) = \text{Cat}(z_t \mid \pi_k)$$

Thus, $\pi_{k\ell}$ is the probability that $z_t = \ell$ given $z_{t-1} = k$. We assume that z_1 is uniformly distributed across the K possible switch modes, and $p(x_1) = \text{Norm}(x_1 \mid 0, I_d)$.

In general, many different state space models may assign equal likelihood to a given dataset, due to generalized versions of the rotational ambiguities underlying PCA and factor analysis models. For all models we consider, we thus constrain $C = [I_p \ O]$, where O is a $p \times (d - p)$ matrix of zeros, and I_p is a $p \times p$ identity matrix. When $d = p$, we have $C = I_p$.

For some experiments, we define our state space model parameters to take one of a few canonical forms. The *constant position* model takes $d = p$, and

$$A = I_p, \quad Q = \sigma_x^2 I_p, \quad R = \sigma_y^2 I_p,$$

so that the observations are noisy estimates of an underlying random walk. The *constant velocity* model generates more smooth motions by taking $d = 2p$, and

$$A = \begin{bmatrix} I_p & I_p \\ O_p & I_p \end{bmatrix}, \quad Q = \begin{bmatrix} \sigma_x^2 I_p & O_p \\ O_p & \sigma_v^2 I_p \end{bmatrix}, \quad R = \sigma_y^2 I_p,$$

where O_p is a $p \times p$ matrix of zeros. More generally, we will use the expectation maximization (EM) algorithm to learn model parameters $\{A, Q, R\}$ from observation sequences.

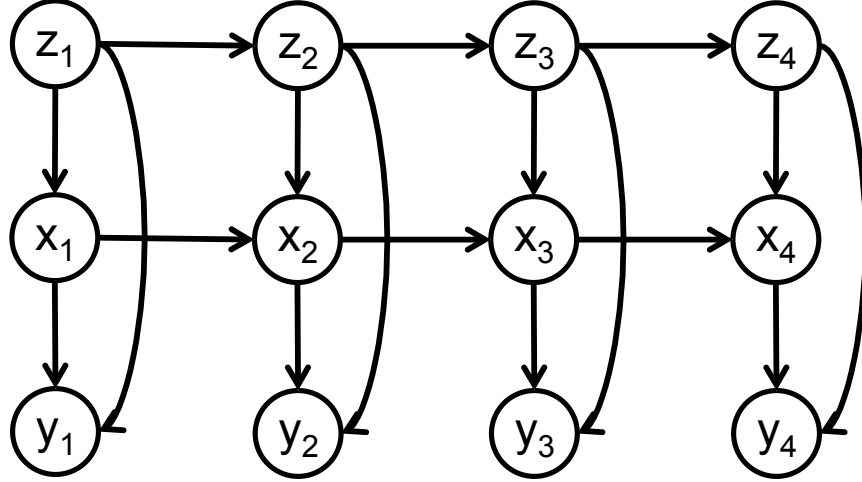


Figure 1: A directed graphical model representing a switching state space model. Discrete switch variables z_t evolve according to a Markov chain. At each time step, the switch z_t selects which state space model determines the continuous state x_t and observations y_t .

Question 1: Tracking, Kalman Filters, & Particle Filters

Throughout this question we assume a non-switching state space model ($K=1$). In the handout code, see `kalman_smoother.m` for an implementation of the Kalman filter. Given a sequence of T observations, this code produces a $d \times T$ matrix of posterior mean estimates, and a $d \times d \times T$ array of posterior covariance estimates.

- Consider the $p = 1$ -dimensional observation sequence we provide in `track`. It was sampled from a constant velocity model with $d = 2$, $\sigma_v^2 = 0.01$, $\sigma_x^2 = 0.01/3$, $\sigma_y^2 = 20$. Apply the Kalman filter code using this correct model. On one set of axes, plot the observation sequence, the posterior mean of the first (position) component of the state, and posterior confidence intervals. The confidence intervals should be determined as the posterior mean plus or minus two times the posterior standard deviation of the first state component.
- Suppose you incorrectly assumed the `track` data was generated by a constant-position model with $d = 1$ and $\sigma_y^2 = 20$. Consider two possible state-transition noise levels, $\sigma_x^2 = 0.01/3$ and $\sigma_x^2 = 10$. For each of these alternative models, apply the Kalman filter code and plot your results as in part (a). Discuss differences from the results in part (a).
- Implement a basic “bootstrap” particle filter, which proposes particles from the state transition distribution $p(x_t | x_{t-1})$, reweights with the observation likelihood $p(y_t | x_t)$, and resamples weighted particles (with replacement) at each time step.
- Apply your particle filter implementation to the model and data from part (a), using 100 particles. At each time step, estimate the posterior mean as the weighted mean (before resampling) of the current particle locations. Repeat this experiment three times, and on a single set of axes, plot the three particle filter mean estimates together with the optimal Kalman filter mean estimates.

- e) Repeat part (d) using only 20 particles. Discuss the robustness of the particle filter, compared to the Kalman filter.
- f) We finish by exploring robustness to outliers. Again consider the `track` data and independently at each time step t , with probability 0.1 replace the true observation by a sample $y_t \sim \text{Norm}(0, 40^2)$. Apply the Kalman filter to this corrupted data, with the (now inaccurate) original observation likelihood. Also apply a particle filter with 100 particles, and an observation likelihood that correctly models the outlier process. Plot the corrupted observation data, as well as the mean estimates for both methods, and discuss.

Question 2: Learning, Kalman Smoothers, & Expectation Maximization (EM)

Throughout this question we assume a non-switching state space model ($K=1$). In the handout code, see `kalman_smoother.m` for an implementation of the Kalman smoother.

- a) Apply the Kalman smoother to the constant velocity model and `track` data from part 1(a). On one set of axes, plot the observation sequence, the posterior mean of the first (position) component of the state, and posterior confidence intervals. Compare to the estimates produced by the Kalman filter.
- b) The marginal log-likelihood of an observation sequence y , integrating over states x for some fixed state space model parameters, can be written as follows:

$$\begin{aligned} \log p(y) &= \log p(y_1) + \sum_{t=1}^{T-1} \log p(y_{t+1} \mid y_1, \dots, y_t) \\ &= \log \int_{\mathcal{X}_1} p(y_1 \mid x_1) p(x_1) dx_1 + \sum_{t=1}^{T-1} \log \int_{\mathcal{X}_t} p(y_{t+1} \mid x_t) p(x_t \mid y_1, \dots, y_t) dx_t \end{aligned}$$

Use these expressions to efficiently evaluate the log-likelihood in closed form, based on the output of the Kalman filter. Write code implementing this log-likelihood computation.

- c) We have provided a partial implementation of the EM algorithm for linear state space models in `em_lds.m`. The Kalman smoother provides the E-step statistics, and we provide the updates for A , Q , and C . You need to add code for computing log-likelihood bounds as in part (b), and for the M-step update of the observation covariance R . Specify the mathematical formula for this M-step update and argue why it is correct, but a detailed derivation is not necessary. **Hint:** Information on EM for state space models can be found in Sec. 24.5.2 of Barber's Bayesian Reasoning and Machine Learning. The update for R is similar to the M-step for the simpler factor analysis model.
- d) Consider the $p = 2$ -dimensional observation sequence we provide in `spiral`, and use the EM algorithm to learn a $d = 2$ -dimensional state space model. Initialize with a constant position model with $\sigma_x^2 = 1, \sigma_y^2 = 1$. Run the EM algorithm for 100 iterations, plot the marginal log-likelihood versus iteration, and report the learned parameters after the final iteration. Plot the observation sequence, and the output of the Kalman smoother after the final iteration, as overlaid 2-dimensional curves.

Question 3: Tracking in Switching State Space Models

We now consider the full switching state space model, and apply it to modeling the tracked motions of dancing honeybees. You can read more about this data at

http://www.cc.gatech.edu/~borg/ijcv_psslds/

Honeybees communicate via dances that can be approximated as switching among $K = 3$ behaviors: *waggle*, *turn left*, *turn right*. Using training data in which behaviors z_t and bee poses y_t are both available, we will learn a switching state space model of bee motion.

We provide three sequences of bee motion. In the questions below we analyze three train-test splits, where each split uses two training sequences and one test sequence. Our bee pose representation y_t is $p = 4$ -dimensional: two position coordinates, and two orientation coordinates. For an orientation angle θ , the orientation coordinates are $\cos(\theta)$, $\sin(\theta)$, and we approximate their behavior-specific conditional distributions as Gaussian to simplify learning and inference. To estimate an angle from inferred angle coordinates, first renormalize those coordinates to lie on the unit circle, and then apply an inverse-trigonometric function.

- a) *Assuming training data in which the active behaviors z_t are observed, specify formulas for maximum likelihood estimation of the 3×3 transition matrix π . Compute and report these ML estimates for each training split.*
- b) *Generalize the EM algorithm implementation from problem 2 to allow learning from multiple sequences independently sampled from a common set of state space model parameters. The E-step should apply the Kalman filter to each sequence independently. The M-step should be based on sufficient statistics collected from all of the training sequences. No formal derivation is required.*
- c) *In general, bees stay in behavior modes for many time steps before switching to other behavior types. Use the provided `load_bees.m` method to extract “segments” corresponding to the contiguous time blocks where $z_t = k$ for $k = 1, 2, 3$, and store these as a set of training sequences for learning the state space model parameters $\{A_k, C_k, Q_k, R_k\}$.*
- d) *Set $d = p = 4$, and use the EM algorithm from part (b) to learn a state space model for each of the three behavior modes, using the data from part (c). In each case, initialize using a constant position model with $\sigma_y^2 = 1$, $\sigma_x^2 = 1$ for location variables, and $\sigma_x^2 = 0.05$ for angular variables. For each train-test split, plot the log-likelihood bounds after each iteration, and report the learned model parameters.*
- e) *We now apply the learned state space model to track bee motion. Simulating the behavior of sometimes-inaccurate vision-based trackers, generate a noisy observation sequence \tilde{y} which depends on the true poses y as follows:*

$$p(\tilde{y}_t | y_t) = 0.9 \text{Norm}(\tilde{y}_t | y_t, 0.1I_4) + 0.1 \text{Norm}(\tilde{y}_t | 0, 5I_4)$$

Derive an explicit formula for $p(\tilde{y}_t | x_t)$.

- f) *Apply a “bootstrap” particle filter to bee tracking, where each particle represents a joint configuration (z_t, x_t) of the hidden variables at a given time point. For each train-test split, use switching state-space model parameters learned in parts (a) and (d), a noisy observation sequence generated as in part (e), and 1000 particles.*

- g) Analyze the results from part (f). For each sequence, plot (i) the inferred probability that each behavior is active over time, together with the true held-out behavior labels; (ii) the posterior mean of the bee's two-dimensional position over time, together with the true position; (iii) the posterior mean of the bee's orientation angle over time (after mapping angle coordinates to the unit circle), together with the true angle. Discuss your results.*