

L08 – Continuous factored models

Barber 24.1,3,4; 27.6 AIMA 13.2.3, 14.4, 14.5.3

What you should know after this lecture

- Kalman filtering
- Particle filter

Linear Gaussian Hidden Markov Models

$$\text{Process step: } \mathbf{x}_t = A\mathbf{x}_{t-1} + \mathbf{w}_{t-1}$$

$$\text{Measurement step: } \mathbf{y}_t = H\mathbf{x}_t + \mathbf{v}_t$$

where

- \mathbf{x}_t — state vector at time step t , a random process
- \mathbf{y}_t — observation received at time step t , a random process
- \mathbf{w}_t — process noise $\sim N(\mathbf{0}, W)$
- \mathbf{v}_t — measurement noise $\sim N(\mathbf{0}, R)$
- A — process model (Note: this is not the same kind of matrix as A in the HMM, although it plays a similar role.)
- H — measurement model
- π — initial distribution $N(\mathbf{x}_0, Q_{0|0})$
- We are ignoring the control term $B\mathbf{u}_{t-1}$ (ignore this comment if it doesn't worry you)
- A and H are assumed known and constant, but could vary
- Continuous time version possible, but hairy

Filtering

Want to compute $P(\mathbf{x}_t \mid \mathbf{y}_{0:t})$.

- We know it's Gaussian because this is a linear Gaussian Bayesian network!
- So $P(\mathbf{x}_t \mid \mathbf{y}_{0:t}) = \mathcal{N}(\hat{\mathbf{x}}_{t|t}, \mathbf{Q}_{t|t})$
- Assume we know parameters of distribution at previous step $\hat{\mathbf{x}}_{t|t-1}, \mathbf{Q}_{t|t-1}$. Note that π is our base case.
- Recursively compute
 1. Transition update finds

$$P(\mathbf{x}_t \mid \mathbf{y}_{0:t-1}) = \mathcal{N}(\hat{\mathbf{x}}_{t|t-1}, \mathbf{Q}_{t|t-1})$$

2. Observation update finds

$$P(\mathbf{x}_t \mid \mathbf{y}_{0:t}) = \mathcal{N}(\hat{\mathbf{x}}_{t|t}, \mathbf{Q}_{t|t})$$

- Can be understood as sum-product on associated Gaussian factor graph

Transition update

- Current belief $P(\mathbf{x}_{t-1} | \mathbf{y}_{0:t-1}) = \mathcal{N}(\hat{\mathbf{x}}_{t-1|t-1}, \mathbf{Q}_{t-1|t-1})$
- Transition $\mathbf{x}_t = \mathbf{A}\mathbf{x}_{t-1} + \mathbf{w}_t$ where $\mathbf{w}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{W})$
- Construct the joint on \mathbf{x}_{t-1} and \mathbf{x}_t :

$$\boldsymbol{\mu} = \begin{pmatrix} \hat{\mathbf{x}}_{t-1|t-1} \\ \mathbf{A}\hat{\mathbf{x}}_{t-1|t-1} \end{pmatrix} \quad \boldsymbol{\Sigma} = \begin{pmatrix} \mathbf{Q}_{t-1|t-1} & \mathbf{Q}_{t-1|t-1}\mathbf{A}^T \\ \mathbf{A}\mathbf{Q}_{t-1|t-1} & \mathbf{A}\mathbf{Q}_{t-1|t-1}\mathbf{A}^T + \mathbf{W} \end{pmatrix}$$

- Marginalize out \mathbf{x}_{t-1}

$$\begin{aligned} P(\mathbf{x}_t | \mathbf{y}_{0:t-1}) &= \mathcal{N}(\hat{\mathbf{x}}_{t|t-1}, \mathbf{Q}_{t|t-1}) \\ \hat{\mathbf{x}}_{t|t-1} &= \mathbf{A}\hat{\mathbf{x}}_{t-1|t-1} \\ \mathbf{Q}_{t|t-1} &= \mathbf{A}\mathbf{Q}_{t-1|t-1}\mathbf{A}^T + \mathbf{W} \end{aligned}$$

Note that $\text{Var}[\mathbf{A} + \mathbf{B}] = \text{Var}[\mathbf{A}] + \text{Var}[\mathbf{B}]$ when \mathbf{A} and \mathbf{B} are independent. Here \mathbf{x}_{t-1} and \mathbf{w}_t are independent. Also $\text{Var}[\mathbf{C}\mathbf{A} + \mathbf{c}]$, where \mathbf{C} and \mathbf{c} are constant, is

Observation update

- Current belief $P(\mathbf{x}_t | \mathbf{y}_{0:t-1}) = \mathcal{N}(\hat{\mathbf{x}}_{t|t-1}, \mathbf{Q}_{t|t-1})$
- Observation model $\mathbf{y}_t = \mathbf{H}\mathbf{x}_t + \mathbf{v}_t$ where $\mathbf{v}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{R})$
- Construct the joint on \mathbf{x}_t and \mathbf{y}_t

$$\boldsymbol{\mu} = \begin{pmatrix} \hat{\mathbf{x}}_{t|t-1} \\ \mathbf{H}\hat{\mathbf{x}}_{t|t-1} \end{pmatrix} \quad \boldsymbol{\Sigma} = \begin{pmatrix} \mathbf{Q}_{t|t-1} & \mathbf{Q}_{t|t-1}\mathbf{H}^T \\ \mathbf{H}\mathbf{Q}_{t|t-1} & \mathbf{H}\mathbf{Q}_{t|t-1}\mathbf{H}^T + \mathbf{R} \end{pmatrix}$$

- Condition on actual observation $\mathbf{y}_t = y_t$

$$P(\mathbf{x}_t | \mathbf{y}_{0:t}) = \mathcal{N}(\hat{\mathbf{x}}_{t|t}, \mathbf{Q}_{t|t})$$

$$\mathbf{Q}_{t|t} = \mathbf{Q}_{t|t-1} - \mathbf{Q}_{t|t-1}\mathbf{H}^T (\mathbf{H}\mathbf{Q}_{t|t-1}\mathbf{H}^T + \mathbf{R})^{-1} \mathbf{H}\mathbf{Q}_{t|t-1}$$

$$\hat{\mathbf{x}}_{t|t} = \hat{\mathbf{x}}_{t|t-1} + \mathbf{Q}_{t|t-1}\mathbf{H}^T (\mathbf{H}\mathbf{Q}_{t|t-1}\mathbf{H}^T + \mathbf{R})^{-1} (\mathbf{y}_t - \mathbf{H}\hat{\mathbf{x}}_{t|t-1})$$

Observation update: simplified

- Define Kalman gain $K_t = Q_{t|t-1}H^T (HQ_{t|t-1}H^T + R)^{-1}$.
- Use (tricky!) matrix algebra-fu to get useful relationships:

$$K_t = Q_{t|t}H^T R^{-1}$$

$$Q_{t|t} = Q_{t|t-1} - K_t H Q_{t|t-1}$$

$$\hat{x}_{t|t} = \hat{x}_{t|t-1} + K_t (y_t - H\hat{x}_{t|t-1})$$

- Call $y_t - H\hat{x}_{t|t-1}$ the innovation: how surprising is our observation?
- K_t maps y_t into an opinion about x_t : Big if observations are accurate and prior on x_t is weak.
- Intuition-building rewrite:

$$\hat{x}_{t|t} = (I - K_t H) \hat{x}_{t|t-1} + K_t y_t$$

Some important properties of the Kalman filter:

- Transition adds uncertainty: $Q_{t|t-1}$ is always “larger” than $Q_{t-1|t-1}$
- Observation reduces uncertainty: $Q_{t|t}$ is always “smaller” than $Q_{t|t-1}$

Kalman smoothing

Just as in discrete HMMs, we can run a similar belief-propagation pass backward to compute $P(\mathbf{x}_t \mid \mathbf{y}_{0:T})$

In Gaussian systems, the max of the individual marginals is the max of the joint!!!

What if your system isn't conjugate?

- Gaussian errors, but non-linear dynamics: extended Kalman filter
- Somewhat non-Gaussian errors, non-linear dynamics: unscented Kalman filter
- Arbitrary model: particle filter

Extended Kalman filter: optional

- Assume system with non-linearity limited to \mathbf{f} and \mathbf{h}

$$\mathbf{x}_t = \mathbf{f}(\mathbf{x}_{t-1}) + \mathbf{w} \quad \mathbf{w}_t \sim \mathcal{N}(\mathbf{0}, W)$$

$$\mathbf{y}_t = \mathbf{h}(\mathbf{x}_t) + \mathbf{v}_t \quad \mathbf{v}_t \sim \mathcal{N}(\mathbf{0}, R)$$

- Taylor series expansion about the current state estimate $\hat{\mathbf{x}}$:

$$\mathbf{f}(\mathbf{x}_{t-1}) = \mathbf{f}(\hat{\mathbf{x}}_{t-1|t-1}) + \left. \frac{\partial \mathbf{f}}{\partial \mathbf{x}} \right|_{\hat{\mathbf{x}}_{t-1|t-1}} (\mathbf{x}_t - \hat{\mathbf{x}}_{t-1|t-1}) + \dots$$

$$\mathbf{h}(\mathbf{x}_t) = \mathbf{h}(\hat{\mathbf{x}}_{t|t-1}) + \left. \frac{\partial \mathbf{h}}{\partial \mathbf{x}} \right|_{\hat{\mathbf{x}}_t} (\mathbf{x}_t - \hat{\mathbf{x}}_{t|t-1}) + \dots$$

assumes that all partial derivatives exist.

$$A(\hat{\mathbf{x}}_{t-1|t-1}) = \left. \frac{\partial \mathbf{f}}{\partial \mathbf{x}} \right|_{\hat{\mathbf{x}}_{t-1|t-1}} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \dots \\ & & \ddots \end{bmatrix}_{\mathbf{x}_t}, \quad H(\hat{\mathbf{x}}_{t|t-1}) = \left. \frac{\partial \mathbf{h}}{\partial \mathbf{x}} \right|_{\hat{\mathbf{x}}_{t|t-1}}$$

- Note that A and H are now time-varying!

Illustrative Example of EKF Weaknesses

- Shows possible reasons for EKF issues (Julier and Uhlmann [1994])
- Consider motion of vehicle following arc $\mathbf{x} = [x(t)y(t)\psi(t)]^T$, velocity $V(t)$ and radius of curvature $Ra(t)$
 - Velocity disturbed by a zero-mean uncorrelated process.
 - Covariance ellipse is oriented in direction of travel (Fig. 1).
- Later, 1/4 way around circle, Fig. 2 shows true position and uncertainty ellipse at time $t + 1$ — covariance ellipse has been expanded and rotated.

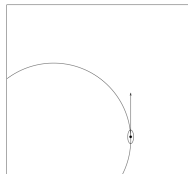


Fig. 1: Mean and covariance of a vehicle at time $t = m$

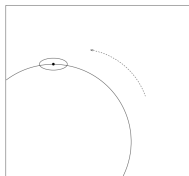


Fig. 2: True mean and covariance prediction to time $t = m + \Delta t$

Illustrative Example of EKF Issues

- EKF predicts covariance using linearized A matrix — equivalent to a constant velocity model tangent to the circle at time t .
- Effect shown in Fig. 3 — mean predicted around arc, but covariance ellipse predicted linearly in direction of travel \Rightarrow EKF loses critical information that largest component of uncertainty is in direction of travel.
- Error can be corrected by adding process noise to system (see Fig. 4), but at the expense of performance.

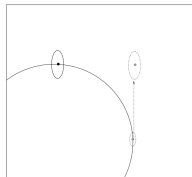


Fig. 3: EKF prediction of mean with linear covariance propagation

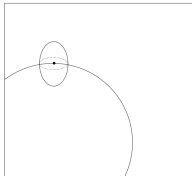


Fig. 4: EKF prediction "adjusted" to compensate for linearisation error

Unscented Kalman filter: optional

- Example of a more general idea: assumed density filtering
- Even if your posterior doesn't have the same parametric form as your prior, find the distribution in the family of the prior that is in some sense closest to the actual posterior
- This is an example of variational inference.
- One strategy is moment matching: estimate mean and covariance of posterior, and pretend it's a Gaussian with those moments. (Can match more moments if that's helpful.)
- UKF is a very simple and even more approximate version of this idea.
- It doesn't stink!

Filtering by sampling

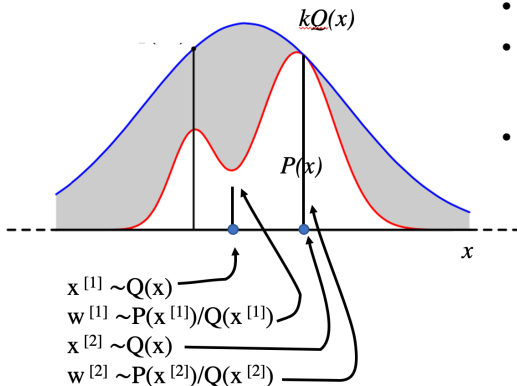
Importance Sampling (AIMA 13.4, Barber 27.6)

Keep all the samples, but weight them! Works in continuous space, too!

- Weight the samples by how much the proposal and target match for the given sample
- When we get a sample in a region where the proposal and target match: weight the sample highly
- When we get a sample in a region where the proposal and target don't match much: give the sample little weight
- Introduce weights

$$w(x) = \frac{P(x)}{Q(x)}$$

Importance Sampling



- If we have unnormalized proposal and target \tilde{p} and \tilde{q} , then we have weights \tilde{w} .
- It's easy to prove that

$$w(x) = \frac{\tilde{w}(x)}{\sum_x \tilde{w}(x)}$$

- Given importance-weighted samples, we can compute the statistics as before, but with weights:

$$E[x] = \sum w(x^{[i]})x^{[i]}$$

$$E[(x - \mu)^2] = \sum w(x^{[i]})(x^{[i]} - \bar{x})^2$$

Sampling Importance Resampling (SIR) (Barber 27.6)

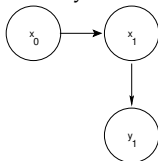
- Imagine that we wanted to infer the distribution $P(\mathbf{x}_{0:T}|\mathbf{y}_{0:T})$
Bad idea:
 - Sample from some multivariate Gaussian proposal distribution $Q(\mathbf{x}_{0:T}) = N(\mathbf{0}, \Sigma)$ where Σ is T dimensional.
 - Weight each sample according to the product of $P(\mathbf{x}_t|\mathbf{x}_{t-1})$ and $P(\mathbf{y}_t|\mathbf{x}_t)$
 - Recover the statistics from the weighted samples.
- If the distance between q and p is large (e.g., Kullback-Leibler divergence), our weights can grow very small
- We may have a lot of samples, but few of them are particularly useful
- Idea: importance **re**-sampling:
 1. Sample from $Q(\mathbf{x})$
 2. Compute weights $w(\mathbf{x}) = P(\mathbf{x})/Q(\mathbf{x})$
 3. Sample from the discrete distribution of sampled $\mathbf{x} \sim w(\mathbf{x})$.

Particle filter

- Use samples as pseudo representation of a distribution (“non-parametric”)
- Constant time per update step
- Weight of samples drops exponentially over time (because they are being generated without dependence on the observations)
- Instead, throw away samples with low weight and generate new ones as we go.

Filtering using Sampling: Sequential Importance Resampling

- Recall during filtering, the distribution we want is $P(\mathbf{x}_t | \mathbf{y}_{0:t})$ for each t
- If $T = 1$, then we have a 3-node Bayes net:



- We can get samples over the two latent variables $\mathbf{x}_{0:1}$, assuming we have a prior over \mathbf{x}_0 .
 1. Sample a value of \mathbf{x}_0 according to its prior
 2. Sample a value of \mathbf{x}_1 according to the sampled value of \mathbf{x}_0 and the noisy dynamics model
 3. Store combined sampled values \mathbf{x}_0 and \mathbf{x}_1 as a single “particle”
 4. Repeat until happy
- Two problems: what do we do about $P(\mathbf{y}_1 | \mathbf{x}_1)$, and how do we marginalize out \mathbf{x}_0 ?

Particle Filtering (Barber 27.6)

- Sampling from $P(\mathbf{y}_1|\mathbf{x}_1)$ doesn't help — we have \mathbf{y}_1 .
- We need to sample from $P(\mathbf{x}_1, \mathbf{x}_0|\mathbf{y}_1)$
- Bayes rule to the rescue!

This is our target:

$$P(\mathbf{x}_0, \mathbf{x}_1|\mathbf{y}_1) = \alpha P(\mathbf{y}_1|\mathbf{x}_1)P(\mathbf{x}_1, \mathbf{x}_0) \quad (\text{Where did } \mathbf{x}_0 \text{ go in the first term?})$$

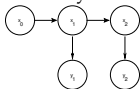
What if this was our proposal?

$$\begin{aligned} Q(\mathbf{x}_0, \mathbf{x}_1) &= P(\mathbf{x}_0, \mathbf{x}_1) \\ \frac{P(\mathbf{x}_0, \mathbf{x}_1|\mathbf{y}_1)}{Q(\mathbf{x}_0, \mathbf{x}_1)} &= \alpha \frac{P(\mathbf{y}_1|\mathbf{x}_1)P(\mathbf{x}_0, \mathbf{x}_1)}{P(\mathbf{x}_0, \mathbf{x}_1)} \\ \Rightarrow w(\mathbf{x}_0, \mathbf{x}_1) &= \alpha P(\mathbf{y}_1|\mathbf{x}_1) \end{aligned}$$

- What do we do about \mathbf{x}_0 ? How do we marginalize it out?
 - We can marginalize out \mathbf{x}_0 by resampling $p(\mathbf{x}_0, \mathbf{x}_1)$ according to the weights, and then dropping \mathbf{x}_0 .

Particle Filtering

- If we move to $k = 2$, we have a 5 node Bayes net



- We could just run the whole process from \mathbf{x}_0 , but recall that $\mathbf{x}_T \perp \mathbf{x}_{0:T-2}, \mathbf{y}_{0:T-1} | \mathbf{x}_{T-1}$.
- This independence means that if we have a distribution over $\mathbf{x}_{T-1} | \mathbf{y}_{0:T-1}$, we can discard the history $\mathbf{x}_{0:T-2}, \mathbf{y}_{0:T-1}$ from the particle.
- Therefore, for each new time step, we run the algorithm one step from \mathbf{x}_{T-1} to get samples over \mathbf{x}_T , weight by the new observation \mathbf{y}_T and resample.
- This gives us the following algorithm:

PARTICLE-FILTER($\mathbf{x}_{T-1}^{[i]}, \mathbf{y}_T$)

for i from 1 to n

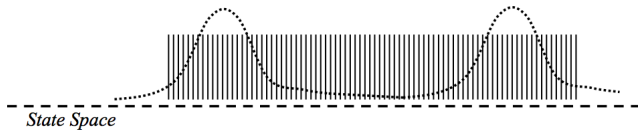
 Sample $\mathbf{x}_T^{[i]} \sim P(\mathbf{x}_T | \mathbf{x}_{T-1}^{[i]})$

 Compute $w^{[i]} = P(\mathbf{y}_T | \mathbf{x}_{T-1}^{[i]}, \mathbf{x}_T^{[i]})$

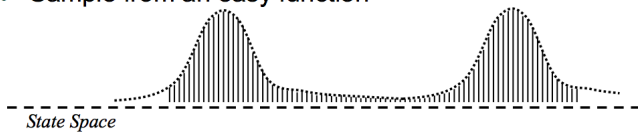
Generate $\mathbf{x}_T^{[1 \dots n]}$ samples from $\{\mathbf{x}_{T-1}^{[i]}, \mathbf{x}_T^{[i]}, w^{[i]}\} \sim w$

return $\mathbf{x}_T^{[1 \dots n]}$

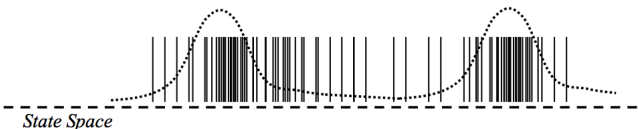
Particle filtering



- Sample from an easy function



- Compute importance weights



- Resample particles from particle set, according to importance weights

Particle filtering – some implementation issues

- You may not want to resample from the weighted particles on every step
 - Resampling may cause important but currently low-probability particles to be lost.
 - One option is to resample only after a certain amount of time when some of the particle weights are consistently very low.
- Need to be careful only to incorporate observations when the dynamics make the observations independent
 - Do not let the particle filter incorporate observations from the same location.
 - This will lead to convergence to a point estimate

Particle filtering – some implementation issues

- Another problem can be accidental particle death, when all the particles are too similar and have very low weight.
- If the measurement likelihood is strongly peaked, only a few particles may have a likely importance weight — these particles will get resampled often, leading to a pool of samples with low diversity \Rightarrow coarsely sampled approximation to the posterior
- To create some particle diversity, after resampling step, may want to add an additional perturbation by sampling a small noise term to be added to the samples.
- May also want to mix in particles from a distribution other than the prior
 - Sample from uniform over the state space : akin to fictitious noise in the Kalman filter: prevents the estimator from becoming too sure due to unmodelled approximations
 - Sample some particles from measurement model, compute importance weights from the dynamics. “Hybrid MCL”

Particle Filtering – Complexity

- There are few useful bounds for sampling techniques, and there is no formal bound on the number of particles required to get good performance
- When there are many local minima in the posterior distribution, need to make sure that you have enough particles
 - Hard to do in general
 - Can use KLD sampling to determine online when more samples are needed [Fox, 2003]
- A “simulated annealing” approach can be used, where the measurements are assumed much noisier than in truth, and the assumed noise is gradually reduced as the distribution converges to a consistent estimate
- Often used for global estimation of the position with no prior knowledge, before “tracking” can begin

So many other important ideas!

- Rao-Blackwellization: particles over some variables, and continuous distributions inside the particles over others.
- Dynamic Bayesian networks: factor states within a time step, and express transitions as a more general Bayes net.