

CS 287: Advanced Robotics Fall 2009

Lecture 21:
HMMs, Bayes filter, smoother, Kalman filters

Pieter Abbeel
UC Berkeley EECS

Overview

- Thus far:
 - Optimal control and reinforcement learning
 - We always assumed we got to observe the state at each time and the challenge was to choose a good action
 - Current and next set of lectures
 - The state is not observed*
 - Instead, we get some sensory information about the state
- Challenge: compute a probability distribution over the state which accounts for the sensory information ("evidence") which we have observed.

Examples

- Helicopter**
 - A choice of state: position, orientation, velocity, angular rate
 - Sensors:
 - GPS : noisy estimate of position (sometimes also velocity)
 - Inertial sensing unit: noisy measurements from (i) 3-axis gyro [=angular rate sensor], (ii) 3-axis accelerometer [=measures acceleration + gravity; e.g., measures (0,0,0) in free-fall], (iii) 3-axis magnetometer
- Mobile robot inside building**
 - A choice of state: position and heading
 - Sensors:
 - Odometry (=sensing motion of actuators): e.g., wheel encoders
 - Laser range finder: measures time of flight of a laser beam between departure and return (return is typically happening when hitting a surface that reflects the beam back to where it came from)

Probability review

For any random variables X, Y we have:

- Definition of conditional probability:**

$$P(X=x | Y=y) = P(X=x, Y=y) / P(Y=y)$$
- Chain rule:** (follows directly from the above)

$$P(X=x, Y=y) = P(X=x) P(Y=y | X=x) = P(Y=y) P(X=x | Y=y)$$
- Bayes rule:** (really just a re-ordering of terms in the above)

$$P(X=x | Y=y) = P(Y=y | X=x) P(X=x) / P(Y=y)$$
- Marginalization:**

$$P(X=x) = \sum_y P(X=x, Y=y)$$

Note: no assumptions beyond X, Y being random variables are made for any of these to hold true (and when we divide by something, that something is not zero)

Probability review

For any random variables X, Y, Z, W we have:

- Conditional probability:** (can condition on a third variable z throughout)

$$P(X=x | Y=y, Z=z) = P(X=x, Y=y | Z=z) / P(Y=y | Z=z)$$
- Chain rule:**

$$P(X=x, Y=y, Z=z, W=w) = P(X=x) P(Y=y | X=x) P(Z=z | X=x, Y=y) P(W=w | X=x, Y=y, Z=z)$$
- Bayes rule:** (can condition on other variable z throughout)

$$P(X=x | Y=y, Z=z) = P(Y=y | X=x, Z=z) P(X=x | Z=z) / P(Y=y | Z=z)$$
- Marginalization:**

$$P(X=x | W=w) = \sum_{y,z} P(X=x, Y=y, Z=z | W=w)$$

Note: no assumptions beyond X, Y, Z, W being random variables are made for any of these to hold true (and when we divide by something, that something is not zero)

Independence

- Two random variables X and Y are independent iff

$$\text{for all } x, y : P(X=x, Y=y) = P(X=x) P(Y=y)$$
- Representing a probability distribution over a set of random variables X_1, X_2, \dots, X_T in its most general form can be expensive.
 - E.g., if all X_i are binary valued, then there would be a total of 2^T possible instantiations and it would require 2^T-1 numbers to represent the probability distribution.
- However, if we assumed the random variables were independent, then we could very compactly represent the joint distribution as follows:
 - $P(X_1=x_1, X_2=x_2, \dots, X_T=x_T) = P(X_1=x_1) P(X_2=x_2) \dots P(X_T=x_T)$
 - Thanks to the independence assumptions, for the binary case, we went from requiring 2^T-1 parameters, to only requiring T parameters!
- Unfortunately independence is often too strong an assumption ...

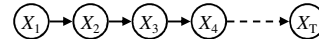
Conditional independence

- Two random variables X and Y are conditionally independent given a third random variable Z iff

$$\text{for all } x, y, z : P(X=x, Y=y | Z=z) = P(X=x | Z=z) P(Y=y | Z=z)$$
- Chain rule (which holds true for all distributions, no assumptions needed):
 - $P(X=x, Y=y, Z=z, W=w) = P(X=x)P(Y=y|X=x)P(Z=z|X=x, Y=y)P(W=w|X=x, Y=y, Z=z)$
 - For binary variables the representation requires $1 + 2^1 + 4^1 + 8^1 = 2^4 - 1$ numbers (just like a full joint probability table)
- Now assume Z independent of X given Y , and assume W independent of X and Y given Z , then we obtain:
 - $P(X=x, Y=y, Z=z, W=w) = P(X=x)P(Y=y|X=x)P(Z=z|Y=y)P(W=w|Z=z)$
 - For binary variables the representation requires $1 + 2^1 + 2^1 + 2^1 = 1 + (4-1) \cdot 2$ numbers --- significantly less!!

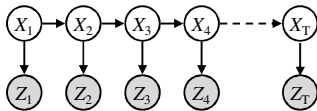
Markov Models

- Models a distribution over a set of random variables X_1, X_2, \dots, X_T where the index is typically associated with some notion of time.
- Markov models make the assumption:
 - X_t is independent of X_1, \dots, X_{t-2} when given X_{t-1}
- Chain rule: (always holds true, not just in Markov models!)
 - $P(X_1 = x_1, X_2 = x_2, \dots, X_T = x_T) = \prod_i P(X_i = x_i | X_{i-1} = x_{i-1}, X_{i-2} = x_{i-2}, \dots, X_1 = x_1)$
- Now apply the Markov conditional independence assumption:
 - $P(X_1 = x_1, X_2 = x_2, \dots, X_T = x_T) = \prod_i P(X_i = x_i | X_{i-1} = x_{i-1})$ (1)
 - in binary case: $1 + 2^*(T-1)$ numbers required to represent the joint distribution over all variables (vs. $2^T - 1$)
- Graphical representation: a variable X_t receives an arrow from the variables appearing in its conditional probability in the expression for the joint distribution (1) [called a Bayesian network or Bayes net representation]



Hidden Markov Models

- Underlying Markov model over states X_t
 - Assumption 1: X_t independent of X_1, \dots, X_{t-2} given X_{t-1}
- For each state X_t there is a random variable Z_t which is a sensory measurement of X_t
 - Assumption 2: Z_t is assumed conditionally independent of the other variables given X_t
- This gives the following graphical (Bayes net) representation:



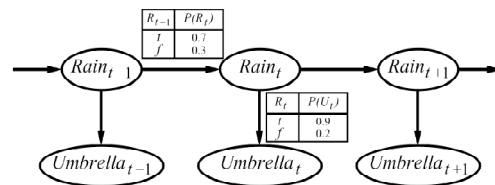
Hidden Markov Models

-
- $P(X_1=x_1, Z_1=z_1, X_2=x_2, Z_2=z_2, \dots, X_T=x_T, Z_T=z_T) =$
- Chain rule: (no assumptions)
 - $P(X_1 = x_1)$
 - $P(Z_1 = z_1 | X_1 = x_1)$
 - $P(X_2 = x_2 | X_1 = x_1, Z_1 = z_1)$
 - $P(Z_2 = z_2 | X_1 = x_1, Z_1 = z_1, X_2 = x_2)$
 - ...
 - $P(X_T = x_T | X_1 = x_1, Z_1 = z_1, \dots, X_{T-1} = x_{T-1}, Z_{T-1} = z_{T-1})$
 - $P(Z_T = z_T | X_1 = x_1, Z_1 = z_1, \dots, X_{T-1} = x_{T-1}, Z_{T-1} = z_{T-1}, X_T = x_T)$
 - HMM assumptions:
 - $P(X_1 = x_1)$
 - $P(Z_t = z_t | X_t = x_t)$
 - $P(X_t = x_t | X_{t-1} = x_{t-1})$
 - $P(Z_t = z_t | X_t = x_t)$
 - ...
 - $P(X_T = x_T | X_{T-1} = x_{T-1})$
 - $P(Z_T = z_T | X_T = x_T)$

Mini quiz

- What would the graph look like for a Bayesian network with no conditional independence assumptions?
- Our particular choice of ordering of variables in the chain rule enabled us to easily incorporate the HMM assumptions. What if we had chosen a different ordering in the chain rule expansion?

Example



- The HMM is defined by:
 - Initial distribution: $P(X_1)$
 - Transitions: $P(X_t | X_{t-1})$
 - Observations: $P(Z_t | X_t)$

Real HMM Examples

- Robot localization:
 - Observations are range readings (continuous)
 - States are positions on a map (continuous)

- Speech recognition HMMs:
 - Observations are acoustic signals (continuous valued)
 - States are specific positions in specific words (so, tens of thousands)

- Machine translation HMMs:
 - Observations are words (tens of thousands)
 - States are translation options

Filtering / Monitoring

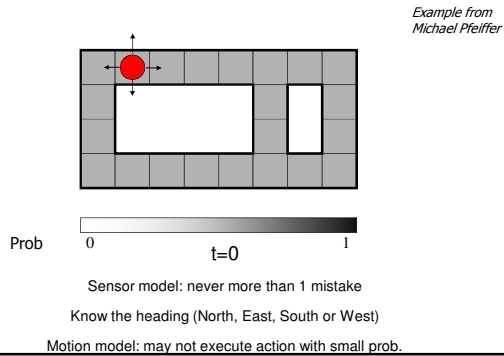
- Filtering, or monitoring, is the task of tracking the distribution $P(X_t | Z_1 = z_1, Z_2 = z_2, \dots, Z_t = z_t)$ over time. This distribution is called the belief state.

- We start with $P(X_0)$ in an initial setting, usually uniform

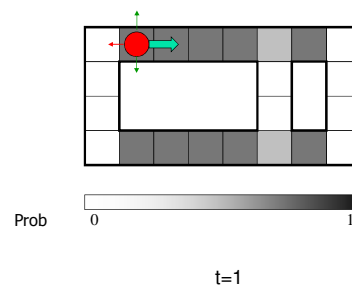
- As time passes, or we get observations, we update the belief state.

- The Kalman filter was invented in the 60's and first implemented as a method of trajectory estimation for the Apollo program. [See course website for a historical account on the Kalman filter. "From Gauss to Kalman"]

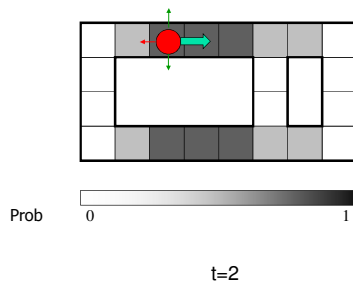
Example: Robot Localization



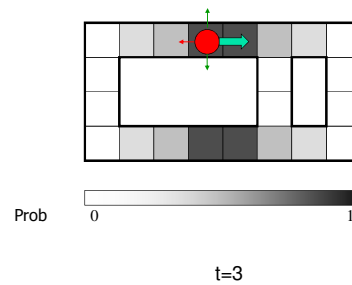
Example: Robot Localization



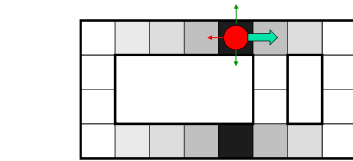
Example: Robot Localization



Example: Robot Localization



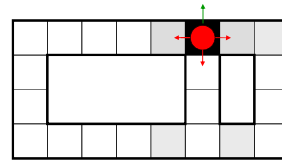
Example: Robot Localization



Prob 0 1

t=4

Example: Robot Localization



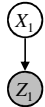
Prob 0 1

t=5

Inference: Base Cases

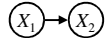
Incorporate observation

Time update



$$P(X_1|z_1)$$

$$P(x_1|z_1) = P(x_1, z_1) / P(z_1) \\ \propto P(x_1, z_1) \\ = P(x_1)P(z_1|x_1)$$



$$P(X_2)$$

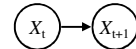
$$P(x_2) = \sum_{x_1} P(x_1, x_2) \\ = \sum_{x_1} P(x_1)P(x_2|x_1)$$

Time update

- Assume we have current belief $P(X | \text{evidence to date})$

$$P(x_t|e_{1:t})$$

- Then, after one time step passes:



$$P(x_{t+1}|e_{1:t}) = \sum_{x_t} P(x_{t+1}|x_t)P(x_t|e_{1:t})$$

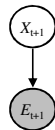
Observation update

- Assume we have:

$$P(x_{t+1}|e_{1:t})$$

- Then:

$$P(x_{t+1}|e_{1:t+1}) \propto P(e_{t+1}|x_{t+1})P(x_{t+1}|e_{1:t})$$



Algorithm

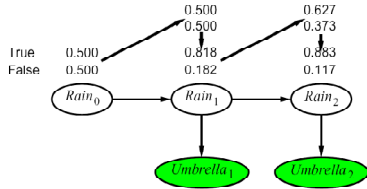
- Init $P(x_t)$ [e.g., uniformly]
- Observation update for time 0:

$$P(x_1|z_1) \propto P(z_1|x_1)P(x_1)$$
- For $t = 1, 2, \dots$
 - Time update

$$P(x_{t+1}|z_{1:t}) = \sum_{x_t} P(x_{t+1}|x_t)P(x_t|z_{1:t})$$
 - Observation update

$$P(x_{t+1}|z_{1:t+1}) \propto P(z_{t+1}|x_{t+1})P(x_{t+1}|z_{1:t})$$
- For continuous state / observation spaces: simply replace summation by integral

Example HMM



R_{t-1}	$P(R_t)$	R_t	$P(U_t)$
T	0.7	T	0.9
F	0.3	F	0.2

The Forward Algorithm

- Time/dynamics update and observation update in one:

$$\begin{aligned}
 P(x_t, z_{1:t}) &= \sum_{x_{t-1}} P(x_{t-1}, x_t, z_{1:t}) \\
 &= \sum_{x_{t-1}} P(x_{t-1}, z_{1:t-1}) P(x_t | x_{t-1}) P(z_t | x_t) \\
 &= P(z_t | x_t) \sum_{x_{t-1}} P(x_t | x_{t-1}) P(x_{t-1}, z_{1:t-1})
 \end{aligned}$$

- recursive update
- Normalization:
 - Can be helpful for numerical reasons
 - However: lose information!
- Can renormalize (for numerical reasons) + keep track of the normalization factor (to enable recovering all information)

The likelihood of the observations

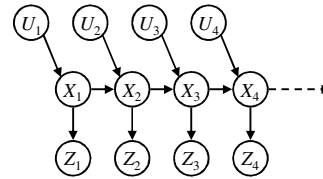
$$P(z_{1:t}) = \sum_{x_1, x_2, \dots, x_t} P(x_{1:t}, z_{1:t}) = \sum_{x_1, x_2, \dots, x_t} \prod_{k=1}^{t-1} P(x_{k+1} | x_k) P(z_k | x_k) P(z_t | x_t)$$

- The forward algorithm first sums over x_1 , then over x_2 and so forth, which allows it to efficiently compute the likelihood at all times t , indeed:

$$P(z_{1:t}) = \sum_{x_t} P(x_t, z_{1:t})$$

- Relevance:
 - Compare the fit of several HMM models to the data
 - Could optimize the dynamics model and observation model to maximize the likelihood
 - Run multiple simultaneous trackers --- retain the best and split again whenever applicable (e.g., loop closures in SLAM, or different flight maneuvers)

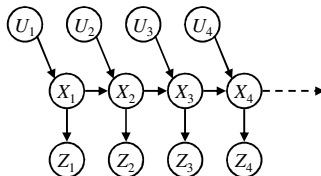
With control inputs



We know track:

$$P(x_t | z_{1:t}, u_{1:t})$$

With control inputs



- Control inputs known:
 - They can be simply seen as selecting a particular dynamics function
- Control inputs unknown:
 - Assume a distribution over them
- Above drawing assumes open-loop controls. This is rarely the case in practice. [Markov assumption is rarely the case either. Both assumptions seem to have given quite satisfactory results.]

Smoothing

- Thus far, filtering, which finds:
 - The distribution over states at time t given all evidence until time t :

$$P(x_t | z_{1:t}, u_{1:t})$$

- The likelihood of the evidence up to time t :

$$P(z_{1:t} | u_{1:t})$$

- How about?

$$P(x_t | z_{1:T}, u_{1:T})$$

- $T < t$: can simply run the forward algorithm until time t , but stop incorporating evidence from time $T+1$ onwards
- $T > t$: need something else

Smoothing

$$\begin{aligned}
 P(x_t | z_{1:T}) &\propto P(x_t, z_{1:T}) \\
 &= \sum_{x_1, x_2, \dots, x_{t-1}, x_{t+1}, \dots, x_T} P(x_{1:T}, z_{1:T}) \\
 &= \sum_{x_1, x_2, \dots, x_{t-1}, x_{t+1}, \dots, x_T} \prod_{k=1}^T P(x_k | x_{k-1}) P(z_k | x_k)
 \end{aligned}$$

- Sum as written has a number of terms exponential in T
- Key idea: order in which variables are being summed out affects computational complexity
 - Forward algorithm** exploits summing out x_1, x_2, \dots, x_{t-1} in this order
 - Can similarly run a **backward algorithm**, which sums out $x_T, x_{T-1}, \dots, x_{t+1}$ in this order

Smoothing

$$\begin{aligned}
 P(x_t, z_{1:T}) &= \sum_{x_1, x_2, \dots, x_{t-1}, x_{t+1}, \dots, x_T} P(x_{1:T}, z_{1:T}) \\
 &= \sum_{x_1, x_2, \dots, x_{t-1}, x_{t+1}, \dots, x_T} \prod_{k=1}^T P(x_k | x_{k-1}) P(z_k | x_k) \\
 &= \left(\sum_{x_1, x_2, \dots, x_{t-1}} \prod_{k=1}^{t-1} P(x_k | x_{k-1}) P(z_k | x_k) \right) \left(\sum_{x_{t+1}, x_{t+2}, \dots, x_T} \prod_{k=t+1}^T P(x_k | x_{k-1}) P(z_k | x_k) \right)
 \end{aligned}$$

Forward algorithm computes this
Backward algorithm computes this

- Can be easily verified from the equations:
 - The factors in the right parentheses only contain x_{t+1}, \dots, x_T , hence they act as a constant when summing out over x_1, \dots, x_{t-1} and can be brought outside the summation
- Can also be read off from the Bayes net graph / conditional independence assumptions:
 - x_1, \dots, x_{t-1} are conditionally of x_{t+1}, \dots, x_T given x_t

Backward algorithm

- Sum out x_t :

$$\begin{aligned}
 P(x_t, e_{1:T}) &= \sum_{x_1, x_2, \dots, x_{t-1}, x_{t+1}, \dots, x_T} \prod_{k=1}^T P(x_k | x_{k-1}) P(e_k | x_k) \\
 &= \sum_{x_1, x_2, \dots, x_{t-1}, x_{t+1}, \dots, x_T} \prod_{k=1}^T P(x_k | x_{k-1}) P(e_k | x_k) \\
 &= \sum_{x_1, x_2, \dots, x_{t-1}, x_{t+1}, \dots, x_T} \prod_{k=1}^{T-1} P(x_k | x_{k-1}) P(e_k | x_k) \sum_{x_T} P(x_T | x_{T-1}) P(e_T | x_T) \\
 &= \sum_{x_1, x_2, \dots, x_{t-1}, x_{t+1}, \dots, x_{T-1}} \prod_{k=1}^{T-1} P(x_k | x_{k-1}) P(e_k | x_k) f_{T-1}(x_{T-1})
 \end{aligned}$$

- Can recursively compute for $l=T, T-1, \dots$:

$$f_{l-1}(x_{l-1}) = \sum_{x_l} P(x_l | x_{l-1}) P(e_l | x_l) f_l(x_l)$$

Smoother algorithm

- Run forward algorithm, which gives
 - $P(x_t, z_1, \dots, z_t)$ for all t
- Run backward algorithm, which gives
 - $f_t(x_t)$ for all t
- Find
 - $P(x_t, z_1, \dots, z_T) = P(x_t, z_1, \dots, z_t) f_t(x_t)$
 - If desirable, can renormalize and find $P(x_t | z_1, \dots, z_T)$

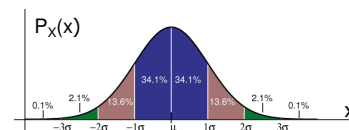
Bayes filters

- Recursively compute
 - $P(x_t, z_{1:t-1}) = \sum_{x_{t-1}} P(x_t | x_{t-1}) P(x_{t-1} | z_{1:t-1})$
 - $P(x_t, z_{1:t}) = P(x_t, z_{1:t-1}) P(z_t | x_t)$
- Tractable cases:
 - State space finite and sufficiently small (what we have in some sense considered so far)
 - Systems with linear dynamics and linear observations and Gaussian noise
 - Kalman filtering

Univariate Gaussian

- Gaussian distribution with mean μ , and standard deviation σ :

$$N(\mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$



Properties of Gaussians

$$N(\mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

- Mean:

$$EX = \int x \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) = \mu$$

- Variance:

$$E((X-\mu)^2) = \int (x-\mu)^2 \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) = \sigma^2$$

Central limit theorem (CLT)

- Classical CLT:

- Let X_1, X_2, \dots be an infinite sequence of *independent* random variables with $E X_i = \mu$, $E(X_i - \mu)^2 = \sigma^2$
- Define $Z_n = ((X_1 + \dots + X_n) - n\mu) / (\sigma n^{1/2})$
- Then for the limit of n going to infinity we have that Z_n is distributed according to $N(0, 1)$

- Crude statement: things that are the result of the addition of lots of small effects tend to become Gaussian.

Multi-variate Gaussians

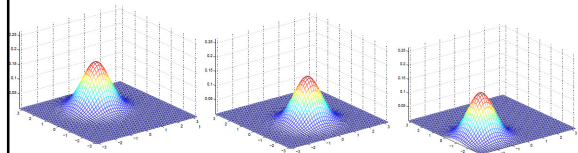
$$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x-\mu)^\top \Sigma^{-1}(x-\mu)\right)$$

$$EX = \int xp(x; \mu, \Sigma) = \mu$$

$$E[(X_i - \mu_i)(X_j - \mu_j)] = \int (x_i - \mu_i)(x_j - \mu_j) p(x; \mu, \Sigma) = \Sigma_{ij}$$

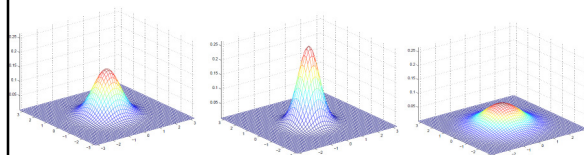
$$\int \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x-\mu)^\top \Sigma^{-1}(x-\mu)\right) = 1$$

Multi-variate Gaussians: examples



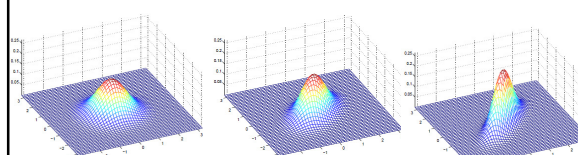
- $\mu = [1; 0]$
- $\Sigma = [1 \ 0; 0 \ 1]$
- $\mu = [-.5; 0]$
- $\Sigma = [1 \ 0; 0 \ 1]$
- $\mu = [-1; -1.5]$
- $\Sigma = [1 \ 0; 0 \ 1]$

Multi-variate Gaussians: examples

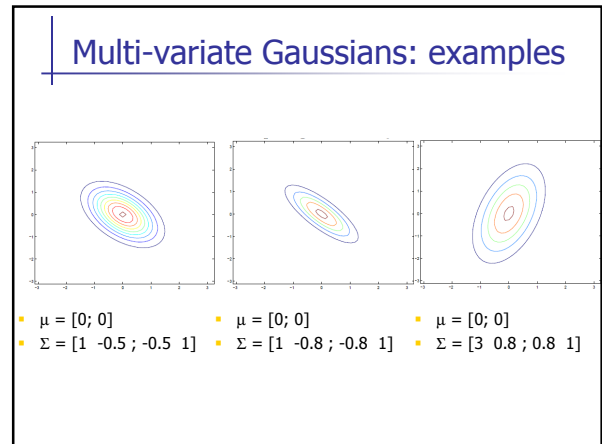
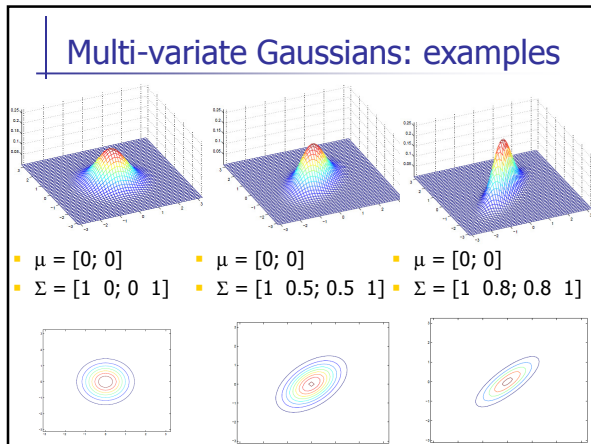


- $\mu = [0; 0]$
- $\Sigma = [1 \ 0; 0 \ 1]$
- $\mu = [0; 0]$
- $\Sigma = [.6 \ 0; 0 \ .6]$
- $\mu = [0; 0]$
- $\Sigma = [2 \ 0; 0 \ 2]$

Multi-variate Gaussians: examples



- $\mu = [0; 0]$
- $\Sigma = [1 \ 0; 0 \ 1]$
- $\mu = [0; 0]$
- $\Sigma = [1 \ 0.5; 0.5 \ 1]$
- $\mu = [0; 0]$
- $\Sigma = [1 \ 0.8; 0.8 \ 1]$



Discrete Kalman Filter

Estimates the state x of a discrete-time controlled process that is governed by the linear stochastic difference equation

$$x_t = A_t x_{t-1} + B_t u_t + \varepsilon_t$$

with a measurement

$$z_t = C_t x_t + \delta_t$$

47

Components of a Kalman Filter

- A_t Matrix ($n \times n$) that describes how the state evolves from t to $t-1$ without controls or noise.
- B_t Matrix ($n \times l$) that describes how the control u_t changes the state from t to $t-1$.
- C_t Matrix ($k \times n$) that describes how to map the state x_t to an observation z_t .
- ε_t Random variables representing the process and measurement noise that are assumed to be independent and normally distributed with covariance R_t and Q_t respectively.

48

Linear Gaussian Systems: Initialization

- Initial belief is normally distributed:

$$bel(x_0) = N(x_0; \mu_0, \Sigma_0)$$

53

Linear Gaussian Systems: Dynamics

- Dynamics are linear function of state and control plus additive noise:

$$x_t = A_t x_{t-1} + B_t u_t + \varepsilon_t$$

$$p(x_t | u_t, x_{t-1}) = N(x_t; A_t x_{t-1} + B_t u_t, R_t)$$

$$\overline{bel}(x_t) = \int p(x_t | u_t, x_{t-1}) \overline{bel}(x_{t-1}) dx_{t-1}$$

$$\downarrow \qquad \qquad \qquad \downarrow$$

$$\sim N(x_t; A_t x_{t-1} + B_t u_t, R_t) \quad \sim N(x_{t-1}; \mu_{t-1}, \Sigma_{t-1})$$

54

Linear Gaussian Systems: Dynamics

$$\begin{aligned} \overline{bel}(x_t) &= \int p(x_t | u_t, x_{t-1}) \quad \overline{bel}(x_{t-1}) \, dx_{t-1} \\ &\quad \downarrow \qquad \qquad \qquad \downarrow \\ &\sim N(x_t; A_t x_{t-1} + B_t u_t, R_t) \quad \sim N(x_{t-1}; \mu_{t-1}, \Sigma_{t-1}) \\ &\quad \downarrow \\ \overline{bel}(x_t) &= \eta \int \exp\left\{-\frac{1}{2}(x_t - A_t x_{t-1} - B_t u_t)^T R_t^{-1} (x_t - A_t x_{t-1} - B_t u_t)\right\} \\ &\quad \exp\left\{-\frac{1}{2}(x_{t-1} - \mu_{t-1})^T \Sigma_{t-1}^{-1} (x_{t-1} - \mu_{t-1})\right\} dx_{t-1} \\ \overline{bel}(x_t) &= \begin{cases} \overline{\mu}_t = A_t \mu_{t-1} + B_t u_t \\ \overline{\Sigma}_t = A_t \Sigma_{t-1} A_t^T + R_t \end{cases} \end{aligned}$$

55

Proof: completion of squares

- To integrate out x_{t-1} , re-write the integrand in the format:

$$\int f(\mu_{t-1}, \Sigma_{t-1}, x_t) \left(\frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(x_{t-1} - m)^T S^{-1} (x_{t-1} - m)\right) \right) dx_{t-1}$$

- This integral is readily computed (integral of a multivariate Gaussian times a constant = that constant) to be

$$f(\mu_{t-1}, \Sigma_{t-1}, x_t)$$

- Inspection of f will show that it is a multi-variate Gaussian in x_t with the mean and covariance as shown on previous slide.

Properties of Gaussians

- We just showed:

$$\left. \begin{matrix} X \sim N(\mu, \Sigma) \\ Y = AX + B \end{matrix} \right\} \Rightarrow Y \sim N(A\mu + B, A\Sigma A^T)$$

- We stay in the "Gaussian world" as long as we start with Gaussians and perform only linear transformations.
- Now we know this, we could find μ_Y and Σ_Y without computing integrals by directly computing the expected values:

$$E[Y] = E[AX + B] = AE[X] + B = A\mu + B$$

$$\begin{aligned} \Sigma_{YY} &= E[(Y - E[Y])(Y - E[Y])^T] = E[(AX + B - A\mu - B)(AX + B - A\mu - B)^T] \\ &= E[A(X - \mu)(X - \mu)^T A^T] = AE[(X - \mu)(X - \mu)^T]A^T = A\Sigma A^T \end{aligned}$$

Self-quiz

Test your understanding of the completion of squares trick! Let $A \in \mathbf{R}^{n \times n}$ be a positive definite matrix, $b \in \mathbf{R}^n$, and $c \in \mathbf{R}$. Prove that

$$\int_{x \in \mathbf{R}^n} \exp\left(-\frac{1}{2}x^T A x - x^T b - c\right) dx = \frac{(2\pi)^{n/2}}{|A|^{1/2} \exp\left(c - \frac{1}{2}b^T A^{-1} b\right)}$$

Linear Gaussian Systems: Observations

- Observations are linear function of state plus additive noise:

$$z_t = C_t x_t + \delta_t$$

$$p(z_t | x_t) = N(z_t; C_t x_t, Q_t)$$

$$\begin{aligned} bel(x_t) &= \eta \quad p(z_t | x_t) \quad \overline{bel}(x_t) \\ &\quad \downarrow \qquad \qquad \qquad \downarrow \\ &\sim N(z_t; C_t x_t, Q_t) \quad \sim N(x_t; \overline{\mu}_t, \overline{\Sigma}_t) \end{aligned}$$

59

Linear Gaussian Systems: Observations

$$\begin{aligned} bel(x_t) &= \eta \quad p(z_t | x_t) \quad \overline{bel}(x_t) \\ &\quad \downarrow \qquad \qquad \qquad \downarrow \\ &\sim N(z_t; C_t x_t, Q_t) \quad \sim N(x_t; \overline{\mu}_t, \overline{\Sigma}_t) \\ &\quad \downarrow \\ bel(x_t) &= \eta \exp\left\{-\frac{1}{2}(z_t - C_t x_t)^T Q_t^{-1} (z_t - C_t x_t)\right\} \exp\left\{-\frac{1}{2}(x_t - \overline{\mu}_t)^T \overline{\Sigma}_t^{-1} (x_t - \overline{\mu}_t)\right\} \\ bel(x_t) &= \begin{cases} \mu_t = \overline{\mu}_t + K_t (z_t - C_t \overline{\mu}_t) \\ \Sigma_t = (I - K_t C_t) \overline{\Sigma}_t \end{cases} \quad \text{with } K_t = \overline{\Sigma}_t C_t^T (C_t \overline{\Sigma}_t C_t^T + Q_t)^{-1} \end{aligned}$$

60

Proof: completion of squares

- Re-write the expression for $bel(x_t)$ in the format:

$$f(\bar{\mu}_t, \bar{\Sigma}_t, C_t, Q_t) \left(\frac{1}{(2\pi)^{\frac{n}{2}} |\bar{S}|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (x_t - m) S^{-1} (x_t - m) \right) \right)$$

- f is the normalization factor
- The expression in parentheses is a multi-variate Gaussian in x_t . Its parameters m and S can be identified to satisfy the expressions for the mean and covariance on the previous slide.

Kalman Filter Algorithm

Algorithm `Kalman_filter` ($\mu_{t-1}, \Sigma_{t-1}, u_t, z_t$):

Prediction:

$$\begin{aligned} \bar{\mu}_t &= A_t \mu_{t-1} + B_t u_t \\ \bar{\Sigma}_t &= A_t \Sigma_{t-1} A_t^T + R_t \end{aligned}$$

Correction:

$$\begin{aligned} K_t &= \bar{\Sigma}_t C_t^T (C_t \bar{\Sigma}_t C_t^T + Q_t)^{-1} \\ \mu_t &= \bar{\mu}_t + K_t (z_t - C_t \bar{\mu}_t) \\ \Sigma_t &= (I - K_t C_t) \bar{\Sigma}_t \end{aligned}$$

Return μ_t, Σ_t

62

How to derive these updates

- Simply work through the integrals
 - Key "trick": completion of squares
- If your derivation results in a different format \rightarrow apply matrix inversion lemma to prove equivalence

$$(A + UCV)^{-1} = A^{-1} - A^{-1}U(C^{-1} + VA^{-1}U)^{-1}VA^{-1}$$