

PRML:

- 4.4 on Laplace Approx (skip 4.4.1)
- 13 on models for sequential data
 - skim: 13.2.1, 13.2.3
 - skip: 13.2.4, 13.2.6, 13.3.2-4

4.4 The Laplace Approximation

Laplace approx: framework to find a Gaussian approx to a probability density over a set of continuous variables

- consider a single continuous var. z

$$p(z) = \frac{1}{Z} f(z) \quad \text{where } Z = \int f(z) dz \quad \text{"normalization coef."}$$

\hookrightarrow value is unknown

- goal: find a Gaussian approx. $q(z)$ which is centered on a mode of $p(z)$

- start: find a mode of $p(z)$; i.e. find a z_0 s.t. $p'(z_0) = 0$:

$$\left. \frac{df(z)}{dz} \right|_{z=z_0} = 0$$

- the log of a Gaussian is a quadratic fn of the variables so can use Taylor expansion:

$$\ln f(z) \approx \ln f(z_0) - \frac{1}{2} A (z - z_0)^2 \quad A = - \left. \frac{d^2}{dz^2} \ln f(z) \right|_{z=z_0}$$

$$f(z) \approx f(z_0) \exp \left[- \frac{A}{2} (z - z_0)^2 \right]$$

- normalized distribution $q(\mathbf{z})$ from the standard normalization of a Gaussian:

$$q(\mathbf{z}) = \left(\frac{A}{2\pi}\right)^{\frac{1}{2}} \exp\left[-\frac{A}{2}(\mathbf{z}-\mathbf{z}_0)^T\right]$$

- extend Laplace method over M -dimensional space \mathbf{z}

$$p(\mathbf{z}) = \frac{1}{Z} f(\mathbf{z})$$

- stationary point \mathbf{z}_0 where $\nabla f(\mathbf{z}) = 0$

$$\ln f(\mathbf{z}) \approx \ln f(\mathbf{z}_0) - \frac{1}{2}(\mathbf{z}-\mathbf{z}_0)^T \mathbf{A}(\mathbf{z}-\mathbf{z}_0)$$

where \mathbf{A} is a $M \times M$ Hessian matrix: $\mathbf{A} = -\nabla \nabla \ln f(\mathbf{z}) \big|_{\mathbf{z}=\mathbf{z}_0}$

↳ "square matrix of 2nd order partial derivatives of a scalar valued function"

$$f(\mathbf{z}) \approx f(\mathbf{z}_0) \exp\left[-\frac{1}{2}(\mathbf{z}-\mathbf{z}_0)^T \mathbf{A}(\mathbf{z}-\mathbf{z}_0)\right]$$

- $q(\mathbf{z})$ is proportional to $f(\mathbf{z})$
- can find normalization w/ standard result for a normalized multivariate Gaussian

$$q(\mathbf{z}) = \frac{|\mathbf{A}|^{1/2}}{(2\pi)^{M/2}} \exp\left[-\frac{1}{2}(\mathbf{z}-\mathbf{z}_0)^T \mathbf{A}(\mathbf{z}-\mathbf{z}_0)\right] \quad |\mathbf{A}| : \text{determinate of } \mathbf{A}$$

$$= N(\mathbf{z}|\mathbf{z}_0, \mathbf{A}^{-1})$$

13. Sequential Data

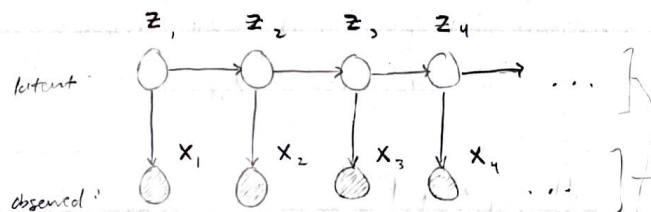
- Markov models: assume that future predictions are independent of all but the most recent observations
- we'll focus on two "state space models": hidden Markov model (HMM), linear dynamical system (LDS)

13.1 Markov models

- relax assumption of i.i.d data
- product rule to express joint distribution for a sequence of observations:

$$p(x_1, \dots, x_N) = \prod_{n=1}^N p(x_n | x_1, \dots, x_{n-1})$$
- assume x_n is independent of all but most recent obs x_{n-1} \rightarrow "first-order Markov chain"

$$p(x_1, \dots, x_N) = \prod_{n=2}^N p(x_n | x_{n-1})$$
- 2nd order MC: $p(x_1, \dots, x_N) = \prod_{n=3}^N p(x_n | x_{n-1}, x_{n-2})$
- autoregressive (AR) models: for continuous variables, use linear Gaussian conditional distributions where each node is a Gaussian dist. w/ mean as a linear function of the parents
- want to build a model w/o a specific Markov order but still limited num. of free parameters
 - need to introduce some latent variables
- for each observation x_n , introduce a corresponding latent var. z_n of lower dim
- form the MC w/ the latent variables to form a "state space model"



$$p(x_1, \dots, x_N, z_1, \dots, z_N) = p(z_1) \left[\prod_{n=2}^N p(z_n | z_{n-1}) \right] \left[\prod_{n=1}^N p(x_n | z_n) \right]$$

13.2 Hidden Markov Models (skim: 13.2.1 & 13.2.3; skip 13.2.4, 13.2.6)

- state space model of discrete latent variables
 - can be interpreted as an extension of a mixture model w/ choice of mixture component dependent on component of previous observation
- latent variables are discrete multinomial variables z_n w/ 1-of-K coding scheme
 - describes which component produced observation x_n
- prob. of z_n is dependent on z_{n-1} : $p(z_n | z_{n-1})$
 - this conditional dist. corresponds to table $A_{K \times K}$ containing "transition prob."

$$A_{jk} = p(z_{nk} = 1 | z_{n-1,j} = 1) \equiv \text{prob. of going from } j \rightarrow k \text{ component}$$

$$p(z_n | z_{n-1}, A) = \prod_{k=1}^K \prod_{j=1}^K A_{jk}^{z_{n-1,j} z_{nk}}$$

special case for first latent node z_1 :

$$p(z_1 | \pi) = \prod_{k=1}^K \pi_k^{z_{1k}} \quad \text{where } \pi \text{ is a vector of probs.: } \pi_k \equiv p(z_{1k} = 1)$$

$\sum_k \pi_k = 1$

- see Fig 13.6 (pg. 611) for a "state transition diagram"
- Fig 13.7 (pg. 612) shows the state transition diagram unfolded over time

- finish prob. models by defining conditional dist. of the observed variables x_n :
 $p(x_n | z_n, \phi)$ where ϕ are the parameters of the dist. ("emission prob.")
 - $p(x_n | z_n, \phi)$ consists of a vector of K numbers corresponding to the K possible states of a binary vector z_n

emission probabilities:

$$p(x_n | z_n, \phi) = \prod_{k=1}^K p(x_n | \phi_k)^{z_{nk}} \quad \leftarrow \text{(recall } z_{nk} \text{ is either 1 or 0)}$$

for a homogeneous model, all latent vars. use same A & all emission dist. use the same ϕ

joint prob. dist. over latent & observed variables:

$$p(X, Z | \theta) = p(z_1 | \pi) \left[\prod_{n=2}^N p(z_n | z_{n-1}, A) \right] \left[\prod_{n=1}^N p(x_n | z_n, \phi) \right]$$

$$X = \{x_1, \dots, x_N\}, \quad Z = \{z_1, \dots, z_N\}, \quad \theta = \{\pi, A, \phi\}$$

13.3 Linear Dynamical System (LDS)

continuous latent variables

the general form of the inference algorithms are the same as for HMM

we will consider a linear-Gaussian state space model

latent vars. $\{z_n\}$ & observed variables $\{x_n\}$ are multivariate Gauss.

means of these dist. are linear fans of their parents in the graph

can interpret LDS as a generalization of continuous latent variable models such as prob. PCA

each pair of nodes $\{z_n, x_n\}$ represents a linear-Gaussian latent var. model

for the observation, but now the latent vars are not independent

because LDS is a linear-Gaussian model, all joint, marginal, & conditional dist. will be Gaussian.

transition dist: $p(z_n | z_{n-1}) = N(z_n | A z_{n-1}, \Gamma)$

emission dist: $p(x_n | z_n) = N(x_n | C z_n, \Sigma)$

initial latent var. dist: $p(z_1) = N(z_1 | \mu_0, V_0)$

alternative formulations of above:

$$z_n = A z_{n-1} + w_n \quad w \sim N(w | 0, \Gamma)$$

$$x_n = C z_n + v_n \quad v \sim N(v | 0, \Sigma)$$

$$z_1 = \mu_0 + u \quad u \sim N(u | 0, V_0)$$

(noise terms)

13.3.1. Inference in LDS

- goal: find marginal dist. for the latent variables conditional on the data
- and make predictions on the next latent state & observation
- accomplish w/ the "sum-product algorithm"
- w.r.t. LDS, results in the "Kalman filter" and "Kalman smoother" equations
- b/c only dealing w/ Gaussians, could use standard results for conditionals & marginals, but "sum-product alg." is more efficient

begin by considering the forward equations

· z_n : root node $h(z_i)$: leaf node

· propagate msg from leaf to root

· propagate messages that are normalized marginal dist.:

$$p(z_n | x_1, \dots, x_n) \rightarrow \hat{\alpha}(z_n) = N(z_n | \mu_n, V_n)$$

recursion equation:

$$c_n \hat{\alpha}(z_n) = p(x_n | z_n) \int \hat{\alpha}(z_{n-1}) p(z_n | z_{n-1}) dz_{n-1}$$

Topics:

1. Laplace Approximation
2. State space models (HMM & linear dynamical systems)

1. Laplace Approx.

- normally w/ MAP estimates, do not get a value of uncertainty
- for Laplace approx., can use the variance as op. estimate of uncertainty
- Laplace approx. is best when the actual prob. dist. is Gaussian-ish
- symmetric w/ long-tails

2. State space models

A. General structure

- when modeling a Markov chain (MC):
 - assuming previous state is all that's needed to know the current state
 - requires being descriptive x
- state space models assume a latent state space z
- the measurements x are noisy observations of z
 - emission model: $p(x_n | z_n)$

2 types of tractable state space models:

1. HMM

- discrete latent space
- $p(z_n | z_{n-1})$ is a transition matrix

2. Kalman filter

- continuous observations x and latent space z
- linear Gaussians $p(z_n | z_{n-1})$ & $p(x_n | z_n)$

B. Filtering (forward pass)

→ "What is my estimate for z_n given all observations $x_{1:n}$?"

assume: $p(z_{n-1} | x_{1:n-1})$

prediction:
of latent
space

$$p(z_n | x_{1:n-1}) = \sum_{z_{n-1}} p(z_n, z_{n-1} | x_{1:n-1})$$

x_{n-1} is not
needed if we
are given z_{n-1}

$$\begin{aligned} &= \sum_{z_{n-1}} p(z_n | z_{n-1}, x_{n-1}) p(z_{n-1} | x_{1:n-1}) \\ &= \sum_{z_{n-1}} p(z_n | z_{n-1}) p(z_{n-1} | x_{1:n-1}) \end{aligned}$$

prediction
of
observation:

$$p(z_n | x_{1:n}) = p(z_n | x_n, x_{1:n-1})$$

(Bayes' Rule)

$$\propto p(x_n | z_n, x_{1:n-1}) p(z_n | x_{1:n-1})$$

C. Smoothing (forward & backward pass)

→ "What is my estimate of some z_n given all observations $x_{1:N}$?"

• guessing at a previous latent space given future data

$$p(z_n | x_{1:N}) = p(z_n | x_{1:n}, x_{n+1:N})$$

(Bayes' Rule)

$$\propto p(x_{n+1} | z_n, x_{1:n}) p(z_n | x_{1:n})$$

if know z_n ,
do not need
 $x_{1:n}$

$$\propto p(x_n | z_n) p(z_n | x_{1:n})$$

"emission model"

\approx
likelihood of x_n

"forward pass" (above) \approx "prior"