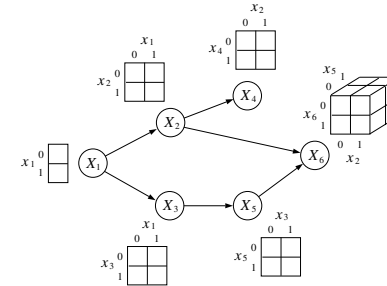


LECTURE 17:

LINEAR GAUSSIAN MODELS

Kevin Murphy  
17 November 2004

We have mostly focused on graphs where all latent nodes are discrete, and all CPDs/potentials are full tables.



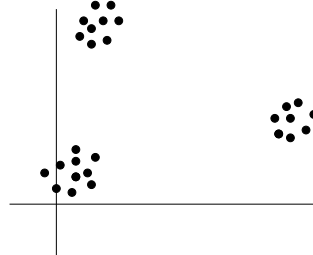
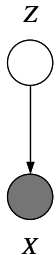
MIXTURES OF GAUSSIANS

- We have also considered discrete latent variables with continuous observed variables.
- e.g., in a mixture of Gaussians:

$$P(Z = i) = \theta_i$$

$$p(X = x | Z = i) = \mathcal{N}(x; \mu_i, \Sigma_i)$$

- This can be used for classification (supervised) and clustering/ vector quantization (unsupervised).



LINEAR GAUSSIAN MODEL

- We now consider the case where all CPDs are linear-Gaussian:

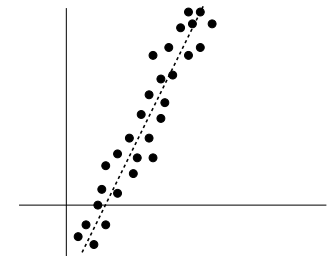
$$p(Z = z) = \mathcal{N}(z; \mu_z, \Sigma_z)$$

$$p(X = x | Z = z) = \mathcal{N}(x; \mu_x + \Lambda z, \Sigma)$$

- i.e., child = linear function of parent plus gaussian noise

$$X = \Lambda Z + \text{noise}, \quad \text{noise} \sim \mathcal{N}(\mu_x, \Sigma_x)$$

- For the supervised case, this is just linear regression:

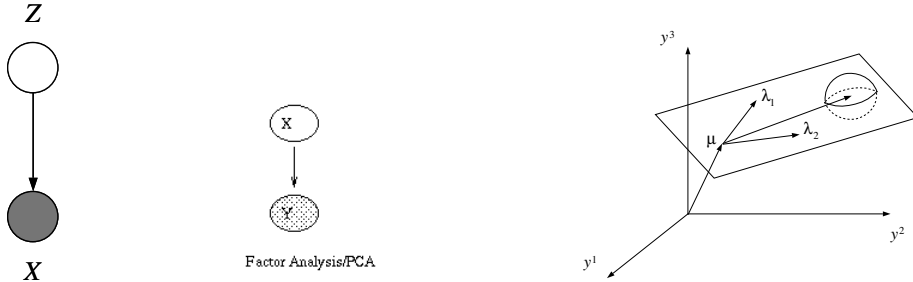


- Unsupervised linear regression is called factor analysis.

$$p(x) = \mathcal{N}(x; 0, I)$$

$$p(y|x) = \mathcal{N}(y; \mu + \Lambda x, \Psi)$$

where  $\Lambda$  is the factor loading matrix and  $\Psi$  is diagonal.



- To generate data, first generate a point within the manifold then add noise. Coordinates of point are components of latent variable.

MARGINAL DATA DISTRIBUTION

- Since a marginal Gaussian times a conditional Gaussian is a joint Gaussian, we can compute the marginal density  $p(y|\theta)$  by integrating out  $x$ .

$$p(y|\theta) = \int_{\mathbf{x}} p(\mathbf{x})p(y|\mathbf{x}, \theta)d\mathbf{x} = \mathcal{N}(y|\mu, \Lambda\Lambda^T + \Psi)$$

which can be done by *completing the square* in the exponent.

- However, since the marginal is Gaussian, we can also just compute its mean and covariance. (Assume noise uncorrelated with data.)

$$E[y] = E[\mu + \Lambda x + \text{noise}] = \mu + \Lambda E[x] + E[\text{noise}]$$

$$= \mu + \Lambda \cdot 0 + 0 = \mu$$

$$\text{Cov}[y] = E[(y - \mu)(y - \mu)^T]$$

$$= E[(\mu + \Lambda x + \text{noise} - \mu)(\mu + \Lambda x + \text{noise} - \mu)^T]$$

$$= E[(\Lambda x + n)(\Lambda x + n)^T] = \Lambda E[\mathbf{x}\mathbf{x}^T]\Lambda^T + E(nn^T)$$

$$= \Lambda\Lambda^T + \Psi$$

- Remember the definition of the mean and covariance of a vector random variable:

$$E[x] = \int_{\mathbf{x}} \mathbf{x}p(\mathbf{x})d\mathbf{x} = \mathbf{m}$$

$$\text{Cov}[x] = E[(\mathbf{x} - \mathbf{m})(\mathbf{x} - \mathbf{m})^T] = \int_{\mathbf{x}} (\mathbf{x} - \mathbf{m})(\mathbf{x} - \mathbf{m})^T p(\mathbf{x})d\mathbf{x} = \mathbf{V}$$

which is the expected value of the outer product of the variable with itself, after subtracting the mean.

- Also, the covariance between two variables:

$$\text{Cov}[\mathbf{x}, \mathbf{y}] = E[(\mathbf{x} - \mathbf{m}_x)(\mathbf{y} - \mathbf{m}_y)^T] = \mathbf{C}$$

$$= \int_{\mathbf{x}\mathbf{y}} (\mathbf{x} - \mathbf{m}_x)(\mathbf{y} - \mathbf{m}_y)^T p(\mathbf{x}, \mathbf{y})d\mathbf{x}d\mathbf{y} = \mathbf{C}$$

which is the expected value of the outer product of one variable with another, after subtracting their means.

Note:  $\mathbf{C}$  is not symmetric.

FA = CONSTRAINED COVARIANCE GAUSSIAN

- Marginal density for factor analysis ( $y$  is  $p$ -dim,  $x$  is  $k$ -dim):

$$p(y|\theta) = \mathcal{N}(y|\mu, \Lambda\Lambda^T + \Psi)$$

- So the effective covariance is the low-rank outer product of two long skinny matrices plus a diagonal matrix:

$$\text{Cov}[y] = \Lambda \begin{matrix} \Lambda^T \\ \hline \end{matrix} + \begin{matrix} \hline \Psi \end{matrix}$$

- In other words, factor analysis is just a constrained Gaussian model. (If  $\Psi$  were not diagonal then we could model any Gaussian and it would be pointless.)

## PROBABILISTIC PRINCIPAL COMPONENT ANALYSIS (PPCA)

---

- In Factor Analysis, we can write the marginal density explicitly:

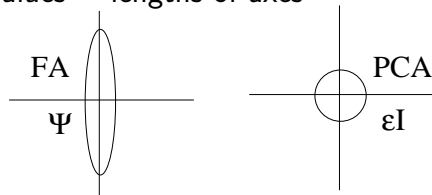
$$p(\mathbf{y}|\theta) = \int_{\mathbf{x}} p(\mathbf{x})p(\mathbf{y}|\mathbf{x}, \theta)d\mathbf{x} = \mathcal{N}(\mathbf{y}|\mu, \Lambda\Lambda^\top + \Psi)$$

- Noise  $\Psi$  must be restricted for model to be interesting. (Why?)
- In Factor Analysis the restriction is that  $\Psi$  is *diagonal* (axis-aligned).
- What if we further restrict  $\Psi = \sigma^2 I$  (ie *spherical*)?
- We get the Probabilistic Principal Component Analysis (PPCA) model: where  $\mu$  is the mean vector, columns of  $\Lambda$  are the *principal components* (usually orthogonal), and  $\sigma^2$  is the *global sensor noise*.

## DIFFERENCE BETWEEN FA AND PPCA

---

- Recall the intuition that Gaussians are hyperellipsoids, where Mean = centre of football, Eigenvectors of covariance matrix = axes of football, Eigenvalues = lengths of axes



- In FA our football is an axis aligned cigar.  
In PPCA our football is a sphere of radius  $\sigma^2$ .
- In FA, the variance of the noise is independent for each dimension; in PPCA, the noise is assumed to be the same.
- PCA looks for directions of large variance,
- Conversely, FA looks for directions of large correlation.

## PCA: ZERO NOISE LIMIT

---

- The traditional PCA model is actually a limit as  $\sigma^2 \rightarrow 0$ .
- Actually PCA is the limit as the ratio of the noise variance on the output to the prior variance on the latent variables goes to zero. We can either achieve this with zero noise or with infinite variance priors.
- (P)PCA is very useful for dimensionality reduction, e.g., for visualizing high-dimensional data.

## MODEL INVARIANCE AND IDENTIFIABILITY

---

- There is *degeneracy* in the FA model.
- Since  $\Lambda$  only appears as outer product  $\Lambda\Lambda^\top$ , the model is invariant to rotation and axis flips of the latent space.
- We can replace  $\Lambda$  with  $\Lambda\mathbf{Q}$  for any unitary matrix  $\mathbf{Q}$  and the model remains the same:  $(\Lambda\mathbf{Q})(\Lambda\mathbf{Q})^\top = \Lambda(\mathbf{Q}\mathbf{Q}^\top)\Lambda^\top = \Lambda\Lambda^\top$ .
- This means that there is no “one best” setting of the parameters. An infinite number of parameters all give the ML score!
- Such models are called un-identifiable since two people both fitting ML params to the identical data will not be guaranteed to identify the same parameters.

- What if we allow the latent variable  $\mathbf{x}$  to have a covariance matrix of its own:  $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|0, \mathbf{P})$ ?
- We can still compute the marginal probability:
 
$$p(\mathbf{y}|\theta) = \int_{\mathbf{x}} p(\mathbf{x})p(\mathbf{y}|\mathbf{x}, \theta)d\mathbf{x} = \mathcal{N}(\mathbf{y}|\mu, \Lambda\mathbf{P}\Lambda^\top + \Psi)$$
- We can always absorb  $\mathbf{P}$  into the loading matrix  $\Lambda$  by diagonalizing it:  $P = EDE^\top$  and setting  $\Lambda = \Lambda ED^{1/2}$ .
- Thus, there is another degeneracy in FA, between  $\mathbf{P}$  and  $\Lambda$ : we can set  $\mathbf{P}$  to be the identity, to be diagonal, whatever we want.
- Traditionally we break this degeneracy by either:
  - set the covariance  $\mathbf{P}$  of the latent variable to be  $I$  (FA) or
  - force the columns of  $\Lambda$  to be orthonormal (PCA)

REVIEW: GAUSSIAN CONDITIONING (JORDAN CH 13)

- Remember the formulas for conditional Gaussian distributions:

$$p\left(\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \mid \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}\right)$$

$$p(\mathbf{x}_1|\mathbf{x}_2) = \mathcal{N}(\mathbf{x}_1|\mathbf{m}_{1|2}, \mathbf{V}_{1|2})$$

$$\mathbf{m}_{1|2} = \mu_1 + \Sigma_{12}\Sigma_{22}^{-1}(\mathbf{x}_2 - \mu_2)$$

$$\mathbf{V}_{1|2} = \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}$$

- Model

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|0, I)$$

$$p(\mathbf{y}|\mathbf{x}, \theta) = \mathcal{N}(\mathbf{y}|\mu + \Lambda\mathbf{x}, \Psi)$$

- Hence the joint distribution of  $\mathbf{x}$  and  $\mathbf{y}$ :

$$p\left(\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \mid \begin{bmatrix} 0 \\ \mu \end{bmatrix}, \begin{bmatrix} I & \Lambda^\top \\ \Lambda & \Lambda\Lambda^\top + \Psi \end{bmatrix}\right)$$

where the corner elements  $\Lambda^\top, \Lambda$  come from  $\text{Cov}[\mathbf{x}, \mathbf{y}]$ :

$$\begin{aligned} \text{Cov}[\mathbf{x}, \mathbf{y}] &= E[(\mathbf{x} - 0)(\mathbf{y} - \mu)^\top] = E[\mathbf{x}(\mu + \Lambda\mathbf{x} + \text{noise} - \mu)^\top] \\ &= E[\mathbf{x}(\Lambda\mathbf{x} + \text{noise})^\top] = \Lambda^\top \end{aligned}$$

- Assume noise is uncorrelated with data or latent variables.

INFERENCE IN FACTOR ANALYSIS

- Apply the Gaussian conditioning formulas to the joint distribution we derived above, where

$$\Sigma_{11} = I$$

$$\Sigma_{12} = \Lambda^\top$$

$$\Sigma_{22}^{-1} = (\Lambda\Lambda^\top + \Psi)^{-1}$$

gives

$$\mathbf{V}_{1|2} = \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}$$

$$= I - \Lambda^\top(\Lambda\Lambda^\top + \Psi)^{-1}\Lambda$$

$$m_{1|2} = m_1 + \Sigma_{12}\Sigma_{22}^{-1}(x_2 - m_2)$$

$$= \Lambda^\top(\Lambda\Lambda^\top + \Psi)^{-1}(\mathbf{y} - \mu)$$

## INFERENCE IN FACTOR ANALYSIS

---

- The previous formula requires inverting a matrix of size  $|y| \times |y|$ .

$$\begin{aligned} p(\mathbf{x}|\mathbf{y}) &= \mathcal{N}(\mathbf{x}|\mathbf{m}, \mathbf{V}) \\ \mathbf{V} &= I - \Lambda^\top(\Lambda\Lambda^\top + \Psi)^{-1}\Lambda \\ \mathbf{m} &= \Lambda^\top(\Lambda\Lambda^\top + \Psi)^{-1}(\mathbf{y} - \mu) \end{aligned}$$

- We now use a good trick for inverting matrices when they can be decomposed into the sum of an easily inverted matrix ( $D$ ) and a low rank outer product. It is called the *matrix inversion lemma*.

$$(D - AB^{-1}A^\top)^{-1} = D^{-1} + D^{-1}A(B - A^\top D^{-1}A)^{-1}A^\top D^{-1}$$

- Apply the matrix inversion lemma we get:

$$\begin{aligned} p(\mathbf{x}|\mathbf{y}) &= \mathcal{N}(\mathbf{x}|\mathbf{m}, \mathbf{V}) \\ \mathbf{V} &= (I + \Lambda^\top\Psi^{-1}\Lambda)^{-1} \\ \mathbf{m} &= \mathbf{V}\Lambda^\top\Psi^{-1}(\mathbf{y} - \mu) \end{aligned}$$

which inverts a matrix of size  $|x| \times |x|$ .

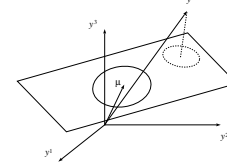
## INFERENCE IS LINEAR PROJECTION

---

- The posterior is

$$\begin{aligned} p(\mathbf{x}|\mathbf{y}) &= \mathcal{N}(\mathbf{x}|\mathbf{m}, \mathbf{V}) \\ \mathbf{V} &= I - \Lambda^\top(\Lambda\Lambda^\top + \Psi)^{-1}\Lambda = (I + \Lambda^\top\Psi^{-1}\Lambda)^{-1} \\ \mathbf{m} &= \Lambda^\top(\Lambda\Lambda^\top + \Psi)^{-1}(\mathbf{y} - \mu) = \mathbf{V}\Lambda^\top\Psi^{-1}(\mathbf{y} - \mu) \end{aligned}$$

- Posterior covariance does not depend on observed data  $y$ ! (Can be precomputed.)
- Also, computing the posterior mean is just a linear operation:  $\mathbf{m} = \beta(\mathbf{y} - \mu)$  where  $\beta$  (hat matrix) can be precomputed.



## INCOMPLETE DATA LOG LIKELIHOOD FUNCTION

---

- Using the marginal density of  $p(y)$ :

$$\begin{aligned} \ell(\theta; \mathcal{D}) &= -\frac{N}{2} \log |\Lambda\Lambda^\top + \Psi| - \frac{1}{2} \sum_n (\mathbf{y}^n - \mu)^\top (\Lambda\Lambda^\top + \Psi)^{-1} (\mathbf{y}^n - \mu) \\ &= -\frac{N}{2} \log |\mathbf{V}| - \frac{1}{2} \text{trace} \left[ \mathbf{V}^{-1} \sum_n (\mathbf{y}^n - \mu)(\mathbf{y}^n - \mu)^\top \right] \\ &= -\frac{N}{2} \log |\mathbf{V}| - \frac{1}{2} \text{trace} \left[ \mathbf{V}^{-1} \mathbf{S} \right] \end{aligned}$$

$\mathbf{V}$  is model covariance;  $\mathbf{S}$  is sample data covariance.

- In other words, we are trying to make the constrained model covariance as close as possible to the observed covariance, where “close” means the trace of the ratio.
- Thus, the sufficient statistics are the same as for the Gaussian: mean  $\sum_n \mathbf{y}^n$  and covariance  $\sum_n (\mathbf{y}^n - \mu)(\mathbf{y}^n - \mu)^\top$ .

## EM FOR FACTOR ANALYSIS

---

- Parameters are coupled nonlinearly in log-likelihood:

$$\ell(\theta; \mathcal{D}) = -\frac{N}{2} \log |\Lambda\Lambda^\top + \Psi| - \frac{1}{2} \text{trace} \left[ (\Lambda\Lambda^\top + \Psi)^{-1} \mathbf{S} \right]$$

- We could use conjugate gradient methods.
- But EM is simpler.
- For E-step we do inference (linear projection)  
For M-step we do linear regression on the expected values.  
More precisely, we maximize the expected complete log likelihood.

$$\mathbf{E} - \text{step} : q_n^{t+1} = p(\mathbf{x}^n | \mathbf{y}^n, \theta^t) = \mathcal{N}(\mathbf{x}^n | \mathbf{m}^n, \mathbf{V}^n)$$

$$\mathbf{M} - \text{step} : \Lambda^{t+1} = \text{argmax}_\Lambda \sum_n \langle \ell_c(\mathbf{x}^n, \mathbf{y}^n) \rangle_{q_n^{t+1}}$$

$$\Psi^{t+1} = \text{argmax}_\Psi \sum_n \langle \ell_c(\mathbf{x}^n, \mathbf{y}^n) \rangle_{q_n^{t+1}}$$

- We know the optimal  $\mu$  is the data mean.  
Assume the mean has been subtracted off  $\mathbf{y}$  from now on.
- The complete likelihood (ignoring mean):

$$\begin{aligned} \ell_c(\Lambda, \Psi) &= \sum_n \log p(\mathbf{x}^n, \mathbf{y}^n) \\ &= \sum_n \log p(\mathbf{x}^n) + \log p(\mathbf{y}^n | \mathbf{x}^n) \\ &= -\frac{N}{2} \log |\Psi| - \frac{1}{2} \sum_n \mathbf{x}^\top \mathbf{x} - \frac{1}{2} \sum_n (\mathbf{y}^n - \Lambda \mathbf{x}^n)^\top \Psi^{-1} (\mathbf{y}^n - \Lambda \mathbf{x}^n) \\ &= -\frac{N}{2} \log |\Psi| - \frac{N}{2} \text{trace}[S \Psi^{-1}] \\ S &= \frac{1}{N} \sum_n (\mathbf{y}^n - \Lambda \mathbf{x}^n)(\mathbf{y}^n - \Lambda \mathbf{x}^n)^\top \end{aligned}$$

M-STEP: OPTIMIZE PARAMETERS

- Derivatives from previous slide:

$$\begin{aligned} \partial \ell_c(\Lambda, \Psi) / \partial \Lambda &= -\Psi^{-1} \sum_n \mathbf{y}_n \mathbf{x}_n^\top + \Psi^{-1} \Lambda \sum_n \mathbf{x}_n \mathbf{x}_n^\top \\ \partial \ell_c(\Lambda, \Psi) / \partial \Psi^{-1} &= +(N/2) \Psi - (N/2) S \end{aligned}$$

- Take the expectation of gradient with respect to  $q^t$  from E-step:

$$\begin{aligned} \langle \ell'_\Lambda \rangle &= -\Psi^{-1} \sum_n \mathbf{y}_n \mathbf{m}_n^\top + \Psi^{-1} \Lambda \sum_n \mathbf{V}_n \\ \langle \ell'_{\Psi^{-1}} \rangle &= +(N/2) \Psi - (N/2) \langle S \rangle \end{aligned}$$

where

$$\begin{aligned} \mathbf{m}_n &= E[X_n | y_n] \\ \mathbf{V}_n &= \text{Var}(X_n | y_n) + E(X_n | y_n) E(X_n | y_n)^\top \\ \langle S \rangle &= \frac{1}{N} \sum_n y_n y_n^\top - y_n m_n^\top \Lambda^\top - \Lambda m_n y_n^\top + \Lambda V_n \Lambda^\top \end{aligned}$$

- Take the derivatives of the complete log likelihood wrt. parameters.
- We will need these identities.

$$\begin{aligned} \frac{\partial}{\partial A} \log |A| &= (A^{-1})^\top \\ \frac{\partial}{\partial A} \text{trace}[B^\top A] &= B \\ \frac{\partial}{\partial A} \text{trace}[B A^\top C A] &= 2C A B \end{aligned}$$

- Hence

$$\begin{aligned} \partial \ell_c(\Lambda, \Psi) / \partial \Lambda &= -\Psi^{-1} \sum_n \mathbf{y}_n \mathbf{x}_n^\top + \Psi^{-1} \Lambda \sum_n \mathbf{x}_n \mathbf{x}_n^\top \\ \partial \ell_c(\Lambda, \Psi) / \partial \Psi^{-1} &= +(N/2) \Psi - (N/2) S \end{aligned}$$

M-STEP: OPTIMIZE PARAMETERS

- Take the expectation of gradient with respect to  $q^t$  from E-step:

$$\begin{aligned} \langle \ell'_\Lambda \rangle &= -\Psi^{-1} \sum_n \mathbf{y}_n \mathbf{m}_n^\top + \Psi^{-1} \Lambda \sum_n \mathbf{V}_n \\ \langle \ell'_{\Psi^{-1}} \rangle &= +(N/2) \Psi - (N/2) \langle S \rangle \end{aligned}$$

- Finally, set the derivatives to zero to solve for optimal parameters:

$$\begin{aligned} \Lambda^{t+1} &= \left( \sum_n \mathbf{y}^n \mathbf{m}_n^\top \right) \left( \sum_n \mathbf{V}_n \right)^{-1} \\ \Psi^{t+1} &= \frac{1}{N} \text{diag} \left[ \sum_n \mathbf{y}_n \mathbf{y}_n^\top - \Lambda^{t+1} \sum_n \mathbf{m}_n \mathbf{y}_n^\top \right] \end{aligned}$$

## FINAL ALGORITHM: EM FOR FACTOR ANALYSIS

- First, set  $\mu$  equal to the sample mean  $(1/N) \sum_n \mathbf{y}_n$ , and subtract this mean from all the data.
- Now run the following iterations:

$$\begin{aligned} \mathbf{E} - \text{step} : q^{t+1} &= p(\mathbf{x}|\mathbf{y}, \theta^t) = \mathcal{N}(\mathbf{x}^n | \mathbf{m}^n, \mathbf{V}^n) \\ \mathbf{V}^n &= (\mathbf{I} + \Lambda^\top \Psi^{-1} \Lambda)^{-1} \\ \mathbf{m}^n &= \mathbf{V}^n \Lambda^\top \Psi^{-1} (\mathbf{y} - \mu) \end{aligned}$$

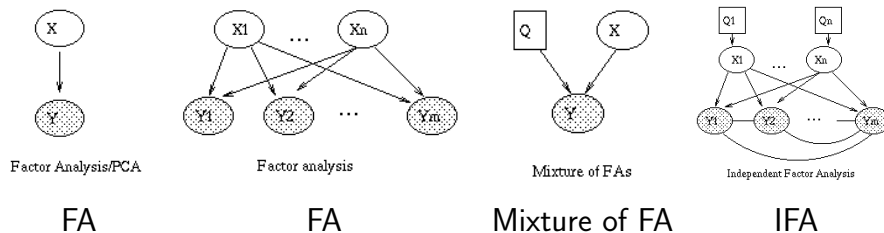
$$\begin{aligned} \mathbf{M} - \text{step} : \Lambda^{t+1} &= \left( \sum_n \mathbf{y}^n \mathbf{m}^{n\top} \right) \left( \sum_n \mathbf{V}^n \right)^{-1} \\ \Psi^{t+1} &= \frac{1}{N} \text{diag} \left[ \sum_n \mathbf{y}^n \mathbf{y}^{n\top} - \Lambda^{t+1} \sum_n \mathbf{m}^n \mathbf{y}^{n\top} \right] \end{aligned}$$

## LEARNING FOR PPCA

- For FA the parameters are coupled in a way that makes it impossible to solve for the ML params directly. We must use EM or other nonlinear optimization techniques.
- But for (P)PCA, the ML params can be solved for directly: The  $k^{\text{th}}$  column of  $\Lambda$  is the  $k^{\text{th}}$  largest eigenvalue of the sample covariance  $\mathbf{S}$  times the associated eigenvector.
- The global sensor noise  $\sigma^2$  is the sum of all the eigenvalues smaller than the  $k^{\text{th}}$  one.
- This technique is good for initializing FA also.
- Note: EM can be used as a fast way to compute eigenvectors!

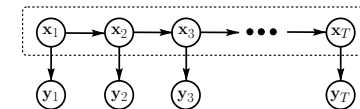
## INDEPENDENT COMPONENTS ANALYSIS (ICA)

- ICA is similar to FA, except it assumes the latent source has *non-Gaussian density*.
- Hence ICA can extract higher order moments (not just second order).
- It is commonly used to solve blind source separation (cocktail party problem).
- Independent Factor Analysis (IFA) is an approximation to ICA where we model the source using a mixture of Gaussians.

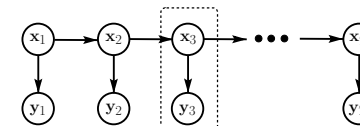


## HIDDEN MARKOV MODELS

- You can think of an HMM as:  
A Markov chain with stochastic measurements.



or a mixture of Gaussians, with latent variables changing over time (lily pad model)



## STATE SPACE MODELS (SSMs)

- An SSM is like an HMM, except the hidden state is continuous valued.
- In general, we can write

$$\begin{aligned}x_t &= f(x_{t-1}) + w_t \\ y_t &= g(x_t) + v_t\end{aligned}$$

where  $f$  is the dynamical (process) model,  $g$  is the observation model.

- $w_t$  is the process noise, embedded inside  $P(X_t|X_{t-1})$ .
- $v_t$  is the observation noise embedded inside  $P(Y_t|X_t)$ .
- For a controlled Markov chain, we condition all quantities on the observed input signals  $u_t$ .

## KALMAN FILTERING AND SMOOTHING

- The Kalman filter is a way to perform exact online inference (sequential Bayesian updating) in an LDS. It is the Gaussian analog of the forwards algorithm for HMMs:

$$P(X_t = i | y_{1:t}) = \alpha_t(i) \propto p(y_t | X_t = i) \sum_j P(X_t = i | X_{t-1} = j) \alpha_{t-1}(j)$$

- The Rauch-Tung-Striavel smoother is a way to perform exact offline inference in an LDS. It is the Gaussian analog of the forwards-backwards algorithm:

$$P(X_t = i | y_{1:T}) \propto \alpha_t(i) \beta_t(i)$$

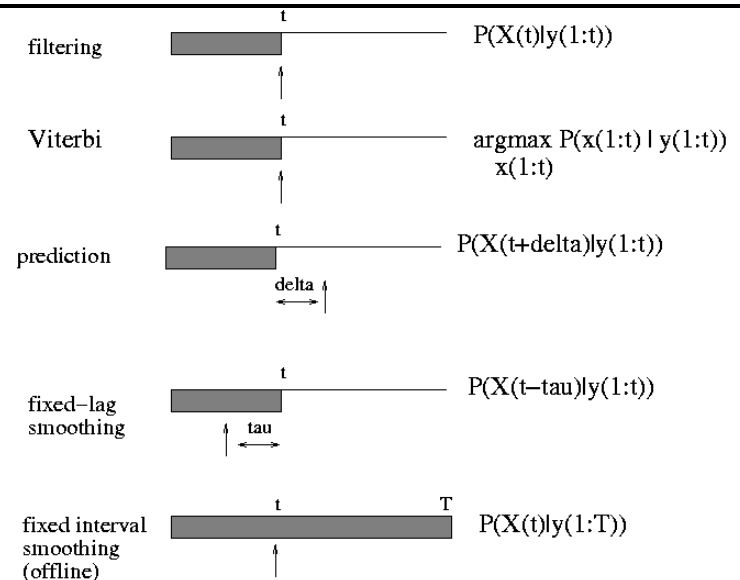
## LINEAR DYNAMIC SYSTEMS (LDS) (JORDAN CH 15)

- An LDS is a special case of a SSM where  $f$  and  $g$  are linear functions and the noise terms are Gaussian.

$$\begin{aligned}x_t &= Ax_{t-1} + w_t \\ y_t &= Bx_t + v_t \\ w_t &\sim \mathcal{N}(0, Q) \\ v_t &\sim \mathcal{N}(0, R)\end{aligned}$$

- An LDS is like factor analysis through time.

## ONLINE VS OFFLINE INFERENCE





## LDS FOR 2D TRACKING

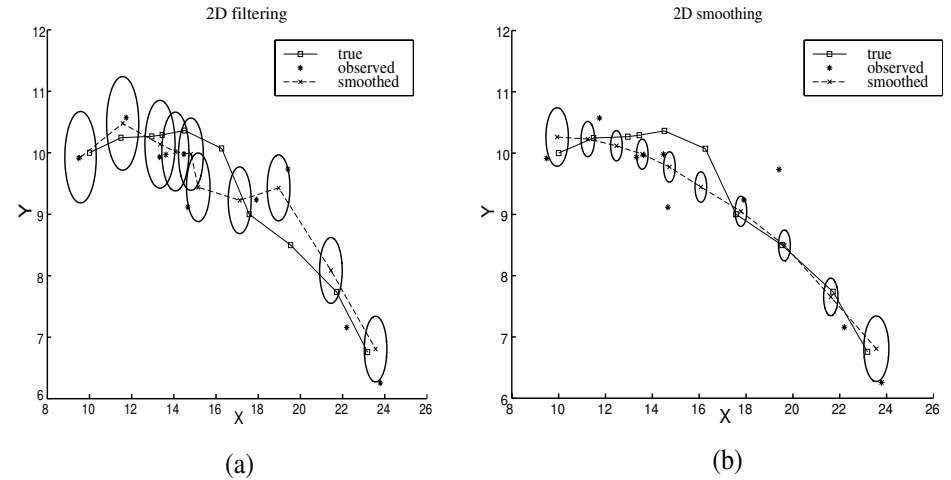
- Dynamics: new position = old position +  $\Delta$  \* velocity + noise  
(constant velocity model, Gaussian acceleration)

$$\begin{pmatrix} x_t^1 \\ x_t^2 \\ \dot{x}_t^1 \\ \dot{x}_t^2 \end{pmatrix} = \Delta \begin{pmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} x_{t-1}^1 \\ x_{t-2}^2 \\ \dot{x}_{t-3}^1 \\ \dot{x}_{t-4}^2 \end{pmatrix} + \text{noise}$$

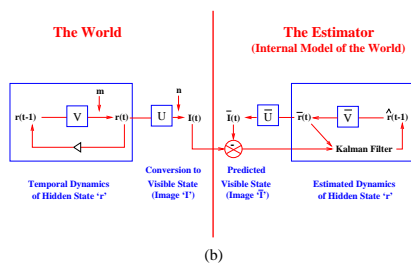
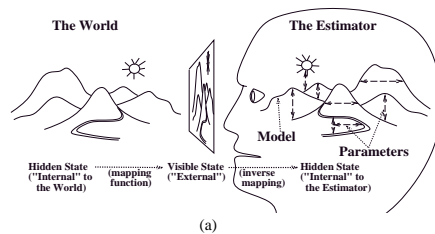
- Observation: project out first two components  
(we observe Cartesian position of object - linear!)

$$\begin{pmatrix} y_t^1 \\ y_t^2 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix} \begin{pmatrix} x_t^1 \\ x_t^2 \\ \dot{x}_t^1 \\ \dot{x}_t^2 \end{pmatrix} + \text{noise}$$

## 2D TRACKING



## KALMAN FILTERING IN THE BRAIN?



## KALMAN FILTERING DERIVATION

- Since all CPDs are linear Gaussian, the system defines a large multi-variate Gaussian.
- Hence all marginals are Gaussian.
- Hence we can represent the belief state  $P(X_t|y_{1:t})$  as a Gaussian with mean  $\hat{x}_{t|t}$  and covariance  $P_{t|t}$ .
- It is common to work with the inverse covariance (precision) matrix  $P_{t|t}^{-1}$ ; this is called *information form*.
- Kalman filtering is a recursive procedure to update the belief state:
  - Predict step: compute  $P(X_{t+1}|y_{1:t})$  from prior belief  $P(X_t|y_{1:t})$  and dynamical model  $P(X_{t+1}|X_t)$ .
  - Update step: compute new belief  $P(X_{t+1}|y_{1:t+1})$  from prediction  $P(X_{t+1}|y_{1:t})$ , observation  $y_{t+1}$  and observation model  $P(y_t|X_t)$ .

## PREDICT STEP

- Dynamical Model:  $x_{t+1} = Ax_t + Gw_t$ ,  $w_t \sim \mathcal{N}(0, Q)$ .
- One step ahead prediction of state:

$$\begin{aligned}\hat{x}_{t+1|t} &= E[X_{t+1}|y_{1:t}] = A\hat{x}_{t|t} \\ P_{t+1|t} &= E[(X_{t+1} - \hat{x}_{t+1|t})(X_{t+1} - \hat{x}_{t+1|t})^T | y_{1:t}] \\ &= AP_{t|t}A^T + GQG^T\end{aligned}$$

- Observation Model:  $y_t = Cy_t + v_t$ ,  $v_t \sim \mathcal{N}(0, R)$ .
- One step ahead prediction of observation:

$$\begin{aligned}E[Y_{t+1}|y_{1:t}] &= E[CX_{t+1} + v_{t+1}|y_{1:t}] = C\hat{x}_{t+1|t} \\ E[(Y_{t+1} - \hat{y}_{t+1|t})(Y_{t+1} - \hat{y}_{t+1|t})^T | y_{1:t}] &= CP_{t+1|t}C^T + R \\ E[(Y_{t+1} - \hat{y}_{t+1|t})(X_{t+1} - \hat{x}_{t+1|t})^T | y_{1:t}] &= CP_{t+1|t}\end{aligned}$$

## KALMAN FILTER EQUATIONS

- Using the conditioning formulas, we get:

$$\begin{aligned}\hat{x}_{t+1|t} &= A\hat{x}_{t|t} \\ P_{t+1|t} &= AP_{t|t}A^T + GQG^T \\ \hat{x}_{t+1|t+1} &= \hat{x}_{t+1|t} + K_{t+1}(y_{t+1} - C\hat{x}_{t+1|t}) \\ P_{t+1|t+1} &= P_{t+1|t} - K_{t+1}CP_{t+1|t}\end{aligned}$$

where  $K_{t+1}$  is the *Kalman gain matrix*

$$K_{t+1} = P_{t+1|t}C^T(CP_{t+1|t}C^T + R)^{-1}$$

- $K_t$  can be precomputed (since it is independent of the data). This rapidly converges to a constant matrix (Ricatti equation).

## UPDATE STEP

- Summarizing results from previous slide, we have  $P(X_{t+1}, Y_{t+1}|y_{1:t}) \sim \mathcal{N}(m_{t+1}, V_{t+1})$ , where

$$m_{t+1} = \begin{pmatrix} \hat{x}_{t+1|t} \\ C\hat{x}_{t+1|t} \end{pmatrix}, \quad V_{t+1} = \begin{pmatrix} P_{t+1|t} & P_{t+1|t}C^T \\ CP_{t+1|t} & CP_{t+1|t}C^T + R \end{pmatrix}$$

- Remember the formulas for conditional Gaussian distributions:

$$\begin{aligned}p\left(\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}\right) &= \mathcal{N}\left(\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \mid \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}\right) \\ p(\mathbf{x}_1|\mathbf{x}_2) &= \mathcal{N}(\mathbf{x}_1|\mathbf{m}_{1|2}, \mathbf{V}_{1|2}) \\ \mathbf{m}_{1|2} &= \mu_1 + \Sigma_{12}\Sigma_{22}^{-1}(\mathbf{x}_2 - \mu_2) \\ \mathbf{V}_{1|2} &= \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}\end{aligned}$$

## EXAMPLE OF KF IN 1D (RUSSELL & NORVIG 2E P554)

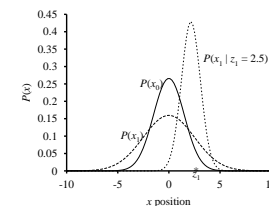
- Consider noisy observations of a 1D particle doing a random walk:  $x_t = x_{t-1} + \mathcal{N}(0, \sigma_x)$ ,  $y_t = x_t + \mathcal{N}(0, \sigma_y)$

- Hence

$$\begin{aligned}P_{t+1|t} &= AP_{t|t}A^T + GQG^T = \sigma_t + \sigma_x \\ K_{t+1} &= P_{t+1|t}C^T(CP_{t+1|t}C^T + R)^{-1} = (\sigma_t + \sigma_x)(\sigma_t + \sigma_x + \sigma_z)^{-1}\end{aligned}$$

$$\hat{x}_{t+1|t+1} = \hat{x}_{t+1|t} + K_{t+1}(y_{t+1} - C\hat{x}_{t+1|t}) = \frac{(\sigma_t + \sigma_x)z_t + \sigma_z\mu_t}{\sigma_t + \sigma_x + \sigma_z}$$

$$P_{t+1|t+1} = P_{t+1|t} - K_{t+1}CP_{t+1|t} = \frac{(\sigma_t + \sigma_x)\sigma_z}{\sigma_t + \sigma_x + \sigma_z}$$



- The KF update of the mean is

$$\hat{x}_{t+1|t+1} = A\hat{x}_{t|t} + K_{t+1}(y_{t+1} - CA\hat{x}_{t|t}) = \frac{(\sigma_t + \sigma_x)z_t + \sigma_z\mu_t}{\sigma_t + \sigma_x + \sigma_z}$$

- The term  $(y_{t+1} - CA\hat{x}_{t|t})$  is called the *innovation*.
- New belief is convex combination of updates from prior and observation, weighted by Kalman Gain matrix:

$$K_{t+1} = P_{t+1|t}C^T(CP_{t+1|t}C^T + R)^{-1}$$

- If the observation is unreliable,  $\sigma_z$  is large so  $K_{t+1}$  is small, so we pay more attention to the prediction.
- If the old prior is unreliable (large  $\sigma_t$ ) or the process is very unpredictable (large  $\sigma_x$ ), we pay more attention to the observation.

- The KF update of the mean is

$$\hat{x}_{t+1|t+1} = A\hat{x}_{t|t} + K_{t+1}(y_{t+1} - CA\hat{x}_{t|t})$$

- Consider the special case where the hidden state is a constant,  $x_t = \theta$ , but the “observation matrix”  $C$  is a time-varying vector,  $C = x_t^T$ .
- Hence  $y_t = x_t^T\theta + v_t$  as in linear regression.
- We can estimate  $\theta$  recursively using the Kalman filter:

$$\hat{\theta}_{t+1} = \hat{\theta}_t + P_{t+1}R^{-1}(y_{t+1} - x_t^T\hat{\theta}_t)x_t$$

This is called the *recursive least squares* (RLS) algorithm.

- We can approximate  $P_{t+1}R^{-1} \approx \eta_{t+1}$  by a scalar constant. This is called the *least mean squares* (LMS) algorithm.
- We can adapt  $\eta_t$  online using stochastic approximation theory.