

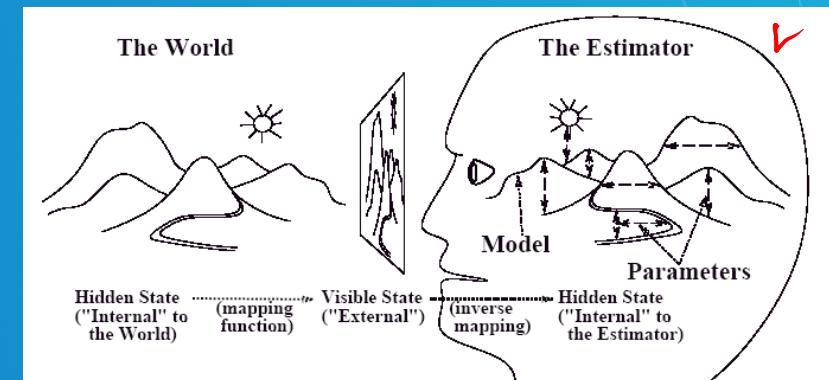
# Probabilistic Graphical Models

## Sequential models

Eric Xing

Lecture 10, February 18, 2019

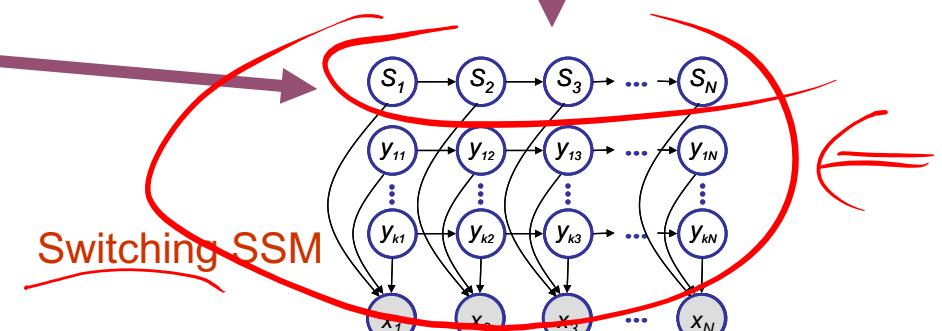
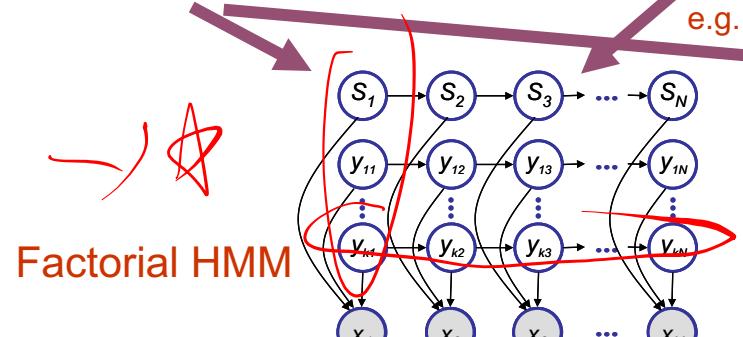
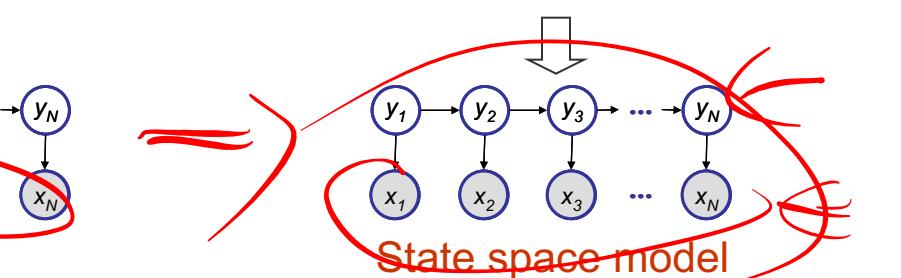
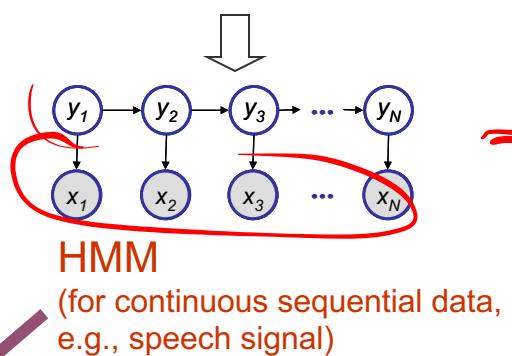
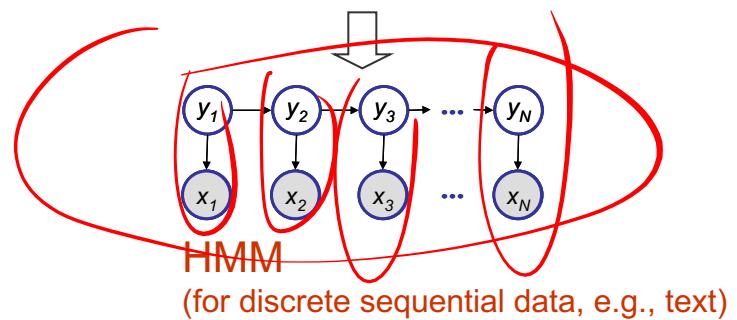
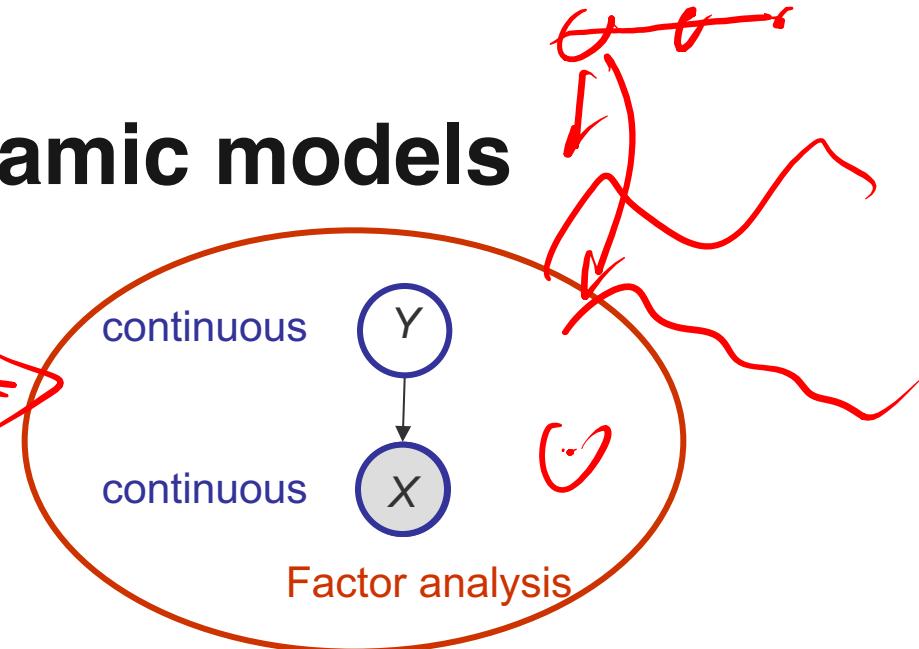
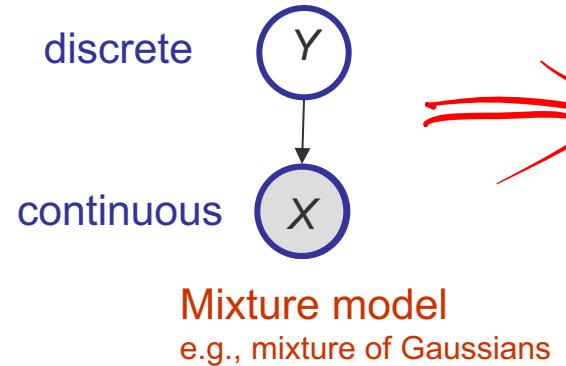
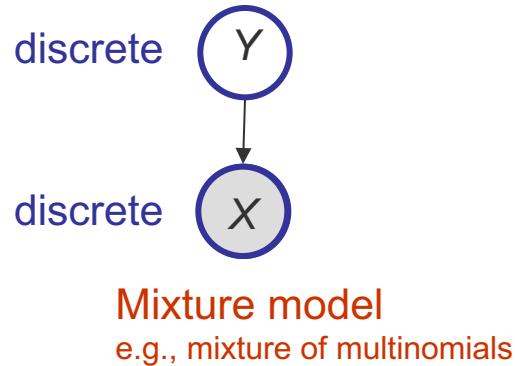
Reading: see class homepage





# A road map to more complex dynamic models

PC(X|Y)





# Recall multivariate Gaussian

- Multivariate Gaussian density:

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

- A joint Gaussian:

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

- How to write down  $p(\mathbf{x}_1)$ ,  $p(\mathbf{x}_1|\mathbf{x}_2)$  or  $p(\mathbf{x}_2|\mathbf{x}_1)$  using the block elements in  $\mu$  and  $\Sigma$ ?
- Formulas to remember:

$$\underbrace{p(\mathbf{x}_2)}_{\mathbf{m}_2^m = \mu_2} = \mathcal{N}(\mathbf{x}_2 | \mathbf{m}_2^m, \mathbf{V}_2^m)$$

$$\mathbf{m}_2^m = \mu_2$$

$$\mathbf{V}_2^m = \Sigma_{22}$$

$$\underbrace{p(\mathbf{x}_1|\mathbf{x}_2)}_{\mathbf{m}_{1|2} = \mu_1 + \Sigma_{12}\Sigma_{22}^{-1}(\mathbf{x}_2 - \mu_2)} = \mathcal{N}(\mathbf{x}_1 | \mathbf{m}_{1|2}, \mathbf{V}_{1|2})$$

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

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





# Review: The matrix inverse lemma

- Consider a block-partitioned matrix:

$$M = \begin{bmatrix} E & F \\ G & H \end{bmatrix}$$

- First we diagonalize  $M$

$$\begin{bmatrix} I & -FH^{-1} \\ 0 & I \end{bmatrix} \begin{bmatrix} E & F \\ G & H \end{bmatrix} \begin{bmatrix} I & 0 \\ -H^{-1}G & I \end{bmatrix} = \begin{bmatrix} E-FH^{-1}G & 0 \\ 0 & H \end{bmatrix}$$

- Schur complement:  $M/H = E-FH^{-1}G$
- Then we inverse, using this formula:  $XYZ = W \Rightarrow Y^{-1} = ZW^{-1}X$

$$\begin{aligned} M^{-1} &= \begin{bmatrix} E & F \\ G & H \end{bmatrix}^{-1} = \begin{bmatrix} I & 0 \\ -H^{-1}G & I \end{bmatrix} \begin{bmatrix} (M/H)^{-1} & 0 \\ 0 & H^{-1} \end{bmatrix} \begin{bmatrix} I & -FH^{-1} \\ 0 & I \end{bmatrix} \\ &= \begin{bmatrix} (M/H)^{-1} & -(M/H)^{-1}FH^{-1} \\ -H^{-1}G(M/H)^{-1} & H^{-1} + H^{-1}G(M/H)^{-1}FH^{-1} \end{bmatrix} = \begin{bmatrix} E^{-1} + E^{-1}F(M/E)^{-1}GE^{-1} & -E^{-1}F(M/E)^{-1} \\ -(M/E)^{-1}GE^{-1} & (M/E)^{-1} \end{bmatrix} \end{aligned}$$

- Matrix inverse lemma

$$(E-FH^{-1}G)^{-1} = E^{-1} + E^{-1}F(H-GE^{-1}F)^{-1}GE^{-1}$$





# Review: Some matrix algebra

- Trace and derivatives
  - Cyclical permutations

$$\text{tr}[A] \stackrel{\text{def}}{=} \sum_i a_{ii}$$

$$\text{tr}[ABC] = \text{tr}[CAB] = \text{tr}[BCA]$$

$$\text{tr}(x^T A x) = x^T A x$$

- Derivatives

$$\frac{\partial}{\partial A} \text{tr}[BA] = B^T$$

$$-\frac{\partial}{\partial A} \text{tr}[x^T A x] = \frac{\partial}{\partial A} \text{tr}[xx^T A] = xx^T$$

- Determinants and derivatives

$$\frac{\partial}{\partial A} \log|A| = A^{-1}$$



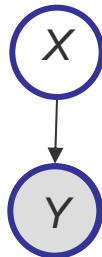


# Factor analysis

$$x \in \mathbb{R}^L$$

$$\mu \sim \Lambda$$

- An unsupervised linear regression model



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

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

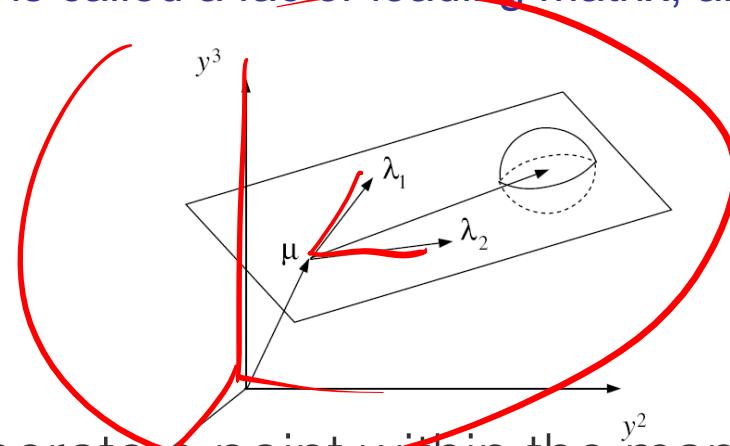
where  $\Lambda$  is called a factor loading matrix, and  $\Psi$  is diagonal.

$$y \in \mathbb{R}^M$$

$$\mu = \Lambda x$$



- Geometric interpretation



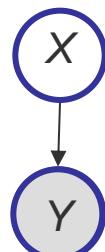
- To generate data, first generate a point within the manifold then add noise. Coordinates of point are components of latent variable.





# Marginal data distribution

- A marginal Gaussian (e.g.,  $p(x)$ ) times a conditional Gaussian (e.g.,  $p(y|x)$ ) is a joint Gaussian
- Any marginal (e.g.,  $p(y)$ ) of a joint Gaussian (e.g.,  $p(x,y)$ ) is also a Gaussian
- Since the marginal is Gaussian, we can determine it by just computing its mean and variance. (Assume noise uncorrelated with data.)



$$\begin{aligned}
 M &= E[Y] = E[\mu + \Lambda X + W] \quad \text{where } W \sim \mathcal{N}(0, \Psi) \\
 &= \mu + \Lambda E[X] + E[W] \\
 &= \mu + 0 + 0 = \mu
 \end{aligned}$$
  

$$\begin{aligned}
 \Sigma_{xy} &= \text{Var}[Y] = E[(Y - \mu)(Y - \mu)^T] \\
 &= E[(\mu + \Lambda X + W - \mu)(\mu + \Lambda X + W - \mu)^T] \\
 &= E[(\Lambda X + W)(\Lambda X + W)^T] \\
 &= \Lambda E[XX^T]\Lambda^T + E[WW^T] \\
 &= \Lambda\Lambda^T + \Psi
 \end{aligned}$$

$\xrightarrow{x \sim p(x)}$   $x \leftarrow \text{f.c. } \tilde{x}$   $\xrightarrow{p(x_{-j}) \sim N(\cdot)}$   
 $y \sim p(y|x)$   $g \sim N(\cdot, \tilde{\sigma}^2)$   $\xrightarrow{p(y|x)} p(y|y) \sim N(\cdot, \tilde{\sigma}^2)$   
 $\Rightarrow p(x|y)$

$$\begin{aligned}
 \checkmark M_x &= \mu, \\
 \checkmark M_y &= M, \\
 \checkmark \Sigma_{xx} &= \mathbb{I} \\
 \checkmark \Sigma_{xy} &= E[(x - \mu_x)(y - \mu_y)^T] \\
 &= E[(x - \mu_x)(\mu + \Lambda x + w - \mu)^T] \\
 &= E[x(\Lambda x + w)^T] \\
 &= E[X(\Lambda X)^T] + 0 \\
 &= \Lambda E[X X^T] + 0
 \end{aligned}$$





# FA joint distribution

- Model

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

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

$$\gamma(x|y)$$

- Covariance between  $\mathbf{x}$  and  $\mathbf{y}$

$$\begin{aligned} \text{Cov}[\mathbf{X}, \mathbf{Y}] &= E[(\mathbf{X} - \mathbf{0})(\mathbf{Y} - \boldsymbol{\mu})^T] = E[\mathbf{X}(\boldsymbol{\mu} + \Lambda \mathbf{X} + \mathbf{W} - \boldsymbol{\mu})^T] \\ &= E[\mathbf{X}\mathbf{X}^T \Lambda^T + \mathbf{X}\mathbf{W}^T] \\ &= \Lambda^T \end{aligned}$$

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

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

- Assume noise is uncorrelated with data or latent variables.





$$\begin{bmatrix} (M/H)^{-1} & -(M/H)^{-1}FH^{-1} \\ -H^{-1}G(M/H)^{-1} & H^{-1} + H^{-1}G(M/H)^{-1}FH^{-1} \end{bmatrix} = \begin{bmatrix} E^{-1} & E^{-1}F(M/E)^{-1}GE^{-1} & -E^{-1}F(M/E)^{-1} \\ -M/E)^{-1}GE^{-1} & (M/E)^{-1} \end{bmatrix}$$

# Inference in Factor Analysis

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

$$\Sigma_{11} = I$$

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

$$\Sigma_{22} = (\Lambda\Lambda^T + \Psi)$$

$$\begin{aligned} E &= \bar{I} \\ F &= \bar{\Lambda}^T \\ G &= \bar{\Lambda} \\ H &= \bar{\Psi} \end{aligned}$$

$$M = \begin{bmatrix} \bar{G} & \bar{F} \\ \bar{F} & \bar{H} \end{bmatrix}$$

we can now derive the posterior of the latent variable  $\mathbf{x}$  given observation  $\mathbf{y}$ , where

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

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

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

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

$$\text{Applying the matrix inversion lemma } = \bar{\Lambda}(\bar{\Lambda}\bar{\Lambda}^T + \bar{\Psi})^{-1}(\mathbf{y} - \mu)$$

- Here we only need to invert a matrix of size  $|\mathbf{x}|'|\mathbf{x}|$ , instead of  $|\mathbf{y}|'|\mathbf{y}|$ .

$$\Rightarrow \mathbf{V}_{1|2} = (I + \bar{\Lambda}^T \bar{\Psi}^{-1} \bar{\Lambda})^{-1}$$

$$(E - EH^{-1}G)^{-1} = (E^{-1} + E^{-1}F(H^{-1}GE^{-1}F)^{-1}GE^{-1})^{-1}$$

$$\mathbf{m}_{1|2} = \mathbf{V}_{1|2} \bar{\Lambda}^T \bar{\Psi}^{-1}(\mathbf{y} - \mu)$$





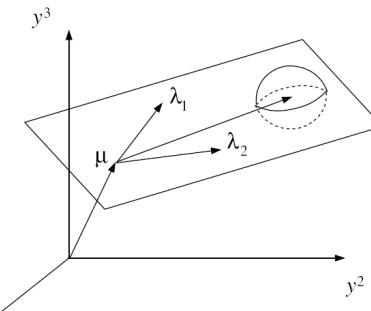
# FA = Constrained-Covariance Gaussian

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

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

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

$$\text{Cov}[\mathbf{y}] = \boldsymbol{\Lambda} \begin{pmatrix} \boldsymbol{\Lambda}^T \\ \boldsymbol{\Psi} \end{pmatrix}$$



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





# Geometric interpretation: inference is linear projection

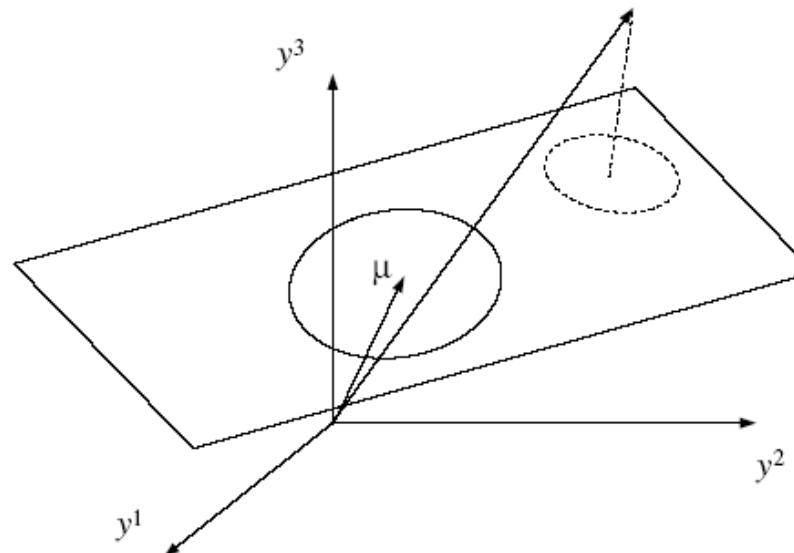
- The posterior is:

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

$$\mathbf{V}_{1|2} = (I + \Lambda^T \Psi^{-1} \Lambda)^{-1}$$

$$\mathbf{m}_{1|2} = \mathbf{V}_{1|2} \Lambda^T \Psi^{-1} (\mathbf{y} - \boldsymbol{\mu})$$

- Posterior covariance does not depend on observed data  $\mathbf{y}$ !
- Computing the posterior mean is just a linear operation:





# Learning FA

- Now, assume that we are given  $\{y_n\}$  (the observation on high-dimensional data) only
- We have derived how to estimate  $x_n$  from  $P(X|Y)$
- How can we learn the model?
  - Loading matrix  $\Lambda$
  - Manifold center  $\mu$
  - Variance  $\Psi$

$$\begin{aligned} \{\Lambda^*, \mu^*, \Psi^*\} &= \arg \max L(\Lambda, \mu, \Psi) \\ &= \arg \max P(Y|\Lambda, \mu, \Psi) \end{aligned}$$





# EM for Factor Analysis

$\mu_{\mathcal{Y}_n}$ ,  $\Sigma_{\mathcal{Y}}$   
 $k(\gamma)$

- Incomplete data log likelihood function (marginal density of  $y$ )

$$\ell(\theta, D) = -\frac{N}{2} \log |\Lambda \Lambda^T + \Psi| - \frac{1}{2} \sum_n (y_n - \mu)^T (\Lambda \Lambda^T + \Psi)^{-1} (y_n - \mu)$$

$$= -\frac{N}{2} \log |\Lambda \Lambda^T + \Psi| - \frac{1}{2} \text{tr} [(\Lambda \Lambda^T + \Psi)^{-1} S],$$

where  $S = \sum_n (y_n - \mu)(y_n - \mu)^T$

- Estimating  $\mu$  is trivial:
- Parameters  $\Lambda$  and  $\Psi$  are coupled nonlinearly in log-likelihood
- Complete log likelihood

$$\ell_c(\theta, D) = \sum_n \log p(x_n, y_n) = \sum_n \log p(x_n) + \log p(y_n | x_n)$$

$$= -\frac{N}{2} \log |I| - \frac{1}{2} \sum_n x_n^T x_n - \frac{N}{2} \log |\Psi| - \frac{1}{2} \sum_n (y_n - \Lambda x_n)^T \Psi^{-1} (y_n - \Lambda x_n)$$

$$= -\frac{N}{2} \log |\Psi| - \frac{1}{2} \sum_n \text{tr}[x_n x_n^T] - \frac{N}{2} \text{tr}[S \Psi^{-1}],$$

where  $S = \frac{1}{N} \sum_n (y_n - \Lambda x_n)(y_n - \Lambda x_n)^T$





# E-step for Factor Analysis

- Compute

$$\langle \ell_e(\theta, D) \rangle_{p(x|y)}$$

$$\underbrace{\langle \ell_e(\theta, D) \rangle}_{\text{red}} = -\frac{N}{2} \log |\Psi| - \frac{1}{2} \sum_n \text{tr} [\langle X_n X_n^T \rangle] - \frac{N}{2} \text{tr} [\langle S \rangle \Psi^{-1}]$$

$P(X|y)$

$\langle x \rangle$

$\langle X X^T \rangle$

$$\langle S \rangle = \frac{1}{N} \sum_n (y_n y_n^T - y_n \langle X_n^T \rangle \Lambda^T - \Lambda \langle X_n^T \rangle y_n^T + \Lambda \langle X_n X_n^T \rangle \Lambda^T)$$

$$\underbrace{\langle X_n \rangle}_{\text{red}} = E[X_n | y_n]$$

$$\underbrace{\langle X_n X_n^T \rangle}_{\text{red}} = \text{Var}[X_n | y_n] + E[X_n | y_n] E[X_n | y_n]^T$$

- Recall that we have derived:

$$V_{1|2} = (I + \Lambda^T \Psi^{-1} \Lambda)^{-1}$$

$$m_{1|2} = V_{1|2} \Lambda^T \Psi^{-1} (\mathbf{y} - \mu)$$

$$\Rightarrow \underbrace{\langle X_n \rangle}_{\text{red}} = m_{x_n | y_n} = V_{1|2} \Lambda^T \Psi^{-1} (y_n - \mu)$$

and

$$\underbrace{\langle X_n X_n^T \rangle}_{\text{red}} = \underbrace{V_{1|2}}_{\text{red}} + \underbrace{m_{x_n | y_n} m_{x_n | y_n}^T}_{\text{red}}$$





# M-step for Factor Analysis

- Take the derivatives of the expected complete log likelihood wrt. parameters.
    - Using the trace and determinant derivative rules:

$$\begin{aligned} \cancel{\frac{\partial}{\partial \Psi^{-1}} \langle \ell_c \rangle} &= \frac{\partial}{\partial \Psi^{-1}} \left( -\frac{N}{2} \log |\Psi| - \frac{1}{2} \sum_n \text{tr} [\langle X_n X_n^T \rangle] - \frac{N}{2} \text{tr} [\langle S \rangle \Psi^{-1}] \right) \\ &= \frac{N}{2} \Psi - \frac{N}{2} \langle S \rangle \quad \Rightarrow \quad \Psi^{t+1} = \langle S \rangle \end{aligned}$$

$$\begin{aligned}
\underbrace{\frac{\partial}{\partial \Lambda} \langle \ell_c \rangle}_{=} &= \frac{\partial}{\partial \Lambda} \left( -\frac{N}{2} \log |\Psi| - \frac{1}{2} \sum_n \text{tr} [\langle X_n X_n^T \rangle] - \frac{N}{2} \text{tr} [\langle S \rangle \Psi^{-1}] \right) = -\frac{N}{2} \Psi^{-1} \frac{\partial}{\partial \Lambda} \langle S \rangle \\
&= -\frac{N}{2} \Psi^{-1} \frac{\partial}{\partial \Lambda} \left( \frac{1}{N} \sum_n (Y_n Y_n^T - Y_n \langle X_n^T \rangle \Lambda^T - \Lambda \langle X_n^T \rangle Y_n^T + \Lambda \langle X_n X_n^T \rangle \Lambda^T) \right) \\
&= \Psi^{-1} \sum_n Y_n \langle X_n^T \rangle - \Psi^{-1} \Lambda \sum_n \langle X_n X_n^T \rangle \quad \Rightarrow \quad \Lambda^{t+1} = \left( \sum_n Y_n \langle X_n^T \rangle \right) \left( \sum_n \langle X_n X_n^T \rangle \right)^{-1}
\end{aligned}$$

$$\begin{aligned}
 & M_{\text{max}} \left( \frac{\partial \ell(\theta)}{\partial \theta} \right) = 0 \\
 & \ell(\mathbf{x}) = \prod_{i=1}^n p(x_i | \theta) = \prod_{i=1}^n \frac{f(x_i)}{f(x_i | \theta)} = \frac{1}{f(\mathbf{x} | \theta)} \prod_{i=1}^n f(x_i) \\
 & \text{likelihood wrt. } \theta = \frac{1}{f(\mathbf{x} | \theta)} \prod_{i=1}^n f(x_i) \\
 & \text{enforce } \left[ \frac{1}{f(\mathbf{x} | \theta)} \right] \left[ \prod_{i=1}^n f(x_i) \right]
 \end{aligned}$$

Lang. E(x)  
p(x|y)





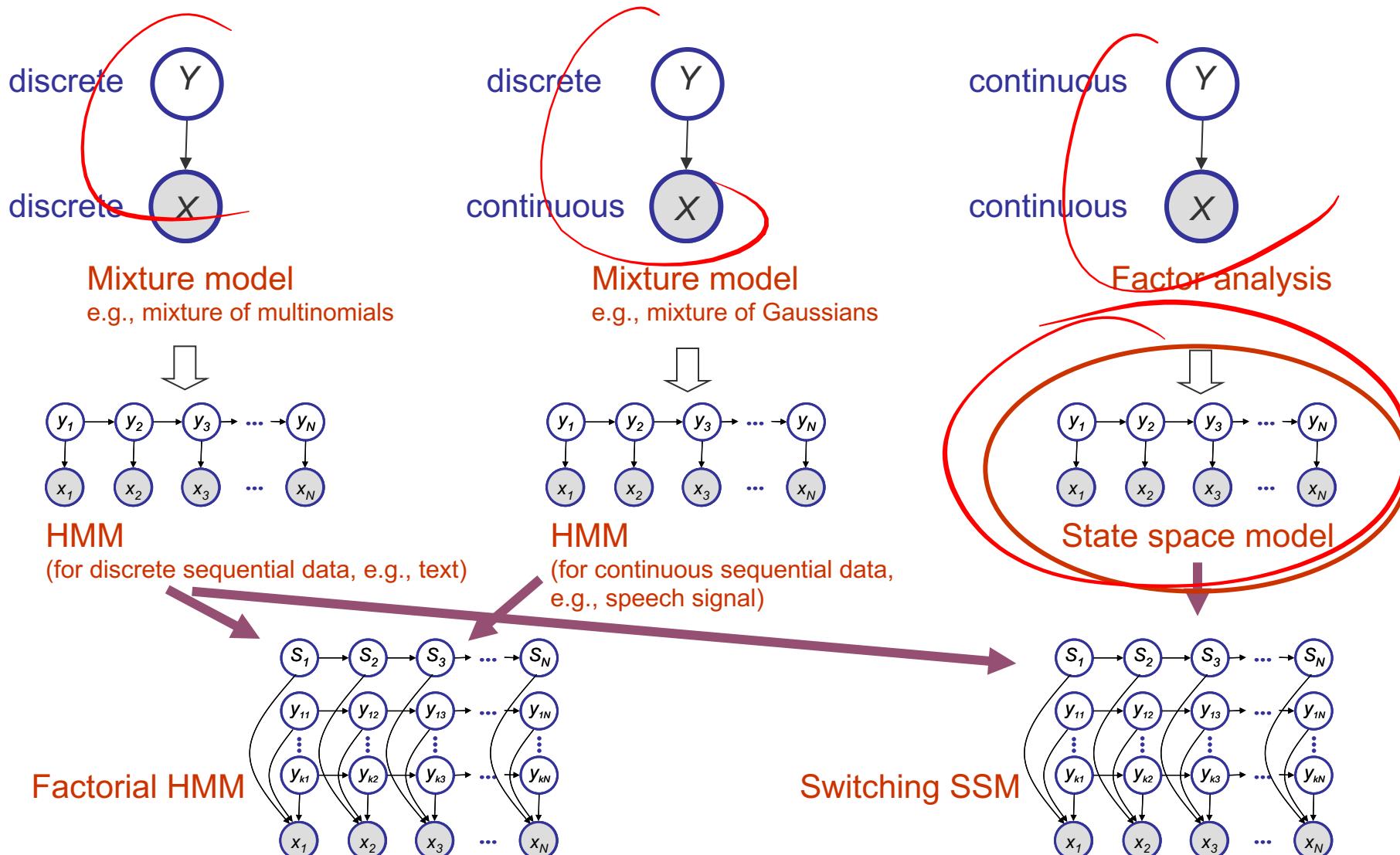
# Model Invariance and Identifiability

- ❑ There is *degeneracy* in the FA model.
- ❑ Since  $\Lambda$  only appears as outer product  $\Lambda\Lambda^T$ , the model is invariant to rotation and axis flips of the latent space.
- ❑ We can replace  $\Lambda$  with  $\Lambda Q$  for any orthonormal matrix  $Q$  and the model remains the same:  $(\Lambda Q)(\Lambda Q)^T = \Lambda(QQ^T)\Lambda^T = \Lambda\Lambda^T$ .
- ❑ 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 parameters to the identical data will not be guaranteed to identify the same parameters.





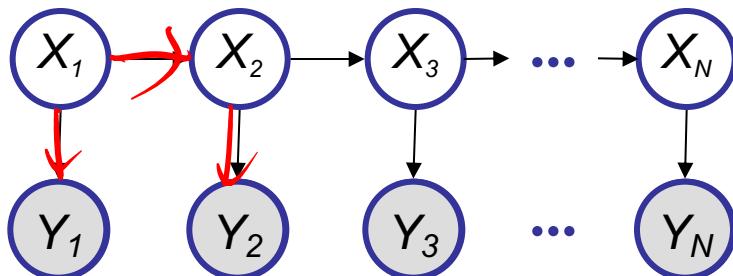
# A road map to more complex dynamic models





# State space models (SSM)

- A sequential FA or a continuous state HMM



$$\begin{aligned} \mathbf{x}_t &= A\mathbf{x}_{t-1} + Gw_t \\ \mathbf{y}_t &= C\mathbf{x}_{t-1} + v_t \\ w_t &\sim \mathcal{N}(0; Q), \quad v_t \sim \mathcal{N}(0; R) \\ \mathbf{x}_0 &\sim \mathcal{N}(0; \Sigma_0), \end{aligned}$$

This is a linear dynamic system.

- In general,

$$\begin{aligned} \mathbf{x}_t &= f(\mathbf{x}_{t-1}) + Gw_t \\ \mathbf{y}_t &= g(\mathbf{x}_{t-1}) + v_t \end{aligned}$$

where  $f$  is an (arbitrary) dynamic model, and  $g$  is an (arbitrary) observation model





# LDS for 2D tracking

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

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

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

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

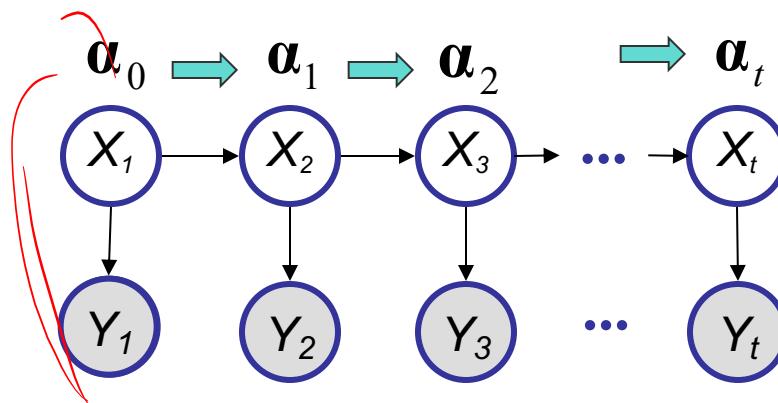




# The inference problem 1

- Filtering → given  $y_1, \dots, y_t$ , estimate  $x_t$ :  $P(x_t | y_{1:t})$ 
  - The Kalman filter is a way to perform exact online inference (sequential Bayesian updating) in an LDS.
  - It is the Gaussian analog of the forward algorithm for HMMs:

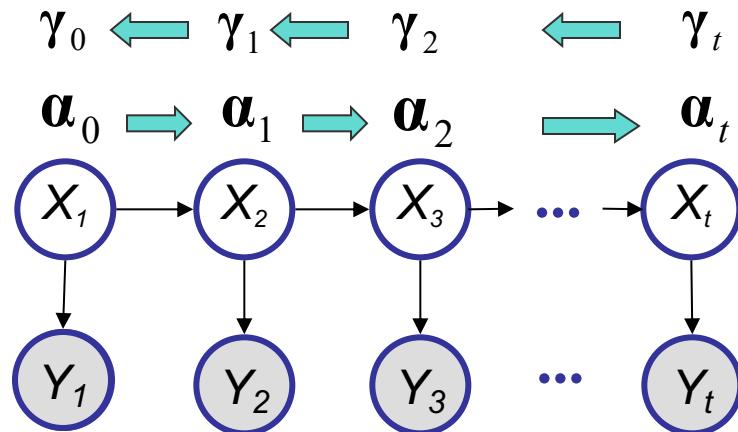
$$p(X_t = i | y_{1:t}) = \alpha_t^i \propto p(y_t | X_t = i) \sum_j p(X_t = i | X_{t-1} = j) \alpha_{t-1}^j$$





# The inference problem 2

- Smoothing → given  $\mathbf{y}_1, \dots, \mathbf{y}_T$ , estimate  $\mathbf{x}_t (t < T)$ 
  - The Rauch-Tung-Strievel smoother is a way to perform exact off-line inference in an LDS. It is the Gaussian analog of the forwards-backwards (alpha-gamma) algorithm:

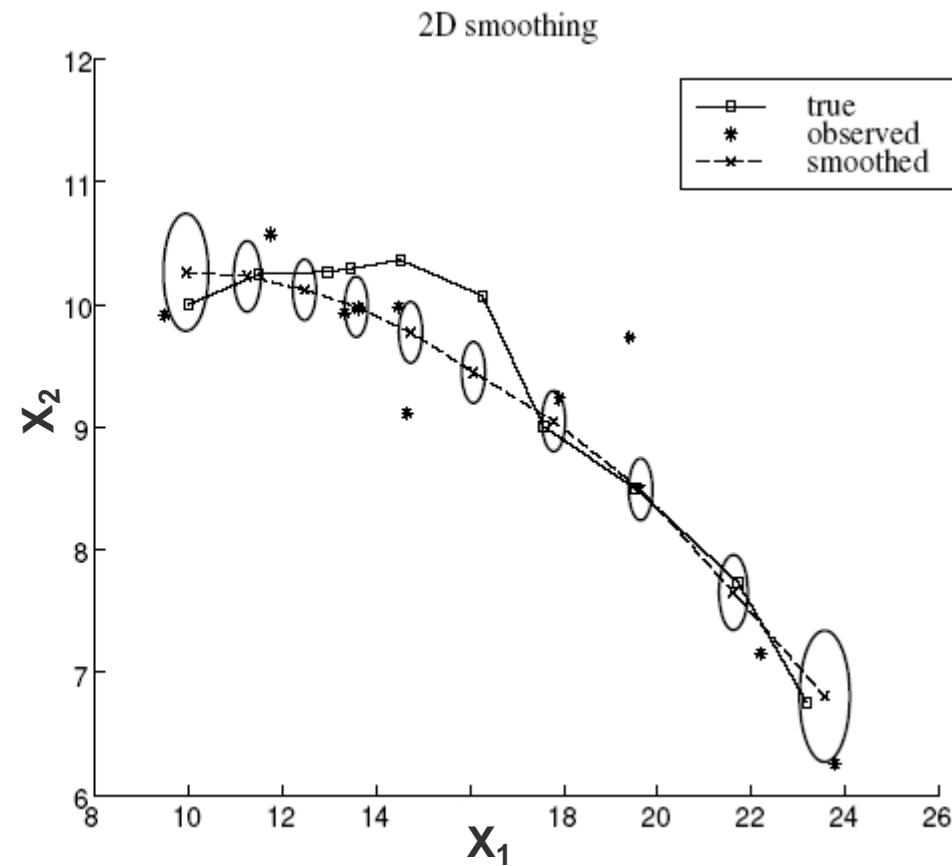
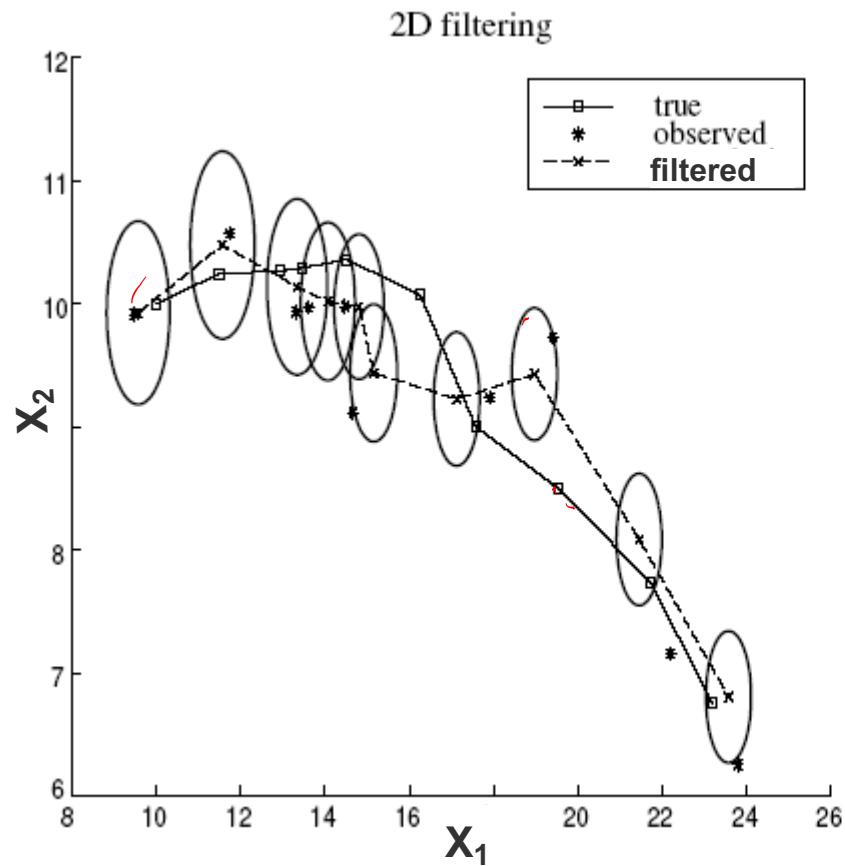


$$p(\mathbf{X}_t = i | \mathbf{y}_{1:T}) = \gamma_t^i \propto \sum_j \alpha_t^i P(X_{t+1}^j | X_i^j) \gamma_{t+1}^j$$



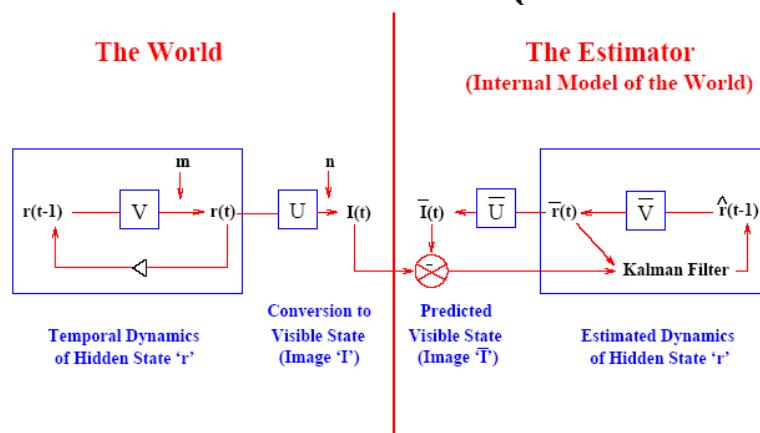
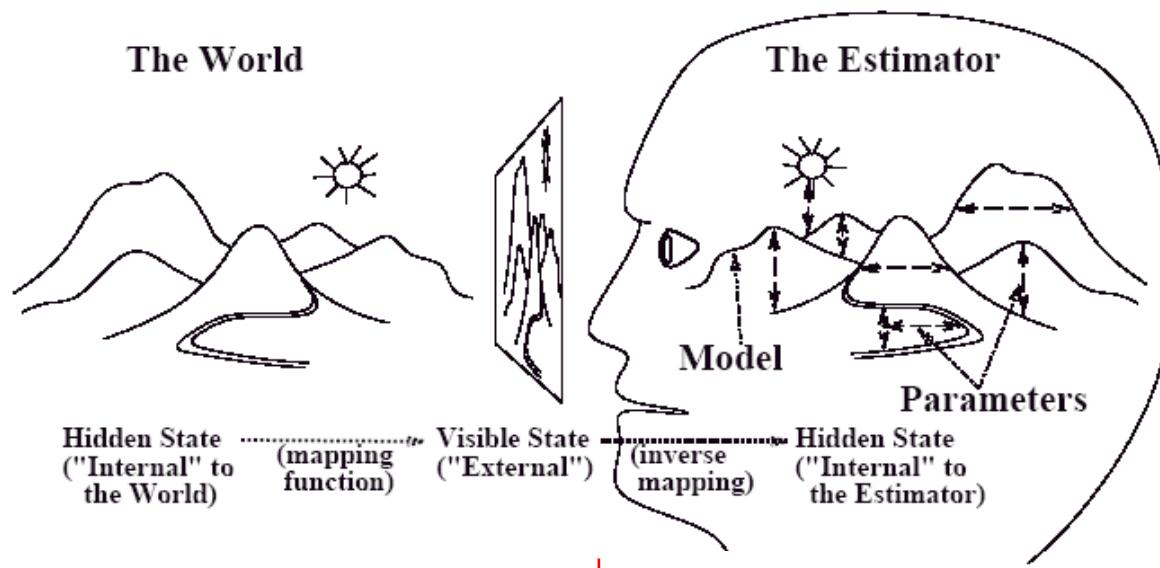


# 2D tracking





# Kalman filtering in the brain?

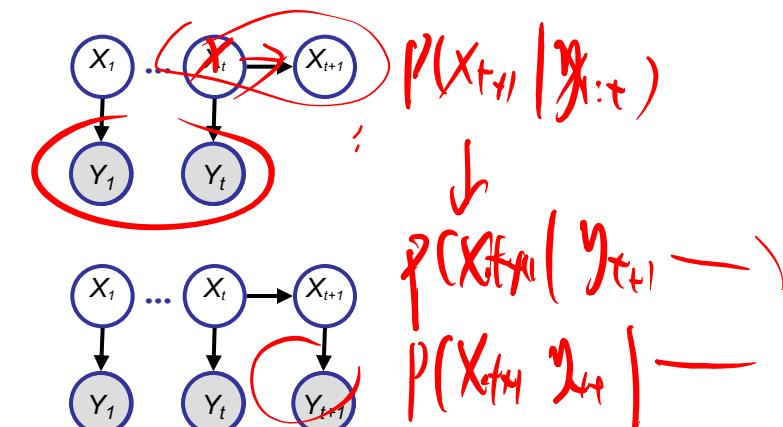




# Kalman filtering derivation

- Since all CPDs are linear Gaussian, the system defines a large multivariate Gaussian. ✓
  - Hence all marginals are Gaussian.
  - Hence we can represent the belief state  $p(\mathbf{X}_t | \mathbf{y}_{1:t})$  as a Gaussian with mean and covariance:  $\mathbb{E}(\mathbf{x}_t | \mathbf{y}_{1:t}) = \mathbf{M}_{1:t}$  ,  $\mathbb{E}((\mathbf{x}_t - \mathbf{M}_{1:t})(\mathbf{x}_t - \mathbf{M}_{1:t})^T) = \mathbf{P}_{1:t}$
  - It is common to work with the inverse covariance (precision) matrix, this is called information form.
  - Kalman filtering is a recursive procedure to update the belief state:
    - Predict step: compute  $p(\mathbf{X}_{t+1} | \mathbf{y}_{1:t})$  from prior belief  $p(\mathbf{X}_t | \mathbf{y}_{1:t})$  and dynamical model  $p(\mathbf{X}_{t+1} | \mathbf{X}_t)$  --- time update
    - Update step: compute new belief  $p(\mathbf{X}_{t+1} | \mathbf{y}_{1:t+1})$  from prediction  $p(\mathbf{X}_{t+1} | \mathbf{y}_{1:t})$ , observation  $\mathbf{y}_{t+1}$  and observation model  $p(\mathbf{y}_{t+1} | \mathbf{X}_{t+1})$  --- measurement update

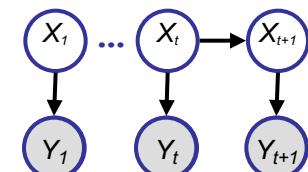
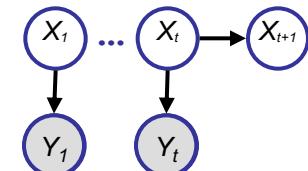
$$\begin{aligned}
 & P(x_t | y_{1:t}) \quad ? \\
 & \frac{P(x_t | y_{t-1})}{P(x_t | x_{t-1})} \quad ) \\
 & \downarrow \quad \downarrow \\
 & \text{The system defines a large } \\
 & \propto P(x_t, \dot{x}_t | y_{t-1}) \\
 & P(x_t, \dot{x}_t | y_{t-1}, y_t) \\
 & \text{as a Gaussian with mean } \\
 & e^{-\lambda_t} (Ae(\mu_t) - P_{0|0}) \\
 & \text{and covariance matrix } P_{0|0} \text{ is called information}
 \end{aligned}$$





# Kalman filtering derivation

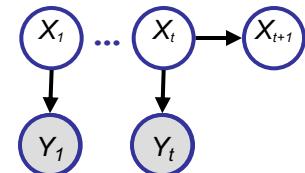
- Kalman filtering is a recursive procedure to update the belief state:
  - Predict step: compute  $p(\mathbf{X}_{t+1}|\mathbf{y}_{1:t})$  from prior belief  $p(\mathbf{X}_t|\mathbf{y}_{1:t})$  and dynamical model  $p(\mathbf{X}_{t+1}|\mathbf{X}_t)$  --- time update
  - Update step: compute new belief  $p(\mathbf{X}_{t+1}|\mathbf{y}_{1:t+1})$  from prediction  $p(\mathbf{X}_{t+1}|\mathbf{y}_{1:t})$ , observation  $\mathbf{y}_{t+1}$  and observation model  $p(\mathbf{y}_{t+1}|\mathbf{X}_{t+1})$  --- measurement update



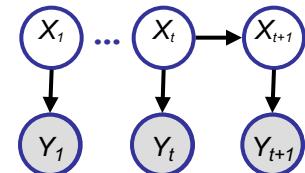


# Predict step

- Dynamical Model:  $\mathbf{x}_{t+1} = A\mathbf{x}_t + Gw_t, \quad w_t \sim \mathcal{N}(0; Q)$ 
  - One step ahead prediction of state:



- Observation model:  $\mathbf{y}_t = C\mathbf{x}_t + v_t, \quad v_t \sim \mathcal{N}(0; R)$ 
  - One step ahead prediction of observation:





# Predict step

- Dynamical Model:  $\mathbf{x}_{t+1} = A\mathbf{x}_t + Gw_t, \quad w_t \sim \mathcal{N}(0; Q)$

- One step ahead prediction of state:

$$\hat{\mathbf{x}}_{t+1|t} = E(\mathbf{X}_{t+1} | \mathbf{y}_1, \dots, \mathbf{y}_t) = A\hat{\mathbf{x}}_{t|t}$$

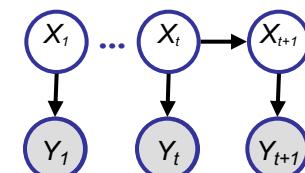
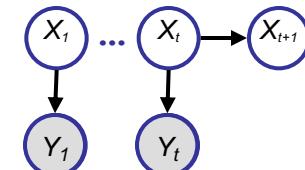
$$P_{t+1|t} = E(\mathbf{X}_{t+1} - \hat{\mathbf{x}}_{t+1|t})(\mathbf{X}_{t+1} - \hat{\mathbf{x}}_{t+1|t})^T | \mathbf{y}_1, \dots, \mathbf{y}_t$$

$$= E(AX_t + Gw_t - \hat{\mathbf{x}}_{t+1|t})(AX_t + Gw_t - \hat{\mathbf{x}}_{t+1|t})^T | \mathbf{y}_1, \dots, \mathbf{y}_t$$

$$= AP_{t|t}A + GQG^T$$

- Observation model:  $\mathbf{y}_t = C\mathbf{x}_t + v_t, \quad v_t \sim \mathcal{N}(0; R)$

- One step ahead prediction of observation:



$$E(\mathbf{y}_{t+1} | \mathbf{y}_1, \dots, \mathbf{y}_t) = E(C\mathbf{X}_{t+1} + v_{t+1} | \mathbf{y}_1, \dots, \mathbf{y}_t) = C\hat{\mathbf{x}}_{t+1|t}$$

$$E(\mathbf{y}_{t+1} - \hat{\mathbf{y}}_{t+1|t})(\mathbf{y}_{t+1} - \hat{\mathbf{y}}_{t+1|t})^T | \mathbf{y}_1, \dots, \mathbf{y}_t = CP_{t+1|t}C^T + R$$

$$E(\mathbf{y}_{t+1} - \hat{\mathbf{y}}_{t+1|t})(\mathbf{X}_{t+1} - \hat{\mathbf{x}}_{t+1|t})^T | \mathbf{y}_1, \dots, \mathbf{y}_t = CP_{t+1|t}$$





# Update step

- Summarizing results from previous slide, we have  $p(\mathbf{X}_{t+1}, \mathbf{Y}_{t+1} | \mathbf{y}_{1:t}) \sim \mathcal{N}(m_{t+1}, V_{t+1})$ , 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:

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

$$p(\mathbf{x}_2) = \mathcal{N}(\mathbf{x}_2 | \mathbf{m}_2^m, \mathbf{V}_2^m)$$

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

$$\mathbf{m}_2^m = \mu_2$$

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

$$\mathbf{V}_2^m = \Sigma_{22}$$

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





# Kalman Filter

- Measurement updates:

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

$$P_{t+1|t+1} = P_{t+1|t} - KCP_{t+1|t}$$

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

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

- Time updates:

$$\hat{\mathbf{x}}_{t+1|t} = A\hat{\mathbf{x}}_{t|t}$$

$$P_{t+1|t} = AP_{t|t}A + GQG^T$$

- $K_t$  can be pre-computed (since it is independent of the data).





# Example of KF in 1D

- Consider noisy observations of a 1D particle doing a random walk:

$$x_{t|t-1} = x_{t-1} + w, \quad w \sim \mathcal{N}(0, \sigma_x) \quad z_t = x_t + v, \quad v \sim \mathcal{N}(0, \sigma_z)$$

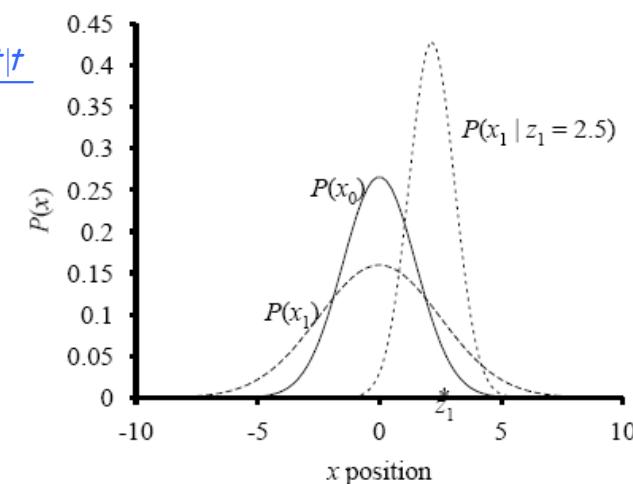
- KF equations:

$$P_{t+1|t} = AP_{t|t}A^T + GQG^T = \sigma_t + \sigma_x, \quad \hat{x}_{t+1|t} = A\hat{x}_{t|t} = \hat{x}_{t|t}$$

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

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

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





# KF intuition

- The KF update of the mean is

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

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

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

- If the observation is unreliable,  $\sigma_z$  (i.e.,  $R$ ) 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.





# Complexity of one KF step

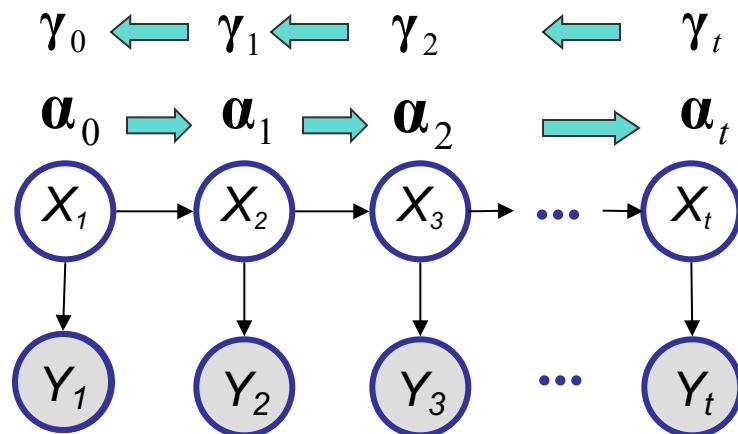
- Let  $X_t \in \mathbb{R}^{N_x}$  and  $Y_t \in \mathbb{R}^{N_y}$ ,
- Computing  $P_{t+1|t} = AP_{t|t}A^T + GQG^T$  takes  $O(N_x^2)$  time, assuming dense  $P$  and dense  $A$ .
- Computing  $K_{t+1} = P_{t+1|t}C^T(CP_{t+1|t}C^T + R)^{-1}$  takes  $O(N_y^3)$  time.
- So overall time is, in general,  $\max\{N_x^2, N_y^3\}$





# The inference problem 2

- ❑ Smoothing → given  $\mathbf{y}_1, \dots, \mathbf{y}_T$ , estimate  $\mathbf{x}_t (t < T)$ 
  - ❑ The Rauch-Tung-Strievel smoother is a way to perform exact off-line inference in an LDS. It is the Gaussian analog of the forwards-backwards (alpha-gamma) algorithm:



$$p(X_t = i | y_{1:T}) = \gamma_t^i \propto \sum_j \alpha_t^j P(X_{t+1}^j | X_i^j) \gamma_{t+1}^j$$

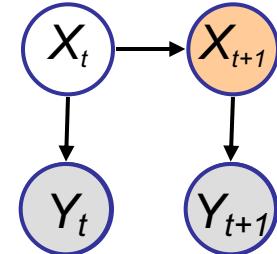




# Rauch-Tung-Strievel smoother

$$\hat{\mathbf{x}}_{t|T} = \hat{\mathbf{x}}_{t|t} + L_t (\hat{\mathbf{x}}_{t+1|T} - \hat{\mathbf{x}}_{t+1|t})$$

$$P_{t|T} = P_{t|t} + L_t (P_{t+1|T} - P_{t+1|t}) L_t^T \quad L_t = P_{t|t} A^T P_{t+1|t}^{-1}$$



- General structure: KF results + the difference of the "smoothed" and predicted results of the next step
- Backward computation: Pretend to know things at  $t+1$  -- such conditioning makes things simple and we can remove this condition finally

□ The difficulty:  $\mathcal{X}_t | \mathcal{Y}_1, \dots, \mathcal{Y}_T$

□ The trick:  $E[\mathcal{X} | \mathcal{Z}] = E[E[\mathcal{X} | \mathcal{Y}, \mathcal{Z}] | \mathcal{Z}]$  (Hw!)

$$Var[\mathcal{X} | \mathcal{Z}] = Var[E[\mathcal{X} | \mathcal{Y}, \mathcal{Z}] | \mathcal{Z}] + E[Var[\mathcal{X} | \mathcal{Y}, \mathcal{Z}] | \mathcal{Z}]$$

$$\begin{aligned}\hat{\mathbf{x}}_{t|T} &\stackrel{\text{def}}{=} E[\mathcal{X}_t | \mathcal{Y}_1, \dots, \mathcal{Y}_T] = E[E[\mathcal{X}_t | \mathcal{X}_{t+1}, \mathcal{Y}_1, \dots, \mathcal{Y}_T] | \mathcal{Y}_1, \dots, \mathcal{Y}_T] \\ &= E[E[\mathcal{X}_t | \mathcal{X}_{t+1}, \mathcal{Y}_1, \dots, \mathcal{Y}_T] | \mathcal{Y}_1, \dots, \mathcal{Y}_T] \\ &= E[\mathcal{X}_t | \mathcal{X}_{t+1}, \mathcal{Y}_1, \dots, \mathcal{Y}_T]\end{aligned}$$

Same for  $P_{t|T}$





# RTS derivation

- Following the results from previous slide, we need to derive  $p(\mathbf{X}_{t+1}, \mathbf{X}_t | \mathbf{y}_{1:t}) \sim \mathcal{N}(\mathbf{m}, \mathbf{V})$ , where

$$\mathbf{m} = \begin{pmatrix} \hat{\mathbf{x}}_{t|t} \\ \hat{\mathbf{x}}_{t+1|t} \end{pmatrix},$$

$$\mathbf{V} = \begin{pmatrix} P_{t|t} & P_{t|t} \mathbf{A}^T \\ \mathbf{A} P_{t|t} & P_{t+1|t} \end{pmatrix},$$

- all the quantities here are available after a forward KF pass
- Remember the formulas for conditional Gaussian distributions:

$$p\left(\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \mid \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix}\right),$$

$$\begin{aligned} p(\mathbf{x}_2) &= \mathcal{N}(\mathbf{x}_2 \mid \mathbf{m}_2^m, \mathbf{V}_2^m) & p(\mathbf{x}_1 \mid \mathbf{x}_2) &= \mathcal{N}(\mathbf{x}_1 \mid \mathbf{m}_{12}, \mathbf{V}_{12}) \\ \mathbf{m}_2^m &= \boldsymbol{\mu}_2 & \mathbf{m}_{12} &= \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} (\mathbf{x}_2 - \boldsymbol{\mu}_2) \\ \mathbf{V}_2^m &= \boldsymbol{\Sigma}_{22} & \mathbf{V}_{12} &= \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21} \end{aligned}$$

- The RTS smoother

$$\begin{aligned} \hat{\mathbf{x}}_{t|T} &= E[\mathbf{X}_t \mid \mathbf{X}_{t+1}, \mathbf{y}_1, \dots, \mathbf{y}_T] \\ &= \hat{\mathbf{x}}_{t|t} + L_t (\hat{\mathbf{x}}_{t+1|T} - \hat{\mathbf{x}}_{t+1|t}) \end{aligned}$$

$$\begin{aligned} P_{t|T} &\stackrel{\text{def}}{=} \text{Var}[\hat{\mathbf{x}}_{t|T} \mid \mathbf{y}_{1:T}] + E[\text{Var}[\mathbf{X}_t \mid \mathbf{X}_{t+1}, \mathbf{y}_{1:t}] \mid \mathbf{y}_{1:T}] \\ &= P_{t|t} + L_t (P_{t+1|T} - P_{t+1|t}) L_t^T \end{aligned}$$





# Learning SSMs

- Complete log likelihood

$$\begin{aligned}\ell_c(\theta, \mathcal{D}) &= \sum_n \log p(x_n, y_n) = \sum_n \log p(x_1) + \sum_n \sum_t \log p(x_{n,t} | x_{n,t-1}) + \sum_n \sum_t \log p(y_{n,t} | x_{n,t}) \\ &= f_1(X_1; \Sigma_0) + f_2(\langle X_t X_{t-1}^T \rangle, \langle X_t X_t^T \rangle, \langle X_t \rangle : \forall t; A, Q, G) + f_3(\langle X_t X_t^T \rangle, \langle X_t \rangle : \forall t; C, R)\end{aligned}$$

- EM

- E-step: compute

$$\langle X_t X_{t-1}^T \rangle, \langle X_t X_t^T \rangle, \langle X_t \rangle \mid y_1, \dots, y_T$$

these quantities can be inferred via KF and RTS filters, etc.,

e.g.,

- M-step: MLE using  $\langle X_t X_t^T \rangle \equiv \text{var}(X_t X_t^T) + \text{E}(X_t)^2 = P_{t|T} + \hat{x}_{t|T}^2$

c.f., M-step in factor analysis

$$\langle \ell_c(\theta, \mathcal{D}) \rangle = f_1(\langle X_1 \rangle; \Sigma_0) + f_2(\langle \langle X_t X_{t-1}^T \rangle, \langle X_t X_t^T \rangle, \langle X_t \rangle : \forall t \rangle; A, Q, G) + f_3(\langle \langle X_t X_t^T \rangle, \langle X_t \rangle : \forall t \rangle; C, R)$$





# Nonlinear systems

- In robotics and other problems, the motion model and the observation model are often nonlinear:

$$x_t = f(x_{t-1}) + w_t, \quad y_t = g(x_t) + v_t$$

- An optimal closed form solution to the filtering problem is no longer possible.
- The nonlinear functions  $f$  and  $g$  are sometimes represented by neural networks (multi-layer perceptrons or radial basis function networks).
- The parameters of  $f$  and  $g$  may be learned offline using EM, where we do gradient descent (back propagation) in the M step, c.f. learning a MRF/CRF with hidden nodes.
- Or we may learn the parameters online by adding them to the state space:  $x_t^i = (x_t, \theta)$ . This makes the problem even more nonlinear.





# Extended Kalman Filter (EKF)

- The basic idea of the EKF is to linearize  $f$  and  $g$  using a second order Taylor expansion, and then apply the standard KF.
- i.e., we approximate a stationary nonlinear system with a non-stationary linear system.

$$\begin{aligned}x_t &= f(\hat{x}_{t-1|t-1}) + A_{\hat{x}_{t-1|t-1}}(x_{t-1} - \hat{x}_{t-1|t-1}) + w_t \\y_t &= g(\hat{x}_{t|t-1}) + C_{\hat{x}_{t|t-1}}(x_t - \hat{x}_{t|t-1}) + v_t\end{aligned}$$

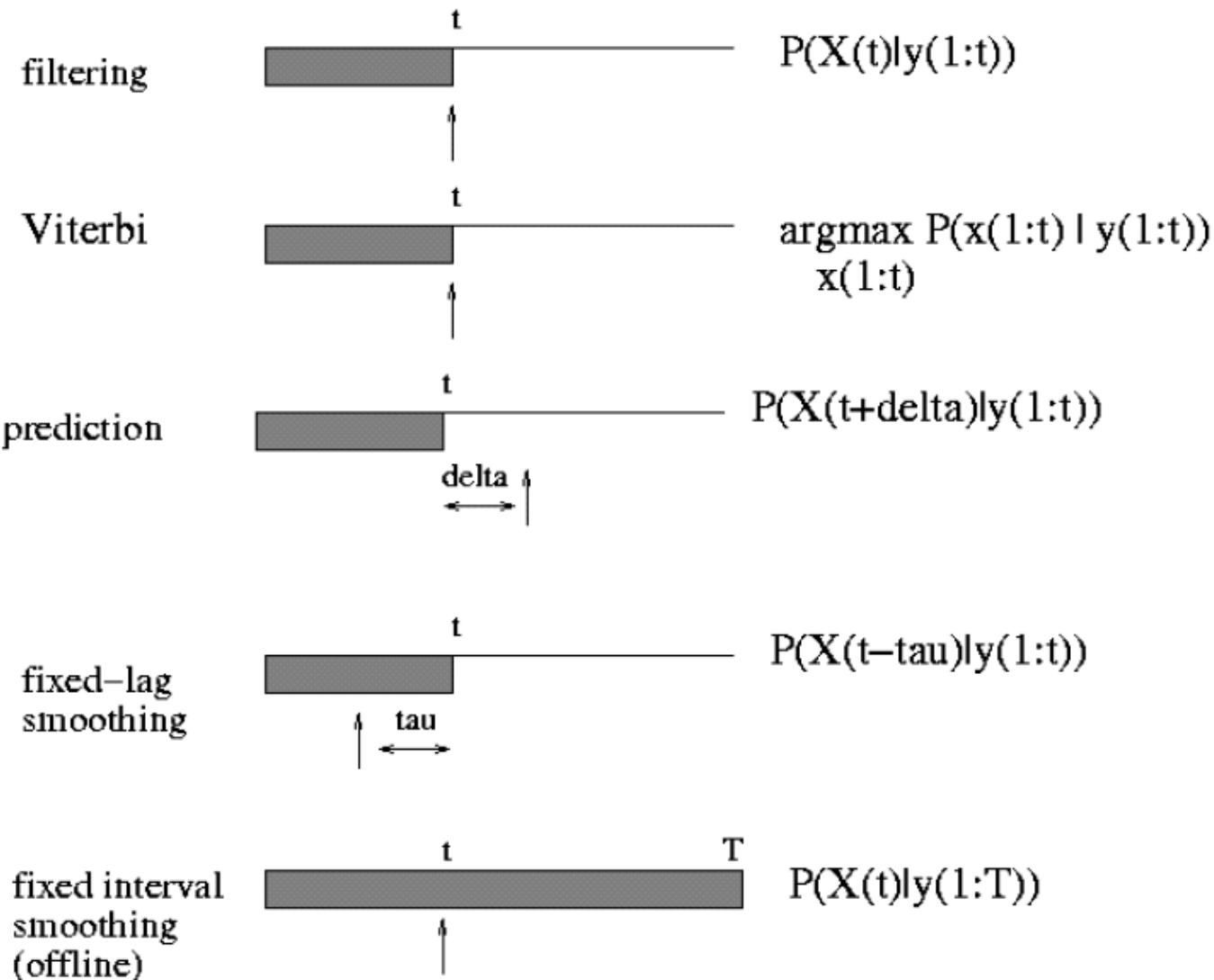
where  $\hat{x}_{t|t-1} = f(\hat{x}_{t-1|t-1})$  and  $A_{\hat{x}} = \frac{\partial f}{\partial x} \Big|_{\hat{x}}$  and  $C_{\hat{x}} = \frac{\partial g}{\partial x} \Big|_{\hat{x}}$

- The noise covariance ( $Q$  and  $R$ ) is not changed, i.e., the additional error due to linearization is not modeled.





# Online vs offline inference





# KF, RLS and LMS

- The KF update of the mean is

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

- 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 the observation model at each time slide,  $y_t = x_t^T \theta + v_t$ , is a linear regression
- We can estimate 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.

