

TS : KALMAN FILTERING.

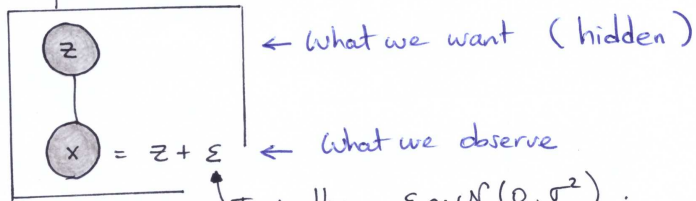
So far, we have considered i.i.d. observations.

- + : likelihood can easily be expressed as a product
- : unrealistic in many situations.

In this chapter, we depart from the i.i.d. assumption, and consider an important model for sequential data.

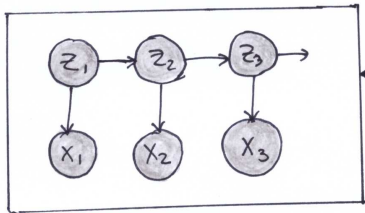
- Ex: (i) Rainfall measurements on successive days at a particular location.
 (ii) Position of a moving object.
 (iii) Daily values of the EUR-RUB exchange rate.

Usually, the quantity of interest (denoted z) is measured using a noisy sensor that returns an observation x representing z plus some noise ϵ



Typically, $\epsilon \sim \mathcal{N}(0, \sigma^2)$.

'Best' guess for z is x
 Better guess = repeat the measurement n times and average the x s.
 This is OK if z is not changing over time.



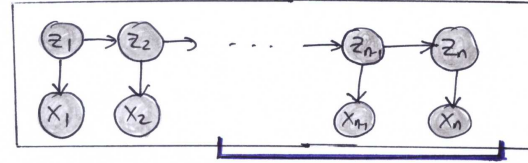
How would you estimate z_1, z_2, z_3, \dots now?

$\hat{z}_1 = x_1, \hat{z}_2 = x_2, \hat{z}_3 = x_3$?

↳ Maybe average the x s to reduce the noise.

⚠ By doing so, you introduce a bias.

Idea: Use a moving window of length L



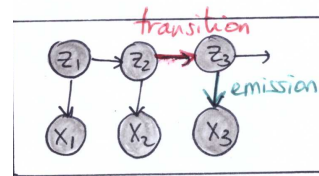
Take the L most recent observations and average them.

- ↗ If z is slowly varying + σ^2 large \Rightarrow Take a long window
- ↘ If z is quickly varying + σ^2 small \Rightarrow Take a short window

↑ Situations in between?

In any case, we need a better way to decide on the value of $L \Rightarrow$ we need a probabilistic model.

• Model # 1



- The z s are discrete & evolve according to a Markov Chain (MC).
- Link between z and x is probabilistic (arbitrary e.g. Gaussian / Binomial / ...)

$P(z_2 = j | z_1 = i) = a_{ij}$ = transition probability

$P(x_2 = k | z_2 = j) = \pi_{jk}$ = emission probability

\Rightarrow End up with Hidden Markov Models (HMM) (not studied here)

• Model # 2

- The z s are continuous & normally distributed.
- Transition probability is also Gaussian: $p(z_j | z_{j-1}) = \mathcal{N}(z_j | A z_{j-1}, \Gamma)$
- Emission probability is Gaussian: $p(x_j | z_j) = \mathcal{N}(x_j | C z_j, \Sigma)$

\Rightarrow End up with equations of Kalman Filtering.
 +: likelihood can still be easily written.

variable name \uparrow
 mean \uparrow
 cov. \uparrow

The transition + emission distributions are commonly expressed in an equivalent way in terms of linear equations: (3)

$$\begin{cases} z_1 = \mu_0 + u & , & u \sim \mathcal{N}(u | 0, \underline{\Sigma}_0) \\ z_j = A z_{j-1} + w_j & , & w_j \sim \mathcal{N}(w_j | 0, \underline{\Gamma}) \\ x_j = C z_j + v_j & , & v_j \sim \mathcal{N}(v_j | 0, \underline{\Sigma}) \end{cases}$$

Q: (i) Estimation of model parameters $\Theta = \{A, \Gamma, C, \Sigma, \mu_0, \Sigma_0\}$
 ↳ use MLE / EM algorithm

(ii) Predict z_n (and x_n) given x_1, \dots, x_{n-1} (and x_n).

We turn our attention to question (ii), and postpone (i) for later.

⇒ We are interested in the posterior distribution

$$p(z_n | x_1, \dots, x_n) =: \hat{\alpha}(z_n)$$

↖ To end up with an efficient algorithm, we should obtain $\hat{\alpha}(z_n)$ easily from $\hat{\alpha}(z_{n-1})$. These should also have the same functional form. Our Gaussian assumption will make everything work.

- $\hat{\alpha}(z_{n-1}) = p(z_{n-1} | x_1, \dots, x_{n-1})$
- $\hat{\alpha}(z_{n-1}) \underbrace{p(z_n | z_{n-1})}_{\text{transition probability}} = p(z_{n-1} | x_1, \dots, x_{n-1}) p(z_n | z_{n-1})$
 $= p(z_{n-1}, z_n | x_1, \dots, x_{n-1})$
 (conditioned on z_{n-1} , z_n is independent of x_1, \dots, x_{n-1})
- $\hat{\alpha}(z_{n-1}) \underbrace{p(z_n | z_{n-1}) p(x_n | z_n)}_{\text{emission probability}} = p(z_{n-1}, z_n, x_n | x_1, \dots, x_{n-1})$
 (conditioned on z_n , x_n is independent of $z_{n-1}, x_1, \dots, x_{n-1}$)

$$\begin{aligned} & \int \hat{\alpha}(z_{n-1}) p(z_n | z_{n-1}) p(x_n | z_n) dz_{n-1} & (4) \\ & = \int p(z_{n-1}, z_n, x_n | x_1, \dots, x_{n-1}) dz_{n-1} \\ & = p(z_n, x_n | x_1, \dots, x_{n-1}) \\ & = p(z_n | x_1, \dots, x_n) p(x_n | x_1, \dots, x_{n-1}) \\ & = \hat{\alpha}(z_n) C_n \end{aligned}$$

We obtained the recursion equation:

$$(*) \quad C_n \hat{\alpha}(z_n) = p(x_n | z_n) \int \hat{\alpha}(z_{n-1}) p(z_n | z_{n-1}) dz_{n-1} \quad n \geq 2$$

Remarks

(i) The initial z_1 is Gaussian. Since z_j is expressed as a linear combination of Gaussian variables, it is also Gaussian. Likewise, the distribution $\hat{\alpha}(z_n)$ can be seen to be Gaussian. We write

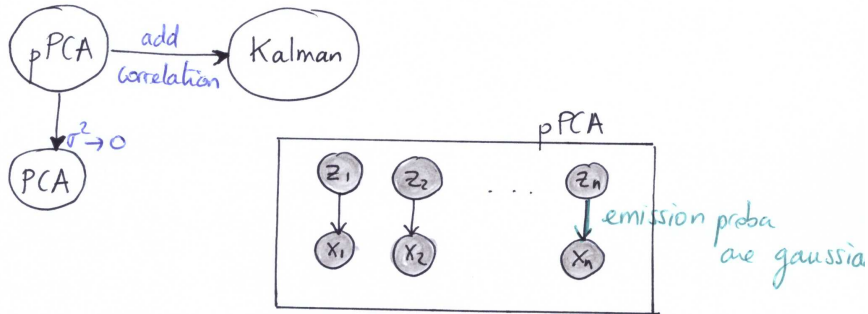
$$\hat{\alpha}(z_n) = \mathcal{N}(z_n | \mu_n, \nu_n)$$

↖ Use the recursion equation (*) to express μ_n, ν_n in terms of μ_{n-1}, ν_{n-1} and the model parameters.
 ⇒ prediction of z_n using μ_n (mean).
 C.I obtained from ν_n .

(ii) Kalman filtering \equiv extension of pPCA & Factor Analysis:
 $\{x_n, z_n\} =$ linear-Gaussian latent variable model
 ↖ latent variable

Recall: $x_n = W z_n + \mu + \epsilon_n$ (page 17
 Chp PCA)

Main difference: the z_n are no longer treated as independent variables: we allow sequential correlation in the data. (5)



OK, time to do the maths & express recurrence equations for μ_n, Σ_n .

Starting point:

$$(*) : c_n \hat{\alpha}(z_n) = \underbrace{p(x_n | z_n)}_{\mathcal{N}(x_n | Cz_n, \Sigma)} \int \underbrace{\hat{\alpha}(z_{n-1})}_{\mathcal{N}(z_{n-1} | \mu_{n-1}, \Sigma_{n-1})} \underbrace{p(z_n | z_{n-1})}_{\mathcal{N}(z_n | Az_{n-1}, \Gamma)} dz_{n-1}$$

→ Evaluation of such an integral is standard.

Toolbox: If $p(x) = \mathcal{N}(x | \mu, \Lambda^{-1})$
 $p(y|x) = \mathcal{N}(y | Ax + b, L^{-1})$

Then

$$p(y) = \mathcal{N}(y | A\mu + b, L^{-1} + A\Lambda^{-1}A^t)$$

$$p(x|y) = \mathcal{N}(x | S\{A^tL(y-b) + \Lambda\mu\}, S)$$

where

$$S = (\Lambda + A^tLA)^{-1}$$

(left as an exercise)

where $p(y) = \int p(y|x)p(x)dx$

$$\Rightarrow \int \mathcal{N}(z_n | Az_{n-1}, \Gamma) \mathcal{N}(z_{n-1} | \mu_{n-1}, \Sigma_{n-1}) dz_{n-1}$$

$$= \mathcal{N}(z_n | A\mu_{n-1}, \underbrace{\Gamma + A\Sigma_{n-1}A^t}_{\text{Call this } \Sigma_{n-1}})$$

$$\Rightarrow c_n \hat{\alpha}(z_n) = \mathcal{N}(x_n | Cz_n, \Sigma) \mathcal{N}(z_n | A\mu_{n-1}, \Sigma_{n-1})$$

$\mathcal{N}(z_n | \mu_n, \Sigma_n) = p(z_n | x_1, \dots, x_n)$

$p(x_n | x_1, \dots, x_{n-1})$

Variable z_n appears on both terms. We need to reshuffle things out to identify the RHS and LHS

$$\underbrace{p(y|x)}_{\mathcal{N}(x_n | Cz_n, \Sigma)} \underbrace{p(x)}_{\mathcal{N}(z_n | A\mu_{n-1}, \Sigma_{n-1})} = p(y|x)p(x)$$

$$= p(x, y)$$

$$= p(x|y)p(y)$$

$$= \mathcal{N}(x_n | \underbrace{CA\mu_{n-1}}_{\mu}, \underbrace{\Sigma + CP_{n-1}C^t}_{\Lambda^{-1}A^tA^t}) \underbrace{p(y)}_{\mathcal{N}(z_n | S\{C^t\Sigma^{-1}(x_n - 0) + P_{n-1}^{-1}A\mu_{n-1}\}, S)}$$

where

$$S = (P_{n-1}^{-1} + C^t\Sigma^{-1}C)^{-1}$$

$$= c_n \hat{\alpha}(z_n)$$

$$\Rightarrow c_n = \mathcal{N}(x_n | CA\mu_{n-1}, CP_{n-1}C^t + \Sigma)$$

We simplify terms in the expression of $\alpha^*(z_n | \dots, \dots)$ (7) by making use of two convenient matrix inverses identities:

this one is easy to recall, right?

$$(P^{-1} + B^t R^{-1} B)^{-1} B^t R^{-1} = P B^t (B P B^t + R)^{-1} \quad (1)$$

$$(A + B D^{-1} C)^{-1} = A^{-1} - A^{-1} B (D + C A^{-1} B)^{-1} C A^{-1} \quad (2)$$

Woodbury

(1a) • $S = (P_{n-1}^{-1} + \frac{C^t \Sigma^{-1} C}{B^t R^{-1} B})^{-1}$
 $= P_{n-1} - P_{n-1} C^t (\Sigma + C P_{n-1} C^t)^{-1} C P_{n-1}$ (2)

(2a) • $(\frac{P_{n-1}^{-1}}{P^{-1}} + \frac{C^t \Sigma^{-1} C}{B^t R^{-1} B})^{-1} C^t \Sigma^{-1} = \frac{P_{n-1} C^t (C P_{n-1} C^t + \Sigma)^{-1}}{B^t R^{-1} B}$ (1)

$$\Rightarrow S \{ C^t \Sigma^{-1} x_n + P_{n-1}^{-1} A \mu_{n-1} \}$$

$$= \underbrace{P_{n-1} C^t (\Sigma + C P_{n-1} C^t)^{-1}}_{(2a)} x_n + \underbrace{\{ P_{n-1} - P_{n-1} C^t (-)^{-1} C P_{n-1} \}}_{(1a)} \times P_{n-1}^{-1} A \mu_{n-1}$$

$$= A \mu_{n-1} + \underbrace{P_{n-1} C^t (\Sigma + C P_{n-1} C^t)^{-1}}_{K_n} (x_n - C A \mu_{n-1})$$

$$= A \mu_{n-1} + K_n (x_n - C A \mu_{n-1})$$

Our μ_n

Also, $S = P_n - K_n C P_n$
 $= (I - K_n C) P_n$

Our V_n

SUMMARY

$$\mu_n = A \mu_{n-1} + K_n (x_n - C A \mu_{n-1})$$

$$V_n = (I - K_n C) P_n$$

$$c_n = \mathcal{W}(x_n | C A \mu_{n-1}, C P_{n-1} C^t + \Sigma)$$

$$\begin{cases} K_n = P_{n-1} C^t (\Sigma + C P_{n-1} C^t)^{-1} \\ P_n = \Gamma + A V_n A^t \end{cases} \quad n \geq 2$$

(8) \Rightarrow Given μ_{n-1}, V_{n-1} and a new observation x_n , we can evaluate the distribution of z_n with mean μ_n and covariance matrix V_n . (+ normalization coefficient c_n)

• $K_n = P_{n-1} C^t (\Sigma + C P_{n-1} C^t)^{-1}$ is known as the KALMAN GAIN.

Initial conditions: $c_1 \hat{\alpha}(z_1) = p(z_1) p(x_1 | z_1)$
 $p(x_1) \quad p(z_1 | x_1) \quad \mathcal{W}(z_1 | \mu_0, P_0) \quad \mathcal{W}(x_1 | C z_1, \Sigma)$

\Rightarrow Proceed as before to obtain:

[Main difference: replace term $A \mu_{n-1}$ by μ_0 :]

INITIAL CONDITIONS

$$\mu_1 = \mu_0 + K_1 (x_1 - C \mu_0)$$

$$V_1 = (I - K_1 C) P_0$$

$$c_1 = \mathcal{W}(x_1 | C \mu_0, C P_0 C^t + \Sigma)$$

where

$$K_1 = P_0 C^t (\Sigma + C P_0 C^t)^{-1}$$

Interpretation of $\mu_n = A \mu_{n-1} + K_n (x_n - C A \mu_{n-1})$

Mean of z_{n-1} updated using the transition matrix A
 $z_n = A z_{n-1} + w_n$

Predicted observation obtained by applying the matrix C to the predicted hidden-state μ_{n-1}

$\Rightarrow x_n - C A \mu_{n-1} = \text{error}$

\Rightarrow posterior mean $\mu_n =$ predicted mean $A \mu_{n-1} +$ correction proportional to the error $(x_n - C A \mu_{n-1})$

We make this precise now

Alternative view of Kalman equations:

(8a)

- In a filtering problem, we are interested in finding the 'best' predictor of z_n given observations x_1, \dots, x_n .

↳ denote it by $\hat{z}_n = \varphi^*(x_1, \dots, x_n)$ = function of the data

'Best' is commonly measured by the square loss between the prediction and the true value:

$$\hat{z}_n = \varphi^*(x_1, \dots, x_n) = \underset{\varphi}{\operatorname{argmin}} \mathbf{E} (z_n - \varphi(x_1, \dots, x_n))^2$$

↑ we know for a long time now that φ^* is given by the CE

$$\varphi^*(x_1, \dots, x_n) = \mathbf{E}(z_n | x_1, \dots, x_n).$$

← page 7

↳ We derived that $p(z_n | x_1, \dots, x_n) = \mathcal{N}(z_n | \mu_n, V_n)$, **BEST**

so that our best estimate of z_n is $\hat{z}_n = \mathbf{E}(z_n | x_1, \dots, x_n) = \mu_n$.

↳ Moreover, the quality of this estimate is provided by the error covariance matrix,

$$\begin{aligned} & \mathbf{E}\{(z_n - \hat{z}_n)(z_n - \hat{z}_n)^t | x_1, \dots, x_n\} \\ &= \mathbf{E}\{(z_n - \mathbf{E}(z_n | x_1, \dots, x_n))(z_n - \mathbf{E}(z_n | x_1, \dots, x_n))^t | x_1, \dots, x_n\} \\ &= V_n. \end{aligned}$$

- We may now interpret equations derived at the bottom of page 7:

↳ Denote by $\hat{z}_{n|n}$ the best estimate at time n , given observations up to and including time n .

↳ Denote by $P_{n|n}$ the error covariance matrix of $\hat{z}_{n|n}$.

- A time $(n-1)$, the best estimate of z_{n-1} given observations up to and including time $(n-1)$ is $\mu_{n-1} = \hat{z}_{n-1|n-1}$.

⇒ Predicted (a priori) latent estimate is $\hat{z}_{n|n-1} = A \hat{z}_{n-1|n-1}$

↑ why 'a priori'?

Because we are looking for an estimate of z_n based solely on observations x_1, \dots, x_{n-1} : x_n is not observed.

⇒ Also, the predicted (a priori) covariance estimate is the covariance matrix of $A \hat{z}_{n-1|n-1} + w_n$, which is $(A P_{n-1|n-1} A^t + \Gamma) =: P_{n|n-1}$

← $\mathcal{N}(0, \Gamma)$

Remark: Expression $\hat{z}_{n|n-1} = A \hat{z}_{n-1|n-1}$ makes sense intuitively: the best we can, given we have not observed x_n , is to make use of the equation $z_n = A z_{n-1} + w_n$ (our model), and to set the noise term to 0.

But how would you prove this formally?

Well, the 'best' estimate of z_n , given x_1, \dots, x_{n-1} is given by the conditional mean $\mathbf{E}(z_n | x_1, \dots, x_{n-1})$, and the predicted covariance matrix is given by

$$\mathbf{E}\{(z_n - \mathbf{E}(z_n | x_1, \dots, x_{n-1}))(z_n - \mathbf{E}(z_n | x_1, \dots, x_{n-1}))^t | x_1, \dots, x_{n-1}\}$$

⇒ We need to derive $p(z_n | x_1, \dots, x_{n-1})$:

$$\begin{aligned} p(z_n | x_1, \dots, x_{n-1}) &= \int p(z_n, z_{n-1} | x_1, \dots, x_{n-1}) dz_{n-1} \\ &= \int p(z_n | z_{n-1}) p(z_{n-1} | x_1, \dots, x_{n-1}) dz_{n-1} \\ &= \int \mathcal{N}(z_n | A z_{n-1}, \Gamma) \mathcal{N}(z_{n-1} | \mu_{n-1}, V_{n-1}) dz_{n-1} \end{aligned}$$

$\Rightarrow p(z_n | x_{1:n-1}, x_{n-1}) = \mathcal{N}(z_n | A \hat{z}_{n-1|n-1}, \Gamma + A V_{n-1} A^t)$ (8c)
 So that indeed $\hat{z}_{n|n-1} = A \hat{z}_{n-1|n-1}$
 $\bullet \text{Cov} = P_{n|n-1} = \Gamma + A P_{n-1|n-1} A^t$

But then, what would be your 'best' estimate of z_{n+1} , given x_1, \dots, x_{n-1} ?

$\Rightarrow \text{Compute } p(z_{n+1} | x_{1:n}, x_{n+1}) = \iint p(z_{n+1} | z_n) p(z_n | z_{n-1}) p(z_{n-1} | x_{1:n-1}, x_{n-1}) dz_{n-1} dz_n$

where

$$\int p(z_{n+1} | z_n) p(z_n | z_{n-1}) dz_n = \int \mathcal{N}(z_{n+1} | A z_n, \Gamma) \mathcal{N}(z_n | A z_{n-1}, \Gamma) dz_n$$

$$= \mathcal{N}(z_{n+1} | A^2 z_{n-1}, \Gamma + A \Gamma A^t)$$

$\Rightarrow p(z_{n+1} | x_{1:n}, x_{n+1}) = \int \mathcal{N}(z_{n+1} | A^2 \hat{z}_{n-1|n-1}, \Gamma + A \Gamma A^t) \mathcal{N}(z_{n-1} | \mu_{n-1}, V_{n-1}) dz_{n-1}$

$$= \mathcal{N}(z_{n+1} | A^2 \hat{z}_{n-1|n-1}, \Gamma + A \Gamma A^t + A V_{n-1} A^t)$$

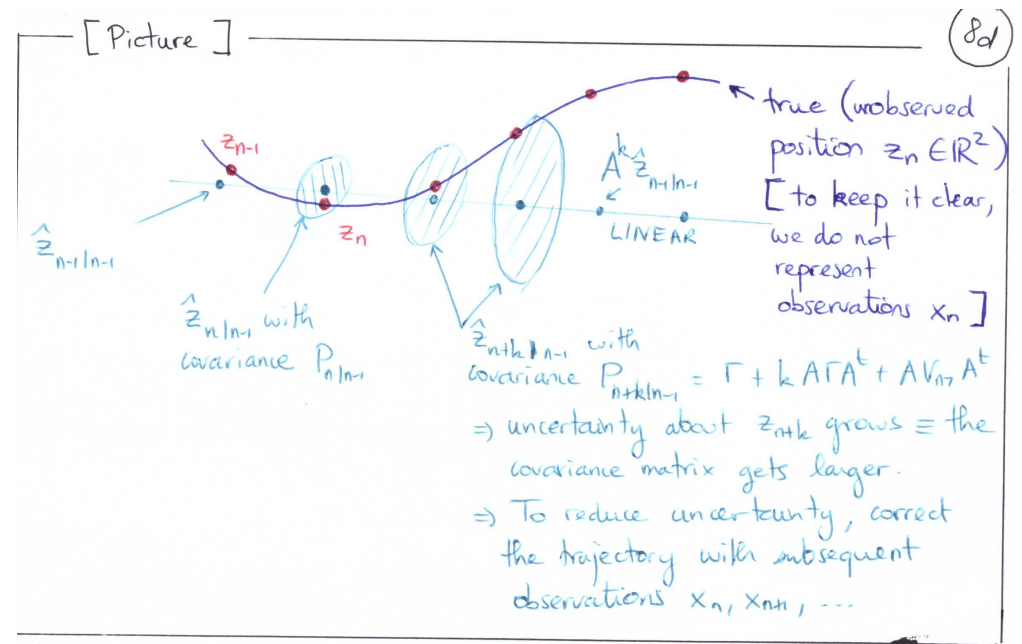
Thus,

- $\hat{z}_{n+1|n-1} = A^2 \hat{z}_{n-1|n-1}$ = propagate $\hat{z}_{n-1|n-1}$ linearly
- $P_{n+1|n-1} = \Gamma + \boxed{A \Gamma A^t} + A V_{n-1} A^t$
uncertainty increases!

More generally, you can convince yourself that

$$p(z_{n+k} | x_{1:n}, x_{n+1}) = \mathcal{N}(z_{n+k} | A^{k+1} \hat{z}_{n-1|n-1}, \Gamma + k A \Gamma A^t + A V_{n-1} A^t)$$

$\Rightarrow \hat{z}_{n+k|n-1} = A^{k+1} \hat{z}_{n-1|n-1}$ k ≥ 0
 $P_{n+k|n-1} = \Gamma + \boxed{k A \Gamma A^t} + A V_{n-1} A^t$
 increases with k



- Update prediction, given you observe x_n :
 - \Rightarrow The residual error (aka INNOVATION) is

$$i_n = x_n - C \hat{z}_{n|n-1} = x_n - C A \hat{z}_{n-1|n-1}$$

↑
Use evolution equation, $x_n = C z_n + v_n$.
 - \Rightarrow The residual covariance is $C P_{n|n-1} C^t + \Sigma =: S_n$

↑
The predicted a priori covariance matrix
 - \Rightarrow The optimal Kalman gain is $K_n = P_{n|n-1} C^t S_n^{-1}$.

⇒ The updated (a posteriori) latent estimate:

(8e)

$$\hat{z}_{n|n} = \hat{z}_{n|n-1} + K_n z_n$$

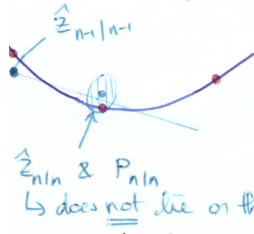
$\hat{z}_{n|n-1}$ is the *a priori estimate*
 K_n is the *optimal 'weight'*
 z_n is the *error*
 $=$ *correction to your linear prediction.*

Compare with the expression of μ_n page 7

⇒ The updated (a posteriori) covariance estimate is:

$$P_{n|n} = (I - K_n C) P_{n|n-1}$$

$P_{n|n} \ll P_{n|n-1}$
 We are indeed reducing uncertainty since $P_{n|n} \prec P_{n|n-1}$.
 lock: $K_n C P_{n|n-1}$
 $= P_{n|n-1} C^T S_n^{-1} C P_{n|n-1}$ is positive definite.



- A closely related problem is that of SMOOTHING (the FILTERING problem aims at recovering at time n some information about z_n given x_1, \dots, x_n). In the smoothing problem, measurements derived later than time n can be used in obtaining information about z_n .

Ex: the way the human brain tackles the problem of reading written handwriting: words are read sequentially. When a word is difficult to read, words before and after it may be used to attempt to deduce the word.

Formally, we are interested in the 'best' estimate of z_j given x_1, \dots, x_n , $1 \leq j \leq n$.

↳ Under a square loss, the best estimate is given by the conditional mean of z_j given x_1, \dots, x_n ; and it is therefore of interest to derive $p(z_j | x_1, \dots, x_n)$.

↳ Section II.1 p.10

Remarks.

(9)

(i) Back to our discussion on the top of page 2: Suppose that the measurement noise σ^2 is small compared to the rate at which the latent variable is changing; i.e. suppose $\Sigma \approx 0$. Take $C = I$, so that observation x corresponds to z + something small.

Bottom of page 7: $\mu_n = A \mu_{n-1} + K_n (x_n - A \mu_{n-1})$
 $K_n = P_{n-1} (\Sigma + P_{n-1})^{-1} = I$
 $\Rightarrow \mu_n = x_n$

⇒ Predict z_n using x_n , in agreement with our intuition.

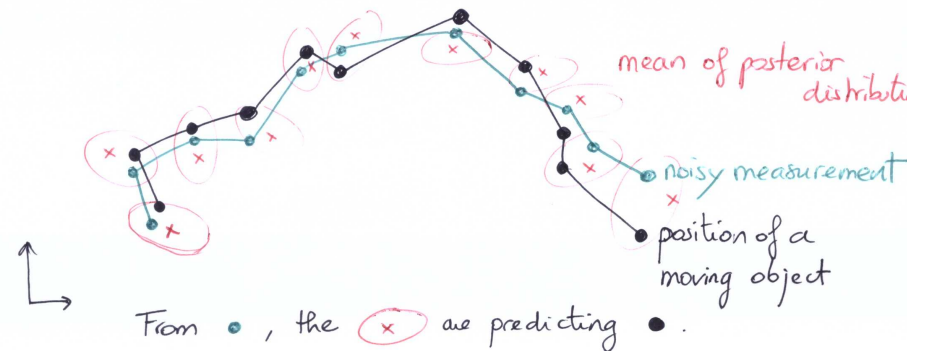
(ii) What if the latent variable is evolving slowly relative to the observation noise level?

In the extreme case, we may assume

- $A = C = I$
 - $\Gamma = 0$
 - $P_0 \rightarrow \infty$
- $\left. \begin{array}{l} \text{ } \\ \text{ } \end{array} \right\} z_j \text{ stays constant over time}$
 (initial distrib. unimportant)

Then it is possible to show that the posterior mean for z_n is determined by the average of the x_1, \dots, x_n .

(iii) Application: tracking a moving object.



I. LEARNING THE MODEL PARAMETERS

(10)

In the previous section we derived the equations of Kalman filter, assuming the parameters are known. In this section, we derive a procedure for estimating $\Theta = \{A, \Gamma, C, \Sigma, \mu_0, \Sigma_0\}$.
 \Rightarrow We do this using ML + EM algorithm (due to the presence of hidden variables).

• First, we need preliminary results:

(i) \rightarrow Compute $p(z_j | x_1, \dots, x_n)$, $1 \leq j \leq n$

(ii) \rightarrow Compute $p(z_{j-1}, z_j | x_1, \dots, x_n)$, $2 \leq j \leq n$

We address (i) + (ii) in the next section.

II.1. Preliminary results.

[To simplify notation, we write $\underline{x} = (x_1, \dots, x_n)$; so that $p(z_j | x_1, \dots, x_n)$ is rewritten $p(z_j | \underline{x})$.]

\rightarrow We are interested in calculating (efficiently) the posterior probability $p(z_j | \underline{x}) =: \gamma(z_j)$, for $1 \leq j \leq n$.

Note that for $j = n$, we have $\gamma(z_n) = \hat{\alpha}(z_n)$ (page 3)

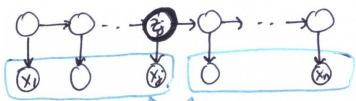
$$\text{Bayes} \Rightarrow \gamma(z_j) = \frac{p(\underline{x} | z_j) p(z_j)}{p(\underline{x})}, \quad j \leq n-1$$

Conditional independence property (left as an exercise)

$$= \frac{p(x_1, \dots, x_j | z_j) p(x_{j+1}, \dots, x_n | z_j) p(z_j)}{p(\underline{x})}$$

$$= \frac{p(x_1, \dots, x_j, z_j) p(x_{j+1}, \dots, x_n | z_j)}{p(\underline{x})}$$

$$=: \frac{\alpha(z_j) \beta(z_j)}{p(\underline{x})}$$



independent, conditionally on z_j^- .

Where we have defined

(11)

$$\alpha(z_j) = p(x_1, \dots, x_j, z_j)$$

$$\beta(z_j) = p(x_{j+1}, \dots, x_n | z_j)$$

Remark: we usually work with scaled versions of $\alpha(z_j)$ and $\beta(z_j)$, for numerical issues. Specifically,

$$\hat{\alpha}(z_j) = p(z_j | x_1, \dots, x_j) = \frac{\alpha(z_j)}{p(x_1, \dots, x_j)}$$

\leftarrow introduced on page 3 already.

We introduce as well

$$c_j = p(x_j | x_1, \dots, x_{j-1}) \quad [\text{defined page 4}]$$

so that

$$p(x_1, \dots, x_j) = p(x_j | x_1, \dots, x_{j-1}) p(x_{j-1} | x_1, \dots, x_{j-2}) \dots p(x_2 | x_1) p(x_1)$$

$$= \prod_{m=1}^j c_m$$

so that

$$\alpha(z_j) = \left(\prod_{m=1}^j c_m \right) \hat{\alpha}(z_j)$$

• Likewise, we defined the scaled variables $\hat{\beta}(z_j)$ as:

$$\beta(z_j) = \left(\prod_{m=j+1}^n c_m \right) \hat{\beta}(z_j)$$

$$\text{Now, } \left(\prod_{m=1}^n c_m \right) = p(\underline{x}) = \left(\prod_{m=1}^j c_m \right) \left(\prod_{m=j+1}^n c_m \right)$$

$$\Rightarrow \prod_{m=j+1}^n c_m = \frac{p(\underline{x})}{p(x_1, \dots, x_j)} = p(x_{j+1}, \dots, x_n | x_1, \dots, x_j)$$

$$\text{so that } \hat{\beta}(z_j) = \frac{\beta(z_j)}{p(x_{j+1}, \dots, x_n | x_1, \dots, x_j)}$$

Recursion equations for $\alpha / \hat{\alpha} / \beta / \hat{\beta}$. (12)

→ We already established on page 4 that

scaled version $\left[c_j \hat{\alpha}(z_j) = p(x_j | z_j) \int \hat{\alpha}(z_{j+1}) p(z_j | z_{j+1}) dz_{j+1} \right] \text{--- (A)}$

Equivalently,

$$c_j \frac{\alpha(z_j)}{\left(\prod_{m=1}^j C_m\right)} = p(x_j | z_j) \int \frac{\alpha(z_{j+1})}{\left(\prod_{m=1}^{j+1} C_m\right)} p(z_j | z_{j+1}) dz_{j+1}$$

unscaled version $\left[\alpha(z_j) = p(x_j | z_j) \int \alpha(z_{j+1}) p(z_j | z_{j+1}) dz_{j+1} \right] \text{--- (B)}$
 ↑ we have a FORWARD message passing: from z_{j+1} to z_j .

→ We establish a similar recursion equation for $\beta / \hat{\beta}$:

$$\begin{aligned} \beta(z_j) &= p(x_{j+1}, \dots, x_n | z_j) \\ &= \int p(x_{j+1}, \dots, x_n, z_{j+1} | z_j) dz_{j+1} \quad (\text{cond. independence}) \\ &= \int p(x_{j+1}, \dots, x_n | z_{j+1}, z_j) p(z_{j+1} | z_j) dz_{j+1} \\ &= \int p(x_{j+2}, \dots, x_n | z_{j+1}, z_j) p(x_{j+1} | z_{j+1}) p(z_{j+1} | z_j) dz_{j+1} \end{aligned}$$

unscaled version $\left[\beta(z_j) = \int \beta(z_{j+1}) p(x_{j+1} | z_{j+1}) p(z_{j+1} | z_j) dz_{j+1} \right] \text{--- (C)}$

Equivalently,

$$\left(\prod_{m=j+1}^n C_m\right) \hat{\beta}(z_j) = \int \left(\prod_{m=j+2}^n C_m\right) \hat{\beta}(z_{j+1}) p(x_{j+1} | z_{j+1}) p(z_{j+1} | z_j) dz_{j+1}$$

scaled version $\left[c_{j+1} \hat{\beta}(z_j) = \int \hat{\beta}(z_{j+1}) p(x_{j+1} | z_{j+1}) p(z_{j+1} | z_j) dz_{j+1} \right] \text{--- (D)}$
 ↑ We have a BACKWARD message passing: from z_{j+1} to z_j .

For these reasons, $\alpha / \hat{\alpha}$ are referred to as the FORWARD variables (13)

$\beta / \hat{\beta}$ — " — the BACKWARD variables.

(You will meet these again when discussing HMM)

Recall from page 10,

$$\begin{aligned} p(z_j | x_1, \dots, x_n) &= \gamma(z_j) = \frac{\alpha(z_j) \beta(z_j)}{p(x_1, \dots, x_n)} \\ &= \frac{\alpha(z_j) \beta(z_j)}{\left(\prod_{m=1}^n C_m\right)} \\ &= \frac{\alpha(z_j)}{\left(\prod_{m=1}^j C_m\right)} \frac{\beta(z_j)}{\left(\prod_{m=j+1}^n C_m\right)} \\ &= \hat{\alpha}(z_j) \hat{\beta}(z_j). \end{aligned}$$

$$\Rightarrow \gamma(z_j) = \frac{\alpha(z_j) \beta(z_j)}{p(x_1, \dots, x_n)} = \hat{\alpha}(z_j) \hat{\beta}(z_j)$$

This representation, together with the recurrence equations satisfied by $\hat{\alpha} / \hat{\beta}$, will allow us to derive the distribution of $z_j | x_1, \dots, x_n$.

Starting point: recurrence equation (D)

$$c_{j+1} \hat{\beta}(z_j) = \int \hat{\beta}(z_{j+1}) p(x_{j+1} | z_{j+1}) p(z_{j+1} | z_j) dz_{j+1} \quad \times \hat{\alpha}(z_j)$$

$$c_{j+1} \hat{\alpha}(z_j) \hat{\beta}(z_j) = \int \hat{\alpha}(z_j) \hat{\beta}(z_{j+1}) p(x_{j+1} | z_{j+1}) p(z_{j+1} | z_j) dz_{j+1}$$

$$\gamma(z_j) = \text{Gaussian} = \mathcal{N}(z_j | \hat{\mu}_j, \hat{\Sigma}_j)$$

We are now looking for recurrence relations on $\hat{\beta}_j / \hat{V}_j$. (14)

$$c_{j+1} \mathcal{N}(z_j | \hat{\beta}_j, \hat{V}_j) = \int \hat{\beta}(z_{j+1}) \hat{\alpha}(z_j) p(x_{j+1} | z_{j+1}) p(z_{j+1} | z_j) dz_{j+1}$$

where

$$\begin{aligned} \hat{\alpha}(z_j) p(z_{j+1} | z_j) &= \mathcal{N}(z_j | \mu_j, V_j) \mathcal{N}(z_{j+1} | A z_j, \Gamma) \\ &\stackrel{\text{notation from bottom of page 5}}{=} p(x) p(y | x) \\ &= p(x | y) p(y) \\ &= \mathcal{N}(z_j | \underbrace{M_j}_{S} (A^t \Gamma^{-1} z_{j+1} + \underbrace{V_j^{-1} \mu_j}_{P_j}), \underbrace{M_j}_{S}) \mathcal{N}(z_{j+1} | \underbrace{A \mu_j}_{Y}, \underbrace{P_j}_{L^t + A A^t A}) \end{aligned}$$

where $M_j = (V_j^{-1} + A^t \Gamma^{-1} A)^{-1}$

Put $m_j = M_j (A^t \Gamma^{-1} z_{j+1} + V_j^{-1} \mu_j)$

$$\Rightarrow \hat{\alpha}(z_j) p(z_{j+1} | z_j) = \mathcal{N}(z_j | m_j, M_j) \mathcal{N}(z_{j+1} | A \mu_j, P_j)$$

Now, $M_j = (V_j^{-1} + A^t \Gamma^{-1} A)^{-1}$

$$\begin{aligned} &= V_j - V_j A^t (\Gamma + A V_j A^t)^{-1} A V_j \quad \text{Relation (2) page 7} \\ &= V_j - V_j A^t P_j^{-1} A V_j \quad (\diamond) \\ &= (I - \underbrace{V_j A^t P_j^{-1}}_{= J_j} A) V_j \\ &= (I - J_j A) V_j \end{aligned}$$

$$\Rightarrow \begin{cases} M_j = (I - J_j A) V_j \\ J_j = V_j A^t P_j^{-1} \end{cases}$$

$$\Rightarrow c_{j+1} \mathcal{N}(z_j | \hat{\beta}_j, \hat{V}_j) = \int \hat{\beta}(z_{j+1}) \underbrace{p(x_{j+1} | z_{j+1}) \mathcal{N}(z_{j+1} | A \mu_j, P_j)}_{c_{j+1} \hat{\alpha}(z_{j+1})} \times \mathcal{N}(z_j | m_j, M_j) dz_{j+1} \quad (15)$$

(third line on page 6)

$$\begin{aligned} &= c_{j+1} \int \gamma(z_{j+1}) \mathcal{N}(z_j | m_j, M_j) dz_{j+1} \\ &\quad \uparrow \text{(depends on } z_{j+1}) \\ &= c_{j+1} \int \mathcal{N}(z_{j+1} | \hat{\beta}_{j+1}, \hat{V}_{j+1}) \mathcal{N}(z_j | m_j, M_j) dz_{j+1} \\ &\quad \stackrel{\text{notation bottom of page 5}}{=} p(x) p(y | x) \\ &\quad \quad m_j = M_j A^t \Gamma^{-1} z_{j+1} + M_j V_j^{-1} \mu_j \\ &= \int p(x) p(y | x) dx \\ &= p(y) \\ &= \mathcal{N}(z_j | \underbrace{M_j A^t \Gamma^{-1} \hat{\beta}_{j+1}}_A + \underbrace{M_j V_j^{-1} \mu_j}_P, \underbrace{M_j + M_j A^t \Gamma^{-1} \hat{V}_{j+1} \Gamma^{-1} A M_j}_{L^{-1} A^{-1} A^t}) \\ &= c_{j+1} \mathcal{N}(z_j | M_j (A^t \Gamma^{-1} \hat{\beta}_{j+1} + V_j^{-1} \mu_j), M_j + M_j A^t \Gamma^{-1} \hat{V}_{j+1} \Gamma^{-1} A M_j) \\ &\quad (M_j = \text{symm.}) \end{aligned}$$

$$\Rightarrow \begin{cases} \hat{\beta}_j = M_j (A^t \Gamma^{-1} \hat{\beta}_{j+1} + V_j^{-1} \mu_j) \\ \hat{V}_j = M_j + M_j A^t \Gamma^{-1} \hat{V}_{j+1} \Gamma^{-1} A M_j \end{cases}$$

Let's simplify these expressions.

Note that $M_j A^t \Gamma^{-1} = \underbrace{(V_j - V_j A^t P_j^{-1} A V_j)}_{= M_j \text{ (relation } \blacklozenge \text{ page 14)}} A^t \Gamma^{-1}$ (16)

$$= V_j A^t (\mathbf{I} - \underbrace{P_j^{-1} A V_j A^t}_{= P_j^{-1} \Gamma \text{ by definition of } P_j, \text{ page 7}}) \Gamma^{-1}$$

$$= V_j A^t (P_j^{-1} \Gamma) \Gamma^{-1}$$

$$= V_j A^t P_j^{-1}$$

$$= J_j \text{ (defined bottom of page 14)}$$

$$\Rightarrow \bullet \hat{\mu}_j^- = M_j (A^t \Gamma^{-1} \hat{\mu}_{j+1} + V_j^{-1} \mu_j)$$

$$= J_j \hat{\mu}_{j+1} + M_j V_j^{-1} \mu_j$$

$$= J_j \hat{\mu}_{j+1} + \underbrace{(\mathbf{I} - J_j A) V_j}_{\text{bottom of page 14}} V_j^{-1} \mu_j$$

$$\boxed{\hat{\mu}_j^- = \mu_j + J_j (\hat{\mu}_{j+1} - A \mu_j)}$$

$$\bullet \hat{V}_j^- = \underbrace{M_j}_{= V_j - V_j A^t P_j^{-1} A V_j} + \underbrace{M_j A^t \Gamma^{-1} \hat{V}_{j+1} \Gamma^{-1} A M_j}_{= J_j \hat{V}_{j+1} J_j^t}$$

$$= V_j + J_j \hat{V}_{j+1} J_j^t - \underbrace{J_j A V_j}_{\substack{J_j = V_j A^t P_j^{-1} \\ J_j P_j = V_j A^t \\ P_j^t J_j^t = A V_j^t \\ P_j J_j^t = A V_j^t \rightarrow V_j, P_j \text{ symm}}}$$

$$= V_j + J_j \hat{V}_{j+1} J_j^t - J_j P_j J_j^t$$

$$\boxed{\hat{V}_j^- = V_j + J_j (\hat{V}_{j+1} - P_j) J_j^t}$$

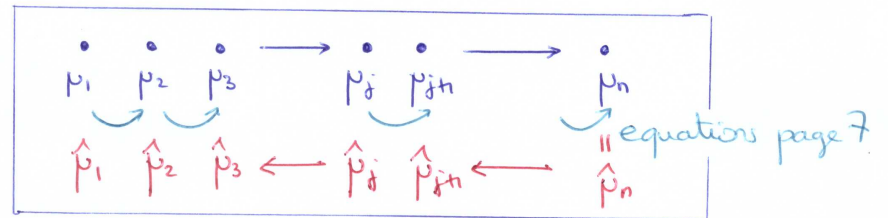
Summary: $\gamma(z_j) = p(x_j | x_1, \dots, x_n)$ (17)

$$= \mathcal{N}(z_j | \hat{\mu}_j^-, \hat{V}_j^-),$$

where $\bullet \hat{\mu}_j^- = \mu_j + J_j (\hat{\mu}_{j+1} - A \mu_j)$

$$\bullet \hat{V}_j^- = V_j + J_j (\hat{V}_{j+1} - P_j) J_j^t, \quad J_j = V_j A^t P_j^{-1}$$

To compute $\hat{\mu}_j^-$ and \hat{V}_j^- , we need the value of μ_j and V_j
 \Rightarrow A FORWARD PASS must be completed first.



\Rightarrow Then, using the 'boundary conditions' $\hat{\mu}_n = \mu_n$ and $\hat{V}_n = V_n$, a BACKWARD PASS yields $\hat{\mu}_j^-$ and \hat{V}_j^- , $j = n-1, \dots, 1$.

Important: Meaning of update $\hat{\mu}_j^- = \mu_j + J_j (\hat{\mu}_{j+1} - A \mu_j)$.
 Well,

$\bullet \mu_j^-$ = our best guess of z_j , computed from x_1, \dots, x_j , and independently of x_{j+1}, \dots, x_n .

\Rightarrow To get our 'best' guess of z_j given the full dataset, one must somehow correct the estimate μ_j^- , using the additional information contained in x_{j+1}, \dots, x_n :
 $\hat{\mu}_j^- = \mu_j^- + \text{some correction}$

- If you only observe x_1, \dots, x_j , you would predict z_{j+1} using $A\mu_j$ (see page 8b)
- But you know how to predict z_{j+1} more efficiently, using the full dataset - it is given by $\hat{\mu}_{j+1}$.

Wisdom: learn from your errors.
Here, learn from $\hat{\mu}_{j+1} - A\mu_j$.

⇒ $\hat{\mu}_j = \mu_j + \text{something proportional to } \hat{\mu}_{j+1} - A\mu_j$

And this 'something' is exactly given by J_j^{-1} , we computed it. Good.

⇒ Similarly for the covariance matrix: update V_j using a term involving $\hat{V}_{j+1} - P_j = \text{error}$.

→ The next quantity of interest is $p(z_{j-1}, z_j | x_1, \dots, x_n) =: \xi(z_{j-1}, z_j)$

We have

$$\xi(z_{j-1}, z_j) = \frac{p(x | z_{j-1}, z_j) p(z_{j-1}, z_j)}{p(x)}$$

$$= \frac{p(x_1, \dots, x_{j-1} | z_{j-1}) p(x_j | z_j) p(x_{j+1}, \dots, x_n | z_j) p(z_j | z_{j-1}) p(z_j)}{p(x)}$$

unscaled version

$$= \frac{\alpha(z_{j-1}) p(x_j | z_j) p(z_j | z_{j-1}) \beta(z_j)}{p(x)}$$

scaled version

$$= \bar{c}_j^{-1} \hat{\alpha}(z_{j-1}) p(x_j | z_j) p(z_j | z_{j-1}) \hat{\beta}(z_j)$$

α/β were introduced on page 11

(7a)

$$\xi(z_{j-1}, z_j) = \frac{\mathcal{N}(z_{j-1} | \mu_{j-1}, V_{j-1}) p(x_j | z_j) p(z_j | z_{j-1}) \gamma(z_j)}{c_j \hat{\alpha}(z_j)}$$

$\hat{\alpha}(z_j) = p(x_j | z_j) \mathcal{N}(z_j | A\mu_{j-1}, P_{j-1})$ (page 5/6)

$$= \frac{\mathcal{N}(z_{j-1} | \mu_{j-1}, V_{j-1}) \mathcal{N}(z_j | A z_{j-1}, \Gamma) \mathcal{N}(z_j | \hat{\mu}_j, \hat{V}_j)}{\mathcal{N}(z_j | A\mu_{j-1}, P_{j-1})}$$

Same calculation as on page 14, with j replaced by $j-1$.

$$= \mathcal{N}(z_{j-1} | m_{j-1}, M_{j-1}) \mathcal{N}(z_j | A\mu_{j-1}, P_{j-1})$$

$$= \frac{\mathcal{N}(z_{j-1} | m_{j-1}, M_{j-1}) \mathcal{N}(z_j | A\mu_{j-1}, P_{j-1}) \mathcal{N}(z_j | \hat{\mu}_j, \hat{V}_j)}{\mathcal{N}(z_j | A\mu_{j-1}, P_{j-1})}$$

⇒ $\xi(z_{j-1}, z_j) = \mathcal{N}(z_{j-1} | m_{j-1}, M_{j-1}) \mathcal{N}(z_j | \hat{\mu}_j, \hat{V}_j)$

From this expression, we may deduce the mean of the vector $\begin{pmatrix} z_{j-1} \\ z_j \end{pmatrix}$, and $\text{cov}(z_{j-1}, z_j)$:

Toolbox: If $p(x) = \mathcal{N}(x | \mu, \Lambda^{-1})$
 $p(y|x) = \mathcal{N}(y | Ax + b, L^{-1})$

Then $(x \ y)^t$ is gaussian with

$$\text{cov} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} \Lambda^{-1} & \Lambda^{-1} A^t \\ A \Lambda^{-1} & L^{-1} + A \Lambda^{-1} A^t \end{pmatrix}$$

$$E \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} \mu \\ A\mu + b \end{pmatrix}$$

(18)

With the notation at the bottom of page 18,

$$\begin{aligned} X &= z_j & Y &= z_{j-1} \\ \mu &= \hat{\mu}_j & A &= M_{j-1}^t A^t \Gamma^{-1} \\ \Lambda &= \hat{\Lambda}_j & b &= M_{j-1}^t V_{j-1}^{-1} \mu_{j-1} \\ & & L^{-1} &= M_{j-1} \end{aligned}$$

$$\Rightarrow \mathbf{E} \begin{pmatrix} z_j \\ z_{j-1} \end{pmatrix} = \begin{pmatrix} \hat{\mu}_j \\ \underbrace{M_{j-1}^t A^t \Gamma^{-1} \hat{\mu}_j + M_{j-1}^t V_{j-1}^{-1} \mu_{j-1}}_b \end{pmatrix} = \begin{pmatrix} \hat{\mu}_j \\ \hat{\mu}_{j-1} \end{pmatrix}$$

(recurrence equation, bottom of page 15)

$$\begin{aligned} \bullet \text{cov}(z_j, z_{j-1}) &= \hat{\Lambda}_j^{-1} \frac{A^t}{\Gamma^{-1} A M_{j-1}} = \hat{\Lambda}_j^{-1} J_{j-1}^t \\ &= J_{j-1}^t \text{ from top of page 16.} \end{aligned}$$

Summary: $p(z_{j-1}, z_j | x_1, \dots, x_n) = \xi(z_{j-1}, z_j)$ is a gaussian with

- $\mathbf{E} \begin{pmatrix} z_{j-1} \\ z_j \end{pmatrix} = \begin{pmatrix} \hat{\mu}_{j-1} \\ \hat{\mu}_j \end{pmatrix}$
- $\text{cov}(z_j, z_{j-1}) = \hat{\Lambda}_j^{-1} J_{j-1}^t$
equivalently,
 $\mathbf{E} \begin{pmatrix} z_j \\ z_{j-1} \end{pmatrix} = \hat{\Lambda}_j^{-1} J_{j-1}^t + \begin{pmatrix} \hat{\mu}_j \\ \hat{\mu}_{j-1} \end{pmatrix}$
- $\text{cov}(z_j, z_j) = \hat{\Lambda}_j$
or
 $\mathbf{E} \begin{pmatrix} z_j \\ z_j \end{pmatrix} = \hat{\Lambda}_j + \begin{pmatrix} \hat{\mu}_j \\ \hat{\mu}_j \end{pmatrix}$.

II.2. EM algorithm.

Step I: Complete log-likelihood.

$$\begin{aligned} \mathcal{L}_c &= \log p(x_1, \dots, x_n, z_1, \dots, z_n) \\ &= \log \left\{ \underbrace{p(z_1)}_{\mathcal{N}(z_1; \mu_0, \Sigma_0)} \prod_{j=2}^n \underbrace{p(z_j | z_{j-1})}_{\mathcal{N}(z_j; A z_{j-1}, \Gamma)} \prod_{j=1}^n \underbrace{p(x_j | z_j)}_{\mathcal{N}(x_j; C z_j, \Sigma)} \right\} \\ &= \log \left\{ p(z_1; \mu_0, \Sigma_0) \right\} + \sum_{j=2}^n \log p(z_j | z_{j-1}; A, \Gamma) \\ &\quad + \sum_{j=1}^n \log p(x_j | z_j; C, \Sigma) \\ &= -\frac{1}{2} \log |\Sigma_0| - \frac{1}{2} (z_1 - \mu_0)^t \Sigma_0^{-1} (z_1 - \mu_0) \end{aligned}$$

we omit terms independent of the parameters)

$$\begin{aligned} &- \frac{n-1}{2} \log |\Gamma| - \frac{1}{2} \sum_{j=2}^n (z_j - A z_{j-1})^t \Gamma^{-1} (z_j - A z_{j-1}) \\ &- \frac{n}{2} \log |\Sigma| - \frac{1}{2} \sum_{j=1}^n (x_j - C z_j)^t \Sigma^{-1} (x_j - C z_j) \end{aligned}$$

Step II: E-step.

Keeping the observed variables x_1, \dots, x_n fixed, we need to compute the expected value of \mathcal{L}_c with respect to $z | x_1, \dots, x_n$.
 \Rightarrow we only need information about $z_j | x_1, \dots, x_n$, which are $z_{j-1}, z_j | x_1, \dots, x_n$, which were derived in the previous section.

We denote $\langle z_j \rangle$ for $E(z_j | x_1, \dots, x_n)$, given by $\hat{\mu}_j$ (bottom of page 19) (21)

$$\text{Similarly, } \begin{cases} \langle z_j z_{j-1}^t \rangle = \hat{\Sigma}_j J_{j-1}^t + \hat{\mu}_j \hat{\mu}_{j-1}^t \\ \langle z_j^t z_j \rangle = \hat{\Sigma}_j + \hat{\mu}_j \hat{\mu}_j^t \end{cases}$$

Step III: M-step.

We treat the green / red / blue terms separately:

$$\begin{aligned} Q_1 &= -\frac{1}{2} \log |P_0| - \frac{1}{2} E_{z_1|X} \left(z_1^t P_0^{-1} z_1 - z_1^t P_0^{-1} \mu_0 - \mu_0^t P_0^{-1} z_1 + \mu_0^t P_0^{-1} \mu_0 \right) \\ &= -\frac{1}{2} \log |P_0| - \frac{1}{2} E_{z_1|X} \text{Tr} \left(P_0^{-1} z_1 z_1^t - P_0^{-1} \mu_0 z_1^t - P_0^{-1} z_1 \mu_0^t + P_0^{-1} \mu_0 \mu_0^t \right) \\ &= \frac{1}{2} \left\{ \log |P_0^{-1}| - \text{Tr} \left\{ P_0^{-1} \left(\langle z_1 z_1^t \rangle - \mu_0 \langle z_1 \rangle^t - \langle z_1 \rangle \mu_0^t + \mu_0 \mu_0^t \right) \right\} \right\} \end{aligned}$$

→ Derivatives with respect to μ_0 :

$$\frac{\partial Q_1}{\partial \mu_0} = -2 P_0^{-1} \langle z_1 \rangle + 2 P_0^{-1} \mu_0 = 0$$

⇒ $\mu_0^{\text{new}} = \langle z_1 \rangle$
computed with 'old' values of the parameters.

$$\frac{\partial}{\partial x} x^t a = \frac{\partial}{\partial x} a^t x = a$$

→ Derivatives with respect to P_0^{-1} : (22)

$$\frac{\partial Q_1}{\partial P_0^{-1}} = \frac{1}{2} \left(P_0 - \langle z_1 z_1^t \rangle + \langle z_1 \rangle \mu_0^t + \mu_0 \langle z_1 \rangle^t - \mu_0 \mu_0^t \right) = 0$$

$$\begin{aligned} \frac{\partial}{\partial A} \text{Tr}(AB) &= B^t \\ \frac{\partial}{\partial A} \ln |A| &= (A^{-1})^t \end{aligned} \Rightarrow P_0^{\text{new}} = \langle z_1 z_1^t \rangle - \langle z_1 \rangle \mu_0^{\text{new}t} - \mu_0^{\text{new}} \langle z_1 \rangle^t + \mu_0^{\text{new}} \mu_0^{\text{new}t} = \langle z_1 z_1^t \rangle - \langle z_1 \rangle \langle z_1 \rangle^t$$

Summarising:

$$\begin{aligned} \mu_0^{\text{new}} &= \langle z_1 \rangle \\ P_0^{\text{new}} &= \langle z_1 z_1^t \rangle - \langle z_1 \rangle \langle z_1 \rangle^t \end{aligned}$$

$$\begin{aligned} Q_2 &= \frac{n-1}{2} \log |\Gamma^{-1}| - \frac{1}{2} E_{z_1|X} \sum_{j=2}^n (z_j - A z_{j-1})^t \Gamma^{-1} (z_j - A z_{j-1}) \\ &= \frac{n-1}{2} \log |\Gamma^{-1}| - \frac{1}{2} E_{z_1|X} \sum_{j=2}^n \left(z_j^t \Gamma^{-1} z_j - z_j^t \Gamma^{-1} A z_{j-1} - z_{j-1}^t A^t \Gamma^{-1} z_j + z_{j-1}^t A^t \Gamma^{-1} A z_{j-1} \right) \\ &= \frac{n-1}{2} \log |\Gamma^{-1}| - \frac{1}{2} \sum_{j=2}^n E_{z_1|X} \text{Tr} \left(\Gamma^{-1} z_j z_j^t - \Gamma^{-1} A z_{j-1} z_j^t - \Gamma^{-1} z_j z_{j-1}^t A^t + \Gamma^{-1} A z_{j-1} z_{j-1}^t A^t \right) \end{aligned}$$

$$Q_2 = \frac{n-1}{2} \log |\Gamma^{-1}| - \frac{1}{2} \sum_{j=2}^n \text{Tr} \left\{ \Gamma^{-1} \left(\begin{aligned} &\langle z_j z_j^t \rangle \\ &- A \langle z_{j-1} z_j^t \rangle \\ &- \langle z_j z_{j-1}^t \rangle A^t \\ &+ A \langle z_{j-1} z_{j-1}^t \rangle A^t \end{aligned} \right) \right\} \quad (23)$$

→ Derivative with respect to A:

$$\frac{\partial \text{Tr}(\Gamma^{-1} A \langle z_{j-1} z_{j-1}^t \rangle A^t)}{\partial A} = 2 \Gamma^{-1} A \langle z_{j-1} z_{j-1}^t \rangle$$

$$\frac{\partial \text{Tr}(\Gamma^{-1} A \langle z_j z_j^t \rangle)}{\partial A} = \Gamma^{-1} \langle z_j z_j^t \rangle$$

$$\frac{\partial \text{Tr}(\Gamma^{-1} \langle z_j z_{j-1}^t \rangle A^t)}{\partial A} = \Gamma^{-1} \langle z_j z_{j-1}^t \rangle$$

Toolbox:

$$\frac{\partial}{\partial A} \text{Tr}(BAC A^t) = B^t A C^t + BAC$$

$$\frac{\partial}{\partial A} \text{Tr}(BAC) = B^t C^t$$

$$\frac{\partial}{\partial A} \text{Tr}(BA^t) = B$$

$$\Rightarrow A^{\text{new}} = \left(\sum_{j=2}^n \langle z_j z_{j-1}^t \rangle \right) \left(\sum_{j=2}^n \langle z_{j-1} z_{j-1}^t \rangle \right)^{-1}$$

$$\begin{aligned} \frac{\partial Q_2}{\partial A} &= \sum_{j=2}^n \Gamma^{-1} A \langle z_{j-1} z_{j-1}^t \rangle \\ &\quad - \Gamma^{-1} \langle z_j z_{j-1}^t \rangle \\ &= 0 \end{aligned}$$

→ Derivative with respect to Γ^{-1} :

$$\frac{\partial Q}{\partial \Gamma^{-1}} = \frac{n-1}{2} \Gamma - \frac{1}{2} \frac{\partial}{\partial \Gamma^{-1}} \left(\sum_{j=2}^n \text{Tr} \left\{ \Gamma^{-1} \left(\begin{aligned} &\langle z_j z_j^t \rangle \\ &- A \langle z_{j-1} z_j^t \rangle \\ &- \langle z_j z_{j-1}^t \rangle A^t \\ &+ A \langle z_{j-1} z_{j-1}^t \rangle A^t \end{aligned} \right) \right\} \right)$$

$$\frac{\partial \text{Tr}(AB)}{\partial A} = B^t$$

$$= \frac{n-1}{2} \Gamma - \frac{1}{2} \sum_{j=2}^n \left(\begin{aligned} &A \langle z_{j-1} z_{j-1}^t \rangle A^t \\ &- A \langle z_{j-1} z_j^t \rangle \\ &- \langle z_j z_{j-1}^t \rangle A^t \\ &+ \langle z_j z_j^t \rangle \end{aligned} \right) = 0$$

$$\Rightarrow \Gamma^{\text{new}} = \frac{1}{n-1} \sum_{j=2}^n \left\{ \begin{aligned} &\langle z_j z_j^t \rangle - A^{\text{new}} \langle z_{j-1} z_j^t \rangle \\ &- \langle z_j z_{j-1}^t \rangle A^{\text{new}^t} \\ &+ A^{\text{new}} \langle z_{j-1} z_{j-1}^t \rangle A^{\text{new}^t} \end{aligned} \right\}$$

Summarizing:

$$A^{\text{new}} = \left(\sum_{j=2}^n \langle z_j z_{j-1}^t \rangle \right) \left(\sum_{j=2}^n \langle z_{j-1} z_{j-1}^t \rangle \right)^{-1}$$

$$\Gamma^{\text{new}} = \frac{1}{n-1} \sum_{j=2}^n \left\{ \begin{aligned} &\langle z_j z_j^t \rangle - A^{\text{new}} \langle z_{j-1} z_j^t \rangle \\ &- \langle z_j z_{j-1}^t \rangle A^{\text{new}^t} \\ &+ A^{\text{new}} \langle z_{j-1} z_{j-1}^t \rangle A^{\text{new}^t} \end{aligned} \right\}$$

The blue term can be treated as the red term :

(25)

wherever z_j appears \rightarrow replace with x_j

$z_{j-1} \xrightarrow{A} z_j$

$A \xrightarrow{C} C$

& Summation starts with $j=1$.

$$C^{new} = \left(\sum_{j=1}^n x_j \langle z_j \rangle^t \right) \left(\sum_{j=1}^n \langle z_j z_j^t \rangle \right)^{-1}$$

$$\Sigma^{new} = \frac{1}{n} \sum_{j=1}^n \left\{ \begin{aligned} &x_j x_j^t - C^{new} \langle z_j \rangle x_j^t \\ &- x_j \langle z_j \rangle^t C^{new} \\ &+ C^{new} \langle z_j z_j^t \rangle C^{new} \end{aligned} \right\}$$

EM algorithm for Kalman Filtering.

(i) Initialize $\mu_0^{old}, P_0^{old}, A, \Gamma^{old}, C^{old}, \Sigma^{old}$

(ii) Repeat

* Using old parameter values, compute

\rightarrow FWD variables $\mu_1, \dots, \mu_n, V_1, \dots, V_n$ (p.7/8)

\rightarrow then BWD variables $\hat{\mu}_n, \dots, \hat{\mu}_1, \hat{V}_n, \dots, \hat{V}_1$ (p.17)

\rightarrow Calculate $\langle z_j \rangle, \langle z_j z_{j-1}^t \rangle, \dots$ (page 19)

* Update parameters

$\rightarrow \mu_0^{new}, P_0^{new}, A^{new}, \Gamma^{new}, C^{new}, \Sigma^{new}$ (p.22/24/25)

* $\mu_0^{old} \leftarrow \mu_0^{new} \quad A^{old} \leftarrow A^{new} \quad C^{old} \leftarrow C^{new}$
 $P_0^{old} \leftarrow P_0^{new} \quad \Gamma^{old} \leftarrow \Gamma^{new} \quad \Sigma^{old} \leftarrow \Sigma^{new}$

Summary of Notation + Useful formula

(26)

(i) NOTATION.

$\alpha(z_j) = p(x_1, \dots, x_j, z_j) =$ FWD variable

$\beta(z_j) = p(x_{j+1}, \dots, x_n | z_j) =$ BWD variable.

$c_j = p(x_j | x_1, \dots, x_{j-1})$

$\alpha(z_j) = \left(\prod_{m=1}^j c_m \right) \hat{\alpha}(z_j), \quad \hat{\alpha}(z_j) = p(z_j | x_1, \dots, x_j)$

$\beta(z_j) = \left(\prod_{m=j+1}^n c_m \right) \hat{\beta}(z_j), \quad \hat{\beta}(z_j) = \frac{\beta(z_j)}{p(x_{j+1}, \dots, x_n | x_1, \dots, x_j)}$

$\gamma(z_j) = p(z_j | x_1, \dots, x_n) = \frac{\alpha(z_j) \beta(z_j)}{p(x_1, \dots, x_n)} = \hat{\alpha}(z_j) \hat{\beta}(z_j)$

$\hat{\alpha}(z_j) = \mathcal{N}(z_j | \mu_j, V_j)$

$\gamma(z_j) = \mathcal{N}(z_j | \hat{\mu}_j, \hat{V}_j)$

$\begin{cases} \mu_j = A \mu_{j-1} + K_j (x_j - C A \mu_{j-1}) \\ V_j = (I - K_j C) P_{j-1} \end{cases}$

where $K_j = P_{j-1} C^t (\Sigma + C P_{j-1} C^t)^{-1}$

$P_j = \Gamma + A V_j A^t$

$c_j = \mathcal{N}(x_j | C A \mu_{j-1}, C P_{j-1} C^t + \Sigma)$

$\begin{cases} \hat{\mu}_j = \mu_j + J_j (\hat{\mu}_{j+1} - A \mu_j) \\ \hat{V}_j = V_j + J_j (\hat{V}_{j+1} - P_j) J_j^t \end{cases}$

where $J_j = V_j A^t P_j^{-1}$

(ii) USEFUL FORMULA.

(27)

- If $p(x) = \mathcal{N}(x | \mu, \Lambda^{-1})$
 $p(y|x) = \mathcal{N}(y | Ax + b, L^{-1})$

Then

$$\rightarrow \begin{pmatrix} x \\ y \end{pmatrix} \text{ is Gaussian with } E \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} \mu \\ A\mu + b \end{pmatrix}$$

$$\text{and } \text{Cov} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} \Lambda^{-1} & \Lambda^{-1}A^t \\ A\Lambda^{-1} & L^{-1} + A\Lambda^{-1}A^t \end{pmatrix}$$

$$\rightarrow p(y) = \mathcal{N}(y | A\mu + b, L^{-1} + A\Lambda^{-1}A^t)$$

$$p(x|y) = \mathcal{N}(x | S \{ A^t L(y - b) + \Lambda \mu \}, S)$$

where

$$S = (\Lambda + A^t L A)^{-1}$$

• Woodbury:

$$(A + B D^{-1} C)^{-1} = A^{-1} - A^{-1} B (D + C A^{-1} B)^{-1} C A^{-1}$$

$$\bullet (P^{-1} + B^t R^{-1} B)^{-1} B^t R^{-1} = P B^t (B P B^t + R)^{-1}$$

• Matrix derivatives:

$$\rightarrow \frac{\partial}{\partial x} x^t a = \frac{\partial}{\partial x} a^t x = a$$

$$\rightarrow \frac{\partial}{\partial A} \text{Tr}(AB) = B^t$$

$$\frac{\partial}{\partial A} \ln |A| = (A^{-1})^t$$

$$\rightarrow \frac{\partial}{\partial A} \text{Tr}(BAC A^t) = B^t A C^t + B A C$$

$$\frac{\partial}{\partial A} \text{Tr}(BAC) = B^t C^t$$

$$\frac{\partial}{\partial A} \text{Tr}(B A^t) = B$$