

110 - sequential models

(*) explore other properties of graphical models with Gaussian distri.

(*) should already be familiar with mixture models (graphically)

- factor analysis just codes both Y and X continuous (cont./discrete)

(*) can better understand embeddings

(*) HMMs as a time-sequenced MM.

(*) HMM inference skipped; focus on variable elim; message passing interpretations

EX: what is used for inference in HMM?

①: - viterbi algorithm? (MAP assigned)

- forward-backwards / α - β recursion / BP \rightarrow param estimation

- inference tasks: - Baum-Welch

- factorial HMM, switching SSM (see diagram).

(*) today's extrapolation is to FA, SSMS (see diagram).

- algorithms prev. studied all have counterparts

- spirit \rightarrow break problem down locally into a subproblem

(*) factorial HMM, switching SSM.

EX: ML is not about equations, maths; use a story/motivating problem to anchor the mathematics in a solid mental model!

② Factorial HMM \rightarrow engine-switching time-series?

- 3 dealers with dishonest casino

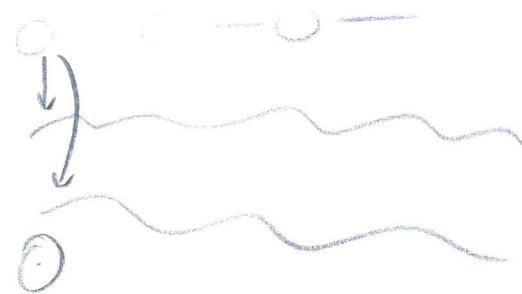
\rightarrow aircraft radar measurements (obs.)
physical locations (states)

SSM

- envisage S as ~~obs~~ generating which aircraft appears on your screen/radar.
- multiple aircrafts

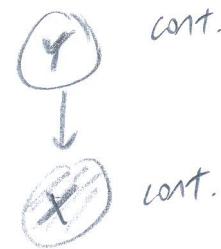
switching SSM

\rightarrow



- Ex: think of what kind of stories can be told with these equations

(lit. correspond with real). This is art of modelling; way more important than the mathematics. latter is a language to express ideas, rather than an end in itself.



cont.

(*) MVG

- covered in previous lectures

(*) M.I.L.

- inclusion of matrices in terms of sub-partitions.

Ex: Remembered this

(*) Matrix algebra:

- convert matrix \rightarrow no. via trace

(1) Review results.

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

(*) Factor Analysis

- as an example; imagine $\underline{x} \in \mathbb{R}^2$ (e.g. a plane).
- imagine different orientations of the plane \rightarrow can be described with 3D world.
- plane is a manifold' subspace
- \underline{y} corresponds to points in 3d space e.g. $\underline{y} \in \mathbb{R}^3$

(*) Relation between \underline{x} and \underline{y}

- orientation of subspace affects how points in manifold assigned to coordinates in 3d space.

$$\underline{y} = \mu + \Lambda \underline{x} \quad (\text{convert } \underline{x} \in \mathbb{R}^2 \rightarrow \underline{y} \in \mathbb{R}^3 \text{ (projection?!)})$$

(*) Λ - factor loading matrix

$\underline{\Lambda}$ - diagonal

- (*) (2) - geometric story review - important ✓
- idea of a latent space (lower dimensional)
(but not all the same)

- similar ideas in PCA

(*) Marginal data distn

- we have a latent factor \underline{x} in a low-dim space that is observed
- we do observe y , whose components/points we do know.

$$\textcircled{X} \quad p(\underline{x}) = N(\underline{x} | 0, I)$$

w-noise term
(like ϵ)

$$\textcircled{Y} \quad p(y|\underline{x}) = N(y | \mu + \Lambda \underline{x}, \Psi) \quad y = \mu + \Lambda \underline{x} + w \quad w \sim N(0, \Psi)$$

goal: infer $p(\underline{x}|y)$ i.e. observed given observations.

- latent space is not necessarily low dimensional (e.g. 10d \rightarrow physical)
2d 3d

ex: A procedure for achieving goal.

- we know $p(\underline{x})$ and $p(y|\underline{x})$ is Gaussian.
- Here $p(\underline{x}, y)$ is jointly Gaussian, $p(\underline{x}|y)$ is Gaussian

$$\begin{pmatrix} \underline{x} \\ y \end{pmatrix} \sim N\left(\begin{pmatrix} \mu_{\underline{x}} \\ \mu_y \end{pmatrix}, \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix}\right)$$

(A3): use Gaussian results to derive $p(y)$ and $p(\underline{x}|y)$

(E1): This is the key idea for inference in factor analysis

(*) Being sloppy
with \underline{x}, y, X, Y
(all same)
i.e. capitalised.

$$\mu_x = 0$$

$$\underline{\mu} = E[y] = E[\mu + \Lambda \underline{x} + w] = \underbrace{\mu}_{=0} + \Lambda E[\underline{x}] + E[w] = \mu$$

$$\Sigma_{xx} = I$$

$$\Sigma_{yy} = var[Y] = E[(Y - \mu)(Y - \mu)^T]$$

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

$$= E[(\Lambda \underline{x} + w)(\Lambda \underline{x} + w)^T] = E[(\Lambda \underline{x} \underline{x}^T \Lambda^T + \Lambda \underline{x} w^T + w \underline{x}^T \Lambda^T + w w^T)]$$

$$= E[(\Lambda \underline{x} \underline{x}^T \Lambda^T + \Lambda E[\underline{x} \underline{x}^T] + E[w \underline{x}^T] \Lambda^T + w w^T)] \textcircled{A4} \cdot \text{review } \checkmark$$

$$= \Lambda E[\underline{x} \underline{x}^T] \Lambda^T + E[w w^T] = \Lambda \Lambda^T + \Psi$$

$$\begin{aligned}
 \Sigma_{xy} &= E[(\underline{x} - \mu_x)(\underline{y} - \mu_y)^T] \\
 &= E[(\underline{x} - \mu_x)(\mu + \Delta x + \underline{w} - \mu)^T] \\
 &= E[\underline{x}(\Delta x + \underline{w})^T] \quad \mu x = 0 \\
 &= E[\underline{x}\underline{x}^T]\Delta^T + E[\underline{x}]E[\underline{w}^T] \\
 &= \Delta^T
 \end{aligned}$$

⑮ Review calc.

(*) FA joint distn

- yielding:

model $p(\underline{x}) = N(0, I) \quad p(y|\underline{x}) = N(\mu + \Delta \underline{x}, \Psi)$

cov. $\text{cov}(\underline{x}, y) = E[(\underline{x} - 0)(y - \mu)^T] = E[\underline{x}(\mu + \Delta \underline{x} + \underline{w} - \mu)^T]$

$$\begin{aligned}
 &= E[\underline{x}\underline{x}^T]\Delta^T + \underline{x}\underline{w}^T \\
 &= E[\underline{x}\underline{x}^T]\Delta^T + E[\underline{x}]\widetilde{\underline{w}}^T \\
 &\quad \underset{=0}{=} \quad \textcircled{ab} \\
 &= \Delta^T
 \end{aligned}$$

Joint distn:- $p\left[\begin{array}{c} \underline{x} \\ y \end{array}\right] = N\left(\left[\begin{array}{c} \underline{x} \\ y \end{array}\right] \mid \left[\begin{array}{c} 0 \\ \mu \end{array}\right], \left[\begin{array}{cc} I & \Delta^T \\ 0 & \Delta^T + \Psi \end{array}\right]\right)$

- (*) Assume noise uncorrelated with data or latent variables

(*) Inferred FA

- Review derivation of post.

- ex: posterior distn that is derived \rightarrow satisfactory?

Q: Yes, ~~and~~ given that invertibility is satisfied

- apply Gaussian condit. formulae.

- set $\Sigma_{11} = I$, $\Sigma_{12} = \Sigma_{21}^T = A^T$ $\Sigma_{22} = (\Delta A^T + \Psi)$

- posterior of latent \mathbf{z} given obs \mathbf{y}

$$p(\mathbf{z} | \mathbf{y}) = N(\mathbf{z} | M_{1|2}, V_{1|2})$$

$$\begin{aligned} M_{1|2} &= \mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (\mathbf{y} - \mu_2) & V_{1|2} &= \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \\ &= A^T (\Delta A^T + \Psi)^{-1} (\mathbf{y} - \mu) & &= I - A^T (\Delta A^T + \Psi)^{-1} A \end{aligned} \quad (*)$$

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

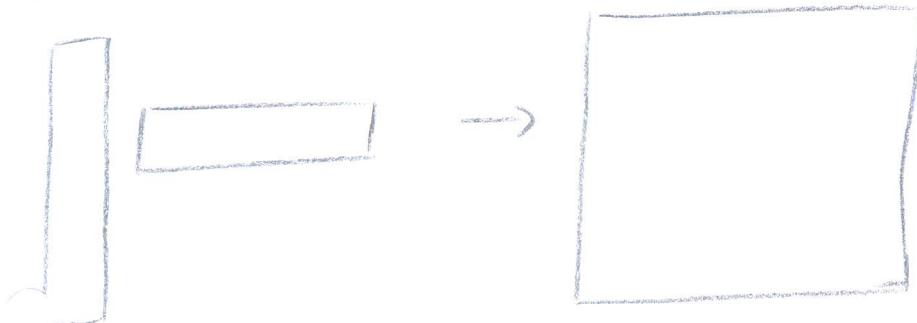
$$\Rightarrow V_{1|2} = (I + A^T \Psi^{-1} A)^{-1} \quad M_{1|2} = V_{1|2} A^T \Psi^{-1} (\mathbf{y} - \mu)$$

(*) computationally \rightarrow examine ΔA^T and Ψ

- inverting Ψ is trivial as diagonal
- ΔA^T (matrix product of factor loadings)

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

(*) Almost one-to-one correspondence



(*) $I - I A^T (\Delta A^T + \Psi)^{-1} A I$

- MIL allows us to re-express $V_{1|2}$ and $M_{1|2}$ in a different form

- Have different computational implications

(*) ex. exp. of computational savings of MIL through projection, dimensionality of Ψ requires clarity on your part. (inversion of smaller matrix). methodologically * focus on this

i) get joint distri for Gaussians.

ii) compute condit. mean/covariance (accounting for comp.)

iii) using MIL to reduce dimensionality

(*) FA - constrained cov. Gaussian

- review slides $\textcircled{A9}$

(*) Geometric interp.

- review $\textcircled{A10}$

(*) Estimating F.A.

$\textcircled{?}$ - derivation of $p(\mathbf{y}|\theta)$..

- from earlier we add mixture
- focus on estimation of params. given $\{\mathbf{y}_n\}_{n=1}^N$
 - loading matrix Λ
 - manifold center μ
 - variance Ψ

Ex: what statistical paradigm is appropriate for estimation?

→ (i.e. what procedure is suitable for the construction of estimators?)

- under MLE (at a cursory level):-

$$[\Lambda^*, \mu^*, \Psi^*] = \underset{\theta}{\operatorname{argmax}} \ell(\theta; \mathbf{y}) = \underset{\theta}{\operatorname{argmax}} p(\mathbf{y})$$

$\textcircled{?}$

(*) EM for Factor Analysis

$$\begin{aligned} \text{- Incomplete data log likelihood} \quad \ell(\theta, \mathbf{D}) &= -\frac{N}{2} \log |\Lambda \Lambda^T + \Psi| - \frac{1}{2} \sum_{n=1}^N (\mathbf{y}_n - \mu)^T (\Lambda \Lambda^T + \Psi)^{-1} (\mathbf{y}_n - \mu) \\ &= -\frac{N}{2} \log |\Lambda \Lambda^T + \Psi| - \frac{1}{2} \operatorname{tr} [(\Lambda \Lambda^T + \Psi)^{-1} \underline{S}] \end{aligned}$$

$$\text{where } \underline{S} = \sum_{n=1}^N (\mathbf{y}_n - \mu)(\mathbf{y}_n - \mu)^T$$

$$\text{Estimating } \mu: \hat{\mu}_{\text{MLE}} = \frac{1}{N} \sum_{n=1}^N \mathbf{y}_n$$

However, est. 1 and 4 tricky as there is a non-linear coupling in log-like complete log-like.

(AII)

$$\begin{aligned} \ell_c(\theta; \Omega) &= \sum_{n=1}^N \log p(x_n, y_n) = \sum_{n=1}^N \log p(x_n) + \log p(y_n | x_n) \\ &= -\frac{N}{2} \log |\Omega| - \frac{1}{2} \sum_{n=1}^N x_n^T x_n - \frac{N}{2} \log |\Psi| - \frac{1}{2} \sum_{n=1}^N (y_n - \Omega x_n)^T \Psi^{-1} (y_n - \Omega x_n) \\ &= -\frac{N}{2} \log |\Psi| - \frac{1}{2} \sum_{n=1}^N \text{tr}[x_n x_n^T] - \frac{N}{2} \text{tr}[S \Psi^{-1}] \end{aligned}$$

where $S = \sum_{n=1}^N (y_n - \mu)(y_n - \mu)^T$

(*) E-step for factor analysis

(*) imagine x is observed; can do inference on x given y ; and can always compute sufficient statistics of x . (?)

(*) $p(x|y) = \langle x \rangle, \langle xx^T \rangle$

(*) M-step for factor analysis

(AII)- review derivation of EM for FA (Jordan 2003)

- A counterpart of MM ^{with} continuous latent r.v.s.

Summary

1. MM has discrete latent state } some topology graphically
FA has continuous latent state }

2. MM $\rightarrow p(x|y) = \frac{p(x,y)}{p(y)} = \frac{p(x,y)}{\sum_x p(x,y)}$ Inference

3. FA $\rightarrow p(x), p(y|x) \rightarrow p(x|y)$

generative

$\hookrightarrow p(x,y) \rightarrow p(x|y) + \text{M.I.L}$

4. Param est.

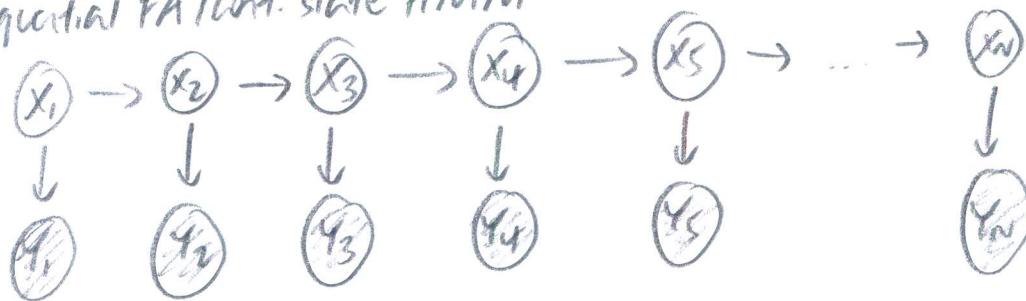
- ML estimation
- Hidden r.v.
- Benefit from inference sol. to comp. exp/sufficient } - use EM

(*) Model variance and identifiability

(iii) - review

(*) SSMs (HMM counterpart), or LDS

- sequential FA / cont. state HMM



$$\begin{aligned}x_t &= \Lambda x_{t-1} + g w_t & w_t &\sim N(0, Q) & z_0 &\sim N(0, \Sigma_0) \\y_t &= (x_t + v_t) & v_t &\sim N(0, R)\end{aligned}$$

- An LDS

- In general:

$$\begin{aligned}x_t &= f(x_{t-1}) + g w_{t-1} & f \text{- arbitrary dyn. model} \\y_t &= g(x_t) + v_t & g \text{- arb. obs. model}\end{aligned}$$

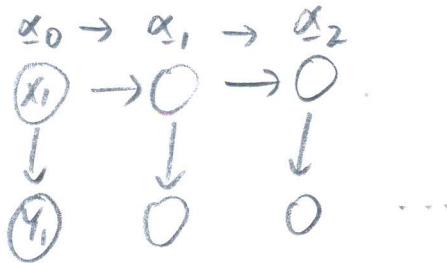
(*) use HMMs as a reference point for similarity/difference

(*) LDS - 2D tracking

(*) inference problem I.

- filtering: given $y_1, y_2, y_3, \dots, y_t$ estimate $x_t = P(x_t | y_1:t)$
- given all previous observations; estimate current latent state (the pos.)
(algo)
- Kalman filter → exact online inference/sequential Bayesian inference.
in an LDS
- Gaussian analog of forward algorithm for HMMs, in continuous space

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



- A non-trivial collection

- HMM as collection of sequential MMs.

- forward algorithm is a recursive algorithm

(*) Inference problem 2 (not emphasised)

- smoothing \rightarrow given y_1, \dots, y_T , estimate x_t ($t < T$)

- Rauch-Tung-Striebel smoother

- \hookrightarrow exact offline inference as an LDS

b) Gaussian analog of forwards-backwards (alpha-gamma recursion)

(†) Kalman filtering derivation

- given $y_{1:t}$ i.e. $y_1, y_2, y_3, y_4, \dots, y_t$

- question $p(x_t|y_{1:t})$

- you already have $\nearrow p(x_{t-1}|y_{1:t})$

$\backslash \nearrow p(x_t|x_{t-1})$

$p(y_t|x_t)$

(cond.)

- due to Gaussian property :- only need mean and cov. of $p(x_t|y_{1:t})$

$$\mathbb{E}[x_t|y_{1:t}] = \mu_{t|t} \quad \mathbb{E}[(x_t - \mu_{t|t})^2] = P_{t|t} \quad \text{①}$$

- inference \rightarrow need to compute cond. mean and covariance.

- Kalman filtering is a recursive procedure to update belief state

- split into

i) prediction step

ii) update step

(*) A13

- Review SSM, RF derivation slides
- continued next lecture

VI - variational inference - loopy belief prop.(*) continue VI - derivation of KF for SSMSKalman filtering derivation

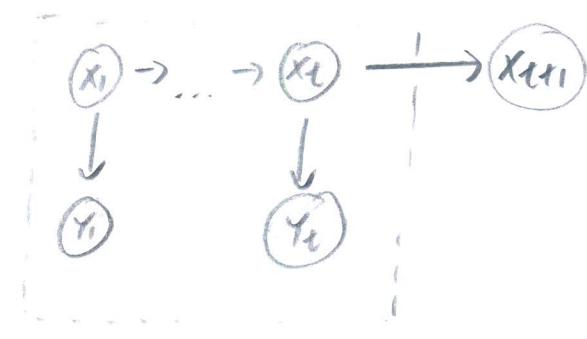
- recall filtering problem
- inference of $p(x_t | y_{1:t})$ (i.e. cond. dist. of last hidden state given all observations up to now)
- use diagram as reference point for visualizing queries.
- we require:-

$p(x_t | y_{1:t})$ from a recursive procedure involving $p(x_{t+1} | y_{1:t})$, y_{t+1}

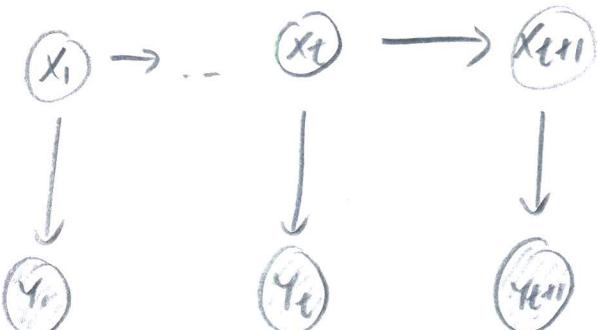
- break up problem into 2 steps:-

1) Time update:

- compute $p(x_{t+1} | y_{1:t})$ from prior belief $p(x_t | y_{1:t})$ and dyn. model $p(x_{t+1} | x_t)$

2) Measurement update:

- 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 obs. model $p(y_{t+1} | x_{t+1})$

(*) Make use of Gaussian properties everywhere

- inference problem \rightarrow dealing with ^{new} means and covariances

(*) same technique as in FA.

(*) Schematic of the strategy:-
(gaussian manip.)

M.I.L.

$$\begin{bmatrix} z_1 \\ z_2 \end{bmatrix} \sim N \left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \right) \quad z_1 \rightarrow z_1 z_2 \rightarrow z_1/z_2$$

corresp.

$$z_1 \sim x_t | y_{1:t} \quad x_{t+1} = Ax_t + Gw_t \rightarrow p(x_{t+1}, y_{1:t} | y_{1:t}) \rightarrow p(x_{t+1} | y_{1:t}, y_{1:t})$$

$$= f(x_t)$$

$$p(x_{t+1} | y_{1:t})$$

④ μ_1

$$p(y_{1:t} | x_{t+1})$$

↑

$$y_{1:t} = Cx_{1:t} + v_{1:t}$$

μ_2 / μ_1 ④

① - tidy this up
(makes some sense; but want
full clarity)

- ex: conditioning is not scary; just
some constants in eq.; focus on
LHS of cond.

(*) Compute means and covariances within the schematic

② review

(*) predict step

- asymmetrical model $x_{t+1} = Ax_t + Gw_t \quad w_t \sim N(0, Q)$

$$\hat{x}_{t+1|t} = E[x_{t+1} | y_1, \dots, y_t] = A\hat{x}_{t|t}$$

$$P_{t+1|t} = E[(x_{t+1} - \hat{x}_{t+1|t})(x_{t+1} - \hat{x}_{t+1|t})^T | y_1, \dots, y_t]$$

$$= E[(Ax_t + Gw_t - \hat{x}_{t+1|t})(Ax_t + Gw_t - \hat{x}_{t+1|t})^T | y_1, \dots, y_t]$$

$$= AP_{t|t}A^T + GP^T G$$

$\hat{x}_{t+1|t}$ and $P_{t+1|t}$ are conditional means and covariances (*)

③ notation

$$\text{observation model} \quad y_{t+1} = Cx_t + v_t \quad v_t \sim N(0, R)$$

$$E[y_{t+1}|y_1, \dots, y_t] = E[(x_{t+1} + v_{t+1})|y_1, \dots, y_t] = \hat{x}_{t+1|t}$$

$$E[(y_{t+1} - \hat{x}_{t+1|t})(y_{t+1} - \hat{x}_{t+1|t})^T | y_1, \dots, y_t] = P_{t+1|t} C^T + R$$

(*) update step

- From previous slide:-

we have joint $p(x_{t+1}, y_{t+1}|y_1, \dots, y_t) = N(\underline{m}_{t+1}, \underline{V}_{t+1})$ where:-

$$\underline{m}_{t+1} = \begin{pmatrix} \hat{x}_{t+1|t} \\ \hat{C}_{t+1|t} \end{pmatrix} \quad \underline{V}_{t+1} = \begin{pmatrix} P_{t+1|t} & P_{t+1|t} C^T \\ C P_{t+1|t} & C P_{t+1|t} C^T + R \end{pmatrix}$$

- Now use conditional Gaussian formulae for partitioned matrices to get $p(x_{t+1}|y_{t+1}, y_{1:t})$

(*) Kalman Filter

- Has closed form measured and time updates:-

measurement update :- $\hat{x}_{t+1|t+1} = \hat{x}_{t+1|t} + K_{t+1}(y_{t+1} - \hat{x}_{t+1|t})$

$$P_{t+1|t+1} = P_{t+1|t} - K_{t+1} P_{t+1|t} \quad ?$$

- K_{t+1} - Kalman gain matrix

time updates: $\hat{x}_{t+1|t} = A\hat{x}_{t|t}$

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

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

(*) Review notation and intuition.

Ex. Be sensitive to conditioning in terms of info being used at each point.

(*) Example of KF in 1D

- just apply KF equations

(**) - intuition behind these

(*) KF intuition (**)

- KF update of mean:- $\hat{x}_{\text{new}} = \hat{x}_{\text{old}} + K_{11} (\underline{z}_{11} - \hat{x}_{\text{old}})$

$$= \frac{(\sigma_t + \sigma_x) \underline{z}_t + \sigma_z \hat{x}_{\text{old}}}{\sigma_t + \sigma_x + \sigma_z}$$

- innov.

@ new belief is convex combo of updates from prior and observation, weighted by Kalman gain matrix

$$K_{11} = P_{\text{old}} C^T (C P_{\text{old}} C^T + R)^{-1}$$

- observation unreliable, σ_z (i.e. R) large so K_{11} is small \rightarrow more att. to prediction

- old prior is reliable, large σ_t , or process is unpredictable, large σ_x , we pay more attention to observation.

(*) Where do A, G, C matrices come from?

- A parameter estimation (as opposed to filtering, which is an inference problem)
 - ↳ i.e. compact p.d. of latent var, given data in this context.
- inference assumes learning / param est. has already occurred.

(*) Complexity of one KF step., Inf prob 2, RTS smoother, RTS deriv.

- Review # in OM time

(*) learning SSMS

- (*) use EM. inference from ICF or RTS filtering forms E-step!
- (*) Form complete log-like.
- (*) from E-step \rightarrow compute suffi. statistics using inference results to give you a posterior.

M-step \rightarrow MLE

Ex: You can and should derive m your spare time.

- (*) very principle \rightarrow inference can be treated as a subroutine for learning in a partially obs. setting

(*) nonlinear systems

- use derivatives, Taylor exp., approxim., linearization.

- (*) L.S thinking \rightarrow you can approximate functions to the order you can afford
- (*) Ex A lot of seemingly difficult, obtuse algebra can be viewed/reduced through an understanding of principles motivating them (often new).

