

110 - sequential models

(*) explore other properties of graphical models with Gaussian distri.
 (*) should already be familiar with mixture models (graphically)

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

(*) can better understand embeddings

(*) HMMs as a time-sequenced MM.

(*) HMM inference skipped; ex focused 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

(*) 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.

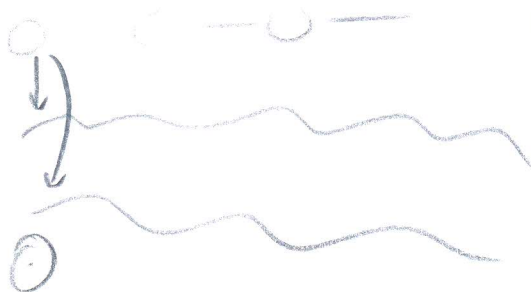
②: Factorial HMM \rightarrow regime-switching time-series?

- 3 dealers with dishonest casino

SSM

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

switching SSM \rightarrow



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

ex: think of what kind of stories can be told with these equations
 (i.e. corresp. 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.

- (*) MVG
- covered in previous lectures

- (*) M.I.L.
- inclusion of matrices in terms of sub-partitions.
Ex: Remembered this

- (*) Matrix algebra.
- invert matrix \rightarrow no. via trace

(A): Review results.

$$\frac{1}{2}(\underline{x}^T \underline{A} \underline{x}) = \underline{x}^T \underline{A} \underline{x}$$

- (*) 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 co-ord.
- 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 ^{are} assigned co-ordinates in 3d space. (subspace)¹

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

- (*) $\underline{\Lambda}$ - factor loading matrix

$\underline{\Lambda}$ - diagonal

- (*) (A) - 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 x in a low-dim space that is unobserved
- we do observe y , whose components/points we do know.

w -noise term
(like ϵ)

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



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

GOAL: infer $p(x|y)$ i.e. unobserved given observations.

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

EX: A procedure for achieving goal.

- we know $p(x)$ and $p(y|x)$ is Gaussian.
- Hence $p(x, y)$ is jointly Gaussian, $p(x|y)$ is Gaussian

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

- (A3): use Gaussian results to derive $p(y)$ and $p(x|y)$

(A4): This is the key idea for inference in factor analysis

(*) being sloppy
with x, y, μ, Σ
(all same) ~~the~~
i.e. capitalised.

$$\mu_x = 0$$

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

$$\Sigma_{xx} = I$$

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

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

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

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

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

$$\begin{aligned}
 \Sigma_{xy} &= E[(X - \mu_x)(Y - \mu_y)^T] \\
 &= E[(X - \mu_x)(\mu + \Lambda X + W - \mu)^T] \\
 &= E[X(\Lambda X + W)^T] \quad \mu_x = 0 \\
 &= E[XX^T]\Lambda^T + E[X]E[W^T] \\
 &= \Lambda^T
 \end{aligned}$$

Ⓐ Review calc.

(*) FA joint distri

- yielding:

model $p(x) = N(0, I) \quad p(y|x) = N(\mu + \Lambda x, \Psi)$

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

$$\begin{aligned}
 &= E[xx^T\Lambda^T + xW^T] \\
 &= E[xx^T]\Lambda^T + E[x]W^T \\
 &\quad \underbrace{E[x] = 0}_{(46)} \\
 &= \Lambda^T
 \end{aligned}$$

Joint distri:- $p\left[\begin{bmatrix} x \\ y \end{bmatrix}\right] = N\left(\begin{bmatrix} x \\ y \end{bmatrix} \middle| \begin{bmatrix} 0 \\ \mu \end{bmatrix}, \begin{bmatrix} I & \Lambda^T \\ 0 & \Lambda\Lambda^T + \Psi \end{bmatrix}\right)$

(*) Assume noise uncorrelated with data or latent variables

(*) Inference in FA

- Ⓐ - Review derivation of post.

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

(N): yes, ~~and~~ given that identifiability is satisfied

- Apply Gaussian condit. formulae.

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

- posterior of latent \underline{x} given obs \underline{y}

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

$\underline{m}_{1|2} = \underline{\mu}_1 + \Sigma_{12} \Sigma_{22}^{-1} (\underline{y} - \underline{\mu}_2)$

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

$= \underline{\Lambda}^T (\underline{\Lambda}\underline{\Lambda}^T + \underline{\Psi})^{-1} (\underline{y} - \underline{\mu})$

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

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

(A7)

$\Rightarrow \underline{V}_{1|2} = (\underline{I} + \underline{\Lambda}^T \underline{\Psi}^{-1} \underline{\Lambda})^{-1}$ $\underline{m}_{1|2} = \underline{V}_{1|2} \underline{\Lambda}^T \underline{\Psi}^{-1} (\underline{y} - \underline{\mu})$

(*) Computationally \rightarrow examine $\underline{\Lambda}\underline{\Lambda}^T$ and $\underline{\Psi}$

- Inverting $\underline{\Psi}$ is trivial as diagonal

- $\underline{\Lambda}\underline{\Lambda}^T$ (matrix product of factor loadings)

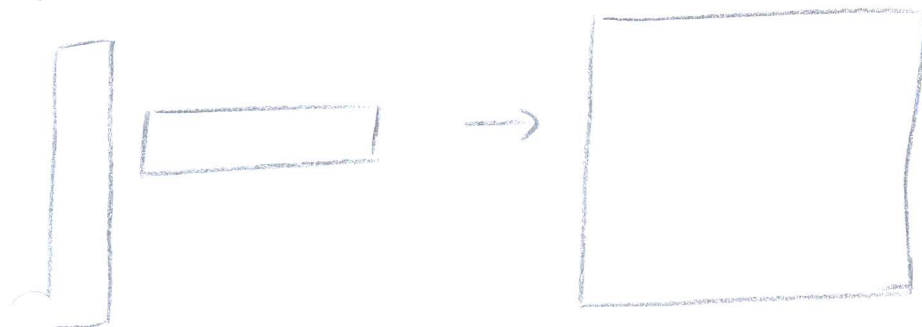
$\underline{E} = \underline{I}$

$\underline{F} = \underline{\Lambda}^T$

$\underline{G} = \underline{\Lambda}$

$\underline{H} = \underline{\Psi}$

(*) Almost one-to-one correspondence



(*) $\underline{I} - \underline{\Lambda}^T (\underline{\Lambda}\underline{\Lambda}^T + \underline{\Psi})^{-1} \underline{\Lambda}$

- MIL allows us to re-express $\underline{V}_{1|2}$ and $\underline{m}_{1|2}$ in a different form

- Have different computational implications

(A8) ex exp. of computational savings of MIL through projection, dimensionality of $\underline{\Psi}$ requires clarity on your part. (reduction 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 (A9)

(*) Geometric interp.

- Review (A10)

Estimating F.A.

- ^② ~~look~~ earlier we did inference - derivation of $p(\mathbf{z}|\mathbf{y})$..

- focus on estimation of params. given $\{\mathbf{y}_n\}_{n=1}^N$

- loading matrix $\mathbf{\Lambda}$

- 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):-

$$\underbrace{[\mathbf{\Lambda}^*, \mu^*, \Psi^*]}_{\underline{\theta}} = \underset{\underline{\theta}}{\operatorname{argmax}} \ell(\underline{\theta}; \mathbf{y}) = \underset{\underline{\theta}}{\operatorname{argmax}} \underbrace{p(\mathbf{y})}_{\textcircled{2}}$$

(*) EM for Factor Analysis

- Incomplete data log like. (marginal) of \mathbf{y}

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

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

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

However, est. $\underline{\Lambda}$ and $\underline{\Psi}$ tricky as there is a non-linear coupling in log-like.
complete log-like. (A11)

$$\begin{aligned} \ell(\theta; D) &= \sum_{n=1}^N \log p(\mathbf{x}_n, \mathbf{y}_n) = \sum_{n=1}^N \log p(\mathbf{x}_n) + \log p(\mathbf{y}_n | \mathbf{x}_n) \\ &= -\frac{N}{2} \log |\underline{\Sigma}| - \frac{1}{2} \sum_{n=1}^N \mathbf{x}_n^T \mathbf{x}_n - \frac{N}{2} \log |\underline{\Psi}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{y}_n - \underline{\Lambda} \mathbf{x}_n)^T \underline{\Psi}^{-1} (\mathbf{y}_n - \underline{\Lambda} \mathbf{x}_n) \\ &= -\frac{N}{2} \log |\underline{\Psi}| - \frac{1}{2} \sum_{n=1}^N \text{tr}[\mathbf{x}_n \mathbf{x}_n^T] - \frac{N}{2} \text{tr}[\underline{S} \underline{\Psi}^{-1}] \end{aligned}$$

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

(*) E-step for factor analysis

(*) Imagine \mathbf{x} is observed; can do inference on \mathbf{x} given \mathbf{y} ; and can always compute sufficient statistics of \mathbf{x} . (?)

(*) $p(\mathbf{x} | \mathbf{y}) = \langle \mathbf{x} \rangle, \langle \mathbf{x} \mathbf{x}^T \rangle$

(*) M-step for factor analysis

(A11) - review derivation of EM for FA (Jordan 2003)
 - A counterpart of MM ^{with} continuous latent: C.V.S.

* Summary

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

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

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

generative

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

4. Param est.

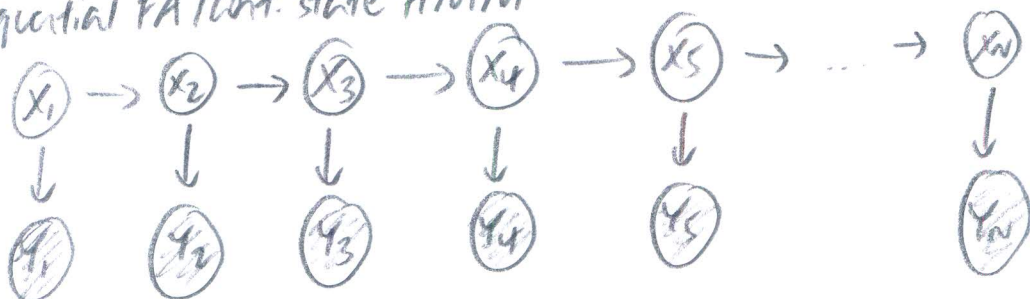
- ML estimation
- HMM i.v.
- Benefit from inference sol. to comp. exp/sufficient } - use EM

(*) Model variance and identifiability

(AII) - review

(*) SSMS (HMM counterpart). or LDS

- sequential FA / cont. state HMM



$$\begin{aligned} x_t &= Ax_{t-1} + Gw_t & w_t &\sim N(0, Q) & x_0 &\sim N(0, \Sigma_0) \\ y_t &= Cx_{t-1} + v_t & v_t &\sim N(0, R) \end{aligned}$$

- An LDS

- In general:

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

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

(*) LDS - 2D tracking

(*) Inference problem 1.

- 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 latest state (the pos.)
- Kalman filter ^(algo) \rightarrow exact online inference/sequential Bayesian inference.
- Gaussian analog of forward algorithm for HMMs, in continuous space

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

$$\alpha_0 \rightarrow \alpha_1 \rightarrow \alpha_2$$

$$(x_1) \rightarrow \bigcirc \rightarrow \bigcirc$$

$$\downarrow \quad \downarrow \quad \downarrow$$

$$(y_1) \quad \bigcirc \quad \bigcirc \quad \dots$$

- A non-trivial connection

- 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 in an LDS

\hookrightarrow Gaussian analog of forwards-backwards (alpha-gamma recursion)

(*) Kalman filtering derivation

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

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

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

$$\begin{array}{l} \swarrow p(x_t | x_{t-1}) \\ \searrow p(y_t | x_t) \end{array} \quad (\text{cond.})$$

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

$$\mathbb{E}[x_t | y_{1:t}] = \mu_{t|t} \quad \mathbb{E}[(x_t - \mu_{t|t}) :] = P_{t|t} \quad \textcircled{A12} \textcircled{2}$$

- 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, KF derivation slides
- continued next lecture