



STA 4103/5107

Computational Methods

in Statistics II

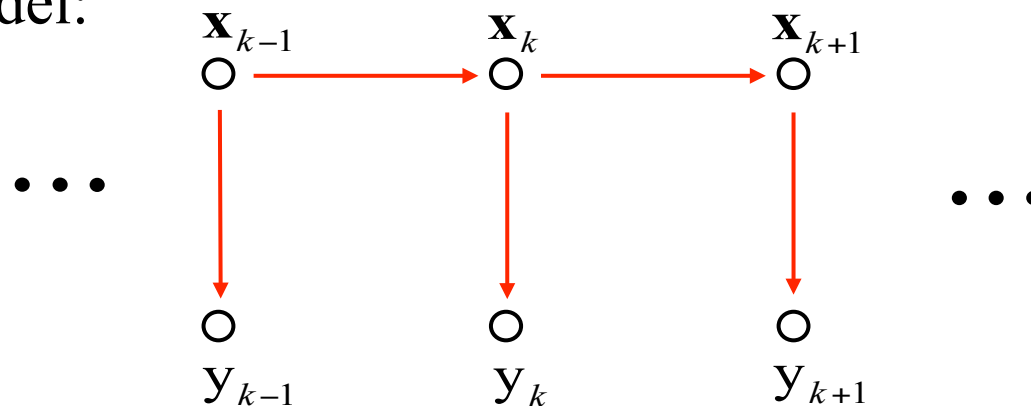
Department of Statistics
Florida State University

Class 11
February 14, 2017



Review: Nonlinear Filtering Problem

- State vector $x_k \in \mathbf{R}^d$, observation vector $y_k \in \mathbf{R}^c$:
- State Equation (prior): $x_k = F(x_{k-1}) + w_k$
 Observation equation (likelihood): $y_k = G(x_k) + q_k$
- Independence Assumptions:
 $f(x_k | x_1, \dots, x_{k-1}) = f(x_k | x_{k-1})$, $f(y_k | y_1, \dots, y_{k-1}, x_k) = f(y_k | x_k)$.
- Graphical model:





Review: Prediction and Update Equations

- **Goal: Estimate**

$$f(x_k | y_1, y_2, \dots, y_k)$$

- **Prediction equation**

$$\begin{aligned} f(x_k | y_1, y_2, \dots, y_{k-1}) &= \int_{x_{k-1}} f(x_k, x_{k-1} | y_1, y_2, \dots, y_{k-1}) dx_{k-1} \\ &= \int_{x_{k-1}} f(x_k | x_{k-1}) f(x_{k-1} | y_1, y_2, \dots, y_{k-1}) dx_{k-1} \end{aligned}$$

- **Update equation**

$$\begin{aligned} f(x_k | y_1, y_2, \dots, y_k) &= \frac{f(x_k, y_1, y_2, \dots, y_k)}{f(y_1, y_2, \dots, y_k)} \\ &= \frac{f(y_k | x_k) f(x_k | y_1, y_2, \dots, y_{k-1})}{f(y_k | y_1, y_2, \dots, y_{k-1})} \end{aligned}$$



Review: Notation in Kalman Filter

- $\mathbf{x}_k \in \mathbb{R}^d$: internal state at k th frame (hidden random variable, e.g. position of the object in the image).

$\mathbf{X}_k = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k]^T$: history up to time step k

- $y_k \in \mathbb{R}^c$: measurement at k th frame (observable random variable, e.g. the given image).

$\mathbf{Y}_k = [y_1, y_2, \dots, y_k]^T$: history up to time step k

- Goal:

Estimating the posterior probability $p(\mathbf{x}_k | \mathbf{Y}_k)$



Kalman Filter: Likelihood Model

- *Generative model* for the observation:

$$y_k = H_k x_k + q_k$$

$$H_k \in \mathbb{R}^{c \times d}, \quad q_k \sim N(0, Q_k), \quad Q_k \in \mathbb{R}^{c \times c}$$

- The likelihood model is equivalent to that

$$y_k | x_k \sim N(H_k x_k, Q_k)$$

- The conditional probability has explicit form:

$$p(y_k | x_k) = \frac{1}{((2\pi)^c \det(Q_k))^{1/2}} \exp\left(-\frac{1}{2}(y_k - H_k x_k)^T Q_k^{-1} (y_k - H_k x_k)\right)$$



Kalman Filter: Prior Model

- *Temporal prior* of the state:

$$\mathbf{x}_k = \mathbf{A}_k \mathbf{x}_{k-1} + \mathbf{w}_k$$

$$\mathbf{A}_k \in \mathbb{R}^{d \times d}, \quad \mathbf{w}_k \sim N(0, \mathbf{W}_k), \quad \mathbf{W}_k \in \mathbb{R}^{d \times d}$$

- The prior model is equivalent to that

$$\mathbf{x}_k | \mathbf{x}_{k-1} \sim N(\mathbf{A}_k \mathbf{x}_{k-1}, \mathbf{W}_k)$$

- The conditional probability has explicit form:

$$p(\mathbf{x}_k | \mathbf{x}_{k-1}) = \frac{1}{((2\pi)^d \det(\mathbf{W}_k))^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x}_k - \mathbf{A}_k \mathbf{x}_{k-1})^T \mathbf{W}_k^{-1}(\mathbf{x}_k - \mathbf{A}_k \mathbf{x}_{k-1})\right)$$



Kalman Filter Model

Definition:

System Equation:

$$\mathbf{x}_k = \mathbf{A}_k \mathbf{x}_{k-1} + \mathbf{w}_k, \quad \mathbf{w}_k \in N(0, \mathbf{W}_k)$$
$$k=2,3,\dots$$

Measurement Equation:

$$\mathbf{y}_k = \mathbf{H}_k \mathbf{x}_k + \mathbf{q}_k, \quad \mathbf{q}_k \in N(0, \mathbf{Q}_k)$$
$$k=1,2,\dots$$

Assumption:

All random variables have Gaussian distributions and they are linearly related.



Learning Kalman Model

- In practice, the parameters in the model need to be estimated from training data. (In training data, we know both hidden states and measurements.)
- Common simplification: A_k, H_k, W_k, Q_k are constant over time (independent of k).
- The A, H, W, Q can be estimated by maximizing the joint probability $p(X_M, Y_M)$. In fact,

$$\begin{aligned} p(X_M, Y_M) &= p(X_M) p(Y_M | X_M) \\ &= [p(x_1) \prod_{k=2}^M p(x_k | x_{k-1})] [\prod_{k=1}^M p(y_k | x_k)] \end{aligned}$$



Splitting the Joint Distribution

$$\arg \max_{A,W,H,Q} p(X_M, Y_M)$$

$$= \{\arg \max_{A,W} p(X_M), \arg \max_{H,Q} p(Y_M | X_M)\}$$

$$= \{\arg \min_{A,W} f(A,W), \arg \min_{H,Q} g(H,Q)\}$$

where

$$f(A,W) = -\alpha \log p(X_M) = \sum_{k=2}^M [\log(\det W) + (x_k - Ax_{k-1})^T W^{-1} (x_k - Ax_{k-1})],$$

$$g(H,Q) = -\beta \log p(Y_M | X_M) = \sum_{k=1}^M [\log(\det Q) + (y_k - Hx_k)^T Q^{-1} (y_k - Hx_k)].$$

How to optimize functions with matrix variables ?



Matrix Calculus

- Definition: assume $\mathbf{X} = (x_{ij})_{mn}$, then

$$d/d\mathbf{X} = \begin{pmatrix} d/dx_{11} & \cdots & d/dx_{1n} \\ \vdots & \ddots & \vdots \\ d/dx_{m1} & \cdots & d/dx_{mn} \end{pmatrix}$$

- Quadratic Products:

$$d/d\mathbf{X} ((\mathbf{X}\mathbf{a}+\mathbf{b})^T \mathbf{C} (\mathbf{X}\mathbf{a}+\mathbf{b})) = (\mathbf{C}+\mathbf{C}^T)(\mathbf{X}\mathbf{a}+\mathbf{b})\mathbf{a}^T$$

- Determinant: $d/d\mathbf{X} (\log(\det(\mathbf{X}))) = \mathbf{X}^{-T}$
- Inverse: $d/d\mathbf{X} (\mathbf{a}^T \mathbf{X}^{-1} \mathbf{b}) = -\mathbf{X}^{-T} \mathbf{a} \mathbf{b}^T \mathbf{X}^{-T}$

(upper case: matrix, lower case: column vector)

<http://www.ee.ic.ac.uk/hp/staff/dmb/matrix/calculus.html>



Detailed Steps

Show one example on prior probability:

$$f(A, W) = \sum_{k=2}^M [\log(\det W) + (x_k - Ax_{k-1})^T W^{-1} (x_k - Ax_{k-1})]$$

$$\text{i) } \frac{\partial}{\partial A} f(A, W) = \sum_{k=2}^M (W^{-1} + W^{-1})(x_k - Ax_{k-1})(-x_{k-1})^T = 0$$

$$\Rightarrow \sum_{k=2}^M x_k x_{k-1}^T = A \sum_{k=2}^M x_{k-1} x_{k-1}^T \Rightarrow A = \left(\sum_{k=2}^M x_k x_{k-1}^T \right) \left(\sum_{k=2}^M x_{k-1} x_{k-1}^T \right)^{-1}$$

$$\text{ii) } \frac{\partial}{\partial W} f(A, W) = \sum_{k=2}^M (W^{-1} - W^{-1}(x_k - Ax_{k-1})(x_k - Ax_{k-1})^T W^{-1}) = 0$$

$$\Rightarrow W = \sum_{k=2}^M (x_k - Ax_{k-1})(x_k - Ax_{k-1})^T / (M - 1)$$



Closed-form Solutions

$$\begin{aligned} A &= \left(\sum_{k=2}^M \mathbf{x}_k \mathbf{x}_{k-1}^T \right) \left(\sum_{k=2}^M \mathbf{x}_{k-1} \mathbf{x}_{k-1}^T \right)^{-1}, \\ W &= \frac{1}{M-1} \left(\sum_{k=2}^M \mathbf{x}_k \mathbf{x}_k^T - A \sum_{k=2}^M \mathbf{x}_{k-1} \mathbf{x}_k^T \right), \\ H &= \left(\sum_{k=1}^M y_k \mathbf{x}_k^T \right) \left(\sum_{k=1}^M \mathbf{x}_k \mathbf{x}_k^T \right)^{-1}, \\ Q &= \frac{1}{M} \left(\sum_{k=1}^M y_k y_k^T - H \sum_{k=1}^M \mathbf{x}_k y_k^T \right). \end{aligned}$$



Recursive Estimation

$$p(\mathbf{x}_k | \mathbf{Y}_k) = \kappa p(y_k | \mathbf{x}_k) \int p(\mathbf{x}_k | \mathbf{x}_{k-1}) p(\mathbf{x}_{k-1} | \mathbf{Y}_{k-1}) d\mathbf{x}_{k-1}$$

Time update:

posterior at previous step:

$$p(\mathbf{x}_{k-1} | \mathbf{Y}_{k-1})$$

temporal prior:

$$p(\mathbf{x}_k | \mathbf{x}_{k-1})$$

prior distribution:

$$p(\mathbf{x}_k | \mathbf{Y}_{k-1}) = \int p(\mathbf{x}_k | \mathbf{x}_{k-1}) p(\mathbf{x}_{k-1} | \mathbf{Y}_{k-1}) d\mathbf{x}_{k-1}$$

Measurement update:

prior distribution:

$$p(\mathbf{x}_k | \mathbf{Y}_{k-1})$$

likelihood:

$$p(y_k | \mathbf{x}_k)$$

posterior distribution:

$$p(\mathbf{x}_k | \mathbf{Y}_k) = \kappa p(y_k | \mathbf{x}_k) p(\mathbf{x}_k | \mathbf{Y}_{k-1})$$



Basic Properties of Normal Distributions

- $x, y \in \mathbb{R}^d$ independent random vectors,
 $x \sim N(0, A), y \sim N(0, B)$,
then, i) for any matrix $C \in \mathbb{R}^{d \times d}$, $Cx \sim N(0, CAC^T)$;
ii) $x + y \sim N(0, A + B)$.

These two properties can be derived by the following rules:

- i)
$$\text{Cov}(Cx) = C(\text{Cov}(x))C^T = CAC^T.$$
- ii)
$$\text{Cov}(x + y) = \text{Cov}(x) + \text{Cov}(y) = A + B.$$



Kalman Filtering, Step I: Time Update

Assume: $\mathbf{x}_{k-1} \mid \mathbf{Y}_{k-1} \sim N(\hat{\mathbf{x}}_{k-1}, \mathbf{P}_{k-1})$

$$\Leftrightarrow \mathbf{x}_{k-1} = \hat{\mathbf{x}}_{k-1} + \mathbf{e}_{k-1}, \quad \mathbf{e}_{k-1} \sim N(0, \mathbf{P}_{k-1})$$

System equation: $\mathbf{x}_k = \mathbf{A}_k \mathbf{x}_{k-1} + \mathbf{w}_k, \quad \mathbf{w}_k \sim N(0, \mathbf{W}_k)$

$$\Rightarrow \mathbf{x}_k = \mathbf{A}_k \hat{\mathbf{x}}_{k-1} + \mathbf{A}_k \mathbf{e}_{k-1} + \mathbf{w}_k$$

Use properties i) and ii):

$$\mathbf{A}_k \mathbf{e}_{k-1} + \mathbf{w}_k \sim N(0, \mathbf{A}_k \mathbf{P}_{k-1} \mathbf{A}_k^T + \mathbf{W}_k)$$

Let

$$\hat{\mathbf{x}}_k^- = \mathbf{A}_k \hat{\mathbf{x}}_{k-1}, \quad \mathbf{P}_k^- = \mathbf{A}_k \mathbf{P}_{k-1} \mathbf{A}_k^T + \mathbf{W}_k,$$

then,

$$\mathbf{x}_k \mid \mathbf{Y}_{k-1} \sim N(\hat{\mathbf{x}}_k^-, \mathbf{P}_k^-)$$



Step II: Measurement Update

Time update:

$$p(\mathbf{x}_k | \mathbf{Y}_{k-1}) \propto \exp(-(\mathbf{x}_k - \hat{\mathbf{x}}_k^-)^T (\mathbf{P}_k^-)^{-1} (\mathbf{x}_k - \hat{\mathbf{x}}_k^-) / 2)$$

Measurement equation:

$$p(y_k | \mathbf{x}_k) \propto \exp(-(\mathbf{y}_k - \mathbf{H}_k \mathbf{x}_k)^T \mathbf{Q}_k^{-1} (\mathbf{y}_k - \mathbf{H}_k \mathbf{x}_k) / 2)$$

Recursive update:

$$\begin{aligned} p(\mathbf{x}_k | \mathbf{Y}_k) &\propto p(y_k | \mathbf{x}_k) p(\mathbf{x}_k | \mathbf{Y}_{k-1}) \\ &\propto \exp\left(-\frac{1}{2} (\mathbf{x}_k - \hat{\mathbf{x}}_k)^T (\mathbf{P}_k)^{-1} (\mathbf{x}_k - \hat{\mathbf{x}}_k)\right) \end{aligned}$$

(details omitted)

where, $\hat{\mathbf{x}}_k = \hat{\mathbf{x}}_k^- + \mathbf{K}_k (\mathbf{y}_k - \mathbf{H}_k \hat{\mathbf{x}}_k^-)$, $\mathbf{P}_k = (\mathbf{I} - \mathbf{K}_k \mathbf{H}_k) \mathbf{P}_k^-$,

$$\mathbf{K}_k = \mathbf{P}_k^- \mathbf{H}_k^T (\mathbf{H}_k \mathbf{P}_k^- \mathbf{H}_k^T + \mathbf{Q}_k)^{-1}.$$

That is,

$$\mathbf{x}_k | \mathbf{Y}_k \sim N(\hat{\mathbf{x}}_k, \mathbf{P}_k)$$



Kalman Filter Algorithm

Time Update

Prior estimate

$$\hat{\mathbf{x}}_k^- = \mathbf{A}_k \hat{\mathbf{x}}_{k-1}$$

Error covariance

$$\mathbf{P}_k^- = \mathbf{A}_k \mathbf{P}_{k-1} \mathbf{A}_k^T + \mathbf{W}_k$$

Measurement Update

Posterior estimate

$$\hat{\mathbf{x}}_k = \hat{\mathbf{x}}_k^- + \mathbf{K}_k (\mathbf{y}_k - \mathbf{H}_k \hat{\mathbf{x}}_k^-)$$

Error covariance

$$\mathbf{P}_k = (\mathbf{I} - \mathbf{K}_k \mathbf{H}_k) \mathbf{P}_k^-$$

Kalman gain

$$\mathbf{K}_k = \mathbf{P}_k^- \mathbf{H}_k^T (\mathbf{H}_k \mathbf{P}_k^- \mathbf{H}_k^T + \mathbf{Q}_k)^{-1}$$

previous estimate of $\hat{\mathbf{x}}_{k-1}$ and \mathbf{P}_{k-1}

(initially: we can let $\hat{\mathbf{x}}_1 = 0, \mathbf{P}_1 = 0$)

Welch & Bishop, *An Introduction to the Kalman Filter*, 2006



Estimation Accuracy

- R^2 Error is commonly-used as a criterion to measure the estimation accuracy:

Let x_k denote the true state and \hat{x}_k denote the estimate. Then

$$R^2 = 1 - \frac{\sum_k \|x_k - \hat{x}_k\|^2}{\sum_k \|x_k - \bar{x}\|^2}$$

- R^2 Error can also be measured component-wise. That is, for the i -th component, we have

$$R_i^2 = 1 - \frac{\sum_k (x_{k,i} - \hat{x}_{k,i})^2}{\sum_k (x_{k,i} - \bar{x}_{\cdot,i})^2}$$