

LECTURE 7

In the approximated NP-PP problem:

$$Y_i \sim \mathcal{P}(\lambda(\hat{t}_i | \mathcal{H}_{t_i}, \theta) \Delta t)$$

$$\log \lambda(t | \mathcal{H}_t, \theta) = g(t, \mathcal{H}_t, \theta)$$

$$i = 1, 2, 3, \dots, m$$

the parameter vector θ is time-invariant and must be estimated from the entire realization (Y_1, Y_2, \dots, Y_m) of the sequence of RVs $Y_1, Y_2, \dots, Y_m \Rightarrow$ What if $\theta = \theta_t$, i.e., the parameter vector is a function of time and perhaps past history? The case is pertinent since λ is a function of t and $\mathcal{H}_t \Rightarrow$ We consider the new problem:

$$Y_i \sim \mathcal{P}(\lambda(\hat{t}_i | \mathcal{H}_{t_i}^i, \theta_i) \Delta t)$$

$$\log \lambda(t | \mathcal{H}_t, \theta_t) = g(t, \mathcal{H}_t, \theta_t)$$

$$\theta_i \triangleq \theta_{t_i} \quad i = 1, 2, 3, \dots, m$$

In this case, we have one new vector θ_i for each new observation $y_i \Rightarrow$ A ML method could not be viable \Rightarrow An efficient and recursive solution can be obtained under the additional assumptions:

- i) $\theta_{i+1} = F_i \theta_i + \eta_i \quad i = 0, 1, 2, \dots, m$
 where: F_i - known matrix $\forall i$
 $\eta_i \sim N(0, Q_i) \quad \forall i$ - white noise
 $Q_i \triangleq$ covariance matrix

- ii) $\theta_0 \triangleq$ initial condition, $\theta_0 \sim N(\mu_0, W_0)$
 μ_0, W_0 - known

Assumption (i) means that there is an evolution model for $\theta_i \quad \forall i$

\Rightarrow

Assumptions (i) and (ii) mean that each vector θ_i is a Gaussian random vector with its own mean and covariance matrix

②

Since ϑ_i is a random vector $\forall i$, we can consider the conditional probability $P(\vartheta_i | Y_i, \mathcal{H}_{\hat{t}_i})$ and the correspondent pdf: $f_{\vartheta_i}(\vartheta_i | Y_i, \mathcal{H}_{\hat{t}_i})$.

For sake of simplicity, let us use the notation: $\mathcal{H}_i = \mathcal{H}_{\hat{t}_i}$ and $f_{\vartheta} = f_{\vartheta_i}$ (the latter makes sense since the random vectors ϑ_i are all Gaussian). By applying the Bayes' theorem, we have:

$$P(\vartheta_i | Y_i, \mathcal{H}_i) = \frac{P(\vartheta_i, Y_i, \mathcal{H}_i)}{P(Y_i, \mathcal{H}_i)} = \frac{P(Y_i | \mathcal{H}_i, \vartheta_i) P(\vartheta_i | \mathcal{H}_i)}{P(Y_i | \mathcal{H}_i)}$$

We can note:

- $P(Y_i | \mathcal{H}_i, \vartheta_i) \cong \exp \left\{ y_i \log(\lambda(\hat{t}_i | \mathcal{H}_i, \vartheta_i) \Delta t) - \lambda(\hat{t}_i | \mathcal{H}_i, \vartheta_i) \Delta t \right\}$
↑
Formulation of the NP-PP problem

- $f_{\vartheta}(\vartheta_i | Y_i, \mathcal{H}_i) = \frac{P(Y_i | \mathcal{H}_i, \vartheta_i)}{P(Y_i | \mathcal{H}_i)} \cdot f_{\vartheta}(\vartheta_i | \mathcal{H}_i) \quad (*)$

- $P(Y_i | \mathcal{H}_i)$ is a normalization factor that assures the integral of f_{ϑ} over ϑ_i is equal to 1

- $f_{\vartheta}(\vartheta_i | \mathcal{H}_i) = \int f(\vartheta_i, \vartheta_{i-1} | \mathcal{H}_i) d\vartheta_{i-1} = \int f_{\vartheta}(\vartheta_i | \vartheta_{i-1}, \mathcal{H}_i) f_{\vartheta}(\vartheta_{i-1} | \mathcal{H}_i) d\vartheta_{i-1}$
↑
joint-probability function

$$= \int f_{\vartheta}(\vartheta_i | \vartheta_{i-1}, \mathcal{H}_i) f_{\vartheta}(\vartheta_{i-1} | Y_{i-1}, \mathcal{H}_{i-1}) d\vartheta_{i-1}$$

$f_{\theta}(\vartheta_i | Y_i, \mathcal{H}_i) \triangleq$ Posterior density of ϑ_i at time t_i
 $f_{\theta}(\vartheta_i | \mathcal{H}_i) \triangleq$ Prediction density of ϑ_i at time t_i with observations up to t_{i-1}

$$f_{\theta}(\vartheta_i | Y_i, \mathcal{H}_i) \propto P(Y_i | \mathcal{H}_i, \vartheta_i) \int f_{\theta}(\vartheta_i | \vartheta_{i-1}, \mathcal{H}_i) f_{\theta}(\vartheta_{i-1} | Y_{i-1}, \mathcal{H}_{i-1}) d\vartheta_{i-1}$$

\uparrow posterior density at time t_i \uparrow It is related to the prediction density \uparrow posterior density at time t_{i-1}

Assumptions (i) and (ii) also determine the following facts:

a) $P(\vartheta_i | \vartheta_{i-1}, \mathcal{H}_i) = P(F_i \vartheta_{i-1} + \eta_i | \vartheta_{i-1}, \mathcal{H}_i) = P(\eta_i)$ since η_i is white noise

\uparrow evolution model \uparrow for any given value $\vartheta_{i-1}, \mathcal{H}_i$

$\Rightarrow f_{\theta}(\vartheta_i | \vartheta_{i-1}, \mathcal{H}_i)$ - Gaussian function of the var.: $\vartheta_i - F_i \vartheta_{i-1}$

b) Because of (a), if $f_{\theta}(\vartheta_{i-1} | Y_{i-1}, \mathcal{H}_{i-1})$ is Gaussian, then $f_{\theta}(\vartheta_i | Y_i, \mathcal{H}_i)$ is Gaussian too, i.e., the posterior density function preserves its characteristics in time \Rightarrow It is therefore plausible to focus on the quantities:

$$\vartheta_{i|i-1} \triangleq E_{\vartheta}[\vartheta_i | \mathcal{H}_i] \quad W_{i|i-1} \triangleq E_{\vartheta}[(\vartheta_i - \vartheta_{i|i-1})^2 | \mathcal{H}_i] = \text{var}[\vartheta_i | \mathcal{H}_i]$$

$$\vartheta_{i|i} \triangleq E_{\vartheta}[\vartheta_i | Y_i, \mathcal{H}_i] \quad W_{i|i} \triangleq E_{\vartheta}[(\vartheta_i - \vartheta_{i|i})^2 | Y_i, \mathcal{H}_i] = \text{var}[\vartheta_i | Y_i, \mathcal{H}_i]$$

Note: - $\vartheta_{i|i-1}$ and $W_{i|i-1}$ are mean and covariance matrix of the one-step prediction density $f_{\theta}(\vartheta_i | \mathcal{H}_i)$
 - $\vartheta_{i|i}$ and $W_{i|i}$ are mean and covariance matrix of the posterior density function $f_{\theta}(\vartheta_i | Y_i, \mathcal{H}_i)$

④

Because of the definition, we have:

$$1) \quad \theta_{i|i-1} = E[\theta_i | \mathcal{H}_i] = F_i E[\theta_{i-1} | \mathcal{H}_i] = F_i \theta_{i-1|i-1}$$

$$2) \quad W_{i|i-1} = \text{var}[\theta_i | \mathcal{H}_i] = \text{var}[F_i \theta_{i-1} + \eta_i | \mathcal{H}_i] = \text{var}[F_i \theta_{i-1} | \mathcal{H}_i] + \text{var}[\eta_i]$$

↑
since η_i is white noise

$$= F_i W_{i-1|i-1} F_i^T + Q_i$$

Because of assumptions (i) and (ii) we also have:

$$f_{\theta}(\theta_i | \mathcal{H}_i) = \int f_{\theta}(\theta_i | \theta_{i-1}, \mathcal{H}_i) f_{\theta}(\theta_{i-1} | \mathcal{Y}_{i-1}, \mathcal{H}_{i-1}) d\theta_{i-1}$$

$$= f(\eta_i) \int f_{\theta}(\theta_{i-1} | \mathcal{Y}_{i-1}, \mathcal{H}_{i-1}) d\theta_{i-1} = f(\eta_i) - \text{Gaussian function}$$

Therefore, if we replace $\eta_i = \theta_i - F_i \theta_{i-1}$ with the prediction error $\varepsilon_i \triangleq \theta_i - \theta_{i|i-1}$ we have:

$$f_{\theta}(\theta_i | \mathcal{H}_i) \cong \frac{1}{\sqrt{(2\pi)^r \det(W_{i|i-1})}} \exp \left\{ -\frac{1}{2} (\theta_i - \theta_{i|i-1})^T W_{i|i-1}^{-1} (\theta_i - \theta_{i|i-1}) \right\}$$

where $r \triangleq$ size of vector θ_i

Therefore, the conditional pdf $f_{\theta}(\theta_i | \mathcal{Y}_i, \mathcal{H}_i)$ in (*) can be approximated with a function that is proportional to:

$$\exp \left\{ \gamma_i \log(\lambda(\varepsilon_i | \mathcal{H}_i, \theta_i) \Delta t) - \lambda(\varepsilon_i | \mathcal{H}_i, \theta_i) \Delta t + \right. \\ \left. - \frac{1}{2} (\theta_i - \theta_{i|i-1})^T W_{i|i-1}^{-1} (\theta_i - \theta_{i|i-1}) \right\} \quad (**)$$

(5)

On the other hand, because of (b), $f_{\theta}(\theta_i | Y_i, \mathcal{H}_i)$ is Gaussian too, and therefore it should be:

$$f_{\theta}(\theta_i | Y_i, \mathcal{H}_i) \propto \exp \left\{ -\frac{1}{2} (\theta_i - \theta_{i|c})^T W_{i|c}^{-1} (\theta_i - \theta_{i|c}) \right\} \quad (***)$$

Beside a constant scaling factor to be calculated, we can therefore combine (**) and (***) and write:

$$\exp \left\{ y_i \log (\lambda(\hat{t}_i | \mathcal{H}_i, \theta_i) \Delta t) - \lambda(\hat{t}_i | \mathcal{H}_i, \theta_i) \Delta t - \frac{1}{2} (\theta_i - \theta_{i|c-1})^T W_{i|c-1}^{-1} (\theta_i - \theta_{i|c-1}) \right\}$$

$$= A \exp \left\{ -\frac{1}{2} (\theta_i - \theta_{i|c})^T W_{i|c}^{-1} (\theta_i - \theta_{i|c}) \right\}$$

coefficient
to be determined

↓

$$y_i \log (\lambda(\hat{t}_i | \mathcal{H}_i, \theta_i) \Delta t) - \lambda(\hat{t}_i | \mathcal{H}_i, \theta_i) \Delta t - \frac{1}{2} (\theta_i - \theta_{i|c-1})^T W_{i|c-1}^{-1} (\theta_i - \theta_{i|c-1})$$

$$= -\frac{1}{2} (\theta_i - \theta_{i|c})^T W_{i|c}^{-1} (\theta_i - \theta_{i|c}) + \log A$$

↓ By differentiating
with respect to θ_i

$$c) \quad W_{i|c}^{-1} (\theta_i - \theta_{i|c}) = W_{i|c-1}^{-1} (\theta_i - \theta_{i|c-1}) - \left(\frac{y_i}{\lambda(\hat{t}_i | \mathcal{H}_i, \theta_i)} - \Delta t \right) \nabla_{\theta} \lambda(\hat{t}_i | \mathcal{H}_i, \theta_i)$$

$$\text{where: } \nabla_{\theta} \lambda \triangleq \left[\frac{\partial \lambda}{\partial \theta_{i,1}}, \frac{\partial \lambda}{\partial \theta_{i,2}}, \dots, \frac{\partial \lambda}{\partial \theta_{i,r}} \right]^T - \text{gradient and } \theta_i \triangleq [\theta_{i,1}, \theta_{i,2}, \dots, \theta_{i,r}]^T$$

Note also the relationship that is usually used:

$$\nabla_{\theta} \log \lambda(\hat{t}_i | \mathcal{H}_i, \theta_i) (y_i - \lambda(\hat{t}_i | \mathcal{H}_i, \theta_i) \Delta t) = \left(\frac{y_i}{\lambda(\hat{t}_i | \mathcal{H}_i, \theta_i)} - \Delta t \right) \nabla_{\theta} \lambda(\hat{t}_i | \mathcal{H}_i, \theta_i)$$

6

Under assumptions (i), (ii), and (b), formula (c) should hold $\forall \theta_i \Rightarrow$ By choosing $\theta_i = \theta_{i|i-1}$ we have:

$$W_{i|i}^{-1} (\theta_{i|i-1} - \theta_{i|i}) = -\nabla_{\theta} \left(\log \lambda(\hat{t}_i | \mathcal{H}_i, \theta_{i|i-1}) \right) (y_i - \lambda(\hat{t}_i | \mathcal{H}_i, \theta_{i|i-1}) \Delta t)$$

\Downarrow

$$3) \quad \theta_{i|i} = \theta_{i|i-1} + W_{i|i} \nabla_{\theta} \left(\log \lambda(\hat{t}_i | \mathcal{H}_i, \theta_{i|i-1}) \right) (y_i - \lambda(\hat{t}_i | \mathcal{H}_i, \theta_{i|i-1}) \Delta t)$$

Moreover, if we differentiate (c) with respect to θ_i :

$$\begin{aligned} d) \quad W_{i|i}^{-1} &= W_{i|i-1}^{-1} - \nabla_{\theta} \left\{ \nabla_{\theta} \left(\log \lambda(\hat{t}_i | \mathcal{H}_i, \theta_i) \right) (y_i - \lambda(\hat{t}_i | \mathcal{H}_i, \theta_i) \Delta t) \right\} \\ &= W_{i|i-1}^{-1} - H_{\theta} \log \lambda (y_i - \lambda \Delta t) + (\nabla_{\theta} \log \lambda)^T (\nabla_{\theta} \log \lambda) \lambda \Delta t \end{aligned}$$

where $H_{\theta} \log \lambda \triangleq$ Hessian of $\log \lambda(\hat{t}_i | \mathcal{H}_i, \theta_i)$ - Evaluation of formula (d) for $\theta_i = \theta_{i|i-1}$ leads to the equation:

$$4) \quad W_{i|i}^{-1} = W_{i|i-1}^{-1} - \left[H_{\theta} \log \lambda (y_i - \lambda \Delta t) + (\nabla_{\theta} \log \lambda)^T (\nabla_{\theta} \log \lambda) (\lambda \Delta t) \right]_{\theta_i = \theta_{i|i-1}}$$

The combination of equations (1)-(4) defines an adaptive filter that recursively updates the estimation of the expected value of parameter vector θ_i at each time step t_i :

$$\left. \begin{aligned} (1) + (2) \\ \theta_{i|i-1} &= F_i \theta_{i-1|i-1} \\ W_{i|i-1} &= F_i W_{i-1|i-1} F_i^T + Q_i \end{aligned} \right\} \begin{array}{l} \text{Prediction} \\ \text{Step} \end{array}$$

(3) + (4)

$$\left. \begin{aligned} W_{i|i}^{-1} &= W_{i|i-1}^{-1} + I(\theta_{i|i-1}) \\ \theta_{i|i} &= \theta_{i|i-1} + W_{i|i} \nabla_{\theta}(\log \lambda)(y_i - \lambda \Delta t) \end{aligned} \right\} \begin{array}{l} \text{Correction} \\ \text{Step} \end{array}$$

$$\text{where: } I(\theta_{i|i-1}) \triangleq -H_0 \log \lambda (y_i - \lambda \Delta t) + (\nabla_{\theta} \log \lambda)^T (\nabla_{\theta} \log \lambda) (\lambda \Delta t)$$

$$\lambda \triangleq \lambda(\hat{t}_i | \mathcal{H}_i, \theta_{i|i-1})$$

Note: $(y_i - \lambda \Delta t)$ provides the difference between the probability of having an event in $(t_{i-1}, t_i]$ given the model and whether the event actually occurs \Rightarrow It is the INNOVATION of the filter

It is the error signal used to update the model parameters

The vector $W_{i|i} \nabla_{\theta}(\log \lambda)$ can be seen as a learning gain in eq (3) \Rightarrow It is the equivalent of the Kalman gain for standard Kalman filters and it depends on the statistical properties of θ_i and the point process

The adaptive filter (1) - (4) is relevant for the following reasons:

* If $\eta_i = 0 \forall i$ (i.e., θ_i is a deterministic process) $\Rightarrow W_{i|i-1} = F_i W_{i-1|i-1} F_i^T \Rightarrow$

$W_{i|i}^{-1} = (F_i W_{i-1|i-1} F_i^T)^{-1} + I(\theta_{i|i-1}) \Rightarrow$ The filter is analogous to the recursive least-square filter (RLS) for deterministic models

$\{\theta_i\}$ is a Gaussian Random Walk

* If $F_i = I_r \forall i$ (identity matrix) $\left. \begin{array}{l} \\ W_{i|i} = \Sigma - \text{constant} \end{array} \right\} \Rightarrow \theta_{i|i} = \theta_{i-1|i-1} + \Sigma \nabla_{\theta}(\log \lambda)(y_i - \lambda \Delta t)$

with $\lambda = \lambda(\hat{t}_i | \mathcal{H}_i, \theta_{i-1|i-1})$

8

⇒ The filter is analogous to the steepest descent formula

* The formulation of the filter (1)-(4) exploits the fact that θ_i $i=1,2,3,\dots,m$ are (Gaussian) RVs and that an evolution model is given (assumptions i-ii)

↓

θ_i could be replaced by a vector of explanatory variables $X_i \triangleq [X_{1,i} X_{2,i} \dots X_{q,i}]^T$ and the value of X_i at each time step could be recursively estimated from the realization of the NP-PP

↓

The filter (1)-(4) can be used for solving the following problem:

? $(x_1, x_2, \dots, x_m) : Y_i \sim P(\lambda(\hat{\theta}_i | \mathcal{H}_i, X_i, \theta) \Delta t)$

DECODING
PROBLEM

$$\log \lambda(t | \mathcal{H}_t, X, \theta) = g(t, \mathcal{H}_t, X, \theta)$$

$$X_i = x_i \quad i=1,2,3,\dots,m$$



Ex.: Let us assume that: $\log \lambda = \alpha_0 + \left\{ \sum_{j=1}^k \beta_j \Delta N_{(t-j\Delta t, t-(j-1)\Delta t)} \right\} x$

(scalar case)

↓

$$\nabla_x \log \lambda = \sum_{j=1}^k \beta_j \Delta N_{(t-j\Delta t, t-(j-1)\Delta t)} - \text{const. for a given time } t$$

Let us have: $x_i = x_{i-1} + \eta_i \Rightarrow x_{i|i-1} = x_{i-1|i-1}$

$$\eta_i \sim N(0, \sigma^2) \Rightarrow W_{i|i-1} = W_{i-1|i-1} + \sigma^2$$

$$W_{i|i}^{-1} = W_{i-1|i-1}^{-1} + (\nabla_x \log \lambda)^T \lambda (\hat{t}_i | \mathcal{H}_i, x_{i-1|i-1}, \vartheta) \Delta t$$

$$x_{i|i} = x_{i-1|i-1} + W_{i|i} (\nabla_x \log \lambda) (y_i - \lambda (\hat{t}_i | \mathcal{H}_i, x_{i-1|i-1}, \vartheta) \Delta t)$$

where $\vartheta \triangleq [\alpha_0 \beta_1 \beta_2 \dots \beta_k]^T$

□

Ex.: Let us assume that: $\log \lambda = \alpha_0 + \sum_{j=1}^k x_{j,t} \Delta N_{(t-j\Delta t, t-(j-1)\Delta t)}$

(vector, log-linear case)

$$X_i = [x_{1,i} \ x_{2,i} \ \dots \ x_{k,i}]^T$$

$$\nabla_x \log \lambda = [\Delta N_{(t-\Delta t, t)} \ \Delta N_{(t-2\Delta t, t-\Delta t)} \ \dots]^T \triangleq \Gamma_i$$

$$(\nabla_x \log \lambda)^T (\nabla_x \log \lambda) = \begin{bmatrix} \Delta N_{(t-\Delta t, t)}^2 & \Delta N_{(t-\Delta t, t)} \Delta N_{(t-2\Delta t, t-\Delta t)} & \dots \\ \vdots & \ddots & \dots \\ \dots & \dots & \dots \end{bmatrix}$$

This matrix is updated at each time step t_i based on the previous history

→ $\underbrace{\hspace{15em}}_{\Lambda_i}$

Let us assume: $X_i = X_{i-1} + \eta_i \Rightarrow X_{i|i-1} = X_{i-1|i-1}$

$\eta_i \sim N(0, \Sigma_i) \Rightarrow W_{i|i-1} = W_{i-1|i-1} + \Sigma_i$

$W_{i|i}^{-1} = W_{i-1|i-1}^{-1} + \Lambda_i \lambda_i \Delta t$

$X_{i|i} = X_{i-1|i-1} + W_{i|i} \Gamma_i (y_i - \lambda_i \Delta t)$

where $\lambda_i \triangleq e^{\alpha_0} \cdot e^{\eta_i^T X_{i-1|i-1}}$

□

Ex.: Let us assume that $X_{t+\tau}$ is the vector of kinematic variables (e.g., position and velocity variables) of a point in the Cartesian 3D space at time $t+\tau$, with $\tau > 0$ - fixed and known

(kinematics decoding)

(10)

Let us also suppose that the vector X_t is an explanatory variable for $M > 1$ point processes (e.g., M spike trains from individual neurons that are collected simultaneously), each one with its own CIF: $\lambda_w(t | \mathcal{H}_t^w, X_{t+\tau}, \theta^w)$

$\theta^w \triangleq$ parameter vector for the w -th point process

$\mathcal{H}_t^w \triangleq$ history up to time t for the w -th point process

In this case, we have to replace the factor $P(Y_i | \mathcal{H}_i, \theta_i)$ in the original formula (*) with the joint probability distribution:

$$P(Y_{i,1} = y_{i,1}, Y_{i,2} = y_{i,2}, \dots, Y_{i,M} = y_{i,M} | \mathcal{H}_i^1, \mathcal{H}_i^2, \dots, \mathcal{H}_i^M, X_{\hat{t}_i+\tau}, \theta^1, \theta^2, \dots, \theta^M)$$

where $Y_{i,w}$ is the Bernoulli RV associated with the w -th point process in the interval $(t_{i-1}, t_i]$, $w=1, 2, 3, \dots, M \Rightarrow$ Because dependencies are encompassed in the histories \mathcal{H}_t^w , $w=1, 2, 3, \dots, M$, we can write:

$$P(Y_{i,1} = y_{i,1}, Y_{i,2} = y_{i,2}, \dots) = \prod_{w=1}^M P(Y_{i,w} | \mathcal{H}_i^w, X_{\hat{t}_i+\tau}, \theta^w)$$

$$\cong \exp \left\{ \sum_{w=1}^M y_{i,w} \log(\lambda_w(\hat{t}_i | \mathcal{H}_i^w, X_{\hat{t}_i+\tau}, \theta^w) \Delta t) - \sum_{w=1}^M \lambda_w(\hat{t}_i | \mathcal{H}_i^w, X_{\hat{t}_i+\tau}, \theta^w) \Delta t \right\} \quad (\nabla \Delta)$$

Also, if $X_{t+\tau}$ is a continuous process (e.g., representation of kinematic variables), it is usually common to assume:

$$X_{t+\tau} = \mu_x + F X_{t+\tau-1} + \eta_{t+\tau} \quad (\nabla \Delta \nabla)$$

where μ_x provides the average position and velocity and F is time-invariant \Rightarrow By replacing $(\nabla \Delta)$ in $(**)$ and using $(\nabla \Delta T)$ we have the multivariate version of the filter:

$$1') \quad \tilde{X}_{i|i-1} = \mu_x + F \tilde{X}_{i-1|i-1}$$

$$2') \quad W_{i|i-1} = F W_{i-1|i-1} F^T + Q_i$$

$$3') \quad \tilde{X}_{i|i} = \tilde{X}_{i|i-1} + W_{i|i} \sum_{w=1}^M \nabla_x (\log \lambda_w) (y_{i,w} - \lambda_w \Delta t)$$

$$4') \quad W_{i|i}^{-1} = W_{i|i-1}^{-1} - \sum_{w=1}^M H_x \log \lambda_w (y_{i,w} - \lambda_w \Delta t) + \sum_{w=1}^M (\nabla_x \log \lambda_w)^T (\nabla_x \log \lambda_w) \lambda_w \Delta t$$

where: $\tilde{X}_i \triangleq X_{\hat{t}_i + \tau} \Rightarrow \tilde{X}_{i|i-1} = X_{\hat{t}_i + \tau / \hat{t}_i + \tau - 1}$
 $\tilde{X}_{i|i} = X_{\hat{t}_i + \tau / \hat{t}_i + \tau}$

$$\lambda_w \triangleq \lambda_w(\hat{t}_i | \mathcal{H}_i^w, \tilde{X}_{i|i-1}, \theta^w)$$

$H_x \triangleq$ hessian matrix with respect to the components of $X_{t+\tau}$

$\nabla_x \triangleq$ gradient with respect to the components of $X_{t+\tau}$

□

* General Properties of the Adaptive Filter

Filters (1)-(4) and (1')-(4') and their variations share a few key properties:

- The original problem is formulated as:

$$\left. \begin{aligned} \gamma_i &\sim \mathcal{P}(\gamma_i | \theta_i) \\ \theta_i &\sim \mathcal{N}(\mu_i, \Sigma_i) \end{aligned} \right\} \rightarrow$$

The first RV is conditioned on the second RV, which is conditioned on parameters μ_i, Σ_i coming from the evolution model

where γ_i is obtained from combining $\lambda(\cdot)$ and $g(\cdot)$

12

⇒ It is a Conditional HIERARCHICAL Model and it requires the calculation of $f_Y(\eta | \theta)$ (⇒ conditional pdf given by the CIF in our case) and $f_\theta(\theta | \mu, \Sigma)$ (⇒ we implicitly used when we computed the prediction steps)

- The original problem has an observable RV (i.e., Y_i) and a latent RV (i.e., θ_i) which cannot be observed and whose value must be inferred from the value of the observable RV ⇒ It is a LATENT VARIABLE Estimation problem
- The adaptive filter can be envisioned as a generalization of the Kalman filter to the Point Processes (but with important differences too):

	<u>Kalman Filter</u>	<u>Adaptive Filter</u>
Evolution Model	$\theta_i = F\theta_{i-1} + \eta_i$ $Y_i = B\theta_i + \varepsilon_i$ $\eta_i \sim N(0, Q) \quad \varepsilon_i \sim N(0, R)$ $\theta_0 \sim N(\mu_0, Q) \quad \forall i$	$\theta_i = F_i\theta_{i-1} + \eta_i$ $Y_i \sim \mathcal{P}(\lambda(\xi_i H_i, \theta_i) \Delta t)$ $\eta_i \sim N(0, Q_i)$ $\theta_0 \sim N(\mu_0, Q_0)$
Prediction Step	$\theta_{i i-1} = F\theta_{i-1 i-1}$ $W_{i i-1} = FW_{i-1 i-1}F^T + Q$	$\theta_{i i-1} = F_i\theta_{i-1 i-1}$ $W_{i i-1} = F_iW_{i-1 i-1}F_i^T + Q_i$
Correction Step	$W_{i i} = W_{i i-1} + W_{i i-1}B^T \dots$ $(BW_{i i-1}B^T + R)^{-1}BW_{i i-1}$	$W_{i i}^{-1} = W_{i i-1}^{-1} + I(\theta_{i i-1})$
	$\theta_{i i} = \theta_{i i-1} + W_{i i-1}B^T \dots$ $(BW_{i i-1}B^T + R)^{-1}(Y_i - B\theta_{i i-1})$	$\theta_{i i} = \theta_{i i-1} + W_{i i} \nabla_{\theta} (\log \lambda)$ $(Y_i - \lambda \Delta t)$
	<p>↑ Error signals used to drive the estimation ↓</p>	

- The formulation of the filter is based on formula (*):

$$f_{\theta}(\theta_i | Y_i, \mathcal{H}_i) = \frac{P(Y_i | \mathcal{H}_i, \theta_i) f_{\theta}(\theta_i | \mathcal{H}_i)}{P(Y_i | \mathcal{H}_i)}$$

where $P(Y_i | \mathcal{H}_i, *)$ substitutes for $f_Y(Y_i | \mathcal{H}_i, *)$ given the discrete nature of $Y_i \Rightarrow$ Let us drop the dependency on \mathcal{H}_i and note:

$f_{\theta}(\theta_i | Y_i) = f_{\theta|Y}(\theta_i | Y_i)$ - Posterior density of θ_i conditioned on Y_i

$f_Y(Y_i | \theta_i) = f_{Y|\theta}(\theta_i | Y_i)$ - Prediction density of the observable variable Y_i

$f_Y(Y_i) = \int f_{Y|\theta}(Y_i | \theta) f_{\theta}(\theta) d\theta$ - Marginal probability density of Y_i

$f_{\theta}(\theta_i) = \int f_{\theta|\theta_{i-1}}(\theta_i | \theta_{i-1}) f_{\theta}(\theta_{i-1}) d\theta_{i-1}$ - We expressed the a priori distribution of θ_i at each step in term of the a priori distribution at the previous step and the prediction density

↓

The original problem is a Bayesian Estimation problem where:

- likelihood function of the point process determines the prediction density of Y_i
- evolution model and prediction density of θ_i determine the a priori density of θ_i
- evolution model provides a formula for the posterior density of θ_i

- In general, the formula:

$$f_{\theta|Y}(\theta | Y) = \frac{f_{Y|\theta}(Y | \theta) f_{\theta}(\theta)}{\int f_{Y|\theta}(Y | \theta) f_{\theta}(\theta) d\theta}$$

can be used to estimate the posterior density regardless of whether $\theta \sim N(\mu, \Sigma)$ for any pair (μ, Σ) or not

For instance, we can have:

(14)

$$f_{\theta}(\theta) = 1 \text{ (Uniform)} \\ f_{Y|\theta}(Y|\theta) = \theta^Y (1-\theta)^{n-Y} \text{ (Binomial)} \Rightarrow f_{\theta|Y}(\theta|Y) = \frac{\theta^Y (1-\theta)^{n-Y}}{\int \theta^Y (1-\theta)^{n-Y} d\theta}$$

$$\Rightarrow f_{\theta|Y}(\theta|Y) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} \text{ - Beta distribution } B(\alpha, \beta)$$

$$\text{with: } \alpha \triangleq Y+1 \quad \beta \triangleq n-Y+1 \quad \Gamma(n) \triangleq \int_0^{\infty} x^{n-1} e^{-x} dx$$

Because $U(0,1) = B(1,1)$ we have: $f_{\theta}(\theta) \sim B \Rightarrow f_{\theta|Y} \sim B$ - The a priori distribution is called a "conjugate" prior distribution. Because the same happens in case of Gaussian priors, we have that the adaptive filter is developed by exploiting conjugate priors. \square

References:

Eden et al. (2004), Neural Comput., vol. 16, pp. 971-998 \Rightarrow A copy is on Husky CT

Truccolo et al. (2005), J. Neurophysiol., vol. 93, pp. 1074-89 \Rightarrow A copy is on Husky CT

Textbook: ch. 16 (sections. 16.1 - 16.1.1 - 16.1.2 - 16.1.3 - 16.2 - 16.2.1 - 16.2.4 - 16.2.5)