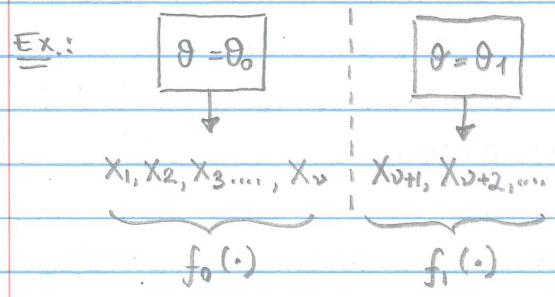


LECTURE 8

* Stochastic Models with Changes

Let us consider the stochastic process $\{X_k\}_k$ satisfying the condition:

- $\exists v \geq 0$: X_1, X_2, \dots, X_v are RVs that satisfy one distribution function $f_0(\cdot)$
 X_{v+1}, X_{v+2}, \dots are RVs that satisfy another distribution function $f_1(\cdot)$
 with $f_1(\cdot) \neq f_0(\cdot)$



The probability function of X_k depends on a parameter vector θ $\forall k$. The vector abruptly changes at time $v \geq 0$, i.e.:

$$f_0(x) \triangleq f_x(x | \theta = \theta_0)$$

$$f_1(x) \triangleq f_x(x | \theta = \theta_1)$$

Ex.: $X_1, X_2, \dots, X_k, \dots$ are observations of a latent RV (i.e., latent state) and the model (i.e., the probability function) of this RV changes over time, i.e.:

$$f_0(x) = f_{\theta_0}(x)$$

$$f_1(x) = f_{\theta_1}(x)$$

Ex.: A variation of the examples reported above can be obtained when the RVs X_1, X_2, \dots are not independent and the functions $f_0(\cdot), f_1(\cdot)$ are either conditional probability functions:

$$f_0 = f_0(X_k | X_{k-1}, X_{k-2}, \dots, X_1)$$

$$f_1 = f_1(X_k | X_{k-1}, X_{k-2}, \dots, X_1)$$

or they are joint probability functions.

In each and every case considered here, $v \geq 0$ is called CHANGE POINT

For a generic process $\{X_k\}$ with changepoint, the following two problems are typically of interest:

2

1) Offline Hypotheses Testing:

We have observations (x_1, x_2, \dots, x_N) of $\{X_k\}$ collected up to $k=N>0$ and we must decide which one of the following two hypotheses is true:

$$H_0) f_x(x_k | x_{k-1}, \dots, x_1) = f_0(x_k | x_{k-1}, \dots, x_1) \quad \forall k \in \{1, 2, 3, \dots, N\}$$

$$H_1) \exists \nu \geq 0 : \begin{cases} f_x(x_k | x_{k-1}, \dots, x_1) = f_0(x_k | x_{k-1}, \dots, x_1) & k=1, 2, \dots, \nu \\ f_x(x_k | x_{k-1}, \dots, x_1) = f_1(x_k | x_{k-1}, \dots, x_1) & k=\nu+1, \nu+2, \dots, N \end{cases}$$

The problem here does NOT require to estimate ν , but just to determine if it exists

2) Sequential Hypothesis Testing:

We have an ongoing collection of observations of $\{X_k\}_k$ and, for any newly collected observation x_n , we must decide whether the sequence of observations (x_1, x_2, \dots, x_n) is enough to be assigned to one of two (or more) possible classes, i.e., we must determine the value:

$$\hat{n} = \min \{n : \exists i \geq 0 \text{ s.t. } (x_1, x_2, \dots, x_n) \sim f_i(\cdot)\}$$

In this case, we need a sequential test that, for every new n , determines if it is enough to make a terminal decision \Rightarrow There will be a stopping time

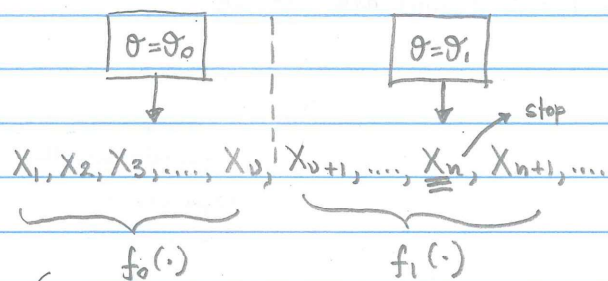
The solution is a policy that makes a tradeoff between average observation time and quality of the decision

The number of observations that is needed can be envisioned as a RV that must be estimated

3) Quickest Change point Detection:

We have an ongoing collection of observations of $\{X_k\}_k$ and, for any newly

collected observation x_n , we must decide whether a change point has occurred or not \Rightarrow The stopping time is now the time when we have enough observations to conclude that a change has occurred



Regardless of the value of v , the stopping time is when we have enough data to decide in favour of a change $\Rightarrow n \geq v$, i.e., a delay is intrinsically related to the problem

The solution is a policy that makes a tradeoff between detection delay and possibility of a FALSE POSITIVE (i.e., $n < v$ or $n \in \mathbb{Z}_0^+$ and $v = \infty$)

Note: Problem 2 and 3 are solved by choosing a decision policy \Rightarrow A decision policy is a pair $\mathcal{S} = (T, d)$ defined by:
 $T \triangleq$ stopping time
 $d \triangleq$ terminal decision

In Problem 3, we can have two scenarios: (i) there are only $N=2$ classes for the observations $(x_1, x_2, \dots) \Rightarrow$ "d" is trivial and follows the choice of T - or the scenario. (ii) there are $N > 2$ classes \Rightarrow "d" is the decision about which post-change hypothesis is satisfied (i.e., which one of the $N-1$ classes the observations x_{v+1}, x_{v+2}, \dots belong to)

Note: The decision policy \mathcal{S} is supposed to be built (or updated) at every new observation (\Rightarrow online implementation) and the management of the tradeoff can be eventually done by minimizing a cost function \Rightarrow An optimization problem can be formulated

We will specifically focus on algorithms that solve the Quickest Change-point Detection problem.

64

* Change point Models

A change point model describes the probabilistic structure of the process $\{X_k\}_k$ and of the change point ν . Numerous options are possible:

ν — RV (with a distribution that is known a priori)

unknown but non-random number

X_k — i.i.d.
independent but nonidentically distributed
dependent observations

Let us consider the history $\mathcal{H}_n \triangleq (X_1 = x_1, X_2 = x_2, \dots, X_{n-1} = x_{n-1})$ and let us denote with $f_{\nu,0}(x)$ and $f_{\nu,1}(x)$ the pdf of X_n under condition $\Theta = \Theta_0$ (i.e., before change point) and $\Theta = \Theta_1$ (i.e., after change point), respectively. We have:

$$f_{\nu}(X_1, X_2, X_3, \dots, X_n) = \prod_{k=1}^{\nu} f_{k,0}(X_k | \mathcal{H}_k) \prod_{k=\nu+1}^n f_{k,1}(X_k | \mathcal{H}_k)$$

joint probability function
when the change point is ν

if $0 \leq \nu \leq n$

In general, if X_k are dependent on the change point (e.g., X_k does not depend on \mathcal{H}_k but only on $(X_{\nu+1}, X_{\nu+2}, \dots, X_{k-1})$ after the change point), we can assume that $f_{k,1}(\cdot)$ is affected by ν and write:

$$f_{\nu}(X_1, X_2, \dots, X_n) = \prod_{k=1}^{\nu} f_{k,0}(X_k | \mathcal{H}_k) \prod_{k=\nu+1}^n f_{k,1}^{(\nu)}(X_k | \mathcal{H}_{\nu+1}^k) \quad (*)$$

Note: X_k - i.i.d. $\Rightarrow f_{k,0}(X_k | \mathcal{H}_k) = f_0(X_k)$

$\mathcal{H}_{\nu+1}^k \triangleq (X_{\nu+1}, \dots, X_{k-1})$

$f_{k,1}^{(\nu)}(X_k | \mathcal{H}_k) = f_1(X_k)$

Also note that, in realistic scenarios, $f_{k,0}(\cdot)$ are known (or they belong to a known class) while $f_{k,1}(\cdot)$ may be unknown and their estimation may be unnecessary in order to choose $S \Rightarrow$ The estimation of the post-change distribution functions is necessary if there are $N-1 > 1$ possible classes of post-change distribution functions (CHANGE DETECTION AND ISOLATION PROBLEM)

Finally, it is possible that, in a time series, there are multiple changepoints \Rightarrow We typically formulate as many changepoint detection problems as the number of changepoints \Rightarrow A reset of \mathcal{S} is assumed at every new detection

Let us consider the changepoint ν . If ν is a RV, its probability function may depend on the observations $X_1, X_2, X_3, \dots \Rightarrow$ We introduce:

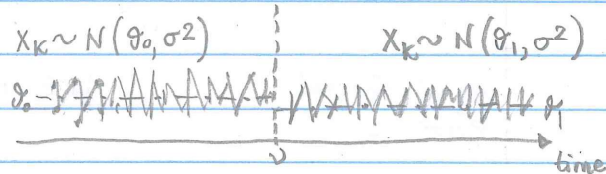
$$q \triangleq P(\nu < 0) \quad (**)$$

$$\pi_k \triangleq P(\nu < k \mid X_k, \mathcal{H}_k)$$

Note that a change point model is defined when a model is given for (*) (at least, for $f_{k,0}(\cdot)$) and (**). By looking at the models for (*), two classes of changes have usually practical interest:

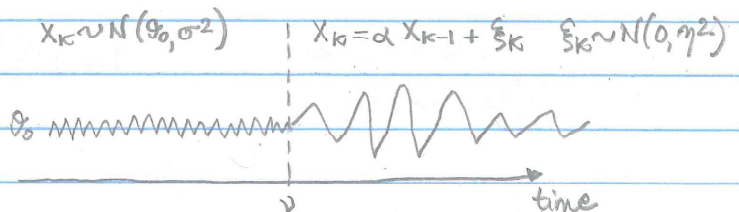
Additive Change \uparrow

- change in the mean value of $\{X_k\}_k$



Non additive Change \uparrow

- change in the mechanism of generation of $\{X_k\}_k$



Additive changes do not affect the probabilistic structure of the process $\{X_k\}$ (i.e., if X_k are observations of a dynamical system, the system dynamics is not affected). Nonadditive changes, instead, are changes that affect the variance, correlations or other nonadditive features of $\{X_k\}_k \Rightarrow$ The system dynamics is affected.

By looking at the models for (**), we can quantify the possibility (a.k.a. "risk") of a false positive via: $E(T \mid \nu = \infty) \triangleq$ mean time to false positive

6

Analogously, we can measure the average detection delay via:

$$E_T(T - \nu | T > \nu) \quad \nu = 0, 1, 2, \dots$$

* Algorithms for Changepoint Detection

Let us assume, for sake of simplicity, the following conditions:

- $\{X_k\}_k$ is made of independent RVs with a probability density $f_X(\cdot)$ conditioned on a scalar $\theta \Rightarrow$ Let us define: $f_\theta(x) \triangleq f_X(x|\theta)$

- $\theta = \theta_0$ for $k \leq \nu$ } with $\nu \geq 0$ changepoint. θ_0 is known, θ_1 may be unknown
 $\theta = \theta_1$ for $k > \nu$ } \Rightarrow Let us define: $f_i(x) \triangleq f_{\theta_i}(x)$

Let us also the following functions:

Log-likelihood Ratio (LLR): $LLR(x) \triangleq \ln \frac{f_1(x)}{f_0(x)} \quad \forall x$

Kullback-Leibler Divergence (KL): $KL(\theta_1, \theta_0) \triangleq E_{\theta_1} \left[\ln \frac{f_1(X)}{f_0(X)} \right] = E_{\theta_1} [LLR]$

A typical detection algorithm requires a function that tracks the probabilistic structure of $\{X_k\}_k$ as the number of observations increases and a threshold to decide when the change has occurred \Rightarrow A simple solution is:

- Function: $S_i^j \triangleq \sum_{k=i}^j LLR(x_k) \quad \forall i, j$

Chart-based Algorithms

- Hypotheses: $H_0) \theta = \theta_0 \quad H_1) \theta = \theta_1$

- Decision: $d = \begin{cases} 0 & \text{if } S_1^N < h \Rightarrow H_0 \text{ is chosen} \\ 1 & \text{if } S_1^N \geq h \Rightarrow H_1 \text{ is chosen} \end{cases}$

with h -threshold to be chosen and $N \geq 1$ fixed and known

In this case, we need N observations to obtain one value of $S_1^N \Rightarrow$ The decision policy is updated every N samples and the stopping time is given by the rule:

$$T = N \cdot \min \{ k : d_k = 1 \}$$

where d_k is the decision taken by using the k -th batch of N observations

Ex.: $X_k \sim N(\theta, \sigma^2)$ - We have:

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\theta)^2}{2\sigma^2}} \Rightarrow \text{LLR}(x) = -\frac{(x-\theta_1)^2}{2\sigma^2} + \frac{(x-\theta_0)^2}{2\sigma^2} = \frac{\theta_1 - \theta_0}{\sigma^2} \left(x - \frac{\theta_0 + \theta_1}{2} \right)$$

$$\text{Hence: } S_1^N = \frac{\theta_1 - \theta_0}{\sigma^2} \sum_{k=1}^N \left(x_k - \frac{\theta_0 + \theta_1}{2} \right) = \frac{\theta_1 - \theta_0}{\sigma^2} \left[\sum_{k=1}^N x_k - \frac{N}{2} (\theta_0 + \theta_1) \right]$$

At the generic batch k , we have:

$$S_1^N(k) = S_{N(k-1)+1}^{Nk}$$

In this example, we need that both θ_0 and θ_1 are known. The decision function is typically formulated as:

$$\left. \begin{array}{l} b \triangleq \frac{\theta_1 - \theta_0}{\sigma} \\ \gamma \triangleq \theta_1 - \theta_0 \end{array} \right\} \Rightarrow S_1^N = \frac{b}{\sigma} N \left(\bar{x} - \theta_0 - \frac{\gamma}{2} \right) \quad \bar{x} \triangleq \frac{1}{N} \sum_{k=1}^N x_k \text{ - sample mean}$$

The threshold is empirically chosen as $h \triangleq \lambda b \sqrt{N} - \frac{N}{2} \left(\frac{\theta_1 - \theta_0}{\sigma} \right)^2$ with λ to be set. In this way:

$$S_1^N(k) \geq h \Leftrightarrow \bar{x}_k \geq \theta_0 + \lambda \frac{\sigma}{\sqrt{N}} \quad \text{Shewhart Control Chart}$$

with $\bar{x}_k \triangleq \frac{1}{N} \sum_{i=1}^N x_{N(k-1)+i}$ and the assumption: $\theta_1 > \theta_0$ □

8

An alternative to the previous algorithm can be conceived to avoid the need for batches of N observations per iteration and to give higher relevance to the latest observations:

$$g_k \triangleq \sum_{i=0}^{\infty} \gamma_i \text{LLR}(x_{k-i}) \quad \Leftarrow \text{New decision function}$$

$$\gamma_i \triangleq \alpha (1-\alpha)^i \quad 0 < \alpha \leq 1 \quad \Leftarrow \text{Exponential weights}$$

In this case, we have: $g_k = g_{k-1} (1-\alpha) + \alpha \text{LLR}(x_k)$

$$T \triangleq \min \{ k : g_k \geq h \}$$

Geometric
Moving Average
Chart

Ex: $x_k \sim N(\mu, \theta^2) \Rightarrow$ We can write:

$$\text{LLR}(x_k) = \ln \frac{\theta_0}{\theta_1} + \left(\frac{1}{\theta_0^2} - \frac{1}{\theta_1^2} \right) \frac{(x_k - \mu)^2}{2}$$

$$\begin{aligned} \Rightarrow g_k &= \ln \frac{\theta_0}{\theta_1} \sum_{i=0}^{\infty} \gamma_i + \left(\frac{1}{\theta_0^2} - \frac{1}{\theta_1^2} \right) \sum_{i=0}^{\infty} \gamma_i \frac{(x_{k-i} - \mu)^2}{2} \\ &= (1-\alpha) g_{k-1} + \alpha \text{LLR}(x_k) \end{aligned}$$

$$\text{If we define: } \tilde{g}_k \triangleq \frac{2\theta_0^2\theta_1^2}{\theta_1^2 - \theta_0^2} g_k - \frac{2\theta_0^2\theta_1^2}{\theta_1^2 - \theta_0^2} \ln \frac{\theta_0}{\theta_1}$$

$$\begin{aligned} \text{we have: } \tilde{g}_k &= \frac{2\theta_0^2\theta_1^2}{\theta_1^2 - \theta_0^2} \left[(1-\alpha) g_{k-1} + \alpha \text{LLR}(x_k) - \ln \frac{\theta_0}{\theta_1} \right] \\ &= \frac{2\theta_0^2\theta_1^2}{\theta_1^2 - \theta_0^2} \left[(1-\alpha) \left(g_{k-1} - \ln \frac{\theta_0}{\theta_1} \right) \right] + \alpha (x_k - \mu)^2 \\ &= (1-\alpha) \tilde{g}_{k-1} + \alpha (x_k - \mu)^2 \end{aligned}$$

Note that the Geometric Moving Average consists of running a infinite-memory filter on the values of LLR \Rightarrow Variants can be defined with finite-memory filters (i.e., finite number of weights γ_i)

By following the same approach, one can use the discrete derivative

$$\nabla g_k \triangleq g_k - g_{k-1}$$

as the decision function. In this case, we are typically interested in counting how many times $|\nabla g_k|$ crosses a given threshold, i.e., we define:

$$I_{\nabla g_k} \triangleq \begin{cases} 1 & |\nabla g_k| \geq h \\ 0 & |\nabla g_k| < h \end{cases} \quad \text{- Then, we consider the decision rule:}$$

$$d \triangleq \begin{cases} 1 & \sum_{i=0}^{N-1} I_{\nabla g_{k-i}} \geq \eta \\ 0 & \text{otherwise} \end{cases} \Rightarrow \text{Hence, } T \triangleq \min \left\{ k: \sum_{i=0}^{N-1} I_{\nabla g_{k-i}} \geq \eta \right\}$$

In these rules, h and η are thresholds to be assigned and N defines the observation window that is monitored to make a decision. \square

These algorithms, though, rely on using fixed thresholds \Rightarrow No adaptation to the dynamics of the process. Also, a decision is taken the very first time a quantity derived from LLR is above threshold \Rightarrow There may be sensitivity to noise. An alternative solution that copes with these issues is:

- Function: $g_k \triangleq S_1^k - m_k$, $m_k \triangleq \min_{1 \leq j \leq k} S_1^j$

CUSUM
Algorithm

- Decision: $d = \begin{cases} 0 & \text{if } g_k < h \\ 1 & \text{if } g_k \geq h \end{cases}$ $h \triangleq$ threshold to be chosen

Note: $g_k \geq h \Leftrightarrow S_1^k \geq m_k + h \Rightarrow$ The threshold is adaptive and depends on the evolution of the cumulative sum S_1^j over time $j \leq k$

Also, note: $S_1^k = \min_{1 \leq j \leq k} S_1^j \Rightarrow g_k = 0$

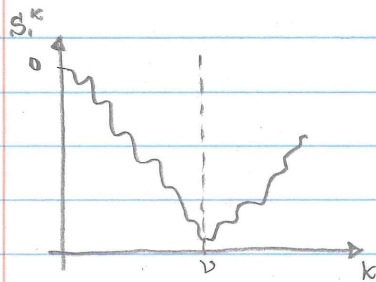
$$S_1^k > \min_{1 \leq j \leq k} S_1^j \Rightarrow m_k = m_{k-1} \text{ and } g_k = g_{k-1} + \text{LLR}(x_k)$$

(10)

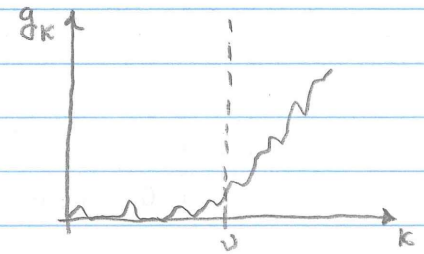
Therefore, we can formulate the decision function in a recursive way:

$$g_k = \sup(0, g_{k-1} + \text{LLR}(x_k)) \leftarrow \begin{array}{l} \text{It is a} \\ \text{RANDOM WALK} \end{array}$$

From a graphical standpoint, one may note: $f_0(x) > f_1(x) \Rightarrow \text{LLR}(x) < 0 \Rightarrow$ The cumulative sum S_1^k has a negative drift for $k \leq \nu$ and eventually changes as $k > \nu$



$\Rightarrow g_k$ captures this trend and amplifies the change in drift



It resets anytime the value of S_k decreases

Note: the negative drift of S_1^k for $k \leq \nu$ can be estimated by computing the expected value of LLR:

$$E_{x_k}(\text{LLR}(x_k) | \nu) = \int_{-\infty}^{\infty} f_0(x) \text{LLR}(x) dx = E_{\theta_0}[\text{LLR}] = -\text{KL}(\theta_0, \theta_1) < 0 \quad \nu \geq k$$

$$E_{x_k}(\text{LLR}(x_k) | \nu) = \int_{-\infty}^{+\infty} f_1(x) \text{LLR}(x) dx = \text{KL}(\theta_1, \theta_0) > 0 \quad \nu < k$$

An interesting generalization can be obtained in the Gaussian case with additive change:

$$\left. \begin{array}{l} X_k \sim N(\theta_0, \sigma) \quad k \leq \nu \\ X_k \sim N(\theta_1, \sigma) \quad k > \nu \end{array} \right\} \Rightarrow \text{LLR}(x) = \frac{b}{\sigma} \left(x - \theta_0 - \frac{\sigma}{2} \right)$$

$$\Rightarrow g_k = \sup\left(0, g_{k-1} + \frac{b}{\sigma} (x_k - \theta_0 - \frac{\sigma}{2})\right)$$

Now, let us suppose that θ_1 can assume either the value: $\theta_1^+ \triangleq \theta_0 + \delta$ or the value: $\theta_1^- \triangleq \theta_0 - \delta$ with $\delta > 0$ fixed and known (e.g., $\delta = \sigma$) $\Rightarrow \text{LLR}(x_k)$ can be one of these two values:

$$LLR(x_k)^+ = \frac{\delta}{\sigma} \left(x_k - \theta_0 - \frac{\delta}{2} \right) \stackrel{\delta=\sigma}{=} x_k - \theta_0 - \frac{\delta}{2} \quad \text{if } \theta_1 = \theta_1^+$$

$$LLR(x_k)^- = \frac{\delta}{\sigma} \left(-x_k + \theta_0 - \frac{\delta}{2} \right) \stackrel{\delta=\sigma}{=} -x_k + \theta_0 - \frac{\delta}{2} \quad \text{if } \theta_1 = \theta_1^-$$

Therefore, we can define two decision functions:

$$\left. \begin{aligned} g_k^+ &= \sup \left(0, g_{k-1}^+ + LLR(x_k)^+ \right) \\ g_k^- &= \sup \left(0, g_{k-1}^- + LLR(x_k)^- \right) \end{aligned} \right\} \Rightarrow \text{The decision rule becomes:}$$

$$d \triangleq \begin{cases} 1 & g_k^+ \geq h \cup g_k^- \geq h \\ 0 & \text{otherwise} \end{cases}$$

Note that such decision rule corresponds to introduce two straight line thresholds for $\sum_{i=1}^k X_i$ (a.k.a. "V-mask") and detect a change when one of them is hit:

$$g_k^+ = \sup \left(0, g_{k-1}^+ + x_k - \theta_0 - \frac{\delta}{2} \right) = \max_{1 \leq j \leq k} \sum_{i=j}^k \left(x_i - \theta_0 - \frac{\delta}{2} \right) \Rightarrow$$

$$g_k^+ \geq h \stackrel{\text{EQUIV}}{\Leftrightarrow} \max_{1 \leq j \leq k} \left[\sum_{i=j}^k x_i - (k-j) \left(\theta_0 - \frac{\delta}{2} \right) \right] \geq h \stackrel{\text{EQUIV}}{\Leftrightarrow} \text{Denoted with } j^* \text{ the}$$

$$\text{arg of the max, we have: } \sum_{i=j^*}^k x_i \geq h + (k-j^*) \left(\theta_0 - \frac{\delta}{2} \right) \leftarrow \text{It's a line} \quad \square$$

A different class of algorithms can be developed under the assumptions:

Bayesian Algorithms	<ul style="list-style-type: none"> • v is a RV • The distribution of v is known A PRIORI 	}	<p>→ The Bayes' Theorem can be invoked</p>
---------------------	--------------------------------------------------------------------------------------------------------------------------------------------	---	--------------------------------------------

In particular, we can introduce:

- Probability of false positive: $P(T \leq v) = \sum_{k=1}^{\infty} P(T \leq v | v=k) P(v=k) \triangleq P^\pi(T \leq v)$

- Average delay to a detection: $E_T(T - v | T > v) = \int_0^{+\infty} (T - v) \frac{f(T - v, T > v)}{P(T > v)} dT$

$$= \int_0^{+\infty} (T - v) \frac{\sum_k f(T - v, T > v | v=k) P(v=k)}{\sum_k P(T > v | v=k) P(v=k)} dT = \frac{1}{P^\pi(T > v)} \sum_k P(v=k) \int_0^{+\infty} (T - v) f(T - v, T > v | v=k) dT$$

(12)

In this formula:

$k > T \Rightarrow T - v < 0$ and $f(T - v, T > v | v = k) = 0$ - Therefore, we can define:

$(T - v)^+ \triangleq \max(0, T - v)$ and write:

$$\int_0^{+\infty} (T - v) f(T - v, T > v | v = k) dv = \mathbb{E}_T \left((T - v)^+ | v = k \right)$$

$$\text{Hence we have: } \mathbb{E}_T \left((T - v)^+ | T > v \right) = \frac{1}{P^\pi(T > v)} \sum_k \underbrace{\mathbb{E}_T \left((T - v)^+ | v = k \right)}_{\mathbb{E}^\pi \left((T - v)^+ \right)} P(v = k) = \frac{\mathbb{E}^\pi \left((T - v)^+ \right)}{P^\pi(T > v)}$$

Let us call: $PFA(T) \triangleq P^\pi(T \leq v)$
 $ADD(T) \triangleq \mathbb{E}^\pi \left((T - v)^+ \right) / P^\pi(T > v)$ } - We can now derive the decision

policy by solving an optimization problem that explicitly minimizes the tradeoff between $PFA(T)$ and $ADD(T) \Rightarrow$ We can represent this tradeoff in a cost function to be minimized.

- What is a good model for $P(v = k)$?

In general, if no prior knowledge is given about v , it is assumed that:

$$v \sim \text{Geom}(p)$$

$$\text{i.e., } P(v = k) = p(1-p)^{k-1} \quad \forall k > 0$$

This assumption simplifies the calculations:

$$P(\theta = \theta_1 \text{ at time } k | \theta = \theta_0 \text{ at time } k-1) = P(v = k-1 | v \geq k-1) = \frac{P(v = k-1)}{P(v \geq k-1)}$$

$$= \frac{p(1-p)^{k-2}}{p(1-p)^{k-2} \sum_{m=0}^{\infty} (1-p)^m} = p \quad \forall k > 0$$

$$P(\theta = \theta_0 \text{ at time } k \mid \theta = \theta_0 \text{ at time } k-1) = P(v \geq k \mid v \geq k-1) = 1-p \quad \forall k \geq 1$$

And similarly:

$$P(\theta = \theta_1 \text{ at time } k \mid \theta = \theta_1 \text{ at time } k-1) = P(v < k \mid v < k-1) = \frac{P(v < k-1)}{P(v < k-1)} = 1$$

$$P(\theta = \theta_0 \text{ at time } k \mid \theta = \theta_1 \text{ at time } k-1) = 0$$

By using the a priori distribution function $P(v=k)$ we can compute the a posteriori probability of a change in (**):

$$\pi_k \triangleq P(v < k \mid X_k, \mathcal{H}_k) = \frac{P(X_k \mid \mathcal{H}_k, v < k) P(v < k \mid \mathcal{H}_k)}{P(X_k \mid \mathcal{H}_k)}$$

$$P(X_k \mid \mathcal{H}_k, v < k) = f_1(X_k) \leftarrow \text{distribution function when } \theta = \theta_1$$

$$\begin{aligned} P(X_k \mid \mathcal{H}_k) &= P(X_k \mid \mathcal{H}_k, v < k) P(v < k \mid \mathcal{H}_k) + P(X_k \mid \mathcal{H}_k, v \geq k) P(v \geq k \mid \mathcal{H}_k) \\ &= f_1(X_k) P(v < k \mid \mathcal{H}_k) + f_0(X_k) P(v \geq k \mid \mathcal{H}_k) \end{aligned}$$

Hence, by denoting: $L(x) \triangleq f_1(x)/f_0(x)$ - likelihood ratio, we can write:

$$\pi_k = \frac{L(X_k) P(v < k \mid \mathcal{H}_k)}{L(X_k) P(v < k \mid \mathcal{H}_k) + P(v \geq k \mid \mathcal{H}_k)}$$

We can further develop the relationship:

$$\begin{aligned} P(v < k \mid \mathcal{H}_k) &= P(v < k-1 \mid X_{k-1}, \mathcal{H}_{k-1}) + P(v = k-1 \mid v \geq k-1, \mathcal{H}_k) P(v \geq k-1 \mid \mathcal{H}_k) \\ &= \pi_{k-1} + p(1-\pi_{k-1}) \end{aligned}$$

$$P(v \geq k \mid \mathcal{H}_k) = P(v \geq k \mid v \geq k-1, \mathcal{H}_k) P(v \geq k-1 \mid \mathcal{H}_k) = (1-p)(1-\pi_{k-1})$$

where we have used the facts:

$$P(v = k-1 | v \geq k-1, \mathcal{H}_k) = P(\theta = \theta_1 \text{ at time } k | \theta = \theta_0 \text{ at time } k-1) = p$$

\uparrow \mathcal{H}_k is the sequence up to $k-1$ and does not provide any further information
 \uparrow Prior distribution is geometric

$$P(v \geq k-1 | \mathcal{H}_k) = 1 - P(v < k-1 | \mathcal{H}_{k-1}, \mathcal{H}_{k-1}) = 1 - \pi_{k-1}$$

$$P(v \geq k | v \geq k-1, \mathcal{H}_k) = P(\theta = \theta_0 \text{ at time } k | \theta = \theta_0 \text{ at time } k-1) = 1-p$$

Therefore, a recursive formulation of π_k in (***) is achieved:

$$\pi_k = \frac{L(x_k) [\pi_{k-1} + p(1-\pi_{k-1})]}{(1-p)(1-\pi_{k-1}) + L(x_k) [\pi_{k-1} + p(1-\pi_{k-1})]} \quad (***)$$

which can be further simplified by defining: $w_k \triangleq \frac{\pi_k}{1-\pi_k}$:

$$w_k = \frac{L(x_k)}{1-p} (w_{k-1} + p) \Rightarrow \text{or, in log-scale: } g_k = \log(e^{g_{k-1} + p}) + \text{LLR}(x_k) - \log(1-p)$$

\uparrow
 $g_k \triangleq \log w_k$

In summary:

Bayesian Approach \Rightarrow We can give a measure for false positives (PFA) and delay in detection (ADD) \Rightarrow We can build a cost function and solve an optimization problem to get the decision policy δ

$v \sim \text{Geom}(p) \Rightarrow$ Calculations simplify and a recursive formulation for π_k is given \Rightarrow The function g_k is not a random walk but it has a specific meaning related to the latent state (it depends on π_k)

References:

Basseville, M. and Nikiforov, I. V. (1993) "Detection of Abrupt Changes: Theory and Applications", Prentice Hall Ed., Englewood Cliffs, NJ

↑

I referred to ch. 1-2 - A copy of the book is on Husky CT