

CHANGE-POINT DETECTION IN A SEQUENCE OF BAGS-OF-DATA

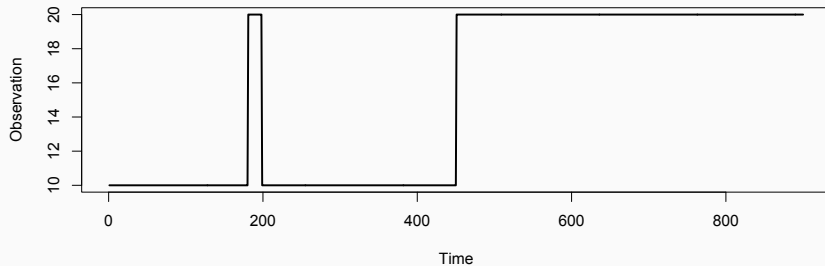
AN EXTENSION OF ANOMALY ANALYSIS

Noboru Murata

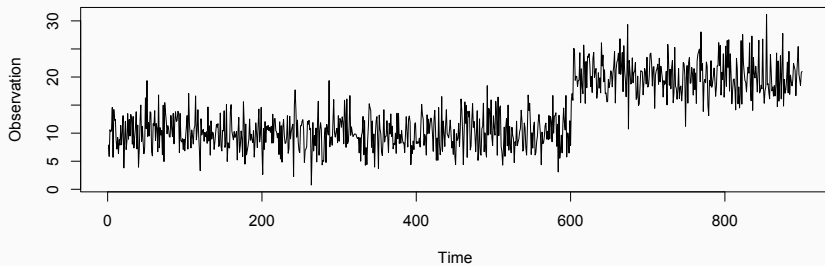
June 20, 2023

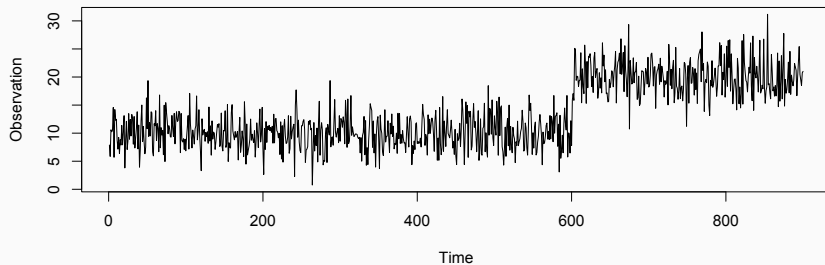
<https://noboru-murata.github.io/>

INTRODUCTION



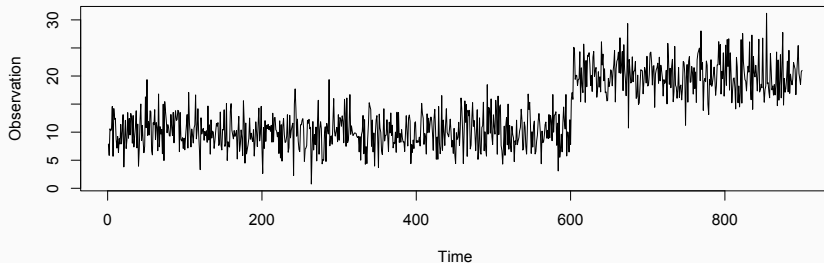
- objective
 - anomaly detection
 - find an outlier of time series
 - change-point detection
 - find a drastic change of time series





- generating mechanism

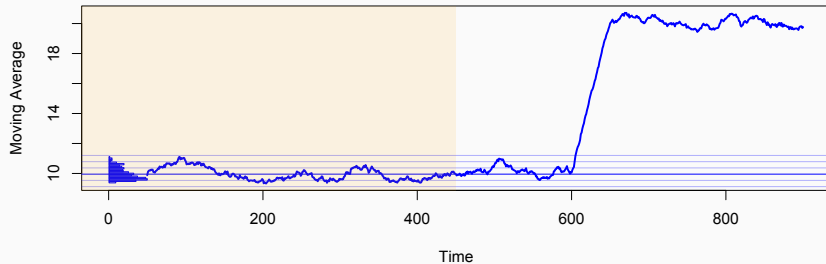
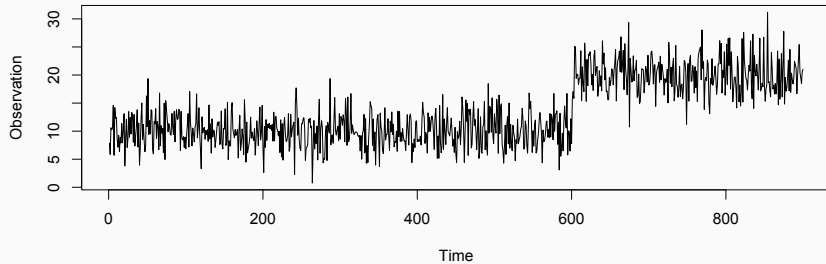
$$X_t = \begin{cases} c_0 + \varepsilon_t, & t < t_0, \\ c_1 + \varepsilon_t, & t \geq t_0, \end{cases} \quad \varepsilon_t \sim P$$

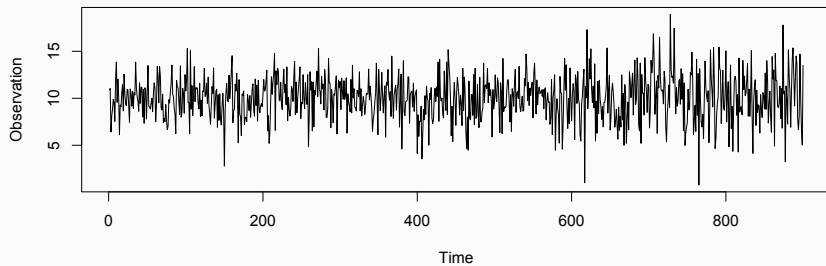


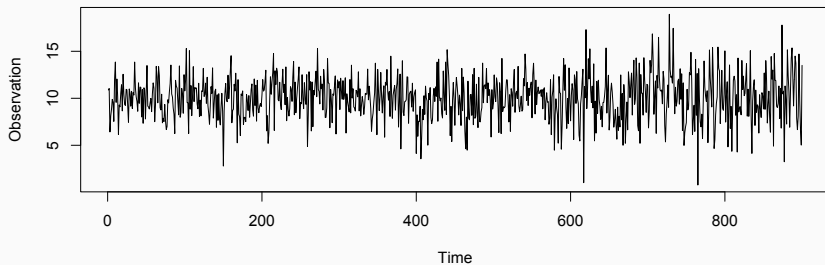
- summary statistics

$$\bar{X}_t = \frac{1}{\tau} \sum_{i=0}^{\tau-1} X_{t-i}$$

estimates of mean values (moving average)

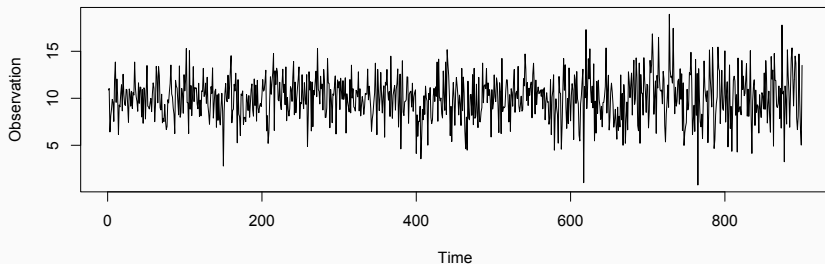






- generating mechanism

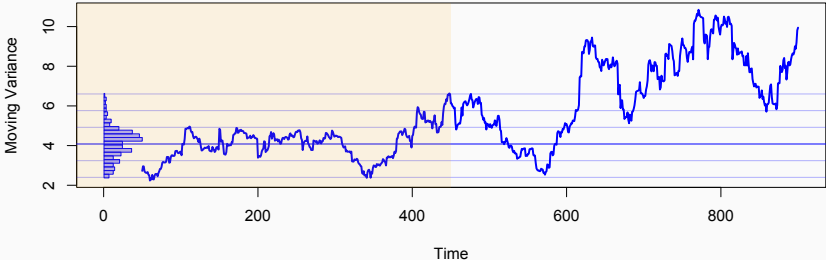
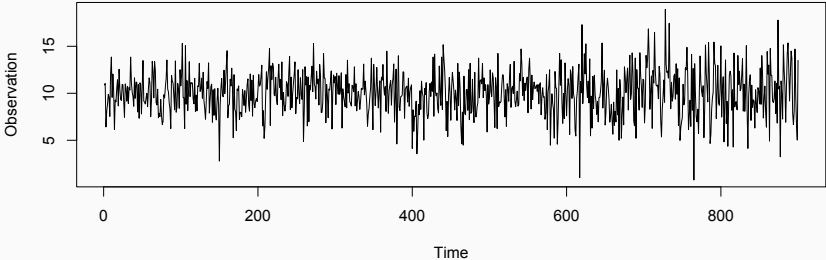
$$X_t = \begin{cases} c_0 + \varepsilon_t, & t < t_0, & \varepsilon_t \sim P \\ c_0 + \xi_t, & t \geq t_0, & \xi_t \sim Q \end{cases}$$

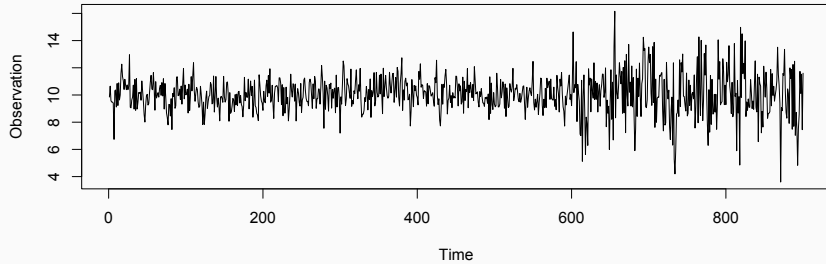


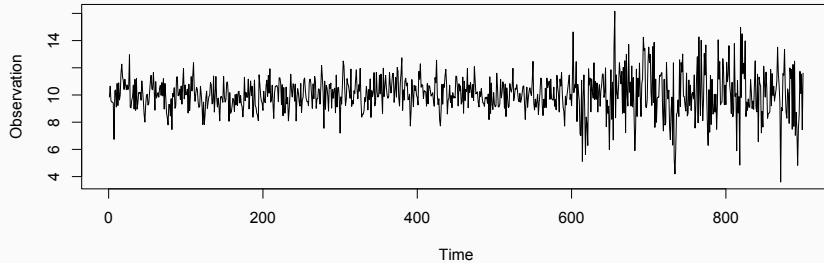
- summary statistics:

$$V_t = \frac{1}{\tau'} \sum_{i=0}^{\tau'-1} (X_{t-i} - \bar{X}_t)^2$$

estimates of variances

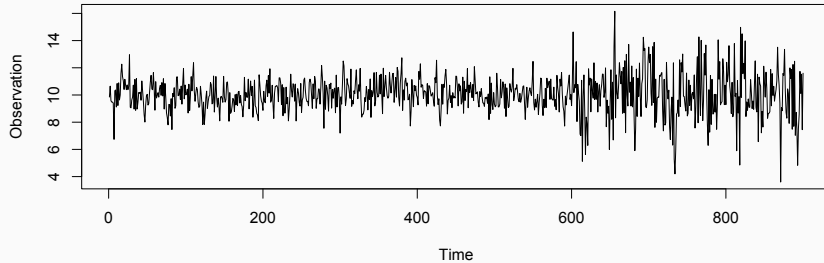






- generating mechanism

$$X_t = aX_{t-1} + bX_{t-2} + \varepsilon_t, \quad \varepsilon_t \sim \begin{cases} P, & t < t_0, \\ Q, & t \geq t_0 \end{cases}$$

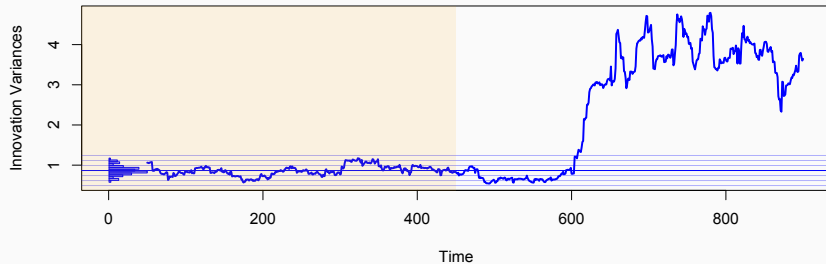
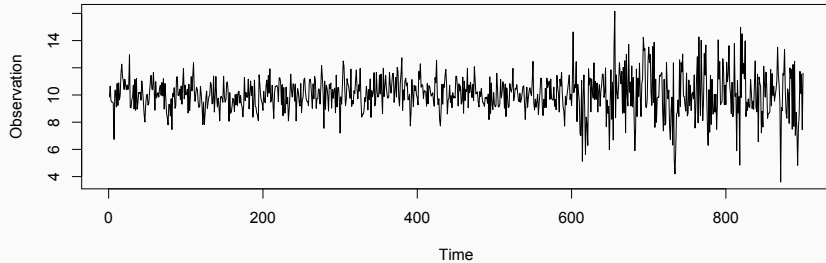


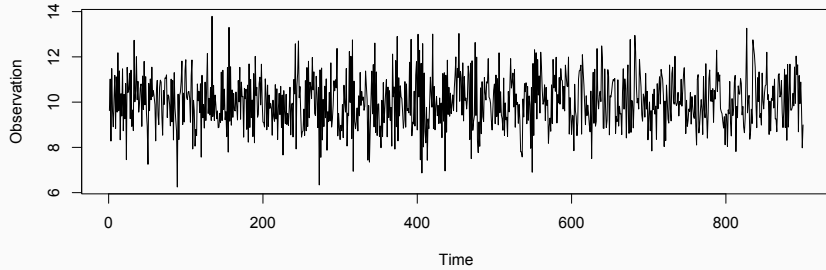
- summary statistics

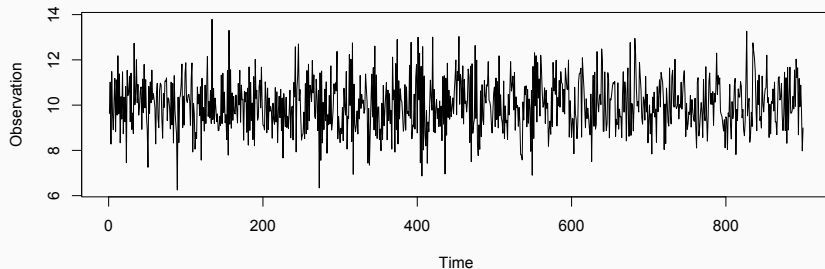
$\text{Var}(\hat{\varepsilon}_t)$ (estimated from X_t, X_{t-1}, \dots)

estimates of innovation variances

$$\hat{\varepsilon}_t = X_t - \hat{X}_t = X_t - (\hat{a}X_{t-1} + \hat{b}X_{t-2})$$

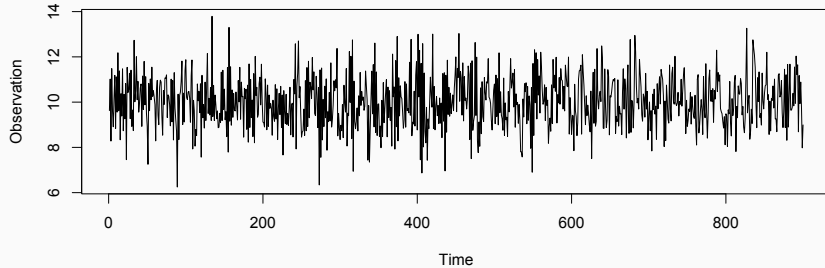






- generating mechanism

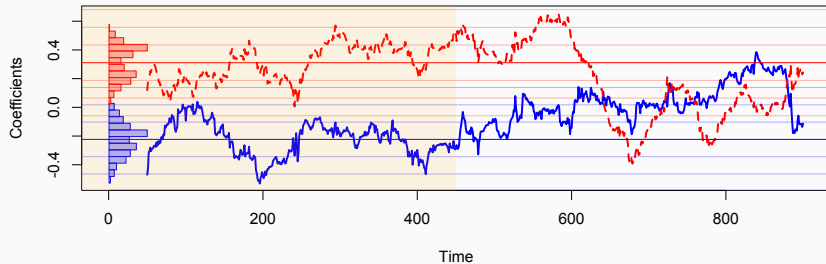
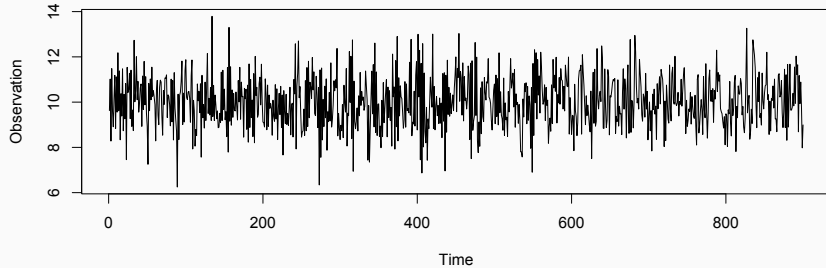
$$X_t = \begin{cases} a_0 X_{t-1} + b_0 X_{t-2} + \varepsilon_t, & t < t_0, \\ a_1 X_{t-1} + b_1 X_{t-2} + \varepsilon_t, & t \geq t_0, \end{cases} \quad \varepsilon_t \sim P$$



- summary statistics

\hat{a}_t, \hat{b}_t (estimated from X_t, X_{t-1}, \dots)
estimates of coefficients

note: multi-dimensional problem

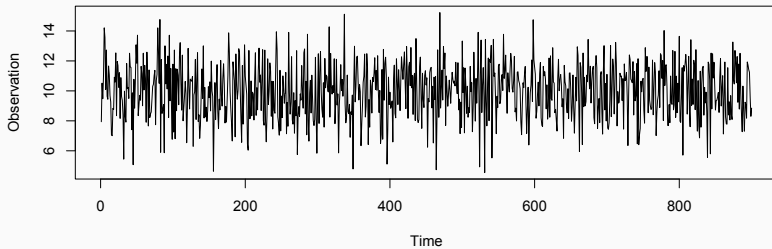


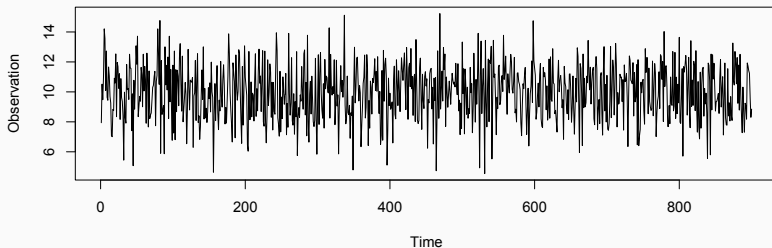
Problem

find time points at which the generating mechanism of time series suddenly changes

- applications
 - intrusion detection in computer networks
 - irregular-motion detection in vision systems
 - signal segmentation in data stream
 - fraud detection in cellular systems
 - fault detection in engineering systems
 - etc.

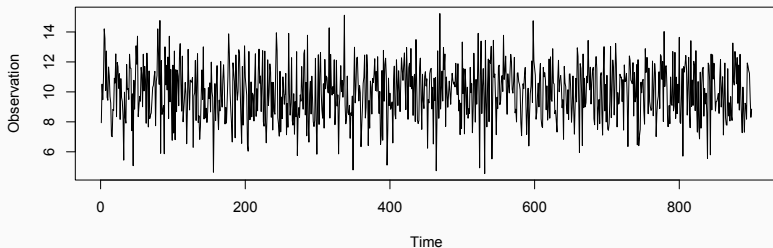
- representative algorithms
 - Singular Spectrum Analysis (Moskvina & Zhigljavskya, 2003)
 - ChangeFinder (Takeuchi & Yamanishi, 2006)
 - Kullback-Leibler Importance Estimation Procedure (Sugiyama et al. 2007)
- differences of these approaches
 - generative models of time series
 - computational costs
 - scalability of data size
 - sensitivity to change of regularity





- generating mechanism

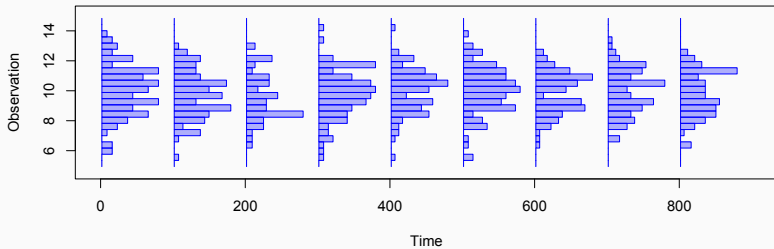
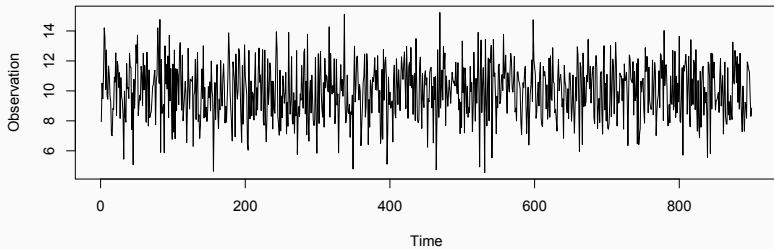
$$X_t = \begin{cases} c_0 + \varepsilon_t, & t < t_0, & \varepsilon_t \sim P \\ c_0 + \xi_t, & t \geq t_0, & \xi_t \sim Q \end{cases}$$



- summary statistics

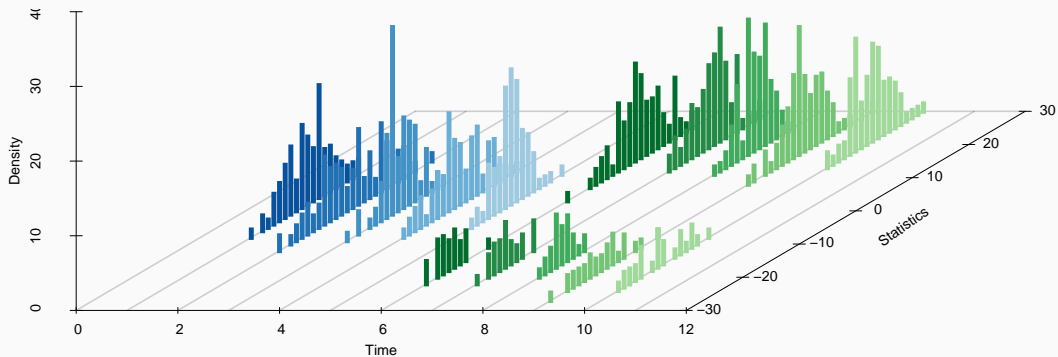
$\hat{P}_t =$ (density estimates of X_t, X_{t-1}, \dots)

i.e. histogram, kernel density estimate, etc.



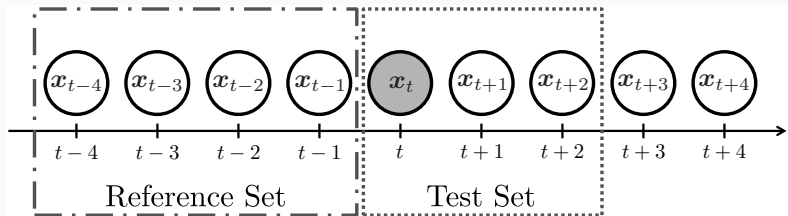
PROBLEM FORMULATION

- framework
 - datum at time t : $B_t = \{X_i; i = 1, \dots, n_t\}$
a set of random variables, i.e. a bag of data
size of bag can be different in time
 - objective:
examine whether B_t, B_{t+1}, \dots differ from B_{t-1}, B_{t-2}, \dots
in statistical setup:
examine whether $\Pr(B_t)$ is predictable from $\Pr(B_{t-1}), \Pr(B_{t-2}), \dots$

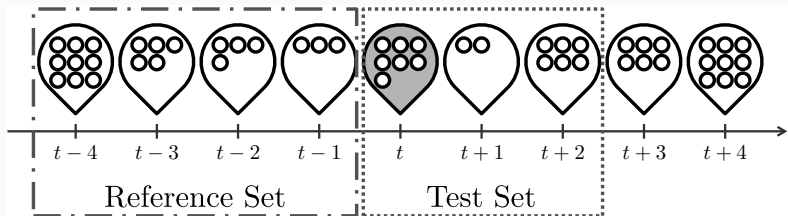


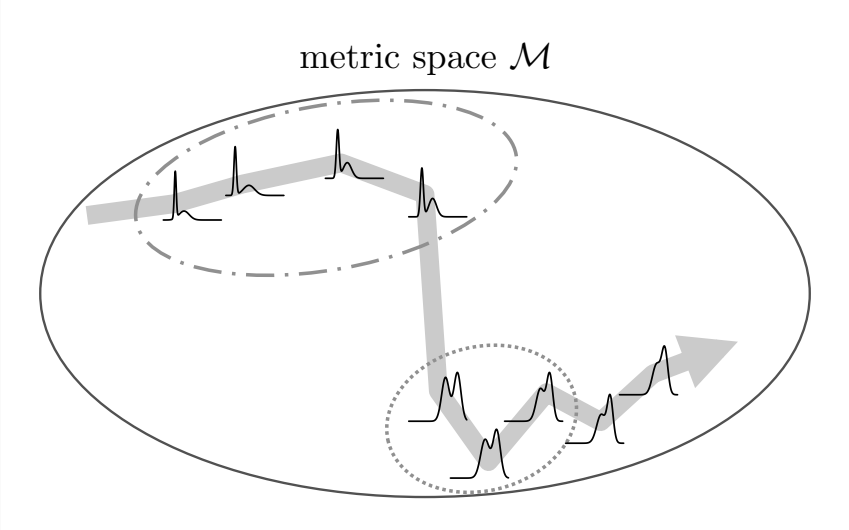
detect a change of distributions behind bags

- standard problem setting

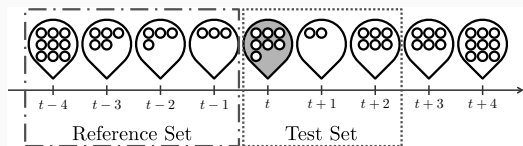


- our problem setting

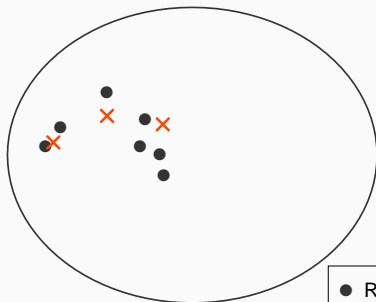




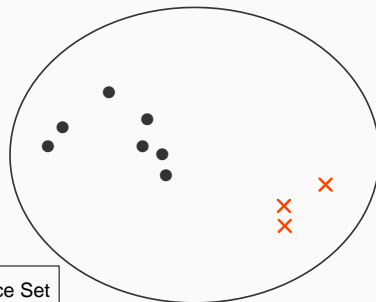
detect a significant change by following a path of bags



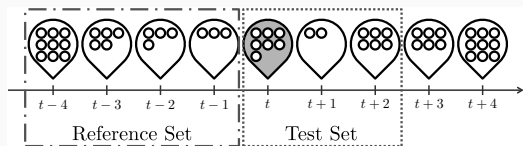
No Change



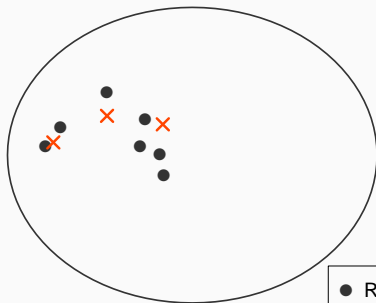
Change



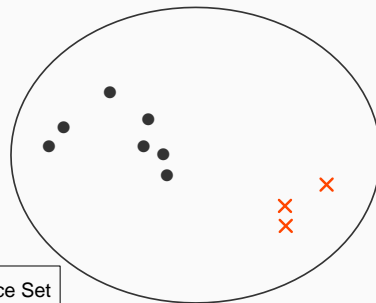
● Reference Set
 × Test Set



No Change



Change



● Reference Set
 × Test Set

- Bayesian bootstrap: Bayesian analogue of the bootstrap
instead of resampling from an empirical distribution, weighted samples are used where weights are sampled from the Dirichlet distribution

$$(N_1, \dots, N_k) \sim \text{Mult}(n; \rho_1, \dots, \rho_k) \quad (\text{resampling})$$

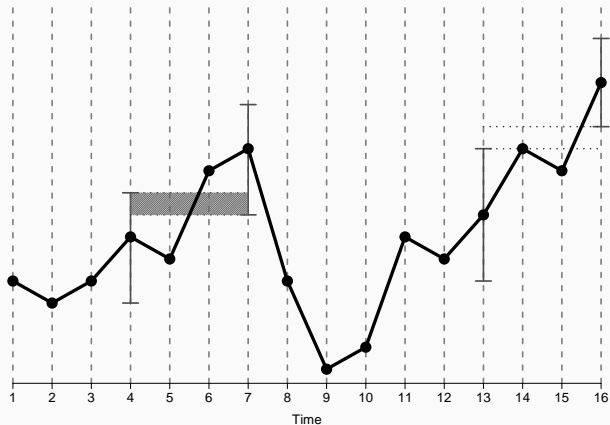
$$(W_1, \dots, W_k) \sim \text{Dir}(\alpha_1, \dots, \alpha_k) \quad (\text{reweighting})$$

- if we let $\alpha_j = n\rho_j$:

$$\mathbb{E}[N_j] = \mathbb{E}[W_j] = \rho_j$$

$$\text{Var}[N_j] = \text{Var}[W_j] \cdot \frac{n+1}{n} = \frac{\rho_j(1-\rho_j)}{n}$$

- confidence interval with Bayesian bootstrap on weights of bags
 - **regular**: intervals intersect each other
 - **anomalous**: otherwise

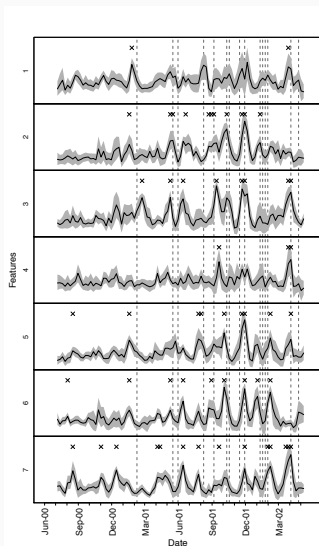


NUMERICAL EXAMPLES

Enron Email Dataset (Cohen, 2009)

email transmission data from about 150 users, mostly senior management of Enron

- duration: 2000/6 – 2002/5 (accounting scandal: 2001)
- time window size of bags: 1 week
- size of reference datasets: 5 weeks
- size of test datasets: 3 weeks
- statistics in bags: 7 stats of bipartite graphs
 - degree of sender / receiver
 - 2nd order degree of sender-sender / receiver-receiver
 - number of messages from sender / to receiver
 - number of messages between sender and receiver
- confidence interval: 0.95



Date	Proposed	GS	Event
February 12, 2001	X	X	Jeff Skilling becomes chief executive of Enron.
May 19, 2001	X		Congress begins implementing President Bush's energy plan into legislation.
June 5, 2001	X	X	Rove divests his stocks in energy.
August 14, 2001	X	X	Skilling resigns abruptly citing personal reasons. Kenneth Lay returns to CEO.
September 11, 2001	X		Four terrorist attacks launched by al-Qaeda.
October 16, 2001	X		Enron reports a \$618 million loss and a \$1.2 billion reduction in shareholder equity.
October 19, 2001	X		Securities and Exchange Commission launches inquiry into Enron finances.
November 19, 2001	X	X	Enron restates its third-quarter earnings and says a \$690 million debt is due Nov. 27.
November 29, 2001	X	X	Dynergy deal collapses.
December 2, 2001	X		Enron files for bankruptcy, the biggest in US history, and lays off 4,000 employees.
January 9, 2002	X	X	The justice department opens a criminal investigation of Enron.
January 17, 2002			Enron fires Andersen blaming the auditor for destroying Enron documents.
January 23, 2002		X	Kenneth Lay resigns as chairman and chief executive of Enron.
January 30, 2002	X	X	Enron names Stephen F. Cooper new CEO.
February 4, 2002	X	X	Kenneth Lay resigns from the board.
April 9, 2002	X		David Duncan, Andersen's former top Enron auditor, pleads guilty to obstruction.
April 24, 2002		X	House passes accounting reform package.

CONCLUSION
