

Takeuchi, J., and Yamanishi, K.:
A Unifying Framework for Detecting
Outliers and Change Points
from Time Series,
*IEEE Trans. on Knowledge and Data
Engineering*, 18(4), pp.482–492, 2006.

2013年度論文ゼミ#9

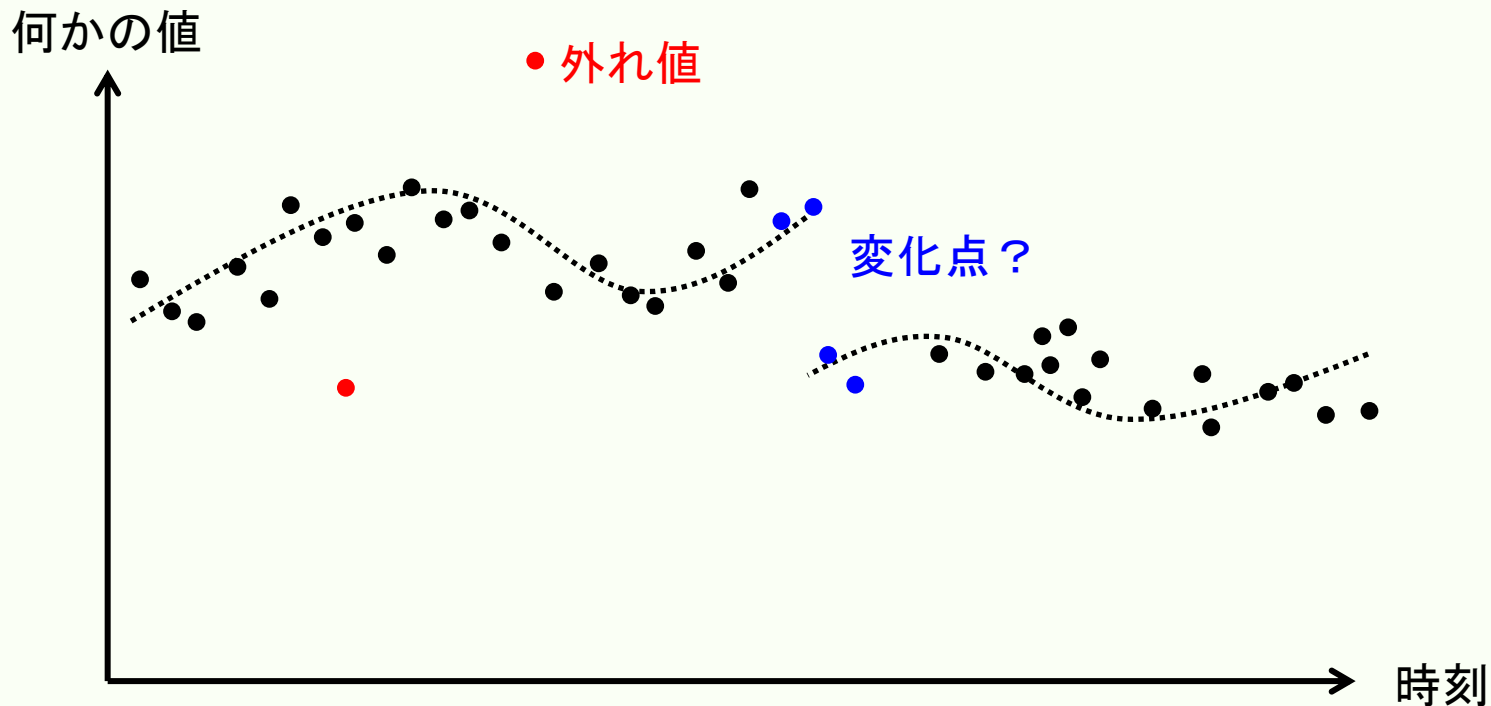
20130621

中西 航

1. Introduction
2. Problem Setting
::本論文の主題”Change Point Detection”の位置づけ
3. Two-Stage Learning for Change Point Detection
::Change Point Detectionの定式化
4. Other Methods
::以下での比較用に別の手法を紹介
5. Simulation
::人工データに対する適用と分析
6. Experiment with Real Data
::実データに対する適用
7. Concluding Remarks

1. 時系列データからの「検出」

- ◆ 時系列データの変化には、
 - ◆ 外れ値の検出
 - ◆ 変化点の検出の2種類があり、ともに興味を引く問題となってきた



- ◆ 時系列データの変化には、
 - ◆ 外れ値の検出
 - ◆ 変化点の検出の2種類があり、ともに興味を引く問題となってきた

- ◆ 例) ネットワークのアクセスログ
 - ◆ 通常状態をデータから学習し、新たに入ってくるデータがそこからどの程度はずれているか＝外れ値＝異常検出
 - ◆ 一方で、膨大な外れ値データが入ってきた場合は、統計的に通常状態自体が変化しているにとらえたい

- ◆ これまで、外れ値検出と変化点検出は別個に行われてきた
→これを統一した枠組みで記述したい

1. 著者らの既往研究との関係

- ◆ 従来：統計的な外れ値検出手法(時系列データに未対応)
→これに対し、
 - ◆ 時系列データに対応
 - ◆ 変化点の検出にも対応
- の2点が本研究の手法。具体的には、
- ◆ Auto Regressionモデルを従来手法に導入する
 - ◆ 外れ値検出と変化点検出を、
2段階の「検出」ステップに分離し、逐次同時に処理する
 - ◆ より具体的には、データ単体の外れ度合いの計算と、
データの移動平均の外れ度合いの計算を行う
 - ◆ この手法をChangeFinderと命名する

1. 関連する研究(の課題)

- ◆ 外れ値検出は時系列データへの対応が不完全
- ◆ 時系列を考慮したARモデルはデータの定常性を仮定

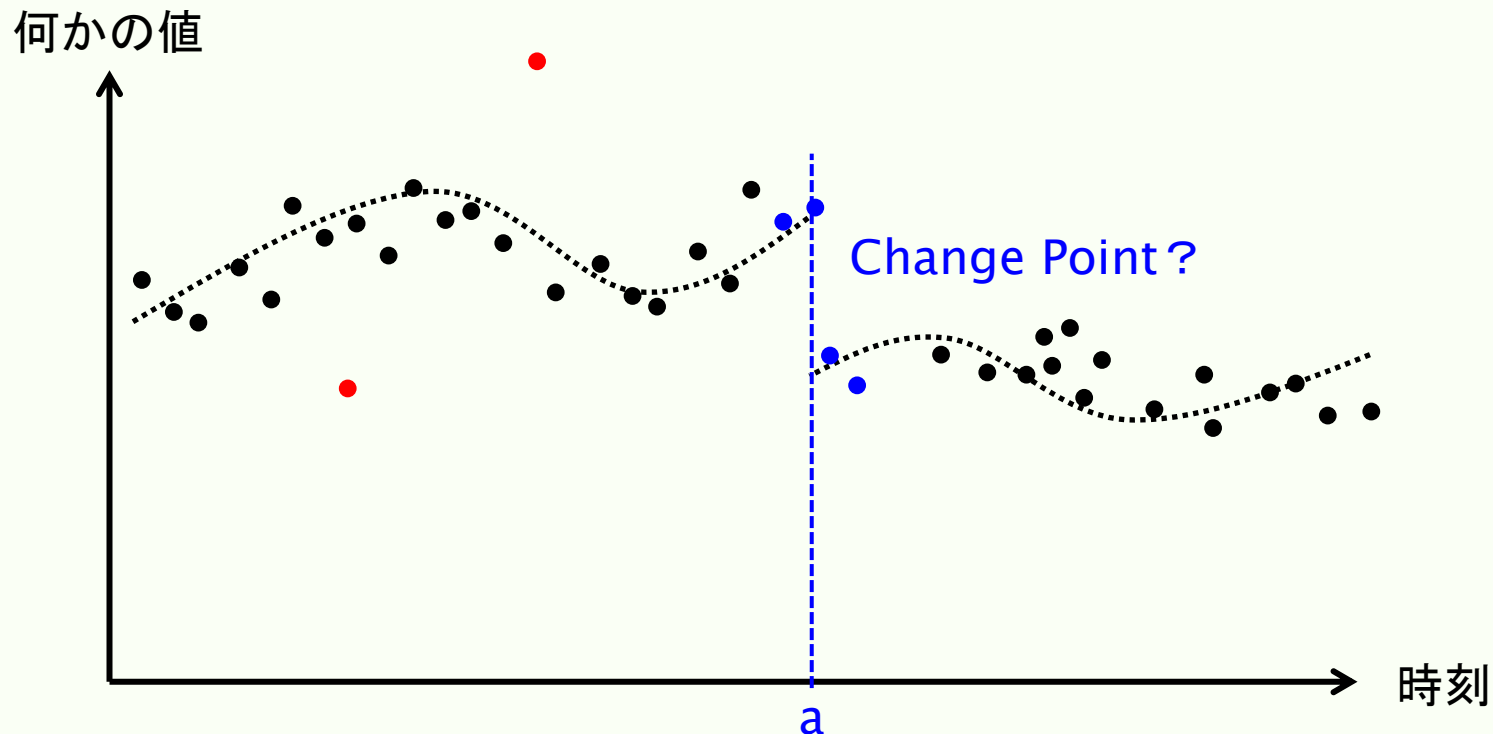
→ 非定常な時系列データを扱う方法として、

- ◆ ARモデルの(共)分散に時間変化を持たせる方法
 - ◆ 割引率の導入による一種の最尤推定(本研究のアプローチ)
-
- ◆ 事前に変化点の個数が分かっていることを前提
 - ◆ 変化後のデータがある程度蓄積されてから変化点が検出

→ より一般化した適用可能性の高い手法の構築

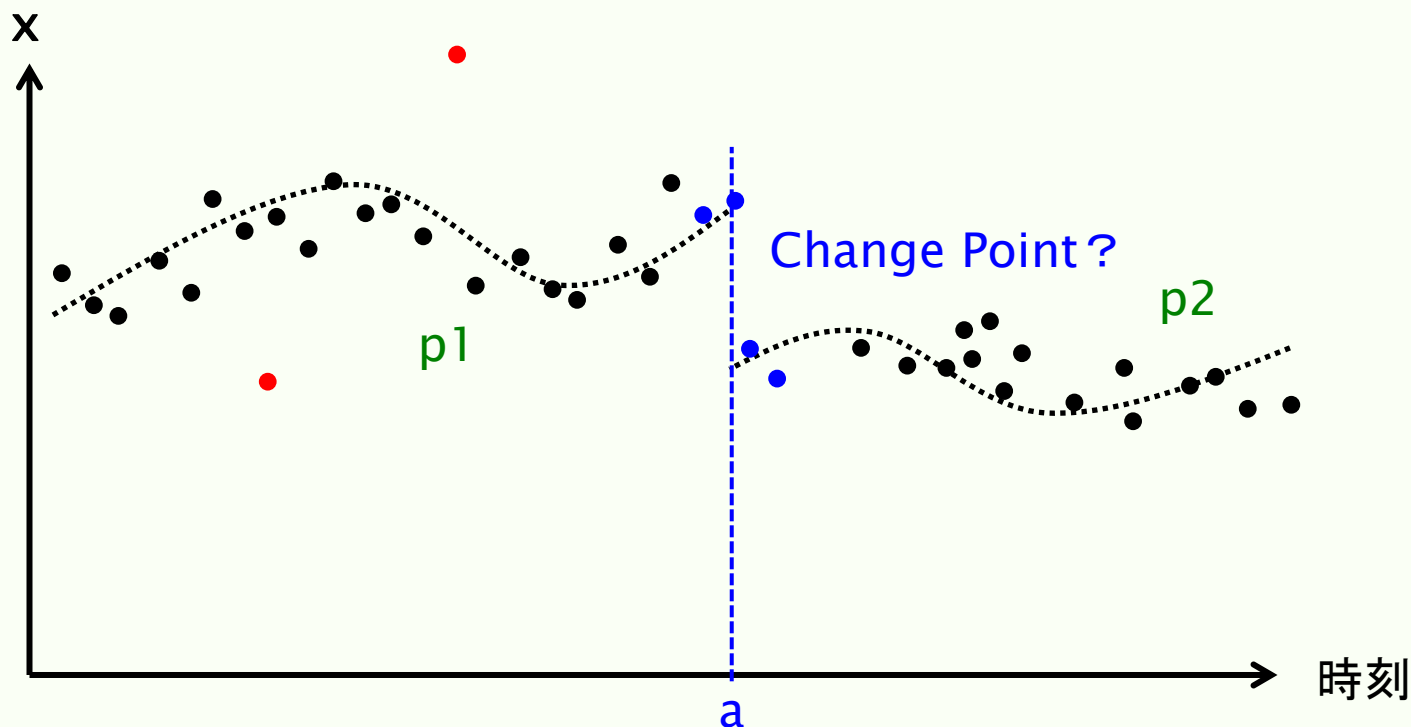
2.Change Point Detectionでは何を行うか

- ◆ Change Pointとは、データの属性が急激に変化する時間軸上の点(時刻)のことで、緩やかな変化は対象としない
- ◆ この時刻 $t=a$ を検出するのがChange Point Detection
 - ◆ (中西注：緩やかな変化に対しては、動的パラメータの逐次推定やファジィシステムの適用が行われているように思われる)



2.Change Point Detectionでは何を行うか

- ◆ ある確率過程 p に基づいて、 d 次元の実数ベクトル \mathbf{x}_t が $t=1$ から順次得られる
- ◆ $p(\mathbf{x}_t|\mathbf{x}^{t-1})$ は \mathbf{x}_1 から \mathbf{x}_{t-1} が得られたもとでの \mathbf{x}_t の分布
- ◆ このとき、ある時刻 $t=a$ がChange Pointならば、 p がある時刻 $t=a$ を境に $p^{(1)}$ と $p^{(2)}$ という2つの(異なる)確率密度関数に分離されるという想定が出来る



2.Change Point Detectionでは何を行うか

- ◆ このときカルバック・ライブラー情報量

$$D(p^{(2)} \parallel p^{(1)}) \equiv \lim_{n \rightarrow \infty} \frac{1}{n} E_{p^{(2)}} \ln \frac{p^{(2)}(\mathbf{x}^n)}{p^{(1)}(\mathbf{x}^n)} \quad E_p : p \text{の期待値}$$

- ◆ たとえば、有意なChange Pointである $t=a$ では D の値がとても大きい、というような想定
- ◆ もし $p^{(1)}$ と $p^{(2)}$ が i.i.d. を満たす確率密度関数ならば、 $p^{(1)}$ のもとで \mathbf{x}_{a+1} から \mathbf{x}_{a+m} までの一連の \mathbf{x} が得られる確率は $p^{(2)}$ のもとで得られる確率の $\exp(-mD)$ 倍に比例する
- ◆ ところで、上記の方法では、唯一の変化点 $t=a$ の存在と、その前後での $p^{(1)}$ と $p^{(2)}$ の定常性を仮定している
- ◆ これに対し本研究は、複数の変化点に対応し、変化点を出来る限りリアルタイムで検出することを目的としている

- ◆ 具体的には以下の2つを扱う

- ◆ Jumping mean(平均の変動)

- $p^{(i)}$ がi.i.d.の1次元正規分布で、平均 $\mu^{(i)}$ 、分散ともに σ^2 のとき

$$D(p^{(2)} \parallel p^{(1)}) = \frac{(\mu^{(1)} - \mu^{(2)})^2}{2\sigma^2}$$

- $|\mu^{(1)} - \mu^{(2)}|$ が大きいときJumping mean型の変化点

- ◆ Jumping variance(分散の変動)

- $p^{(i)}$ がi.i.d.の1次元正規分布で平均ともに0、分散が $\sigma_{(i)}^2$ のとき

$$D(p^{(2)} \parallel p^{(1)}) = \frac{1}{2} \left(\frac{\sigma_{(2)}^2}{\sigma_{(1)}^2} - 1 - \ln \frac{\sigma_{(2)}^2}{\sigma_{(1)}^2} \right)$$

- $\sigma_{(1)}^2$ と $\sigma_{(2)}^2$ が大きく異なるときJumping variance型の変化点

3. ChangeFinder

◆ 2段階学習を行う

x_t : 時系列データ
 p : x の時系列変化を表す(条件付き)確率密度関数
 y_t : p_{t-1} のもとでの x_t の合致度の移動平均
 q : y の時系列変化を表す(条件付き)確率密度関数

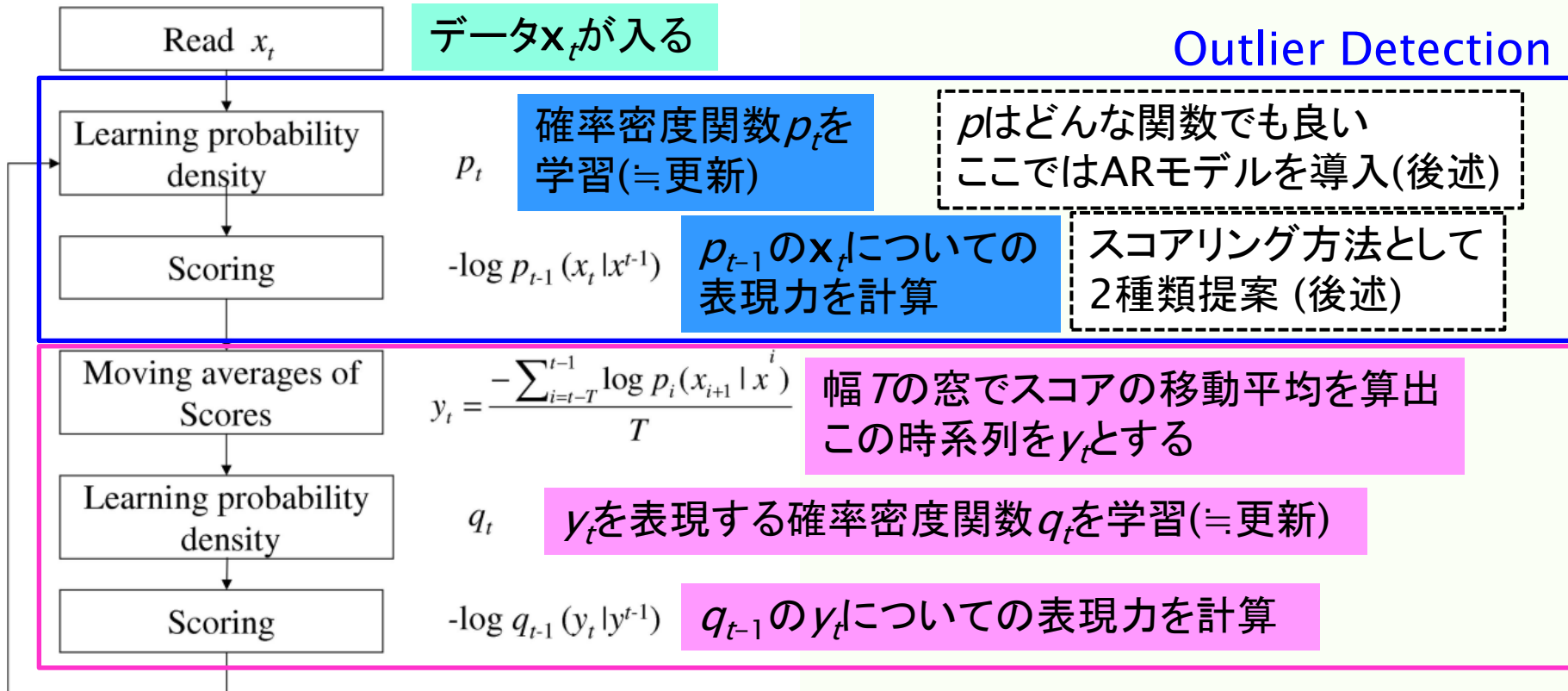
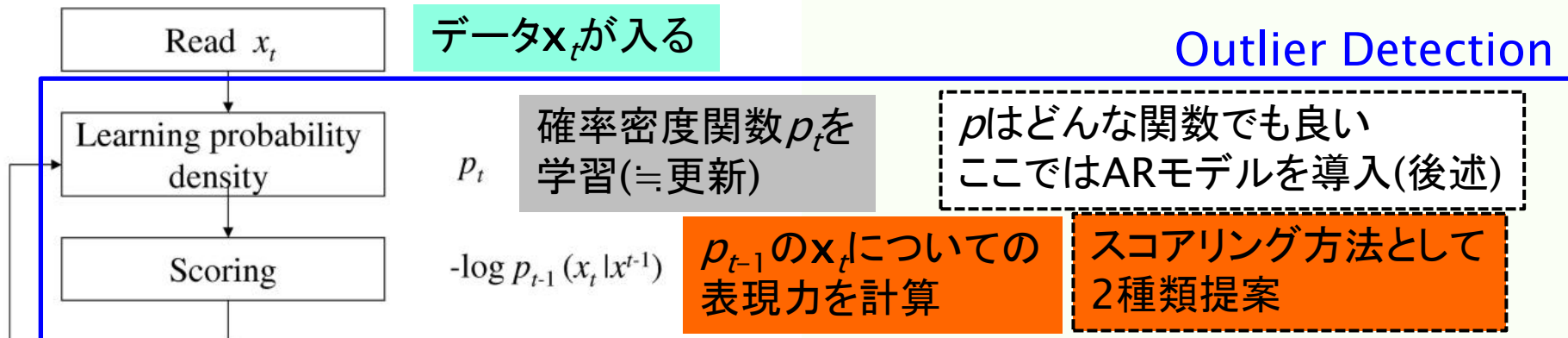


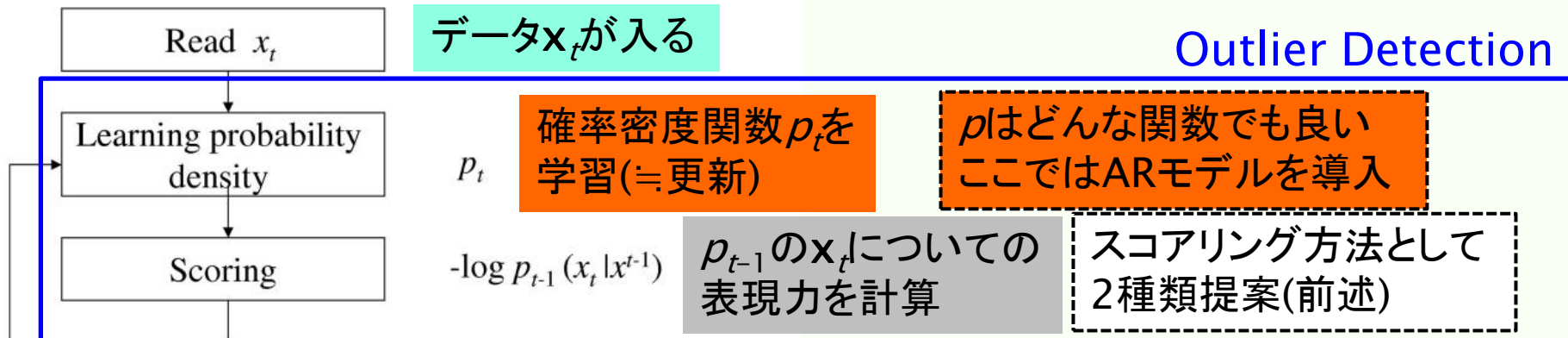
Fig. 1. Flow of ChangeFinder.



◆ Scoringの方法

1. 対数損失 $Score(\mathbf{x}_t) = -\log p_{t-1}(\mathbf{x}_t | \mathbf{x}^{t-1})$
情報理論の立場では、 p_{t-1} に従って生成されるバイナリデータをもとに \mathbf{x}_t をエンコードするのに必要な記述長
2. 二次損失 $Score(\mathbf{x}_t) = (\mathbf{x}_t - \hat{\mathbf{x}}_t)^2, \hat{\mathbf{x}}_t = \int \mathbf{x} p_{t-1}(\mathbf{x}_t | \mathbf{x}^{t-1}) d\mathbf{x}$

★ $Score(\mathbf{x})$ が大きいとき外れ値の可能性が高いということ



- ◆ p にARモデルを導入
初期値の平均が0であるような d 次元ベクトル \mathbf{z}_t が k 次のARモデルに従う

$$\mathbf{z}_t = \sum_{i=1}^k \mathbf{A}_i \mathbf{z}_{t-i} + \boldsymbol{\varepsilon}, \boldsymbol{\varepsilon} \sim N(0, \boldsymbol{\Sigma})$$

$$\mathbf{x}_t = \mathbf{z}_t + \boldsymbol{\mu}$$

- ◆ \mathbf{x}_{t-k} から \mathbf{x}_{t-1} までの系列データを \mathbf{x}_{t-k}^{t-1} と表記すれば、

$$p(\mathbf{x}_t | \mathbf{x}_{t-k}^{t-1}; \boldsymbol{\theta}) = \frac{1}{(2\pi)^{k/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x}_t - \mathbf{w})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}_t - \mathbf{w})\right)$$

$$\mathbf{w} = \sum_{i=1}^k \mathbf{A}_i (\mathbf{x}_{t-i} - \boldsymbol{\mu}) + \boldsymbol{\mu}, \boldsymbol{\theta} = (\mathbf{A}_1, \dots, \mathbf{A}_k, \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

- ◆ さきほどの式展開

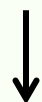
$$\mathbf{x}_t = \mathbf{z}_t + \boldsymbol{\mu}$$

$$\mathbf{w} = \sum_{i=1}^k \mathbf{A}_i (\mathbf{x}_{t-i} - \boldsymbol{\mu}) + \boldsymbol{\mu} = \mathbf{A}_1 (\mathbf{x}_{t-1} - \boldsymbol{\mu}) + \mathbf{A}_2 (\mathbf{x}_{t-2} - \boldsymbol{\mu}) + \dots + \mathbf{A}_k (\mathbf{x}_{t-k} - \boldsymbol{\mu}) + \boldsymbol{\mu}$$



$$\mathbf{x}_t - \mathbf{w} = \mathbf{z}_t - (\mathbf{A}_1 \mathbf{z}_{t-1} + \mathbf{A}_2 \mathbf{z}_{t-2} + \dots + \mathbf{A}_k \mathbf{z}_{t-k}) = \boldsymbol{\varepsilon} \sim N(0, \boldsymbol{\Sigma})$$

\mathbf{z}_t についての回帰式



$$p(\mathbf{x}_t | \mathbf{x}_{t-k}^{t-1} : \boldsymbol{\theta}) = \frac{1}{(2\pi)^{k/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x}_t - \mathbf{w})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_t - \mathbf{w})\right)$$

- ◆ Sequential Discounting AR learningによる $\boldsymbol{\theta}$ の推定：
割引率 r を導入した以下の式を最大化する $\boldsymbol{\theta}$

$$\sum_{i=1}^t (1-r)^{t-i} \log p(\mathbf{x}_i | \mathbf{x}^{i-1}, \boldsymbol{\theta})$$

- ◆ Sequential Discounting AR learningによる θ の推定量
「:=」は代入(更新)であり、逐次推定を行うという意味

$$\begin{aligned}\hat{\mu} &:= (1-r)\hat{\mu} + rx_t, \\ C_j &:= (1-r)C_j + r(x_t - \hat{\mu})(x_{t-j} - \hat{\mu})^T \\ &\quad (j = 0, \dots, k).\end{aligned}$$

μ の推定量

C: 計算上設定する変数
「自己共分散関数」と呼ばれる

Solve the following equation (C_{-i} denotes C_i^T):

$$C_j = \sum_{i=1}^k A_i C_{j-i} \quad (j = 1, \dots, k).$$

← Yule-Walkerの方程式
これを解くと
Aの推定量が得られる

Letting the solution to (14) be $\hat{A}_1, \dots, \hat{A}_k$, then calculate

$$\begin{aligned}\hat{x}_t &:= \sum_{i=1}^k \hat{A}_i (x_{t-i} - \hat{\mu}) + \hat{\mu}. \\ \hat{\Sigma} &:= (1-r)\hat{\Sigma} + r(x_t - \hat{x}_t)(x_t - \hat{x}_t)^T.\end{aligned}$$

Aの推定量を代入すると
 x と Σ の推定量が得られる

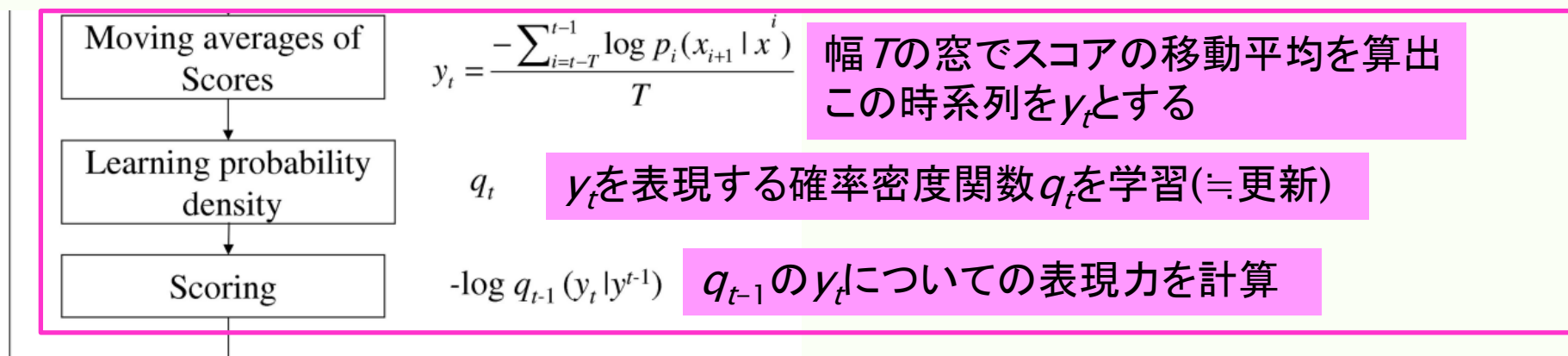


Fig. 1. Flow of ChangeFinder.

Change Point Detection

- ◆ $Score(\mathbf{x})$ を元に新たな時系列データ y を生成 $y_t = \frac{1}{T} \sum_{i=t-T+1}^t Score(\mathbf{x}_i)$
 - ◆ いわゆる移動平均
- ◆ 外れ値 \mathbf{x} の前後での y の値の変化は小さく、変化点 t の前後での y の値の変化は大きいことが想定される
- ◆ $\mathbf{x}_t \rightarrow p_t$ と同じ要領で、 y_t から q_t を生成する
- ◆ p_t と同じ要領で q_t について Scoring し、 $Score(t)$ とする
- ◆ T : 小 \rightarrow 変化点検出高速、外れ値との峻別能力低い
- ◆ T : 大 \rightarrow 変化点検出低速、検出精度高い

- ◆ GS: Guralnik and Srivastava(1999)に基づく手法
 - ◆ 誤差(残差)の平方和の最小化
 - ◆ GSでは多項式による近似も行っているが今回は線形近似
- ◆ SC: GSにARモデルを導入した手法
 - ◆ 確率的コンプレキシティ(Stochastic Complexity)最小化
 - ◆ 確率密度関数 $p(x_t | x_{t-k}^{t-1}, \theta)$ と系列データ x_u^t が与えられたとき、
確率的コンプレキシティ $I(x_u^t) \equiv -\sum_{i=u}^t \log p(x_i | x_{i-k}^{i-1}, \theta(x_u^{i-1}))$
 - ◆ θ は与えられた x (時刻 u から $t-1$ まで) のもとでの最尤推定量
 - ◆ 変化点 $t=a$ のあと、毎時刻 $t=b$ で $\frac{1}{t_b - t_a + 1} \left(I(x_{t_a}^{t_b}) - \min_{i:t_a < i < t_b} (I(x_{t_a}^i) + I(x_{i+1}^{t_b})) \right)$ を計算し、閾値以上となれば $t=i$ を変化点として検出
 - ◆ 変化点が検出されたら、 i を a として繰り返す
 - ◆ $t=i$ で区間を分割し、異なる θ による p を用いた方が確率的コンプレキシティが低下するとき、表現力が高いモデルだということ

- ◆ x_t : 2次のARモデルに従って生成

$$x_t = a_1 x_{t-1} + a_2 x_{t-2} + \varepsilon_t$$

$$a_1 = 0.6, a_2 = -0.5, \varepsilon \sim N(0, \sigma^2)$$

- ◆ y_t : 3次のARモデルに従って生成

- ◆ $Score(x)$: 対数損失で計算

- ◆ $Score(y)$: 二次損失で計算

- ◆ 移動平均の幅 $T=5$

- ◆ 割引率 $r=0.02$

- ◆ 変化点は $t=1000x$ ($x=1, 2, \dots, 9$) に設ける
→ 3種類のパターンで試行

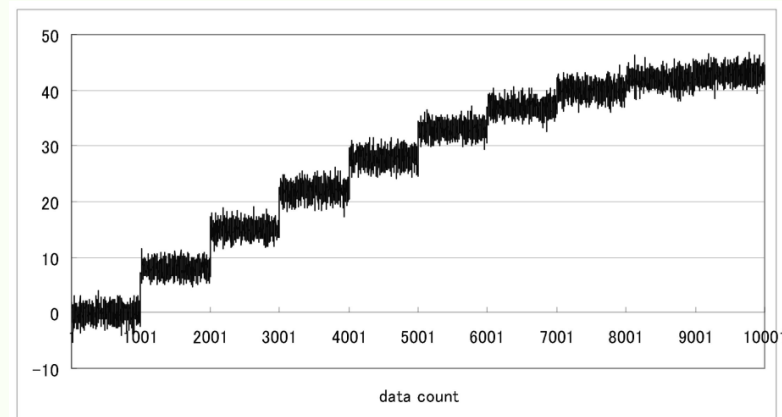
5. Simulation: Jumping mean with constant variance

- ◆ $t=1000x(x=1,2,\dots,9)$ の9回、平均を $\Delta x=10-x$ 変化(加算)
 - ◆ Δx をchange sizeと呼ぶ
- ◆ $\sigma^2=1$ (一定)
- ◆ x 回目の変化前後でのKL情報量は

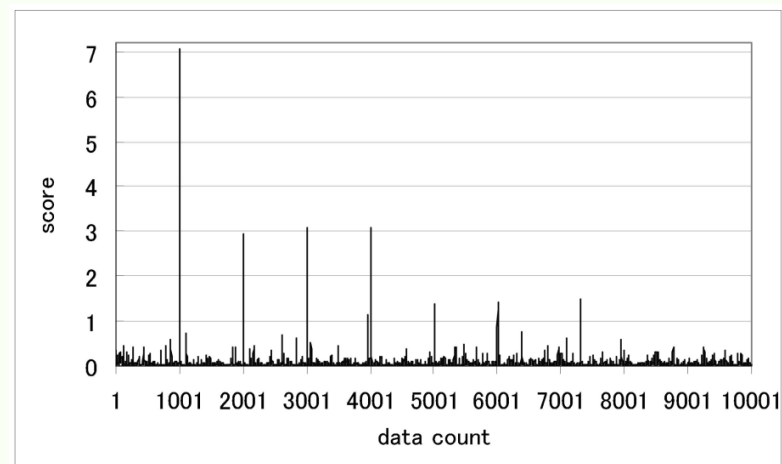
$$K(x) = \frac{\Delta(x)^2 \left(1 - \sum_{i=1}^x a_i\right)^2}{2\sigma^2} = 0.405\Delta(x)^2 \approx 0.4(10-x)^2$$

- ◆ $t=6000$ まで検出できている
このとき $K(x) \doteq 6.4$

元データ



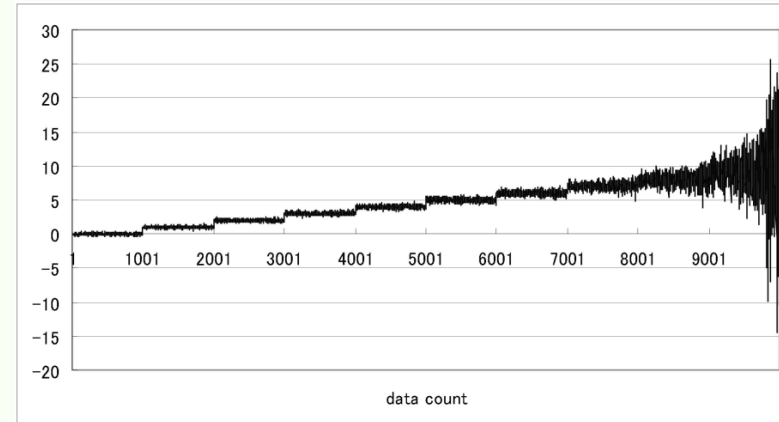
Score



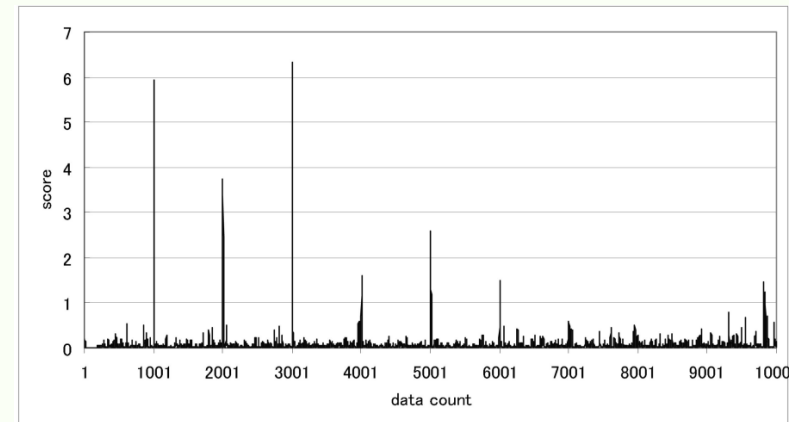
5. Simulation: Jumping mean with varying variance

- ◆ $t=1000x(x=1,2,\dots,9)$ の9回、平均を $\Delta x=1$ 変化(加算)
- ◆
$$\sigma^2 = \frac{0.1}{(0.01 + (10000 - t) / 10000)}$$
 - ◆ 時間経過とともに $0.1 \rightarrow 10$ に増加
- ◆ x 回目の変化前後でのKL情報量は $K(x) \approx 0.4(10 - x)^2$
- ◆ $t=6000$ まで検出できているこのとき $K(x) \doteq 6.4$

元データ

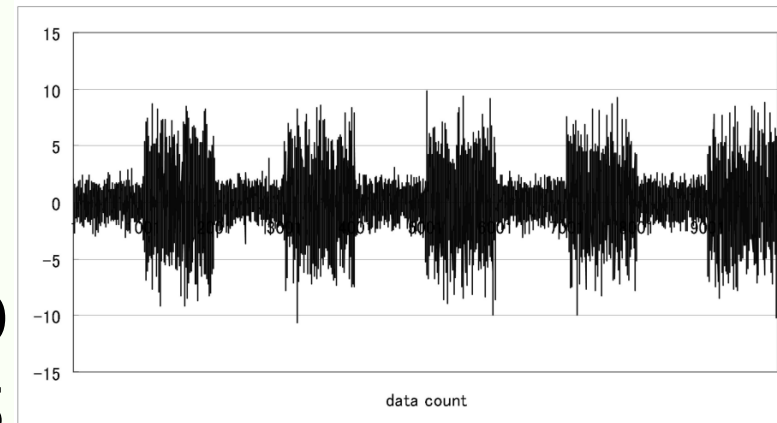


Score

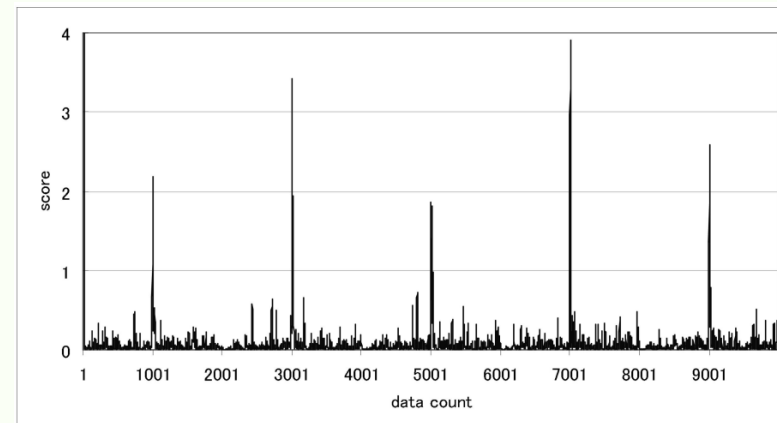


- ◆ 平均は0で一定
- ◆ $t=1000x(x=1,2,\dots,9)$ の9回、 σ^2 を1.0と9.0とで行き来させる
- ◆ 1.0→9.0のときのKL情報量 ≈ 2.90
- ◆ 9.0→1.0のときのKL情報量 ≈ 0.65
- ◆ 分散が増大するときには変化点が検出できているが、減少するときには検出できていない

元データ



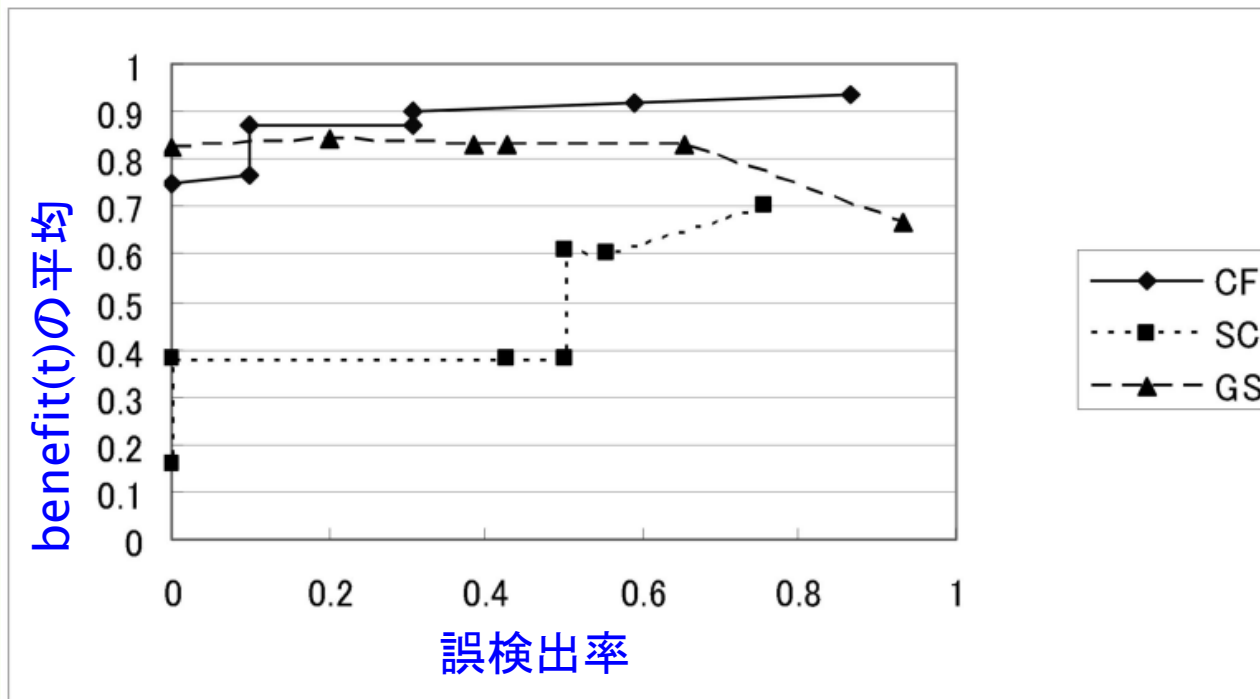
Score



- ◆ Activity Monitoring[Fawcett and Provost, 1999]に基づき、検出力を検出速度と誤検出率の観点から評価
- ◆ Scoreに対して設ける閾値によって検出力が変動する
→ Effective AlarmとFalse Alarm Rateとを軸としたROC曲線のようなことを考える
 - ◆ Effective Alarm : 変化点から20タイムステップ以内に検出
 t^* を真の変化時刻、 t を検出時刻としたとき、 $benefit(t) \equiv 1 - (t - t^*) / 20$
 - ◆ Alarm Policy :
最新のAlarmから20タイムステップ以内のAlarmは無視
 - ◆ SCやGSの場合は t^* より前の時点で検出することがあるので、
 $benefit(t) \equiv 1 - |t - t^*| / 10$ とする
 - ◆ False Alarm Rate : False Alarm / 全Alarm

◆ Jumping mean with constant variance

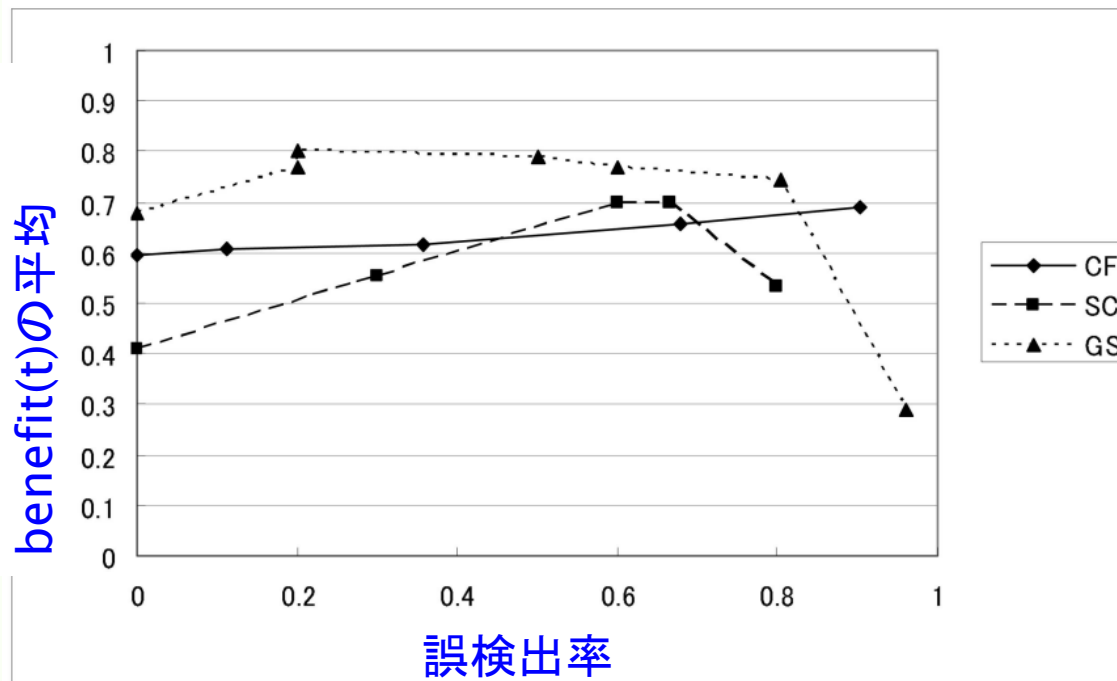
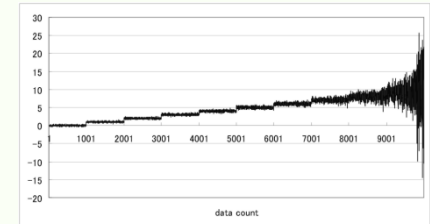
- ◆ データセットは前出同様、
ただしchange sizeは5で固定している



★CFの性能が高い

◆ Jumping mean with varying variance

- ◆ データセットは前出同様

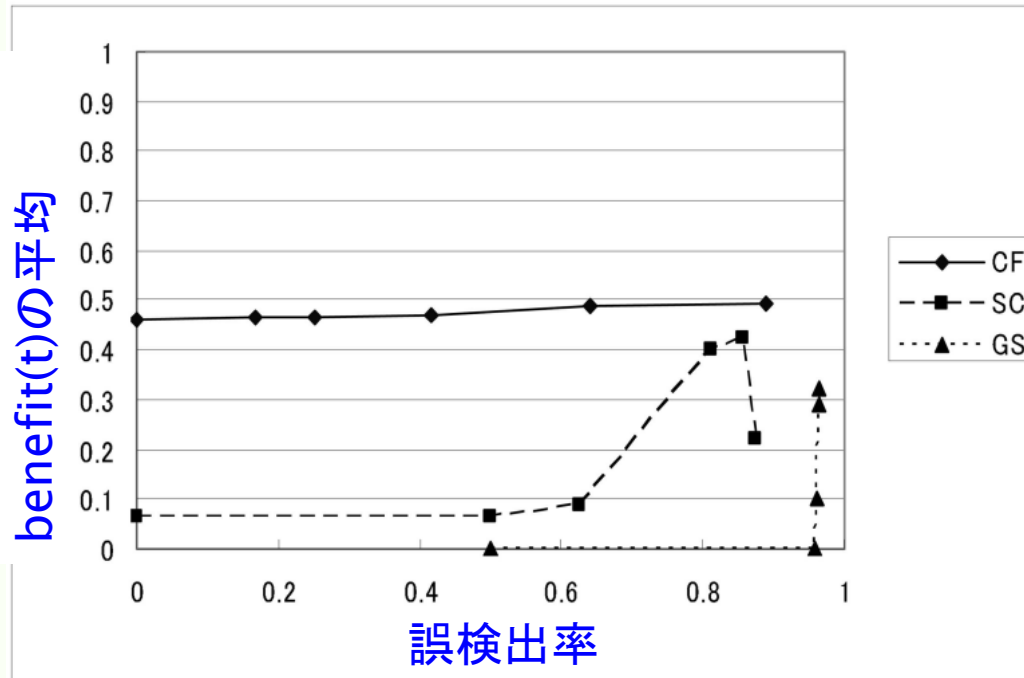
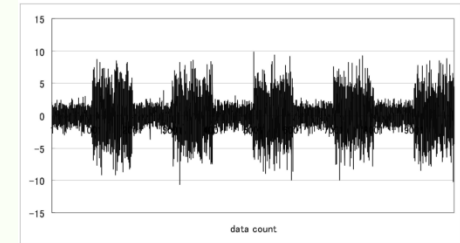


★CFよりもGSの性能が高い

- ◆ ただし、分散が徐々に縮小していくパターンを試すと、CFのみが変化点を検出できた(と書いてあるがグラフは無い)

◆ Jumping variance with constant mean

- ◆ データセットは前出同様



★CFの性能が高い

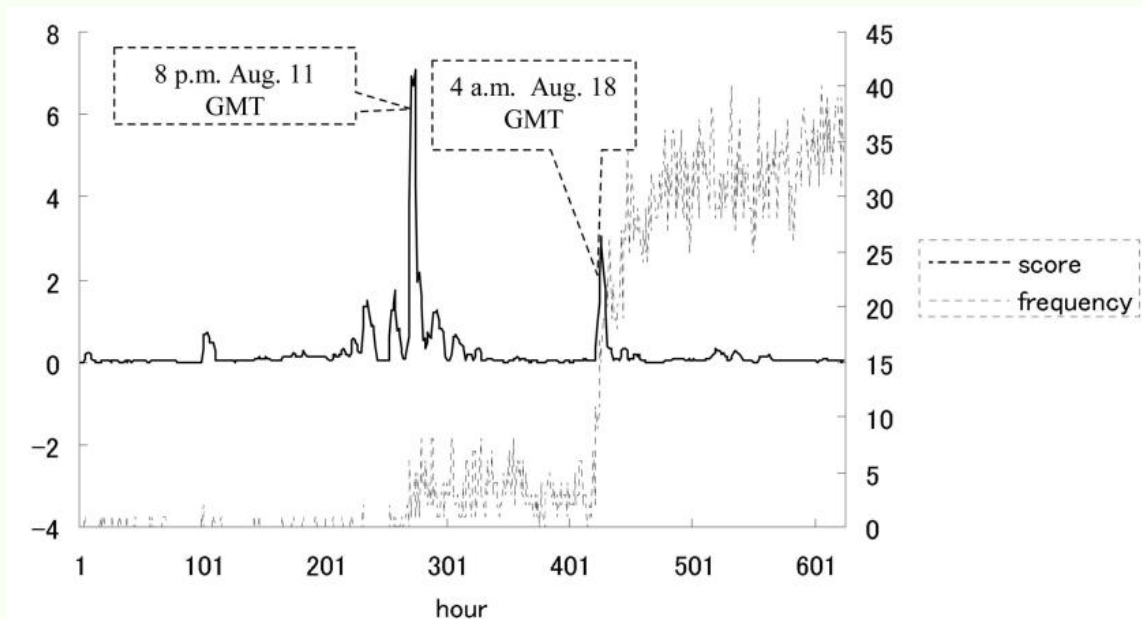
- ◆ ただし、分散が小さくなる(KL情報量が小さい)変化点を検出できていないのでbenefitが低い
- ◆ そもそもKL情報量に依存する手法であることの課題

◆ 考察

- ◆ 定性的には、CFはSCやGSよりも変化点検出に遅れの出にくい方法といえる。
- ◆ SCやGSで、False Alarm RateをあげてもBenefitが低下することがある。これは、Alarmの発生が真の変化点よりも早すぎるが多くなるため。
- ◆ CFはSCやGSに匹敵する検出力を持っている。特に3番目のデータセットと、2番目のデータセットの分散を小さくしていくパターンでは、相対的に高性能といえる。
- ◆ 計算のオーダーは、SCやGSが $O(n^2)$ であるのに対し、CFは $O(n)$ であるため、リアルタイム分析にはCFが適している。

◆ DoS攻撃の実例に適用

- ◆ パラメータ設定はSimulationと同様、Scoringは対数損失



◆ CFにより、変化点を正しく検出できた。

- ◆ 8月11日のAlarmは、現実のトレンドマイクロの発表より1時間44分早く、8/18のAlarmも約1日の遅れに留まっている
- ◆ SCやGSでも変化点は検出できるが、CFより検出時刻が遅い

◆ 成果

- ◆ 2段階の学習プロセスからなるChangeFinderの開発
- ◆ ChangeFinderへのARモデルとSDARモデルの適用
- ◆ 外れ値検出と変化点検出の統合的な枠組みの構築
- ◆ いくつかのデータによる挙動の検証

◆ 課題

- ◆ T の最適化、データセットからの T 自体の学習
- ◆ 分散が減少する変化点の検出法
- ◆ ARIMAモデルや状態空間モデルへの拡張
- ◆ 隠れマルコフモデル等との統合