



時系列データからの 因果関係の発見

NTTコミュニケーション科学基礎研究所

研究員

近原 鷹一

自己紹介

- Name: 近原 鷹一 (ちかはらよういち)
 - ✓ 2013.03: 慶應義塾大学理工学部生命情報学科 (舟橋研)
 - ✓ 2015.03: 東京大学大学院情報理工学系研究科 コンピュータ科学専攻 (宮野研)
 - システム生物学、バイオインフォマティクス
 - ✓ 2015.04-: NTT コミュニケーション科学基礎研究所 (CS研)
 - ✓ 2019.10-: 京都大学情報学研究科知能情報学専攻 (鹿島研; 社会人博士課程)
 - 機械学習・因果推論(因果探索・因果効果推定)

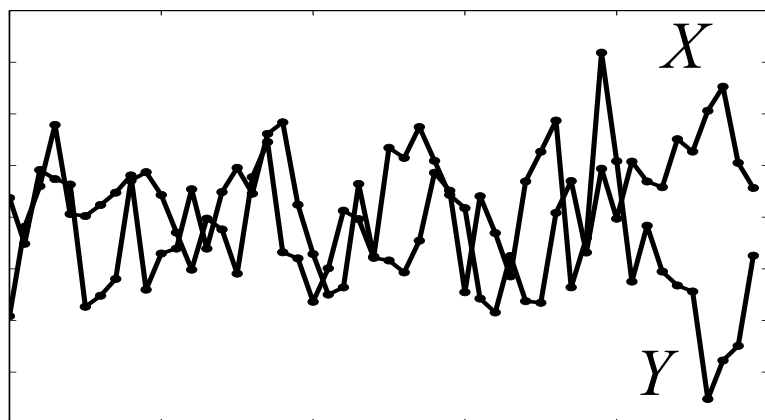


本日の内容

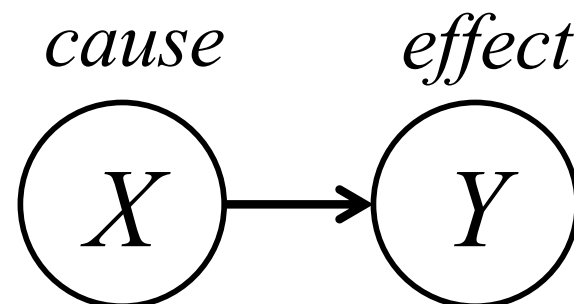
- 時系列データからの因果発見とは
 - ✓ 応用例
- Granger causalityの推定
 - ✓ 因果関係とは – Granger causalityの定義 –
 - ✓ 教師あり学習に基づく Granger causalityの推定
[Chikahara+; IJCAI2018(AI分野最難関会議)]
- Granger causality研究の最前線
 - ✓ イベント間の因果関係 – 連続時間系のGranger –

時系列データからの因果関係の発見

- 与えられた時系列データから変数間の**因果関係**を推定するタスク



入力: 時系列データ

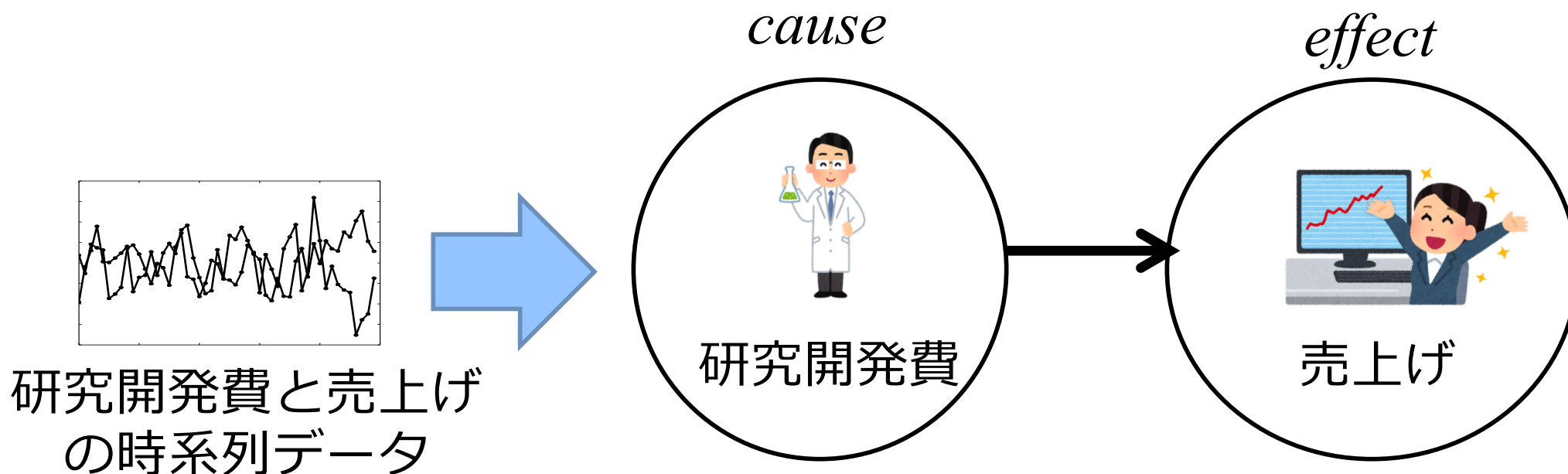


出力: **因果関係**

目的: 因果関係に関する**知識の発見**

Application 1: Economics

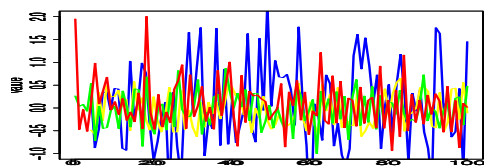
- 「研究開発費は売上金に影響を与える」
という知識は企業にとって有用



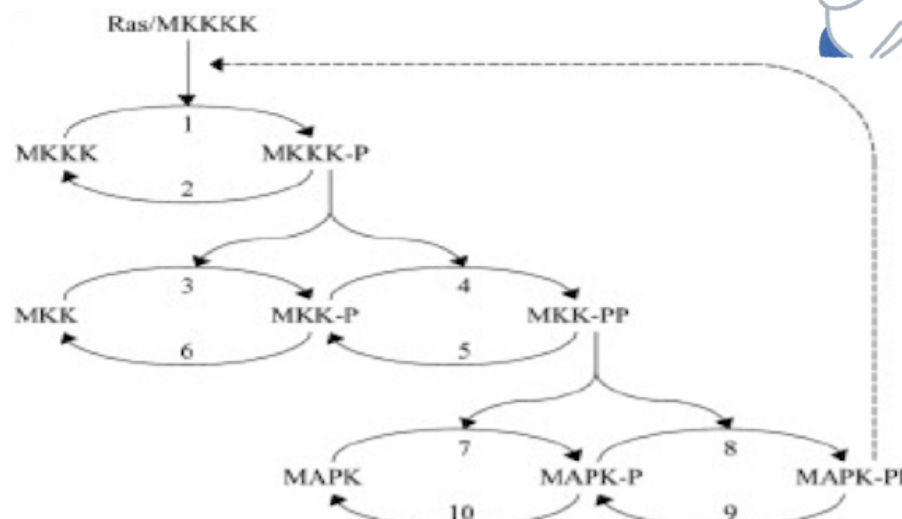
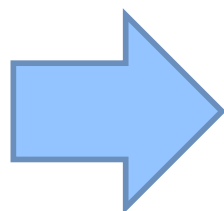
※逆ではない
(売上げは研究開発費に影響を与えない)

Application 2: Bioinformatics

- 遺伝子の発現量の時系列データから「遺伝子の制御関係」がわかると新規薬剤の探索に役立つ



遺伝子発現量
に関するデータ



遺伝子の制御関係

“因果関係”とは何か？

どう定義されるのか？

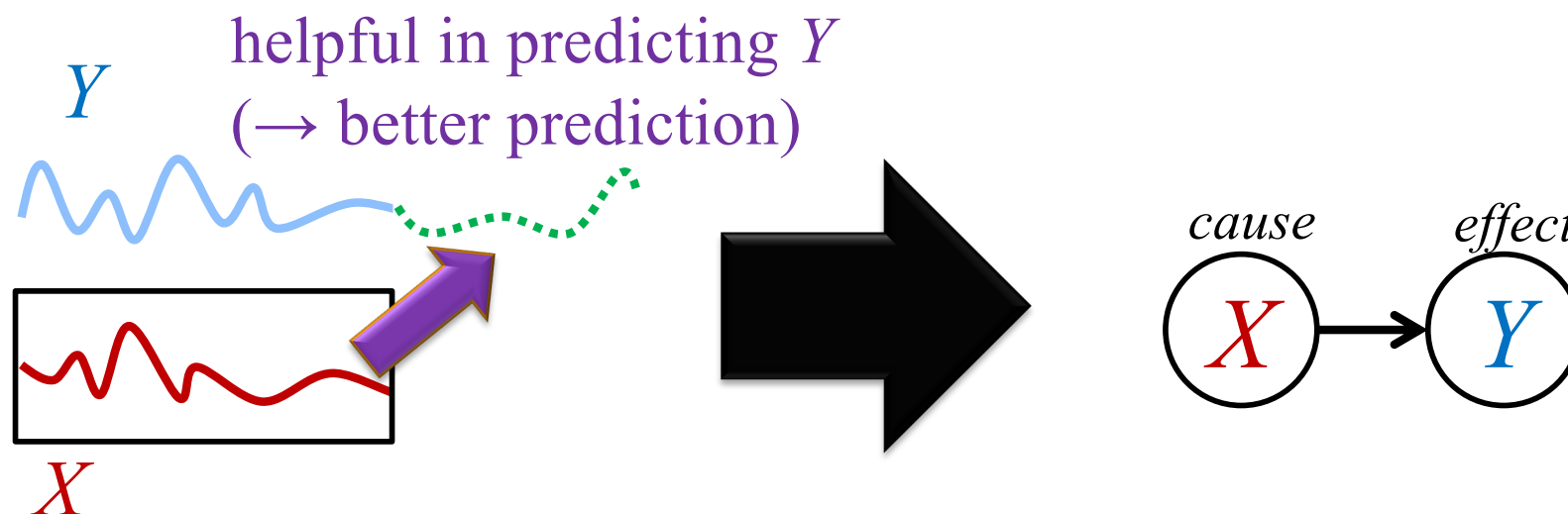
時系列における因果関係の定義

Granger causality [Granger1969]

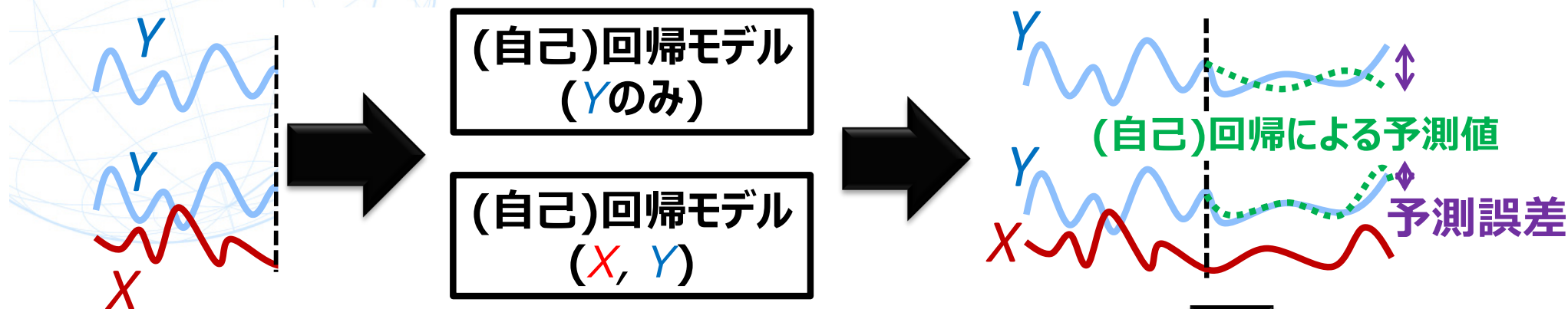


X is the cause of Y

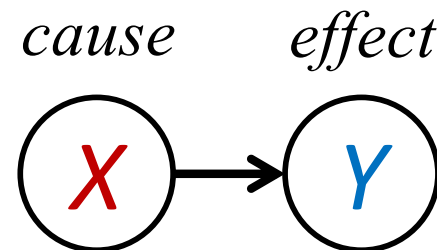
if the past values of X are **helpful in predicting**
the future values of Y



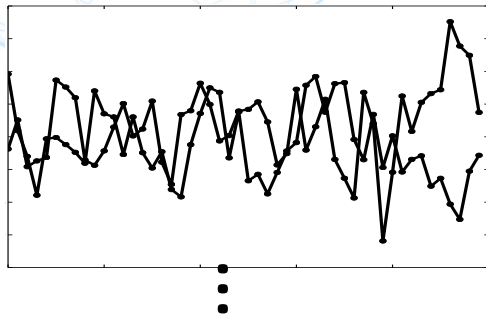
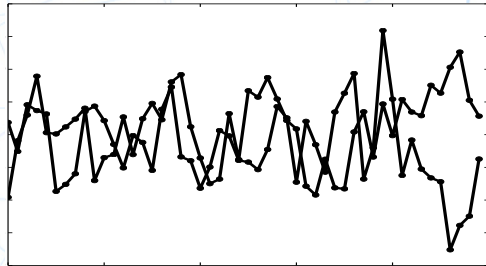
既存手法によるGranger causalityの推定



過去の X の値を用いることで
予測誤差が有意に小さくなれば,
 X は Y の原因



既存手法の問題点



(自己)回帰モデル



どの回帰モデルを使えばいいの？

VAR
モデル

GAM

ガウス過程

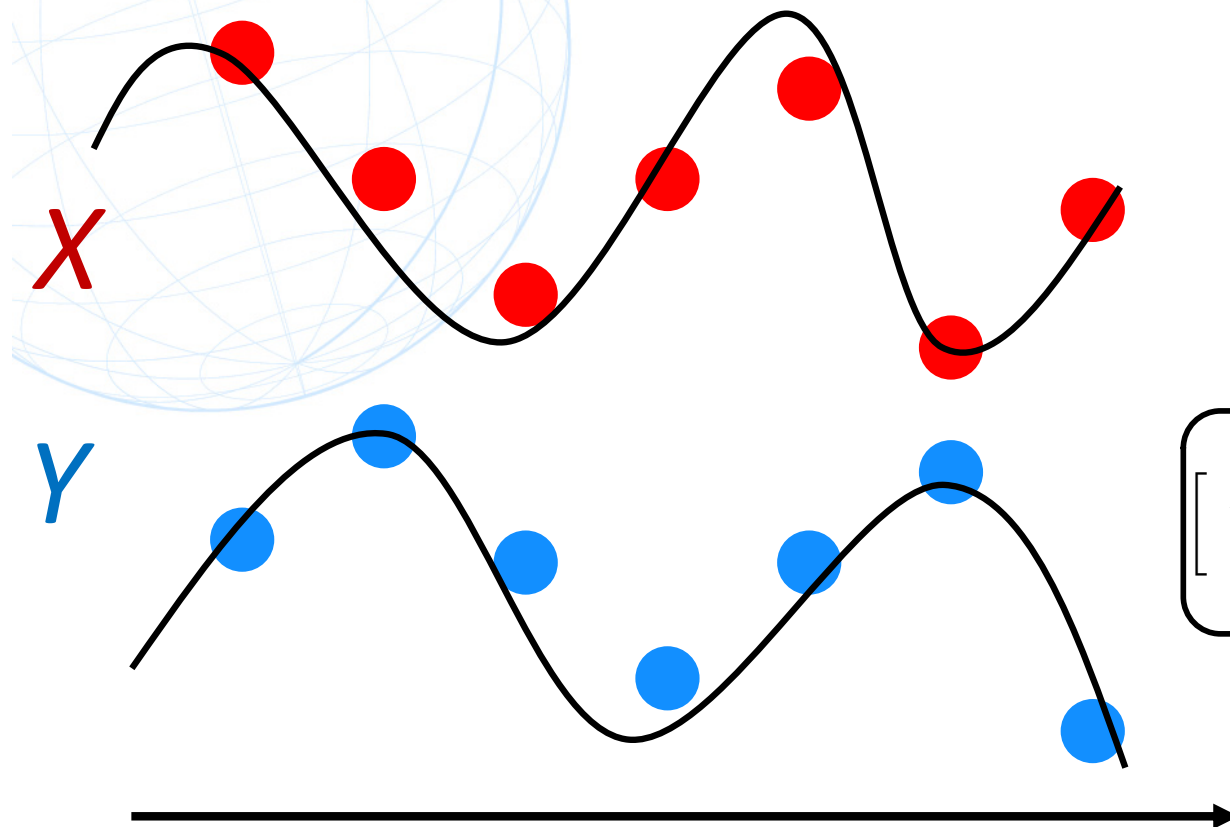
⋮

$$\begin{bmatrix} X_t \\ Y_t \end{bmatrix} = \frac{1}{P} \sum_{\tau=1}^P A_{\tau} \begin{bmatrix} X_{t-\tau} \\ Y_{t-\tau} \end{bmatrix} + \begin{bmatrix} E_{X_t} \\ E_{Y_t} \end{bmatrix}$$

(自己)回帰モデルを各データに対して適切に選ばなければ
正しくGranger causalityの方向・有無を推定できない

本研究のゴール：回帰モデルの選択が不要なアプローチの提案

(補足) 既存のGranger causality推定



VARモデルの場合

$$\begin{bmatrix} X_t \\ Y_t \end{bmatrix} = \frac{1}{P} \sum_{\tau=1}^P A_{\tau} \begin{bmatrix} X_{t-\tau} \\ Y_{t-\tau} \end{bmatrix} + \begin{bmatrix} E_{X_t} \\ E_{Y_t} \end{bmatrix}$$

↑
モデルの係数

データ点に合わせて回帰モデルを学習する (最小二乗法)

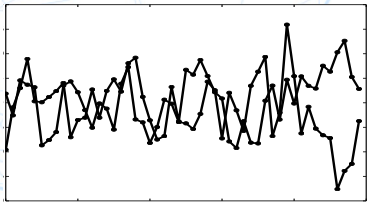
モデルの係数のゼロ/非ゼロで
判定

予測誤差に関する統計量で
有意差判定(仮説検定)

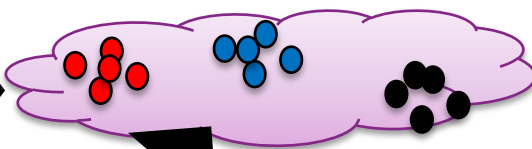
提案手法： 教師あり学習に基づく Granger causality の推定

テストデータ

(因果関係を知りたいデータ)



特徴ベクトル



分類器

$X \rightarrow Y$

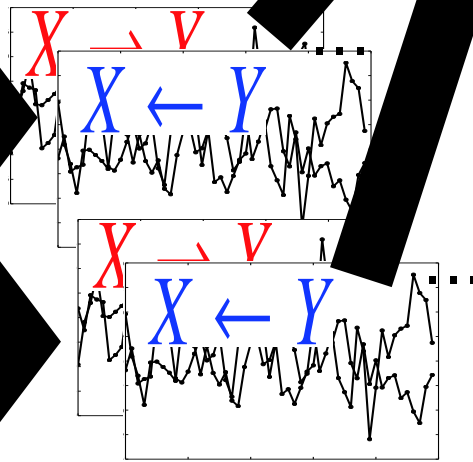
$X \leftarrow Y$

因果なし

サンプリング

線形
モデル
(VAR)

非線形
モデル



訓練データ

(因果関係が既知のデータ)

さまざまな(人工)データ
を活用

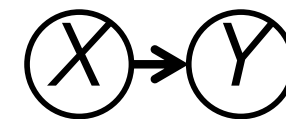
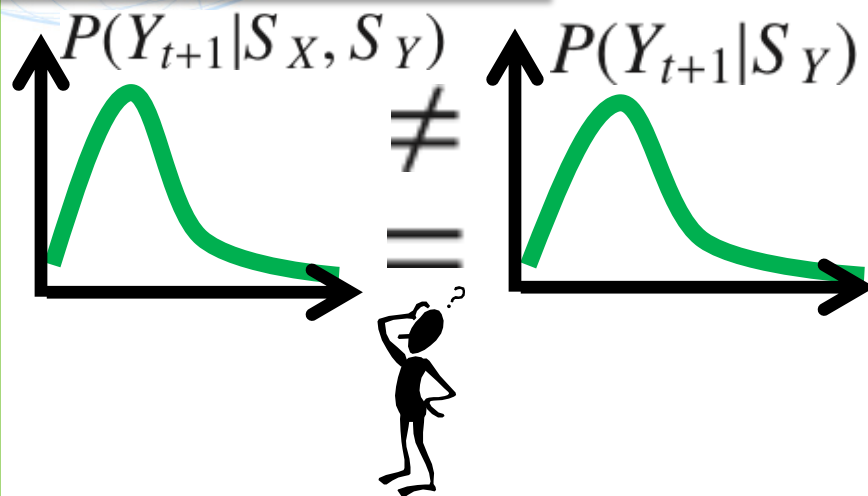
さまざまな
確率モデル

どのように特徴ベクトルを得るか？

– Granger causalityの定義から考える –

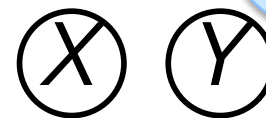
Granger因果性の有無・方向は
2つの(条件付き)分布が等しいか否かで決まる

Granger因果性の
フォーマルな定義



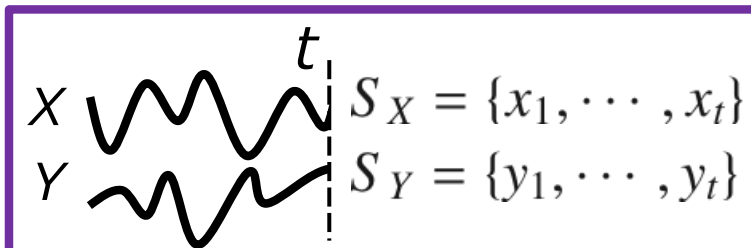
if $P(Y_{t+1}|S_X, S_Y) \neq P(Y_{t+1}|S_Y)$

Xを用いると
Yの予測が異なる

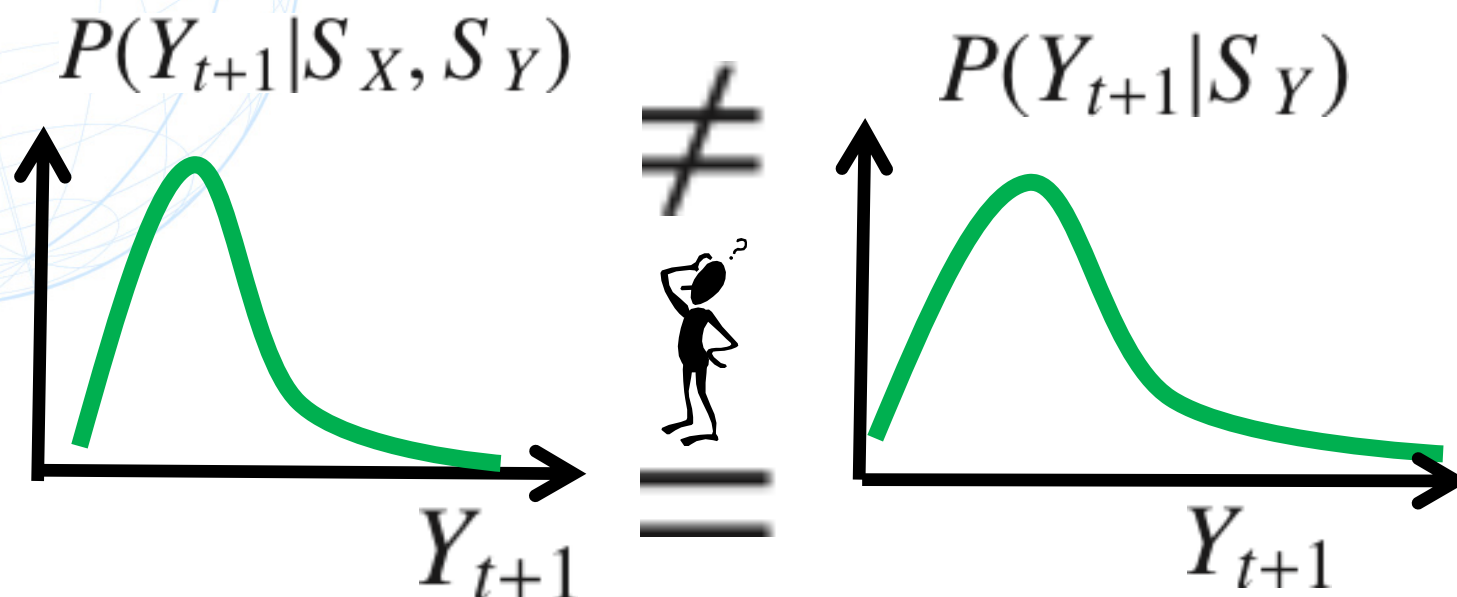


if $P(Y_{t+1}|S_X, S_Y) = P(Y_{t+1}|S_Y)$

Xがあってもなくても
Yの予測は変わらない



2つの(条件付き)分布が等しいか否かで Granger因果性の有無はわかる

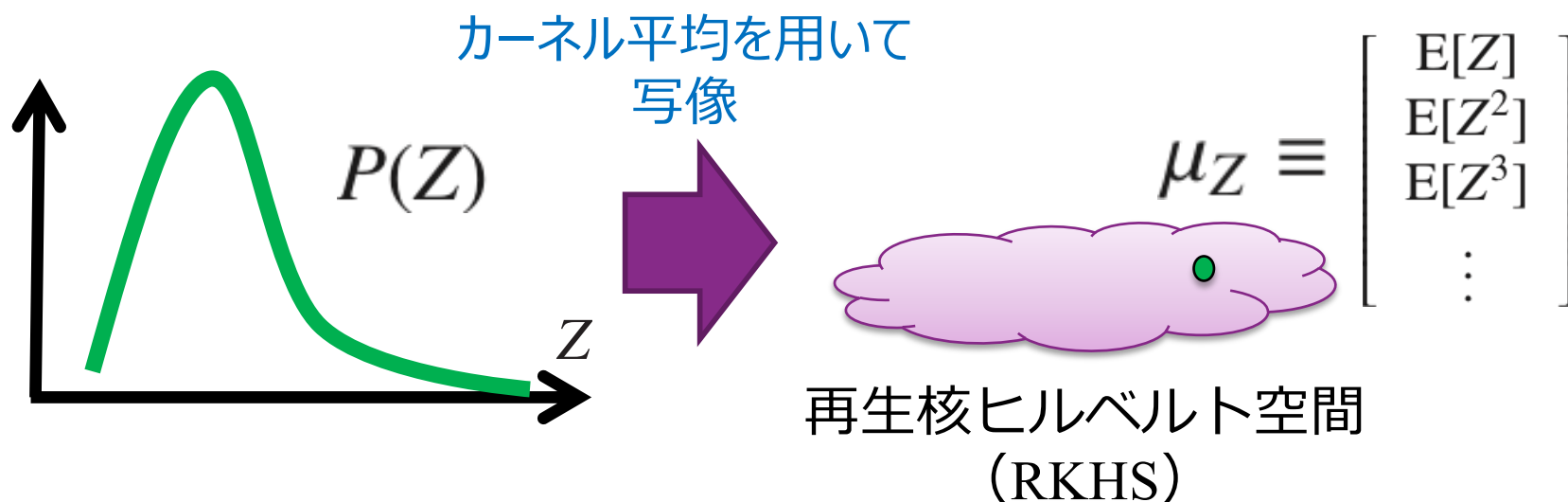


2つの分布が等しいか否かを
判定する特徴量が必要

分布間距離尺度MMDを使う

MMD: 条件付き確率分布の違いをkernel mean間の距離として定量

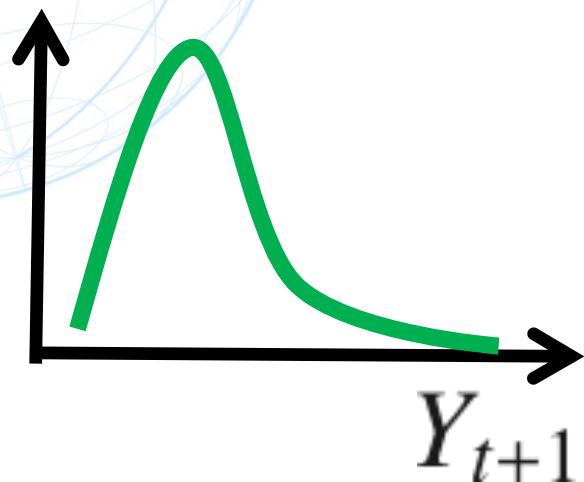
- kernel meanとは：分布の平均, 分散, ...を返す関数



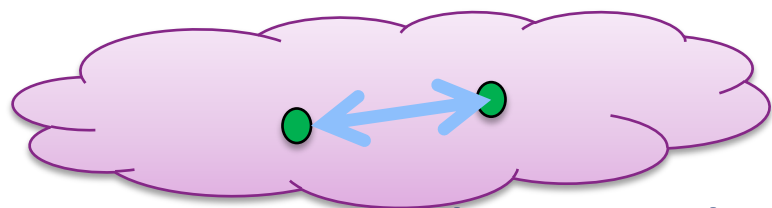
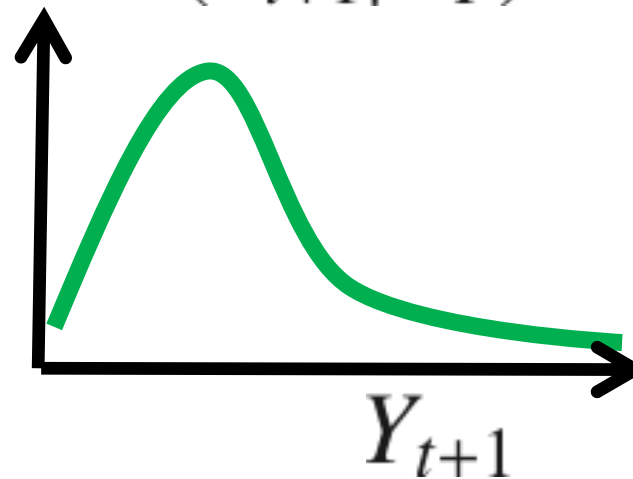
※特殊なカーネル(Gaussian kernel等)を用いれば
異なる2つの分布は同一の点に写像されない(写像による情報損失がない)

2つの(条件付き)分布が等しいか否かで Granger因果性の有無はわかる

$$P(Y_{t+1}|S_X, S_Y)$$

 \neq  \equiv

$$P(Y_{t+1}|S_Y)$$



Distance (MMD)

(maximum mean discrepancy
[Gretton+; NIPS2007])

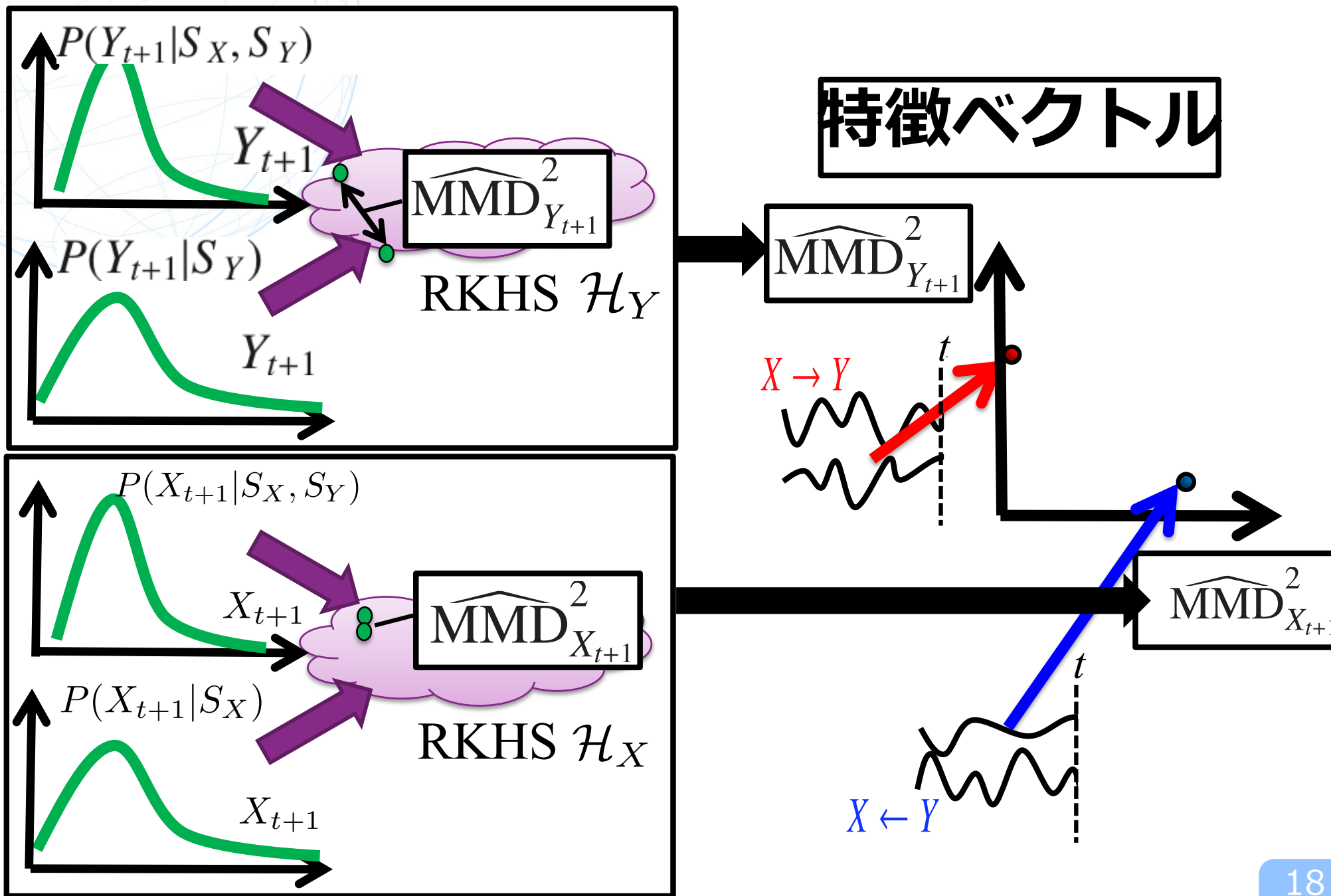
 \neq

0

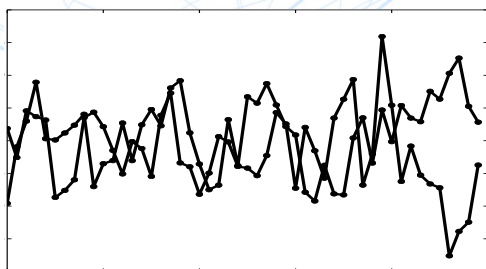
 \equiv

0

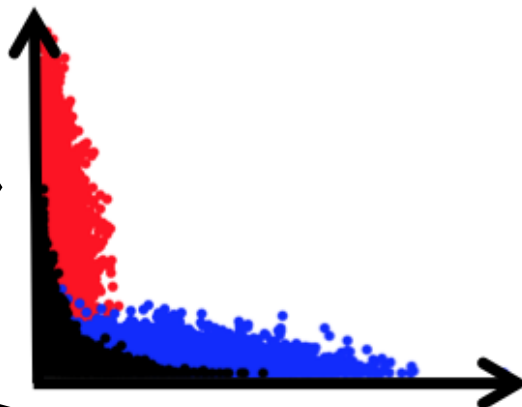
Kernel meanを用いて確率分布の特徴を抽出



テストデータ



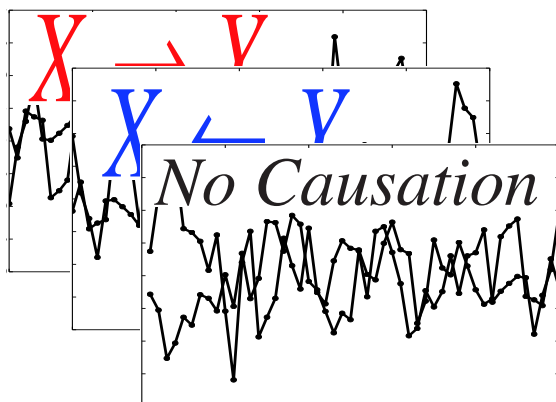
分類器
(e.g., Random Forest)



$X \rightarrow Y$

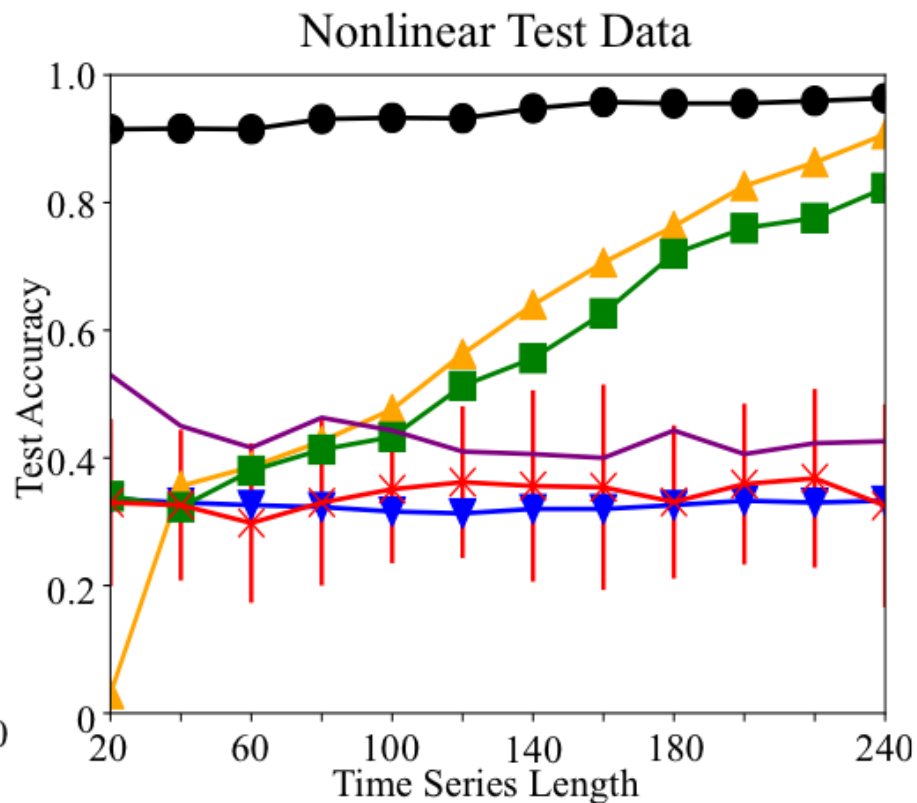
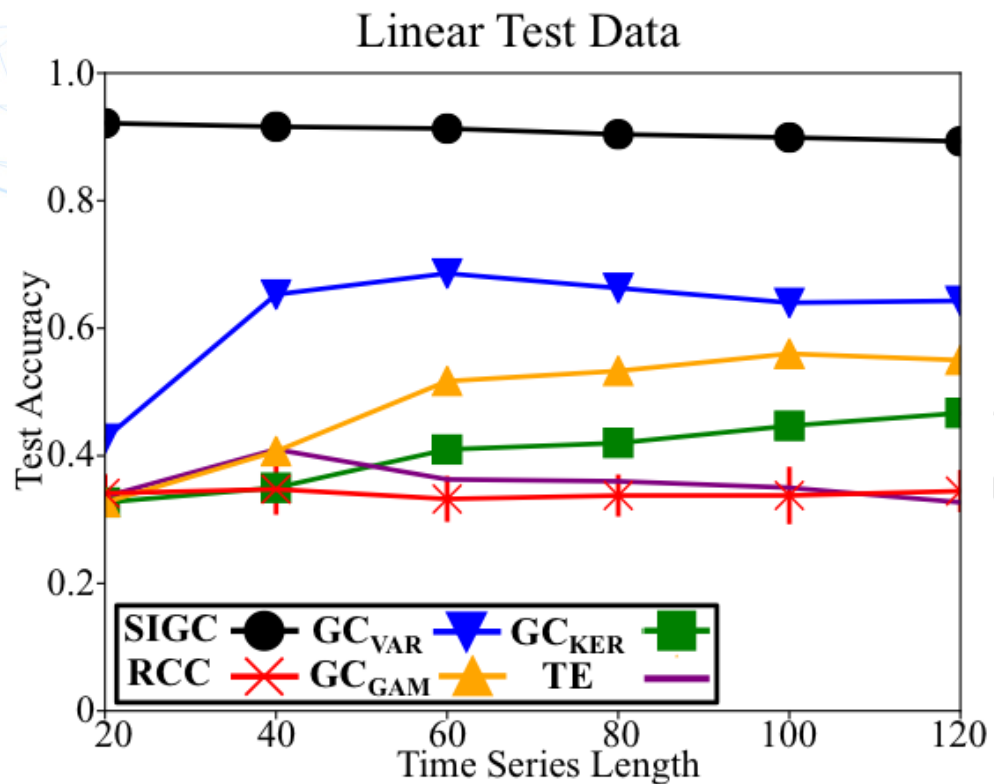
$X \leftarrow Y$

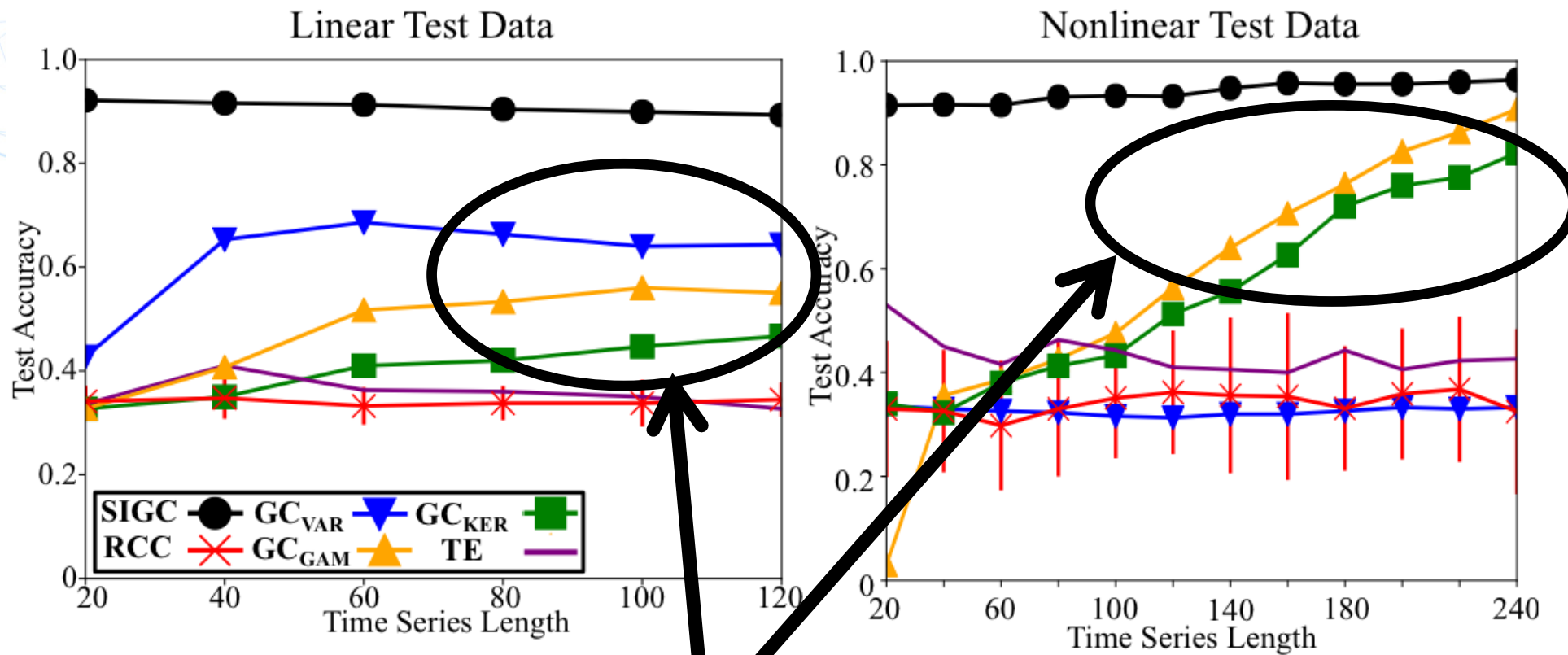
訓練データ
(人工データ)



No Causation

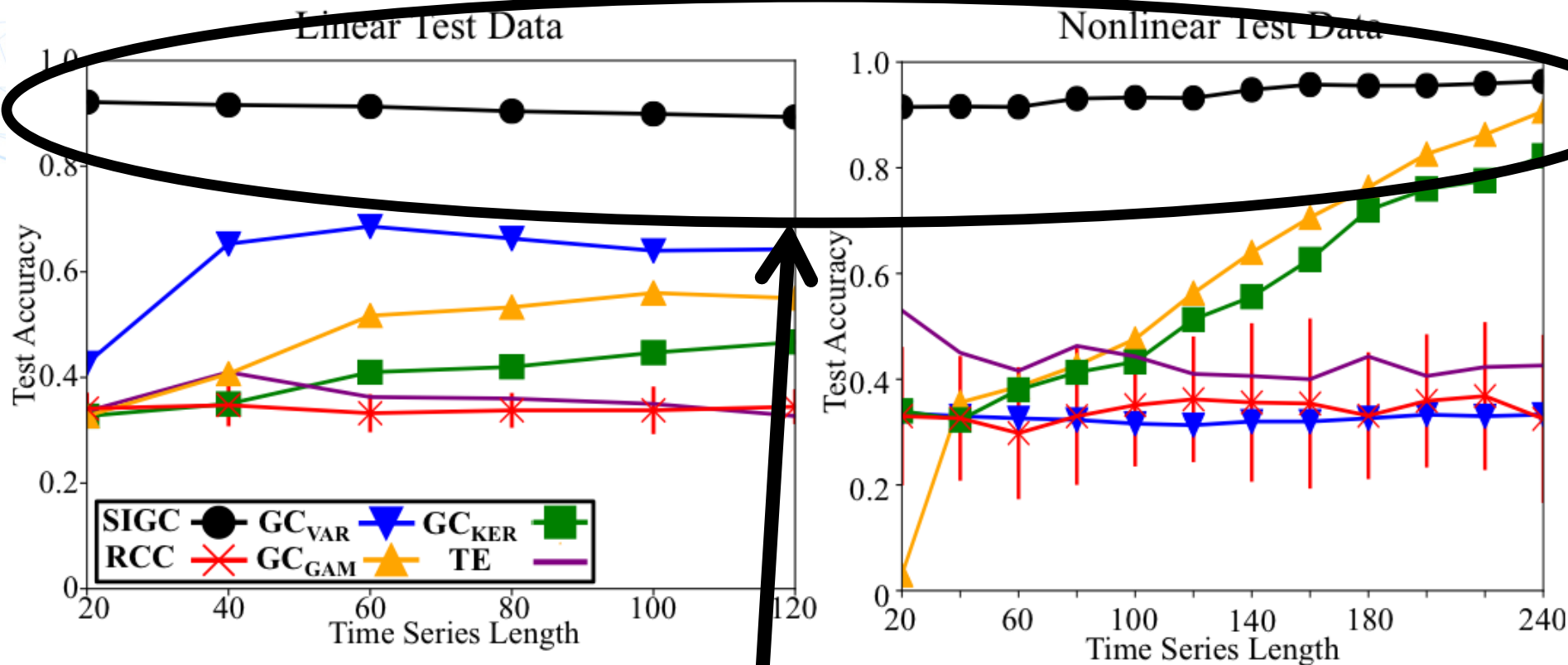
人工データ実験





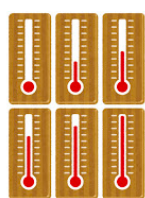
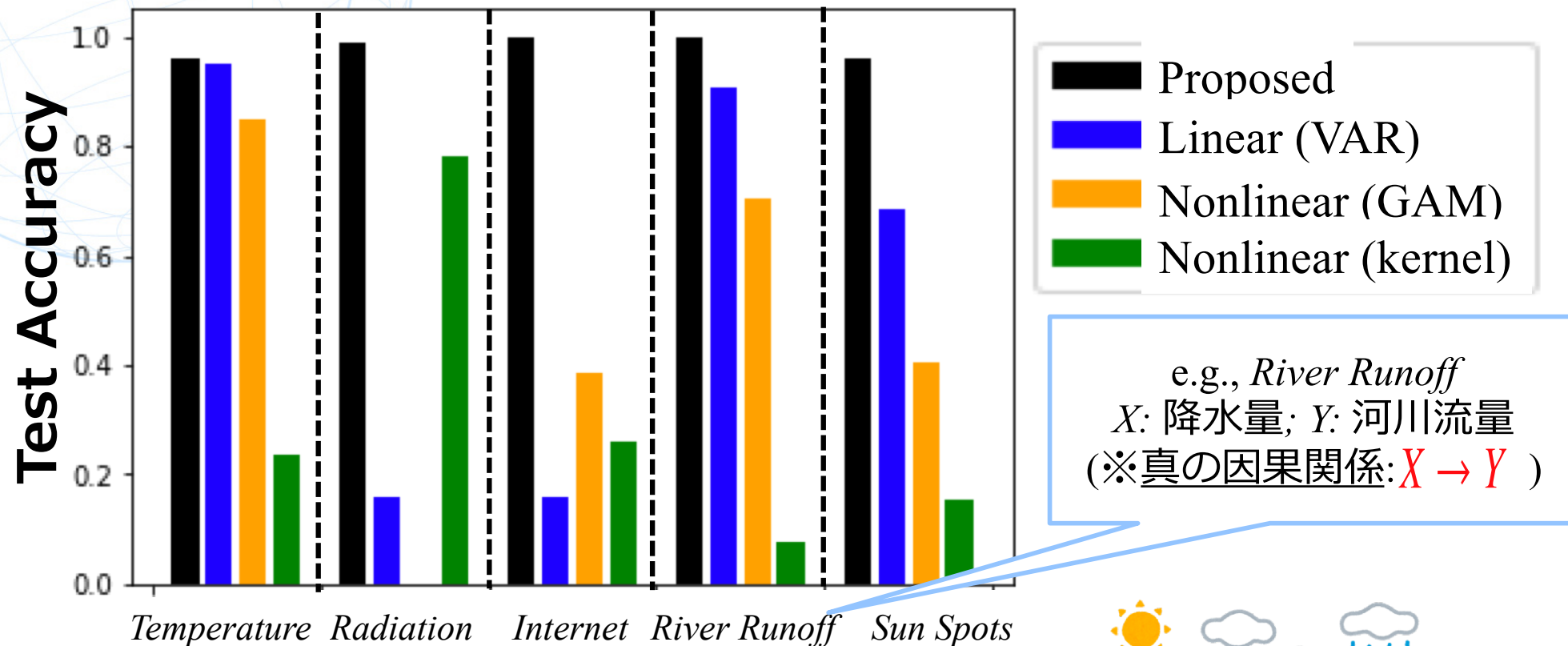
既存手法

推定精度は選んだ回帰モデルの当てはまり度合いに強く依存



提案手法

線形・非線形データの両方で高い推定精度

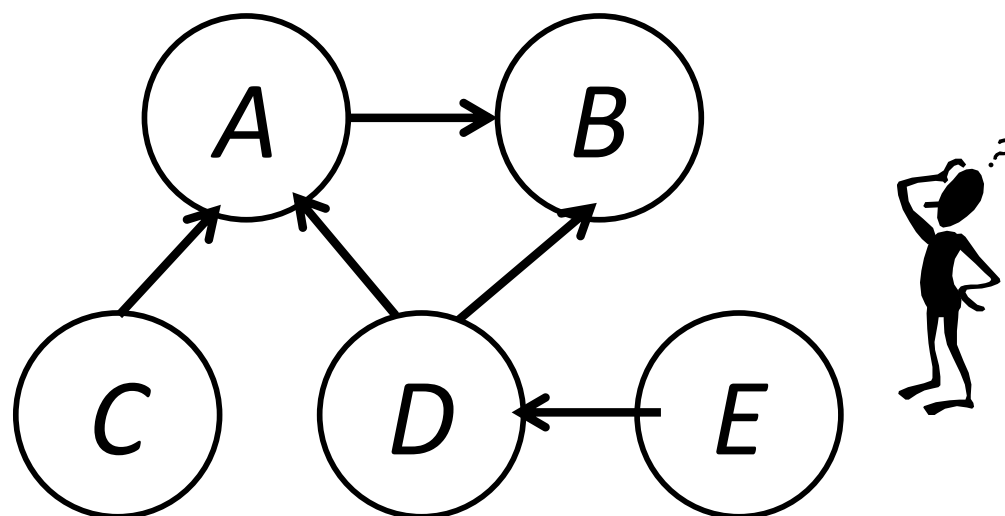


5つの実データセット

多変数(n変数)への拡張

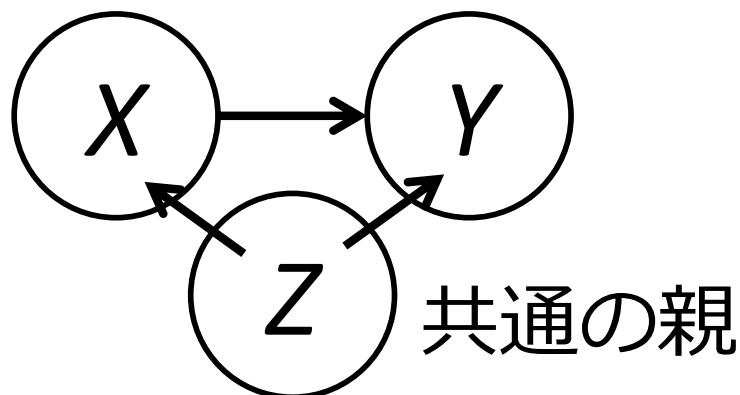
- ${}_n C_2$ 種類の各ペアに因果関係のラベルを割り当てる

$X \rightarrow Y$ $X \leftarrow Y$ *No Causation*



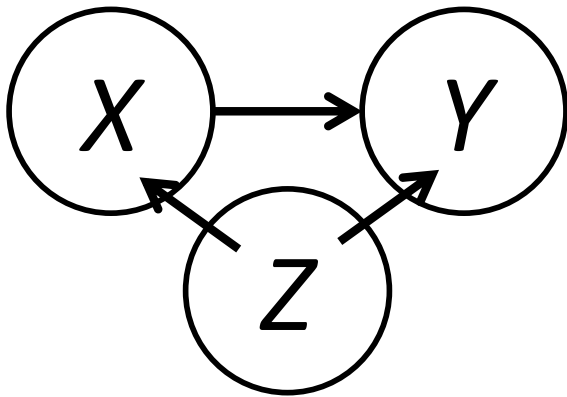
(補足) 多変数(n 変数)への拡張

- 仮定：各ペアの共通の親は高々1つ
→ 第3の変数の影響のみ考えて3つ組をとる

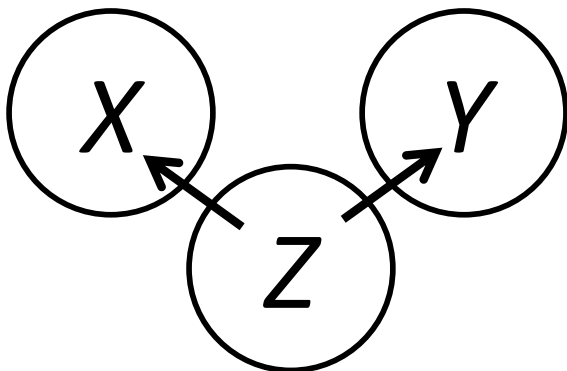


- ペア(X, Y)に対するラベルの割り当て方：
 - ✓ $(n-2)$ 種類の各3つ組 (X, Y, Z_v) ($v \in \{1, \dots, n-2\}$) に対し
 1. 特徴ベクトルを計算
 2. ラベル確率を計算
 - ✓ ラベル確率の平均が最大のラベルを出力

- **Conditional Granger causality** [Geweke JASA1984]:
第3の変数 Z で条件付けた2つの条件付き分布を比較



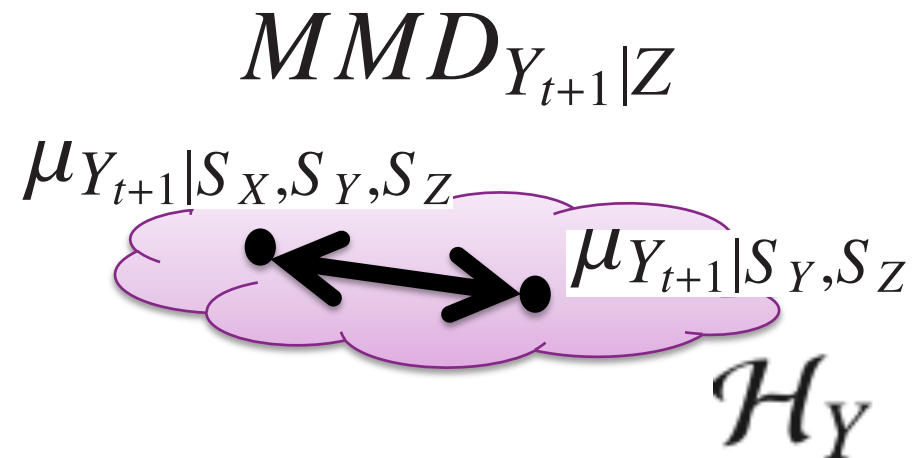
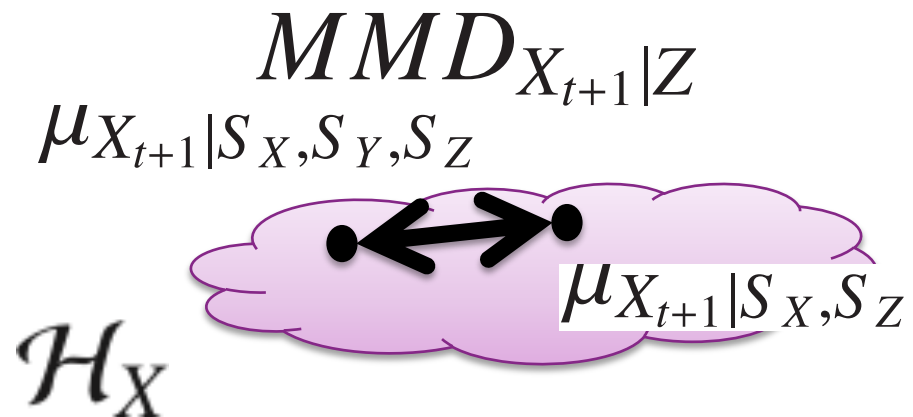
if $P(Y_{t+1} | \underline{S_X}, S_Y, \underline{S_Z}) \neq P(Y_{t+1} | S_Y, \underline{S_Z})$



if $P(Y_{t+1} | \underline{S_X}, S_Y, \underline{S_Z}) = P(Y_{t+1} | S_Y, \underline{S_Z})$

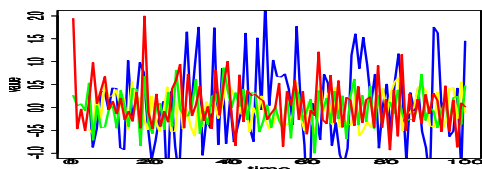
特徴ベクトルの計算

- 同様に求めた条件付き分布間の距離を使って分類の特徴ベクトルを得る

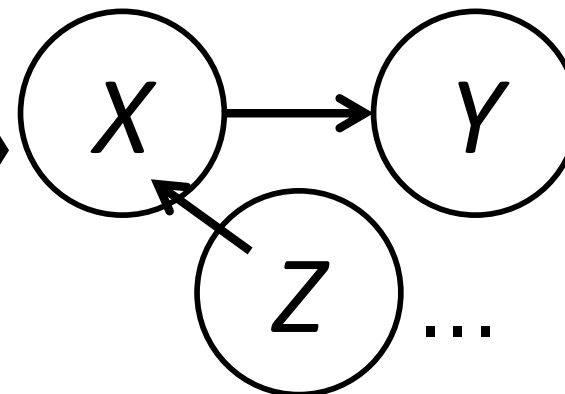


テストデータ

Yeast cell cycle gene expression data
[Spellman+ 1998]
14 variables (genes)



分類器



訓練データ
(大量の3つ組人工時系列)

※真の因果関係 : パスウェイデータベースKEGGを参照

結果(macro F1 score, micro F1 scoreで評価)

	SIGC_{tri}	GC _{VAR}	GC _{KER}	SIGC _{bi}	GC _{GAM}	TE	RCC
macro-averaged F1	0.483 (0.0)	0.351	0.437	0.431 (0.007)	0.457	0.430	0.407 (0.096)
micro-averaged F1	0.637 (0.0)	0.436	0.513	0.578 (0.011)	0.567	0.449	0.567 (0.161)

※Higher is better

既存手法よりは高いが**十分な推定精度ではない**

1. 仮定（共通の親が高々1つ）が強すぎる
2. **時系列が短すぎる** ($T = 57$), etc.

データ取得コストが大きいライフサイエンス分野において
高精度に因果関係の有無・方向を推定するのは未だ難しい

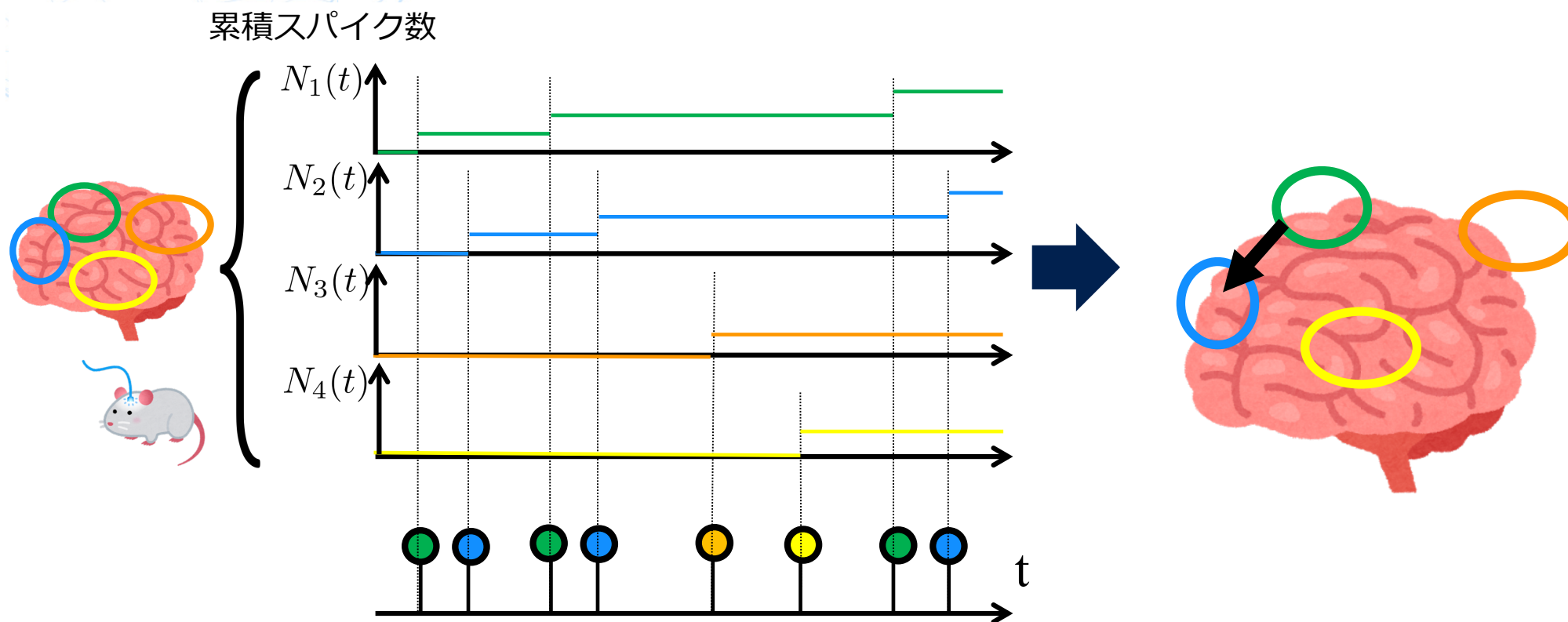


本日の内容

- 時系列データからの因果発見とは
 - ✓ 応用例
- Granger causalityの推定
 - ✓ 因果関係とは – Granger causalityの定義 –
 - ✓ 教師あり学習に基づくGranger causalityの推定
[Chikahara+; IJCAI2018(AI分野最難関会議)]
- Granger causality研究の最前線
 - ✓ イベント間の因果関係 – 連続時間系のGranger –

イベント時系列間のGranger causality

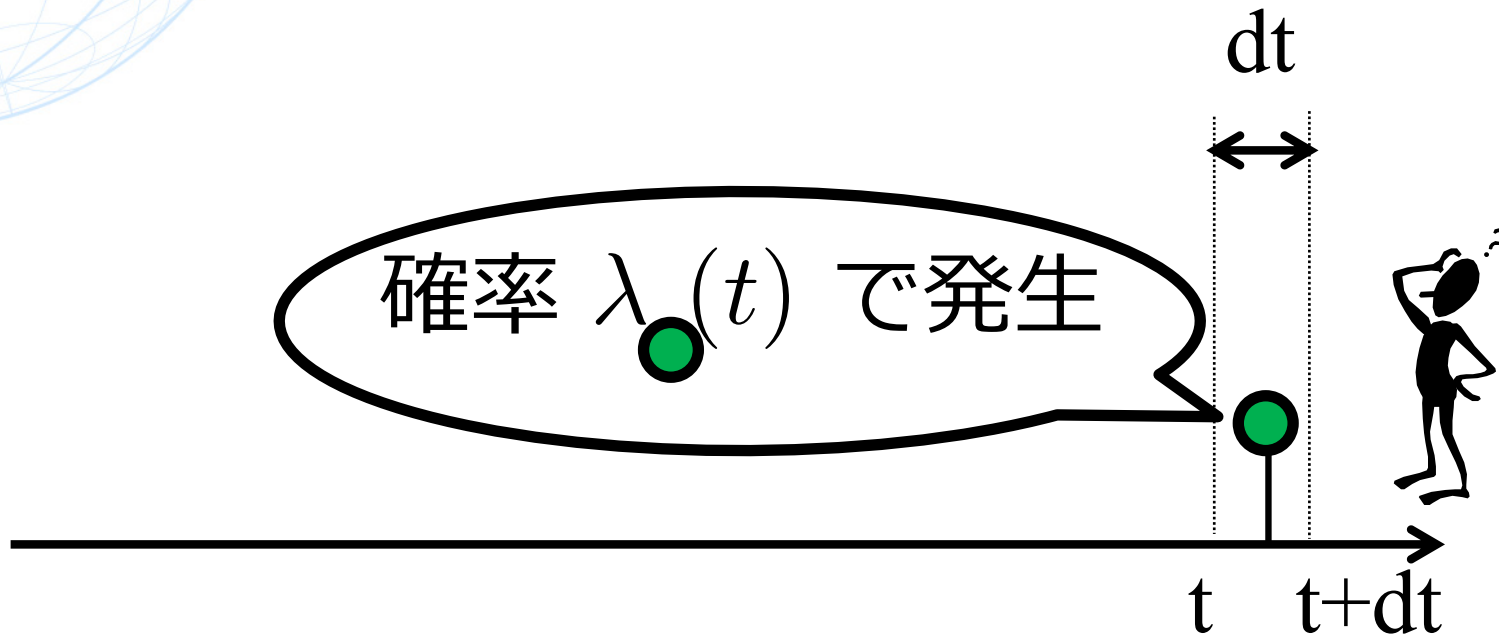
- 不均一な時間間隔で生じる離散的なイベントを考える
 - ✓ (例) イベント = 特定の脳領域での脳波スパイク



イベント間のGranger causalityは
点過程モデルを用いて表せる

点過程モデル(Point Processes)

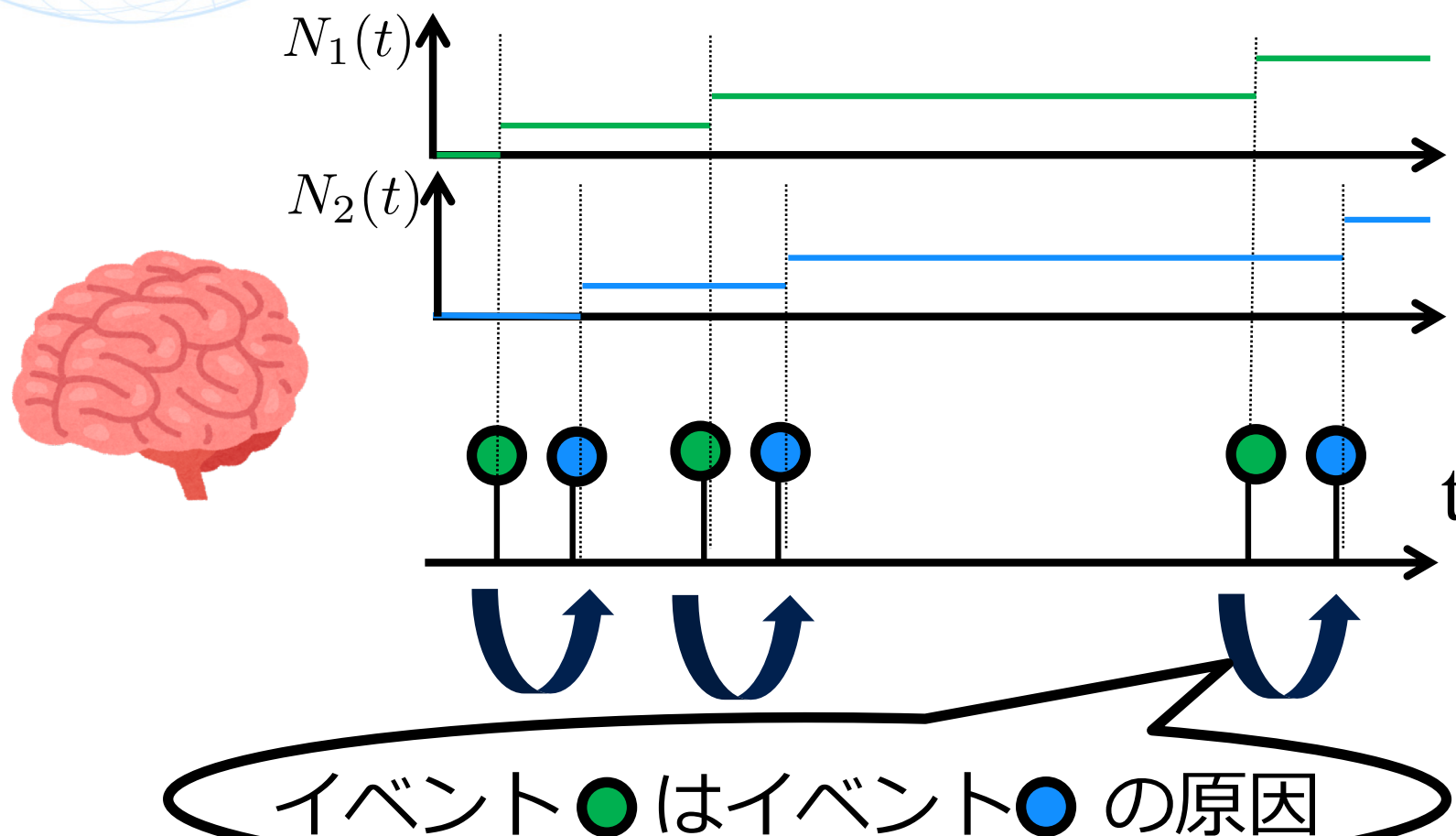
- 点過程モデル：微小時間 dt の間にイベントが生じる確率 $\lambda(t)$ をモデル化




Hawkes過程

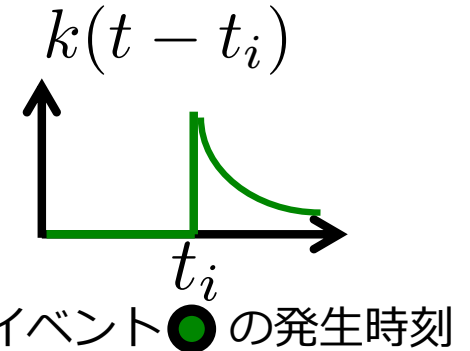
- Hawkes過程：イベント発生確率が他のイベントの発生に影響される点過程モデル

累積スパイク数



Hawkes過程


- イベント  の発生確率 :



$$\lambda_{\bullet}(t) = \mu_{\bullet} + \alpha_{\bullet} \sum_{t_i \in \mathcal{H}_{\bullet}(t)} k(t - t_i) + \alpha_{\bullet} \sum_{t_i \in \mathcal{H}_{\bullet}(t)} k(t - t_i)$$

時間によらず
一定確率で発生



過去の自分自身の
イベント発生に影響

過去のイベント  の
発生に影響

Hawkes過程によるGranger causalityの定義

- イベント  の発生確率 :

$$\lambda_{\bullet}(t) = \mu_{\bullet} + \alpha_{\bullet} \sum_{t_i \in \mathcal{H}_{\bullet}(t)} k(t - t_i) + \alpha_{\bullet} \sum_{t_i \in \mathcal{H}_{\bullet}(t)} k(t - t_i)$$

$\alpha_{\bullet} \neq 0$ ならば
イベント  はイベント  の原因
[Eichler+; 2016]

- **推定** : 最尤推定でパラメータ μ, α を推定

- 時系列の因果関係の定義の一つとして Granger causality をご紹介
 - ✓ 予測の「有用性」で定義
 - ✓ 離散的なイベントどうしの因果関係も定義可能
- 多種多様な応用が期待される一方、依然として多くの課題がある
 - ✓ 予測モデルの与え方
 - ✓ **非観測変数が影響を及ぼす場合（非常に難しいが重要）**

