

学習数理情報学研究室の紹介

東京大学 大学院情報理工学系研究科

創造情報学専攻/数理情報学専攻

山西健司

<http://www.ibis.t.u-tokyo.ac.jp/>

2014/6/7

学習数理情報学研究室

- 機械知能の本質を数理で捉える -



教授 山西健司 (H25より創造学専攻本務)

専門: 情報論的学習理論、Latent Dynamics
データマイニング応用、異常検知

講義: 確率数理工学(夏:学部)
情報論的学習理論(冬:大学院)
データマイニングによる異常検知(冬:大学院)



助教:松島慎

専門: 機械学習(大規模分類問題)
講義: プログラミング演習



特任助教: 森野佳生

専門: 機械学習, 非線型振動子
講義: 数理情報工学実験第二

大学院生:PD: 寺園, 田村

木村(D3), 梶野(D2),

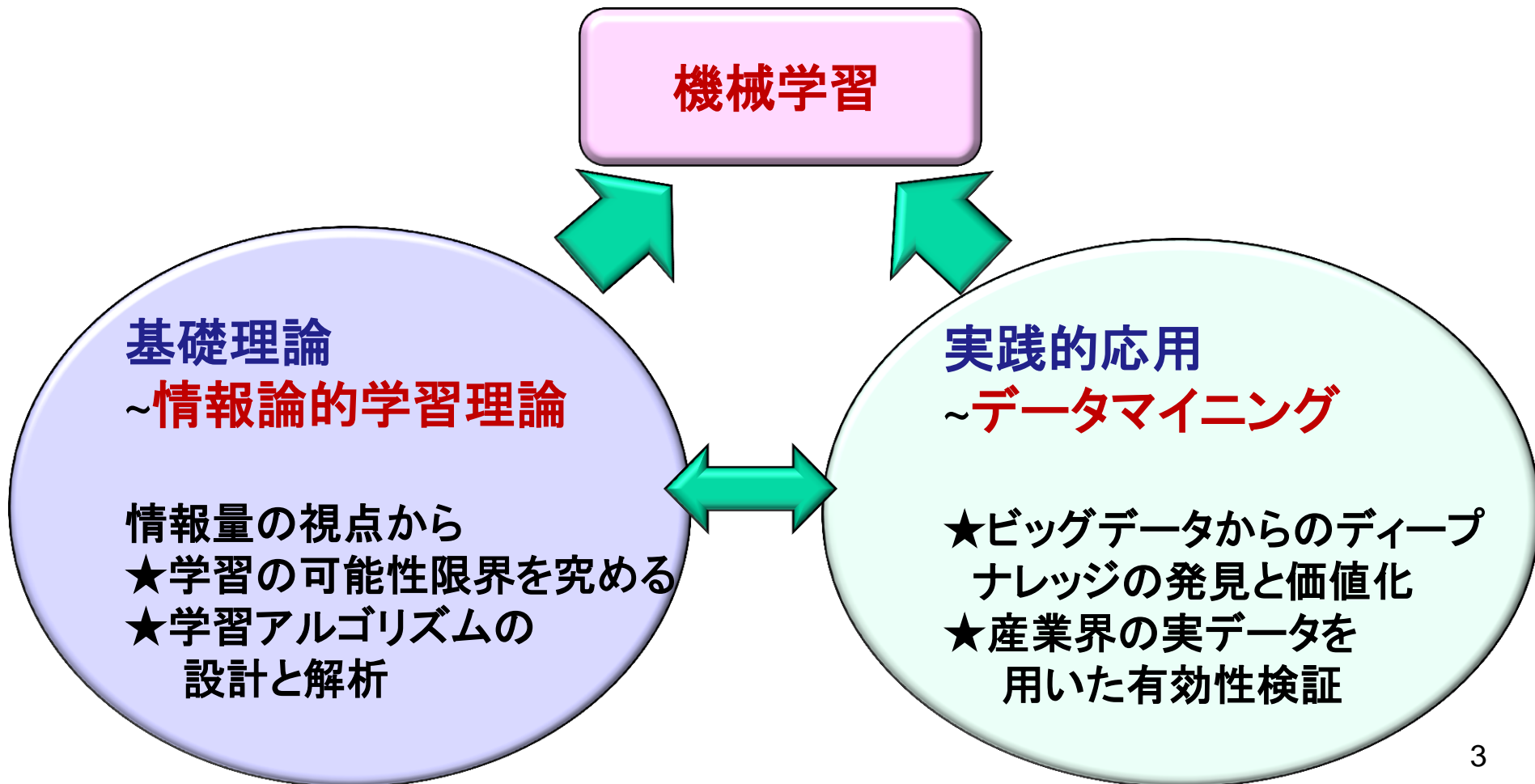
家入(M2), 伊藤(M2), 貝ヶ石(M2), 真矢(M2), 渡邊(M2),

石井(M1), 江田(M1), 尾亦(M1), 友田(M1), 中村(M1),

宮口(M1), 要名本(M1)

学習数理情報学における両輪

情報論的学習理論とデータマイニング実践を両輪とする



研究室の特徴

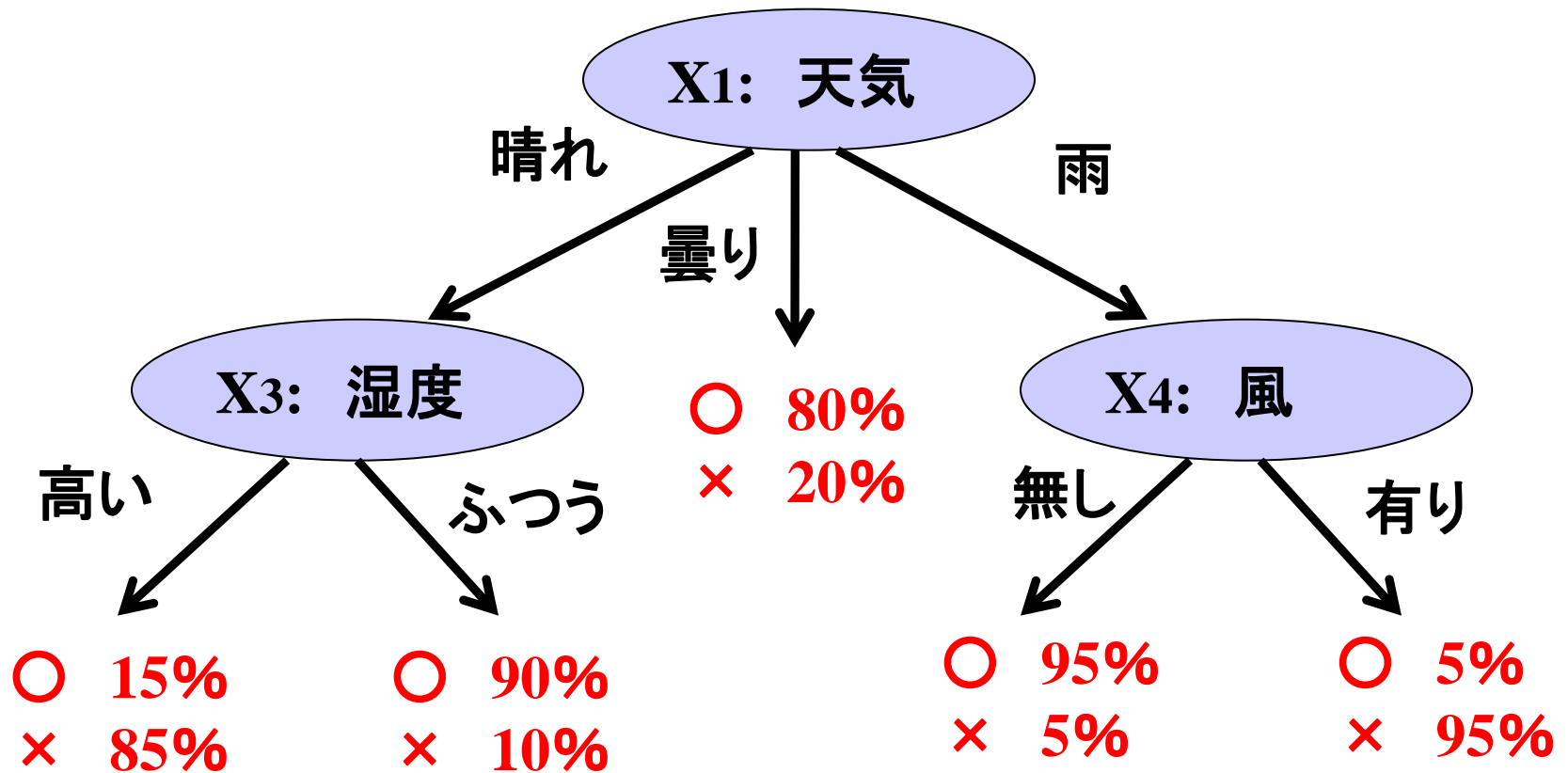
- しっかりとした基礎理論に立脚しつつ、現場に生きる機械学習・データマイニング研究をめざします
- JST CREST ビッグデータ基盤「複雑データからのディープナレッジの発見と価値化」の研究拠点です
- 幅広い産業界や公的研究機関との連携のもとで、研究を進めます
- 顧客視点に立って独善的ではないデータの利活用を心掛けています。

情報論的学習とは

X1 天気	X2 温度	X3 湿度	X4 風	Y プレイ
晴れ	暑い	高い	無し	×
晴れ	暑い	高い	有り	×
曇り	暑い	高い	無し	○
雨	暖かい	高い	無し	○
雨	涼しい	ふつう	無し	○
雨	涼しい	ふつう	有り	×
曇り	涼しい	ふつう	有り	○
晴れ	暖かい	高い	無し	×
晴れ	涼しい	ふつう	無し	○
雨	暖かい	ふつう	無し	○
晴れ	暖かい	ふつう	有り	○
曇り	暖かい	高い	有り	○
曇り	暑い	ふつう	無し	○
雨	暖かい	高い	有り	×
.....	

知識発見

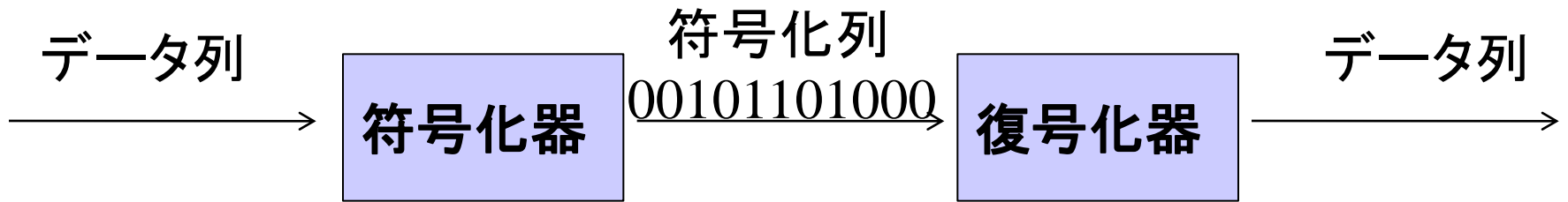
$X=(X_1, X_2, X_3, X_4)$ $Y=O$ or \times



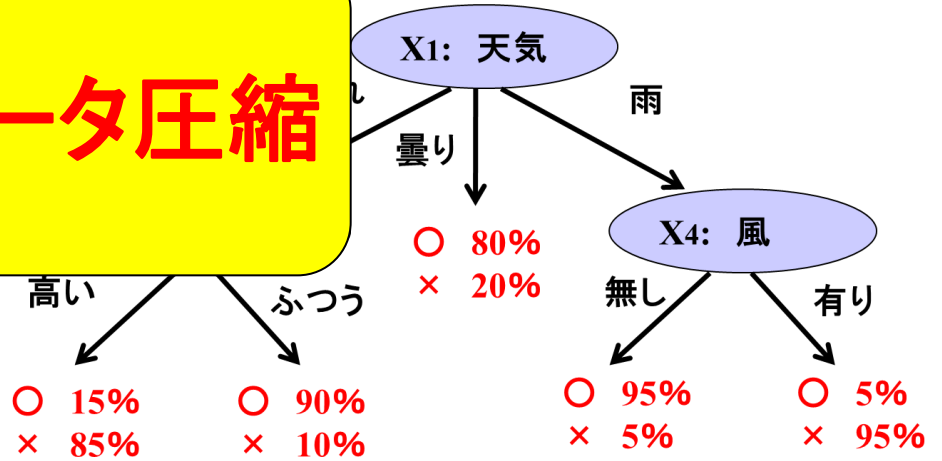
$P(Y|X)$ 知識

機械学習とは データから $P(Y|X)$ を導き出すこと

MDL原理と学習



X1	天気	X2	温度	X3	湿度	X4	風	Y	プレイ
	晴れ		暑い		高い		無し		×
	晴れ		暑い		高い		有り		×
	曇り		暑い		高い		無し		○
	雨		暖かい		高い		無し		○
	雨		涼しい		ふつう		無し		○
	雨		涼しい		ふつう		有り		×
	曇り		涼しい		ふつう		有り		○
	晴れ		暖かい		高い		無し		×
	晴れ		涼しい		ふつう		無し		○
	雨		暖かい		ふつう		無し		○
	晴れ		暖かい		ふつう		有り		○
	曇り		暖かい		高い		有り		○
	曇り		暑い		ふつう		無し		○
	雨		暖かい		高い		有り		×
.....



$$\mathcal{L}(D^n | M) + \mathcal{L}(M) \Rightarrow \text{最小化 w.r.t. } M$$

[データ符号長] [モデル符号長]

確率的コンプレキシティ

MDL (Minimum Description Length) 原理

情報論的学習理論のシナリオ

学習を通じて様々な学問領域が融合

古典的学習問題

推定問題

予測問題

識別問題

発展的学習問題

集団学習
転移学習

潜在的
ダイナミクス

関係データ
統合予測

データ圧縮

1つの視点

学習とは何か？

統計学

情報理論

統計物理

計算理論

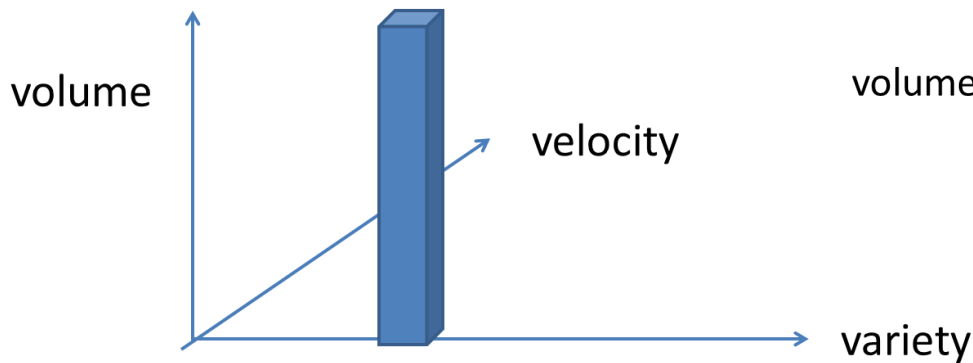
最適化

AI

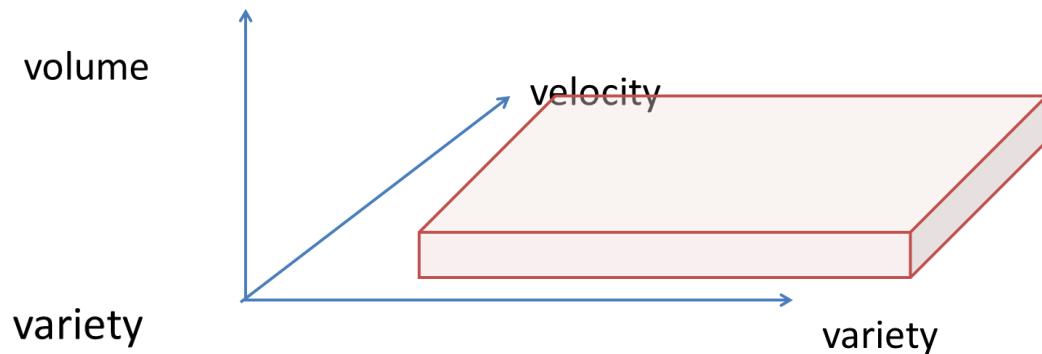
ビッグデータへの取り組み

5V

- **Volume** (大量性、大次元)
- **Variety** (多様性、複雑性(潜在性))
- **Velocity** (動的、非定常、変化)
- **Value** (付加価値)
- **Veracity** (真実性)



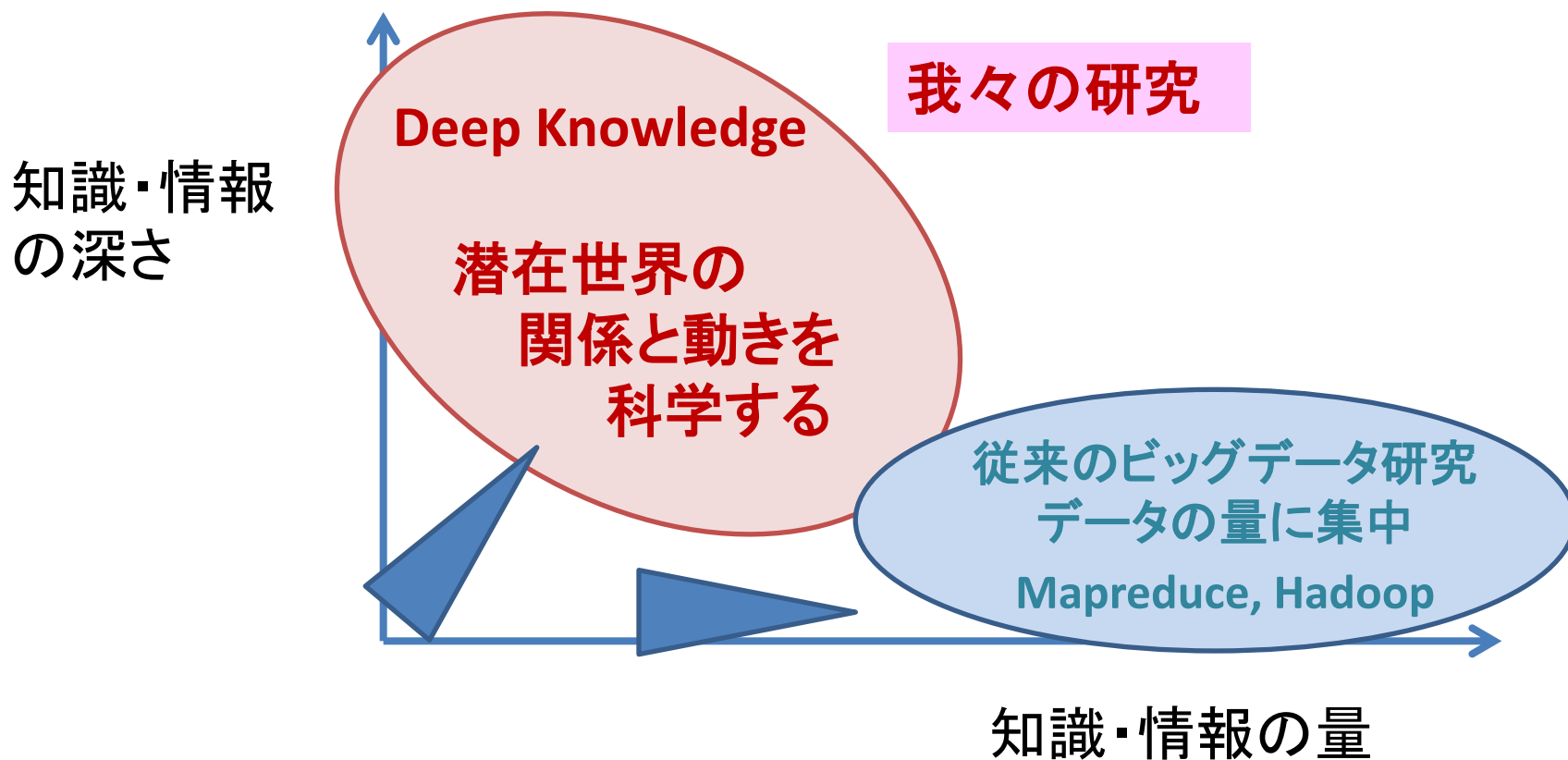
世間的イメージ



実際のビッグデータ

「複雑データからのディープナレッジの発見」 の研究拠点

知識の深さに重点を置き、データに潜在するDeep Knowledgeを発見し、valueを引き出す新しい方法論を開発し、体系化する



学習数理情報学研究室の キーワード

情報論的学習理論
IBIS(Information-Based Induction Sciences)



ディープ
ナレッジ

潜在的
ダイナミクス
Latent
Dynamics

関係データ
統合予測

学習数理情報学研究室のテーマ

■ 情報論的学習理論

- 潜在的ダイナミクス
- モデル選択、MDL原理
- 関係データ統合予測、潜在変数モデル
- 正則化理論、テンソル型学習
- マルチタスク学習、転移学習

■ データマイニング（ディープナレッジの発見）

- 異常検知・変化検知
- 緑内障進行予測
- Social Networkマイニング
- 教育データマイニング
- 交通・マーケティング向けデータマイニング

潜在的ダイナミクス(1/4)

クラスタリング構造変化検知

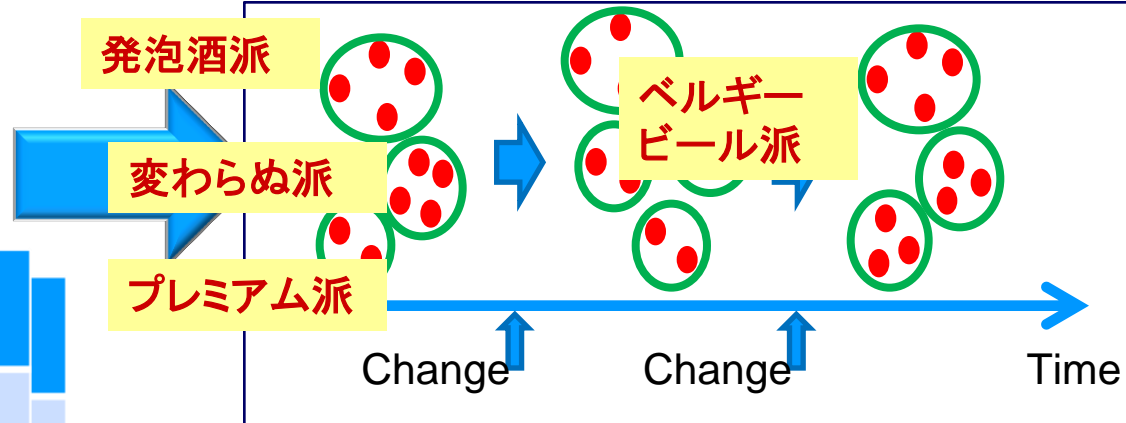
購買データからのユーザ層変化検知への応用

[Hirai Yamanishi KDD2012,
ISIT2011, IEEE IT2013]

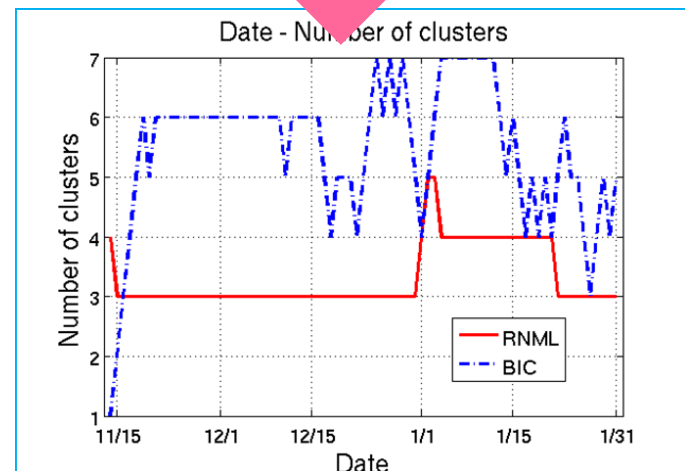
多変数時系列

	Beer 1	Beer 2	...
User 1	350	700	...
User 2	1050	350	...
...

購買層クラスター系列



データ圧縮に基づく
変化点解析

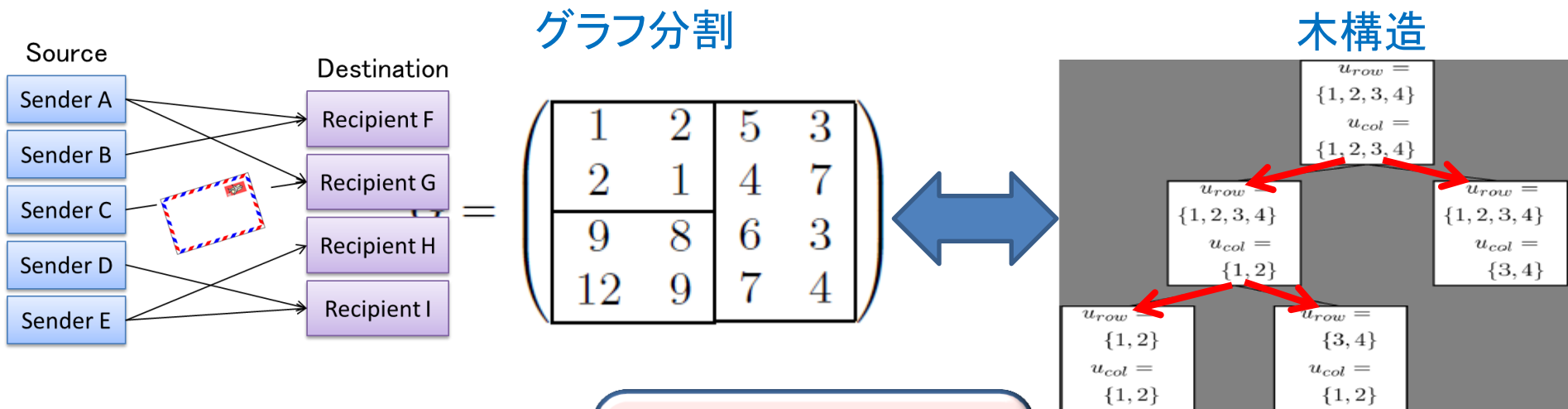


Inc.

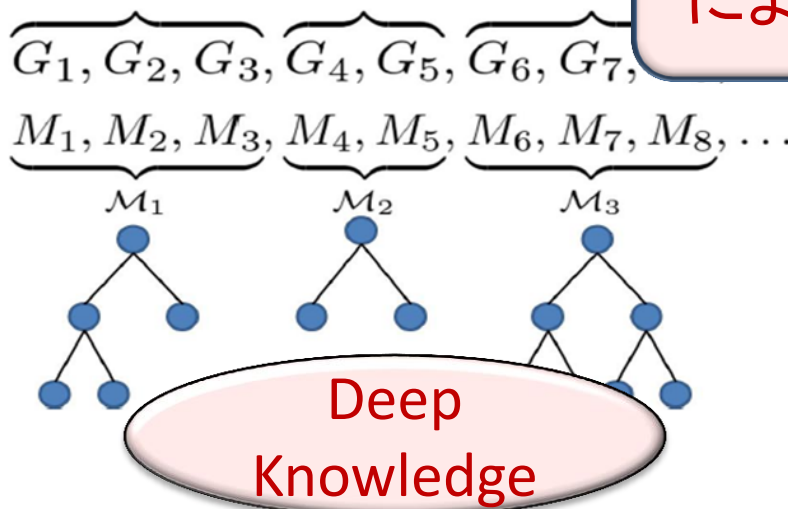
潜在的ダイナミクス(2/4)

木構造に基づくグラフ分割構造変化検知

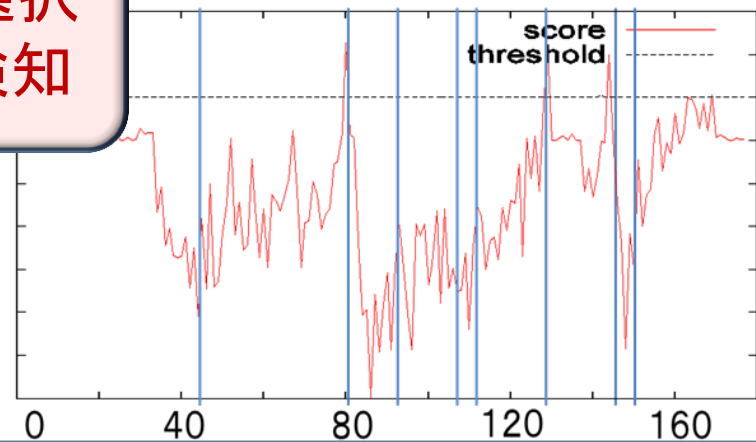
[Sato Yamanishi ICDM2013]



グラフ分割時系列



動的モデル選択
による変化検知



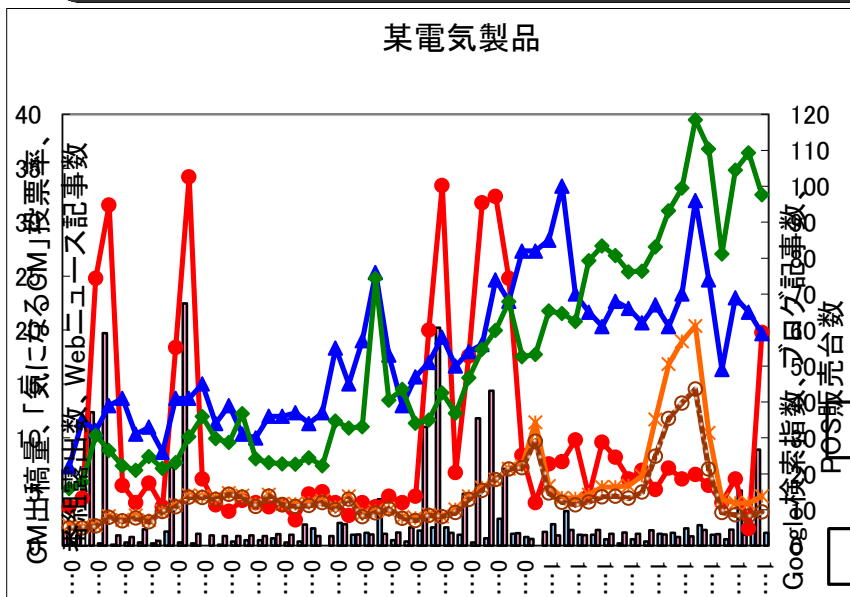
イベント予兆検知

潜在的ダイナミクス(3/4)

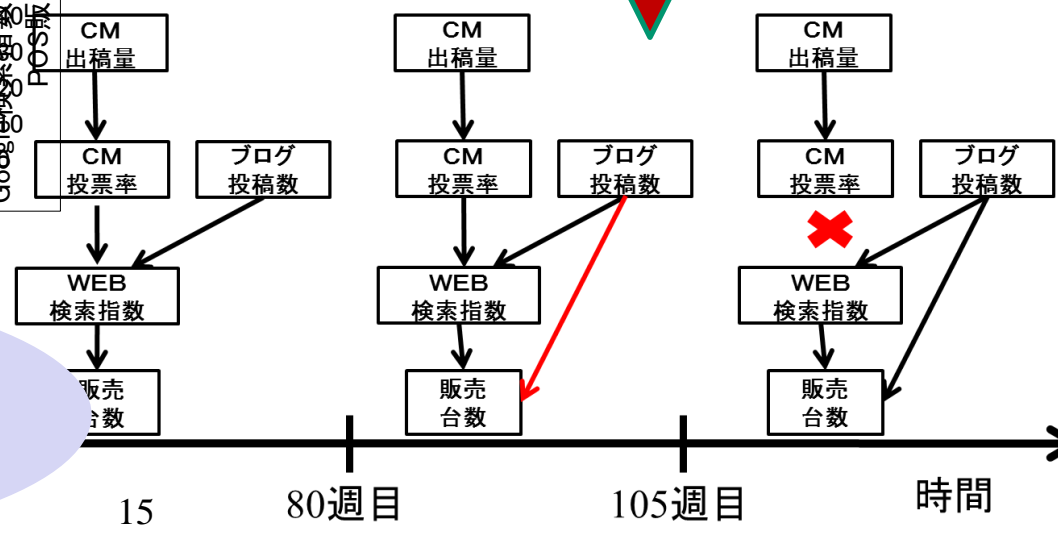
ベイジアンネット構造変化検知に基づく広告効果測定

[Hayashi Yamanishi ICDM2013, DAMI2013]

多次元時系列データからベイジアンネットワークの
構造変化検知により広告効果を測定



広告出稿



広告出稿により
ブログでの評判が高まり
売り上げに結び付いた

潜在的ダイナミクス(4/4)

教育データマイニング

[Oeda Yamanishi EDM2013]

試験時系列データからのスキル構造変化の検知



Many examination results



		Skills		
		s_1	s_2	s_3
Items	i_1	0	1	0
	i_2	0	0	1
	i_3	1	0	0
	i_4	0	0	1

A made Q-matrix automatically

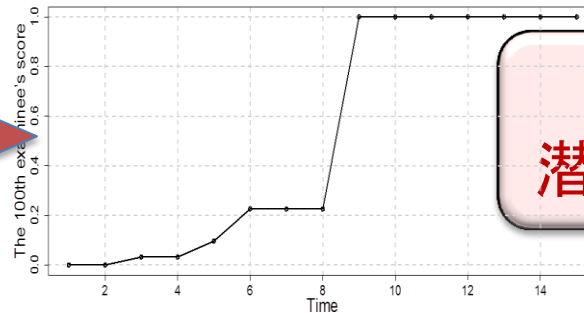
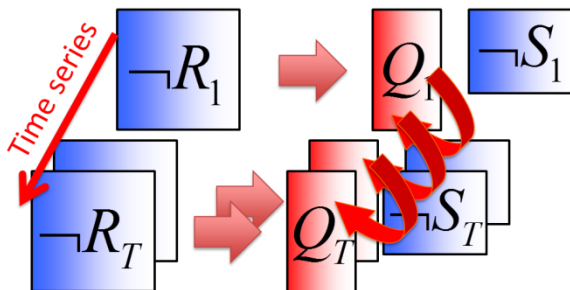


Analyzing new examination results by the Q-matrix

Non-negative matrix factorization(NMF)

$$\begin{matrix} & \text{Examinees} \\ \text{Items} & \begin{pmatrix} 0 & 0 & 1 & 0 & 1 \\ 0 & 1 & 0 & 1 & 1 \\ 1 & 1 & 1 & 0 & 1 \\ 0 & 1 & 0 & 1 & 1 \end{pmatrix}
 \end{matrix}
 =
 \begin{matrix} & \text{Skills} \\ \text{Items} & \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix}
 \end{matrix}
 \begin{matrix} & \text{Examinees} \\ \text{Skills} & \begin{pmatrix} 1 & 1 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 \\ 0 & 1 & 0 & 1 & 1 \end{pmatrix}
 \end{matrix}$$

• Online NMF:



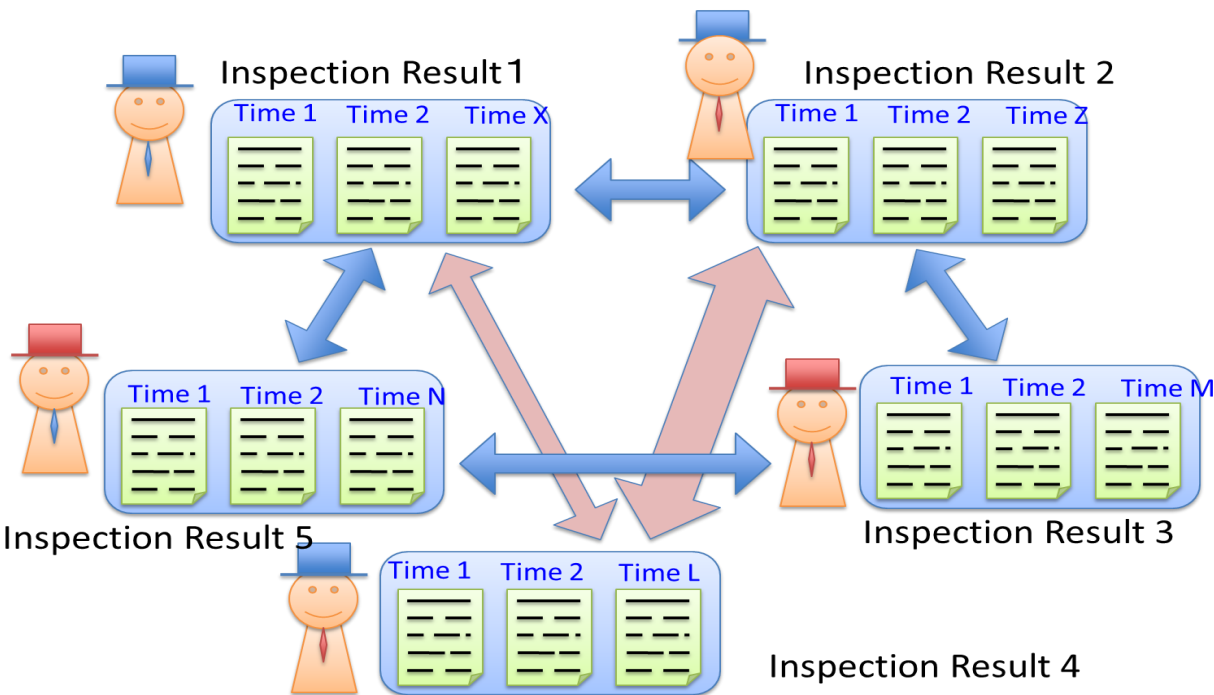
Deep Knowledge
潜在スキル構造の変化

関係データ統合予測(1/2)

緑内障進行予測

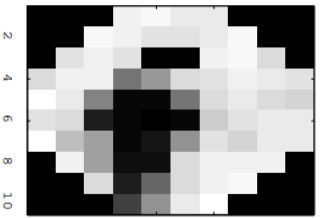
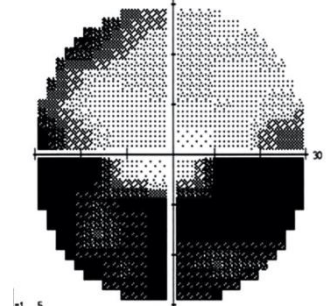
[Liang, Tomioka, Murata, Asaoka, Yamanishi ICDM013]

病状の時空間パタンの似ている患者の
情報を利用して進行予測精度を増強

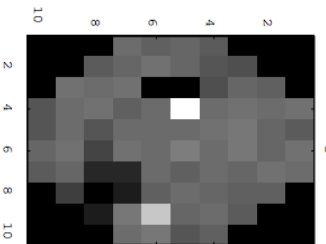


Deep Knowledge
= 患者ネットワーク
クラスタリング

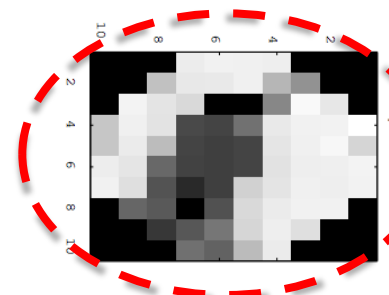
Correct



Prediction
without using
relational data



Prediction
aggregating
relational
data

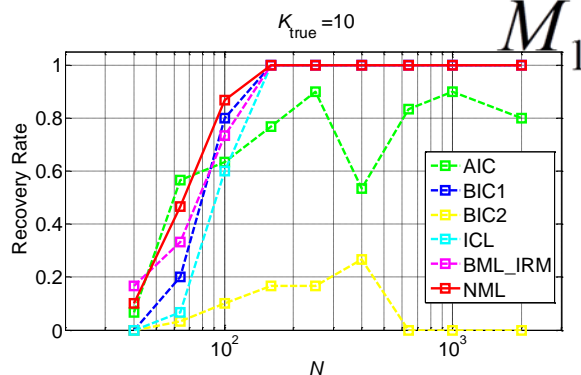
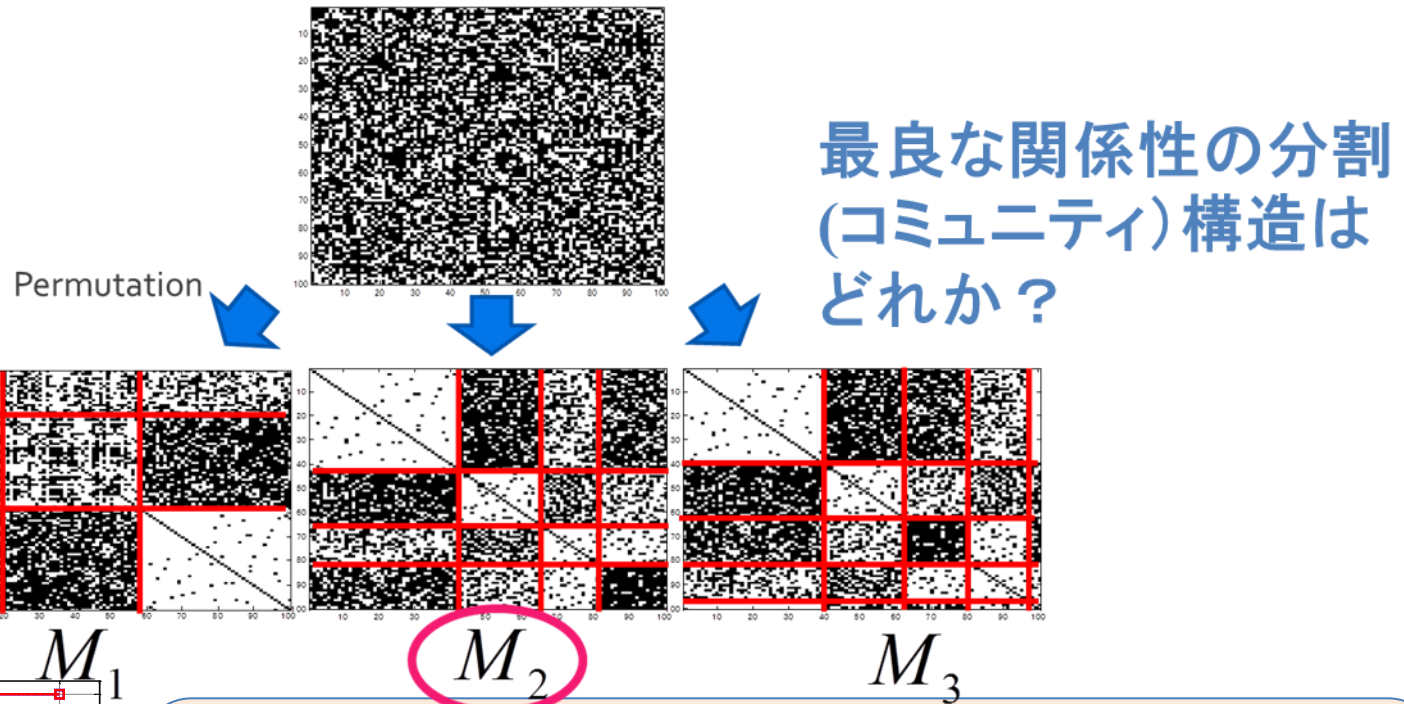


関係データ統合予測(2/2)

関係データモデル選択

[Sakai, Yamanishi BigData2013]

NML符号長によりヘテロな関係データの最適分割構造を推定



$\ln C_{\text{SBM}}(N, K)$ 一般関係モデルのモデル選択規準を開発

$$= \frac{K-1+K(K+1)}{2} \ln N - \frac{K-1+K(K+1)/2}{2} \ln 2\pi - \frac{K}{2} \ln 2$$

$$+ K \ln \Gamma((K+2)/2) - \ln \Gamma(K(K+2)/2) + \frac{K(K+1)}{2} \ln \pi$$

SNSマイニング

リンク情報に基づくTwitterからの話題出現検知

[Takahashi, Tomioka, Yamanishi TKDE2014]

従来の話題検出法

GW 原発
母の日

単語に注目

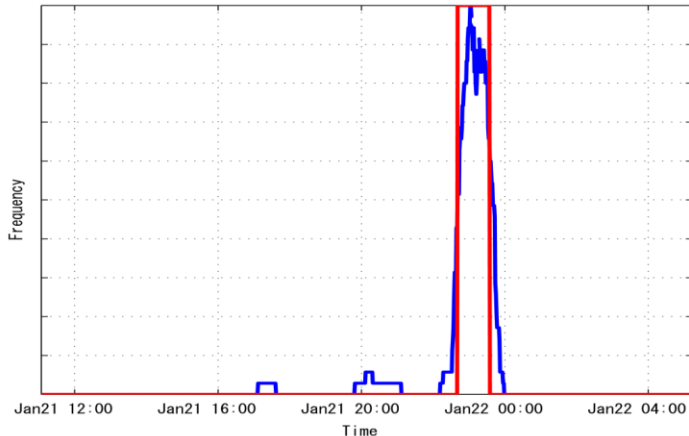
問題点: 言い換えに弱い

提案する話題検出法

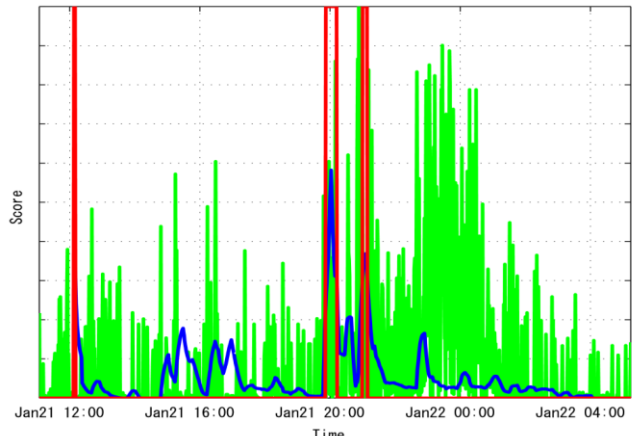
ユーザ間のリンクに注目

従来手法(キーワード頻度)より速い話題検出が可能

比較手法(キーワード頻度)



提案手法(変化点検出)



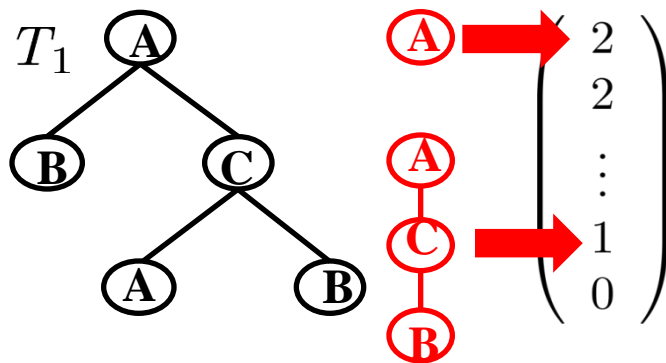
木構造データの高速カーネル関数の設計

木構造をもつカーネルの高速計算で複雑な対象を分類

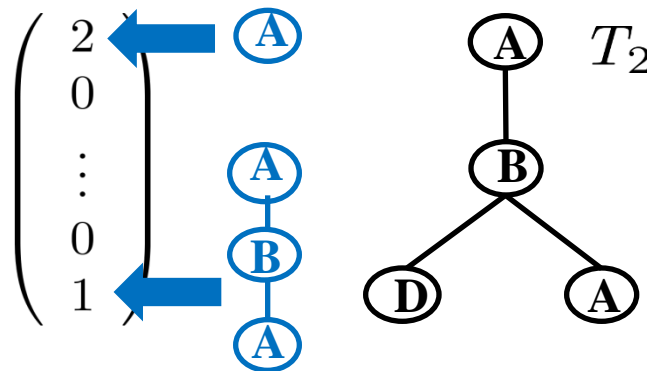
[Kimura Kashima ICML 2012]

$$k(T_1, T_2) = \sum_{p \in P} \underbrace{num(T_{1p})}_{\text{パス } p \text{ の出現回数}} \underbrace{num(T_{2p})}$$

パスに基づく木構造間の類似度を定義



内積
(類似度)



i	SA	LCP	S_i
1	1	0	A
2	3	0	BA
3	2	1	CA
4	4	-1	CBA

拡張接尾辞配列を用いることで2つの木に共通するパス数を高速に計算可能

人-機械・協調計算のプライバシー保護

情報理論・計算理論に基づき情報漏洩、品質保証精度を定量化

[Kajino, Arai, Kashima DAMI2014]



問題: 処理対象の情報漏洩

解: 必要以外の情報を排除



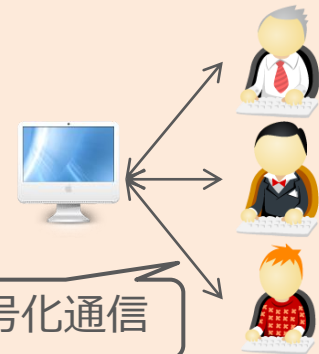
顔は塗りつぶせるが、誰かわからない

品質尺度: $E_I [\text{KL}(p(R|I) \parallel p'(R|I))]$
(R : 成果物, I : 対象物) 保護前 保護後

漏洩尺度: $E_I [\text{KL}(q(S|I) \parallel q(S))]$
(S : 秘匿情報, I : 対象物) 保護前 保護後

問題: 処理結果の情報漏洩

解: 統合化した統計量を返す
統合化には暗号化



単一マシンを用いた大規模最適化

学習最適化問題をDual Cashed Loopsを用いて高速処理

[Matsushima, Vishwanathan, Smola, KDD2012]

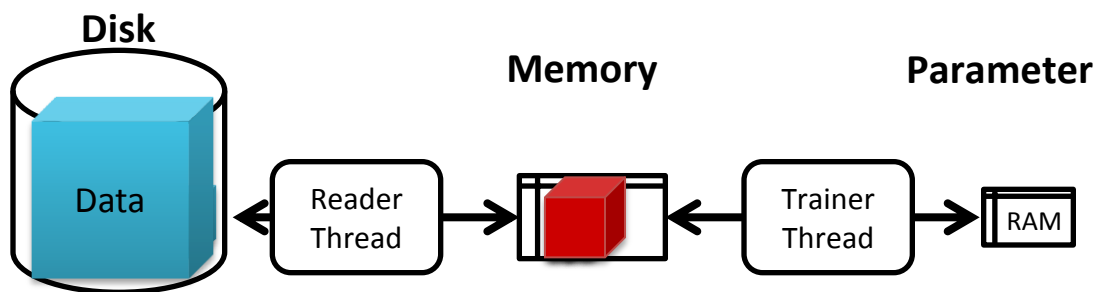
$$\min_{w \in \mathbf{R}^d} \underbrace{\sum_{j=1}^d \phi(w_j)}_{\text{regularizer}} + \underbrace{\sum_{i=1}^n \ell(w \cdot y_i x_i)}_{\text{loss}}$$

Dual Problem:

$$\min_{\alpha \in \mathbf{R}^n} \sum_{j=1}^d \phi^* \left(\sum_{i=1}^n \alpha_i y_i x_{ij} \right) + \sum_{i=1}^n \ell^*(-\alpha_i)$$

Application: 線形回帰、ロジスティック回帰、サポートベクターマシン (SVM) など

Solution: Dual Cashed Loops



- 読込/訓練スレッドを同時に非同期的に動作
- 訓練スレッドは高速にアクセス可能なメモリにのみアクセス
- “uninformative” なデータを優先的にメモリから削除する

近年の修論テーマ

- ・非定常情報源に対するパラメータ推定とモデル選択の研究
- ・正規化最尤符号に基づくクラスタリングの研究
- ・ネットワーク構造変化検知の研究
- ・木構造データに対する高速カーネル関数の設計
- ・Convex Formulation for Learning from Crowds
- ・木構造に基づくグラフ構造変化検知の研究
- ・区間的定常情報源の学習に関する研究
- ・独立成分分析に基づく変化検知の研究
- ・行列補完問題に対する転移学習アプローチ
- ・非一様データからの統合的学習
- ・グループ最適化に基づくネットワーククラスタリングの研究
- ・Dueling Bandit問題に対する漸近最適アルゴリズム