

第2回「ネットワーク生態系と空間デザイン」シンポジウム

2004年3月5日 於 ATR人間情報科学研究所

# Webサイトの競合ダイナミクス

木村 昌弘 齊藤 和巳 上田 修功

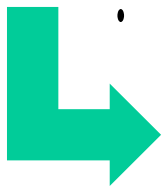
NTT コミュニケーション科学基礎研究所

# World Wide Web

- ・ ハイパーリンクで繋がれたサイト群の巨大なネットワーク
- ・ 巨大な情報空間であり、コミュニケーションの新重要メディア  
→ ネットワーク生態系の最重要例の一つ

社会学や経済学の観点では、

- ・ サイト管理者がユーザに情報商品を提供する  
グローバルなマーケット
- ・ あるサイトへの訪問者数はそのサイトの成功指標



「類似したサービスを提供するWebサイト群は、各自のサイトへの訪問者数を増やすために、互いに競合している。」

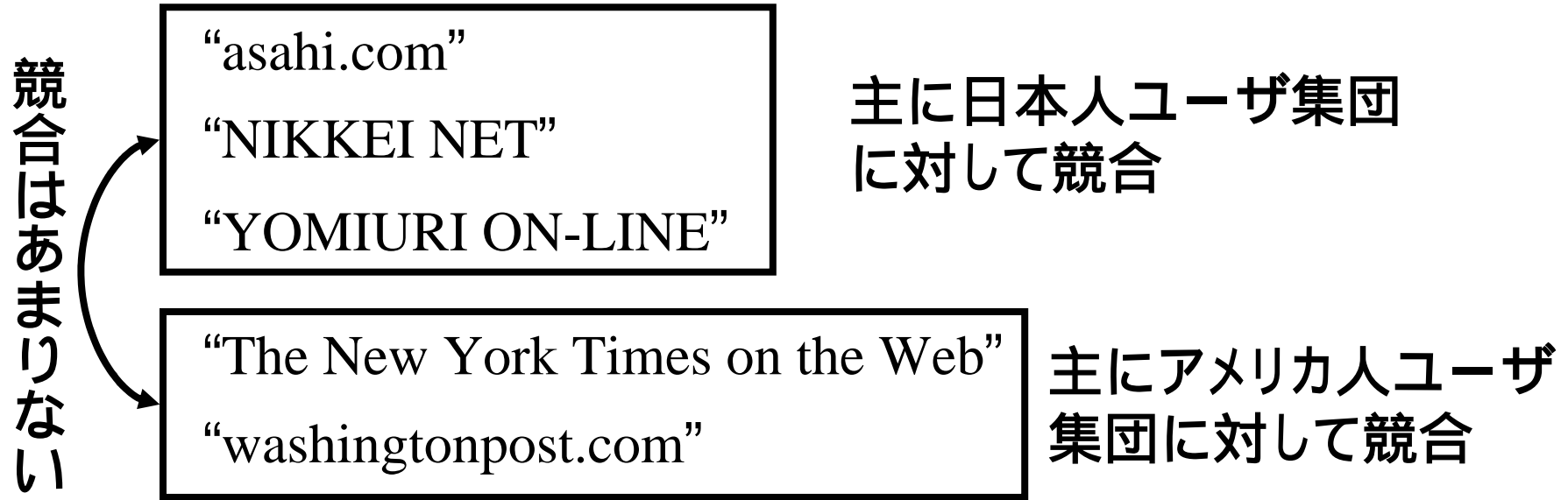
重要

マーケットを形成するWebサイト群に対して、

各サイトの**訪問者数の変動** → モデリング & 解析  
**競合ダイナミックスの観点**

# Webマーケットの競合構造

最新ニュースを提供するWebサイト集合の一例



マーケットを形成するサイト群 → **競合グループ**に分類される

**重要な課題**

サイト集合の競合構造を, サイト訪問者数の変動データから抽出

(サイト管理者のビジネス戦略立案への貢献)

**必要性は高い**

サイト競合ダイナミックスのモデリング ← **競合構造を組み込む**

# 関連研究

Webサーファのアクセスパターンに関する統計的規則性の発見  
例) サイト毎の訪問者数分布がべき分布

↓  
定性的説明

Webユーザの振る舞いに関して:

- ・統計的理論 Huberman *et al* (1998), Adamic & Huberman (2000)
- ・エージェントモデル Liu *et al* (to appear)

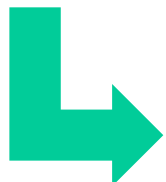


類似したサービスを提供するWebサイト群に関して:

サイト訪問者数ダイナミックスのモデル Maurer & Huberman (2003)

↓  
定性的説明

サイトの競合関係がそのマーケットの性質に及ぼす影響



観測データの予測モデル(定量的モデル)ではない

# 関連研究(続き)

あるWebマーケットに対して、  
各サイトの訪問者数：

観測データ → 近未来予測

考えられる

ブラックボックスモデルの適用

(ARモデル, ニューラルネットなど)

現象の根本的な構造やメカニズムを陽に表現しない

例えば：

- ・ マーケットの競合構造
- ・ サイトの競合ダイナミクス

# 研究の目的

マーケットを形成するサイト集合に対して、

各サイトへの訪問者数の短期変動現象のモデリング

観測時系列データ 予測モデル(定量的モデル) :

- ・ サイト間の相互作用の構造とメカニズムを、  
サイト集合の競合構造と競合ダイナミックスの見地から表現
- ・ 観測時系列データを定量的に説明し、さらに、  
各サイトの近未来のシェアを予測



- ・ 新たな確率的ダイナミクスモデルの提案
- ・ その学習アルゴリズムの導出

# サイト訪問者数変動のモデリング

$S = \{s_1, \dots, s_N\}$  : マーケットを形成するサイト集合

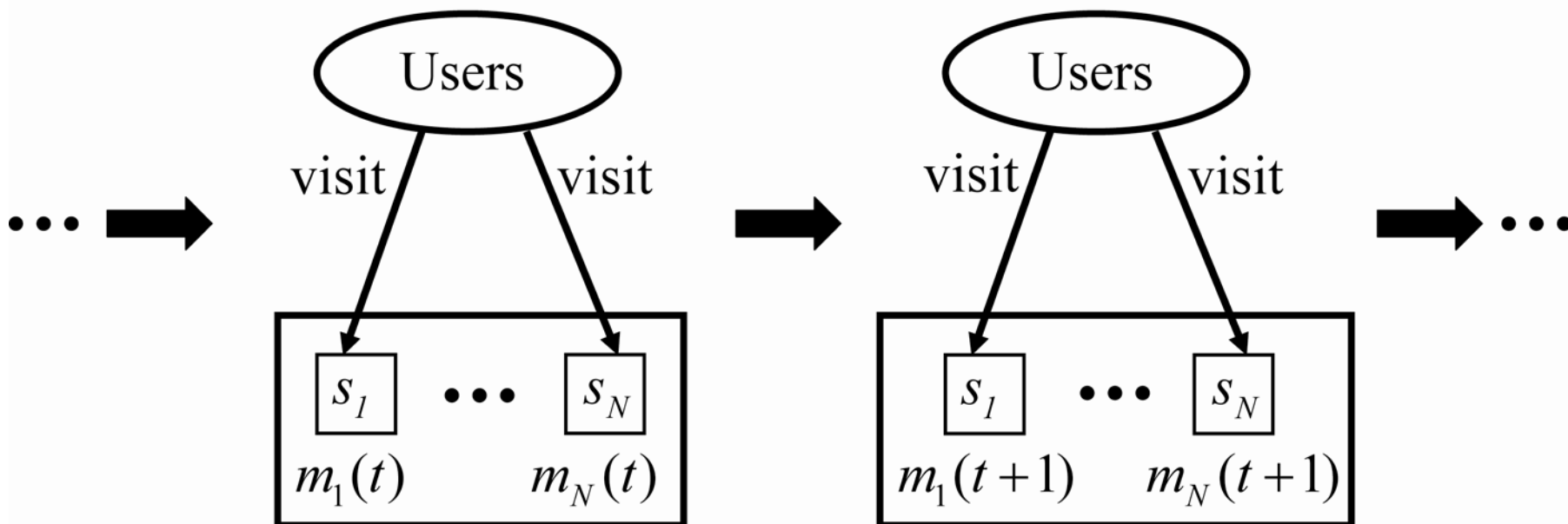
$m_i(t)$  : サイト  $s_i$  への時間ステップ  $t$  における訪問者数

$\mathbf{m}(t) = (m_1(t), \dots, m_N(t))$  : 時間ステップ  $t$  における  
マーケット  $S$  の訪問ベクトル

 **ダイナミクス  
モデリング**

Time-step  $t$

Time-step  $t+1$



# シェアダイナミクス

$M(t) = \sum_{i=1}^N m_i(t)$  :  $S$  への時間ステップ  $t$  における総訪問者数

ダイナミクス

  $S$  を含むもっと大きい社会システムを調べて、  
あらかじめ別にモデリングすべき。

  $M(t)$  は与えられると**仮定**し、  
各サイト  $s_i$  の**シェア**を予測する問題を考える。

$x_i(t) = m_i(t)/M(t)$  :  $s_i$  の時間ステップ  $t$  における**シェア**

$S$  の時間ステップ  $t$  におけるシェアベクトル

$$\mathbf{x}(t) = (x_1(t), \dots, x_N(t)) \in \Delta^{N-1}$$

 **単体上**の確率ダイナミクスのモデリングが必要



# S の競合構造

仮定:

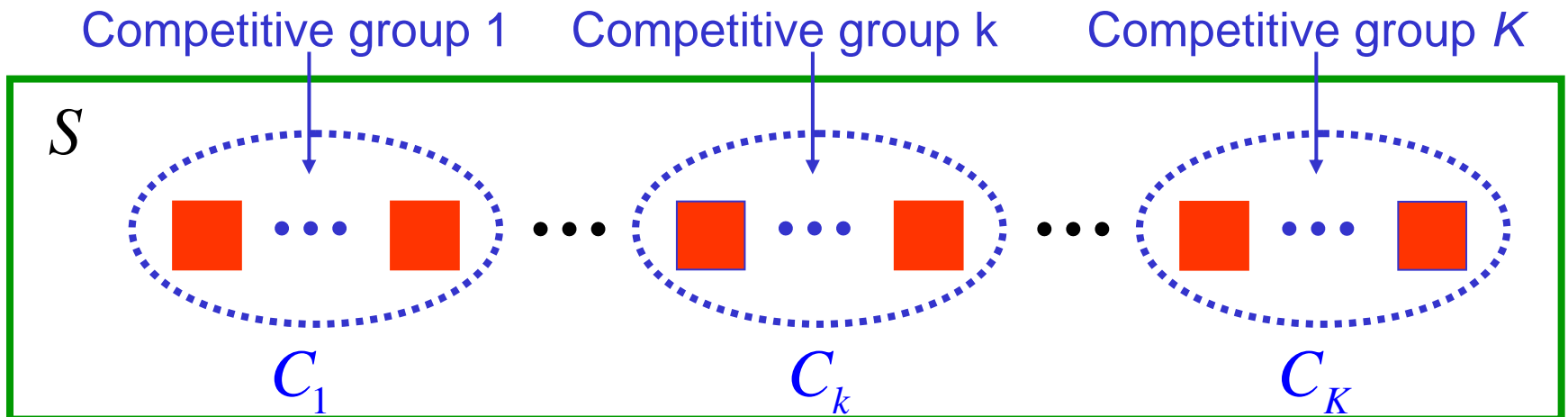
S は K 個の競合グループに分割される

$$S = \bigcup_{k=1}^K C_k, \quad (\text{disjoint union})$$



$C_k = \{s_i; \lambda_i = k\}$ : 競合グループ(ジャンル)  
( $k = 1, \dots, K$ )

$$\lambda = (\lambda_1, \dots, \lambda_N); \quad \lambda_i \in \{1, \dots, K\}, i = 1, \dots, N$$



# 提案モデル

## 確率過程

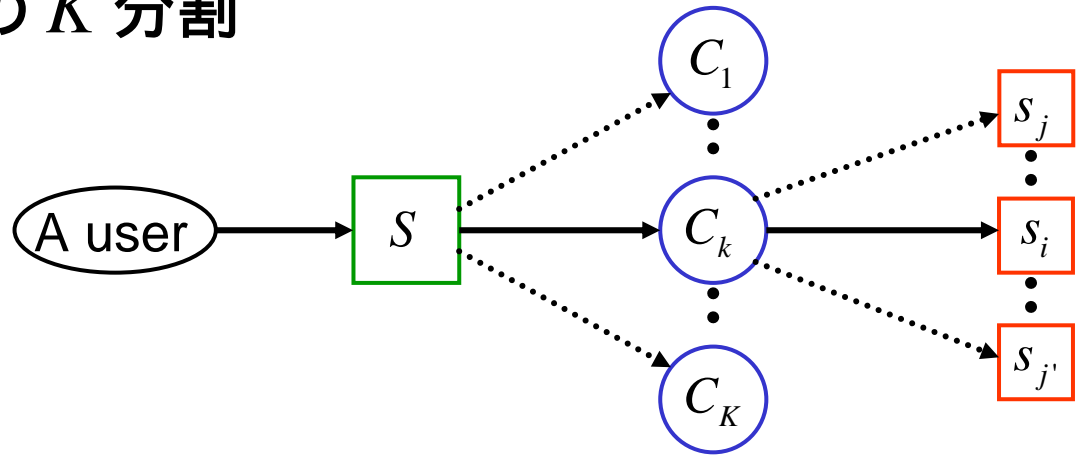
$$P(\mathbf{m}(t+1) | \mathbf{m}(t), M(t+1), \Theta, \lambda) \propto \prod_{i=1}^N \{P(s_i | \mathbf{m}(t), \Theta, \lambda)\}^{m_i(t+1)}$$

(多項分布)

$\Theta = ((\theta_1, \dots, \theta_K), \Phi_1, \dots, \Phi_K)$ : パラメータベクトル

$\lambda = (\lambda_1, \dots, \lambda_N)$ :  $S$  の  $K$  分割

$S$  の一般ユーザが  
 $s_i \in C_k$  を訪問する行動



$$P(s_i | \mathbf{m}(t), \Theta, \lambda) = P(s_i | C_k, \mathbf{m}(t), \Theta, \lambda) P(C_k | \mathbf{m}(t), \Theta, \lambda)$$

ただし,

$$P(C_k | \mathbf{m}(t), \Theta, \lambda) \triangleq \theta_k \quad \text{時間に依らない定数}$$

# $P(s_i | C_k, \mathbf{m}(t), \Theta, \lambda)$ の定義

$$P(s_i | C_k, \mathbf{m}(t), \Theta, \lambda) \triangleq f_{k,i}(\mathbf{x}_k(t); \Phi_k) \text{ シェアダイナミクス}$$

||

$C_k$  内でのシェアベクトル

$$\overline{x_{k,i}}(t+1) \quad x_{k,i}(t) : s_i \text{ の } C_k \text{ 内でのシェア}$$

## 定義

$$f_{k,i}(\mathbf{x}_k(t); \Phi_k) = (1 - \xi_k) \underbrace{g_k(\mathbf{x}_k(t); \alpha_k, \beta_k, \eta_k)}_{\text{Replicator dynamics}} + \xi_k \frac{1}{\underbrace{N_k}_{\text{Uniform dynamics}}} \quad (N_k = |C_k|)$$

Replicator dynamics

Uniform dynamics

$$\frac{\overline{x_{k,i}}(t+1) - x_{k,i}(t)}{x_{k,i}(t)} = \underbrace{A_{k,i}(\mathbf{x}_k(t))}_{\text{魅力度}} - \underbrace{\overline{A_k}(\mathbf{x}_k(t))}_{\text{平均魅力度}}$$

$s_i$  の  $C_k$  内での**魅力度**

$C_k$  の**平均魅力度**

# サイト魅力度の定義

$s_i$  の  $C_k$  内での魅力度

$$A_{k,i}(\mathbf{x}_k(t)) = (1 - \eta_k) \alpha_{k,i} x_{k,i}(t) + \eta_k \beta_{k,i} (1 - x_{k,i}(t))$$

一般受けする魅力度  
(一般価値的魅力度)

専門家受けする魅力度  
(希少価値的魅力度)

パラメータベクトル  $\Theta$ :

魅力度:  $0 < \alpha_{k,i}, \beta_{k,i} < 1; \sum_i \alpha_{k,i} = \sum_i \beta_{k,i} = 1$

混合度:  $0 < \eta_k < 1, 0 < \xi_k < 1$

$$\rightarrow \Phi_k = (\alpha_k, \beta_k, \eta_k, \xi_k)$$

選択率:  $0 < \theta_k < 1; \sum_{k=1}^K \theta_k = 1$

$$\rightarrow \Theta = ((\theta_1, \dots, \theta_K), \Phi_1, \dots, \Phi_K)$$

# 学習アルゴリズム

  $\{m(0), \dots, m(T)\}$ :  $S$  の訪問ベクトルの観測時系列データ

最適な確率ダイナミクスモデル

$$P\left(m(t) \mid m(t-1), M(t), \hat{\Theta}, \hat{\lambda}, \hat{K}\right), \quad (t \geq 1)$$

を推定する。

パラメータベクトル

分割ベクトル

分割数

戦略:

- ・ 最尤推定  $\rightarrow \hat{\Theta}, \hat{\lambda}$
- ・ サイトの1期先シェアの平均予測誤差の最小化  $\rightarrow \hat{K}$

# $\Theta$ の推定

$\{m(0), \dots, m(T)\}$ ,  $K$ ,  $\lambda$  が与えられたとき,  
最適パラメータベクトル  $\hat{\Theta}$  は

$$L(\Theta) = -\log P(m(0), \dots, m(T)) - \log P(\Theta)$$


を最小化することにより計算する.

$\Theta \in \Delta^{K-1} \times \Delta^{2N_1} \times \dots \times \Delta^{2N_K}$  であるので,

$$P(\Theta) \propto \prod_j \Theta_j, \quad (\text{Dirichlet 分布の独立積})$$

と仮定  $\rightarrow$  *Laplace smoothing*

$$\therefore L(\Theta) = -\sum_i a_i \log \left( \sum_j b_{ij} \Theta_j \right), \quad a_i, b_{ij} > 0.$$

- 
- ・ 極小解 = 大域最適解
  - ・ EMアルゴリズムで効率的に計算できる

# $\lambda$ の推定

$\{m(0), \dots, m(T)\}$ ,  $K$  が与えられたとき,

最適分割ベクトル  $\hat{\lambda}$  は

$$\lambda \mapsto L(\Theta(\lambda)),$$

を最小化することにより計算する。ただし,

$\Theta(\lambda)$  : 分割  $\lambda$  に対する  $\Theta$  の最適値

戦略:

$K$  分割  $\lambda = (\lambda_1, \dots, \lambda_N)$ ,  $(\lambda_i \in \{1, \dots, K\}, i = 1, \dots, N)$

を局所的に最適な方向へ変更



$\hat{\lambda}$  を探索

# 予測性能の評価尺度

$\{m(0), \dots, m(T)\}$  : 訪問ベクトルの観測時系列データ

$T_0 > 0$  を固定し,

実験  $\rightarrow (T = 50, T_0 = 10)$

$1 \leq \forall \Delta T \leq T_0$  に対して,

$\{m(0), \dots, m(T - \Delta T)\}$  を学習データとして,

時間ステップ  $T - \Delta T + 1$  でのシェアベクトル

$$x(T - \Delta T + 1) = \frac{m(T - \Delta T + 1)}{M(T - \Delta T + 1)} \quad \text{を予測}$$

予測誤差

$$\mathcal{E}(\Delta T) = \frac{1}{2} \sum_{i=1}^N \left| \frac{m_i(T - \Delta T + 1)}{M(T - \Delta T + 1)} - \hat{x}_i(T - \Delta T + 1) \right|.$$

観測値

予測値



平均予測誤差

$$\bar{\mathcal{E}} = \frac{1}{T_0} \sum_{\Delta T=1}^{T_0} \mathcal{E}(\Delta T)$$

で性能評価



# 予測タスクにおける比較対象

- ・ ガウスノイズに基づいた確率過程モデル  

(∴)

必ずしも  $\Delta^{N-1} \rightarrow \Delta^{N-1}$  とは写像しない

→ あるサイトのシェアを負にし得る



- ARモデルやニューラルネットの単純な適用  

- ・ 単純多項分布法 (NMM)

- ・ 平行移動法 (PDM)

比較すべき単純な従来法  
として採用

# 予測タスクにおける比較対象 (続き)

## ・ 単純多項分布法 (NMM)

仮定: 「訪問ベクトル  $m(t)$  は, ある多項分布にしたがって  
時間  $t$  とは独立に生成される .」

学習データ  $\{m(0), \dots, m(T - \Delta T)\}$  から

その多項分布を最尤推定し,

$\omega = (\omega_1, \dots, \omega_N)$  : 推定された多項分布のパラメータ

時間ステップ  $T - \Delta T + 1$  でのシェアベクトルの予測値を

$$\hat{x}(T - \Delta T + 1) = \omega$$

とする .

# 予測タスクにおける比較対象 (続き)

- ・ 平行移動法 (PDM)

学習データ  $\{m(0), \dots, m(T - \Delta T)\}$  から

時間ステップ  $T - \Delta T + 1$  でのシェアベクトルを予測するのに

時間ステップ  $T - \Delta T$  でのシェアベクトルの観測値で予測する。

すなわち,

$$\hat{\mathbf{x}}(T - \Delta T + 1) = \mathbf{x}(T - \Delta T)$$

# 分類性能の評価尺度

$\lambda = \{C_k\}_{k=1}^K$  :  $S$  の真の  $K$  分割

$\hat{\lambda} = \{\hat{C}_l\}_{l=1}^L$  :  $\{m(0), \dots, m(T - \Delta T)\}$  から推定された  $S$  の  $K$  分割

各  $\hat{C}_l$  に対して,

$C_{k[l]}$  : 期間  $[0, T - \Delta T]$  にサイト集合  $C_k \cap \hat{C}_l$  を訪問した  
ユーザ数が最も多かった  $C_k$

$\hat{\lambda}$  の  $\lambda$  に対する micro-averaged precision  $AP(\Delta T)$

$$\frac{\text{期間 } [0, T - \Delta T] \text{ に } \hat{C}_l \cap C_{k[l]} \text{ を訪問したユーザ数}}{\text{期間 } [0, T - \Delta T] \text{ に訪問した総ユーザ数}} \times 100$$

# 分類タスクにおける比較対象

Webサイト集合を関連するサイト集合に分類



多くの研究

→ しばしば要求される

テキストやリンクの静的情報を使って、  
トピック的に関連するサイト群やコミュニティを抽出する研究



適当ではない

ユーザに関する動的情報に基づいた競合サイトのグループ化

一方、経済物理学の分野で、Mantegnaは:

「株価の時系列データから同業種の会社群を抽出することに成功」

同一視

サイトの訪問者数変動      会社の株価変動

比較対象として、  
Mantegnaの手法  
を考える



# 分類タスクにおける比較対象 (続き)

## 比較対象とする従来法

1. Mantegnaの手法を用いて,  
学習データ  $\{m(0), \dots, m(T - \Delta T)\}$  に基づいて,  
各サイトを球面上の点と同一視する.

$$s_i \leftrightarrow w_i = (w_i(1), \dots, w_i(T - \Delta T)), \quad (i = 1, \dots, N),$$

ここに,

$$w_i(t) = (y_i(t) - \mu_i) / \sigma_i,$$

$$y_i(t) = \log m_i(t) - \log m_i(t - 1).$$

2. Spherical  $K$ -means法を用いて,

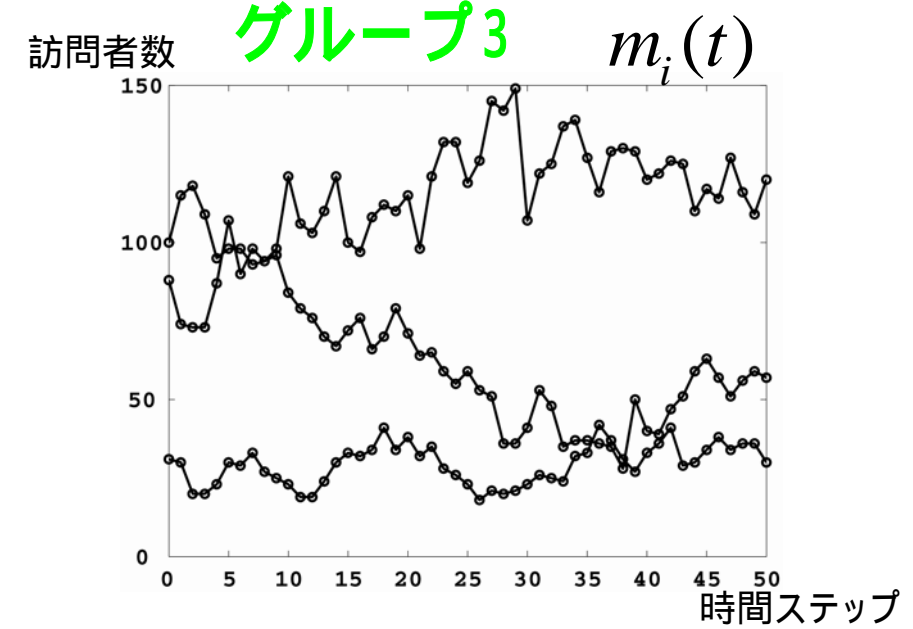
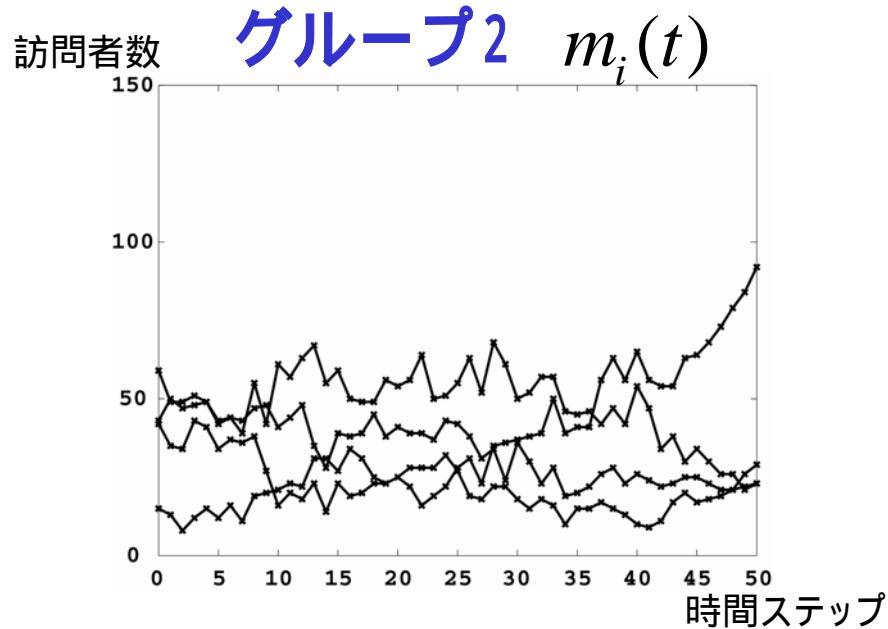
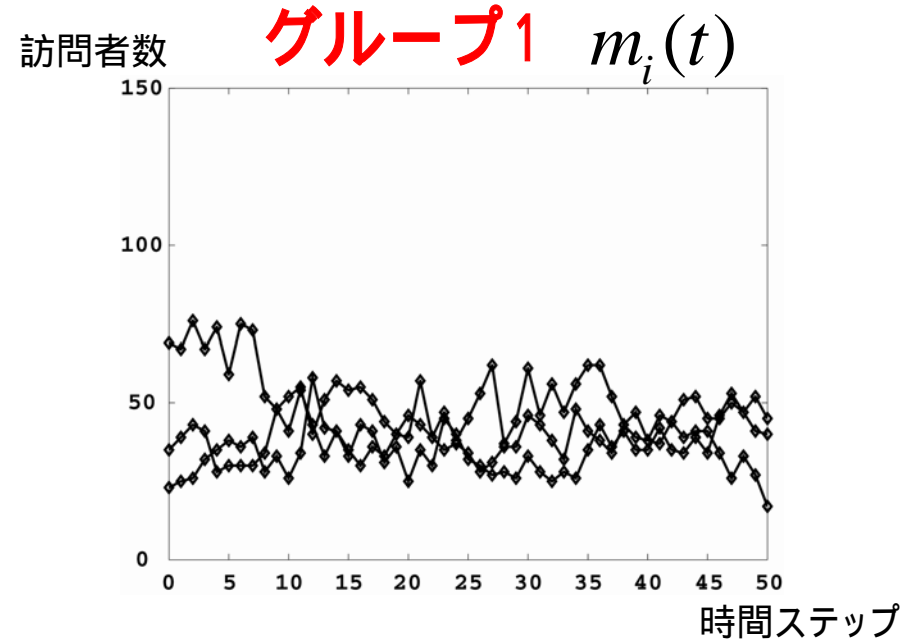
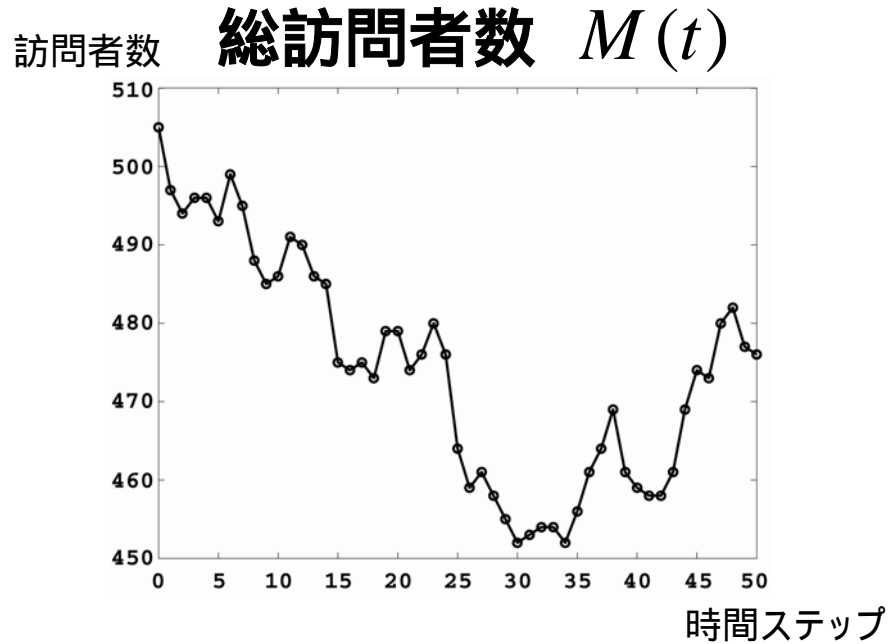
$w_1, \dots, w_N$  を  $K$  分割する.

$\Leftrightarrow s_1, \dots, s_N$  を  $K$  分割する.

注)

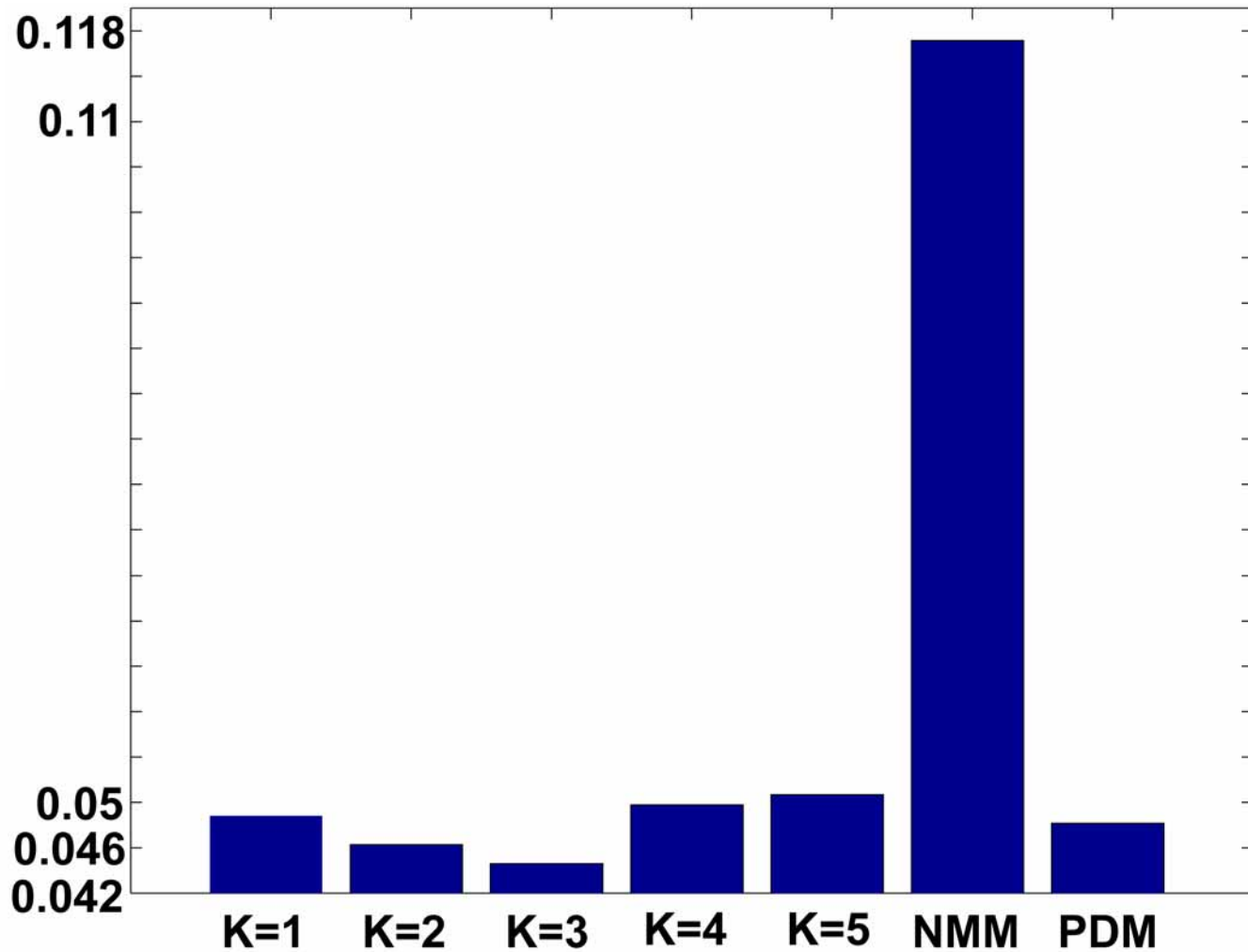
Mantegnaのオリジナル:  
Euclid距離に基づいて  
デンドログラムを構築

# 人工データによる実験評価(10サイト)



# 予測性能の比較(人工データ)

$\bar{\epsilon}$   
平均予測誤差

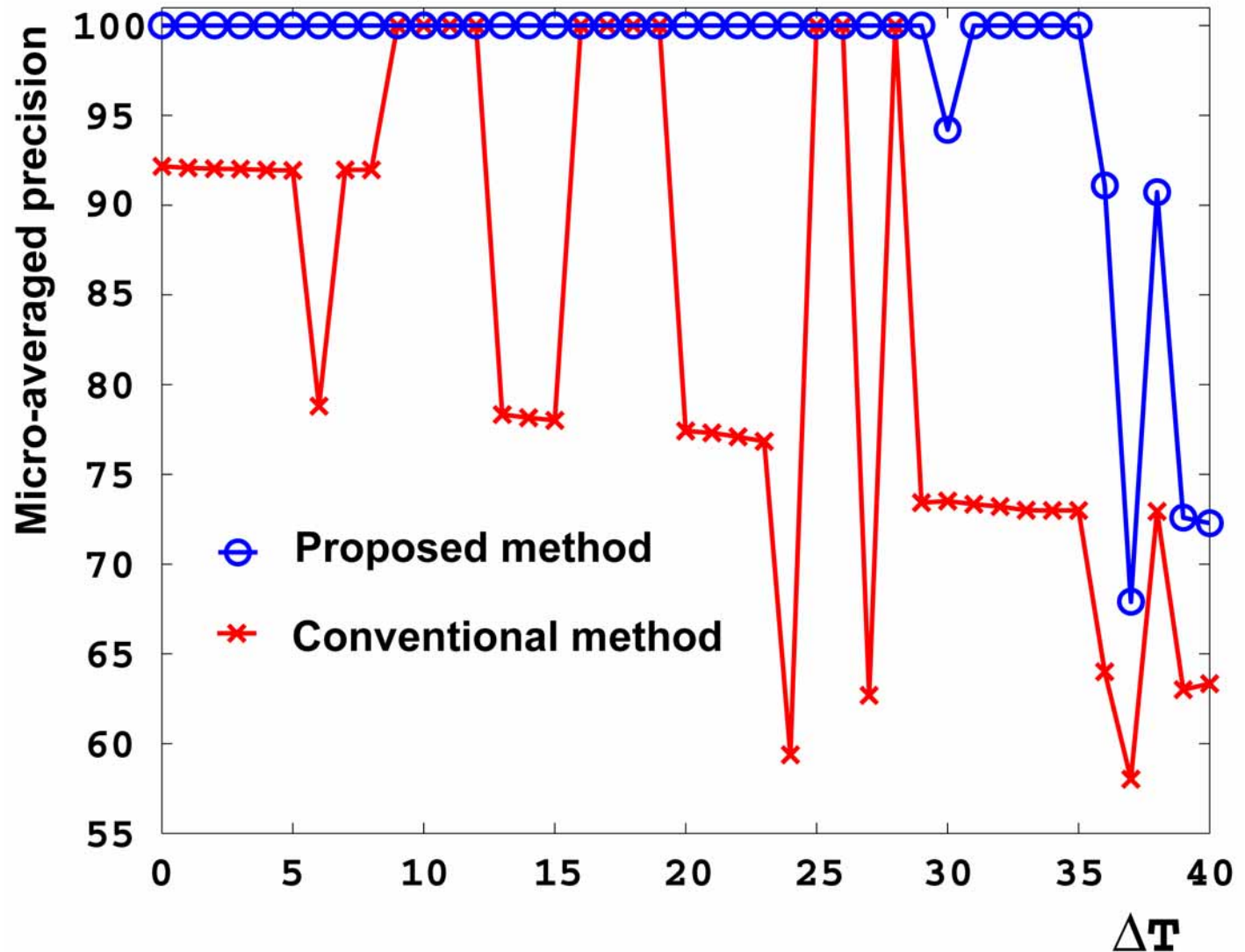




# 分類性能の比較(人工データ)

$K = 3$

$AP(\Delta T)$



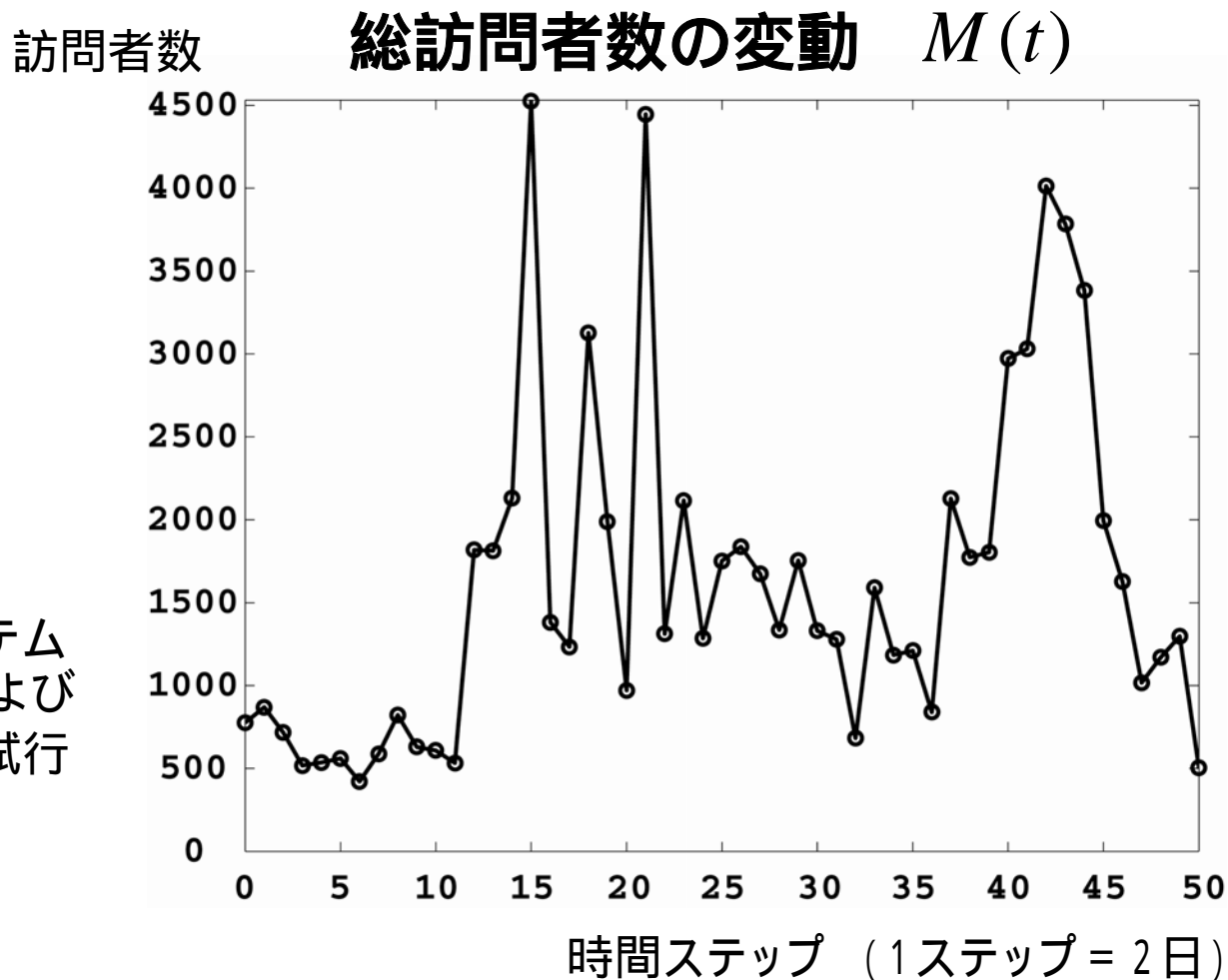
注)  $\{m(0), \dots, m(T - \Delta T)\} \rightarrow AP(\Delta T)$

# 実Webデータによる実験評価

ストリーミングビデオコンテンツを提供する，日本の20 Webサイトの利用ログのデータ

NTTが東京ニュース通信社と共同で2002年に実施した「ブロードバンド番組ガイド」の試行サービス実験におけるデータ

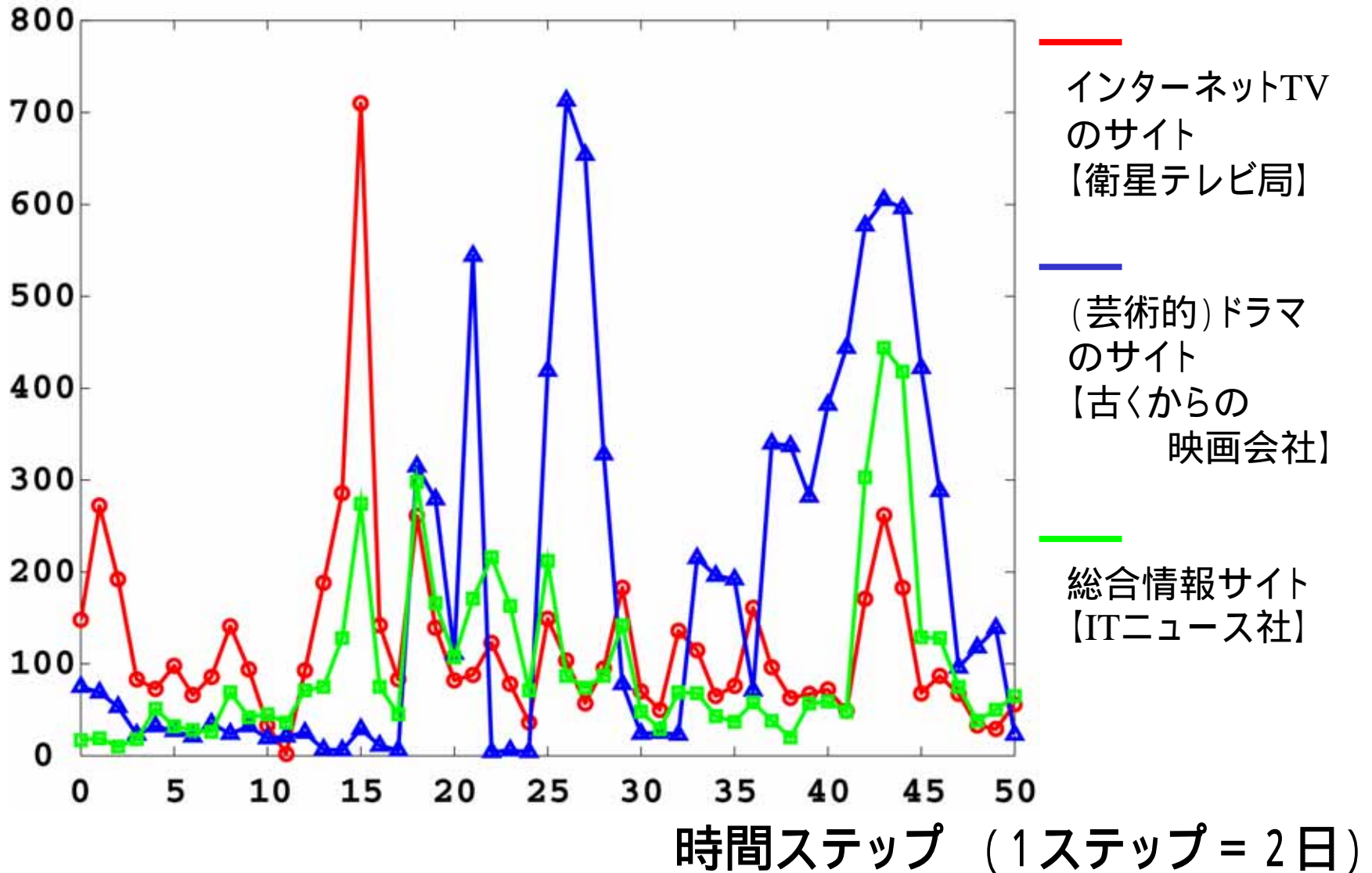
安部，宮原，林，外村，  
散策型コンテンツガイドシステム  
「AssociaGuide」：システムおよび  
「ブロードバンド番組ガイド」試行  
サービス実験概要，  
映像メディア学会技術報告，  
26(81)，1 - 4，2002.



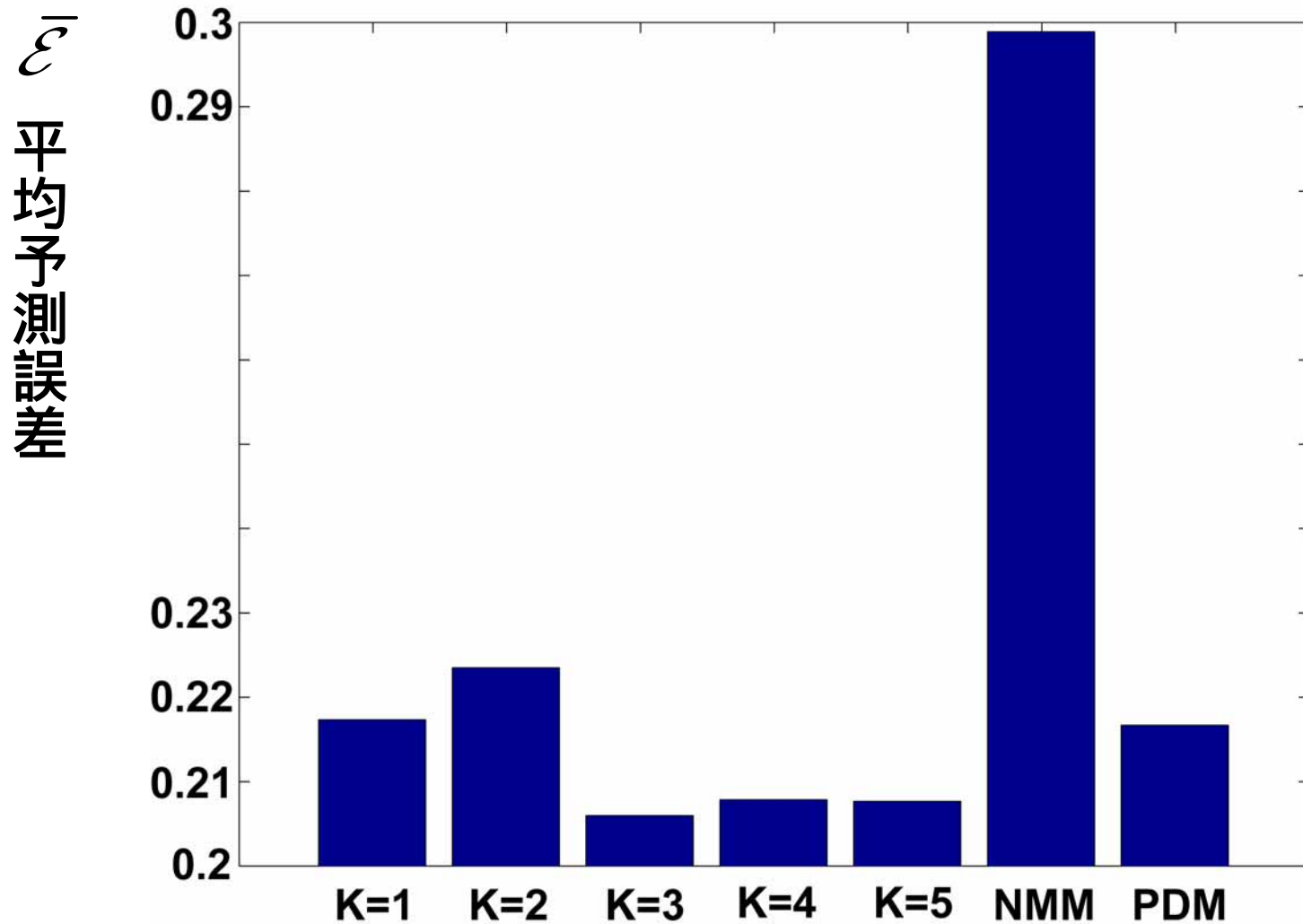
# サイト訪問者数の変動例(実データ)

訪問者数 (人)

$$m_i(t)$$



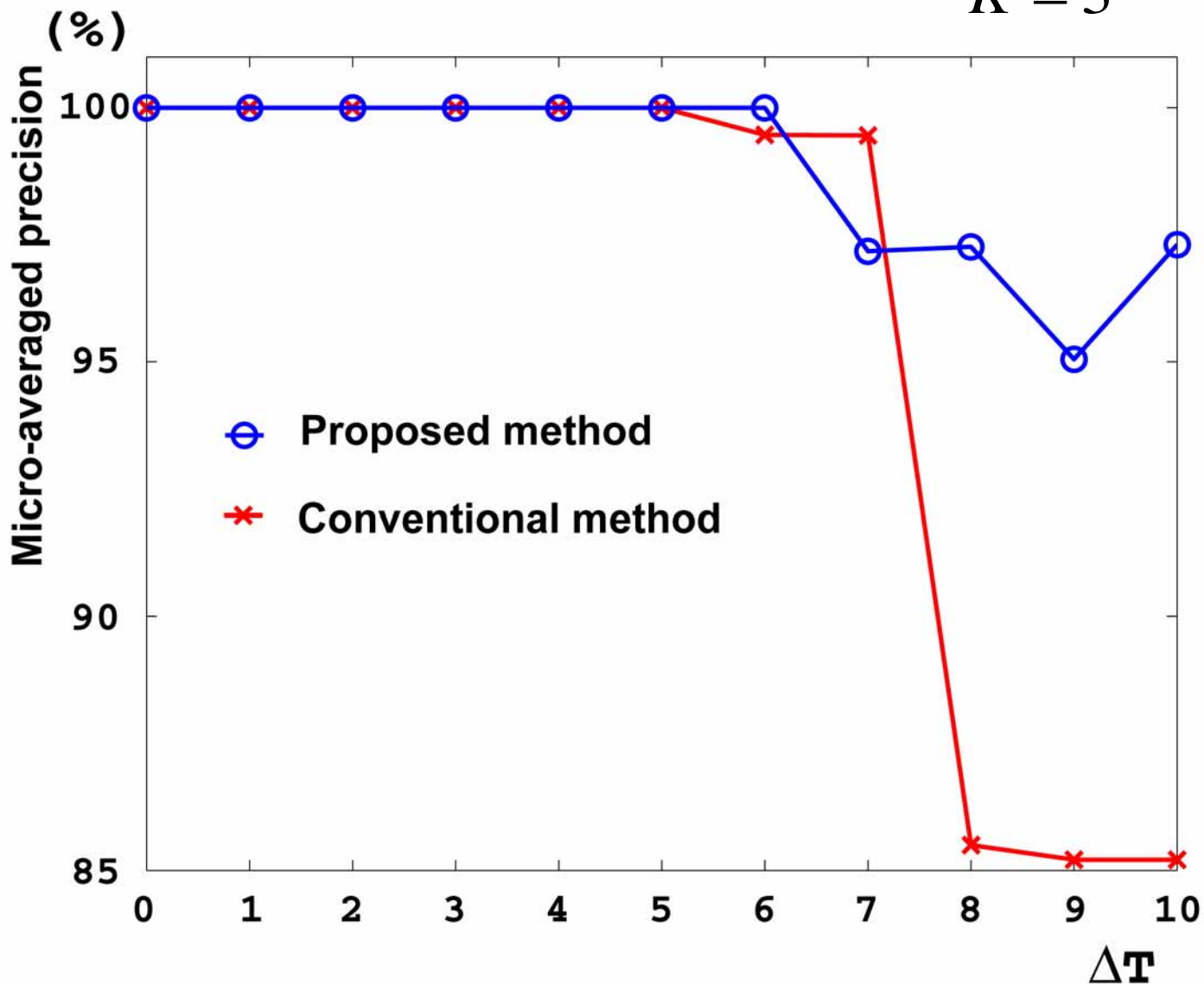
# 予測性能の比較(実データ)



# 分類性能の比較(実データ)

$K = 3$

$AP(\Delta T)$



注)  $\{m(0), \dots, m(T - \Delta T)\} \rightarrow AP(\Delta T)$

# 定性的評価（提案法による分類結果）

## (1) ポピュラーエンターテインメントサイト群

- 1) ポピュラー音楽ビデオサイト
- 2) 衛星テレビ局運営のインターネットTVサイト
- 3) 渋谷ライフスタイル・ショッピング情報サイト
- 4) イベント情報サイト
- 5) スポーツエンターテインメントサイト
- 6) 卓球情報サイト
- 7) インタラクティブドラマサイト

## (2) 芸術的趣味サイト群

- 8) 映像制作会社運営の次世代スポーツと古典芸能情報サイト
- 9) エコ商品情報サイト
- 10) 音楽情報サイト
- 11) ラジオ局運営のインターネット放送サイト
- 12) 地方のインターネット放送サイト
- 13) 古くからの映画会社運営の芸術的ドラマサイト
- 14) ミステリードラマサイト

## (3) 知識提供サイト群

- 15) 新聞社運営の総合情報サイト
- 16) ITニュース社運営の総合情報サイト
- 17) 通信会社運営の総合情報サイト
- 18) 地方のIT会社が運営の総合情報サイト
- 19) 子供の教育サイト
- 20) 育毛情報サイト

# まとめ

1. マーケットを形成するWebサイト集合における訪問者数の変動を、競合ダイナミクスの観点からモデリング
  - ・ 競合構造に基づいたダイナミクスモデルの提案：  
多項分布の確率的混合モデル
    - ~ 多項分布のパラメータ値が  
レプリケータ方程式で時間発展
  - ・ 学習アルゴリズムの導出：  
観測時系列データ →  
競合サイトのグループを同定  
各グループのレプリケータ方程式を推定
2. 予測モデル(定量的モデル)としての有効性の検証