

統計的機械学習による 都市インテリジェンス研究

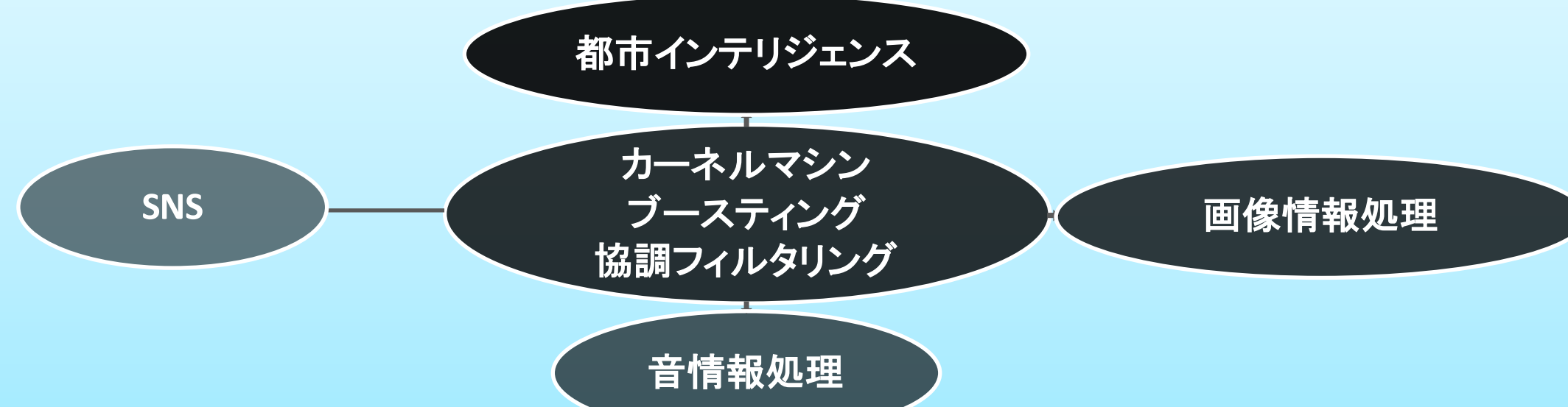
松井 知子 モデリング研究系 教授

【概要】

本研究室では統計的学習機械を用いて、音声/音楽/画像/SNSなどを処理する方法について研究しています。具体的にはカーネルマシン、ブースティング、協調フィルタリングの手法を用いて、

1. 音声・話者認識
2. 音楽情報処理
3. 画像識別
4. SNS解析
5. WEBユーザビリティ評価
6. 都市インテリジェンス など

の研究課題に取り組んでいます。



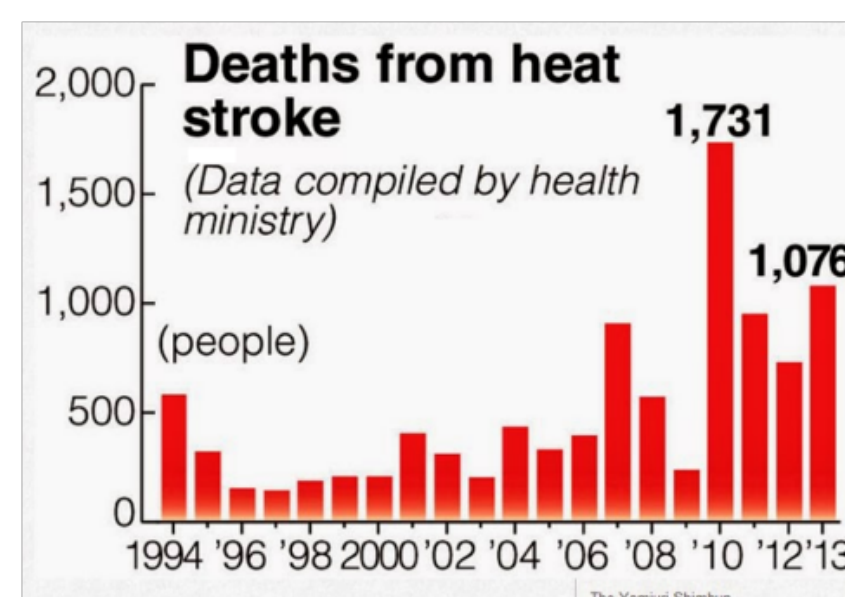
本研究室では統計的機械学習とその応用研究に興味のある学生さんを募集しています！

【統計的機械学習】

- 統計科学を用いて、
 - データから、内在する数学的な構造を発見する。
 - その数学的な構造に基づいて、予測や判別などの情報処理を行う。
- 帰納的アプローチ
 - V.S.
- 自然科学でよく見られる演繹的アプローチ
 - 仮説をたて、推論し、実験的または理論的に検証する。
- カーネルマシン
 - 自動的な特徴(ノモデル)選択機構を含む。
 - 非線形の扱いに優れている。
 - サポートベクターマシン(SVM)、罰金付ロジスティック回帰マシン
- いろいろな確率モデルによる方法
 - 混合ガウス分布モデル
 - 隠れマルコフモデル
- ガウス過程状態空間モデル など

【都市リスク管理のためのツイートデータと異常気象事象の時空間解析】

Introduction



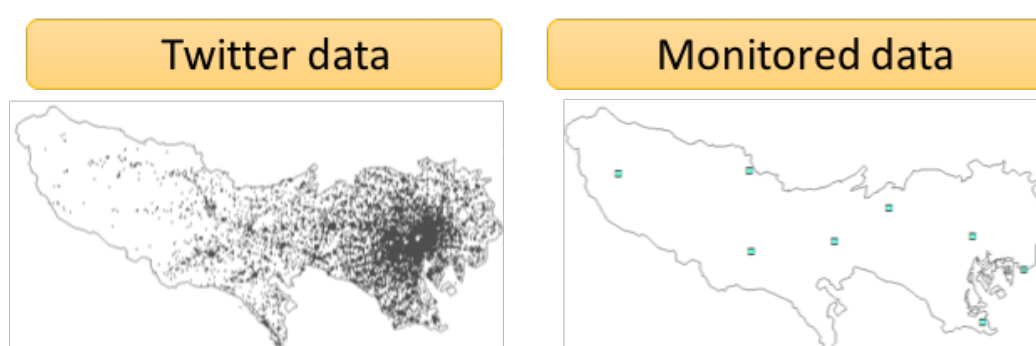
To cope with heatwave risks, it is important to know the risk in a local (e.g., district-level) and real-time manner.

Available temperature-related data

	Official monitoring data	Remote sensing data
Spatial resolution	Poor (8 monitoring stations)	High (1km grids)
Temporal resolution	Medium (every 1 hour)	Poor (10:30, 13:30, 20:30, 23:30)
Image	<p>Temperature</p> <p>□ 24-26 ■ 26-28 ■ 28-30 ■ 30-</p>	<p>Ground Temperature</p> <p>40.0 - 40.0 37.5 - 37.5 35.0 - 35.0 32.5 - 32.5 30.0 - 30.0 27.5 - 27.5 25.0 - 25.0 22.5 - 22.5 20.0 - 20.0 17.5 - 17.5</p>

Twitter data

- Geo-tagged tweets in Tokyo in August 2012
 - They consist of 1% of random samples
 - Sample size (after a cleaning): 130,331



Spatial resolution **High** Noise **Large**

Extraction of heat-tweets

Heat-tweets : Tweets including any of synonyms of "heat"		
あつい	Hot	蒸し
暑い	Hot	水分補給
暑	Heat	Rehydration
猛暑	Heat wave	体調管理
炎天下	Blazing sun	Health management
真夏日	Hot day	死に
残暑	Lingering summer heat	異常
熱中症	Heat illness	不快感
バテ	Faint	Discomfort
寝苦しい	Cannot sleep well	嫌
夏本番	Midsummer	Unpleasant
日差し	Sunlight	不快
照り	Reflected heat	クソ
湿度	Humid	Orz
湿気	Moisture	Expletive
汗	Sweat	きつい
ジメジメ	Damp	辛い
ムシムシ	Humid	大変
ベタベタ	Sticky	しんどい
		Severe
		苦手
		Weak

Objective

As a first step of using tweets for a real-time and local heatwave management, we examine whether tweets are useful for local real-time temperature estimation.

- Research questions
 - Do tweets explain temperature?
→ A correlation analysis between heat-tweets and temperatures is conducted.
 - How can we utilize tweets as a participatory sensing data that complement unobserved temperatures?
→ The spatial best linear unbiased estimator (S-BLUE) is applied.

Model

- Temperature:** $y(\mathbf{x}_i)$
 - Temperatures are describes by a latent spatiotemporal process, $f(\mathbf{x}_i)$, and a noise process, $v(\mathbf{x}_i)$.

$$y(\mathbf{x}_i) = f(\mathbf{x}_i) + v(\mathbf{x}_i)$$

$$f(\mathbf{x}_i) \sim N(0, c(\mathbf{x}_i, \mathbf{x}_j))$$

$$v(\mathbf{x}_i) \sim N(0, \sigma^2)$$

\mathbf{x}_i : Spatiotemporal coordinates of monitored temperatures
 $c(\mathbf{x}_i, \mathbf{x}_j)$: Covariance between sites \mathbf{x}_i and \mathbf{x}_j
 σ^2 : Variance parameter

- Heat-tweets:** $y_{ht}(\mathbf{x}_i)$
 - It assumes that the probability of being a heat-tweet increases if $f(\mathbf{x}_i)$ exceeds a threshold temperature T .

$$P(y_{ht}(\mathbf{x}_i) = 1 | f(\mathbf{x}_i)) = \sum_{p=0}^1 P(y_{ht}(\mathbf{x}_i) = 1 | y_{ht_0}(\mathbf{x}_i) = p) P(y_{ht_0}(\mathbf{x}_i) = p | f(\mathbf{x}_i))$$

$$y_{ht_0}(\mathbf{x}_i) = \begin{cases} 1 & \text{if } f(\mathbf{x}_i) \geq T \\ 0 & \text{otherwise} \end{cases}$$

\mathbf{x}_i : Spatiotemporal coordinates of i -th tweet
 $y_{ht}(\mathbf{x}_i)$: 1 if i -th tweet is a heat-tweet, and 0 otherwise

Spatiotemporal temperature estimation equation

- Spatial best linear unbiased estimator (S-BLUE)**
 - It estimates the temperature at \mathbf{x}_* , by minimizing the error variance.
- $$\hat{f}_* = \hat{\alpha} + \hat{\mathbf{B}}\mathbf{Y} = \arg \min_{\alpha, \mathbf{B}} E[(f_* - (\alpha + \mathbf{B}\mathbf{Y}))^2] = E[f_* \mathbf{Y}'] E^{-1}[\mathbf{Y} \mathbf{Y}'] \mathbf{Y}$$
- f_* : The value of the temperature process at site \mathbf{x}_*
 α : Constant
 \mathbf{Y} : A vector of monitored temperatures, $y(\mathbf{x}_i)$ and heat-tweets, $y_{ht}(\mathbf{x}_i)$
 \mathbf{B} : A matrix of unknown parameters

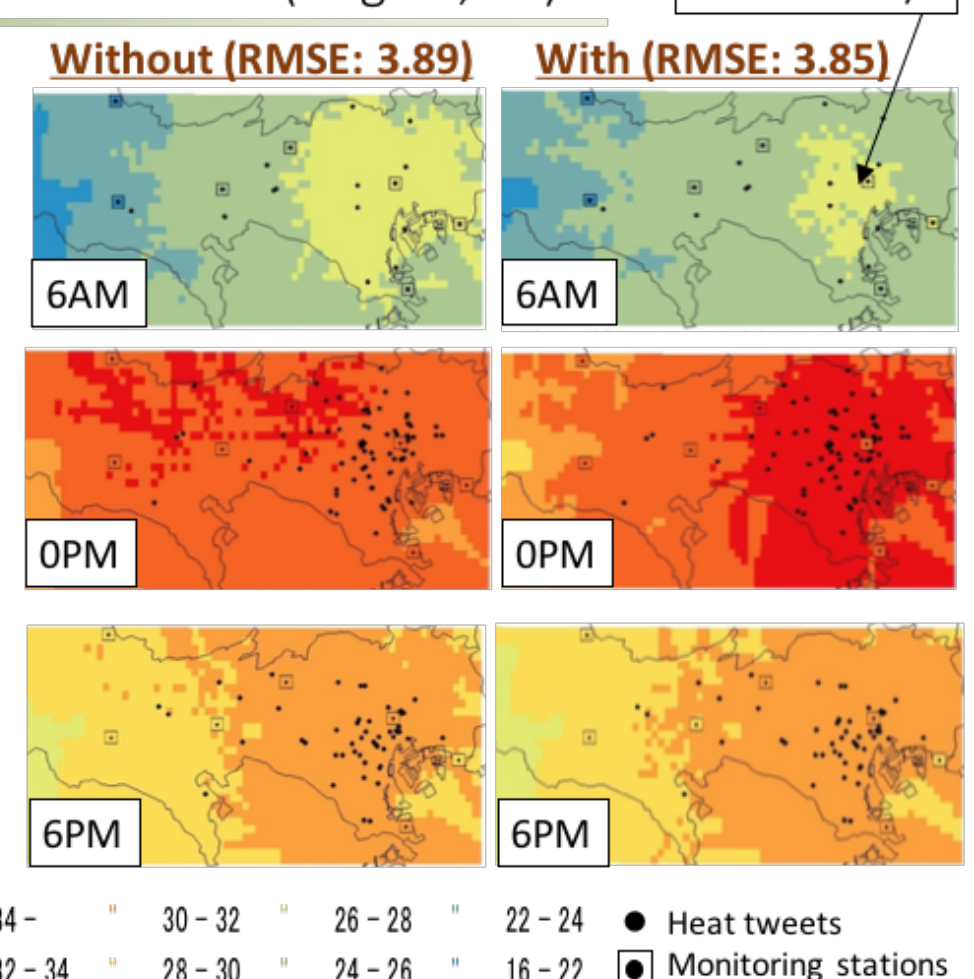
- $c(\mathbf{x}_i, \mathbf{x}_j)$ in $E[f_* \mathbf{Y}']$ and $E[\mathbf{Y} \mathbf{Y}']$ is usually given by a kernel function. However, the **kernel approach** is computationally expensive.
- Hence, we developed an EM-algorithm-based approach, which lighten the problem by defining $c(\mathbf{x}_i, \mathbf{x}_j)$ using radial **basis functions**.

Temperature estimation result (August, 25)

RMSE is improved by using tweets.

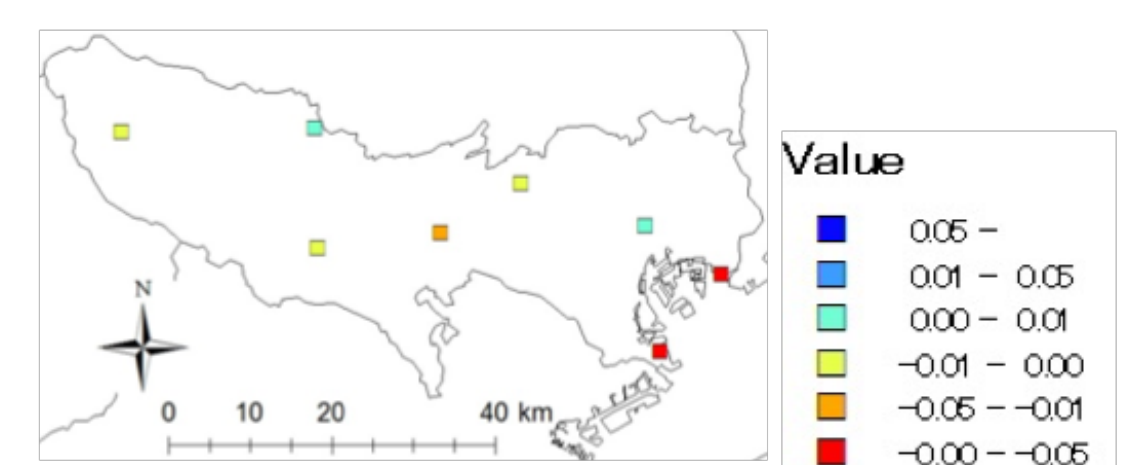
“With” approach successfully captures the heat island effect in the central area at the noon.

The strange result of “Without” approach at 0PM would be due to the small number of the monitoring stations (8).



Gap of RMSEs (With minus Without)

- Red : With twitter approach is better
- Blue: Without twitter approach is better



- The accuracy is improved in **the bayside area**
 - It might be due to the fact that heat-tweets capture the heat island effects in this area around noon.

Remaining issues

- Improvement of the heat-tweet extraction approach**
 - Some words might explain heat stronger than others
 - Some tweets might not comment about the current situation, but about future prediction.
- Difference between feeling temperature and actual temperature**
 - To analyze heatwave risks, it is needed to analyze feeling temperature, which depends on humidity, sunlight, mood/sentiment, etc.
 - Use of vital sensors would be needed.
- Fusion with other data**
 - Remo sensing data, data of vital sensors etc.
- The computational complexity of S-BLUE**
 - S-BLUE still requires the inversion of a covariance matrix, which becomes large if large samples are utilized.

共同研究者:
山形与志樹、村上大輔 (国立環境研究所)
Gareth W. Peters (UCL)