

# 十分定着した単語は使われ方は日々どのくらいいつつ変化しているか？

## — 様々な言語や媒体に共通してみられる単語出現頻度の対数拡散と非整数微積分 —

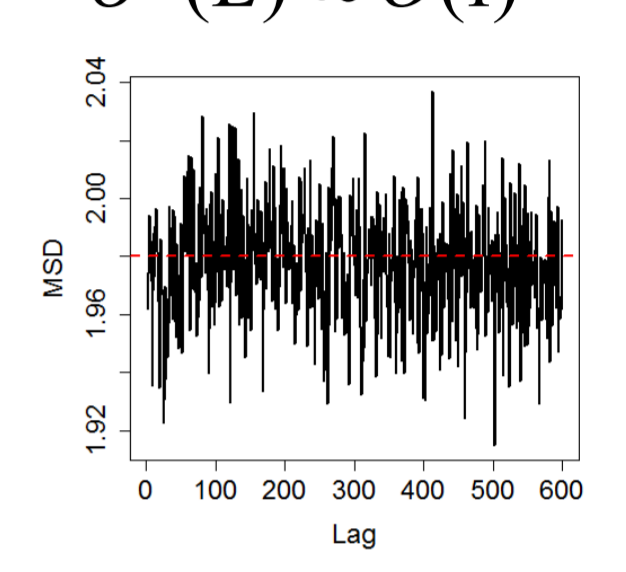
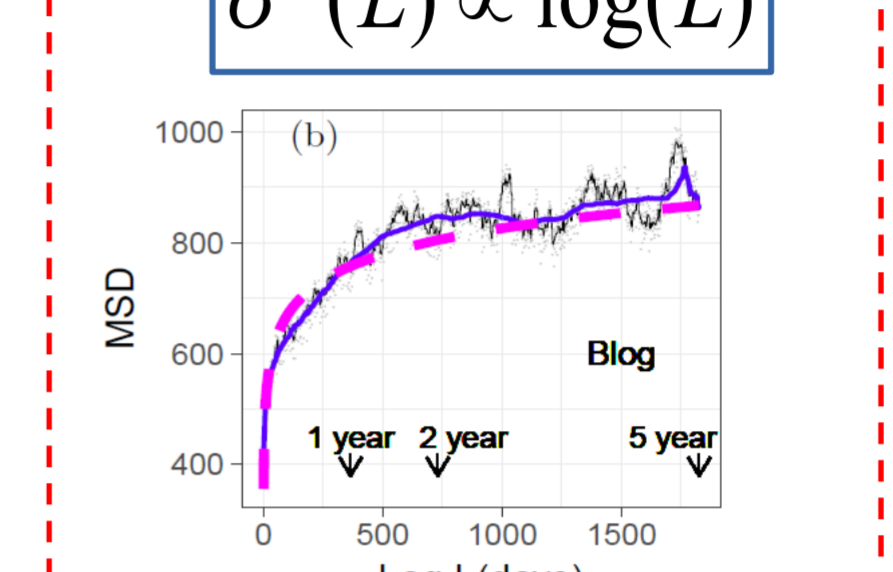
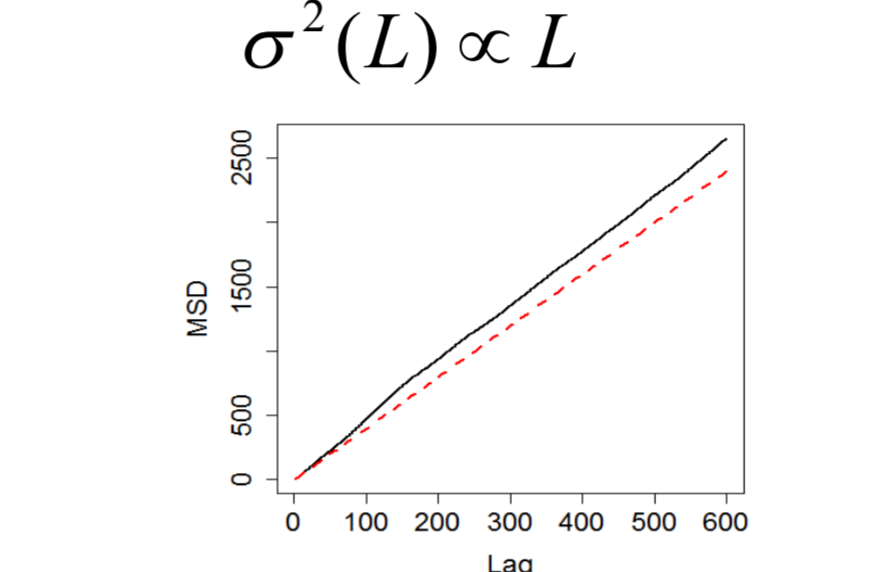
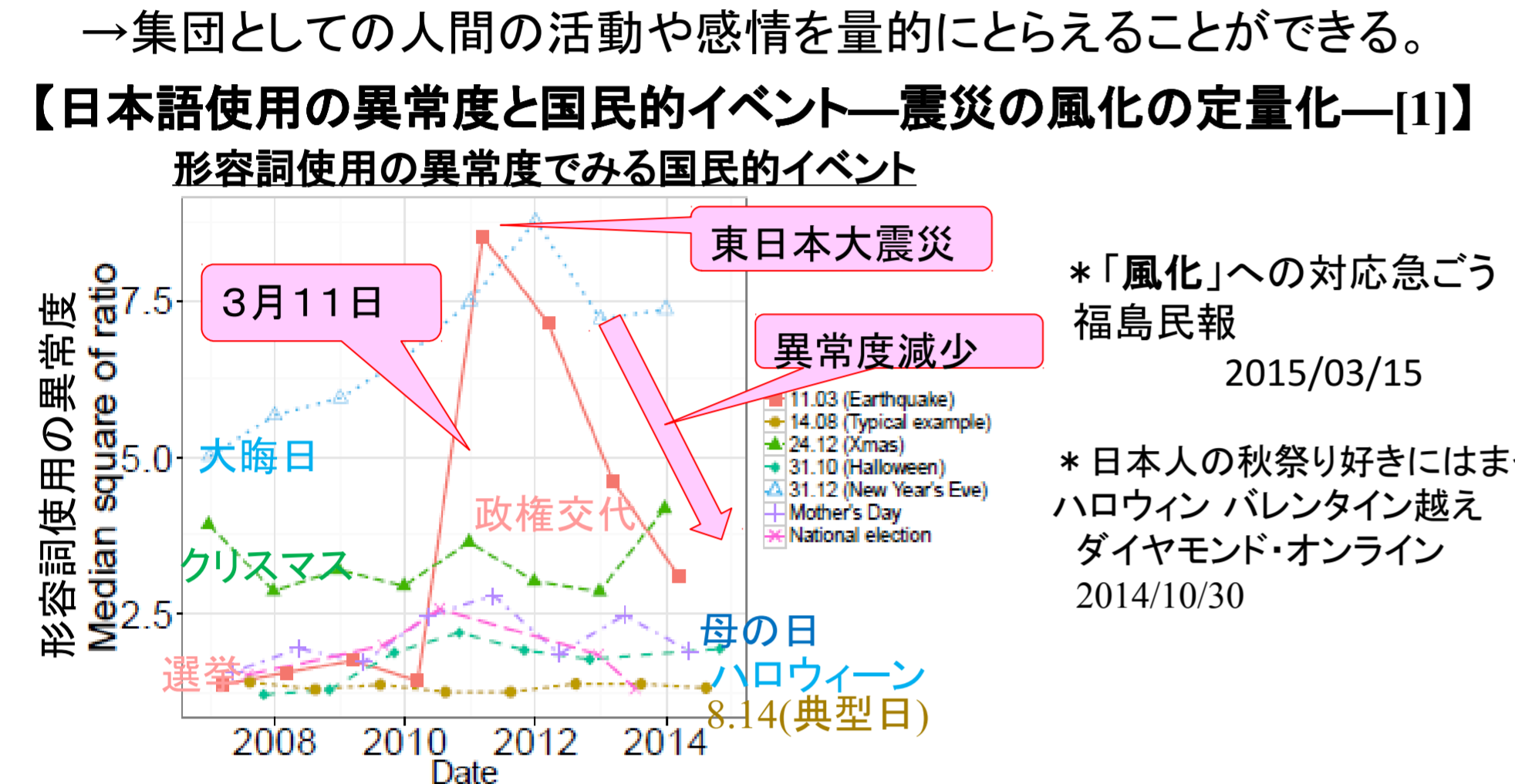
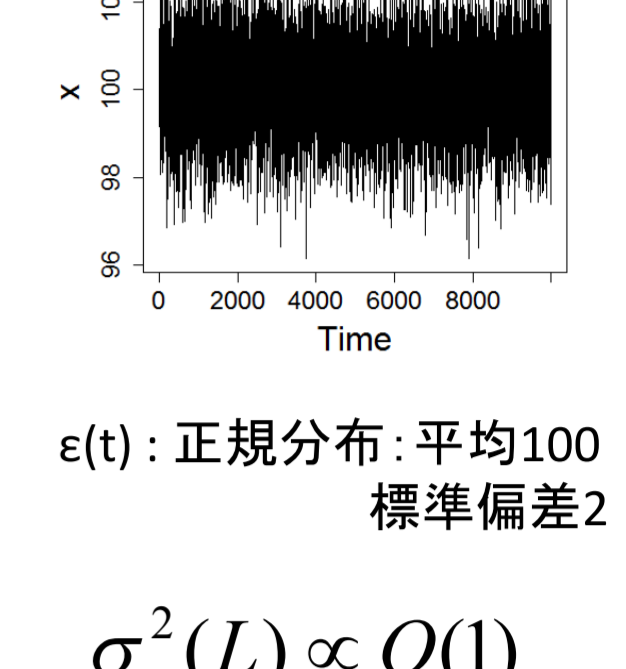
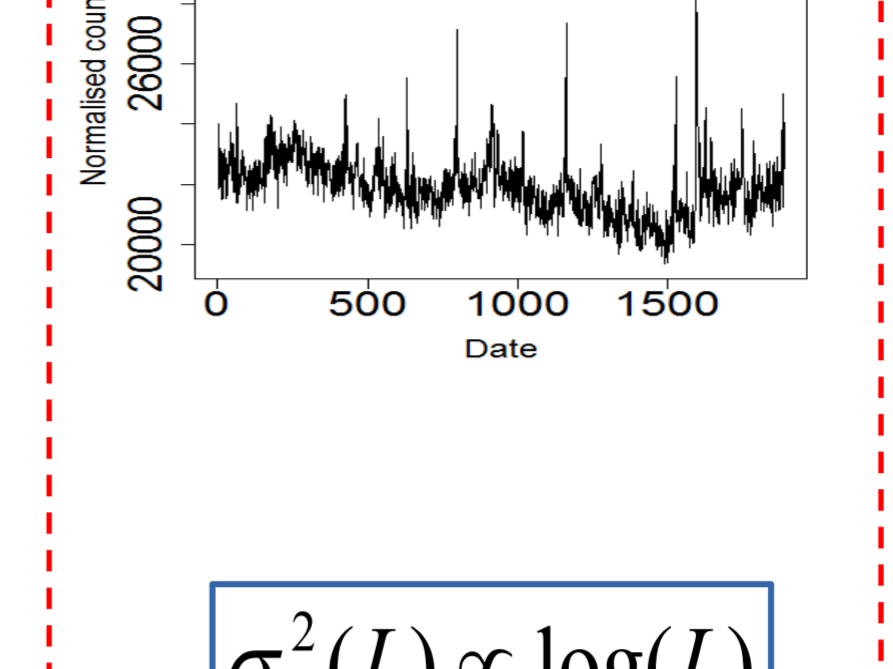
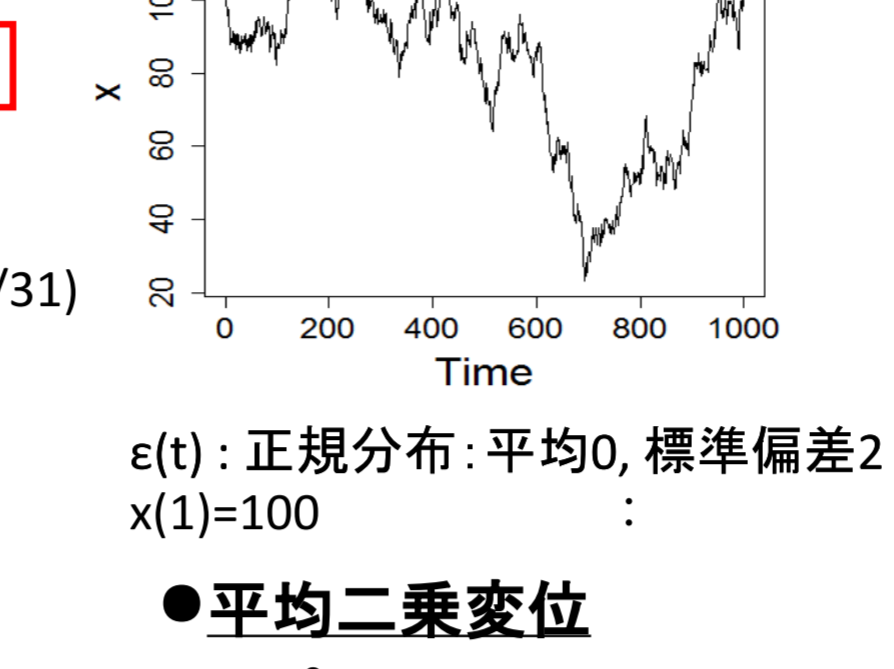
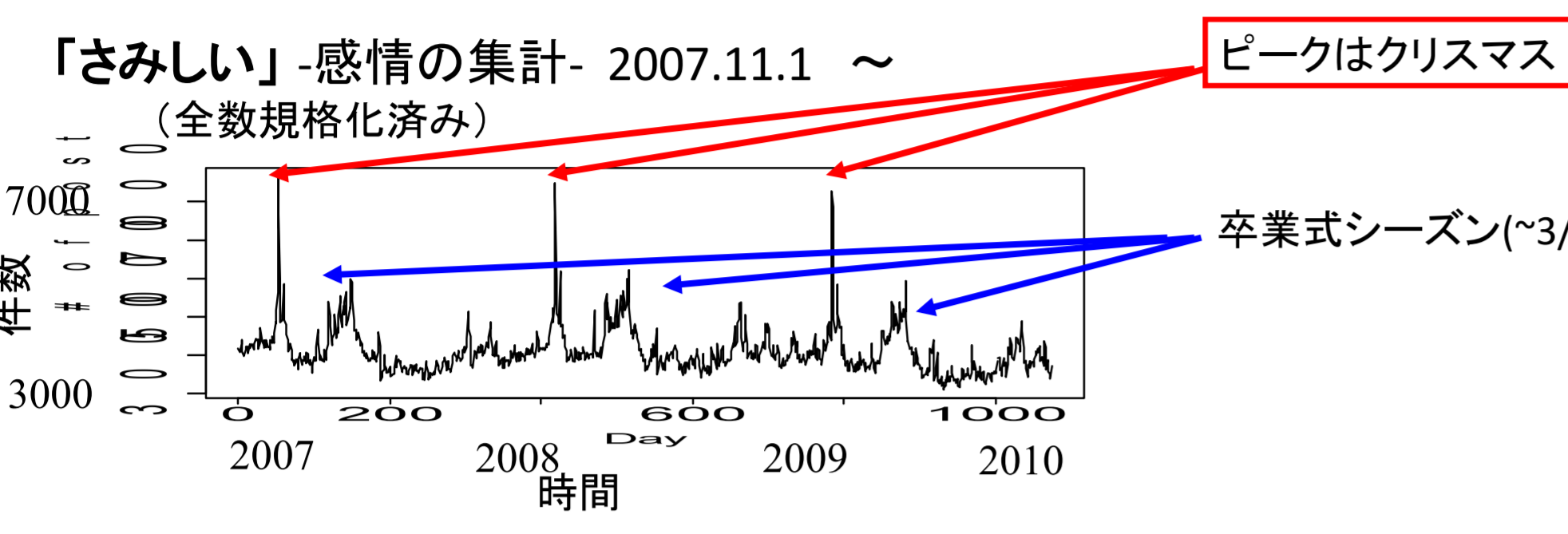
渡邊隼史 リスク解析戦略研究センター(統数研), 社会データ構造化センター(情・シス機構) 特任助教 h-wata@ism.ac.jp

【ブログデータ】

- Web上の日記のようなもの。日付きのテキストデータ
- データ収集は、ホットリンク社のロコミ係長を利用
- 2006年11月1日～までの日本語ブログ記事30億記事
- 基本的な量である着目キーワードの国内ブログ上での出現頻度時系列に着目

【結果1(データの解析): 十分定着した単語の拡散特性】

- 時系列
- ランダムウォーク(熱拡散、為替)  $x(t+1) = x(t) + \varepsilon(t)$
- 平均二乗変位  $\sigma^2(L) \propto L$
- 対数拡散  $\sigma^2(L) \propto \log(L)$
- 分布  $\varepsilon(t)$ : 正規分布: 平均0, 標準偏差2  $x(1)=100$
- 無相関定常ノイズ (毎日独立にずっと同じ)  $x(t) = \varepsilon(t)$
- 平均二乗変位  $\sigma^2(L) \propto O(1)$



【研究の目的】

- 社会の安定性を定量的・数学的に調べたい。
  - 社会は常に変化している部分もあるが、「コア部分」は安定しているように見える
  - 例えば、「こんにちは」という挨拶はここ30年くらいはあまりかわっていない。
  - では、どのくらい安定してどのように安定してるかを定量的に調べてみたい。
    - 「変わらないや安定」は社会の基盤なので目立つ大きな動きと同様に社会の動的な性質の理解に重要と考えられる。
- 今研究では、そのひとつの例として、言葉において「十分定着した単語」ほどのくらい安定しているかをデータを用いて精密に測定。
  - 例えば、「だから」、「重い」などの十分定着している語は日本のブログでの書き込み件数がほぼ定常で安定しているように見える。
  - しかし、少しづつは変化しているはずで…
  - 実際にはどのくらい変化してるかを日次の解像度で調べてみたい。

不安定 → 安定

対数拡散

(1) 定常ではないがとてもゆっくりした拡散  
→ 言葉は安定してるが徐々に変化?

(2) 自然現象ではほとんど観測されない。  
→ 物理で理論予測はあったが実観測がなかった

(3) 新聞記事数データや英語・仏語・中国語・日本語  
○ Wikipedia page view でも同等な対数拡散を観測

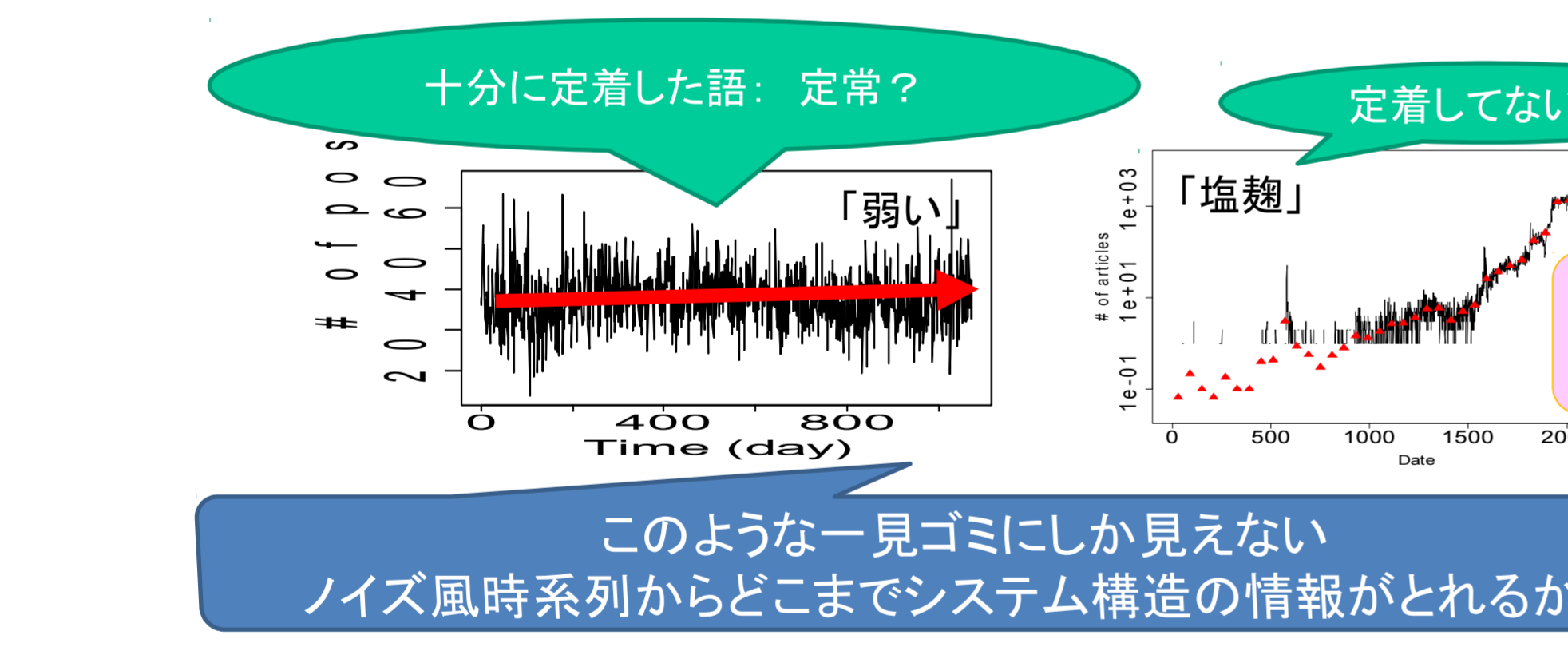
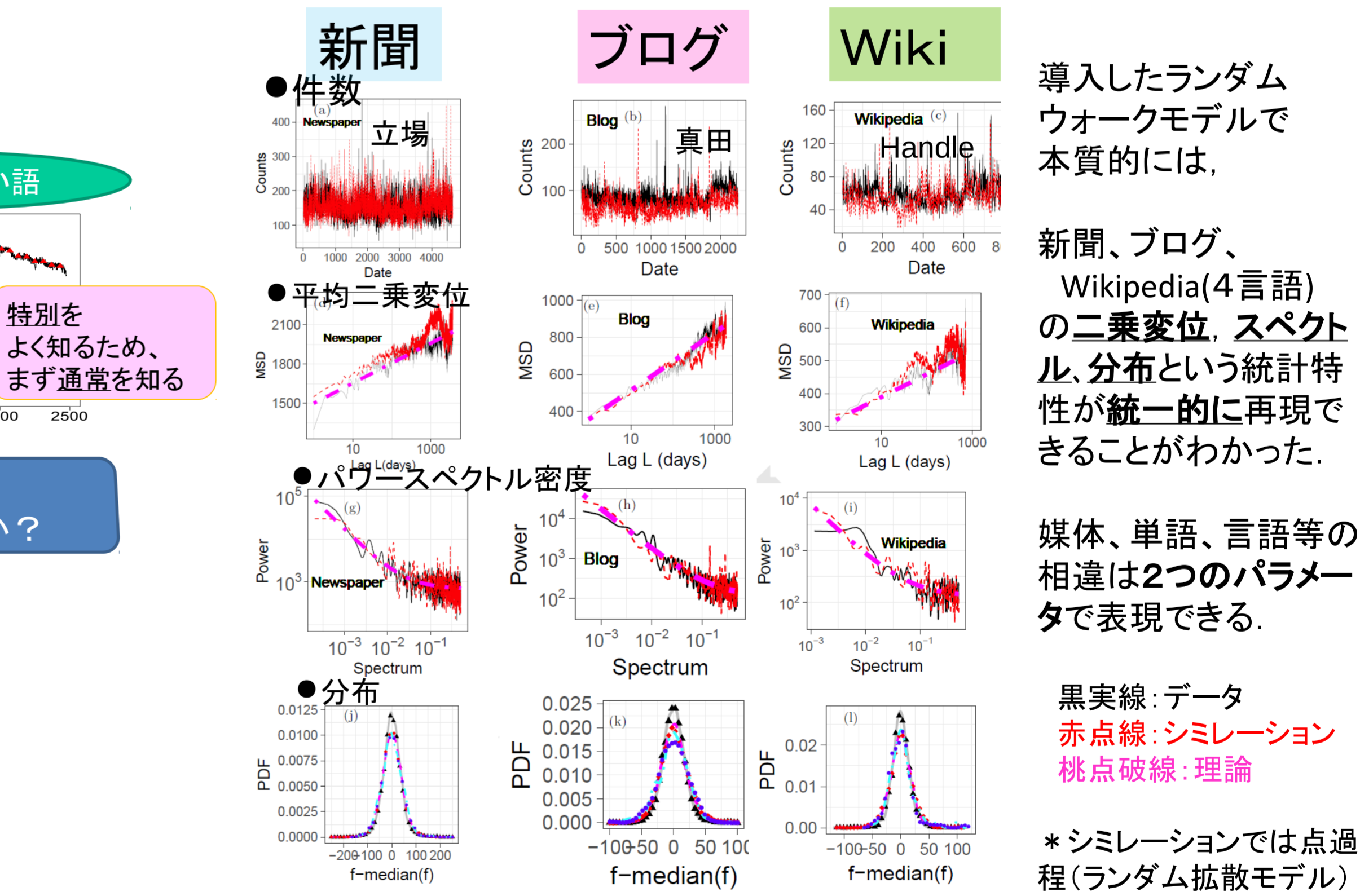
【結果2(数理のモデル解析): 非整数微分との関係】

- べき忘却ランダムウォークモデルが性質をよく再現する

$$r_j(t) = r_j(0) + \sum_{s=1}^t (t-s+1)^{-\beta} \eta_j(s)$$

→ とくにブログのような対数拡散は忘却速度  $\beta=0.5$

\* 時系列は「べき忘却+点過程(RDモデル)」でよく再現



【手法: 平均二乗変異】

- 今回の研究では「平均二乗変位」という時系列量を測定した
- 「平均二乗変位」は、着目日のL日語にどのくらい書き込み数が平均的に変化するかをあらわした量(拡散特性)

$$\sigma^2(L) = \langle (x(t+L) - x(t))^2 \rangle$$

t+L日の書き込み数が ある時間tの書き込み数に対して

どれくらい変化しているか? 平均

$$\sigma^2(L) \approx \frac{\sum_{t=1}^{T-L} (x(t+L) - x(t))^2}{T-L}$$

色々なLについて計算

【考察: 忘却速度  $\beta=0.5$  の意味】

(1) 定常と非定常の境界 (非常に微妙な立ち位置のパラメータ)

$$\langle (r_j(t+L) - r_j(t))^2 \rangle \propto \begin{cases} L^{1-2\beta} & (0 \leq \beta < 0.5) \text{ 非定常} \\ \log(L) & (\beta = 0.5) \text{ 単語カウント} \\ O(1) & (\beta > 0.5) \text{ 定常} \end{cases}$$

忘れ方が  $\beta=0.5$  より少しでも大きければ「定常」に  
忘れ方が  $\beta=0.5$  より少しでも小さければ「非定常」になる。

• べき忘却ランダムウォーク(積分形)  $r_j(t) = \frac{1}{\Gamma(1-\beta)} \sum_{s=0}^t (t-s)^{-\beta} \cdot \eta_j(t-s),$

• 微分形  $\frac{d^{1-\beta}}{dt^{1-\beta}} r(t) = \eta(t),$

$\left(\frac{d^{1/3}}{dt^{1/3}}\right)^3 x(t) = \frac{d^{1/3}}{dt^{1/3}} \frac{d^{1/3}}{dt^{1/3}} \frac{d^{1/3}}{dt^{1/3}} x(t) = \frac{d}{dt} x(t).$

The of

• ランダムウォーク (1階微分)  $\frac{d}{dt} r_j(t) = \eta(t)$

• ブログなど(1/2階微分)  $\sqrt{\frac{d}{dt}} r(t) \approx \eta(t),$

• IIDノイズ(0階微分)  $\frac{d^0}{dt^0} r_j(t) = r_j(t) = \eta(t)$