

- 潜在変数モデルなどの確率モデルのパラメタ推定に事前知識を活用したい (正則化, 制約, etc.)
- 特に**専門家の知識**を有効利用したいが, データに関する知識を直接記述するのは高コスト or 不可能
- 一方, データ生成プロセスの**間接的・あいまいな関係性**は比較的容易に記述できる場合がある
  - (例1) データ内の変数同士の正確な関係性は複雑/未知だが, ありうる依存関係の有無は簡単に列挙できる
  - (例2) センサの示す物理量の間依存関係は複雑/未知だが, センサ搭載機器間の関係は簡単に列挙できる

### 直接的・詳細な関係性

- $\omega = \text{const. if } V > V_{th}$
- $P \propto V^3$
- $\omega f(\varepsilon) \propto P$
- $\varepsilon = g(P, \omega, T) \dots\dots$

→ モデリングにそのまま利用できるが高コスト/未知のため記述しにくい



### データ

- 属性①  $V$  風速
- 属性②  $T$  気温
- 属性③  $P$  発電量
- 属性④  $\omega$  回転数
- 属性⑤  $\varepsilon$  翼歪
- ⋮

### 間接的・あいまいな関係性

- ① ←(関係あり)→ ③
- ④ ←(関係あり)→ ⑤
- ② ←(関係なし)→ ①
- ⑤ ←(未知)→ ② ⋯⋯

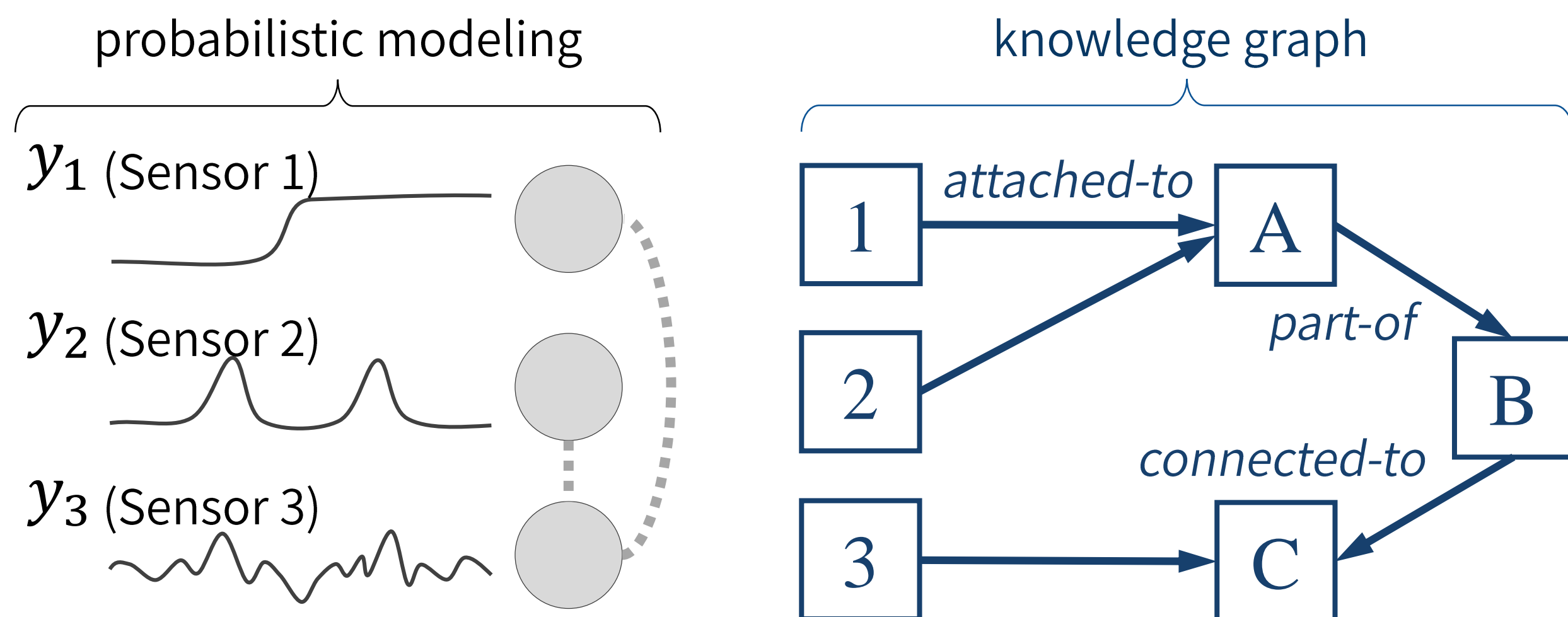
→ 多くの場合低コストで記述できるが既存手法ではそのまま利用しにくい



## Knowledge-Based Distant Regularization

### 想定する状況

- 潜在変数モデルなどの確率モデルを学習したい
- 利用したい間接的な知識が**ナレッジグラフ**で表せる
- データ内の(少なくとも一部の)オブジェクトや属性に対応するエンティティがナレッジグラフ中にある



### 提案手法の概要

- ナレッジグラフのエンティティの**連続値埋め込み**<sup>[1]</sup>を利用して**確率モデルパラメタの事前分布**を用意

### 関連研究

- グラフ疎性正則化<sup>[2]</sup>, グラフラプラシアン正則化<sup>[3]</sup>
- 統計的関係学習の方法による論理ルールの利用<sup>[4]</sup>等

## 例: Distantly Regularized PCA

### 観測モデル (確率的PCA<sup>[5]</sup>)

$$p(Y | X, W, \mu, \sigma^2) = \prod_i \prod_j \mathcal{N}(y_{i,j} | w_j^T x_i, \sigma^2)$$

- $y_{i,j} \in \mathbb{R}$  オブジェクト  $i$  の属性  $j$  の観測値
- $x_i \in \mathbb{R}^d$  オブジェクト  $i$  の隠れ変数 (主成分スコア)
- $w_j \in \mathbb{R}^d$  属性  $j$  のパラメタ (因子負荷行列の一部)

### 提案手法によるパラメタの事前分布

$$p(w_j | e_j) = \mathcal{N}(w_j | f_\xi(e_j), g_\xi(e_j))$$

where  $f_\xi: \mathbb{R}^k \rightarrow \mathbb{R}^d, g_\xi: \mathbb{R}^k \rightarrow \mathbb{R}_+$

- 属性  $j = 1, \dots, m'$  に対応するエンティティがナレッジグラフ中にあるとする
- $e_j \in \mathbb{R}^k$  属性  $j$  に対応するエンティティの埋め込み

### 最終的な目的関数

$$\mathcal{L}(W, \mu, \sigma^2, \xi, e_{1:m'}, e_{KG}, r_{KG}) = (\text{PCA+W事前分布の周辺尤度}) + \psi(e_{1:m'}, e_{KG}, r_{KG})$$

- $\psi(e, r)$  ナレッジグラフ埋め込みのスコア

今後, 理論的な解析・保証が必要!

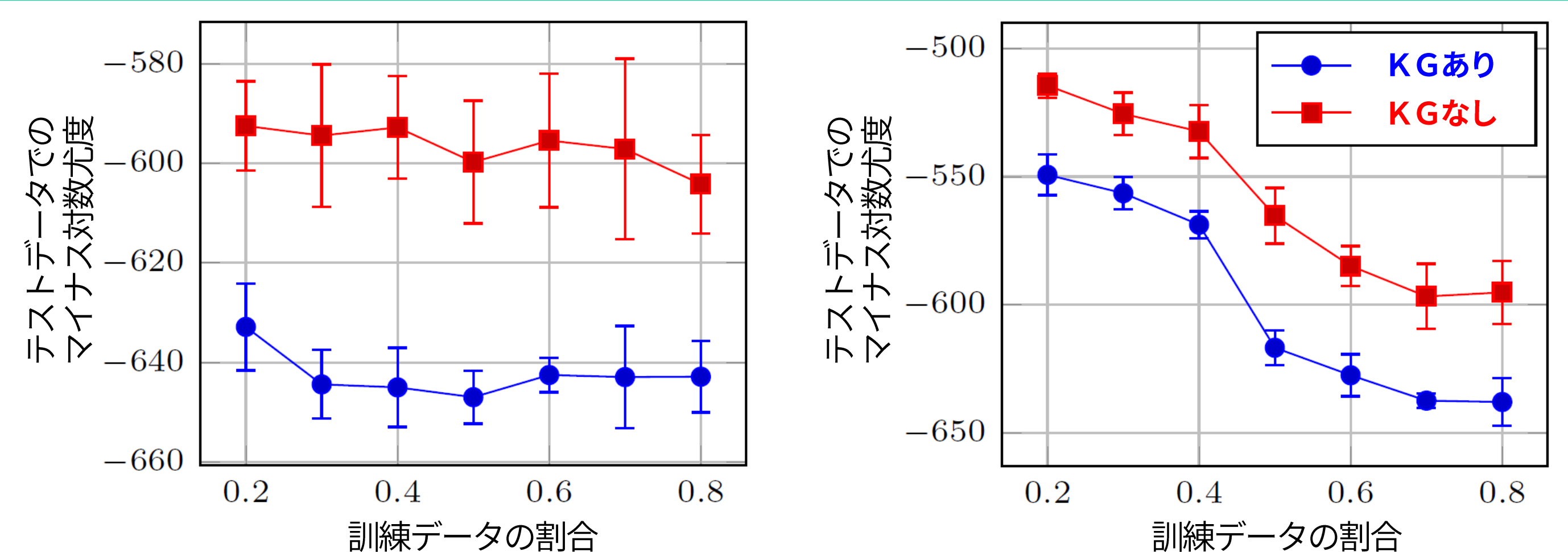
## 予備実験: Distantly Regularized PCA の汎化性能

### データ

- 227の国・地域の月平均降水量 (1901–2015年)<sup>[6]</sup>

### ナレッジグラフ

- 国・地域の地理的知識<sup>[7]</sup>をナレッジグラフとして利用
  - (Japan, is-inside, Asia), (Norway, is-neighbor, Sweden) ⋯
- しかし, これは間接的な知識しか与えない
  - 国の地理的関係性は必ずしも気象の関係性を与えない
  - is-neighborは陸の隣接のみ記述 → グラフ正則化は適用し難い



ナレッジグラフによる事前分布を**与える場合**と**与えない場合**について, テストデータに関するマイナスイタナ対数尤度 (小さい値の方がよい) を比較.

ナレッジグラフ埋め込みのスコア  $\psi$  としてシンプルな DistMult<sup>[8]</sup> を使用.

(左) 訓練/テスト分割前にデータをランダムに入替 (右) ランダム入替なし