

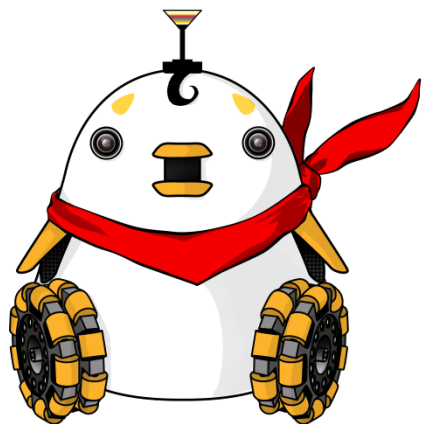
人工知能

第11章 学習と認識(1)

クラスタリングと教師なし学習

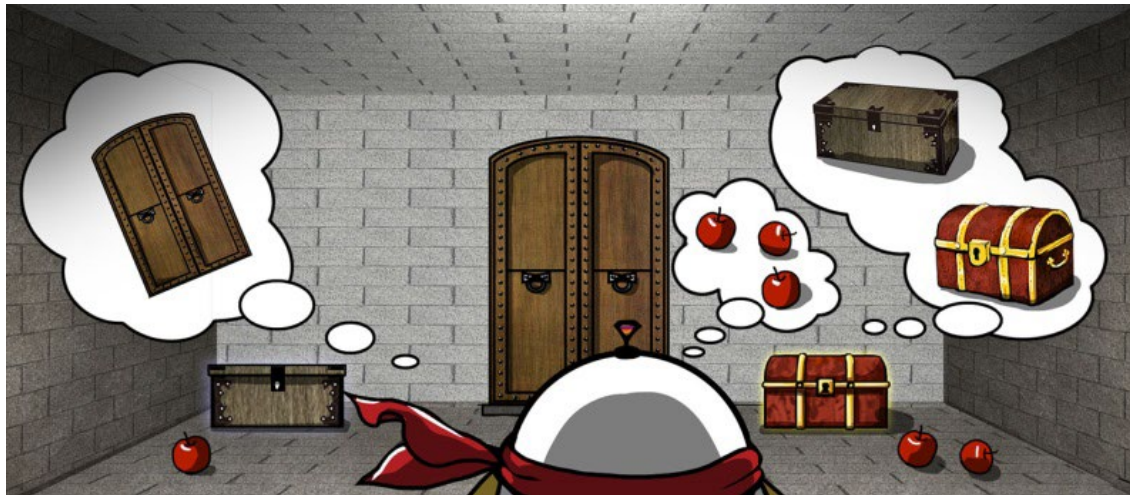
立命館大学 情報理工学部

谷口彰



STORY クラスタリングと教師なし学習

- 迷路を探索し，通り抜ける方法もわかった．自分の位置を見失っても自己位置推定で思い出すことができる．ホイールダック2号はこれで大丈夫だと思った．「さあ，お宝にとってゴールに向かうぞ！」
- しかし，ちょっと待てよ．「お宝」や「ゴール」って何だろう．
「**お宝**」とは**どんなもの**で「**ゴール**」って**どんな見た目**なんだろう．
ホイールダック2号は地図はわかるが，目の前に「お宝」や「ゴール」があったとしても，それが「お宝」や「ゴール」であることを**認識**することができない．まずは，「お宝」や「ゴール」とは**どんなものなのか，学習**していないと話にならない．



仮定 クラスタリングと教師なし学習

- ホイールダック 2 号は適切な画像特徴量を**有限次元ベクトル**として取得できるものとする.

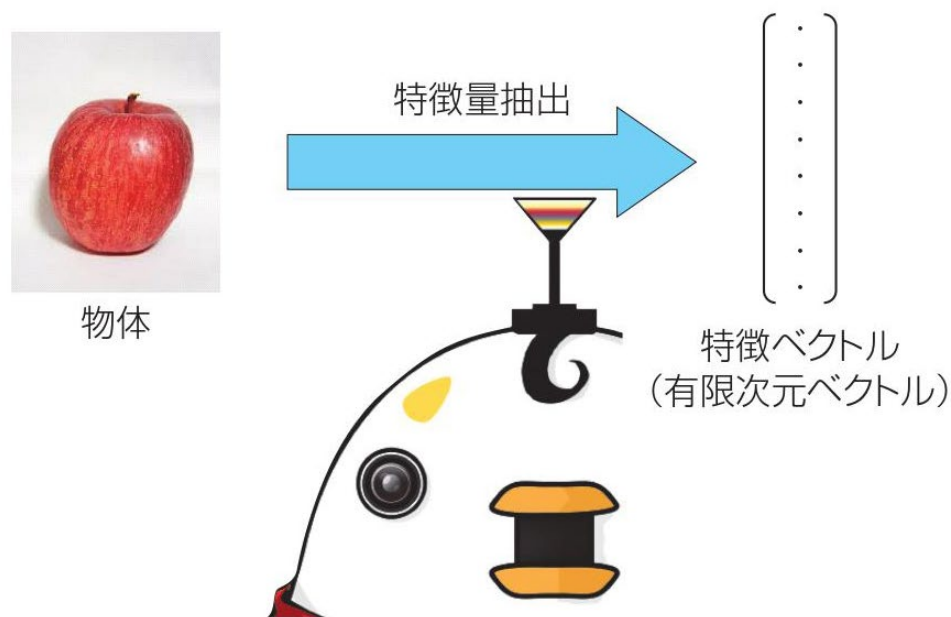


図 11.4 対象物体からの特徴ベクトル抽出

Contents

□11.1 クラスタリング

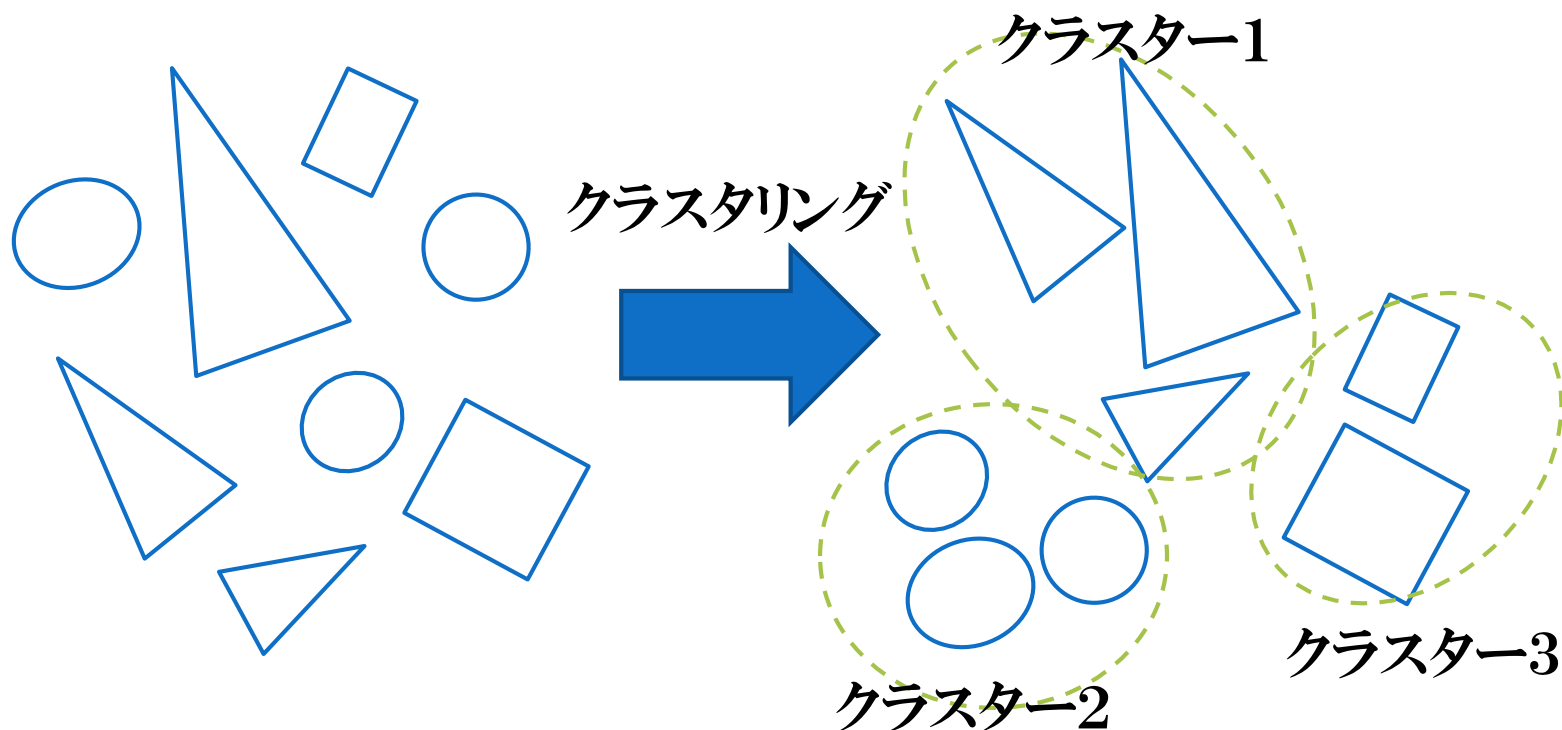
□11.2 k-means法

□11.3 混合分布モデル

□11.4 表現学習

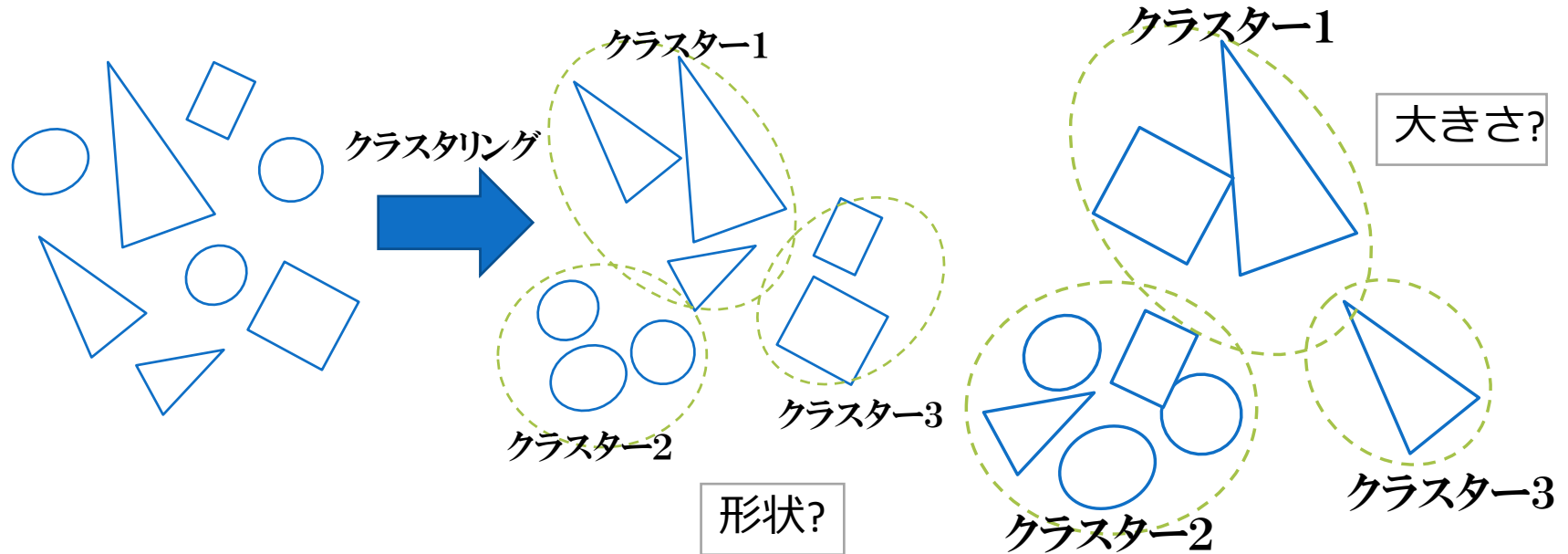
11.1.1 クラスタリングとは何か？ (clustering)

- データの集まりをデータ間の類似度にしたがっていくつかのグループに分類することを**クラスタリング**という。
- この作業を自動化するのが機械学習におけるクラスタリングという種類に属する手法
- 自ら概念を獲得するロボットをつくらうとする場合にはクラスタリングは重要な要素技術になる。



11.1.2 特徴抽出

「自然な」クラスタリングとは？



- ロボットにとってこのグループ分けが「自然な」ものであるかどうかは、ロボットにどのような基準を与えるかに依存する。
- そのような類似性を定義するために、**特徴量**(feature value)や**特徴ベクトル**(feature vector)によって張られる**特徴空間**(feature space)の設計が重要になる。

特徴量抽出とクラスタリング

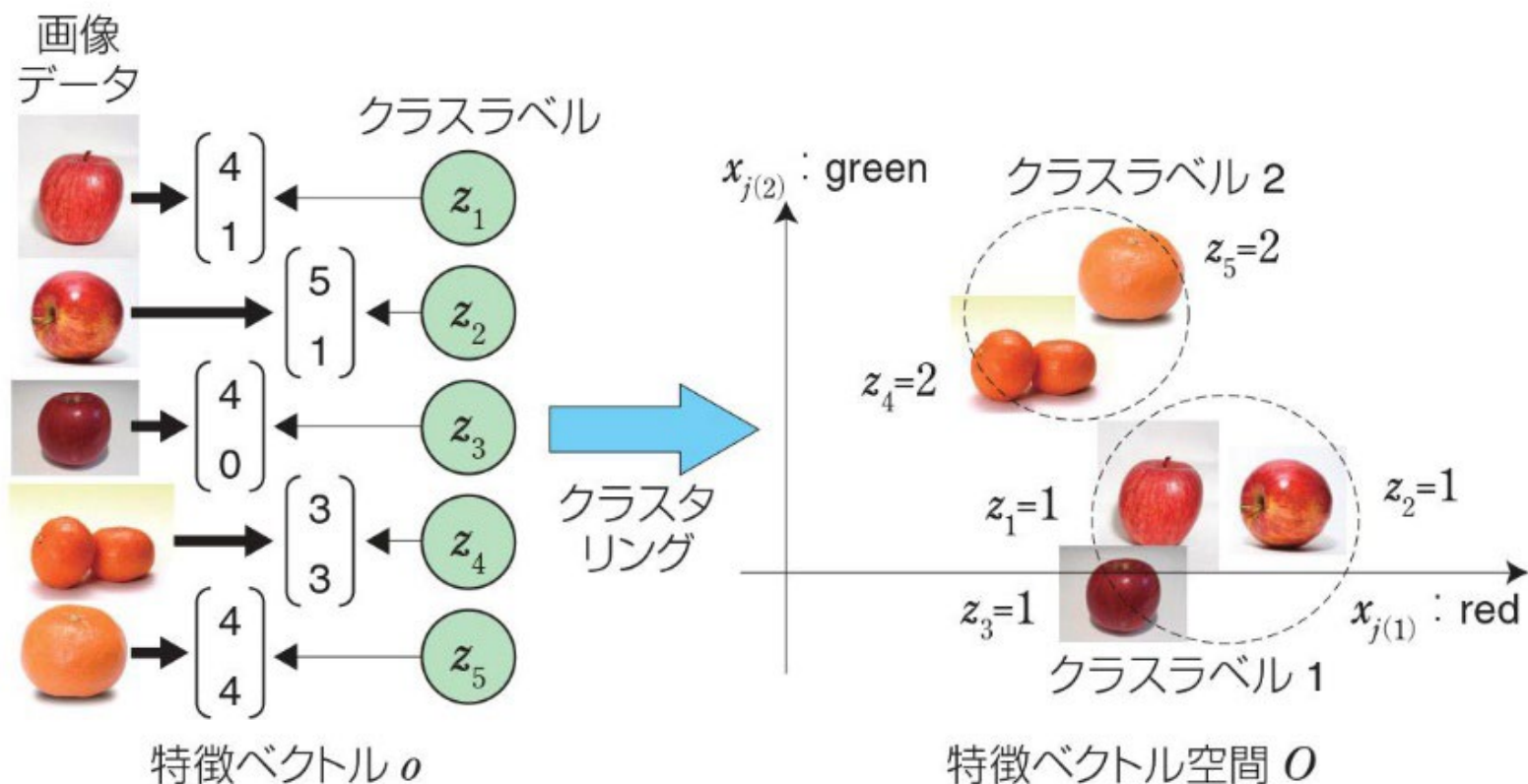


図 11.5 特徴ベクトルの抽出とクラスタリング

対象が特徴空間上の点として表されると、クラスタリングは特徴空間上の点をグループ分けする数学的な問題になる。
(潜在的なクラスラベルを推定することに等しい)

教師なし学習

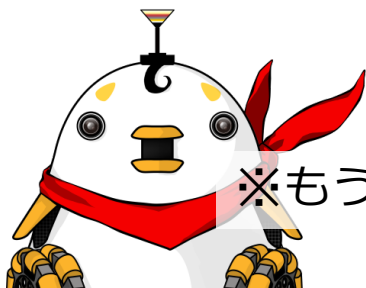
□入力として与えられたデータに潜む知識を発見する方法

□クラスタリング

- 大量のデータを幾つかのグループに自動的に分類する.
- 分類問題を教師データを用いずに行う.

□表現学習

- 高次元のデータをより潜在空間（多くの場合は低次元）に写像することで、データを説明する少数のパラメータを発見する. または、可視化する.



※もう少し詳細な機械学習の分類は次回行います。

Contents

□11.1 クラスタリング

□11.2 k-means法

□11.3 混合分布モデル

□11.4 表現学習

11.2.1 k-means法のア​​ルゴリズム

Algorithm 11.1 k-means 法

- ① k 個のクラスタの代表点 $(\mu_1, \mu_2, \dots, \mu_k)$ を初期化する.
- ② repeat
- ③ 各データ点 x_i ($i \in \{1, 2, \dots, N\}$) について, x_i と μ_j の距離を $d(x, y) = \|x - y\|^2$ で測り, x_i のクラスラベルである z_i を x_i と最も近いクラスタ代表点 μ_j の添字 j に更新する.

$$z_i \leftarrow \underset{j}{\operatorname{argmin}} d(x_i, \mu_j) \quad (11.1)$$

- ④ μ_j を各クラスタに含まれるデータの重心値で更新する.
- ⑤ until すべてのクラスタの割り当て z_i が変化しなくなる.

□このアルゴリズムでコスト関数 J を単調減少させられる.

$$J = \sum_{j=1}^k \sum_{\forall i, z_i=j} \|x_i - \mu_j\|^2$$

k-means法の実行例

- $S=\{2,4,6,10,12,14\}$ という6個の一次元データがあったとする。これをk-means法を用いてクラスタリングする。
- 初期クラスターを $S_1=\{2,4,10\}$, $S_2=\{6,12,14\}$ とした際に, k-means法のアルゴリズムを実行する。

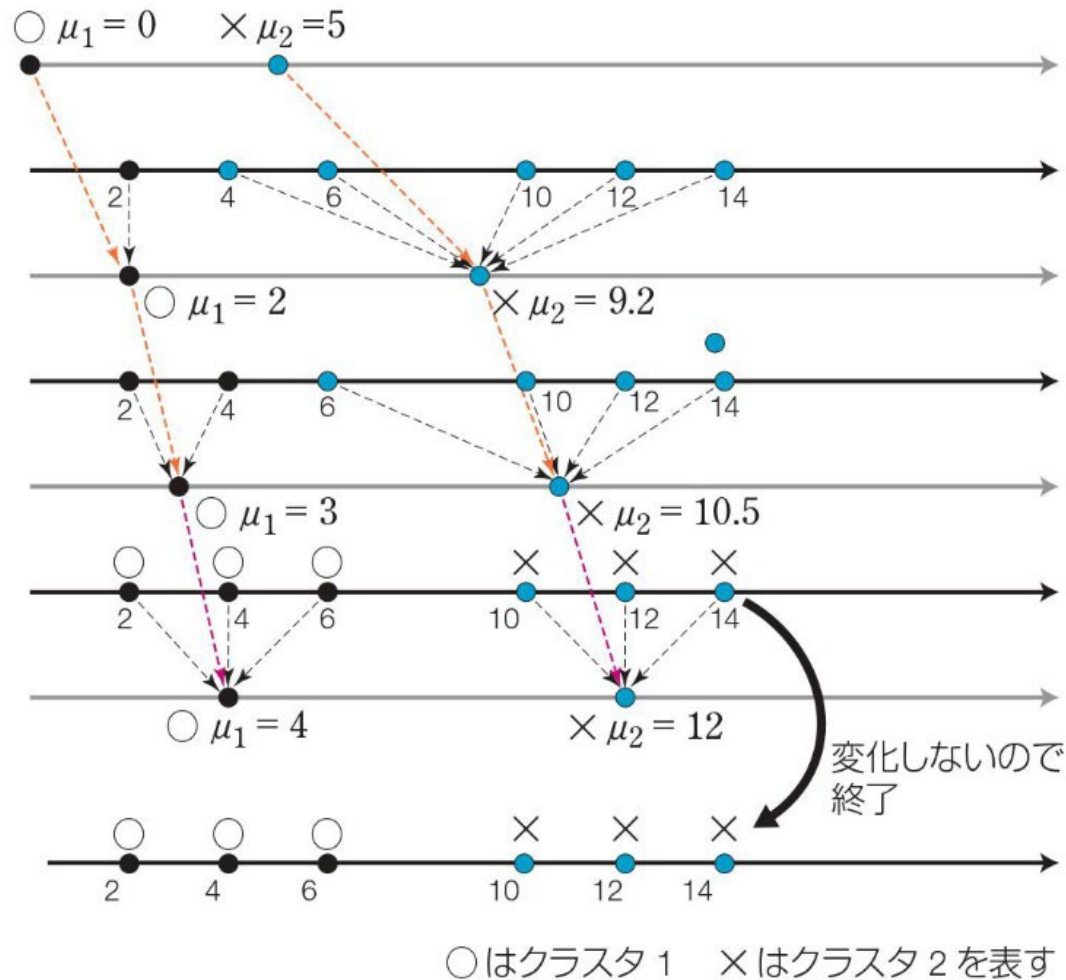
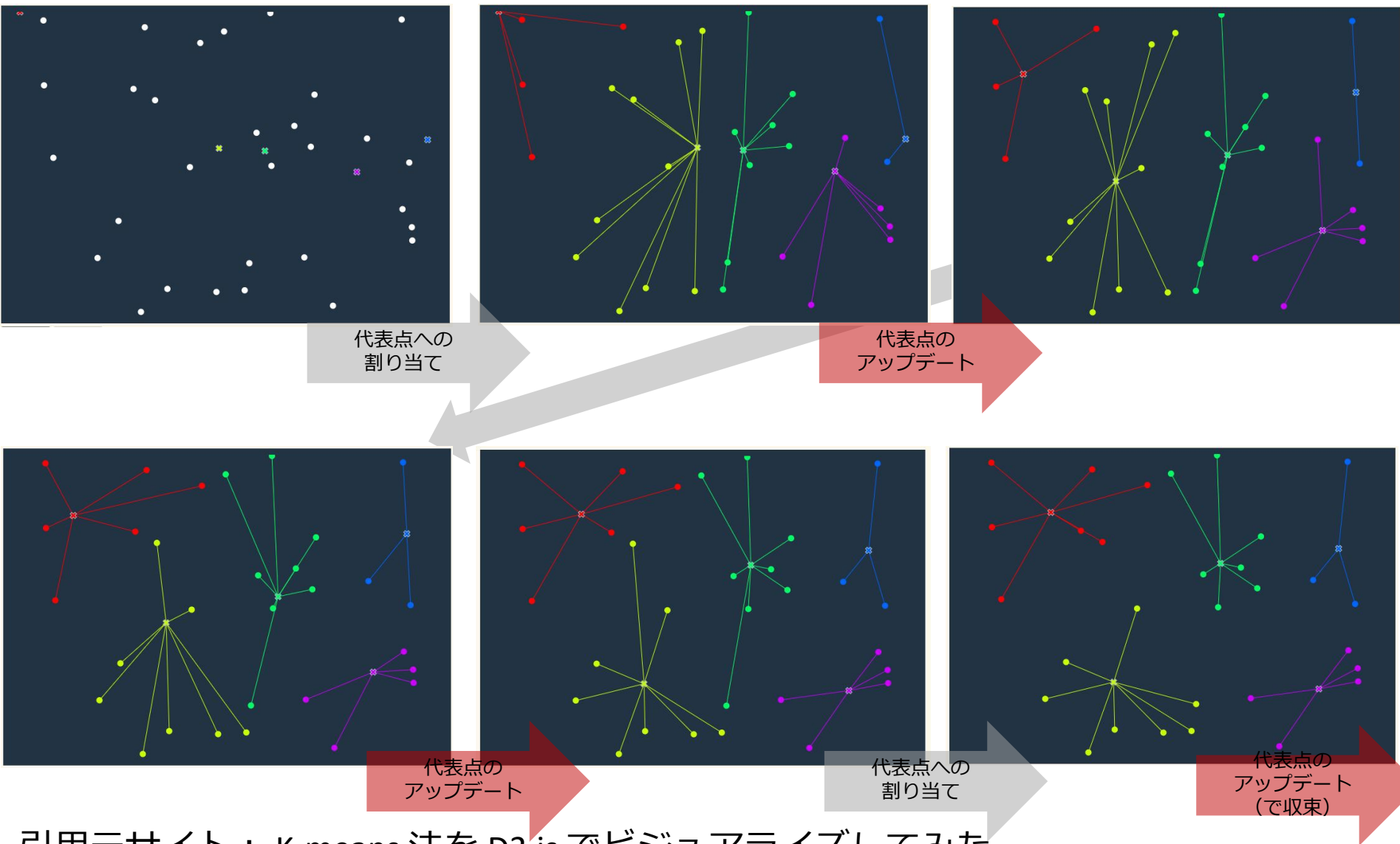


図 11.6

k-means 法によるクラスタリングの過程

2次元での事例



引用元サイト： K-means 法を D3.js でビジュアライズしてみた

<http://tech.nitoyon.com/ja/blog/2013/11/07/k-means/>

演習11-1 k-means法とは？

- k-means 法の説明として最も**不適切**なものを選べ.
 - ① データを最も近いクラスタに帰属させ、その後にクラスタの代表点を更新する.
 - ② クラスタ内のデータとクラスタの代表点の距離の和を減少させる.
 - ③ クラスタの代表点を更新する際にはデータの重心値をとるのであって中央値をとるのではない.
 - ④ k 個の方法を組み合わせることで学習を進行させる.

演習11-2 k-means法

- 2次元平面上に $\{(0,0), (0,1), (0,2), (4,0), (4,1), (4,2)\}$ の6点の点集合がある. これらに対してk-means法を適用しクラスタリングを行え.
- 初期のグループ分けはランダムに行うこと.
- クラスタ数は $k=2$ とせよ.

Contents

□11.1 クラスタリング

□11.2 k-means法

□11.3 混合分布モデル

□11.4 表現学習

11.3.1 確率モデルに基づくクラスタリング

- ❑ k-meansでは境界が確定的なので、クラスタへの帰属度合いなどが議論しにくい。
- ❑ また、データがどのクラスタに属するかの判定が距離のみで判断されるために、クラスタごとにデータ分布の広がり異なるようなデータを適切に分けることができない。

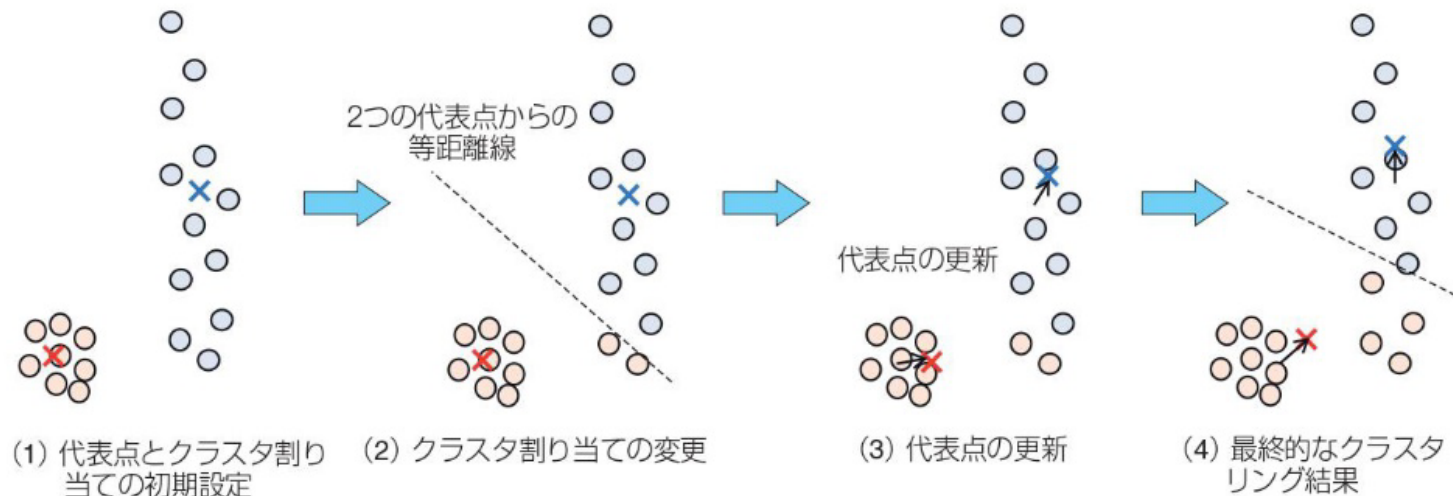


図 11.7 異なる形状の分布からなるデータに k-means 法を適用した場合に起きる典型的な問題

裏でデータが生成される確率を明示的に考える

➡ 混合分布モデル（確率的生成モデル）に基づくアプローチ

11.3.2 混合分布モデルのデータ生成過程

□ 混合分布モデルでは、データが、元々どのようにして生成されたデータであるか、というモデルを考えて、その生成過程をベイズの定理を用いて逆方向に推定することでクラスタリングを行う。

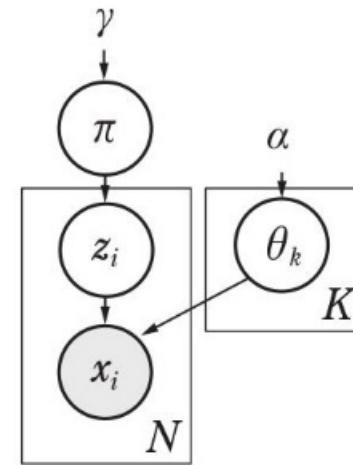


図 11.8

混合分布モデルを表すグラフィカルモデル

$$P(z_i = k)$$

要素分布の
選択確率

$z_i = 1$

$z_i = 2$

$z_i = 3$

$$P(x|\theta_{k=z_1})$$

要素分布

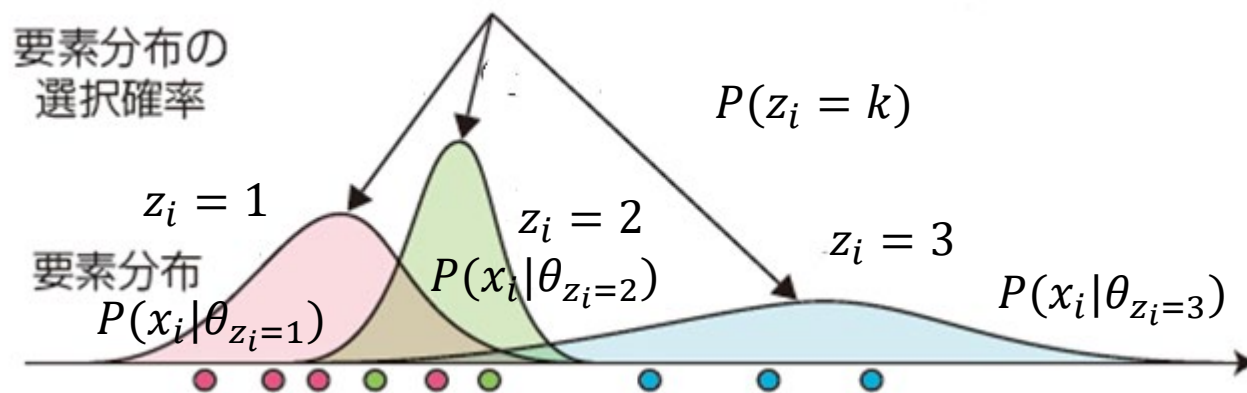
$$P(x|\theta_{k=z_2})$$

$$P(x|\theta_{k=z_3})$$



混合ガウス分布

$$P(x_i) = \sum_{k'} P(z_i = k') P(x_i | \theta_{z_i=k'})$$



• 混合ガウス分布

- 混合分布モデルで要素分布がガウス分布であるもの.
- 各要素分布が**平均パラメータ**と**分散パラメータ**を持つ.
- パラメータ更新がk-means法の重心の更新に相当する.

• EMアルゴリズム

- **最尤推定**にもとづいて混合ガウス分布を学習するためのアルゴリズム

推論方法：EMアルゴリズム

- 混合ガウス分布の学習はEMアルゴリズムを用いることが多い。EMアルゴリズムは平均パラメータについては以下のようなアルゴリズムになる。
- Eステップ
 - ガウス分布の平均値パラメータを固定した上で、全ての観測 x_i に対して、 $P(z_i = k|x_i)$ を計算する。
 - $P(z_i = k|x_i)$ はデータ x_i のクラスタ k への帰属度を与えていると考えられる。
- Mステップ
 - k 番目のガウス分布について全てのデータ x_i を $P(z_i = k|x_i)$ で重みづけて平均をとり、平均値パラメータを更新する。

k-means法はEMアルゴリズムの近似になっている。

他にマルコフ連鎖モンテカルロ法(MCMC)の一種であるギブスサンプリングや変分推論（変分ベイズ）による推論法などがある。



⑪ 須山敦志（著），杉山将（監修）：
ベイズ推論による機械学習入門，講談社，2018。

11.3.5 確率的生成モデルによる発展的なクラスタリング*

□ LDA (Latent Dirichlet Allocation) 潜在ディリクレ配分法

□ Bleiらによって2003年に提案されて以降文章クラスタリングの標準的手法として用いられている。

□ 多項分布の混合モデル。 **トピックモデル**とも呼ばれる。

□ HMM (Hidden Markov Model) 隠れマルコフモデル

□ 音声認識や自然言語処理を始めとしたさまざまな時系列データのクラスタリングや分節化, 認識, 系列ラベリングに用いられる。

□ 混合モデルに時間方向の依存性を導入

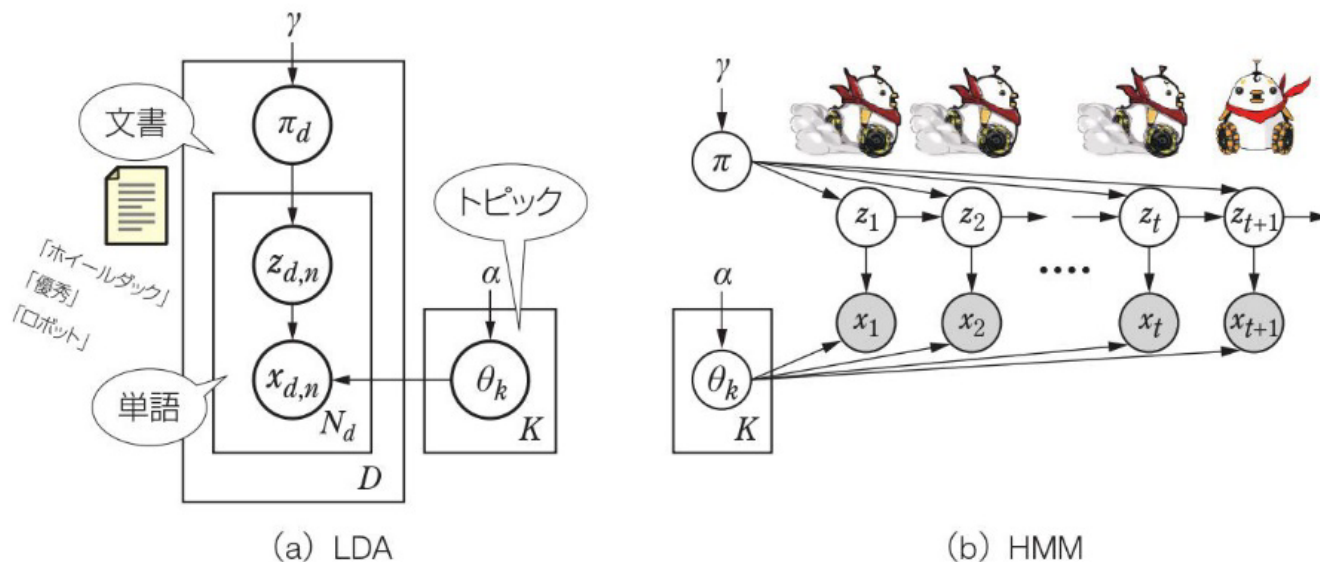


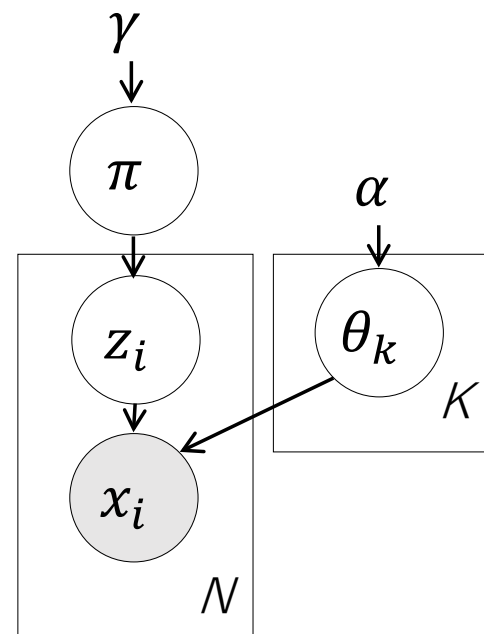
図 11.9

LDA と HMM のグラフィカルモデル

演習11-3 混合分布

□右のグラフィカルモデルは混合ガウス分布（GMM）を表すグラフィカルモデルである。
以下のそれぞれを表す変数はどれか？答えよ。

1. どのガウス分布が選ばれやすいか？
2. 各データがどのガウス分布に当てはめられるか？
3. 各ガウス分布のパラメータ
4. 観測データ



Contents

□11.1 クラスタリング

□11.2 k-means法

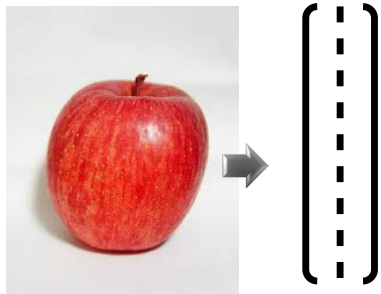
□11.3 混合分布モデル

□11.4 表現学習

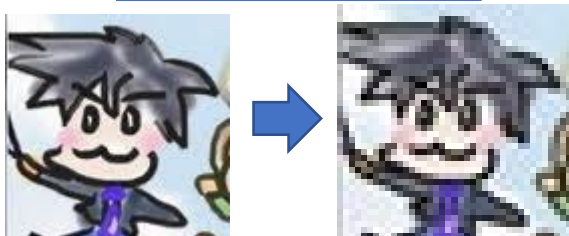
11.4.1 クラスタリングと表現学習

- クラスタリングと並ぶ教師なし学習の手法.
- 表現学習は多くの場合高次元のデータをより低次元のベクトルで表現する. (より高次元に表現することもある)
- 多様体学習などと類似した概念.

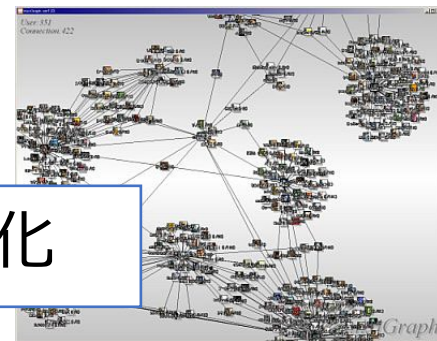
特徴ベクトル抽出



データ圧縮



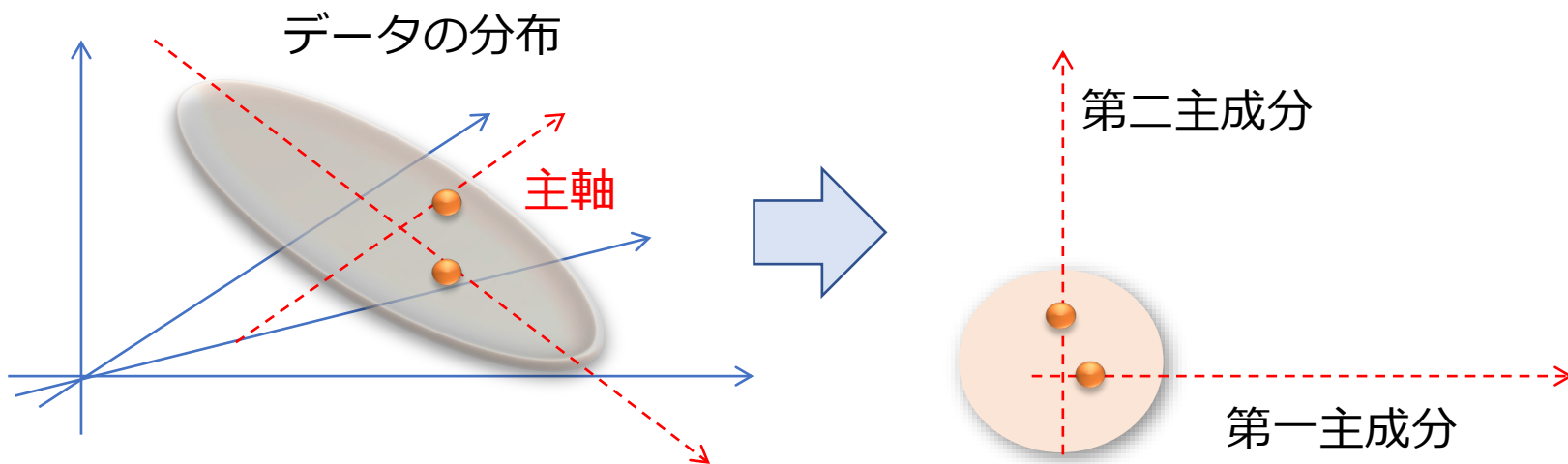
可視化



ソーシャルネットワークグラフ
[twitter mention map](#)

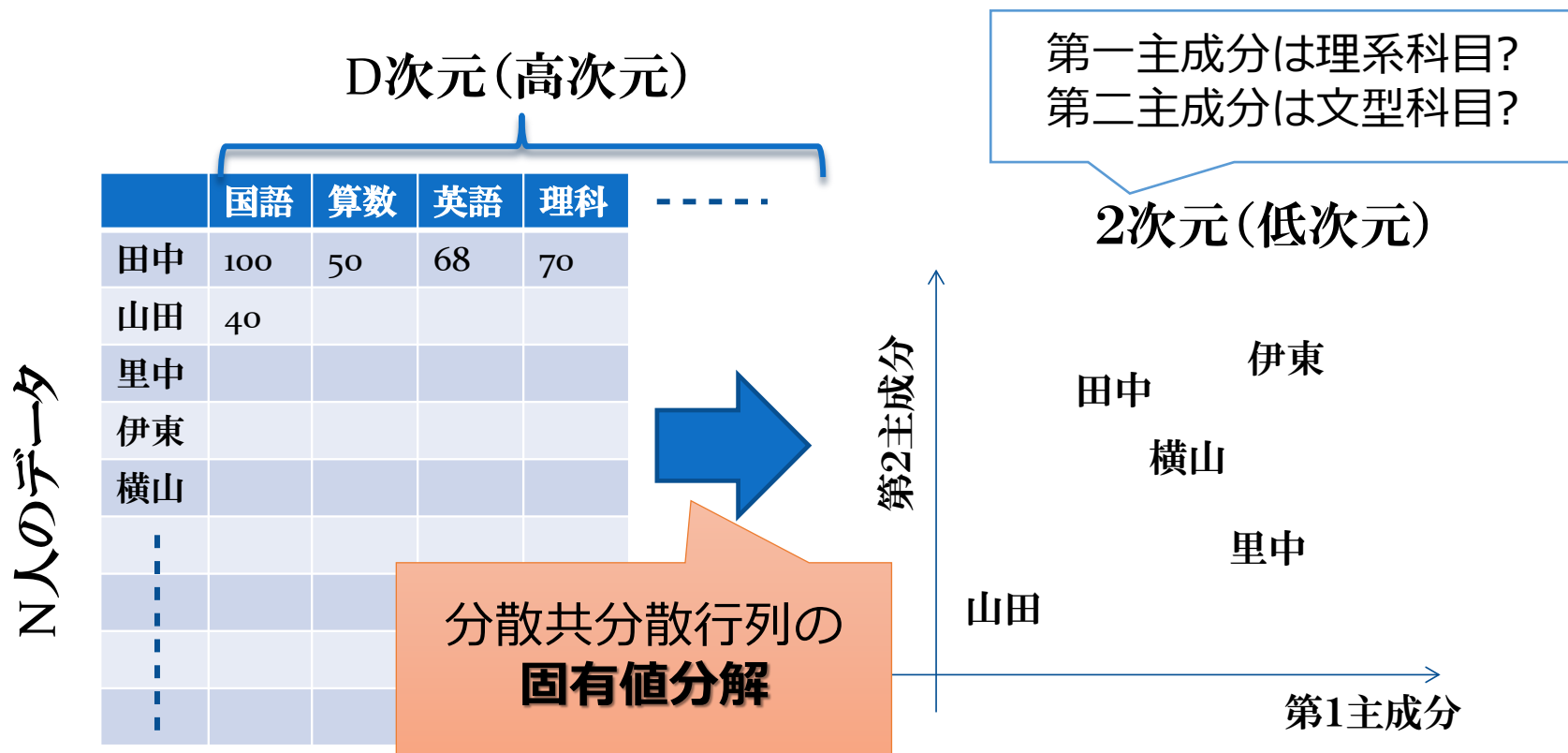
主成分分析

□主成分分析は具体的にはデータが高次元空間上でガウス分布をしていると仮定して、その分布の主軸方向（最も分散の大きい方向）を発見し、それを第1主成分とする。その後、その次に分散の大きい軸をとるというように、順次、軸をとっていくことで、低次元空間を得ていく。



主成分分析の例

- $N = 1000$ 人の学生が $D = 30$ 科目の授業の履修を終えて、それぞれに100点満点の成績を得たとする.
- 30次元のデータを最も上手く表現できるような低次元の表現を得る.



11.4.5 自己符号化器：オートエンコーダ

- ニューラルネットワークを用いた古典的な表現学習手法である**自己符号化器（オートエンコーダ：Auto-Encoder）**。砂時計型のニューラルネットワークによって、入力データ自身を予測する学習を行う（図10.10）。このとき、隠れ層には入力データの情報を保持する特徴ベクトルが得られる。
- さらにこの自己符号化器を多段階に積み上げたのが**積層自己符号化器（Stacked Auto Encoder）**

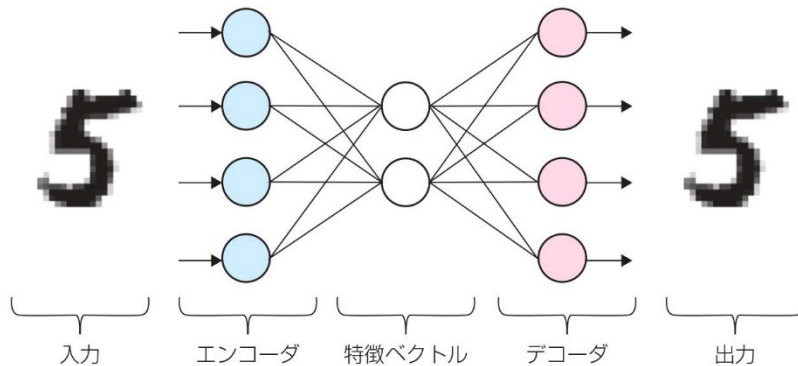


図 11.12 オートエンコーダによる表現学習

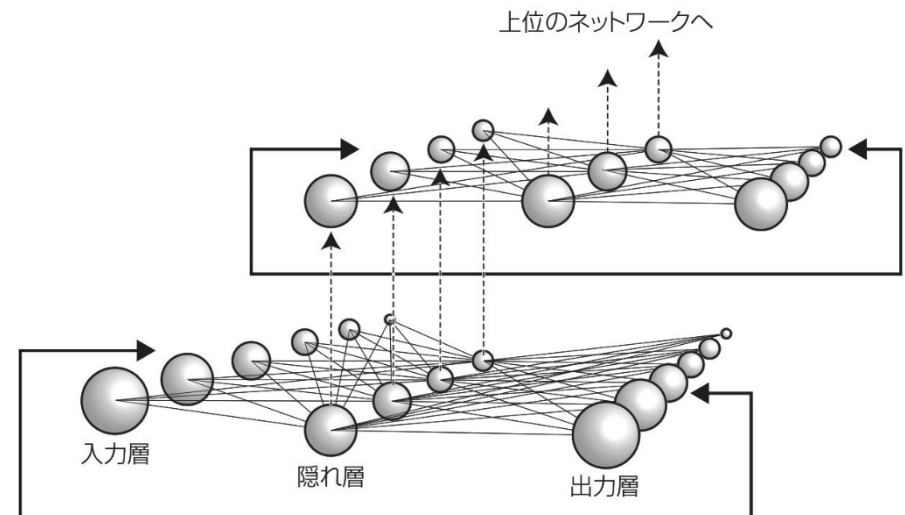
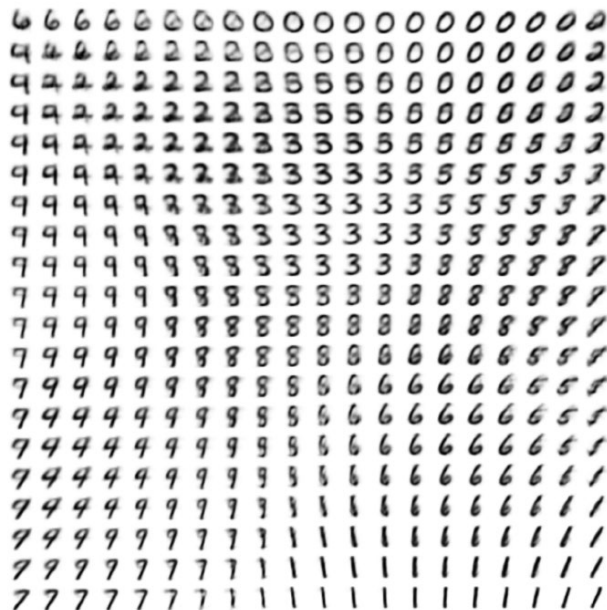


図 11.13 積層オートエンコーダ

11.4.6 変分自己符号化器: VAE

- 変分自己符号化器(VAE: Variational Auto-Encoder) は確率的生成モデルとニューラルネットワークに基づく表現学習の手法である.



(b) Learned MNIST manifold

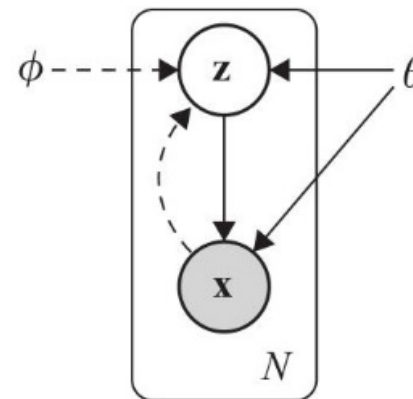


図 11.14

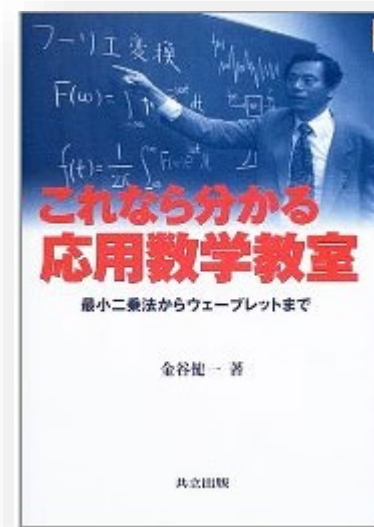
変分オートエンコーダのグラフィカルモデル

Kingma, Diederik P., and Max Welling. "Auto-encoding variational bayes." *arXiv preprint arXiv:1312.6114* (2013).

様々な表現学習（低次元化）手法

- 主成分分析
- 独立成分分析
- カーネル主成分分析
- MDS (多次元尺度法)
- 自己組織化マップ(SOM)
- GPLVM
- 自己符号化器（オートエンコーダ）
- 変分自己符号化器
（変分オートエンコーダ）： VAE

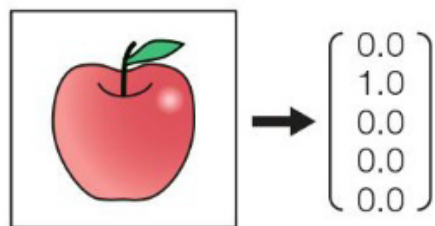
主成分分析を学ぶなら
とりあえず、これなど・・・



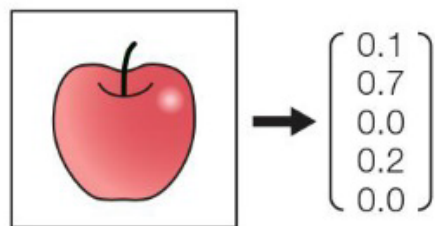
これなら分かる応用数学教室—最小二乗法からウェーブレットまで，金谷 健一

クラスタリングと表現学習の関係

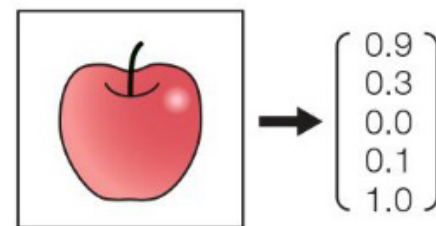
- 確定的であるにせよ確率的であるにせよクラスタリングではある観測データは一つの潜在的な特徴表現, つまり, クラスタに対応する要素分布のパラメータから生成されると考えた.
- 表現学習はいくつかの潜在的な特徴表現に分解して表現しようとする.
- しかしいずれもK次元ベクトルへと情報表現する点では類似.



(a) 確定的クラスタリング



(b) 確率的クラスタリング



(c) 表現学習

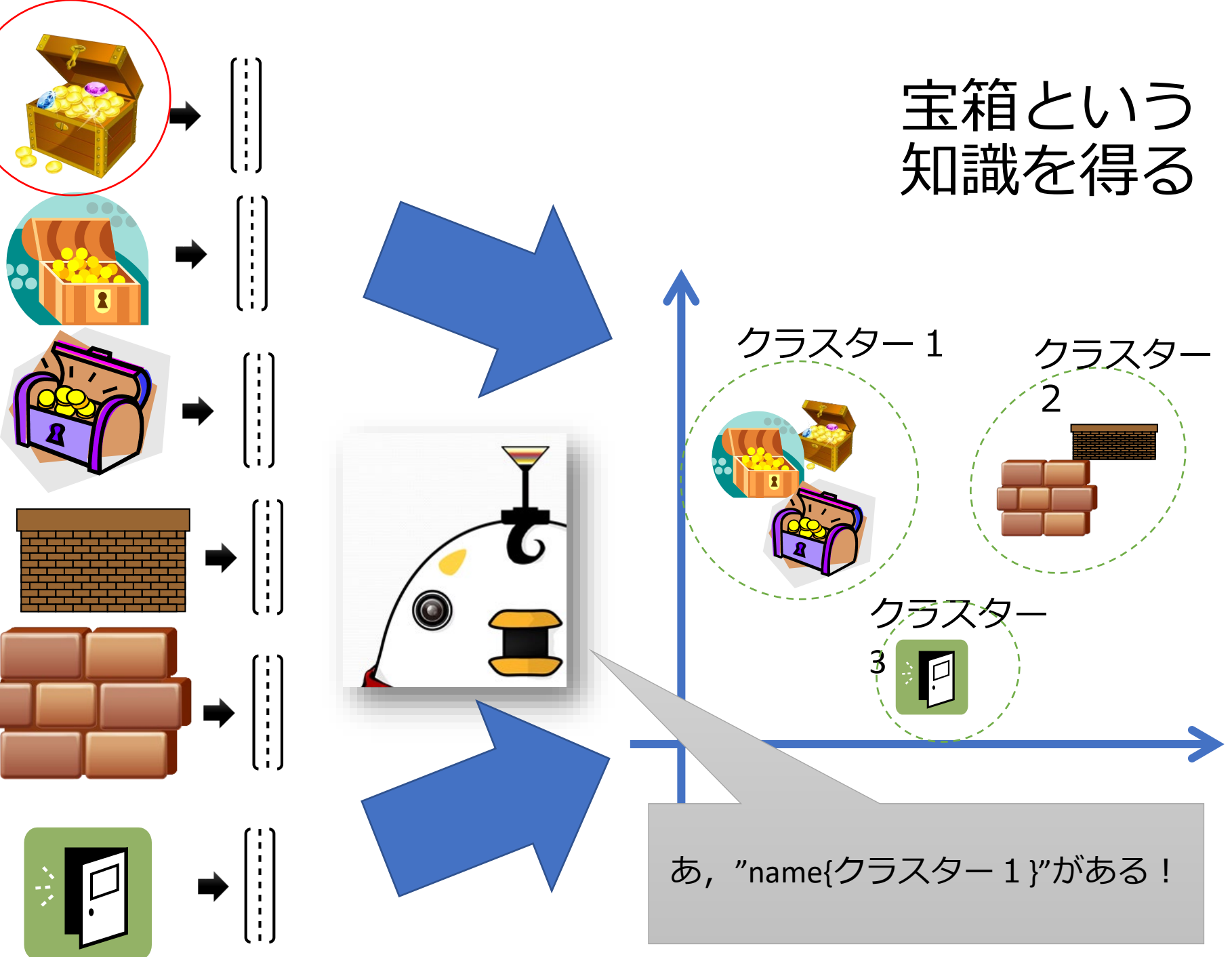
図 11.10

クラスタリングと表現学習の関係

演習11-4 表現学習

- 以下のアルゴリズムの中で表現学習の例として適切でないものを選べ
 1. 主成分分析
 2. 変分オートエンコーダ
 3. GPLVM
 4. 重回帰分析

宝箱という 知識を得る



まとめ

- クラスタリングの基礎について学んだ.
- k-means 法のアルゴリズムを学び, 簡単な数値例を通じてその動作を確認した.
- 混合ガウス分布の生成過程に関して学んだ.
- 混合ガウス分布を用いてクラスタリングを行うアルゴリズムに関して学んだ.
- 表現学習の概要について学び, その代表的な手法である主成分分析, 独立成分分析, カーネル主成分分析, オートエンコーダなどの概要を知った.