

人工知能A

Topic 6: Introduction to clustering and PCA

クラスタリングと主成分分析（の入門）

6

クラスタリング

- 講義の内容
 - 概念学習とクラスタリング
 - K-means法
 - 階層的クラスタリング
 - クラスタリングの応用例
 - 主成分分析
 - 混合（ガウス）分布（概要→最適化の講義）
- 目標：
 - クラスタリングの基本、主にk-means法と凝集法について学ぶ。主成分分析（次元圧縮）について学ぶ

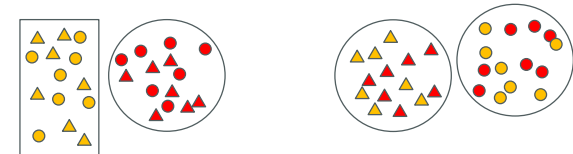
2

3

概念学習とクラスタリング

クラスタリング

- クラスタリング(clustering) とは、データの集合を、その表現（通常は数値のベクトル）の類似性に従って分類すること。クラスタ分析 (cluster analysis) とも言う。
- 分類した一つ一つのまとまりを **クラスタ (cluster)** という。
- 特に正解データや報酬があるわけではなく、データだけから分類するので、教師なし学習の一つ。（どのように分類すべきかは未知）
- 分類したものはまとめて概念として保存すると有効なこともあり、概念学習の一つでもある。
- その分類の観点（あるいは意味）は複数有り、何が適切かは問題に応じて異なる



4

特徴の表現

- あるデータ（対象）を区別するには、その特徴を数値化して、表現する。ただし、特徴は多様であるので、複数の数値（ベクトル）として表現。
- これを**特徴ベクトル (feature vector)** あるいは**特徴量 (feature value)** と言う。
- 各特徴を表す数値（ベクトルの要素）の範囲は、実数全体、整数、正数、 $[0, 1]$ 区間などそれぞれによって決まる。表現可能な特徴量の空間を**特徴空間 (feature space)** と言う。
- たとえば、みかんであれば、形、色、上に緑の点がある、などなど多くの特徴があり、これで（完全ではないが）表現できる。
- クラスタリングは、ベクトル空間（の部分集合）の点を、適当な軸で分割・分類することとも言える。



5

6 基本的なクラスタリング法

クラスタリングの基本 (K-means法)

- K 個の距離的に近いグループに分けること。
- アルゴリズム：データの集合を Ω とする。
 - ここからランダムに K 個の点を取り、それを c_1, \dots, c_K とする。
 - 全ての $d \in \Omega$ について c_i との距離 $\|d - c_i\|^2$ を求め、一番近い c_i のインデックス i を d に付与する。これを $I(d) \in \{1, \dots, K\}$ と表す。
 - すべてのインデックス i について、そのデータの集合の重心（平均）を新たな c_i とする。
 - 2.に戻る。この繰り返しを、インデックスの付け替えが無くなるまで行う。
- 上記は以下の評価値 $J(c_1, \dots, c_K)$ を最小化する $\forall c_j$ と (Ω_i) を求める。

$$J(c_1, \dots, c_K) = \sum_{i=1}^K \sum_{d \in \Omega_i} \|d - c_i\|^2$$

ただし、 $\Omega_i = \{d \in \Omega : I(d) = i\}$ である。

7

K-means法の特徴（利点と欠点）

- 利点
 - 単純で効率がよい。分かりやすい。
 - 多次元空間における外れ値（異常値）を容易に見つけることができる。
- 欠点
 - 複数の相関関係の影響を受ける。また高次元では固まりやすい（分割・分類しにくくなる）。
 - 基本的に距離に基づいているので、超球の形状のみに制限される。
 - 結果は初期値（ c_1, \dots, c_K の取り方）に依存し、安定的ではない。何回か繰り返し良いものを選ぶ（ J を最小化）。

8

階層的クラスタリング（凝集型）

- K-meansとは逆に、個々の点から近いものを集めてクラスタリングする方法。
- 徐々に統合により大きくなるため階層性がある。
 - 凝集型クラスタリング (agglomerative clustering)あるいは凝集型階層クラスタリング (agglomerative hierarchical clustering) という
- アルゴリズム： $\Omega = \{d_1, \dots, d_n\}$ をデータの集合とする
 - 初期値として $\forall i$ に対し、 $c_i = \{d_i\}$ とおく。 c_i はクラスタ。
 - 各クラスタ間で距離 $\text{dist}(c_i, c_j)$ を測り、最小となるペアを統合する。
$$(c_{i_0}, c_{j_0}) = \arg \min_{(c_i, c_j) \in \Omega \times \Omega, i \neq j} \text{dist}(c_i, c_j)$$
 - 上記をクラスタの数が K （あるいは1）になるまで繰り返す。
- 議論： c_i は集合である。その距離の定義の仕方は？
 - Centroid法、ward法などがあり、よく使われる。

9

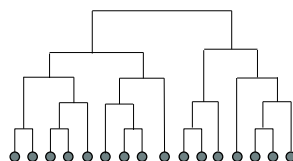
凝集型階層クラスタリングの種類

- 最短距離法あるいは単連結法 (single-link method)
 - 二つの集合のうち、最も近い点同士の距離を集合の距離とする。
 - 同様に最長距離法 (=完全連結法, complete-link method) もある。
- 群平均法 (group-average method)
 - 二つの集合の点のペアの距離の平均値を距離とする。
- 重心法 (Centroid method)
 - それぞれの集合の重心を求めて、それらの距離とする。
- Ward 法 (Ward's method)
 - 集合 C の重心 b と C の各点からの距離の2乗和を $L(C)$ と置く。このとき、 $\text{dist}(c_i, c_j) = L(c_i \cup c_j) - L(c_i) - L(c_j)$ と定義する。
 - $L(C)$ は、 C の全ての点ではなく、サンプルを利用する方法もある。

10

デンドログラム

- 集約すると右下のグラフのように階層的な木構造が生まれる。
 - このグラフをデンドログラム (dendrogram) という。



- 課題 6-1（発展）
この他に潜在的ディリクレ配分法
LDA (Latent Dirichlet allocation)
などもあるので、興味がある学生は調べてみるとよいでしょう
(特に言語解析やトピック (話題) モデルに興味があれば)。

11

12

適用例

例：画像圧縮

- 画像は多数の点でできていて、各点は色 (RGBなど) で表されている
- たとえば、RGBのそれぞれの色の強度を8ビット (0-255) の256段階で表現すると24ビット分の情報が必要。各点の色をこの3次元ベクトルで表現。
- たとえば、ある写真について、上記のベクトル空間に $K=32$ としてK-means法を適用し、その重心を c_1, \dots, c_K をとして、同一クラスに属するものをすべて一色の c_i で描く(c_i にも色が対応する)。
- 使われる色は、32色に。近似している色に変わる。これにより画像情報がかなり下がる。
 - 問題点：3次元空間での色の近さと人間の認知はかならずしも一致しない。
 - 問題点：人間はたとえば顔など自然に注目する箇所があり、感覚的に不自然に感じる。
 - 画像圧縮については、いろいろな研究があり、それらは専門の講義で。

13

例：言語処理

- 簡単な例としてbag-of-words (BOW, bagは多重集合のことでもある)
- 単語 (全単語でもよいが普通は着目する単語。それでも数万以上はある) の出現数をベクトル表示。たとえば、
“全単語でもよいが普通は着目する単語” \Rightarrow
全/単語/でも/よい/が/普通/は/着目/する/単語/
- もし単語、普通、集合、出現、着目、数 だけを考えれば、
[2, 1, 0, 0, 1 0]
これを文や文書ごとに集計し、ベクトル表示。
- 文章を数百から数万次元の空間に配置し、クラスタリングすると、文章を分類できる。
- この他にも、関連する項目として、連続BOW (continuous bag-of-words, N-gram, word2vec等があるが、これらは自然言語処理の講義で述べられるはず)

14

マーケティング (顧客分類)

- $\Omega = \{x_1, \dots, x_n\}$ を顧客の集合とする。ネット販売やポイント (たとえばTポイントなど) などで購入実績が分かる。
- x_i の履歴 (たとえば最近1年とか1月とか) を調査し、以下のベクトルで表現。
 - v_1 : 購入総額、 v_2 : 雑貨の購入額、 v_3 : 食品の購入額、 v_4 : 電気製品の購入額
 - v_5 : 1月の購入額、 v_6 : 2月の購入額、 v_7 : 3月の購入額、 v_8 : 4月の購入額 などなど。力まかせで、とにかくいろいろ調べる。
 - 重みを均一にしたいなら正規化 (例えば0-10の範囲に縮小・拡大) する (たとえば、購入総額は大きいので、その強い影響を消し、バランスをとる。もちろん、解析者の意図もある)。
- この表現を用いて顧客を分類し、クラス内の顧客は同じ傾向があると考える。
 - たとえば、購入の多い月の直前にメールを送る、
 - 類似した購入指向があれば、他の顧客の購入活動から購入推薦を行う、など

15

16 数学の復習

主成分分析 (PCA) の準備として

数学の復習(内積)

- n 次元ベクトル空間のベクトル $\mathbf{a} = (a_1, \dots, a_n)$, $\mathbf{b} = (b_1, \dots, b_n)$ の内積は、

$$\mathbf{a} \cdot \mathbf{b} = a_1 b_1 + \dots + a_n b_n$$

と定義。これは内積の公式から、

$$\mathbf{a} \cdot \mathbf{b} = |\mathbf{a}| |\mathbf{b}| \cos \theta$$

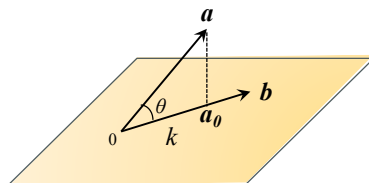
\mathbf{a} の \mathbf{b} への射影 (垂線を下ろしたときの) の長さ (あるいは \mathbf{a} を点と考えて、その \mathbf{b} へ垂線を下ろした点 \mathbf{a}_0 の原点からの長さ) を k とすると、 $k = |\mathbf{a}| \cos \theta$ 。したがって、

$$k = \frac{\mathbf{a} \cdot \mathbf{b}}{|\mathbf{b}|}$$

である。したがって、

$$\mathbf{a}_0 = k \frac{\mathbf{b}}{|\mathbf{b}|} = \frac{\mathbf{a} \cdot \mathbf{b}}{|\mathbf{b}|^2} \mathbf{b}$$

- 特に \mathbf{b} が単位ベクトルなら、 $k = \mathbf{a} \cdot \mathbf{b}$ 、 $\mathbf{a}_0 = k \mathbf{b}$ である。



17

数学の復習 (不偏推定量)

- n 個の観測データ X_1, \dots, X_n から、母集団の平均と分散を推定する不偏推定量はそれぞれ、

$$\bar{X} = \frac{X_1 + \dots + X_n}{n}$$

と、

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

である。初めの係数を $1/n$ とすると、 $1/n$ だけ期待値が小さくなり、ややずれる (勿論 n が大きければ差は小さい)。

- n 個の観測データ X_1, \dots, X_n と Y_1, \dots, Y_n から共分散を推定する不偏推定量は、

$$s_{XY} = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})$$

である。

18

数学の復習 (ラグランジュの未定乗数法)

- ラグランジュの未定乗数法 (Lagrange multiplier method)

- $f(x_1, \dots, x_n)$, $g(x_1, \dots, x_n)$ は関数。 $\mathbf{x} = (x_1, \dots, x_n)$ とおく。条件 $g(x_1, \dots, x_n) = 0$ のもと、 $f(x_1, \dots, x_n)$ を最大化 (最小化) する解 $\mathbf{c} = (c_1, \dots, c_n)$ は、

$$L(x_1, \dots, x_n, \lambda) = f(x_1, \dots, x_n) - \lambda g(x_1, \dots, x_n) \text{ としたとき、}$$

$$\frac{\partial L(\mathbf{x}, \lambda)}{\partial x_1} = \dots = \frac{\partial L(\mathbf{x}, \lambda)}{\partial x_n} = \frac{\partial L(\mathbf{x}, \lambda)}{\partial \lambda} = 0 \quad (1)$$

の解である (つまり L の停留点)。 λ をラグランジュ乗数という

- Maximize (minimize) $f(x_1, \dots, x_n)$
subject to $g(x_1, \dots, x_n) = 0$
- $f(\mathbf{x})$ を (あるいは $g(\mathbf{x}) = 0$ なので $L(\mathbf{x}, \lambda)$ を) を最大化する \mathbf{c} を求めればよい。

19

数学の復習 (ラグランジュの未定乗数法)

- 証明 (幾何学的な概要。「最適化理論」等の講義で詳しくやると思う)

- $\frac{\partial L(\mathbf{x}, \lambda)}{\partial \lambda} = 0$ は、 $g(x_1, \dots, x_n) = 0$ を意味する。
- $\lambda \neq 0$ として、残りの式は $\frac{\partial f(\mathbf{x})}{\partial x_i} = \lambda \frac{\partial g(\mathbf{x})}{\partial x_i}$ ($i=1, \dots, n$) となる。
- これをまとめて書くと

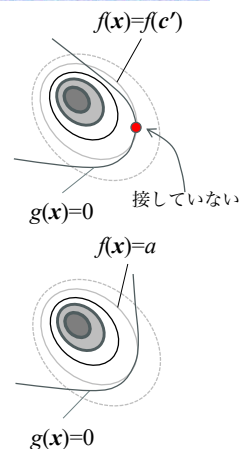
$$\begin{pmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_n} \end{pmatrix} = \lambda \begin{pmatrix} \frac{\partial g(\mathbf{x})}{\partial x_1} \\ \vdots \\ \frac{\partial g(\mathbf{x})}{\partial x_n} \end{pmatrix}$$

$\lambda \neq 0$ なので、これらのベクトルは平行。つまり $f(\mathbf{x})$, $g(\mathbf{x})$ を偏微分したベクトルは勾配 (つまり曲線の法線ベクトル (前出: Topic2, スライド113で) であるが、これが同じ向きであることを示している。

20

数学の復習（ラグランジュの未定乗数法）

4. これは、 $f(x_1, \dots, x_n) = \exists a$ と $g(x_1, \dots, x_n) = 0$ が接していることを意味する。
5. 他方、勾配が同じではない点 $c' = (c'_1, \dots, c'_n)$ があり、それが元の制約を満たす解（最大値）とする。
6. その点の近傍では、 $f(x) = f(c')$ と $g(x) = 0$ は（接していないので）交差し、実数の連続性から、 $f(x) > f(c')$ かつ $g(x) = 0$ となる点が存在する。これは、 c' が制約を満たす解であることと矛盾。
（上の図のようなことはない）
7. 従って、式(1)を満たす解で、 f の値が最大となるもの $c = (c_1, \dots, c_n)$ を選ばばよい。



21

数学の復習（ラグランジュの未定乗数法）

- ラグランジュの未定乗数法のやや一般化
- $f(x_1, \dots, x_n), g_k(x_1, \dots, x_n)$ は関数 ($1 \leq k \leq m$)。 $x = (x_1, \dots, x_n)$ とおく。条件 $g_k(x_1, \dots, x_n) = 0$ のもと、 $f(x_1, \dots, x_n)$ を最大化（最小化）する解 $c = (c_1, \dots, c_n)$ は、

$$L(x_1, \dots, x_n, \lambda_1, \dots, \lambda_m) = f(x_1, \dots, x_n) - \sum_{k=1}^m \lambda_k g_k(x_1, \dots, x_n)$$

としたとき、

$$\frac{\partial L(x, \lambda)}{\partial x_1} = \dots = \frac{\partial L(x, \lambda)}{\partial x_n} = \frac{\partial L(x, \lambda)}{\partial \lambda_1} = \dots = \frac{\partial L(x, \lambda)}{\partial \lambda_m} = 0 \quad (2)$$

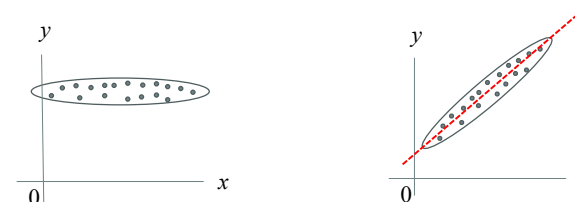
の解（つまり停留点）でもある。

22

23 主成分分析 (PCA)

主成分分析

- これまでの例では次元はとても高そうである。一方で、使われない変数もありそう。
- 仮に使われている変数でも、データを区分する（クラスタリングする）のに重要（不要）な変数はありそう。
- イメージ：たとえば、左下では、 x 軸の値の方が重要そう。右の方は、 x, y 軸共に重要そうだが、新しく軸を赤線（点線）のようにとれば、一つの軸で区別がある程度できそうである。



24

主成分分析

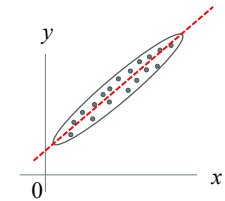
- n 変数 (n 次元)のデータが N 個ある。これを少ない次元 m なるべくデータの区分が可能ないように表現したい。ただし、 $m \leq n$ である。
- このように少ない変数で表現し直すことを低次元化という。
- (必要であれば簡単な一次変換を施して) なるべく有効な変数を残し、重要で無さそうなものは消して、簡単に (低次元で) 表したい。
- 通常は新しい直交座標を (w_1, \dots, w_m) として、重要なものから順に第1成分、第2成分...と呼ぶ。この順に区分けが明確になるよう決める。
- クラスタリングとは別の意味の区分になるので、教師無し学習とも考えられる。
- 初めに主成分分析で次元を落としてからクラスタリングをする、逆にクラスタリングののち各クラスタの特徴を探るためにそのデータのみに着目して主成分分析を行う、両者平行して行い比較するなど、組み合わせで用いることもある。

25

主成分分析 (principle component analysis, PCA)

- 基本方針
 - データを「ある軸」に射影したときに、なるべく区分けができるように分散が大きい軸を「ある軸」 w_1 としてを選択する (第1主成分)。
 - 次に、第1主成分と直交する軸で、同様にその軸に射影したときに分散値が大きくなるものを第2主成分 w_2 として選択する。
 - これまで選んだ第 $k-1$ 主成分まですべてと直交し、射影したときに分散値が大きくなる軸を第 k 主成分 w_k として選択する。
 - 上記を第 m 主成分 w_m を得るまで繰り返す。
 - なお軸 $w_i = (w_{i1}, w_{i2}, \dots, w_{im})$ の大きさは1、つまり

$$\|w_i\|^2 = w_i \cdot w_i^T = 1$$
 としておく (一般性を失わない)。



26

主成分の求め方(1)

- 各データ x_i は、 n 個の値のベクトルで表現でき、このようなデータが N 個あるとする。これを縦に並べて以下のように表す。

$$X = \begin{pmatrix} x_1 \\ \vdots \\ x_N \end{pmatrix} = \begin{pmatrix} x_{11} & \cdots & x_{1n} \\ \vdots & \ddots & \vdots \\ x_{N1} & \cdots & x_{Nn} \end{pmatrix}$$

- まず求めたい第1主成分を $w_1 = (w_{11}, w_{12}, \dots, w_{1n})$ と表す。また、ベクトルの各要素 x_{1j}, \dots, x_{Nj} の値の標本平均 μ_j を求め、平均値が0となるように $\bar{x}_{ij} \leftarrow x_{ij} - \mu_j$ と変換し、以下のようにおく。

$$\bar{X} = \begin{pmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_N \end{pmatrix} = \begin{pmatrix} \bar{x}_{11} & \cdots & \bar{x}_{1n} \\ \vdots & \ddots & \vdots \\ \bar{x}_{N1} & \cdots & \bar{x}_{Nn} \end{pmatrix}$$

- 点 \bar{x}_i の w_1 への射影は、内積で表せるので、

$$p_i = \bar{x}_i \cdot w_1^T \quad \left(= \sum_{k=1}^n \bar{x}_{ik} w_{1k} \right)$$

28

主成分の求め方(2)

- これをまとめれば、

$$\begin{pmatrix} p_1 \\ \vdots \\ p_N \end{pmatrix} = \bar{X} \cdot w_1^T$$

- データを射影した p_i の不偏標本分散 σ^2 を求めると標準化してあるので

$$\sigma^2 = \frac{1}{N-1} \sum_{i=1}^N p_i^2 = w_1 \cdot \left(\frac{1}{N-1} \bar{X}^T \bar{X} \right) \cdot w_1^T$$

(標本分散を使っても良いこととする、そのときは $1/N$ をつかう)。

- ここで、 $\frac{1}{N-1} \bar{X}^T \bar{X}$ の ij 成分を $\bar{\sigma}_{ij}$ とおけば、これは

$$\bar{\sigma}_{ij} = \frac{1}{N-1} \sum_{k=1}^N \bar{x}_{ki} \bar{x}_{kj} = \frac{1}{N-1} \sum_{k=1}^N (x_{ki} - \mu_i)(x_{kj} - \mu_j)$$

となり、共分散行列の不偏推定 (統計の講義) となる。この行列を S とおく。

29

主成分の求め方(3)

- $\|\mathbf{w}_1\| - 1 = 0$ のもと、 $\sigma^2 = \mathbf{w}_1^T \mathbf{S} \mathbf{w}_1$ を最大化する \mathbf{w}_1 を求めればよい。
- ラグランジュ未定乗数法を適用する。
$$L(\mathbf{w}_1, \lambda) = \mathbf{w}_1^T \mathbf{S} \mathbf{w}_1 - \lambda(\mathbf{w}_1^T \mathbf{w}_1 - 1)$$
- 上記より、
$$L(\mathbf{w}_1, \lambda) = \left(\sum_{i=1}^n w_{1i} \sigma_{i1}, \dots, \sum_{i=1}^n w_{1i} \sigma_{in} \right) \cdot \mathbf{w}_1^T - \lambda(\mathbf{w}_1^T \mathbf{w}_1 - 1)$$
$$= \sum_{j=1}^n \sum_{i=1}^n w_{1i} \sigma_{ij} w_{1j} - \lambda(\mathbf{w}_1^T \mathbf{w}_1 - 1)$$
- $\frac{\partial L(\mathbf{w}_1, \lambda)}{\partial \mathbf{w}_1} = 0$ を考えると、
$$\frac{\partial L(\mathbf{w}_1, \lambda)}{\partial \mathbf{w}_1} = 2(\mathbf{S} - \lambda \mathbf{I}) \mathbf{w}_1^T = 0$$
つまり、 $\mathbf{S} \mathbf{w}_1^T = \lambda \mathbf{w}_1^T$ となり、 \mathbf{S} の固有値を求めればよい。

30

主成分の求め方(4)とPCAの意味

- 固有値は、 $\det(\mathbf{S} - \lambda \mathbf{I}) = 0$ を解けばよいが、一般に固有値は複数個ある。どれを選ぶか？これまでは第1主成分を考えてきたが、その他についても同様で、 \mathbf{S} の固有値問題に帰着できる。
- \mathbf{S} の固有値を $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$ とおき、対応する固有値ベクトルを $\omega_1, \dots, \omega_m$ とする。
$$\sigma^2 = \mathbf{w}_1^T \mathbf{S} \mathbf{w}_1 = \sigma^2 = \mathbf{w}_1^T \lambda_k \mathbf{w}_1 = \lambda_k$$
- 以上から第1主成分は $\mathbf{w}_1 = \omega_1$ 、第2主成分は $\mathbf{w}_2 = \omega_2$ などとなる。
- 主成分分析の固有値ベクトルからの意味
 - ω_1 に射影したとき分散が最大となり、データを広く区分けできる。
 - 一方、それと直交する軸のデータは失われる。失われたデータの中でデータを最も広く区分するの軸が ω_2 （以下同様）。

31

PCAの利点と欠点

- 利点
 - 相関のある変数の除去。
 - 次元を下げるができる。これによる単純化と高速化の利益は大きい
- 欠点
 - 線形性に強く依存している。これに合わないものは適用できない。

32

応用例：画像圧縮

- 画像は多数の点でできていて、各点は色(RGBなど)で表されている。
- たとえば、 $1024 \times 640 = 655360$ の点(pixel)でできているが、これを 64×64 に分割し、 16×10 の画像が4096集まっていると考える。各点の色はRGBで表現(3次元)。分割した画像を $16 \times 10 \times 3 = 480$ 次元の点で表現。これが4096ある。
- この4096の480次元データをPCAで次元圧縮する。
- $\mathbf{x}_i = (x_{i1}, \dots, x_{i480})$ (480次元) $\Rightarrow \mathbf{y}_i = (y_{i1}, \dots, y_{im})$ (m 次元)
- $\mathbf{w}_1, \dots, \mathbf{w}_m$ と合わせて、 $y_{i1} \mathbf{w}_1 + \dots + y_{im} \mathbf{w}_m$ として元の次元に戻す(色に戻すということ。もちろん結果は前とは違う)
 - たくさんの色が使われている領域では、大きく劣化するかもしれない。
 - 他にも方法がある。画像圧縮については、いろいろな研究があり、それらは専門の講義で。

33

課題 6-1

- 以下の手順に従って、主成分分析に関するプログラムを書き、実験を行いなさい
 - 正規分布 $N(0, 4)$ に従う乱数を適当な数 M 個 (M は50~100ぐらい) 生成し、それを $a_i (i=1, \dots, M)$ とおく。
 - $N(3, 2)$ に従う乱数を M 個生成し、これを b_i とおく。
 - 0~4の範囲で一様に発生する乱数を M 個用意し、これを c_i とおく。
- 1. $x_i=(a_i, b_i, c_i)$ として M 個のデータを作る。このデータから不偏共分散行列を求めよ。
- 2. この行列の固有値を求めよ。
- 3. PCAの第1主成分、第2主成分をもとめ、その結果について考察せよ。
- なお、正規分布の平均値を μ 分散を σ^2 とすると、 $N(\mu, \sigma^2)$ と表します。

34

課題 6-2

- 以下の手順に従って、K-means方でクラスタリングするプログラムを書き、実験を行いなさい
 - 0~5の範囲で一様に発生する乱数を数 M 個 (M は50~100ぐらい) 生成し、それを $a_i (i=1, \dots, M)$ とおく。
 - 正規分布 $N(7, 4)$ と $N(-3, 3)$ に従う乱数を同数、合わせて M 個生成し、ランダムに並べ替える。これを b_i とおく。
 - $N(5, 5)$ と $N(-1, 4)$ に従う乱数を同数、合わせて M 個生成し、ランダムに並べ替える。これを c_i とおく。
- 1. $x_i=(a_i, b_i, c_i)$ として M 個のデータを作る。このデータを K 個のクラスに分けてみよ。 $K=2, 3, 4, 5 \dots$ として試してその重心をもとめてみなさい。
- 2. K はいくつが良さそうか、考えなさい。

35

36 混合分布

ここはやや難しいので概略だけ。たぶん3年の後期か大学院で。
人工知能の概論と言うより他の応用に結びつけて。

準備 (EMアルゴリズム)

- Expectation maximization algorithm : 基本的には隠れ変数 (値が未知の変数) とある統計量について、隠れ変数の推定と統計量の推定 (最尤推定) を交互に行い、値を徐々に更新する方法。
- ある統計量 θ (たとえば平均値とか確率とか、何か知りたいもの) を最尤推定する。最尤推定は、観測 $X=(x_1, \dots, x_n)$ に対して、その観測が最大となる確率 $P(X|\theta)$ が最大となる $\hat{\theta}$ を求めれば (つまり微分して0となるときの値を調べる) それが θ の推定値であるとするもの。
- ただし、観測できない変数もあり (隠れ変数という)、これを z_1, \dots, z_k とする。
- $\hat{\theta}$ を適当に決め、そこから隠れ変数 z_1, \dots, z_k を推定する。次に決めた z_1, \dots, z_k に基づいて θ の推定値 $\hat{\theta}$ を求める。その $\hat{\theta}$ を利用して、隠れ変数 z_1, \dots, z_k を推定。これを収束条件が満たすまで繰り返す。

37

EMアルゴリズム (1)

- 少し確率・統計の復習。
- 観測できる変数 \mathbf{x} と観測できない変数 \mathbf{z} を仮定し、確率 $p(\mathbf{x}, \mathbf{z})$ を考える。 \mathbf{x} の周辺分布を求める。

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}) = \sum_{\mathbf{z}} p(\mathbf{z})p(\mathbf{x}|\mathbf{z}) (= \sum_{\mathbf{z}} p(\mathbf{x})p(\mathbf{z}|\mathbf{x}))$$

- これに推定したいパラメータ θ を表現に加えて、 $p(\mathbf{x}|\theta)$ と置換え。これは推定したいパラメータを θ としたときの \mathbf{x} が観測できる確率。

$$p(\mathbf{x}|\theta) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}|\theta) = \sum_{\mathbf{z}} p(\mathbf{z}|\theta)p(\mathbf{x}|\mathbf{z}, \theta)$$

- なお各確率の値は θ に依存するものとする。

38

EMアルゴリズム (2)

- EMアルゴリズムを形式的に。
 1. $\theta(0)$ を初期値として選ぶ。
 2. $E: p(\mathbf{z}|\mathbf{x}, \theta(t))$ を計算する。($E = \text{expectation}$)
 3. $M: J = \sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}, \theta(t)) \ln p(\mathbf{x}, \mathbf{z}|\tilde{\theta})$ を最大化する $\tilde{\theta}$ を求めこれを $\theta(t+1)$ とする。($M = \text{maximization}$)
 4. EとMを収束するまで繰り返す。
- なお、 θ の関数 $Q(\theta; \theta(t)) = \sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}, \theta(t)) \ln p(\mathbf{x}, \mathbf{z}|\theta)$ において、 $\tilde{\theta} = \arg \max_{\theta} Q(\theta; \theta(t))$ と表現することもある。
- M の式で自然対数が使われているが、最尤推定ではよく使われる。同時確率は、それらの積となるので対数をとると和となり簡易化
- M の式で自然対数値の期待値が使われているが、これは \mathbf{z} を直接観測できないため、 E で求めた \mathbf{z} の確率を利用して求めている。

39

EMアルゴリズムの適用例

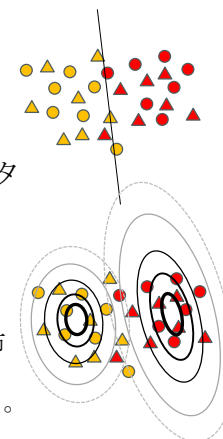
- K-means法
 - $\Omega = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, $\theta = (c_1, \dots, c_K)$. 隠れ変数は $\gamma_i = (\gamma_{i1}, \dots, \gamma_{iK})$, ただし

$$\gamma_{ik} = \begin{cases} 1 & \text{if } \mathbf{x}_i \text{ が } c_k \text{ に一番近いとき} \\ 0 & \text{その他} \end{cases}$$
- 当然隠れ変数は分からない（どちらかと言えば、これを求めたい）
- 1. ランダムに $\theta(0) = (c_1(0), \dots, c_K(0))$ を選ぶ。
- 2. $E: \theta(t)$ を使って $\gamma_{ik}(t)$ を求める。これが隠れ変数の推定値。具体的には、 \mathbf{x}_i に一番近い $c_k(t)$ を求めることに相当。
- 3. $M: c_k(t+1)$ を求める。 $Q(c_k; c_k(t)) = \sum_{i=1}^n \gamma_{ik}(t) \|\mathbf{x}_i - c_k\|^2$ が最小となる c_k を $c_k(t+1)$ とする。なおここでは最小値だが、これが、クラスタに属する確率を最大化するものと定義している。具体的には、 $Q(c_k; c_k(t))$ を c_k で微分して0と置けば新しい $c_k(t+1)$ が平均値（重心）として求められる。

40

混合ガウス分布（混合正規分布）

- これまでの方法の欠点
 - 実データでは境界付近が綺麗に分かれるとは限らない。これを強制的に一方に分けていた。
 - どこのクラスタに属するかを確率的に表現し、境界の部分は曖昧性をもって表現する。
- このためにデータありきではなく、データが生成された過程を考える。
 - そもそも異なるクラスタのデータは、異なるメカニズム（ここでは確率分布）にしたがって生成されていると仮定。
 - クラスタが K 個あるとすると、 K 個の異なる分布 $f_k(\mathbf{x})$ が確率 α_k に従って選択され($1 \leq k \leq K$)、それに基づいて m 次元データ \mathbf{x} が生成されると考える。



41

混合ガウス分布（混合正規分布）

- ここで分布 $f_k(\mathbf{x})$ が正規分布 $N(\boldsymbol{\mu}_k, \Sigma_k)$ に従うと仮定する。
 - もちろん、平均 $\boldsymbol{\mu}_k$ と（分散）共分散行列分散 Σ_k は不明なので、これを知る必要がある。
 - ここでK-means法を思い出すと、クラスタの中心 \mathbf{c}_i に相当するのはクラスタの平均 $\boldsymbol{\mu}_k$ に相当する。（K-means法との関係）
- 考え方（最適化の講義で主に扱うので簡単に）
 - データの集合を $\Omega = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ とする。 \mathbf{x}_i は、 f_1, \dots, f_K のどれかに従うので、最も尤もらしいものを対応づける→最尤推定法 (EM)
 - \mathbf{x}_i が分布 f_j に従って生成される確率は、 $P(\mathbf{x}_i, f_j) = P(\mathbf{x}_i | f_j)P(f_j) = P(\mathbf{x}_i | f_j)\alpha_j$ である。 f_j は正規分布を仮定しているので、 $P(\mathbf{x}_i | f_j)$ は正規分布にしたがった確率となる。つまり $f_j(\mathbf{x}_i)$ 。

$$P(\mathbf{x}_i) = \prod_{j=1}^K \alpha_j \cdot f_j(\mathbf{x}_i)$$

42

混合ガウス分布（混合正規分布）

- 一方ベイズの定理から、

$$P(f_j | \mathbf{x}_i) = \frac{P(\mathbf{x}_i, f_j)}{\sum_{k=1}^K P(f_k)P(\mathbf{x}_i | f_k)} = \frac{\alpha_j f_j(\mathbf{x}_i)}{\sum_{k=1}^K \alpha_k f_k(\mathbf{x}_i)}$$
 これを γ_{ij} とおく。ただし

$$f_k(\mathbf{x}_i) = \frac{1}{(2\pi)^{\frac{m}{2}} \sqrt{|\Sigma_k|}} \exp\left(-\frac{1}{2}(\mathbf{x}_i - \boldsymbol{\mu}_k) \Sigma_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T\right) \quad (\in [0, 1])$$
- ここから $\alpha_k, \boldsymbol{\mu}_k, \Sigma_k$ を求める。
- しかし \mathbf{x}_i がどの正規分布に従ったものかは分からない。これを隠れ変数とする。つまり、
- $\mathbf{z}_i = (z_{i1}, \dots, z_{iK})$ であり、ただ一つの要素のみ1で、残りは0。もちろん分からない。

44

混合ガウス分布（混合正規分布）

- この先は最適化の講義で詳しくやと思うが、

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{i=1}^n \gamma_{ik} \cdot \mathbf{x}_i$$

ただし、 $N_k = \sum_{i=1}^n \gamma_{ik}$

- $\Sigma_k = \frac{1}{N_k} \sum_{i=1}^n \gamma_{ik} \cdot (\mathbf{x}_i - \boldsymbol{\mu}_k)^2$
- $\alpha_k = \frac{N_k}{N}$
- となり、EMアルゴリズムで解く
- 混合分布は最適化問題やパターン認識でよく使われるので、そちらの講義で詳しくあるはず（大学院かも）。

45