

パターン認識と学習

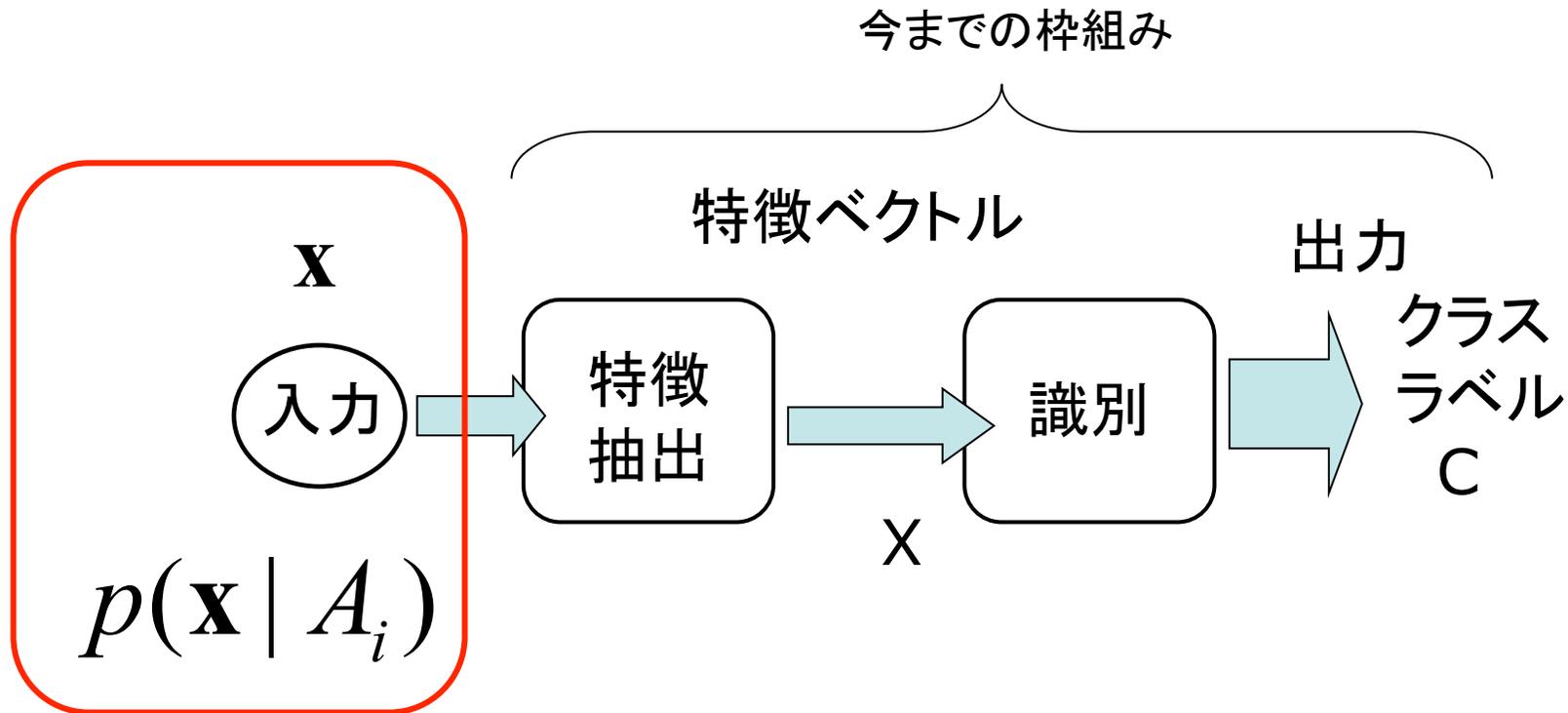
岡山大学大学院 竹内孔一

本日の内容

- 識別部の設計

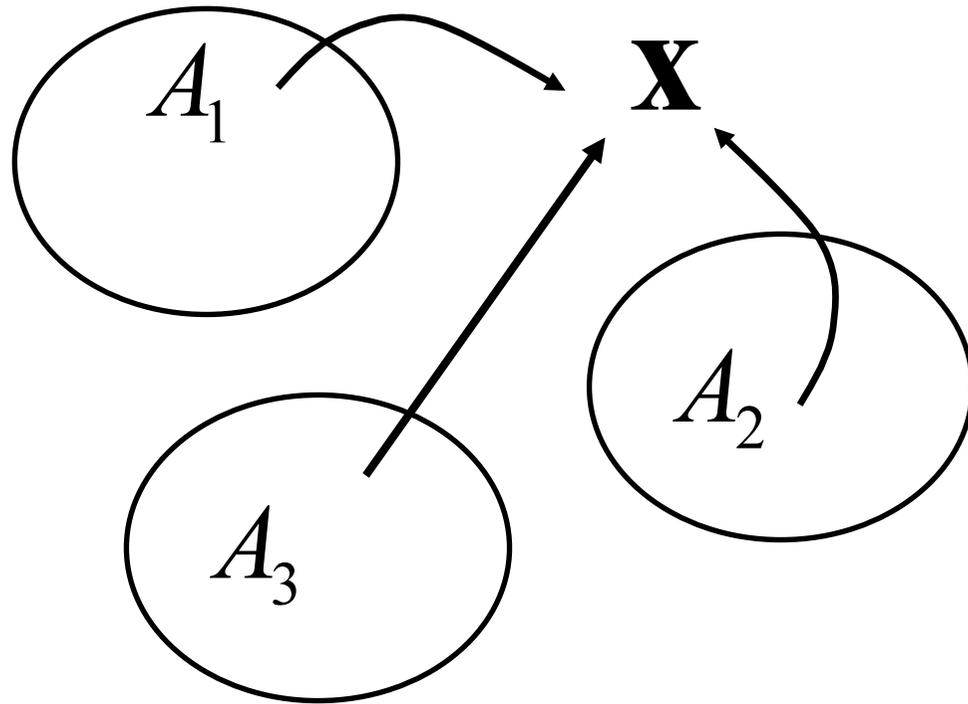
情報源の仮定

- 事例が統計的な分布に従う



前提: 事例 \mathbf{x} は各クラス A_i に対してある確率密度分布に従って生起していると仮定する

情報源の仮定



$P(A_i)$ 事前確率

$P(A_i | \mathbf{x})$ 事後確率

$p(\mathbf{x})$ \mathbf{x} の生起確率

$p(\mathbf{x} | A_i)$ 各クラスからの
 \mathbf{x} の生起分布

$$P(A_i | \mathbf{x}) = \frac{p(\mathbf{x} | A_i)}{p(\mathbf{x})} P(A_i)$$

Bayes
theorem

識別関数

- 事後確率を最大にするクラスを判別結果とする

$$\max_i \{P(A_i | \mathbf{x})\} = P(A_l | \mathbf{x})$$

未知の入力 \mathbf{x} はクラス A_l に分類される

→ つまり, $P(A_i | \mathbf{x})$ を識別関数と同じと見なす

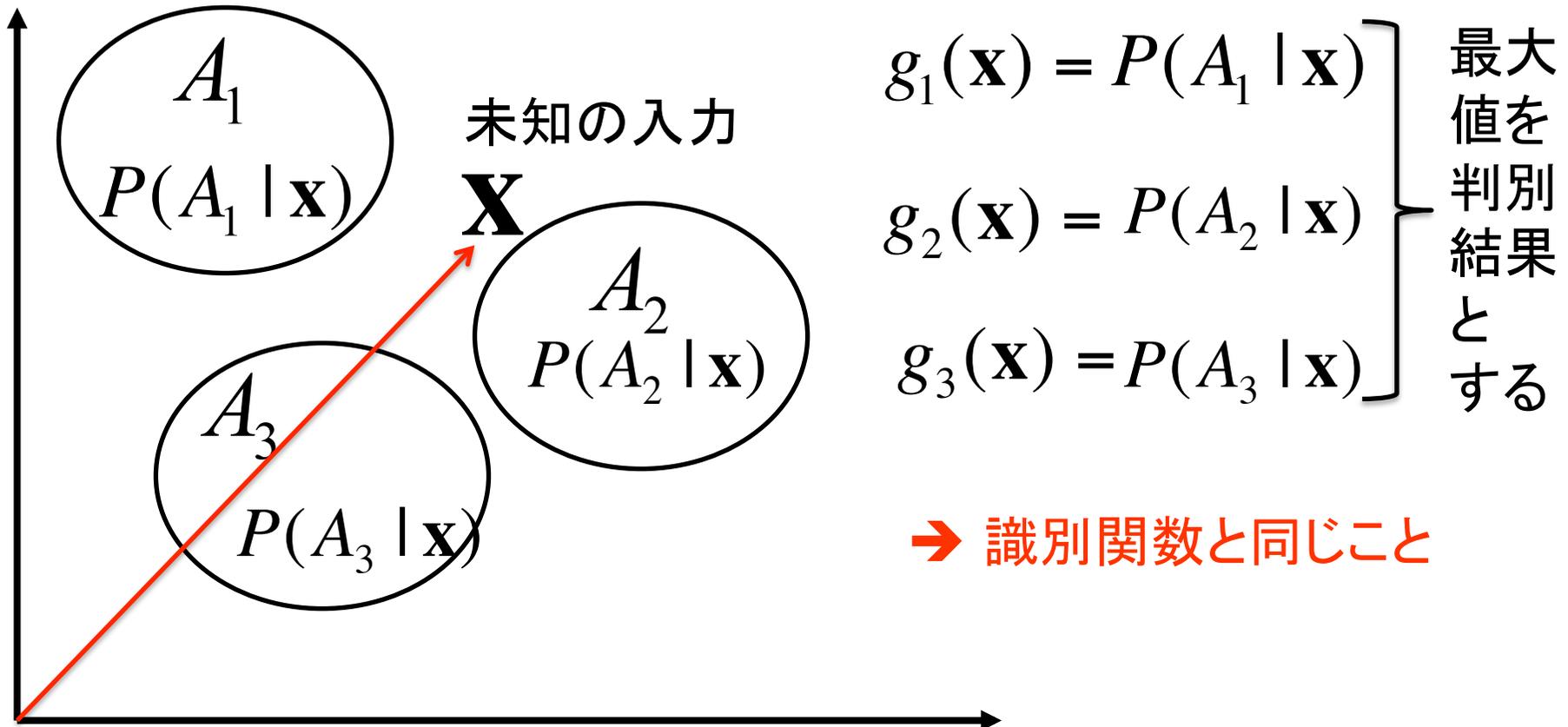
$$g_i(\mathbf{x}) = \frac{p(\mathbf{x} | A_i)}{p(\mathbf{x})} P(A_i)$$

def

$$= p(\mathbf{x} | A_i) P(A_i)$$

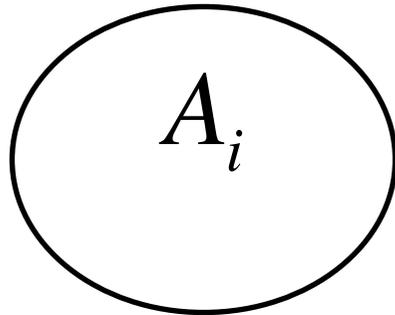
$p(\mathbf{x})$ はクラス A_i を
選択するときは
固定なので省略

事後確率の最大化とは？



練習

- 各クラス A_i 内である特徴量 x が生起する確率は正規分布であるとする
 - 平均ベクトル m_i
 - 共分散行列 Σ_iとするとき, 下記の図にその様子を書き込め (m_i, Σ_i を使うこと)



- 共分散行列 $\Sigma_i = \Sigma_0$ とはどのような状況か説明せよ
- 共分散行列 $\Sigma_0 = I$ (単位行列) とはどのような状況か説明せよ

例題

- 入力が各クラスの正規分布で発生していると仮定したとき, NN法(最小距離識別法)はどのような特殊な条件であるか答えよ

識別部の設計

- 確認
 - 前提としてある事例 \mathbf{x}_p はクラス A_i からどのように生成されていると考えられるか
 - 共分散行列 (covariance matrix) はどういうものか説明せよ
 - マハラノビス距離 (Mahalanobis distance) はどのような特徴を測ったもの?

生成モデルと判別モデル

- 仮定と事後確率

判別モデル $p(A_i | \mathbf{x})$ 事後確率
-> 識別関数

生成モデル $p(\mathbf{x} | A_i)$ 仮説としておいた
分布

やろうと思えば
学習データから
直接どちらも求まる

Bayes の定理を利用して $p(\mathbf{x} | A_i)$ を通して $p(A_i | \mathbf{x})$ を求める方法を生成モデルと呼ぶ。Hidden Markov Model など。

直接学習データから $p(A_i | \mathbf{x})$ を求める方法もある。事後確率を直接求めると分母が疎なので推定が(普通)よくないが、Conditional Random Field などでは $\exp()$ を分布として仮定して疎なデータに対して良い予測を与える。判別モデルと呼ばれる。

パラメータ推定

- 下記の式の意味を答えなさい

$$p(X; \phi) = \prod_{j=1}^n p(\mathbf{x}_j; \phi)$$

$$p(X | A_i; \phi_i) = \prod_{j=1}^n p(\mathbf{x}_j | A_i; \phi_i)$$

- パラメトリックな学習とは何か？

最尤推定

- 全学習データ X に対して下記の式が最大になるようにパラメータ ϕ を学習する

– そもそもなぜ最大にするのだろうか? $p(X; \phi) = \prod_{j=1}^n p(\mathbf{x}_j; \phi)$

$p(X; \phi)$ は ϕ によって $p(x)$ の形が決まり, 全学習データ X の確率値を書けた確率値(いつものパターン!)

$p(x; \phi)$ は未知だけど, 実際にデータ $X = \{x_1, x_2, \dots, x_n\}$ は出現した \rightarrow これが出てきたのだからデータ X を $p(x; \phi)$ はうまく説明して欲しい

今, x は分布 $p(x; \phi)$ に依存して出力してると考えている.
ならば, $p(x; \phi)$ から X が出力したと考えると, その p で全データ X に対して確率値 $p(X; \phi)$ を求めれば, 最大の値になるはず

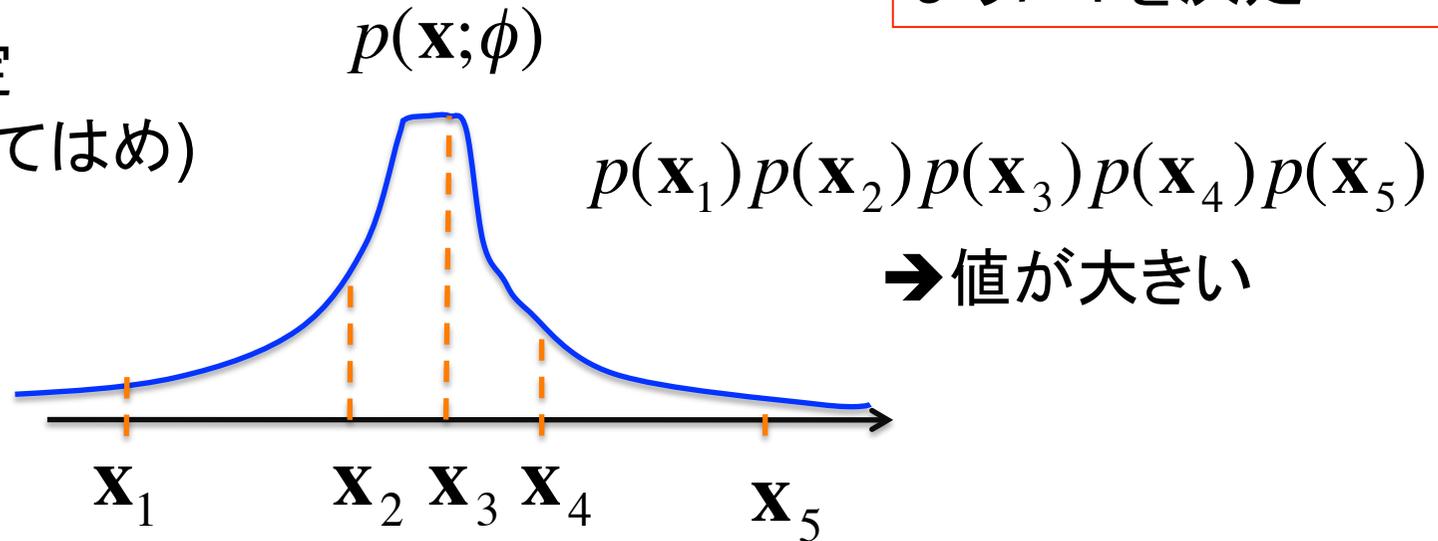
具体例

学習データ $X = \{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_5\}$
が下記のように分布しているとする

最尤推定 →
密なところに大きな
確率値を割り当てる
ように ϕ を決定

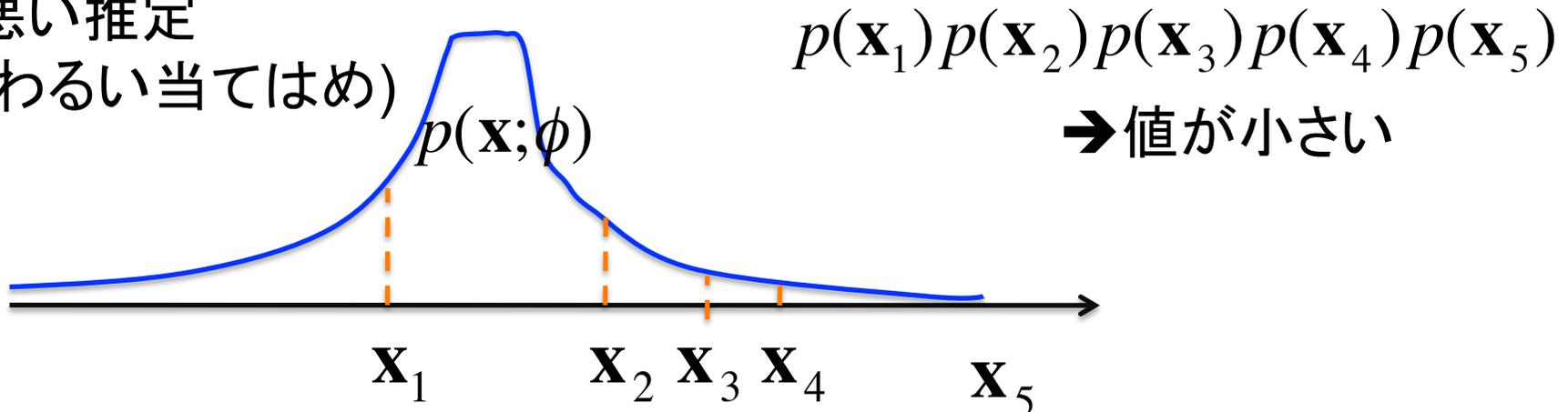
良い推定

(良い当てはめ)



悪い推定

(わるい当てはめ)



最尤推定(確率最大化→数を数える)

- 推定対象の確率変数

$\{\theta_1, \dots, \theta_N\}$

- N個の事象

$\{N_1, \dots, N_n\}$

- 事象 N の出現回数が n_i

- 事象の生起確率

$\{P(\theta_1), \dots, P(\theta_N)\} \rightarrow$ これを最尤推定

次の式を以下の制約のもとに最大化する

$$L = \sum_{i=1}^N n_i \log P(\theta_i)$$

制約

$$\sum_{i=1}^N P(\theta_i) = 1$$

Lanrange 未定乗数法

$P(s_i|s_{i-1})$ の場合

θ_i は 例えば $s_i|s_{i-1}$

N_i は s_{i-1} から s_i の遷移

n_i は s_{i-1} から s_i の遷移の回数

$C(s_{i-1} \rightarrow s_i)$

最尤推定

Lagrange 関数

$$F = \sum_{i=1}^N n_i \log P(\theta_i) + \lambda \sum_{i=1}^N P(\theta_i)$$

これを θ_i で偏微分する

$$\frac{\partial F}{\partial P(\theta_i)} = \frac{n_i}{P(\theta_i)} + \lambda = 0 \quad \text{よって} \quad P(\theta_i) = -\frac{n_i}{\lambda} \quad \textcircled{1}$$

制約式に代入して

$$1 = \sum_{i=1}^N P(\theta_i) = -\frac{1}{\lambda} \sum_{i=1}^N n_i \quad \lambda = -\sum_{i=1}^N n_i$$

よって λ を①に代入して

$$P(\theta_i) = \frac{n_i}{\sum_{k=1}^N n_k}$$

参考:奈良先端大:
音声情報処理講義録
鹿野先生 1994年

ラグランジュ未定乗数法

- 制約の中で式を解く
 - M個の制約条件 $g_i(x) = 0$ ($i=1,2,\dots,M$)
 - $f(X)$ の極値をとりたい (極大or極小)

ラグランジュ乗数 $\alpha = [\alpha_1, \dots, \alpha_M]^t$ を使って

$$L(X, \alpha) = f(X) - \sum_{i=1}^M \alpha_i g_i(X)$$

$$\frac{\partial L}{\partial x_k} = 0 \quad \frac{\partial L}{\partial \alpha_i} = 0$$

$(k = 1, 2, \dots, N) \quad (i = 1, 2, \dots, M)$

を解けばよい.