

情報意味論(13)

(簡単に)事例ベースアプローチ

櫻井彰人
慶應義塾大学理工学部

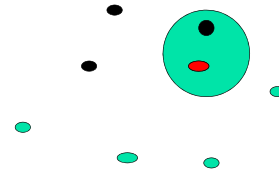
事例ベース学習

- キーアイデア
 - 訓練データ $\langle x_i, f(x_i) \rangle$ を全て憶えていよう(とりあえずは、何も、または、あまりしない)
 - 問い合わせがあつたら、その時点で、しよう
- この類に属する方法
 - 最近傍法 (Nearest neighbor)
 - k -Nearest neighbor
 - Locally weighted regression
 - Radial basis functions
- Lazy 対 eager

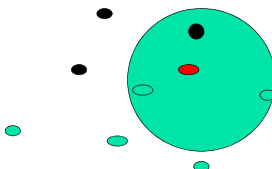
最近傍法

- 最近傍法 (Nearest neighbor)
 - 問合せ x_q に対し、最近接の x_n を見つけ、 $f(x_q) \leftarrow f(x_n)$ とする
- k -Nearest neighbor
 - k 個の最近接データの間で、多数決
 - k 個の最近接データの間で、平均値

1-Nearest Neighbor

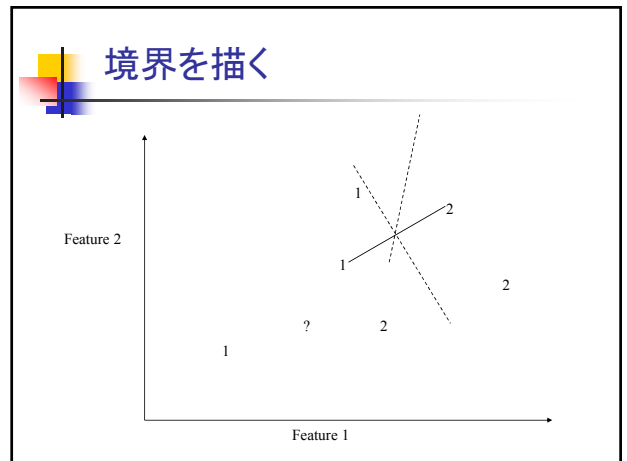
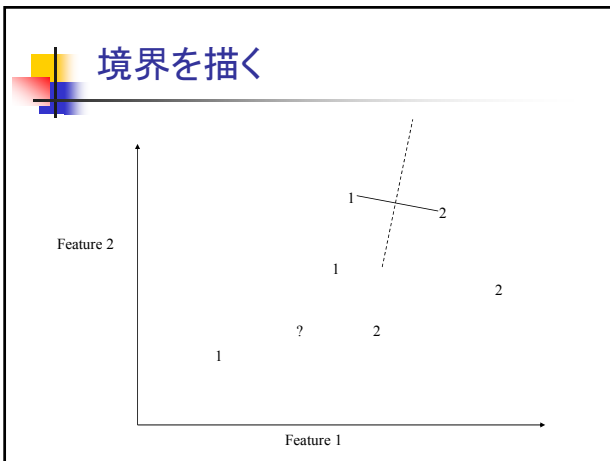
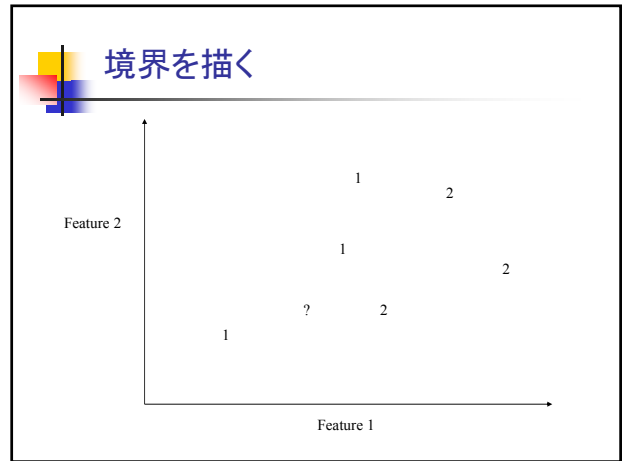
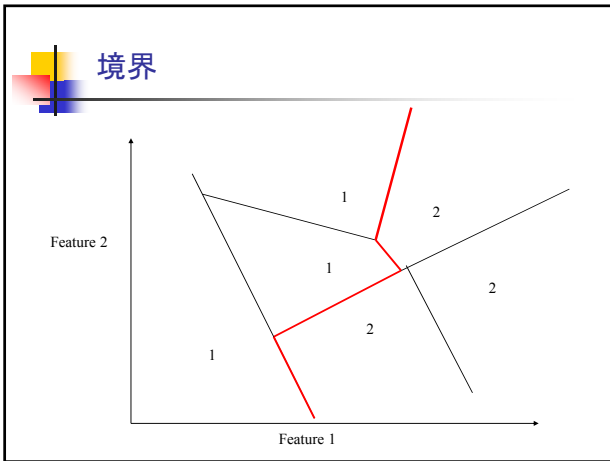
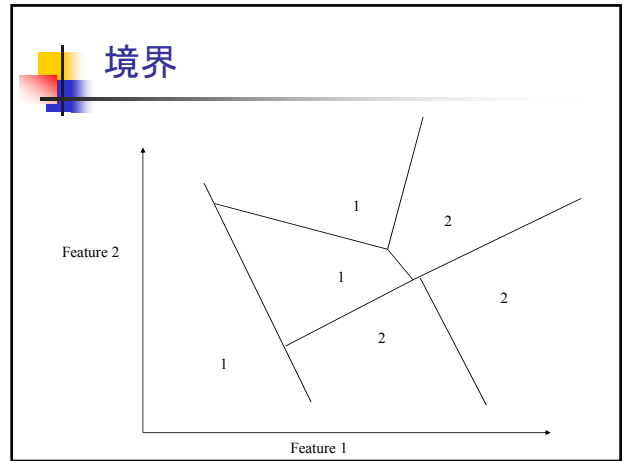
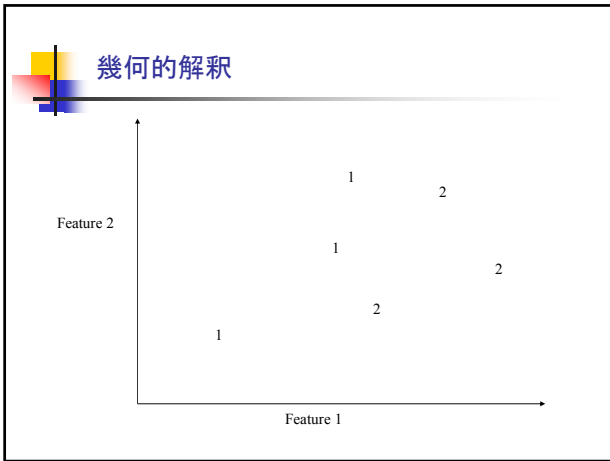


3-Nearest Neighbor



最近傍法の特徴

- いつ使うか
 - 属性が R^n の点とみなせる
 - 属性数はあまり多くない(数十個?)
 - 大量の訓練データ
- 長所
 - 学習が速い
 - 複雑な目標関数も表現可能
 - (訓練データがもつ)情報を失うことがない
- 短所
 - 問合せ時、遅い
 - 無関係な属性によって、簡単に、ごまかされる



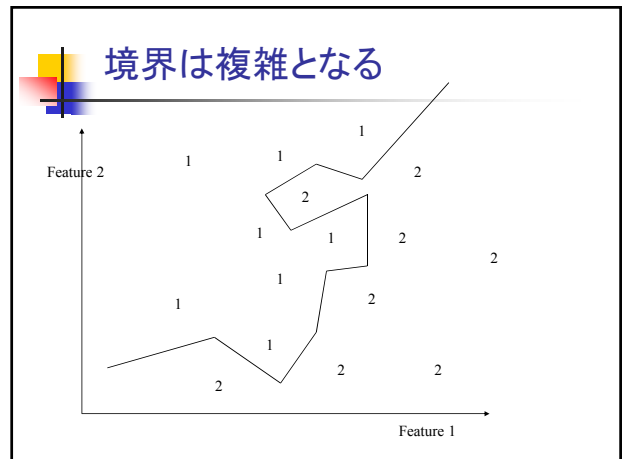
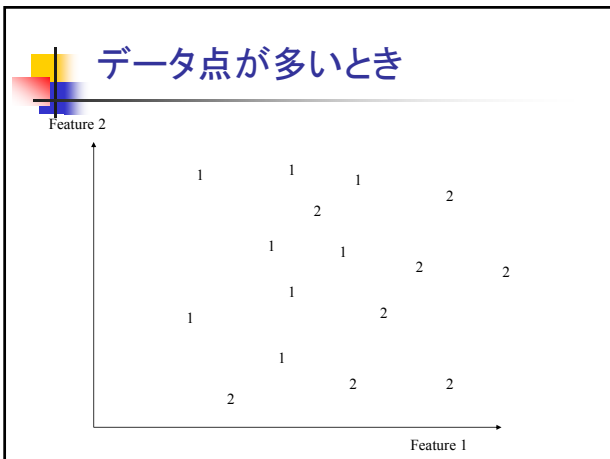
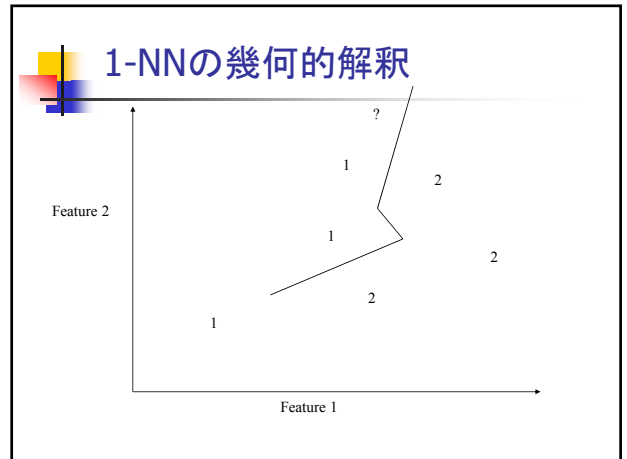
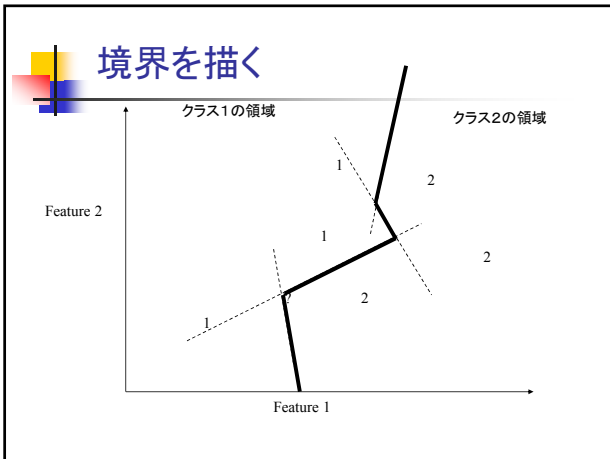
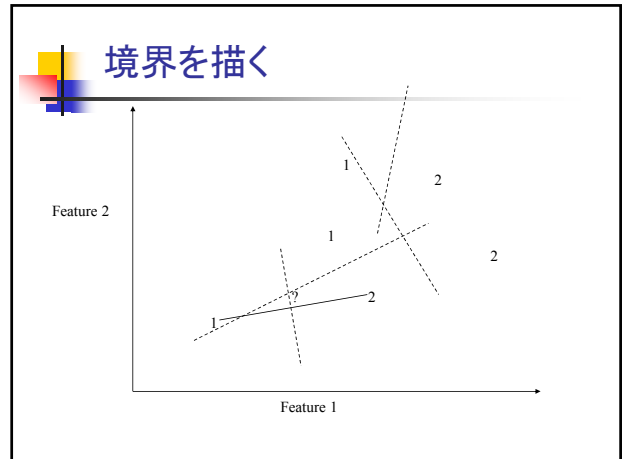
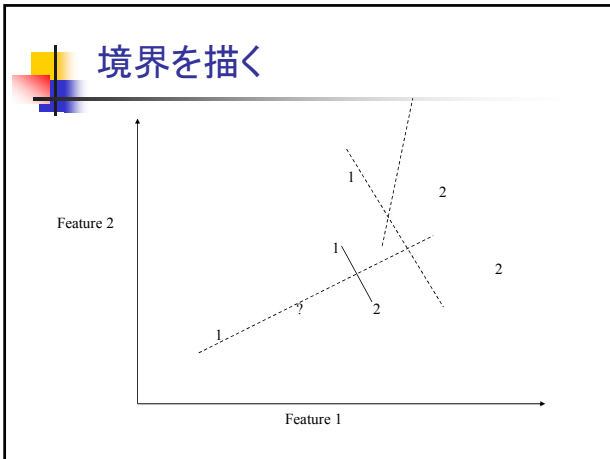


Table VI. Comparative Results Among Different Classifiers Obtained on Five Different Versions of Reuters. (Columns of maximum recall, circles indicate the macro-averaged breakdown point, with its parentheses. 'M' indicates macro-entropy and 'F1' indicates use of the F1 measure; boldface indicates the best performer on the selection.)

System	Type	Results reported by	F1	F2	F3	F4	F5
		# of documents	21,459	14,347	19,272	12,300	12,592
		# of training documents	14,704	10,687	8,210	10,603	10,041
		# of test documents	6,746	3,660	3,662	3,259	3,259
		# of categories	135	63	62	56	73
Word	rule-based	Yang (1999)	169	110	207	152	313
ParaRank	probabilistic	(Dumais et al. 1998)					
	probabilistic	(Joachims 1998)					729
	probabilistic	(Lee et al. 1997)	443 (M _F)				
	probabilistic	(Liu 1992a)	650				
Na	probabilistic	(Li and Yamashita 1996)				747	773
	probabilistic	(Yang and Liu 1999)				795	
C4.5	decision tree	(Dumais et al. 1998)					884
	decision tree	(Joachims 1998)					
Iaa	decision tree	(Liu and Ringuette 1994)	670			794	
	decision tree	(Liu et al. 1991)		805			
Naive Bayes	decision rule	(Cohen and Singer 1999)	683	811		820	
	decision rule	(Cohen and Singer 1999)	753	759		827	829
Support Vector	decision rule	(Li and Yamashita 1996)					829
	decision rule	(Mollinari and Giacinto 1998)					
Classifier	decision rule	(Mollinari et al. 1998)					
	decision rule	(Yang 1999)					
Log	regression	(Yang and Liu 1999)		855	810		
	regression	(Yang et al. 1997)	747 (M)	843 (M)			849
BALANCE WISCONSIN	co-linear linear	(Dumais et al. 1998)					822
	co-linear linear	(Liu and He 1998)					
Rocchio	batch linear	(Cohen and Singer 1999)	690	745		745	648
	batch linear	(Dumais et al. 1998)					747
Rocchio	batch linear	(Joachims 1998)					799
	batch linear	(Lee and He 1998)					784
Rocchio	batch linear	(Li and Yamashita 1996)					825
	batch linear	(Yang et al. 1997)					
Clique	neural network	(Yang and Liu 1999)		802			
	neural network	(Wu et al. 1995)					838
GmW	example-based	(Liu and He 1998)			820		
	example-based	(Joachims 1998)					809
k-NN	example-based	(Liu and He 1998)					825
	example-based	(Yang 1999)	690	852	820		829
k-NN	example-based	(Yang and Liu 1999)					856
	example-based	(Dumais et al. 1998)					879
SvmLearn	SVM	(Liu and He 1998)					864
	SVM	(Li and Yamashita 1996)					841
SvmLearn	SVM	(Yang and Liu 1999)					859
	SVM	(Wu et al. 1995)					
AmiBoostM1	committee	(Schapire and Singer 2000)		800			
	committee	(Wu et al. 1995)					878
Bayesian net	Bayesian net	(Dumais et al. 1998)					800
	Bayesian net	(Lee et al. 1997)	542 (M _F)				

Fabrizio Sebastiani, Machine learning in automated text categorization, ACM Computing Surveys, vol.34, no.1, 1-47 (2002)

極限における振り舞い

- $p(x)$: 事例 x がラベル1(正)をもつ事後確率
- Nearest neighbor:
 - 事例数 $\rightarrow\infty$ のとき, Gibbsアルゴリズムに漸近
 - Gibbs: 確率 $p(x)$ で1を予測
- k -Nearest neighbor
 - 事例数 $\rightarrow\infty$ かつ k が大きくなると, Bayes最適
 - Bayes最適: $p(x) > 0.5$ なら1, それ以外0

注: Gibbs の期待誤差はBayesの倍以下

復習

Bayes 最適な分類器

$$\arg \max_{c_j \in \{+,-\}} \sum_{h_j \in H} P(c_j | h_j) P(h_j | D)$$

注: Bayes 最適な分類器は H に含まれるとは限らない
 注: 偏置にはうまくいくと報告されているのだが、試してみるとMAPやMLと変わらない場合がある。どのような場合にそうなるか、興味のあるところである
 注: 実行可能か? 見るからに時間がかりそう

Gibbs 分類器 - 速度向上

- 仮説を $P(h|D)$ に従ってランダムに選ぶ
- 新事例をこれに従って分類する

復習: もし仮説を事前分布 $P(h)$ に従ってランダムに選ぶと,
 $E[\text{error}_{\text{Gibbs}}] \leq 2E[\text{error}_{\text{BayesOptimal}}]$

(詳細は "Mitchell Machine Learning Chap. 6.8")
 復習の回数が増えてくると、ベイズ最適な分類器が計算できないと表に実行

距離荷重つき k-NN

- 近い事例の判断を重視したい

$$\hat{f}(x_q) \leftarrow \frac{\sum_{i=1}^k w_i f(x_i)}{\sum_{i=1}^k w_i}, \quad w_i \equiv \frac{1}{d(x_q, x_i)^2}$$

但し, $d(x_q, x_i)$ は, x_q と x_i の間の距離

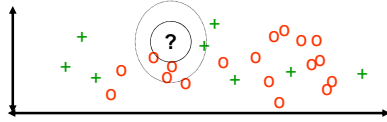
- これにより, k 個のみならず全データを使うことに意味がでてくる \Rightarrow Shepardの方法

K-NN と不要な特徴

+++ 00 0(070) ++0+ 0 0000+00000 +

K-NN と不要な特徴

K-NN と不要な特徴



距離の問題

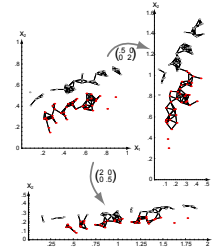
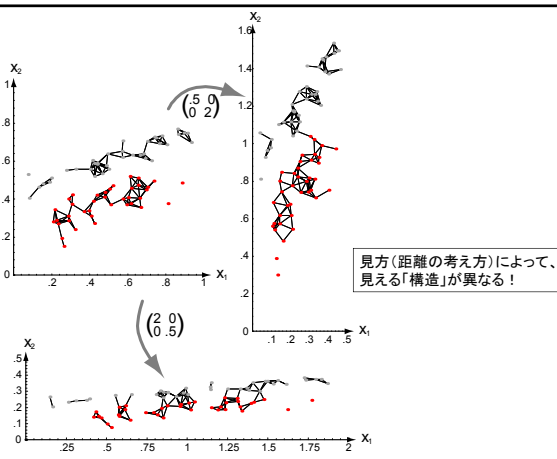


FIGURE 10.8 Scaling axes affects the clusters in a minimum distance cluster method. The original data and minimum-distance clusters are shown in the upper left; points in one cluster are shown in red, while the others are shown in gray. When the vertical axis is expanded by a factor of 2.0 and the horizontal axis shrunk by a factor of 0.5, the clustering is altered (as shown at the right). Alternatively, if the vertical axis is shrunk by a factor of 0.5 and the horizontal axis is expanded by a factor of 2.0, smaller more numerous clusters result (shown at the bottom). In both these scaled cases the assignment of points to clusters differs from that in the original space. From: Richard O. Duda, Peter E. Hart, and David G. Stork, Pattern Classification, Copyright © 2001 by John Wiley & Sons, Inc.



次元の呪い

- 20個の属性で記述されるが、その内、たった2属性のみが意味ある場合を考える
- 次元の呪い:
 - k -NNなら、他の18属性の値でどんな結論も出うる
- 解決方法
 - j 番目の属性に z_j の荷重を。 z_j は予測誤差最小となるように選択
 - cross-validationを用いて自動的に z_j を決定

Locally weighted regression

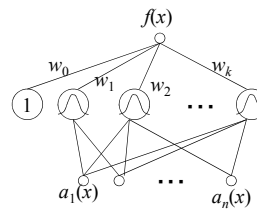
- k -NN は各問合せ x_q で f の局所近似を構成していた
- x_q の周囲で $f(x)$ の近似関数を明示的に構成したらどうだろうか?
 - k -NNに線型回帰したら?
 - 2次回帰では?
 - 区分回帰したら?
- 最小化すべき誤差にもいくつかの候補が

$$E_1(x_q) = \frac{1}{2} \sum_{x \in x_q \text{ の } k\text{-NN}} (f(x) - \hat{f}(x_q))^2$$

$$E_2(x_q) = \frac{1}{2} \sum_{x \in D} (f(x) - \hat{f}(x_q))^2 K(d(x_q, x))$$

Radial Basis Function Network

- 局所近似の線型結合による大域近似
- 神経回路網の一種
- distance-weighted regression に類似
 - lazy ではなく eager であるが



$$f(x) = w_0 + \sum_{u=1}^k w_u K_u(d(x_u, x))$$

$K_u(d(x_u, x))$ の一例

$$K_u(d(x_u, x)) \equiv e^{-\frac{1}{2\sigma^2}d(x_u, x)^2}$$

RBFの学習

- $K_u(d(x_u, x))$ の x_u の定め方
 - 事例空間に一様にばら撒く
 - 事例を使用(事例の分布が反映)
- 荷重の学習 (K_u は正規分布とする)
 - 各 K_u の分散(と平均)を定める
 - 例えば、EMを使用
 - K_u を固定したまま、線型出力部分を学習
 - 線型回帰で高速に

Lazy 対 eager

- Lazy: 事例からの一般化をしない。問合せがあったときに考える
 - k-Nearest Neighbor
- Eager: 問合せ前に予め一般化しておく
 - 「学習」アルゴリズム、ID3, 回帰, RBF, ...
- 違いはあるか？
 - Eager学習は全域的な近似を作成
 - Lazy学習は局所近似を大量に作成
 - 同じ仮説空間を使うなら、lazyの方が複雑な関数を作成
 - over-fittingの可能性
 - 柔軟(複雑なところと単純なところの組合せ)

まとめ

- 事例ベースアプローチ
 - 大域的な構造を仮定しない
 - どんな場合にも使える
 - 雑音に弱い(大域構造を用いた平滑化ができない)
 - 次元の呪い