

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

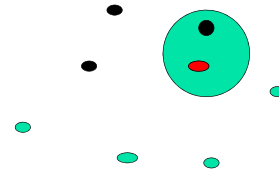
事例ベース学習

- 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

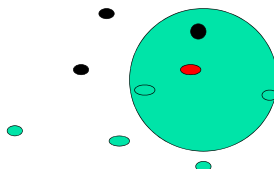
最近傍法

- 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



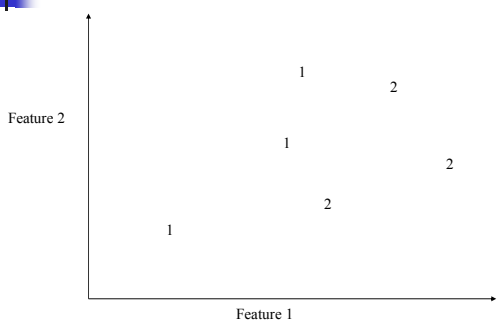
3-Nearest Neighbor



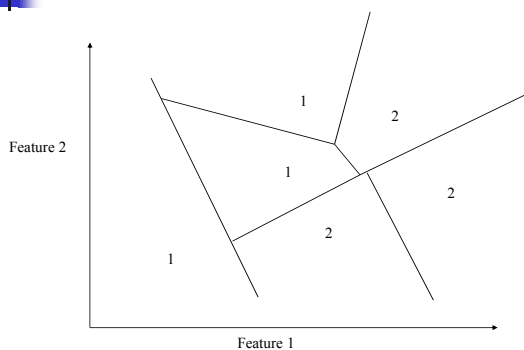
最近傍法の特徴

- 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

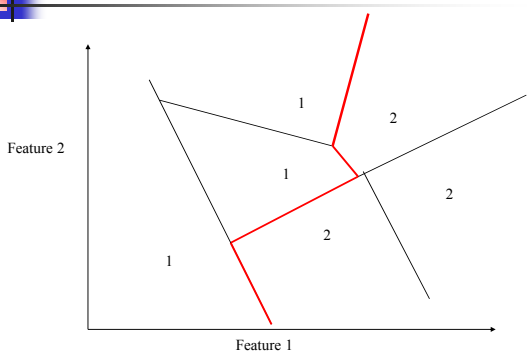
幾何的解釈



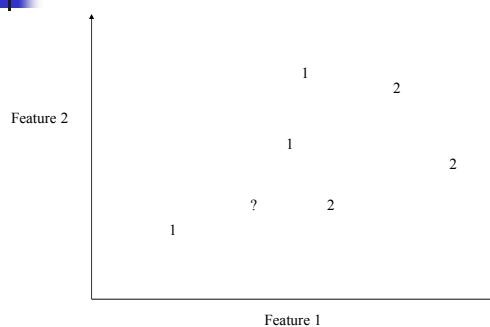
境界



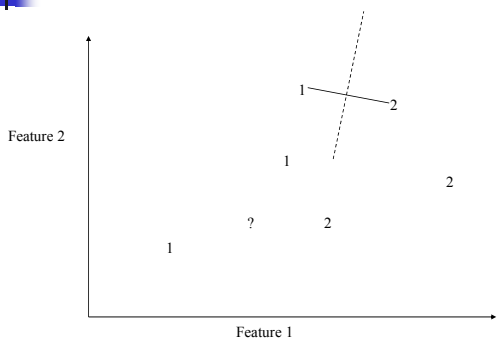
境界



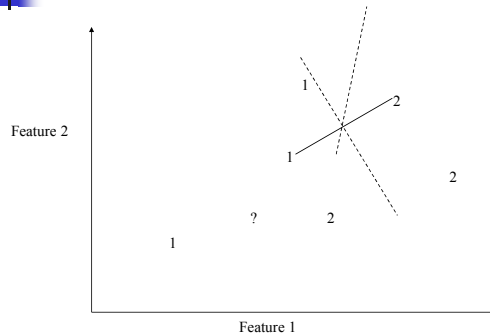
境界を描く

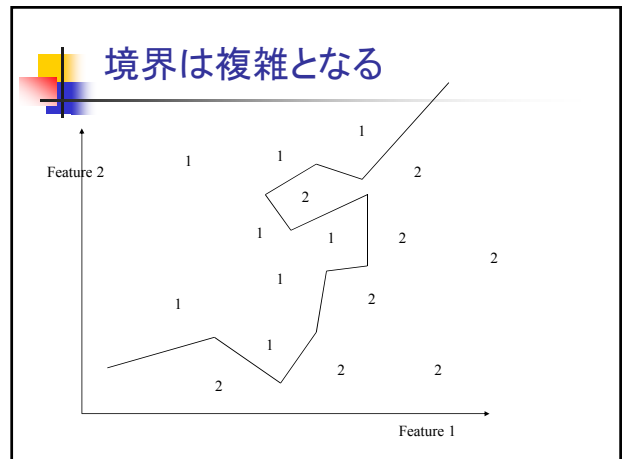
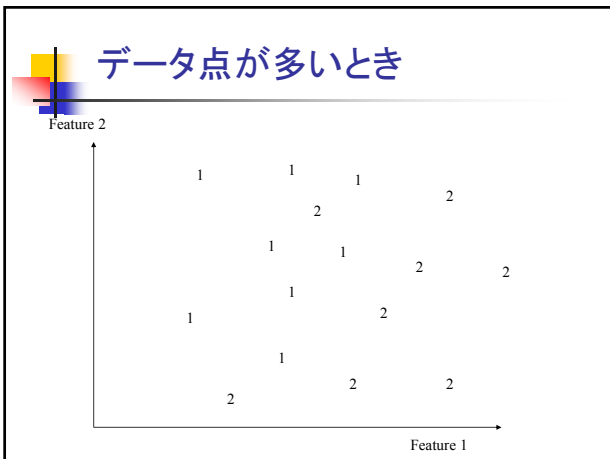
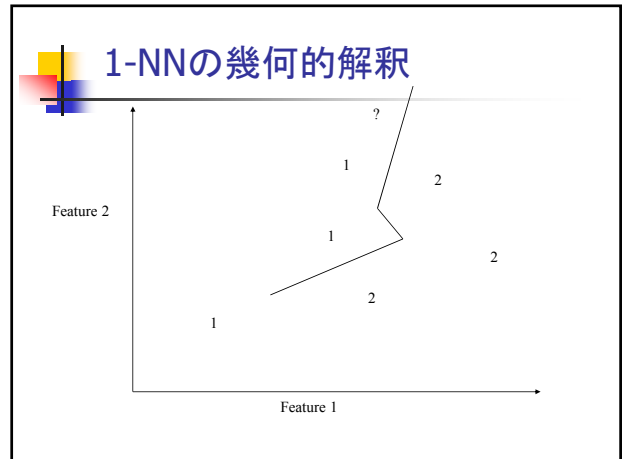
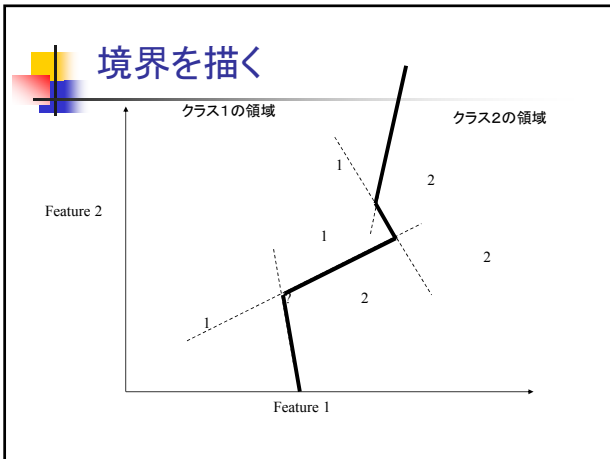
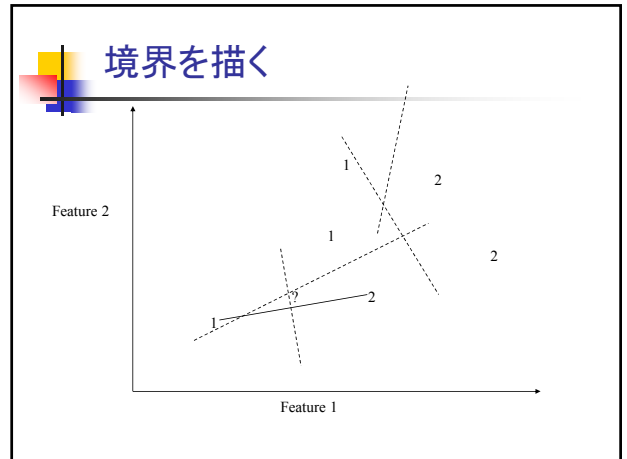
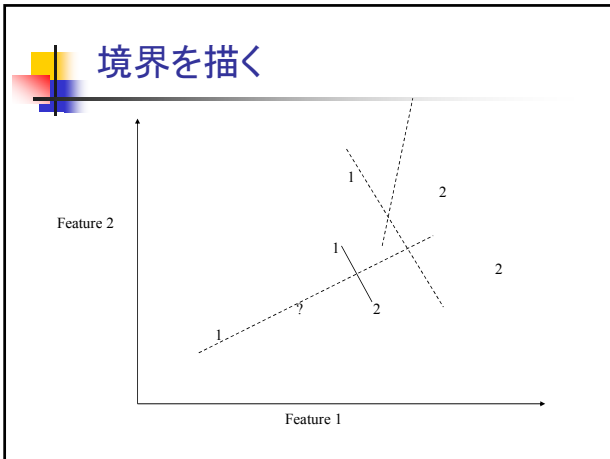


境界を描く



境界を描く





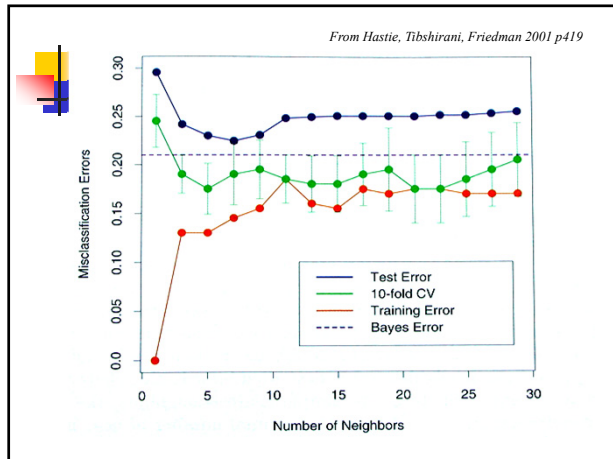
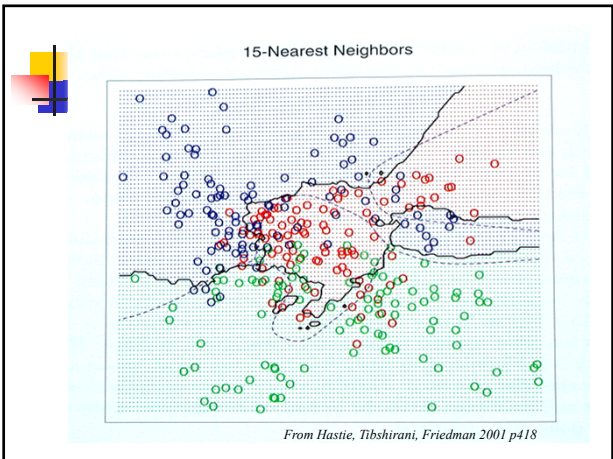
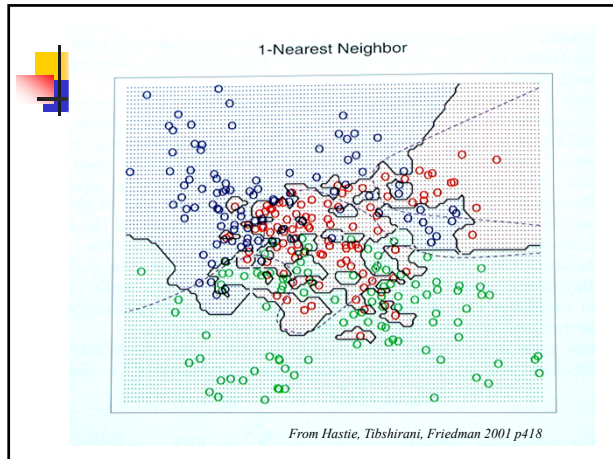


Table 6. Results summary of TC systems on Reuters versions 1-4.

System	Reuters version 1	Reuters version 2	Reuters version 3	Reuters version 4
WORD	—	.15 (Scut)	.31 (Pcut)	.29 (Pcut)
kNN	—	.69 (Scut)	.85 (Scut)	.82 (Scut)
LLSF	—	—	.85 (Scut)	.81 (Scut)
NNets.PARC (perceptron)	—	—	—	.82 (Pcut)
CLASS1 (perceptron)	—	—	.80	—
RIPPER (DNF)	—	.72 (Scut)	.80 (Scut)	—
SWAP-1 (DNF)	—	—	.79	—
DTree IND	—	.67 (Pcut)	—	—
DTree C4.5	—	—	.79 (F1)	—
CHARADE (DNF)	—	—	.78	—
EXPERTS (n-gram)	—	.75 (Scut)	.76 (Scut)	—
Rocchio	—	.66 (Scut)	.75 (Scut)	—
NaiveBayes	—	.65 (Pcut)	.71	—
CONSTRUE (Exp. Sys.)	.90	—	—	—

Yiming Yang, An Evaluation of Statistical Approaches to Text Categorization, Information Retrieval, vol.1, 69-90 (1999)

System	Type	Reuters reported by	#1	#2	#3	#4	#5
Wtd	linear	Yang 1999	100	310	200	752	815
ProbBayes	probabilistic	[Dumais et al 1998]	485 (M _{F1})	—	—	—	720
ES	probabilistic	[Joachims 1995]	—	—	—	—	—
Sa	probabilistic	[Lewis 1992]	—	—	—	—	—
C4.5	decision tree	[Li and Yamashita 1999]	—	—	—	717	775
Idp	decision tree	[Yang and Liu 1999]	—	—	—	765	—
Idp	decision tree	[Dumais et al 1998]	—	—	—	841	794
Idp	decision tree	[Joachims 1995]	.670	—	—	—	—
Idp	decision tree	[Yang et al 1999]	—	.805	—	—	—
Idp	decision tree	[Cohen and Singer 1999]	.683	.811	—	.820	—
Idp	decision tree	[Cohen and Singer 1999]	.784	.795	—	.820	—
Idp	decision tree	[Li and Yamashita 1999]	—	—	738	—	—
Idp	decision tree	[Moulinier and Gnanou 1996]	—	—	763 (F1)	—	—
Idp	decision tree	[Moulinier et al 1996]	—	—	—	—	—
Idp	regression	[Yang 1999]	—	.855	.810	—	—
Idp	regression	[Yang and Liu 1999]	—	—	—	—	—
Idp	regression	[Dumais et al 1998]	.717	.833	—	.822	—
Idp	regression	[Lam and He 1994]	—	—	—	716	—
Idp	regression	[Cohen and Singer 1999]	.620	.748	—	.817	—
Idp	regression	[Dumais et al 1998]	—	—	—	—	.646
Idp	regression	[Joachims 1995]	—	—	—	—	.759
Idp	regression	[Lam and He 1994]	—	—	—	781	—
Idp	regression	[Li and Yamashita 1999]	—	—	—	725	—
Idp	regression	[Yang et al 1999]	—	—	.802	—	—
Idp	regression	[Yang and Liu 1999]	—	—	—	—	.834
Idp	regression	[Witner et al 1995]	—	—	—	.820	—
Idp	regression	[Joachims 1995]	—	—	—	—	.858
Idp	regression	[Lam and He 1994]	—	—	—	—	.825
Idp	regression	[Yang 1999]	.690	.852	.820	—	.856
Idp	regression	[Yang and Liu 1999]	—	—	—	—	.866
Idp	regression	[Dumais et al 1998]	—	—	—	—	.870
Idp	regression	[Joachims 1995]	—	—	—	—	.884
Idp	regression	[Yang and Liu 1999]	—	—	—	—	.889
Idp	regression	[Cohen and Singer 2000]	—	—	.860	—	.878
Idp	regression	[Waltz et al 1999]	—	—	—	—	.890
Idp	regression	[Friedman et al 1996]	—	—	—	—	.890
Idp	regression	[Lam et al 1997]	.542 (M _{F1})	—	—	—	.890

Table 6. Comparative results among different classifiers obtained on five different version of the Reuters collection. Unless otherwise noted, entries indicate the microaveraged breakeven point; within parentheses, “M” indicates macroaveraging and “F1” indicates use of the F1 measure. Boldface indicates the best performer on the collection.

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

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

System	Type	Results reported by	F1	F2	F3	F4	F5
		# of documents	21,150	14,347	17,272	12,300	12,502
		# of training documents	14,704	10,687	8,210	10,603	10,041
		# of test documents	6,746	3,660	9,062	3,298	2,269
		# of categories	135	83	82	106	131
Word	probabilistic	Yang (1999)	168	110	207	152	113
ParaRank	probabilistic	(Dumais et al. 1998)					
	probabilistic	(Joachims 1998)					129
	probabilistic	(Lee et al. 1997)	443 (MF)				147
	probabilistic	(Liu 1992a)	650				173
	probabilistic	(Li and Yamashita 1996)					205
	probabilistic	(Li and Yamashita 1996)					284
	probabilistic	(Yang and Liu 1999)					284
C4.5	decision tree	(Dumais et al. 1998)					584
	decision tree	(Joachims 1998)					
	decision tree	(Lee and Elmaghrabi 1994)	670				794
Naïve Bayes	decision rules	(McLis et al. 1991)		806			820
SVM	decision rules	(Cohen and Singer 1999)	683	811			827
Support Vector	decision rules	(Cohen and Singer 1999)	753	759			829
Naïve Bayes	decision rules	(Li and Yamashita 1996)					829
Classifier	decision rules	(Moulinier and Gauthier 1998)		738			787 (F)
Classifier	decision rules	(Moulinier et al. 1998)		787 (F)			
Log	regression	Yang (1999)		856	810		849
Log	regression	(Yang and Liu 1999)					
HALAKIRI	co-line linear	(Dumais et al. 1997)	747 (M)	833 (M)			822
WordFlow	co-line linear	(Lee and He 1998)					
Reccmo	batch linear	(Lesh and Singer 1999)	690	745			745
Reccmo	batch linear	(Dumais et al. 1998)					747
Reccmo	batch linear	(Joachims 1998)					759
Reccmo	batch linear	(Lee and He 1998)					784
Reccmo	batch linear	(Li and Yamashita 1996)					825
Clean	neural network	(Yang et al. 1997)		802			838
Naïve	neural network	(Yang and Liu 1999)					838
GenM	neural network	(Warner et al. 1995)			820		
GenM	example-based	(Liu and He 1998)			809		821
k-NN	example-based	(Joachims 1998)			821		829
k-NN	example-based	(Lee and He 1998)			829		
k-NN	example-based	(Yang 1999)	690	852	820		856
k-NN	example-based	(Yang and Liu 1999)			856		
SVM	SVM	(Dumais et al. 1998)					879
SVM	SVM	(Joachims 1998)					864
SVM	SVM	(Li and Yamashita 1996)					841
SVM	SVM	(Yang and Liu 1999)					859
AmECorM	committee	(Schapire and Singer 2000)			800		
AmECorM	committee	(Waltos et al. 1998)					878
Bayesian net	Bayesian net	(Dumais et al. 1998)					800
Bayesian net	Bayesian net	(Lee et al. 1997)	542 (MF)				800

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: 全ての仮説を考える
 $p(x) > 0.5$ 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

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

注: Bayes 最適な分類器は H に含まれるとは限らない
注: 疑念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

慶賞: 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")
仮説の個数が増えてくると、ベイズ最適な分類器が計算できないときに有用

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

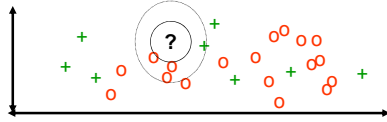
K-NN と不要な特徴

Diagram illustrating feature selection for K-NN. A sequence of features is shown: + + + 00 (circled) + + 0 + 0000 + 00000 +. The circled '00' has a question mark above it, suggesting it might be an unnecessary feature.

K-NN と不要な特徴

Scatter plot illustrating K-NN with unnecessary features. A point (circled with a question mark) is surrounded by several other points. Some are '+' and some are 'o'. The plot shows how distance-based weighting affects classification.

K-NN と不要な特徴



距離の問題

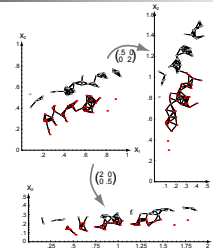
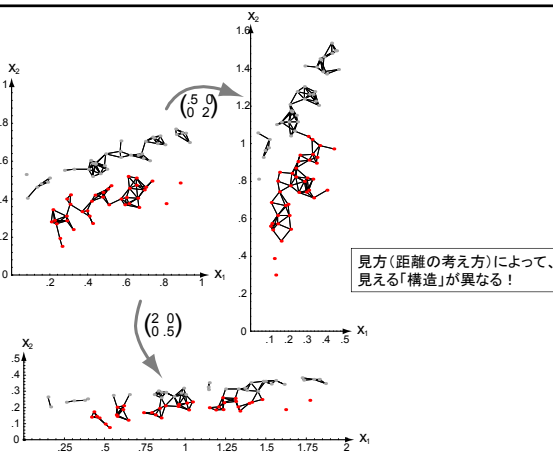


FIGURE 10.8 Scaling axes affects the clustering in a minimum distance cluster method. The original data and minimum-distance clusters are shown in the upper left; points in one cluster are shown in red, while the others are shown in gray. When the vertical axis is expanded by a factor of 2.0 and the 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

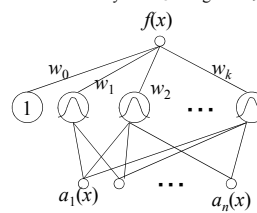
- k -NN は各問合せ x_q で f の局所近似を構成していた
- x_q の周囲で $f(x)$ の近似関数を明示的に構成したらどうだろうか?
 - k -NNに線型回帰したら?
 - 2次回帰では?
 - 区分回帰したら?
- 最小化すべき誤差にもいくつかの候補が

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

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

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