

## 情報意味論(10)

(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

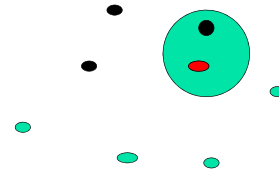
## 事例ベース学習

- 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

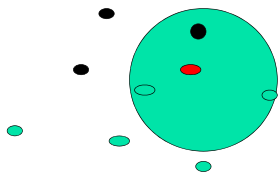
## 最近傍法

- 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



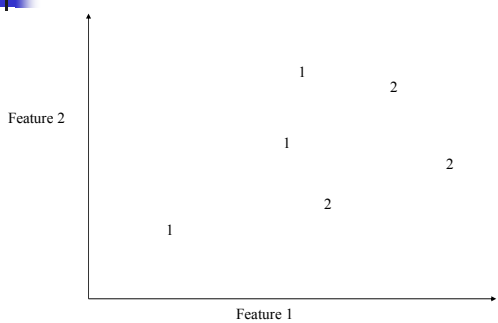
## 3-Nearest Neighbor



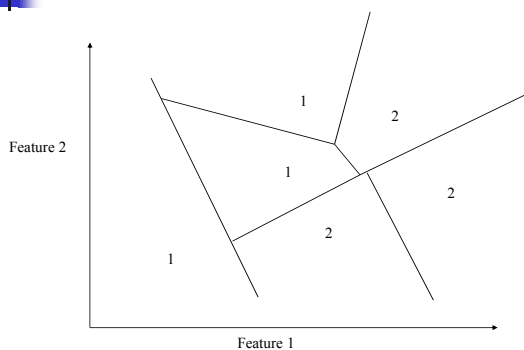
## 最近傍法の特徴

- 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

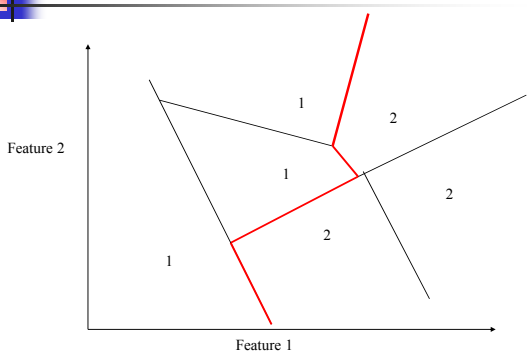
### 幾何的解釈



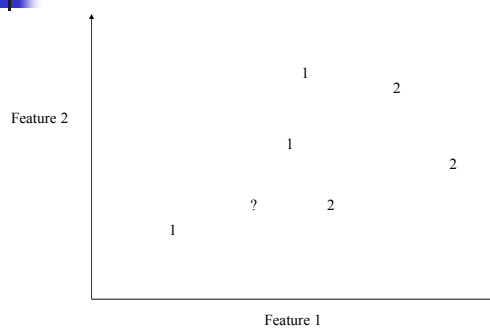
### 境界



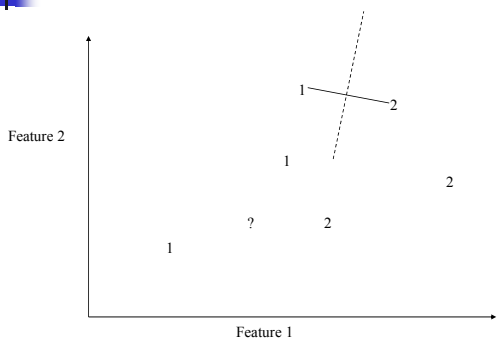
### 境界



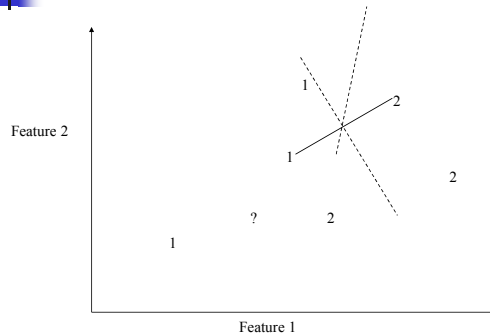
### 境界を描く

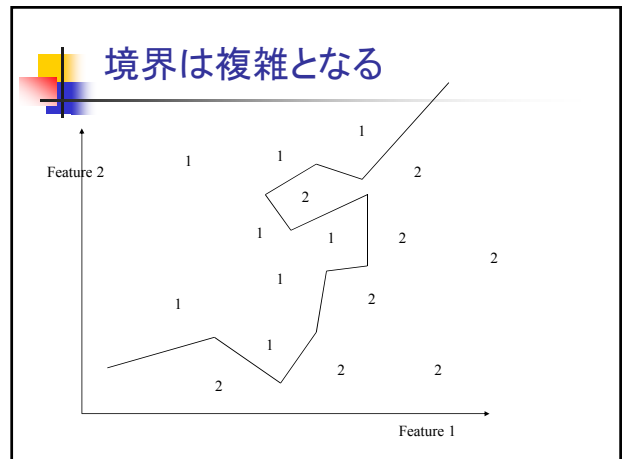
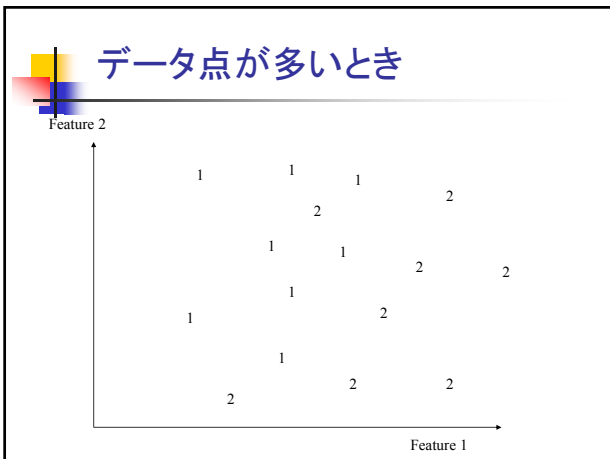
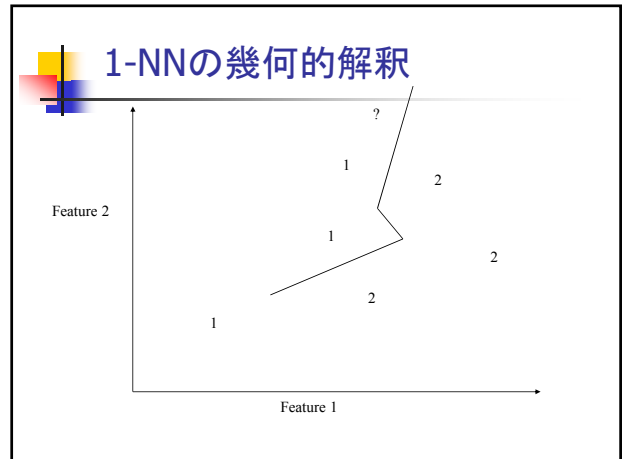
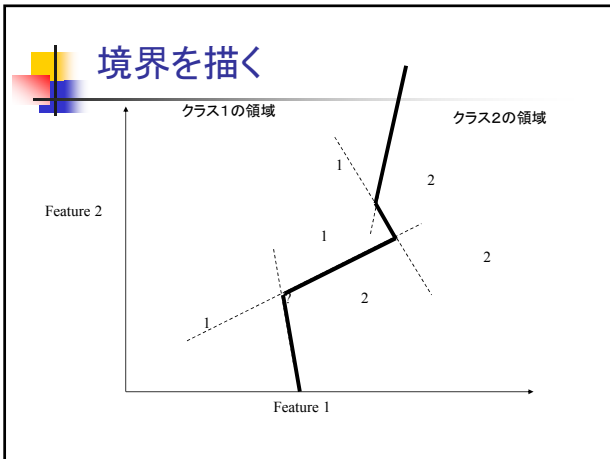
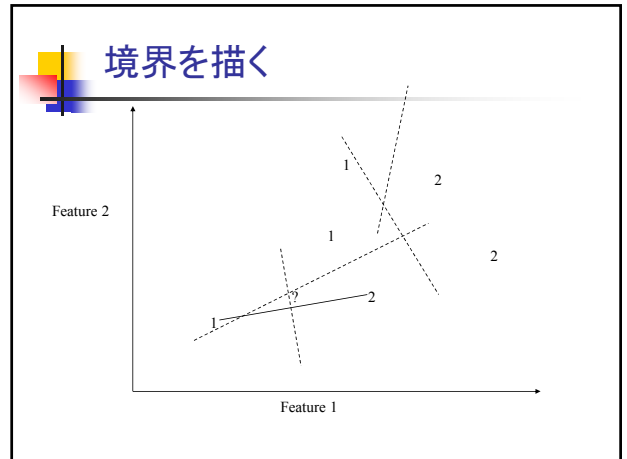
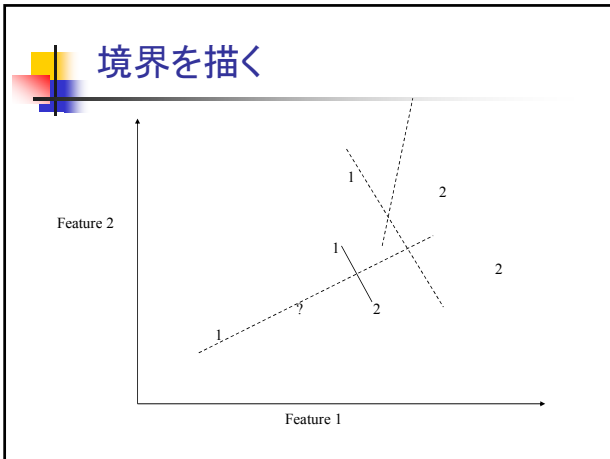


### 境界を描く

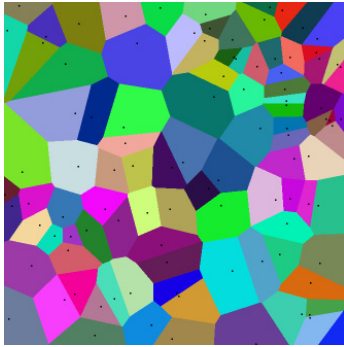


### 境界を描く

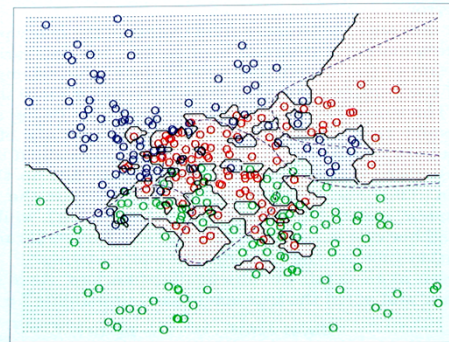




# Voronoi

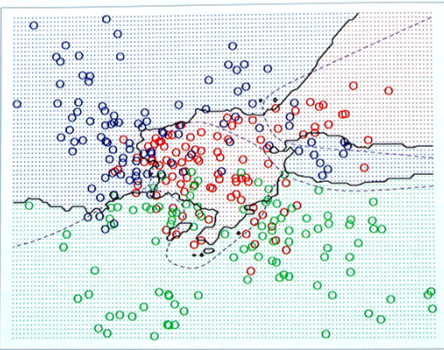


# 1-Nearest Neighbor



From Hastie, Tibshirani, Friedman 2001 p418

# 15-Nearest Neighbors



From Hastie, Tibshirani, Friedman 2001 p418

From Hastie, Tibshirani, Friedman 2001 p419

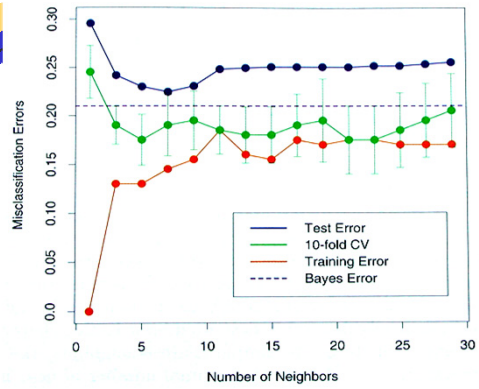


Table 6. Results summary of TC systems on Reuters versions 1-4.

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)
CLASSI (perceptron)	—	—	.80	—
RIPPER (DNF)	—	.72 (Scut)	.80 (Scut)	—
SWAP-1 (DNF)	—	—	.79	—
DTree IND	—	.67 (Pcut)	—	—
DTree C4.5	—	—	.79 ( $F_1$ )	—
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
WORD	linear (naïve)	Yang 1999	189	310	290	752	815
NaiveBayes	probabilistic	[Dumais et al. 1998]	445 (MP <sub>1</sub> )	—	—	—	729
EM	probabilistic	[Joachims 1999]	—	—	—	—	—
Sa	probabilistic	[Lewis 1992]	—	—	—	—	—
C4.5	decision tree	[Li and Yamashita 1999]	—	—	—	—	—
Id3	decision tree	[Yang and Liu 1999]	—	—	—	—	—
Id3	decision tree	[Liu and Elmqvist 1994]	.670	—	—	—	—
Support Vector	decision rule	[Cohen and Singer 1999]	.683	.805	—	.820	—
Support Vector	decision rule	[Cohen and Singer 1999]	.783	.796	—	.820	—
Decision	decision rule	[Li and Yamashita 1999]	—	—	—	—	—
Decision	decision rule	[Moulinier and Guancia 1996]	738	—	—	—	—
Decision	decision rule	[Moulinier et al. 1996]	763 (F <sub>1</sub> )	—	—	—	—
Log	regression	[Yang 1999]	—	.855	.810	—	—
Log	regression	[Yang and Liu 1999]	—	—	—	—	—
BayesNet/Window	on-line linear	[Dagan et al. 1995]	.727	.833	—	.822	—
Window-Map	on-line linear	[Lam and He 1998]	—	—	—	—	—
Boosting	back linear	[Cohen and Singer 1999]	.690	.748	—	.770	.646
Boosting	back linear	[Joachims 1999]	—	—	—	—	—
Boosting	back linear	[Lam and He 1998]	—	—	—	.781	.759
Boosting	back linear	[Li and Yamashita 1999]	—	—	—	—	—
Boosting	back linear	[Yang 1999]	—	—	—	.726	—
Boosting	back linear	[Yang and Liu 1999]	—	—	—	—	—
Boosting	back linear	[Wu et al. 1995]	—	.802	—	—	—
Boosting	back linear	[Yang and Liu 1999]	—	—	—	.820	.838
Boosting	back linear	[Joachims 1999]	—	—	—	—	—
Boosting	back linear	[Lam and He 1998]	—	—	—	.820	.823
Boosting	back linear	[Yang 1999]	.690	.852	.820	.820	.836
Boosting	back linear	[Yang and Liu 1999]	—	—	—	—	—
Boosting	back linear	[Dumais et al. 1998]	—	—	—	.870	.920
Boosting	back linear	[Joachims 1999]	—	—	—	.844	.861
Boosting	back linear	[Yang and Liu 1999]	—	—	—	.839	—
Boosting	back linear	[Cohen and Singer 2000]	—	.860	—	.878	—
Boosting	back linear	[Watts et al. 1999]	—	—	—	.878	—
Boosting	back linear	[Friedman et al. 1998]	—	—	—	.900	.850
Boosting	back linear	[Lam et al. 1997]	.542 (MP <sub>1</sub> )	—	—	—	—

Table 6. Comparative results among different classifiers obtained on five different version of the Reuters collection. Unless otherwise noted, entries indicate the microaveraged break-even point; within parentheses, "M" indicates macroaveraging and "F<sub>1</sub>" indicates use of the F<sub>1</sub> 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 nearest neighbor, circles indicate the macro-averaged breakdown point, with its parentheses. 'M' indicates macro-entropy and 'F' 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	12,902	12,902
		# of training documents	14,704	10,687	8,210	10,670	10,670
		# of test documents	6,746	3,660	3,662	2,269	2,269
		# of categories	135	63	62	106	131
Word	non-linear	Yang (1999)	169	110	207	152	313
PairWise	probabilistic	(Dumais et al. 1998)					
	probabilistic	(Joachims 1998)					729
	probabilistic	(Lee et al. 1997)	443 (M <sub>F</sub> )				747
	probabilistic	(Liu and Yamashita 1999)	650				773
Na	probabilistic	(Lee and Lu 1999)					795
	probabilistic	(Dumais et al. 1998)					884
C1.2	decision tree	(Joachims 1998)					794
	decision tree	(Lee and Bengio 1994)	670				
SVM1	decision rule	(Liu et al. 1999)		805			
	decision rule	(Cohen and Singer 1999)	683	811			820
SupportVector	decision rule	(Cohen and Singer 1999)	753	759			827
	decision rule	(Liu and Yamashita 1999)					829
Chance	decision rule	(Mollinari and Gansner 1998)		718			
	decision rule	(Mollinari et al. 1998)		783 (F)			
Log	regression	Yang (1999)		855	810		
	regression	(Yang and Liu 1999)					849
BALANCEWISER	on-line linear	(Dumais et al. 1997)	747 (M)	843 (M)			822
	on-line linear	(Lee and He 1998)					
Rocchio	batch linear	(Cohen and Singer 1999)	690	745			715
	batch linear	(Dumais et al. 1998)					647
Rocchio	batch linear	(Joachims 1998)					799
	batch linear	(Lee and Lu 1999)					784
Rocchio	batch linear	(Liu and Yamashita 1999)					825
	batch linear	(Liu and Yamashita 1999)					858
Naer	neural network	Yang et al. (1997)		802			
	neural network	(Wang et al. 1995)				820	
GenM	example-based	(Liu and He 1998)					809
	example-based	(Joachims 1998)					825
k-NN	example-based	(Liu and He 1998)					829
	example-based	(Yang 1999)	690	852	820		856
k-NN	example-based	(Yang and Liu 1999)					879
	example-based	(Dumais et al. 1998)					920
SvmLearn	SVM	(Joachims 1998)		864			
	SVM	(Liu and Yamashita 1999)					841
SvmLearn	SVM	(Liu and Lu 1999)					859
	SVM	(Wang et al. 1995)					878
AmbiCoEMH	committee	(Schapire and Singer 2000)		800			
	committee	(Wang et al. 1995)					890
Bayesian net	Bayesian net	(Dumais et al. 1998)					890
	Bayesian net	(Lee et al. 1997)	642 (M <sub>F</sub> )				890

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

## 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

## K-NN と不要な特徴

## K-NN と不要な特徴

## 距離の問題

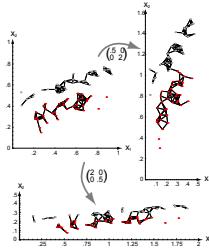
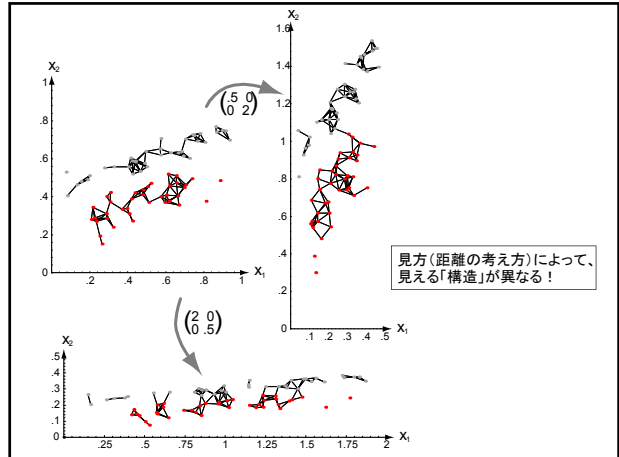


FIGURE 10.8 Scaling affects the clusters in a minimum distance cluster method. The original data and minimum distance clusters are shown in the upper left; points in one cluster are shown in red, while the others are shown in gray. When the vertical axis is expanded by a factor of 2.0 and the horizontal 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 expanded by a factor of 2.0, a smaller number of clusters result (shown at the bottom). In both the rescaled 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

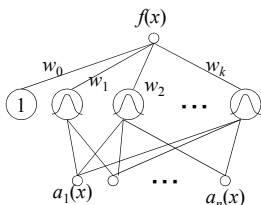
- $k$ -NN は各問合せ $x_q$ で $f$ の局所近似を構成していた
- $x_q$ の周囲で $f(x)$ の近似関数を明示的に構成したらどうだろうか?
  - $k$ -NNに線型回帰したら?
  - 2次回帰では?
  - 区分回帰したら?
- 最小化すべき誤差にもいくつかの候補が

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

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

## 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
  - 事例空間に一樣にばら撒く
  - 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