

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)
 - 問い合わせがあったら、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 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 (Peut)	.29 (Peut)
kNN	—	.69 (Scut)	.85 (Scut)	.82 (Scut)
LLSF	—	—	.85 (Scut)	.81 (Scut)
NNets PARC (perceptron)	—	—	—	.82 (Peut)
CLASSI (perceptron)	—	—	.80	—
RIPPER (DNF)	—	.72 (Scut)	.80 (Scut)	—
SWAP-1 (DNF)	—	—	.79	—
DTree IND	—	.67 (Peut)	—	—
DTree C4.5	—	—	.79 (F_1)	—
CHARADE (DNF)	—	—	.78	—
EXPERTS (n-gram)	—	.75 (Scut)	.76 (Scut)	—
Rocchio	—	.66 (Scut)	.75 (Scut)	—
NaiveBayes	—	.65 (Peut)	.71	—
CONSTRUE (Exp. Sys.)	.90	—	—	—


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

System	Type	Break-even point			
		$\#$ of documents	$\#$ of training documents	$\#$ of test documents	$\#$ of categories
WORD	(naive) averaging	13,450	13,337	13,272	12,502
kNN	probabilistic	14,704	10,667	9,410	9,603
LLSF	probabilistic	14,704	10,667	9,410	9,603
Peut	probabilistic	443 (M^P_1)	400	—	—
CLASSI	probabilistic	14,704	10,667	9,410	9,603
RIPPER	probabilistic	14,704	10,667	9,410	9,603
SWAP-1	probabilistic	14,704	10,667	9,410	9,603
DTree IND	probabilistic	14,704	10,667	9,410	9,603
DTree C4.5	probabilistic	14,704	10,667	9,410	9,603
CHARADE	probabilistic	14,704	10,667	9,410	9,603
EXPERTS	probabilistic	14,704	10,667	9,410	9,603
Rocchio	probabilistic	14,704	10,667	9,410	9,603
NaiveBayes	probabilistic	14,704	10,667	9,410	9,603
CONSTRUE	probabilistic	14,704	10,667	9,410	9,603
WORD	regression	13,450	13,337	13,272	12,502
CHARADE	regression	13,450	13,337	13,272	12,502
NaiveBayes	regression	13,450	13,337	13,272	12,502
WORD	rule learner	13,450	13,337	13,272	12,502
CHARADE	rule learner	13,450	13,337	13,272	12,502
NaiveBayes	rule learner	13,450	13,337	13,272	12,502
WORD	batch linear	13,450	13,337	13,272	12,502
CHARADE	batch linear	13,450	13,337	13,272	12,502
NaiveBayes	batch linear	13,450	13,337	13,272	12,502
WORD	neural network	13,450	13,337	13,272	12,502
CHARADE	neural network	13,450	13,337	13,272	12,502
NaiveBayes	neural network	13,450	13,337	13,272	12,502
WORD	k-NN	13,450	13,337	13,272	12,502
CHARADE	k-NN	13,450	13,337	13,272	12,502
NaiveBayes	k-NN	13,450	13,337	13,272	12,502
WORD	example-based	13,450	13,337	13,272	12,502
CHARADE	example-based	13,450	13,337	13,272	12,502
NaiveBayes	example-based	13,450	13,337	13,272	12,502
WORD	SV-M	13,450	13,337	13,272	12,502
CHARADE	SV-M	13,450	13,337	13,272	12,502
NaiveBayes	SV-M	13,450	13,337	13,272	12,502
WORD	AdaBoost MH	13,450	13,337	13,272	12,502
CHARADE	AdaBoost MH	13,450	13,337	13,272	12,502
NaiveBayes	AdaBoost MH	13,450	13,337	13,272	12,502
WORD	Bayesian net	13,450	13,337	13,272	12,502
CHARADE	Bayesian net	13,450	13,337	13,272	12,502
NaiveBayes	Bayesian net	13,450	13,337	13,272	12,502
WORD	Bayesian net	13,450	13,337	13,272	12,502
CHARADE	Bayesian net	13,450	13,337	13,272	12,502
NaiveBayes	Bayesian net	13,450	13,337	13,272	12,502
WORD	Bayesian net	13,450	13,337	13,272	12,502
CHARADE	Bayesian net	13,450	13,337	13,272	12,502
NaiveBayes	Bayesian net	13,450	13,337	13,272	12,502
WORD	Bayesian net	13,450	13,337	13,272	12,502
CHARADE	Bayesian net	13,450	13,337	13,272	12,502
NaiveBayes	Bayesian net	13,450	13,337	13,272	12,502
WORD	Bayesian net	13,450	13,337	13,272	12,502
CHARADE	Bayesian net	13,450	13,337	13,272	12,502
NaiveBayes	Bayesian net	13,450	13,337	13,272	12,502
WORD	Bayesian net	13,450	13,337	13,272	12,502
CHARADE	Bayesian net	13,450	13,337	13,272	12,502
NaiveBayes	Bayesian net	13,450	13,337	13,272	12,502
WORD	Bayesian net	13,450	13,337	13,272	12,502
CHARADE	Bayesian net	13,450	13,337	13,272	12,502
NaiveBayes	Bayesian net	13,450	13,337	13,272	12,502
WORD	Bayesian net	13,450	13,337	13,272	12,502
CHARADE	Bayesian net	13,450	13,337	13,272	12,502
NaiveBayes	Bayesian net	13,450	13,337	13,272	12,502
WORD	Bayesian net	13,450	13,337	13,272	12,502
CHARADE	Bayesian net	13,450	13,337	13,272	12,502
NaiveBayes	Bayesian net	13,450	13,337	13,272	12,502
WORD	Bayesian net	13,450	13,337	13,272	12,502
CHARADE	Bayesian net	13,450	13,337	13,272	12,502
NaiveBayes	Bayesian net	13,450	13,337	13,272	12,502
WORD	Bayesian net	13,450	13,337	13,272	12,502
CHARADE	Bayesian net	13,450	13,337	13,272	12,502
NaiveBayes	Bayesian net	13,450	13,337	13,272	12,502
WORD	Bayesian net	13,450	13,337	13,272	12,502
CHARADE	Bayesian net	13,450	13,337	13,272	12,502
NaiveBayes	Bayesian net	13,450	13,337	13,272	12,502
WORD	Bayesian net	13,450	13,337	13,272	12,502
CHARADE	Bayesian net	13,450	13,337	13,272	12,502
NaiveBayes	Bayesian net	13,450	13,337	13,272	12,502
WORD	Bayesian net	13,450	13,337	13,272	12,502
CHARADE	Bayesian net	13,450	13,337	13,272	12,502
NaiveBayes	Bayesian net	13,450	13,337	13,272	12,502
WORD	Bayesian net	13,450	13,337	13,272	12,502
CHARADE	Bayesian net	13,450	13,337	13,272	12,502
NaiveBayes	Bayesian net	13,450	13,337	13,272	12,502
WORD	Bayesian net	13,450	13,337	13,272	12,502
CHARADE	Bayesian net	13,450	13,337	13,272	12,502
NaiveBayes	Bayesian net	13,450	13,337	13,272	12,502
WORD	Bayesian net	13,450	13,337	13,272	12,502
CHARADE	Bayesian net	13,450	13,337	13,272	12,502
NaiveBayes	Bayesian net	13,450	13,337	13,272	12,502
WORD	Bayesian net	13,450	13,337	13,272	12,502
CHARADE	Bayesian net	13,450	13,337	13,272	12,502
NaiveBayes	Bayesian net	13,450	13,337	13,272	12,502
WORD	Bayesian net	13,450	13,337	13,272	12,502
CHARADE	Bayesian net	13,450	13,337	13,272	12,502
NaiveBayes	Bayesian net	13,450	13,337	13,272	12,502
WORD	Bayesian net	13,450	13,337	13,272	12,502
CHARADE	Bayesian net	13,450	13,337	13,272	12,502
NaiveBayes	Bayesian net	13,450	13,337	13,272	12,502
WORD	Bayesian net	13,450	13,337	13,272	12,502
CHARADE	Bayesian net	13,450	13,337	13,272	12,502
NaiveBayes	Bayesian net	13,450	13,337	13,272	12,502
WORD	Bayesian net	13,450	13,337	13,272	12,502
CHARADE	Bayesian net	13,450	13,337	13,272	12,502
NaiveBayes	Bayesian net	13,450	13,337	13,272	12,502
WORD	Bayesian net	13,450	13,337	13,272	12,502
CHARADE	Bayesian net	13,450	13,337	13,272	12,502
NaiveBayes	Bayesian net	13,450	13,337	13,272	12,502
WORD	Bayesian net	13,450	13,337	13,272	12,502
CHARADE	Bayesian net	13,450	13,337	13,272	12,502
NaiveBayes	Bayesian net	13,450	13,337	13,272	12,502
WORD	Bayesian net	13,450	13,337	13,272	12,502
CHARADE	Bayesian net	13,450	13,337	13,272	12,502
NaiveBayes	Bayesian net	13,450	13,337	13,272	12,502
WORD	Bayesian net	13,450	13,337	13,272	12,502
CHARADE	Bayesian net	13,450	13,337	13,272	12,502
NaiveBayes	Bayesian net	13,450	13,337	13,272	12,502
WORD	Bayesian net	13,450	13,337	13,272	12,502
CHARADE	Bayesian net	13,450	13,337	13,272	12,502
NaiveBayes	Bayesian net	13,450	13,337	13,272	12,502
WORD	Bayesian net	13,450	13,337	13,272	12,502
CHARADE	Bayesian net	13,450	13,337	13,272	12,502
NaiveBayes	Bayesian net	13,450	13,337	13,272	12,502
WORD	Bayesian net	13,450	13,337	13,272	12,502
CHARADE	Bayesian net	13,450	13,337	13,272	12,502
NaiveBayes	Bayesian net	13,450	13,337	13,272	12,502
WORD	Bayesian net	13,450	13,337	13,272	12,502
CHARADE	Bayesian net	13,450	13,337	13,272	12,502
NaiveBayes	Bayesian net	13,450	13,337	13,272	12,502
WORD	Bayesian net	13,450	13,337	13,272	12,502
CHARADE	Bayesian net	13,450	13,337	13,272	12,502
NaiveBayes	Bayesian net	13,450	13,337	13,272	12,502
WORD	Bayesian net	13,450	13,337	13,272	12,502
CHARADE	Bayesian net	13,450	13,337	13,272	12,502
NaiveBayes	Bayesian net	13,450	13,337	13,272	12,502
WORD	Bayesian net	13,450	13,337	13,272	12,502
CHARADE	Bayesian net	13,450	13,337	13,272	12,502
NaiveBayes	Bayesian net	13,450	13,337	13,272	12,502
WORD	Bayesian net	13,450	13,337	13,272	12,502
CHARADE	Bayesian net	13,450	13,337	13,272	12,502
NaiveBayes	Bayesian net	13,450	13,337	13,272	12,502
WORD					

Table VI. Comparative Results Among Different Classifiers Obtained on Five Different Versions of Reuter. (Unless otherwise noted, entries indicate the macroaveraged break-even point; within parentheses, "M" indicates macroaveraging and " F_1 " indicates use of the F_1 measure; boldface indicates the best performer on the collection.)

復習

Bayes 最適な分類器



注: Bayes 最適な分類器は H に含まれるとは限らない
注: 論文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
 - 新事例をこれに従い分類する

慶賀:もしも仮説を事前分布 $P(h)$ 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” を参照）

K-NN と不要な特徴



(1)

極限における振舞い

- $p(x)$: 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


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⇒Shepardの方法

K-NN と不要な特徴



K-NN と不要な特徴



距離の問題





FIGURE 9.9 Scaling of axes changes minimum distances. Data and cluster centers 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 clusters of points shown at the right. Alternately, if the vertical axis is shrunk by a factor of 0.5 and the horizontal axis expanded by a factor of 2.0, smaller more fine-grained clusters result (shown at the bottom). In both these cases the assignment of points to clusters differ 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

$$E_1(x_q) \equiv \frac{1}{2} \sum_{x \in N_q} (f(x) - \hat{f}(x_q))^2$$

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

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



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でいる。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