

Rish (2001)

①

Empirical Study of Naive Bayes

Independence is generally a poor assumption - Still NB competes well

Bayes Classifier assigns most likely class based on a data instance feature vector

Learning is much simplified by assuming that features are independent per class

$$P(X|C_i) = \prod_{i=1}^n P(X_i|C_i)$$
$$X = (X_1, \dots, X_n), C = \text{Class}$$

Success of NB despite poor assumption =
Optimal classification error is not necessarily related to Quality of fit of distributions
(approx & independence assumption)

Rather the optimal 

(2)

Rather, optimal classifier =
both actual & estimated
distributions agree on the most
probable class

This is vague & the paper
wants to understand the
data characteristic that affect
naive Bayes

MOST studies of NB uses Benchmark
comparison to other classifiers

This paper uses Monte Carlo sim to
allow more systematic study of
classification errors

- on parametric families of random
~~gener~~ generated problems

- Focuses on bias not variance
- by assuming infinite data
- Allows separation of approx error (bias)
& error from sample variance (variance)

"Distribution entropy"

Quantifies the amount of uncertainty associated w/ the distib of a variable

low-entropy dependances

- Attribs are related to each other
- esp strong when you focus on data points that belong to the same class
- these related attributes are therefore ~~related~~ predictable within each class

(relationship)

With naive Bayes, if the dependances are low-entropy (predictable), the classifier works well.

This might be considered to be contradict given the independence assum/structure

error of NB vanishes as entropy approaches 0

4

Bayes works best under two Easterns

- Completely unrelated Independent feature - which we would expect
- functionally Dependant feature-
 - i.e. one feature can be used to deterministicly predict another
 - $y \propto X^k$ then Y
 - conversion cm to mm
- it performs worst in between
- This would be where features are related by not perfectly
- often conditionally related based on other feature

e.g. Weather, temp change effect Humid — but same ΔT does not always cause same change in Humid as other factors @ play

(5)

Also

Another way of stating these finds is that the accuracy of NB is not directly correlated w/ the degree of feature dependencies measured as the class-conditioned mutual information between the features.

- Mutual info = How attrs together provide info about class

Better pred of NB accuracy can be the loss of info that features contain. given as:

$$I(C; X_i, X_j) - \underbrace{I_{NB}(C; X_i, X_j)}_{\text{mutual info}}$$

2 Definitions and Backgrounds

Discriminant function → Select highest class

$$h(x) = \operatorname{argmax}_i \in \{0, \dots, M-1\} f_i(x)$$

$$\text{Bayes Disc} = f_i(x) = P(X=x | c_i) P(c_i)$$

~~Not~~

Class conditional

Probability Distribution

Likelihood

feature space

$P(X=x | c_i)$ is hard to estimate high-Dims

↳ why Approximations are used →

simplify assume to indep features

yields the NB:

$$f_i^{\text{NB}} = \prod_{j=1}^n P(X_j=a_j | c_i) P(c=i)$$