

(1)

Naive Bayes Classifier

The naive Bayes approach makes the assumption that all attributes are independent

This assumption means that the likelihood can be decomposed into a product of dimension-wise probabilities:

$$P(x|c_i) = P(x_1, x_2, \dots, x_d | c_i) = \prod_{j=1}^d P(x_j | c_i)$$

How is this different?

- A full Bayes tries to model the joint Probability of all the features
- This means it takes into account the dependancies & correlations between features
- This is what makes Bayes complex & costly
- in theory with enough Data it is more accurate as it captures true relationships

(2)

How is Naive Bayes different?

- features are conditionally indep given class
- that is the value of one does not affect the value of another, given C_i
- This drastically reduces complexity
- Instead of est a complex joint distrib
- Naive Bayes only est marginal probs of each class feature given class
- Computationally efficient, reg less train
- Generally still performs well partic in domains where independ assumption is valid

(3)

Naive Bayes: Numeric Attrbs

Default assumption that each is normally distributed

Let μ_i & σ^2 be for attribute x_j for class c_i

Note here x_j is a single attribute but for all attributes

The likelihood for c_i given dimension x_j

$$P(x_j | c_i) \propto f(x_j | \mu_{ij}, \sigma_{ij}^2) = \frac{1}{\sqrt{2\pi \sigma_{ij}^2}} \exp \left\{ -\frac{(x_j - \mu_{ij})^2}{2\sigma_{ij}^2} \right\}$$

α = proportionality constant

Prob Dens function

tells us that the ~~conditional~~ probability is proportional to the gaussian PDF upto a constant factor

Why? Bayes claim we often care more about the relative probabilities of a data point belong to different classes rather than the exact prob

(4)

Addition consequence of the independence assumption corresponds to setting all covariances to 0 - (non-diags) Σ_i

$$\Sigma_i = \begin{pmatrix} \sigma_{ij}^2 & 0 & \dots & 0 \\ 0 & \sigma_{i2}^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \\ 0 & 0 & \dots & \sigma_{id}^2 \end{pmatrix}$$

Indep so
no cross
interacts

Determinant of Σ_i : $|I\Sigma_i|$

- the determinant of a matrix is a scalar value that provides info about mat
- for diag mats, it is the product of the diag elements
- This is one of the factors that makes the calc a lot quicker overall

$$|I\Sigma_i| = \det(\Sigma_i) = \sigma_{i1}^2 \cdot \sigma_{i2}^2 \cdot \dots \cdot \sigma_{id}^2 = \prod_{j=1}^d \sigma_{ij}^2$$

Product
of all varB

(5)

Inverse Matrix of covariance matrix

Σ_i^{-1} = Diag Matrix where all values are $\frac{1}{\sigma^2}$

We want the inverse because we can plug it into the multivariate Gaussian PDF (likelihood)

$$(x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i) = \sum_{j=1}^d \frac{(x_j - \mu_{ij})^2}{\sigma_{ij}^2}$$

(RHS of PDF)

use of ~~clamies~~ allows to account for all data.

~~Rows of a attr~~ Attrs

Note a univariate PDF only looks@ a single attribute

plugging values into PDF for num attrs



(6)

this is a univariate
PDF

$$\textcircled{1} \quad f(x) = \frac{1}{(\sqrt{2\pi})^d \sqrt{|\Sigma_i|}} \exp \left\{ -\frac{(x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i)}{2} \right\}$$

$\prod_{j=1}^d \sigma_{ij}^2$ $\sum_{j=1}^d \frac{(x_j - \mu_{ij})^2}{\sigma_{ij}^2}$

\textcircled{2} Plugging gives us:

$$P(x|c_i) = \frac{1}{(\sqrt{2\pi})^d \sqrt{\prod_{j=1}^d \sigma_{ij}^2}} \exp \left\{ -\sum_{j=1}^d \frac{(x_j - \mu_{ij})^2}{2\sigma_{ij}^2} \right\}$$

2 remains as part of gaussian dist

\textcircled{3} use distributions to take \prod out of Brackets

~~$$\prod_{j=1}^d \left(\frac{1}{\sqrt{2\pi} \sigma_{ij}} \exp \left\{ -\frac{(x_j - \mu_{ij})^2}{2\sigma_{ij}^2} \right\} \right)$$~~

$$\textcircled{4} \quad \prod_{j=1}^d P(x_j | c_i)$$

final form is the product of the likelihoods

(7)

To summarise what we did here

the joint probability has been decomposed into a product of the probability along each dimension

\hookrightarrow caused by indep assumption

Naive Bayes uses the sample mean:

$$\hat{\mu}_i = (\hat{\mu}_{i1}, \dots, \hat{\mu}_{id})^T$$

Diag Sample covariance mat $= \hat{\Sigma}_i = \text{diag}(\hat{\sigma}_{i1}^2, \dots, \hat{\sigma}_{id}^2)$
for each class i

Thus the complexity is only ever
 ~~$O(n^d)$~~ $O(nd)$

Pred:

$$\hat{y} = \arg\max \left\{ P(c_i) \prod f(x_i | \hat{\mu}_{ij}, \hat{\sigma}_{ij}^2) \right\}$$

Prior

Likelihood
Multivariate

②

Naive Bayes: Categorical Attrb

Indep assumption leads to a simplification of the joint prob mass function

$$\text{Bayes: } P(x|c_i) = f(v|c_i) = r(X_1 = e_1, \dots, X_d = e_d | c_i)$$

Recall the func is just $\frac{n_i(v)}{n_i}$

this can be rewritten as:

$v = \text{single}$
 not vect

$$P(x|c_i) = \prod_{j=1}^d P(x_j|c_i) = \prod_{j=1}^d f(X_j = e_{j,r_j} | c_i)$$

this prob makes
function for X_j

As w/ bayes we can use pseudo counts if the counts are zero:

$$\hat{f}(v_j | c_i) = \frac{n_i(v_j) + 1}{n_i + n_j}$$