

① Naive Bayes classifiers are collection of classification algorithm based on Bayes Theorem. When we have large training set available then it is moderate to use. Attributes that describes instances are conditionally independent given classification.

Successful Applications &

- Classifying Text Documents
- Diagnosis

Assume target function $f: X \rightarrow V$, where each instance x described by attributes (a_1, a_2, \dots, a_n) . Most probable value of $f(x)$ is:

$$v_{map} = \underset{v_j \in V}{\operatorname{argmax}} P(v_j | a_1, a_2, \dots, a_n)$$

$$= \frac{P(a_1, a_2, \dots, a_n | v_j) P(v_j)}{P(a_1, a_2, \dots, a_n)}$$

$$= P(a_1, a_2, \dots, a_n | v_j) P(v_j)$$

Naive Bayes Assumption which gives

$$P(a_1, a_2, \dots, a_n | v_j) = \prod_i P(a_i | v_j)$$

$$v_{NB} = P(v_j) \prod_i P(a_i | v_j)$$

Algorithm &

N. Bayes (examples)

For each target value v_j

$\hat{P}(v_j) \leftarrow \text{estimate } P(v_j)$

For each attribute value a_i of each attribute a

$\hat{P}(a_i | v_j) \leftarrow \text{estimate } P(a_i | v_j)$

Classify New Instance (n)

$$v_{na} = P(v_i) \prod P(a_i | v_i)$$

Example & Consider Play Tennis

Outlook = sun, Temp = cool, Humid = high, Wind = strong

$$v_i = P(v_i) \prod P(a_i | v_i)$$

$$P(y) P(\text{sun} | y) P(\text{cool} | y) P(\text{high} | y) P(\text{strong} | y) = 0.005$$

$$P(n) P(\text{sun} | n) P(\text{cool} | n) P(\text{high} | n) P(\text{strong} | n) = 0.021$$

i) Conditional independence assumption is often violated $P(a_1, a_2, \dots, a_n | v_i) = \prod P(a_i | v_i)$

• It works surprisingly well anyway. Note don't need estimated posteriors $P(v_i | n)$ to be correct, need only that.

$$\hat{P}(v_i) = \prod \hat{P}(a_i | v_i) = P(v_i) P(a_1, \dots, a_n | v_i)$$

• Naive Bayes posteriors often unrealistically close to 1 or 0.

ii) What if none of the training instance with target value v_j have attribute value a_i ? then

$$\hat{P}(a_i | v_j) = 0$$

$$\hat{P}(v_j) \prod \hat{P}(a_i | v_j) = 0$$

Bayesian Estimation for $\hat{P}(a_i | v_j)$

$$\hat{P}(a_i | v_j) \leftarrow \frac{n_{a_i} + m}{n + m}$$

where, n is no. of training example for $v = v_j$

n_{a_i} " " " examples for $v = v_j$ & $a = a_i$

p is prior estimate for $\hat{P}(a_i | v_j)$

m is weight given to prior