

①

Naïve Bayes Classifier

All previous formulations of Classification =

given data $\{x_{j1}, x_{j2}, \dots, x_{jd}\}, \forall d, d=1, 2, \dots, n$

Predict $\hat{y}_j : \hat{y}_j \in \{c_k, \forall k, k=1, 2, \dots, K\}$

Alternative but equiv formulation:

given data ...

what is the prob of a class given inputs

$p(c_k | x_{j1}, x_{j2}, \dots, x_{jd}), \forall k, k=1, 2, \dots, K$

$\hat{y}_j = c_k$ for $\operatorname{argmax} p(c_k | x_{j1}, x_{j2}, \dots, x_{jd})$

~~Bayes Rules:~~

~~for two events A & B~~

$$\cancel{P(A \cap B) = P(A|B) P(B)} \\ \cancel{\qquad \qquad \qquad P(A)}$$

(2)

Bayes Rule
for two events A & B

$$P(B|A) = \frac{P(A|B)P(B)}{P(A)}$$

$P(B|A)$ = cond prob of seeing B given
observation of A

$P(A)$
 $P(B)$ = Prior probability of B or ~~B~~

$P(A|B)$ = Prob of A occurring given obs of B

Probabilistic Model formulation
Rewrite Bayes rule for a set of possible classes

$$P(Y=c_k | X) = \frac{P(X|Y=c_k)P(Y=c_k)}{P(X)}$$

$P(X)$ can be ignored
as it is the same for all classes

Same as:

$$P(Y=c_k | X) = \frac{P(X_1=x_{1,1}, X_2=x_{2,2}, \dots, X_D=x_{D,p} | Y=c_k)}{P(X=c_k)}$$

(3)

BAYES Classifier

$$P(C_k | X) = P(X_1, X_2, \dots, X_D | C_k) P(C_k)$$

$P(C_k)$ can be estimated from the training data as $\frac{N_{C_k}}{N} \rightarrow$

- $N_{C_k} \rightarrow$ No. of C_k instances
- $N \rightarrow$ n of instance

In eng: the prob of seeing class C_k given the parameters, against the overall prob of seeing the class

$P(X_D | C_k)$ can be estimated from the training data given a distribution e.g gaussian $N(\mu_{C_k}, \Sigma_{C_k})$

Probability of X given the parameters distribution

$$P(X_D | C_k) = E[P(X | \mu_{C_k}, \Sigma_{C_k})]$$

μ_{C_k} = D-Dimension mean vector mean

Σ_{C_k} = DxD covariance SD

$E[\cdot]$ = Expected Value

Find Optimal Params

(4)

Assuming a gaussian distribution $N(\mu_{ck}, \Sigma_{ck})$ the prob dist params are μ_{ck}, Σ_{ck}

The optimal dist params max the k lik. hood w.r.t the train data = $E[P(x | \mu_{ck}, \Sigma_{ck})]$

Recall that opt value of param θ for a function $L = \frac{\partial L}{\partial \theta} = 0$

Log likelihood derivation

$P(x | \mu_{ck}, \Sigma_{ck})$ = Plug in gaussian prob Dens func

(1) take log of both sides

(2) log of product = log of sum

(3) apply log rules

finally Max/opt likelihood of X given params μ_{ck}, Σ_{ck}

$$\frac{\partial}{\partial \mu} \log P(x | \mu_{ck}, \Sigma_{ck}) = 0$$

- Sub in log derivative $\log P(x | \mu_{ck}, \Sigma_{ck})$
- Remove non-linear terms as per Diff rule
- Expand & Simplify again

(5)

- After full simplification, make U_k the subject

$$U_{ck} = \frac{1}{N_{ck}} \sum_{n=1}^{N_{ck}} x_n$$

↳ Ends up being standard mean formula

- Next optimize $\Sigma_{ck} = \frac{\partial^2}{\partial \Sigma} \log P(X|U_{ck}, \Sigma_{ck}) = 0$
- Plug in logs, expand, simplify, set to 0

$$\Sigma_{ck} = \frac{1}{N_{ck}} E((x_n - U_{ck})(x_n - U_{ck})^T)$$

↳ Ends up as standard variance formula

(6)

Now we have optimal μ & Σ values

$P(X|U_{ck}, \Sigma_{ck})$ = Prob Dens func plug in

$$U_{ck} = \frac{1}{N_{ck}} \sum x_n$$

$$\Sigma_{ck} = \frac{1}{N_{ck}} \sum (x_n - U_{ck})^T (x_n - U_{ck})$$

Plug optimal values into PDF = Bayes Class

Naive Bayes classifier assumes input params are independent so can be isolated

Summary

Bayes Classifier gives the posterior probability of each class by updating the prior probability of the class with the likelihood of the observed data