

(1)

Alpaydin

Bayesian Decision Theory.

3/26/2008

Note Title
Chp 3.

Probability theory is the framework for making decisions under uncertainty.

For classification, Bayes' rule is used to calculate the probabilities of classes.

More generally, make decisions by minimizing the expected risk.

Assume data is generated by a random process - or by a deterministic process that we only partially know.

X is a random variable $\begin{cases} X=1 & \text{heads} \\ X=0 & \text{tails.} \end{cases}$

$$P(X=1) = p_0 \quad \& \quad P(X=0) = 1 - P(X=1) = 1 - p_0.$$

If we know $P(x)$.

Then we can predict:

i.e if $p_0 > \frac{1}{2}$, predict heads
otherwise, predict tails

If we do not know $p(x)$, then learn it from training samples

$$\hat{p}_0 = \frac{\# \text{ tosses which are heads}}{\# \text{ tosses}}$$

(2)

Classification

Credit scoring - is the customer a credible risk?

$C=1$ high-risk customer

$C=0$ low-risk customer.

$$X = (x_1, x_2)$$

x_1 - income

x_2 - savings

Suppose we know $P(C|X_1, X_2)$

and get a new customer with (x_1, x_2)

choose $\begin{cases} C=1, & \text{if } P(C=1|x_1, x_2) > 0.5 \\ C=0, & \text{otherwise} \end{cases}$

Equivalently

$$\begin{cases} C=1, & \text{if } P(C=1|x_1, x_2) > P(C=0|x_1, x_2) \\ C=0, & \text{otherwise} \end{cases}$$

Prob of error = $1 - \max\{P(C=1|x_1, x_2), P(C=0|x_1, x_2)\}$

How to calculate $P(C|x)$? $x = (x_1, x_2)$

Bayes Rule $P(C|x) = \frac{p(C)p(x|C)}{p(x)}$

$p(C)$ - prior probability (before observing data)

$p(x|C)$ - class likelihood.

$p(x)$ ~ evidence

(3)

Classification (cont)

$$\text{posterior} = \frac{\text{prior} \times \text{likelihood}}{\text{evidence}}$$

Now, we assume we know the prior and likelihood (later in the course, we will learn it)

In general case, with K classes $C_i : i=1\dots K$,

$$\text{prior}, \quad P(C_i) \geq 0, \quad \sum_{i=1}^K P(C_i) = 1$$

$$\text{likelihood}, \quad P(x|C_i),$$

$$\text{posterior}, \quad P(C_i|x) = \frac{P(x|C_i)P(C_i)}{P(x)}$$

Bayes' classifier pick the class with the highest posterior probability:

$$\text{choose } C_i : \text{if } P(C_i|x) = \max_k P(C_k|x)$$

(4)

Losses & Risks

But decisions may not be equally good or costly.
Need to allow for different gain & loss.

Action α_i is decision to classify input to C_i .

γ_{ik} loss incurred for taking action
 α_i if class is C_k .

Expected risk:

$$R(\alpha_i | x) = \sum_{k=1}^K \gamma_{ik} P(C_k | x)$$

Choose action with minimal risk:

choose α_i if $R(\alpha_i | x) = \min_k R(\alpha_k | x)$.

K actions, α_i s: $i = 1, \dots, K$ α_i assigns to C_i
zero-one loss $\gamma_{ik} = \begin{cases} 0, & \text{if } i=k \\ 1, & \text{if } i \neq k. \end{cases}$

$$R(\alpha_i | x) = \sum_{k=1}^K \gamma_{ik} P(C_k | x)$$

$$= \sum_{k \neq i} P(C_k | x) = 1 - P(C_i | x)$$

To minimize risk (zero-one loss) we choose
the most probable class $P(C_i | x) = \max_k P(C_k | x)$

(5) In some applications - e.g. medical diagnosis wrong decisions may have a very high cost - and we may require a more complex decision.

New action - reject/doubt α_{K+1} .

$$x_{ik} = \begin{cases} 0, & \text{if } i=k \\ \lambda, & \text{if } i=K+1 \\ 1, & \text{otherwise} \end{cases} \quad 0 < \lambda < 1$$

Risk of reject

$$R(\alpha_{K+1}|x) = \sum_{k=1}^K \lambda P(C_k|x) = \lambda$$

Risk of class i

$$R(\alpha_i|x) = \sum_{k \neq i}^K P(C_k|x) = 1 - P(C_i|x).$$

Optimal Decision Rule:

choose C_i if $R(\alpha_i|x) < R(\alpha_k|x) \forall k \neq i$
or $R(\alpha_i|x) < R(\alpha_{K+1}|x)$.

reject if $R(\alpha_{K+1}|x) < R(\alpha_i|x), i=1, K$

Given simple loss function.

choose C_i if $P(C_i|x) > P(C_k|x), \forall k \neq i$
 $P(C_i|x) > 1 - \lambda$

reject, otherwise.

|| 0 ($\lambda < 1$), if $\lambda = 0$, always reject ||
 $\lambda = 1$, never reject. ||

(6)

Discriminant Functions.

Classification can be seen as implementing a set of discriminant functions. $g_i(\underline{x})$, $i=1, \dots, K$ s.t.
choose C_i , if $g_i(\underline{x}) = \max_k g_k(\underline{x})$

$$g_i(\underline{x}) = -R(\alpha_i | \underline{x})$$

maximum discriminant function corresponds to minimum risk.

with zero-one loss function,

$$g_i(\underline{x}) = P(C_i | \underline{x})$$

or $g_i(\underline{x}) = p(\underline{x} | C_i) P(C_i)$

ignoring normalization (constant) term $P(x)$.

These divide the feature space into K decision regions R_1, \dots, R_K

$$R_i := \{ \underline{x} \mid g_i(\underline{x}) = \max_k g_k(\underline{x}) \}$$

separated by decision boundaries.

(where discriminants are tied).

For two classes,

$$\text{single discriminant } g(\underline{x}) = g_1(\underline{x}) - g_2(\underline{x})$$

$$\text{choose } \begin{cases} C_1 & \text{if } g(\underline{x}) > 0 \\ C_2 & \text{otherwise} \end{cases}$$

7)

Utility Theory

Optimistic version of decision theory.

Utility function U_{ik} .

$$\text{Expected Utility } EU(x_i | \underline{x}) = \sum_k U_{ik} P(S_k | \underline{x})$$

States S_k

Maximizing utility - or minimize risk

Value of Information

Medical diagnosis - many tests can be applied, but some are expensive
measuring pulse is cheap.
blood test is costly.

But blood test may give more information.

If we observe \underline{x} , then expected utility is

$$EU(\underline{x}) = \max_i \sum_k U_{ik} P(S_k | \underline{x})$$

If we observe new feature z

$$EU(\underline{x}, z) = \max_i \sum_k U_{ik} P(S_k | \underline{x}, z)$$

If $EU(\underline{x}, z) > EU(\underline{x})$, then z is useful.

Difference is the value of information-

(But need some way to estimate the value of z).

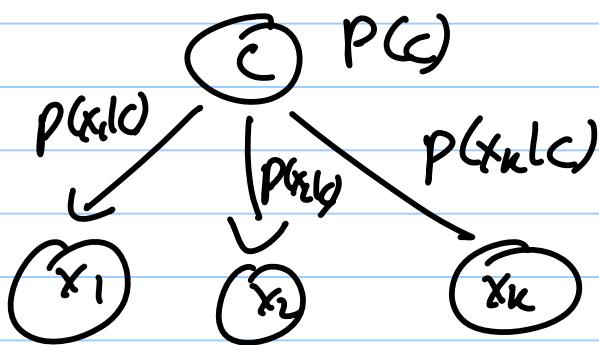
(8)

Bayesian Networks – optional topic

$$p(x_1 \dots x_n) = \prod_i p(x_i | X_{\text{par}})$$

X_{par}

Naive Bayes classifier (ignores dependencies)



$$p(x|C) = \prod_{j=1}^k p(x_j|C)$$

Usually there is more structure

Influence Diagrams

Association Rules:

$$\text{Confidence } (X \rightarrow Y) = P(Y|X) = \frac{P(X, Y)}{P(X)}$$

Support of association rule

$$\text{Support}(X, Y) = P(X, Y)$$

