# Bayesian Classification Scheme

A statistical approach leads to automatic identification of classes.

*Ames Research Center, Moffett Field, California*

A scheme derived via a statistical approach automatically identifies classes in a set of data. The scheme is applicable to real-valued, discrete (e.g., binary) or continuous data in the form of parameter vectors that represent the attributes of objects. It determines the most probable number of classes, the probabilistic descriptions of the classes, and the probability that each object is a member of each class. It incorporates partial prior knowledge of some classes that may be present, striking a balance between the tendency to put each datum in a class by itself and the tendency to put all data in the same class.

The scheme is based on the theory of finite mixtures, in which each datum in a set of $I$ data is assumed to be drawn from one of $J$ classes. The $j$th class, $C_j$, is described by a class distribution function, $p(x_i | x_i \in C_j, \theta_j)$, which gives the probability distribution of the parameters of a datum (components of a vector) that belongs to that class. The $j$th class distribution is described by a class parameter vector, $\theta_j$. For a single-attribute normal distribution, $\theta_j$ consists of the class mean, $\mu_j$, and variance, $\sigma_j^2$.

The probability of an object being drawn from class $j$ is called the class probability $\pi_j$. Thus, the probability of a given datum coming from a set of classes is the sum of the probabilities that it came from each class separately, weighted by the class probabilities.

$$p(x_i | \theta, \pi, J) = \sum_{j=1}^{J} \pi_j \, p(x_i | x_i \in C_j, \theta_j) \quad \text{(a)}$$

Provided that the data are unordered and independent of each other, the likelihood of measuring an entire data base is the product of the probabilities of measuring each datum.

*ive !*

$$p(x | \theta, \pi, J) = \prod_{i=1}^{I} p(x_i | \theta, \pi, J) \quad \text{(b)}$$

For given values of the class parameters, one can calculate the probability that object $i$ belongs to class $j$ by use of Bayes's theorem:

*add*

$$p(x_i \in C_j | x_i, \theta, \pi, J) = \frac{\pi_j \, p(x_i | x_i \in C_j, \theta_j)}{p(x_i | \theta, \pi, J)} \quad \text{(c)}$$

The problem of identifying a finite mixture is broken into two parts: determining the classification parameters for a given number of classes and determining the number of classes. Rather than seek the class parameter vectors, $\theta$, and the class probabilities, $\pi$, one seeks the full posterior probability distribution of them. The posterior distribution is proportional to the product of the prior distribution of the parameters $p(\theta, \pi | J)$ and the likelihood function $p(x | \theta, \pi, J)$.

$$p(\theta, \pi | x, J) = \frac{p(\theta, \pi | J) p(x | \theta, \pi, J)}{p(x | J)} \quad \text{(d)}$$

The pseudolikelihood $p(x | J)$ is simply the normalizing constant of the posterior distribution, obtained by integrating out the classification parameters — in effect, treating them as "nuisance" parameters:

$$p(x | J) = \iint p(\theta, \pi | J) p(x | \theta, \pi, J) d\theta \, d\pi \quad \text{(e)}$$

To solve the second half of the classification problem (determining the number of classes), one calculates the posterior distribution of the number of classes $J$. This is proportional to the product of the prior distribution $p(J)$ and the pseudolikelihood function $p(x | J)$.

$$p(J | x) = \frac{p(J) p(x | J)}{p(x)} \quad \text{(f)}$$

In principle, one can determine the most probable number of classes by evaluating $p(J | x)$ over the range of $J$ for which the prior $p(J)$ is significant. In practice, the multidimensional integrals of equation (e) are computationally intractable, and one must search for the maximum of the function and approximate it about that point. The search and approximate integration are performed by the AutoClass algorithm. The complete problem involves starting with more classes than are believed to be present (as specified by the user), searching to find the best class parameters for that number of classes, approximating the integral to find the relative probability of that number of classes, and then decreasing the number of classes and repeating the procedure. AutoClass uses a Bayesian variant of Dempster's and Laird's EM algorithm to find the best class parameters for a given number of classes [the maximum of equation (d)]. To derive the algorithm, one differentiates the posterior distribution with respect to the class parameters and equates the derivative with zero. This yields a system of nonlinear equations that are solved iteratively.

The number of classes is determined as follows. If a class has negligible posterior probability $\pi_j$, then including that class in the model cannot improve the likelihood of the data at all. At the same time, the prior probability of one class probability being near zero is very low. Thus models in which a class has negligible probability are always less probable than are models that simply omit that class. The user runs AutoClass with $J$ larger than the expected number of classes. If all resulting classes have significant probability, then the user increases $J$ and repeats until some classes are empty. AutoClass then ignores the empty classes, and the populated classes represent an optimal classification of the data, given the assumed class model function.

*(?..)*

*To ARC —*
*Source author should double check just in case vector and scalar quantities were not mixed. We seem to have it OK*
*Thanks*

*Type fonts on eqs b and c need corrections. Also, Σ and Π should be rechecked. Need different fonts.*

**APPROVAL COPY**