

Κάθε γραμμή του συνόλου δεδομένων ανήκει σε ένα πρότυπο. Οι γραμμές καλούνται αντικείμενα, παραδείγματα ή παρατηρήσεις. Οι στήλες αναφέρονται σε μια ιδιότητα των αντικειμένων, όπως πχ η απόσταση των ματιών, το χρώμα των ματιών κλπ. Οι στήλες καλούνται και πεδία (fields), μεταβλητές (variables), γνωρίσματα (attributes) ή χαρακτηριστικά (features). Το γνώρισμα το οποίο περιέχει την απάντηση για το αν ο επιβάτης που ελέγχεται είναι εγκληματίας ή όχι, είναι το γνώρισμα της κλάσης. Στο Σχήμα 3.1 παρουσιάζεται το σύνολο δεδομένων αυτού του παραδείγματος.

Στάδια κατηγοριοποίησης

Η κατηγοριοποίηση περιλαμβάνει τρία στάδια, το στάδιο της εκμάθησης, το στάδιο του ελέγχου του μοντέλου και το στάδιο της εφαρμογής.

- **Εκμάθηση (Learning).** Στο στάδιο αυτό, μια μέθοδος κατηγοριοποίησης αναλύει ένα σύνολο δεδομένων προκειμένου να σχηματιστεί ένα μοντέλο. Η κατασκευή ή εκπαίδευση του μοντέλου καθοδηγείται από τις τιμές του γνωρίσματος της κλάσης και για τον λόγο αυτό η διαδικασία ονομάζεται επιβλεπόμενη μάθηση. Το σύνολο δεδομένων, το οποίο χρησιμοποιείται για την εκπαίδευση του μοντέλου, ονομάζεται σύνολο εκπαίδευσης (training data set). Η επιλογή του συνόλου εκπαίδευσης είναι καθοριστικής σημασίας, γιατί το μοντέλο που θα προκύψει θα αποτυπώνει σχέσεις που υπάρχουν στο σύνολο εκπαίδευσης.
- **Έλεγχος (Testing).** Στο στάδιο αυτό ακολουθεί η διαδικασία ελέγχου της ακρίβειας του μοντέλου, η ικανότητα του δηλαδή να προβλέπει σωστά. Το μοντέλο τροφοδοτείται με αντικείμενα, των οποίων η κλάση είναι γνωστή. Αναλύοντας τα στοιχεία των ανεξάρτητων γνωρισμάτων κάθε αντικειμένου, το μοντέλο προβλέπει την κλάση του αντικειμένου και στη συνέχεια το συγκρίνει η πρόβλεψη του μοντέλου με την πραγματική τιμή της κλάσης. Αν το μοντέλο προβλέψει σωστά την κλάση ενός ικανοποιητικού ποσοστού παρατηρήσεων, τότε θεωρείται επιτυχημένο και μπορεί να χρησιμοποιηθεί για τη διατύπωση προβλέψεων.
- **Εφαρμογή (Application).** Στο στάδιο αυτό το μοντέλο αφού έχει εκπαιδευτεί και έχει κριθεί αποδεκτό, χρησιμοποιείται για τη διατύπωση προβλέψεων. Εφαρμόζεται πλέον σε νέα δεδομένα ίδιου τύπου των οποίων η κατηγοριοποίηση είναι άγνωστη και το μοντέλο καλείται να τα κατατάξει.