

Statistical classification

In machine learning and statistics, **classification** is the problem of identifying to which of a set of categories (sub-populations) a new observation belongs, on the basis of a training set of data containing observations (or instances) whose category membership is known. Examples are assigning a given email to the "spam" or "non-spam" class, and assigning a diagnosis to a given patient based on observed characteristics of the patient (sex, blood pressure, presence or absence of certain symptoms, etc.). Classification is an example of pattern recognition.

In the terminology of machine learning,^[1] classification is considered an instance of supervised learning, i.e., learning where a training set of correctly identified observations is available. The corresponding unsupervised procedure is known as clustering, and involves grouping data into categories based on some measure of inherent similarity or distance.

Often, the individual observations are analyzed into a set of quantifiable properties, known variously as explanatory variables or *features*. These properties may variously be categorical (e.g. "A", "B", "AB" or "O", for blood type), ordinal (e.g. "large", "medium" or "small"), integer-valued (e.g. the number of occurrences of a particular word in an email) or real-valued (e.g. a measurement of blood pressure). Other classifiers work by comparing observations to previous observations by means of a similarity or distance function.

An algorithm that implements classification, especially in a concrete implementation, is known as a **classifier**. The term "classifier" sometimes also refers to the mathematical function, implemented by a classification algorithm, that maps input data to a category.

Terminology across fields is quite varied. In statistics, where classification is often done with logistic regression or a similar procedure, the properties of observations are termed explanatory variables (or independent variables, regressors, etc.), and the categories to be predicted are known as outcomes, which are considered to be possible values of the dependent variable. In machine learning, the observations are often known as *instances*, the explanatory variables are termed *features* (grouped into a feature vector), and the possible categories to be predicted are *classes*. Other fields may use different terminology: e.g. in community ecology, the term "classification" normally refers to cluster analysis, i.e., a type of unsupervised learning, rather than the supervised learning described in this article.

Contents

Relation to other problems

Frequentist procedures

Bayesian procedures

Binary and multiclass classification

Feature vectors

Linear classifiers

Algorithms

Evaluation

Application domains

See also

References

Relation to other problems

Classification and clustering are examples of the more general problem of pattern recognition, which is the assignment of some sort of output value to a given input value. Other examples are regression, which assigns a real-valued output to each input; sequence labeling, which assigns a class to each member of a sequence of values (for example, part of speech tagging, which assigns a part of speech to each word in an input sentence); parsing, which assigns a parse tree to an input sentence, describing the syntactic structure of the sentence; etc.

A common subclass of classification is probabilistic classification. Algorithms of this nature use statistical inference to find the best class for a given instance. Unlike other algorithms, which simply output a "best" class, probabilistic algorithms output a probability of the instance being a member of each of the possible classes. The best class is normally then selected as the one with the highest probability. However, such an algorithm has numerous advantages over non-probabilistic classifiers:

- It can output a confidence value associated with its choice (in general, a classifier that can do this is known as a *confidence-weighted classifier*).
- Correspondingly, it can *abstain* when its confidence of choosing any particular output is too low.
- Because of the probabilities which are generated, probabilistic classifiers can be more effectively incorporated into larger machine-learning tasks, in a way that partially or completely avoids the problem of *error propagation*.

Frequentist procedures

Early work on statistical classification was undertaken by Fisher,^{[2][3]} in the context of two-group problems, leading to Fisher's linear discriminant function as the rule for assigning a group to a new observation.^[4] This early work assumed that data-values within each of the two groups had a multivariate normal distribution. The extension of this same context to more than two-groups has also been considered with a restriction imposed that the classification rule should be linear.^{[4][5]} Later work for the multivariate normal distribution allowed the classifier to be nonlinear:^[6] several classification rules can be derived based on slight different adjustments of the Mahalanobis distance, with a new observation being assigned to the group whose centre has the lowest adjusted distance from the observation.

Bayesian procedures

Unlike frequentist procedures, Bayesian classification procedures provide a natural way of taking into account any available information about the relative sizes of the different groups within the overall population.^[7] Bayesian procedures tend to be computationally expensive and, in the days before Markov chain Monte Carlo computations were developed,

approximations for Bayesian clustering rules were devised.^[8]

Some Bayesian procedures involve the calculation of group membership probabilities: these can be viewed as providing a more informative outcome of a data analysis than a simple attribution of a single group-label to each new observation.

Binary and multiclass classification

Classification can be thought of as two separate problems – binary classification and multiclass classification. In binary classification, a better understood task, only two classes are involved, whereas multiclass classification involves assigning an object to one of several classes.^[9] Since many classification methods have been developed specifically for binary classification, multiclass classification often requires the combined use of multiple binary classifiers.

Feature vectors

Most algorithms describe an individual instance whose category is to be predicted using a feature vector of individual, measurable properties of the instance. Each property is termed a feature, also known in statistics as an explanatory variable (or independent variable, although features may or may not be statistically independent). Features may variously be binary (e.g. "on" or "off"); categorical (e.g. "A", "B", "AB" or "O", for blood type); ordinal (e.g. "large", "medium" or "small"); integer-valued (e.g. the number of occurrences of a particular word in an email); or real-valued (e.g. a measurement of blood pressure). If the instance is an image, the feature values might correspond to the pixels of an image; if the instance is a piece of text, the feature values might be occurrence frequencies of different words. Some algorithms work only in terms of discrete data and require that real-valued or integer-valued data be *discretized* into groups (e.g. less than 5, between 5 and 10, or greater than 10)

Linear classifiers

A large number of algorithms for classification can be phrased in terms of a linear function that assigns a score to each possible category k by combining the feature vector of an instance with a vector of weights, using a dot product. The predicted category is the one with the highest score. This type of score function is known as a linear predictor function and has the following general form:

$$\text{score}(\mathbf{X}_i, k) = \boldsymbol{\beta}_k \cdot \mathbf{X}_i,$$

where \mathbf{X}_i is the feature vector for instance i , $\boldsymbol{\beta}_k$ is the vector of weights corresponding to category k , and $\text{score}(\mathbf{X}_i, k)$ is the score associated with assigning instance i to category k . In discrete choice theory, where instances represent people and categories represent choices, the score is considered the utility associated with person i choosing category k .

Algorithms with this basic setup are known as linear classifiers. What distinguishes them is the procedure for determining (training) the optimal weights/coefficients and the way that the score is interpreted.

Examples of such algorithms are

- Logistic regression and Multinomial logistic regression

- [Probit regression](#)
- [The perceptron algorithm](#)
- [Support vector machines](#)
- [Linear discriminant analysis](#).

Algorithms

In [unsupervised learning](#), classifiers form the backbone of cluster analysis and in [supervised](#) or semi-supervised learning, classifiers are how the system characterizes and evaluates unlabeled data. In all cases though, classifiers have a specific set of dynamic rules, which includes an interpretation procedure to handle vague or unknown values, all tailored to the type of inputs being examined.^[10]

Since no single form of classification is appropriate for all data sets, a large toolkit of classification algorithms have been developed. The most commonly used include:^[11]

- [Linear classifiers](#)
 - [Fisher's linear discriminant](#)
 - [Logistic regression](#)
 - [Naive Bayes classifier](#)
 - [Perceptron](#)
- [Support vector machines](#)
 - [Least squares support vector machines](#)
- [Quadratic classifiers](#)
- [Kernel estimation](#)
 - [k-nearest neighbor](#)
- [Boosting \(meta-algorithm\)](#)
- [Decision trees](#)
 - [Random forests](#)
- [Neural networks](#)
- [Learning vector quantization](#)

Evaluation

Classifier performance depends greatly on the characteristics of the data to be classified. There is no single classifier that works best on all given problems (a phenomenon that may be explained by the [no-free-lunch theorem](#)). Various empirical tests have been performed to compare classifier performance and to find the characteristics of data that determine classifier performance. Determining a suitable classifier for a given problem is however still more an art than a science.

The measures [precision](#) and [recall](#) are popular metrics used to evaluate the quality of a classification system. More recently, [receiver operating characteristic](#) (ROC) curves have been used to evaluate the tradeoff between true- and false-positive rates of classification algorithms.

As a performance metric, the uncertainty coefficient has the advantage over simple accuracy in that it is not affected by the relative sizes of the different classes.^[12] Further, it will not penalize an algorithm for simply *rearranging* the classes.

Application domains

Classification has many applications. In some of these it is employed as a data mining procedure, while in others more detailed statistical modeling is undertaken.

- Computer vision
 - Medical imaging and medical image analysis
 - Optical character recognition
 - Video tracking
- Drug discovery and development
 - Toxicogenomics
 - Quantitative structure-activity relationship
- Geostatistics
- Speech recognition
- Handwriting recognition
- Biometric identification
- Biological classification
- Statistical natural language processing
- Document classification
- Internet search engines
- Credit scoring
- Pattern recognition
- Micro-array classification

See also

- Artificial intelligence
- Binary classification
- Class membership probabilities
- Classification rule
- Compound term processing
- Data mining
- Data warehouse
- Fuzzy logic
- Information retrieval
- List of datasets for machine learning research
- Machine learning
- Recommender system

References

1. Alpaydin, Ethem (2010). *Introduction to Machine Learning* (<https://books.google.com/books?id=7f5bBAAQBAJ&printsec=frontcover#v=onepage&q=classification&f=false>). MIT Press. p. 9. ISBN 978-0-262-01243-0.
2. Fisher R.A. (1936) "The use of multiple measurements in taxonomic problems (<https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.1469-1809.1936.tb02137.x>)", *Annals of Eugenics*, 7, 179–188
3. Fisher R.A. (1938) "The statistical utilization of multiple measurements (<https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.1469-1809.1938.tb02189.x>)", *Annals of Eugenics*, 8, 376–386
4. Gnanadesikan, R. (1977) *Methods for Statistical Data Analysis of Multivariate Observations*, Wiley. ISBN 0-471-30845-5 (p. 83–86)
5. Rao, C.R. (1952) *Advanced Statistical Methods in Multivariate Analysis*, Wiley. (Section 9c)
6. Anderson, T.W. (1958) *An Introduction to Multivariate Statistical Analysis*, Wiley.
7. Binder, D.A. (1978) "Bayesian cluster analysis", *Biometrika*, 65, 31–38.
8. Binder, D.A. (1981) "Approximations to Bayesian clustering rules (<https://academic.oup.com/biomet/article-abstract/68/1/275/237691>)", *Biometrika*, 68, 275–285.
9. Har-Peled, S., Roth, D., Zimak, D. (2003) "Constraint Classification for Multiclass Classification and Ranking." In: Becker, B., Thrun, S., Obermayer, K. (Eds) *Advances in Neural Information Processing Systems 15: Proceedings of the 2002 Conference*, MIT Press. ISBN 0-262-02550-7
10. "What is a Classifier in Machine Learning?" (<https://deepai.org/machine-learning-glossary-and-terms/classifier>).
11. "A Tour of The Top 10 Algorithms for Machine Learning Newbies" (<https://towardsdatascience.com/a-tour-of-the-top-10-algorithms-for-machine-learning-newbies-dde4edffae11>). *Towards Data Science*. 2018-01-20. Retrieved 2018-10-10.
12. Peter Mills (2011). "Efficient statistical classification of satellite measurements". *International Journal of Remote Sensing*. arXiv:1202.2194 (<https://arxiv.org/abs/1202.2194>). doi:10.1080/01431161.2010.507795 (<https://doi.org/10.1080%2F01431161.2010.507795>).

Retrieved from "https://en.wikipedia.org/w/index.php?title=Statistical_classification&oldid=879646071"

This page was last edited on 22 January 2019, at 15:20.

Text is available under the [Creative Commons Attribution-ShareAlike License](#); additional terms may apply. By using this site, you agree to the [Terms of Use](#) and [Privacy Policy](#). Wikipedia® is a registered trademark of the [Wikimedia Foundation, Inc.](#), a non-profit organization.