

# Байесовский классификатор (Bayesian classifier)

Синонимы: Простой байесовский классификатор, Байесовская классификация, Наивный байесовский классификатор, Naïve Bayes

Разделы: [Алгоритмы](#)

В [машинном обучении](#) — семейство простых вероятностных классификаторов, основанных на использовании [теоремы Байеса](#) и «наивном» предположении о независимости признаков классифицируемых объектов.

Анализ на основе байесовской классификации активно изучался и использовался начиная с 1950-х годов в области классификации документов, где в качестве признаков использовались частоты слов. Алгоритм является масштабируемым по числу признаков, а по точности сопоставим с другими популярными методами, такими как [машины опорных векторов](#).

Как и любой классификатор, байесовский присваивает [метки классов](#) наблюдениям, представленным векторами признаков. При этом предполагается, что каждый признак независимо влияет на вероятность принадлежности наблюдения к классу. Например, объект можно считать яблоком, если он имеет округлую форму, красный цвет и диаметр около 10 см. Наивный байесовский классификатор «считает», что каждый из этих признаков независимо влияет на вероятность того, что этот объект является яблоком, независимо от любых возможных корреляций между характеристиками цвета, формы и размера.

Простой байесовский классификатор строится на основе [обучения с учителем](#). Несмотря на малореалистичное предположение о независимости признаков, простые байесовские классификаторы хорошо зарекомендовали себя при решении многих практических задач. Дополнительным преимуществом метода является небольшое число [примеров](#), необходимых для обучения.

По сути, байесовский классификатор представляет собой вероятностную модель. Пусть задано множество наблюдений, каждое из которых представлено вектором признаков  $x = (x_1, x_2, \dots, x_n)$ . Модель присваивает каждому наблюдению условную вероятность  $p(C_k | x_1, x_2, \dots, x_n)$ ,  $C_k$  — класс.

Используя теорему Байеса, можно записать:

$$p(C_k | x) = \frac{p(C_k)p(x|C_k)}{p(x)}$$

В этой формуле интерес с точки зрения классификации представляет только числитель, поскольку знаменатель от метки классов не зависит и является константой. При условии, что признаки независимы, можно показать, что

$$p(C_k|x_1, x_2, \dots, x_n) = p(C_k)p(x_1|C_k)p(x_2|C_k)\dots p(x_n|C_k) = \prod_n p(x_i|C_k).$$

Тогда простой байесовский классификатор можно рассматривать как функцию, которая каждому выходному значению модели присваивает метку класса, т.е.  $y = C_k$  следующим образом:

$$y = \arg_k \max_{1\dots k} \prod_n p(x_i|C_k)$$

Таким образом, выбирается класс  $C_k$ , который максимизирует функцию правдоподобия, представляющую собой произведение условных вероятностей значений признака  $x_i$  по каждому классу  $C_k$ .

Вероятностный классификатор предсказывает класс с самой большой условной вероятностью для заданного вектора признаков  $x$ .