Bayesian learning methods are relevant to study of machine learning for two different reasons.

• First, Bayesian learning algorithms that calculate explicit probabilities for hypotheses, such as the naive Bayes classifier, are among the most practical approaches to certain types of learning problems

• The second reason is that they provide a useful perspective for understanding many learning algorithms that do not explicitly manipulate probabilities.

Features of Bayesian Learning Methods

- Each observed training example can incrementally decrease or increase the estimated probability that a hypothesis is correct. This provides a more flexible approach to learning than algorithms that completely eliminate a hypothesis if it is found to be inconsistent with any single example
- Prior knowledge can be combined with observed data to determine the final probability of a hypothesis. In Bayesian learning, prior knowledge is provided by asserting (1) a prior probability for each candidate hypothesis, and (2) a probability distribution over observed data for each possible hypothesis
- Bayesian methods can accommodate hypotheses that make probabilistic predictions
- New instances can be classified by combining the predictions of multiple hypotheses, weighted by their probabilities.
- Even in cases where Bayesian methods prove computationally intractable, they can provide a standard of optimal decision making against which other practical methods can be measured.

Practical difficulty in applying Bayesian methods

- One practical difficulty in applying Bayesian methods is that they typically require initial knowledge of many probabilities. When these probabilities are not known in advance they are often estimated based on background knowledge, previously available data, and assumptions about the form of the underlying distributions.
- A second practical difficulty is the significant computational cost required to determine the Bayes optimal hypothesis in the general case. In certain specialized situations, this computational cost can be significantly reduced.

BAYES THEOREM

Bayes theorem provides a way to calculate the probability of a hypothesis based on its prior probability, the probabilities of observing various data given the hypothesis, and the observed data itself.

Notations

- P(h) prior probability of h, reflects any background knowledge about the chance that h is correct
- P(D) prior probability of D, probability that D will be observed
- P(D|h) probability of observing D given a world in which h holds
- P(h|D) posterior probability of h, reflects confidence that h holds after D has been observed.

Bayes theorem is the cornerstone of Bayesian learning methods because it provides a way to calculate the posterior probability P(h|D), from the prior probability P(h), together with P(D) and P(D(h).

Bayes Theorem:

$$P(h|D) = \frac{P(D|h)P(h)}{P(D)}$$

P(h|D) increases with P(h) and with P(D|h) according to Bayes theorem.

P(h|D) decreases as P(D) increases, because the more probable it is that D will be observed independent of h, the less evidence D provides in support of h.

Maximum a Posteriori (MAP) Hypothesis

• In many learning scenarios, the learner considers some set of candidate hypotheses H and is interested in finding the most probable hypothesis $h \in H$ given the observed data D. Any such maximally probable hypothesis is called a maximum a posteriori (MAP) hypothesis.

- Bayes theorem to calculate the posterior probability of each candidate hypothesis is h_{MAP} is a MAPhypothesis provided

$$h_{MAP} = \underset{h \in H}{\operatorname{argmax}} P(h|D)$$
$$= \underset{h \in H}{\operatorname{argmax}} \frac{P(D|h)P(h)}{P(D)}$$
$$= \underset{h \in H}{\operatorname{argmax}} P(D|h)P(h)$$

• P(D) can be dropped, because it is a constant independent of h

Maximum Likelihood (ML) Hypothesis

In some cases, it is assumed that every hypothesis in H is equally probable a priori (P(hi) = P(hj) for all hi and hj in H).

In this case the below equation can be simplified and need only consider the term P(D|h) to find the most probable hypothesis.

$$h_{MAP} = \underset{h \in H}{\operatorname{argmax}} P(D|h)P(h)$$

the equation can be simplified

$$h_{ML} = \underset{h \in H}{\operatorname{argmax}} P(D|h)$$

P(D|h) is often called the likelihood of the data D given h, and any hypothesis that maximizes P(D|h) is called a maximum likelihood (ML) hypothesis.

Example

Consider a medical diagnosis problem in which there are two alternative hypotheses

- The patient has a particular form of cancer (denoted by cancer)
- The patient does not (denoted by cancer)

The available data is from a particular laboratory with two possible outcomes: + (positive) and - (negative).

$$\begin{aligned} P(cancer) &= .008 & P(\neg cancer) = 0.992 \\ P(\oplus | cancer) &= .98 & P(\ominus | cancer) = .02 \\ P(\oplus | \neg cancer) &= .03 & P(\ominus | \neg cancer) = .97 \end{aligned}$$

- Suppose a new patient is observed for whom the lab test returns a positive (+) result.
- Should we diagnose the patient as having cancer or not ?

$$\begin{split} P(\oplus|cancer)P(cancer) &= (.98).008 = .0078\\ P(\oplus|\neg cancer)P(\neg cancer) = (.03).992 = .0298\\ &\Rightarrow h_{MAP} = \neg cancer \end{split}$$

