

SL09. Bayesian Learning

Introduction:

- We're trying to learn the best (most probable) hypothesis H given input data and domain knowledge.

Best == Most probable

- It's the probability of some hypothesis h given input data D :

$$P_r(h | D)$$

- We're trying to find hypothesis h with the highest probability P_r :

$$\operatorname{argmax}_{h \in H} (P_r(h | D))$$

Bayes Rule:

- Bayes Rule for probability states that:

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

- $P_r(h | D) \rightarrow$ The probability of a specific hypothesis given input data (Posterior probability).
- $P_r(D | h) \rightarrow$ The probability of data given the hypothesis. It's the (likelihood) of seeing some particular labels associated with input points, given a world where some hypothesis h is true.
- $P_r(h) \rightarrow$ The prior probability of a particular hypothesis. This value encapsulates our prior belief that one hypothesis is likely or unlikely compared to other hypotheses. This is basically the domain knowledge.
- $P_r(D) \rightarrow$ The likelihood of the data under all hypotheses (A normalizing term).
- Bayes Rule is derived from the chain rule:

$$P_r(a, b) = P_r(a, b)P_r(b)$$

$$P_r(a, b) = P_r(b, a)P_r(a)$$

$$\text{then } \rightarrow P_r(a, b)P_r(b) = P_r(b, a)P_r(a)$$

$$P_r(a, b) = \frac{P_r(b, a)P_r(a)}{P_r(b)}$$

Bayesian Learning:

- Bayesian Learning algorithm:

For each $h \in H$:

$$\text{Calculate } P_r(h | D) = \frac{P_r(D | h)P_r(h)}{P_r(D)}$$

Output:

$$h = \operatorname{argmax}_{h \in H} (P_r(h | D))$$



- Using this approximate probability, we can calculate the Maximum a Posteriori (MAP), which is the maximum probability hypothesis given the data across all hypotheses:

$$h_{map} = \operatorname{argmax}_{h \in H} (P_r(h | D))$$

$$h_{map} = \operatorname{argmax}_{h \in H} \left(\frac{P_r(D | h) P_r(h)}{P_r(D)} \right)$$

- Since we're interested in finding the hypothesis with the highest probability, not the exact probability value for each hypothesis, our prior on the data isn't exactly relevant. That is, we don't care about the $P_r(D)$ term in the denominator as it affects all computations equally:

$$h_{map} = \operatorname{argmax}_{h \in H} (P_r(D | h) P_r(h))$$

- We can also assume that our prior belief is uniform over all the hypotheses $h \in H$ (we equally believe in every $h \in H$), then we can drop $P_r(h)$ from the equation, ending up with the Maximum Likelihood:

$$h_{ml} = \operatorname{argmax}_{h \in H} (P_r(D | h))$$

- The problem with Bayes Learning is that it's not practical to perform direct computations for large hypotheses spaces, because you have to look into every single hypothesis.

Bayesian Learning in Action:

- Assume:
 - Given noise-free training data $\{\langle x_i, d_i \rangle\}$ as examples of c .
 - $c \in H$
 - Uniform prior.
- We need to calculate $P_r(h | D)$:

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

$$P_r(h) = \frac{1}{|H|}, \quad \text{because we have a uniform prior}$$

$$P_r(D | h) = \begin{cases} 1 & \text{if } d_i = h(x_i) \forall x_i, d_i \in D \\ 0 & \text{otherwise} \end{cases}$$

- This equation basically means that $P_r(D | h) = 1$ if $h \in VS(D)$

$$P_r(D) = \sum_{h_i \in H} P_r(D | h_i) P_r(h_i) = \sum_{h_i \in VS_{H,D}} 1 \cdot \frac{1}{|H|} = \frac{|VS|}{|H|}$$

$$P_r(h | D) = \frac{1 \cdot \frac{1}{|H|}}{\frac{|VS|}{|H|}} = \frac{1}{|VS|}$$

- This means that given data D , the probability of h to be a correct hypothesis is a uniform over all the hypotheses in the version space.

Bayesian Learning with Noise:

- Assume:
 - Given $\{x_i, d_i\}$
 - $d_i = f(x_i) + \varepsilon_i$
 - $\varepsilon_i \sim N(0, \sigma^2) \rightarrow$ IID (Independent and Identically Distributed)
- We need to calculate $P_r(h | D)$:

$$h_{ml} = \operatorname{argmax}_{h \in H} (P_r(D | h))$$

- To find $P_r(D | h)$ for IID, we find the product of the probability of each data point given the hypothesis:

$$h_{ml} = \operatorname{argmax}_{h \in H} \prod_i P_r(d_i | h)$$

- Given a Gaussian noise:

$$h_{ml} = \operatorname{argmax}_{h \in H} \prod_i \frac{1}{\sqrt{2\pi\sigma^2}} e^{\left(\frac{-1}{2} \cdot \frac{(d_i - h(x_i))^2}{\sigma^2}\right)}$$

- Since we're looking for the maximum:
 - We can remove $\frac{1}{\sqrt{2\pi\sigma^2}}$ since we're looking for the maximum.
 - We can take the natural log \ln to remove the exponential. Since \ln of a product is equal to the sum of the terms, we end up with the following function.

$$h_{ml} = \operatorname{argmax}_{h \in H} \sum_i \frac{-1}{2} \cdot \frac{(d_i - h(x_i))^2}{\sigma^2}$$

- Again, since we're calculating the maximum, we can remove the $\frac{1}{2}$ and the σ^2 :

$$h_{ml} = \operatorname{argmax}_{h \in H} - \sum_i (d_i - h(x_i))^2$$

- Maximizing a negative value is the same as minimizing the positive sum of this value:

$$h_{ml} = \operatorname{argmin}_{h \in H} \sum_i (d_i - h(x_i))^2$$

- This means: If you're looking for the maximum likelihood hypothesis, you should minimize the sum of squared error.
- This model will not work if the data is corrupted with any sort of noise other than Gaussian noise.

Minimum Description Length:

$$h_{map} = \operatorname{argmax}_{h \in H} (P_r(D | h) P_r(h))$$

$$h_{map} = \operatorname{argmax}_{h \in H} [\log P_r(D | h) + \log P_r(h)]$$

$$h_{map} = \operatorname{argmin}_{h \in H} [-\log P_r(D | h) - \log P_r(h)]$$

- Information theory: The optimal code for some event w with probability P_r has a length of $-\log P_r$.
- This means that in order to maximize the Maximum a Posteriori hypothesis, we need to minimize two terms that can be described as length:
 - $\log P_r(h) \rightarrow$ This is the length of the hypothesis, which is the number of bits needed to represent this hypothesis.
 - $\log P_r(D | h) \rightarrow$ This is the length of the data given a particular hypothesis. If the hypothesis perfectly describes the data, so we don't need any points. But if the hypothesis labels some points wrong, so we need the correct labels for these points to be able to come up with a better hypothesis. So basically this term captures the error.
- This is always a trade of, a more complex hypothesis will drive down error, while a simple hypothesis will have some error.
- We need to find the best hypothesis, which is the simplest hypothesis that minimizes error. This hypothesis is called the Minimum Description.

Bayesian Classification:

- The question in classification is "What is the best label?" not the best hypothesis.
- To find the best label, we need to do a weighted vote for every single hypothesis in the hypotheses set, where the weight is the probability $P_r(h | D)$.
- Now we end up trying to maximize v_{map} :

$$V_{map} = \operatorname{argmax}_{v_j \in V} \sum_{h_i \in H} P_r(v_j | h_i) P_r(h_i | D)$$

- The Bayes optimal classifier is computationally very costly. This is because the posterior probability $P_r(h | D)$ must be computed for each hypothesis $h \in H$ and combined with the prediction $P_r(v | h)$ before V_{map} can be computed.