

Machine Learning

Summary by Ari Feiglin (ari.feiglin@gmail.com)

Contents

1	Bayes Decision Theory	1
1.1	Formalizing the Theory	1
1.2	Minimizing Risk	1
1.3	Parameter Estimation	2

1 Bayes Decision Theory

1.1 Formalizing the Theory

Let us imagine the following scenario: you, a computer science student, are preparing to leave your house for the first time in a while. Should you or should you not take your umbrella? If you take your umbrella and it doesn't rain then you've inconvenienced yourself, but if you don't and it rains then you'll end up getting wet. You look outside, see that clouds are grey, and decide to take your umbrella.

Here you are trying to decide between some possible actions, each with their own cost. These actions are based on the state of the world, of which you have some set of observations.

The problem can be formalized as follows, you have

- (1) a set of world states: $\{\omega_i\}_{i \in I}$, which are disjoint and exhaustive: $\Omega = \bigcup_{i \in I} \omega_i$ (where Ω is the entire space).
- (2) a set of observations: $S_n = \{x_1, \dots, x_n\}$.
- (3) a probabilistic model: conditionals $\mathbb{P}(S_n | \omega)$ and priors $\mathbb{P}(\omega)$.
- (4) a set of possible actions $A = \{\alpha_1, \dots, \alpha_k\}$.
- (5) a set of cost functions $\Lambda = \{\lambda(\alpha_k | \omega_j)\}_{j,k}$.

So in our example above, we have two world states: raining and not raining, our observation is that the clouds are grey, we have some prior belief about the probabilities of each world state as well as the probability of observing our observation given a world state, our actions are to take and to not take an umbrella, each has an associated cost.

So suppose we've made an observation x , we'd like to compute the new probability of a world state ω given this observation. This can be done via Bayes's law:

$$\mathbb{P}(\omega | x) = \frac{\mathbb{P}(x | \omega)}{\mathbb{P}(x)} \cdot \mathbb{P}(\omega)$$

$\mathbb{P}(\omega | x)$ is called the *posterior*. But $\mathbb{P}(x)$ is not known, so how do we compute it? Using total probability:

$$\mathbb{P}(x) = \sum_{i \in I} \mathbb{P}(x | \omega_i) \mathbb{P}(\omega_i)$$

As we make more and more observations, we can employ Baye's law over and over to refine the posterior.

1.2 Minimizing Risk

Our goal in Bayes decision theory is to define a strategy $\alpha(S_n)$ which determines which action α to take so that it minimizes our expected costs, given observations S_n .

Definition 1.2.1

Given an action α , we define its **conditional risk** given an observation x to be

$$R(\alpha | x) = \sum_{i \in I} \lambda(\alpha | \omega_i) \mathbb{P}(\omega_i | x)$$

Where the posterior $\mathbb{P}(\omega_i | x)$ is computed as above using Bayes's law. If we define the random variable $\lambda(\alpha)$ to be $\lambda(\alpha | \omega)$ under the state ω , then this is just $\mathbb{E}[\lambda(\alpha) | x]$.

Example 1.2.2

Suppose our actions $\alpha_1, \dots, \alpha_k$ are to guess the world state $\omega_1, \dots, \omega_k$. The cost we pay is 1 if we are wrong and 0 if we are correct, meaning $\lambda(\alpha_k | \omega_j) = 1 - \delta_{kj}$ where $\delta_{kj} = 1$ when $k = j$ and 0 otherwise.

Then the conditional risk is

$$R(\alpha_k | x) = \sum_{j=1}^k \lambda(\alpha_k | \omega_j) \mathbb{P}(\omega_j | x) = \sum_{j \neq k} \mathbb{P}(\omega_j | x) = 1 - \mathbb{P}(\omega_k | x)$$

So we minimize the conditional risk when we take α_k such that the posterior $\mathbb{P}(\omega_k | x)$ is maximal.

Example 1.2.3

Suppose we have two world states ω_1, ω_2 and two actions α_1, α_2 . Then our costs form a 2×2 matrix: $\lambda_{kj} = \lambda(\alpha_k | \omega_j)$. And by definition

$$R(\alpha_i | x) = \lambda_{i1} \mathbb{P}(\omega_1 | x) + \lambda_{i2} \mathbb{P}(\omega_2 | x)$$

We choose α_1 if $R(\alpha_1 | x) < R(\alpha_2 | x)$, which is equivalent to

$$(\lambda_{12} - \lambda_{22}) \mathbb{P}(\omega_2 | x) < (\lambda_{21} - \lambda_{11}) \mathbb{P}(\omega_1 | x)$$

which is equivalent to

$$\iff \frac{\mathbb{P}(\omega_1 | x)}{\mathbb{P}(\omega_2 | x)} > \frac{\lambda_{12} - \lambda_{22}}{\lambda_{21} - \lambda_{11}} \iff \frac{\mathbb{P}(x | \omega_1)}{\mathbb{P}(x | \omega_2)} > \frac{\mathbb{P}(\omega_2)}{\mathbb{P}(\omega_1)} \cdot \frac{\lambda_{12} - \lambda_{22}}{\lambda_{21} - \lambda_{11}}$$

The left-hand side is called the **likelihood ratio** and the right-hand side is called the **decision boundary**.

If we have many observations x_1, \dots, x_n then this becomes

$$\begin{array}{cc} \text{Likelihood ratio} & \text{Decision boundary} \\ \frac{\mathbb{P}(x_1, \dots, x_n | \omega_1)}{\mathbb{P}(x_1, \dots, x_n | \omega_2)} > \frac{\mathbb{P}(\omega_2)}{\mathbb{P}(\omega_1)} \cdot \frac{\lambda_{22} - \lambda_{12}}{\lambda_{11} - \lambda_{21}} = \Theta \end{array}$$

If the observations x_1, \dots, x_n are independent then this becomes

$$\frac{\prod_i \mathbb{P}(x_i | \omega_1)}{\prod_i \mathbb{P}(x_i | \omega_2)} > \Theta$$

taking the log of both sides gives

$$\sum_i \log \frac{\mathbb{P}(x_i | \omega_1)}{\mathbb{P}(x_i | \omega_2)} > \log \Theta = \Theta'$$

the left-hand side is called the **log-likelihood ratio**.

1.3 Parameter Estimation

Suppose we know the distribution of some random variable X up to some parameter θ , i.e. we know the function $\mathbb{P}(X | \theta)$. For example, X is the number of heads in n coin tosses where the coin has a bias of θ , then $X | \theta \sim \text{Bin}(n, \theta)$.

We use Bayes decision theory to estimate θ . Suppose we have a prior $\mathbb{P}(\theta)$, we use this to estimate θ . Our actions will be choosing some prediction of θ , $\hat{\theta}$. Define a cost function $\lambda(\hat{\theta}, \theta)$. Then given a sequence of observations $S_n = \{x_1, \dots, x_n\}$, we want to minimize the expected cost

$$\mathbb{E}[\lambda(\hat{\theta}) | S_n] = \int \lambda(\hat{\theta}, \theta) \mathbb{P}(\theta | S_n) d\theta$$

We then define the *Bayes estimator* to be the prediction θ^* which minimizes this:

$$\theta^* = \text{argmin}_{\hat{\theta}} \mathbb{E}[\lambda(\hat{\theta}) | S_n]$$

To find θ^* we differentiate $\mathbb{E}[\lambda(\hat{\theta}) | S_n]$ and compare it with zero. Hand-waving away all the technical details because this is computer science, we can swap the order of differentiation and integration and so

$$\frac{d}{d\hat{\theta}} \mathbb{E}[\lambda(\hat{\theta}) | S_n] = \frac{d}{d\hat{\theta}} \int \lambda(\hat{\theta}, \theta) \mathbb{P}(\theta | S_n) d\theta = \int \frac{d}{d\hat{\theta}} \lambda(\hat{\theta}, \theta) \mathbb{P}(\theta | S_n) d\theta$$

Example 1.3.1 (Square Loss)

Suppose we use a **square loss** cost function: $\lambda(\hat{\theta}, \theta) = (\theta - \hat{\theta})^2$. Then the Bayes estimator is

$$\mathbb{E}[\lambda(\hat{\theta})] = \int \lambda(\hat{\theta}, \theta) \mathbb{P}(\theta | S_n) d\theta = \int (\theta - \hat{\theta})^2 \mathbb{P}(\theta | S_n) d\theta$$

Differentiating gives that we want

$$\int \theta \mathbb{P}(\theta | S_n) d\theta = \hat{\theta} \int \mathbb{P}(\theta | S_n) d\theta = \hat{\theta}$$

(since $\int \mathbb{P}(\theta | S_n) = 1$.) So we get that our estimator is

$$\hat{\theta}_{SE} = \int \theta \mathbb{P}(\theta | S_n) d\theta = \mathbb{E}[\theta | S_n]$$

Example 1.3.2 (Maximum A posteriori)

Suppose we use a **zero-one** cost function: $\lambda(\hat{\theta}, \theta) = 1 - \delta_{\hat{\theta}\theta}$. We have already showed that to minimize the cost, we must maximize $\text{argmax} \mathbb{P}(\theta | S_n)$. By Bayes, this is just $\text{argmax} \mathbb{P}(S_n | \theta) \frac{\mathbb{P}(\theta)}{\mathbb{P}(S_n)} = \text{argmax} \mathbb{P}(S_n | \theta) \mathbb{P}(\theta)$. Since the observations in S_n are independent, we see that this is then just equal to $\text{argmax} \prod_i \mathbb{P}(x_i | \theta) \mathbb{P}(\theta)$. Finally we take the log, which is monotonic, to get that the estimator

$$\hat{\theta}_{MAP} = \text{argmax} \sum_i \log \mathbb{P}(x_i | \theta) + \log \mathbb{P}(\theta)$$

as n grows, the last term becomes negligible and this just becomes the maximum log likelihood function.

Exercise 1.3.3

Compute the MAP estimator for $X \sim \text{Exp}(\theta)$, i.e. $f_X(x) = \theta \exp(-\theta x)$, with the prior $\mathbb{P}(\theta) = \exp(-\theta)$. We want to compute

$$\sum_i \log \mathbb{P}(x_i | \theta) + \log \mathbb{P}(\theta) = \sum_i \log(\theta \exp(-\theta x_i)) + \log \exp(-\theta) = n \log \theta - \theta \sum_i x_i - \theta$$

Differentiating with respect to θ gives

$$\frac{n}{\theta} - \sum_i x_i - 1 = 0 \Rightarrow \hat{\theta}_{MAP} = \frac{n}{\sum_i x_i + 1}$$

Here the prior behaves similarly to a sample, it can be viewed as a “pseudo-sample”.

Example 1.3.4 (Maximum Likelihood)

We know that

$$\hat{\theta}_{MAP} = \text{argmax} \sum_i \log \mathbb{P}(x_i | \theta) + \log \mathbb{P}(\theta)$$

A common approach is to approximate this by dropping the prior, giving

$$\hat{\theta}_{ML} = \operatorname{argmax}_{\theta} \sum_i \log \mathbb{P}(x_i \mid \theta)$$