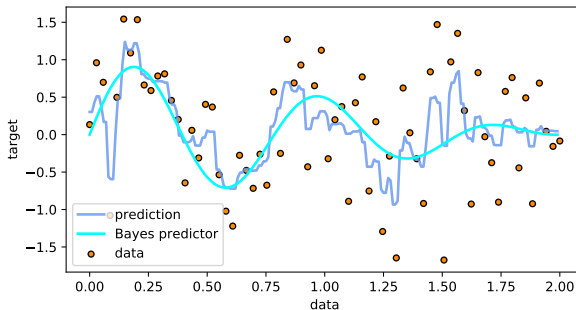


Fondamentaux théoriques du machine learning

Random forest regression
number of estimators: 20
max depth: 5
test error: 8.09E-01
Bayes risk: 7.00E-01



Step 1

We have seen that

$$R(f_n) - R(f_a) \leq 2 \sup_{h \in F} |R(h) - R_n(h)| \quad (1)$$

As a consequence, for all $t \geq 0$:

$$P\left(R(f_n) - R(f_a) \geq t\right) \leq P\left(2 \sup_{h \in F} |R(h) - R_n(h)| \geq t\right) \quad (2)$$

Step 2

The fact that

$$2 \sup_{h \in F} |R(h) - R_n(h)| \geq t \quad (3)$$

is equivalent to :

$$\cup_{h \in F} \left(2 |R(h) - R_n(h)| \geq t \right) \quad (4)$$

Step 3

Boole's inequality shows that :

$$P\left(\bigcup_{h \in F} \left(2|R(h) - R_n(h)| \geq t\right)\right) \leq \sum_{h \in F} P(2|R(h) - R_n(h)| \geq t) \quad (5)$$

Step 4

For each $h \in F$, we need to bound

$$P(2|R(h) - R_n(h)| \geq t) \tag{6}$$

Step 4

For each $h \in F$, we need to bound

$$P(2|R(h) - R_n(h)| \geq t) \quad (7)$$

With Hoeffding's inequality we get

$$P(2|R(h) - R_n(h)| \geq t) \leq 2 \exp\left(-\frac{nt^2}{2(b-a)^2}\right) \quad (8)$$

Step 5

With Hoeffding's inequality we get

$$P(2|R(h) - R_n(h)| \geq t) \leq 2 \exp\left(-\frac{nt^2}{2(b-a)^2}\right) \quad (9)$$

Finally, putting everything together :

$$\begin{aligned} P(R(f_n) - R(f_a) \geq t) &\leq \sum_{h \in F} P(2|R(h) - R_n(h)| \geq t) \\ &\leq \sum_{h \in F} 2 \exp\left(-\frac{nt^2}{2(b-a)^2}\right) \\ &= 2|F| \exp\left(-\frac{nt^2}{2(b-a)^2}\right) \end{aligned} \quad (10)$$

Conclusion

$$P\left(R(f_n) - R(f_a) \geq t\right) \leq 2|F| \exp\left(-\frac{nt^2}{2(b-a)^2}\right) \quad (11)$$

Conclusion

We write

$$\delta = 2|F| \exp\left(-\frac{nt^2}{2(b-a)^2}\right) \quad (12)$$

We assume that $b - a = 1$. Then, with probability $1 - \delta$, we can compute and show that

$$R(f_n) \leq R(f_a) + 2\sqrt{\frac{\log(|F|) + \log(\frac{2}{\delta})}{2n}} \quad (13)$$