



# SNS COLLEGE OF TECHNOLOGY

Coimbatore-37.

An Autonomous Institution



**COURSE NAME : 19CST301 & INTRODUCTION TO MACHINE LEARNING**

**III YEAR/ V SEMESTER**

**UNIT – 3 DEEP LEARNING**

**Topic: Interpretation**

Mrs.S.R.Janani

Assistant Professor

Department of Computer Science and Engineering



# Interpretation

- Single decision trees are highly interpretable.
- The entire model can be completely represented by a simple two-dimensional graphic (binary tree) that is easily visualized.
- Linear combinations of trees lose this important feature, and must therefore be interpreted in a different way.
  - Relative Importance of Predictor Variables
  - Partial Dependence Plots



# Relative Importance of Predictor Variables



- In data mining applications the **input predictor variables** are seldom equally relevant.
- Often only a few of them have substantial influence on the response; the **vast majority are irrelevant** and could just as well have not been included.
- It is often useful **to learn** the relative importance or contribution of each input variable in predicting the response



For a single decision tree  $T$ , Breiman et al. (1984) proposed

$$\mathcal{I}_\ell^2(T) = \sum_{t=1}^{J-1} \hat{v}_t^2 I(v(t) = \ell) \quad (10.42)$$

as a measure of relevance for each predictor variable  $X_\ell$ .

The sum is over the  $J - 1$  internal nodes of the tree.

At each such node  $t$ , one of the input variables  $X_{v(t)}$  is used to partition the region associated with that node into two subregions; within each a separate constant is fit to the response values.

The particular variable chosen is the one that gives maximal estimated improvement  $\hat{v}_t^2$  in **squared error risk** over that for a constant fit over the entire region.

The squared relative importance of variable  $X_\ell$  is the sum of such squared improvements over all internal nodes for which it was chosen as the splitting variable.



This importance measure is easily generalized to additive tree expansions. it is simply averaged over the trees

$$\mathcal{I}_\ell^2 = \frac{1}{M} \sum_{m=1}^M \mathcal{I}_\ell^2(T_m). \quad (10.43)$$

Due to the **stabilizing effect of averaging**, this measure turns out to be more reliable than is its counterpart for a single tree.

Also, because of shrinkage the masking of important variables by others with which they are **highly correlated is much less of a problem.**

Note that (10.42) and (10.43) refer to squared relevance; the actual relevances are their respective square roots.

Since these measures are relative, it is customary to assign the largest a value of 100 and then scale the others accordingly.

Figure 10.6 shows the relevant importance of the 57 inputs in predicting spam versus email.



- For K-class classification, K separate models  $f_k(x), k = 1, 2, \dots, K$  are induced, each consisting of a sum of trees

$$f_k(x) = \sum_{m=1}^M T_{km}(x). \quad (10.44)$$

In this case (10.43) generalizes to

$$\mathcal{I}_{\ell k}^2 = \frac{1}{M} \sum_{m=1}^M \mathcal{I}_{\ell}^2(T_{km}). \quad (10.45)$$

Here  $\mathcal{I}_{\ell k}$  is the relevance of  $X_{\ell}$  in separating the class  $k$  observations from the other classes. The overall relevance of  $X_{\ell}$  is obtained by averaging over all of the classes

$$\mathcal{I}_{\ell}^2 = \frac{1}{K} \sum_{k=1}^K \mathcal{I}_{\ell k}^2. \quad (10.46)$$



# Partial Dependence Plots

- After the most relevant variables have been identified, the next step is to attempt to understand the nature of the dependence of the approximation  $f(X)$  on their joint values.
- Graphical renderings of the  $f(X)$  as a function of its arguments provides a comprehensive summary of its dependence on the joint values of the input variables.
- Unfortunately, such visualization is limited to low-dimensional views.
- We can easily display functions of one or two arguments, either continuous or discrete (or mixed), in a variety of different ways; this book is filled with such displays.
- Functions of slightly higher dimensions can be plotted by conditioning on particular sets of values of all but one or two of the arguments, producing a trellis of plots (Becker et al., 1996).<sup>1</sup>



- For more than two or three variables, viewing functions of the corresponding higher-dimensional arguments is more difficult.
- A useful alternative can sometimes be to view a collection of plots, each one of which shows the partial dependence of the approximation  $f(X)$  on a selected small subset of the input variables.
- Although such a collection can seldom provide a comprehensive depiction of the approximation, it can often produce helpful clues, especially when  $f(x)$  is dominated by low-order interactions





Consider the subvector  $X_S$  of  $\ell < p$  of the input predictor variables  $X^T = (X_1, X_2, \dots, X_p)$ , indexed by  $S \subset \{1, 2, \dots, p\}$ .

Let  $C$  be the complement set, with  $S \cup C = \{1, 2, \dots, p\}$ .

A general function  $f(X)$  will in principle depend on all of the input variables:  $f(X) = f(X_S, X_C)$ .

One way to define the average or partial dependence of  $f(X)$  on  $X_S$  is

$$f_S(X_S) = E_{X_C} f(X_S, X_C). \quad (10.47)$$



- This is a marginal average of  $f$ , and can serve as a useful description of the effect of the chosen subset on  $f(X)$  when,
  - for example, the variables in  $X_S$  do not have strong interactions with those in  $X_C$ .
  - Partial dependence functions can be used to interpret the results of any “black box” learning method.
  - They can be estimated by

$$\bar{f}_S(X_S) = \frac{1}{N} \sum_{i=1}^N f(X_S, x_{iC}), \quad (10.48)$$

- where  $\{x_{1C}, x_{2C}, \dots, x_{NC}\}$  are the values of  $X_C$  occurring in the training data. This requires a pass over the data for each set of joint values of  $X_S$  for which  $\bar{f}_S(X_S)$  is to be evaluated.



- This can be computationally intensive, even for moderately sized data sets. Fortunately with decision trees,  $\tilde{f}_S(X_S)$  (10.48) can be rapidly computed from the tree itself without reference to the data .
- It is important to note that partial dependence functions defined in (10.47) represent the effect of  $X_S$  on  $f(X)$  after accounting for the (average) effects of the other variables  $X_C$  on  $f(X)$ .
- They are not the effect of  $X_S$  on  $f(X)$  ignoring the effects of  $X_C$ . The latter is given by the conditional expectation

$$\tilde{f}_S(X_S) = E(f(X_S, X_C)|X_S), \quad (10.49)$$



- and is the best least squares approximation to  $f(X)$  by a function of  $X_S$  alone.
- The quantities  $\tilde{f}_S(X_S)$  and  $\bar{f}_S(X_S)$  will be the same only in the unlikely event that  $X_S$  and  $X_C$  are independent.

- For example, if the effect of the chosen variable subset happens to be purely additive,

$$f(X) = h_1(X_S) + h_2(X_C). \quad (10.50)$$

- Then (10.47) produces the  $h_1(X_S)$  up to an additive constant. If the effect is purely multiplicative,

$$f(X) = h_1(X_S) \cdot h_2(X_C), \quad (10.51)$$

- then (10.47) produces  $h_1(X_S)$  up to a multiplicative constant factor.
- On the other hand, (10.49) will not produce  $h_1(X_S)$  in either case.
- In fact, (10.49) can produce strong effects on variable subsets for which  $f(X)$  has no dependence at all.



- Viewing plots of the partial dependence of the boosted-tree approximation (10.28) on selected variables subsets can help to provide a qualitative description of its properties.
- Illustrations are shown in Sections 10.8 and 10.14. Owing to the limitations of computer graphics, and human perception, the size of the subsets  $X_S$  must be small ( $|S| \approx 1, 2, 3$ ).
- There are of course a large number of such subsets, but only those chosen from among the usually much smaller set of highly relevant predictors are likely to be informative.
- Also, those subsets whose effect on  $f(X)$  is approximately additive (10.50) or multiplicative (10.51) will be most revealing.



- For K-class classification, there are K separate models (10.44), one for each class. Each one is related to the respective probabilities (10.21) through

$$f_k(X) = \log p_k(X) - \frac{1}{K} \sum_{l=1}^K \log p_l(X). \quad (10.52)$$

- Thus each  $f_k(X)$  is a monotone increasing function of its respective probability on a logarithmic scale.
- Partial dependence plots of each respective  $f_k(X)$  (10.44) on its most relevant predictors (10.45) can help reveal how the log-odds of realizing that class depend on the respective input variables.



# References

- AlpaydinEthem, “Introduction to Machine Learning”, MIT Press, Second Edition, 2010.
- Trevor Hastie, Robert Tibshirani, Jerome Friedman, “The Elements of Statistical Learning: Data Mining, Inference, and Prediction”, Springer; Second Edition, 2009.

