

STUDY OF STRUCTURAL AND PROBABILISTIC MODELLING AND MACHINE LEARNING

Sapna Thakur and Dr. Priyanka Bhalerao

Department of Mathematics

Dr. A. P. J. Abdul Kalam University, Indore (M. P.) – 452010

Corresponding Author Email : sapana.thakur15@yahoo.in

Abstract:

Probabilistic modelling provides a frame-work for understanding what learning is, and has therefore emerged as one of the principal theoretical and practical approaches for designing machines that learn from data acquired through experience. The probabilistic framework, which describes how to represent and manipulate uncertainty about models and predictions, plays a central role in scientific data analysis, machine learning, robotics, cognitive science, and artificial intelligence. This article provides an introduction to this probabilistic framework, and reviews some state-of-the-art advances in the field, namely, probabilistic programming, Bayesian optimisation, data compression, and automatic model discovery.

1. Introduction

The key idea behind the probabilistic framework to machine learning is that learning can be thought of as inferring plausible models to explain observed data. A machine can use such models to make predictions about future data, and decisions that are rational given these predictions. Uncertainty plays a fundamental role in all of this. Observed data can be consistent with many models, and therefore which model is appropriate given the data is uncertain. Similarly, predictions, about future data and the future consequences of actions, are uncertain. Probability theory provides a framework for modelling uncertainty. This article starts with an introduction to the probabilistic approach to machine learning and Bayesian inference, and then reviews some of the state-of-the-art in the field. The central thesis is that many aspects of learning and intelligence depend crucially on the careful probabilistic representation of uncertainty. Probabilistic approaches have only recently become a main-stream paradigm in artificial intelligence, robotics and machine learning [1]. Even now, there is controversy in these fields about how important it is to fully represent uncertainty. For example, recent advances using deep neural networks to solve challenging pattern recognition problems such as speech recognition, image classification and prediction of words in text [2], do not overtly represent the uncertainty in the structure or parameters of those neural networks. However, my focus will not be on these types of pattern recognition problems, characterised by the availability of large amounts of data, but rather on problems where uncertainty is really a key ingredient, for example where a decision may depend on the amount of uncertainty.

2. Probabilistic Modelling and the Representation of Uncertainty

At a most basic level, machine learning seeks to develop methods for computers to improve their performance at certain tasks based on observed data. Typical examples of such tasks might include detecting pedestrians in images taken from an autonomous vehicle, classifying gene-expression patterns from leukaemia patients into subtypes by clinical outcome, or translating English sentences into French. However, as we will see, the scope of machine learning tasks is even broader than these pattern classification or mapping tasks, and can include optimisation and decision making, compressing data, and automatically extracting interpretable models from data.

Data are the key ingredient of all machine learning systems. But data, even so called "Big Data", is useless on its own until one extracts knowledge or inferences from it. Almost all machine learning tasks can be formulated as making inferences about missing or latent data from the observed data. We will variously use the terms inference, prediction or forecasting to refer to this general task. Elaborating the example mentioned above, consider classifying leukaemia patients into one of the four main subtypes of this disease, based on each patient's measured gene-expression patterns. Here the observed data are pairs of gene-expression patterns and labelled subtypes, and the unobserved or missing data to be inferred are the subtypes for new patients.

There are many forms of uncertainties in modelling. At the lowest level, model uncertainty is introduced from measurement noise, e.g., pixel noise or blur in images. At higher levels, a model may have many parameters, such as the linear regression, and there is uncertainty about which values of these parameters will be good at predicting new data. Finally, at the highest levels, there is often uncertainty about even the general structure of the model: is linear regression appropriate or a neural network, if the latter, how many layers, etc.

The probabilistic approach to modelling uses probability theory to express all forms of uncertainty. Probability theory is the mathematical language for representing and manipulating uncertainty, in much the same way as calculus is the language for representing and manipulating rates of change. Fortunately, the probabilistic approach to modelling is conceptually very simple: probability distributions are used to represent all the uncertain unobserved quantities in a model (including structural, parametric, and noise-related) and how they relate to the data. Then the basic rules of probability theory are used to infer the unobserved quantities given the observed data. Learning from data occurs through the transformation of the prior probability distributions (defined before observing the data), into posterior distributions (after observing data). The application of probability theory to learning from data is called Bayesian learning.

Probabilistic modelling also has some conceptual advantages over alternatives as a normative theory for learning in artificially intelligent (AI) systems. How should an AI system represent and update its beliefs about the world in light of data? The Cox axioms define some desiderata for representing beliefs; a consequence of these axioms is that 'degrees of belief', ranging from "impossible" to "absolutely certain", must follow all the rules of probability theory [3]. This justifies the use of subjective Bayesian probabilistic representations in AI. An argument for Bayesian representations in AI that is motivated by decision theory is given by the Dutch-Book theorems. The argument rests on the idea that the strength of beliefs of an agent can be assessed

by asking the agent whether it would be willing to accept bets at various odds (ratios of payoffs). The Dutch-Book theorems state that unless an AI system's (or human's, for that matter) degrees of beliefs are consistent with the rules of probability it will be willing to accept bets that are guaranteed to lose money [4]. Because of the force of these and many other arguments on the importance of a principled handling of uncertainty for intelligence, Bayesian probabilistic modelling has emerged not only as the theoretical foundation for rationality in AI systems but also as a model for normative behaviour in humans and animals and much research is devoted to understanding how neural circuitry may be implementing Bayesian inference [5].

3. Flexibility through Non-Parametrics

The best way to understand non-parametric models is through comparison to parametric ones. In a parametric model, there are a number of parameters, and no matter how much training data are observed, all the data can do is set these parameters that control future predictions. In contrast, nonparametric approaches have predictions that grow in complexity with the amount of training data, either by considering a nested sequence of parametric models with increasing numbers of parameters or by starting out with a model with infinitely many parameters.

For example, in a classification problem, whereas a linear (i.e., parametric) classifier will always predict using a linear boundary between classes, a nonparametric classifier can learn a nonlinear boundary whose shape becomes more complex with more data. Many nonparametric models can be derived starting from a parametric model and considering that happens as the model grows to the limit of infinitely many parameters [6]. Clearly, fitting a model with infinitely many parameters to finite training data would result in "over fitting", in the sense that the model's predictions might reflect quirks of the training data rather than regularities that can be generalised to test data.

3.1 Probabilistic Programming

The basic idea in probabilistic programming is to use computer programs to represent probabilistic models [2],[7]. One way to do this is for the computer program to define a generator for data from the probabilistic model, i.e., a simulator. This simulator makes calls to a random number generator in such a way that repeated runs from the simulator would sample different possible data sets from the model. This simulation framework is more general than the graphical model framework described previously since computer programs can allow constructs such as recursion (functions calling themselves) and control flow statements (e.g., if statements resulting in multiple paths a program can follow) which are difficult or impossible to represent in a finite graph. In fact, for many of the recent probabilistic programming languages that are based on extending Turing-complete languages (a class that includes almost all commonly-used languages), it is possible to represent any computable probability distribution as a probabilistic program [8].

The full potential of probabilistic programming comes from automating the process of inferring unobserved variables in the model conditioned on the observed data. Conceptually, conditioning needs to compute input states of the program that generate data matching the observed data. Whereas normally we think of programs running from inputs to outputs, conditioning involves solving the inverse problem of inferring the inputs (in particular the random number calls) that match a certain program output. Such conditioning is performed by a universal inference engine, usually implemented by Monte Carlo sampling over possible

executions over the simulator program that are consistent with the observed data. The fact that defining such universal inference algorithms for computer programs is even possible is somewhat surprising, but it is related to the generality of certain key ideas from sampling such as rejection sampling, sequential Monte Carlo and "approximate Bayesian computation" [9].

3.2 Bayesian optimisation

Consider the very general problem of finding the global maximum of an unknown function which is expensive to evaluate (say, evaluating the function requires performing lots of computation, or conducting an experiment). Mathematically, for a function f on a domain X , the goal is to find a global maximiser x^* :

$$x^* = \arg \max_{x \in \mathcal{X}} f(x).$$

Bayesian optimisation poses this as a problem in sequential decision theory: where should one evaluate next so as to most quickly maximize f , taking into account the gain in information about the unknown function f [10]? For example, having evaluated at three points measuring the corresponding values of the function at those points, $f(x_1)$; $f(x_2)$; $f(x_3)$, which point x should the algorithm evaluate next, and where does it believe the maximum to be? This is a classic machine intelligence problem with a wide range of applications in science and engineering, e.g., from drug design to robotics where the function could be the drug's efficacy or the speed of a robot's gait respectively. Basically, it can be applied to any problem involving the optimisation of expensive functions; the qualifier "expensive" comes because Bayesian optimisation might use substantial computational resources to decide where to evaluate next, and these resources have to be traded with the cost of function evaluations.

3.3 Data Compression

Consider the problem of compressing data so as to communicate it or store it in as few bits as possible, in such a manner that the original data can be recovered exactly from the compressed data. Methods for such lossless data compression are ubiquitous in information technology, from computer hard drives to data transfer over the internet. Data compression and probabilistic modelling are two sides of the same coin, and Bayesian machine learning methods are increasingly advancing the state of the art in compression. The connection between compression and probabilistic modelling was established in Shannon's seminal work on the source coding theorem [11] which states that the number of bits required to compress data is bounded by the entropy of the probability distribution of the data. All commonly used lossless data compression algorithms (e.g., gzip, etc) can be viewed as probabilistic models of sequences of symbols.

4. Automatically Discovering Interpretable Models from Data

One of the grand challenges of machine learning is to fully automate the process of learning and explaining statistical models from data. This is the goal of the Automatic Statistician, a system that can automatically discover plausible models from data, and explain what it has discovered in plain English [12]. This could be useful to almost any field of endeavour that is reliant on extracting knowledge from data. In contrast to much of the machine learning literature which has been focused on extracting increasing performance improvements on pattern recognition problems using techniques such as kernel methods, random forests, or deep learning, the Automatic Statistician needs to build models that are composed of interpretable

components, and to have a principled way of representing uncertainty about model structures given data. It also needs to be able to give reasonable answers not just for big data sets but also for small ones. Bayesian approaches provide an elegant way of trading off the complexity of the model and the complexity of the data, and probabilistic models are compositional and interpretable as described previously [13].

Probabilistic Model for Attributes A probabilistic relational model O specifies a probability distribution over all instantiations K of the relational schema. It consists of the qualitative dependency structure, P , and the parameters associated with it, Q^*R . The dependency structure is defined by associating with each attribute A a set of parents $Pa(A)$. Each parent has the form $S \rightarrow T$ where S is either empty or a single slot $\%$. (PRMs also allow dependencies on longer slot chains, but we have chosen to omit those for simplicity of presentation.) To understand the semantics of this dependence, note that $U \rightarrow S$ is a multi-set of values V in S . We use the notion of aggregation from database theory to define the dependence on a multi-set; thus, U will depend probabilistically on some aggregate property $W \rightarrow X$. In this paper, we use the median for ordinal attributes, and the mode (most common value) for others. When V is single-valued, both reduce to a dependence on the value of $U \rightarrow S$.

The quantitative part of the PRM specifies the parameterization of the model. Given a set of parents for an attribute, we can define a local probability model by associating with it a conditional probability distribution (CPD). For each attribute we have a CPD that specifies $Y[Z \mid Pa(A)]$.

The quantitative part of the PRM specifies the parameterization of the model. Given a set of parents for an attribute, we can define a local probability model by associating with it a conditional probability distribution (CPD). For each attribute we have a CPD that specifies $Y[Z \mid Pa(A)]$.

Definition 1: A probabilistic relational model (PRM) O for a relational schema P is defined as follows. For each class and each descriptive attribute, we have a set of parents $Pa(A)$, and a conditional probability distribution (CPD) that represents $Y_{\%} \mid [Z \mid Pa(A)]$. Given a relational skeleton L , a PRM O specifies a distribution over a set of instantiations K consistent with L :

$$P(I \mid \sigma_r, \Pi) = \prod_{x \in \sigma_r(X)} \prod_{A \in \mathcal{A}(x)} P(x.A \mid Pa(x.A))$$

where L are the objects of each class as specified by the relational skeleton L (in general we will use the notation L to refer to the set of objects of each class as defined by any type of domain skeleton). For this definition to specify a coherent probability distribution over instantiations, we must ensure that our probabilistic dependencies are acyclic, so that a random variable does not depend, directly or indirectly, on its own value. Moreover, we want to guarantee that this will be the case for any skeleton. For this purpose, we use a class dependency graph, which describes all possible dependencies among attributes. In this graph, we have an (intra-object) edge $S \rightarrow T$ if S is a parent of T . If $S \rightarrow T$ is a parent of U , and $U \rightarrow V$, we have an (inter-object) edge $S \rightarrow V$. If the dependency graph of O is acyclic, then it defines a legal model for any relational skeleton L [14].

Definition 2: A probabilistic relational model O with reference uncertainty has the same components as in Definition 1. In addition, for each reference slot $\% \rightarrow Y$, we have:

a set of attributes Y, \dots

$$P(\mathcal{I} | \sigma_o, \Pi) = \prod_{x \in \sigma_o(X)} \prod_{A \in \mathcal{A}(x)} P(x.A | \text{Pa}(x.A)) \prod_{\rho \in \mathcal{R}(x), \text{Range}[\rho]=Y} \frac{P(x.S_\rho = v[x.\rho] | \text{Pa}(x.S_\rho))}{|Y_v|}$$

a new selector attribute V_6 within which takes on values in the cross-product space Y_v , a set of parents and a CPD for V_6 . To define the semantics of this extension, we must define the probability of reference slots as well as descriptive attributes:

5. Conclusion

We evaluated the methods on several real-life data sets, comparing standard PRMs, PRMs with reference uncertainty (RU), and PRMs with existence uncertainty (EU). Our experiments used the Bayesian score with a uniform Dirichlet parameter prior with equivalent sample size, and a uniform distribution over structures. We first tested whether the additional expressive power allows us to better capture regularities in the domain. Toward this end, we evaluated the likelihood of test data given our learned models. Unfortunately, we cannot directly compare likelihoods, since the PRMs involve different sets of probabilistic events. Instead, we compare the two variants of PRMs with structural uncertainty, EU and RU, to “baseline” models which incorporate link probabilities, but make the “null” assumption that the link structure is uncorrelated with the descriptive attributes. For reference uncertainty, the baseline has 7-8% for each slot. For existence uncertainty, it forces U_6 to have no parents in the model.

References

1. Murphy, K. P. Machine learning: a probabilistic perspective (MIT press, 2012).
2. Hinton, G. et al. Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. *Signal Processing Magazine, IEEE* 29, 82 (2012).
3. Krizhevsky, A., Sutskever, I. & Hinton, G. E. ImageNet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, 1097-1105 (2012).
4. Sermanet, P. et al. Overfeat: Integrated recognition, localization and detection using convolutional networks. In *International Conference on Learning Representations (ICLR 2014)* (arXiv preprint arXiv:1312.6229, 2014).
5. Ghahramani, Z. Bayesian nonparametrics and the probabilistic approach to modelling. *Philosophical Transactions of the Royal Society A* 371, 20110553 (2013).
6. Marcus, G. F. & Davis, E. How robust are probabilistic models of higher-level cognition? *Psychological Science* 24, 2351-2360 (2013).
7. Goodman, N. D. et al. Relevant and robust a response to Marcus and Davis (2013). *Psychological Science* 0956797614559544 (2015).
8. Doya, K., Ishii, S., Pouget, A. & Rao, R. P. N. *Bayesian Brain: Probabilistic Approaches to Neural Coding* (The MIT Press, 2007).
9. Neal, R. M. MCMC using Hamiltonian dynamics. In S. Brooks A. Gelman, G. J. & Meng, X.-L. (eds.) *Handbook of Markov Chain Monte Carlo* (Chapman & Hall / CRC Press, 2010).

10. Girolami, M. & Calderhead, B. Riemann manifold Langevin and Hamiltonian Monte Carlo methods. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 73, 123{214 (2011).
11. Sutskever, I., Vinyals, O. & Le, Q. V. Sequence to sequence learning with neural networks. In Ghahramani, Z., Welling, M., Cortes, C., Lawrence, N. & Weinberger, K. (eds.) *Advances in Neural Information Processing Systems* 27, 3104{3112 (Curran Associates, Inc., 2014).
12. Orbanz, P. & Teh, Y. W. Bayesian nonparametric models. In *Encyclopedia of Machine Learning* (Springer, 2010).
13. Hjort, N., Holmes, C., Müller, P. & Walker, S. (eds.) *Bayesian Nonparametrics* (Cambridge University Press, 2010).
14. Lu, C. & Tang, X. Surpassing human-level face verification performance on LFW with GaussianFace. In *Proceedings of the 29th AAAI Conference on Artificial Intelligence (AAAI)*(2015).