



Lecture notes on Bayesian machine learning

What is BML, why use it, and how to implement it.

Rémi Bardenet and Julyan Arbel



Copyright © 2021 Rémi Bardenet and Julyan Arbel

This template is adapted from Mathias Legrand's and Vel's *Orange Book* template v2.4, licensed under CC BY-NC-SA 3.0 (<http://creativecommons.org/licenses/by-nc-sa/3.0/>).

Contents

I	What is Bayesian machine learning?	
1	The example of penalized linear regression	9
1.1	Fisher does linear regression	9
1.2	Wald does linear regression	10
1.3	Savage does linear regression	11
1.4	Hoerl and Kennard do linear regression	11
2	Maximizing expected utility	13
2.1	ML problems are decision problems	13
2.2	Bayesians maximize expected utility	14
2.3	Specifying a joint model	14
2.4	The Bayes decision rule for common ML problems	14
II	Foundations	
3	Many (incompatible) reasons to be a Bayesian	17
3.1	Because you abide by the likelihood principle	17
3.1.1	The formal LP	17
3.1.2	SEU satisfies the LP	18
3.1.3	The stopping rule principle	18
3.1.4	Pros and cons of the LP	18

3.2	Because you place coherence above all things: subjective Bayes	19
3.2.1	A closer look at the axioms	19
3.2.2	Major criticisms	21
3.3	Because you like coherence and consensus: objective Bayes	22
3.4	Because you are a Waldian frequentist in disguise	22
3.4.1	On the consistency of Bayesian estimators	22
3.4.2	Complete class theorems	22
3.4.3	PAC-Bayes statistical learning	22

III

Implementing Bayesian machine learning

4	Markov chain Monte Carlo	25
4.1	Basic Monte Carlo	25
4.2	The Metropolis-Hastings algorithm	26
4.3	Gibbs sampling	27
4.4	Hamiltonian Monte Carlo	27
4.4.1	An abstract variant of Metropolis-Hastings	27
4.4.2	An augmented target	27
4.4.3	Hamiltonian dynamics	27
4.4.4	Ideal HMC and numerical HMC	28
4.4.5	On the ergodicity of HMC	28
4.4.6	An ubiquitous variant: NUTS	28
4.5	MCMC practice: convergence diagnostics	28
5	Variational Bayes	29
5.1	The evidence lower bound (ELBO)	29
5.2	Mean-field inference	30
5.2.1	Mean-field VB for LDA	31
5.2.2	Mean-field VB for marginal LDA	32
5.2.3	VB generalizes the EM algorithm	32
5.3	Gradient-based algorithms	33
5.3.1	VB for deep networks	33
5.4	Theoretical guarantees	34
5.5	Alternatives to the KL divergence	34

IV

Nonparametric Bayes

6	Random functions: Gaussian processes	37
7	Random probability measures: Dirichlet processes and the like	39
8	Asymptotic frequentist properties	41

V	Bayesian deep learning	
9	Bayesian neural networks	45
10	Modern questions	47
	Bibliography	49



What is Bayesian machine learning?

1	The example of penalized linear regression	9
1.1	Fisher does linear regression	
1.2	Wald does linear regression	
1.3	Savage does linear regression	
1.4	Hoerl and Kennard do linear regression	
2	Maximizing expected utility	13
2.1	ML problems are decision problems	
2.2	Bayesians maximize expected utility	
2.3	Specifying a joint model	
2.4	The Bayes decision rule for common ML problems	

1. The example of penalized linear regression

In this chapter, we introduce different historical approaches to learning and inference on a running example, penalized linear regression. We shall exaggerate the stances of some famous scientists, for the sake of illustration.

1.1 Fisher does linear regression

Say we have a data set (x_i, y_i) , $1 \leq i \leq N$, where $x_i \in \mathbb{R}^d$ are the *features* and $y_i \in \mathbb{R}$ the *response*. We want to study the influence of features on the response, and we ask British statistician Ronald Fisher (1890–1962) for help. He recommends positing a simple statistical model, i.e. a parametrized collection of PDFs for $y|x$. Calling the parameter θ and abusively denoting the parametrized distributions by $y|x, \theta$, we posit

$$p(y_i|x_i, \theta) = \mathcal{N}(y_i|x_i^T \theta, \sigma^2), 1 \leq i \leq N,$$

i.i.d., with known σ for simplicity. To characterize the influence of features on the response, Fisher recommends that we estimate θ , along with a confidence interval to communicate our uncertainty. Fisher was the one to introduce the maximum likelihood estimator

$$\hat{\theta}_{\text{MLE}} \triangleq \arg \max_{\theta} p(\mathbf{y}|X, \theta) = \arg \max_{\theta} \prod_{i=1}^N \mathcal{N}(y_i|x_i^T \theta, \sigma^2).$$

Now

$$\log \prod_{i=1}^N \mathcal{N}(y_i|x_i^T \theta, \sigma^2) = \log \mathcal{N}(\mathbf{y}|X\theta, \sigma^2 I) \propto -\|\mathbf{y} - X\theta\|^2,$$

where the proportionality sign allows us to drop additive terms that do not depend on θ . This entails that $\hat{\theta}_{\text{MLE}}$ is nothing but the least squares estimator

$$\hat{\theta}_{\text{MLE}} = (X^T X)^{-1} X^T \mathbf{y},$$

where we assumed that X has full rank. There is nothing inherently good in using the MLE rather than another estimator, but it often has good *frequentist* properties, i.e., properties established by integrating over the posited data generation model. For instance, it is easy to prove the following.

Proposition 1.1.1 Under $\mathbf{y} \sim \mathcal{N}(X\theta, \sigma^2 I)$, and assuming X has full rank,

$$\hat{\theta}_{\text{MLE}} \sim \mathcal{N}(\theta, \sigma^2 (X^T X)^{-1}). \quad (1.1)$$

Proof. Left as an exercise. ■

Fisher then argues that (1.1) gives you a wealth of properties that make $\hat{\theta}_{\text{MLE}}$ an interesting estimator. For starters, $\hat{\theta}_{\text{MLE}}$ is unbiased ($\mathbb{E}\hat{\theta}_{\text{MLE}} = \theta$) and $\hat{\theta}_{\text{MLE}} \rightarrow \theta$ in probability. Moreover, the (random) ellipsoid

$$\mathcal{E}_\alpha = \{\theta \in \mathbb{R}^d \text{ such that } \sigma^{-2}(\theta - \hat{\theta}_{\text{MLE}})^T X^T X (\theta - \hat{\theta}_{\text{MLE}}) \leq \alpha\}$$

contains θ with known probability $1 - \delta(\alpha)$. It thus makes sense, to Fisher, to find the smallest value $\alpha_{0.99}$ such that $1 - \delta(\alpha) \geq 0.99$, and report the *confidence region* $\mathcal{E}_{\alpha_{0.99}}$ to quantify his uncertainty about θ . The property that, under repetition of the data generation, the (random) ellipsoid $\mathcal{E}_{\alpha_{0.99}}$ contains θ about 99% of the time, is called *coverage*.

1.2 Wald does linear regression

Imagine that you had asked Hungarian-American statistician Abraham Wald (1902–1950) for help instead of Fisher. He would have argued that estimating θ , or outputting a region of \mathbb{R}^d that encodes your uncertainty, are two different actions to be taken under uncertainty.

Consider estimation. Wald would ask you to specify the loss $L(\theta, \hat{\theta})$ that you incur by estimating θ by $\hat{\theta} \in \mathbb{R}^d$. You might answer $L_\alpha(\theta, \hat{\theta}) = \|\theta - \hat{\theta}\|^\alpha$ for some $\alpha > 0$. Wald would then say that the accuracy of an estimator $\hat{\theta} = \hat{\theta}(\mathbf{y})$ is characterized by its risk function

$$R(\cdot, \hat{\theta}) : \theta \mapsto \mathbb{E}_{\mathbf{y}|X, \theta} L(\theta, \hat{\theta}(\mathbf{y})). \quad (1.2)$$

If the risk function of an estimator is smaller than that of another estimator, we say that the former dominates the latter. Wald's minimum requirement for an estimator is that the estimator is *admissible*, i.e., that it is not dominated. Maybe surprisingly, the MLE for regression with the squared loss L_2 is *not* admissible! This is an extension by **TBC** of an important theorem known as the James-Stein theorem (**TBC**). We leave the original James-Stein theorem as Exercise ???. We shall see an even simpler estimator that dominates the MLE in Section ??.

So what should we do, according to Wald, if the MLE is no option? It would be natural to find an estimator with as small as possible a risk function. Unfortunately, the risk function being a function, many pairs of estimators are incomparable. Wald might recommend to sum up a risk function by a single number, and look, for instance, for an estimator minimizing the worst-case risk, a so-called *minimax* estimator

$$\hat{\theta}_{\text{minimax}} \in \arg \inf_{\hat{\theta}} \sup_{\theta} R(\theta, \hat{\theta}). \quad (1.3)$$

Another solution to sum up the risk function is to integrate it against a measure for which the risk is integrable. This leads to an estimator that we call the *Bayes* estimator, anticipating over Chapter 2,

$$\hat{\theta}_{\text{Bayes}} \in \arg \min_{\hat{\theta}} \mathbb{E}_{\theta} R(\theta, \hat{\theta}) = \mathbb{E}_{\theta} \mathbb{E}_{\mathbf{y}|X, \theta} R(\theta, \hat{\theta}). \quad (1.4)$$

We now have a joint distribution over θ, \mathbf{y} . As long as the loss $L(\theta, \hat{\theta}(\mathbf{y}))$ is integrable w.r.t. this joint, the towering property of the expectation yields

$$\mathbb{E}_{\theta} \mathbb{E}_{\mathbf{y}|X, \theta} L(\theta, \hat{\theta}) = \mathbb{E}_{\mathbf{y}} \mathbb{E}_{\theta|X, \mathbf{y}} L(\theta, \hat{\theta}). \quad (1.5)$$

Since $\hat{\theta}$ is only a function \mathbf{y} , minimizing (1.5) boils down to setting

$$\hat{\theta}_{\text{Bayes}} = \hat{\theta}_{\text{Bayes}}(\mathbf{y}) = \arg \inf_{\hat{\theta}} \mathbb{E}_{\theta|\mathbf{X},\mathbf{y}} L(\theta, \hat{\theta}).$$

For the squared loss $L = L_2$, the Bayes estimator is thus the mean of the *posterior* distribution, i.e. $\mathbb{E}_{\theta|\mathbf{X},\mathbf{y}} \theta$. For the one-loss $L = L_1$, the Bayes estimator is a generalized median of the same posterior distribution. In particular, note that the Bayes estimator depends on the loss function, and is not necessarily the posterior expectation.

1.3 Savage does linear regression

1.4 Hoerl and Kennard do linear regression

2. Maximizing expected utility

After formally defining decision problems, we show that basic machine learning problems such as classification, regression, model choice, etc. are decision problems. Then we introduce subjective expected utility, the single unique guideline of all Bayesians, and go over its consequences for ML decisions. Justifying subjective expected utility shall wait until Chapter II.

For this chapter, we mostly used two references. (PaIn09) focuses on ideas and is a formidable entry point to (Bayesian) decision theory, while (Sch12) is a great textbook-level reference for advanced reading, with mathematical details.

2.1 ML problems are decision problems

Formally, a decision problem is defined as

1. a set \mathcal{S} of *states*. For technical reasons, we also require a σ -algebra $\Sigma_{\mathcal{S}}$ that makes $(\mathcal{S}, \Sigma_{\mathcal{S}})$ a Borel space (Sch12).
2. a set of *rewards* \mathcal{R} . We also require a σ -algebra $\Sigma_{\mathcal{R}}$, which should contain all singletons.
3. a set \mathcal{A} of measurable functions from \mathcal{S} to \mathcal{R} , called *actions*.
4. a utility function $u : \mathcal{R} \rightarrow \mathbb{R}$.

We think of states as encoding all information about the situation at hand. Actions are what we are tasked to choose, and picking action a while the situation is described by a given state s leads to reward $a(s)$. Note that we assume here that the same set of actions \mathcal{A} remains available in every state s .¹

Most basic ML problems are of this kind; see Figure 2.1 for a few classical formalizations. Note that most choices made in this table are arbitrary, and correspond to the simplest variant of each problem. For instance, in classification, one might penalize false negatives and false positives differently; see Exercises. Note also that, since Wal50 and as done in Section ??, it is also customary to define loss functions instead of utilities, as $L(a, s) = -u(a(s))$. At this point of the document, both notations are as expressive, and we shall use them interchangeably. The distinction

¹In future versions of the course, we might make the framework more general, to include, e.g. Markov decision processes.

	\mathcal{S}	\mathcal{R}	\mathcal{A}	$u(r)$
Regression	$(\mathcal{X} \times \mathcal{Y})^{n+1}$	$\mathcal{Y} = \mathbb{R}$	$\{a_g : s \mapsto y - g(x; x_{1:n}, y_{1:n})\}$	$-\ r\ ^2$
Classification	$(\mathcal{X} \times \mathcal{Y})^{n+1}$	$\mathcal{Y} = \{0, 1\}$	$\{a_g : s \mapsto y - g(x; x_{1:n}, y_{1:n})\}$	$\mathbb{1}_{\{r=0\}}$
Point estimation	$\mathcal{Y}^n \times \Theta$	Θ	$\{a_g : s \mapsto \theta - g(y_{1:n})\}$	$-\ r\ ^2$
Interval estimation	$\mathcal{Y}^n \times \Theta$	$\{0, 1\} \times \mathbb{R}_+$	$\{a_g : s \mapsto (\mathbb{1}_{\{\theta \in g(y_{1:n})\}}, g(y_{1:n}))\}$	$r_1 + \gamma r_2$
Model choice	$\mathcal{Y}^n \times (\cup_{m=1}^M \{m\} \times \Theta_m)$	$\{0, 1\}$	$\{a_g : s \mapsto \mathbb{1}_{\{m=g(y_{1:n})\}}\}$	r

Figure 2.1: Some classical formalizations of ML problems as decision problems. Actions are labeled by functions g (“predictors”), the domain and codomain of which should be obvious from the definition; for instance g outputs a $\{0, 1\}$ label in classification, and a Borel subset of Θ in “interval” estimation.

will come later in Chapter II, when we discuss state-dependent utilities.

2.2 Bayesians maximize expected utility

By definition, a Bayesian is someone who, facing a decision problem, picks a joint distribution p over $(\mathcal{S}, \Sigma_{\mathcal{S}})$, and chooses action

$$a^* \in \arg \max_{a \in \mathcal{A}} \mathbb{E}_{s \sim p} u(a(s)). \quad (2.1)$$

The principle of decision-making encoded by (2.1) is called *expected utility*, or *subjective expected utility*, to insist on the fact that p is an arbitrary choice from the decision maker. At this stage, we have not discussed how Bayesians choose p , and there are many ways to do so; see Section ???. Finally, to give an example of Bayesian decision, the ridge regression estimator is the Bayes action for a particular decision problem and joint distribution over states; see Section ???.

In ML, it is customary to split the state variable into (s_O, s_U) , where $O, U \subset \{1, \dots, \dim \mathcal{S}\}$ are disjoint subsets that respectively index observed states (“data”) and unknown states. In particular, we can rewrite (2.1) as

$$a^* \in \arg \max_{a \in \mathcal{A}} \mathbb{E}_{s_O} \mathbb{E}_{s_U | s_O} u(a(s)). \quad (2.2)$$

Now assume that actions are labeled by a “predictor”, which maps s_O to one or several of the unknown variables of interest, say s_I for some $I \subset U$. This is the case for all rows in Figure 2.1. In classification or regression, for instance, the variable of interest is the new label y , and actions are labeled by measurable predictors of this variable of interest: evaluating $g(x; x_{1:n}, y_{1:n})$ is thought of as training a given algorithm (say, an SVM) over $\{(x_i, y_i), 1 \leq i \leq n\}$ and evaluating the corresponding predictor at the new feature vector x . Now, to maximize (2.2) over \mathcal{A} , it is enough maximize it over g . And since g is a function of s_O only, the optimal g is

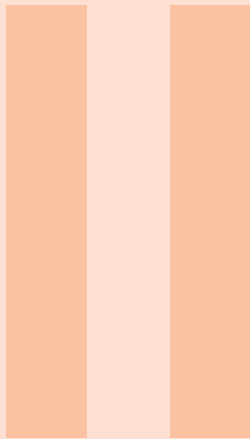
$$g^* : s_O \mapsto \arg \max_g \mathbb{E}_{s_U | s_O} u(a_g(s)). \quad (2.3)$$

Indeed, by maximizing the innermost expectation in (2.1) for each fixed value of s_O , we maximize the whole expectation. The resulting g^* is called the *Bayes decision rule*.

2.3 Specifying a joint model

We assume here familiarity with probabilistic graphical models, to the point of telling from a graph whether two sets of nodes are independent given a third one. The reader needing a recap is referred to (Mur12).

2.4 The Bayes decision rule for common ML problems



Foundations

3	Many (incompatible) reasons to be a Bayesian	17
3.1	Because you abide by the likelihood principle	
3.2	Because you place coherence above all things: subjective Bayes	
3.3	Because you like coherence and consensus: objective Bayes	
3.4	Because you are a Waldian frequentist in disguise	

3. Many (incompatible) reasons to be a Bayesian

So far, the main appeal of subjective expected utility (SEU) has been its conceptual simplicity, and the fact that it answers all decision problems in the same manner. In this chapter, we review some attempts at justifying SEU more formally, i.e. show that SEU logically follows from some simple, consensual principle. When you hear someone say that “Non-Bayesians are incoherent”, that “a prior encapsulates the information available *before* an experiment is made, so that the prior cannot depend on data”, or that “the posterior is the experimenter’s updated belief after the experiment”, or that “Non-Bayesians violate the likelihood principle”, they are all referring to one or the other of the formal justifications below. Not all these justifications are compatible, and none can really claim to be uniformly superior, so it is important to know upon what arguments you are resting your interpretation of the Bayesian procedure.

3.1 Because you abide by the likelihood principle

The “formal” likelihood principle (**BeWo88**) is a half-formal justification that deals primarily with the estimation problem. It is only half-formal because it deals with notions that are hard to rigorously define.

3.1.1 The formal LP

Consider two statistical experiments

$$E_i = (\mathcal{Y}_i, \mathbf{s}, \{p_i(\cdot|\vartheta), \vartheta \in \Theta\}), \quad i = 1, 2.$$

Assume that for some realizations \mathbf{y}_1 and \mathbf{y}_2 ,

$$p_1(\mathbf{y}_1|\cdot) \propto p_2(\mathbf{y}_2|\cdot).$$

Note that we are using bold characters to insist on \mathbf{y}_1 and \mathbf{y}_2 being vectors concatenating (possibly many, of arbitrary dimension) observations, the label $i = 1, 2$ is only there to indicate which experiment we consider. In particular, \mathbf{y}_1 and \mathbf{y}_2 may differ in dimension.

Now, assuming that there exists a quantity $\text{Ev}(E, x)$ that encapsulates the “evidence on θ arising from E and x ”, the formal LP principle is the requirement that

$$\text{Ev}(E_1, \mathbf{y}_1) = \text{Ev}(E_2, \mathbf{y}_2).$$

As a corollary, $\text{Ev}(E, \mathbf{y})$ can depend on x solely through $p(\mathbf{y}|\cdot)$.

3.1.2 SEU satisfies the LP

Letting $\mathcal{S} = \mathcal{Y}^n \times \Theta$, SEU satisfies the LP as long as the joint distribution over states has either p_1 or p_2 as its conditional of \mathbf{y} given θ . Indeed, let

$$p_i(s_i) = p_i(\mathbf{y}_i, \theta) = p_i(\mathbf{y}_i|\theta)p(\theta) = \mathcal{Z}p_i(\theta|\mathbf{y}_i), \quad i = 1, 2.$$

Note how we use a common prior. Then for $a : \mathcal{S} \rightarrow \mathcal{X}$,

$$\int L(a, s_1) \frac{p_1(\mathbf{y}_1|\theta)p(\theta)}{\mathcal{Z}} d\theta \propto \int L(a, s_2) \frac{p_2(\mathbf{y}_2|\theta)p(\theta)}{\mathcal{Z}} d\theta,$$

so that the posterior expected losses are the same in both experiments, and Bayes actions coincide. However, note that full expected utilities are different in general,

$$\int L(a, s_1) p_1(\mathbf{y}_1|\theta)p(\theta) d\mathbf{y}_1 d\theta \neq \int L(a, s_2) p_2(\mathbf{y}_2|\theta)p(\theta) d\mathbf{y}_2 d\theta.$$

3.1.3 The stopping rule principle

The same kind of computations shows that SEU with a particular choice of joint distribution is immune to data-dependent stopping rules. This can also be seen as a consequence of the LP (BeWo88), but we stick to SEU with some conditions on its joint distribution of states for simplicity.

Assume that we want to model the following inference problem. We collect data one item at a time, independently from some distribution $y_i|\theta$, until the first $n \in \mathbb{N}$ such that $y_1, \dots, y_n \in A_n$, and then you want to estimate θ . We model this by $\mathcal{S} = \Theta \times \cup_{n \geq 1} \mathcal{Y}^n$, and decide to take a joint distribution p such that y_1, y_2, \dots are independent given θ , just like we assume the data generating mechanism works. Then the Bayes action $a^* = a_{g^*}$ minimizes

$$\begin{aligned} \mathbb{E}L(a, s) &= \mathbb{E} \left[L(a, s) \sum_n \mathbb{1}_{\{N=n\}} \right] \\ &= \sum_n \mathbb{E} [L(a, s) \mathbb{1}_{\{N=n\}}] \\ &= \sum_n \int L(a, (\theta, y_{1:n})) \mathbb{1}_{\{y_{1:n} \in A_n\}} \prod_{k < n} \mathbb{1}_{\{y_{1:k} \notin A_k\}} p(y_{1:n}|\theta) p(\theta) dy_{1:n} d\theta. \\ &= \sum_n \int dy_{1:n} \mathbb{1}_{\{y_{1:n} \in A_n\}} \prod_{k < n} \mathbb{1}_{\{y_{1:k} \notin A_k\}} \int L(a, (\theta, y_{1:n})) p(y_{1:n}|\theta) p(\theta), \end{aligned}$$

where we used the monotone convergence theorem and Fubini’s theorem (assuming, e.g., that the loss is bounded). So, to find the minimizer g^* defined on $\cup_n \mathcal{Y}^n$ of the overall expected loss, it is enough, for each n , to define $g^*(y_{1:n})$ as the usual Bayes rule for fixed n , i.e. as the minimizer of the inner integral. In other words, as long as the prior $p(\theta)$ does not depend on data, the Bayes decision is immune to data-dependent stopping rules: just act as if there were no stopping rule.

3.1.4 Pros and cons of the LP

- The LP is compelling to many (BeWo88), but it has its downsides.

- Being Bayesian is not the only way to abide by the LP.
- I am personally uncomfortable with the stopping rule principle, probably because my frequentist intuition is still too strong.
- It is hard to make fully formal: is $\text{Ev}(E, x)$ even meaningful? See answer by LeCam to (BeWo88).
- It assumes we want to specify a likelihood, this prevents model-free Bayesianism.
- It separates the roles of the likelihood and the prior. For LP-abiding Bayesians, **the prior is not allowed to depend on data.**

3.2 Because you place coherence above all things: subjective Bayes

The literature on the foundations of subjective Bayes is rich, and we refer to (PaIn09) for entry points. A major milestone was obtained by Sav53, building on work of VoMo and Ram. Savage gave a list of properties (the so-called *Savage axioms* of coherence) that a binary relation \succ on the action space \mathcal{A} should satisfy in order for it to have an essentially unique expected utility representation. More precisely, \succ satisfies the Savage axioms if and only if there is a bounded utility function $u : \mathcal{S} \rightarrow \mathbb{R}_+$ and a (finitely additive) probability measure p on \mathcal{S} such that for $a, b \in \mathcal{A}$,

$$a \succ b \Leftrightarrow \int u(a(s))dp(s) > \int u(b(s))dp(s);$$

u is unique up to affine transformations. The pair (u, p) thus characterizes the behaviour of a Savage-abiding decision maker: p can be interpreted as the DM's degree of belief in the subsets of \mathcal{S} , and the DM's preferred action is the one maximizing SEU according to the (u, p) pair. There is no constraint on the probability p : any pair (u, p) corresponds to a coherent decision-maker.

Savage's result is beautiful, in that it builds a probability measure p and a utility function u over states from a coherent (in the sense of satisfying the Savage axioms) ranking of actions. Following Savage, there is no need to assume that the phenomenon we are studying is probabilistic, or to disentangle different forms of uncertainty: a coherent decision maker *must* have a utility and a probability in mind!

3.2.1 A closer look at the axioms

Sav54 derives a unique (finitely additive) probability on $2^{\mathcal{S}}$ and utility function directly from a preference relation over acts $a : \mathcal{S} \rightarrow \mathcal{Z}$. He relies on NM, but without compound acts (considered as artificial) and no physical chance mechanism (unlike Ramsey or Anscombe & Aumann). Since axiom numbers 'Px' in the original paper of Sav54 are often referred to, I use them as labels. Yet I introduce axioms in the same order as PaIn09.

Axiom P1 (Preference). \succ is complete and transitive.

Axiom P5 (No total indifference). $\exists z_1, z_2 \in \mathcal{Z}$ such that $z_1 \succ z_2$.

Now for $a, b \in \mathcal{A}$, $T \in \mathcal{S}$, define action a_T^b by $a_T^b(\theta) = a(\theta)1_{\theta \in T} + b(\theta)1_{\theta \in T^c}$.

Axiom P2 (Sure Thing principle). $\forall a, b, h_1, h_2 \in \mathcal{A}$ and $\forall T \subset \mathcal{S}$,

$$a_T^{h_1} \succ b_T^{h_1} \Leftrightarrow a_T^{h_2} \succ b_T^{h_2}.$$

Here, Savage directly formulates that if two acts coincide on part of \mathcal{S} , then preference should only depend on their values where they differ. As a side remark, we can now define a new family of preference relations, called *conditional preferences*. Let $T \subset \mathcal{S}$, and define $a \succ b|T$ by

$$c_T^a \succ c_T^b \tag{3.1}$$

for some (equivalently all) $c \in \mathcal{A}$. This family of conditional preferences will be discussed later, when we consider how to rank actions after some T has been observed.

Define now a null state as a $T \subset \mathcal{S}$ such that $a \sim b|T$ for all $a, b \in \mathcal{A}$. That is, preferences are insensitive to T obtaining.

Axiom P3 (No reduction to conditional indifference). *If T is not null, then $a_T^{z_1} \succ a_T^{z_2}|T$ iff $z_1 \succ z_2$.*

In **P3 (No reduction to conditional indifference)**, we have identified $z \in \mathcal{Z}$ with the constant action $z : s \mapsto z$. Note that by **P2 (Sure Thing principle)**, the values of a and b outside T do not matter. Axiom **P3 (No reduction to conditional indifference)** makes sure that no preference among consequences can be reduced to indifference conditional on a non-null event. I do not fully grasp this axiom, but it seems harmless.

Axiom P4 (Separation). *Assume that $z_1 \succ z_2$ and $z'_1 \succ z'_2$. Let $T_1, T_2 \subset \mathcal{S}$, and*

$$a = z_{2T_1}^{z_1}, b = z_{2T_2}^{z_1}, a' = z_{2T_1}^{z'_1}, b' = z_{2T_2}^{z'_1}.$$

Then $a \succ b \Leftrightarrow a' \succ b'$.

In words, for two actions that are both piecewise constant taking the same set of values, I can change the constant values without altering the preference, as long as I preserve the order between outcomes. Intuitively, if outcomes were money, your willingness to bet on T_1 obtaining rather than T_2 obtaining does not depend on how much money you make/lose in each of these two bets.

We can now define a binary relation $T_1 \succ T_2$ over subsets of \mathcal{S} that we could interpret as “ T_1 is more likely to me than T_2 ”. We say $T_1 \succ T_2$ if for all $z_1, z_2 \in \mathcal{Z}$ such that $z_1 \succ z_2$,

$$z_{2T_1}^{z_1} \succ z_{2T_2}^{z_1}.$$

In the betting metaphor, you’re more willing to bet on T_1 obtaining than T_2 for the same rewards.

These first five axioms are enough to get the following representation.

Theorem 3.2.1 — Qualitative representation. If **P1 (Preference)**, **P2 (Sure Thing principle)**, **P3 (No reduction to conditional indifference)**, **P4 (Separation)**, and **P5 (No total indifference)** hold, then \succ on \mathcal{S} is a qualitative probability, that is

- \succ is negatively transitive: $T_1 \preceq T_2 \preceq T_3$ implies $T_1 \preceq T_3$,
- $\forall R \subset \mathcal{S}, T \succ \emptyset$.
- If $T_1 \cap U = T_2 \cap U = \emptyset$, then $T_1 \succ T_2$ iff $T_1 \cup U \succ T_2 \cup U$.

We could stop here and work with beliefs specified by pairwise rankings of events. But **Sav54** goes forward, and keeps adding axioms to get a more usual *quantitative* probability. In particular, we need an additional structural axiom on \mathcal{S} . There are variations of this “partition axiom”, and Savage chooses to embed it in its Archimedean axiom.

Axiom P6 (Archimedean). $\forall a, b \in \mathcal{A}$ such that $a \succ b$, and $\forall z \in \mathcal{Z}$, there exists a finite partition of $\mathcal{S} = T_1 \cup \dots \cup T_M$ such that for all T_i , either $a_{T_i}^z \succ b$ or $a \succ b_{T_i}^z$.

An important consequence of **P6 (Archimedean)** is that \mathcal{S} must be rich enough to be splittable into tiny pieces suiting any pair of actions. This is in general not possible for a discrete space, and we will generally have to use a continuous state space.

Theorem 3.2.2 — Qualitative representation. Under **P1 (Preference)**, **P2 (Sure Thing principle)**, **P3 (No reduction to conditional indifference)**, **P4 (Separation)**, **P5 (No total indifference)**, and **P6 (Archimedean)**, there exists a unique finitely additive probability measure Π on $2^{\mathcal{S}}$ such that

- $T_1 \succ T_2$ iff $\Pi(T_1) > \Pi(T_2)$,

- $\forall T_1 \subset \mathcal{S}$ and $k \in [0, 1]$, there exists $T_2 \subset \mathcal{S}$ such that $\Pi(T_2) = k\Pi(T_1)$.

Now to obtain a full expected utility representation, we need a final axiom.

Axiom P7 (No aversion to risk). $\forall T \subset \mathcal{S}$,

- If $\forall s \in T$, $a \succ b(s)|T$, then $a \succ b|T$.
- If $\forall s \in T$, $a(s) \succ b|T$, then $a \succ b|T$.

I call this *no aversion to risk*, since it intuitively means that you do not make a difference between preferring a over b to preferring a over any sure consequence of b , and vice versa. I guess this plays a role in paradoxes like Ellsberg's where people can revert preferences in favour of more sure outcomes.

Theorem 3.2.3 — Qualitative representation. Under P1 (Preference), P2 (Sure Thing principle), P3 (No reduction to conditional indifference), P4 (Separation), P5 (No total indifference), P6 (Archimedean), and P7 (No aversion to risk), there exists a unique finitely additive probability measure satisfying the results of Theorem 3.2.2. Furthermore, there is a unique (up to positive affine transformations) bounded utility function $u : \mathcal{Z} \rightarrow \mathbb{R}$ such that

$$a \succ b \Leftrightarrow \int u(a(s))d\pi > \int u(b(s))d\pi.$$

It is crucial to note that the probability measure in Theorem 3.2.3 is finitely additive, and defined on the Boolean algebra $2^{\mathcal{S}}$ of all subsets of \mathcal{S} . In particular, the integrals in Theorem 3.2.3 are *not* Lebesgue integrals, see e.g. **Kre88**.

3.2.2 Major criticisms

Yet there are several points that can be raised against using Savage's result to justify Bayesian learning, which we roughly rank by increasing seriousness. First, the resulting probability measure is only finitely additive: it is hard to bring a σ -algebra like $\Sigma_{\mathcal{S}}$ into the picture with a consensual coherence axiom that does not involve an infinite number of choices from the DM. This is a minor inconvenience, as we can always restrict actions to a set of measurable functions and choose a *bona fide* probability distribution. The price is that we might be unable to represent all coherent behaviour. A second objection is the boundedness of the utility function, first noticed by (**Fis70**). Common losses such as the squared loss need to be trimmed to fit into the picture without modifying the axioms. At the cost of some less natural axioms, one can however accomodate unbounded loss functions (**Wak93**). A third objection is the pivotal role played by constant actions in Savage's (and actually von Neumann and Morgenstern's) proof. Constant actions are those that assign the same reward $r \in \mathcal{R}$ to all states $s \in \mathcal{S}$. The utility function in approaches derived from **VoMo** is defined by finding the weight in a convex combination of two extreme constant actions; see e.g. (**PaIn09**). But constant actions are not part of the action sets that we considered in Section ?? . In binary classification, for instance, there are two constant actions: the ideal classifier a^* that always predicts the right label, and the worst possible classifier a_* , which consistently predicts the wrong label. The utility function of any reward r is defined in reference to these two idealized classifiers. It would be more satisfying to have a set of axioms that does not give such a key role to *non-physical* classifiers. A fourth and related objection is that the utility in Savage's result is independent of the state: $u(a(s))$ only depends on the state through $a(s)$. This means that the notation $L(a, s) = -u(a(s))$ is inappropriate, as it lets the user think that the loss of action a can depend on the state s . In textbook ML tasks, this is not much of a problem, but in real applications, this can prove limiting; see (**PaIn09**).

Maybe the most important weakness of axiomatic constructions like Savage's is that they do not prescribe what action to choose after some part of the state is observed. Using the conditional SEU is only a natural candidate for what to do, not more. We should not be fooled by the name "conditional expectation": call it "strange mathematical construction" and then try to think if this is how you want to change your behaviour once you observe an event. There are ways to justify conditioning by means of sequential Dutch books (), but none that is not controversial in the general case of an infinite state space. Worse: it is doubtful that this can be done, since in practice, we all refrain from conditioning from time to time. For instance, if you observe that your classifier has poor test error, you change the family \mathcal{G} of classifiers that you consider. Unless you included the choice of \mathcal{G} in your state space, this kind of "external" move, which looks at some calibration property of a subjective model, is not constrained by Savage-like constructions. Modern attempts to build Bayesian statistics with more focus on calibration include the prequential view of **Daw**. One could also argue that PAC-Bayes goes along those lines, by abandoning conditioning in favour of better frequentist properties; more about this in Section ??.

3.3 Because you like coherence and consensus: objective Bayes

3.4 Because you are a Waldian frequentist in disguise

3.4.1 On the consistency of Bayesian estimators

Restrict to parametric Bernstein von-Mises and its consequences.

3.4.2 Complete class theorems

3.4.3 PAC-Bayes statistical learning



Implementing Bayesian machine learning

4	Markov chain Monte Carlo	25
4.1	Basic Monte Carlo	
4.2	The Metropolis-Hastings algorithm	
4.3	Gibbs sampling	
4.4	Hamiltonian Monte Carlo	
4.5	MCMC practice: convergence diagnostics	
5	Variational Bayes	29
5.1	The evidence lower bound (ELBO)	
5.2	Mean-field inference	
5.3	Gradient-based algorithms	
5.4	Theoretical guarantees	
5.5	Alternatives to the KL divergence	

4. Markov chain Monte Carlo

Maximizing expected utility requires computing integrals. Numerical integration consists in finding T nodes x_t and weights w_t , such that

$$\mathcal{E}_T(f) \triangleq \int f(x) d\mu(x) - \sum_{t=1}^T w_t f(\theta_t) = o_{T \rightarrow \infty}(1), \quad \forall f: \mathbb{R}^d \rightarrow \mathbb{R} \in \mathcal{C},$$

where \mathcal{C} is a large class of functions. For any smooth f , Riemann-like (i.e., grid-based) integration leads to $\mathcal{E}_T(f) \sim \frac{\sqrt{d}}{T^{1/d}}$, so that grids become essentially useless beyond $d = 3$. Monte Carlo methods are randomized constructions of nodes and weights. Since $\mathcal{E}_T(f)$ is then random, statements on the error shall be with high probability, in expectation, about the asymptotic fluctuations, etc.

We shall see that Monte Carlo methods (i) can accommodate settings where the target μ in (??) is only available through the evaluation of an unnormalized density $d\mu(x) = \pi(x)dx = \pi_u(x)/Zdx$, as is almost always the case in Bayesian learning; and that (ii) some Monte Carlo methods have an error that scales only polynomially with the dimension of the support of the integrand.

References.

The reference book for basic Monte Carlo methods is (RoCa04). For theoretical results on MCMC, we refer to (DoMoSt14). For more recent methods, we shall refer to papers. For demos of MCMC samplers, see <https://chi-feng.github.io/mcmc-demo/>. Finally, we shall not cover deterministic alternatives to Monte Carlo methods in large dimension, see quasi-Monte Carlo methods (DiPi10).

4.1 Basic Monte Carlo

If we knew how to sample from π , we could take $x_t \sim \pi$ i.i.d., $w_t = 1/T$. Chebyshev's inequality would lead to

$$\mathbb{P}\left(\mathcal{E}_T(f) \geq \alpha \frac{\sigma(f)}{\sqrt{T}}\right) \leq \frac{1}{\alpha^2}, \quad \forall \alpha > 0,$$

as soon as $\sigma(f)^2 := \mathbb{E}_{X \sim \pi}[f(X) - \int f(x)\pi(x)dx]^2 < +\infty$.

In practice, we never have access to a sampler of π , so we choose an instrumental distribution $q(x)dx$, sample x_t i.i.d. from q . If we can evaluate π , then $w_t = \pi(x_t)/q(x_t)$ leads to an unbiased estimator called the *importance sampling*. Its error can also be controlled by the Chebyshev inequality. But more often than not, we can only evaluate an unnormalized density π_u . An alternative is to then set $w_t \propto \pi_u(x_t)/q(x_t)$ and normalize the weights so that $\sum_t w_t = 1$. This leads to the *self-normalized importance sampling estimator*

$$\hat{I}_T^{\text{NIS}} = \frac{\sum_{t=1}^T \frac{\pi_u(\theta_t)}{q(\theta_t)} f(x_t)}{\sum_{t=1}^T \frac{\pi_u(\theta_t)}{q(\theta_t)}} \rightarrow \frac{\int f(x) \pi_u(x) dx}{\int \pi_u(x) dx} = \int f(x) dx,$$

where the convergence is almost sure (apply the strong law of large numbers to both the numerator and the denominator).

Proposition 4.1.1 — CLT for NIS. The NIS estimator satisfies

$$\sqrt{T} \left(\hat{I}_T^{\text{NIS}} - \int f(x) \pi(x) dx \right) \rightarrow_d \mathcal{N}(0, \sigma_{\text{NIS}}^2(f))$$

where f is such that

$$\sigma_{\text{NIS}}^2(f) \triangleq TBC < \infty.$$

Proof. See exercise sheet. ■

Unfortunately, while NIS does accomodate unnormalized targets, it does not solve the curse of dimensionality: even when π and q are both Gaussian, only with different covariance matrices, one can show that

$$\log \sigma_{\text{NIS}}^2(f) = \Theta(d).$$

4.2 The Metropolis-Hastings algorithm

Metropolis-Hastings (MH) is the archetypal MCMC algorithm, and is still the main building block of modern MCMC algorithms. The idea is to take the nodes as the truncated history of a Markov chain (X_t) , which we build so as to guarantee that $\mathcal{E}_T(f) \rightarrow 0$. To see how to build (X_t) , remember first the law of large numbers for Markov chains.

Proposition 4.2.1 — LLN for Markov chains; see e.g. DoMoSt14. Let $(X_t)_{t \in \mathbb{N}}$ be a Markov chain with Markov kernel P . If

1. There exists π s.t.

$$\int d\pi(x) P(x, B) = \pi(B).$$

2. For any A with $\pi(A) > 0$, for any $\theta \in \Theta$,

$$\mathbb{P}_x \left(\sum_{t=0}^{\infty} 1_{\theta_t \in A} = +\infty \right) = 1,$$

then for any initial distribution μ_0 of X_0 , almost surely

$$\frac{1}{T} \sum_{t=1}^T f(\theta_t) \rightarrow \int f d\pi,$$

where $f \in L^1(\pi)$.

The first condition states that π is an *invariant distribution* of the chain: if $X_t \sim \pi$, then $X_{t+1} \sim \pi$. The second condition is called *Harris recurrence*: starting from any x , i.e. $X_0 \sim \delta_x$, then the chain returns an infinite number of times to A almost surely, as soon as A is charged by π . Intuitively, this makes sure that there is an infinity of nodes on A , so that the integral of f on A is accurately estimated.

4.3 Gibbs sampling

4.4 Hamiltonian Monte Carlo

We closely follow **BoSa18**, who give a clear introduction to the ingredients of HMC. To facilitate going back and forth between the notes and (**BoSa18**), we temporarily adopt their notational conventions: the variable over which we wish to integrate is $q \in \mathbb{R}^d$, while $x \in \mathcal{X}$ denotes a generic variable, later taken to be $x = (p, q)$. Note that vanilla HMC is limited to continuous variables.

4.4.1 An abstract variant of Metropolis-Hastings

Let S be a linear involution of $\mathcal{X} \subset \mathbb{R}^{2d}$, such that $\eta \circ S = \eta$ for some (possibly unnormalized) PDF η . Let further $\Phi : \mathbb{R}^{2d} \rightarrow \mathbb{R}^{2d}$ be a C^1 -diffeomorphism that is reversible w.r.t. to S , that is, $S \circ \Phi = \Phi^{-1} \circ S$. Now let

$$\alpha(x) \triangleq 1 \wedge \frac{\eta(\Phi(x))}{\eta(x)} |\Phi'(x)|, \quad (4.1)$$

and consider the Markov kernel

$$P_{aHMC}(x, A) = \alpha(x) 1_{\Phi(x) \in A} + (1 - \alpha(x)) 1_{S(x) \in A}.$$

Algorithmically, this “abstract” HMC kernel corresponds to accepting $\Phi(x)$ with probability $\alpha(x)$, and otherwise setting the new Markov state to $S(x)$.

Proposition 4.4.1 P_{aHMC} leaves η invariant.

Proof. Left as an exercise. ■

4.4.2 An augmented target

Consider the PDF on \mathbb{R}^{2d} defined by $\tilde{\pi}(q, p) = \frac{1}{2} \mathcal{N}(p|0, M) \times \pi(q)$. Clearly, the p -marginal is Gaussian, while the q -marginal is π . If we manage to obtain an MCMC chain for $\tilde{\pi}$, i.e. a chain with a Markov kernel that leaves $\tilde{\pi}$ invariant, then simply discarding the p -component of every realization will yield a chain that is invariant w.r.t. π . This is an example of *augmentation* of the state space: unlike for the collapsed Gibbs sampling of Section ??, we augment the dimensionality of the problem in the hope to make sampling easier. Our hope is justified here by the fact that we know how to efficiently move in \mathbb{R}^{2d} along the level lines of the augmented density $\tilde{\pi}$: this is where Hamiltonian dynamics come into play.

4.4.3 Hamiltonian dynamics

Let $H(q, p) = \log \tilde{\pi}(q, p)$, and consider the differential equation

$$\frac{d}{dt} \begin{pmatrix} q \\ p \end{pmatrix} = J \nabla H(q, p), \quad \text{where } J = \begin{pmatrix} 0_d & -I_d \\ I_d & 0_d \end{pmatrix}. \quad (4.2)$$

In particular, if ∇H is Lipschitz, which we shall always assume in this section, the Cauchy-Lipschitz theorem yields a unique solution to (4.2) passing through (q_0, p_0) at $t = 0$, which we denote by $t \mapsto \phi_t(q_0, p_0)$. We shall further assume that ϕ_t is well defined for all t , and call ϕ_t the *Hamiltonian flow*. By definition, the Hamiltonian flow preserves the Hamiltonian, since

$$\frac{d}{dt} H(\phi_t(q_0, p_0)) = \nabla H(\phi_t(q_0, p_0))^T J \nabla H(\phi_t(q_0, p_0)) = 0,$$

J being skew-symmetric. In other words, the flow ϕ_t follows level lines of H .

4.4.4 Ideal HMC and numerical HMC

The ideal HMC is the concatenation of two kernels. Given (q_n, p_n) , first resample p from its conditional under $\tilde{\pi}$; i.e. $p' \sim \mathcal{N}(0, M)$. This obviously leaves $\tilde{\pi}$ invariant, as in Gibbs sampling. Then set $(q_{n+1}, p_{n+1}) = \phi_T(q_n, p')$. In words, one step of the corresponding Markov chain consists of sampling a random momentum variable from the corresponding conditional, and then following the Hamiltonian flow ϕ_t up to time $T > 0$. Note that, unlike Gibbs sampling, this second step changes both variables.

One can formulate the intuition that, since the second step just follows a level line of H , the ideal HMC kernel leaves $\tilde{\pi}$ invariant. This is indeed the case (BoSa18), but since ϕ_t is usually not available in closed form, the ideal HMC kernel cannot be implemented, and we will skip the proof of its invariance.

In practice, one has to approximate the Hamiltonian flow ϕ_t , and there is a large literature in numerical analysis on the subject, with integrators showcasing many interesting properties. In terms of notation, denote by h a stepsize parameter, and $n = \lfloor T/h \rfloor$, so that we think of the numerical integrator $\psi_h^n(x, p)$ as an approximation to $\phi_T(x, p)$. There exists numerical integrators ϕ_h^n that are (i) C^1 diffeomorphisms, (ii) are reversible w.r.t. to momentum flip, and are volume-preserving, i.e. $|\det(\psi_h^n)'(q, p)| = 1$. Since the momentum flip preserves $\tilde{\pi}$, we can replace the second step of the ideal HMC algorithm by the abstract HMC algorithm of Section 4.4.1, with $\eta = \tilde{\pi}$, $\Phi = \phi_h^n$ and S the momentum flip. Because the numerical integrator is volume-preserving, the acceptance probability becomes suprisingly simple, as

$$\alpha_{HMC}((q, p), (q', p')) = 1 \wedge e^{-H(q, p) - H(q', p')}.$$

Intuitively, the acceptance step compensates for the fact that we did not exactly follow the level lines of H .

One common numerical integrator satisfying all the required properties is the leapfrog (aka velocity Verlet in (BoSa18)) integrator defined as $\psi_h^n = \psi_h \circ \dots \circ \psi_h$, where $(p', q') = \psi_h(p, q)$ is defined by

$$\begin{aligned} p_{1/2} &= p + \frac{h}{2} \nabla \log \pi(q) \\ q' &= q + hM^{-1} p_{1/2} \\ p' &= p_{1/2} + \frac{h}{2} \nabla \log \pi(q'); \end{aligned}$$

see (BoSa18) for more information, and <https://chi-feng.github.io/mcmc-demo/> for an interactive demo.

4.4.5 On the ergodicity of HMC

4.4.6 An ubiquitous variant: NUTS

NUTS for auto-tuning, etc.

4.5 MCMC practice: convergence diagnostics

Convergence diagnostics. Discuss the output of pymc.

5. Variational Bayes

While Monte Carlo methods are randomized numerical quadratures, *variational Bayes* (VB) stands for finding an approximation to the target distribution within some prespecified set \mathcal{Q} of distributions. This is done by minimizing some notion of distance between the target π and the so-called variational approximation $q \in \mathcal{Q}$. For instance, a common choice is to solve

$$q^* \in \arg \min_{q \in \mathcal{Q}} \text{KL}(q, \pi) := \mathbb{E}_q \log \frac{q(\theta)}{\pi(\theta)} = \int q(\theta) \log \frac{q(\theta)}{\pi(\theta)} d\theta, \quad (5.1)$$

where π and q are PDFs w.r.t. some common measure $d\theta$. Once (5.1) has been solved, one can either use the resulting q^* as a plug-in replacement for the target π , or, say, use $q^*(\theta)d\theta$ as a proposal distribution in importance sampling.

Compared to MCMC, the big plus of VB is that some modification of the optimization problem (5.1) can often be implemented in large-scale settings where either the number of data items or the dimension of the problem are large, e.g. in deep learning. The major downside of VB is the difficulty to provide theoretical guarantees on its results. One reason for this is that the distance to minimize, such as the reverse KL divergence in (5.1), is often chosen for computational convenience rather than for its guarantees on integrating functions of interest. Another reason is that approximating complex target distributions requires a large set \mathcal{Q} of variational approximations, which usually makes (5.1) a difficult optimization problem that has to be further modified before an efficient implementation is possible.

5.1 The evidence lower bound (ELBO)

Remember that the density π is often known only through the evaluation of an unnormalized version π_u , i.e., $\pi_u(\theta) = Z\pi(\theta), \forall \theta$. If we are to carry out (5.1), we thus need to know how to express $\text{KL}(q, \pi)$ using only π_u .

Lemma 1. Let $J(q) := \int q(\theta) \log \frac{q(\theta)}{\pi_u(\theta)} d\theta$. Then

$$J(q) = \text{KL}(q, \pi) - \log Z. \quad (5.2)$$

Proof. Left as an exercise. ■

Two remarks are in order. First, since Z does not depend on q , the optimization problem in (5.1) is equivalent to $\min J(q)$. Letting $L(q) = -J(q)$, (5.1) is further equivalent to $\max L(q)$. Second and the nonnegativity of the KL divergence implies that

$$L(q) \leq \log Z. \quad (5.3)$$

In Bayesian inference,

$$\pi_u(\theta) = p(y_{1:N}|\theta)p(\theta),$$

so that $Z = p(y_{1:N})$. Furthermore, (5.3) says that $L(q)$ is a lower bound for the (logarithm of the) evidence, shortened in ELBO. Most VB algorithms in the literature are cast as maximizing the ELBO.

5.2 Mean-field inference

The most common variational family is the so-called *mean-field approximation*. If you need to approximate a posterior over parameters $\theta \in \mathbb{R}^d$ and latent variables z_1, \dots, z_N , this means taking

$$\mathcal{Q} = \left\{ \theta \mapsto \prod_{d=1}^D q_d(\theta_d) \prod_{i=1}^N q_i(z_i) \right\}. \quad (5.4)$$

In other words, we approximate π with a separable PDF. Note that (5.4) only specifies the structure of the variational approximations. This is enough to derive the abstract form of the VB updates, which we do in the remainder of this section. In practice, though we further make explicit parametric choices for the individual factors, as we shall see in Section 5.2.1.

The whole motivation of the mean-field variational family is that if your target has simple conditionals, coordinate-wise optimization of the ELBO is easy. Indeed, write $x = (\theta, z) \in \mathbb{R}^p$ and let $1 \leq i \leq p$. Writing $q(x) = q_i(x_i)q_{\setminus i}(x_{\setminus i})$, and keeping track only of the additive terms that depend on q_i , it comes

$$\begin{aligned} L(q) &= \iint q_i(x_i)q_{\setminus i}(x_{\setminus i}) [\log \pi_u(x) - (\log q_i(x_i) + \log q_{\setminus i}(x_{\setminus i}))] dx_i dx_{\setminus i} \\ &\propto \int q_i(x_i) \left[\int q_{\setminus i}(x_{\setminus i}) \log \pi_u(x) dx_{\setminus i} \right] dx_i - \int q_i(x_i) \log q_i(x_i) dx_i \\ &= -KL(q_i, \phi_i), \end{aligned}$$

where

$$\phi_i(x_i) = \exp \left[\int q_{\setminus i}(x_{\setminus i}) \log \pi_u(x) dx_{\setminus i} \right] \quad (5.5)$$

is an unnormalized PDF. By a fundamental property of the KL, the ELBO L is thus maximized by setting $q_i \propto \phi_i$. The bottleneck is thus to be able to compute ϕ_i in (5.5). Like in deriving conditionals in Gibbs sampling, this is the part where conjugate distributions play a role. In practice, the choice of \mathcal{Q} is often made so that this step is easy, as we shall see in Section 5.2.1.

Finally, note that taking the variational family to be (5.4) is akin to assuming independence of all variables under the posterior. Combined with the fact that reverse KL (5.1) penalizes q^* putting a lot of mass where π does not, this often implies a gross underestimation of the support of the target (and thus of posterior uncertainty), along with the built-in ignorance of posterior correlations; see Figure ?? . While separability makes algorithmic derivations easier, we thus usually rather aim for “as separable as required by computation”. In other words, if, for modeling reasons, you believe

that there is correlation under π of some subset of the variables, say $x_i, i \in I$, you should try to keep these variables correlated in your variational approximation, by rather defining

$$\mathcal{Q} = \left\{ \theta \mapsto q_I(x_I) \prod_{i \notin I} q_i(x_i) \right\},$$

for some nonseparable q_I . **Mur12** calls this *structured mean-field*.

5.2.1 Mean-field VB for LDA

Recall the latent Dirichlet allocation model from Section ??, for which

$$\log p(y, z, \pi, B) \tag{5.6}$$

$$\begin{aligned} &= \sum_{i=1}^N \left[\log p(\pi_i | \alpha) + \sum_{\ell=1}^{L_i} \left(\log p(z_{i\ell} | \pi_i) + \log p(y_{i\ell} | z_{i\ell}, B) \right) \right] + p(B | \gamma) \\ &\propto \sum_{i=1}^N \left[\sum_{k=1}^K \alpha_k \log \pi_{ik} + \sum_{\ell=1}^{L_i} \left(\sum_{k=1}^K 1_{z_{i\ell}=k} \log \pi_{ik} + \sum_{v=1}^V \sum_{k=1}^K 1_{y_{i\ell}=v} 1_{z_{i\ell}=k} \log b_{kv} \right) \right] \\ &\quad + \sum_{k=1}^K \sum_{v=1}^V \gamma_v \log b_{kv}. \end{aligned} \tag{5.7}$$

We want to fit a mean-field approximation

$$\mathcal{Q} = \left\{ \prod_{i=1}^N \left[\text{Dir}(\pi_i | \tilde{\pi}_i) \prod_{\ell=1}^{L_i} \text{Cat}(z_{i\ell} | \tilde{z}_{i\ell}) \right] \prod_{k=1}^K \text{Dir}(B_k | \tilde{B}_k) \right\}.$$

Tilded variables parametrize the variational approximation q , and optimizing over q will thus be implemented as an optimization over these parameters. As we shall see, the Dirichlet distributions are chosen to make the following computations easy thanks to conjugacy.

To implement VB, we need to compute (5.5) for every coordinate, that is, we need to integrate the log joint distribution (5.7) with respect to all variables but one, for every choice of that singled out variable.

Singling out π_i .

We start by singling out $\pi_i \in \Delta_K$ for some $1 \leq i \leq N$, denoting the corresponding expectation by $\mathbb{E}_{\setminus \pi_i}$. We are confident that we shall be able to identify the functional form (and thus the normalization constant) of the resulting distribution, and thus we do not keep track of additive variable that do not imply π_i . This yields

$$\begin{aligned} \mathbb{E}_{\setminus \pi_i} \log p(y, z, \pi, B) &\propto \sum_{k=1}^K \alpha_k \log \pi_{ik} + \sum_{\ell=1}^{L_i} \mathbb{E}_{z_{i\ell}} \sum_{k=1}^K 1_{z_{i\ell}=k} \log \pi_{ik} \\ &= \sum_{k=1}^K \alpha_k \log \pi_{ik} + \sum_{\ell=1}^{L_i} \sum_{k=1}^K \tilde{z}_{i\ell k} \log \pi_{ik} \end{aligned}$$

and we recognize the log PDF of a Dirichlet distribution in π_i , with parameters

$$\tilde{\pi}_i \triangleq \left(\alpha_k + \sum_{\ell=1}^{L_i} \tilde{z}_{i\ell k} \right)_{1 \leq k \leq K}.$$

Singling out $z_{i\ell}$.

To compute $\mathbb{E}_{z_{i\ell}} \log p(y, z, \pi, B)$, we need to be able to compute expectation of log weights w.r.t. a Dirichlet distribution.

Lemma 2. Let $\Psi(\cdot) := \Gamma'(\cdot)/\Gamma(\cdot)$ be the digamma function. Let $\tilde{\eta} \in \Delta_M$ be a probability distribution over $\{1, \dots, M\}$. Then, for $m \in \{1, \dots, M\}$,

$$\mathbb{E}_{\text{Dir}(\eta|\tilde{\eta})} \log \eta_m = \Psi(\tilde{\eta}_m) - \Psi(\|\tilde{\eta}\|_1) \triangleq \Psi_m(\tilde{\eta}).$$

Proof. Left as an exercise. ■

Now we derive

$$\begin{aligned} \mathbb{E}_{z_{i\ell}} \log p(y, z, \pi, B) &\propto \sum_{k=1}^K 1_{z_{i\ell}=k} \mathbb{E}_{\pi_i} \log \pi_{ik} + \sum_{v=1}^V \sum_{k=1}^K 1_{y_{i\ell}=v} 1_{z_{i\ell}=k} \mathbb{E}_{B_{k:}} \log b_{kv} \\ &= \sum_{k=1}^K 1_{z_{i\ell}=k} \Psi_k(\tilde{\pi}_i) + \sum_{v=1}^V \sum_{k=1}^K 1_{y_{i\ell}=v} 1_{z_{i\ell}=k} \Psi_v(\tilde{B}_{k:}) \end{aligned}$$

We recognize a categorical distribution with parameters

$$\tilde{z}_{i\ell} \propto \left(\exp \left[\Psi_k(\tilde{\pi}_i) + \Psi_{y_{i\ell}}(\tilde{B}_{k:}) \right] \right)_{1 \leq k \leq K}.$$

Once again, the normalization constant can be guessed after doing the computation, since necessarily

$$\tilde{z}_{i\ell} = \left(\frac{\exp \left[\Psi_k(\tilde{\pi}_i) + \Psi_{y_{i\ell}}(\tilde{B}_{k:}) \right]}{\sum_{k=1}^K \exp \left[\Psi_k(\tilde{\pi}_i) + \Psi_{y_{i\ell}}(\tilde{B}_{k:}) \right]} \right)_{1 \leq k \leq K}.$$

Singling out $B_{k:}$.

In the same vein,

$$\mathbb{E}_{B_{k:}} \log p(y, z, \pi, B) \propto \sum_{i=1}^N \sum_{\ell=1}^{L_i} \sum_{v=1}^V 1_{y_{i\ell}=v} \mathbb{E}_{z_{i\ell}} 1_{z_{i\ell}=k} \log b_{kv} + \sum_{v=1}^V \gamma_v \log b_{kv}$$

and we recognize a Dirichlet with parameters

$$\tilde{B}_{k:} \triangleq \left(\gamma_v + \sum_{i=1}^N \sum_{\ell=1}^{L_i} 1_{y_{i\ell}=v} \tilde{z}_{i\ell k} \right)_{1 \leq v \leq V}.$$

This concludes the derivation of VB for LDA.

5.2.2 Mean-field VB for marginal LDA

As an exercise, derive the updates for the marginalized LDA model of Section ??; see **Mur12** for the solution.

5.2.3 VB generalizes the EM algorithm

TBC

5.3 Gradient-based algorithms

An alternate approach to finding a simple \mathcal{Q} leading to closed-form updates is to directly run a gradient algorithm on the ELBO (5.3). Take $\mathcal{Q} = \{q(\cdot|\phi), \phi \in \Phi\}$, where for all x , $\phi \mapsto q(x|\phi)$ is differentiable. Then, assuming the necessary regularity conditions, **PaBlJo12** note that gradient of the ELBO can be rewritten using the so-called *score function trick* as

$$\begin{aligned}\nabla_{\phi} L(q) &= \nabla_{\phi} \mathbb{E}_{x \sim q(\cdot|\phi)} \log \frac{\pi_u(x)}{q(x|\phi)} \\ &= \int \log \pi_u(x) \nabla_{\phi} q(x|\phi) dx + \nabla_{\phi} H[q(\cdot|\phi)]. \\ &= \mathbb{E}_{x \sim q(\cdot|\phi)} [\log \pi_u(x) \nabla_{\phi} \log q(x|\phi)] + \nabla_{\phi} H[q(\cdot|\phi)].\end{aligned}$$

The first term can be estimated by vanilla Monte Carlo. The entropy term can usually be differentiated in closed form; if not, it can be estimated by vanilla Monte Carlo as well. Overall, we can plug an unbiased estimator of the gradient of the ELBO in any stochastic gradient algorithm. More often than not, the ELBO as a function of ϕ is not convex, though, and one has to be happy with searching for local optima. Moreover, vanilla Monte Carlo estimators of (??) have been reported to have high variance even in simple models **PaBlJo12**.

Variance reduction for ELBO gradients has been a field of active research. **PaBlJo12** propose to use control variates, while **KiWe14** and a large body of follow-up work propose *reparametrization tricks* that work as follows. Assume that there exists a (deterministic) smooth and invertible function f such that $f(\varepsilon, \phi) \sim q(\cdot|\phi)$ whenever $\varepsilon \sim p(\varepsilon)$, with ε easy to sample. Now rewrite

$$\nabla_{\phi} L(q) = \nabla_{\phi} \mathbb{E}_{\varepsilon \sim p} \log \pi_u(f(\varepsilon, \phi)) + \nabla_{\phi} H[q(\cdot|\phi)].$$

This time the gradient can be passed under the integral in the first term, without relying on the score function trick, and we obtain

$$\nabla_{\phi} L(q) = \mathbb{E}_{\varepsilon \sim p} \nabla_{\phi} \log \pi_u(f(\varepsilon, \phi)) + \nabla_{\phi} H[q(\cdot|\phi)].$$

As long as we can compute gradients of $\log \pi_u$, we can compute the gradient in the expectation using the chain rule. This suggests a second vanilla Monte Carlo estimator, drawing $\varepsilon_i \sim p$ i.i.d. In practice, the resulting estimator has been found to have much lower variance (**ReMoWi14**), like in variational auto-encoders (**KiWe14**). I haven't seen a completely convincing explanation why and when variance reduction happens with the reparametrization trick in general, though. Finally, note that we again assumed that the entropy of q could be differentiated in closed form, but the entropy term can also be treated using the reparametrization trick if needed. We shall do so in the next section, for the sake of illustration.

5.3.1 VB for deep networks

One of the hot applications of gradient-based VB is for Bayesian deep learning, which has generated a huge literature in a short amount of time; see e.g. recent NeurIPS tutorials and workshops for pointers. For instance, **BCKW15** proceed as follows. We consider networks as generative models, so consider the softmax (classification) or squared (regression) loss. A network thus corresponds to a likelihood $p(\mathbf{y}|w)$. We take a prior $p(w)$ for the weights, and want to fit $q(w|\phi)$ to the posterior $\pi(w) \propto p(\mathbf{y}|w)p(w) = \pi_u(w)$. The gradient of the reparametrized ELBO writes

$$\begin{aligned}\nabla_{\phi} L(q(\cdot|\phi)) &= \mathbb{E}_{\varepsilon} \nabla_{\phi} \log \frac{\pi_u(f(\varepsilon, \phi))}{q(f(\varepsilon, \phi)|\phi)} \\ &\approx \frac{1}{N_{\varepsilon}} \sum_{i=1}^{N_{\varepsilon}} \nabla_{\phi} \log \frac{\pi_u(f(\varepsilon_i, \phi))}{q(f(\varepsilon_i, \phi)|\phi)}.\end{aligned}$$

Now notice that $\log \pi_u(w) = \log p(w) + \sum_{i=1}^{N_y} \log p(y_i|w)$, so that one can further uniformly draw (with or without replacement) a minibatch of data points B , and further obtain an unbiased estimator

$$\nabla_{\phi} L(q(\cdot|\phi)) \approx \frac{N_y}{N_{\epsilon}|B|} \sum_{i=1}^{N_{\epsilon}} \sum_{y \in B} \nabla_{\phi} \left[\frac{1}{N_y} \log p(f(\epsilon_i, \phi)) + \log p(y|f(\epsilon_i, \phi)) - \frac{1}{N_y} \log q(f(\epsilon_i, \phi)|\phi) \right].$$

Note that following **BCKW15**, we do not assume that the entropy can be differentiated in closed form, but the method applies *mutatis mutandis*. Note also that it is not obvious that replacing the entropy by its closed-form would reduce the variance of the estimator.

Now the key argument is that the gradient inside the sum can be computed using the chain rule, backpropagation, and the (assumed known) gradient of q . As an example, assume $w \in \mathbb{R}^d$, $\phi = (\mu, \sigma) \in \mathbb{R}^{d+1}$, so that $f(\phi, \epsilon) = \mu + \sigma \epsilon \sim \mathcal{N}(\mu, \sigma I_d)$. For $y \in B$, Let $F(w) = \log p(y|w) + \log p(w)$. The gradient of F is provided by backpropagation. Now the chain rule yields

$$\begin{aligned} \nabla_{\phi} (F(f(\epsilon, \cdot)))(\phi_0) &= J_{f(\epsilon, \cdot)}(\phi_0)^T \nabla F(f(\epsilon, \phi_0)) \\ &= \begin{pmatrix} I_d & \epsilon \end{pmatrix}^T \nabla F(f(\epsilon, \phi_0)). \end{aligned}$$

5.4 Theoretical guarantees

5.5 Alternatives to the KL divergence

IV

Nonparametric Bayes

6	Random functions: Gaussian processes	37
7	Random probability measures: Dirichlet processes and the like	39
8	Asymptotic frequentist properties	41



6. Random functions: Gaussian processes



7. Random probability measures: Dirichlet process

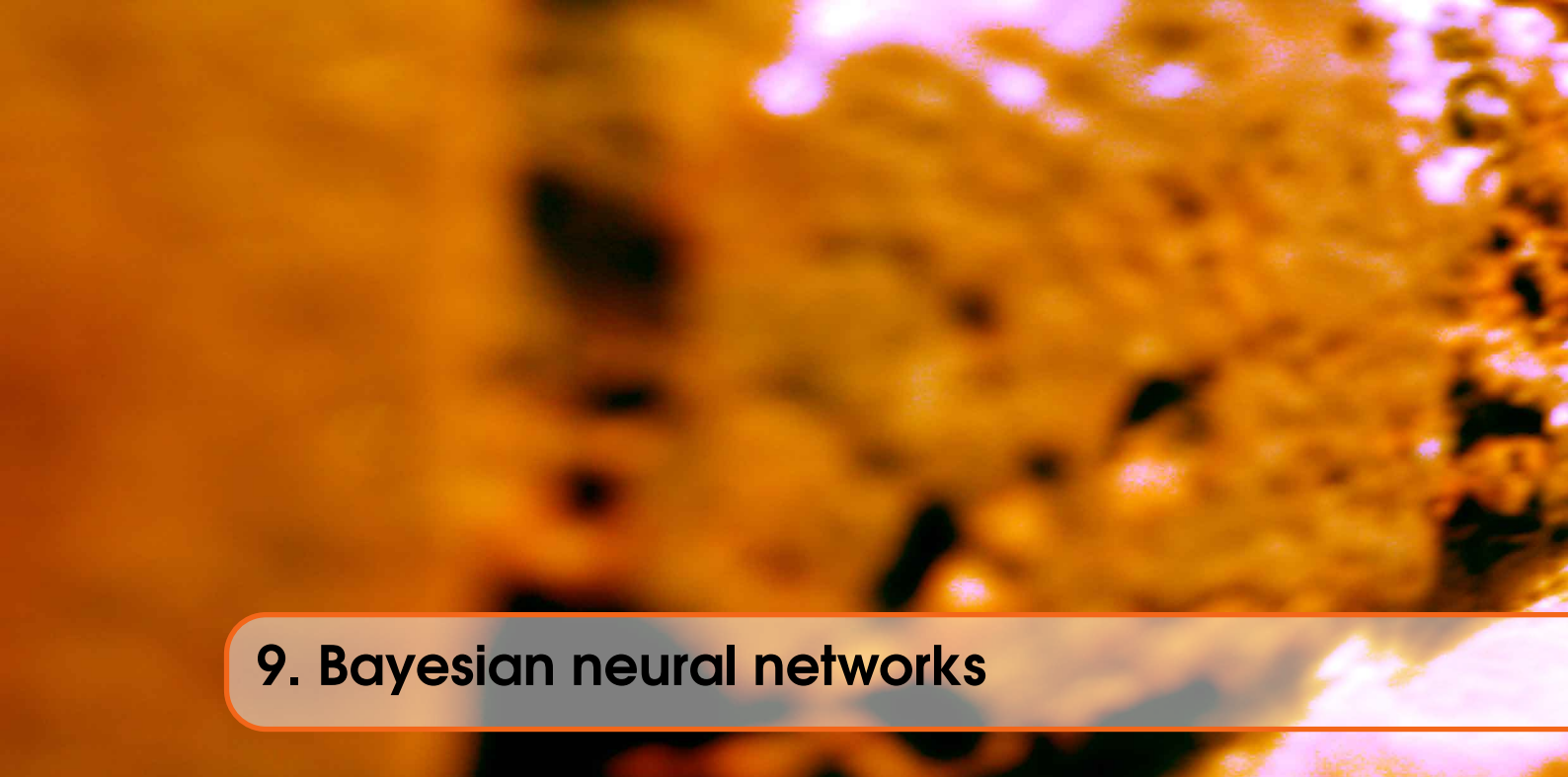


8. Asymptotic frequentist properties




Bayesian deep learning

9	Bayesian neural networks	45
10	Modern questions	47
	Bibliography	49



9. Bayesian neural networks



10. Modern questions



Bibliography

