
Double-Bayesian Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Contemporary machine learning methods will try to approach the Bayes error, as
2 it is the lowest possible error any model can achieve. This paper postulates that
3 any decision is composed of not one but two Bayesian decisions and that decision-
4 making is, therefore, a double-Bayesian process. The paper shows how this duality
5 implies intrinsic uncertainty in decisions and how it incorporates explainability.
6 The proposed approach understands that Bayesian learning is tantamount to finding
7 a base for a logarithmic function measuring uncertainty, with solutions being fixed
8 points. Furthermore, following this approach, the golden ratio describes possible
9 solutions satisfying Bayes' theorem. The double-Bayesian framework suggests
10 using a learning rate and momentum weight with values similar to those used in
11 the literature to train neural networks with stochastic gradient descent.

12 1 Introduction

13 Despite the progress in machine learning, several problems stand out for which convincing solutions
14 have yet to be found. With massive training sets, enormously sized networks, and immense computing
15 power, training machine learning models has become a brute force approach, arguably more concerned
16 with memorization than generalization. However, quoting from a post by Y. LeCun (Nov. 23, 2023),
17 we know that

18 *Animals and humans get very smart very quickly with vastly smaller amounts of training data than*
19 *current AI systems. Current large language models (LLMs) are trained on text data that would take*
20 *20,000 years for a human to read. And still, they haven't learned that if A is the same as B, then B*
21 *is the same as A. Humans get a lot smarter than that with comparatively little training data. Even*
22 *corvids, parrots, dogs, and octopuses get smarter than that very, very quickly, with only 2 billion*
23 *neurons and a few trillion "parameters."*

24 This raises the question of whether modern training techniques and principles are actually biologically
25 implemented in the human brain and, if not, what alternative methods could save resources. More
26 efficient methods would be better at generalizing with smaller amounts of training data, which almost
27 certainly would also improve the explainability and interpretability of neural networks.

28 This paper investigates what it takes for a classifier to be optimal. The starting point is Bayes' theorem,
29 which is the foundation of the Bayes classifier. The Bayes classifier is considered optimal because
30 it minimizes the Bayes risk, meaning it has the smallest probability of misclassification among all
31 classifiers. However, applying the Bayes classifier directly is often impossible because of the difficulty
32 in computing the posterior probabilities. For this reason, most classifiers are trying to approximate
33 the Bayes classifier, like the naïve Bayes classifier, for instance. The information-theoretical analysis
34 presented in this paper splits the decision of a Bayes classifier into two decisions, each following
35 Bayes' theorem, where one decision can serve as an explanation or verification of the other. Each of
36 the two decision processes faces intrinsic uncertainty, as its decision depends on the output of the
37 other process. The paper will investigate the theoretical ramifications of this approach. As a practical

38 result, it will discuss the consequences for two hyperparameters of stochastic gradient descent used
39 in the training process of a neural network: learning rate and momentum weight.

40 The structure of the paper is as follows: After this introduction, Section 2 motivates one of the main
41 ideas, namely that learning to make a decision involves solving two sub-problems and, thus, two
42 decisions. Section 3 discusses Bayes' theorem, which is central to statistical decision-making and
43 is the starting point of the theoretical approach outlined in the following. Section 4 then introduces
44 the double-Bayesian model as the key concept of the paper. The next section, Section 5, shows
45 how to represent possible solutions of the double-Bayesian decision model. Section 6 discusses the
46 golden ratio, including its functional equations and how it defines a solution to the double-Bayesian
47 model. Then, Section 7 discusses the theoretical implications for training double-Bayesian networks
48 with stochastic gradient descent. Finally, Section 8 summarizes the key concepts, followed by a
49 conclusion.

50 2 Dual decisions

Suppose a sender transmits the image on the left-hand side of Figure 1 to a receiver. This image

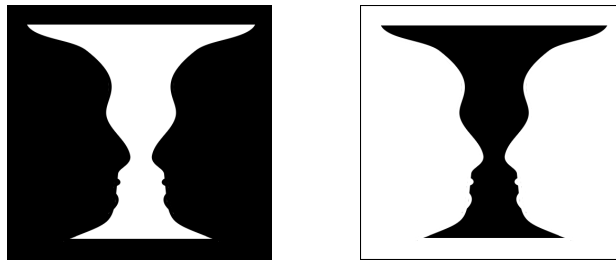


Figure 1: An image of Rubin's vase (left) and its inverted counterpart (right) - (Rubin, 1915)

51
52 depicts *Rubin's vase* by the Danish psychologist Edgar Rubin (Rubin, 1915), which shows a vase
53 or two faces looking at each other, depending on the receiver's perception. The receiver then faces
54 an unsolvable conundrum: 1) If the receiver thinks the image represents a vase, the receiver cannot
55 be certain that the vase is indeed the intended message the sender wanted to convey. Maybe the
56 sender wanted to send the faces. 2) If the receiver is expecting a picture of a vase (or faces) and
57 thus knows the intended message, there is no certainty that an image of a vase has been transmitted.
58 After all, the image could show faces. Therefore, two decisions are involved in making the final
59 interpretation of the image: 1) a decision about the perception of the image (vase or faces), and 2)
60 a decision about whether the perceived image coincides with the intended message, meaning the
61 image transmitted. Both decisions together are fraught with intrinsic uncertainty because deciding the
62 ultimate interpretation of Rubin's vase, a vase or faces, is impossible. Therefore, neither the sender
63 nor the receiver can make both decisions without uncertainty. Instead, the knowledge is distributed.
64 The sender knows the intended message (a vase or faces) but not the receiver's perceived image. On
65 the other hand, the receiver knows the perceived image (a vase or faces) but not the intended message.
66 Therefore, the sender and the receiver must collaborate to get the true interpretation across their
67 communication channel.

68 Let the sender and receiver perceive Rubin's vase differently, with contrary opinions about the
69 foreground and background color (black or white), where the foreground represents the perceived
70 image, either a vase or faces. Furthermore, let the sender and the receiver both be able to send
71 an image of Rubin's vase to each other so that both become senders and receivers alike and can
72 share their knowledge about the perceived image and intended message. The image that the sender
73 perceives is then the inverted image that the sender perceives. The goal is to collaborate so that the
74 perceived image (foreground) equals the intended message on both ends.

75 A sender can either send the image of Rubin's vase on the left-hand side of Figure 1 or send the
76 image with colors inverted, as shown on the right-hand side of Figure 1, depending on the perceived
77 image or intended message, respectively. On the other end, the receiver has two options: 1) accept
78 the received image if it is identical to the image expected, or 2) tell the sender to invert the image if it
79 is different. After this feedback, the image on the receiver end will be the same as the image on the

80 sender side. By making the images on both sides the same, the receiver has completed half of the
81 decision process without making a mistake and has thus behaved optimally. The receiver has ensured
82 that both sides see the same image. It is now up to the sender to make the final, second decision about
83 what image needs to be inverted to arrive at the final interpretation, either the image of the sender
84 or the image of the receiver. Thus, the first process tries to make the images identical, whereas the
85 second process tries to make the images different on both ends to reflect the different perceptions of
86 the sender and receiver.

87 Although described as a sequential process, the two dual decision processes leading to the final
88 interpretation are running in parallel. The sender is also a receiver, and the receiver is also a sender.
89 One of them conveys the correct foreground information (black or white), while the other conveys the
90 message. Note that neither the sender nor the receiver will ever see the true interpretation of the image.
91 The receiver in the example above will never know whether the received image needs to be inverted
92 after making the images identical because this would mean the receiver knows the true interpretation
93 of the image, which is not possible according to the uncertainty principle described above. A similar
94 statement can be made for the sender. The sender and the receiver can be considered dual and
95 complementary forces because of their different interpretations of foreground and background. They
96 make two binary decisions, deciding on the correct foreground color (black or white) and on the
97 message (a vase or faces). They decide whether Rubin's vase should be interpreted as a white vase, a
98 black vase, white faces, or black faces.

99 3 Bayes theorem

100 Bayes' theorem is a fundamental law in probability theory that describes the probability of an event
101 given prior knowledge. The theorem is of central importance in machine learning, where it guides the
102 training of machines for decision-making, such as in Bayesian inference or naïve Bayes classification.
103 For two events A and B , with prior probabilities $P(A)$ and $P(B)$, and $P(B) \neq 0$, Bayes' theorem
104 states the following:

$$P(A|B) = \frac{P(A) \cdot P(B|A)}{P(B)}, \quad (1)$$

105 where $P(A|B)$ and $P(B|A)$ are the conditional or posterior probabilities. Thus, $P(A|B)$ is the
106 probability of event A occurring when B is true, and analogously, $P(B|A)$ is the probability of B
107 given that A is true.

108 For a machine learning application, A would be the class of an observed input pattern B . The
109 probability $P(A)$ is then the prior probability of class A , and $P(B)$ is the prior probability of seeing
110 pattern B . Consequently, $P(A|B)$ is the posterior probability of class A when seeing pattern B , and
111 $P(B|A)$ is the posterior probability of B within A . According to Bayes' theorem, three probabilities
112 are needed to compute the probability $P(A|B)$ that class A is observed when seeing pattern B : $P(A)$,
113 $P(B)$, and $P(B|A)$. However, several obstacles prevent Bayes' theorem from being applied in this
114 way. No particular method can help determine the prior probabilities, which are often unknown.
115 Furthermore, the posterior probability is often not readily available and is approximated by making
116 assumptions about the distribution of B given A , for example, assuming a normal distribution.

117 To cope with these limitations, the next section describes decision-making as a dual process based on
118 Bayes' theorem, with uncertainty intrinsically involved.

119 4 Double-Bayesian framework

120 The Bayes Theorem is typically stated as in Eq. 1. However, restating the theorem in the following
121 equivalent form highlights the two decision processes for the two subproblems involved, as motivated
122 in Section 2:

$$\frac{P(A|B)}{P(B|A)} = \frac{P(A)}{P(B)} \quad (2)$$

123 The left-hand side of Eq. 2 features a fraction of the posterior probabilities, whereas the right-hand
124 side shows the prior probabilities. Following the motivation in Section 2, the posterior probabilities,
125 $P(A|B)$ and $P(B|A)$, can be understood as the probability that A or B is the intended message,
126 respectively. Then, the prior probabilities, $P(A)$ and $P(B)$, would express the probabilities that A or
127 B is in the foreground.

128 With only one equation for four parameters, Eq. 2 is underdetermined. However, it is fair to assume
 129 that $1 - (P(A|B) = P(B|A)$ and $1 - P(A) = P(B)$, which leaves one equation with one parameter
 130 on each side. This is possible because either A or B can be the message or foreground, not both of
 131 them at the same time, following again the reasoning in Section 2. Therefore, the intrinsic uncertainty
 132 in Bayes' theorem can be described as follows: if the true foreground is known, then whether the
 133 message needs to be swapped is unknown; on the other hand, if the message is known, then whether
 134 the foreground needs to be swapped is unknown. The fractions on both sides of Eq. 2 are thus
 135 "cognitively entangled."

136 The two remaining unknown parameters can be computed using two separate processes, each adding
 137 a constraint to handle the uncertainty. To illustrate this, Eq. 3 restates Bayes' theorem in yet another
 138 way:

$$1 = \frac{P(A)}{P(B)} \cdot \frac{P(B|A)}{P(A|B)} \quad (3)$$

139 Assuming that $P(B) = P(B|A)$, Eq. 3 simplifies to $P(A|B) = P(A)$. This assumption of B being
 140 independent of A is fair because, according to the motivation in Section 2, the decisions about the
 141 message and the foreground are independent of each other. Under this assumption, only one unknown
 142 remains, either $P(A|B)$ or $P(A)$, which follows directly from either $P(A)$ or $P(A|B)$, depending
 143 on which is input and which is output.

144 A similar, symmetric statement can be made when using the reciprocals on both sides of Eq. 2, which
 145 leads to the following equation:

$$1 = \frac{P(B)}{P(A)} \cdot \frac{P(A|B)}{P(B|A)} \quad (4)$$

146 Here, assuming that A is independent of B simplifies Eq. 4 to $P(B|A) = P(B)$.

147 Solving Eq. 2, Eq. 3, or Eq. 4 will be referred to as solving the outer Bayes equation. On the other
 148 hand, making both multiplicands on the right-hand side of Eq. 3 or Eq. 4 identical will be referred to
 149 as solving the inner Bayes equation, or simply solving the inner equation of Eq. 3 or Eq. 4. For Eq. 3,
 150 the inner Bayes equation thus states as follows:

$$\frac{P(A)}{P(B)} = \frac{P(B|A)}{P(A|B)} \quad (5)$$

151 Accordingly, the inner Bayes equation for Eq. 4 is obtained by using the reciprocals of the fractions
 152 on both sides of Eq. 5:

$$\frac{P(B)}{P(A)} = \frac{P(A|B)}{P(B|A)} \quad (6)$$

153 Consequently, the inner Bayes equations can be derived by inverting a fraction on one side of Bayes'
 154 theorem, as stated in Eq. 2. The inner Bayes equations are thus "entangled" versions of Bayes'
 155 theorem.

156 The two independent decision processes motivated above are solving the inner and outer Bayes equa-
 157 tions. To further formalize these processes, the following section will add a logarithmic expression
 158 to Eq. 3 and Eq. 4. Adding a logarithm offers several advantages: 1) using information theory to
 159 measure uncertainty; 2) using a reciprocal becomes equivalent to changing the sign of a logarithm;
 160 and 3) solving the equation in Bayes' theorem is reduced to finding a suitable base for a logarithm.

161 5 Fixpoint solutions

162 Using a logarithmic expression in Eq. 3 and Eq. 4 is possible when solutions become fixed points
 163 of a logarithmic function. To illustrate this, let $\log_b(x)$ be the logarithm for an input x and a base b .
 164 By definition, the logarithm is the inverse function of taking the power. Therefore, the following
 165 equation holds:

$$x = \log_b(b^x) \quad (7)$$

166 For the base b of a logarithm, any positive real number can be used so long as $b \neq 1$. A logarithm
 167 computed for base b can be converted into a logarithm for base b' as follows:

$$\log'_b(x) = \log_b(x) / \log_b(b') \quad (8)$$

168 Therefore, the simple term \log is used for the logarithm in the following.

169 By applying the logarithm to probabilities, they become information. For the two dual processes
 170 above, the information of one process will be its counterpart's information with a different sign. To
 171 achieve this, the following identity is required:

$$\log(x) = x \quad (9)$$

172 The following lemma states that this requirement can be met for general input values.

173 *Lemma:* For every $x \in \mathbb{R}^+ \setminus \{1\}$, there exists a base λ so that $\log_\lambda(x) = x$.

174 *Proof:* Let $b \in \mathbb{R}^+ \setminus \{1\}$ be an arbitrary basis for which $\log_b(x) = y$. Furthermore, let k be
 175 a multiplier so that $ky = x$. Then, $\log_\lambda(x) = x$ for $\lambda = b^{1/k}$. This follows from Eq.8, with
 176 $\log_\lambda(x) = \log_b(x)/\log_b(\lambda) = \log_b(x)/\log_b(b^{1/k}) = \log_b(x) \cdot k = x$. \square

177 Note that the common logarithmic rules apply for a fixed λ . However, when requiring a λ that always
 178 satisfies $\log_\lambda(x) = x$, computations become ambiguous, as seen here: $-\log_\lambda(x) = -x \neq 1/x =$
 179 $\log_\lambda(1/x)$. The base λ should be understood as a dynamic parameter that a learning system can
 180 modify over time so that $\log_\lambda(x)$ converges to the input x .

181 Using the \log_λ expression of the above Lemma, the Bayes' equation in Eq. 3 can be written as
 182 follows:

$$1 = \frac{P(A)}{P(B)} \cdot \log_\lambda \left(\frac{P(B|A)}{P(A|B)} \right) \quad (10)$$

183 Then, the following sequence of transformations can be derived from Eq. 10:

$$P(A|B) = \frac{P(A)}{P(B)} \cdot \log_\lambda \left(\frac{P(B|A)}{1} \right) \quad (11)$$

$$= \frac{1 - P(B)}{P(B)} \cdot \log_\lambda (P(B|A)) \quad (12)$$

$$= (1 - P(B)) \cdot \log_\lambda (P(B)^2) \quad (13)$$

$$= P(B) \cdot \log_\lambda (1 - P(B)^2) \quad (14)$$

$$= 2 \cdot P(B) \cdot \log_\lambda (\sqrt{1 - P(B)^2}) \quad (15)$$

$$= 2 \cdot \sin(\phi) \cdot \log_\lambda (\cos(\phi)), \quad (16)$$

184 where the last expression holds for an angle $\phi \in [0; \frac{\pi}{2}]$. The reasoning behind these transformations
 185 is as follows:

186 The first step, Eq. 11, moves the posterior probability $P(A|B)$ back to the left-hand side of the
 187 equation. The result is Bayes' theorem in its original form, as shown in Equation 1.

188 The next step, Eq. 12, replaces $P(A)$ with $1 - P(B)$, removing one degree of freedom as motivated
 189 above.

190 In the same way, Eq. 13 reformulates Eq. 12, assuming that $P(B) = P(B|A)$ and that the two
 191 multipliers on the right-hand side of the equation are equal to meet the inner Bayes equation.

192 Then, Eq. 14 rewrites the right-hand side of Eq. 13, transforming $1 - P(B) = P(B)^2$ into the
 193 equivalent $P(B) = 1 - P(B)^2$, which must hold true to satisfy the inner Bayes equation.

194 Finally, Eq. 15 extracts a factor of two from the \log_λ expression to get a radical input expression for
 195 the logarithm, following the standard rules for logarithms. The new input term to the \log_λ expression
 196 in Eq. 15 allows visualizing all possible solutions to the outer and inner Bayes equations.

197 To illustrate this further, Eq. 16 rewrites Eq. 15 using trigonometric functions and the Pythagorean
 198 relationship between \sin and \cos : $\sin^2 \phi + \cos^2 \phi = 1$, and thus $\sin \phi = \pm \sqrt{1 - \cos^2 \phi}$ and
 199 $\cos \phi = \pm \sqrt{1 - \sin^2 \phi}$. Solutions to the outer and inner Bayes equations then correspond to an
 200 angle ϕ in Equation 16, depending on the base λ . Thus, solutions are points on the unit circle.
 201 By changing the angle ϕ in Equation 16, all the possible solutions to the outer and inner Bayes

202 equations can be visualized. Following the reasoning above, the right-hand side of Eq. 16 represents
 203 the inner Bayes equation. Accordingly, after bringing the factor 2 on the other side of Eq. 16,
 204 the inner Bayes equation is satisfied when $\sin(\phi) = \cos(\phi)$, which is the case for $\phi = \pi/4$, with
 205 $\sin(\pi/4) = \cos(\pi/4) = 1/\sqrt{2}$.

206 For the dual process, the \log_λ expression can be used in combination with the other term of the inner
 207 Bayes equation in Eq. 3, as shown here:

$$1 = \log_\lambda \left(\frac{P(A)}{P(B)} \right) \cdot \frac{P(B|A)}{P(A|B)} \quad (17)$$

208 Note that the \log_λ expression has moved to the left compared to the right-hand side of Eq. 10. From
 209 this equation, the following sequence of transformations can be derived similar to the transformations
 210 above.

$$P(B) = \log_\lambda \left(\frac{P(A)}{1} \right) \cdot \frac{P(B|A)}{P(A|B)} \quad (18)$$

$$= \log_\lambda (P(A)) \cdot \frac{1 - P(A|B)}{P(A|B)} \quad (19)$$

$$= \log_\lambda (P(A|B)^2) \cdot (1 - P(A|B)) \quad (20)$$

$$= \log_\lambda (1 - P(A|B)^2) \cdot P(A|B) \quad (21)$$

$$= 2 \cdot \log_\lambda (\sqrt{1 - P(A|B)^2}) \cdot P(A|B) \quad (22)$$

$$= 2 \cdot \log_\lambda (\sin(\phi)) \cdot \cos(\phi) \quad (23)$$

211 During this sequence, assumptions similar to the ones in Eq. 12 and Eq. 13 are made. In Eq. 19,
 212 $P(B|A)$ was replaced by $1 - P(A|B)$, and Eq. 20 assumes that $P(A) = P(A|B)$. Again, all
 213 transformations assume that both multiplicands on the right-hand side are equal to satisfy the inner
 214 Bayes equation.

215 The intrinsic uncertainty for the dual processes can again be seen in Eq. 16 and Eq. 23, where it
 216 manifests like this: if the base λ is known, then the angle ϕ is unknown; and vice versa, if ϕ is known,
 217 then λ is unknown. Each process contributes knowledge about λ and ϕ , which the other process does
 218 not know.

219 The process knowledge about λ and ϕ does not need to be "all-or-nothing." The uncertainty ranges
 220 continuously between two extremes, and both dual processes can be somewhat knowledgeable about
 221 both parameters. When $\sin(\phi) = \cos(\phi)$, with $\phi = \pi/4$, one process has no or full knowledge
 222 about one parameter. With ϕ approaching 0 or $\pi/2$, where $\sin(\phi)$ and $\cos(\phi)$ become different, this
 223 knowledge increases or decreases, respectively.

224 6 Golden ratio

225 The solution to the inner Bayes equation is connected to the golden ratio (Livio, 2002), which becomes
 226 evident from the transformations of equations above and the assumptions made for both processes.
 227 Based on their right-hand equations, both dual processes must meet the same requirement to satisfy
 228 the inner Bayes equation, assuming that $\log_\lambda(x)$ produces x . For Eq. 12, with $P(B) = P(B|A)$,
 229 and for the corresponding Eq. 19 of the dual process, with $P(A) = P(A|B)$, this requirement can be
 230 written as

$$p = \frac{1 - p}{p}, \quad (24)$$

231 where the variable p is a placeholder for one of the probabilities. Eq. 24 holds true if p is the golden
 232 ratio, which is defined by the equivalent quadratic equation,

$$p^2 + p - 1 = 0, \quad (25)$$

233 which has two irrational solutions p_1 and p_2 :

$$p_1 = \frac{\sqrt{5} - 1}{2} \approx 0.618, \quad (26)$$

234 and

$$p_2 = \frac{-\sqrt{5} - 1}{2} \approx -1.618 \quad (27)$$

235 A key observation is that the complement of both solutions, $1 - p$, equals their square:

$$1 - p = p^2 \quad (28)$$

236 Alternatively, another quadratic equation that may be more frequently encountered in textbooks can
237 be used to arrive at the golden ratio. This equation is obtained by substituting $-p$ for p in Eq. 25:

$$p^2 - p - 1 = 0 \quad (29)$$

238 The alternative equation also possesses two irrational solutions, namely the negations of p_1 and p_2 :

$$-p_1 \approx -0.618 \quad \text{and} \quad -p_2 \approx 1.618 \quad (30)$$

239 For these solutions, the complement $1 - p$ is the negative reciprocal:

$$1 - p = -\frac{1}{p} \quad (31)$$

240 Computing the complement of the golden ratio allows changing viewpoints and switching between
241 the solutions to the inner and outer Bayes equations. This will become important in the next section
242 for training neural networks.

243 The golden ratio is sometimes represented by the letter φ in the literature. It is often defined as a
244 single value, usually $\varphi \approx 1.618$, and negative values are not considered (Livio, 2002; Huntley, 1970).
245 However, each of the four solutions to the aforementioned quadratic equations will be referred to as
246 the golden ratio in the context of this paper.

247 7 Theoretical implications

248 Supervised training methods first present a teaching input to a neural network and then try to make
249 the network's output the same as the input by adjusting the network weights. This equalizing of
250 input and output can be related to equalizing multiplicands to satisfy the inner Bayes equation. For
251 example, in Eq. 18, the term $P(B|A)/P(A|B)$ can be considered as input and the term $P(A)$ in
252 the lambda expression as output. The task of the lambda expression is then to make both terms the
253 same to satisfy the inner Bayes equation. Moreover, the lambda expression $\log_\lambda(P(A))$ becomes
254 the gradient of a linear function for the outer Bayes equation. These relationships help to determine
255 the optimal learning rate and momentum weight for training based on backpropagation and stochastic
256 gradient descent (SGD).

257 A training method based on backpropagation estimates the gradient of a loss function with respect to
258 each network weight, where the loss function measures the difference between input and network
259 output. Backpropagation methods try to minimize the loss by following the gradient and updating the
260 network weights accordingly (LeCun et al., 2012). They accomplish this for one network layer at
261 a time, iteratively propagating the gradient back from the output layer to the input layer. To move
262 along the gradient towards the minimum of the loss function, a delta is added to each weight, which
263 often has the following form, including a momentum term:

$$\Delta w_{ij}(t) = -\eta \frac{\partial L}{\partial w_{ij}(t)} + \alpha \cdot \Delta w_{ij}(t-1) \quad (32)$$

264 In (32), L is the loss function, and $\Delta w_{ij}(t)$ denotes the delta added to each weight w_{ij} between a
265 node i and a node j in the network at training iteration (or time) t . The term $\partial L / \partial w_{ij}(t)$ is the partial
266 derivative of the loss function with respect to w_{ij} , at time t , which is multiplied with the learning
267 rate η . The sign of $\Delta w_{ij}(t)$ is negative, so the loss function approaches its minimum. In practice, a
268 momentum term describing the weight change at time $t-1$, $\Delta w_{ij}(t-1)$, is commonly added. This
269 term is typically multiplied by a weighting factor α , as seen in (32).

270 The traditional understanding is that the momentum term improves stochastic gradient descent by
271 dampening oscillations. However, the dual process model offers another explanation for the per-
272 formance improvement brought about by the momentum term. As of yet, a conclusive theory for
273 the optimal values of the learning rate η and the momentum weight α has been lacking. Although

274 second-order methods (Bengio, 2012; Sutskever et al., 2013; Spall, 2000) as well as adaptive meth-
 275 ods (Jacobs, 1988; Kingma and Ba, 2014; Duchi et al., 2011; Tieleman and Hinton, 2012) have been
 276 tried with various degrees of success, an ultimate answer has still to be found. Both parameters are
 277 usually determined heuristically through empirical experiments or systematic search (Bergstra and
 278 Bengio, 2012). Training results can be very sensitive to the value of the learning rate. For example, a
 279 small learning rate may result in slow convergence, whereas a larger learning rate may result in the
 280 search passing over the minimum loss. Negotiating this delicate trade-off in the regularization of the
 281 training process can be time-consuming in practical applications. The literature seems to prefer initial
 282 learning rates around 0.01 or smaller for SGD, although reported values differ by several orders of
 283 magnitude. For the momentum weight, higher initial values around 0.9 are more common (Li et al.,
 284 2020; Krizhevsky et al., 2012; Simonyan and Zisserman, 2014; He et al., 2016).

285 As shown in the following, the proposed dual process model allows deriving theoretical values for
 286 both regularization parameters: learning rate η and momentum weight α . In the weight adjustment
 287 given by Eq. 32, each summand represents a gradient of one of the two dual processes. These are
 288 the partial derivative $\partial L/\partial w_{ij}(t)$ and the momentum term $\Delta w_{ij}(t-1)$. The momentum weight α
 289 follows from the results above, where the lambda expression can be considered as the gradient of
 290 the current iteration at time t . The other multiplicand of the inner Bayes equation corresponds to the
 291 gradient of the other dual process at time $t-1$, assuming that both dual processes are interleaved, if
 292 not in parallel.

293 The previous sections showed that the inner Bayes equation is met when both summands are equal to
 294 $\sin(\pi/4) = \cos(\pi/4) = 1/\sqrt{2}$ and when they are equal to the golden ratio. Therefore, the delta at
 295 $t-1$, $\Delta w_{ij}(t-1)$, needs to be multiplied by a constant to obtain the golden ratio. This constant is
 296 the momentum weight α , which needs to satisfy $\alpha/\sqrt{2} = p_1$, and can thus be computed as follows.

$$\alpha = \sqrt{2} \cdot p_1 \approx 0.874, \quad (33)$$

297 where p_1 is the value of the golden ratio in Eq. 26. So, this logic provides the value of the first
 298 regularization term, namely the momentum weight α , with $\alpha \approx 0.874$.

299 The learning rate η can be derived from the momentum weight α by converting the latter to the
 300 corresponding value for the dual process. The dual process does not aim to satisfy the inner Bayes
 301 equation with $\phi = \pi/4$. Instead, it aims to satisfy the outer Bayes equation, with $\phi = 0$ or $\phi = \pi/2$,
 302 and thus $\sin(\phi) = 0$ and $\cos(\phi) = 1$, or $\sin(\phi) = 1$ and $\cos(\phi) = 0$. By moving in the opposite
 303 direction of the gradient of its dual counterpart, the first process can minimize its loss in satisfying
 304 the inner Bayes equation. Accordingly, taking the complement of the momentum weight α twice
 305 results in the learning rate η for the gradient change at time t . Taking the complement of α twice can
 306 be understood as looking at the same process from a dual point of view. Mathematically, this can be
 307 achieved by squaring the simple complement, $1 - \alpha$. Squaring the complement follows the functional
 308 equation of the golden ratio described by Eq.28. Squaring also means bringing the multiplier 2 back
 309 in, which was extracted from the lambda expression in Eq. 15 and Eq. 22 to represent all solutions
 310 graphically. Applying these steps to the momentum weight α then results in the following equation
 311 for the learning rate η :

$$\eta = (1 - \alpha)^2 \approx 0.016 \quad (34)$$

312 So, this computation provides the value for the second regularization term, learning rate η , with $\eta \approx$
 313 0.016.

314 8 Discussion

315 Starting from Bayes' theorem, this paper develops a theoretical framework that describes any decision
 316 of a machine classifier as the result of two processes. The first decision process determines the input
 317 message; specifically, it decides whether the input is encoded according to its true value or needs to
 318 be inverted. On the other hand, the second decision process decides whether the output should be
 319 equal to the input or needs to be inverted. Although both decision processes run simultaneously, they
 320 are independent processes, with each possessing knowledge not accessible to the other process. What
 321 is uncertain for one process is certain for the other, and vice versa. The first process does not know
 322 whether the input should be equal to the output, and conversely, the second process does not know
 323 whether the input needs to be inverted. This means a binary decision always involves two bits, one
 324 indicating the encoding of the input and the other defining the relationship between input and output.

325 However, practically, only one of the two processes can be performed at a time, leaving one bit of
326 uncertainty for one of the processes.

327 Theoretically, the framework proposed here formulates this duality with two processes having
328 different perceptions of zero and one (black and white). The output of one process is the input to
329 the other process. While one process tries to make its output equal to its input, the other aims for
330 the opposite and tries to make its output as different as possible. The mathematical definitions of
331 these processes are defined by the outer and inner Bayes equation, the latter of which is an entangled
332 version of the original Bayes' theorem. By introducing the logarithm, each process is given a control
333 parameter, namely the base of the logarithm, to achieve its goal. This parameter, which is essentially
334 a multiplier, allows each process to control the magnitude of the input/output.

335 The solution space of the proposed double-Bayesian decision framework can be visualized with the
336 trigonometric functions \sin and \cos . Furthermore, the golden ratio defines solutions to the inner
337 Bayes equation. Connecting these two observations leads to specific values for momentum weight
338 and learning rate for stochastic gradient descent, which tries to minimize the difference between
339 training input and output during training.

340 The supplemental material to this paper contains experiments for the MNIST dataset (LeCun et al.,
341 accessed May 21, 2024), where the proposed double-Bayesian learning framework is practically
342 evaluated. The theoretical parameters found in this paper did, in fact, provide the best performance
343 for a network trained with stochastic gradient descent in a large grid search for learning rate and
344 momentum weight.

345 9 Conclusion

346 Three primary characteristics define the work presented in this paper: First, a double-Bayesian
347 approach that understands learning as a process involving two Bayesian decisions instead of a single
348 decision, like in contemporary approaches. Second, solving a Bayesian decision problem is equivalent
349 to finding a fixed point for a logarithmic function measuring uncertainty. Third, the golden ratio
350 defines solutions to a Bayesian decision problem. These three characteristics make the proposed
351 approach novel and unique.

352 The double-Bayesian framework leads to new theoretical results for training neural networks, particu-
353 larly specific hyperparameter values for backpropagation and gradient descent. These results are in
354 contrast with other gradient descent heuristics in the literature that either use dynamic hyperparam-
355 eters or second-order methods for adjusting parameters during training. It will be interesting to see how
356 this conceptual difference will be resolved in the future. The proposed framework offers new ways to
357 understand how neural networks make decisions and may thus contribute to the interpretability and
358 explainability of neural networks, an actively investigated research area.

359 The proposed framework may also help build bridges to other disciplines like neuroscience or
360 physics. For example, representing all possible solutions to a double-Bayesian decision by means
361 of trigonometric functions, as done in this paper, introduces waves. Incorporating brain waves into
362 machine learning, a feature that traditional machine learning approaches are arguably lacking, would
363 likely entail a better understanding of learning in general. This better understanding could mean
364 training methods for smaller networks that could achieve the same performance with less training
365 data, as motivated at the beginning of this paper.

366 Another example of a discipline that could be related to this work is quantum mechanics. One of the
367 fundamental concepts in quantum mechanics is Heisenberg's uncertainty principle, which states that
368 certain pairs of physical properties, such as the position and momentum of an electron, cannot be
369 measured with absolute certainty. The more accurately one property is measured, the less is known
370 about the other property. The proposed double-Bayesian framework incorporates such an intrinsic
371 uncertainty and makes a connection to Bayesian decision theory, which could lead to new insights.

372 Although empirical evidence in the literature supports the theoretical hyperparameter values derived
373 in this paper, and the experiments in the supplemental material show that these values outperform
374 other value pairs, more practical experiments are needed to corroborate these values. To address this
375 limitation, future work will validate the practicality of the derived hyperparameter values in additional
376 experiments across different domains and compare their performance with the performance of other
377 values and other optimization strategies.

378 References

- 379 Y. Bengio. Practical recommendations for gradient-based training of deep architectures. In *Neural*
380 *networks: Tricks of the trade*, pages 437–478. Springer, 2012.
- 381 J. Bergstra and Y. Bengio. Random search for hyper-parameter optimization. *Journal of machine*
382 *learning research*, 13(2), 2012.
- 383 L. Buitinck, G. Louppe, M. Blondel, F. Pedregosa, A. Mueller, O. Grisel, V. Niculae, P. Prettenhofer,
384 A. Gramfort, J. Grobler, R. Layton, J. VanderPlas, A. Joly, B. Holt, and G. Varoquaux. API
385 design for machine learning software: experiences from the scikit-learn project. In *ECML PKDD*
386 *Workshop: Languages for Data Mining and Machine Learning*, pages 108–122, 2013.
- 387 J. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and stochastic
388 optimization. *Journal of machine learning research*, 12(7), 2011.
- 389 K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings*
390 *of the IEEE Conference on computer vision and pattern recognition*, pages 770–778, 2016.
- 391 H. Huntley. *The Divine Proportion*. Dover Publications, 1970.
- 392 R. Jacobs. Increased rates of convergence through learning rate adaptation. *Neural networks*, 1(4):
393 295–307, 1988.
- 394 D. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*,
395 2014.
- 396 A. Krizhevsky, I. Sutskever, and G. Hinton. Imagenet classification with deep convolutional neural
397 networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- 398 Y. LeCun, L. Bottou, G. Orr, and K. Müller. Efficient backprop. In *Neural networks: Tricks of the*
399 *trade*, pages 9–48. Springer, 2012.
- 400 Y. LeCun, C. Cortes, and C. Burges. *The MNIST Database*, accessed May 21, 2024. URL <http://yann.lecun.com/exdb/mnist/>.
401
- 402 H. Li, P. Chaudhari, H. Yang, M. Lam, A. Ravichandran, R. Bhotika, and S. Soatto. Rethinking the
403 hyperparameters for fine-tuning. *arXiv preprint arXiv:2002.11770*, 2020.
- 404 M. Livio. *The Golden Ratio*. Random House, Inc., 2002.
- 405 F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Pretten-
406 hofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and
407 E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*,
408 12:2825–2830, 2011.
- 409 E. Rubin. *Rubin Vase*. Wikimedia Commons (last accessed March 4, 2022, CC BY-SA 3.0), 1915.
410 URL <https://commons.wikimedia.org/wiki/File:Facevase.png>.
- 411 Scikit-learn developers (BSD License). *Scikit-learn machine learning library*, accessed May 21,
412 2024. URL [https://scikit-learn.org/stable/modules/generated/sklearn.model_](https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.StratifiedShuffleSplit.html)
413 [selection.StratifiedShuffleSplit.html](https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.StratifiedShuffleSplit.html).
- 414 K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition.
415 *arXiv preprint arXiv:1409.1556*, 2014.
- 416 J. Spall. Adaptive stochastic approximation by the simultaneous perturbation method. *IEEE transac-*
417 *tions on automatic control*, 45(10):1839–1853, 2000.
- 418 I. Sutskever, J. Martens, G. Dahl, and G. Hinton. On the importance of initialization and momentum
419 in deep learning. In *International Conference on Machine Learning*, pages 1139–1147, 2013.
- 420 T. Tieleman and G. Hinton. Lecture 6.5-rmsprop: Divide the gradient by a running average of its
421 recent magnitude. *COURSERA: Neural networks for machine learning*, 4(2):26–31, 2012.

422 NeurIPS Paper Checklist

423 1. Claims

424 Question: Do the main claims made in the abstract and introduction accurately reflect the
425 paper's contributions and scope?

426 Answer: [Yes]

427 Justification: **This is a theoretical paper that tries to explain hyperparameter values that**
428 **have been successfully used in the literature. The paper investigates what it takes for a**
429 **classifier to be optimal, as stated in the introduction. Although the literature and the**
430 **practical experiments provided in the supplemental material support the theoretical**
431 **results, providing more practical experiments would be desirable.**

432 Guidelines:

- 433 • The answer NA means that the abstract and introduction do not include the claims
434 made in the paper.
- 435 • The abstract and/or introduction should clearly state the claims made, including the
436 contributions made in the paper and important assumptions and limitations. A No or
437 NA answer to this question will not be perceived well by the reviewers.
- 438 • The claims made should match theoretical and experimental results, and reflect how
439 much the results can be expected to generalize to other settings.
- 440 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
441 are not attained by the paper.

442 2. Limitations

443 Question: Does the paper discuss the limitations of the work performed by the authors?

444 Answer: [Yes]

445 Justification: **The limitations are discussed at the very end of the paper in the con-**
446 **clusion. A comparison with other hyperparameter optimization strategies would be**
447 **desirable to corroborate the theoretical results even more. Specifically, a systematic**
448 **comparison with second-order methods and other methods that dynamically adapt**
449 **hyperparameters during training should shed more light on the performance of this**
450 **approach.**

451 Guidelines:

- 452 • The answer NA means that the paper has no limitation while the answer No means that
453 the paper has limitations, but those are not discussed in the paper.
- 454 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 455 • The paper should point out any strong assumptions and how robust the results are to
456 violations of these assumptions (e.g., independence assumptions, noiseless settings,
457 model well-specification, asymptotic approximations only holding locally). The authors
458 should reflect on how these assumptions might be violated in practice and what the
459 implications would be.
- 460 • The authors should reflect on the scope of the claims made, e.g., if the approach was
461 only tested on a few datasets or with a few runs. In general, empirical results often
462 depend on implicit assumptions, which should be articulated.
- 463 • The authors should reflect on the factors that influence the performance of the approach.
464 For example, a facial recognition algorithm may perform poorly when image resolution
465 is low or images are taken in low lighting. Or a speech-to-text system might not be
466 used reliably to provide closed captions for online lectures because it fails to handle
467 technical jargon.
- 468 • The authors should discuss the computational efficiency of the proposed algorithms
469 and how they scale with dataset size.
- 470 • If applicable, the authors should discuss possible limitations of their approach to
471 address problems of privacy and fairness.
- 472 • While the authors might fear that complete honesty about limitations might be used by
473 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
474 limitations that aren't acknowledged in the paper. The authors should use their best

475 judgment and recognize that individual actions in favor of transparency play an impor-
476 tant role in developing norms that preserve the integrity of the community. Reviewers
477 will be specifically instructed to not penalize honesty concerning limitations.

478 3. Theory Assumptions and Proofs

479 Question: For each theoretical result, does the paper provide the full set of assumptions and
480 a complete (and correct) proof?

481 Answer: [Yes]

482 Justification: **All assumptions are discussed in detail, and one proof has been included.**

483 Guidelines:

- 484 • The answer NA means that the paper does not include theoretical results.
- 485 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
486 referenced.
- 487 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 488 • The proofs can either appear in the main paper or the supplemental material, but if
489 they appear in the supplemental material, the authors are encouraged to provide a short
490 proof sketch to provide intuition.
- 491 • Inversely, any informal proof provided in the core of the paper should be complemented
492 by formal proofs provided in appendix or supplemental material.
- 493 • Theorems and Lemmas that the proof relies upon should be properly referenced.

494 4. Experimental Result Reproducibility

495 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
496 perimental results of the paper to the extent that it affects the main claims and/or conclusions
497 of the paper (regardless of whether the code and data are provided or not)?

498 Answer: [Yes]

499 Justification: **Experimental results are listed in the supplemental material, with infor-
500 mation to reproduce the results, including the code itself.**

501 Guidelines:

- 502 • The answer NA means that the paper does not include experiments.
- 503 • If the paper includes experiments, a No answer to this question will not be perceived
504 well by the reviewers: Making the paper reproducible is important, regardless of
505 whether the code and data are provided or not.
- 506 • If the contribution is a dataset and/or model, the authors should describe the steps taken
507 to make their results reproducible or verifiable.
- 508 • Depending on the contribution, reproducibility can be accomplished in various ways.
509 For example, if the contribution is a novel architecture, describing the architecture fully
510 might suffice, or if the contribution is a specific model and empirical evaluation, it may
511 be necessary to either make it possible for others to replicate the model with the same
512 dataset, or provide access to the model. In general, releasing code and data is often
513 one good way to accomplish this, but reproducibility can also be provided via detailed
514 instructions for how to replicate the results, access to a hosted model (e.g., in the case
515 of a large language model), releasing of a model checkpoint, or other means that are
516 appropriate to the research performed.
- 517 • While NeurIPS does not require releasing code, the conference does require all submis-
518 sions to provide some reasonable avenue for reproducibility, which may depend on the
519 nature of the contribution. For example
 - 520 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
521 to reproduce that algorithm.
 - 522 (b) If the contribution is primarily a new model architecture, the paper should describe
523 the architecture clearly and fully.
 - 524 (c) If the contribution is a new model (e.g., a large language model), then there should
525 either be a way to access this model for reproducing the results or a way to reproduce
526 the model (e.g., with an open-source dataset or instructions for how to construct
527 the dataset).

528 (d) We recognize that reproducibility may be tricky in some cases, in which case
529 authors are welcome to describe the particular way they provide for reproducibility.
530 In the case of closed-source models, it may be that access to the model is limited in
531 some way (e.g., to registered users), but it should be possible for other researchers
532 to have some path to reproducing or verifying the results.

533 5. Open access to data and code

534 Question: Does the paper provide open access to the data and code, with sufficient instruc-
535 tions to faithfully reproduce the main experimental results, as described in supplemental
536 material?

537 Answer: [Yes]

538 Justification: **Please see the supplemental material for the code and information about**
539 **reproducing the experimental results. The publicly available MNIST database has**
540 **been used for the experiments.**

541 Guidelines:

- 542 • The answer NA means that paper does not include experiments requiring code.
- 543 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
544 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 545 • While we encourage the release of code and data, we understand that this might not be
546 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
547 including code, unless this is central to the contribution (e.g., for a new open-source
548 benchmark).
- 549 • The instructions should contain the exact command and environment needed to run to
550 reproduce the results. See the NeurIPS code and data submission guidelines ([https:
551 //nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 552 • The authors should provide instructions on data access and preparation, including how
553 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 554 • The authors should provide scripts to reproduce all experimental results for the new
555 proposed method and baselines. If only a subset of experiments are reproducible, they
556 should state which ones are omitted from the script and why.
- 557 • At submission time, to preserve anonymity, the authors should release anonymized
558 versions (if applicable).
- 559 • Providing as much information as possible in supplemental material (appended to the
560 paper) is recommended, but including URLs to data and code is permitted.

561 6. Experimental Setting/Details

562 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
563 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
564 results?

565 Answer: [Yes]

566 Justification: **Please see the information in the supplemental material.**

567 Guidelines:

- 568 • The answer NA means that the paper does not include experiments.
- 569 • The experimental setting should be presented in the core of the paper to a level of detail
570 that is necessary to appreciate the results and make sense of them.
- 571 • The full details can be provided either with the code, in appendix, or as supplemental
572 material.

573 7. Experiment Statistical Significance

574 Question: Does the paper report error bars suitably and correctly defined or other appropriate
575 information about the statistical significance of the experiments?

576 Answer: [NA]

577 Justification: **The paper provides theoretical results. For the experimental results in**
578 **the supplemental material, only the relative performance to other hyperparameter**
579 **combinations was investigated, significant or not, to see whether the proposed values**

580 **define the optimum or are at least close to it. To compare the proposed method and**
581 **values with other optimization methods, future experiments may require significance**
582 **tests.**

583 Guidelines:

- 584 • The answer NA means that the paper does not include experiments.
- 585 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
586 dence intervals, or statistical significance tests, at least for the experiments that support
587 the main claims of the paper.
- 588 • The factors of variability that the error bars are capturing should be clearly stated (for
589 example, train/test split, initialization, random drawing of some parameter, or overall
590 run with given experimental conditions).
- 591 • The method for calculating the error bars should be explained (closed form formula,
592 call to a library function, bootstrap, etc.)
- 593 • The assumptions made should be given (e.g., Normally distributed errors).
- 594 • It should be clear whether the error bar is the standard deviation or the standard error
595 of the mean.
- 596 • It is OK to report 1-sigma error bars, but one should state it. The authors should
597 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
598 of Normality of errors is not verified.
- 599 • For asymmetric distributions, the authors should be careful not to show in tables or
600 figures symmetric error bars that would yield results that are out of range (e.g. negative
601 error rates).
- 602 • If error bars are reported in tables or plots, The authors should explain in the text how
603 they were calculated and reference the corresponding figures or tables in the text.

604 8. Experiments Compute Resources

605 Question: For each experiment, does the paper provide sufficient information on the com-
606 puter resources (type of compute workers, memory, time of execution) needed to reproduce
607 the experiments?

608 Answer: [Yes]

609 Justification: **Please see the supplemental material for more information.**

610 Guidelines:

- 611 • The answer NA means that the paper does not include experiments.
- 612 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
613 or cloud provider, including relevant memory and storage.
- 614 • The paper should provide the amount of compute required for each of the individual
615 experimental runs as well as estimate the total compute.
- 616 • The paper should disclose whether the full research project required more compute
617 than the experiments reported in the paper (e.g., preliminary or failed experiments that
618 didn't make it into the paper).

619 9. Code Of Ethics

620 Question: Does the research conducted in the paper conform, in every respect, with the
621 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

622 Answer: [Yes]

623 Justification: **There is no violation of the NeurIPS Code of Ethics.**

624 Guidelines:

- 625 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 626 • If the authors answer No, they should explain the special circumstances that require a
627 deviation from the Code of Ethics.
- 628 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
629 eration due to laws or regulations in their jurisdiction).

630 10. Broader Impacts

631 Question: Does the paper discuss both potential positive societal impacts and negative
632 societal impacts of the work performed?

633 Answer: [NA]

634 Justification: **The paper proposes a generic method to find hyperparameter values**
635 **for optimizing the performance of neural networks. Its societal impacts, therefore,**
636 **correlate with the risks of machine learning in general, which does not need to be**
637 **pointed out in particular according to the guidelines below.**

638 Guidelines:

- 639 • The answer NA means that there is no societal impact of the work performed.
- 640 • If the authors answer NA or No, they should explain why their work has no societal
641 impact or why the paper does not address societal impact.
- 642 • Examples of negative societal impacts include potential malicious or unintended uses
643 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
644 (e.g., deployment of technologies that could make decisions that unfairly impact specific
645 groups), privacy considerations, and security considerations.
- 646 • The conference expects that many papers will be foundational research and not tied
647 to particular applications, let alone deployments. However, if there is a direct path to
648 any negative applications, the authors should point it out. For example, it is legitimate
649 to point out that an improvement in the quality of generative models could be used to
650 generate deepfakes for disinformation. On the other hand, it is not needed to point out
651 that a generic algorithm for optimizing neural networks could enable people to train
652 models that generate Deepfakes faster.
- 653 • The authors should consider possible harms that could arise when the technology is
654 being used as intended and functioning correctly, harms that could arise when the
655 technology is being used as intended but gives incorrect results, and harms following
656 from (intentional or unintentional) misuse of the technology.
- 657 • If there are negative societal impacts, the authors could also discuss possible mitigation
658 strategies (e.g., gated release of models, providing defenses in addition to attacks,
659 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
660 feedback over time, improving the efficiency and accessibility of ML).

661 11. Safeguards

662 Question: Does the paper describe safeguards that have been put in place for responsible
663 release of data or models that have a high risk for misuse (e.g., pretrained language models,
664 image generators, or scraped datasets)?

665 Answer: [NA]

666 Justification: **There is no risk of misusing the proposed method beyond misusing**
667 **machine learning in general.**

668 Guidelines:

- 669 • The answer NA means that the paper poses no such risks.
- 670 • Released models that have a high risk for misuse or dual-use should be released with
671 necessary safeguards to allow for controlled use of the model, for example by requiring
672 that users adhere to usage guidelines or restrictions to access the model or implementing
673 safety filters.
- 674 • Datasets that have been scraped from the Internet could pose safety risks. The authors
675 should describe how they avoided releasing unsafe images.
- 676 • We recognize that providing effective safeguards is challenging, and many papers do
677 not require this, but we encourage authors to take this into account and make a best
678 faith effort.

679 12. Licenses for existing assets

680 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
681 the paper, properly credited and are the license and terms of use explicitly mentioned and
682 properly respected?

683 Answer: [Yes]

684 Justification: **The main paper cites relevant references for the scientific content and the**
685 **supplemental material provides more details about the data and software sources.**

686 Guidelines:

- 687 • The answer NA means that the paper does not use existing assets.
- 688 • The authors should cite the original paper that produced the code package or dataset.
- 689 • The authors should state which version of the asset is used and, if possible, include a
690 URL.
- 691 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 692 • For scraped data from a particular source (e.g., website), the copyright and terms of
693 service of that source should be provided.
- 694 • If assets are released, the license, copyright information, and terms of use in the
695 package should be provided. For popular datasets, `paperswithcode.com/datasets`
696 has curated licenses for some datasets. Their licensing guide can help determine the
697 license of a dataset.
- 698 • For existing datasets that are re-packaged, both the original license and the license of
699 the derived asset (if it has changed) should be provided.
- 700 • If this information is not available online, the authors are encouraged to reach out to
701 the asset’s creators.

702 13. New Assets

703 Question: Are new assets introduced in the paper well documented and is the documentation
704 provided alongside the assets?

705 Answer: [Yes]

706 Justification: **The paper provides new assets in the form of knowledge about hyperpa-**
707 **rameter values to train neural networks with gradient descent and software to find**
708 **the best combination of momentum weight and learning rate with a grid search. Each**
709 **asset is documented in the paper and supplemental material, respectively.**

710 Guidelines:

- 711 • The answer NA means that the paper does not release new assets.
- 712 • Researchers should communicate the details of the dataset/code/model as part of their
713 submissions via structured templates. This includes details about training, license,
714 limitations, etc.
- 715 • The paper should discuss whether and how consent was obtained from people whose
716 asset is used.
- 717 • At submission time, remember to anonymize your assets (if applicable). You can either
718 create an anonymized URL or include an anonymized zip file.

719 14. Crowdsourcing and Research with Human Subjects

720 Question: For crowdsourcing experiments and research with human subjects, does the paper
721 include the full text of instructions given to participants and screenshots, if applicable, as
722 well as details about compensation (if any)?

723 Answer: [NA]

724 Justification: **The paper does not involve crowdsourcing nor research with human**
725 **subjects.**

726 Guidelines:

- 727 • The answer NA means that the paper does not involve crowdsourcing nor research with
728 human subjects.
- 729 • Including this information in the supplemental material is fine, but if the main contribu-
730 tion of the paper involves human subjects, then as much detail as possible should be
731 included in the main paper.
- 732 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
733 or other labor should be paid at least the minimum wage in the country of the data
734 collector.

735 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human**
736 **Subjects**

737 Question: Does the paper describe potential risks incurred by study participants, whether
738 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
739 approvals (or an equivalent approval/review based on the requirements of your country or
740 institution) were obtained?

741 Answer: [NA]

742 Justification: **The paper does not involve crowdsourcing nor research with human**
743 **subjects.**

744 Guidelines:

- 745 • The answer NA means that the paper does not involve crowdsourcing nor research with
746 human subjects.
- 747 • Depending on the country in which research is conducted, IRB approval (or equivalent)
748 may be required for any human subjects research. If you obtained IRB approval, you
749 should clearly state this in the paper.
- 750 • We recognize that the procedures for this may vary significantly between institutions
751 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
752 guidelines for their institution.
- 753 • For initial submissions, do not include any information that would break anonymity (if
754 applicable), such as the institution conducting the review.

755 **A Appendix / supplemental material**

756 Two grid searches for the publicly available MNIST dataset were performed to corroborate the
757 learning rate and momentum weight derived in the main paper (LeCun et al., accessed May 21, 2024).
758 The MNIST dataset contains gray-scale images of handwritten digits and is one of the prominent
759 datasets used to evaluate machine learning methods. It is split into a training and a test set, where the
latter serves as a standard of comparison. Figure 2 shows an example of the MNIST data.

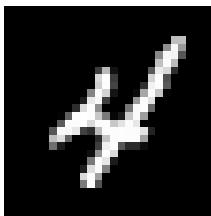


Figure 2: A slightly enlarged example from the MNIST dataset showing a handwritten digit (4).

760

761 **A.1 Experiments**

762 The grid searches were performed on the full-size MNIST dataset and a smaller version of MNIST
763 containing only 50% of the training data. In the latter case, a stratified sampling method named
764 *StratifiedShuffleSplit* was used to create a stratified random subset of the training samples (Scikit-
765 learn developers, BSD License; Pedregosa et al., 2011; Buitinck et al., 2013). This ensured that the
766 class distribution in the training subset was the same as in the original full-size training set. The
767 degradation in dataset size allowed observing how each optimizer performed under varying amounts
768 of training data, assuming that providing less training data posed a harder problem.

769 A deep learning model was trained based on a convolutional neural network (CNN). The model
770 consisted of two convolutional layers, each followed by a ReLU activation function and a max
771 pooling operation. The first convolutional layer had a single-channel input (grayscale image) and
772 applied 16 filters, followed by a second convolutional layer that expanded the channel size to 32.
773 Both convolutional layers used a 3x3 kernel size, a stride of one, and a padding of one. After each
774 convolution, a ReLU activation function introduced non-linearity, and a max pooling operation with

775 a 2x2 kernel and stride reduced the spatial dimensions by half. A dropout layer with a rate of 0.25
 776 was applied after flattening the output to prevent overfitting. The network concluded with two fully
 777 connected layers with a final output of 10 classes, where the maximum output value determined the
 778 class of an input image. The number of parameters was around two hundred thousand for an MNIST
 779 input image of size 28x28. A weight initialization was performed using the Kaiming uniform method.
 780 No data augmentation techniques were applied; however, the input was normalized to the range [-1,1].
 781 The training used a batch size of 64 and was conducted over 30 epochs, employing cross entropy
 782 as the loss function. The sizes of the training, validation, and test datasets were 54,000, 6,000, and
 783 10,000, respectively. Finally, the model’s performance was assessed through 10-fold cross-validation.

784 **A.2 Results**

785 The results of both grid searches are shown in Figure 3 for the full-size training set and in Figure 3
 for the smaller training set with 50% of the size. The following values were used as momentum

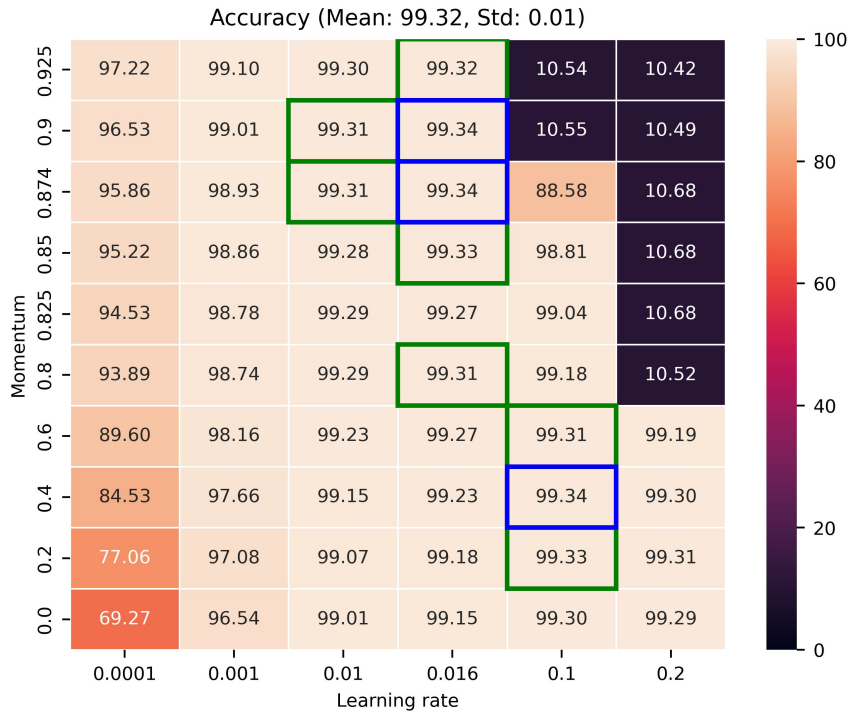


Figure 3: Grid search results for MNIST

786 weights for each grid search: 0, 0.2, 0.4, 0.6, 0.8, 0.825, 0.85, 0.874, 0.9, and 0.925. On the other
 787 hand, the following values were used as learning rates: 0.0001, 0.001, 0.01, 0.016, 0.1, 0.2. These
 788 values included the momentum weight derived in the paper ($\alpha \approx 0.874$) and the derived learning
 789 rate ($\eta \approx 0.016$). Other values were chosen based on their use in the literature or to increase the
 790 resolution around the derived theoretical values. All possible combinations of values span a 6x10
 791 grid. The color of each square in the grids of Figure 3 and Figure 4 represent the performance of the
 792 corresponding pair of momentum weight and learning rate, with lighter colors representing higher
 793 performance. Green rectangles indicate the top ten performing pairs, whereas blue rectangles show
 794 the best-performing pair. Note that more than one pair can share the best performance, as in Figure 3.
 795

796 Figure 3 shows that no pair of momentum weight and learning rate provides better performance
 797 on the full-size MNIST set than the pair derived in the paper, (0.016, 0.874), although this pair has to
 798 share its first place with other pairs. The classification accuracies for the reduced training set size
 799 are slightly lower in the table of Figure 4, as one would expect for a problem with less training data.

Nevertheless, the theoretical values derived in the paper for momentum weight and learning rate show

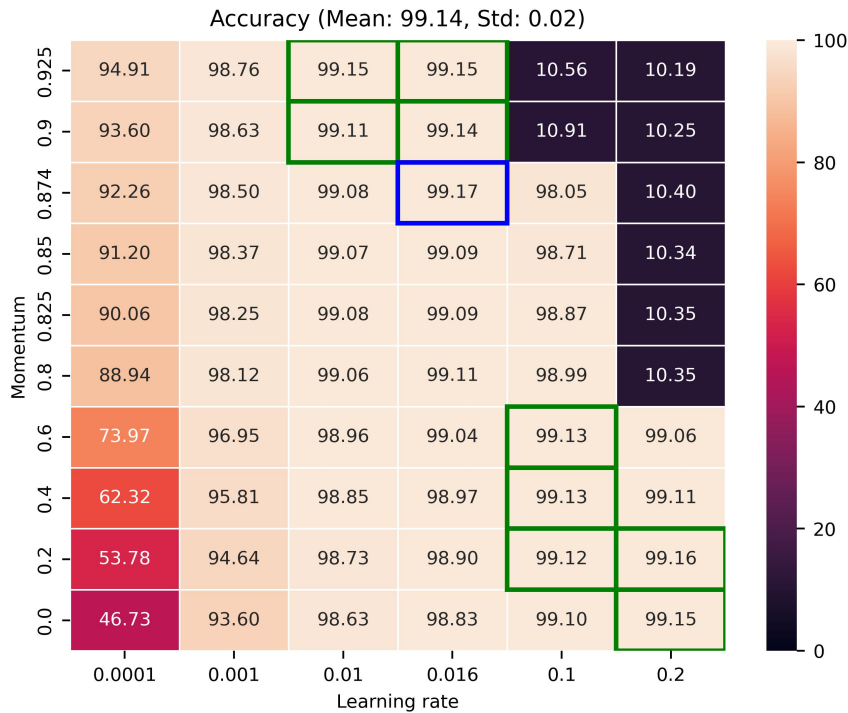


Figure 4: Grid search results for MNIST using only 50% of the training data

800
801 again the best performance.

802 A.3 Computational environment and runtime

803 The software was developed using Python 3.10, and the Convolutional Neural Network (CNN) model
804 was implemented in Pytorch 2.2.2. For each combination of learning rate and momentum weight (60
805 combinations in total), the training time was approximately three hours for 100% of the training set
806 size and about 1.5 hours for 50% of the training set. Consequently, the cumulative GPU time for all
807 experiments was approximately $(3 + 1.5) \times 60$ hours, which is 270 hours. The average memory usage
808 was roughly 1 GB for each combination. For more information about the software requirements and
809 workflow, see the Readme file uploaded as supplemental material together with the code.

810 A.4 Computing cluster

811 Figure 5 shows an overview of the GPU computing cluster that was available for the experiments,
812 including the type of GPUs among which the processing was distributed.

GPU nodes	Processor cores per node	Memory	Network
36	32 x 2.8 GHz (AMD Epyc 7543p) hyperthreading enabled 256 MB level 3 cache 4 x NVIDIA A100 GPUs (80 GB VRAM, 6912 cores, 432 Tensor cores) NVLINK	256 GB	200 Gb/s HDR Infiniband (1:1)
56	36 x 2.3 GHz (Intel Gold 6140) hyperthreading enabled 25 MB secondary cache 4 x NVIDIA V100-SXM2 GPUs (32 GB VRAM, 5120 cores, 640 Tensor cores) NVLINK	384 GB	200 Gb/s HDR Infiniband (1:1)
8	28 x 2.4 GHz (Intel E5-2680v4) hyperthreading enabled 35 MB secondary cache 4 x NVIDIA V100 GPUs (16 GB VRAM, 5120 cores, 640 Tensor cores)	128 GB	56 Gb/s FDR Infiniband (1.11:1)
48	28 x 2.4 GHz (Intel E5-2680v4) hyperthreading enabled 35 MB secondary cache 4 x NVIDIA P100 GPUs (16 GB VRAM, 3584 cores)	128 GB	56 Gb/s FDR Infiniband (1.11:1)
72	28 x 2.4 GHz (Intel E5-2680v4) hyperthreading enabled 35 MB secondary cache 2 x NVIDIA K80 GPUs with 2 x GK210 GPUs each (24 GB VRAM, 4992 cores)	256 GB	56 Gb/s FDR Infiniband (1.11:1)

Figure 5: GPU computing cluster