

A BAYESIAN PERSPECTIVE ON GENERALIZATION AND STOCHASTIC GRADIENT DESCENT

Samuel L. Smith* & Quoc V. Le

Google Brain,

{slsmith, qvl}@google.com

ABSTRACT

This paper tackles two related questions at the heart of machine learning: how can we predict if a minimum will generalize to the test set, and why does stochastic gradient descent find minima that generalize well? Our work is inspired by Zhang et al. (2016), who showed deep networks can easily memorize randomly labeled training data, despite generalizing well when shown real labels of the same inputs. We show here that the same phenomenon occurs in small linear models. These observations are explained by evaluating the Bayesian evidence, which penalizes sharp minima but is invariant to model parameterization. We also explore the “generalization gap” between small and large batch training, identifying an optimum batch size which maximizes the test set accuracy. Interpreting stochastic gradient descent as a stochastic differential equation, we identify a “noise scale” $g = \epsilon(\frac{N}{B} - 1) \approx \epsilon N/B$, where ϵ is the learning rate, N training set size and B batch size. Consequently the optimum batch size is proportional to the learning rate and the training set size, $B_{opt} \propto \epsilon N$. We verify these predictions empirically.

1 INTRODUCTION

Zhang et al. (2016) argued that deep learning cannot be understood without rethinking our understanding of generalization. To support this claim, they trained deep convolutional networks on both ImageNet and CIFAR10, achieving excellent predictive accuracy on both the training and test sets. They then took the same input images, but randomized the labels, and showed that while the model is now completely unable to generalize to the test set, it still successfully memorizes the training set labels. They demonstrated that these results cannot be understood within the framework of classical theories of statistical learning like Rademacher complexity (Bartlett & Mendelson, 2002), VC-dimension (Vapnik, 1998) or uniform stability (Mukherjee et al., 2006).

Many authors have suggested that “broad minima” whose curvature is small may generalize better than “sharp minima” whose curvature is large (Chaudhari et al., 2016; Hochreiter & Schmidhuber, 1997). Following this intuition, Krueger et al. (2017) observed that the magnitude of the largest Hessian eigenvalue increased significantly when training on random labels. However Dinh et al. (2017) argued that curvature alone cannot explain generalization, since the curvature of a minimum can be arbitrarily increased by changing the model parameterization. Keskar et al. (2016) invoked curvature to explain why stochastic gradient descent (SGD) often generalizes better with smaller training batches. Goyal et al. (2017) observed a linear scaling rule between batch size and learning rate in a deep ResNet, while Wilson et al. (2017) found SGD consistently generalized better to the test set than second-order methods. Our work makes two main contributions:

- The results of Zhang et al. (2016) are not unique to deep learning; we observe the same phenomenon in a small linear model. We demonstrate these observations are explained by evaluating the Bayesian evidence in favor of each model, which penalizes sharp minima.
- SGD integrates a stochastic differential equation whose “noise scale” $g \approx \epsilon N/B$, where ϵ is the learning rate, N training set size and B batch size. Noise drives SGD away from sharp minima, and there is an optimal noise scale which maximizes the test accuracy. Consequently the optimal batch size is proportional to both learning rate and training set size.

*Work done as a member of the Google Brain Residency Program (g.co/brainresidency)

We introduce Bayesian model comparison in section 2. In section 3 we replicate the observations of Zhang et al. (2016) in a linear model, and show they are explained by evaluating the Bayesian evidence. We explore the “generalization gap” in section 4 and show that there is an optimum batch size. Finally in section 5 we derive scaling rules between the optimum batch size, learning rate and training set size. Throughout this work, generalization gap refers to the gap in test accuracy between small and large batch training with SGD, not the gap in accuracy between training and test sets.

2 BAYESIAN MODEL COMPARISON

Here we give a brief tutorial on Bayesian model comparison and the Gaussian approximation to the evidence, closely following the seminal work of David MacKay (1992). We can infer a probability distribution over all possible parameter values by applying Bayes theorem. Considering a classification model M with a single parameter ω , training inputs $\{x_i\}_{i=1}^N$ and training labels $\{y_i\}_{i=1}^N$,

$$P(\omega|\{y\}, \{x\}; M) = \frac{P(\{y\}|\omega, \{x\}; M)P(\omega; M)}{P(\{y\}|\{x\}; M)} \quad (1)$$

$$\propto P(\{y\}|\omega, \{x\}; M)P(\omega; M). \quad (2)$$

The likelihood, $P(\{y\}|\omega, \{x\}; M) = \prod_i P(y_i|\omega, x_i; M) = e^{-H(\omega; M)}$, where we identify the cross-entropy for one-hot labels, $H(\omega; M) = -\sum_i \ln(P(y_i|\omega, x_i; M))$. We typically use an uncorrelated Gaussian prior, $P(\omega; M) = \sqrt{\lambda/2\pi}e^{-\lambda\omega^2/2}$, and we discuss the Bayesian motivation for feature whitening in appendix A. Thus the probability density of the parameter given the training data, $P(\omega|\{y\}, \{x\}; M) \propto \sqrt{\lambda/2\pi}e^{-C(\omega; M)}$, where $C(\omega; M) = H(\omega; M) + \lambda\omega^2/2$ denotes the L2 regularized cross entropy, hereafter referred to as the “cost function”, and λ is the regularization coefficient. The value ω_0 which minimizes this cost function, maximizes the probability density of the parameter. To predict an unknown label y_t of a new input x_t , we should compute the integral,

$$P(y_t|x_t, \{x\}, \{y\}; M) = \int d\omega P(y_t|\omega, x_t; M)P(\omega|\{y\}, \{x\}; M) \quad (3)$$

$$= \frac{\int d\omega P(y_t|\omega, x_t; M)e^{-C(\omega; M)}}{\int d\omega e^{-C(\omega; M)}}. \quad (4)$$

However these integrals are dominated by the region near ω_0 , and since $P(y_t|\omega, x_t; M)$ is smooth we usually approximate $P(y_t|x_t; M) \approx P(y_t|\omega_0, x_t; M)$. Having minimized $C(\omega; M)$ to find ω_0 , we now wish to compare two or more different models and select the best one. The probability ratio,

$$\frac{P(M_1|\{y\}, \{x\})}{P(M_2|\{y\}, \{x\})} = \frac{P(\{y\}|\{x\}; M_1) P(M_1)}{P(\{y\}|\{x\}; M_2) P(M_2)}. \quad (5)$$

The first factor on the right hand side is the evidence ratio, while the second is the prior ratio. The prior ratio describes which model is most plausible. To avoid unnecessary subjectivity, we usually set this to 1. Meanwhile the evidence ratio decides how much the training data should change our prior beliefs. Usually the evidence ratio of two models is exponentially large, such that any reasonable prior leads to the same conclusions. In statistics the evidence is also called the “marginal likelihood”. To compute it, we evaluate the normalizing constant in the denominator of equation 1,

$$P(\{y\}|\{x\}; M) = \int d\omega P(\{y\}|\omega, \{x\}; M)P(\omega; M) \quad (6)$$

$$= \sqrt{\frac{\lambda}{2\pi}} \int d\omega e^{-C(\omega; M)}. \quad (7)$$

Once again, this integral is dominated by the region surrounding the minimum ω_0 , so we can estimate the evidence by Taylor expanding $C(\omega; M) \approx C(\omega_0) + C''(\omega_0)(\omega - \omega_0)^2/2$,

$$P(\{y\}|\{x\}; M) \approx e^{-C(\omega_0)} \sqrt{\frac{\lambda}{2\pi}} \int d\omega e^{-C''(\omega_0)(\omega - \omega_0)^2/2} \quad (8)$$

$$= e^{-C(\omega_0)} \sqrt{\frac{\lambda}{C''(\omega_0)}} \quad (9)$$

$$= \exp \left\{ - \left(C(\omega_0) + \frac{1}{2} \ln(C''(\omega_0)/\lambda) \right) \right\}. \quad (10)$$

Within this ‘‘Gaussian’’ approximation, the evidence in favor of a model is controlled by the value of the cost function at the minimum, and by the logarithm of the ratio of the curvature about this minimum compared to the regularization constant. Thus far we have considered models of a single parameter; in realistic models with many parameters, $P(\{y\}|\{x\}; M) = \lambda^{\frac{P}{2}} e^{-C(\omega_0)} / |\nabla\nabla C(\omega)|_{\omega_0}^{1/2}$, where $|\nabla\nabla C(\omega)|_{\omega_0}$ is the determinant of the Hessian, and P denotes the number of model parameters. The determinant of a matrix is simply the product of its eigenvalues, $(\prod_{i=1}^P \lambda_i)$, and thus,

$$P(\{y\}|\{x\}; M) \approx \exp \left\{ - \left(C(\omega_0) + \frac{1}{2} \sum_{i=1}^P \ln(\lambda_i/\lambda) \right) \right\}. \quad (11)$$

The contribution $(\lambda^{\frac{P}{2}} / |\nabla\nabla C(\omega)|_{\omega_0}^{1/2})$ is often called the ‘‘Occam factor’’, because it enforces Occam’s razor; *when two models describe the data equally well, the simpler model is usually better* (Gull, 1988). Minima with low curvature are simple, because the parameters do not have to be fine-tuned to fit the data. Intuitively, the Occam factor describes the fraction of the prior parameter space consistent with the data. Since this fraction is always less than one, we propose to approximate equation 11 away from local minima by only performing the summation over eigenvalues $\lambda_i \geq \lambda$.

In this work, we will compare the evidence against a null model which assumes the labels are entirely random, assigning equal probability to each class. This unusual model has no parameters, and so the evidence is controlled by the likelihood alone, $P(\{y\}|\{x\}; NULL) = (1/n)^N = e^{-N \ln(n)}$, where n denotes the number of model classes and N the number of training labels. Thus the evidence ratio,

$$\frac{P(\{y\}|\{x\}; M)}{P(\{y\}|\{x\}; NULL)} = e^{-E(\omega_0)}, \quad (12)$$

Where $E(\omega_0) = C(\omega_0) + (1/2) \sum_i \ln(\lambda_i/\lambda) - N \ln(n)$ is the log evidence ratio in favor of the null model. Clearly, we should only assign any confidence to the predictions of our model if $E(\omega_0) < 0$.

The Bayesian evidence supports the intuition that broad minima generalize better than sharp minima, but unlike the curvature it does not depend on the model parameterization. Dinh et al. (2017) showed one can increase the Hessian eigenvalues (curvature) by rescaling the parameters, but they must simultaneously rescale the regularization parameter, otherwise the model changes. Since Occam’s factor arises from the log ratio, $\ln(\lambda_i/\lambda)$, these two effects cancel out. Unfortunately it is difficult to evaluate the evidence for deep networks, since we cannot compute the Hessian of models with millions of parameters. Neural networks also exhibit many equivalent minima, since we can permute the hidden units without changing the model. To compute the evidence of deep networks we must carefully account for this ‘‘degeneracy’’. No degeneracy arises in linear models. The Bayesian evidence can be reframed in the language of information theory, whereby Occam’s factor penalizes the amount of information the model must learn about the parameters to accurately model the training data (Hinton & Van Camp, 1993; Achille & Soatto, 2017; Shwartz-Ziv & Tishby, 2017).

3 BAYES THEOREM AND GENERALIZATION

We now turn to the experimental observations of Zhang et al. (2016); namely that deep neural networks generalize well when shown training data with real labels, and yet the same model can drastically overfit when the labels are randomized; perfectly memorizing the training set. To demonstrate that these observations are not unique to large or deep networks, let’s consider a far simpler model; logistic regression. We define a toy task, whereby our inputs are randomly generated 200D vectors x , and each component x_i is drawn from the Gaussian distribution. We normalize the inputs, to ensure $\langle x_i \rangle = 0$ and $\langle x_i^2 \rangle = 1$. Each training example is assigned to one of two classes, $y = \{0, 1\}$, and we train using the L2-regularized sigmoid cross-entropy loss. There are two tasks. In the first task, the labels of both the training and test sets are random. In the second task, the label $y = 1$ if $\sum_i x_i > 0$, while $y = 0$ otherwise. Our model has just 200 weights and 1 bias. For both tasks, our training set contained 200 examples, while our test set contained 10000 examples.

We show the accuracy of our model predictions on both the training and test sets in figure 1. When trained on the informative labels, the model generalizes well to the test set, so long as it is sufficiently weakly regularized. However the model also perfectly memorizes the random labels, replicating the observations of Zhang et al. (2016) in deep networks. No improvement in model performance is observed as the regularization coefficient is increased.

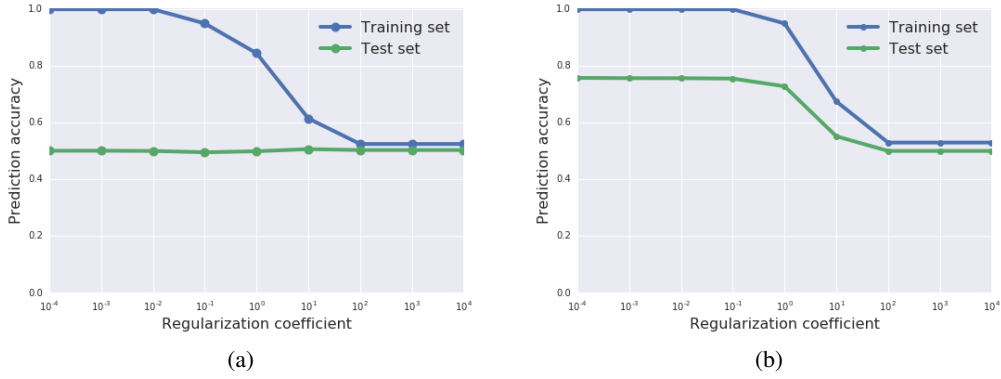


Figure 1: Prediction accuracy as a function of regularization coefficient, for a logistic regression trained on random (a) and informative (b) labels. Exactly as observed by Zhang et al., weakly regularized logistic regression generalizes well on informative labels but memorizes random labels.

Zhang et al. (2016) argued these results are inconsistent with classical statistical learning theory. We show here that they are explained by Bayesian model comparison. Consider figure 2, where we plot the mean cross-entropy of the model predictions, evaluated on both training and test sets, as well as the Bayesian log evidence ratio defined in the previous section. Looking first at figure 2a, while the cross-entropy on the training set vanishes when the model is weakly regularized, the cross-entropy on the test set explodes. Not only does the model make random predictions, but it is extremely confident in those predictions. As the regularization coefficient is increased the test set cross-entropy falls, settling at $\ln 2$, the cross-entropy of assigning equal probability to both classes (The model is now so heavily regularized that it cannot learn from the training data). Now consider the Bayesian evidence, which we evaluate *on the training set*. The log evidence ratio is large and positive when the model is weakly regularized, indicating that the model is exponentially less plausible than assigning equal probabilities to each class. As the regularization parameter is increased, the log evidence ratio falls, but it is always positive, indicating that the model can never be expected to generalize well.

Now consider figure 2b. Once again, the training cross-entropy falls to zero when the model is weakly regularized, while the test cross-entropy is high. Even though the model is making accurate predictions, those predictions are overconfident. As the regularization is increased, the test cross-entropy falls below $\ln 2$, indicating that the model is successfully generalizing to the test set. Now consider the Bayesian evidence. The log evidence ratio is large and positive when the model is weakly regularized, but as the regularization coefficient is increased, the log evidence ratio drops *below zero*, indicating that the model is exponentially more plausible than assigning equal probabilities to each class (when properly regularized). As we further increase the regularization, the log

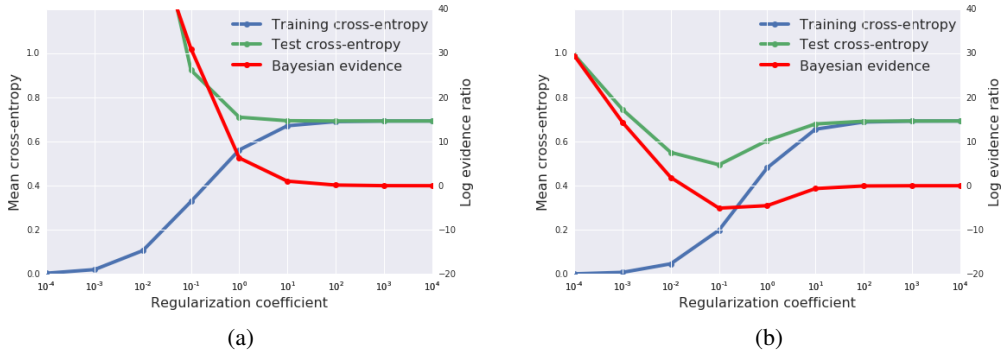


Figure 2: The cross-entropy and log evidence ratio, evaluated on random (a) or informative (b) labels. The evidence, evaluated on the training set, is highly correlated with the test cross-entropy.

evidence ratio rises to zero while the test cross-entropy rises to $\ln 2$. Test cross-entropy and Bayesian evidence are strongly correlated, with minima at the same regularization strength.

Thus the Bayesian evidence has successfully explained the generalization of our logistic regression. Meanwhile in deep networks, Krueger et al. (2017) showed that the largest Hessian eigenvalue increased when training on random labels; implying that the Occam factor decreases, indicative of a decrease in the model evidence. We conclude that Bayesian model comparison is quantitatively consistent with the results of Zhang et al. (2016) in linear models where we can compute the evidence, and qualitatively consistent with their results in deep networks where we cannot.

4 BAYES THEOREM AND STOCHASTIC GRADIENT DESCENT

In the proceeding section, we showed that generalization is strongly correlated with the Bayesian evidence, which is a weighted combination of the depth of a minimum (the cost function) and its breadth (the Occam factor). Therefore, when we train deep networks, we should not minimize the cost function; we should try to find the local minimum which maximizes the evidence. To achieve this, Bayesians add noise to the gradient updates (Mandt et al., 2017; Welling & Teh, 2011). We propose that the same principles account for the “generalization gap” (Keskar et al., 2016), whereby the test accuracy often falls as the SGD batch size is increased. Small batch training introduces noise to the gradients, and this noise drives the SGD away from sharp minima, thus enhancing generalization. Since the gradient drives the SGD towards deep minima, while noise drives the SGD towards broad minima; we expect the test set performance to be maximized at an optimal batch size, which introduces the right amount of noise to balance these competing contributions to the evidence.

We were unable to observe a generalization gap in linear models. This is consistent with Bayesian intuition; since a linear model has a single convex minimum there are no sharp minima to avoid. Consequently we introduce a new model; a shallow neural network with 800 hidden units, trained on MNIST without regularization. We use SGD with momentum, with a momentum parameter of 0.9. In our initial experiments, we use a constant learning rate of 1.0 which does not depend on the batch size or decay during training. Furthermore, we train on just 1000 images, selected at random from the MNIST training set, which enables us to compare small batch to full batch training. We emphasize that we are not trying to achieve optimal performance, but to explore the simplest possible model which shows a generalization gap between small and large batch training.

In figure 3, we exhibit the evolution of the test accuracy and test cross-entropy during training. Our small batches are composed of 30 images, randomly sampled from the training set. Looking first at figure 3a, small batch training takes longer to converge, but after a thousand gradient updates a clear generalization gap in model accuracy emerges between small and large training batches. Now consider figure 3b. While the test cross-entropy for small batch training is lower at the end of training; the cross-entropy of both small and large training batches is increasing, indicative of

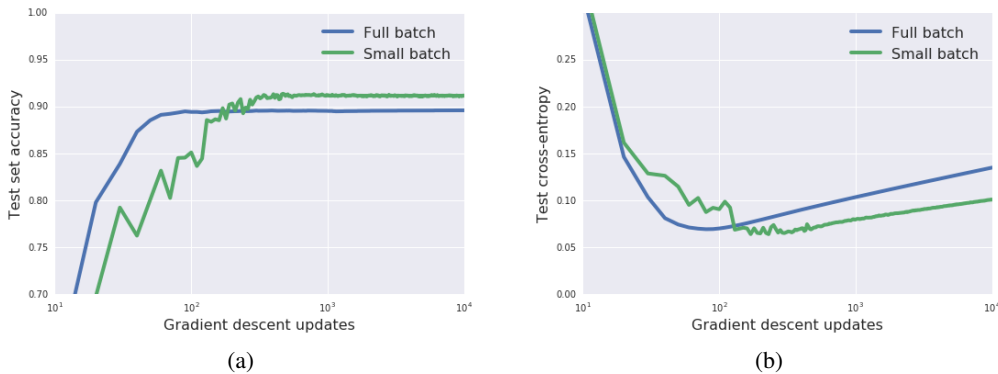


Figure 3: The generalization gap between full batch and small batch training. Full batches are composed of 1000 images, while small batches comprise 30 images. In figure a, we exhibit the evolution of the test accuracy during training. In figure b, we exhibit the test cross-entropy.

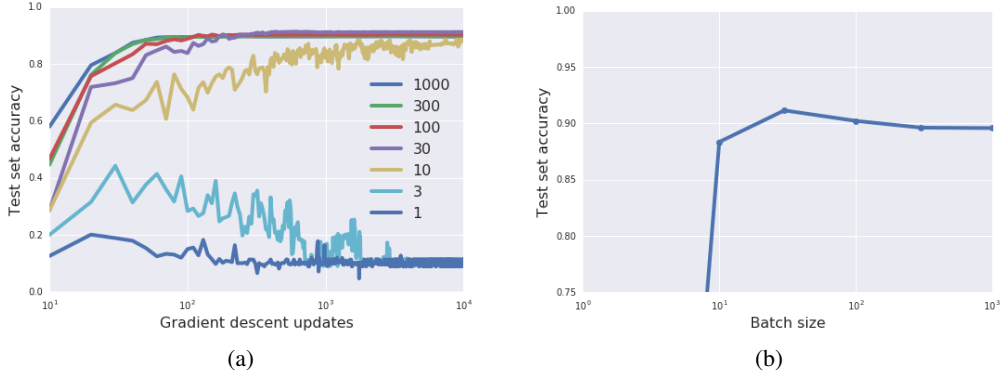


Figure 4: The test accuracy for a range of batch sizes, during training (a) and after 10000 steps (b).

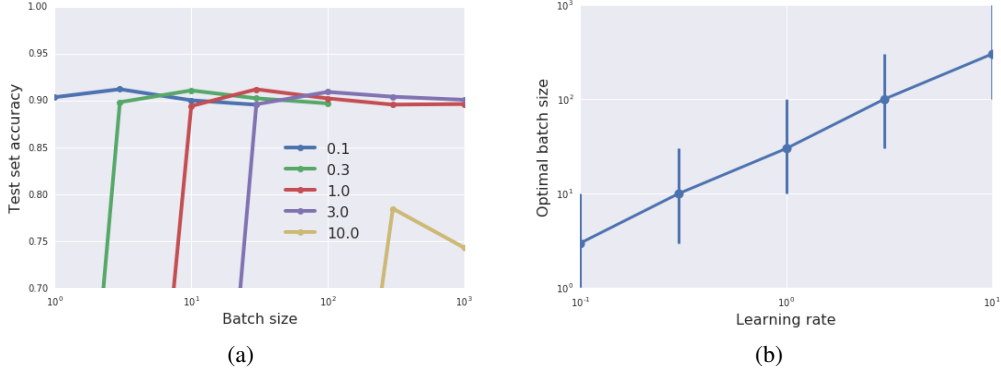


Figure 5: a) The test set accuracy as a function of batch size, for a range of learning rates ϵ . The performance peak shifts to the right as we increase ϵ , but the overall performance falls once $\epsilon \gtrsim 3$. b) The best observed batch size is proportional to the learning rate across two orders of magnitude.

over-fitting. This is another reminder that model accuracy and calibration may not be strongly correlated. Both models exhibit a minimum test cross-entropy, although after different numbers of gradient updates. Small batch training reaches a lower minimum but also exhibits larger fluctuations. Intriguingly, we show in appendix B that the generalization gap between small and large batch training shrinks significantly when we introduce L2 regularization.

From now on we focus on the test set accuracy (since this converges as the number of gradient updates increases). In figure 4a, we exhibit training curves for a range of batch sizes between 1 and 1000. We find that the model cannot train when the batch size $B \lesssim 10$. In figure 4b we plot the mean test set accuracy after 10000 training steps. A clear peak emerges, indicating that there is indeed an optimum batch size which maximizes the test accuracy, as predicted by Bayesian intuitions.

5 STOCHASTIC DIFFERENTIAL EQUATIONS AND THE SCALING RULES

We now predict three simple scaling rules, which describe how the optimal batch size depends on the learning rate, training set size and momentum coefficient. We verify these scaling rules on MNIST using the shallow neural network introduced in section 4. We recall that SGD introduces noise to the gradient updates, and that the magnitude of this noise grows as the batch size shrinks. For clarity, we consider vanilla SGD here, and SGD with momentum in appendix C. A gradient update,

$$\Delta\omega = \frac{\epsilon}{N} \left(\frac{dC}{d\omega} + \left(\frac{d\hat{C}}{d\omega} - \frac{dC}{d\omega} \right) \right), \quad (13)$$

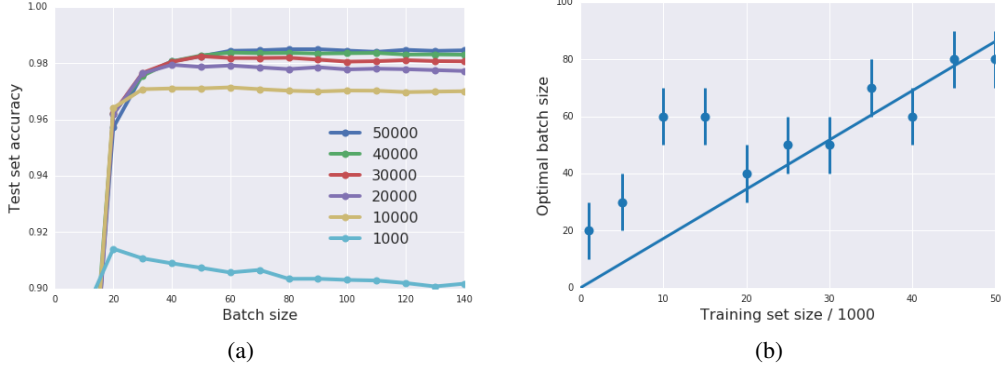


Figure 6: a) The test accuracy as a function of batch size, for a range of training set sizes. To reduce noise, we average each curve over five experiments. The performance peak shift to the right as we increase the size of the training set. Unsurprisingly, the overall model performance also improves. b) The best observed batch size is proportional to the size of the training set once $N \gtrsim 20000$.

Where ϵ is the learning rate, N is the training set size, $\frac{dC}{d\omega} = \sum_{i=1}^N \frac{dC_i}{d\omega}$ is the true gradient, and $\frac{d\hat{C}}{d\omega} = \frac{N}{B} \sum_{i=1}^B \frac{dC_i}{d\omega}$ is the gradient evaluated on a batch. The expected gradient of a single training example, $\langle \frac{dC_i}{d\omega} \rangle = \frac{1}{N} \frac{dC}{d\omega}$, while $\langle \frac{dC_i}{d\omega} \frac{dC_j}{d\omega} \rangle = (\frac{1}{N} \frac{dC}{d\omega})^2 + F(\omega) \delta_{ij}$. $F(\omega)$ is a matrix describing the gradient covariances between different parameters, which are themselves a function of the current parameter values. To proceed, we adopt the central limit theorem and model the gradient error $\alpha = \left(\frac{d\hat{C}}{d\omega} - \frac{dC}{d\omega} \right)$ with Gaussian random noise. It is easy to show that $\langle \alpha \rangle = 0$, while $\langle \alpha^2 \rangle = N(\frac{N}{B} - 1)F(\omega)$. Typically $N \gg B$, such that $\langle \alpha^2 \rangle \approx N^2 F(\omega)/B$. To continue, we interpret equation 13 as a discrete update of the stochastic differential equation (Gardiner, 1985),

$$\frac{d\omega}{dt} = \frac{dC}{d\omega} + \eta(t), \quad (14)$$

Where t is a continuous variable, $\eta(t)$ represents noise, $\langle \eta(t) \rangle = 0$ and $\langle \eta(t)\eta(t') \rangle = gF(\omega)\delta(t-t')$. The constant g controls the scale of random fluctuations in the dynamics. To relate this differential equation to the SGD, we compute a single gradient update $\Delta\omega = \int_0^{\epsilon/N} \frac{d\omega}{dt} dt = \frac{\epsilon}{N} \frac{dC}{d\omega} + \int_0^{\epsilon/N} \eta(t) dt$. Finally, to measure g , we equate the variance in this gradient update to the variance in equation 13,

$$\begin{aligned} \left(\frac{\epsilon}{N} \right)^2 \langle \alpha^2 \rangle &= \epsilon^2 \left(\frac{N}{B} - 1 \right) F(\omega) / N \\ &= \langle \left(\int_0^{\epsilon/N} dt \eta(t) \right)^2 \rangle = \int_0^{\epsilon/N} dt \int_0^{\epsilon/N} dt' \langle \eta(t)\eta(t') \rangle = \epsilon g F(\omega) / N. \end{aligned} \quad (15)$$

Rearranging, $g = \epsilon(\frac{N}{B} - 1) \approx \epsilon N/B$. SGD simulates a stochastic differential equation, and the scale of random fluctuations are inversely proportional to the batch size. Consequently, an optimal batch size emerges when the underlying scale of random fluctuations is also optimal.

To reach the same performance when changing the learning rate, we should ensure g stays constant; which implies the optimal batch size is proportional to the learning rate, $B_{opt} \propto \epsilon$. In figure 5a, we plot the test accuracy as a function of batch size after $(10000/\epsilon)$ training steps, for a range of learning rates. Exactly as predicted, the peak moves to the right as ϵ increases (Unsurprisingly, the model accuracy drops dramatically when ϵ is too large). In figure 5b, we plot the best observed batch size as a function of learning rate, observing a clear linear trend. The error bars indicate the distance from the best observed batch size to the next batch size sampled in our experiments. We note that Goyal et al. (2017) used this linear scaling rule to train a deep ResNet on ImageNet in one hour.

Furthermore, it is likely the optimal scale of random fluctuations does not change when the training set size increases. Holding g constant again, we predict that the optimal batch size is also proportional to the training set size, $B_{opt} \propto N$. To confirm this scaling rule, in figure 6a we exhibit the test set accuracy as a function of batch size for a range of training set sizes after 10000 steps ($\epsilon = 1$

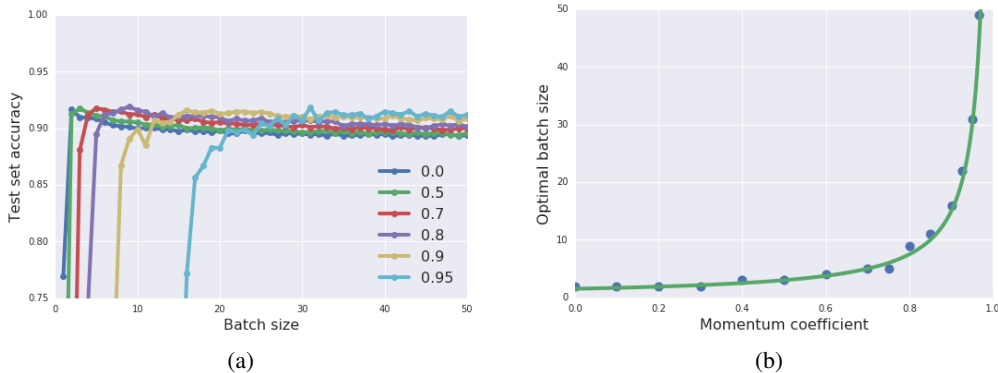


Figure 7: a) The test set accuracy as a function of batch size for a range of momentum coefficients. As expected, the peak moves to the right as the momentum coefficient increases. b) The best observed batch size for a range of momentum coefficients. The green curve exhibits the scaling rule.

everywhere). Once again, the peak shifts to the right as the training set size rises, although the generalization gap is less pronounced when the training set size increases. In figure 6b, we plot the best observed batch size as a function of training set size; observing another linear trend.

These scaling rules have two consequences. First, we can significantly increase the batch size with no loss in performance and no increase in computational cost by also increasing the learning rate. This enables us to exploit increased parallelism across multiple GPUs, reducing training time. Second, we also expect the size of many commercial datasets to grow considerably over time, and consequently we expect large batch training to become increasingly common. We note that Hoffer et al. (2017) proposed an alternative scaling rule for the learning rate; $\epsilon \propto \sqrt{B_{opt}}$. They claim that scaling the learning rate alone is sufficient to close the generalization gap; however we showed in figure 5 that this is not the case, since the model performance drops considerably once the learning rate is too large, imposing a maximum achievable batch size.

Finally, in appendix C we extend our analysis to SGD with momentum, identifying the noise scale, $g \approx \frac{\epsilon N}{B(1-m)}$, where m denotes the momentum coefficient. Notice that this reduces to the noise scale of vanilla SGD as $m \rightarrow 0$. When $m > 0$, we obtain the additional scaling rule $B_{opt} \propto 1/(1-m)$. This scaling rule predicts that the optimal batch size will increase when the momentum parameter is increased, providing an alternative strategy for achieving large batch training which has not yet been explored in the literature. In figure 7a we plot the test set performance as a function of batch size after 10000 gradient updates ($\epsilon = 1$ everywhere), for a range of momentum coefficients. As predicted, the batch size required to train the model increases when the momentum parameter increases. Finally in figure 7b, we plot the best observed batch size as a function of momentum, and fit our results to the new scaling rule above; obtaining remarkably good agreement. We propose a simple heuristic for tuning the batch size, learning rate and momentum coefficient in appendix D.

6 CONCLUSIONS

Just like deep neural networks, linear models which generalize well on informative labels can memorize random labels of the same inputs. These observations contradict classical statistical learning theory, but they are explained by evaluating the Bayesian evidence, or “marginal likelihood”. This Bayesian evidence penalizes sharp minima, but it is invariant to changes in model parameterization. We explore the generalization gap between small and large batch training, and show that there is an optimum batch size which maximizes the test set accuracy. This is consistent with Bayesian intuition, since gradient noise drives SGD away from sharp minima. Interpreting the SGD as the discretization of a stochastic differential equation, we predict that this optimum batch size should scale linearly with both the learning rate and the size of the training set, $B_{opt} \propto \epsilon N$. For SGD with momentum, we derive an additional scaling rule, $B_{opt} \propto 1/(1-m)$ between the optimum batch size and momentum coefficient. We verify these scaling rules empirically and discuss their implications.

ACKNOWLEDGMENTS

We thank Pieter-Jan Kindermans, Prajit Ramachandran, Jascha Sohl-Dickstein, Jon Shlens, Kevin Murphy, Samy Bengio, Yasaman Bahri and Saeed Saremi for helpful comments on the manuscript.

REFERENCES

- Alessandro Achille and Stefano Soatto. On the emergence of invariance and disentangling in deep representations. *arXiv preprint arXiv:1706.01350*, 2017.
- Peter L Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3:463–482, 2002.
- Pratik Chaudhari, Anna Choromanska, Stefano Soatto, and Yann LeCun. Entropy-SGD: Biasing gradient descent into wide valleys. *arXiv preprint arXiv:1611.01838*, 2016.
- Laurent Dinh, Razvan Pascanu, Samy Bengio, and Yoshua Bengio. Sharp minima can generalize for deep nets. *arXiv preprint arXiv:1703.04933*, 2017.
- Crispin W Gardiner. *Handbook of Stochastic Methods*, volume 4. Springer Berlin, 1985.
- Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch SGD: Training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017.
- Stephen F Gull. Bayesian inductive inference and maximum entropy. In *Maximum-entropy and Bayesian methods in science and engineering*, pp. 53–74. Springer, 1988.
- Geoffrey E Hinton and Drew Van Camp. Keeping the neural networks simple by minimizing the description length of the weights. In *Proceedings of the sixth annual conference on Computational learning theory*, pp. 5–13. ACM, 1993.
- Sepp Hochreiter and Jürgen Schmidhuber. Flat minima. *Neural Computation*, 9(1):1–42, 1997.
- Elad Hoffer, Itay Hubara, and Daniel Soudry. Train longer, generalize better: closing the generalization gap in large batch training of neural networks. *arXiv preprint arXiv:1705.08741*, 2017.
- Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang. On large-batch training for deep learning: Generalization gap and sharp minima. *arXiv preprint arXiv:1609.04836*, 2016.
- David Krueger, Nicolas Ballas, Stanislaw Jastrzebski, Devansh Arpit, Maxinder S Kanwal, Tegan Maharaj, Emmanuel Bengio, Asja Fischer, and Aaron Courville. Deep nets don’t learn via memorization. *ICLR Workshop*, 2017.
- David JC MacKay. A practical bayesian framework for backpropagation networks. *Neural computation*, 4(3):448–472, 1992.
- Stephan Mandt, Matthew D Hoffman, and David M Blei. Stochastic gradient descent as approximate bayesian inference. *arXiv preprint arXiv:1704.04289*, 2017.
- Sayan Mukherjee, Partha Niyogi, Tomaso Poggio, and Ryan Rifkin. Learning theory: stability is sufficient for generalization and necessary and sufficient for consistency of empirical risk minimization. *Advances in Computational Mathematics*, 25(1):161–193, 2006.
- Ravid Shwartz-Ziv and Naftali Tishby. Opening the black box of deep neural networks via information. *arXiv preprint arXiv:1703.00810*, 2017.
- Vladimir Vapnik. *Statistical learning theory*, volume 1. Wiley New York, 1998.
- Max Welling and Yee W Teh. Bayesian learning via stochastic gradient langevin dynamics. In *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, pp. 681–688, 2011.

Ashia C Wilson, Rebecca Roelofs, Mitchell Stern, Nathan Srebro, and Benjamin Recht. The marginal value of adaptive gradient methods in machine learning. *arXiv preprint arXiv:1705.08292*, 2017.

Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. *arXiv preprint arXiv:1611.03530*, 2016.

A FEATURE WHITENING AND BAYESIAN REGULARIZATION

As we showed in the main text, L2 regularization introduces an uncorrelated Gaussian prior centered on the origin. Most authors use a single regularization coefficient λ for all parameters. In the absence of any information about the relevance of each parameter, one could ensure this prior is plausible by checking when we initialize the model that $\langle \frac{\partial y}{\partial \omega_i} \rangle = 0$ and $\langle \frac{\partial y}{\partial \omega_i} \frac{\partial y}{\partial \omega_j} \rangle = \delta_{ij}$ for all parameters $\{\omega_i, \omega_j\}$. Here y denotes the predictions of a regression task or the logits of a classification. These considerations provide a Bayesian justification for feature whitening in linear models. We suspect that similar considerations also arise when batch normalization is used in deep networks.

B THE EFFECT OF REGULARIZATION ON THE GENERALIZATION GAP

In the experiments of section 4 of the main text, the L2 regularization coefficient $\lambda = 0$. In figure 8, we plot the evolution of the training curves when $\lambda = 0.1$, for both small batch and full batch training. Excluding the regularization parameter, these experiments are identical to figure 3. To our surprise, regularized full batch training took longer to converge than small batch training. In another surprise, regularization significantly reduced the size of the generalization gap. While large batch regularized training achieves slightly lower test set accuracy than unregularized small batch training, it also achieves lower test cross-entropy. The test cross-entropy of our regularized models does not degrade after many gradient updates, removing the need for early stopping.

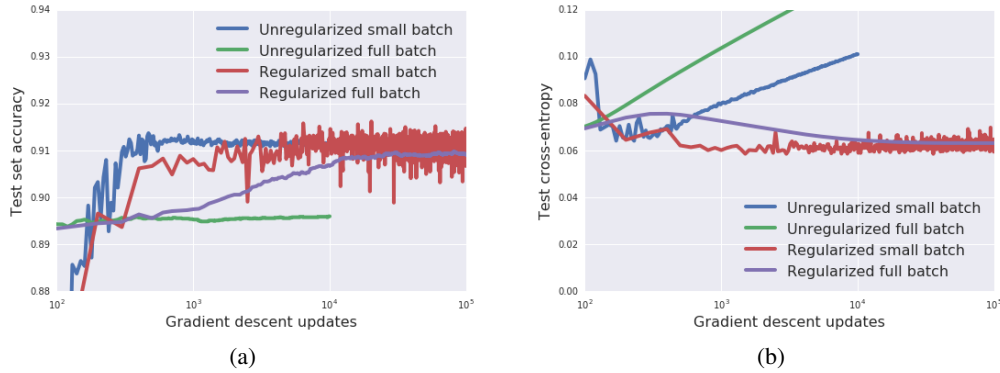


Figure 8: The mean test accuracy (a) and the mean test cross-entropy (b) of a regularized model during training. While full batch training takes longer to converge, it achieves similar performance at long times. The noise inherent in small batch training causes the performance to fluctuate.

C DERIVING THE SCALING RULES FOR SGD WITH MOMENTUM

Momentum simulates a Langevin equation (with structured fluctuations),

$$\frac{d^2\omega}{dt^2} = -\lambda \frac{d\omega}{dt} - \frac{dC}{d\omega} + \eta(t). \quad (16)$$

λ is the “damping coefficient” and $\eta(t)$ describes Gaussian noise, whose statistics $\langle \eta(t) \rangle = 0$ and $\langle \eta(t)\eta(t') \rangle = g\lambda F(\omega)\delta(t - t')$. As before, the coefficient g describes the scale of random fluctuations in the dynamics, and $F(\omega)$ describes the gradient covariances between parameters. We include

a factor of λ in the noise variance to satisfy the fluctuation-dissipation theorem, which states that we can vary the damping coefficient without changing the probability of sampling any particular configuration of parameters in the limit $t \rightarrow \infty$, if we proportionally increase the noise variance.

To relate this Langevin equation to the usual momentum equations, we first re-express it as two coupled first order differential equations,

$$\frac{dp}{dt} = -\lambda p - \frac{dC}{d\omega} + \eta(t), \quad (17)$$

$$\frac{d\omega}{dt} = p. \quad (18)$$

Integrating over a single step $\Delta t/N$,

$$\Delta p = -(\lambda \Delta t/N)p - \frac{\Delta t}{N} \frac{dC}{d\omega} + \eta, \quad (19)$$

$$\Delta \omega = p \Delta t/N. \quad (20)$$

Where now $\langle \eta \rangle = 0$ and $\langle \eta^2 \rangle = g \Delta t \lambda F(w)/N$. We define the accumulation $A = p/\Delta t$,

$$\Delta A = -(\lambda \Delta t/N)A - \frac{1}{N} \frac{dC}{d\omega} + \frac{\eta}{\Delta t}, \quad (21)$$

$$\Delta \omega = (\Delta t)^2 A/N. \quad (22)$$

These equations can be compared to the TensorFlow update equations for momentum,

$$\Delta A = (m-1)A - \frac{1}{N} \left(\frac{dC}{d\omega} + \alpha \right), \quad (23)$$

$$\Delta \omega = \epsilon A. \quad (24)$$

Where $\alpha = \left(\frac{d\hat{C}}{d\omega} - \frac{dC}{d\omega} \right)$ denotes the error in the gradient update. As discussed in the main text, we can approximate this error as Gaussian noise with statistics $\langle \alpha \rangle = 0$ and $\langle \alpha^2 \rangle \approx N^2 F(\omega)/B$. Equations 21 and 22 match equations 23 and 24 if the step size $\epsilon = (\Delta t)^2/N$, and the momentum parameter $m = (1 - \lambda \Delta t/N)$. Finally we equate the noise by setting $\langle \alpha^2 \rangle/N^2 = \langle \eta^2 \rangle/(\Delta t)^2$, and solve for the noise scale g to obtain,

$$g \approx \frac{\Delta t N}{B \lambda} = \frac{(\Delta t)^2}{B(1-m)} \quad (25)$$

$$= \frac{\epsilon N}{B(1-m)}. \quad (26)$$

As observed in the main text, if we wish to keep the scale of random fluctuations constant, then we should scale the batch size $B \propto \epsilon N$. We also predict an additional scaling relation between the batch size and the momentum parameter, $B \propto 1/(1-m)$. Taken together, our results demonstrate that the influence of four different hyper-parameters (ϵ, m, N, B) can be summarized through a single effective coordinate g which measures the random fluctuation scale of the dynamics. One can also interpret $\epsilon_{eff} = \epsilon/(1-m)$ as the ‘‘effective learning rate’’.

D HOW TO TUNE THE BATCH SIZE

Here we propose a simple heuristic for tuning the batch size, learning rate and momentum parameter:

1. Set the learning rate to 0.1 and the momentum coefficient to 0.9. Run experiments at a range of batch sizes on a logarithmic scale, and identify the optimal batch size which maximizes the validation set accuracy.
2. Repeatedly increase the batch size by a factor of 3, while scaling the learning rate $\epsilon \propto B$, until the validation set accuracy starts to fall. Then repeatedly increase the batch size by a factor of 3, while scaling the momentum coefficient $(1-m) \propto 1/B$, until either the validation set accuracy falls or the batch size reaches the limits of your hardware.
3. Having identified the final learning rate and momentum parameter, retune the batch size on a linear scale in the local neighborhood of the current batch size.

We believe that this simple procedure will increase the test accuracy, reduce the cost of tuning hyper-parameters, and significantly reduce the final number of gradient updates required to train a model.