

A Perturbative Analysis of Stochastic Gradient Descent

Samuel C. Tenka

MIT, CSAIL
C O L I @ M I T . e d u

Abstract

We quantify how gradient noise shapes the dynamics of stochastic gradient descent (SGD) by Taylor-expanding with respect to the learning rate. We present a new diagram-based technique that permits re-summation of that series to convergent results on short timescales or near local minima. We physically interpret the resulting terms as entropic corrections to deterministic descent. Our theory predicts that gradient noise biases SGD’s trajectory toward low-curvature, low-noise regions of the loss landscape. We contrast our theory against popular approximations of SGD such as SDE and experimentally verify our theory’s predictions.


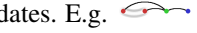
Introduction

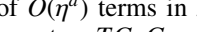
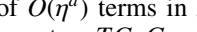

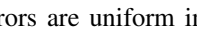
Gradient estimates, measured on minibatches and thus noisy, form the primary learning signal in deep learning. While users of deep learning benefit from the intuition that *stochastic gradient descent* (SGD) approximates deterministic descent (GD) (Bottou 1991; LeCun, Bengio, and Hinton 2015), SGD’s gradient noise alters training dynamics and testing losses (Goyal et al. 2018; Wu et al. 2020). We analyze these dynamics on short timescales or near minima. We apply our theory to find that **gradient noise biases learning** toward low-curvature, low-noise regions of the loss landscape.

Specifically, we study the expectation $\mathbb{E}[l(\theta_T)]$ over training sets of the post-training testing loss by Taylor expanding that loss in the learning rate η . By induction on T , the loss changes by $-TG\eta G + o(\eta)$ after T training steps, where $G = \mathbb{E}_x[\nabla l_x(\theta_0)]$ is the expectation over testing samples x of the gradient at initialization. This estimate is exact when $\nabla l_x(\theta)$ depends on neither x nor θ , i.e., for deterministic linear loss landscapes. We compute how noise and curvature correct this estimate.

A Taylor series analysis of SGD presents three challenges. *First*, the terms explode in variety. Even the main correction to $-TG\eta G$ represents the diverse ways that some past update may affect a future weight θ_t and thus a future update involving $\nabla l_{x_t}(\theta_t)$. That is, updates **interact**. *Second*, SGD’s gradient noise is correlated between timesteps: the same training sample reappears in each epoch. Such **finite-sample** effects complicate the simple induction that led to

$-TG\eta G$. *Third*, the series’ d th order truncation **diverges** as T grows. Indeed, on landscapes such as least squares linear regression, the loss grows exponentially with time for $\eta < 0$. So on no neighborhood of $\eta = 0$ does any Taylor truncation suffer an error uniform in T .

We address the three challenges by using diagrams such as  to organize and evaluate many terms at once, including correlation effects. Roughly, the analytical procedure we develop is this. We interpret each diagram as depicting a class of interactions or “**histories**” between updates. E.g.  depicts an update’s (red’s) double effect on a future update (green) that in turn affects the testing loss (blue). The rightmost (“root”) node always represents a post-training measurement. Then, up to $o(\eta^d)$ error, *the testing loss is a sum — over all histories of all diagrams with $\leq d$ edges — of certain diagram-dependent tensor expressions*.

More formally, each d -edged diagram represents a class of $O(\eta^d)$ terms in $\mathbb{E}[l(\theta_T)]$ ’s Taylor series. E.g. ’s terms sum to $-TG\eta G$, which diverges as T grows. However, we find that each diagram’s (e.g. ) divergence is countered by those of higher-order topologically related diagrams (e.g. , , \dots). These “**re-summed**” expressions’ errors are uniform in T for quadratic landscapes (with non-Gaussian gradient noise) and are provably finite and empirically small for convolutional landscapes.

Paradigmatic example

We illustrate our theory by deriving a flatness-seeking tendency of SGD while avoid the next section’s generality.

Notation and assumptions, I

We formalize the loss — suffered by a fixed architecture on a random datapoint — as a distribution \mathcal{D} over functions from a space \mathcal{M} of weights. The *testing loss* $l : \mathcal{M} \rightarrow \mathbb{R}$ is \mathcal{D} ’s mean. We write $\theta \in \mathcal{M}$, $l_x \sim \mathcal{D}$ for generic elements. We consider training sequences $(l_n : 0 \leq n < N) \sim \mathcal{D}^N$. We call n and l_n *training points*. Each initialization $\theta_0 \in \mathcal{M}$ then induces via SGD a distribution over trajectories $(\theta_t : 0 \leq t \leq T)$. Specifically, SGD runs T steps of η -steepest descent:

$$\theta_{t+1}^\mu := \theta_t^\mu - \sum_v \eta^{\mu\nu} \nabla_\nu l_{n_t}(\theta_t)$$

Each sequence $(n_t : kN \leq t < kN+N)$ is a permutation of $(n : 0 \leq n < N)$. Our Greek indices name components of θ, η, ∇ w.r.t. a fixed basis. We view η as a bilinear form so that the only type-correct expressions have geometric significance.

We heavily use the **gradient** $G_\mu = \mathbb{E}[\nabla_\mu l_x(\theta)]$, **hessian** $H_{\mu\nu} = \mathbb{E}[\nabla_\mu \nabla_\nu l_x(\theta)]$, **jerk** $J_{\mu\nu\xi} = \mathbb{E}[\nabla_\mu \nabla_\nu \nabla_\xi l_x(\theta)]$, **covariance** $C_{\mu\nu} = \mathbb{E}[(\nabla l_x(\theta) - G)^{\otimes 2}]_{\mu\nu}$, and **skew** $S_{\mu\nu\xi} = \mathbb{E}[(\nabla l_x(\theta) - G)^{\otimes 3}]_{\mu\nu\xi}$, typically evaluated at initialization $(\theta = \theta_0)$. So G, H, J, C, S respectively have 1, 2, 3, 2, 3 many axes, each of which transforms under change-of-basis like a covector.¹

To illustrate our notation we quote Nesterov (2004), §2.1:



Prop 0. $G = \nabla l(\theta_0)$ controls the loss to leading order. Precisely, $\mathbb{E}[l(\theta_T) - l(\theta_0)]^\mu = -T \sum_{\mu\nu} G_\mu \eta^{\mu\nu} G_\nu + o(\eta^1)$.

Our work identifies how noise and curvature correct Prop 0.

Informal Statement of Result

Definition 1. A **diagram** is a rooted tree equipped with an equivalence relation on (i.e. a partition of) its non-root nodes. Convention: we orient the tree left-to-right so that children precede parents; the root is thus rightmost. We draw the partition with fuzzy outlines. **Colors** in diagrams lack formal meaning but help us refer to diagram parts. \diamond

Valid diagrams include , , , . And

 depicts the same diagram as .

Definition 2. A **history** of a diagram is an assignment of non-root nodes to (n, t) pairs such that: the n th training point participates in the t th batch; parents' t s strictly exceed their children's t s; and any two nodes' n s are equal if the nodes are in the same part of the partition. \diamond

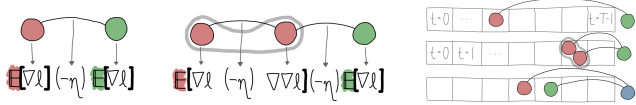
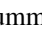
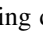
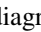
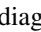


Figure 1: **Left:** Diagrams evaluate to tensor expressions: a degree- d node gives a d th derivative; an edge gives a $(-\eta)$; expectation brackets enclose each fuzzy group. **Right:** Each diagram depicts a class of histories; here is an example history for each of , , . Summing diagram values over all possible histories gives the expected testing loss.

Theorem (informal). *The expected testing loss $\mathbb{E}[l(\theta_T)]$ is a sum over all diagram histories, where each diagram is evaluated by a procedure as in Figure 1. If we sum only diagrams with at most d edges, we suffer only $o(\eta^d)$ error.*

For example, the only diagram with 1 edge is . This diagram evaluates to $-G\eta G$ and, since its red node can be any update while its green node represents testing, describes

¹Gradients (covectors, dollars-per-mile) and displacements (vectors, miles) have geometrically distinct types. We respect this distinction throughout; that's why our η is a tensor. (Misner, Thorne, and Wheeler 1973) (§2.5) visualizes this geometry; (McCullagh 1987) (§1.4) discusses vectors vs covectors in statistics.

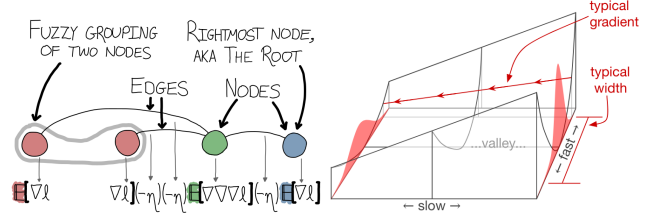
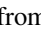



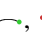
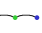



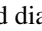
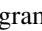
Figure 2: **Left:** a diagram consists of nodes, edges, and fuzzy groupings that dictate the diagram's corresponding tensor expression. **Right:** Gradient noise pushes SGD toward minima flat w.r.t. C . A 2D loss near a valley of minima. Red densities show typical θ s, perturbed by noise, in two slices of the valley. The hessian changes across the valley: $J \neq 0$.


T many histories. We thus recover Prop 0's (un-resummed value) $-TG\eta G$. The re-summed value, which turns out to be $-TG \frac{1 - \exp(-T\eta H)}{H} G$, tempers Prop 0's large- T behavior by including contributions from , , , \dots .

SGD descends on a C -smoothed landscape

We sketched above how a sum over diagrams gives $\mathbb{E}[l(\theta_T)]$ for the testing loss l . We may likewise compute $\mathbb{E}[s(\theta_T)]$ for other $s : \mathcal{M} \rightarrow \mathbb{R}$; we need only replace the l corresponding to each diagram's root by s . For instance,  ordinarily evaluates to $-\eta^3 C J(\nabla l)$ (see Figure 2), but more generally evaluates to $-\eta^3 C J(\nabla s)$.

Now let's study the displacement $\mathbb{E}[\theta_T - \theta_0]$ of one-epoch, batch-size-one SGD initialized at a testing loss minimum. To do this, we let $s(\theta) = w \cdot \theta$ be linear. Then $\mathbb{E}[s(\theta_T)]$ reveals the displacement's w -component. Since s 's higher derivatives vanish, any diagram whose root has degree > 1 will evaluate to zero. So only diagrams with degree-1 root are relevant to computing displacements (e.g. , , ).

At a minimum, $G = 0$ so we may ignore yet more diagrams: those with a factor of $\mathbb{E}[\nabla l]$, i.e., those with a degree-1 non-root node that's not fuzzily grouped. So we rule out , , \dots . Only one fewest-edged diagram remains:

, which above we evaluated as $-\eta^3 (C \nabla H) \cdot (\nabla s)$. So to leading order the displacement is some T -dependent number of histories times $-\eta^3 (C \nabla H)$.

Since $-\eta^3 C \nabla H$ points toward small H , **SGD moves toward flat minima** (Figure 2). This effect scales with C . Re-summation more precisely quantifies this 'entropic force'.²³

Corollary 1. *Start SGD at a minimum of l with $H > 0$, $N = T$. Use an eigenbasis of $K = \eta H$. For any T and with*

²Thermal systems tend toward disorder as if pushed by an 'entropic force'. So arises the tension of rubber bands: their polymers can wreathe in many ways but be straight in only one. Such 'forces' characteristically scale with temperature (the noise intensity C).

³Our result ($T \gg 1$) is $\Theta(\eta^2)$; (Yaida 2019a)'s ($T = 2$) is $\Theta(\eta^3)$. We integrate noise over time, amplifying C 's effect.

$\mathcal{P}_T(s) = (1 - \exp(-Ts))/s$, the expected displacement is

$$-\sum_{\substack{\mu\nu \\ \delta\pi p}} C_{\mu\delta} \mathcal{P}_T(K_{\mu\mu} + K_{\delta\delta}) \eta^{\mu\nu} \eta^{\delta\pi} J_{\nu\pi\xi} \mathcal{P}_T(K_{\xi\xi}) \eta^{\xi\rho} / 2 + o(\eta^3)$$

For instance, consider a 1D valley of near-minima wherein ηH has spectrum $\lambda_0 \ll 1/T \ll \lambda_1 \leq \dots$ for eigenvectors v_i . Let's ignore diffusion along the valley: $(\eta C)v_0 = 0$. Then $\mathcal{P}_T(\lambda_0) \approx T$ and every $C_{ij}\mathcal{P}_T(\lambda_i + \lambda_j)$ is $O(T^0)$. So $\mathbb{E}[\theta_T - \theta_0]$ scales linearly with T . We thus expect SGD to move with velocity $\approx -\eta^2 C \nabla H / 2$ per timestep toward flat minima. Observe that $\nabla(\mathbb{G}_C \star l) = \nabla l + C \nabla H / 2 + o(C)$, where $\mathbb{G}_C \star$ denotes convolution with a *fixed* centered C -shaped Gaussian; we conclude with the intuition that *SGD descends on a C -smoothed landscape that changes as C does*.

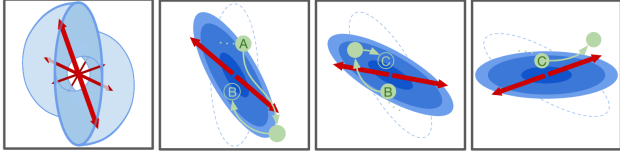


Figure 3: ■■■: HELIX is defined on an $\mathcal{M} = \mathbb{R}^3$ extending into the page. A helical level surface S (blue) of l winds around a 1D valley of minima orthogonal to the page. Gradient noise (red bi-arrows), parallel to the page, twists out of phase with l . ■■■: SGD's trajectory (green) over cross sections of HELIX descend progressively into the page. Dotted curves help compare adjacent panes. Gradient noise kicks θ from A; θ then falls (■■■) to B in ■■■. At C, noise kicks θ uphill (■■■); θ thus never settles and the descent persists.

To test Corollary 1's C -dependence, we construct a landscape, HELIX, on whose valley of global minima C varies (Figure 3). As in Rock-Paper-Scissors, each θ has a neighbor that is more attractive (flatter) w.r.t. $C(\theta)$. This induces eternal motion into the page despite HELIX's discrete translation symmetry. Corollary 1 predicts a velocity of $+\eta^2/6$ per timestep, while (Chaudhari and Soatto 2018)'s SDE-based analysis predicts¹ a constant velocity of 0 (Figure 6). Wrapping HELIX in a loop makes SGD circulate. This is possible because $C \nabla H$, unlike $\nabla(CH)$, is not a total derivative. In avoiding (Wei and Schwab 2019)'s constant- C assumption, we find that SGD's velocity field can have curl.

Perturbative theory of SGD

After explaining the diagram-based Taylor series approach to SGD dynamics, we state our main result: that diagram-based analysis is correct.

Notation and assumptions, II

REGULARITY CONDITIONS — We assume the following throughout. **Derivative bounds:** there are compact sets $(K_k : k \geq 0)$ so that $\nabla^k l_x(\theta) \in K_k$ for all θ, l_x . Here $\nabla^k l_x(\theta)$ is a k th derivative, a k -axis tensor. **Analytic moments:** any polynomial p of l_x and its higher derivatives induces a random variable so that $\mathbb{E}[p] : \mathcal{M} \rightarrow \mathbb{R}$ (exists and) is analytic in θ ;

¹For, HELIX' velocity is η -perpendicular to the image of $(\eta C)_v^\mu$.

moreover, $\mathbb{E}[p]$'s radii of convergence are strictly bounded from 0, even as θ varies. Consequently, $\nabla \mathbb{E}[p] = \mathbb{E}[\nabla p]$.

BATCHES AND EPOCHS — Generalizing the previous section, our theory describes SGD with any number N of training points, \mathbf{T} of updates, and \mathbf{B} of points per batch. SGD runs T updates (hence $\mathbf{E} = TB/N$ epochs or $\mathbf{M} = T/N$ updates per training point) of the form

$$\theta_{t+1}^\mu := \theta_t^\mu - \sum_v \eta^{\mu\nu} \nabla_v \sum_{n \in \mathcal{B}_t} l_n(\theta) / B$$

where in each epoch we sample the t th batch \mathcal{B}_t without replacement from the training sequence.

COMPARING TENSORS — A (potentially tensor) quantity q vanishes to order η^d when for some homogeneous degree- d polynomial p $\lim_{\eta \rightarrow 0} q/p(\eta) = 0$; we then say $q \in o(\eta^d)$. We write $A \leq B$ ($<$) for symmetric bilinear forms A, B when $A(v, v) \leq B(v, v)$ ($<$) for all $v \neq 0$.

CONTINUOUS TIME — Ordinary and stochastic differential equations (ODE, SDE) are popular models of SGD (Liao et al. 2018; Barrett and Dherin 2021). They correspond to continuous-time limits of large-training-set SGD with independent noise ($E = B = 1, N = T = kT_0, \eta = \eta_0/k, k \rightarrow \infty$) ODE descends on a noiseless version of the landscape $l_x(\theta)$, while SDE descends on a version whose gradient noise is independent gaussian of shape $C(\theta) = \mathbb{E}[\nabla l_x \nabla l_x] - \mathbb{E}[\nabla l_x] \mathbb{E}[\nabla l_x]$, scaled as in a Wiener process:

$$l_x^{\text{ODE}}(\theta) - l(\theta) = 0 \quad l_x^{\text{SDE}}(\theta) - l(\theta) \sim \mathcal{N}(0, k \cdot \theta^T C(\theta) \theta)$$

In this paper we shall take k to be large but finite and express results about ODE, SDE with error terms such as $o(1/k)$. We emphasize that the constants in these little- o s are permitted to depend on the loss landscape and optimization parameters such as η, T . It is physical intuition and experiment that determine when such error terms are negligible.

Diagrams arise from and organize Taylor series

Structure of the Taylor Expansion We discuss how to analyze SGD by expanding in powers of η . We begin by proving Prop 0 (c.f. (Nesterov 2004; Roberts 2018)).

Proof. By our gradient bound assumption: $\theta_T - \theta_0$ is $O(\eta^1)$. We claim that $(\theta_T - \theta_0)^\mu = -\sum_t \sum_v \eta^{\mu\nu} \nabla_v l_{n_t}(\theta_0) + o(\eta^1)$. The claim holds when $T = 0$. Say the claim holds for \tilde{T} -step SGD with $T = \tilde{T} + 1$. The displacement $(\theta_T - \theta_{\tilde{T}})^\mu$ is:

$$\begin{aligned} & -\sum_v \eta^{\mu\nu} \nabla_v l_{n_{\tilde{T}}}(\theta_{\tilde{T}}) \\ &= -\sum_v \eta^{\mu\nu} \nabla_v \left(l_{n_{\tilde{T}}}(\theta_0) + \sum_\xi \nabla_\xi l_{n_{\tilde{T}}}(\theta_0) (\theta_{\tilde{T}} - \theta_0)^\xi + o(\theta_{\tilde{T}} - \theta_0) \right) \\ &= -\sum_v \eta^{\mu\nu} \nabla_v \left(l_{n_{\tilde{T}}}(\theta_0) + \nabla l_{n_{\tilde{T}}}(\theta_0) \cdot O(\eta^1) + o(O(\eta^1)) \right) \\ &= -\sum_v \eta^{\mu\nu} \nabla_v l_{n_{\tilde{T}}}(\theta_0) + o(\eta^1) \end{aligned}$$

Applying the induction hypothesis proves the claim. We plug the claim into l 's Taylor series:

$$\begin{aligned} \mathbb{E}[l(\theta_T) - l(\theta_0)] &= \sum_\mu \nabla_\mu l(\theta_0) \mathbb{E}[\theta_T - \theta_0]^\mu + \mathbb{E}[o(\theta_T - \theta_0)] \\ &= \sum_\mu \nabla_\mu l(\theta_0) (-T \eta G + o(\eta^1)) + o(O(\eta^1)) \\ &= -\sum_{\mu\nu} T G_{\mu\nu} \eta^{\mu\nu} G_\nu + o(\eta^1) \end{aligned}$$

Indeed, due our assumption of analytic moments, the above expectations of $o(\eta^1)$ terms are still $o(\eta^1)$. \square

The above proof gives an order-1 result. At higher order, higher derivatives correct \dots and higher moments augment \dots . Whereas above the displacement is a sum over \tilde{T} 's of \dots 's, due to \dots 's corrections the displacement at higher order is a sum over *tuples* of times with summands such as $\nabla l_{n_T} \nabla l_{n_{\tilde{T}}}$ instead of ∇l_{n_T} . When we then take expectations of \dots to evaluate \dots as \dots , some summands (e.g. $\mathbb{E}[\nabla l_5 \nabla l_2] = \mathbb{E}[\nabla \nabla l_5] \mathbb{E}[\nabla l_2]$) are uncorrelated and thus factor; others (e.g. $\mathbb{E}[\nabla \nabla l_5 \nabla l_5]$) do not. This is how ∇l_x 's higher cumulants such as C, S appear in our analysis.

Overall, a general summand in \dots has the following form (evaluated at $\theta = \theta_0$):

$$\sum_{\text{all Greek indices}} \left(\prod_{j \in J} \eta^{\mu_j \nu_j} \right) \left(\prod_{i \in I} \left(\prod_{k \in K_i} \nabla_{\xi_{i,k}} \right) l_{x_i} \right) \left(\prod_{k \in K_*} \nabla_{\xi_{*,k}} \right) l$$

We represent such a summand's expectation as a diagram with edges indexed by $j \in J$, nodes indexed by $i \in I \sqcup \{\star\}$, an edge j incident to a node i when $\{\xi_{i,k} : k \in K_i\}$ meets $\{\mu_j, \nu_j\}$, and nodes i, i' grouped in the partition when $x_i = x_{i'}$.

Definition 3. A diagram D 's *un-resummed value* (**uvalue**) is a product with one factor of l_x 's d th derivative for each degree- d node, grouped under cumulant symbols \mathbb{C} (think: expectation symbols \mathbb{E})¹ per D 's fuzzy groups, and tensor-contracted via a factor $\eta^{\mu\nu}$ for each edge. \diamond

E.g. $\text{uvalue}(\text{---})$ is:

$$\sum_{\mu\nu\xi\delta\sigma\rho} \eta^{\mu\xi} \eta^{\nu\delta} \eta^{\sigma\rho} \mathbb{C}[(\nabla_\mu l_x) \cdot (\nabla_\nu l_x)] \mathbb{C}[(\nabla_\xi \nabla_\delta \nabla_\sigma l_x)] \mathbb{C}[(\nabla_\rho l_x)]$$

There are dozens of small diagrams. In many analyses, only a few diagrams are relevant. Examples: for fixed T , **to order d we may neglect diagrams with more than d edges**; if $E = B = 1$ (§), each diagram with an ancestor-descendant pair in the same part contributes zero to $\mathbb{E}[l(\theta_T)]$; for θ_0 a minimum of l , all diagrams vanish that contain a leaf node not fuzzily grouped.

SGD as a Sum over Histories Having expressed terms in $\mathbb{E}[l(\theta_T)]$'s Taylor expansion as uvalues of diagrams, we seek the coefficient for each uvalue. Intuitively, a diagram represents a process (as in Figure 9) and a diagram's contribution scales with the number of ways that process may occur. Specifically, the Key Lemma in § establishes that a diagram's coefficient in the expansion is the number of its *histories* (weighted by symmetry factors to counter overcounting), where we define *histories* below with respect to given SGD hyperparameters N, E, B .²

E.g. --- has just one non-root node. It has as many histories as there are (n, t) cells where n participates in the t th update, i.e., $B \cdot T$ histories. Since --- is the only 1-edged diagram, it gives the full η^1 contribution to final testing loss.

Diagrams streamline analysis of SGD because it is in practice straightforward to count a diagram's histories. Also,

¹Inconsequential technicality: uvalues are products of *cumulants* such as C , not of un-centered moments such as $GG + C$. The symbol $\mathbb{C}[a]$ gives a 's mean; $\mathbb{C}[a \cdot b]$, a, b 's covariance; we center higher cumulants $\mathbb{C}[\prod_i a_i]$ with respect to lower cumulants.

²and w.r.t. a deterministic selector of the t th batch. One may take expectations over such algorithms — §.

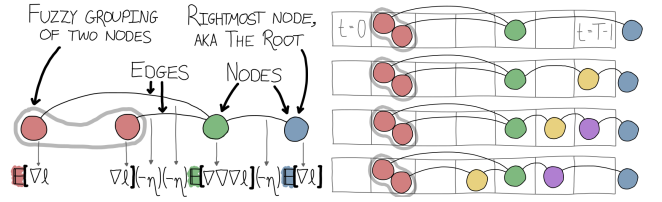


Figure 4: --- : a diagram consists of nodes, edges, and fuzzy groupings that dictate the diagram's corresponding tensor expression. --- shows one of a diagram's histories. Re-summation evaluates a whole class of topologically related histories at once (---).

the topology of diagrams has dynamical significance: the t^{dT-p} -th order correction³ to the ODE approximation of SGD is given by diagrams with d edges and $d + 1 - p$ many fuzzy groups (counting each node not fuzzily grouped as its own group). Likewise, if we seek to isolate the effect, say, of C or H or S or J , we may consider only those diagrams that contain the corresponding subgraph.

Re-summation cures large- T divergences So far, diagrams have been a convenient but dispensable book-keeping tool; so far, §'s polynomial divergence remains in (3). We now show how diagrams enable us to tame this divergence.

Let us collect similar diagrams, where our notion of 'similar' permits chains to grow or shrink (see Definition 4). We obtain lists (each conveniently represented by its smallest member) such as --- , --- , --- , \dots . We will express in closed form the total contribution to (3) of all diagrams in such a list. The idea is that the uvalues of chains are powers of Hessians — e.g. $\text{uvalue}(\text{---}) = (\eta)^4 GH^3 G$ — so we may sum over chain lengths via geometric series.

We define an embedded diagram's *resummed value* or **rvalue** as we defined the uvalue, except that we use $(I - \eta H)^{\Delta t - 1} \eta$ instead of η to contract a pair of tensors embedded Δt timesteps apart. For example, take Figure 7's history of --- (topmost of four). The associated uvalue is $\sum_{\mu\nu} G_\mu \eta^{\mu\nu} G_\nu$: a G for each degree-one node and an η for each edge. By contrast, the associated rvalue is $\sum_{\mu\nu} G_\mu ((I - \eta H)^{11-1} \eta)^{\mu\nu} G_\nu$ since the edge spans 11 timesteps. Distributing out this expression reveals uvalues for histories of --- , etc.

Definition 4. A **link** is a degree-2 non-root node that is not fuzzily grouped. E.g. --- has one link (green). To *reduce* at a link, we replace the link by a black edge connecting the link's two neighbors. E.g. $\text{---} \rightsquigarrow \text{---}$. Reduction generates an equivalence relation on diagrams. Each equivalence class contains exactly one **linkless** diagram. \diamond

Main result

Our main result is abstract but specializes to several concrete corollaries.

³We compare ODE integrated to time t to T steps of SGD with $\eta = \eta_{*,t}/T$ and $E = B = 1$, and we assume $p \neq 0$.

diagram	#embed.s	interpretation
	T	naïve descent
	$\binom{T+1}{2}$	θ -dependent loss
	T	gradient noise
	0	correl'd batches

Theorem 1. $\forall T : \exists \eta_0 > 0 : \forall 0 \leq \eta < \eta_0$: the final testing loss is a sum over linkless diagrams:

$$\mathbb{E}[l(\theta_T)] = \sum_{D \text{ a linkless diagram}} \sum_{f \text{ an embed-ding of } D} \frac{1}{|\text{Aut}_f(D)|} \frac{\text{rvalue}_f(D)}{(-B)^{|\text{edges}(D)|}}$$

Here, $|\text{Aut}_f(D)|$ counts the graph automorphisms of D that preserve f . (Typically $|\text{Aut}_f(D)| = 1$.)

Remark 1. A diagram with d edges scales as $O(\eta^d)$, so the Theorem expresses a series in η . In practice, we truncate to small d (thus focusing on few-edged diagrams) and we replace sums over histories by integrals over t ; $(I - \eta H)^t$, by $\exp(-\eta H t)$, thus reducing to a routine integration of exponentials at the cost an error factor $1 + o(\eta)$. \diamond

Remark 2. Theorem 1 gives us the expectation of the testing loss. Straightforward variations of the theorem permit us to compute variances instead of expectations, training statistics instead of testing statistics, and weight displacements instead of losses. See §appendix:solve-variants. \diamond

Theorem 2. For constant- $M = P$ or constant- $N = P$ SGD: $\forall \theta_*, P : (G(\theta) = 0 \wedge H(\theta) > 0) \implies \exists U \ni \theta_* \text{ open} : \forall \theta_0 \in U : \text{Theorem 1's } d\text{-th-order truncation converges as } T \rightarrow \infty$.

QUADRATIC EXACTNESS!

Remark 3. Thm 2 claims only that a limit $\lim_{T \rightarrow \infty} L_d(T, \eta)$ of the d -th-order truncation $L_d(T, \eta)$ exists. It does not compare $\lim_{d \rightarrow \infty} \lim_{T \rightarrow \infty} L_d(T, \eta)$ with $\lim_{T \rightarrow \infty} \lim_{d \rightarrow \infty} L_d(T, \eta)$ (however, we note that we have not in practice observed pathologies of non-commuting limits). Our theory suggests but does not guarantee that when d, T are large (and finite), then computation of $L_d(T, \eta)$ by Taylor methods gives insight into SGD's behavior. It is by empirical tests and physical intuition that we decide whether in a given situation our theory's little- o error terms may be ignored.

Example Computation A We improve on Prop 0 for SGD with $E = B = 1$. The one 1-edged diagram () embeds in T ways (one for each timestep) and contributes (let $K_v^\mu = \sum_\xi \eta^{\mu\xi} H_{\xi v}$):

$$\sum_{0 \leq t < T} \sum_{\mu\nu} G_\mu \left[(I - K)^{T-t-1} \eta \right]^{\mu\nu} G_\nu = \sum_{\mu\nu} G_\mu \left[\frac{I - K^T}{I - K} \eta \right]^{\mu\nu} G_\nu$$

to the loss. This is the re-summed Prop.

Example Computation B There are three 2-edged linkless diagrams (see table for intuitive interpretations). The two diagrams involving higher cumulants (i.e., gray outlines) give the leading order effect of gradient noise. Since $E = 1$, has no histories; so only contributes. Up

to a $1 + o(\eta)$ factor, its rvalues sum to (a sum over all indices of):

$$\int_t C_{\mu\nu} [\exp(-(T-t)(K \otimes I + I \otimes K))]_{\xi\phi}^{\mu\nu} \eta^{\xi\pi} \eta^{\phi\rho} H_{\pi\rho} \\ = C_{\mu\nu} \left[\frac{I - \exp(-T(K \otimes I + I \otimes K))}{(K \otimes I + I \otimes K)} \right]_{\xi\phi}^{\mu\nu} \eta^{\xi\pi} \eta^{\phi\rho} H_{\pi\rho}$$

We used Remark 1 to approximate $(I - K)_\xi^\mu (I - K)_\phi^\nu$ by $\exp(-K \otimes I - I \otimes K)_{\xi\phi}^{\mu\nu}$. The above is the leading contribution of the gradient covariance C to SGD's final testing loss. We have derived an $E = B = 1$ variant of Corollary 6's $E = T$, $B = N$ result (§).

Consequences of the theory

By Cor.s 2 and 7, gradient noise repels SGD. By Cor. 5, SGD senses changes in H more than SDE; in fact, (Cor. 1) SGD seeks small- H weights. Cor. 6 relates C and H to overfitting.

Gradient noise repels SGD

Physical intuition suggests that noise repels SGD: if two neighboring regions of weight space have high and low levels of gradient noise, respectively, then the rate at which θ jumps from the former to the latter exceeds the opposite rate. There is thus a net movement toward regions of small C .¹ Our theory makes this precise; θ drifts in the direction $-\nabla C$, and the effect is weaker when gradient noise is averaged out by large batch sizes:

Corollary 2 (Computed from). SGD with $E = B = 1$ avoids high- C regions more than GD: $\mathbb{E}[\theta_{GD} - \theta_{SGD}]^p = T \cdot \frac{N-1}{4N} \sum_{\mu\nu\xi} \eta^{\mu\rho} \eta^{\nu\xi} \nabla_\mu C_{\nu\xi} + o(\eta^2)$.

(Roberts 2019) obtained a version of this Corollary with a nearly equal error of $O(\eta^2/N) \vee o(\eta^2)$. The Corollary's proof implies that if \hat{l}_c is a smooth unbiased estimator of $\frac{N-1}{4N} C_\nu^\nu$, then GD on $l + \hat{l}_c$ has an expected testing loss that agrees with SGD's to order η^2 . We call this method **GDC**.

An analogous form of averaging occurs over multiple epochs.

Time discretization penalizes sloped regions

The following corollary recovers (Barrett and Dherin 2021)'s main dynamical result:

Corollary 3 (). On a noiseless landscape, SGD prefers small- G^2 regions more than ODE: $\mathbb{E}[\theta_{SGD} - \theta_{ODE}]^p = -\frac{T}{4} \sum_{\mu\nu} \eta^{\xi\rho} \eta^{\mu\nu} \nabla_\xi (G_\mu G_\nu) + o(\eta^2) + o(1/k)$.

Due to time discretization, in the presence of curvature SGD's response to gradient noise 'overshoots' more than SDE. The following corollary makes this precise and separates SDE from SGD, even on landscapes obeying SDE's assumption of gaussian noise:

Corollary 4 (). The covariance of gradient noise contributes $\frac{T}{2} \sum_{\mu\nu\xi\phi} C_{\mu\nu} \eta^{\mu\xi} \eta^{\nu\phi} H_{\xi\phi} + o(\eta^2) + o(1/k)$ to $E = B = 1$ SGD's final testing loss excess over SDE's.

¹This is the same mechanism by which sand on a vibrating plate accumulates in quiet regions (Chladni 1787). We thus dub the SGD phenomenon the Chladni drift.

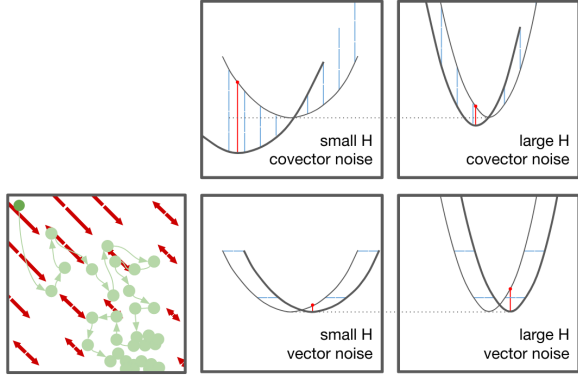



Figure 5: Left. Chladni drift on $\mathcal{M} = \mathbb{R}^2$. Red bi-arrows depict $C(\theta)$'s major axis. SGD updates (green) tend toward small C . Both curvature and noise affect overfitting. In each pane, the \leftrightarrow axis represents weight space and the \updownarrow axis represents loss. Noise (blue) transforms the testing loss (thin curve) into the observed loss (thick curve). Red dots mark the testing loss at the arg-min of the observed loss. $\square\square$: covector-perturbed landscapes favor large H s. $\square\square$: vector-perturbed landscapes favor small H s.

Jerk distinguishes SDE and SGD

SDE differs from SGD in ways beyond time-discretization effects. For instance, the inter-epoch noise correlations in multi-epoch SGD measurably affect SGD's final testing loss (Corollary 7), but SDE assumes uncorrelated gradient updates. Even if we restrict to single-epoch SGD, non-Gaussian noise lead SGD and SDE to respond differently to changes in curvature:


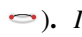
Corollary 5 (). For $E = B = 1$ SGD and up to error $o(\eta^3) + o(1/k)$: the skewness of gradient noise contributes

$$-\frac{\eta^3}{3!} \sum_{\mu\nu\lambda} S_{\mu\nu\lambda} \frac{1 - \exp(-T\eta(H_{\mu\mu} + H_{\nu\nu} + H_{\lambda\lambda}))}{\eta(H_{\mu\mu} + H_{\nu\nu} + H_{\lambda\lambda})} J_{\mu\nu\lambda}$$

to the excess final testing loss over SDE (in an eigenbasis of ηH).

Both flat and sharp minima overfit less

Intuitively, sharp minima are robust to slight changes in the average *gradient* and flat minima are robust to slight *displacements* in weight space (Figure 5 $\square\square$). However, as SGD by definition equates displacements with gradients, it may be unclear how to reason about overfitting in the presence of curvature. Our theory accounts for the implicit regularization of fixed- T descent and shows that both effects play a role. In fact, by routine calculus on Corollary 6, overfitting is maximized for medium minima with curvature $H \sim (\eta T)^{-1}$.

Corollary 6 (from  , ). Initialize GD at a non-degenerate test minimum θ_* . The overfitting (testing loss minus $l(\theta_*)$) and generalization gap (testing minus training

loss) due to training are:

$$\sum_{\mu\nu\rho\lambda} (C/(2NH))_{\mu\nu}^{\rho\lambda} \left((I - \exp(-\eta TH))^{\otimes 2} \right)_{\rho\lambda}^{\mu\nu} + o(\eta^2)$$

and

$$\sum_{\mu\nu\rho\lambda} (C/(2NH))_{\mu\nu}^{\mu\lambda} (I - \exp(-\eta TH))_{\lambda}^{\nu} + o(\eta)$$

The generalization gap tends to $C_{\mu\nu}(H^{-1})^{\mu\nu}/N$ as $T \rightarrow \infty$. For maximum likelihood (ML) estimation in well-specified models near the “true” minimum, $C = H$ is the Fisher metric, so we recover the AIC: (model dimension)/ N . Unlike AIC, our more general expression is descendably smooth, may be used with MAP or ELBO tasks instead of just ML, and does not assume a well-specified model.

Experiments

Our theory does not control our Taylor series' rates of convergence. We thus test our theory by experiment. We perceive support for our theory in drastic rejections of the null hypothesis. For instance, in Figure 3, (Chaudhari and Soatto 2018) predict a velocity of 0 while we predict a velocity of $\eta^2/6$. Likewise, published intuitions (§) suggest that Figure 6 $\square\square$ overfitting (test loss minus test minimum) is monotonic in a landscape's hessian, whereas we do not. Here, I bars, + signs, and shaded regions all mark 95% confidence intervals based on the standard error of the mean. § describes neural architectures, artificial landscapes, sample sizes, and further plots.

Discussion

Related work

Kiefer and Wolfowitz (1952) united gradient descent (Cauchy 1847) with stochastic approximation (Robbins and Monro 1951) to invent SGD. Since the development of back-propagation (Werbos 1974), SGD has been used to train connectionist models, e.g. neural networks (Bottou 1991), recently to remarkable success (LeCun, Bengio, and Hinton 2015). Several research programs treat overfitting of SGD-trained networks. (Neyshabur et al. 2017a). Bartlett, Foster, and Telgarsky (2017) controls the Rademacher complexity of deep hypothesis classes, leading to optimizer-agnostic generalization bounds. Yet SGD-trained networks generalize despite their ability to shatter large sets (Zhang et al. 2017), so generalization must arise from not only architecture but also optimization (Neyshabur et al. 2017b).

Some analyses of implicit regularization use a Langevin or SDE approximation (e.g. Chaudhari and Soatto (2018); Zhu et al. (2019)), but, per Yaida (2019b), such continuous-time or uncorrelated-noise analyses treat SGD noise incorrectly. We avoid these pitfalls by Taylor expanding around $\eta = 0$ as in Roberts (2018). Unlike that work, we generalize beyond order η^1 and $T = 2$. Our interpretation of the resulting terms offers a new qualitative picture of SGD as a superposition of simpler information-flow processes. Other research focuses on *double descent* and suggests that some highly overparameterized models share implicit regularization properties with linear least-squares models (Belkin

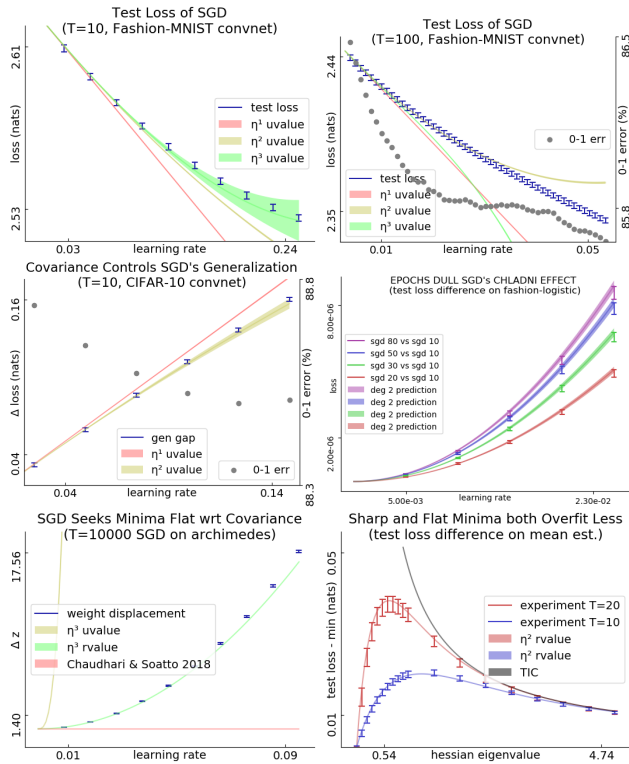


Figure 6: ■■■: Fashion-MNIST convnet’s testing loss vs learning rate. For all initializations tested (1 shown, 11 unshown), the order 3 prediction agrees with experiment through $\eta T \approx 10^0$, corresponding to a decrease in 0-1 error of $\approx 10^{-3}$. □□□: Fashion-MNIST convnet’s testing loss. For large ηT , our predictions break down. Here, the order 3 prediction holds until the 0-1 error improves by $5 \cdot 10^{-3}$. Beyond this, 2nd order agreement with experiment is coincidental. □□□: CIFAR-10 convnet generalization gaps. For all initializations tested (1 shown, 11 unshown), the degree-2 prediction agrees with experiment through $\eta T \approx 5 \cdot 10^{-1}$. □□□: For MEAN ESTIMATION with fixed C and a range of H s, initialized at the truth, the testing losses after fixed- T GD are smallest for very sharp and very flat H . Near $H = 0$, our predictions improve on AIC, TIC (Dixon and Ward 2018). **Right:** Predictions near minima excel for large ηT .

et al. 2019), for example by bounding log-determinants (and hence the effective dimensions) of feature matrices and weight spaces (Mei and Montanari 2020).¹ Our work reveals new dynamics toward and within valleys of minima, dynamics that may also reduce the effective dimension of model space. However, our focus on the structure of gradient noise may be overspecific, since recent work finds that GD and SGD may both converge to the same global minima (Zou et al. 2020) or that noise covariance but not higher moments are relevant to regularization (Wu et al. 2020).

Our predictions are vacuous for large η . Other work

¹(Mei and Montanari 2020)’s eq. 75 bounds a log-determinant defined in eq. 61 of a transformed feature matrix. C.f. to linear Representer Theorems (Mohri, Rostamizadeh, and Talwalkar 2018).

treats large- η learning phenomenologically, whether by finding empirical correlates of the generalization gap (Liao et al. 2018), by showing that *flat* minima generalize (Hof-fer, Hubara, and Soudry 2017; Keskar et al. 2017; Wang et al. 2018), or by showing that *sharp* minima generalize (Stein 1956; Dinh et al. 2017; Wu, Ma, and E 2018). SGD’s implicit regularization mediates between these seemingly clashing intuitions (§??).

Prior work analyzes SGD perturbatively: (Dyer and Gur-Ari 2019) perturb in inverse network width, using ‘t Hooft diagrams to correct the Gaussian Process approximation for specific nets. Perturbing to order η^2 , (Chaudhari and Soatto 2018) and (Li, Tai, and E 2017) assume uncorrelated Gaussian noise while (Barrett and Dherin 2021) compares GD to ODE. By contrast, we use Penrose diagrams (Penrose 1971) to compute testing losses and to compare to ODE and SDE to *arbitrary order* in η . We allow correlated, non-Gaussian noise and thus *any* smooth architecture.

Conclusion

This paper studies stochastic optimization on short timescales or near minima. Generalizing (Liao et al. 2018; Wei and Schwab 2019; Zhu et al. 2019; Barrett and Dherin 2021), we model correlated, non-gaussian, non-isotropic, non-constant gradient noise and find qualitative differences in dynamics. For example, we construct a non-pathological loss landscape on which SGD’s trajectory *ascends*. We verify our theory on convolutional CIFAR-10 and Fashion-MNIST landscapes. Corollaries 1 and 6 together show that SGD avoids curvature and noise, which to leading order control generalization.

Our theory offers a new physics-inspired perspective of SGD as a superposition of concurrent processes in which data influence weights. Notating such processes with such diagrams, we show how to compute the effect of each process and that summing the finitely many processes with d or fewer edges suffices to answer dynamical questions to error $o(\eta^d)$. We thus factor the analysis of SGD into the analyses of individual processes, a technique that may power future theoretical inquiries.

Since our predictions depend only on loss data near initialization, they break down after the weight moves far from initialization. Our theory thus best applies to small-movement contexts, whether for long times (large ηT) near an isolated minimum or for short times (small ηT) in general. Thus, the theory might aid future analysis of fine-tuners such as (Finn, Abbeel, and Levine 2017)’s MAML.

Much as meteorologists understand the dance of warm and cold fronts despite long-term forecasting’s intractability, we quantify how curvature and noise contribute to counter-intuitive dynamics governing each short-term interval of SGD’s trajectory. Equipped with our theory, users of deep learning may refine intuitions — e.g. that SGD descends on the training loss — to account for noise.

moo

moo

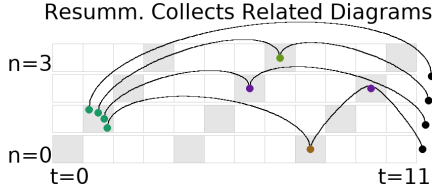


Figure 7: **Resummation propagates information damped by curvature.** Each resummed value (here, for \rightarrow) represents many un-resummed values, four shown here, each modulated by the Hessian (\checkmark) in a different way.

G_μ	$= \mathbb{E}[\nabla_\mu l_x(\theta)]$	\leftrightarrow	
$H_{\mu\nu}$	$= \mathbb{E}[\nabla_\mu \nabla_\nu l_x(\theta)]$	\leftrightarrow	
$J_{\mu\nu\xi}$	$= \mathbb{E}[\nabla_\mu \nabla_\nu \nabla_\xi l_x(\theta)]$	\leftrightarrow	
$C_{\mu\nu}$	$= \mathbb{E}[(\nabla l_x(\theta) - G)^{\otimes 2}]_{\mu\nu}$	\leftrightarrow	
$S_{\mu\nu\xi}$	$= \mathbb{E}[(\nabla l_x(\theta) - G)^{\otimes 3}]_{\mu\nu\xi}$	\leftrightarrow	

Figure 8: **Named tensors**, typically evaluated at initialization ($\theta = \theta_0$). Def. ?? explains how diagrams depict tensors.

For a tight comparison, we scale the learning rates appropriately so that, to leading order, few-epoch and many-epoch SGD agree. Then few-epoch and many-epoch SGD differ, to leading order, in their sensitivity to ∇C :

Corollary 7 (\rightarrow). *SGD with $E = B = 1$, $\eta = \eta_0$ avoids high- C regions more than SGD with $E = E_0$, $B = 1$, $\eta = \eta_0/E_0$. Precisely: $\mathbb{E}[\theta_{E=E_0} - \theta_{E=1}]^\mu = \left(\frac{E_0-1}{4E_0}\right) N \eta^{\mu\rho} \eta^{\nu\xi} \nabla_\mu C_{\nu\xi} + o(\eta^2)$.*

In sum, high- C regions repel small- (E, B) SGD more than large- (E, B) SGD. We thus extend the $T = 2$ result of (Roberts 2018) and resolve some questions posed therein.

Figure 9 shows an instance of the process¹

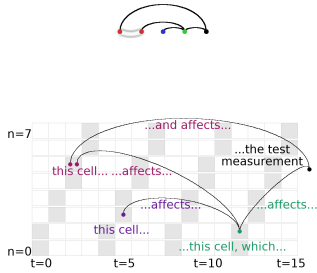


Figure 9: **A sub-process of SGD.** Timesteps index columns; training data index rows. The 5th datum participates in the 2nd SGD update. This $(n = 5, t = 2)$ event affects the testing loss both directly and via the $(1, 12)$ event, which is itself modulated by the $(2, 5)$ event.

Taylor series: method and challenges

Let's study $\mathbb{E}[\theta_T]$, $\mathbb{E}[l(\theta_T)]$. To warm up, we'll prove Prop 0 (c.f. (Nesterov 2004; Roberts 2018)).

¹Throughout, colors help us refer to parts of diagrams; colors lack mathematical meaning.

Proof. By gradient bounds: $\theta_T - \theta_0$ is $O(\eta^1)$. We **claim** that $(\theta_T - \theta_0)^\mu = -T \sum_\nu \eta^{\mu\nu} G_\nu + o(\eta^1)$. The claim holds when $T = 0$. Say the claim holds for \tilde{T} -step SGD with $T = \tilde{T} + 1$. The displacement $(\theta_T - \theta_{\tilde{T}})^\mu$ evaluates to:

$$\begin{aligned}
 & - \sum_\nu \eta^{\mu\nu} \nabla_\nu l_{n_{\tilde{T}}}(\theta_{\tilde{T}}) \\
 &= - \sum_\nu \eta^{\mu\nu} \nabla_\nu \left(l_{n_{\tilde{T}}}(\theta_0) + \sum_\xi \nabla_\xi l_{n_{\tilde{T}}}(\theta_0) (\theta_{\tilde{T}} - \theta_0)^\xi + o(\theta_{\tilde{T}} - \theta_0) \right) \\
 &= - \sum_\nu \eta^{\mu\nu} \nabla_\nu \left(l_{n_{\tilde{T}}}(\theta_0) + \nabla l_{n_{\tilde{T}}}(\theta_0) \cdot O(\eta^1) + o(O(\eta^1)) \right) \\
 &= - \sum_\nu \eta^{\mu\nu} \nabla_\nu l_{\tilde{T}}(\theta_0) + o(\eta^1)
 \end{aligned}$$

Applying the induction hypothesis proves the claim. We plug the claim into l 's Taylor series:

$$\begin{aligned}
 \mathbb{E}[l(\theta_T) - l(\theta_0)] &= \sum_\mu \nabla_\mu l(\theta_0) \mathbb{E}[\theta_T - \theta_0]^\mu + \mathbb{E}[o(\theta_T - \theta_0)] \\
 &= \sum_\mu \nabla_\mu l(\theta_0) (-T \eta G + o(\eta^1)) + o(O(\eta^1)) \\
 &= - \sum_{\mu\nu} T G_\mu \eta^{\mu\nu} G_\nu + o(\eta^1)
 \end{aligned}$$

Indeed, due to the assumption of analytic moments, the above expectations of $o(\eta^1)$ terms are still $o(\eta^1)$. \square

What happens when we keep higher order terms?

MULTIPLE MOMENTS — We used above that, to order η^1 , $\mathbb{E}[l(\theta_T)]$ depends on the training data only through the first moment $\mathbb{E}[\theta_T - \theta_0]$. But to compute $\mathbb{E}[l(\theta_T)]$ to higher order, we'd also need the k th moments $M_k^{\mu_0\mu_1\cdots} = \mathbb{E}[\prod_i (\theta_T - \theta_0)^{\mu_i}]$. We may achieve this by inductively proving multiple **claims**, one for each moment.

TUPLES OF TIMES — Complications arise even as we compute M_1 to order η^2 . We may not neglect the gradient correction $\nabla(\nabla l_{n_{\tilde{T}}}(\theta_0) \cdot (\theta_{\tilde{T}} - \theta_0))$ at the \tilde{T} th induction step. As the displacement $\theta_{\tilde{T}} - \theta_0$ contains (to order η^1) \tilde{T} terms, so will the correction. Totalling the correction over time thus yields $\sum_{0 \leq \tilde{T} < T} \tilde{T} = \binom{T}{2}$ summands, each (e.g. $\nabla l_5 \nabla l_2$) involving a pair of times. Order- d corrections represent the joint influence of d -tuples of times. Prop 0's result $\sum_{\tilde{T}} (-\eta \nabla l_{\tilde{T}}(\theta_0))$ is degree 1 in T ; but the order- d displacement is a degree d polynomial — very divergent — in T .

FACTORING'S FAILURE — To obtain $-TG\eta G$, we multiplied l 's derivatives by the expectations of such summands. In contrast to Prop 0, these expectations, even those of a fixed degree in η , now vary in form due to noise: some (e.g. $\nabla l_5 \nabla l_2$) have statistically independent factors that permit expectations to factor; others (e.g. $\nabla \nabla l_5 \nabla l_5$) do not. This is how ∇l_x 's higher cumulants (such as the covariance and skew of the gradient distribution) appear in our analysis.

DIVERSE DERIVATIVES — At order η^3 , a hessian correction $\nabla((\theta_{\tilde{T}} - \theta_0) \cdot \nabla l_{n_{\tilde{T}}}(\theta_0) \cdot (\theta_{\tilde{T}} - \theta_0)/2)$ augments the gradient correction. Then M_1 's order- η^3 summands vary in form, even when all expectations factor (as happens on noiseless landscapes). For instance, the hessian and gradient corrections respectively induce order- η^3 summands of $\mathbb{E}[l(\theta_T)]$ such as

$$\sum_{\mu\nu\xi} \eta^{\mu\phi} \eta^{\nu\pi} \eta^{\xi\rho} (\nabla_\mu l_x) (\nabla_\phi \nabla_\nu l_y) (\nabla_\pi \nabla_\xi l_z) (\nabla_\rho l)$$





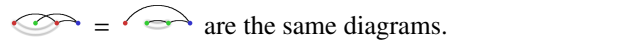
And M_2, M_3 's terms are yet more diverse. In short, a Taylor expansion even to low degrees yields a combinatorial explosion of terms. Our paper develops tools to organize and interpret these terms.

Since \curvearrowright has only one non-rightmost node, it has as many histories as there are timesteps. Since two nodes have degree 1, each history's contribution $G(-\eta(I - \eta H)^{\Delta-1})G$ involves two 1st derivatives: G . We use geometric series to sum from $\Delta = 1$ through $\Delta = T$ to find the testing loss change up to $o(\eta^1)$ error: $-G((I - (I - \eta H)^T)/H)G$. This is a more precise form of the above $-TG\eta G$ result.

The displacement $\mathbb{E}[\theta_T - \theta_0]$ contains many order- η^3 summands, including those of the form

$$\Delta_{xyz}^\xi \propto -\sum_{\delta\pi\rho}^{\mu\nu} \eta^{\mu\delta} \eta^{\nu\pi} \eta^{\xi\rho} \mathbb{E}[(\nabla_\mu l_x)(\nabla_\nu l_y)(\nabla_\delta \nabla_\pi \nabla_\rho l_z)]$$

where $0 \leq x, y, z < N$ label datapoints. Let $\Delta_o = \mathbb{E}[\Delta_{xxz} - \Delta_{xyz}]$ for x, y, z distinct. l_x, l_y, l_z are i.i.d., so: $\Delta_o^\xi \propto -\sum \dots \eta^{\mu\delta} \eta^{\nu\pi} \eta^{\xi\rho} C_{\mu\nu} J_{\delta\pi\rho}$ or, schematically, $\Delta_o^\xi \propto -\eta^3 C \nabla H$. Here, $C_{\mu\nu} = \mathbb{E}_x[\nabla_\mu l_x \nabla_\nu l_x] - G_\mu G_\nu$ is the covariance of gradients, $H_{\pi\rho} = \nabla_\pi \nabla_\rho l$ is l 's hessian, and $J_{\delta\pi\rho} = \nabla_\delta H_{\pi\rho}$ is l 's 'jerk'.¹

Valid diagrams include ,  but not . Since a diagram is just a rooted tree and partition,  =  are the same diagrams.

References

Absil, P.-A.; Mahony, R.; and Sepulchre, R. 2007. Optimization Algorithms on Matrix Manifolds, Chapter 4. *Princeton University Press*.

Amari, S.-I. 1998. Natural Gradient Works Efficiently. *Neural Computation*.

Barrett, D.; and Dherin, B. 2021. Implicit Gradient Regularization. *ICLR*.

Bartlett, P.; Foster, D.; and Telgarsky, M. 2017. Spectrally-Normalized Margin Bounds for Neural Networks. *NeurIPS*.

Belkin, M.; Hsu, D.; Ma, S.; and Mandal, S. 2019. Reconciling Modern Machine Learning Practice and the Bias-Variance Trade-off. *PNAS*.

Bottou, L. 1991. Stochastic Gradient Learning in Neural Networks. *Neuro-Nîmes*.

Cauchy, A.-L. 1847. Méthode générale pour la résolution des systèmes d'équations simultanées. *Comptes rendus de l'Académie des Sciences*.

Chaudhari, P.; and Soatto, S. 2018. SGD performs variational inference, converges to limit cycles for deep networks. *ICLR*.

Chladni, E. 1787. Entdeckungen über die Theorie des Klages. *Leipzig*.

Comon, P. 2014. An Introduction to Tensors. *IEEE Signal Processing Magazine*.

Dinh, L.; Pascanu, R.; Bengio, S.; and Bengio, Y. 2017. Sharp Minima Can Generalize For Deep Nets. *ICLR*.

Dixon, M.; and Ward, T. 2018. Takeuchi Information as a form of Regularization. *Arxiv Preprint*.

Dyer, E.; and Gur-Ari, G. 2019. Asymptotics of Wide Networks from Feynman Diagrams. *ICML Workshop*.

Dyson, F. 1949. The Radiation Theories of Tomonaga, Schwinger, and Feynman. *Physical Review*.

Finn, C.; Abbeel, P.; and Levine, S. 2017. Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks. *ICML*.

Gauss, C. 1823. Theoria Combinationis Obsevationum Erroribus Minimis Obnoxiae, section 39. *Proceedings of the Royal Society of Gottingen*.

Goyal, P.; Dollár, P.; Girshick, R.; Noordhuis, P.; Wesolowski, L.; Kyrola, A.; Tulloch, A.; Jia, Y.; and He, K. 2018. Accurate, Large Minibatch SGD. *Data @ Scale*.

Hoffer, E.; Hubara, I.; and Soudry, D. 2017. Train Longer, Generalize Better. *NeurIPS*.

Keskar, N.; Mudigere, D.; Nocedal, J.; Smelyanskiy, M.; and Tang, P. 2017. On Large-Batch Training for Deep Learning: Generalization Gap and Sharp Minima. *ICLR*.

Kiefer, J.; and Wolfowitz, J. 1952. Stochastic Estimation of the Maximum of a Regression Function. *Annals of Mathematical Statistics*.

Kolář, I.; Michor, P.; and Slovák, J. 1993. Natural Operations in Differential Geometry. *Springer*.

Krizhevsky, A. 2009. Learning Multiple Layers of Features from Tiny Images. *UToronto Thesis*.

Landau, L.; and Lifshitz, E. 1951. The Classical Theory of Fields. *Addison-Wesley*.

Landau, L.; and Lifshitz, E. 1960. Mechanics. *Pergamon Press*.

LeCun, Y.; Bengio, Y.; and Hinton, G. 2015. Deep Learning. *Nature*.

Li, Q.; Tai, C.; and E, W. 2017. Stochastic Modified Equations and Adaptive Stochastic Gradient Algorithms I. *PMLR*.

Liao, Q.; Miranda, B.; Banburski, A.; Hidary, J.; and Poggio, T. 2018. A Surprising Linear Relationship Predicts Test Performance in Deep Networks. *Center for Brains, Minds, and Machines Memo 91*.

McCullagh, P. 1987. Tensor Methods in Statistics. *Chemical Rubber Company Press*.

Mei, S.; and Montanari, A. 2020. The Generalization Error of Random Features Regression. *Arxiv Preprint*.

Misner, C.; Thorne, K.; and Wheeler, J. 1973. Gravitation. *W.H. Freeman and Company*.

¹'Jerk' names 3rd derivatives in dynamical systems: ISO 2041.

- Mohri, M.; Rostamizadeh, A.; and Talwalkar, A. 2018. Foundations of Machine Learning, Section 6.3.2. *MIT Press* .
- Nesterov, Y. 2004. Lectures on Convex Optimization: Minimization of Smooth Functions. *Springer Applied Optimization* 87, Section 2.1 .
- Neyshabur, B.; Bhojanapalli, S.; McAllester, D.; and Srebro, N. 2017a. Exploring Generalization in Deep Learning. *NeurIPS* .
- Neyshabur, B.; Tomioka, R.; Salakhutdinov, R.; and Srebro, N. 2017b. Geometry of Optimization and Implicit Regularization in Deep Learning. *Chapter 4 from Intel CRI-CI: Why and When Deep Learning Works Compendium* .
- Nickel, M.; and Kiela, D. 2017. Poincaré Embeddings for Learning Hierarchical Representations. *ICML* .
- Paszke, A.; Gross, S.; Massa, F.; Lerer, A.; Bradbury, J.; Killeen, T.; Lin, Z.; Gimelshein, N.; Antiga, L.; Desmaison, A.; Kopf, A.; Yang, E.; DeVito, Z.; Raison, M.; Tejani, A.; Chilamkurthy, S.; Steiner, B.; Fang, L.; Bai, J.; and Chintala, S. 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library. *NeurIPS* .
- Penrose, R. 1971. Applications of Negative Dimensional Tensors. *Combinatorial Mathematics and its Applications* .
- Robbins, H.; and Monro, S. 1951. A Stochastic Approximation Method. *Pages 400-407 of The Annals of Mathematical Statistics* .
- Roberts, D. 2018. SGD Implicitly Regularizes Generalization Error. *NeurIPS: Integration of Deep Learning Theories Workshop* .
- Roberts, D. 2019. SGD. *Personal communication* .
- Rota, G.-C. 1964. Theory of Möbius Functions. *Zeitschrift für Wahrscheinlichkeitstheorie und Verwandte Gebiete* .
- Sidiropoulos, N.; Lathauwer, L.; Fu, X.; Huang, K.; Papalexakis, E.; and Faloutsos, C. 2017. Tensor Decomposition for Signal Processing and Machine Learning. *IEEE Transactions on Signal Processing* .
- Stein, C. 1956. Inadmissibility of the Usual Estimator for the Mean of a Multivariate Normal Distribution. *Berkeley Symposium on Mathematical Probability* .
- Wang, H.; Keskar, N.; Xiong, C.; and Socher, R. 2018. Identifying Generalization Properties in Neural Networks. *Arxiv Preprint* .
- Wei, M.; and Schwab, D. 2019. How Noise Affects the Hessian Spectrum in Overparameterized Neural Networks. *Arxiv Preprint* .
- Werbos, P. 1974. Beyond Regression: New Tools for Prediction and Analysis. *Harvard Thesis* .
- Wu, J.; Hu, W.; Xiong, H.; Huan, J.; Braverman, V.; and Zhu, Z. 2020. On the Noisy Gradient Descent that Generalizes as SGD. *ICML* .
- Wu, L.; Ma, C.; and E, W. 2018. How SGD Selects the Global Minima in Over-Parameterized Learning. *NeurIPS* .
- Xiao, H.; Rasul, L.; and Vollgraf, R. 2017. Fashion-MNIST: a Novel Image Dataset for Benchmarking Machine Learning Algorithms. *Arxiv Preprint* .
- Yaida, S. 2019a. A First Law of Thermodynamics for SGD. *Personal Communication* .
- Yaida, S. 2019b. Fluctuation-Dissipation Relations for SGD. *ICLR* .
- Zhang, C.; Bengio, S.; Hardt, M.; Recht, B.; and Vinyals, O. 2017. Understanding deep learning requires rethinking generalization. *ICLR* .
- Zhang, H.; Reddi, S.; and Sra, S. 2016. Fast stochastic optimization on Riemannian manifolds. *NeurIPS* .
- Zhu, Z.; Wu, J.; Yu, B.; and Ma, J. 2019. The Anisotropic Noise in Stochastic Gradient Descent. *ICML* .
- Zou, D.; Cao, Y.; Zhou, D.; and Gu, Q. 2020. Stochastic Gradient Descent Optimizes Over-parameterized Deep ReLU Networks. *MLJ* .

Organization of the appendices

The following three appendices serve three respective functions:

- to explain how to calculate using diagrams;
- to prove our results (and pose a conjecture);
- to specify our experimental methods and results.

In more detail, we organize the appendices as follows.

A Tutorial: how to use diagrams	page 12
A.1 An example calculation: the effect of epochs	12
A.2 How to identify the relevant grid	17
A.3 How to identify the relevant diagram histories	17
A.4 How to evaluate each history	18
A.5 How to sum the histories' values	19
A.6 How to solve variant problems	20
A.7 Do diagrams streamline computation?	20
B Mathematics of the theory	page 23
B.1 Setting and assumptions	23
B.2 A key lemma à la Dyson	23
B.3 From Dyson to diagrams	24
B.4 Proof of Theorem 1	25
B.5 Proof of Theorem 2	26
B.6 Proofs of corollaries	26
B.7 Future topics	27
C Experimental methods	page 28
C.1 What artificial landscapes did we use?	28
C.2 What image-classification landscapes did we use?	28
C.3 Measurement process	29
C.4 Implementing optimizers	29
C.5 Software frameworks and hardware	29
C.6 Unbiased estimators of landscape statistics	29
C.7 Additional figures	30
D Review of Tensors	page ??
D.1 Vectors versus covectors	??
D.2 What is a tensor?	??
D.3 Tensors of type u_d	??
D.4 Contraction of tensors	??

Tutorial: how to use diagrams

This paper presents a new technique for calculating the expected learning curves of SGD in terms of statistics of the loss landscape near initialization. Here, we explain this technique. There are **four steps** to computing the expected testing loss, or other quantities of interest, after a specific number of gradient updates:

- **Specify, as a grid**, the batch size, training set size, and number of epochs.
- **Draw histories**, of diagrams into the grid, as needed for the desired precision.
- **Evaluate each diagram history**, whether exactly (via r values) or roughly (via u values).
- **Sum the histories' values** to obtain the quantity of interest as a function of η .

After presenting two example calculations that follow these four steps, we detail each step individually. Though we focus on the computation of expected testing losses, we describe how the four steps may give us other quantities of interest: variances instead of expectations, training statistics instead of testing statistics, or weight displacements instead of losses.

Two example calculations

We illustrate the four step procedure above by using it to answer the following two questions.

Our first example calculation reproduces Prop 0. In other words, it answers the question:

Question 1 (Leading order effect of gradients). *What's the leading order loss decrease $\mathbb{E}[l(\theta_T) - l(\theta_0)]$? We seek an answer expressed in terms of the landscape statistics at initialization: G, H, C, \dots . We expect only G to be relevant.*

Our second example is (an illustrative case of) Corollary 7.

Question 2 (Leading order effect of epochs). *How does multi-epoch SGD differ from single-epoch SGD? Specifically, what is the difference between the final testing losses of the following two versions of SGD?*

- SGD over $T = M_0 \times N$ time steps, learning rate η_0/M , and batch size $B = 1$
- SGD over $T = N$ time steps, learning rate η_0 , and batch size $B = 1$

We seek an answer expressed in terms of the landscape statistics at initialization: G, H, C, \dots .

To make our discussion concrete, we will set $M_0 = 2$; our analysis generalizes directly to larger M_0 .

We scaled the above two versions of SGD deliberately, to create an interesting comparison. Specifically, on a noiseless linear landscape $l_x = l \in (\mathbb{R}^n)^*$, the versions attain equal testing losses, namely $l(\theta_0) - T l_\mu \eta^{\mu\nu}$. So Question 2's answer will be second-order (or higher-order) in η .

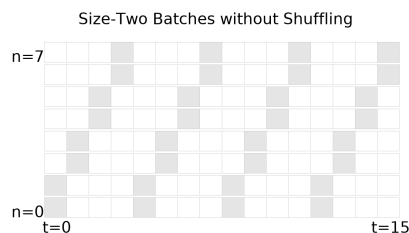
Computations: Grids We begin by asking a question about the final testing loss of some form of SGD. We specify the batch size, training set size, and number of epochs of the setting under analysis by drawing an appropriate grid. That is, we

- draw an $N \times T$ grid and
- shade its cells, shading the (n, t) th cell **when the t th batch includes the n th data point.**

Thus, each column contains B (batch size) many shaded cells and each row contains E (epoch number) many shaded cells.

EFFECT OF GRADIENTS (QUESTION 1)

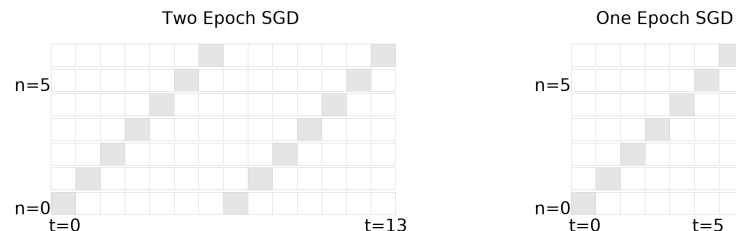
Question 1 does not specify a batch size, epoch number, or training set size and so does not specify a grid. In fact, we wish to answer the Question for any choice of those hyperparameters. E.g. we'll answer the Question for SGD with hyperparameters $B, E, N = 2, 4, 8$:



A grid for SGD with batch size $B = 2$ run for $E = 4$ epochs on $N = 8$ training points for a total of $T = 16$ timesteps.

EFFECT OF EPOCHS (QUESTION 2)

Two grids are relevant to Question 2: one for multi-epoch sgd and another for single-epoch SGD. See below.




Grids for single-epoch and multi-epoch SGD. Both grids depict $N = 7$ training points and batch size $B = 1$. **Left:** SGD with $M = 2$ update per training sample for a total of $T = MN = 2N$ many updates. **Right:** SGD with $M = 1$ update per training sample for a total of $T = MN = N$ many updates.

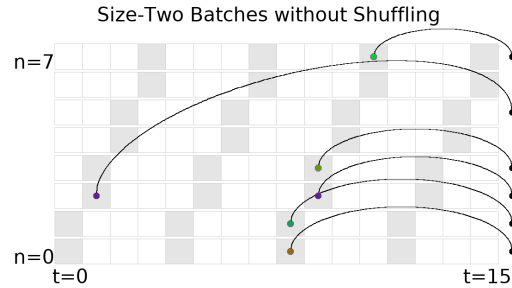
Computations: Embeddings of diagrams into grids Say we permit order- d errors. We draw all relevant diagrams with d or fewer edges and then characterize the histories of those diagrams in §'s grid. An *history* of a diagram D in a grid is an assignment of D 's non-root nodes to shaded cells (n, t) obeying the following criteria:


- **time-ordering condition:** the times t strictly increase along each path from leaf to root; and
- **correlation condition:** if two nodes are in the same part of D 's partition, then they are assigned to the same datapoint n .

We draw histories by placing nodes in their assigned shaded (n, t) cells; we draw the root nodes outside the grids (at arbitrary positions).



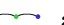


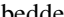




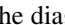

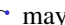
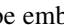
EFFECT OF GRADIENTS (QUESTION 1)

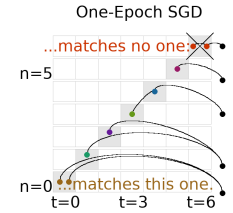
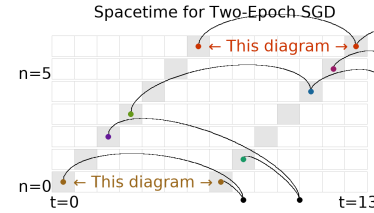
We seek an order 1 result and thus consider one-edged diagrams; there is only one: . We now describe the histories of this diagram:





 has one history for each shaded cell. An history of a non-root node at cell (n, t) represents the influence of the datapoint n on the testing loss due to the t th update. To first order (i.e. one-edged diagrams), the influences of different timesteps do not interact. The combinatorics of histories is thus straightforward.

EFFECT OF EPOCHS (QUESTION 2)

We seek an order 2 result and thus consider two-edged diagrams; there are four: , , , and . The figure below shows some histories of order-1 and order-2 diagrams (i.e. one-edged and two-edged diagrams) into the grid relevant to Question 2. Specifically, from top to bottom in each grid, the five diagrams embedded are  (or ) , , , and  (or ). The diagram  may be embedded wherever the diagram  may be embedded, but not vice versa. Likewise for  and .





Here,  embeds into the multi-epoch but not single-epoch grid.

Left:  embeds into the multi-epoch grid. **Right:**  cannot embed into the single-epoch grid: the correlation condition forces both red nodes into the same row and thus the same cell; the time-ordering condition forces the red nodes into distinct columns and thus distinct cells.

Computations: Evaluating each diagram history We evaluate §’s diagrams. We choose here to compute uvalues (apt for fixed T), not rvalues. These rules translate diagrams to numbers:

- **Node rule:** Replace each degree d node by $\nabla^d l_x$.
- **Outline rule:** surround the nodes in each part of the partition by a “cumulant bracket”. If a part contains one node x , the cumulant bracket is the expectation: $\mathbb{E}[x]$. If the part contains two nodes x, y , the cumulant bracket is the covariance: $\mathbb{E}[xy] - \mathbb{E}[x]\mathbb{E}[y]$.¹
- **Edge rule:** insert a $\eta^{\mu\nu}$ for each edge. The indices μ, ν should match the corresponding indices of the two nodes incident to the edge.

EFFECT OF GRADIENTS (QUESTION 1)

In § we determined the histories of . Now we evaluate uvalue(). The node rule suggests that we begin with

$$\nabla_\mu l_x \nabla_\nu l_x$$




We have two factors, each with one derivative because the diagram has two nodes, each of degree one. Note that the number of indices (here, two) is the total degree over all nodes and thus also twice the number (here, one) of edges. The outline rule transforms this to



$$\mathbb{E}[\nabla_\mu l_x] \mathbb{E}[\nabla_\nu l_x] = G_\mu G_\nu$$

since all parts in the diagram’s partition have size one. The edge rule inserts a factor $\eta^{\mu\nu}$ to yield:

$$\text{uvalue}(\text{red arrow}) = G_\mu G_\nu \eta^{\mu\nu} = G_\mu G^\mu$$

EFFECT OF EPOCHS (QUESTION 2)

In § we saw that  embeds similarly into multi-epoch and single-epoch grids: its multi-epoch histories correspond by a $M_0^2 : 1$ map to its single-epoch histories. Since we scaled the learning rate of the two SGD versions by a factor of M_0 , and since  (being two-edged) scales as η^2 , *the total uvalue of its multi-epoch histories will match the total uvalue of its single-epoch histories*. So we need not compute ’s contribution.

We see that this cancellation happens for all of the order-2 diagrams *except* for . Therefore, we must only compute uvalue().

The node rule suggests that we begin with $\nabla_\mu l_x \nabla_\nu \nabla_\lambda l_x \nabla_\rho l_x$. The outline rule transforms this to

$$\left(\mathbb{E}[\nabla_\mu l_x \nabla_\nu \nabla_\lambda l_x] - \mathbb{E}[\nabla_\mu l_x] \mathbb{E}[\nabla_\nu \nabla_\lambda l_x] \right) \mathbb{E}[\nabla_\rho l_x] = (\nabla_\nu C_{\mu\lambda}/2) G_\rho$$

The edge rule inserts a factor $\eta^{\mu\lambda} \eta^{\nu\rho}$ to yield:

$$\text{uvalue}(\text{red arrow}) = (\nabla_\nu C_{\mu\lambda}/2) G_\rho \eta^{\mu\lambda} \eta^{\nu\rho} = G^\nu \nabla_\nu C_\mu^\mu / 2$$

¹The general pattern is that the cumulant bracket $\mathbb{C}[\prod_{i \in I} x_i]$ of a product indexed by I is (here, P ranges over partitions of I with at least two parts; $I = \sqcup_{p \in P} p$):

$$\mathbb{C}[\prod_{i \in I} x_i] = \mathbb{E}[\prod_{i \in I} x_i] - \sum_{\text{partition } P} \prod_{p \in P} \mathbb{C}[\prod_{i \in p} x_i]$$

Computations: Summing the histories' values Our Key Lemma('s restatement) says that to compute a testing loss,

- we sum §'s uvalues, each weighted by the number of ways its diagram embeds in the grid,
- where histories with s many symmetries count only $1/s$ much toward the total number of histories.

A symmetry of an history f of a diagram D , i.e. an element of $\text{Aut}_f(D)$, is defined to be a relabeling of D 's nodes that simultaneously preserves D 's rooted tree structure, D 's partition structure, and f 's assignment of nodes to (n, t) cells of the grid. This is a strong constraint, so there will typically be no symmetries except for the identity, meaning that $s = 1$.

EFFECT OF GRADIENTS (QUESTION 1)

Referring again to §, we see that $D = \text{---} \text{---} \text{---}$ has TB many histories (B many histories per column for T many columns). Since D has no non-trivial automorphisms (i.e. no non-trivial relabeling of nodes that preserves the root, the graph structure, and the equivalence relation on non-root nodes), D has no non-trivial automorphisms that preserve any given history. Thus $|\text{Aut}_f(D)| = 1$ for each history of D . We conclude that the Restated Key Lemma's expression (3)

$$\sum_{\substack{D \text{ a} \\ \text{diagram}}} \sum_{\substack{f \text{ an embed-} \\ \text{-ding of } D}} \frac{(-B)^{-|\text{edges}(D)|}}{|\text{Aut}_f(D)|} \text{uvalue}(D)$$

has as its contribution from $D = \text{---} \text{---} \text{---}$ the value

$$(\# \text{of histories } f) \cdot \frac{(-B)^{-1}}{1} G_\mu G^\mu = TB \cdot (-G_\mu G^\mu / B)$$

Prop 0's expression $-TG_\mu G^\mu$ follows.

EFFECT OF EPOCHS (QUESTION 2)

Referring again to §, we see that $\text{---} \text{---} \text{---}$ has $\binom{M_0}{2} N$ many histories into the multi-epoch grid (one history per pair of distinct epochs, per row) — and no histories into the single-epoch grid. Moreover, each history of $\text{---} \text{---} \text{---}$ has $|\text{Aut}_f(D)| = 1$. We conclude that the testing loss of $M = M_0$ SGD exceeds the testing loss of $M = 1$ SGD by this much:

$$\binom{M_0}{2} N \cdot \frac{(-1)^2}{1} \cdot (\nabla_\nu C_{\mu\lambda}/2) G^\rho \eta^{\mu\lambda} \eta^{\nu\rho} + o(\eta^2)$$

Since Question 2 defines $\eta^2 = \eta_0^2 / M_0^2$, we can rewrite our answer as:

$$l(\theta_{M=M_0, \eta=\eta_0/M_0}) - l(\theta_{M=1, \eta=\eta_0}) = \frac{M_0 - 1}{4M_0} N \cdot G^\nu (\nabla_\nu C_\mu^\mu) + o(\eta_0^2)$$

where we use η_0 to raise indices. This completes the example problem.

How to identify the relevant grid

Diagrams tell us about the loss landscape but not about SGD’s batch size, number of epochs, and training set size. We encode this SGD data as a set of pairs (n, t) , where we have one pair for each participation of the n th datapoint in the t th update. For instance, full-batch GD has NT many pairs, and singleton-batch SGD has T many pairs. We will draw these (n, t) pairs as shaded cells in an $N \times T$ grid; we will call the shaded grid the SGD’s **grid**. See Figure 10.

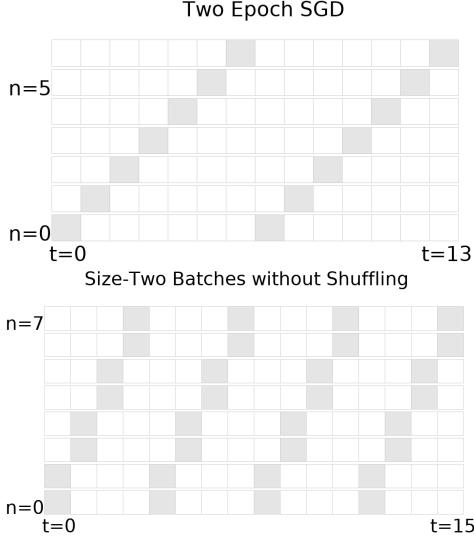


Figure 10: **The grids of two SGD variants.** Shaded cells show (n, t) pairs (see text). **Left:** Two epoch SGD with batch size one. **Right:** Four epoch SGD with batch size two.

WHEN USING THE diagram method to solve a problem relating to SGD (with batch size B and E many training samples), one shades the cells of an $N \times T$ grid with B shaded cells per column and E shaded cells per row.

Note: A grid may also depict the inter-epoch permuting of training sets due to which the b th batch in one epoch differs from the b th batch in a different epoch. For instance, see the grid to the right. Since each grid commits to a concrete sequence of training set permutations, we may analyze SGD with randomized permutations by taking expectations over multiple grids. However, the corollaries in this text are invariant to inter-epoch training set permutations, so we will not focus on this point.


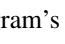
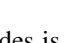
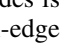
How to identify the relevant diagram histories

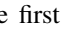
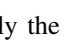
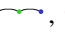

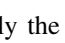


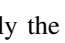


We explain below what the words mean in this green summary box.

WHEN USING THE diagram method to compute SGD’s final testing loss to edges and that have a non-zero number of histories into the relevant grid.

If we seek the isolated contribution due to a landscape statistic (e.g. d) that subgraph. If we are in a setting where a certain landscape statistic may neglect all diagrams that contain that subgraph.

If we are using r values (see next section for discussion of r values and d values). For each diagram, we must enumerate the histories, i.e. the assignments of time-ordering condition and correlation condition.

A *diagram* is a finite rooted tree equipped with a partition of its nodes, such that the root node occupies a part of size 1. For example, there are four diagrams with two edges: , , , and . As always, we specify a diagram’s root by drawing it rightmost.

A diagram is *linkless* when each of its degree-2 nodes is in a part of size one. Intuitively, this rules out multi-edge chains unadorned by fuzzy ties. Thus, only the first diagram in the list , , ,  is linkless. Only the first diagram in the list , ,  is linkless. Only the first diagram in the list , ,  is linkless.

An *history* of a diagram D into a grid is an assignment of D ’s non-root nodes to shaded cells (n, t) that obeys the following two criteria:

- **time-ordering condition:** the times t strictly increase along each path from leaf to root; and
- **correlation condition:** if two nodes are in the same part of D ’s partition, then they are assigned to the same datapoint n .

We may conveniently draw histories by placing nodes in the shaded cells to which they are assigned. Then, the time-ordering condition forbids (among other things) intra-cell

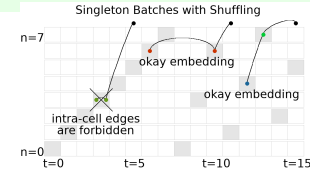
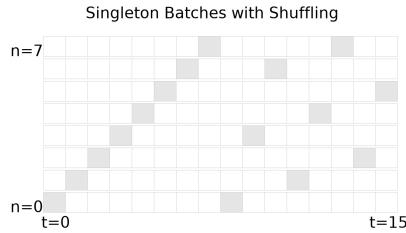
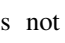

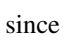
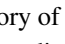
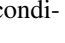





Figure 11: Embeddings, legal and illegal. **Left:** illegal history of , since the time-ordering condition is not obeyed. For the same reason, not a legal history of . **Middle:** an history of . Also an history of , since the correlation condition is obeyed. **Right:** a legal history of . Not an history of , since the correlation condition is not obeyed.

In principle, the relevant diagrams for a calculation with error $o(\eta^d)$ are the diagrams with at most d edges. For d greater than 2, there will be many such diagrams. However, in practice we gain insight even from considering one diagram at a time:

Remark. In this paper’s corollaries, we seek to extract the specific effect of a specific landscape or optimization feature such as skewed noise (Prop ??) or multiple epochs (§). In these cases, it is usually the case that most diagrams are irrelevant. For example, because a diagram evaluates to a product of its components, the only way the skewness of gradient noise can appear in our calculations is through diagrams such as  that have a part of size 3. Likewise, in § we argued by considering which histories that the only diagram relevant to Question 2 is . \diamond

Here are some further examples. Table 1 shows the 6 diagrams that may embed into the grid of $E = B = 1$. It shows each diagram in multiple ways to underscore that diagrams are purely topological and to suggest the ways in which these diagrams may embed into a grid.







$\Theta((\eta N)^3 N^{-0})$	$\Theta((\eta N)^3 N^{-1})$	$G \oplus \left(\frac{\mathbb{E}_x[\nabla \nabla l_x(\theta)]}{\mathbb{E}_x[\nabla l_x(\theta)]} \right)$	$C \triangleq \mathbb{E}_x[(\nabla l_x(\theta) - G)^2]$
		$H \triangleq \mathbb{E}_x[\nabla \nabla l_x(\theta)] \triangleq \checkmark$	$S \triangleq \mathbb{E}_x[(\nabla l_x(\theta) - G)^3]$
		$J \triangleq \mathbb{E}_x[\nabla \nabla \nabla l_x(\theta)] \triangleq \checkmark$	$\mathbb{E}_x[(\nabla l_x(\theta) - G)^4] - 3C^2$
		$\mathbb{E}_x[(\nabla l_x(\theta) - G)(\nabla \nabla l_x(\theta) - H)] \triangleq \checkmark$	
		$\mathbb{E}_x[(\nabla \nabla l_x(\theta) - H)(\nabla \nabla l_x(\theta) - H)] \triangleq \checkmark$	
		$\mathbb{E}_x[(\nabla l_x(\theta) - G)(\nabla \nabla \nabla l_x(\theta) - J)] \triangleq \checkmark$	$\mathbb{E}_x[(\nabla l_x(\theta) - G)^5] - 10CS$

Table 1: **Multiple ways to draw the 6 distinct degree-3 diagrams for $B = E = 1$ SGD’s testing loss.** Because the grid of $B = E = 1$ SGD has only one cell per row and one cell per column, the only diagrams that have a non-zero number of histories are the diagrams such that each ancestor-descendant pair in the rooted tree occupies two different parts of the partition. We show $(4 + 2) + (2 + 2 + 3) + (1)$ ways to draw the 6 diagrams. In fact, these drawings show all of the time-orderings of the diagrams’ nodes that are consistent with the time-ordering condition. **Organization:** We organize the diagrams into columns by the number of parts in their partitions. Because partitions (fuzzy outlines) indicate correlations between nodes (i.e. noise), diagrams with fuzzy outlines show deviations of SGD away from deterministic ODE. The big- Θ notation that heads the columns gives the asymptotics of the sum-over-histories of each diagram’s uvalues (for N large and η small even relative to $1/N$). **Left:** Diagrams for ODE behavior. **Center:** 1st order deviation of SGD away from ODE. **Right:** 2nd order deviation of SGD from ODE with appearance of non-Gaussian statistics.

How to evaluate each history

We will discuss how to compute both rvalues and uvalues. Both are ways of turning a diagram history into a number. The paper body mainly mentions rvalues. uvalues are simpler to calculate, since they depend only on a diagram’s topology, not on the way it is embedded. Physical intuition suggests that rvalues are more accurate; in particular, when we initialize near a non-degenerate local minimum, rvalues do not diverge to $\pm\infty$ as $T \rightarrow \infty$.

We will explain the following green summary box.

TURN AN HISTORY OF A DIAGRAM into its uvalue or rvalue by applying the fo

- **Node rule:** Replace each degree d node by $\nabla^d l_x$.
- **Outline rule:** surround the nodes in each part of the partition by a “cumulant bracket” is the expectation: $\mathbb{E}[x]$. If the part contains two nodes x and y , the cumulant bracket is $\mathbb{C}[\prod_{i \in I} x_i] = \mathbb{E}[\prod_{i \in I} x_i] - \sum_{\text{partition } P} \prod_{p \in P} \mathbb{C}[\prod_{i \in p} x_i]$. (The general pattern is that the cumulant bracket $\mathbb{C}[\prod_{i \in I} x_i]$ of a product indexed by I with at least two parts and $I = \sqcup_{p \in P} p$: $\mathbb{C}[\prod_{i \in I} x_i] = \mathbb{E}[\prod_{i \in I} x_i] - \sum_{\text{partition } P} \prod_{p \in P} \mathbb{C}[\prod_{i \in p} x_i]$ fine partitions, i.e., those partitions whose parts have size one.)
- If we wish to compute a uvalue, then we apply the **Edge rule for uvalue** match the corresponding indices of the two nodes incident to the edge.
- If we wish to compute an rvalue, then we apply the **Edge rule for rvalue** insert a factor of $K^{|\ell| - \ell - 1} \eta$, where $K \triangleq (I - \eta H)$. Here, we consider the ro

Un-resummed values: uvalue(D) Each part in a diagram’s partition looks like one of the following fragments (or one of the infinitely many analogous fragments):

The above examples illustrate the **Node rule**: each degree d node evaluates to $\nabla^d l_x$.

Fuzzy outlines dictate how to collect the $\nabla^d l_x$ s into expectation brackets. For example, we could collect the nodes within each part (of the partition) into a pair of expectation brackets $\mathbb{E}_x[\cdot]$ — call the result the **moment value**. However, this would yield (un-centered) moments such as $\mathbb{E}_x[(\nabla l_x(\theta))^2]$ instead of cumulants such as $C = \mathbb{E}_x[(\nabla l_x(\theta) - G)^2]$. For technical reasons (§?? and §), cumulants will be easier to work with than moments, so we will choose to define the values of diagrams slightly differently as follows.¹

Outline rule: surround the nodes in each part of the partition by a “cumulant bracket”. The cumulant bracket $\mathbb{C}[\prod_{i \in I} x_i]$ of a product indexed by I is (here, P ranges over partitions of I with at least two parts; $I = \sqcup_{p \in P} p$):

$$\mathbb{C}[\prod_{i \in I} x_i] = \mathbb{E}[\prod_{i \in I} x_i] - \sum_{\text{partition } P} \prod_{p \in P} \mathbb{C}[\prod_{i \in p} x_i]$$

Thus, a cumulant bracket of a diagram is the moment bracket of that diagram minus other terms. Those other terms are obtained by considering diagrams with the same graph structure but strictly more parts in their partition. The recursive definition of \mathbb{C} grounds out because the maximal number of parts in a partition of a finite set is finite.

¹ This is just the standard Möbius recursion for defining cumulants (see (Rota 1964)).

For example, if a part contains one node x , the cumulant bracket is the expectation: $\mathbb{E}[x]$. If the part contains two nodes x, y , the cumulant bracket is the covariance: $\mathbb{E}[xy] - \mathbb{E}[x]\mathbb{E}[y]$. If a part contains three nodes x, y, z , then the cumulant bracket is

$$\mathbb{E}[xyz] - \mathbb{E}[x]\mathbb{E}[yz] - \mathbb{E}[y]\mathbb{E}[xz] - \mathbb{E}[z]\mathbb{E}[xy] + 2\mathbb{E}[x]\mathbb{E}[y]\mathbb{E}[z]$$

We visualize the above in the following example:

Example 1. For example, if we denote moment values by solid gray fuzzy ties (instead of fuzzy outlines), then:

We will use the concept of “moment values” again in §??.

Finally, we come to edges. **Edge rule:** insert a factor of $\eta^{\mu\nu}$ for each edge. The indices μ, ν should match the corresponding indices of the two nodes incident to the edge.

Example 2 (Un-resummed value). Remember that = $C_{\mu\nu}$ and = $H_{\lambda\rho}$, so that = $C_{\mu\nu}H_{\lambda\rho}$. Then

$$\text{uvalue}(\text{red-green edge}) = C_{\mu\nu}H_{\lambda\rho}\eta^{\mu\lambda}\eta^{\nu\rho}$$

Here, has two edges, which correspond in this example to the tensor contractions via $\eta^{\mu\lambda}$ and via $\eta^{\nu\rho}$, respectively. \diamond

Resummed values: $\text{rvalue}_f(D)$ The only difference between rvalues and uvalues is in their rule for evaluating edges.

Edge rule: if an edge’s endpoints are embedded to times t, t' , insert a factor of $K^{t'-t-1}\eta$, where $K \triangleq (I - \eta H)$. Here, we consider the root node as embedded to the time T .

Example 3 (Re-summed value). Recall as in Example 2 that = $C_{\mu\nu}$ and = $H_{\lambda\rho}$, so that = $C_{\mu\nu}H_{\lambda\rho}$. Then

if f is an history of that sends the diagram’s red part to a time t (and its green root to T), we have:

$$\text{rvalue}_f(\text{red-green edge}) = C_{\mu\nu}H_{\lambda\rho} \left(K^{T-t-1}\eta \right)^{\mu\lambda} \left(K^{T-t-1}\eta \right)^{\nu\rho}$$

Here, has two edges, which correspond in this example to the tensor contractions via $(K^{\dots}\eta)^{\mu\lambda}$ and via $(K^{\dots}\eta)^{\nu\rho}$, respectively. \diamond

Overall In sum, we evaluate an history of a diagram by using the **node**, **outline**, and **edge** rules to build an expression of $\nabla^d l_{\text{xs}}$, \mathbb{E}_{xs} and η s. The difference between uvalues and rvalues lies only in their edge rule.

How to sum the histories’ values

We give examples of automorphism groups and we illustrate the integration mentioned in this green summary box:

WE OBTAIN OVERALL final testing loss expressions by adding together automorphism-group sizes as in 1.

If we are using rvalues instead of uvalues, we approximate sums over histories $\exp(-\eta H t)$, and we apply:rule

$$\int_{0 \leq u < T} du \exp(-uA) = (I - \exp(-TA))/A$$

When written in an eigenbasis of ηH , this A ’s coefficients are sums of edge involved in the relevant degrees of freedom over which we integrate

The Restated Key Lemma and Theorem 1 together say

Theorem. For any T : for η small enough, SGD has expected testing loss

$$\sum_{D \text{ a linkless diagram}} \sum_{f \text{ an embedding of } D} \frac{(-B)^{-|\text{edges}(D)|}}{|\text{Aut}_f(D)|} \text{rvalue}_f(D)$$

which is the same as

$$\sum_{D \text{ a diagram}} \sum_{f \text{ an embedding of } D} \frac{(-B)^{-|\text{edges}(D)|}}{|\text{Aut}_f(D)|} \text{uvalue}(D)$$

Here, B is the batch size.

How do we evaluate the above sum? Summing uvalues reduces to counting histories, which in all the applications reported in this text is a routine combinatorial exercise. However, when summing rvalues, it is often convenient to replace a sum over histories by an integral over times, and the power $(I - \eta H)^{\Delta t-1}$ by the exponential $\exp(-\Delta t \eta H)$. This incurs a term-by-term $1 + o(\eta)$ error factor, meaning that it preserves leading order results.

Example 4. Let us return to $D = \text{red-green edge}$, embedded, say, in the grid of one-epoch one-sample-per-batch SGD. From Example 3, we know that we want to sum the following value over all histories f , i.e. over all $0 \leq t < T$ to which the red part of the diagram’s partition may be assigned:

$$\text{rvalue}_f(\text{red-green edge}) = C_{\mu\nu} \left(K^{T-t-1}\eta \right)^{\mu\lambda} \left(K^{T-t-1}\eta \right)^{\nu\rho} H_{\lambda\rho}$$

Each history has a factor $(-B)^{-|\text{edges}(D)|} / |\text{Aut}_f(D)| = (-B)^{-2}/2$; we will multiply in this factor at the end so we now we focus on the \sum_f . So, using the aforementioned approximation, we seek to evaluate

$$\int_{0 \leq t < T} dt C_{\mu\nu} (\exp(-(T-t)\eta H)) \eta^{\mu\lambda} (\exp(-(T-t)\eta H)) \eta^{\nu\rho} H_{\lambda\rho} = C_{\mu\nu} \left(\int_{0 \leq t < T} dt \exp(-(T-t)((\eta H) \otimes I + I \otimes (\eta H))) \right)_{\pi\sigma}^{\mu\nu} \eta^{\pi\lambda} \eta^{\sigma\rho} H_{\lambda\rho}$$


We know from linear algebra and calculus that $\int_{0 \leq u < T} du \exp(-uA) = (I - \exp(-TA))/A$ (when A is a non-singular linear endomorphism). Applying this rule for

$u = T - t$ and $A = (\eta H) \otimes I + I \otimes (\eta H)$, we evaluate the integral as:


$$\dots = C_{\mu\nu} \left(\frac{I - \exp(-T((\eta H) \otimes I + I \otimes (\eta H)))}{(\eta H) \otimes I + I \otimes (\eta H)} \right)^{\mu\nu}_{\pi\sigma} \eta^{\pi\lambda} \eta^{\sigma\rho} H_{\lambda\rho}$$


This is perhaps easier to write in an eigenbasis of ηH :

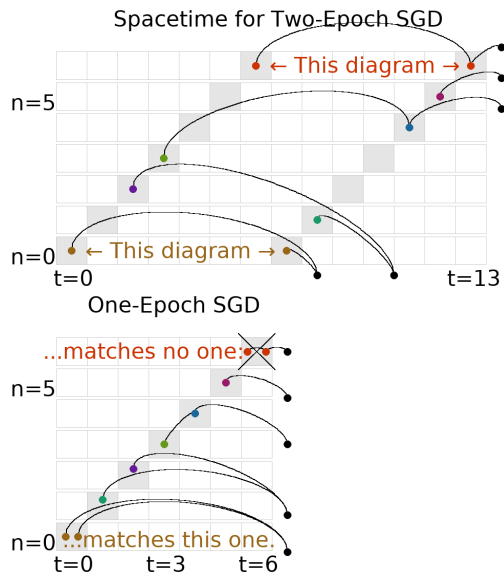
$$\dots = \sum_{\mu\nu} C_{\mu\nu} \frac{1 - \exp(-T((\eta H)_{\mu}^{\mu} + (\eta H)_{\nu}^{\nu}))}{(\eta H)_{\mu}^{\mu} + (\eta H)_{\nu}^{\nu}} (\eta H \eta)^{\mu\nu}$$

Multiplying this expression by the aforementioned $(-B)^{-2}/2$ gives the contribution of  to SGD's test loss. \diamond

An **automorphism of D that preserves a history f of D** is a relabeling of D 's nodes that preserves the root, the graph structure, the equivalence relation on non-root nodes, and that at the same time respects f : f must send a node to the same (n, t) pair to which the node's relabeling is sent. The sizes of automorphism groups appear in the denominators of our main results. They are usually of size one (that is, they usually contain only the identity relabeling).

Example 5 (Automorphisms). We take as examples the diagram histories of §?? (figure reproduced here). All the histories in the **left figure** have $|\text{Aut}_f(D)| = 1$. For instance, the bottom-most history (of ) has a non-trivial relabeling (namely, swap the two non-root nodes) that obeys all of the automorphism conditions *except* the “respects f ” condition. Indeed, the two non-root nodes are assigned to different cells in the grid, so we may not swap them without violating the “respects f ” condition.



By contrast, one of the histories (among the valid histories) shown in the **right figure** has a non-trivial automorphism group. This is the bottom-most history (of , again). Observe that swapping that diagram's two non-root nodes preserves the root, preserves the graph structure, and preserves the equivalence relation on non-root nodes. Moreover, such swapping respects f , since the two swapped nodes embed into the same cell. Thus, in this case, $|\text{Aut}_f(D)| = 2$. \diamond




How to solve variant problems

In §, we briefly discuss second-order methods and natural gradient descent. Here, we briefly discuss modifications. We omit proofs, which would closely follow §'s proof of the expectation-of-test-loss case.



Variance (instead of expectation) To compute variances instead of expectations (with respect to the noise in the training set), one considers generalized diagrams that have “two roots” instead of one. More precisely, to compute, say, the un-centered second moment of testing loss, one uses diagrams whose edge structures are not rooted trees but instead forests consisting of two rooted trees. We require that the set of roots (now a set of size two instead of size one) is a part of the diagram's partition. We draw the two roots rightmost.

For example, the generalized diagrams  or 

may appear in this computation.

Measuring on the training (instead of test) set To compute the training loss, we compute with all the same diagrams as the testing loss, and we also allow all the additional generalized diagrams that violate the constraint that a diagram's root should be in a part of size one. Therefore, to compute the generalization gap (i.e. testing loss minus training loss), we sum over all the diagrams that expressly violate this constraint (and then, since gen. gp is test minus train instead of train minus test, we multiply the whole answer by -1). For example, the generalized diagrams 

or  may appear in this computation.

Weight displacement (instead of loss) To compute displacements instead of losses, one considers generalized diagrams that have a “loose end” instead of a root. For example, the generalized diagrams  or  may appear in this computation.

Do diagrams streamline computation?

Diagram methods from Stueckelberg to Peierls have flourished in physics because they enable swift computations and offer immediate intuition that would otherwise require laborious algebraic manipulation. We demonstrate how our diagram formalism likewise streamlines analysis of descent by comparing direct perturbation¹ to the new formalism on two sample problems.

Aiming for a conservative comparison of derivation ergonomics, we lean toward explicit routine when using diagrams and allow ourselves to use clever and lucky simplifications when doing direct perturbation. For example, while solving the first sample problem by direct perturbation, we structure the SGD and GD computations so that the coefficients (that in both the SGD and GD cases are) called $a(T)$

¹By “direct perturbation”, we mean direct application of our Key Lemma (§).


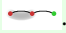
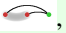


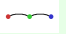
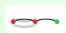
manifestly agree in their first and second moments. This allows us to save some lines.

Despite these efforts, the diagram method yields arguments about *four times shorter* — and strikingly more conceptual — than direct perturbation yields. (We make no attempt to compare the re-summed version of our formalism to direct perturbation because the algebraic manipulations involved for the latter are too complicated to carry out.) These examples specifically suggest that: diagrams obviate the need for meticulous index-tracking, from the start focus one's attention on non-cancelling terms by making visually obvious which terms will eventually cancel, and allow immediate exploitation of a setting's special posited structure, for instance that we are initialized at a test minimum or that the batch size is 1. We regard these examples as evidence that diagrams offer a practical tool for the theorist.

We now compare **Diagram Rules** vs **Direct Perturbation**.

Effect of batch size We compare the testing losses of pure SGD and pure GD. Because pure SGD and pure GD differ in how samples are correlated, their testing loss difference involves a covariance and hence occurs at order η^2 .

DIAGRAM METHOD —

Since SGD and GD agree on noiseless landscapes, we consider only diagrams with fuzzy ties. Since we are working to second order, we consider only two-edged diagrams. There are only two such diagrams,  and . The first diagram, , embeds in GD's space time in N^2 as many ways as it embeds in SGD's spacetime, due to horizontal shifts. Likewise, there are N^2 times as many histories of  in distinct epochs of GD's spacetime as there are in distinct epochs of SGD's spacetime. However, each same-epoch history of  within any one epoch of GD's spacetime corresponds by vertical shifts to an history of  in SGD. There are $MN \binom{N}{2}$ many such histories in GD's spacetime, so GD's testing loss exceeds SGD's by $\frac{MN \binom{N}{2}}{N^2}$ . Reading the diagram's value from its graph structure, we unpack that expression as:

$$\eta^2 \frac{M(N-1)}{4} G \nabla C$$

DIRECT PERTURBATION —

We compute the displacement $\theta_T - \theta_0$ to order η^2 for pure SGD and separately for pure GD. Expanding $\theta_t \in \theta_0 + \eta a(t) + \eta^2 b(t) + o(\eta^2)$, we find:

$$\begin{aligned} \theta_{t+1} &= \theta_t - \eta \nabla l_{n_t}(\theta_t) \\ &\in \theta_0 + \eta a(t) + \eta^2 b(t) - \eta(\nabla l_{n_t} + \eta \nabla^2 l_{n_t} a(t)) + o(\eta^2) \\ &= \theta_0 + \eta(a(t) - \nabla l_{n_t}) + \eta^2(b(t) - \nabla^2 l_{n_t} a(t)) + o(\eta^2) \end{aligned}$$

To save space, we write l_{n_t} for $l_{n_t}(\theta_0)$. It's enough to solve the recurrence $a(t+1) = a(t) - \nabla l_{n_t}$ and $b(t+1) = b(t) - \nabla^2 l_{n_t} a(t)$. Since $a(0), b(0)$ vanish, we have $a(t) = -\sum_{0 \leq t_0 < t_1 < T} \nabla l_{n_{t_1}}$ and $b(t) = \sum_{0 \leq t_0 < t_1 < T} \nabla^2 l_{n_{t_1}} \nabla l_{n_{t_0}}$. We now expand l :

$$\begin{aligned} l(\theta_T) &\in l + (\nabla l)(\eta a(T) + \eta^2 b(T)) \\ &\quad + \frac{1}{2}(\nabla^2 l)(\eta a(T) + \eta^2 b(T))^2 + o(\eta^2) \\ &= l + \eta(\nabla l a(T)) + \eta^2((\nabla l) b(T) + \frac{1}{2}(\nabla^2 l) a(T)^2) + o(\eta^2) \end{aligned}$$

Then $\mathbb{E}[a(T)] = -MN(\nabla l)$ and, since the N many singleton batches in each of M many epochs are pairwise independent,

$$\begin{aligned} \mathbb{E}[(a(T))^2] &= \sum_{0 \leq t < T} \sum_{0 \leq s < T} \nabla l_{n_t} \nabla l_{n_s} \\ &= M^2 N(N-1) \mathbb{E}[\nabla l]^2 + M^2 N \mathbb{E}[(\nabla l)^2] \end{aligned}$$

Likewise,

$$\begin{aligned} \mathbb{E}[b(T)] &= \sum_{0 \leq t_0 < t_1 < T} \nabla^2 l_{n_{t_1}} \nabla l_{n_{t_0}} \\ &= \frac{M^2 N(N-1)}{2} \mathbb{E}[\nabla^2 l] \mathbb{E}[\nabla l] + \\ &\quad \frac{M(M-1)N}{2} \mathbb{E}[(\nabla^2 l)(\nabla l)] \end{aligned}$$

Similarly, for pure GD, we may demand that a, b obey recurrence relations $a(t+1) = a(t) - \sum_n \nabla l_n / N$ and $b(t+1) = b(t) - \sum_n \nabla^2 l_n a(t) / N$, meaning that $a(t) = -t \sum_n \nabla l_n / N$ and $b(t) = \binom{t}{2} \sum_{n_0} \sum_{n_1} \nabla^2 l_{n_0} \nabla l_{n_1} / N^2$. So $\mathbb{E}[a(T)] = -MN(\nabla l)$ and

$$\begin{aligned} \mathbb{E}[(a(T))^2] &= M^2 \sum_{n_0} \sum_{n_1} \nabla l_{n_0} \nabla l_{n_1} \\ &= M^2 N(N-1) \mathbb{E}[\nabla l]^2 + M^2 N \mathbb{E}[(\nabla l)^2] \end{aligned}$$

and

$$\begin{aligned} \mathbb{E}[b(T)] &= \binom{MN}{2} \frac{1}{N^2} \sum_{n_0} \sum_{n_1} \nabla^2 l_{n_0} \nabla l_{n_1} \\ &= \frac{M(MN-1)(N-1)}{2} \mathbb{E}[\nabla^2 l] \mathbb{E}[\nabla l] + \\ &\quad \frac{M(MN-1)}{2} \mathbb{E}[(\nabla^2 l)(\nabla l)] \end{aligned}$$

We see that the expectations for a and a^2 agree between pure SGD and pure GD. So only b contributes. We conclude that pure GD's testing loss exceeds pure SGD's by


$$\begin{aligned} &\eta^2 \left(\frac{M(MN-1)(N-1)}{2} - \frac{M^2 N(N-1)}{2} \right) \mathbb{E}[\nabla^2 l] \mathbb{E}[\nabla l]^2 \\ &+ \eta^2 \left(\frac{M(MN-1)N}{2} - \frac{M(M-1)N}{2} \right) \mathbb{E}[(\nabla^2 l)(\nabla l)] \mathbb{E}[\nabla l] \\ &= \eta^2 \frac{M(N-1)}{2} \mathbb{E}[\nabla l] \left(\mathbb{E}[(\nabla^2 l)(\nabla l)] - \mathbb{E}[\nabla^2 l] \mathbb{E}[\nabla l] \right) \end{aligned}$$

Since $(\nabla^2 l)(\nabla l) = \nabla((\nabla l)^2)/2$, we can summarize this difference as

$$\eta^2 \frac{M(N-1)}{4} G \nabla C$$

Effect of non-Gaussian noise at a minimum. We consider vanilla SGD initialized at a local minimum of the testing loss. One expects θ to diffuse around that minimum according to gradient noise. We compute the effect on testing loss of non-Gaussian diffusion. Specifically, we compare SGD testing loss on the loss landscape to SGD testing loss on a different loss landscape defined as a Gaussian process whose every covariance agrees with the original landscape's. We work to order η^3 because at lower orders, the Gaussian landscapes will by construction match their non-Gaussian counterparts.

DIAGRAM METHOD —

Because $\mathbb{E}[\nabla l]$ vanishes at initialization, all diagrams with a degree-one vertex that is a singleton vanish. Because we work at order η^3 , we consider 3-edged diagrams. Finally, because all first and second moments match between the two landscapes, we consider only diagrams with at least one partition of size at least 3. The only such test diagram is . This embeds in T ways (one for each spacetime cell of vanilla SGD) and has symmetry factor $1/3!$ for a total of

$$\frac{T\eta^3}{6} \mathbb{E}[\nabla^3 l] \mathbb{E}[\nabla l_{n_a} \nabla l_{n_b} \nabla l_{n_c}]$$

DIRECT PERTURBATION —

We compute the displacement $\theta_T - \theta_0$ to order η^3 for vanilla SGD. Expanding $\theta_t \in \theta_0 + \eta a_t + \eta^2 b_t + \eta^3 c_t + o(\eta^3)$, we find:

$$\begin{aligned} \theta_{t+1} &= \theta_t - \eta \nabla l_{n_t}(\theta_t) \\ &\in \theta_0 + \eta a_t + \eta^2 b_t + \eta^3 c_t \\ &\quad - \eta \left(\nabla l_{n_t} + \nabla^2 l_{n_t}(\eta a_t + \eta^2 b_t) + \frac{1}{2} \nabla^3 l_{n_t}(\eta a_t)^2 \right) + o(\eta^3) \\ &= \theta_0 + \eta (a_t - \nabla l_{n_t}) \\ &\quad + \eta^2 (b_t - \nabla^2 l_{n_t} a_t) \\ &\quad + \eta^3 \left(c_t - \nabla^2 l_{n_t} b_t - \frac{1}{2} \nabla^3 l_{n_t} a_t^2 \right) + o(\eta^3) \end{aligned}$$

We thus have the recurrences $a_{t+1} = a_t - \nabla l_{n_t}$, $b_{t+1} = b_t - \nabla^2 l_{n_t} a_t$, and $c_{t+1} = c_t - \nabla^2 l_{n_t} b_t - \frac{1}{2} \nabla^3 l_{n_t} a_t^2$ with solutions: $a_t = -\sum_i \nabla l_{n_i}$ and $\eta^2 b_t = +\eta^2 \sum_{t_0 < t_1} \nabla^2 l_{n_{t_1}} \nabla l_{n_{t_0}}$. We do not compute c_t because we will soon see that it will be

multiplied by 0. To third order, the testing loss of SGD is

$$\begin{aligned} l(\theta_T) &\in l(\theta_0) + (\nabla l)(\eta a_T + \eta^2 b_T + \eta^3 c_T) \\ &\quad + \frac{\nabla^2 l}{2}(\eta a_T + \eta^2 b_T)^2 \\ &\quad + \frac{\nabla^3 l}{6}(\eta a_T)^3 + o(\eta)^3 \\ &= l(\theta_0) + \eta ((\nabla l) a_T) \\ &\quad + \eta^2 \left((\nabla l) b_T + \frac{\nabla^2 l}{2} a_T^2 \right) \\ &\quad + \eta^3 \left((\nabla l) c_T + (\nabla^2 l) a_T b_T + \frac{\nabla^3 l}{6} a_T^3 \right) + o(\eta)^3 \end{aligned}$$

Because $\mathbb{E}[\nabla l]$ vanishes at initialization, we neglect the (∇l) terms. The remaining η^3 terms involve $a_T b_T$, and a_T^3 . So let us compute their expectations:

$$\begin{aligned} \mathbb{E}[a_T b_T] &= - \sum_t \sum_{t_0 < t_1} \mathbb{E}[\nabla l_{n_t} \nabla^2 l_{n_{t_1}} \nabla l_{n_{t_0}}] \\ &= - \sum_{t_0 < t_1} \sum_{t \notin \{t_0, t_1\}} \mathbb{E}[\nabla l_{n_t}] \mathbb{E}[\nabla^2 l_{n_{t_1}}] \mathbb{E}[\nabla l_{n_{t_0}}] \\ &\quad - \sum_{t_0 < t_1} \sum_{t=t_0} \mathbb{E}[\nabla l_{n_t} \nabla l_{n_{t_0}}] \mathbb{E}[\nabla^2 l_{n_{t_1}}] \\ &\quad - \sum_{t_0 < t_1} \sum_{t=t_1} \mathbb{E}[\nabla l_{n_t} \nabla^2 l_{n_{t_1}}] \mathbb{E}[\nabla l_{n_{t_0}}] \end{aligned}$$

Since $\mathbb{E}[\nabla l]$ divides $\mathbb{E}[a_T b_T]$, the latter vanishes.

$$\begin{aligned} \mathbb{E}[a_T^3] &= - \sum_{t_a, t_b, t_c} \mathbb{E}[\nabla l_{n_a} \nabla l_{n_b} \nabla l_{n_c}] \\ &= - \sum_{\substack{t_a, t_b, t_c \\ \text{disjoint}}} \mathbb{E}[\nabla l_{n_a}] \mathbb{E}[\nabla l_{n_b}] \mathbb{E}[\nabla l_{n_c}] \\ &\quad - 3 \sum_{t_a = t_b \neq t_c} \mathbb{E}[\nabla l_{n_a} \nabla l_{n_b}] \mathbb{E}[\nabla l_{n_c}] \\ &\quad - \sum_{t_a = t_b = t_c} \mathbb{E}[\nabla l_{n_a} \nabla l_{n_b} \nabla l_{n_c}] \end{aligned}$$

As we initialize at a test minimum, only the last line remains, at it has T identical summands. When we plug into the expression for SGD testing loss, we get

$$\frac{T\eta^3}{6} \mathbb{E}[\nabla^3 l] \mathbb{E}[\nabla l_{n_a} \nabla l_{n_b} \nabla l_{n_c}]$$

Mathematics of the theory

Assumptions and Definitions


We assume throughout this work the following regularity properties of the loss landscape.

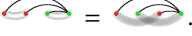
Existence of Taylor Moments — we assume that each finite collection of polynomials of the 0th and higher derivatives of the l_x , all evaluated at any point θ , may be considered together as a random variable insofar as they are equipped with a probability measure upon the standard Borel algebra.

Analyticity Uniform in Randomness — we assume that the functions $\theta \mapsto l_x(\theta)$ — and the expectations of polynomials of their 0th and higher derivatives — exist and are analytic with radii of convergence bounded from 0 (by a potentially θ -dependent function). So expectations and derivatives commute.

Boundedness of Gradients — we also assume that the gradients $\nabla l_x(\theta)$, considered as random covectors, are bounded by some continuous function of θ .¹ A metric-independent way of expressing this boundedness constraint is that the gradients all lie in some subset $S \subseteq TM$ of the tangent bundle of weight space, where, for any compact $C \subseteq M$, we have that the topological pullback — of $S \hookrightarrow TM \rightarrow M$ and $C \hookrightarrow M$ — is compact.

Now we turn to definitions.

Definition 5 (Diagrams). *A diagram is a finite rooted tree equipped with a partition of non-root nodes. We draw the tree using thin “edges”. By convention, we draw each node to the right of its children; the root is thus always rightmost. We draw the partition by connecting the nodes within each part via fuzzy “ties”. For example,  has 2 parts.*

We insist on using as few fuzzy ties as possible so that, if d counts edges and c counts ties, then $d + 1 - c$ counts parts of the partition. There may be multiple ways to draw a single diagram, e.g. .

Definition 6 (Embedding a Diagram into a Grid). *An history of a diagram into a grid is an assignment of that diagram’s non-root nodes to pairs (n, t) such that each node occurs at a time t' strictly after each of its children and such that two nodes occupy the same row n if they inhabit the same part of D ’s partition. In situations of interest, the permissible pairs (n, t) are those pairs such that the n th trainingpoint participates in the t th SGD update. We prefer to draw all such pairs by shading the corresponding cells in a square grid whose rows are indexed by training points n and whose columns are indexed by timesteps t .*

We define $\text{uvalue}(D)$ and $\text{rvalue}_f(D)$ as in §.

A key lemma à la Dyson

Suppose s is an analytic function defined on the space of weights. The following Lemma, reminiscent of (Dyson 1949), helps us track $s(\theta)$ as SGD updates θ :

¹Some of our experiments involve Gaussian noise, which is not bounded and so violates the hypothesis. In practice, Gaussians are effectively bounded in that our predictions vary smoothly with the first few moments of this distribution, so that a ± 12 -clipped Gaussian will yield almost the same predictions.

Key Lemma. *For all T : for η sufficiently small, $s(\theta_T)$ is a sum over tuples of natural numbers:*

$$s(\theta_T) = \sum_{(d_t: 0 \leq t < T) \in \mathbb{N}^T} (-\eta)^{\sum_t d_t} \left(\prod_{0 \leq t < T} \left(\frac{(g\nabla)^{d_t}}{d_t!} \Big|_{g=\sum_{m \in \mathcal{B}_t} \nabla l_m(\theta)/B} \right) \right) (s)(\theta_0) \quad (1)$$

Moreover, the expectation symbol (over training sets) commutes with the outer sum.

Here, we consider each $(g\nabla)^{d_t}$ as a higher order function that takes in a function f defined on weight space and outputs a function equal to the d_t th derivative of f , times g^{d_t} . The above product then indicates composition of $(g\nabla)^{d_t}$ ’s across the different t ’s. In total, that product takes the function s as input and outputs a function equal to some polynomial of s ’s derivatives.

Proof of the Key Lemma. We work in a neighborhood of the initialization so that the tangent space of weight space is a trivial bundle. For convenience, we fix a coordinate system, and with it the induced flat, non-degenerate inverse metric $\tilde{\eta}$; the benefit is that we may compare our varying η against one fixed $\tilde{\eta}$. Henceforth, a “ball” unless otherwise specified will mean a ball with respect to $\tilde{\eta}$ around the initialization θ_0 . Since s is analytic, its Taylor series converges to s within some positive radius ρ ball. By assumption, every l_t is also analytic with radius of convergence around θ_0 at least some $\rho > 0$. Since gradients are x -uniformly bounded by a continuous function of θ , and since in finite dimensions the closed ρ -ball is compact, we have a strict gradient bound b uniform in both x and θ on gradient norms within that closed ball. When

$$2\eta T b < \rho \tilde{\eta} \quad (2)$$

as norms, SGD after T steps on any train set will necessarily stay within the ρ -ball.² We note that the above condition on η is weak enough to permit all η within some open neighborhood of $\eta = 0$.

Condition 2 together with analyticity of s then implies that $(\exp(-\eta g \nabla) s)(\theta) = s(\theta - \eta g)$ when θ lies in the $\tilde{\eta}$ ball (of radius ρ) and its η -distance from that $\tilde{\eta}$ ball’s boundary exceeds b , and that both sides are analytic in η, θ on the same domain — and *a fortiori* when θ lies in the ball of radius $\rho(1 - 1/(2T))$. Likewise, a routine induction through T gives the value of s (after doing T gradient steps from an initialization θ) as

$$\left(\prod_{0 \leq t < T} \exp(-\eta g \nabla) \Big|_{g=\nabla l_t(\theta)} \right) (s)(\theta)$$

for any θ in the $\rho(1 - T/(2T))$ -ball (that is, the $\rho/2$ -ball), and that both sides are analytic in η, θ on that same domain. Note that in each exponential, the ∇_v does not act on the $\nabla_\mu l(\theta)$ with which it pairs.

Now we use the standard expansion of \exp . Because (by analyticity) the order d coefficients of l_t, s are bounded by

²The 2 ensures that SGD initialized at any point within a $\rho/2$ ball will necessarily stay within the ρ -ball.

some exponential decay in d that has by assumption an x -uniform rate, we have absolute convergence and may rearrange sums. We choose to group by total degree:

$$\dots = \sum_{0 \leq d < \infty} (-\eta)^d \sum_{\substack{(d_t: 0 \leq t < T) \\ \sum_t d_t = d}} \left(\prod_{0 \leq t < T} \frac{(g \nabla)^{d_t}}{d_t!} \Big|_{g = \nabla l_t(\theta)} \right) s(\theta) \quad (3)$$

The first part of the Key Lemma is proved. It remains to show that expectations over train sets commute with the above summation.

We will apply Fubini's Theorem. To do so, it suffices to show that

$$|c_d(l_t : 0 \leq t < T)| \triangleq \left| \sum_{\substack{(d_t: 0 \leq t < T) \\ \sum_t d_t = d}} \left(\prod_{0 \leq t < T} \frac{(g \nabla)^{d_t}}{d_t!} \Big|_{g = \nabla l_t(\theta)} \right) s(\theta) \right|$$

has an expectation that decays exponentially with d . The symbol c_d we introduce purely for convenience; that its value depends on the train set we emphasize using function application notation. Crucially, no matter the train set, we have shown that the expansion 3 (that features c_d appear as coefficients) converges to an analytic function for all η bounded as in condition 2. The uniformity of this demanded bound on η implies by the standard relation between radii of convergence and decay of coefficients that $|c_d|$ decays exponentially in d at a rate uniform over train sets. If the expectation of $|c_d|$ exists at all, then, it will likewise decay at that same shared rate.

Finally, $|c_d|$ indeed has a well-defined expected value, for $|c_d|$ is a bounded continuous function of a (finite-dimensional) space of T -tuples (each of whose entries can specify the first d derivatives of an l_t) and because the latter space enjoys a joint distribution. So Fubini's Theorem applies. The Key Lemma follows. \square

From Dyson to diagrams

We now describe the terms that appear in the Key Lemma. The following result looks like Theorem 1, except it has $\text{uvalue}(D)$ instead of $\text{uvalue}_f(D)$, and the sum is over all diagrams, not just linkless ones. In fact, we will use Theorem 3 to prove Theorem 1.

Theorem 3 (Test Loss as a Path Integral). *For all T : for η sufficiently small, SGD's expected test loss is*

$$\sum_D \sum_{\text{histories } f} \frac{1}{|\text{Aut}_f(D)|} \frac{\text{uvalue}(D)}{(-B)^{|\text{edges}(D)|}}$$

Here, D is a diagram whose root r does not participate in any fuzzy edge, f is an history of D into a grid, and $|\text{Aut}_f(D)|$ counts the graph-automorphisms of D that preserve f 's assignment of nodes to cells. If we replace D by $(-\sum_{p \in \text{parts}(D)} (D_{rp} - D)/N)$, where r is D 's root, we obtain the expected generalization gap (testing minus training loss).

Theorem 3 describe the terms that appear in the Key Lemma by matching each term to an history of a diagram

in a grid, so that the infinite sum becomes a sum over all diagram grid configurations. The main idea is that the combinatorics of diagrams parallels the combinatorics of repeated applications of the product rule for derivatives applied to the expression in the Key Lemma. Balancing against this combinatorial explosion are factorial-style denominators, again from the Key Lemma, that we summarize in terms of the sizes of automorphism groups.

Proof of Theorem 3. We first prove the statement about testing losses. Due to the analyticity property established in our proof of the Key Lemma, it suffices to show agreement at each degree d and train set individually. That is, it suffices to show — for each train set $(l_n : 0 \leq n < N)$, grid S , function $\pi : S \rightarrow [N]$ that induces \sim , and natural d — that

$$\begin{aligned} & (-\eta)^d \sum_{\substack{(d_t: 0 \leq t < T) \\ \sum_t d_t = d}} \left(\prod_{0 \leq t < T} \frac{(g \nabla)^{d_t}}{d_t!} \Big|_{g = \nabla l_t(\theta)} \right) l(\theta) = \\ & \sum_{\substack{D \in \text{im}(\mathcal{F}) \\ \text{with } d \text{ edges}}} \left(\sum_{f: D \rightarrow \mathcal{F}(S)} \frac{1}{|\text{Aut}_f(D)|} \right) \frac{\text{uvalue}_\pi(D, f)}{B^d} \quad (4) \end{aligned}$$

Here, uvalue_π is the value of a diagram history before taking expectations over train sets. We have for all f that $\mathbb{E}[\text{uvalue}_\pi(D, f)] = \text{uvalue}(D)$. Observe that both sides of 4 are finitary sums.

Remark 4 (Differentiating Products). *The product rule of Leibniz easily generalizes to higher derivatives of finitary products:*

$$\nabla^{|M|} \prod_{k \in K} p_k = \sum_{v: M \rightarrow K} \prod_{k \in K} (\nabla^{|v^{-1}(k)|} p_k)$$

The above has $|K|^{|M|}$ many term indexed by functions to K from M .

We proceed by joint induction on d and S . The base cases wherein S is empty or $d = 0$ both follow immediately from the Key Lemma, for then the only history is the unique history of the one-node diagram \bullet . For the induction step, suppose S is a sequence of $\mathcal{M} = \min S \subseteq S$ followed by a strictly smaller S and that the result is proven for (\tilde{d}, \tilde{S}) for every $\tilde{d} \leq d$. Let us group by d_0 the terms on the left hand side of desideratum 4. Applying the induction hypothesis with $\tilde{d} = d - d_0$, we find that that left hand side is:

$$\begin{aligned} & \sum_{0 \leq d_0 \leq d} \sum_{\substack{\tilde{D} \in \text{im}(\mathcal{F}) \\ \text{with } d - d_0 \text{ edges}}} \frac{1}{d_0!} \sum_{\tilde{f}: \tilde{D} \rightarrow \mathcal{F}(\tilde{S})} \left(\frac{1}{|\text{Aut}_{\tilde{f}}(\tilde{D})|} \right) \cdot \\ & (-\eta)^{d_0} (g \nabla)^{d_0} \Big|_{g = \nabla l_0(\theta)} \frac{\text{uvalue}_\pi(\tilde{D}, \tilde{f})}{B^{d-d_0}} \end{aligned}$$

Since $\text{uvalue}_\pi(\tilde{D}, \tilde{f})$ is a multilinear product of $d - d_0 + 1$ many tensors, the product rule for derivatives tells us that $(g \nabla)^{d_0}$ acts on $\text{uvalue}_\pi(\tilde{D}, \tilde{f})$ to produce $(d - d_0 + 1)^{d_0}$ terms. In fact, $g = \sum_{m \in \mathcal{M}} \nabla l_m(\theta)/B$ expands to $B^{d_0} (d - d_0 + 1)^{d_0}$ terms, each conveniently indexed by a pair of functions

$\beta : [d_0] \rightarrow \mathcal{M}$ and $\nu : [d_0] \rightarrow \tilde{D}$. The (β, ν) -term corresponds to an history f of a larger diagram D in the sense that it contributes $\text{uvalue}_\pi(D, f)/B^{d_0}$ to the sum. Here, (f, D) is (\tilde{f}, \tilde{D}) with $|(\beta \times \nu)^{-1}(n, \nu)|$ many additional edges from the cell of datapoint n at time 0 to the ν th node of \tilde{D} as embedded by \tilde{f} .

By the Leibniz rule of Remark , this (β, ν) -indexed sum by corresponds to a sum over histories f that restrict to \tilde{f} , whose terms are multiples of the value of the corresponding history of D . Together with the sum over \tilde{f} , this gives a sum over all histories f . So we now only need to check that the coefficients for each $f : D \rightarrow S$ are as claimed.

We note that the (β, ν) diagram (and its value) agrees with the $(\beta \circ \sigma, \nu \circ \sigma)$ diagram (and its value) for any permutation σ of $[d_0]$. The corresponding orbit has size

$$\frac{d_0!}{\prod_{(m,i) \in \mathcal{M} \times \tilde{D}} |(\beta \times \nu)^{-1}(m, i)|!}$$

by the Orbit Stabilizer Theorem of elementary group theory.

It is thus enough to show that

$$|\text{Aut}_f(D)| = |\text{Aut}_{\tilde{f}}(D)| \prod_{(m,i) \in \mathcal{M} \times \tilde{D}} |(\beta \times \nu)^{-1}(m, i)|!$$

We will show this by a direct bijection. First, observe that $f = \beta \sqcup \tilde{f} : [d_0] \sqcup \tilde{D} \rightarrow \mathcal{M} \sqcup \tilde{S}$. So each automorphism $\phi : D \rightarrow D$ that commutes with f induces both an automorphism $\mathcal{A} = \phi|_{\tilde{D}} : \tilde{D} \rightarrow \tilde{D}$ that commutes with \tilde{f} together with the data of a map $\mathcal{B} = \phi|_{[d_0]} : [d_0] \rightarrow [d_0]$ that both commutes with β . However, not every such pair of maps arises from a ϕ . For, in order for $\mathcal{A} \sqcup \mathcal{B} : D \rightarrow D$ to be an automorphism, it must respect the order structure of D . In particular, if $x \leq_D y$ with $x \in [d_0]$ and $y \in \tilde{D}$, then we need

$$\mathcal{B}(x) \leq_D \mathcal{A}(y)$$

as well. The pairs $(\mathcal{A}, \mathcal{B})$ that thusly preserve order are in bijection with the $\phi \in \text{Aut}_f(D)$. There are $|\text{Aut}_{\tilde{f}}(\tilde{D})|$ many \mathcal{A} . For each \mathcal{A} , there are as many \mathcal{B} as there are sequences $(\sigma_i : i \in \tilde{D})$ of permutations on $\{j \in [d_0] : j \leq_D i\} \subseteq [d_0]$ that commute with \mathcal{B} . These permutations may be chosen independently; there are $\prod_{m \in \mathcal{M}} |(\beta \times \nu)^{-1}(m, i)|!$ many choices for σ_i . Claim ?? follows, and with it the correctness of coefficients.

The argument for generalization gaps parallels the above when we use $l - \sum_n l_n/N$ instead of l as the value for s . Theorem 3 is proved. \square

Remark 5 (The Case of $E = B = 1$ SGD). *The grid of $E = B = 1$ SGD permits all and only those histories that assign to each part of a diagram's partition a distinct cell. Such histories factor through a diagram ordering and are thus easily counted using factorials per Prop 1. That prop immediately follows from the now-proven Theorem 3.*

Prop 1. *The order η^d contribution to the expected testing loss of one-epoch SGD with singleton batches is:*

$$\frac{(-1)^d}{d!} \sum_D |\text{ords}(D)| \binom{N}{P-1} \binom{d}{d_0, \dots, d_{P-1}} \text{uvalue}(D)$$

where D ranges over d -edged diagrams. Here, D 's parts have sizes $d_p : 0 \leq p \leq P$, and $|\text{ords}(D)|$ counts the total orderings of D s.t. children precede parents and parts are contiguous.

Proof of Theorem 1

(We say an history is **strict** if it assigns to each part a different datapoint n . Then, by Möbius inversion ((Rota 1964)), a sum over strict histories of moment values (§) matches a sum over all histories of uvalues.)

The diagrams summed in Theorem 1 and 2 may be grouped by their geometric realizations. Each nonempty class of diagrams with a given geometric realization has a unique element with minimally many edges, and in this way all and only linkless diagrams arise.

We encounter two complications: on one hand, that the sizes of automorphism groups might not be uniform among the class of diagrams with a given geometric realization. On the other hand, that the histories of a specific member of that class might be hard to count. The first we handle using Orbit-Stabilizer. The second we address as described by §?? via Möbius sums.

Proof of Theorem 1. We apply Möbius inversion (§??) to Theorem 3 (§).

The difference in loss from the noiseless case is given by all the diagram histories with at least one fuzzy tie, where the fuzzy tie pattern is actually replaced by a difference between noisy and noiseless cases as prescribed by the preceding discussion on Möbius Sums. Beware that even relatively noiseless histories may have illegal collisions of non-fuzzily-tied nodes within a single grid (data) row. Throughout the rest of this proof, we permit such illegal histories of the fuzz-less diagrams that arise from the aforementioned decomposition.

Because the Taylor series for analytic functions converge absolutely in the interior of the disk of convergence, the rearrangement of terms corresponding to a grouping by geometric realizations preserves the convergence result of Theorem 3.

Let us then focus on those diagrams σ with a given geometric realization represented by an linkless diagram ρ . By Theorem 3, it suffices to show that

$$\sum_{f:\rho \rightarrow S} \sum_{\substack{\tilde{f}:\sigma \rightarrow S \\ \exists i_*: f=\tilde{f} \circ i_*}} \frac{1}{|\text{Aut}_{\tilde{f}}(\sigma)|} = \sum_{f:\rho \rightarrow S} \sum_{\substack{\tilde{f}:\sigma \rightarrow S \\ \exists i_*: f=\tilde{f} \circ i_*}} \sum_{i:\rho \rightarrow \sigma} \frac{1}{|\text{Aut}_f(\rho)|} \quad (5)$$

Here, f is considered up to an equivalence defined by precomposition with an automorphism of ρ . We likewise consider \tilde{f} up to automorphisms of σ . And above, i ranges through maps that induce isomorphisms of geometric realizations, where i is considered equivalent to \hat{i} when for some automorphism $\phi \in \text{Aut}_{\tilde{f}}(\sigma)$, we have $\hat{i} = i \circ \phi$. Name as X the set of all such i is under this equivalence relation.

In equation 5, we have introduced redundant sums to structurally align the two expressions on the page; besides this rewriting, we see that equation 5's left hand side matches Theorem 3 resulting formula and that its right hand side is the desired formula of Theorem 1.

To prove equation 5, it suffices to show (for any f, \tilde{f}, i as above) that

$$|\text{Aut}_f(\rho)| = |\text{Aut}_{\tilde{f}}(\sigma)| \cdot |X|$$

We will prove this using the Orbit Stabilizer Theorem by presenting an action of $\text{Aut}_{\tilde{f}}(\rho)$ on X . We simply use precomposition so that $\psi \in \text{Aut}_{\tilde{f}}(\rho)$ sends $i \in X$ to $i \circ \psi$. Since $f \circ \psi = \tilde{f}$, $i \circ \psi \in X$. Moreover, the action is well-defined, because if $i \sim \hat{i}$ by ϕ , then $i \circ \psi \sim \hat{i} \circ \psi$ also by ϕ .

The stabilizer of i has size $|\text{Aut}_{\tilde{f}}(\rho)|$. For, when $i \sim i \circ \psi$ via $\phi \in \text{Aut}_{\tilde{f}}(\rho)$, we have $i \circ \psi = \phi \circ i$. This relation in fact induces a bijective correspondence: every ϕ induces a ψ via $\psi = i^{-1} \circ \phi \circ i$, so we have a map $\text{stabilizer}(i) \leftrightarrow \text{Aut}_{\tilde{f}}(\rho)$ seen to be well-defined and injective because structure set morphisms are by definition strictly increasing and because i must induce isomorphisms of geometric realizations. Conversely, every ψ that stabilizes enjoys *only* one ϕ via which $i \sim i \circ \psi$, again by the same (isomorphism and strict increase) properties. So the stabilizer has the claimed size.


Meanwhile, the orbit is all of $|X|$. Indeed, suppose $i_A, i_B \in X$. We will present $\psi \in \text{Aut}_{\tilde{f}}(\rho)$ such that $i_B \sim i_A \circ \psi$ by $\phi = \text{identity}$. We simply define $\psi = i_A^{-1} \circ i_B$, well-defined by the aforementioned (isomorphisms and strict increase) properties. It is then routine to verify that $f \circ \psi = \tilde{f} \circ i_A \circ i_A^{-1} \circ i_B = \tilde{f} \circ i_B = f$. So the orbit has the claimed size, and by the Orbit Stabilizer Theorem, the coefficients in the expansions of Theorems 1 and 3 match. \square

Proof of Theorem 2

Proof of Theorem 2. Since we assumed Hessians are positive: for any m , the propagator $K^t = ((I - \eta H)^{\otimes m})^t$ exponentially decays to 0 (at a rate dependent on m). Since up to degree d only a finite number of diagrams exist and hence only a finite number of possible m s, the exponential rates are bounded away from 0. Moreover, for any fixed t_{big} , the number of diagrams — involving no exponent t exceeding t_{big} — is eventually constant as T grows. Meanwhile, the number involving at least one exponent t exceeding that threshold grows polynomially in T (with degree d). The exponential decay of each term overwhelms the polynomial growth in the number of terms, and Theorem's first part follows. \square

Proofs of corollaries



Corollary 1

Proof. The relevant linkless diagram is  colored (amputated as in the previous subsection). An history of this diagram into $E = B = 1$ SGD's grid is determined by two durations — t from red to green and \tilde{t} from green to blue — obeying $t + \tilde{t} \leq T$. The automorphism group of each history has size 2: identity or switch the red nodes. So the answer is:

$$C_{\mu\nu} J_{\sigma}^{\rho\lambda} \left(\int_{t+\tilde{t} \leq T} (\exp(-t\eta H)\eta)^{\mu\rho} (\exp(-t\eta H)\eta)^{\nu\lambda} (\exp(-\tilde{t}\eta H)\eta)^{\sigma\pi} \right)$$

Standard calculus then gives the desired result. \square


Corollary 6's first part

Proof. The relevant linkless diagram is  (which equals  because we are at a test minimum). This diagram has one history for each pair of same-row shaded cells, potentially identical, in a grid; for GD, the grid has every cell shaded, so each *non-decreasing* pair of durations in $[0, T]^2$ is represented; the symmetry factor for the case where the cells is identical is $1/2$, so we lose no precision by interpreting a automorphism-weighted sum over the *non-decreasing* pairs as half of a sum over all pairs. Each of these may embed into N many rows, hence the factor below of N . The two integration variables (say, t, \tilde{t}) separate, and we have:

$$\frac{N}{B^{\text{degree}}} \frac{C_{\mu\nu}}{2} \int_t (\exp(-t\eta H))_{\lambda}^{\mu} \int_{\tilde{t}} (\exp(-\tilde{t}\eta H))_{\rho}^{\nu} \eta^{\lambda\sigma} \eta^{\rho\pi} H_{\sigma\pi}$$

Since for GD we have $N = B$ and we are working to degree 2, the prefactor is $1/N$. Since $\int_t \exp(at) = (I - \exp(-aT))/a$, the desired result follows. \square

Corollary 6's second part We apply the generalization gap modification (described in §) to Theorem 1's result about testing losses.

Proof. The relevant linkless diagram is . This diagram has one history for each shaded cell of grid; for GD, the grid has every cell shaded, so each duration from 0 to T is represented. So the generalization gap is, to leading order,

$$+ \frac{C_{\mu\nu}}{N} \int_t (\exp(-t\eta H))_{\lambda}^{\mu} \eta^{\lambda\nu}$$

Here, the minus sign from the gen-gap modification canceled with the minus sign from the odd power of $-\eta$. Integration finishes the proof. \square

Corollaries 7 and 2 Corollary 7 and Corollary 2 follow from plugging appropriate values of M, N, B into the following prop.

Prop 2. *To order η^2 , the testing loss of SGD — on N samples for $T = MN$ timesteps with batch size B dividing N and with any shuffling scheme — has expectation*

$$I - MNG_{\mu}G^{\mu} + MN \left(MN - \frac{1}{2} \right) G_{\mu}H_{\nu}^{\mu}G^{\nu} \\ + MN \left(\frac{M}{2} \right) C_{\mu\nu}H^{\mu\nu} + MN \left(\frac{M - \frac{1}{B}}{2} \right) (\nabla_{\mu}C_{\nu}^{\nu})G^{\mu}/2$$

of Prop 2. To prove Prop 2, we simply count the histories of the diagrams, noting that the automorphism groups are all of size 1 or 2. See Table 2. \square

Corollary 5 The corollary's first part follows immediately from Prop 2.

Proof of second part. Because $\mathbb{E}[\nabla I]$ vanishes at initialization, all diagrams with a degree-one vertex that is a singleton vanish. Because we work at order η^3 , we consider 3-edged diagrams. Finally, because all first and second moments match between the two landscapes, we consider only diagrams with at least one part (in their partition) of size at








diagram	embed.s w/ $ \text{Aut}_f = 1$	embed.s w/ $ \text{Aut}_f = 2$
	1	0
	MNB	0
	$\binom{MN}{2}B^2$	0
	$N\binom{MB}{2}$	0
	$\binom{MNB}{2}$	MNB
	$N\binom{MB}{2}$	MNB

Table 2: Terms used in proof of Prop 2

least 3. The only such test diagram is . This embeds in T ways (one for each grid cell — recall that $E = B = 1$) and has symmetry factor $1/3!$ for a total of

$$\frac{T\eta^3}{6}\mathbb{E}[\nabla^3 l]\mathbb{E}[\nabla l_{n_a}\nabla l_{n_b}\nabla l_{n_c}] = \frac{T\eta^3}{6}S_{\mu\nu\sigma}J^{\mu\nu\sigma}$$

This is the un-resummed expression. To obtain the resummed expression, we replace $\eta^{\mu\nu}$ with $(I - \eta H)^{\Delta t - 1} \eta^{\mu\nu}$. The histories range over T many times uniformly spaced in $[0, T]$. So we may integrate (let's name our variable of integration τ , and let's have it represent $\tau = \Delta t - 1 = T - t - 1$):

$$\int_{0 \leq \tau < T} (\exp(-\tau \eta H))^{\mu\nu} (\exp(-\tau \eta H))^{\pi\sigma} (\exp(-\tau \eta H))^{\lambda\rho} S_{\mu\pi\lambda} J_{\nu\sigma\rho}$$

Observe that

$$(\exp(-\tau \eta H))_{\mu}^{\nu} (\exp(-\tau \eta H))_{\pi}^{\sigma} (\exp(-\tau \eta H))_{\lambda}^{\rho} = (\exp(-\tau Y))_{\mu\nu\pi}^{\sigma\lambda\rho}$$

where $Y = \eta H \otimes I \otimes I + I \otimes \eta H \otimes I + I \otimes I \otimes \eta H$ is a six-index tensor. We finish by recalling that

$$\int_t \exp(At) = \frac{\exp(At)}{A} \Big|_t$$

□

Future topics

Our diagrams invite exploration of Lagrangian formalisms and curved backgrounds:¹

Question 3. *Does some least-action principle govern SGD; if not, what is an essential obstacle to this characterization?*

Lagrange's least-action formalism intimately intertwines with the diagrams of physics. Together, they afford a modular framework for introducing new interactions as new terms or diagram nodes. In fact, we find that some *higher-order* methods — such as the Hessian-based update $\theta \leftarrow \theta - (\eta^{-1} + \lambda \nabla \nabla l_t(\theta))^{-1} \nabla l_t(\theta)$ parameterized by small η, λ — admit diagrammatic analysis when we represent the λ term as a second type of diagram node. Though diagrams suffice for computation, it is Lagrangians that most deeply illuminate scaling and conservation laws.

Our work assumes a flat metric $\eta^{\mu\nu}$, but it might generalize to weight spaces curved in the sense of Riemann.²

¹(Landau and Lifshitz 1960, 1951) review these concepts.

²One may represent the affine connection as a node, thus giving rise to non-tensorial and hence gauge-dependent diagrams.

Such curvature finds concrete application in the *learning on manifolds* paradigm of (Absil, Mahony, and Sepulchre 2007; Zhang, Reddi, and Sra 2016), notably specialized to (Amari 1998)'s *natural gradient descent* and (Nickel and Kiela 2017)'s *hyperbolic histories*. While that work focuses on *optimization* on curved weight spaces, in machine learning we also wish to analyze *generalization*. Starting with the intuition that “smaller” hypothesis classes generalize better and that curvature controls the volume of small neighborhoods, we conjecture that sectional curvature regularizes learning:

Conjecture 1 (Sectional curvature regularizes). *If $\eta(\tau)$ is a Riemann metric on weight space, smoothly parameterized by τ , and if the sectional curvature through every 2-form at θ_0 increases as τ grows, then the gen. gap attained by fixed- T SGD with learning rate $c\eta(\tau)$ (when initialized from θ_0) decreases as τ grows, for all sufficiently small $c > 0$.*

We are optimistic our formalism may resolve conjectures such as above.

Experimental methods

What artificial landscapes did we use?

We define three artificial landscapes, called GAUSS, HELIX, and MEAN ESTIMATION.

GAUSS Consider fitting a centered normal $\mathcal{N}(0, \sigma^2)$ to some centered standard normal data. We parameterize the landscape by $h = \log(\sigma^2)$ so that the Fisher information matches the standard dot product (Amari 1998). More explicitly, the GAUSS landscape is a probability distribution \mathcal{D} over functions $l_x : \mathbb{R}^1 \rightarrow \mathbb{R}$ on 1-dimensional weight space, indexed by standard-normally distributed 1-dimensional datapoints x and defined by the expression:

$$l_x(h) \triangleq \frac{1}{2} (h + x^2 \exp(-h))$$

The gradient at sample x and weight σ is then $g_x(h) = (1 - x^2 \exp(-h))/2$. Since $x \sim \mathcal{N}(0, 1)$, the gradient $g_x(h)$ will be affinely related to a chi-squared, and in particular non-Gaussian.

To measure overfitting, we initialize at the true test minimum $h = 0$, then train and see how much the testing loss increases. At $h = 0$, the expected gradient vanishes, and the testing loss of SGD involves only diagrams that have no leaves of size one.

HELIX The HELIX landscape has chirality, much like Archimedes' screw. Specifically, the HELIX landscape has weights $\theta = (u, v, z) \in \mathbb{R}^3$, data points $x \sim \mathcal{N}(0, 1)$, and loss:

$$l_x(\theta) \triangleq \frac{1}{2} H(\theta) + x \cdot S(\theta)$$

Here,

$$H(\theta) = u^2 + v^2 + (\cos(z)u + \sin(z)v)^2$$

is quadratic in u, v , and

$$S(\theta) = \cos(z - \pi/4)u + \sin(z - \pi/4)v$$

is linear in u, v . Also, since $x \sim \mathcal{N}(0, 1)$, the $x \cdot S(\theta)$ term has expectation 0. In fact, the landscape has a three-dimensional continuous screw symmetry consisting of translation along z and simultaneous rotation in the $u - v$ plane. Our experiments are initialized at $u = v = z = 0$, which lies within a valley of global minima defined by $u = v = 0$.

The paper body showed that SGD travels in HELIX' $+z$ direction. By topologically quotienting the weight space, say by identifying points related by a translation by $\Delta z = 200\pi$, we may turn the line-shaped valley into a circle-shaped valley. Then SGD eternally travels, say, counterclockwise. Alternatively, one may preserve the homotopy type of the underlying weight space by Nash-embedding a flat solid torus

$$[-10^1, +10^1] \times [-10^1, +10^1] \times [-10^3, +10^3] / ((x, y, -10^3) \sim (x, y, +10^3))$$

in a higher-dimensional Euclidean space and extending HELIX from that torus to the ambient space.

Slightly modifying HELIX by adding a linear term $\alpha \cdot z$ to l for $\eta\alpha^2 \ll \eta^2/6$ leads SGD to perpetually ascend.

MEAN ESTIMATION The MEAN ESTIMATION family of landscapes has 1 dimensional weights θ and 1-dimensional datapoints x . It is defined by the expression:

$$l_x(\theta) \triangleq \frac{1}{2} H\theta^2 + xS\theta$$

Here, H, S are positive reals parameterizing the family; they give the hessian and (square root of) gradient covariance, respectively.

For our hyperparameter-selection experiment (Figure 13) we introduce an l_2 regularization term as follows:

$$l_x(\theta, \lambda) \triangleq \frac{1}{2} (H + \lambda)\theta^2 + xS\theta$$

Here, we constrain $\lambda \geq 0$ during optimization using projections; we found similar results when parameterizing $\lambda = \exp(h)$, which obviates the need for projection but necessitates a non-canonical choice of initialization. We initialize $\lambda = 0$.

What image-classification landscapes did we use?

Architectures In addition to the artificial loss landscapes GAUSS, HELIX, and MEAN ESTIMATION, we tested our predictions on logistic linear regression and simple convolutional networks (2 convolutional weight layers each with kernel 5, stride 2, and 10 channels, followed by two dense weight layers with hidden dimension 10) for the CIFAR-10 (Krizhevsky 2009) and Fashion-MNIST datasets (Xiao, Rasul, and Vollgraf 2017). The convolutional architectures used tanh activations and Gaussian Xavier initialization. To set a standard distance scale on weight space, we parameterized the model so that the Gaussian-Xavier initialization of the linear maps in each layer differentially pulls back to standard normal initializations of the parameters.

Some of our experiments involve Gaussian noise, which is not bounded and so violates the our assumptions. In practice, Gaussians are effectively bounded in that our predictions vary smoothly with the first few moments of this distribution, so that a ± 12 -clipped Gaussian will yield almost the same predictions. Even more experiments permit arbitrarily large losses and thus also violate our boundedness assumptions; since in practice SGD with small learning rates does not explore regions of very-large loss, we consider this violation negligible.

Datasets For image classification landscapes, we regard the finite amount of available data as the true (sum of diracs) distribution \mathcal{D} from which we sample testing and training sets in i.i.d. manner (and hence "with replacement"). We do this to gain practical access to a ground truth against which we may compare our predictions. One might object that this sampling procedure would cause testing and training sets to overlap, hence biasing testing loss measurements. In fact, testing and training sets overlap only in reference, in sense: the situation is analogous to a text prediction task in which two training points culled from different corpora happen to record the same sequence of words, say, "Thank you!". In any case, all of our experiments focus on the limited-data regime, e.g. 10^1 datapoints out of $\sim 10^{4.5}$ dirac masses, so overlaps are rare.

Measurement process

Diagram evaluation on real landscapes We implemented the formulae of § in order to estimate diagram values from real data measured at initialization from batch averages of products of derivatives.

Descent simulations We recorded testing and training losses for each of the trials below. To improve our estimation of average differences, when we compared two optimizers, we gave them the same random seed (and hence the same training sets).

We ran $2 \cdot 10^5$ trials of GAUSS with SDE and SGD, initialized at the test minimum with $T = 1$ and η ranging from $5 \cdot 10^{-2}$ to $2.5 \cdot 10^{-1}$. We ran $5 \cdot 10^1$ trials of HELIX with SGD with $T = 10^4$ and η ranging from 10^{-2} to 10^{-1} . We ran 10^3 trials of MEAN ESTIMATION with GD and STIC with $T = 10^2$, H ranging from 10^{-4} to $4 \cdot 10^0$, a covariance of gradients of 10^2 , and the true mean 0 or 10 units away from initialization.

We ran $5 \cdot 10^4$ trials of the CIFAR-10 convnet on each of 6 Glorot-Xavier initializations we fixed once and for all through these experiments for the optimizers SGD, GD, and GDC, with $T = 10$ and η between 10^{-3} and $2.5 \cdot 10^{-2}$. We did likewise for the linear logistic model on the one initialization of 0.

We ran $4 \cdot 10^4$ trials of the Fashion-MNIST convnet on each of 6 Glorot-Xavier initializations we fixed once and for all through these experiments for the optimizers SGD, GD, and GDC with $T = 10$ and η between 10^{-3} and $2.5 \cdot 10^{-2}$. We did likewise for the linear logistic model on the one initialization of 0.

Implementing optimizers

We approximated SDE by refining time discretization by a factor of 16, scaling learning rate down by a factor of 16, and introducing additional noise in the shape of the covariance in proportion as prescribed by the Wiener process scaling.

Our GDC regularizer was implemented using the unbiased estimator

$$\hat{C} \triangleq (I_x - I_y)_\mu I_{xy}/2$$

For our tests of regularization based on Corollary 6, we exploited the low-dimensional special structure of the artificial landscape in order to avoid diagonalizing to perform the matrix exponentiation: precisely, we used that, even on training landscapes, the covariance of gradients would be degenerate in all but one direction, and so we need only exponentiate a scalar.

Software frameworks and hardware

All code and data-wrangling scripts can be found on github.com/???????/perturb. This link will be made available after the period of double-blind review. Our code uses PyTorch 0.4.0 (Paszke et al. 2019) on Python 3.6.7; there are no other substantive dependencies. The code’s randomness is parameterized by random seeds and hence reproducible. We ran experiments on a Lenovo laptop and on our institution’s clusters; we consumed about 100 GPU-hours.

Unbiased estimators of landscape statistics

We use the following method — familiar to some but apparently nowhere described in writing — for obtaining unbiased estimates for various statistics of the loss landscape. The method is merely an elaboration of Bessel’s factor (Gauss 1823). For completeness, we explain it here.

Given samples from a joint probability space $\prod_{0 \leq d < D} X_d$, we seek unbiased estimates of *multipoint correlators* (i.e. products of expectations of products) such as $\langle x_0 x_1 x_2 \rangle \langle x_3 \rangle$. Here, angle brackets denote expectations over the population. For example, say $D = 2$ and from $2S$ samples we’d like to estimate $\langle x_0 x_1 \rangle$. Most simply, we could use $A_{0 \leq s < 2S} x_0^{(s)} x_1^{(s)}$, where A denotes averaging over the sample. In fact, the following also works:

$$S \left(A_{0 \leq s < S} x_0^{(s)} \right) \left(A_{0 \leq s < S} x_1^{(s)} \right) + (1 - S) \left(A_{0 \leq s < S} x_0^{(s)} \right) \left(A_{S \leq s < 2S} x_1^{(s)} \right) \quad (6)$$

When multiplication is expensive (e.g. when each $x_d^{(s)}$ is a tensor and multiplication is tensor contraction), we prefer the latter, since it uses $O(1)$ rather than $O(S)$ multiplications. This in turn allows more efficient use of batch computations on GPUs. We now generalize this estimator to higher-point correlators (and $D \cdot S$ samples).

For uniform notation, we assume without loss that each of the D factors appears exactly once in the multipoint expression of interest; such expressions then correspond to partitions on D elements, which we represent as maps $\mu : [D] \rightarrow [D]$ with $\mu(d) \leq d$ and $\mu \circ \mu = \mu$. Note that $|\mu| := |\text{im}(\mu)|$ counts μ ’s parts. We then define the statistic

$$\{x\}_\mu \triangleq \prod_{0 \leq d < D} A_{0 \leq s < S} x_d^{(\mu(d) \cdot S + s)}$$

and the correlator $\langle x \rangle_\mu$ we define to be the expectation of $\{x\}_\mu$ when $S = 1$. In this notation, 6 says:

$$\langle x \rangle_{\boxed{0} \boxed{1}} = \mathbb{E} \left[S \cdot \{x\}_{\boxed{0} \boxed{1}} + (1 - S) \cdot \{x\}_{\boxed{0} \boxed{1}} \right]$$

Here, the boxes indicate partitions of $[D] = [2] = \{0, 1\}$. Now, for general μ , we have:

$$\mathbb{E} \left[S^D \{x\}_\mu \right] = \sum_{\tau \leq \mu} \left(\prod_{0 \leq d < D} \frac{S!}{(S - |\tau(\mu^{-1}(d))|)!} \right) \langle x \rangle_\tau \quad (7)$$

where ‘ $\tau \leq \mu$ ’ ranges through partitions *finer* than μ , i.e. maps τ through which μ factors. In smaller steps, 7 holds because

$$\begin{aligned} \mathbb{E} \left[S^D \{x\}_\mu \right] &= \mathbb{E} \left[\sum_{(0 \leq s_d < S) \in [S]^D} \prod_{0 \leq d < D} x_d^{(\mu(d) \cdot S + s_d)} \right] \\ &= \sum_{\substack{(0 \leq s_d < S) \\ \in [S]^D}} \mathbb{E} \left[\prod_{0 \leq d < D} x_d^{(\min\{\vec{d} : \mu(\vec{d}) \cdot S + s_{\vec{d}} = \mu(d) \cdot S + s_d\})} \right] \\ &= \sum_{\tau} \left| \left\{ \left(\begin{smallmatrix} (0 \leq s_d < S) \in [S]^D : \\ \mu(d) = \mu(\vec{d}) \\ \wedge s_d = s_{\vec{d}} \end{smallmatrix} \right) \Leftrightarrow \tau(d) = \tau(\vec{d}) \right\} \right| \langle x \rangle_\tau \\ &= \sum_{\tau \leq \mu} \left(\prod_{0 \leq d < D} \frac{S!}{(S - |\tau(\mu^{-1}(d))|)!} \right) \langle x \rangle_\tau \end{aligned}$$

Solving 7 for $\langle x \rangle_\mu$, we find:

$$\langle x \rangle_\mu = \frac{S^D}{S^{|\mu|}} \mathbb{E}[\{x\}_\mu] - \sum_{\tau \prec \mu} \left(\prod_{d \in \text{im}(\mu)} \frac{(S-1)!}{(S - |\tau(\mu^{-1}(d))|)!} \right) \langle x \rangle_\tau$$

This expresses $\langle x \rangle_\mu$ in terms of the batch-friendly estimator $\{x\}_\mu$ as well as correlators $\langle x \rangle_\tau$ for τ *strictly* finer than μ . We may thus (use dynamic programming to) obtain unbiased estimators $\langle x \rangle_\mu$ for all partitions μ . Symmetries of the joint distribution and of the multilinear multiplication may further streamline estimation by turning a sum over τ into a multiplication by a combinatorial factor. For example, in the case of complete symmetry:

$$\langle x \rangle_{\boxed{012}} = S^2 \{x\}_{\boxed{012}} - \frac{(S-1)!}{(S-3)!} \{x\}_{\boxed{0} \boxed{1} \boxed{2}} - 3 \frac{(S-1)!}{(S-2)!} \{x\}_{\boxed{0} \boxed{12}}$$

Additional figures

In the rightmost figure, we add Corollary 6’s generalization gap estimate to l . By descending on this regularized loss, we may tune smooth hyperparameters such as l_2 regularization coefficients for small datasets ($H \ll C/N$) (§). Since matrix exponentiation takes time cubic in dimension, this regularizer is most useful for small models.

Review of Tensors

The linear algebra of tensors plays a role in our analysis similar to the role of types in C. One eschew types entirely, but to do so would hinder interpretation of the raw bitstrings we manipulate and would obscure which operations (e.g. dereferencing of a pointer) are natural and which are not (e.g. dereferencing of a float). The language of tensors supplies our analysis with abstractions and operations appropriate to derivatives and moments in high dimensional loss landscapes.

For example, we recognize $C^{-1} + H$ as ill-formed, even though both C^{-1} and H are square grids of numbers of the same shape. This is because C^{-1} and H are tensors of types $\frac{2}{0}$ and $\frac{0}{2}$, respectively.

So here’s a brief refresher for tensors. We recommend these sources for more details: (Sidiropoulos et al. 2017) §1 for motivational background; (Misner, Thorne, and Wheeler 1973) §2.5 for helpful pictures; (Comon 2014) §2 for examples of how new tensors arise from old; (Kolář, Michor, and Slovák 1993) §7 and §14 for precise formalism; and (McCullagh 1987) §1.4 for statistics-relevant examples.¹²

Vectors versus covectors

Imagine an air-conditioned hallway’s temperature gradient (0.1°K/meter) and length (20 meters). When we switch units from m to cm, the temperature gradient numerically *decreases* (to 0.001) but the length numerically *increases* (to 2000). So temperature gradients and spatial displacements are instances of distinct geometric types. A displacement inhabits a vector space V of primary interest; a gradient inhabits V ’s *dual space* V^* , defined as the set of linear maps from V to \mathbb{R} . When V is clear from context (e.g. throughout our paper V consists of tangent vectors on \mathcal{M}), we call elements of V *vectors* and elements of V^* *covectors*.³⁴

Imagine a smooth real-valued function $f : V \rightarrow \mathbb{R}$. Then we can specify the first derivative $\nabla f(x)$ (say at $x = 0$) by specifying f ’s p many partials with respect to a basis of V . Imagine a (say, compactly supported) probability measure μ

¹Some of these sources contrast the physicists’ and statisticians’ uses of tensors and of linear algebra overall. Physicists often use vector spaces whose elements represent changes through physical space and time; statisticians often use vector spaces whose elements represent mixtures of experimental subjects. Our paper doesn’t commit to or rely on any domain interpretation of our loss landscape and its associated vector spaces, so to us these distinctions are irrelevant.

²Some of these sources refer to *symmetric tensors* and *antisymmetric tensors* (a.k.a. *alternating forms*). Our paper does not use these concepts, so we invite the reader to black-box adjectives such as ‘symmetric’ when consulting those sources. A symmetric tensor (etc) is just a special case of tensors as we define in this appendix.

³This duality is most apparent when we recall that elements in V correspond bijectively with linear maps from \mathbb{R} to V . Then a *vector* (a.k.a. *column* a.k.a. *displacement* a.k.a. *primal* vector) is a linear map of type $\mathbb{R} \rightarrow V$ and a *covector* (a.k.a. *row* a.k.a. *gradient* a.k.a. *dual* vector) is a linear map of type $V \rightarrow \mathbb{R}$.

⁴Here’s why the word *dual* is apt. We may map $V \rightarrow (V^*)^*$ by evaluation: $v \mapsto (f \mapsto f(v))$. In finite dimensions, each $\omega \in (V^*)^*$ arises this way. So we identify $(V^*)^* = V$. So a space’s double dual is the space itself.

on V . We can specify the first moment $\mathbb{E}_\mu[v]$ of μ by specifying $\mathbb{E}_\mu[v]$'s p many projections with respect to a basis of V .

So both $\nabla f(0)$ and $\mathbb{E}_\mu[v]$ are objects of dimension p . But $\mathbb{E}_\mu[v] \in V$ is a vector while $\nabla f(0) \in V^*$ is a covector: the two differ in their geometric properties. There is a canonical way to push μ forward along a linear map $\phi : V \rightarrow W$ and thereby to send $\mathbb{E}_\mu[v]$ to $\mathbb{E}_{\mu[\phi(v)]}$. And there is a canonical way to pull f back along a linear map $\psi : U \rightarrow V$ and thereby to send $\nabla f(0)$ to $\nabla(f \circ \psi)(0)$. But there is no natural non-zero way to pull $\mu, \mathbb{E}_\mu[v]$ backward or to push $f, \nabla f(0)$ forward (try it!).

What is a tensor?

We've seen that the first derivative ∇f of a smooth function $f : V \rightarrow \mathbb{R}$ (evaluated, say, at $x \in V$) inhabits V^* , not V . How about higher derivatives? To answer this question, we introduce **tensors**: multi-axis grids of numbers viewed as linear algebraic objects.

Take the third derivative $\nabla \nabla \nabla f$ evaluated at $x \in V$. We can describe $\nabla \nabla \nabla f(x)$ completely by specifying the p^3 many partial derivatives with respect to a basis $\mathcal{E} = (e^i : 0 \leq i < p)$ of V :

$$[D_{i,j,k}^{\mathcal{E}} f](x) \triangleq \left. \frac{\partial}{\partial t} \frac{\partial}{\partial t'} \frac{\partial}{\partial t''} f(x + t \cdot e^i + t' \cdot e^j + t'' \cdot e^k) \right|_{t=t'=t''=0}$$

The p^3 many components $[D_{i,j,k}^{\mathcal{E}} f](x)$ enjoy key algebraic properties such as symmetry:

$$[D_{i,j,k}^{\mathcal{E}} f](x) = [D_{j,i,k}^{\mathcal{E}} f](x) = [D_{i,k,j}^{\mathcal{E}} f](x)$$

and linearity (which we see from the chain rule for derivatives):

$$[D_{i,j,k}^{\mathcal{F}} f](x) = \sum_i \sum_j \sum_k C_i^i C_j^j C_k^k [D_{i,j,k}^{\mathcal{E}} f](x)$$

Here, \mathcal{F} is a basis expressible in terms of \mathcal{E} via $f^i = \sum_i C_i^i e^i$.

By showing how $\nabla \nabla \nabla f$'s components change coherently as we change basis, the linearity rule helps us interpret $\nabla \nabla \nabla f$ as a geometric object with magnitude and directional information rather than as a mere grid of numbers.

We capture the essence of such linearity rules of by defining what it means for a vector space P to be a tensor product of given vector spaces V, W ; a 'tensor' is then just an element of P . The **DEFINING QUALITY** of P is that **a linear map from P to X is 'as good as' a bilinear map from V, W to X , for any vector space X** .¹ We write $V \otimes W$ for the two-axis tensor P . More generally, we write $U \otimes V \otimes \dots \otimes W$ for

¹Precisely, a tensor product of V, W is a vector space P equipped with a natural family of vector-space isomorphisms

$$\iota_X : \text{Bilin}(V, W; X) \xrightarrow{\sim} \text{Lin}(P; X)$$

from the vector space of bilinear maps to the vector space of linear maps. The word 'natural' just means that for any linear map $\phi : X \rightarrow Y$ and any $b \in \text{Bilin}(V, W; X)$:

$$\phi \circ \iota_X(b) = \iota_Y(\phi \circ b)$$

Naturality captures the intuition that ι just 're-packages' data without 'altering' its content. For instance, if $\kappa : X \rightarrow X$ is a non-trivial

an m -axis tensor P , the maps from which correspond to the m -linear maps from U, V, \dots, W .

We may now interpret the third derivative $\nabla \nabla \nabla f(x)$ as a linear-algebraic object:

$$\nabla \nabla \nabla f(x) \in (V \otimes V \otimes V)^*$$

The defining quality of $V \otimes V \otimes V = P$ says that such an element is specified the moment we specify a trilinear map from V, V, V to \mathbb{R} . This trilinear map is simply

$$(u, v, w) \mapsto \frac{\partial}{\partial t} \frac{\partial}{\partial t'} \frac{\partial}{\partial t''} f(x + t \cdot u + t' \cdot v + t'' \cdot w)$$

So $\nabla \nabla \nabla f(x)$ is a tensor. We now explain what we mean when we say that it is a tensor of type $\frac{0}{3}$.

Tensors of type $\frac{u}{d}$

We've seen two ways to make new finite-dimensional vector spaces from old: we can dualize (to get V^* from V) and we can tensor (to get $U \otimes V \otimes \dots \otimes W$ from U, V, \dots, W).

No matter how we compose these operations, we'll get a space of the standard form

$$(U \otimes V \otimes \dots \otimes W) \otimes (X^* \otimes Y^* \otimes \dots \otimes Z^*)$$

~~In other words, we have~~ a tensor product of u many non-dualized vector spaces and d many dualized vector spaces. We won't provide a proof.²

An example is that $((U^* \otimes V)^* \otimes W^*)^*$, though not of the above form, is naturally isomorphic to $U^* \otimes V \otimes W$, which *is* of the above form. That example illustrates a general algorithm: the factors in an un-standardized expression correspond with the factors in a standardized expression; whether a standardized factor is dualized or not depends simply on the **PARITY** of the number of dual operations that act on the corresponding factor in the original expression. Here, U, V, W appear with 3, 2, 2 layers of duals, respectively, whence by parity we obtain the 1, 0, 0 layers of duals in the standardization.

In the common case where all the underlying spaces U, \dots, Z agree with a vector space V understood from context, we say that elements of the combined space are **tensors of type $\frac{u}{d}$** . We represent such tensors as grids with $u+d$ many axes, where each axis has a length equal to the dimension of V . Thus, elements of $((V^* \otimes V)^* \otimes V^*)^*$ are tensors of type $\frac{2}{1}$. Numerically, they have $2+1=3$ axes: 2 that transform like vectors and 1 that transforms like a covector. The derivative

isomorphism, then defining $\iota'_X(b) = \kappa \circ \iota_X(b)$ and $\iota'_Z = \iota_Z$ for all $Z \neq X$ would turn a natural family ι of isomorphisms into a non-natural family ι' of isomorphisms.

²Here's an special case that contains all the ideas of a general proof. We'll define a natural isomorphism $\phi : U^* \otimes V^* \otimes W^* \rightarrow (U \otimes V \otimes W)^*$. By the defining quality of the domain, to specify ϕ is to specify a trilinear map from U^*, V^*, W^* into the codomain. By the defining quality of the codomain, to specify a map into the codomain is to specify a map into the space of trilinear real-valued maps on U, V, W . We accomplish all this by sending f, g, h to $((u, v, w) \mapsto f(u) \cdot g(v) \cdot h(w))$. We've thus defined a linear map ϕ . It is routine to check that ϕ sends no element except zero to zero, that ϕ 's domain and codomain have equal and finite dimension, and thus that ϕ is an isomorphism.

of θ_T 's covariance matrix with respect to the initialization θ_0 is an instance of such a type- 2_1 tensor.

As special cases, scalars are tensors of type 0_0 , vectors are tensors of type 1_0 , covectors are tensors of type 0_1 , and linear maps from V to V are tensors of type 1_1 .

Whenever we have a tensor of type a_a , we may interpret it as a linear map from the a -fold tensor product $V \otimes V \otimes \dots \otimes V$ to itself. In our work, this linear map will sometimes be invertible, in which case we use standard notation $(\cdot)^{-1}$ to denote the inverse. As a variation on this theme, a tensor (such as C) of type 0_2 we may view as a linear map from V to V^* . This linear map is invertible whenever C is positive definite, in which case we use the notation C^{-1} to denote its inverse. Observe that this inverse, being a linear map from V^* to V , is a tensor of type 2_0 , unlike C .

Contraction

Let's define a map $\epsilon_V : V \otimes V^* \rightarrow \mathbb{R}$ by evaluation: $(v, f) \mapsto f(v)$. We may similarly define

$$\text{id}_{(T \otimes U)} \otimes \epsilon_V : (T \otimes U) \otimes V \otimes V^* \rightarrow (T \otimes U)$$

by $(x, v, f) \mapsto f(v) \cdot x$ for $x \in T \otimes U$. In general, we can map any tensor product of the form

$$(U \otimes V \otimes \dots \otimes W) \otimes (U^* \otimes Y^* \otimes \dots \otimes Z^*)$$

to one of the form

$$(V \otimes \dots \otimes W) \otimes (Y^* \otimes \dots \otimes Z^*)$$

by applying the evaluation map ϵ_U .

Say we have a tensor of type ${}^{u+1}_{d+1}$. If we specify a specific axis among the $u+1$ vector-type axes as well as a specific axis among the $d+1$ covector-type axes THEN we may apply the above linear map to get a tensor of type u_d . We call this operation **contraction**. Operationally, contraction just means summing across the two specified axes over the $\dim(V)$ many 'diagonal' index pairs (i, i) .

For example, a linear map ϕ is the same as a tensor of type 1_1 . We may contract it to get a tensor of type 0_0 , that is, a real number. The number we get is the trace of ϕ .

As another example, consider the learning rate η (a tensor of type 2_0) and the hessian $H = \nabla G$ (a tensor of type 2_2). Then $\eta \otimes H$ is a tensor of type 2_2 and we may contract (in 2×2 many ways, here all equivalent due to η, H 's symmetry) to obtain tensor of type 1_1 : a linear map.

What is this linear map, concretely? It acts on a vector v by sending it to the vector w such that a gradient descent step $\theta \mapsto \theta - \eta \nabla f$ on the non-standard loss function

$$f(x) = G(x) \cdot v \in \mathbb{R}$$

displaces θ by $-w$. We see that the language of tensors concisely expresses otherwise cumbersome conceptual gymnastics. In our experience, this language also guides us through otherwise treacherously brittle algebraic manipulations.

¹Here's an example that our paper doesn't use but that may aid geometric intuition. If V has dimension p , then to specify a linear-algebraic notion of volume in V is to specify a tensor of type 0_p obeying certain antisymmetry properties. Intuitively, such a tensor produces for any given sequence of p input vectors a real number. We interpret this real number as the volume of a parallelepiped spanned by those p vectors.

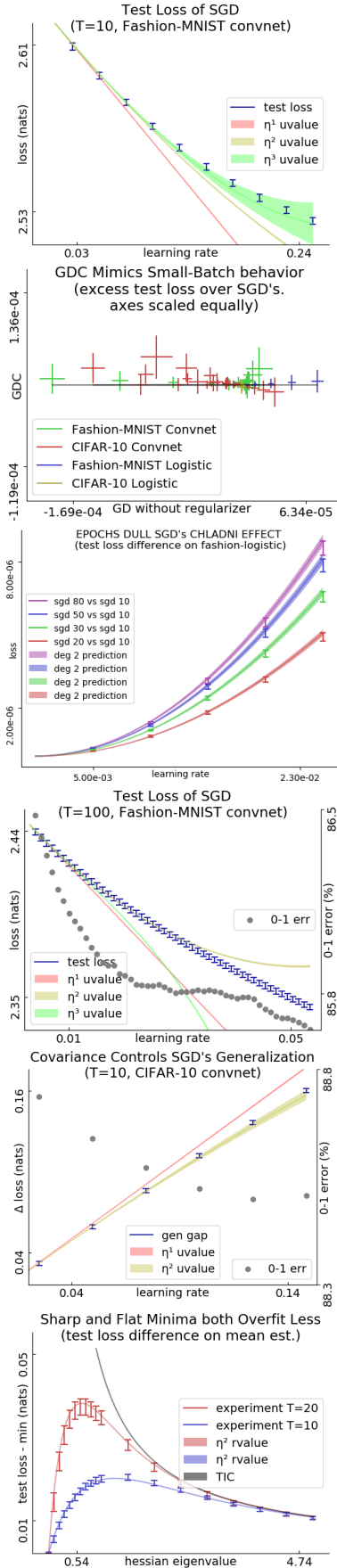


Figure 12: **Experiments on natural and artificial landscapes.** rvalue refers to Theorem 1's predictions, approximated as in Remark 1. uvalues are simpler but (see [\[10\]](#)) less accurate.

Left: Perturbation models SGD for small ηT . Fashion-MNIST convnet's testing-loss vs learning rate. In this

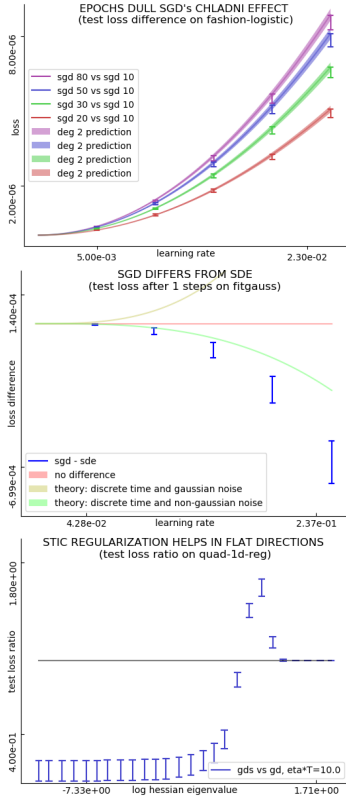


Figure 13: **Further experimental results.** **Left:** SGD with 2, 3, 5, 8 epochs incurs greater test loss than one-epoch SGD (difference shown in I bars) by the predicted amounts (predictions shaded) for a range of learning rates. Here, all SGD runs have $N = 10$; we scale the learning rate for E -epoch SGD by $1/E$ to isolate the effect of inter-epoch correlations away from the effect of larger ηT . **Center:** SGD's difference from SDE after $\eta T \approx 10^{-1}$ with maximal coarseness on GAUSS. Two effects not modeled by SDE — time-discretization and non-Gaussian noise oppose on this landscape but do not completely cancel. Our theory approximates the above curve with a correct sign and order of magnitude; we expect that the fourth order corrections would improve it further. **Right:** Blue intervals show regularization using Corollary 6. When the blue intervals fall below the black bar, this proposed method outperforms plain GD. For MEAN ESTIMATION with fixed C and a range of H s, initialized a fixed distance *away* from the true minimum, descent on an l_2 penalty coefficient λ improves on plain GD for most Hessians. The new method does not always outperform GD, because λ is not perfectly tuned according to STIC but instead descended on for finite ηT .