

1 General motivation

The inductive biases of deep learning optimizers are not well understood. SGD and its variants (e.g. Adam) are discrete update algorithms. At the beginning of training, the model parameters are initialized from some probability distribution (e.g. some gaussian distribution). Given some initial set of parameters, the SGD update has some stochasticity caused by sampling finite batches. These two sources of randomness means that SGD could end up in different minimizers from the same initial set of parameters. We want to characterize the distribution of the parameters during and at the end of training.

The motivation for safety is to be able to sample such distribution of parameters. If we had a good sampler of the distribution of parameters during training, we could systematically sample models that are compatible with minimizing the loss. We would have better guarantees that the sampling process is unbiased and how long it would take to sample the distribution of parameters that are accessible to the sampler. This would for example allows to determine the fraction of model minimizing the loss that are safe or unsafe and give us a lower bound on the probability of sampling a safe model in the set of degenerate solutions of the loss minimization and in particular and also help us understand the robustness of some safety property.

The simplest way to model SGD is to consider the continuous-time limit of SGD. The continuous-time limit of SGD is a stochastic differential equation (SDE) that describes the evolution of the parameters during training. Under Gaussian assumptions of the noise, this continuous-time limit leads to an SDE with a drift term and a Wiener process. Under further assuming that the noise covariance of the noise is constant, we can show that the stationary distribution of SGD obtained from Fokker-Planck is the same as the tempered Bayesian posterior distribution i.e. a Gibbs distribution with a temperature parameter.

This model of SGD has been heavily exploited in the literature, but the assumptions that are made by such models are not always valid. For example, the noise covariance is not constant (anisotropy) and SGD noise has heavy tails. Also, finite learning rate is important because of edge of stability phenomena where finite learning rate allows induced some sensitivity to the curvature in the dynamics inducing SGD to escape the sharpest minima. Subdiffusive phenomena have also been observed which can be modelled at the level of the Fokker-Planck equation by a fractional time derivative operator.

A more general model of SGD would take into account: a central flow term to model edge of stability phenomena from finite learning rate, a Levy process with alpha stable distribution to model heavy tails of the noise, a fractional time derivative operator to model subdiffusive phenomena and a non-constant noise covariance matrix to model anisotropy of the noise.

In this work we propose to bring together a more general SDE model of SGD taking into account these different effects. We study the empirical relevance of such SDE on toy models: deep linear networks, modular addition and toy transformer models. The main hypothesis that we test are whether the distribution of parameters obtained by such model better capture the training dynamics in toy models than the Langevin SDE model of SGD which leads to a tempered posterior distribution.

2 Introduction to Stochastic Gradient Descent

2.1 Problem Setup

Let $(\mathcal{X} \times \mathcal{Y}, P)$ be a probability space, where \mathcal{X} is the input space and \mathcal{Y} is the output space. We observe a dataset $\mathcal{D}_n = \{(X_i, Y_i)\}_{i=1}^n$ consisting of n i.i.d. samples from the distribution P .

Consider a parametric model $f : \Theta \times \mathcal{X} \rightarrow \mathcal{Y}$, where $\Theta \subseteq \mathbb{R}^d$ is the parameter space. For neural networks, $f(\theta; x)$ represents the network's output given parameters $\theta \in \Theta$ and input $x \in \mathcal{X}$.

2.2 Loss Functions

We define a loss function $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$ that measures the discrepancy between predictions and targets. For regression tasks, we typically use the squared loss:

$$\ell(\hat{y}, y) = \frac{1}{2} \|\hat{y} - y\|^2$$

This induces the following hierarchy of loss functions:

1. **Individual sample loss:** For a single data point (x, y) ,

$$\ell(\theta; x, y) := \ell(f(\theta; x), y)$$

2. **Empirical risk** (finite-sample loss): Over the dataset \mathcal{D}_n ,

$$L_n(\theta) := \frac{1}{n} \sum_{i=1}^n \ell(\theta; X_i, Y_i)$$

3. **Population risk** (expected loss): Over the true distribution P ,

$$L(\theta) := \mathbb{E}_{(X,Y) \sim P}[\ell(\theta; X, Y)]$$

In practice, we only have access to $L_n(\theta)$ and use it as a proxy for $L(\theta)$, which we truly wish to minimize.

2.3 Gradients

Define the gradients with respect to parameters θ :

1. **Individual gradient:** $g_i(\theta) := \nabla_{\theta} \ell(\theta; X_i, Y_i)$
2. **Empirical gradient:** $g_n(\theta) := \nabla L_n(\theta) = \frac{1}{n} \sum_{i=1}^n g_i(\theta)$
3. **Stochastic gradient:** For a batch $\mathcal{B} \subset D_n$ of size B sampled uniformly at random,

$$g_{\mathcal{B}}(\theta) := \frac{1}{B} \sum_{i \in \mathcal{B}} g_i(\theta)$$

4. **Population gradient:** $g(\theta) := \nabla L(\theta) = \mathbb{E}_{(X,Y) \sim P}[\nabla_{\theta} \ell(\theta; X, Y)]$

Note that $\mathbb{E}_{\mathcal{B}}[g_{\mathcal{B}}(\theta)] = \nabla L_n(\theta)$, making the stochastic gradient an unbiased estimator of the empirical gradient.

2.4 The SGD Algorithm

Stochastic Gradient Descent (SGD) performs the following iterative update:

where $\eta_k > 0$ is the learning rate (or step size) at iteration k .

2.5 Mini-batch Sampling Strategies

There are two common sampling strategies:

1. **With replacement:** Each element of \mathcal{B}_k is drawn independently and uniformly from $\{1, \dots, n\}$
2. **Without replacement:** Sample B distinct indices uniformly from $\{1, \dots, n\}$

The with-replacement case is often easier to analyze theoretically due to independence.

Algorithm 1 Stochastic Gradient Descent

```
1: Initialize  $\theta_0 \in \Theta$ 
2: for  $k = 0, 1, 2, \dots, K - 1$  do
3:   Sample mini-batch  $\mathcal{B}_k \subseteq D_n$  uniformly at random
4:   Compute stochastic gradient  $g_{\mathcal{B}_k}(\theta_k)$ 
5:   Update:  $\theta_{k+1} = \theta_k - \eta_k g_{\mathcal{B}_k}(\theta_k)$ 
6: end for
```

2.6 Decomposition into Drift and Noise

It is useful to decompose the SGD update as:

$$\theta_{k+1} = \theta_k - \eta_k g_n(\theta_k) + \eta_k \xi(\theta_k)$$

where the **gradient noise** is defined as:

$$\xi(\theta_k) := g_n(\theta_k) - g_{\mathcal{B}_k}(\theta_k)$$

This yields a decomposition into:

- **Drift:** $\eta_k g_n(\theta_k)$ (deterministic gradient)
- **Diffusion:** $\eta_k \xi(\theta_k)$ (stochastic perturbation)

2.7 Special Cases

1. **Full-batch gradient descent:** $\mathcal{B}_k = \{1, \dots, n\}$, so $\xi_k = 0$ (deterministic)
2. **Single-sample SGD:** $|\mathcal{B}_k| = 1$, maximum noise
3. **Constant learning rate:** $\eta_k = \eta$ for all k
4. **Decaying learning rate:** $\eta_k \rightarrow 0$ as $k \rightarrow \infty$ (e.g., $\eta_k = \eta_0 / \sqrt{k+1}$)

““latex

2.8 Properties of the Gradient Noise

The stochasticity in SGD arises from two distinct sources of randomness, which we now analyze.

2.8.1 Two Sources of Gradient Noise

1. Sampling Noise (Empirical vs. Population)

The first source of randomness comes from using a finite dataset. Even if we compute the full empirical gradient $g_n(\theta)$, it differs from the true population gradient $g(\theta)$:

$$\nu_{\text{sampling}}(\theta) := g_n(\theta) - g(\theta) = \frac{1}{n} \sum_{i=1}^n g_i(\theta) - \mathbb{E}_{(X,Y) \sim P}[\nabla_{\theta} \ell(\theta; X, Y)]$$

Since the data points are i.i.d. samples from P , we have:

$$\mathbb{E}_{\mathcal{D}_n}[\nu_{\text{sampling}}(\theta)] = \mathbb{E}_{\mathcal{D}_n} \left[\frac{1}{n} \sum_{i=1}^n g_i(\theta) \right] - g(\theta) = 0 \quad (1)$$

The covariance is:

$$\text{Cov}_{\mathcal{D}_n}[\nu_{\text{sampling}}(\theta)] = \text{Cov}_{\mathcal{D}_n} \left[\frac{1}{n} \sum_{i=1}^n g_i(\theta) \right] \quad (2)$$

$$= \frac{1}{n^2} \sum_{i=1}^n \text{Cov}_{(X_i, Y_i)}[g_i(\theta)] \quad (3)$$

$$= \frac{1}{n} \Sigma_{\text{pop}}(\theta) \quad (4)$$

where $\Sigma_{\text{pop}}(\theta) := \text{Cov}_{(X, Y) \sim P}[\nabla_{\theta} \ell(\theta; X, Y)]$ is the population gradient covariance.

2. Mini-batch Noise (Empirical vs. Batch)

The second source of randomness comes from using mini-batches instead of the full dataset:

$$\nu_{\text{batch}}(\theta) := g_{\mathcal{B}}(\theta) - g_n(\theta) = \frac{1}{B} \sum_{i \in \mathcal{B}} g_i(\theta) - \frac{1}{n} \sum_{i=1}^n g_i(\theta)$$

For uniform sampling with replacement, conditioning on the dataset \mathcal{D}_n :

$$\mathbb{E}_{\mathcal{B}}[\nu_{\text{batch}}(\theta) | \mathcal{D}_n] = \mathbb{E}_{\mathcal{B}} \left[\frac{1}{B} \sum_{i \in \mathcal{B}} g_i(\theta) \right] - g_n(\theta) = 0 \quad (5)$$

The conditional covariance is:

$$\text{Cov}_{\mathcal{B}}[\nu_{\text{batch}}(\theta) | \mathcal{D}_n] = \text{Cov}_{\mathcal{B}} \left[\frac{1}{B} \sum_{i \in \mathcal{B}} g_i(\theta) \middle| \mathcal{D}_n \right] \quad (6)$$

$$= \frac{1}{B} \Sigma_{\text{emp}}(\theta, \mathcal{D}_n) \quad (7)$$

where $\Sigma_{\text{emp}}(\theta, \mathcal{D}_n) := \frac{1}{n} \sum_{i=1}^n [g_i(\theta) - g_n(\theta)][g_i(\theta) - g_n(\theta)]^T$ is the empirical gradient covariance.

2.8.2 Total Gradient Noise

The total noise in SGD is:

$$\xi_{\text{total}}(\theta) := g_{\mathcal{B}}(\theta) - g(\theta) = \underbrace{g_{\mathcal{B}}(\theta) - g_n(\theta)}_{\nu_{\text{batch}}(\theta)} + \underbrace{g_n(\theta) - g(\theta)}_{\nu_{\text{sampling}}(\theta)}$$

However, in the SGD algorithm as defined, we use:

$$\xi(\theta) := g_n(\theta) - g_{\mathcal{B}}(\theta) = -\nu_{\text{batch}}(\theta)$$

This is because we can only compute deviations from the empirical gradient, not the unknown population gradient.

2.8.3 Covariance Structure of SGD Noise

For the noise $\xi(\theta) = g_n(\theta) - g_{\mathcal{B}}(\theta)$ used in SGD:

Proposition 1. *Under uniform mini-batch sampling with replacement, the gradient noise has the following properties:*

1. **Unbiasedness:** $\mathbb{E}_{\mathcal{B}}[\xi(\theta) | \mathcal{D}_n, \theta] = 0$

2. **Covariance:**

$$\text{Cov}_{\mathcal{B}}[\xi(\theta)|\mathcal{D}_n, \theta] = \frac{1}{B} \Sigma_{\text{emp}}(\theta, \mathcal{D}_n)$$

3. **Expected Covariance:** Taking expectation over datasets,

$$\mathbb{E}_{\mathcal{D}_n}[\text{Cov}_{\mathcal{B}}[\xi(\theta)|\mathcal{D}_n, \theta]] = \frac{1}{B} \cdot \frac{n-1}{n} \Sigma_{\text{pop}}(\theta)$$

Proof.

1. Follows immediately from $\mathbb{E}_{\mathcal{B}}[g_{\mathcal{B}}(\theta)] = g_n(\theta)$.
2. For sampling with replacement:

$$\text{Cov}_{\mathcal{B}}[\xi(\theta)|\mathcal{D}_n] = \text{Cov}_{\mathcal{B}}[g_{\mathcal{B}}(\theta)|\mathcal{D}_n] \tag{8}$$

$$= \text{Var}_{\mathcal{B}} \left[\frac{1}{B} \sum_{j=1}^B g_{I_j}(\theta) \middle| \mathcal{D}_n \right] \tag{9}$$

where I_j are i.i.d. uniform on $\{1, \dots, n\}$. Since the I_j are independent:

$$= \frac{1}{B^2} \sum_{j=1}^B \text{Var}_{I_j}[g_{I_j}(\theta)|\mathcal{D}_n] \tag{10}$$

$$= \frac{1}{B} \cdot \frac{1}{n} \sum_{i=1}^n [g_i(\theta) - g_n(\theta)][g_i(\theta) - g_n(\theta)]^T \tag{11}$$

3. Using the identity $\mathbb{E}[\Sigma_{\text{emp}}] = \frac{n-1}{n} \Sigma_{\text{pop}}$ completes the proof.

2.8.4 Key Observations

1. **Batch size scaling:** The noise covariance scales inversely with batch size B . Larger batches reduce noise.
2. **Position dependence:** The covariance $\Sigma_{\text{emp}}(\theta, \mathcal{D}_n)$ depends on the current parameters θ , making the noise *multiplicative* rather than additive.
3. **Finite sample correction:** The factor $\frac{n-1}{n}$ in the expected covariance reflects the finite dataset size. As $n \rightarrow \infty$, we recover $\mathbb{E}[\text{Cov}[\xi]] = \frac{1}{B} \Sigma_{\text{pop}}(\theta)$.
4. **Sampling without replacement:** If we sample without replacement, the covariance becomes:

$$\text{Cov}_{\mathcal{B}}[\xi(\theta)|\mathcal{D}_n] = \frac{n-B}{n-1} \cdot \frac{1}{B} \Sigma_{\text{emp}}(\theta, \mathcal{D}_n)$$

The factor $\frac{n-B}{n-1}$ represents the finite population correction.

2.8.5 Estimating the Noise Covariance

In practice, we can estimate $\Sigma_{\text{emp}}(\theta, \mathcal{D}_n)$ by:

1. Computing multiple stochastic gradients $\{g_{\mathcal{B}_j}(\theta)\}_{j=1}^m$ at the same θ
2. Estimating: $\hat{\Sigma} = \frac{B}{m-1} \sum_{j=1}^m [g_{\mathcal{B}_j}(\theta) - \bar{g}][g_{\mathcal{B}_j}(\theta) - \bar{g}]^T$

where $\bar{g} = \frac{1}{m} \sum_{j=1}^m g_{\mathcal{B}_j}(\theta)$.

This characterization of gradient noise is fundamental for understanding the continuous-time limit of SGD, as the noise covariance $\Sigma(\theta)$ directly determines the diffusion coefficient in the limiting stochastic differential equation.

3 Continuous-time limit of SGD

3.1 Heuristic derivation of the continuous-time limit

Consider the SGG update:

$$\theta_{k+1} = \theta_k - \eta_k g_{\mathcal{B}_k}(\theta_k) \quad (12)$$

where $g_{\mathcal{B}_k}(\theta_k)$ is the gradient at the k -th iteration for batch \mathcal{B}_k . Define the gradient noise as:

$$\xi(\theta_k) := g_n(\theta_k) - g_{\mathcal{B}_k}(\theta_k) \quad (13)$$

where $g_n(\theta_k)$ is the empirical gradient. The new update can be written as:

$$\Delta\theta_{k+1} := \theta_{k+1} - \theta_k = -\eta_k g_{\mathcal{B}_k}(\theta_k) + \eta_k \xi(\theta_k) \quad (14)$$

The 1-sample gradient expectation is:

$$\mathbb{E}[g_i(\theta_k)] = \frac{1}{n} \sum_{i=1}^n g_i(\theta_k) = g_n(\theta_k) \quad (15)$$

Hence the expectation of the noise is zero. The 1-sample noise covariance is:

$$\Sigma(\theta_k) := \frac{1}{n} \sum_{i=1}^n [g_i(\theta_k) - g_n(\theta_k)][g_i(\theta_k) - g_n(\theta_k)]^T \quad (16)$$

where $g_i(\theta_k)$ is the gradient at the i -th data point. The batch noise covariance is:

$$\Sigma_{\text{batch}}(\theta_k) := \frac{1}{B^2} \sum_{i \in \mathcal{B}_k} \text{Cov}(\xi(\theta_k)) = \frac{1}{B} \Sigma(\theta_k) \quad (17)$$

Assume that the learning rate is independent of k . Assume that batch size is large enough to apply CLT and much smaller than the dataset size i.e. $B \ll n$ (**independent batches?**) and $\frac{B}{n} \ll 1$. Assume that number of samples is large enough to apply CLT. Assume that samples are i.i.d. and that the gradient update (and maybe gradient noise?) have finite variance. Under these assumptions, we can apply the CLT to the batch gradient and get:

$$\xi(\theta_k) \sim \mathcal{N}(0, \frac{1}{B} \Sigma(\theta_k)) \quad (18)$$

Proper scaling of the noise is crucial to apply CLT and get the continuous-time limit. Physical intuition: the average displacement of the noise is 0 but the typical displacement of the noise is in \sqrt{t} where $t = N\eta$ for N steps. More specifically

$$\text{Var}\left(\sum_{k=1}^N \eta \xi_k\right) = N\eta^2 \Sigma(\theta_k) = t\eta \Sigma(\theta_k) \quad (19)$$

$$\text{std}\left(\sum_{k=1}^N \eta \xi_k\right) = \sqrt{t\eta} \sqrt{\Sigma(\theta_k)} \quad (20)$$

with $N = \frac{t}{\eta}$ for N steps. We see that as $\eta \rightarrow 0$, the typical displacement of the accumulated noise goes to 0. If we use the scaling $\sqrt{(\eta)}$ for the noise, we get:

$$\text{Var}\left(\sum_{k=1}^N \sqrt{(\eta)} \xi_k\right) = N\eta \Sigma(\theta_k) = t\Sigma(\theta_k) \quad (21)$$

$$\text{std}\left(\sum_{k=1}^N \sqrt{(\eta)} \xi_k\right) = \sqrt{t} \sqrt{\Sigma(\theta_k)} \quad (22)$$

With $N = t/\eta$ for N steps. We see that the noise has proper scaling as the typical displacement is in \sqrt{t} and does not depend on η . The continuous-time limit of SGD is:

$$d\theta_t = -\nabla L_n(\theta_t)dt + \sqrt{\frac{\eta}{B}\Sigma(\theta_t)}dW_t \quad (23)$$

where W_t is a Wiener process. Note that a more rigorous derivation would use the Lindeberg condition and the Donsker's theorem.

3.2 Why do we care: the tempered posterior distribution

Using the machinery of Fokker-Planck equation, we can derive the stationary distribution of the SGD process as the tempered posterior distribution. The Fokker-Planck equation is:

$$\frac{\partial p(\theta, t)}{\partial t} = \nabla \cdot \left(\nabla L_N(\theta)p(\theta, t) + \frac{1}{2} \frac{\eta}{B} \Sigma(\theta) \nabla p(\theta, t) \right) \quad (24)$$

Assume that the noise covariance matrix is positive definite. Let the probability current be:

$$j(\theta, t) = \nabla L_N(\theta)p(\theta, t) + \frac{1}{2} \frac{\eta}{B} \Sigma(\theta) \nabla p(\theta, t) \quad (25)$$

The stationarity condition implies that the probability current is divergence free:

$$\nabla \cdot j(\theta, t) = 0 \quad (26)$$

Assume a detailed balance condition i.e. $j(\theta, t) = 0$. This implies that the stationary distribution satisfies:

$$\nabla L_N(\theta)p(\theta) + \frac{1}{2} \frac{\eta}{B} \Sigma(\theta) \nabla p(\theta) = 0 \quad (27)$$

Assume that the noise covariance matrix is constant $\Sigma(\theta) = \Sigma$, positive and definite. Then solving the latter differential equation gives us the stationary distribution as the tempered posterior distribution:

$$p_\infty(\theta) = \frac{1}{Z} \exp(-\beta \Sigma^{-1} L_N(\theta)) \quad (28)$$

where $\beta = \frac{B}{2\eta}$ is the inverse temperature. In the case where the noise covariance matrix is not constant but still positive and definite, we get:

$$p_\infty(\theta) = \frac{1}{Z} |\Sigma(\theta)|^{-1/2} \exp\left(-\frac{2B}{\eta} L_N(\theta)\right) \quad (29)$$

where Z is the partition function. In that case, the stationary distribution is given by the tempered posterior i.e. up to the temperature parameter SGD is similar to Bayesian inference. This result has been investigated in more details by Mandt. The assumptions in this derivation are not realistic since the noise covariance is singular i.e. not definite positive. In fact, deep neural networks are highly degenerate and we also know that SGD favours flat minima.

3.3 List of assumptions

Assumptions for the continuous-time limit of SGD:

- GN satisfies the Lindeberg condition: necessary for CLT to apply. Unclear how much this matters.

- Gradient noise has finite variance: necessary for CLT to apply. An alternative would be to consider heavy tailed noise. It is unclear how much this matters. This is debated in the literature and it would be useful to understand better how much this matters by reviewing the literature.
- η_k is independent of k : Uncertain but technically should not be necessary as we could have $\eta_k \rightarrow dt$ after some updates and then take the continuous-time limit if we only care about the local behaviour of SGD. Unclear how much this matters, in practice it is fairly common to have time-varying learning rate.
- $\eta \rightarrow dt \ll 1$: makes the learning rate small which is central to take the continuous-time limit. This is an important difference with the practical usage of SGD. Although it might be possible to approximate the discrete time dynamics with a continuous time dynamics by using a central flows which essentially add the average jittered caused by the discrete learning rate to the continuous time dynamics. See for example this paper: Understanding Optimization in Deep Learning with Central Flows .
- $B \ll n$ and $\frac{B}{n} \ll 1$: Have enough batch and data samples to apply CLT
- No auto-correlation in the noise: makes the noise white which enables to take the continuous-time limit. Alternative is having coloured noise. This could happens for example if the batches are not independent, which could be the case, intuitively if the batches, are large. Another example is if the sampling of the batches is not without replacement. In this case it seems that the noise is anti-correlated as argued in this paper. Anti-correlation could be important to generalization as the latter paper argue that it biases SGD further towards flat minima. Momentum can also introduces auto-correlation in the noise. If the noise has some auto-correlation but decaying with time, the functional CLT still can be applied and we can get to the continuous-time with a Wiener process. If not we have to consider a fractional Brownian motion. **Idea**: Could we have a fractional Levy process with heavy tail noise and long-range auto-correlation? Is it the case that in practice long-range auto-correlation matter? **Note**: Fractional Brownian or Levy process could be another connection with fractality in the training process. **Question**: if heavy tailed noise is considered, does it make sense to consider auto-correlation in the noise?
- Samples are i.i.d.: necessary for CLT to apply
- For now we are not considering momentum. Unsure how it changes the analysis. For example our analysis contrasts with Adam which has a momentum and a second order correction of past gradient.

Assumptions for the stationary distribution of SGD being the tempered posterior distribution:

- The noise covariance matrix is positive and definite: seems important to solve the Fokker-Planck equation as we invert the noise covariance matrix to get the stationary distribution. This assumption is violated in practice since deep neural networks are highly degenerate.
- The noise covariance matrix is constant: this is not necessary but makes the calculation easier.
- Detailed balance condition: helps with getting a first ODE for the stationary distribution which enables to get the tempered posterior distribution.
- According to the Helmholtz decomposition, the current can be decomposed into a curl-free part and a divergence-free part. If the noise is anisotropic then we have $\nabla \times j(\theta, t) \neq 0$ In other words, the current is not curl-free. The detailed balanced condition does not hold

in that condition. When the diffusion matrix $\Sigma(\theta)$ depends on θ , the stationary dynamics of SGD is an *out-of-equilibrium* steady state with a non-zero *solenoidal probability current*. Probability circulates in parameter space rather than remaining static, so detailed balance—and the simple Boltzmann form $p \propto e^{-\beta L_N}$ —no longer hold. **Note:** I am still a bit confused about the physical intuition about this. How should i think about this probability mass orbiting along some level set of the loss?

3.4 Ranking assumptions

After a deep research with chatGPT here and with elicit here this is my current ranking of the assumptions:

- Application of CLT from assuming gaussian gradient noise (finite variance in particular, likely that noise is not gaussian but Levy). Also an interesting point is that the heavy tailness could be anisotropic i.e. be sensitive to the degeneracies in the loss landscape.
- Detailed balance condition from constant noise covariance. In particular this means we must consider a curl component to the probability current.
- **Note:** I am a bit confused about the ranking of the next two items.
- Small learning rate: unrealistic but see central flows paper. Would be useful to look more into the edge of stability phenomena literature.
- No auto-correlation in the noise: this will be violated with momentum in practice (but maybe that's ok as long as it's not long-range auto-correlation?) **Note:** I am a bit confused about this one. I'd like to understand better violation of long-range auto-correlation. Modification might include: fractional Levy, generalized Langevin dynamics, etc. Need colored noise models with memory. Sampling without replacement is common introducing some auto-correlation.
- Batch size large enough but small compared with the dataset size seems realistic especially in large scale training.

4 SDE on DLNs

For input-output pairs (x^μ, y^μ) and weight matrices $W^1 \in \mathbb{R}^{d_1 \times d_0}$, $W^2 \in \mathbb{R}^{d_2 \times d_1}$, let the empirical (squared-error) loss be

$$L_N(W^1, W^2) = \frac{1}{N} \sum_{\mu=1}^N \|y^\mu - W^2 W^1 x^\mu\|_2^2.$$

Expanded form. With Einstein summation over repeated indices (i, j, k, l, p, r) ,

$$L_N = \frac{1}{N} \sum_{\mu=1}^N \left(y_i^\mu y_i^\mu - 2 y_i^\mu W_{ik}^2 W_{kj}^1 x_j^\mu + x_l^\mu W_{kl}^1 W_{ki}^2 W_{ip}^1 W_{pr}^1 x_r^\mu \right)$$

Gradient w.r.t. the first layer W^1 .

$$g_{mn}^1 = \frac{\partial L_N}{\partial W_{mn}^1} = \frac{2}{N} \sum_{\mu=1}^N x_n^\mu W_{im}^2 (W_{ik}^2 W_{kj}^1 x_j^\mu - y_i^\mu)$$

Matrix form (with sample matrices $X = [x^1 \dots x^N]$ and $Y = [y^1 \dots y^N]$):

$$\nabla_{W^1} L_N = \frac{2}{N} W^{2\top} (W^2 W^1 X - Y) X^\top$$

Gradient w.r.t. the second layer W^2 .

$$g_{mn}^2 = \frac{\partial L_N}{\partial W_{mn}^2} = \frac{2}{N} \sum_{\mu=1}^N (W_{mk}^2 W_{kj}^1 x_j^\mu - y_m^\mu) (W_{nj}^1 x_j^\mu)$$

Matrix form:

$$\nabla_{W^2} L_N = \frac{2}{N} (W^2 W^1 X - Y) (W^1 X)^\top$$

5 TODO

- Explain why we care about continuous-time limit and connection with the tempered posterior distribution (also mention invertability of the noise covariance)
- Have a better understanding of the detailed balance condition and links with the curl of the probability current (leading to a non stationary distribution but to an orbit instead)
- Explain the rigorous mathematical derivation of the continuous-time limit of SGD.
- Clarify the importance of assuming finite learning rate. Read the central flows paper. Edge of stability implies that naive continuous time limit is not valid.
- Relatedly to a review of edge of stability phenomena literature.
- Check proof of stationary distribution for the case where the noise covariance is not constant.
- Explain connection SGD with SLT and spectrum between continuous limit with all assumptions and discrete SGD
- Include momentum considerations.

6 Questions that we could investigate

We could start with replacing the gaussian noise assumption with a heavy tailed noise assumption. One generic idea would be to test the assumption of heavy tailness during the training dynamics. But it seems the sort of things that has already been done. What would be interesting would be to test heavy tailness if we think of a specific hypothesis about why this is plays an important role in the training dynamics. Under what conditions can we compute the stationary distribution of SGD with heavy tailed noise?

See deep research literature on heavy tailed noise from chatGPT here. Other mention: we could also look at the learning at parity paper and compare heavy tail with non heavy tail noise. One key takeaway is that there is evicence for heavy tailness, that it might be important to escape faster from shallow deep minima and saddle point and find broader flat minima which help with generalization. An implication of Levy flight will be that the probability distribution will be differed from the tempered posterior distribution and will be closer to a non equilibrium distribution.

We could investigate the saddle to saddle dynamics in deep linear networks. We could in particular look at the importance of heavy tailness to favour degenerate minima. Using analytic results from deep linear networks, we could probably compare Gaussian SDE with Levy SDE and see which ones better predicts the training dynamics and the distribution after training. Could we derive some theoretical results about the stationary distribution in deep linear networks using Levy SDE?

We could look at toy model of superposition and look at phase transition between geometric states and amplitude of noise. One hypothesis could be that phase transition happened when some sufficiently high amplitude noise.

More interestingly we could look at grokking in modular addition. One prediction would be that the noise helps the model escape toward broader and flatter basins that generalize better and that this is correlated with the heavy tailness of the noise. Would need to detail the experimental protocol a bit more.

6.1 Experiments on deep linear networks

See chat with Opus here.

7 Questions

How does gradient clipping affects SGD noise? Does not prevent heavy tail because gradient clipping is applied to batch gradient and not to the noise. So noise could be unbounded because of large fluctuations coming from difference between batch gradient and empirical gradient.

What is the difference between the tempered posterior distribution and the Bayesian posterior?

What is the intuition behind the curl of the probability current? How much does it deviate from the tempered posterior distribution/Bayes posterior? Also how much does it matter that the current is not curl-free?

How important is considering SGD vs ADAM?

Is there a tension between having a time-fractional and a parameter-fractional fokker planck operator? Perhaps there are different regime in which subdiffusive and superdiffusive behavior matters more and we don't need to have both operators at the same time.

8 Research Proposal

8.1 Motivation

Stochastic Gradient Descent (SGD) is the central algorithm driving deep learning. While it is often modeled mathematically using simplified stochastic differential equation (SDE) approximations, these standard approximations typically rely on assumptions that are empirically invalid during actual deep neural network training. Specifically, standard assumptions include:

- **Gaussian Gradient Noise:** Empirically violated due to observed heavy-tailed gradient distributions.
- **Detailed Balance (Curl-Free Probability Current):** Empirically violated due to noise anisotropy, leading to a non-zero curl in probability currents.
- **Small Learning Rate Approximation:** Empirically violated by observed edge-of-stability phenomena, including oscillatory dynamics.

These violations imply that the distribution of SGD trajectories significantly deviates from the classical tempered posterior (Gibbs distribution) predicted by naive SDE models.

To address these shortcomings, we propose a refined SDE model incorporating key empirically supported phenomena:

- **Gradient Term:** Standard loss gradient.
- **Central Flow Term:** Capturing oscillatory dynamics and sharpness-induced flows from edge-of-stability phenomena.
- **Levy Noise with Anisotropic Covariance:** Capturing heavy-tailed and anisotropic gradient noise empirically observed during training.

The primary motivation from an AI alignment perspective is to better characterize the inductive biases of SGD. Accurately sampling the SGD trajectory distribution is critical for understanding the robustness of safety properties. Such an approach will aid in determining the likelihood of model trajectories remaining aligned and safe upon further training, as opposed to bifurcating toward deceptive or misaligned models (or being already deceptive if most models sampled on the degenerate set of loss minimizers contain many deceptive models).

8.2 Proposed Experiment

Objective: Validate our refined SDE model by comparing the distributions it generates against actual SGD trajectory distributions and those predicted by naive Gaussian-based SDE approximations.

Setup: Deep Linear Networks (DLNs) We focus on DLNs, as they offer analytically tractable solutions and clear characterizations of degenerate minima. We also know that it is possible to have interesting saddle-to-saddle dynamics in DLNs and the more lazy NTK regime. Furthermore, we know from the saddle to saddle dynamics that SGD favour low rank solutions during training. It would be interesting to understand the influence different terms in the SDE model on the saddle to saddle dynamics.

Experimental Protocol.

Initialization:

- **Saddle to saddle initialization:** Initialize parameters with a small covariance such that the initial parameters are close to a saddle point.
- **More lazy NTK regime:** Initialize parameters with a large covariance such that initial parameters can be close to a global minimum.

Data generation:

- **Generate synthetic data as random matrices:** can try normal or heavy tailed random matrices.

Trajectory Generation:

- **SGD Trajectories:** Simulate true SGD with mini-batches of synthetic data.
- **Langevin SDE:** Standard Brownian-motion-driven Langevin dynamics with isotropic (or anisotropic) Gaussian noise and no central flow.
- **Refined SDE:** Proposed SDE with central flow, Levy-stable noise (heavy tails), and anisotropic covariance (rotational current) reflecting different assumptions behind SGD.

Distribution Analysis:

- Estimate some distance between distributions obtained from SGD and the various SDE models.
- Quantify the shape of distributions (e.g. covariance, skewness and higher moments).

Distance Metrics for Comparison:

- **Wasserstein-2 (Earth Mover’s Distance):** Measure of distributional differences.
- **Maximum Mean Discrepancy (MMD):** Non-parametric kernel-based measure suitable for finite samples.
- **Approximate KL Divergence:** can we estimate some empirical version?
- Pb: not sure if any of this is tractable

Tracking escape dynamics from saddle points. We could test whether heavy tailness of SDE with Levy process (using alpha stable distribution) helps to predict escaping from saddle points compared with a Gaussian SDE. Process:

- Detect saddle points (for at loss plateaus, monitor gradient norm and hessian eigenvalues as one should be negative).
- Collect gradient noise around saddle points from SGD and SDE models (naive vs refined)
- Track tail index indicator e.g. Hill estimator. Hypothesis: Heavy tailness correlates with faster escape from saddle points.
- We could look at escape time from saddle point vs heavy tailness and rotational current. Need to estimate escape time from saddle point but not sure how.

8.3 Hypothesis

Overall main Hypothesis: The refined SDE model sampler produces trajectory distributions that significantly better match actual SGD empirical distributions along degenerate minima than naive Gaussian-based SDE samplers.

Evaluation:

- Compute pairwise distributional distances: Actual SGD \leftrightarrow Naive SDE, Actual SGD \leftrightarrow Refined SDE.
- Statistically assess the hypothesis (e.g., via bootstrapping or significance testing).