

# The long-run behaviour of SGD on nonconvex functions

Jérôme MALICK

CNRS, Grenoble, France

Based on joint work with

Waïss AZIZIAN

Franck IUTZELER Panayotis MERTIKOPOULOS

SIAM Optimization

Edinburgh

June 2026

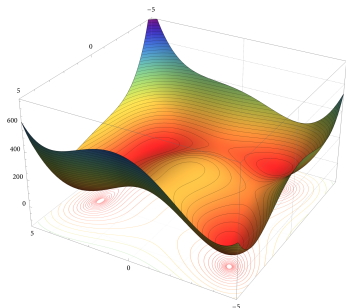


## SGD does not converge...

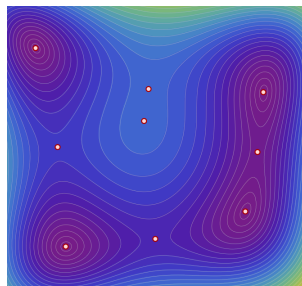
Stochastic Gradient Descent with constant step-size  $\eta > 0$

$$x_{t+1} = x_t - \underbrace{\eta}_{\text{step-size}} \hat{g}_t \quad \text{with} \quad \hat{g}_t = \nabla f(x_t) + \underbrace{Z(x_t; \omega_t)}_{\text{zero-mean noise}}$$

**Illustration:**



Himmelblau function



$$f(x, y) = (x^2 + y - 11)^2 + (x + y^2 - 7)^2$$

**Video:** Run 25 SGD (for each trajectory: plot only last 20 iterates)

**Observation:** SGD ~~converges~~ concentrates to local minimizers

## What is known

SGD with constant step-size:  $x_{t+1} = x_t - \eta \left[ \nabla f(x_t) + Z(x_t; \omega_t) \right]$

- What is known, when  $f$  is convex
  - $f$  strongly convex: SGD converges in a ball around the minimizer [Polyak, 1987]
  - $f$  convex: average of iterates asymptotically optimal [Polyak, Juditsky, 1992]

- What is known, when  $f$  is **non-convex**

- In average, SGD is close to criticality [Lan, 2012] [Ghadimi, Lan, 2013]

$$\mathbb{E} \left[ \frac{1}{T} \sum_{t=0}^{T-1} \|\nabla f(x_t)\|^2 \right] = O\left(\frac{1}{\sqrt{T}}\right)$$

- Almost surely, SGD is not stuck in saddle points [Pemantle, 1990; Benaïm, Hirsch, 1995]

- What we would like to know

**Question:** which points are more likely to be observed — and by how much?

# Our results

## What is the Long-Run Distribution of Stochastic Gradient Descent? A Large Deviations Analysis

Waïss Azizian<sup>1</sup> Franck Iutzeler<sup>2</sup> Jérôme Malick<sup>1</sup> Panayotis Mertikopoulos<sup>3</sup>

### Abstract

In this paper, we examine the long-run distribution of stochastic gradient descent (SGD) in general, non-convex problems. Specifically, we seek to understand which regions of the problem's state space are more likely to be visited by SGD, and by how much. Using an approach based on the theory of large deviations and randomly perturbed dynamical systems, we show that the long-run dis-

architectures – from large language models to reinforcement learning and recommender systems. This phenomenal success is largely owed to the method's simplicity: given a smooth function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  and the associated optimization problem

$$\min_{x \in \mathbb{R}^d} f(x) \quad (\text{Opt})$$

the SGD algorithm is given by the simple update rule

$$x_{n+1} = x_n - \eta \hat{g}_n \quad (\text{SGD})$$

### Theorem (Azizian et al., 2024; very informal)

- 1 The **critical set** of  $f$  is visited exponentially more often than any **non-critical** set
- 2 Critical points are visited according to a **Boltzmann–Gibbs** distribution with temp.  $\eta$
- 3 **Minimizers** are visited exponentially more often than **non-minimizers**
- 4 The **ground state** of the problem is visited exponentially more often than any other min.

Goal of this talk: make these results more precise

## The way to these results

- **Key object:** asymptotic distributions of  $x_t$ 
  - **Occupation measure:** average time spent by the iterates in a set of interest  $\mathcal{B}$

$$\mu_t(\mathcal{B}) = \mathbb{E} \left[ \frac{1}{t} \sum_{s=1}^t \mathbf{1}_{\{x_s \in \mathcal{B}\}} \right]$$

- Its limit is an **invariant measure**:  $x_t \sim \mu_\infty \implies x_{t+1} \sim \mu_\infty$

**Question:** where does  $\mu_\infty$  concentrate?

## The way to these results

- **Key object:** asymptotic distributions of  $x_t$ 
  - **Occupation measure:** average time spent by the iterates in a set of interest  $\mathcal{B}$

$$\mu_t(\mathcal{B}) = \mathbb{E} \left[ \frac{1}{t} \sum_{s=1}^t \mathbf{1}_{\{x_s \in \mathcal{B}\}} \right]$$

- Its limit is an **invariant measure**:  $x_t \sim \mu_\infty \implies x_{t+1} \sim \mu_\infty$

**Question:** where does  $\mu_\infty$  concentrate?

- **Approach:** Randomly perturbed dynamical systems [Freidlin Wenzel 1998, Kifer 1980]  
→ estimate the probability of rare events, such as SGD escaping a local minimum
- **Refinement:** We adapt this theory, with three main challenges:  
Lack of compactness + Noise models (finite sum) + Discrete-time dynamics

## What we are not doing

- **Stochastic approximation:**

$$x_{t+1} = x_t - \eta_t [\nabla f(x_t) + Z(x_t; \omega_t)] \quad \text{with } \eta_t \text{ vanishing}$$

Convergence to local minima [Bertsekas & Tsitsiklis '00, Mertikopoulos et al. '20]  
**but** no information about which one

- **Langevin dynamics:**

$$x_{t+1} = x_t - \eta \nabla f(x_t) + \sqrt{\eta} \mathcal{N}(0, \sigma^2)$$

Convergence of the distribution of the iterates to  $e^{-f/\sigma^2}$  [Raginsky et al. '17]  
**but** scaling of the noise differs from SGD + we want non-gaussian noise

- **Continuous-time dynamics** (Gradient flow, SDE):

$$dX_t = -\nabla f(X_t) dt + \sqrt{\eta \text{cov}(Z(X_t; \cdot))} dW_t$$

Results on stationary distribution + **finite-time** approximation of SGD [Li et al. '17]  
**but** unsuitable for concentration questions

Existing analysis do not carry over

## Blanket assumptions

### Assumptions on the objective

- Smoothness:  $\|\nabla f(x') - \nabla f(x)\| \leq L\|x' - x\|$
- Coercivity:  $f(x) \rightarrow \infty$  as  $\|x\| \rightarrow \infty$  (+ a bit more, like coercive gradient:  $\|\nabla f(x)\| \rightarrow \infty$ )
- Critical set:  $\text{crit}(f) = K_1 \cup \dots \cup K_p$  with smoothly connected components  $K_i$

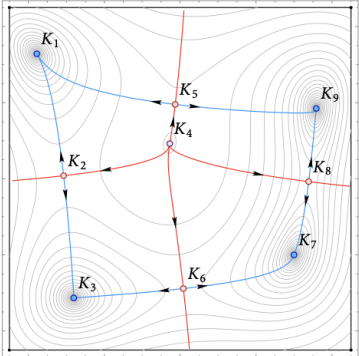
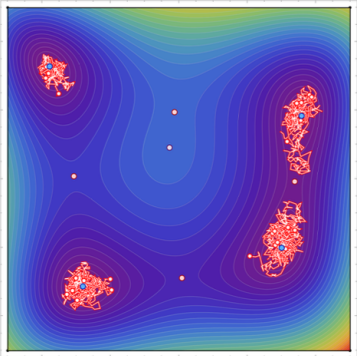
### Assumptions on the noise

- Properness:  $\mathbb{E}[Z(x; \omega)] = 0$  and  $\text{cov}(Z(x; \omega)) \succ 0$
- Smooth growth:  $Z$  is  $C^2$ -smooth and  $\|Z(x; \omega)\| \leq C(1 + \|x\|)$
- Sub-Gaussian:  $\log \mathbb{E}[e^{\langle v, Z(x; \omega) \rangle}] \leq \frac{\sigma^2}{2} \|v\|^2$  ( $\Rightarrow$  bounds  $\mathbb{P}(\|Z(x; \omega)\| > z) \leq \exp(-az^2 + b)$ )

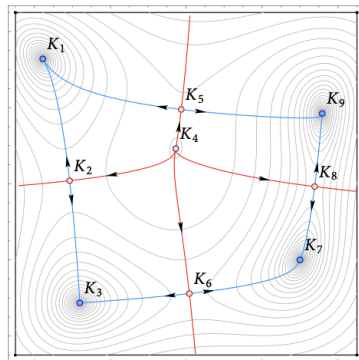
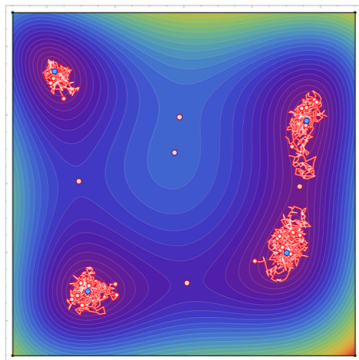
**Example:** Regularized ERM:  $f(x) = \frac{1}{m} \sum_{i=1}^m \ell(x; \xi_i) + \frac{\lambda}{2} \|x\|^2$  (with  $\ell$  smooth, Lipschitz)

Sample one (or several) data-point  $\xi_\omega$ , and undergo  $Z(x; \omega) = \frac{1}{m} \sum_{i=1}^m \nabla \ell(x; \xi_i) - \nabla \ell(x; \xi_\omega)$

# Main idea of the proof



## Main idea of the proof



the long-run distribution of SGD resembles  
a finite-state Markov chain on the critical components

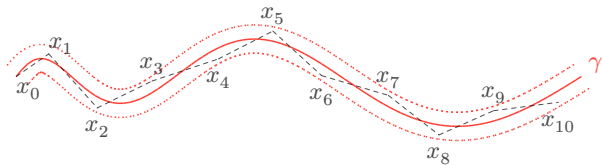
- Make this precise (difficult) – 100 pages of maths in one slide: next slide
- Get the results for SGD from the induced chain (easier) – the following slide

## Key ingredient: large deviations for SGD

For  $T > 0$

Continuous path  $\gamma : [0, T] \rightarrow \mathbb{R}^d$

$$\mathbb{P}(\text{SGD} \approx \gamma) = ?$$



**Proposition (Azizian et al., 2024; informal)**

For small  $\eta$ : 
$$\mathbb{P}(\text{SGD on } [0, T/\eta] \approx \gamma) \approx \exp\left(-\frac{S_T[\gamma]}{\eta}\right)$$

C. Gen. Fct:  $\mathcal{H}(x, v) = \log \mathbb{E} [e^{\langle v, Z(x; \omega) \rangle}]$

Conj.:  $\mathcal{L}(x, v) = \mathcal{H}^*(x, -v - \nabla f(x))$

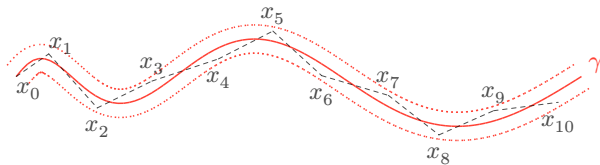
Action:  $S_T[\gamma] = \int_0^T \mathcal{L}(\gamma_t, \dot{\gamma}_t) dt$

## Key ingredient: large deviations for SGD

For  $T > 0$

Continuous path  $\gamma : [0, T] \rightarrow \mathbb{R}^d$

$$\mathbb{P}(\text{SGD} \approx \gamma) = ?$$



**Proposition (Azizian et al., 2024; informal)**

For small  $\eta$ : 
$$\mathbb{P}(\text{SGD on } [0, T/\eta] \approx \gamma) \approx \exp\left(-\frac{S_T[\gamma]}{\eta}\right)$$

Gaussian noise ( $Z(x; \omega) \sim \mathcal{N}(0, \sigma^2 I_d)$ )

C. Gen. Fct:  $\mathcal{H}(x, v) = \log \mathbb{E} [e^{\langle v, Z(x; \omega) \rangle}]$   $\mathcal{H}(x, v) = \frac{\sigma^2}{2} \|v\|^2$

Conj.:  $\mathcal{L}(x, v) = \mathcal{H}^*(x, -v - \nabla f(x))$   $\mathcal{L}(x, v) = \frac{\|v + \nabla f(x)\|^2}{2\sigma^2}$

Action:  $S_T[\gamma] = \int_0^T \mathcal{L}(\gamma_t, \dot{\gamma}_t) dt$   $S_T[\gamma] = \frac{1}{2\sigma^2} \int_0^T \|\dot{\gamma}_t + \nabla f(\gamma_t)\|^2 dt$

## From transitions to energy

We get (...) the **transition cost** from one critical region to another

$$B_{i,j} = \min\{S_T[\gamma] : \gamma(0) \in K_i, \gamma(T) \in K_j, T > 0\}$$

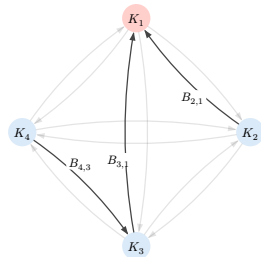
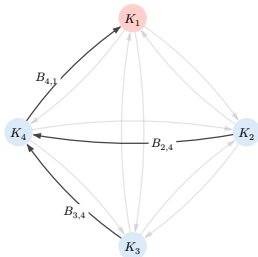
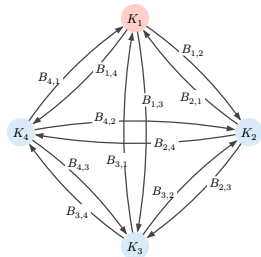
Then, using exact formulas for finite-state Markov chains, we get (...):

### Lemma (very informal)

The invariant measure of SGD restricted to  $\{K_1, \dots, K_p\}$  is, for  $\eta > 0$  small enough

$$\pi(i) \propto \exp\left(-\frac{E_i}{\eta}\right)$$

where **energy** of  $K_j$ :  $E_j = \min\left\{\sum_{j \rightarrow k \in T} B_{j,k} : T \text{ spanning tree pointing to } i\right\}$



## From energy to invariant distribution

### Theorem (Azizian et al., 2024; less informal)

Given  $\mathcal{U}_i$  neighborhoods of  $K_i$ , and  $\eta > 0$  small enough:

- ① **Concentration near critical points:** there is some  $c > 0$  s.t.

$$\mu_\infty \left( \bigcup_{i=1}^p \mathcal{U}_i \right) \geq 1 - e^{-c/\eta} \quad \text{for some } c > 0$$

- ② **Boltzmann–Gibbs distribution:** for all  $i$ ,

$$\mu_\infty(\mathcal{U}_i) \propto \exp \left( -\frac{E_i + O(\varepsilon)}{\eta} \right)$$

- ③ **Saddle-point avoidance:** if  $K_i$  is a saddle, then there is  $K_j$  local minimum with  $E_j < E_i$

$$\frac{\mu_\infty(\mathcal{U}_i)}{\mu_\infty(\mathcal{U}_j)} \approx e^{-(E_i - E_j)/\eta}$$

- ④ **Ground state concentration:**  $\mathcal{U}_0$  neighborhood of the **ground state**  $K_0 = \operatorname{argmin}_i E_i$

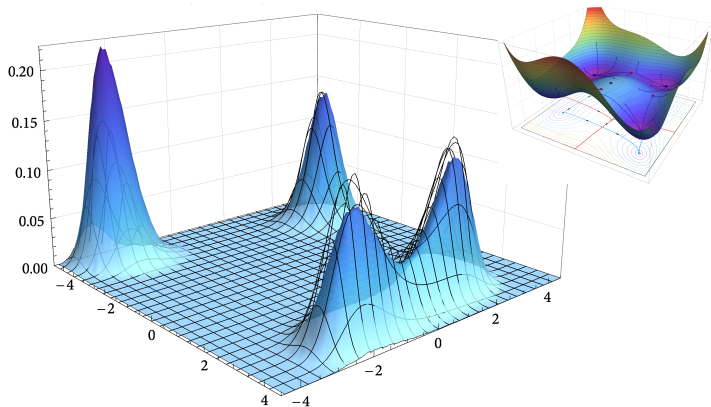
$$\mu_\infty(\mathcal{U}_0) \geq 1 - e^{-c/\eta} \quad \text{for some } c > 0$$

## Illustration: Gaussian noise

Assume  $Z(x; \omega) \sim \mathcal{N}(0, \sigma^2 I_d)$

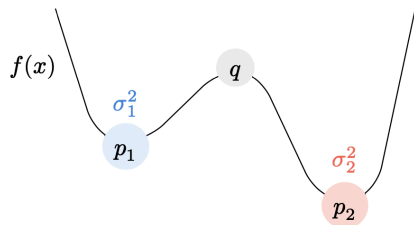
Compute  $E_i = \frac{2f(K_i)}{\sigma^2} \implies \mu_\infty(\mathcal{U}_i) \approx \exp\left(-\frac{2f(K_i)}{\sigma^2 \eta}\right)$  and  $K_0 = \text{global min.}$

**Illustration:** run 1000 SGD on Himmelblau, observe distribution of iterates at  $T = 1000$



Energy landscape (wireframe) vs. empirical distribution of SGD (blue)

Remark: ground state = global minimizer ?



$$E_1 = \frac{f(q) - f(p_2)}{\sigma_2^2} \quad \text{and} \quad E_2 = \frac{f(q) - f(p_1)}{\sigma_1^2}$$

$$\sigma_1 \text{ small enough} \implies E_1 < E_2 \implies \mu_\infty(p_1) \ll \mu_\infty(p_2)$$

In general, minimizer of the energy  $\neq$  minimizer of the function

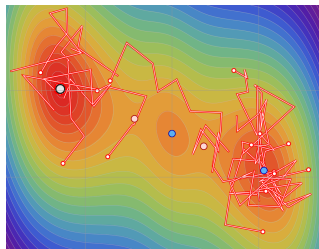
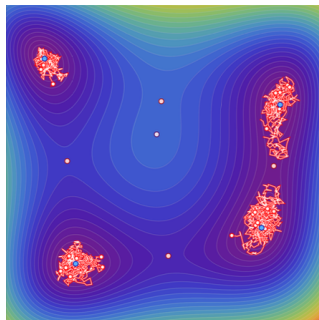
# Let's sum up and look forward

## Main take-aways

- SGD concentrates near local minima  
in particular, near those that minimize the energy
- We developed a new theoretical framework  
to analyze SGD with the help of large deviations

## Perspectives

- Beyond SGD ? (inertia, Adam...)
- Beyond minimization ? (min-max, equilibrium...)
- Beyond distribution ? Eg. Hitting time [Azizian et al '25]



thank you all 😊