

Stochasticity of Deterministic Gradient Descent: Quantitative Local Min Escape in Multiscale Landscape

Molei Tao  School of Mathematics, Georgia Tech, USA

November 12, 2023

Caltech ACM

Machine Learning & Applied Computational Math

Machine Learning & Applied Computational Math

→ ML for problems in computing, sciences & engineering

Machine Learning & Applied Computational Math

- ML for problems in computing, sciences & engineering
- ← how ACM helps design and **analyze**
optimization, sampling, and **deep learning practices**

Machine Learning & Applied Computational Math

→ ML for problems in computing, sciences & engineering

← how ACM helps design and **analyze**

optimization, sampling, and **deep learning practices**

*common belief: **large learning rate** is good*

Machine Learning & Applied Computational Math

→ ML for problems in computing, sciences & engineering

← how ACM helps design and **analyze**

optimization, sampling, and **deep learning practices**

*common belief: **large learning rate** is good*

$\min f(x)$

$$x_{k+1} = x_k - h \nabla f(x_k)$$

Gradient Descent

Machine Learning & Applied Computational Math

→ ML for problems in computing, sciences & engineering

← how ACM helps design and **analyze**

optimization, sampling, and **deep learning practices**

*common belief: **large learning rate** is good*

1st thought:

min $f(x)$

$$x_{k+1} = x_k - h \nabla f(x_k)$$

Gradient Descent

$$\updownarrow x_i \approx x(ih)$$

$$\dot{x} = -\nabla f(x)$$

Gradient Flow

Machine Learning & Applied Computational Math

→ ML for problems in computing, sciences & engineering

← how ACM helps design and **analyze**

optimization, sampling, and **deep learning practices**

*common belief: **large learning rate** is good*

1st thought:

$\min f(x)$

$$x_{k+1} = x_k - h \nabla f(x_k)$$

Gradient Descent

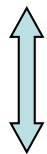
$x(T)$ with fixed T



$[T/h]$ steps



large h ?



$$x_i \approx x(ih)$$

$$\dot{x} = -\nabla f(x)$$

Gradient Flow

Machine Learning & Applied Computational Math

→ ML for problems in computing, sciences & engineering

← how ACM helps design and **analyze**

optimization, sampling, and **deep learning practices**

common belief: large learning rate is good

1st thought:

$$\min f(x)$$

$$x_{k+1} = x_k - h \nabla f(x_k)$$

Gradient Descent

x(T) with fixed T



[T/h] steps



large h?

$$x_i = x_0 + i h \nabla f(x)$$

$$\dot{x} = -\nabla f(x)$$

Gradient Flow

Machine Learning & Applied Computational Math

→ ML for problems in computing, sciences & engineering

← how ACM helps design and **analyze**

optimization, sampling, and **deep learning practices**

*common belief: **large learning rate** is good*

not just faster, but implicitly bias toward (desirable) global structures

???

Machine Learning & Applied Computational Math

→ ML for problems in computing, sciences & engineering

← how ACM helps design and **analyze**

optimization, sampling, and **deep learning practices**

*common belief: **large learning rate** is good*

not just faster, but implicitly bias toward (desirable) global structures

???

→ better training & better test accuracies

training

$$\min_x f(x) := \sum_i d(\text{output}_i, g(x, \text{input}_i))$$

model (e.g., neural network) parameters

training data (known)

training

$$\min_x f(x) := \sum_i d(\text{output}_i, g(x, \text{input}_i))$$

model (e.g., neural network) parameters

training data (known)

trapped in **local min** \leftrightarrow **suboptimal** training accuracy

model (e.g., neural network) parameters

training

$$\min_x f(x) := \sum_i d(\text{output}_i, g(x, \text{input}_i))$$

training data (known)

large h? this talk

trapped in **local min** \leftrightarrow **suboptimal** training accuracy

training

$$\min_x f(x) := \sum_i d(\text{output}_i, g(x, \text{input}_i))$$

model (e.g., neural network) parameters

training data (known)

large h? this talk

trapped in local min \leftrightarrow suboptimal training accuracy

testing

$$d(\widehat{\text{output}}_j, g(x_{\text{trained}}, \widehat{\text{input}}_j))$$

matters

test data

model (e.g., neural network) parameters

training

$$\min_x f(x) := \sum_i d(\text{output}_i, g(x, \text{input}_i))$$

training data (known)

large h? this talk

trapped in local min \leftrightarrow suboptimal training accuracy

generalization



testing

$$d(\widehat{\text{output}}_j, g(x_{\text{trained}}, \widehat{\text{input}}_j))$$

matters

test data

model (e.g., neural network) parameters

training

$$\min_x f(x) := \sum_i d(\text{output}_i, g(x, \text{input}_i))$$

training data (known)

large h? this talk

trapped in **local min** \leftrightarrow **suboptimal** training accuracy

large h? offline

testing

$$d(\widehat{\text{output}}_j, g(x_{\text{trained}}, \widehat{\text{input}}_j))$$

test data

matters

generalization



Stochasticity of Deterministic Gradient Descent: Large Learning Rate for Multiscale Objective Function

NeurIPS 2020, arXiv: 2002.06189

Lingkai Kong¹ and *Molei Tao*¹

1 Georgia Institute of Technology (USA)



$\min f(x)$

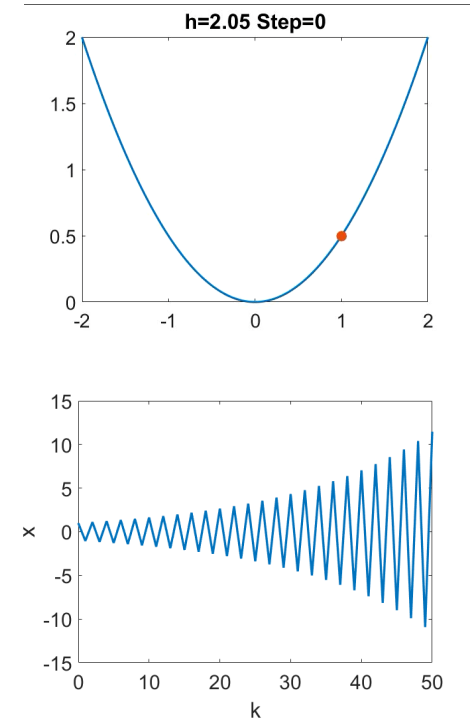
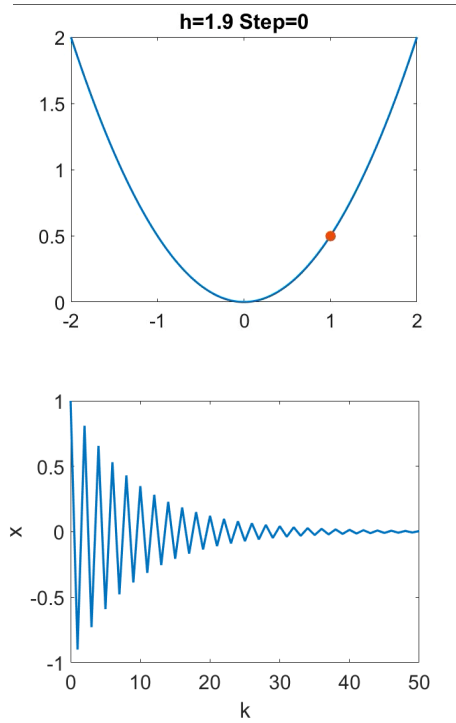
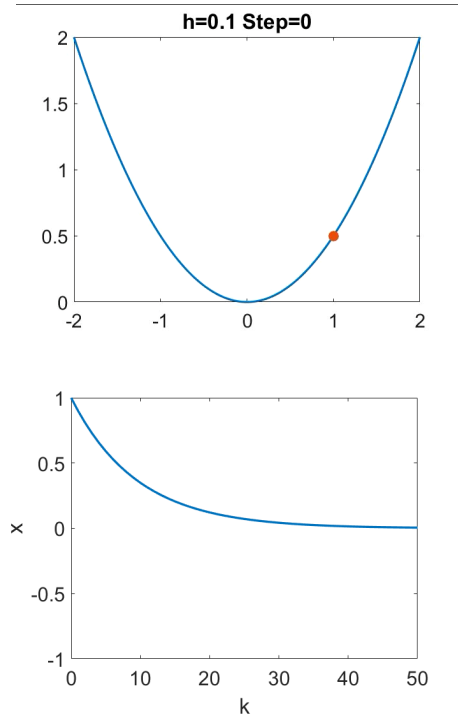
$$x_{k+1} = x_k - h \nabla f(x_k)$$

Gradient Descent

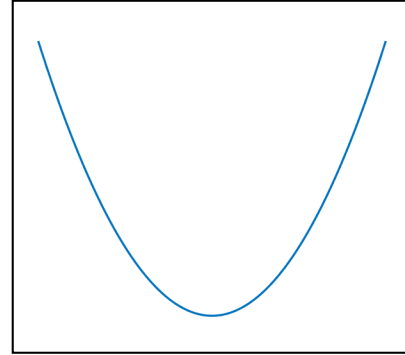
$\min f(x)$

$$x_{k+1} = x_k - h \nabla f(x_k)$$

Gradient Descent

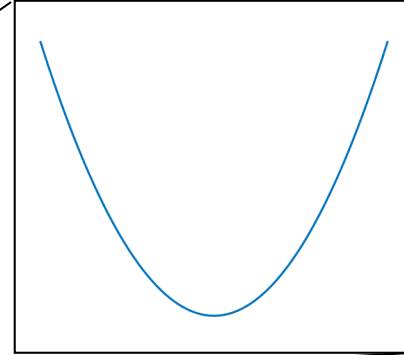
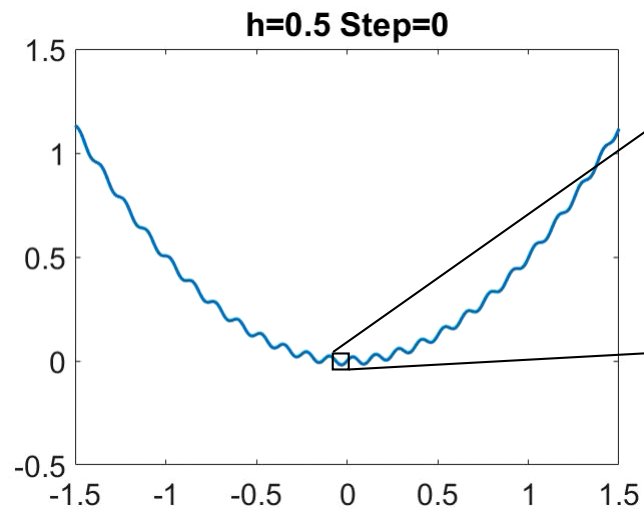


1-2. Introduction: GD with 'large' LR for multiscale objective -- phenomenology



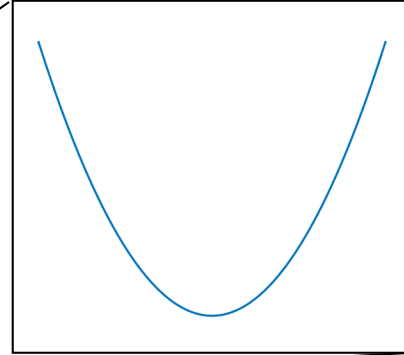
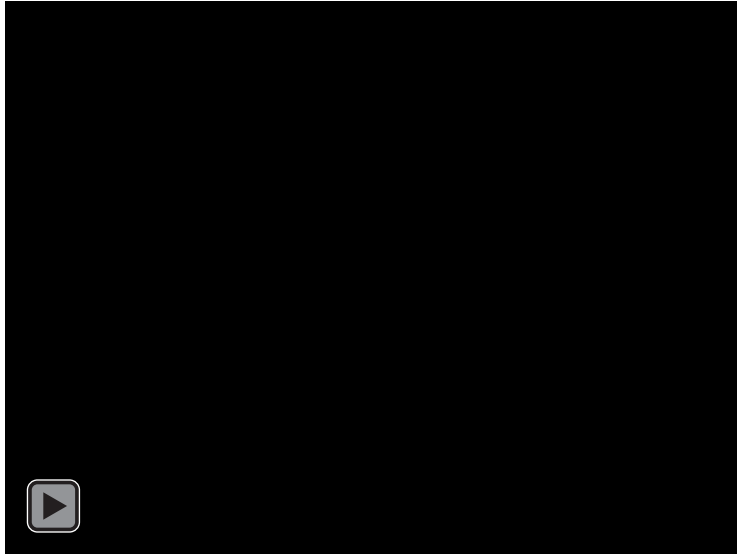
1-2. Introduction: GD with 'large' LR for multiscale objective -- phenomenology

GD animation

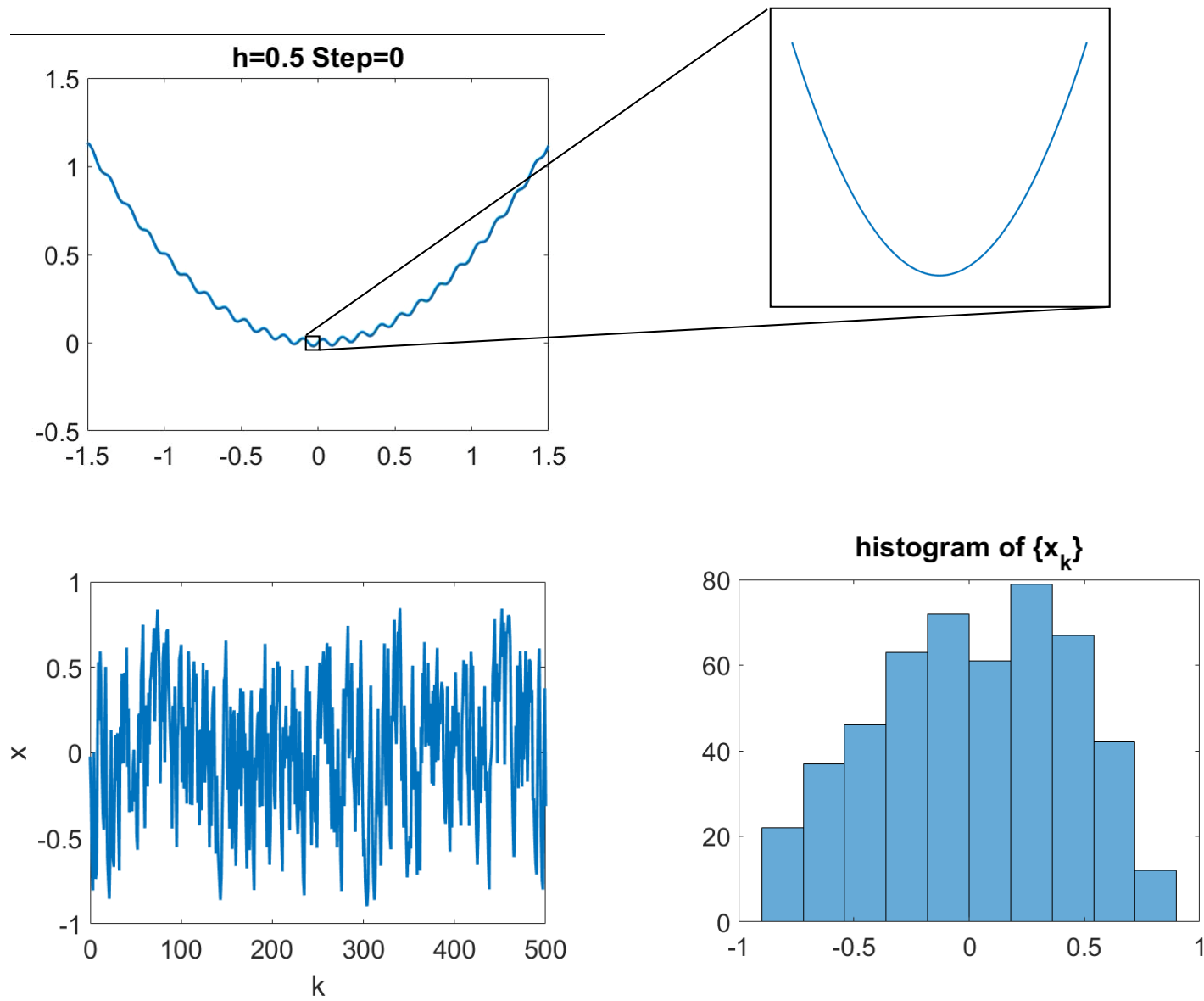


1-2. Introduction: GD with 'large' LR for multiscale objective -- phenomenology

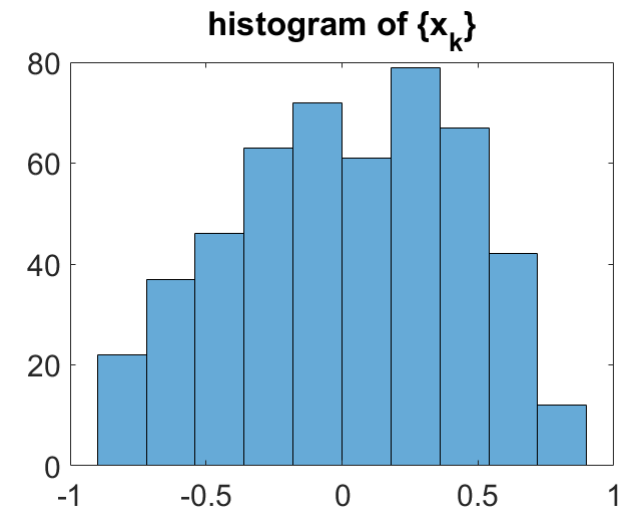
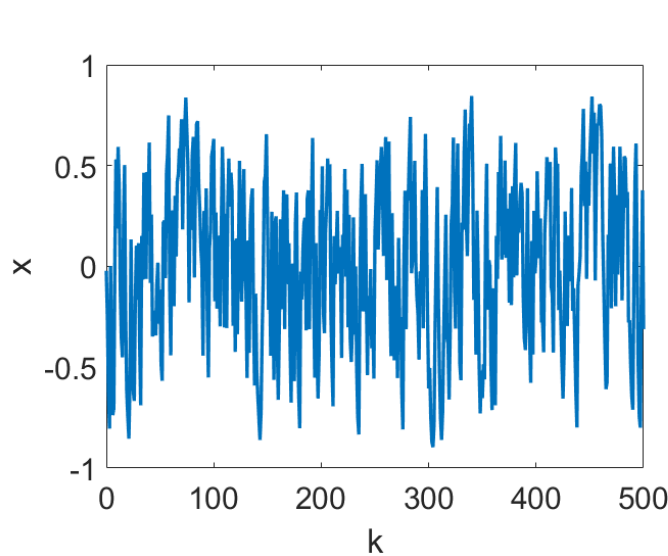
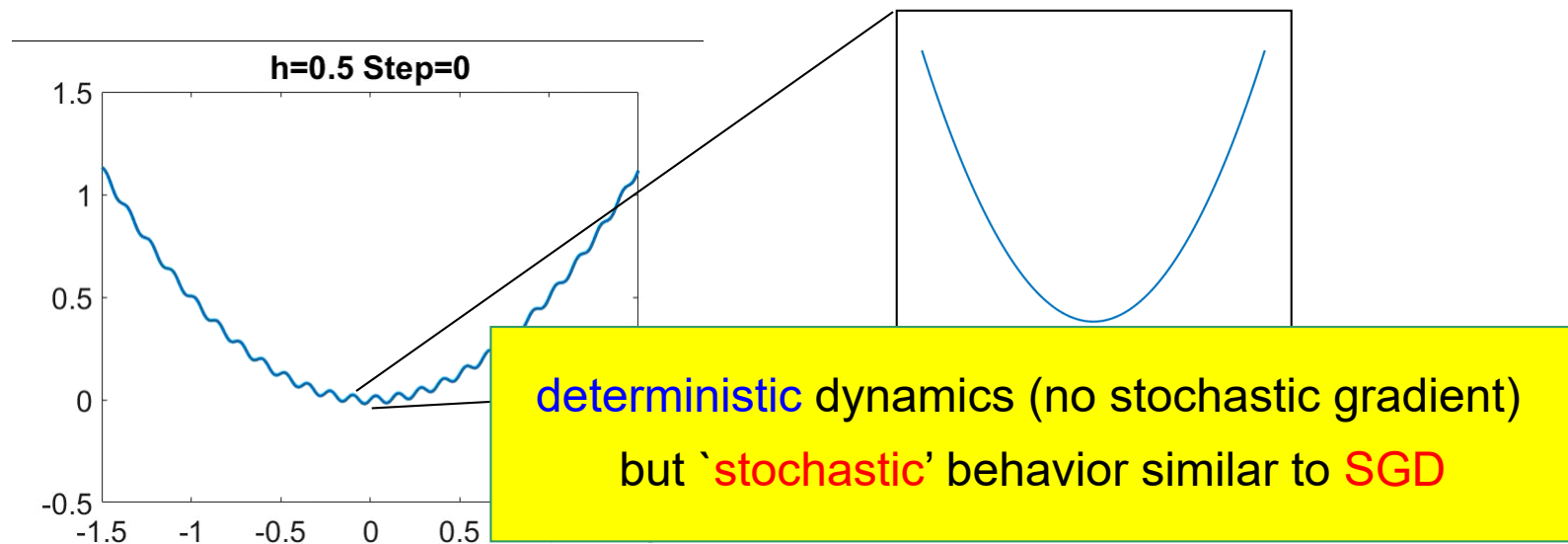
GD animation



1-2. Introduction: GD with 'large' LR for multiscale objective -- phenomenology



1-2. Introduction: GD with 'large' LR for multiscale objective -- phenomenology



2-1. Theory: the setup: multiscale objective function

multiscale objective +
deterministic GD with large LR



'stochastic'
behaviors

2-1. Theory: the setup: multiscale objective function

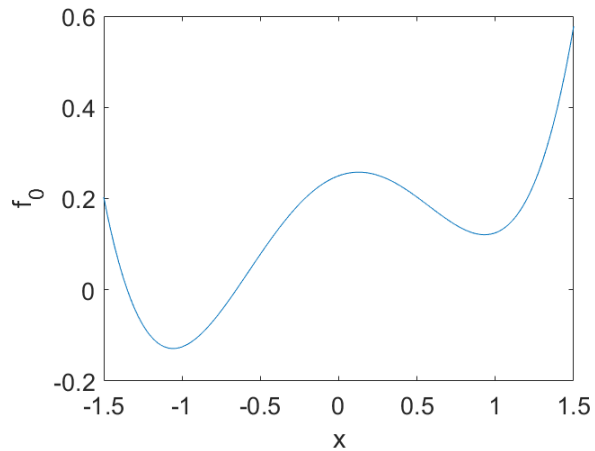
multiscale objective +
deterministic GD with large LR

multiscale objective fct. $f : \mathbb{R}^d \rightarrow \mathbb{R}$. $f(x) := f_0(x) + f_{1,\epsilon}(x)$

2-1. Theory: the setup: multiscale objective function

multiscale objective +
deterministic GD with large LR

multiscale objective fct. $f : \mathbb{R}^d \rightarrow \mathbb{R}$. $f(x) := \overset{\text{macro}}{\boxed{f_0(x)}} + f_{1,\epsilon}(x)$

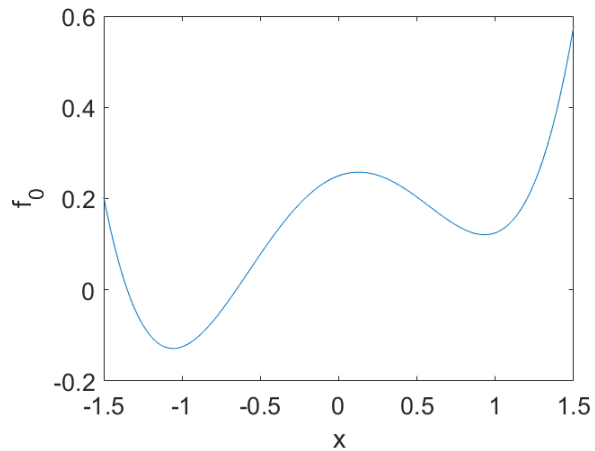


2-1. Theory: the setup: multiscale objective function

multiscale objective +
deterministic GD with large LR

multiscale objective fct. $f : \mathbb{R}^d \rightarrow \mathbb{R}$. $f(x) := \boxed{f_0(x)} + \boxed{f_{1,\epsilon}(x)}$ $\epsilon \ll 1$

macro micro



2-1. Theory: the setup: multiscale objective function

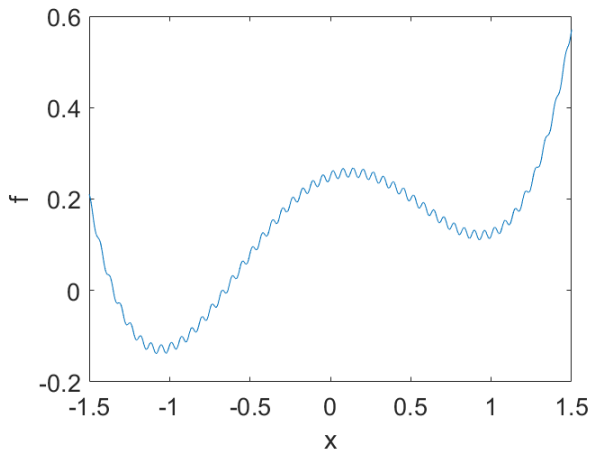
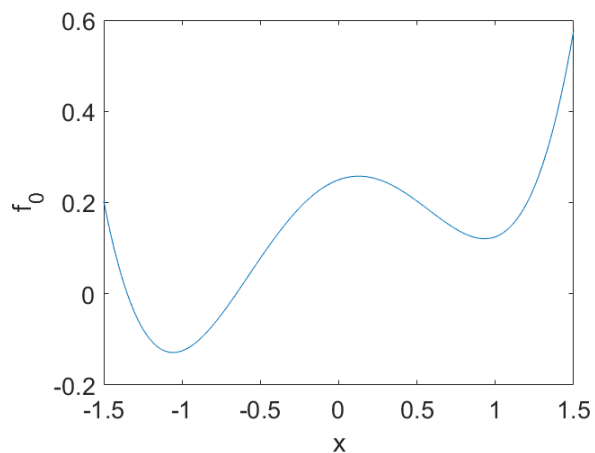
multiscale objective +
deterministic GD with large LR

multiscale objective fct. $f : \mathbb{R}^d \rightarrow \mathbb{R}$. $f(x) := f_0(x) + \boxed{f_{1,\epsilon}(x)}$

micro $\epsilon \ll 1$

EX $f_{1,\epsilon} := \epsilon f_1 \left(\frac{x}{\epsilon} \right)$

f_1 periodic



2-1. Theory: the setup: multiscale objective function

multiscale objective +
deterministic GD with large LR

multiscale objective fct. $f : \mathbb{R}^d \rightarrow \mathbb{R}$. $f(x) := f_0(x) + \boxed{f_{1,\epsilon}(x)}$

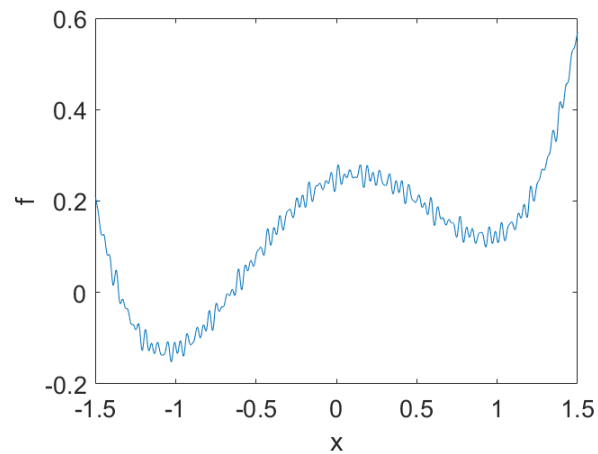
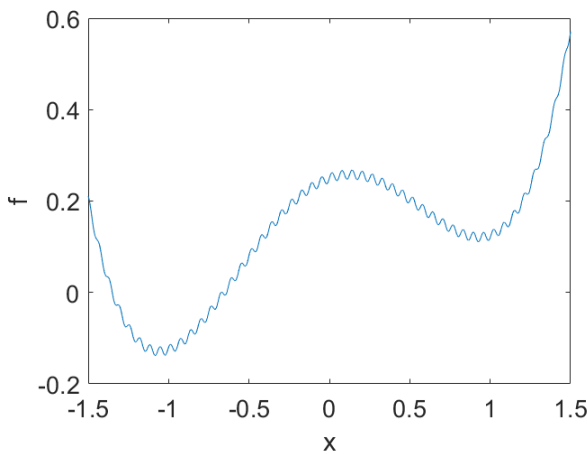
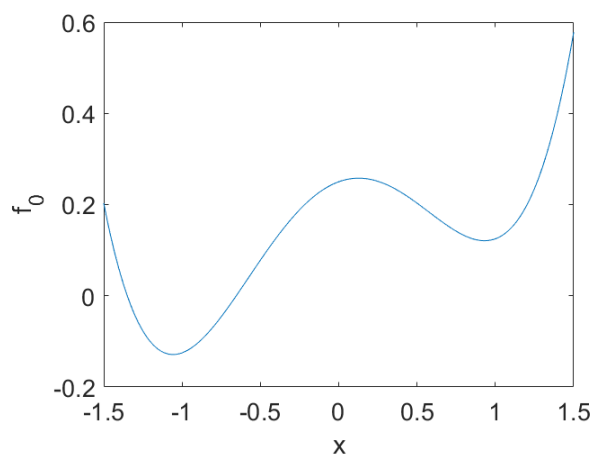
micro $\epsilon \ll 1$

EX $f_{1,\epsilon} := \epsilon f_1 \left(\frac{x}{\epsilon} \right)$

f_1 periodic

EX $f_{1,\epsilon} := \epsilon F_1 \left(\frac{\omega_1 x}{\epsilon}, \frac{\omega_2 x}{\epsilon}, \dots, \frac{\omega_N x}{\epsilon} \right)$

F_1 1-periodic in each argument



2-1. Theory: the setup: multiscale objective function

multiscale objective +
deterministic GD with large LR

multiscale objective fct. $f : \mathbb{R}^d \rightarrow \mathbb{R}$. $f(x) := f_0(x) + \boxed{f_{1,\epsilon}(x)}$

micro $\epsilon \ll 1$

EX $f_{1,\epsilon} := \epsilon f_1 \left(\frac{x}{\epsilon} \right)$

f_1 periodic

EX $f_{1,\epsilon} := \epsilon F_1 \left(\frac{\omega_1 x}{\epsilon}, \frac{\omega_2 x}{\epsilon}, \dots, \frac{\omega_N x}{\epsilon} \right)$

F_1 1-periodic in each argument

more general:

Condition 1. \exists bounded random variable ζ with $P(\zeta = 0) < 1$, s.t. for any bounded rectangle $\Gamma \subset \mathbb{R}^d$ with area $|\Gamma| > 0$, $-\nabla f_{1,\epsilon}(U_\Gamma) \xrightarrow{w} \zeta$ when $\epsilon \rightarrow 0$. Assume without loss of generality that $\mathbb{E}\zeta = 0$ (nonzero mean can be absorbed into f_0).

Condition 2. $\epsilon \nabla^2 f_{1,\epsilon}$ is uniformly bounded as $\epsilon \rightarrow 0$, and $\exists m \in \mathbb{R}$, s.t. for any bounded rectangle $\Gamma \subset \mathbb{R}^d$ whose area $|\Gamma| > 0$, $\mathbb{E} [\ln \|\epsilon \nabla^2 f_{1,\epsilon}(U_\Gamma)\|_2] \rightarrow m$.

2-1. Theory: the setup: multiscale objective function

multiscale objective +
deterministic GD with large LR

multiscale objective fct. $f : \mathbb{R}^d \rightarrow \mathbb{R}$. $f(x) := f_0(x) + \boxed{f_{1,\epsilon}(x)}$

micro $\epsilon \ll 1$

EX $f_{1,\epsilon} := \epsilon f_1 \left(\frac{x}{\epsilon} \right)$

f_1 periodic

EX $f_{1,\epsilon} := \epsilon F_1 \left(\frac{\omega_1 x}{\epsilon}, \frac{\omega_2 x}{\epsilon}, \dots, \frac{\omega_N x}{\epsilon} \right)$

F_1 1-periodic in each argument

more general:

$O(1)$ 1st-derivative with a weak limit

Condition 1. \exists bounded random variable ζ with $P(\zeta = 0) < 1$, s.t. for any bounded rectangle $\Gamma \subset \mathbb{R}^d$ with area $|\Gamma| > 0$, $-\nabla f_{1,\epsilon}(U_\Gamma) \xrightarrow{w} \zeta$ when $\epsilon \rightarrow 0$. Assume without loss of generality that $\mathbb{E}\zeta = 0$ (nonzero mean can be absorbed into f_0).

Condition 2. $\epsilon \nabla^2 f_{1,\epsilon}$ is uniformly bounded as $\epsilon \rightarrow 0$, and $\exists m \in \mathbb{R}$, s.t. for any bounded rectangle $\Gamma \subset \mathbb{R}^d$ whose area $|\Gamma| > 0$, $\mathbb{E} [\ln \|\epsilon \nabla^2 f_{1,\epsilon}(U_\Gamma)\|_2] \rightarrow m$.

2-1. Theory: the setup: multiscale objective function

multiscale objective +
deterministic GD with large LR

multiscale objective fct. $f : \mathbb{R}^d \rightarrow \mathbb{R}$. $f(x) := f_0(x) + \boxed{f_{1,\epsilon}(x)}$

micro $\epsilon \ll 1$

EX $f_{1,\epsilon} := \epsilon f_1 \left(\frac{x}{\epsilon} \right)$

f_1 periodic

EX $f_{1,\epsilon} := \epsilon F_1 \left(\frac{\omega_1 x}{\epsilon}, \frac{\omega_2 x}{\epsilon}, \dots, \frac{\omega_N x}{\epsilon} \right)$

F_1 1-periodic in each argument

more general:

$O(1)$ 1st-derivative with a weak limit

Condition 1. \exists bounded random variable ζ with $P(\zeta = 0) < 1$, s.t. for any bounded rectangle $\Gamma \subset \mathbb{R}^d$ with area $|\Gamma| > 0$, $-\nabla f_{1,\epsilon}(U_\Gamma) \xrightarrow{w} \zeta$ when $\epsilon \rightarrow 0$. Assume without loss of generality that $\mathbb{E}\zeta = 0$ (nonzero mean can be absorbed into f_0).

Condition 2. $\epsilon \nabla^2 f_{1,\epsilon}$ is uniformly bounded as $\epsilon \rightarrow 0$, and $\exists m \in \mathbb{R}$, s.t. for any bounded rectangle $\Gamma \subset \mathbb{R}^d$ whose area $|\Gamma| > 0$, $\mathbb{E} [\ln \|\epsilon \nabla^2 f_{1,\epsilon}(U_\Gamma)\|_2] \rightarrow m$.

$O(1/\epsilon)$ 2nd-derivative with a limiting expectation

2-1. Theory: the setup: large Learning Rate

multiscale objective +
deterministic GD with large LR

$$f(x) := f_0(x) + f_{1,\epsilon}(x)$$

$$\text{e.g., } f_{1,\epsilon} = \epsilon f_1(x/\epsilon)$$

$$x \mapsto x - \eta \nabla f(x), \quad \text{i.e.} \quad x \mapsto x - \eta (\nabla f_0(x) + \nabla f_{1,\epsilon}(x))$$

2-1. Theory: the setup: large Learning Rate

multiscale objective +
deterministic GD with large LR

$$f(x) := f_0(x) + f_{1,\epsilon}(x)$$

$$\text{e.g., } f_{1,\epsilon} = \epsilon f_1(x/\epsilon)$$

$$x \mapsto x - \eta \nabla f(x), \quad \text{i.e.} \quad x \mapsto x - \eta (\nabla f_0(x) + \nabla f_{1,\epsilon}(x))$$

small LR?

2-1. Theory: the setup: large Learning Rate

multiscale objective +
deterministic GD with large LR

$$f(x) := f_0(x) + f_{1,\epsilon}(x)$$

e.g., $f_{1,\epsilon} = \epsilon f_1(x/\epsilon)$

$$x \mapsto x - \eta \nabla f(x), \quad \text{i.e.} \quad x \mapsto x - \eta (\nabla f_0(x) + \nabla f_{1,\epsilon}(x))$$

small LR? $\eta \ll 1/L$ L : Lipschitz const. of ∇f , $L = \mathcal{O}(1/\epsilon)$
i.e., $\eta \ll \epsilon$ resolve the small scale \rightarrow conv. to local min

2-1. Theory: the setup: large Learning Rate

multiscale objective +
deterministic GD with large LR

$$f(x) := f_0(x) + f_{1,\epsilon}(x)$$

e.g., $f_{1,\epsilon} = \epsilon f_1(x/\epsilon)$

$$x \mapsto x - \eta \nabla f(x), \quad \text{i.e.} \quad x \mapsto x - \eta (\nabla f_0(x) + \nabla f_{1,\epsilon}(x))$$

small LR? $\eta \ll 1/L$ L : Lipschitz const. of ∇f , $L = \mathcal{O}(1/\epsilon)$

i.e., $\eta \ll \epsilon$ resolve the small scale \rightarrow conv. to local min

bad LR? $\eta \gg 1$ can't even resolve the large scale

i.e., $x \mapsto x - \eta \nabla f_0(x)$ is unstable

2-1. Theory: the setup: large Learning Rate

**multiscale objective +
deterministic GD with large LR**

$$f(x) := f_0(x) + f_{1,\epsilon}(x)$$

e.g., $f_{1,\epsilon} = \epsilon f_1(x/\epsilon)$

$$x \mapsto x - \eta \nabla f(x), \quad \text{i.e.} \quad x \mapsto x - \eta (\nabla f_0(x) + \nabla f_{1,\epsilon}(x))$$

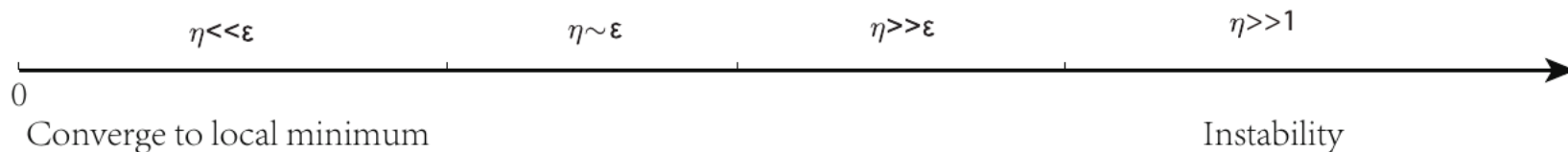
small LR? $\eta \ll 1/L$ L : Lipschitz const. of ∇f , $L = \mathcal{O}(1/\epsilon)$

i.e., $\eta \ll \epsilon$ resolve the small scale \rightarrow conv. to local min

bad LR? $\eta \gg 1$ can't even resolve the large scale

i.e., $x \mapsto x - \eta \nabla f_0(x)$ is unstable

in-between...



2-1. Theory: the setup: large Learning Rate

multiscale objective +
deterministic GD with large LR

$$f(x) := f_0(x) + f_{1,\epsilon}(x)$$

e.g., $f_{1,\epsilon} = \epsilon f_1(x/\epsilon)$

$$x \mapsto x - \eta \nabla f(x), \quad \text{i.e.} \quad x \mapsto x - \eta (\nabla f_0(x) + \nabla f_{1,\epsilon}(x))$$

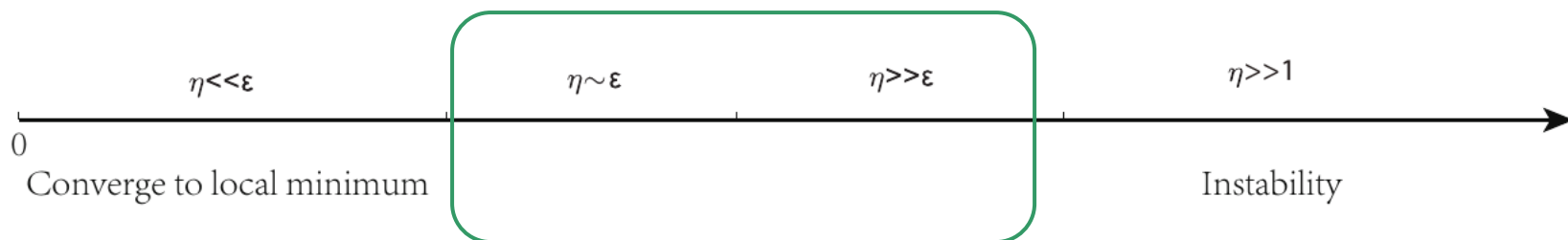
small LR? $\eta \ll 1/L$ L : Lipschitz const. of ∇f , $L = \mathcal{O}(1/\epsilon)$

i.e., $\eta \ll \epsilon$ resolve the small scale \rightarrow conv. to local min

bad LR? $\eta \gg 1$ can't even resolve the large scale

i.e., $x \mapsto x - \eta \nabla f_0(x)$ is unstable

in-between...



2-1. Theory: the setup: large Learning Rate

multiscale objective +
deterministic GD with large LR

$$f(x) := f_0(x) + f_{1,\epsilon}(x)$$

e.g., $f_{1,\epsilon} = \epsilon f_1(x/\epsilon)$

$$x \mapsto x - \eta \nabla f(x), \quad \text{i.e.} \quad x \mapsto x - \eta (\nabla f_0(x) + \nabla f_{1,\epsilon}(x))$$

small LR? $\eta \ll 1/L$ L : Lipschitz const. of ∇f , $L = \mathcal{O}(1/\epsilon)$

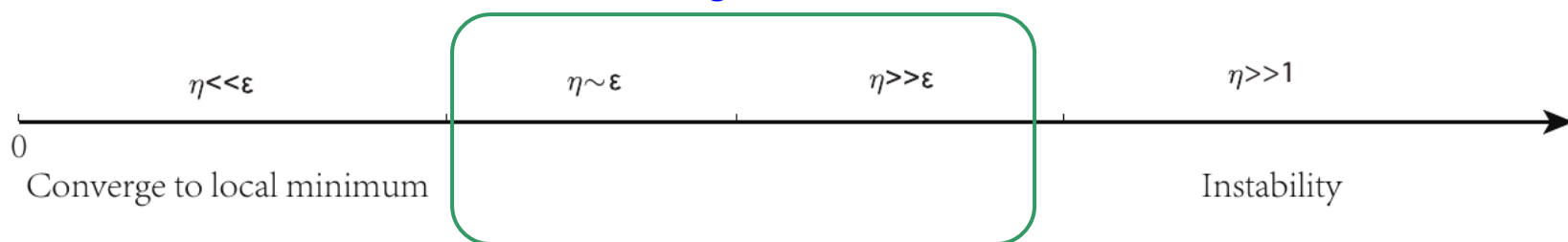
i.e., $\eta \ll \epsilon$ resolve the small scale \rightarrow conv. to local min

bad LR? $\eta \gg 1$ can't even resolve the large scale

i.e., $x \mapsto x - \eta \nabla f_0(x)$ is unstable

in-between...

large LR



2-1. Theory: the setup: large Learning Rate

multiscale objective +
deterministic GD with large LR

$$f(x) := f_0(x) + f_{1,\epsilon}(x)$$

e.g., $f_{1,\epsilon} = \epsilon f_1(x/\epsilon)$

$$x \mapsto x - \eta \nabla f(x), \quad \text{i.e.} \quad x \mapsto x - \eta (\nabla f_0(x) + \nabla f_{1,\epsilon}(x))$$

small LR? $\eta \ll 1/L$ L : Lipschitz const. of ∇f , $L = \mathcal{O}(1/\epsilon)$

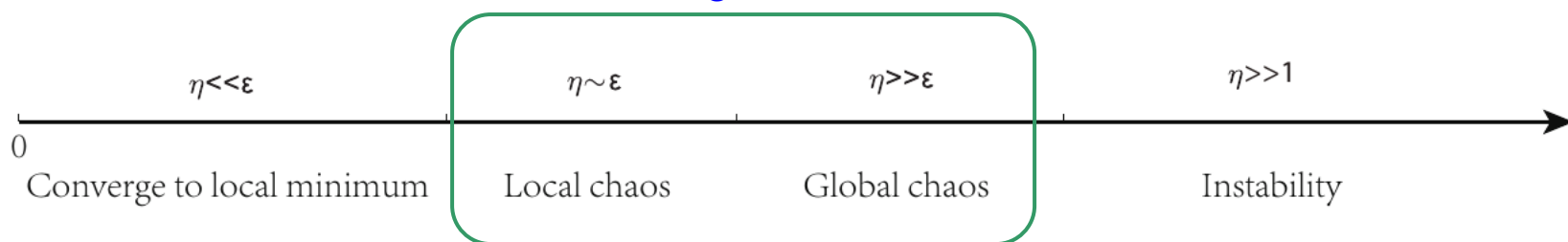
i.e., $\eta \ll \epsilon$ resolve the small scale \rightarrow conv. to local min

bad LR? $\eta \gg 1$ can't even resolve the large scale

i.e., $x \mapsto x - \eta \nabla f_0(x)$ is unstable

in-between...

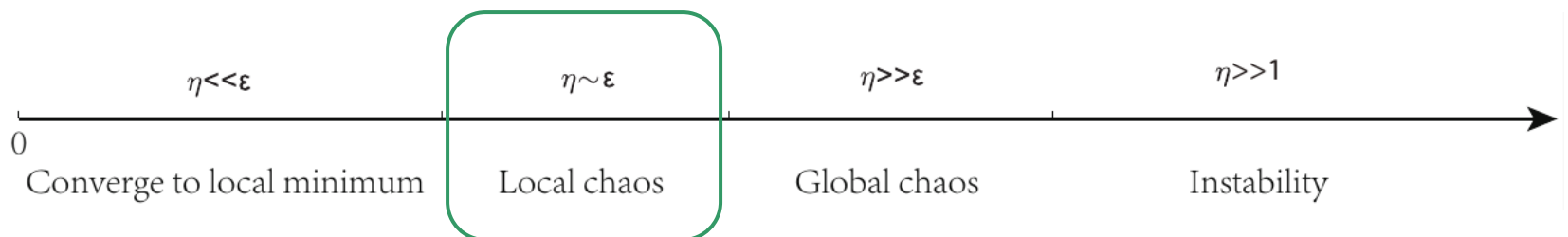
large LR



2-2. Theory: the local chaos regime

deterministic GD map

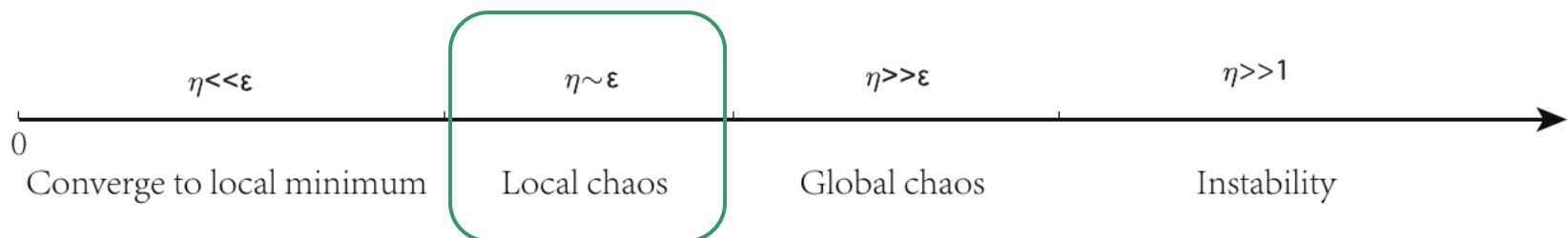
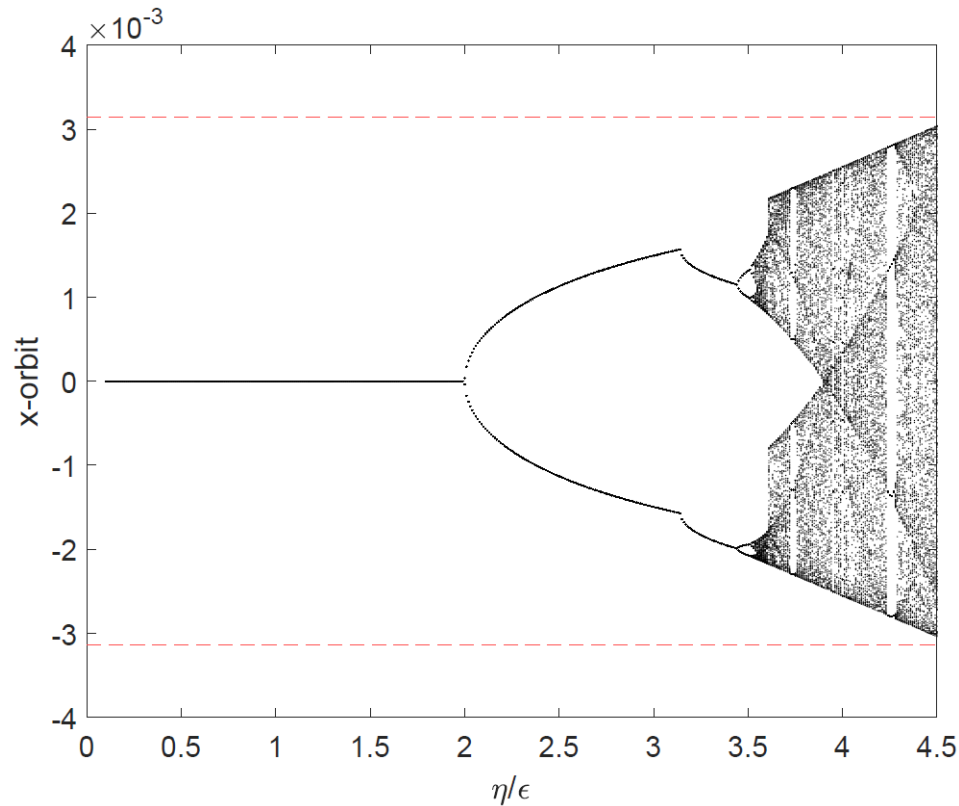
$$\varphi(x) := x - \eta(\nabla f_0(x) + \nabla f_{1,\epsilon}(x))$$



deterministic GD map

$$\varphi(x) := x - \eta(\nabla f_0(x) + \nabla f_{1,\epsilon}(x))$$

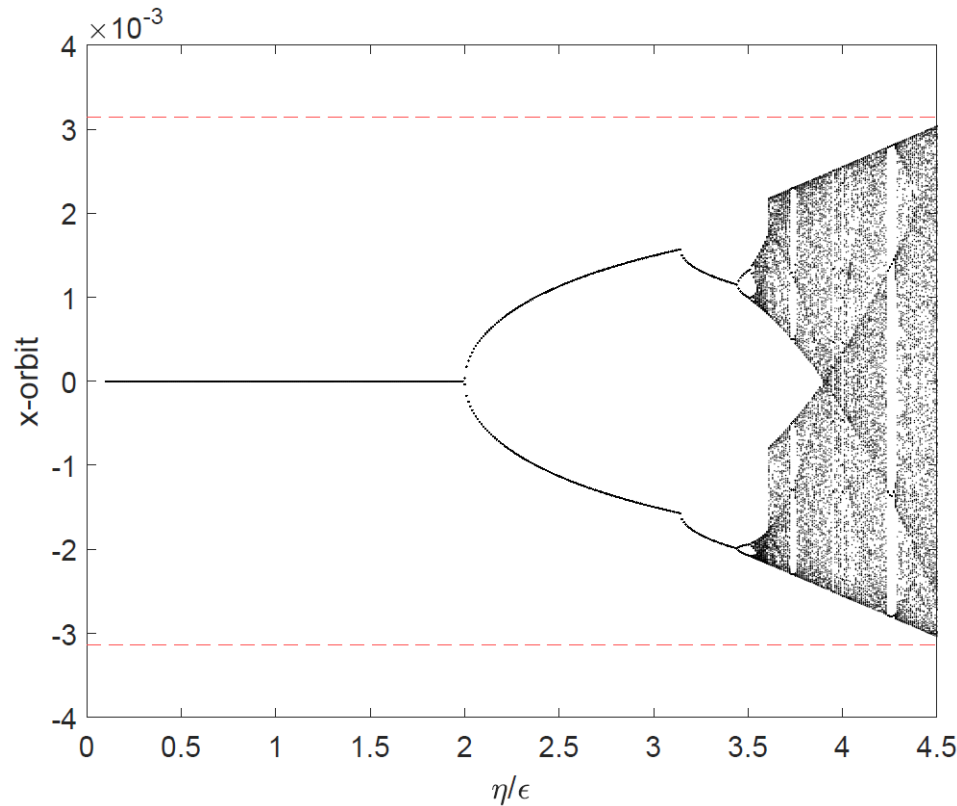
transition
into
chaos
via
period
doublings



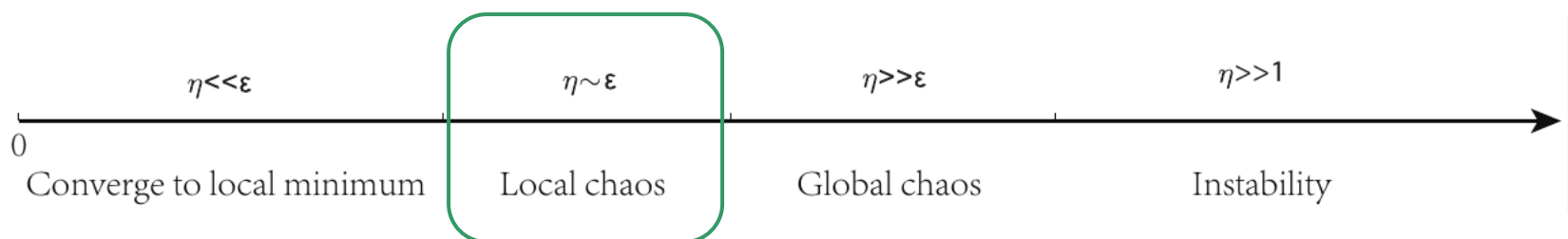
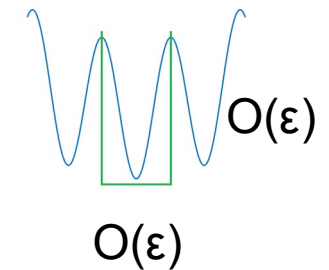
deterministic GD map

$$\varphi(x) := x - \eta(\nabla f_0(x) + \nabla f_{1,\epsilon}(x))$$

transition
into
chaos
via
period
doublings



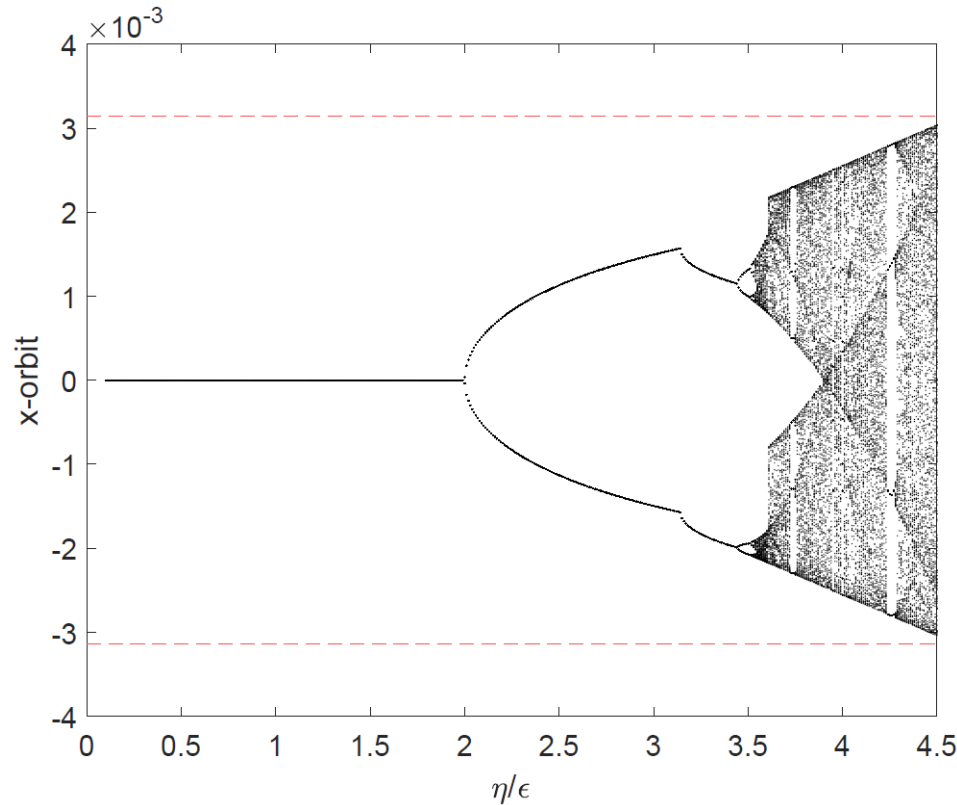
x stays within
a microscopic
potential well



deterministic GD map

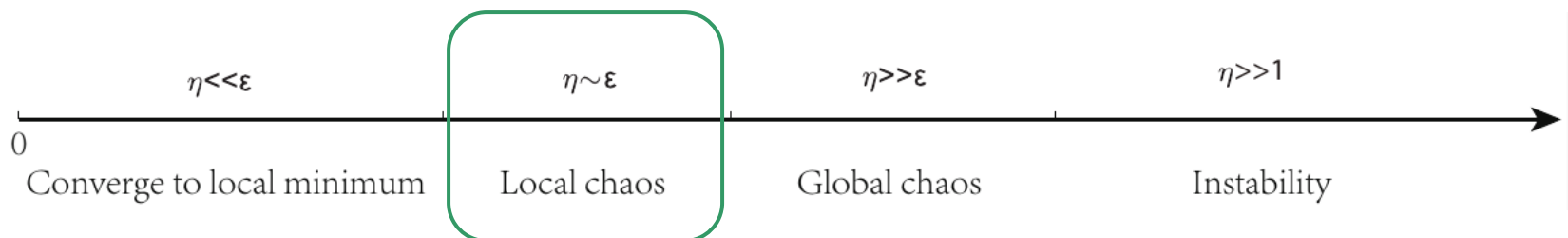
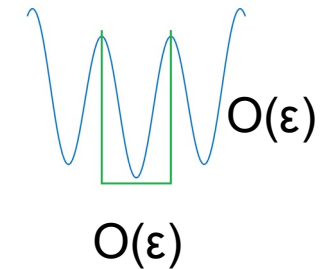
$$\varphi(x) := x - \eta(\nabla f_0(x) + \nabla f_{1,\epsilon}(x))$$

transition
into
chaos
via
period
doublings



unimodal map
within a well

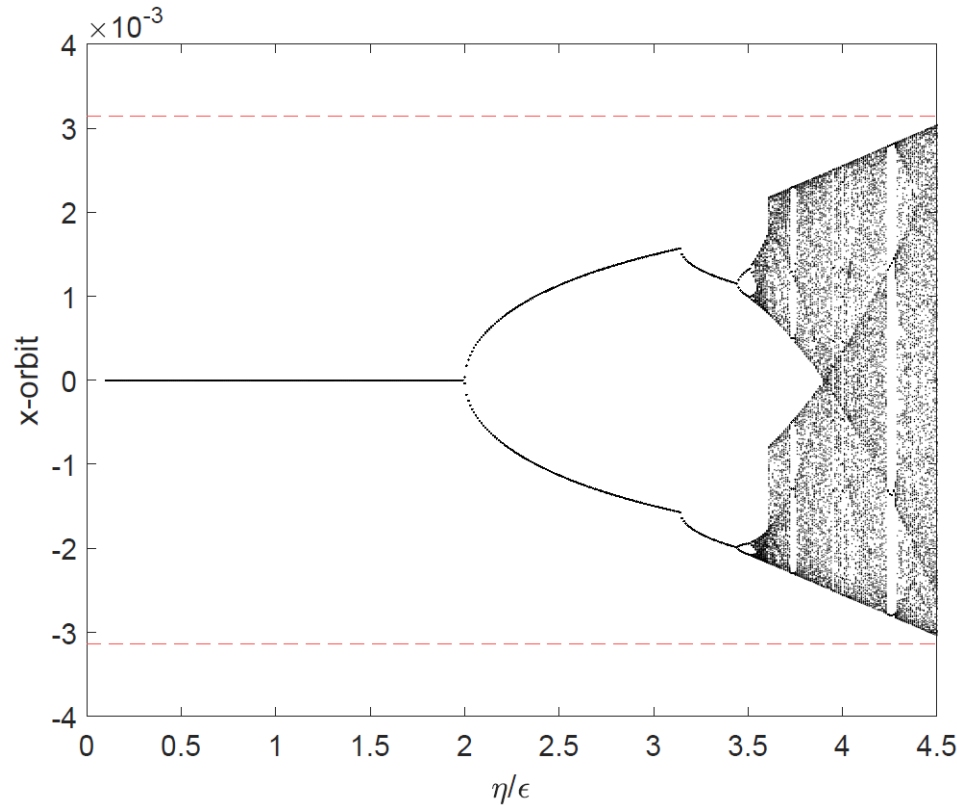
x stays within
a microscopic
potential well



deterministic GD map

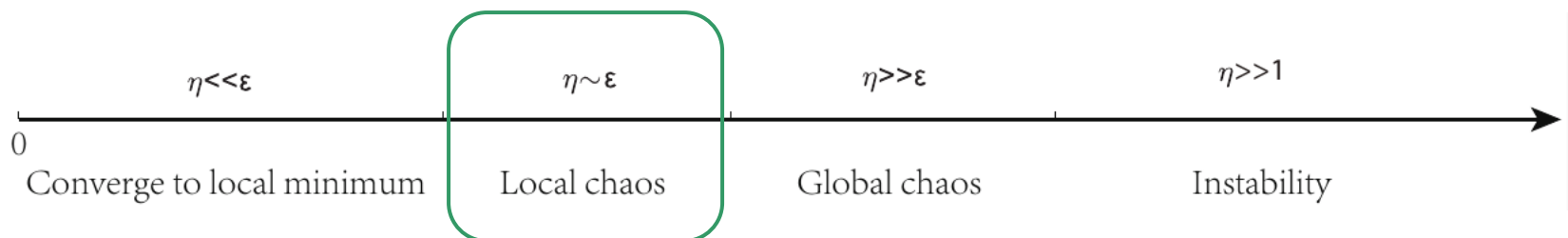
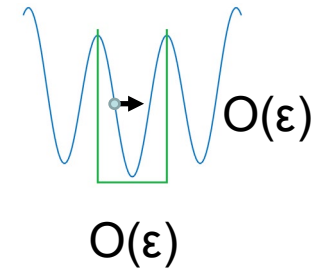
$$\varphi(x) := x - \eta(\nabla f_0(x) + \nabla f_{1,\epsilon}(x))$$

transition
into
chaos
via
period
doublings



unimodal map
within a well

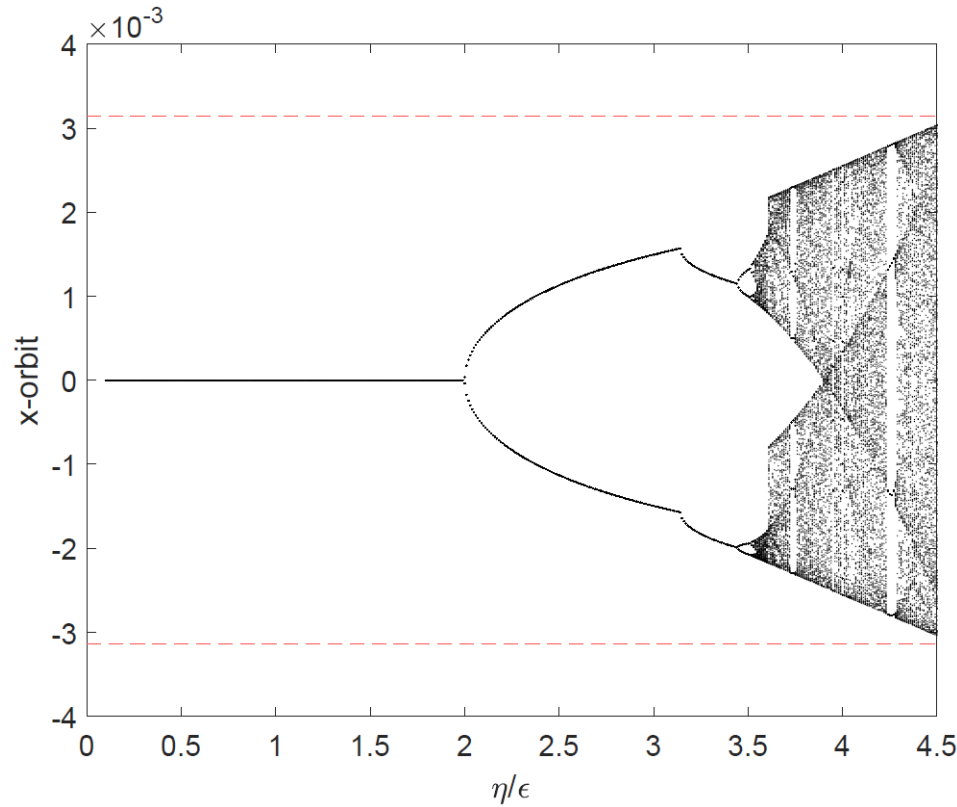
x stays within
a microscopic
potential well



deterministic GD map

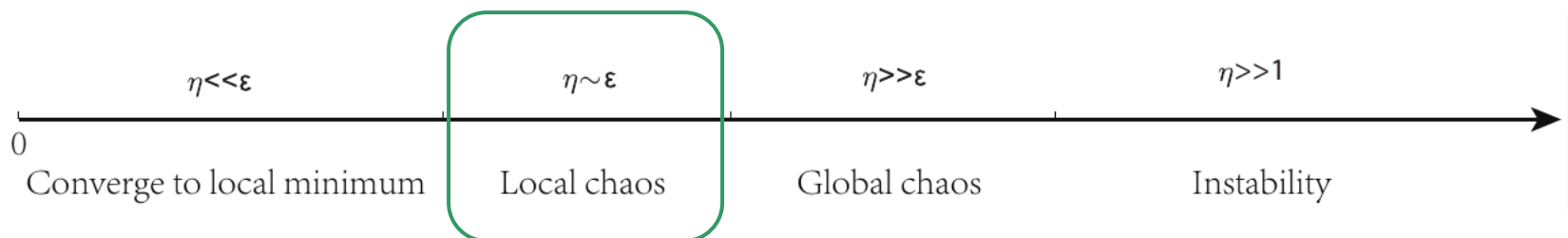
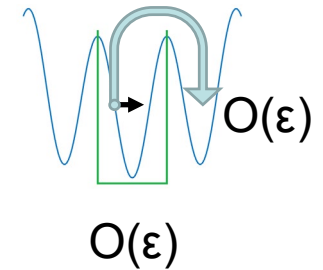
$$\varphi(x) := x - \eta(\nabla f_0(x) + \nabla f_{1,\epsilon}(x))$$

transition
into
chaos
via
period
doublings



unimodal map
within a well

x stays within
a microscopic
potential well

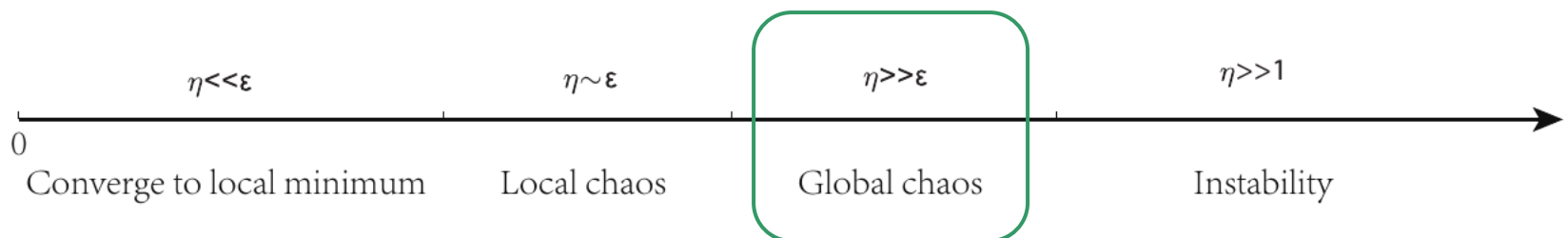


2-3. Theory: the global chaos regime

deterministic GD map

$$\varphi(x) := x - \eta(\nabla f_0(x) + \nabla f_{1,\epsilon}(x))$$

η further increases (i.e. the $\epsilon \rightarrow 0$ regime)



2-3. Theory: the global chaos regime

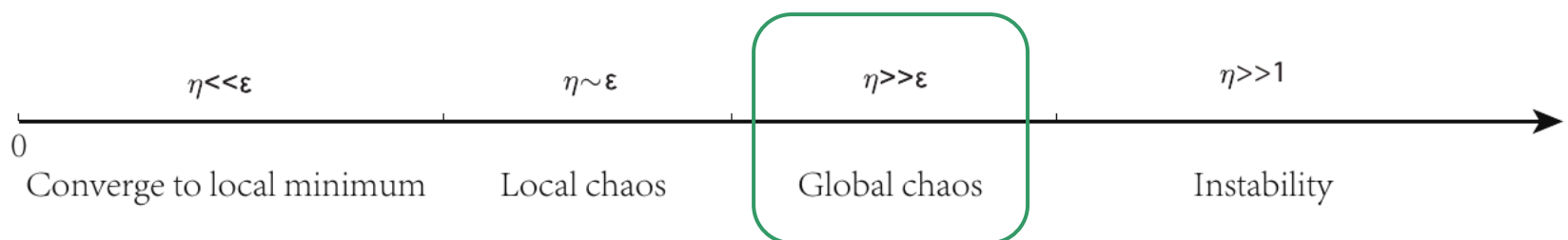
deterministic GD map $\varphi(x) := x - \eta(\nabla f_0(x) + \nabla f_{1,\epsilon}(x))$

η further increases (i.e. the $\epsilon \rightarrow 0$ regime)

- f_0 strongly convex, C^2 , L -smooth \rightarrow

limiting distribution of φ iterates $\stackrel{\text{small } \epsilon}{\approx}$ nearly Gibbs

$$\frac{1}{Z} \exp\left(-\frac{2f_0(x)}{\eta\sigma^2}\right) dx + \mathcal{O}(\eta^2)$$



2-3. Theory: the global chaos regime

deterministic GD map

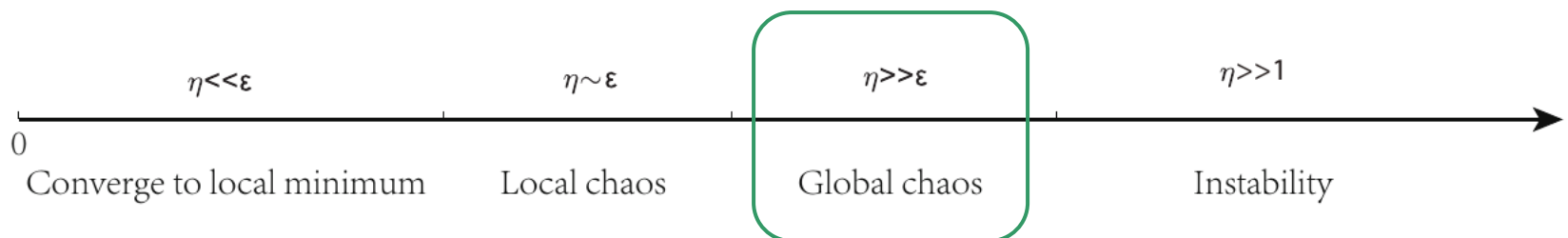
$$\varphi(x) := x - \eta(\nabla f_0(x) + \nabla f_{1,\epsilon}(x))$$

η further increases (i.e. the $\epsilon \rightarrow 0$ regime)

limiting distribution: nearly Gibbs

(f_0 strongly convex)

$$\frac{1}{Z} \exp\left(-\frac{2f_0(x)}{\eta\sigma^2}\right) dx + \mathcal{O}(\eta^2)$$



deterministic GD map

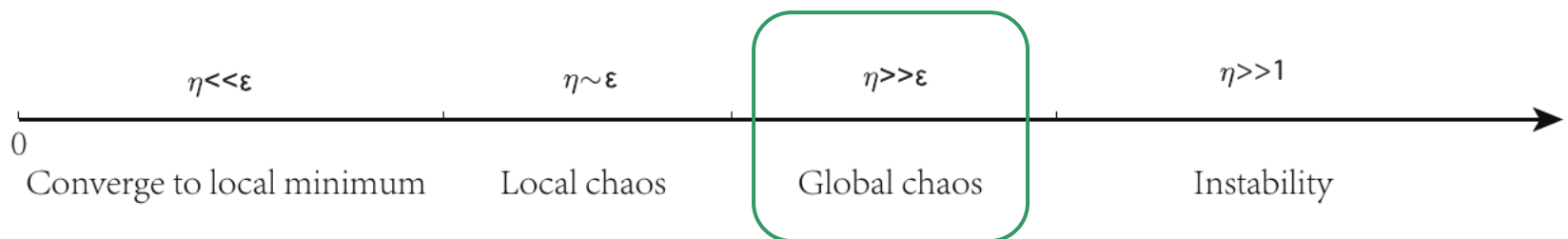
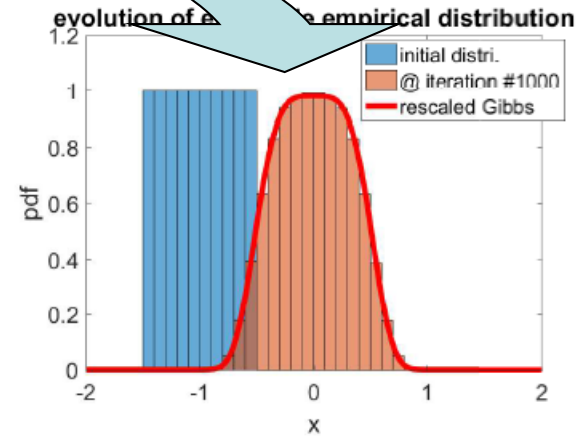
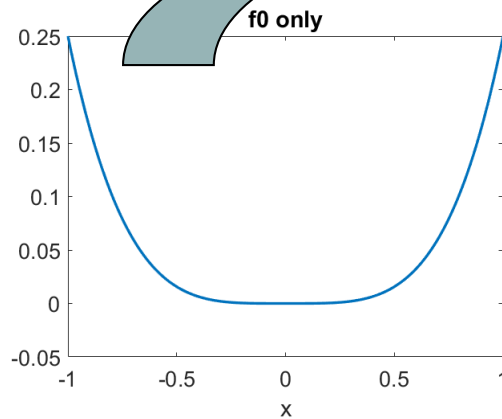
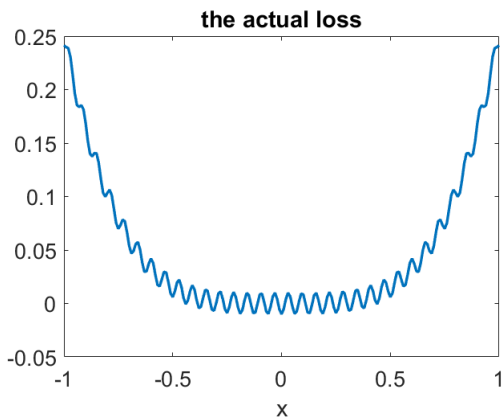
$$\varphi(x) := x - \eta(\nabla f_0(x) + \nabla f_{1,\epsilon}(x))$$

η further increases (i.e. the $\epsilon \rightarrow 0$ regime)

limiting distribution: nearly Gibbs

(f_0 strongly convex)

$$\frac{1}{Z} \exp\left(-\frac{2f_0(x)}{\eta\sigma^2}\right) dx + \mathcal{O}(\eta^2)$$



deterministic GD map

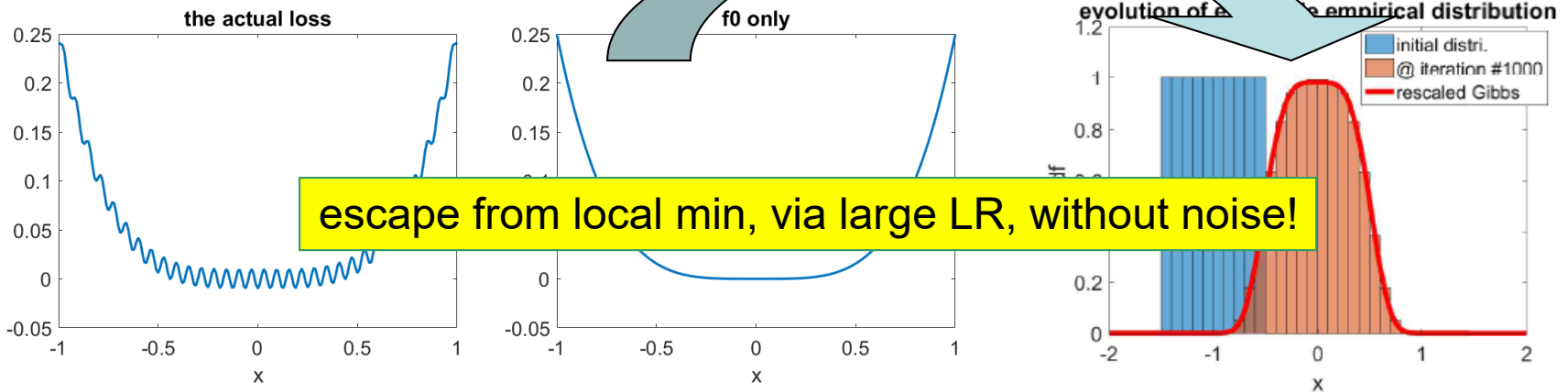
$$\varphi(x) := x - \eta(\nabla f_0(x) + \nabla f_{1,\epsilon}(x))$$

η further increases (i.e. the $\epsilon \rightarrow 0$ regime)

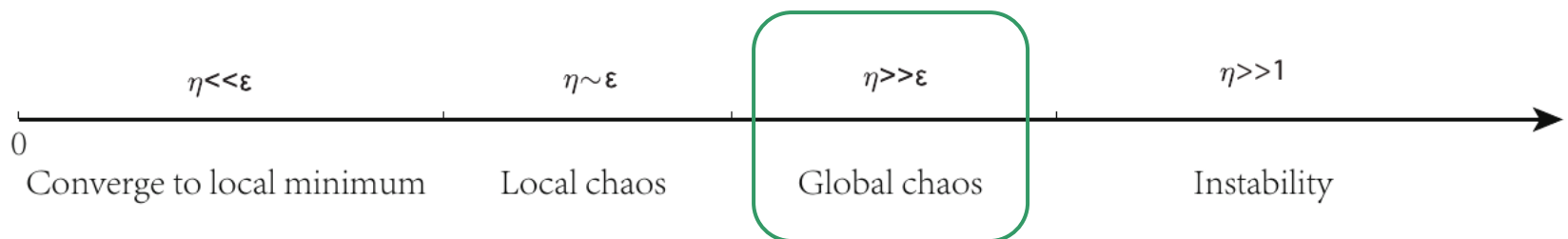
limiting distribution: nearly Gibbs

(f_0 strongly convex)

$$\frac{1}{Z} \exp\left(-\frac{2f_0(x)}{\eta\sigma^2}\right) dx + \mathcal{O}(\eta^2)$$



escape from local min, via large LR, without noise!



2-5. Theory: the global chaos regime: non-convex macroscale

deterministic GD map

$$\varphi(x) := x - \eta(\nabla f_0(x) + \nabla f_{1,\epsilon}(x))$$

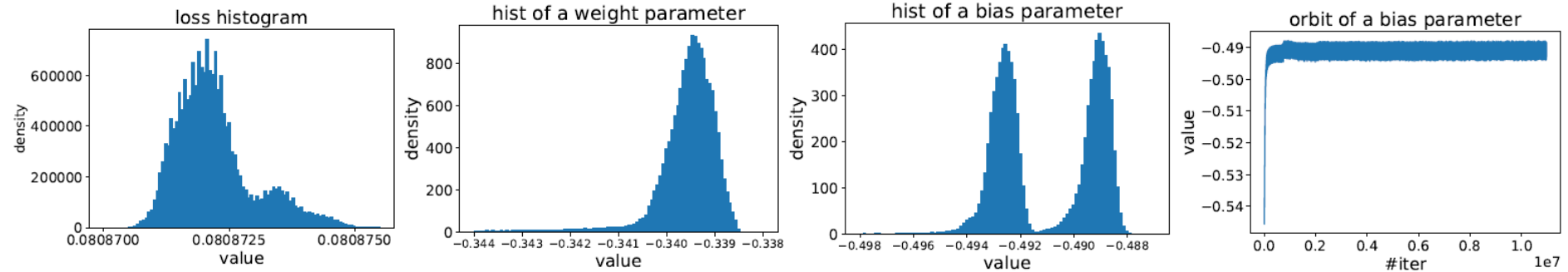
previously discussed
(f_0 strongly convex):

$$\frac{1}{Z} \exp\left(-\frac{2f_0(x)}{\eta\sigma^2}\right) dx + \mathcal{O}(\eta^2)$$

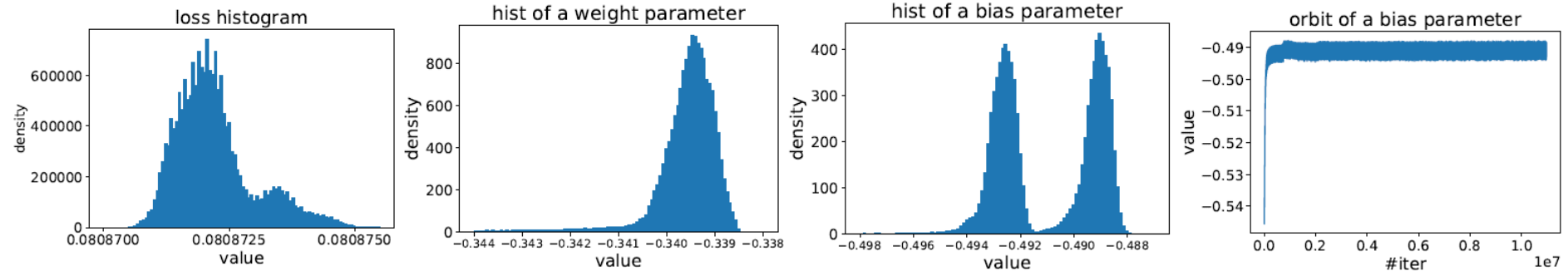
? non-convex f_0

? Is f for **real** problems **multiscale**?

EX 5-16-2 FF neural network, regression, UCI Airfoil Self-Noise dataset, large LR

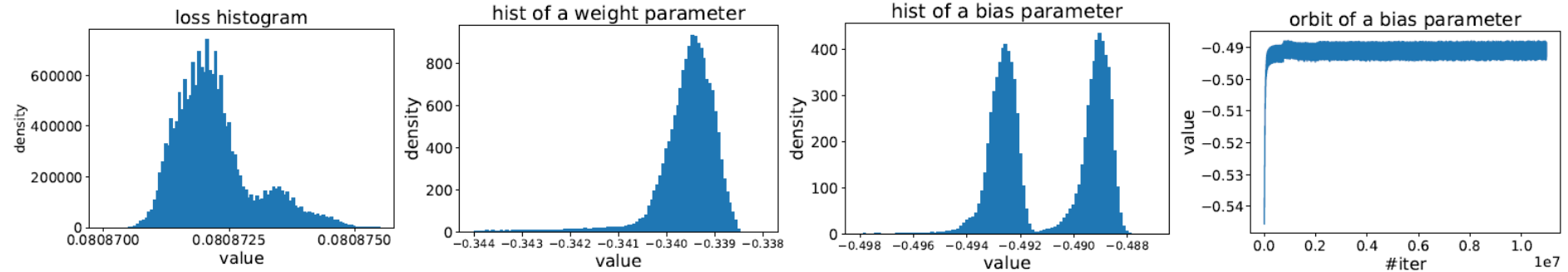


EX 5-16-2 FF neural network, regression, UCI Airfoil Self-Noise dataset, large LR



Theory multiscale data (mean + fluctuation) \rightarrow multiscale loss?

EX 5-16-2 FF neural network, regression, UCI Airfoil Self-Noise dataset, large LR



Theory multiscale data (mean + fluctuation) \rightarrow multiscale loss?

2 layer, periodic activation

chaotic dynamics can help **deterministic** GD
escape **microscopic** local minima
as it ~~optimizes the loss~~ samples a distribution

chaotic dynamics can help **deterministic** GD
escape microscopic local minima
as it ~~optimizes the loss~~ samples a distribution

? why this matters

chaotic dynamics can help **deterministic** GD
escape microscopic local minima
as it ~~optimizes the loss~~ samples a distribution

? why this matters

deeper minimum → better **training** (& thus test) accuracies

Thank *you* for your attention and feedback!

Support:

NSF DMS-1847802, ECCS-1936776

GT Cullen-Peck Scholarship

Emory-GT AI.Humanity Award



itsdynamical



MoleiTaoMath

