

Implicit regularisation of gradient algorithms

Scott Pesme



Loucas Pillaud-Vivien



Nicolas Flammarion



TML lab

EPFL

GRADIENT DESCENT MAXIMIZES THE MARGIN OF HOMOGENEOUS NEURAL NETWORKS

Kaifeng Lyu & Jian Li

Institute for Interdisciplinary Information Sciences

Tsinghua University

Beijing, China

vfleaking@gmail.com, lijian83@mail.tsinghua.edu.cn

The Implicit Bias of Gradient Descent on Separable Data

Daniel Soudry

Elad Hoffer

Mor Shpigel Nacson

Department of Electrical Engineering, Technion

Haifa, 320003, Israel

DANIEL.SOUDRY@GMAIL.COM

ELAD.HOFFER@GMAIL.COM

MOR.SHPIGEL@GMAIL.COM

Suriya Gunasekar

Nathan Srebro

Toyota Technological Institute at Chicago

Chicago, Illinois 60637, USA

SURIYA@TTIC.EDU

NATI@TTIC.EDU

Implicit Bias of Gradient Descent for Wide Two-layer Neural Networks Trained with the Logistic Loss

Lénaïc Chizat

LENAIC.CHIZAT@UNIVERSITE-PARIS-SACLAY.FR

Laboratoire de Mathématiques d'Orsay, CNRS, Université Paris-Saclay, France

Francis Bach

FRANCIS.BACH@INRIA.FR

Research University, Paris, France

Implicit Bias of Gradient Descent on Linear Convolutional Networks

Suriya Gunasekar

TTI at Chicago, USA

suriya@ttic.edu

Jason D. Lee

USC Los Angeles, USA

jasonlee@marshall.usc.edu

Daniel Soudry

Technion, Israel

daniel.soudry@gmail.com

Nathan Srebro

TTI at Chicago, USA

nati@ttic.edu

GRADIENT DESCENT ALIGNS THE LAYERS OF DEEP LINEAR NETWORKS

Ziwei Ji & Matus Telgarsky

Department of Computer Science

University of Illinois at Urbana-Champaign

{ziwei2, mjt}@illinois.edu

Implicit Bias in Deep Linear Classification: Initialization Scale vs Training Accuracy

Edward Moroshko

edward.moroshko@gmail.com

Technion

Blake Woodworth

blake@ttic.edu

TTI Chicago

Suriya Gunasekar

suriya@ttic.edu

Microsoft Research

Jason D. Lee

asonlee@princeton.edu

Princeton University

Nathan Srebro

nati@ttic.edu

TTI Chicago

Daniel Soudry

daniel.soudry@gmail.com

Technion

Implicit Regularization in ReLU Networks with the Square Loss

Gal Vardi

GAL.VARDI@WEIZMANN.AC.IL

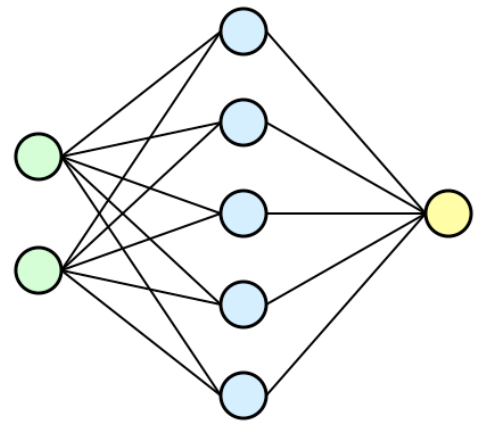
and **Ohad Shamir**

OHAD.SHAMIR@WEIZMANN.AC.IL

Weizmann Institute of Science

Motivation

Training architecture



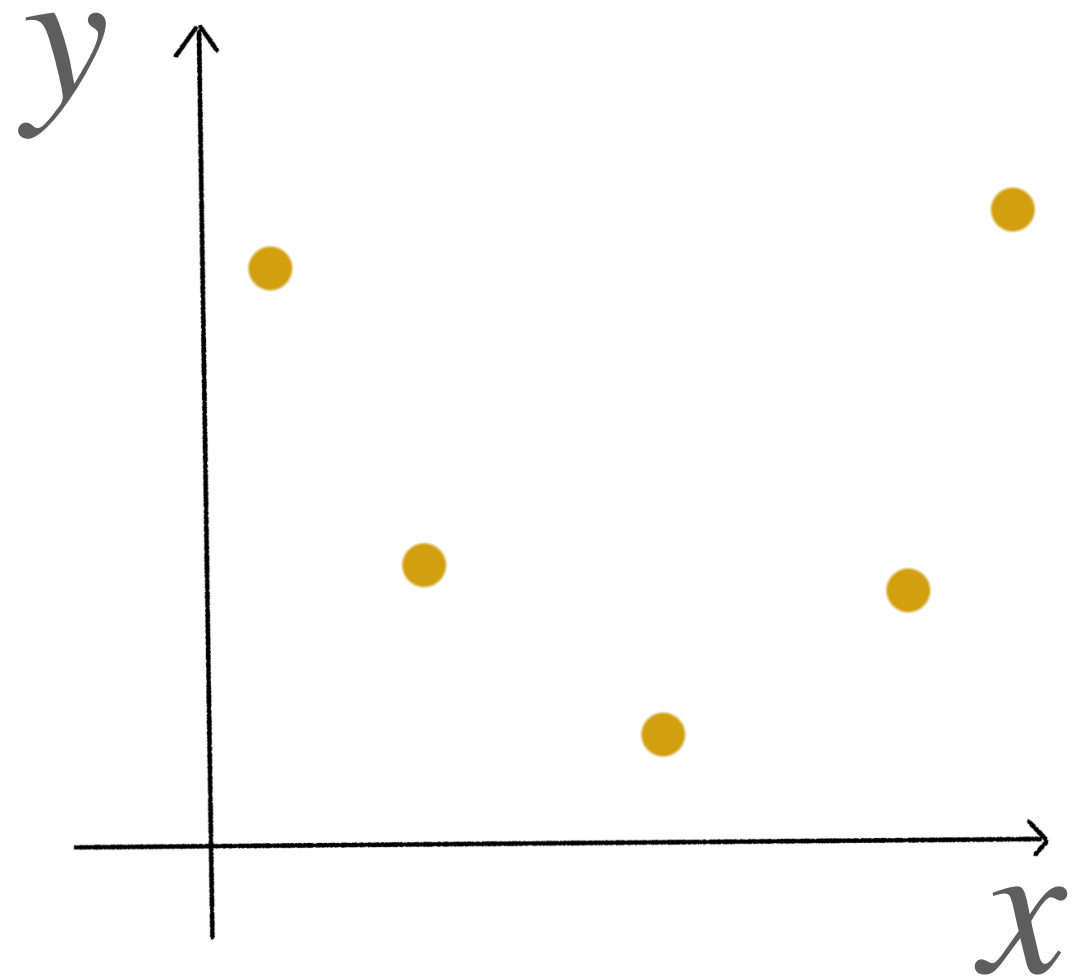
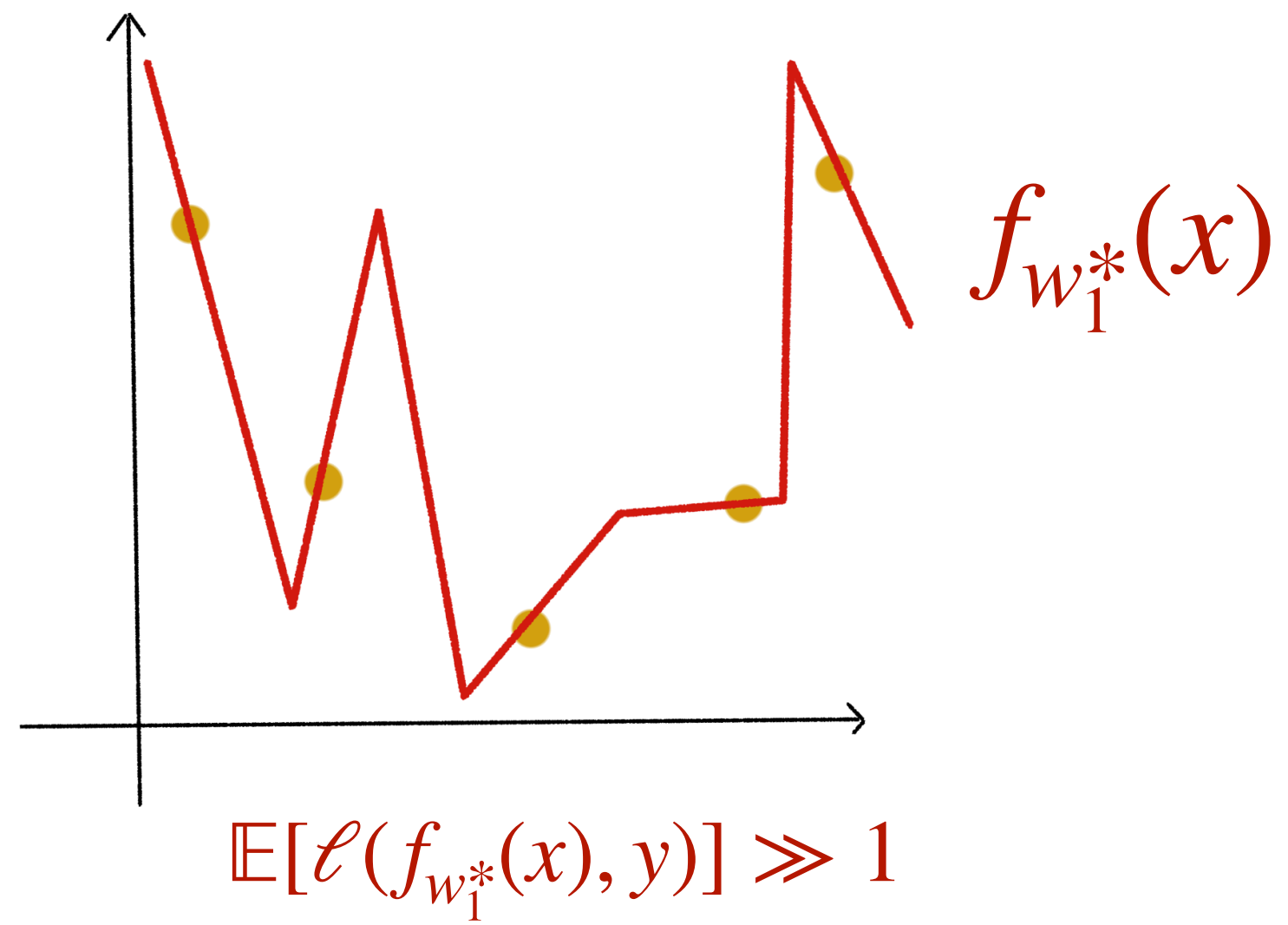
Training dataset



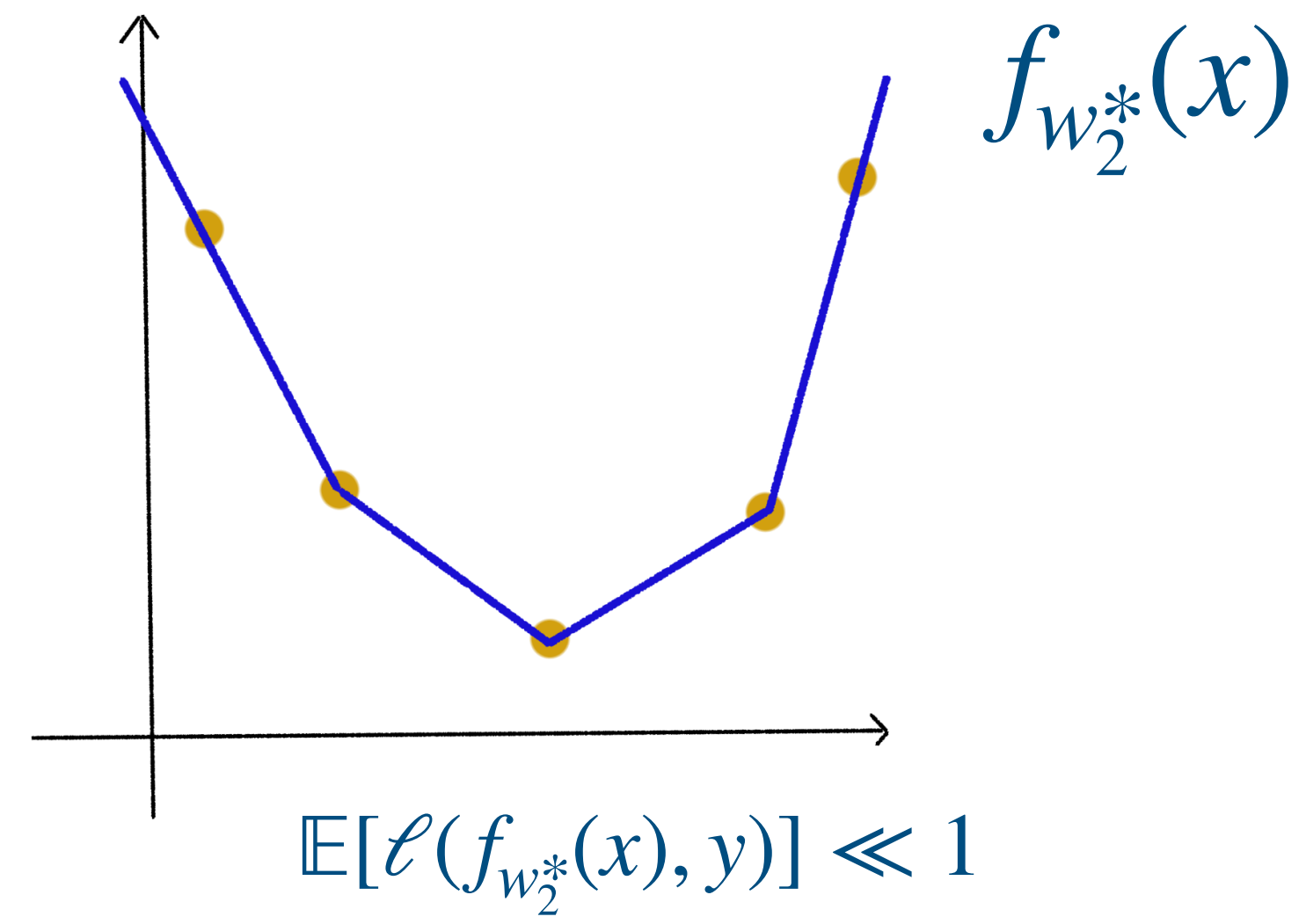
Training algorithm

Some ERM

An infinity of interpolating solutions



GD
SGD
etc.



A panoply of algorithms:

Some
ERM

GD

SGD

+ momentum

Batch
Normalisation

Etc.

All can lead to zero training error but do not generalise the same.

Cifar-10 dataset:

Classification task!

ResNet-18 (with BN):

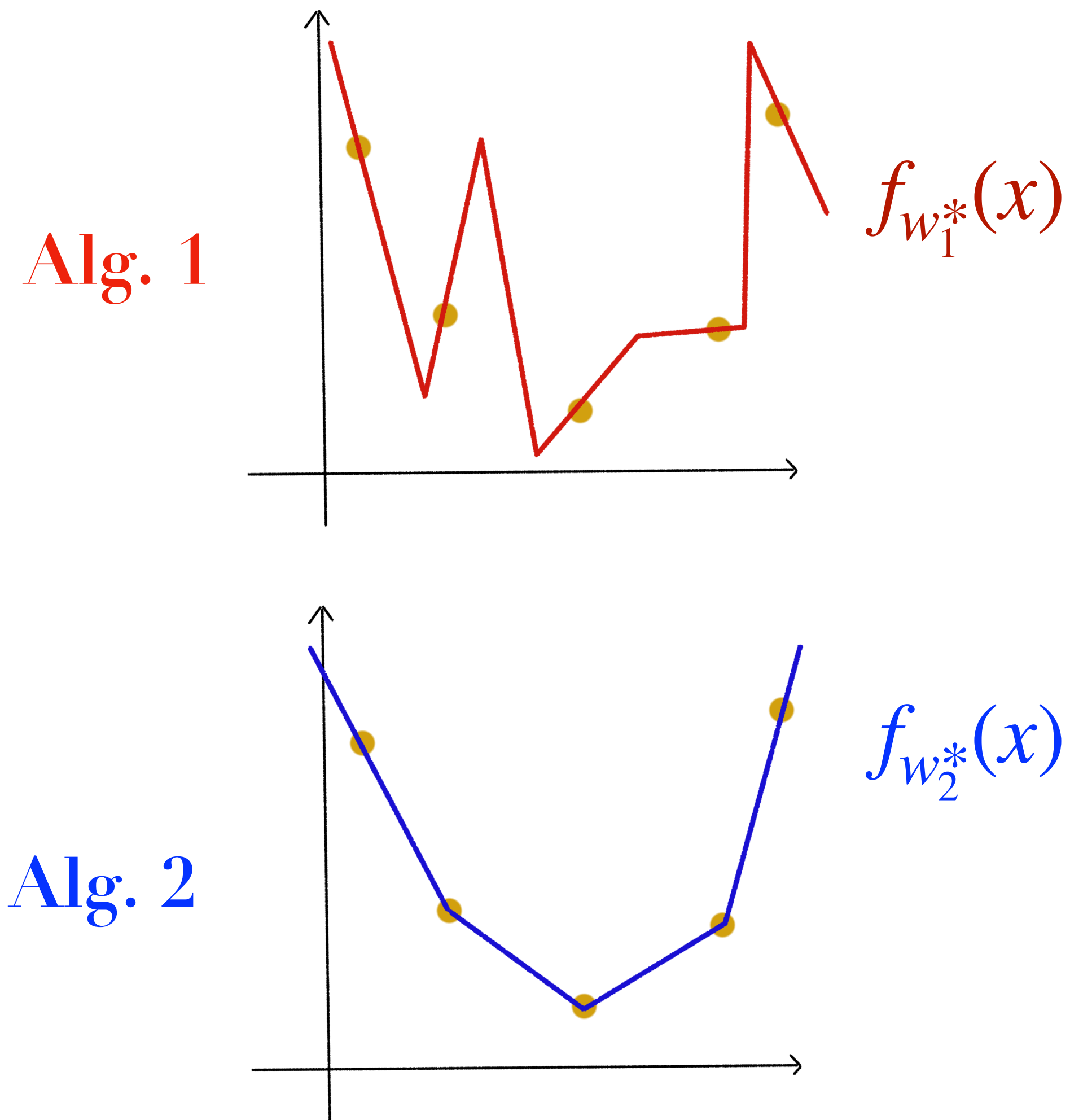
	Train accuracy	Test accuracy
GD	100%	~72% (??)
SGD	100%	85%
SGD + momentum	100%	89%
SGD + mom + DA + ℓ_2	100%	95%
GD + mom + DA + ℓ_2	100%	87%

Stochastic Training is Not Necessary for Generalization, Geiping et al. 2021

Bad Global Minima Exist and SGD Can Reach Them, Liu et al. 2020

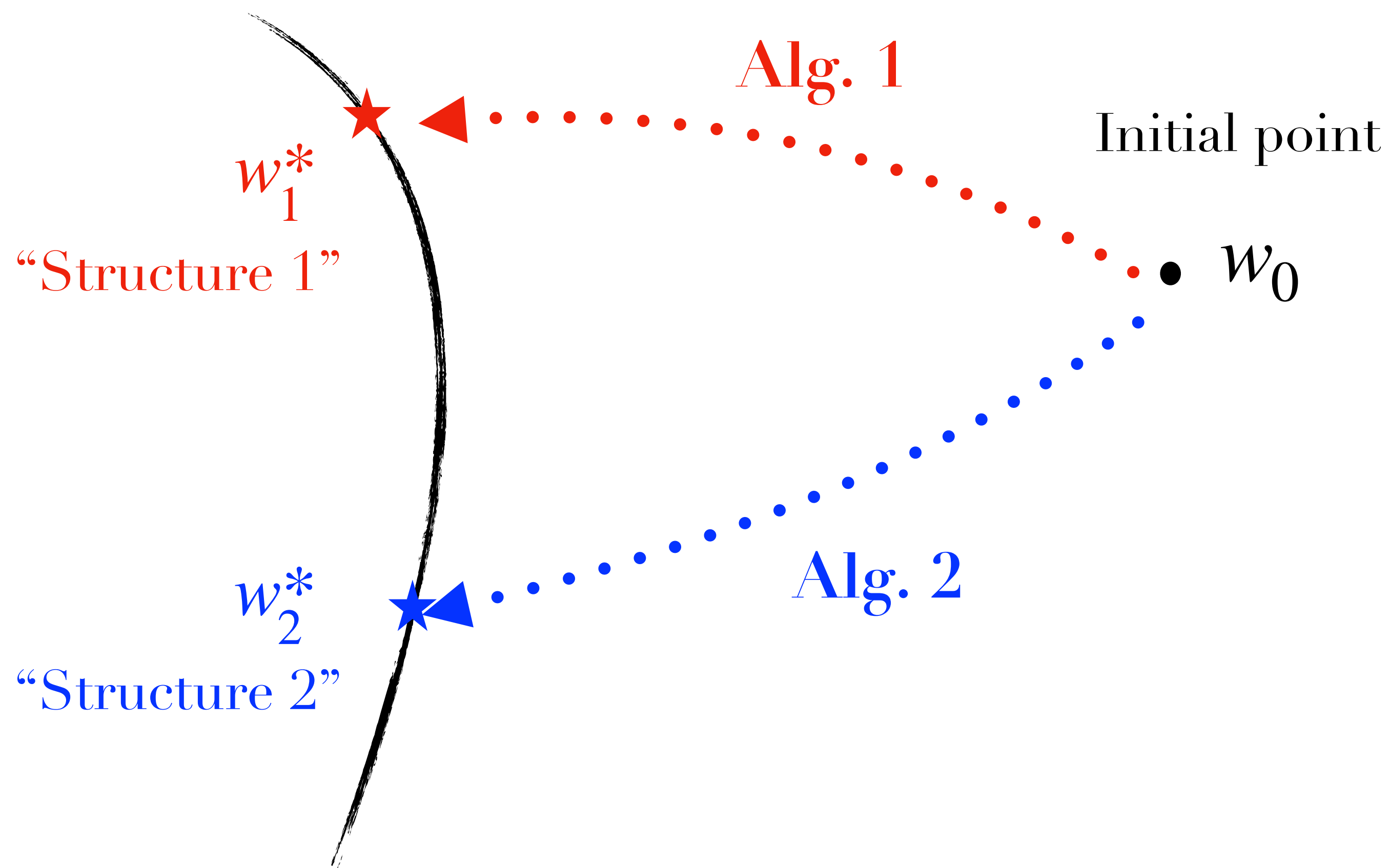
Revisiting Small Batch Training For Deep Neural Networks, Masters and Luschi 2018

What does all this mean ?



Interpolation manifold

$$\{w^* \text{ s.t. } f_{w^*}(x_i) = y_i \forall i\}$$



“Algorithmic implicit bias” : the algorithm “chooses” a particular solution.

Implicit bias and minimal norm solutions (for regression)

Minimise $L(w) = \frac{1}{2n} \sum_{i=1}^n (f_w(x_i) - y_i)^2$ with some algorithm.

“It turns out”:

$$L(w_{\infty}^{alg}) = 0$$

\Leftrightarrow

$$f_{w_{\infty}^{alg}}(x_i) = y_i$$

and

w_{∞}^{alg} enjoys a “nice” structure

\equiv

$$w_{\infty}^{alg} = \arg \min_{w, \forall i, f_w(x_i)=y_i} R_{alg}(w)$$

$\|w\|_2$

$\|w\|_1$

Contrast with explicit regularisation:

$$\min_{w \in \mathbb{R}^d} \text{Reg}L(w) := \frac{1}{n} \sum_{i=1}^n \ell_i(w) + \lambda R(w)$$

(Often unique)
(Not an interpolator !)

A few disclaimers:

We hope to exhibit:

$$w_{\infty}^{alg} = \arg \min_{w, \forall i, f_w(x_i)=y_i} R_{alg}(w)$$

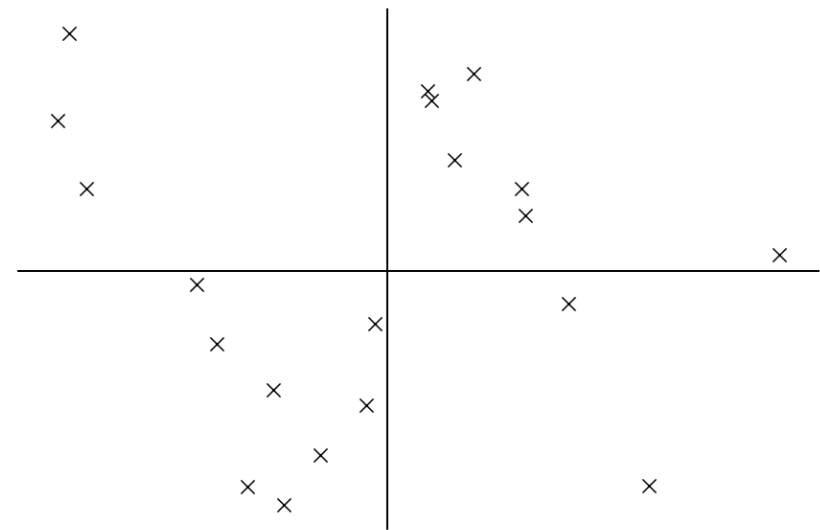
- We will only consider **regression** tasks
- The characterisation of the solution **does not (on its own) say anything about generalisation !**
 - overfitting the training set is not always good (but often works in practice)
 - the benign overfitting literature (partly) covers the generalisation properties of min norm interpolators
 - the generalisation questions **depend on the true distribution**, but not the implicit regularisation problem

Toy examples: implicit bias doesn't explain (on its own) generalisation

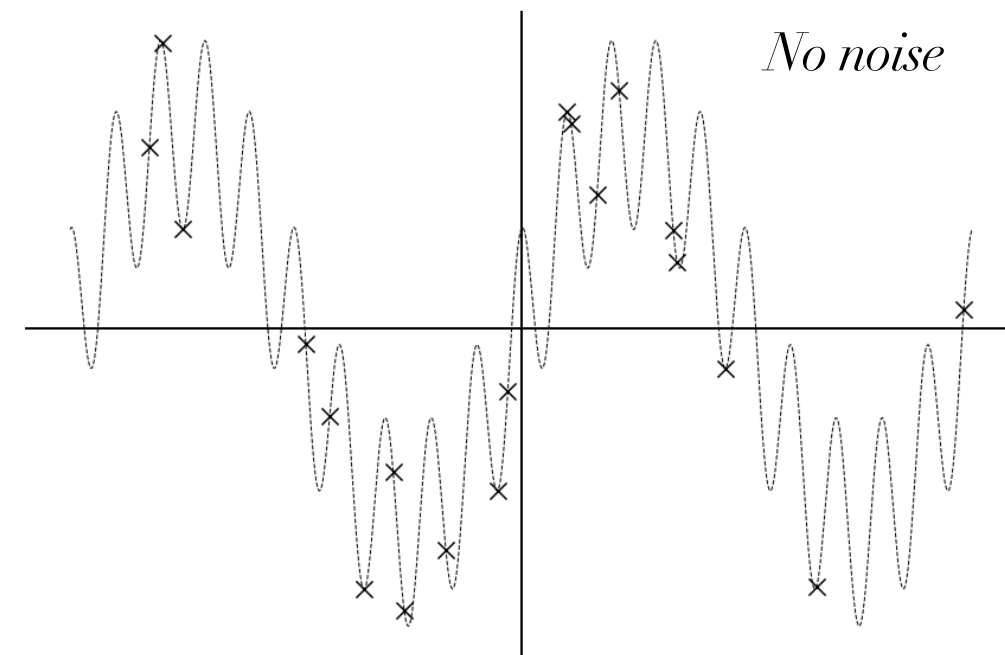
Samples $(x_i, y_i)_{1 \leq i \leq n} \in \mathbb{R} \times \mathbb{R}$ from some distribution \mathcal{D} .

We want to linearly interpolate with feature expansion $\phi(x) = (\frac{1}{i} \cos(2\pi i x), \frac{1}{i} \sin(2\pi i x))_{1 \leq i \leq d/2} \in \mathbb{R}^d$: $f_w(x) = \langle w, \phi(x) \rangle$.

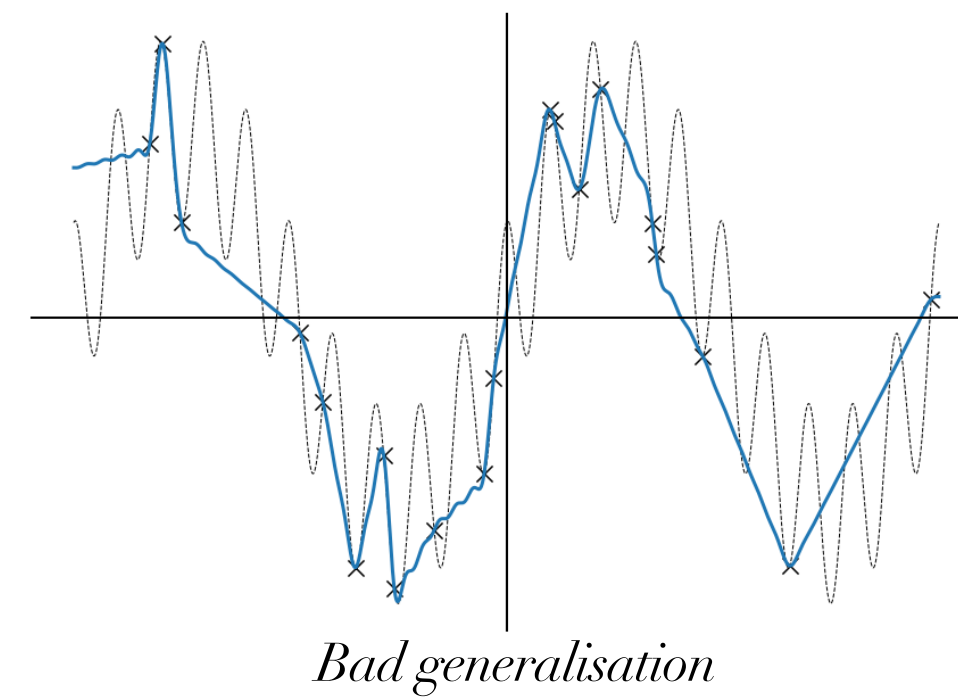
Training set 1



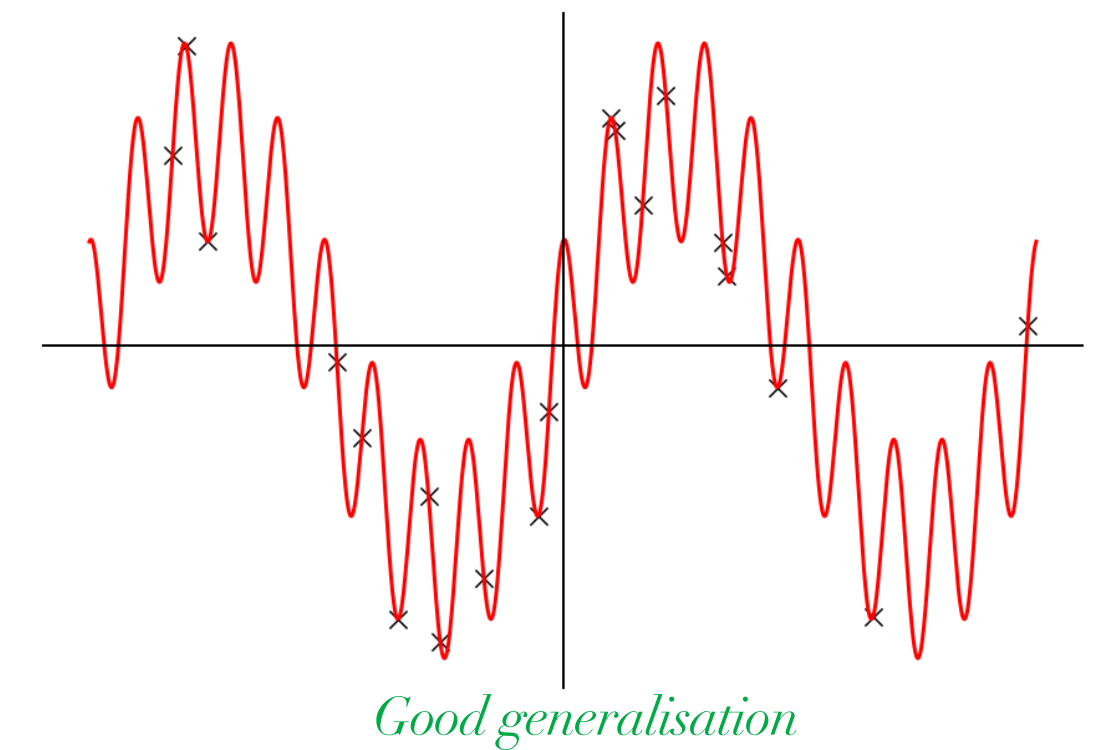
from distribution 1 (sparse)



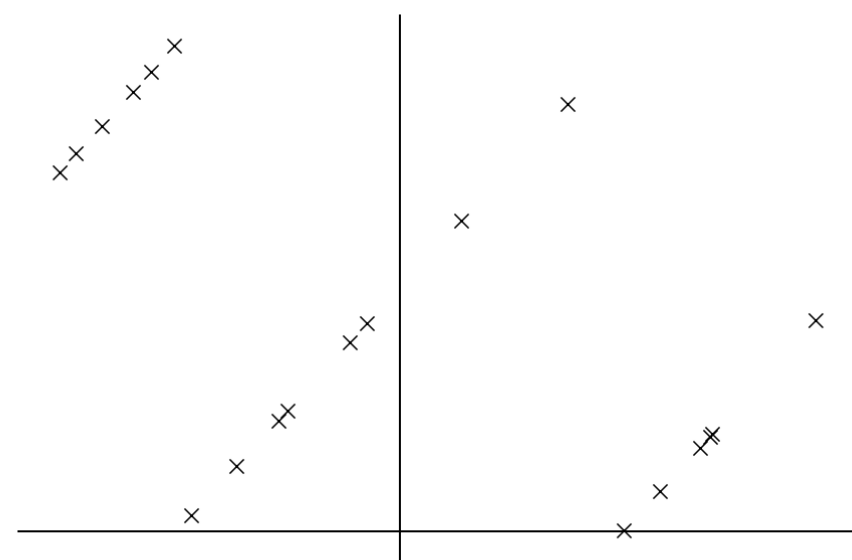
$\arg \min_{w, \forall i, f_w(x_i)=y_i} \|w\|_2$



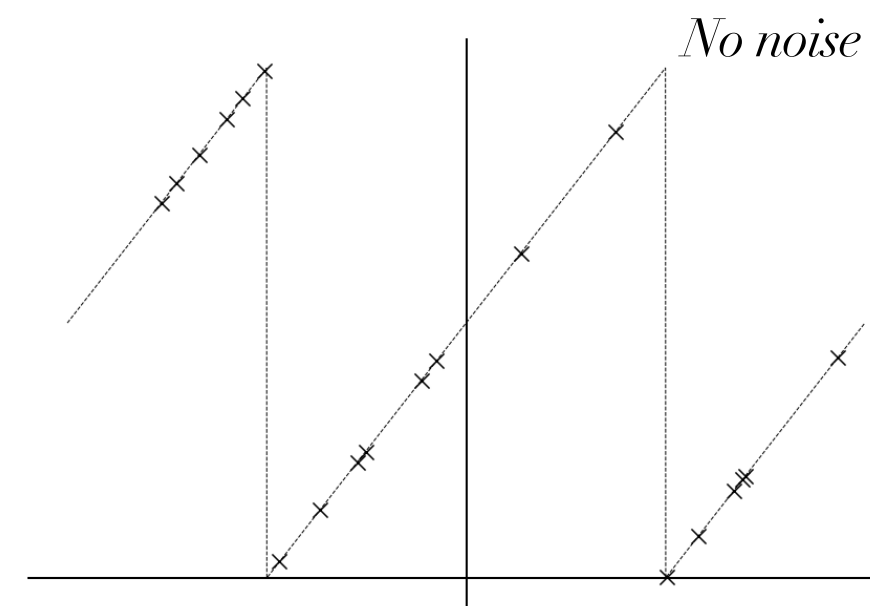
$\arg \min_{w, \forall i, f_w(x_i)=y_i} \|w\|_1$



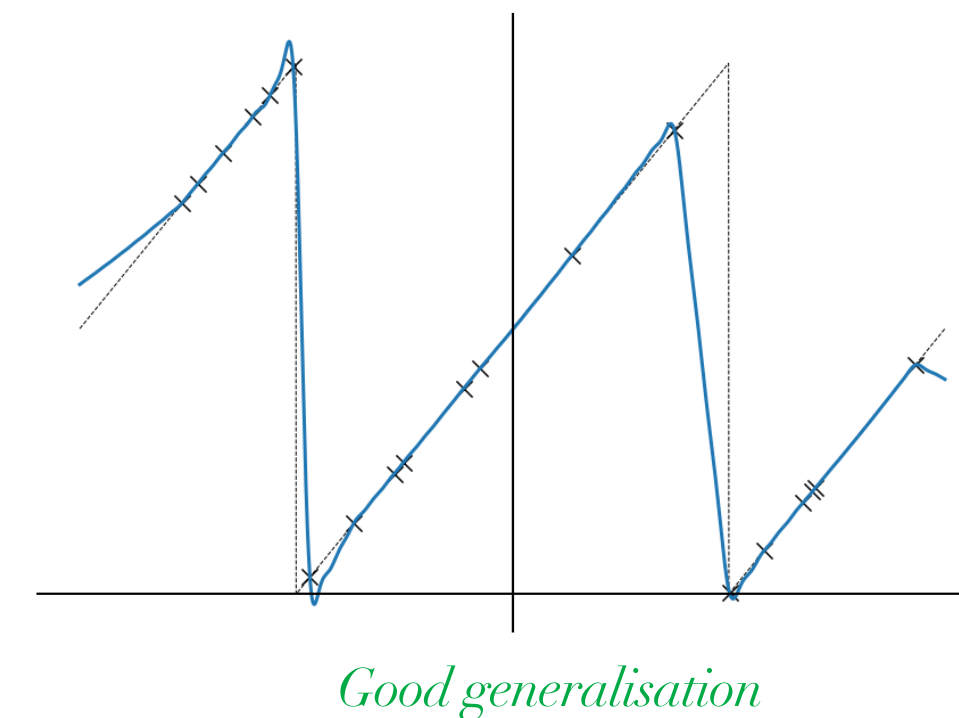
Training set 2



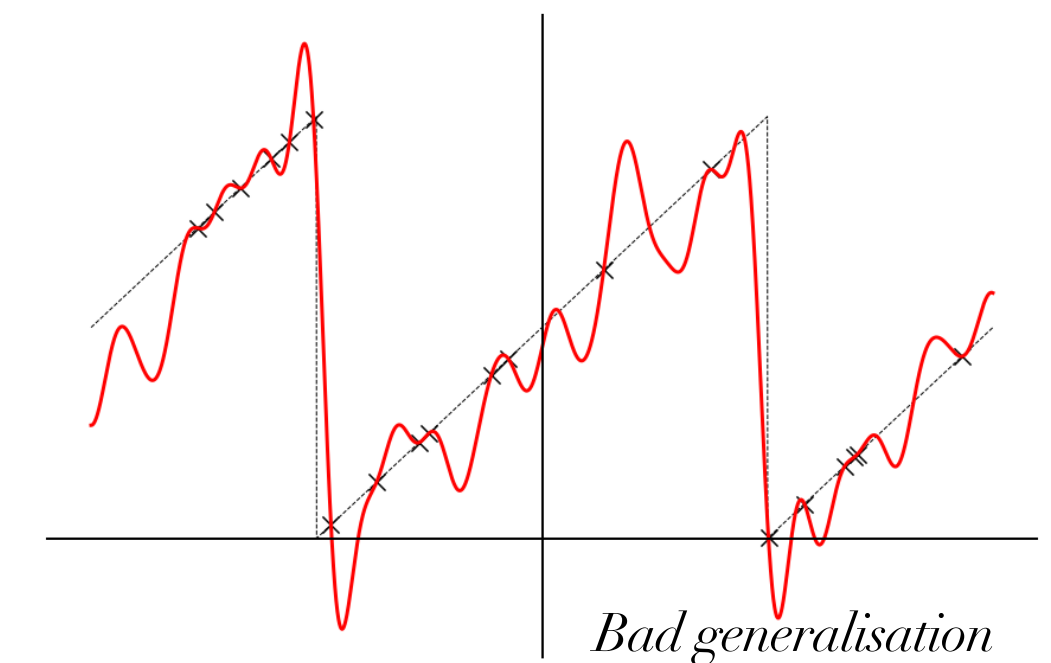
from distribution 2 (dense)



$\arg \min_{w, \forall i, f_w(x_i)=y_i} \|w\|_2$



$\arg \min_{w, \forall i, f_w(x_i)=y_i} \|w\|_1$



Depending on the true data distribution, **“Structure 1”** \preceq **“Structure 2”**.

Back to implicit bias / regularisation:

Can you give me a simple example showing this phenomenon ?

Simplest example: linear regression

Unique root loss (can be non-convex)

$$L(w) = \frac{1}{2n} \sum_{i=1}^n (y_i - \langle w, x_i \rangle)^2$$

But more generally still true with

$$L(w) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, \langle w, x_i \rangle)$$

$$f_w(x) = \langle w, x \rangle$$

$$\text{GD: } w_{t+1} = w_t + \underbrace{\gamma \frac{1}{n} \sum_i (y_i - \langle x_i, w_t \rangle) x_i}_{\in \text{span}(x_1, \dots, x_n)}$$

$$\text{SGD: } w_{t+1} = w_t + \underbrace{\gamma (y_{i_t} - \langle x_{i_t}, w_t \rangle) x_{i_t}}_{\in \text{span}(x_1, \dots, x_n)}$$

Implicit regularisation

$$w_{\infty}^{GD/SGD} \in w_0 + \text{span}(x_1, \dots, x_n)$$

$$\forall i, \langle w_{\infty}^{GD/SGD}, x_i \rangle = y_i$$

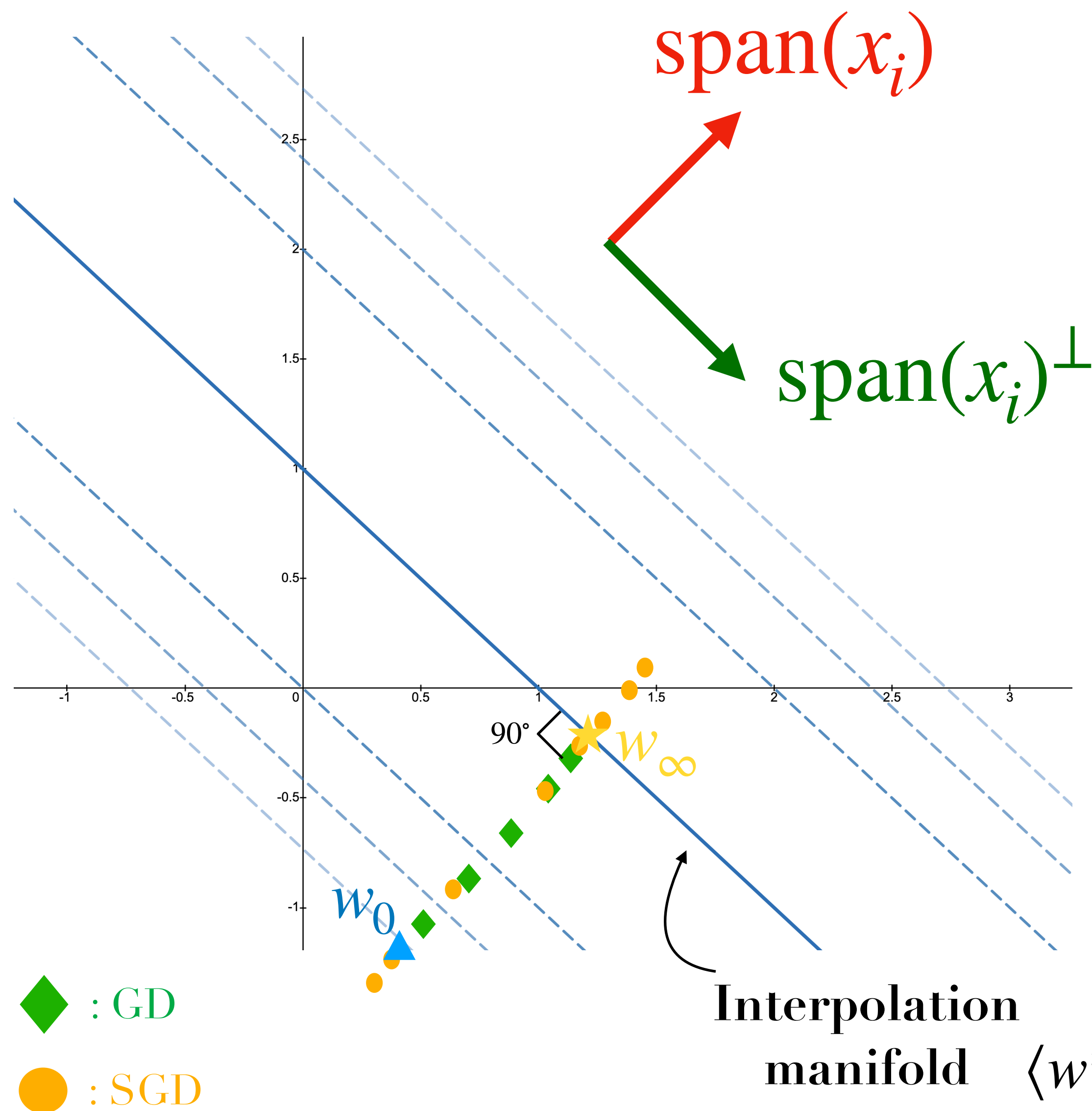
implies

(Pythagorean
Theorem)

Pythagoras et al. 500 BC

$$w_{\infty}^{GD/SGD} = \underset{w, \forall i, \langle w, x_i \rangle = y_i}{\text{argmin}} \|w - w_0\|_2^2$$

Simplest example: linear regression



$$L(w) = \frac{1}{2n} \sum_{i=1}^n (y_i - \langle w, x_i \rangle)^2$$

Implicit regularisation

$$w_{\infty}^{GD/SGD} = \operatorname{argmin}_{w, \forall i, \langle w, x_i \rangle = y_i} \|w - w_0\|_2^2$$

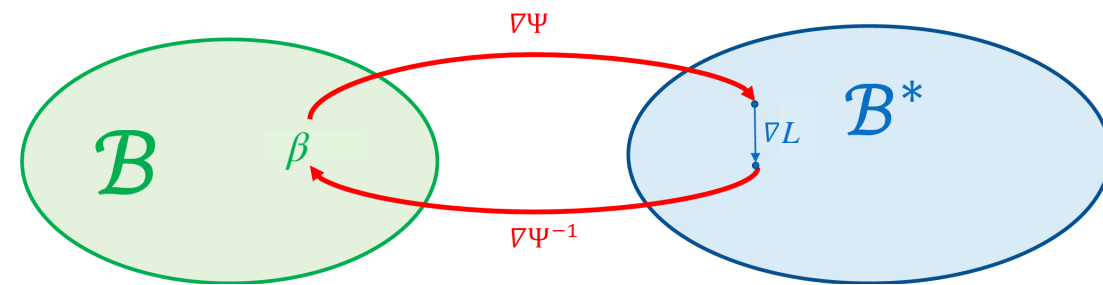
Interpolation manifold $\langle w^*, x_i \rangle = y_i, \forall i$

Second simplest: Mirror descent

$$\nabla \Psi(\beta_{t+1}) = \nabla \Psi(\beta_t) - \gamma \nabla L(\beta_t)$$

Ψ is a convex and differentiable potential.

$$\Psi(\beta) = \|\beta\|_2^2 : \text{back to GD}$$



$$\text{MD: } \nabla \Psi(\beta_{t+1}) = \nabla \Psi(\beta_t) + \underbrace{\gamma \frac{1}{n} \sum_i (y_i - \langle x_i, \beta_t \rangle) x_i}_{\in \text{span}(x_1, \dots, x_n)}$$

$$\text{SMD: } \nabla \Psi(\beta_{t+1}) = \nabla \Psi(\beta_t) + \underbrace{\gamma (y_{i_t} - \langle x_{i_t}, \beta_t \rangle) x_{i_t}}_{\in \text{span}(x_1, \dots, x_n)}$$

$$\beta_\infty^{MD} ? \quad \nabla \Psi(\beta_\infty^{MD}) \in \nabla \Psi(\beta_0) + \text{span}(x_i) \quad \Longrightarrow \quad \beta_\infty^{MD} = \underset{\beta, \langle \beta, x_i \rangle = y_i}{\text{argmin}} D_\Psi(\beta, \beta_0)$$

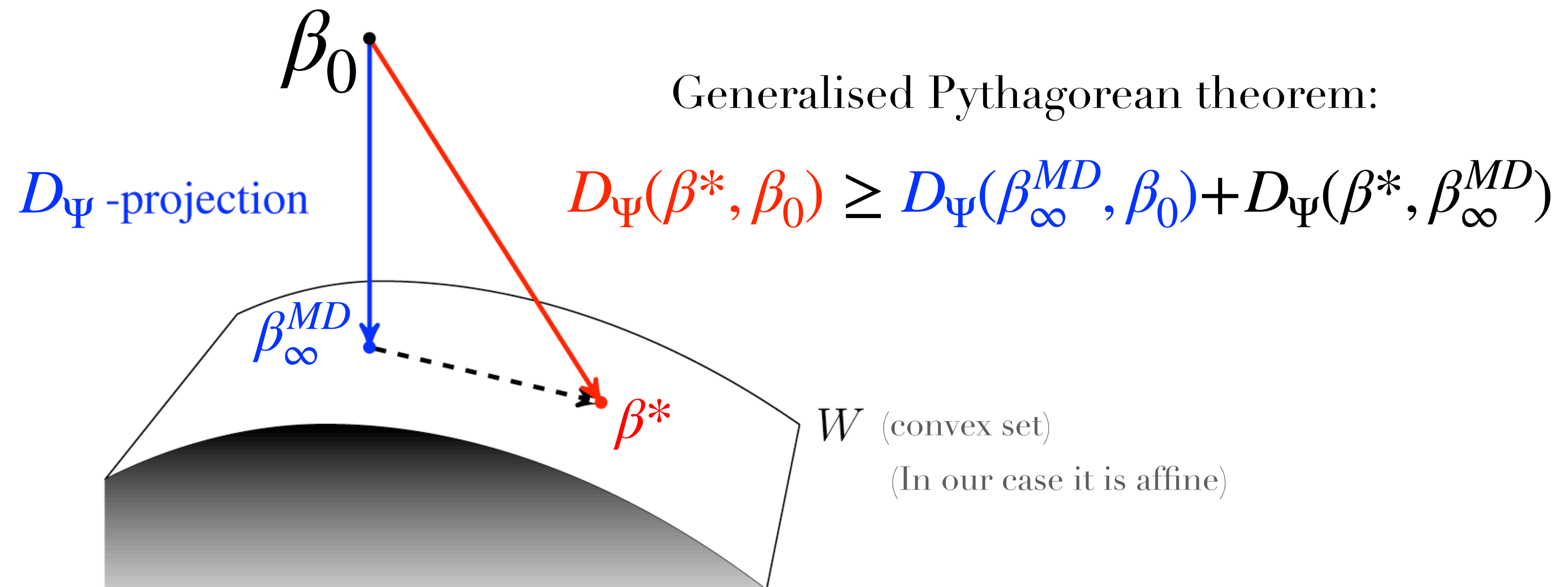
(Pythagorean Theorem)

Regression with linear models, recap.

$$L(w) = \frac{1}{2n} \sum_{i=1}^n (y_i - \langle w, x_i \rangle)^2$$

$$w_{\infty}^{SGD} = w_{\infty}^{GD} = \arg \min_{w, \forall i, \langle w, x_i \rangle = y_i} \|w - w_0\|_2^2 \quad (= w_{\infty}^{mom} = w^{GF} = \dots)$$

$$\beta_{\infty}^{MD} = \beta_{\infty}^{SMD} = \arg \min_{w, \forall i, \langle \beta, x_i \rangle = y_i} D_{\Psi}(\beta, \beta_0) \quad (= \beta_{\infty}^{mom} = \beta^{MF} = \dots) \quad (= \text{Proj}_{\{\langle \beta, x_i \rangle = y_i\}}^{\Psi}(\beta_0))$$



MD on linear models:

$$\beta_{\infty}^{MD} = \beta_{\infty}^{SMD} = \arg \min_{w, \forall i, \langle \beta, x_i \rangle = y_i} D_{\Psi}(\beta, \beta_0)$$

For the intuition, recall that $\beta_{t+1} = \operatorname{argmin}_{\beta \in \mathbb{R}^d} L(\beta_t) + \langle \nabla L(\beta_t), \beta - \beta_t \rangle + \frac{1}{\gamma} D_{\Psi}(\beta, \beta_t)$

➡ Independent of stochasticity

But who (in DL) cares about mirror descent ?

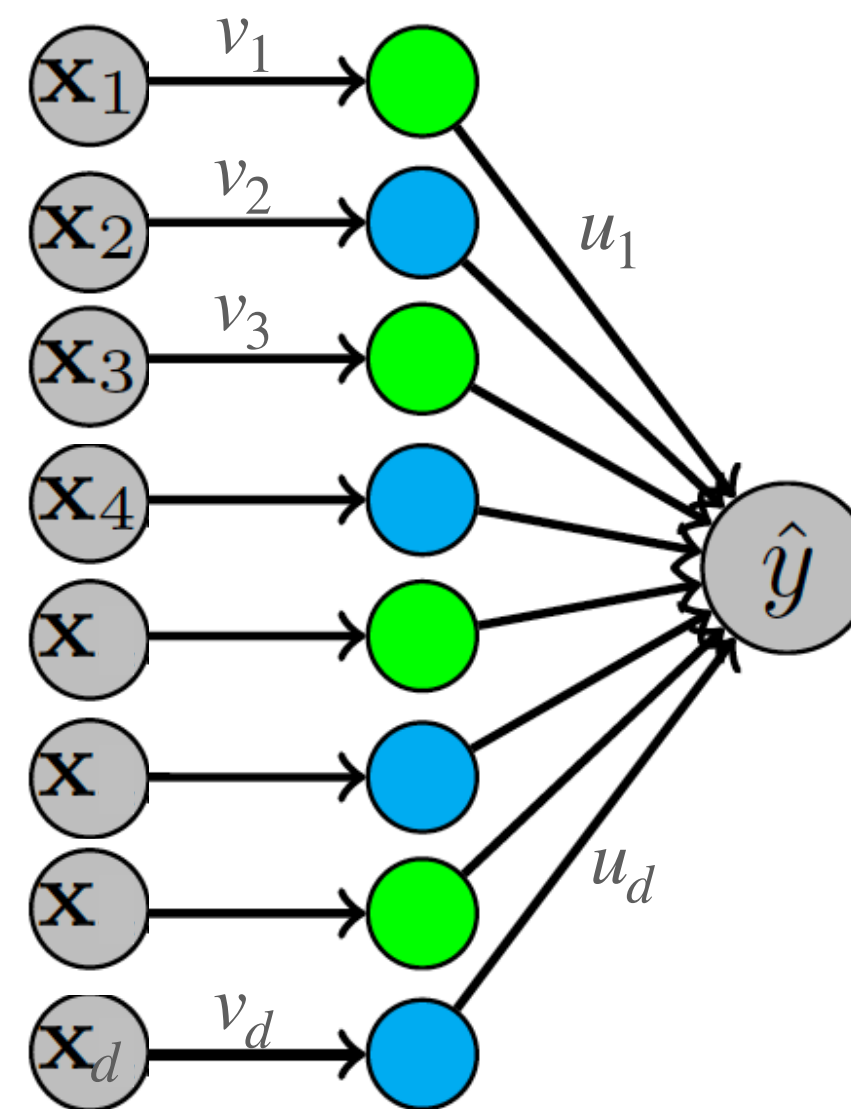
Mirror descent is a framework in which things are easy:

- Convergence of the iterates and of the training loss
- Tight rates
- Implicit bias

A “practical example”: 2-layer diagonal linear network.

Architecture

Diagonal linear network :



$$f_w(x) = \langle u \odot v, x \rangle$$

$$w = \begin{pmatrix} u \\ v \end{pmatrix} \in \mathbb{R}^{2d}$$

Square-loss

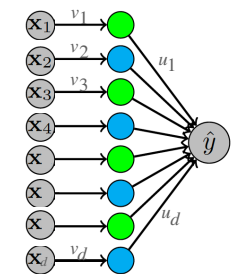
$$\min_{w \in \mathbb{R}^{2d}} L(w) = \frac{1}{4n} \sum_{i=1}^n (y_i - \underbrace{\langle u \odot v, x_i \rangle}_{\beta_w})^2$$

Non-convex in w

Final model is linear: but training is changed.

Hidden mirror descent

Setting $L(w) = \frac{1}{4n} \sum_{i=1}^n (y_i - \langle u \odot v, x_i \rangle)^2$



$w = \begin{pmatrix} u \\ v \end{pmatrix} \in \mathbb{R}^{2d}$ $\beta_w := u \odot v \in \mathbb{R}^d$

Gradient flow on the neurons:

$$\begin{aligned} du_t &= -\nabla_u L(w_t) dt & u_{t=0} &= \alpha \mathbf{1} \in \mathbb{R}^d \\ dv_t &= -\nabla_v L(w_t) dt & v_{t=0} &= \mathbf{0} \in \mathbb{R}^d \end{aligned} \quad \beta_{w_{t=0}} = \mathbf{0} \in \mathbb{R}^d$$

What about the dynamics of $\beta_t := \beta_{w_t} = u_t \odot v_t$?

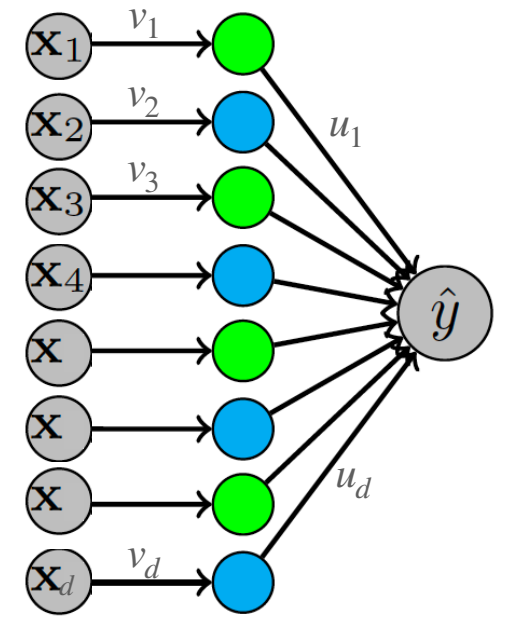
Turns out that:

$$\frac{d \nabla \phi_\alpha(\beta_t)}{dt} = -\nabla_\beta L(\beta_t) \quad \text{i.e. continuous mirror descent with } \Psi = \phi_{\underbrace{\alpha}_{\text{Initialisation scale!}}}$$

Implicit bias of the gradient flow for 2 layer diagonal linear network

Architecture:

$$\min_{w \in \mathbb{R}^{2d}} L(w) = \frac{1}{4n} \sum_{i=1}^n (y_i - \underbrace{\langle u \odot v, x_i \rangle}_{\beta_w})^2$$



Algorithm:

Gradient flow on the neurons w_t starting at $u_{t=0} = \alpha \mathbf{1} \in \mathbb{R}^d$
 $v_{t=0} = \mathbf{0} \in \mathbb{R}^d$

Implicit bias:

$$\beta_\infty^\alpha = \arg \min_{\beta, \langle \beta, x_i \rangle = y_i} \phi_\alpha(\beta) \quad (= D_{\phi_\alpha}(\beta, \beta_0 = 0))$$

where $\phi_\alpha \underset{\alpha \rightarrow \infty}{\sim} \|\cdot\|_2$ and $\phi_\alpha \underset{\alpha \rightarrow 0}{\sim} \|\cdot\|_1$

(+ convergence, rates etc.)

Implicit bias:

$$\beta_\infty^\alpha = \arg \min_{\beta, \langle \beta, x_i \rangle = y_i} \phi_\alpha(\beta)$$

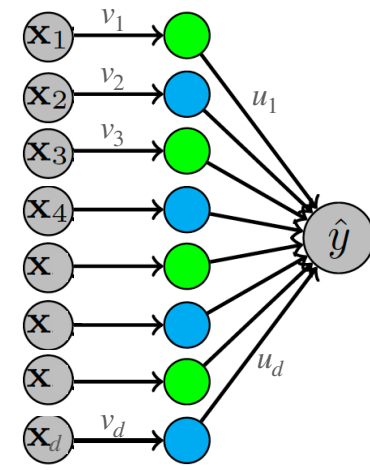
GF on the neurons

$$u_{t=0} = \alpha \mathbf{1} \in \mathbb{R}^d$$

$$v_{t=0} = \mathbf{0} \in \mathbb{R}^d$$

MD on the predictor

$$\beta_t := u_t \odot v_t ?$$



$$\beta, \langle \beta, x_i \rangle = y_i$$

where

$$\phi_\alpha \underset{\alpha \rightarrow 0}{\sim} \|\cdot\|_1$$

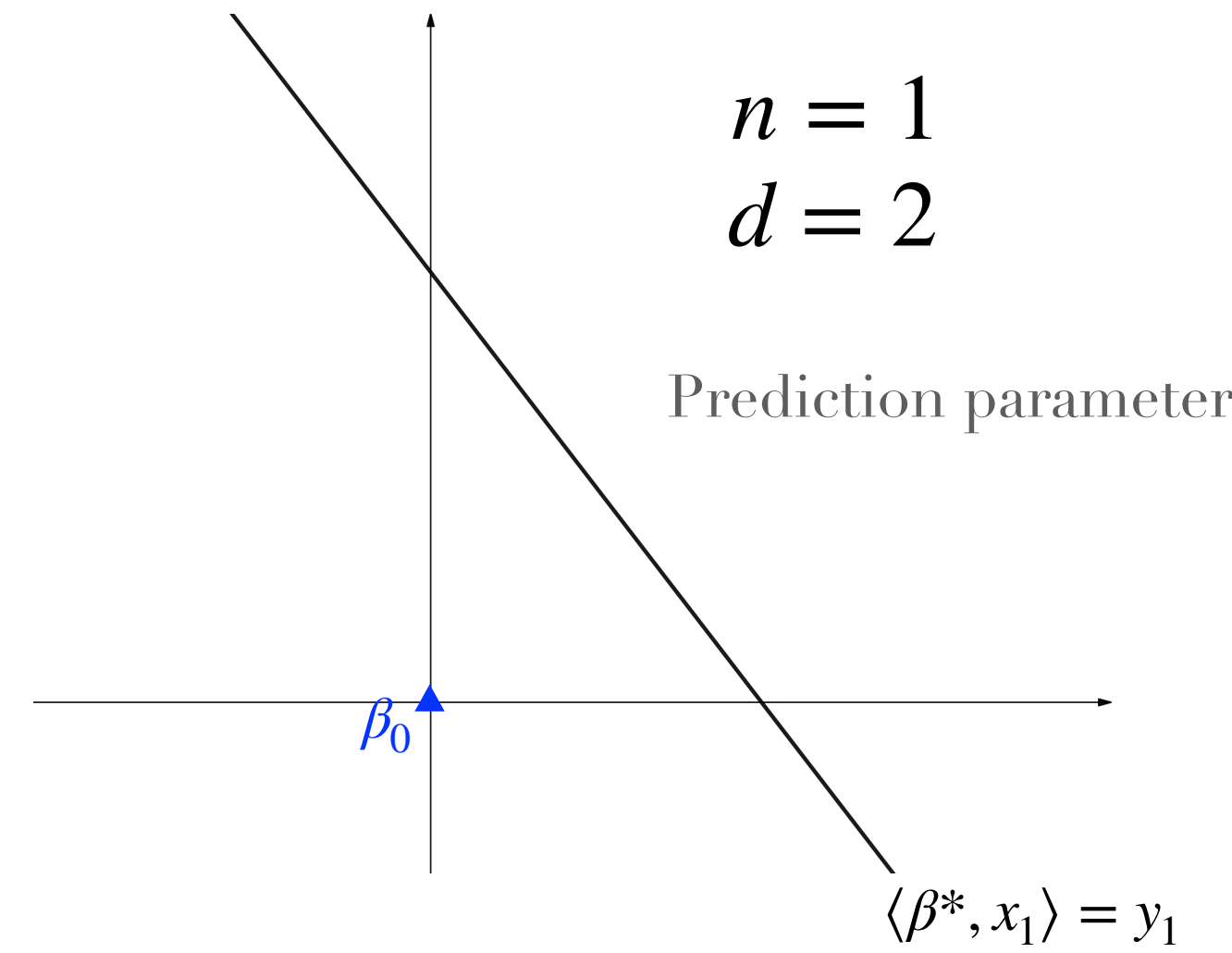
$$\phi_\alpha \underset{\alpha \rightarrow \infty}{\sim} \|\cdot\|_2$$

Toy illustration

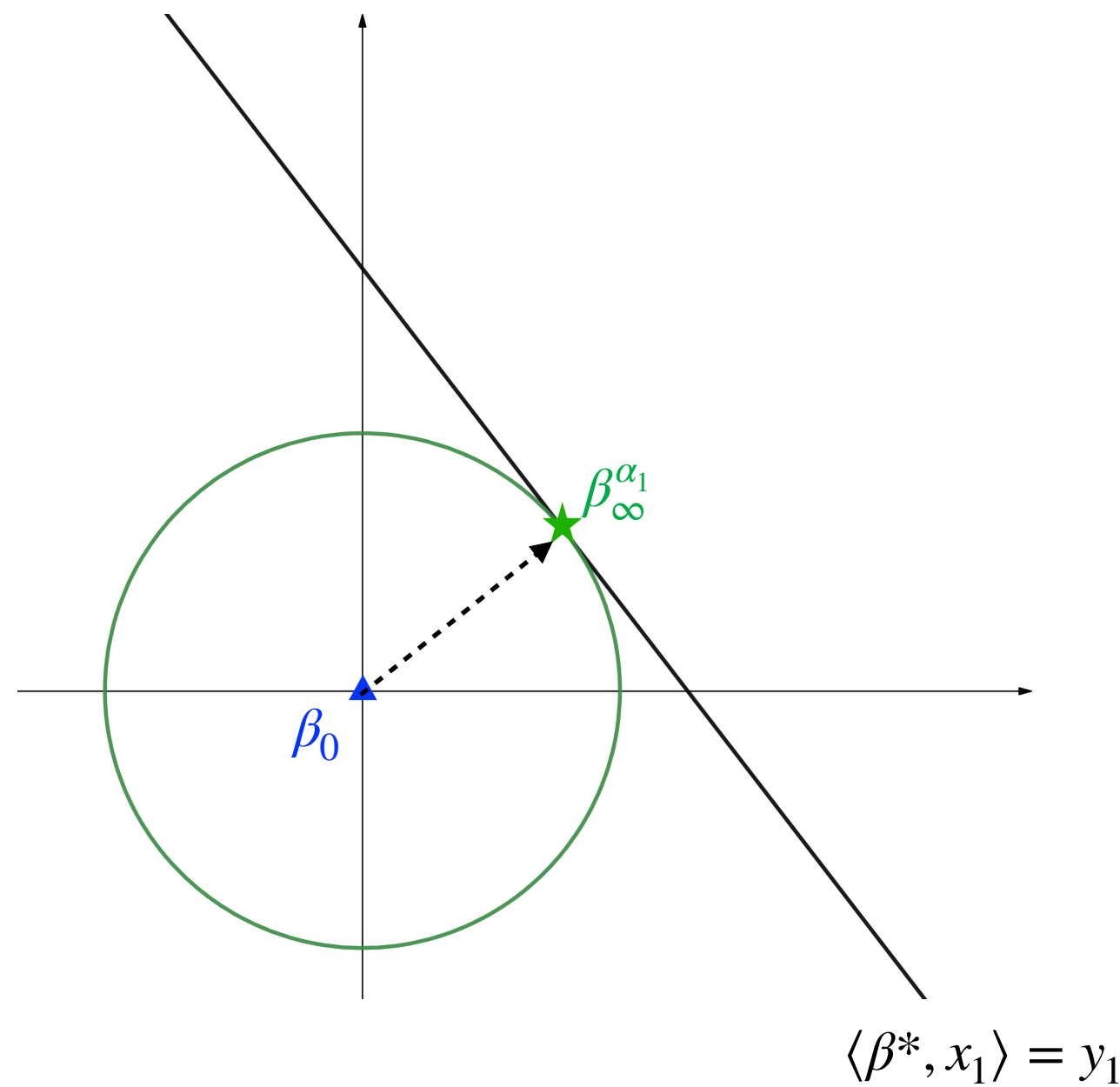
$$n = 1$$

$$d = 2$$

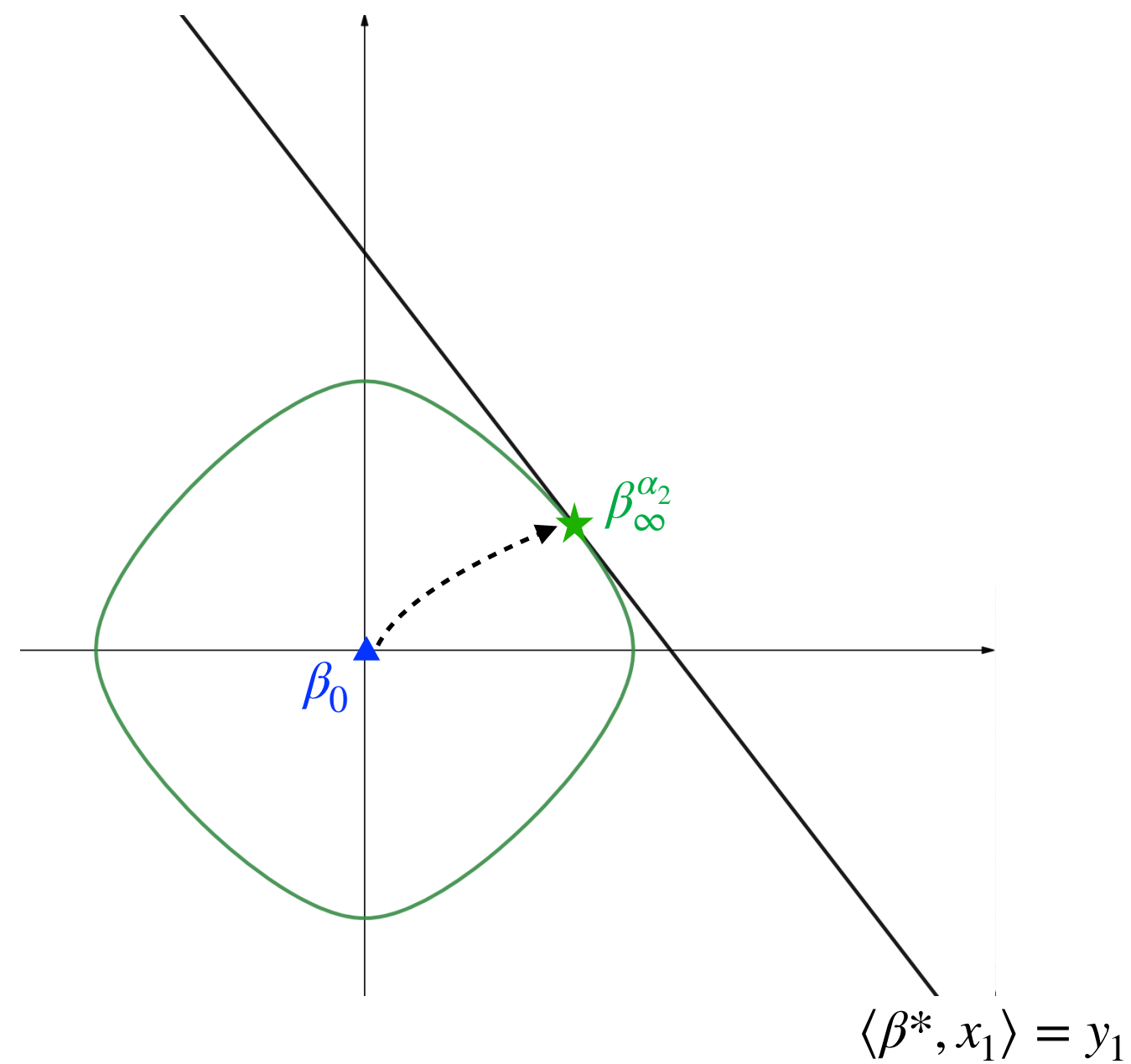
Prediction parameter space



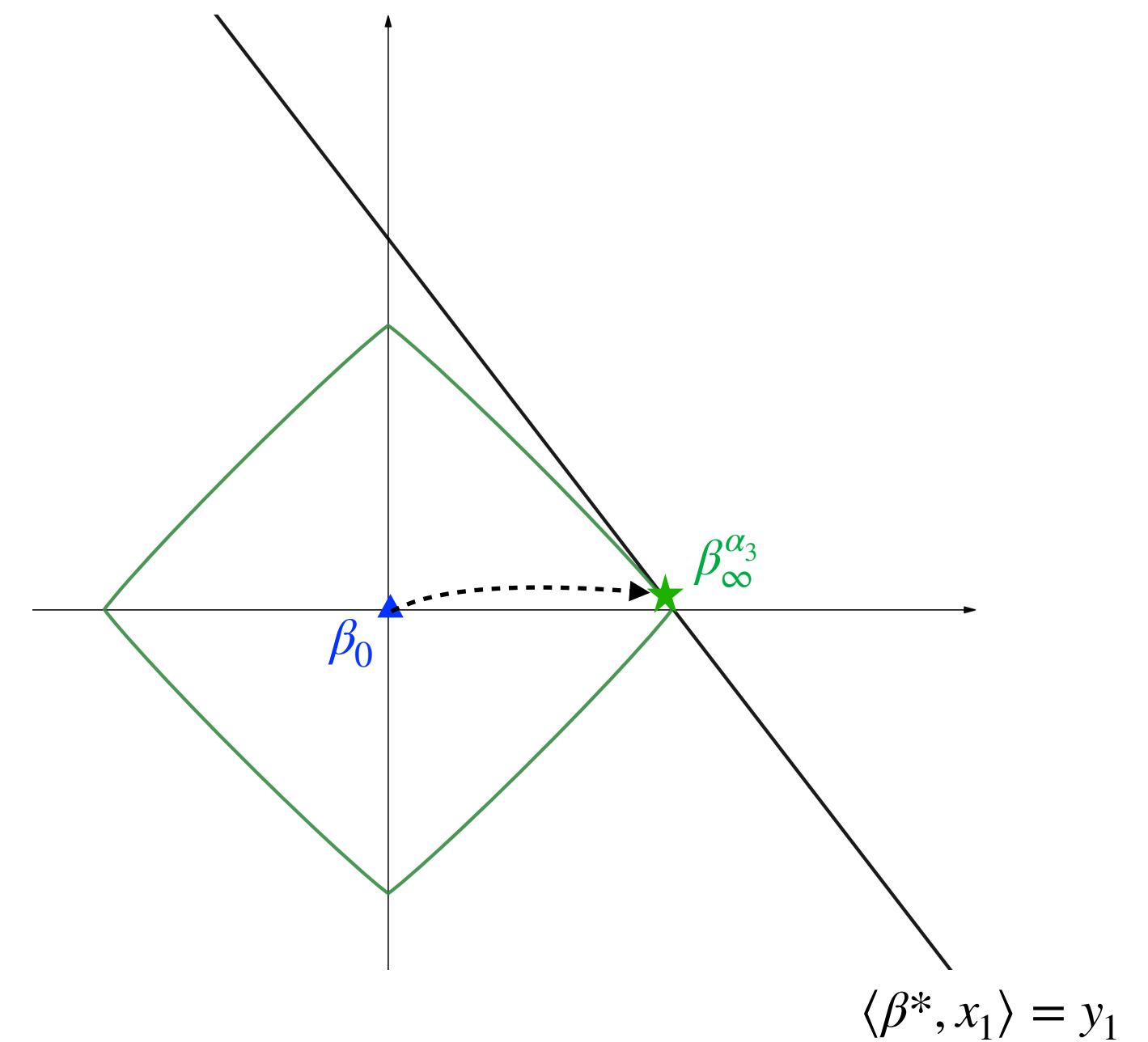
“Big” initialisation α_1



“Intermediate” initialisation α_2



“Small” initialisation α_3



Numerical illustration:

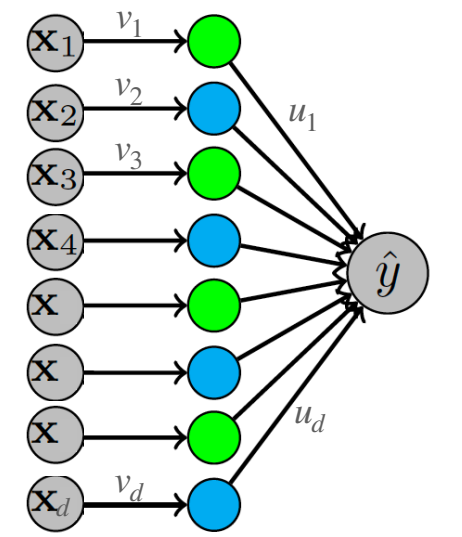
Gradient descent with fixed step-size and $u_{t=0} = \alpha \mathbf{1} \in \mathbb{R}^d$, $v_{t=0} = \mathbf{0} \in \mathbb{R}^d$:

Sparse overparametrised regression with

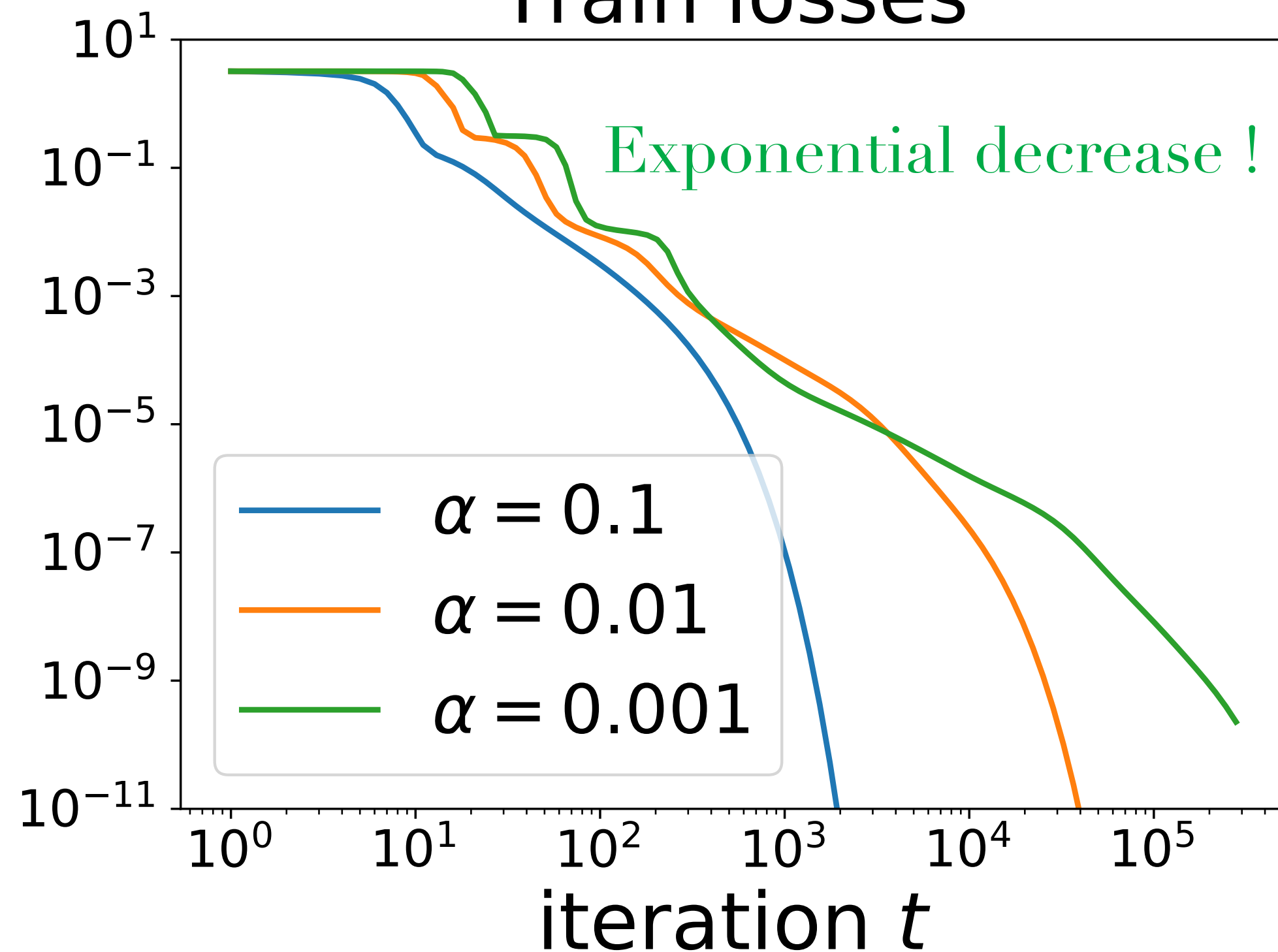
$$x_i \sim \mathcal{N}(0, I_d) \quad y_i = \langle x_i, \beta_{\ell_0}^* \rangle \quad \|\beta_{\ell_0}^*\|_0 = 5$$

$$n = 40 \quad d = 100 \quad d \gg n$$

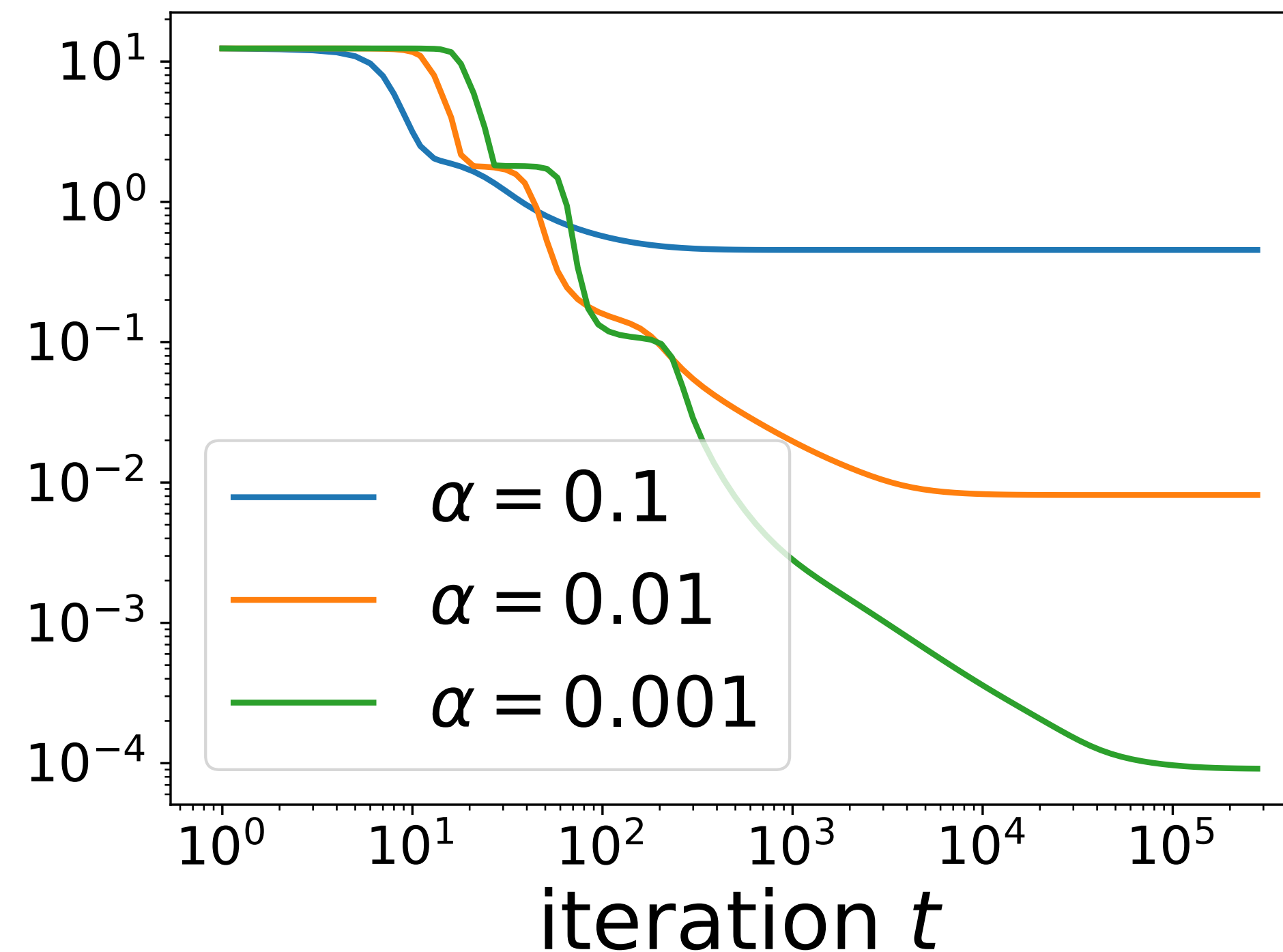
$$\min_{w \in \mathbb{R}^{2d}} L(w) = \frac{1}{4n} \sum_{i=1}^n (y_i - \langle \underbrace{u \odot v}_{\beta_w}, x_i \rangle)^2$$



Train losses



Test losses



$$\beta_{\infty}^{\alpha} = \arg \min_{\beta, \langle \beta, x_i \rangle = y_i} \phi_{\alpha}(\beta)$$

Initialisation gets smaller

$$\beta_{\infty}^{\alpha} \xrightarrow{\alpha \rightarrow 0} \arg \min_{\langle x_i, \beta \rangle = y_i} \|\beta\|_1$$

(= $\beta_{\ell_0}^*$ due to ℓ_1 magic!)

Can we intuitively understand where the ℓ_1 norm comes from ?



$$L(w) = \frac{1}{4n} \sum_{i=1}^n (y_i - \langle u \odot v, x_i \rangle)^2$$

$$w = \begin{pmatrix} u \\ v \end{pmatrix} \in \mathbb{R}^{2d} \quad \beta_w := u \odot v \in \mathbb{R}^d$$

Recall that we do **gradient flow on the neurons** $w = (u, v) \in \mathbb{R}^{2d}$.

This leads to a **mirror descent on the prediction parameter** $\beta \in \mathbb{R}^d$.

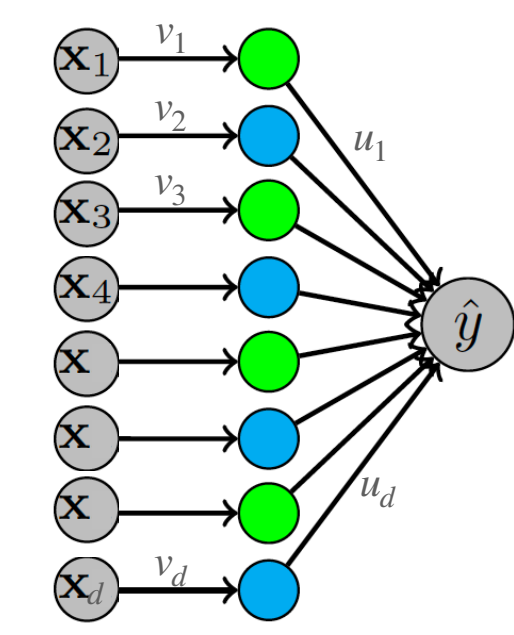
As in the linear case, we “would expect” that $w_\infty^{\alpha \rightarrow 0}$ is the solution of $\min_{w \in \mathbb{R}^{2d}, \langle u \odot v, x_i \rangle = y_i} \|w\|_2^2$

$$\begin{aligned} \text{And:} \quad \min_{w \in \mathbb{R}^{2d}, \langle u \odot v, x_i \rangle = y_i} \|w\|_2^2 &= \min_{(u, v), \langle u \odot v, x_i \rangle = y_i} \|u\|_2^2 + \|v\|_2^2 \\ &= 2 \min_{\beta \in \mathbb{R}^d, \langle \beta, x_i \rangle = y_i} \|\beta\|_1 \end{aligned}$$

$$\text{However:} \quad \min_{(u, v) \in \mathbb{R}^{2d}, \langle u \odot v, x_i \rangle = y_i} \|u - \alpha \mathbf{1}\|_2^2 + \|v\|_2^2 \neq \min_{\beta \in \mathbb{R}^d, \langle \beta, x_i \rangle = y_i} \phi_\alpha(\beta) !$$

Recap.

Training architecture



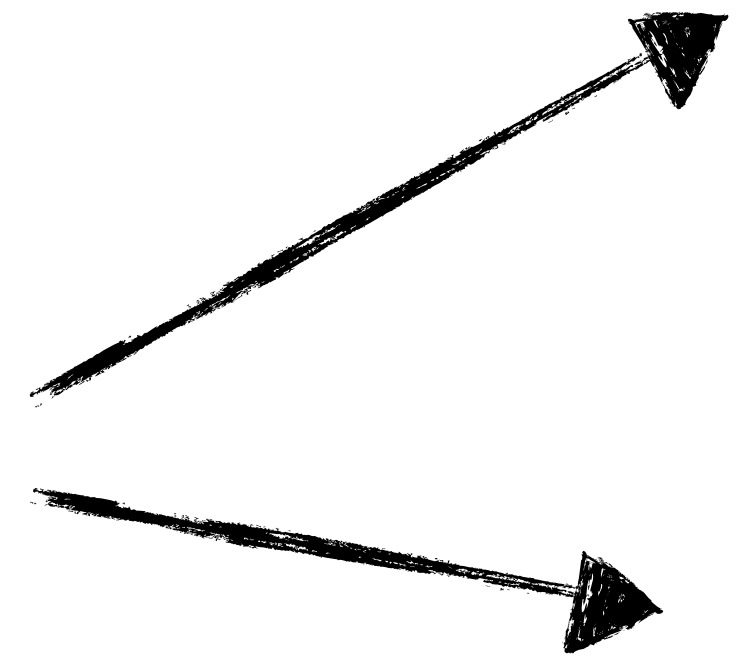
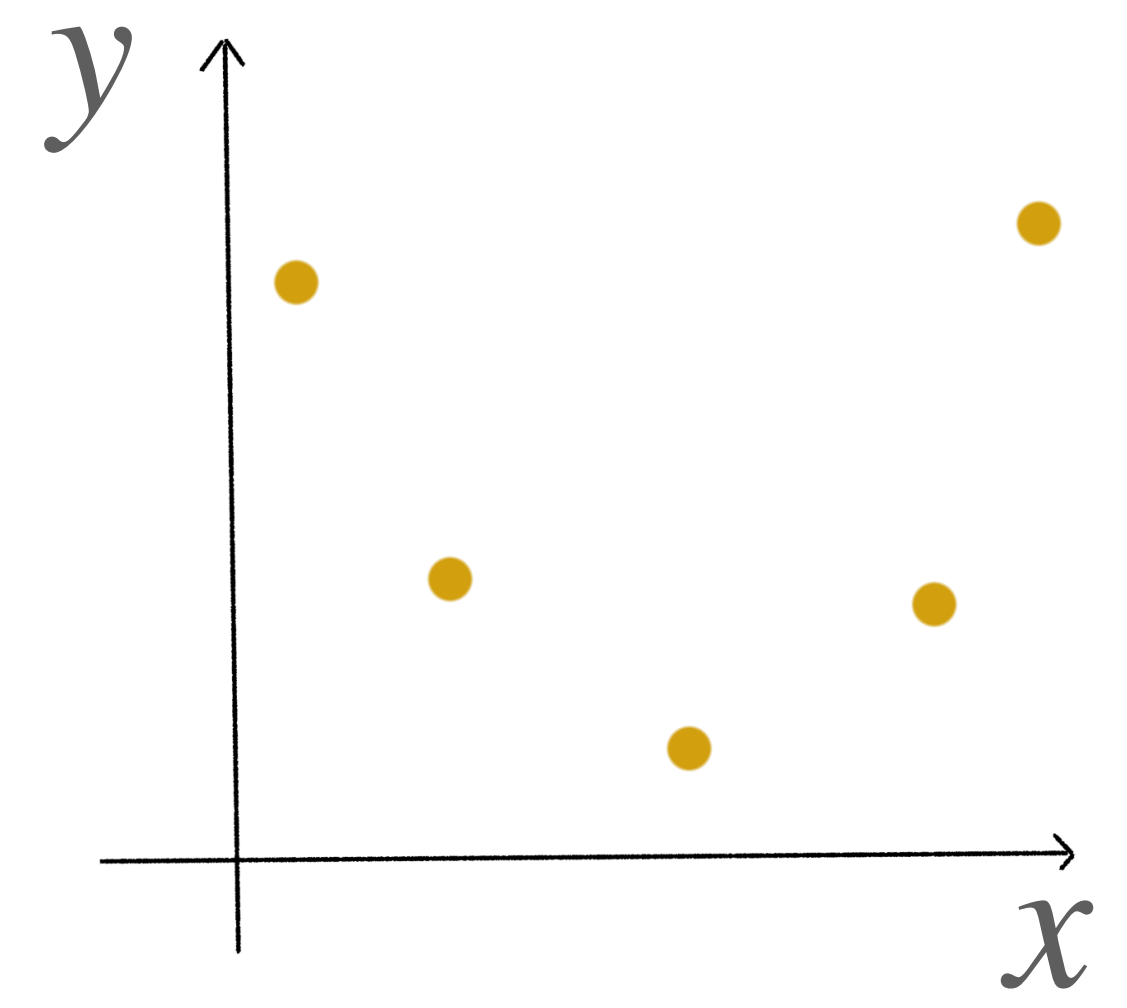
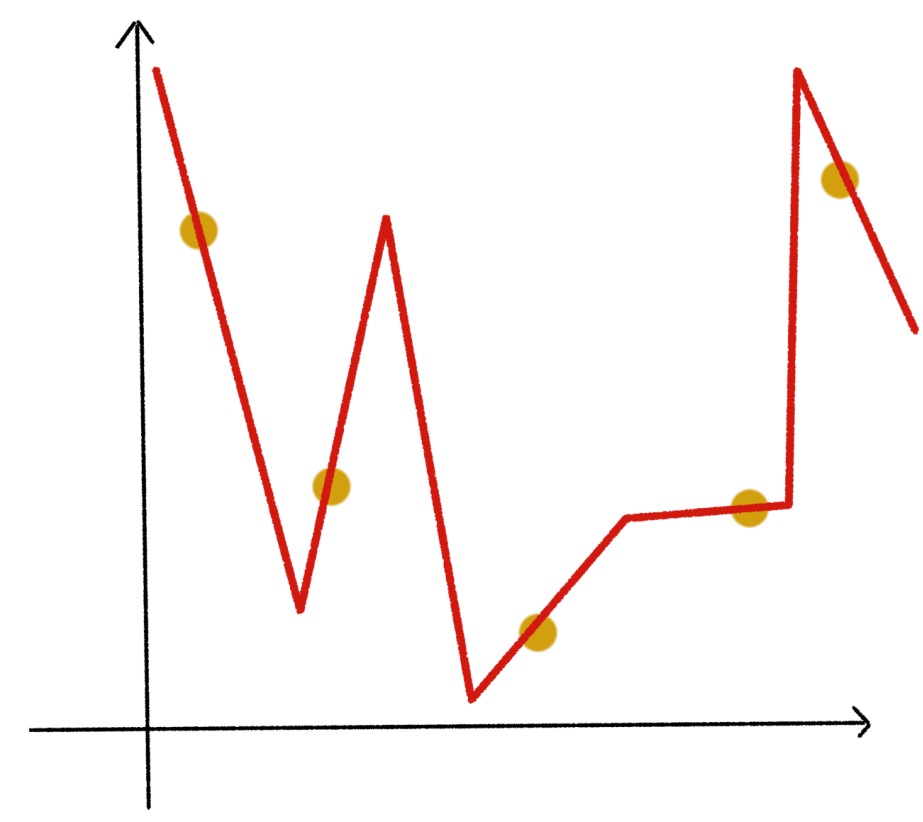
Training dataset



Training algorithm

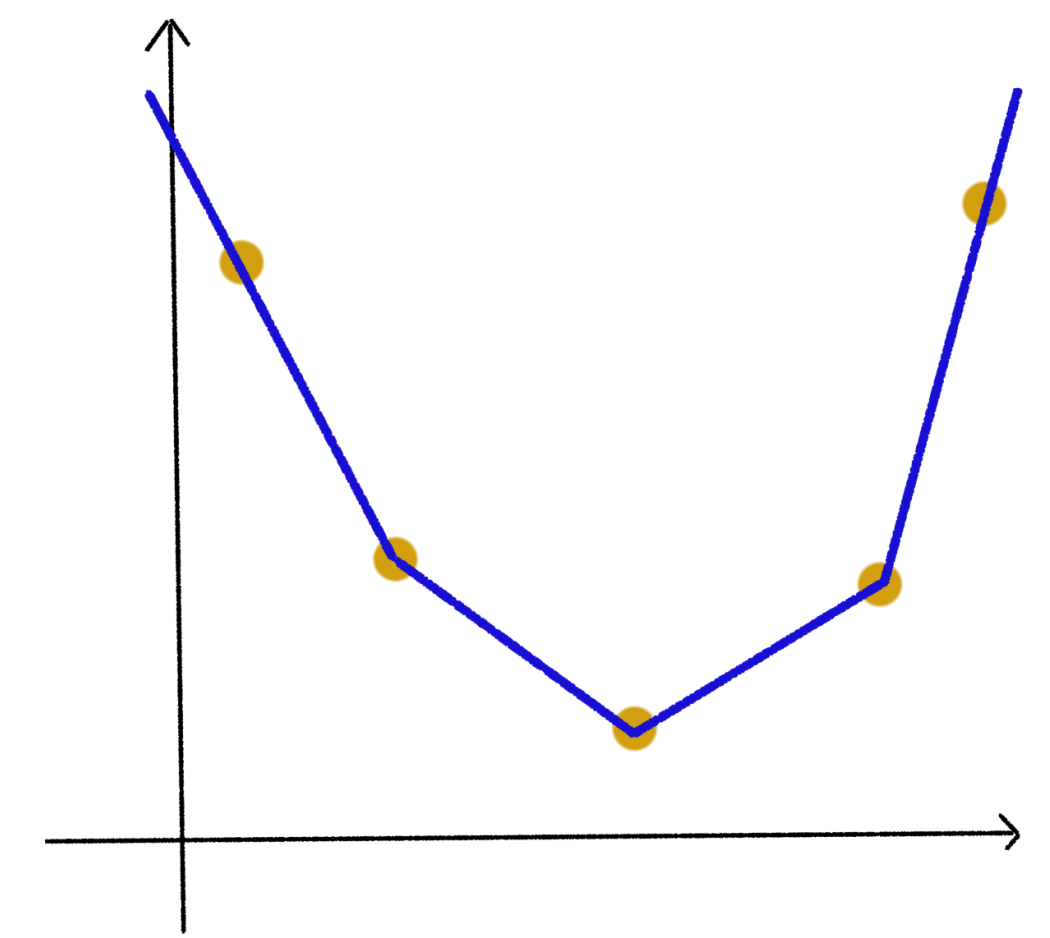
Some ERM

An infinity of interpolating solutions



GF

Initialisation scale α



$$\beta_{\infty}^{\alpha} = \arg \min_{\beta, \langle \beta, x_i \rangle = y_i} \phi_{\alpha}(\beta)$$

where

$$\phi_{\alpha} \underset{\alpha \rightarrow \infty}{\sim} \|\cdot\|_2 \quad \text{and} \quad \phi_{\alpha} \underset{\alpha \rightarrow 0}{\sim} \|\cdot\|_1$$

Second part of the talk, lets talk noise.

Main question : is there a difference of implicit bias between **SGD** and **GD** ?

- We already saw that it is **not the case** when training linear models
- What about non-linear ? Neural-networks ?

Some empirical evidence that **SGD** often outputs models which generalise better than **GD**:

Efficient BackProp 1998

Yann LeCun¹, Leon Bottou¹, Genevieve B. Orr², and Klaus-Robert Müller³

Advantages of Stochastic Learning

1. Stochastic learning is usually *much* faster than batch learning.
2. Stochastic learning also often results in better solutions.
3. Stochastic learning can be used for tracking changes.

ON LARGE-BATCH TRAINING FOR DEEP LEARNING: GENERALIZATION GAP AND SHARP MINIMA

Keskar et al. 2017

Model Name	Testing Accuracy	
	SB (Small batch)	LB (Large batch)
F_1	98.03% ± 0.07%	97.81% ± 0.07%
F_2	64.02% ± 0.2%	59.45% ± 1.05%
C_1	80.04% ± 0.12%	77.26% ± 0.42%
C_2	89.24% ± 0.12%	87.26% ± 0.07%
C_3	49.58% ± 0.39%	46.45% ± 0.43%
C_4	63.08% ± 0.5%	57.81% ± 0.17%

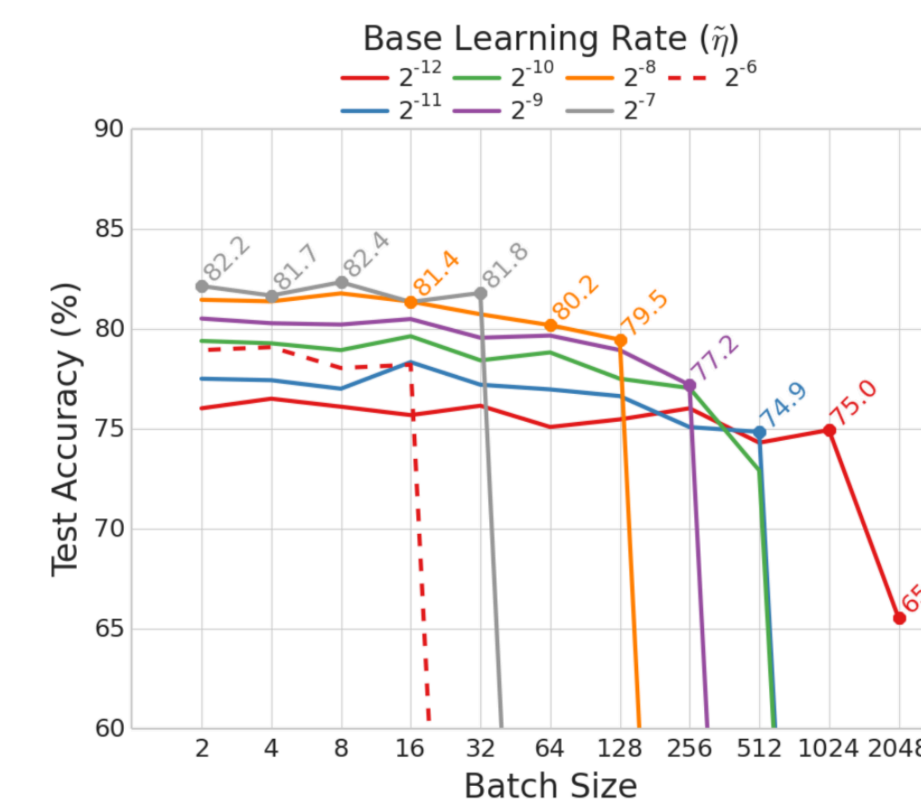
Classification task !
Train ~ 100%
For Adam !

≈ -2%

“Experiments with other optimizers for the large-batch experiments, including SGD, led to similar results.”

REVISITING SMALL BATCH TRAINING FOR DEEP NEURAL NETWORKS 2018

Dominic Masters and Carlo Luschi



Vanilla SGD, AlexNet, CIFAR 10

What about our toy neural network ?

Back to our toy neural network.

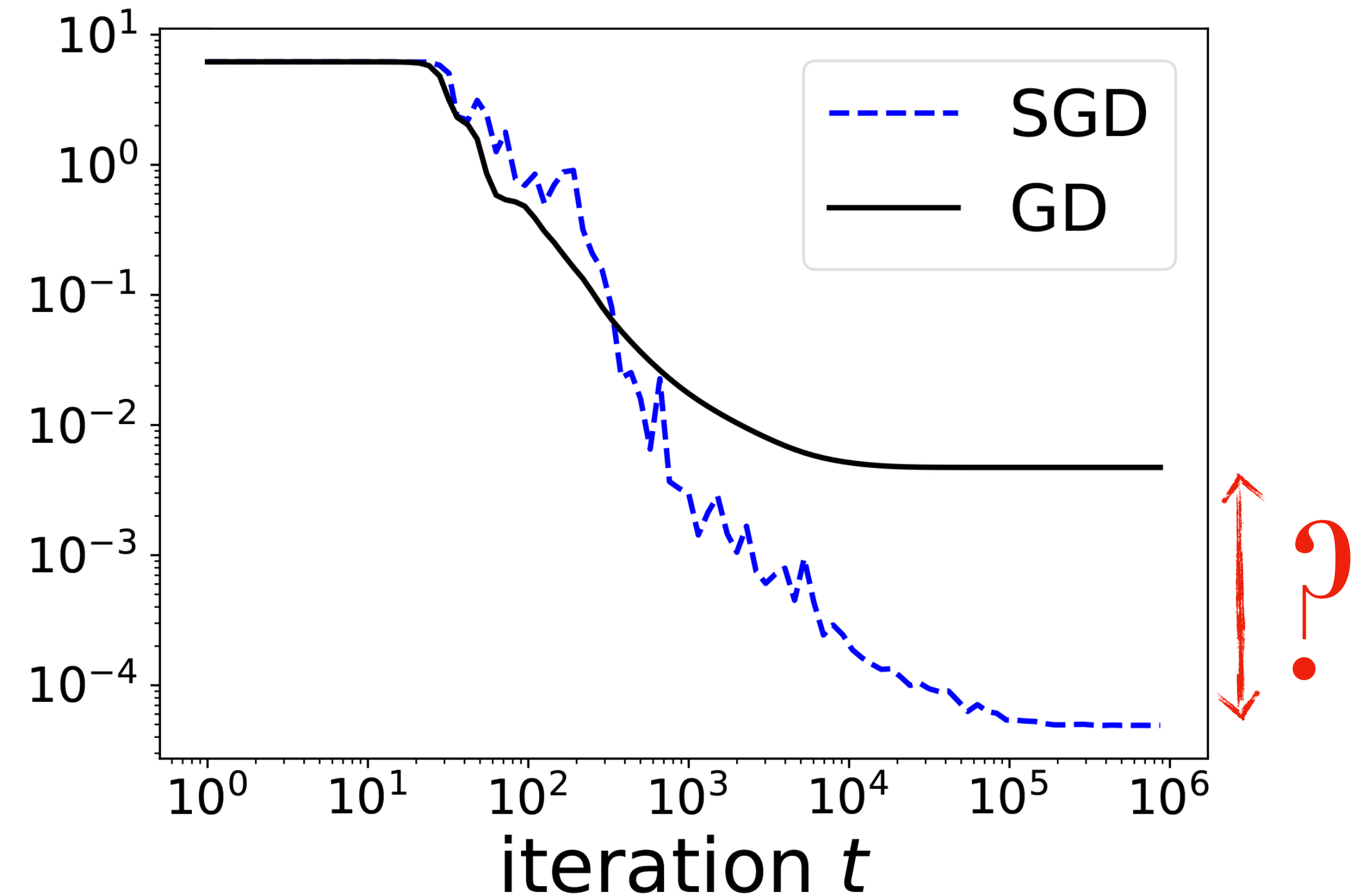
Sparse overparametrised regression with

$x_i \sim \mathcal{N}(0, I_d)$ $y_i = \langle x_i, \beta_{\ell_0}^* \rangle$ $\|\beta_{\ell_0}^*\|_0 = 5$

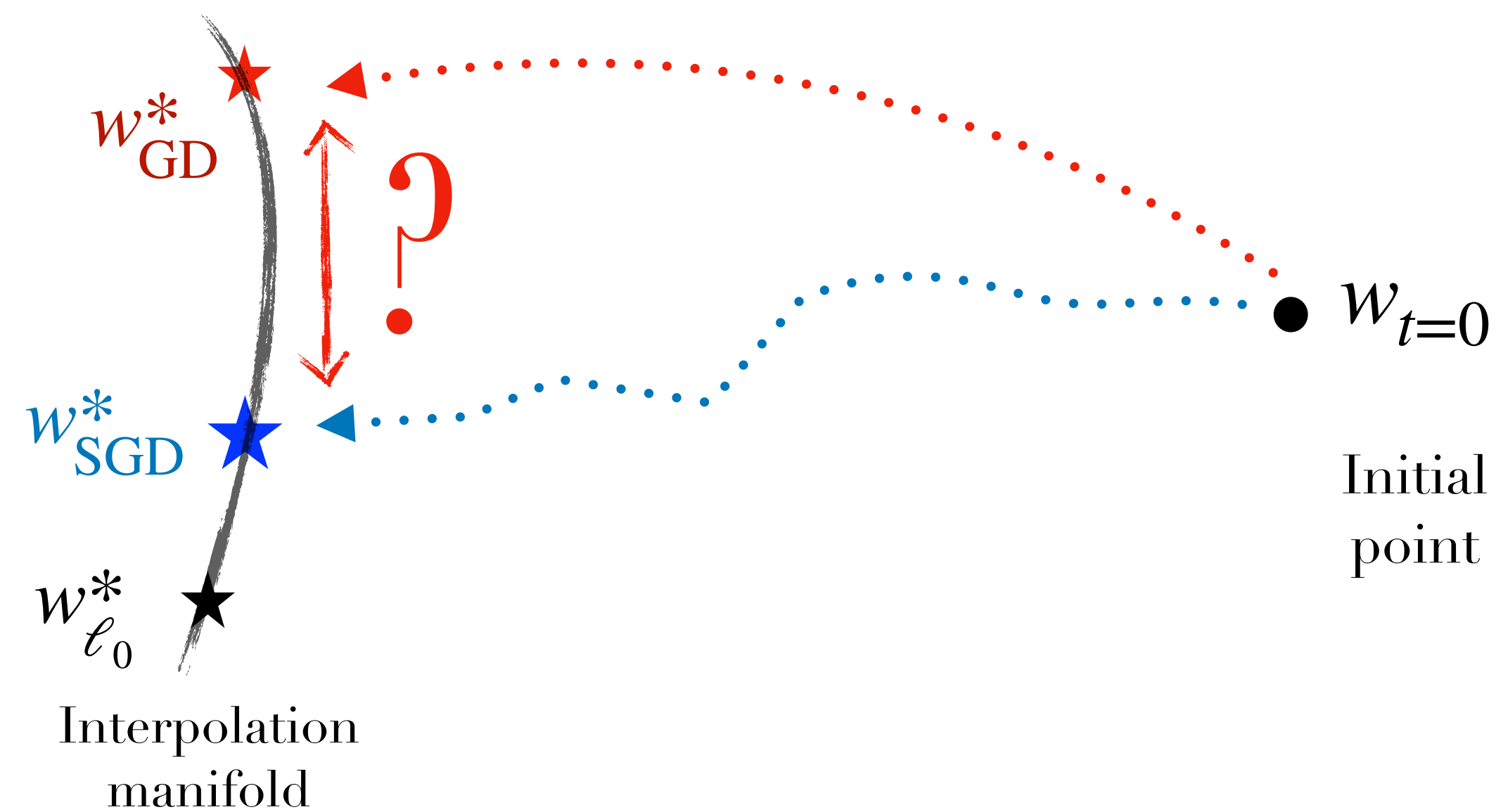
$n = 40$ $d = 100$ $d \gg n$

$$\min_{w \in \mathbb{R}^{2d}} L(w) = \frac{1}{4n} \sum_{i=1}^n (y_i - \langle \underbrace{u \odot v}_{\beta}, x_i \rangle)^2$$

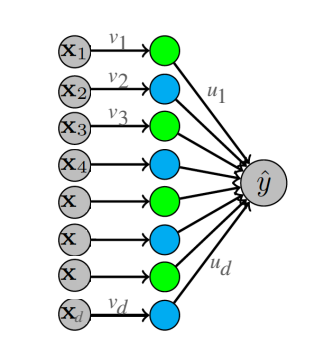
Test losses



GD and SGD with fixed step-size and $u_{t=0} = \alpha \mathbf{1} \in \mathbb{R}^d$, $v_{t=0} = \mathbf{0} \in \mathbb{R}^d$.



SGD to Stochastic Gradient Flow



$$L(w) = \frac{1}{4n} \sum_{i=1}^n (y_i - \langle u \odot v, x_i \rangle)^2$$

$$w = \begin{pmatrix} u \\ v \end{pmatrix} \in \mathbb{R}^{2d} \quad \beta_w := u \odot v \in \mathbb{R}^d$$

SGD: $w_{t+1} = w_t - \gamma \nabla L_{i_t}(w_t) \implies u_{t+1} = u_t - \gamma \langle \beta_{w_t} - \beta^*, x_{i_t} \rangle x_{i_t} \odot v_t$ (Similar for v_t)

What is the “best” continuous version of this recursion ? $du_t = -\nabla_u L(w_t) dt + \Sigma(w_t) dB_t$

Crucial part is to correctly model the noise’s structure.

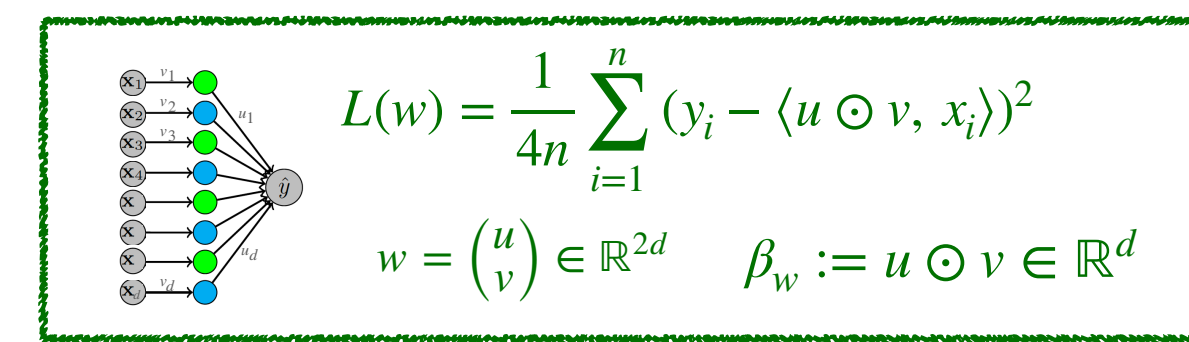
Re-writing SGD: $u_{t+1} = u_t - \gamma \nabla_u L(w_t) + \underbrace{\gamma v_t \odot [X^\top \xi_{i_t}(w_t)]}_{\text{Zero mean, state dependent, vanishing, sampling noise}}$

Zero mean, state dependent, vanishing, sampling noise

Two key properties of the noise: (i) belongs to $\text{span}(x_1 \odot v, \dots, x_n \odot v)$

(ii) has covariance $\Sigma_{SGD}(w) := \gamma^2 \text{diag}(v) X^\top \text{Cov}_{i_t}(\xi_{i_t}(\beta)) X \text{diag}(v) \in \mathbb{R}^{d \times d}$

SGD to Stochastic Gradient Flow



$$L(w) = \frac{1}{4n} \sum_{i=1}^n (y_i - \langle u \odot v, x_i \rangle)^2$$

$$w = \begin{pmatrix} u \\ v \end{pmatrix} \in \mathbb{R}^{2d} \quad \beta_w := u \odot v \in \mathbb{R}^d$$

What is the “best” continuous version of the SGD recursion ? $du_t = -\nabla_u L(w_t)dt + \Sigma(w_t)dB_t$

Crucial part is to correctly model the noise’s structure.

Re-writing SGD: $u_{t+1} = u_t - \gamma \nabla_u L(w_t) + \gamma v_t \odot [X^\top \underbrace{\xi_{i_t}(w_t)}]$

Zero mean, state dependent, vanishing, sampling noise

Two key properties of the noise: (i) belongs to $\text{span}(x_1 \odot v, \dots, x_n \odot v)$

(ii) has covariance $\Sigma_{SGD}(w) := \gamma^2 \text{diag}(v) X^\top \text{Cov}_{i_t}(\xi_{i_t}(\beta)) X \text{diag}(v) \in \mathbb{R}^{d \times d}$

We consider the following stochastic differential equation:

$$du_t = -\nabla_u L(w_t)dt + 2\sqrt{\gamma n^{-1} L(w_t)} v_t \odot [X^\top dB_t]$$

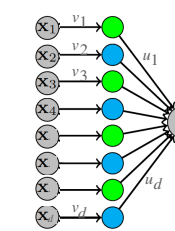
Because it conserves the two key properties: (i) structure

(ii) (nearly) matching covariance

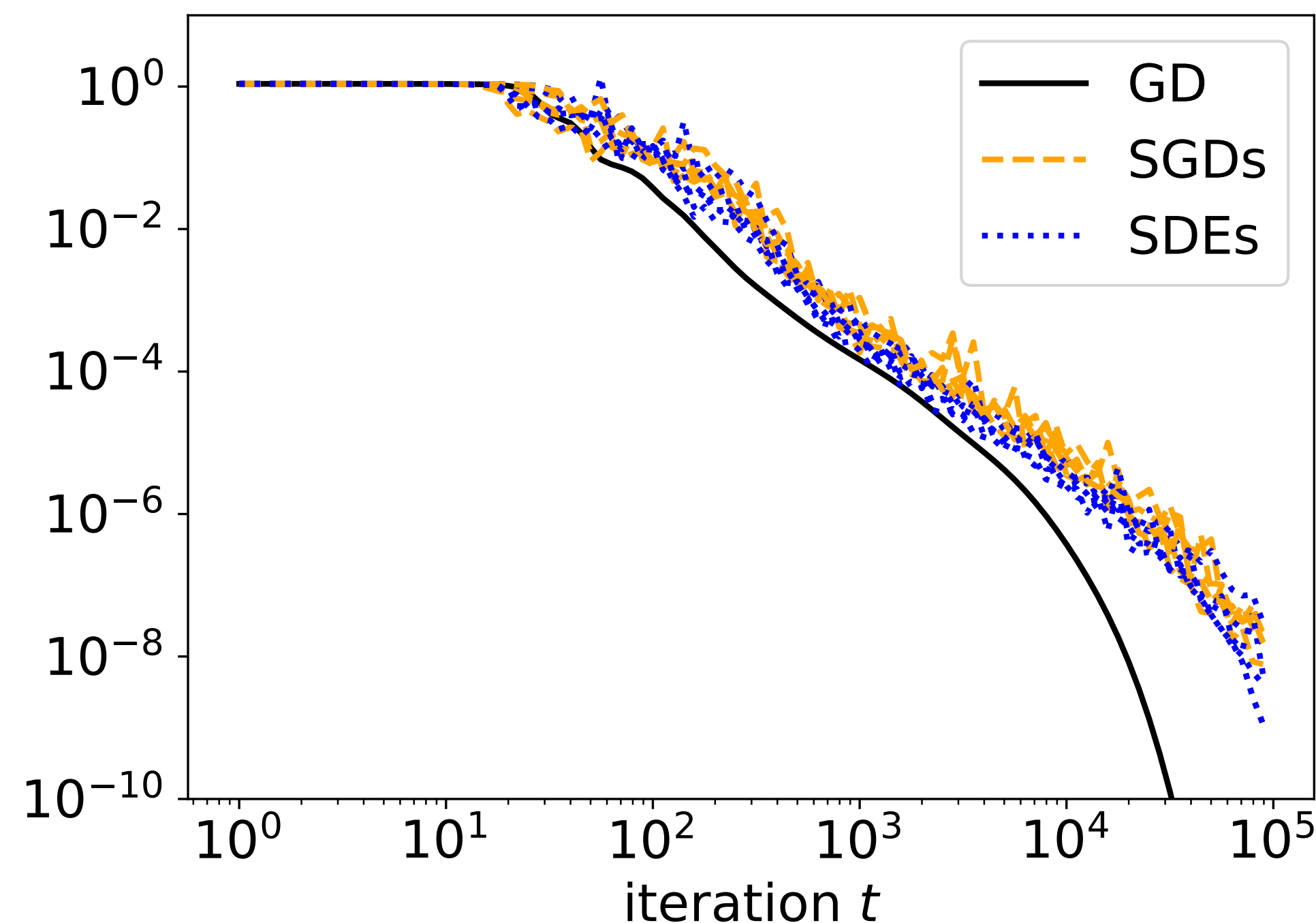
Numerical “validation”

$$x_i \sim \mathcal{N}(0, I_d) \quad y_i = \langle x_i, \beta_{\ell_0}^* \rangle \quad \|\beta_{\ell_0}^*\|_0 = 5$$

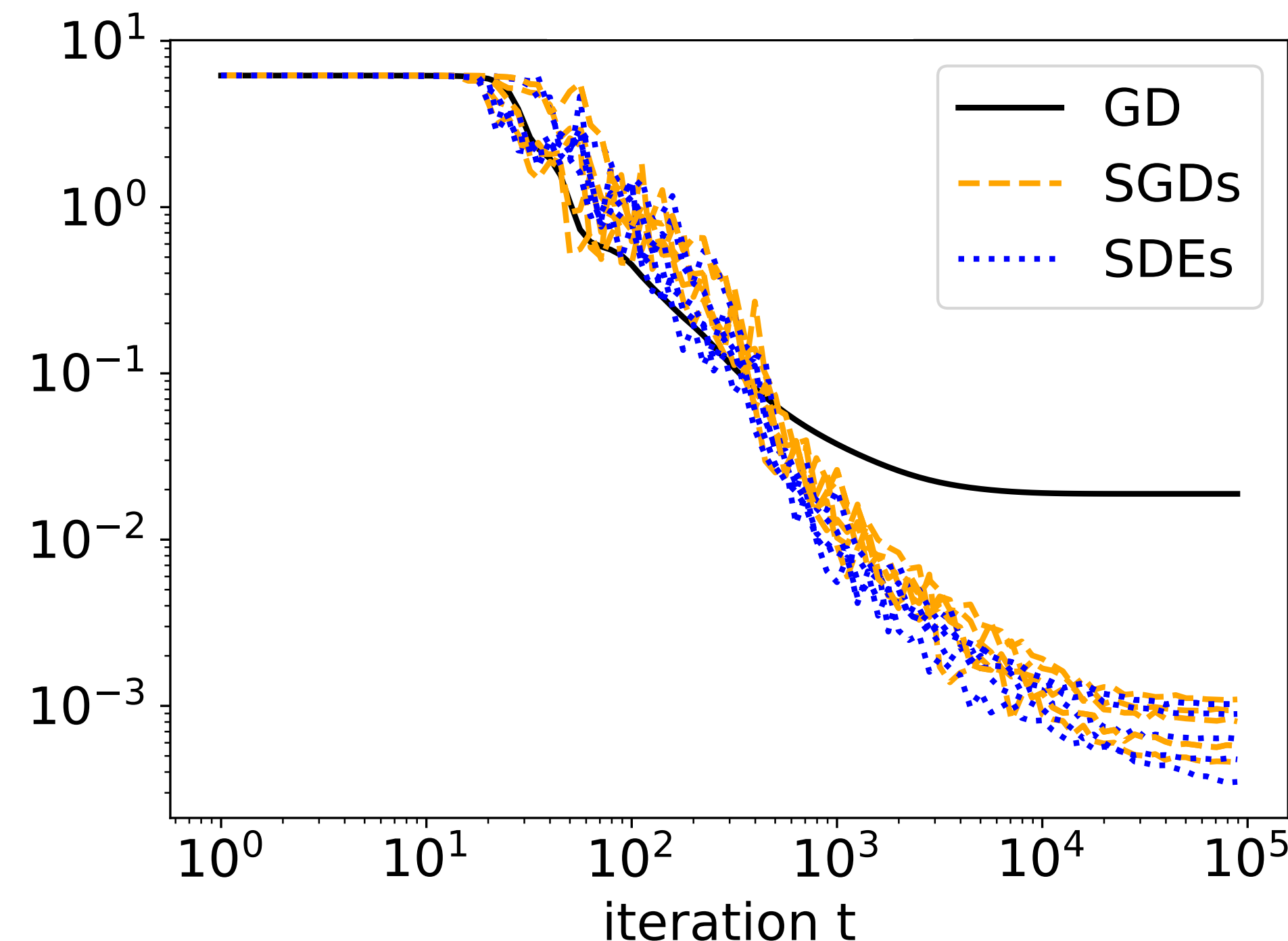
$$\min_{w \in \mathbb{R}^{2d}} L(w) = \frac{1}{4n} \sum_{i=1}^n (y_i - \langle u \odot v, x_i \rangle)^2$$



Train losses



Test losses

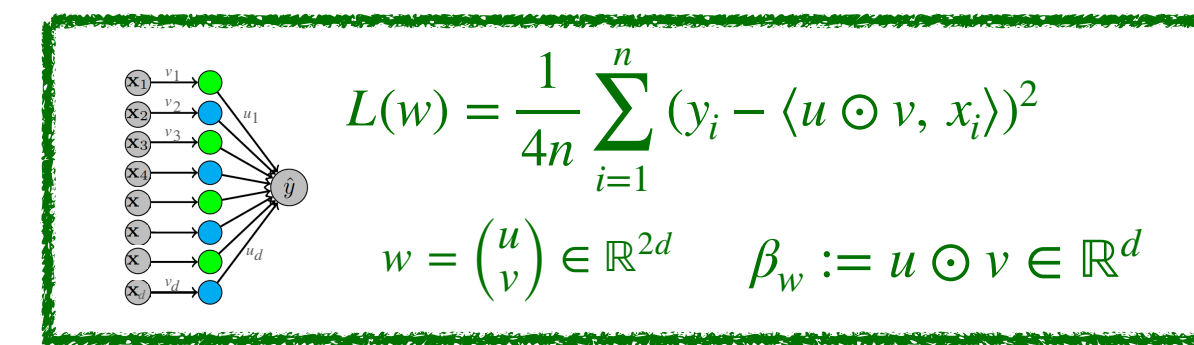


The SDE seems to faithfully capture SGD’s behaviour for macroscopic step-sizes !

Keep in mind: • This is a model !

- There are unfortunately no theoretical guarantees for macroscopic step-sizes (as for GF!)
- However it captures the key ingredients to understand the implicit bias of SGD.

Implicit bias of the stochastic gradient flow



$$L(w) = \frac{1}{4n} \sum_{i=1}^n (y_i - \langle u \odot v, x_i \rangle)^2$$

$$w = \begin{pmatrix} u \\ v \end{pmatrix} \in \mathbb{R}^{2d} \quad \beta_w := u \odot v \in \mathbb{R}^d$$

Assumptions: probability $p \in (0,1)$ and initialisation $u_{t=0} = \alpha \in \mathbb{R}^d, v_{t=0} = 0$.

Step-size $\gamma \leq \tilde{O}\left(\frac{1}{\ln(4/p)\lambda_{\max}\|\beta_{\ell_1}^*\|_1}\right)$ where

$$\left\{ \begin{array}{l} \lambda_{\max} = \lambda_{\max}(X^\top X/n) \\ \beta_{\ell_1}^* = \operatorname{argmin}_{\beta \text{ s.t. } X\beta=y} \|\beta\|_1 \end{array} \right.$$

Result: With probability $1 - p$, the **Stochastic Gradient Flow** (u_t, v_t) is such that:

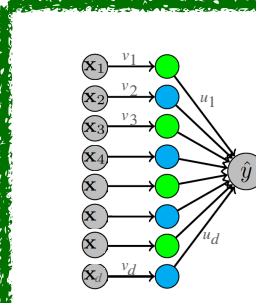
Convergence \rightarrow • The flow $(\beta_t)_{t \geq 0} = (u_t \odot v_t)_{t \geq 0}$ converges towards a zero-training error solution $\beta_\infty^{\alpha, \text{SGF}}$

Implicit Bias \rightarrow • This solution $\beta_\infty^{\alpha, \text{SGF}}$ satisfies

$$\beta_\infty^{\alpha, \text{SGF}} = \operatorname{arg min}_{\beta \in \mathbb{R}^d, \langle \beta, x_i \rangle = y_i} \phi_{\alpha_\infty}(\beta) \text{ where } \underbrace{\alpha_\infty = \alpha \odot \exp\left(-2\gamma \operatorname{diag}\left(\frac{X^\top X}{n}\right)\right)}_{\text{“effective” initialisation}} \underbrace{\int_0^{+\infty} L(\beta_s) ds}_{\text{training loss}} \underbrace{< \alpha}_{\text{initialisation scale}}$$

stochastic !

What does this mean ?



$$L(w) = \frac{1}{4n} \sum_{i=1}^n (y_i - \langle u \odot v, x_i \rangle)^2$$

$$w = \begin{pmatrix} u \\ v \end{pmatrix} \in \mathbb{R}^{2d} \quad \beta_w := u \odot v \in \mathbb{R}^d$$

$$\beta_{\infty}^{\alpha, \text{SGF}} = \arg \min_{\beta \in \mathbb{R}^d, \langle \beta, x_i \rangle = y_i} \phi_{\alpha_{\infty}}(\beta) \quad \text{where} \quad \underbrace{\alpha_{\infty}}_{\text{“effective”}} = \alpha \odot \exp\left(-2\gamma \text{diag}\left(\frac{X^{\top} X}{n}\right)\right) \underbrace{\int_0^{+\infty} L(\beta_s) ds}_{\text{training loss}} < \underbrace{\alpha}_{\text{initialisation scale}}$$

GF vs SGF: Recall that: $\beta_{\infty}^{\alpha, \text{GF}} = \arg \min_{\beta \in \mathbb{R}^d, \langle \beta, x_i \rangle = y_i} \phi_{\alpha}(\beta)$

- Implicit bias of SGF is the same as GF but with an **effective initialisation**
- Since $\alpha_{\infty} < \alpha$, $\phi_{\alpha_{\infty}}$ is closer to the ℓ_1 norm than ϕ_{α} and $\beta_{\infty}^{\alpha, \text{SGF}}$ is “sparser” than $\beta_{\infty}^{\alpha, \text{GF}}$

The slower the convergence, the “better” the bias:

$$\int_0^{+\infty} L(\beta_s) ds \gg 1 \implies \alpha_{\infty} \ll \alpha$$

The bigger the step-size, the “better” the bias

Convergence for fixed step-size !

What does this mean ?

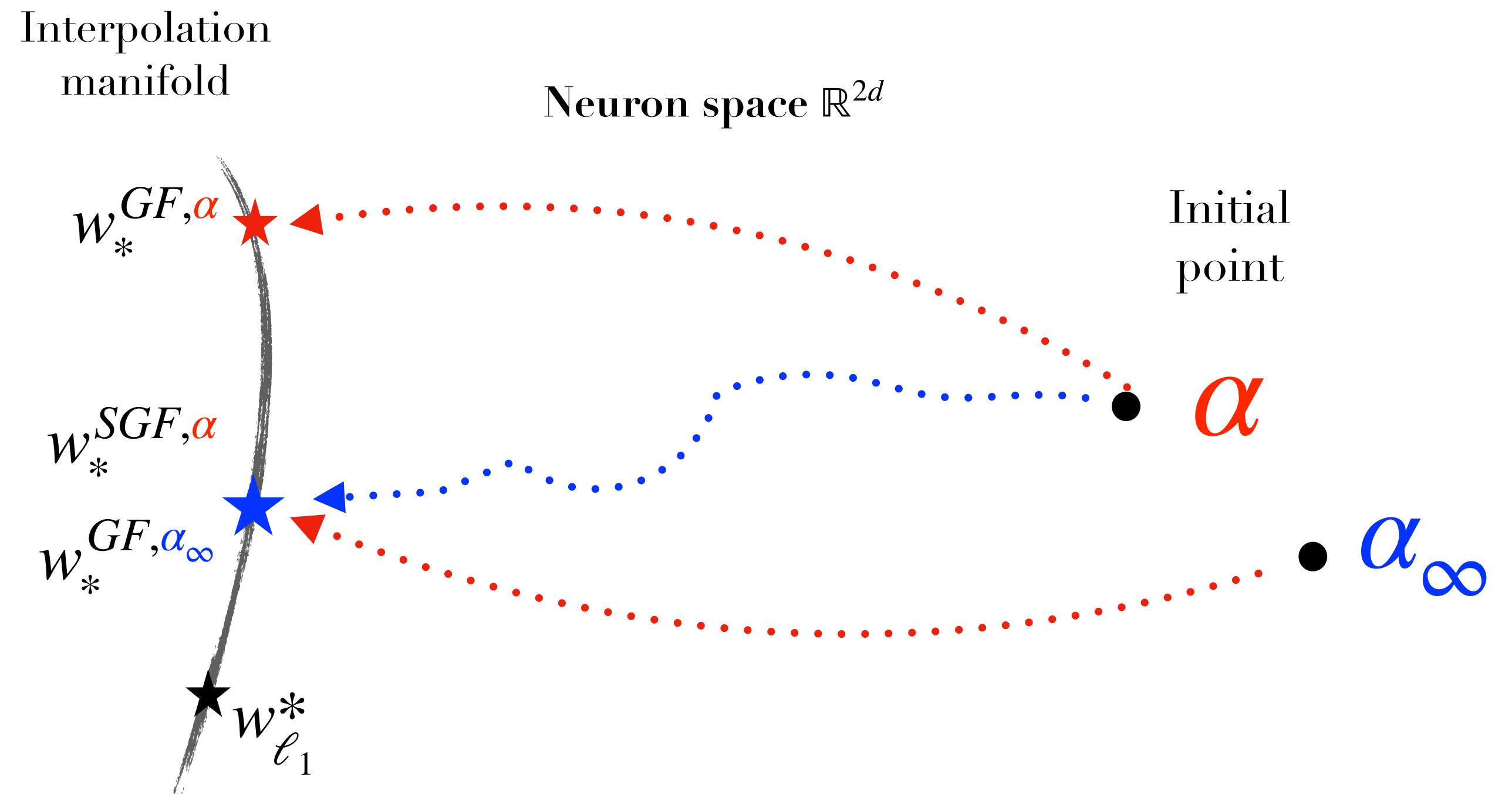
$$\beta_{\infty}^{\alpha, \text{SGF}} = \arg \min_{\beta \in \mathbb{R}^d, \langle \beta, x_i \rangle = y_i} \phi_{\alpha_{\infty}}(\beta)$$

$$\beta_{\infty}^{\alpha, \text{GF}} = \arg \min_{\beta \in \mathbb{R}^d, \langle \beta, x_i \rangle = y_i} \phi_{\alpha}(\beta)$$

$$\underbrace{\alpha_{\infty}}_{\text{"effective" initialisation}} = \alpha \odot \exp \left(-2\gamma \text{diag} \left(\frac{X^{\top} X}{n} \right) \underbrace{\int_0^{+\infty} L(\beta_s) ds}_{\text{training loss}} \right) < \alpha \quad \text{Neuron initialisation scale}$$

stochastic !

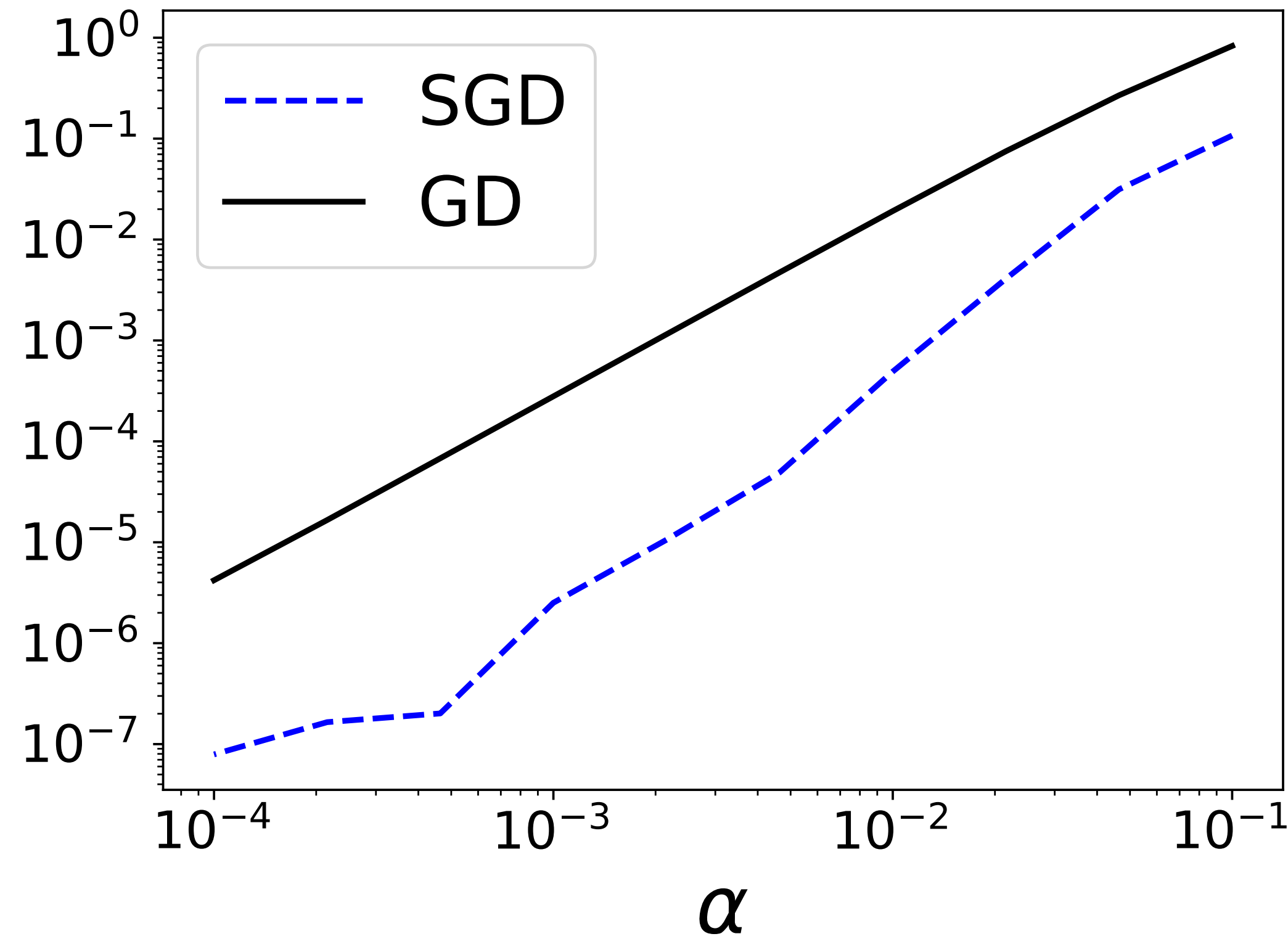
$$\beta_{\infty}^{\alpha, \text{SGF}} = \beta_{\infty}^{\alpha_{\infty}, \text{GF}}$$



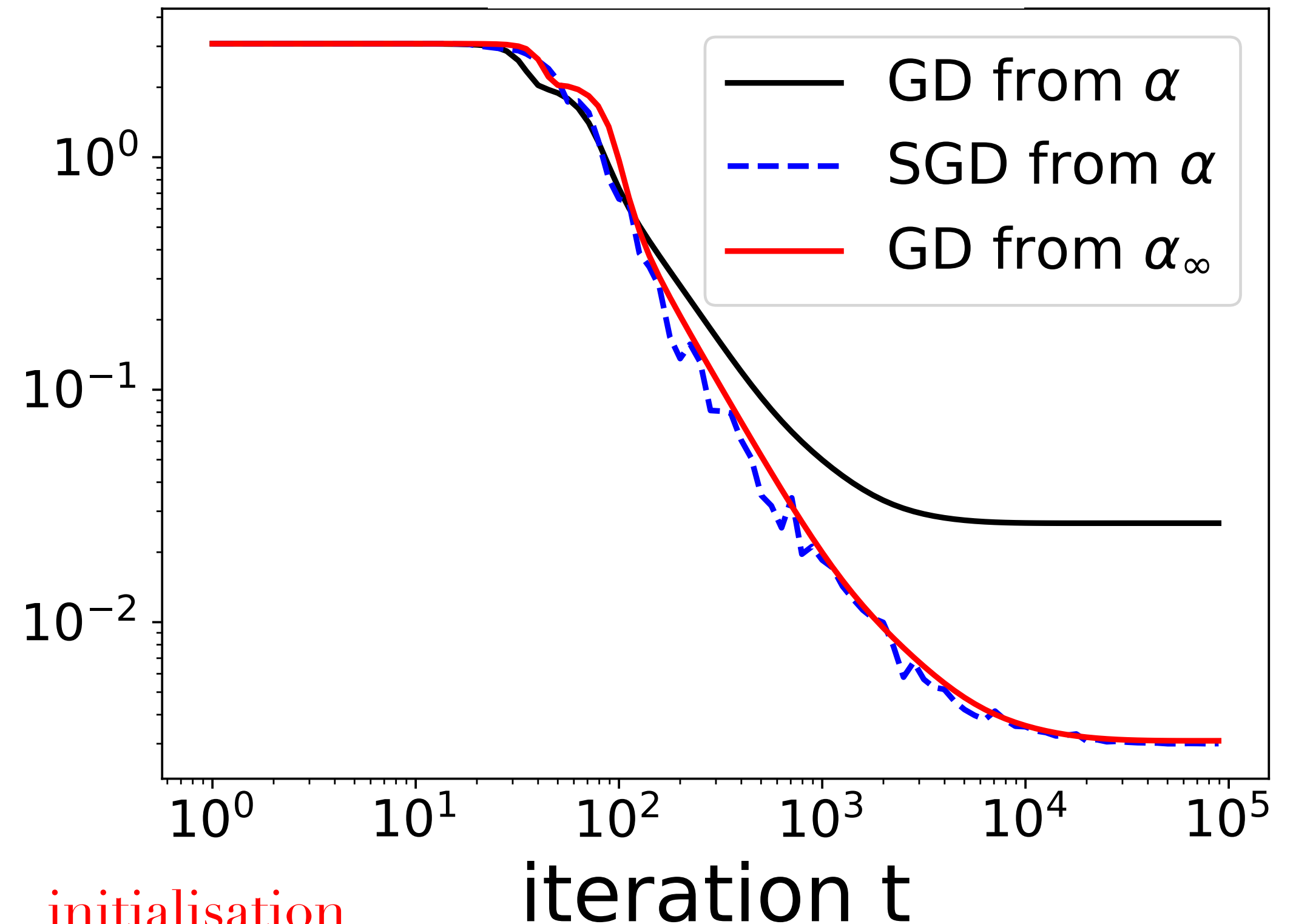
Benefit of stochasticity

Setting: $n = 40$ $d = 100$ $\|\beta_{\ell_0}^*\|_0 = 5$
 $x_i \sim \mathcal{N}(0, I)$ $y_i = \langle x_i, \beta_{\ell_0}^* \rangle$ $\alpha = 0.1$

Test losses at convergence



Test losses



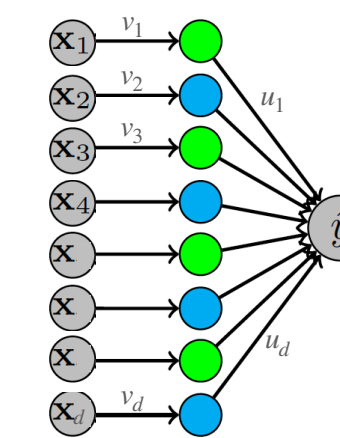
$$\underbrace{\alpha_\infty}_{\text{“effective” initialisation}} = \underbrace{\alpha}_{\text{initialisation scale}} \odot \exp\left(-2\gamma \text{diag}\left(\frac{X^\top X}{n}\right) \underbrace{\int_0^{+\infty} L(\beta_s) ds}_{\text{training loss}}\right) < \underbrace{\alpha}_{\text{initialisation scale}}$$

stochastic !

Benefit of stochasticity

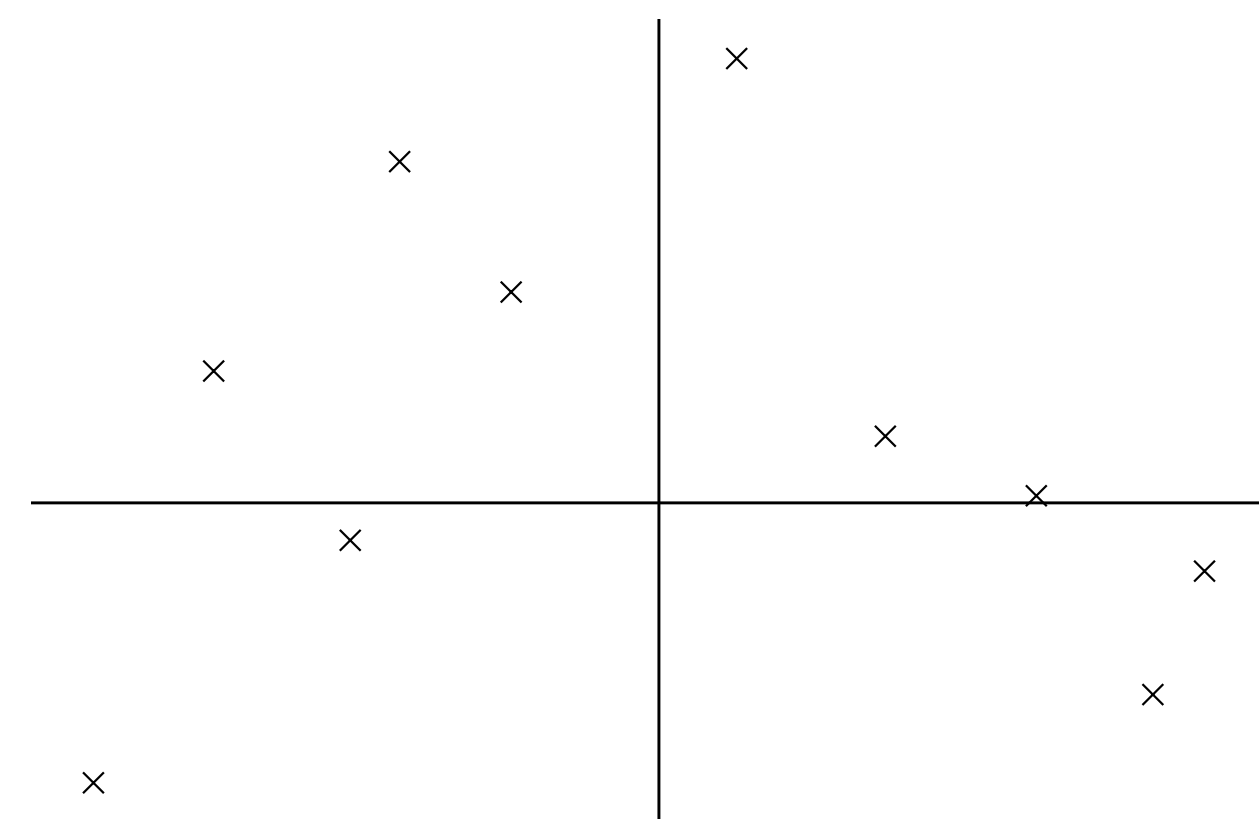
Samples $(x_i, y_i)_{1 \leq i \leq n} \in \mathbb{R} \times \mathbb{R}$ from some distribution \mathcal{D} .

We want to linearly interpolate with feature expansion $\phi(x) = (\cos(2\pi i x), \sin(2\pi i x))_{1 \leq i \leq d/2} \in \mathbb{R}^d$.

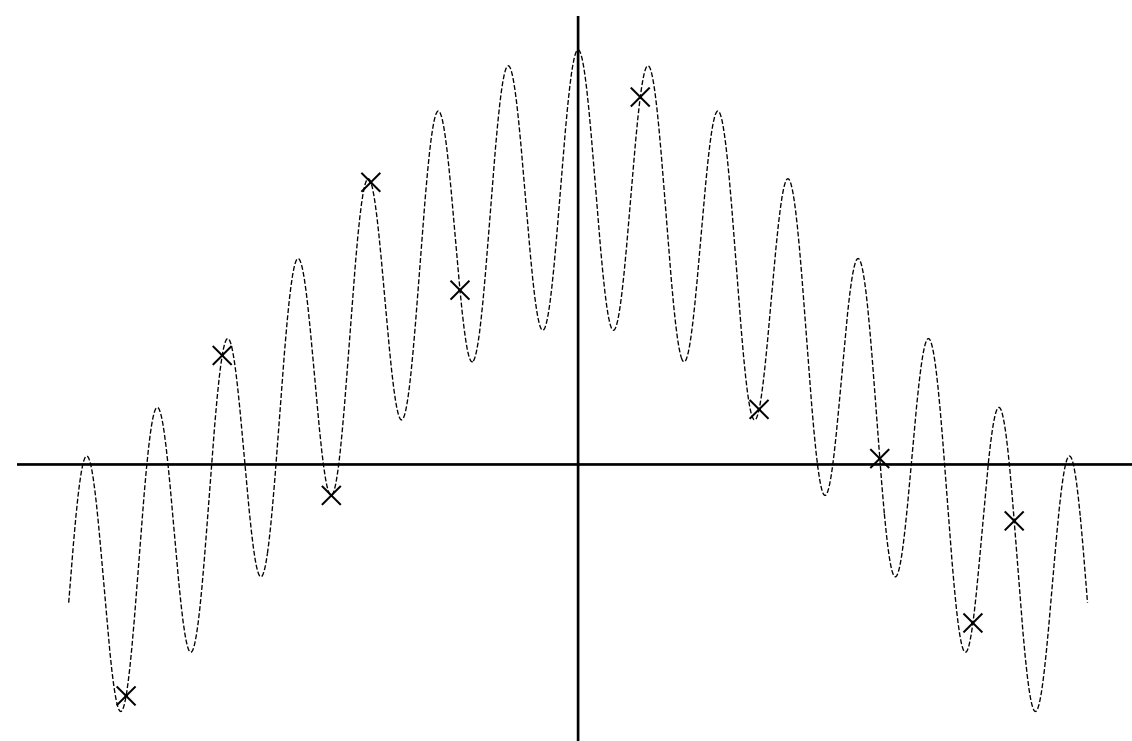


$$f_w(x) = \langle u \odot v, \phi(x) \rangle$$

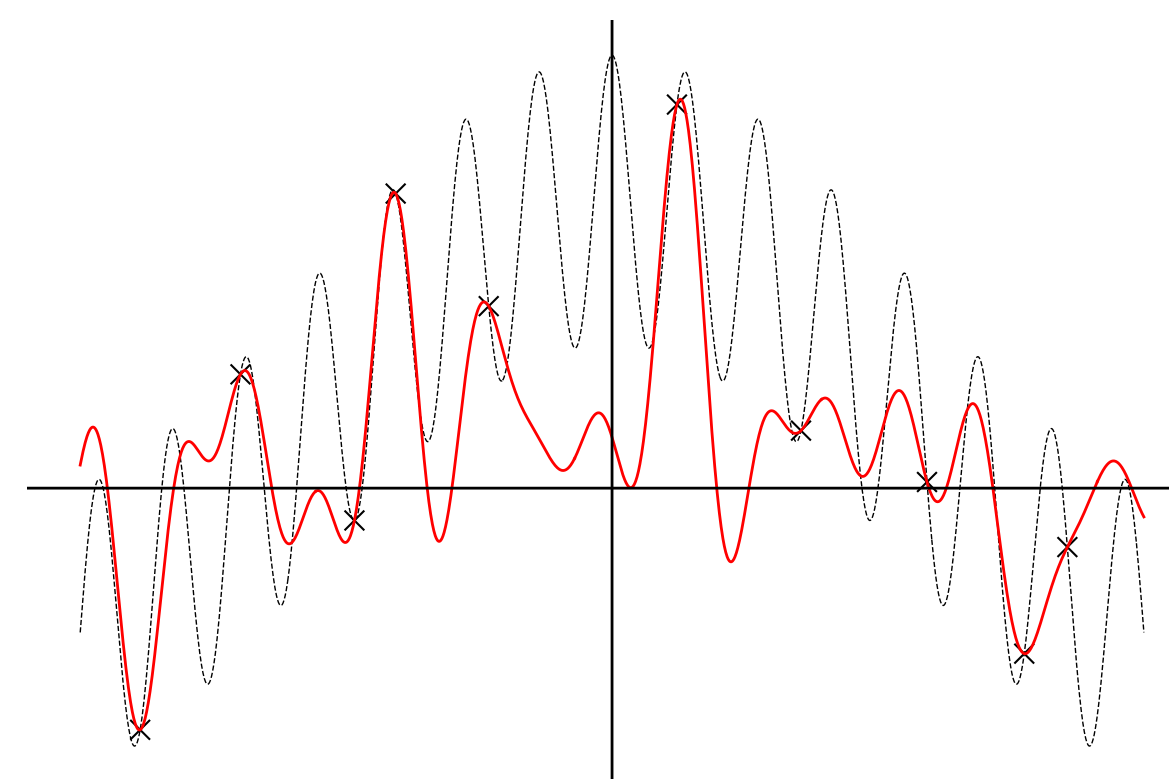
Training set 1



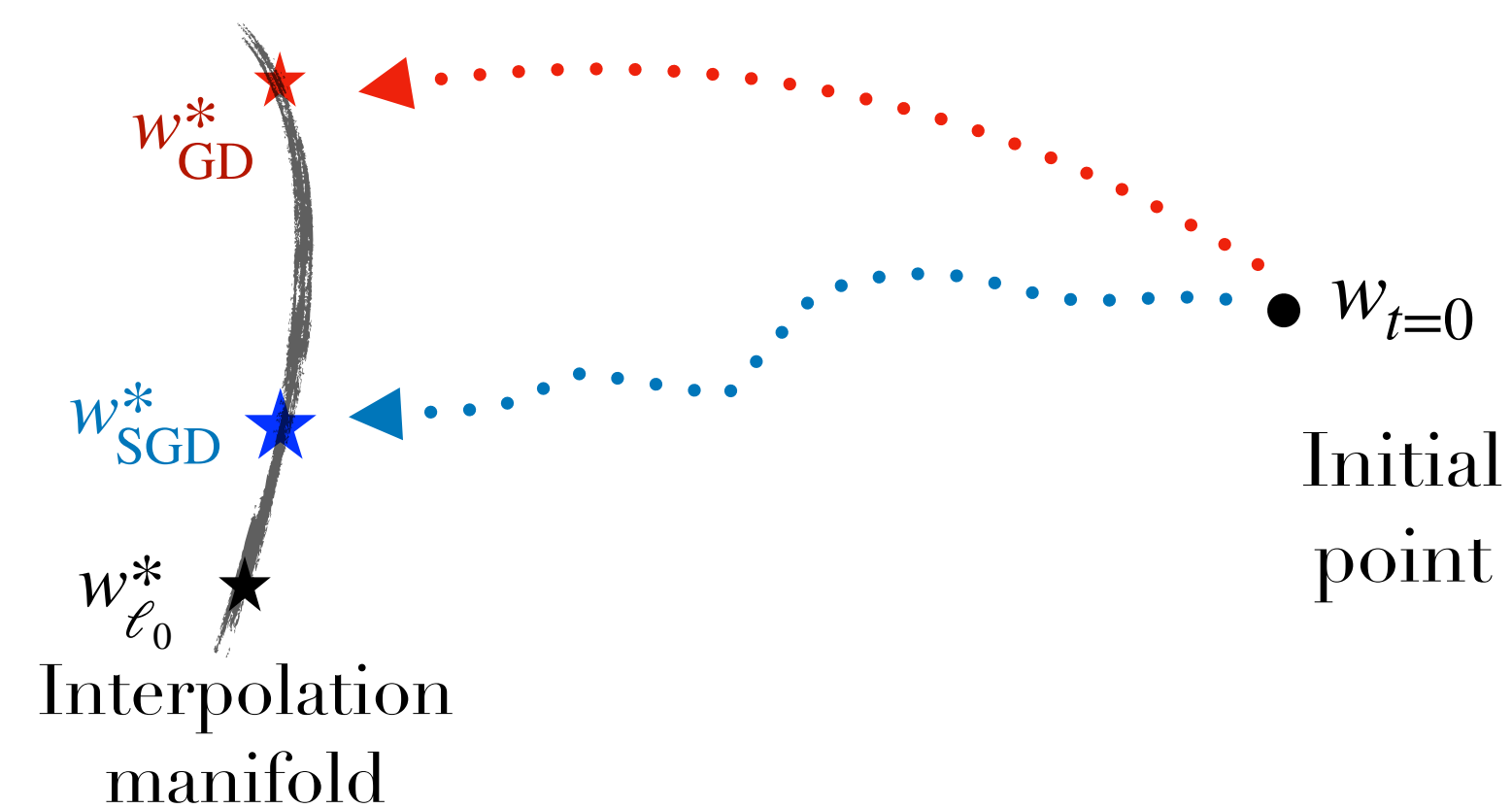
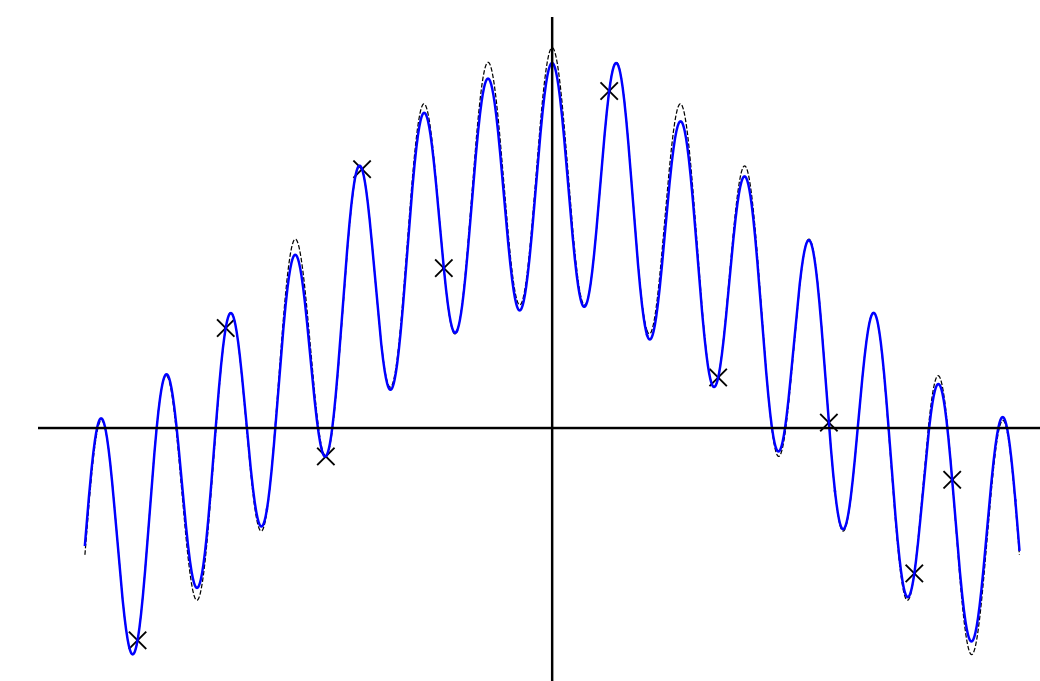
From sparse distribution



GD from initialisation α



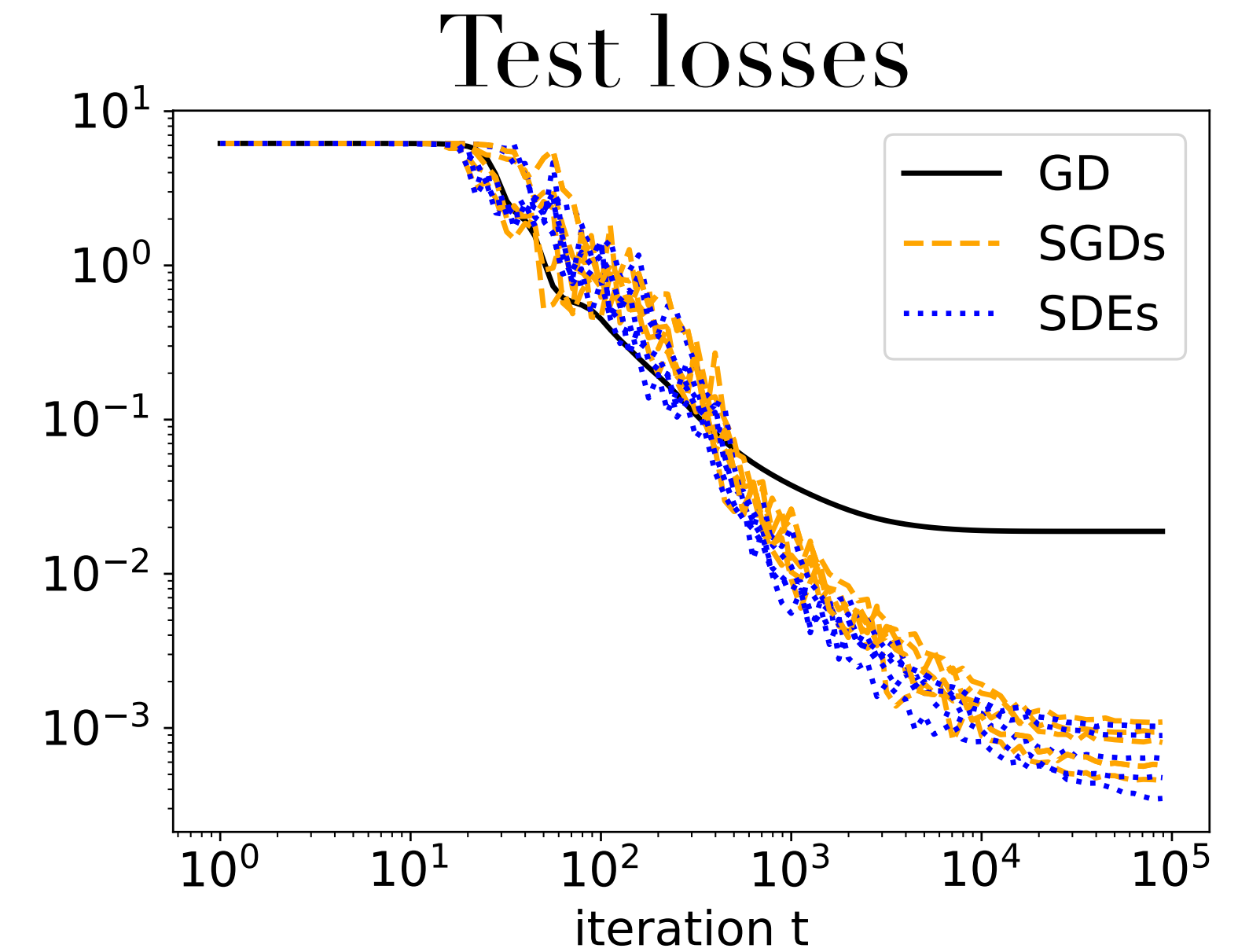
SGD from initialisation α
(+label noise)



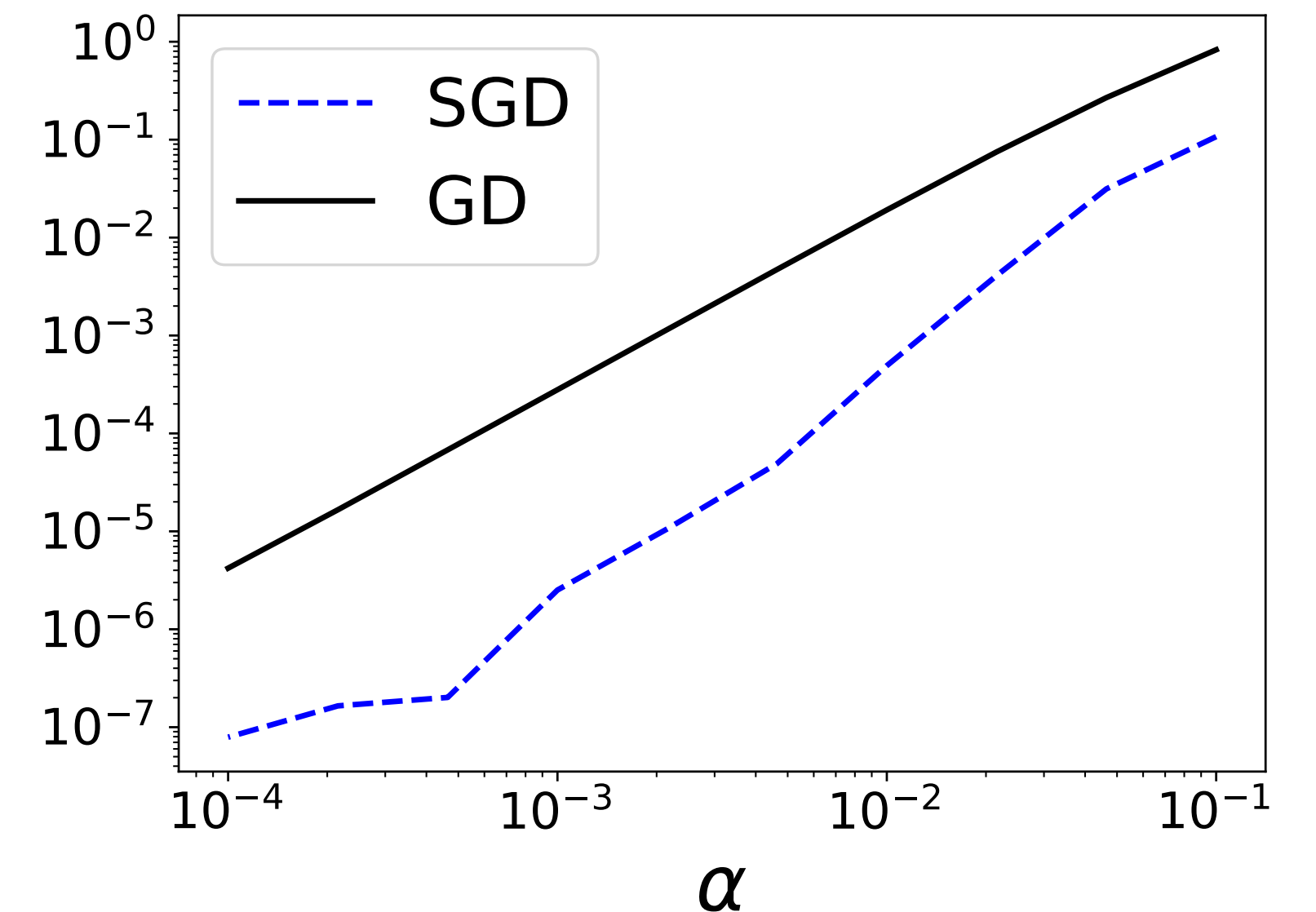
Take home messages

1. Considering (appropriate) stochastic gradient flows can lead to interesting and pertinent results

2. For a very specific toy problem, the noise inherent to SGD's stochasticity helps recover a solution which has better sparsity properties than that of GD.



Test losses at convergence



Implicit regularisation of gradient algorithms

Scott Pesme



Loucas Pillaud-Vivien



Nicolas Flammarion



TML lab

EPFL

Bonus 0: other SDE modelisations

SGD: $u_{t+1} = u_t - \gamma \nabla_u L(w_t) + \gamma v_t \odot [X^\top \xi_{i_t}(w_t)]$

Our SDE:

$$du_t = -\nabla_u L(w_t)dt + 2\sqrt{\gamma n^{-1} L(w_t)} v_t \odot [X^\top dB_t]$$

Overdamped Langevin:

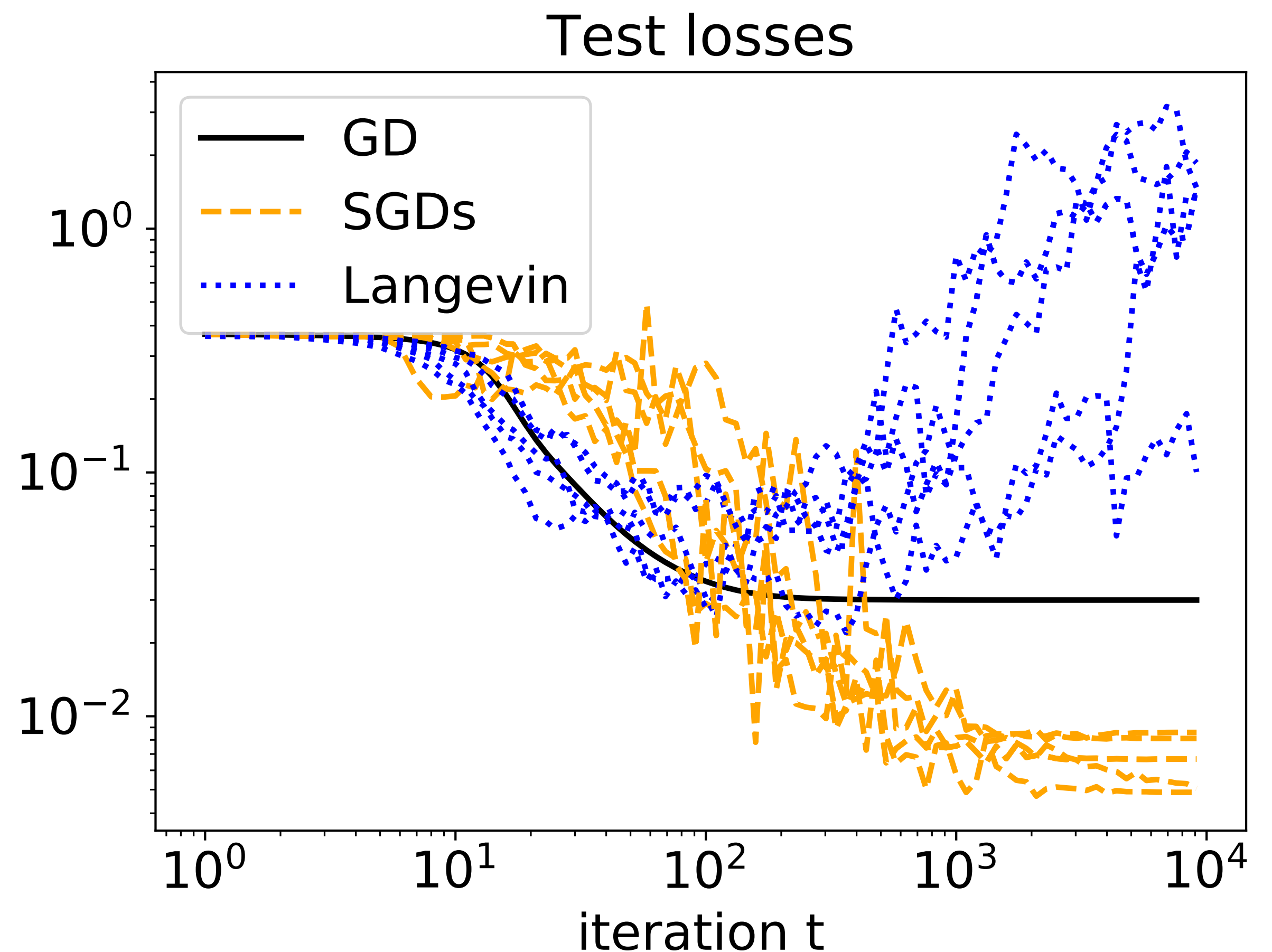
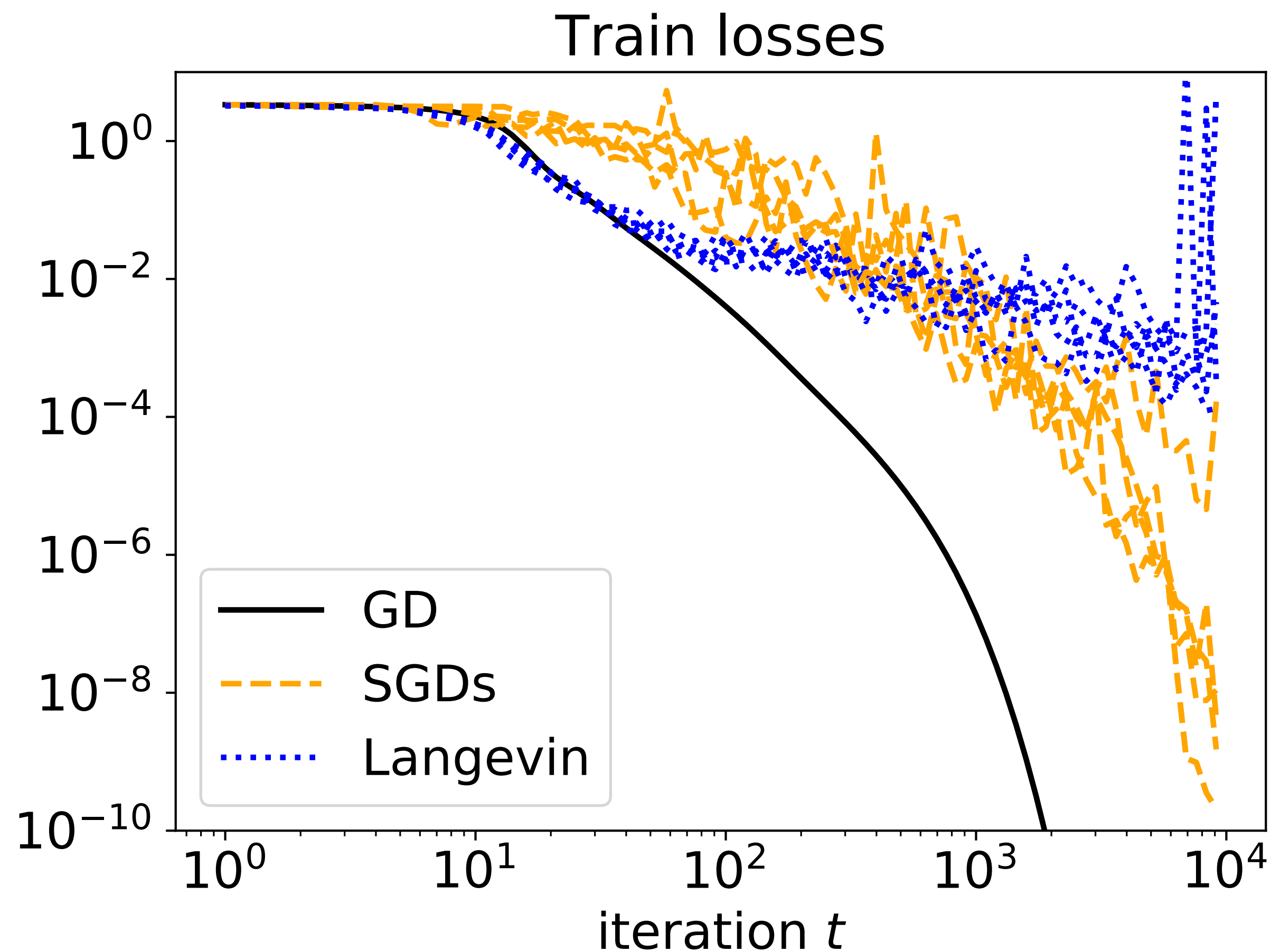
$$du_t = -\nabla_u L(w_t)dt + \sqrt{2\eta^{-1}} d\tilde{B}_t$$

“Wrong SDE”:

$$du_t = -\nabla_u L(w_t)dt + 2\sqrt{\gamma n^{-1} L(w_t)} v_t \odot d\tilde{B}_t$$

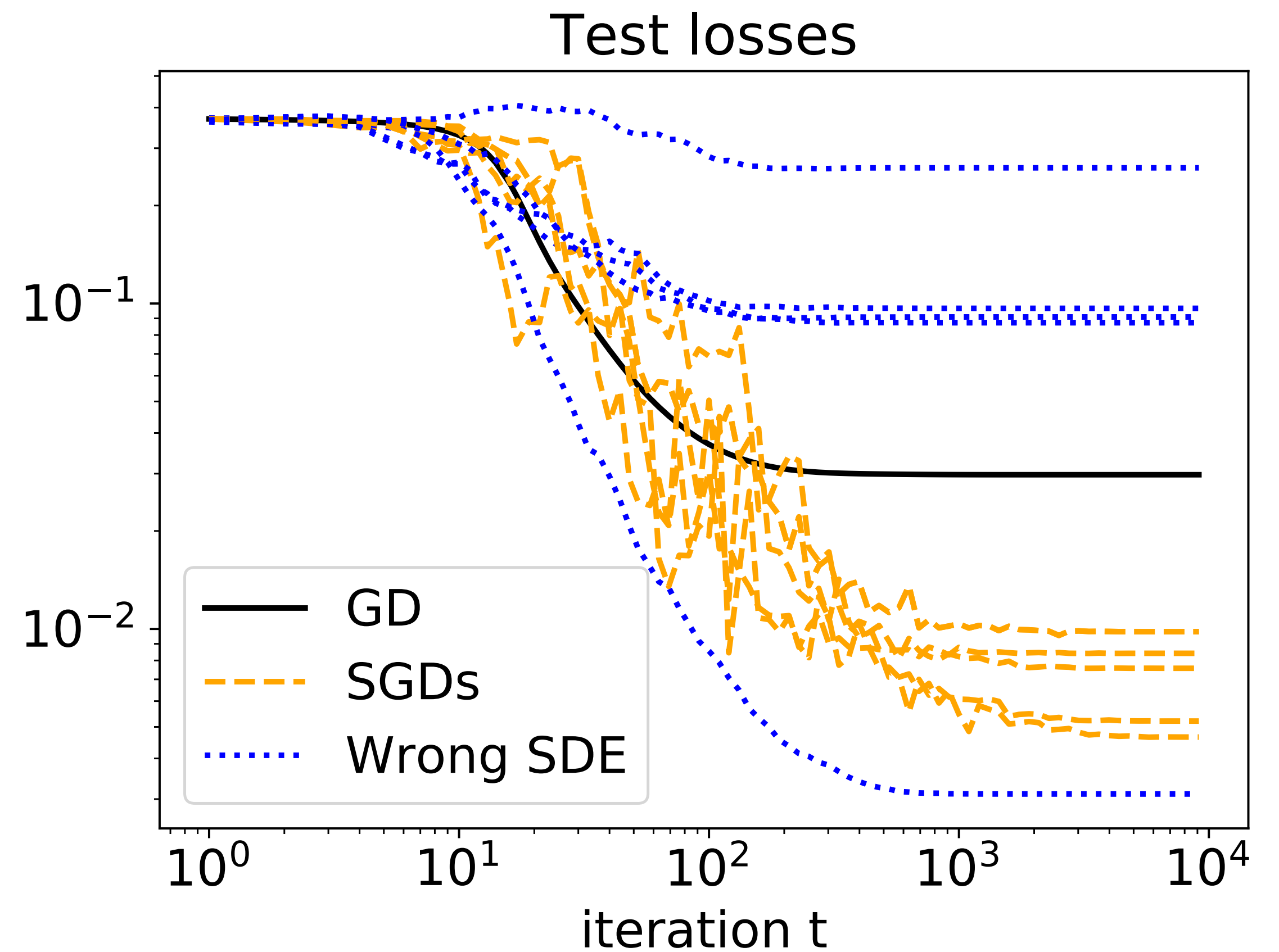
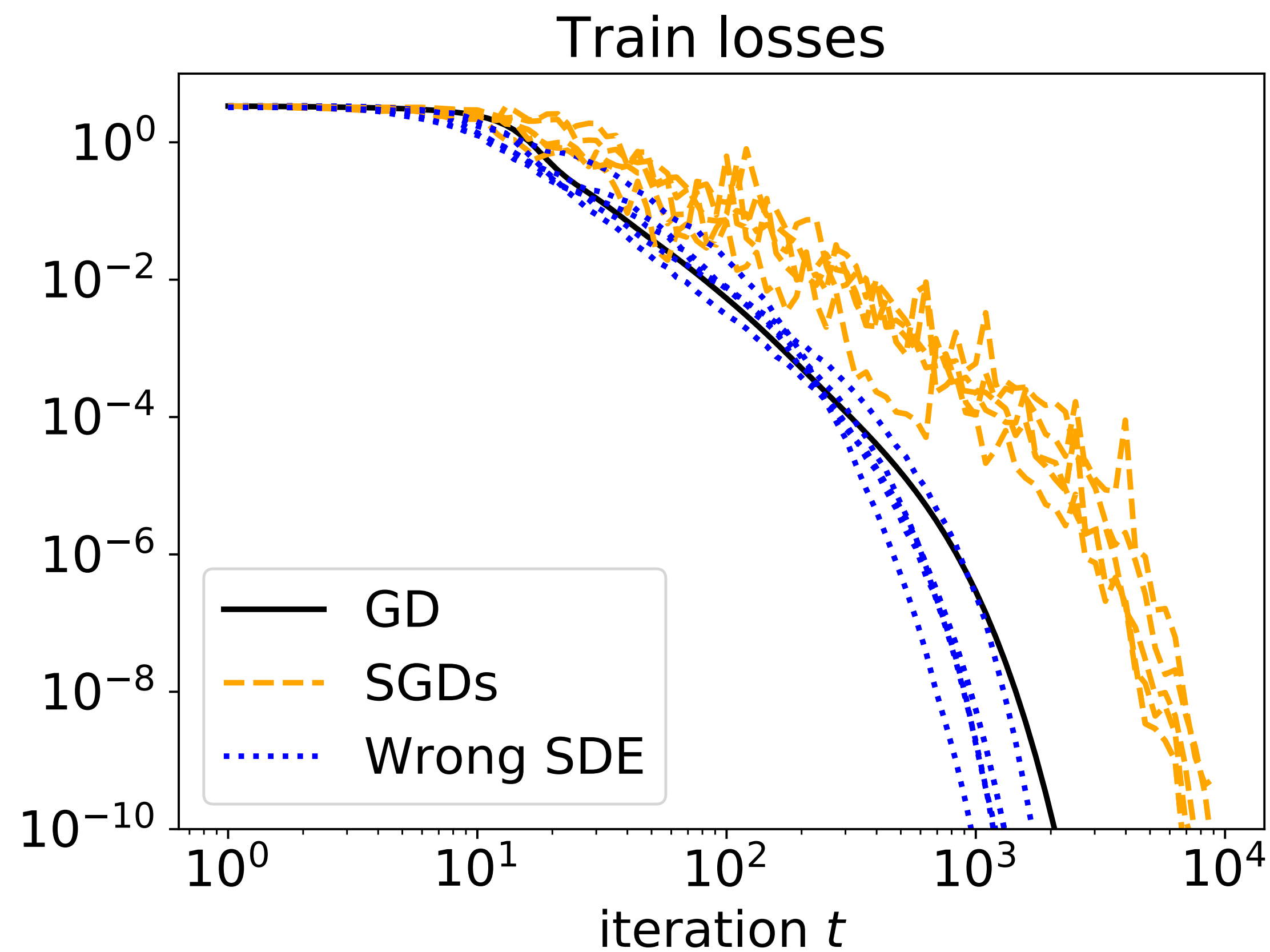
Bonus 0: other SDE modelisations

Overdamped Langevin: $du_t = -\nabla_u L(w_t)dt + \sqrt{2\eta^{-1}}d\tilde{B}_t$



Bonus 0: other SDE modelisations

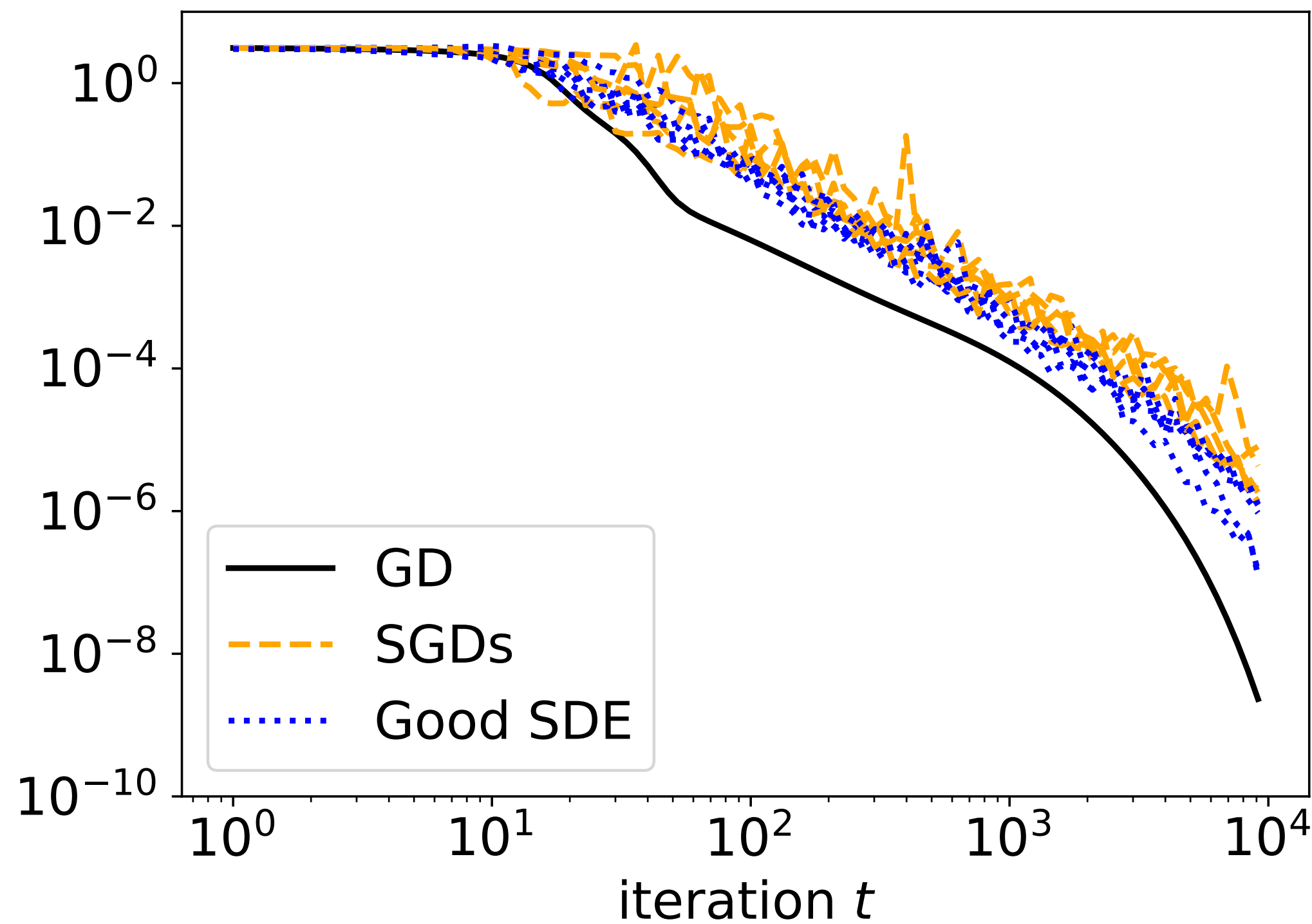
“Wrong SDE”: $du_t = -\nabla_u L(w_t)dt + 2\sqrt{\gamma n^{-1} L(w_t)}v_t \odot d\tilde{B}_t$



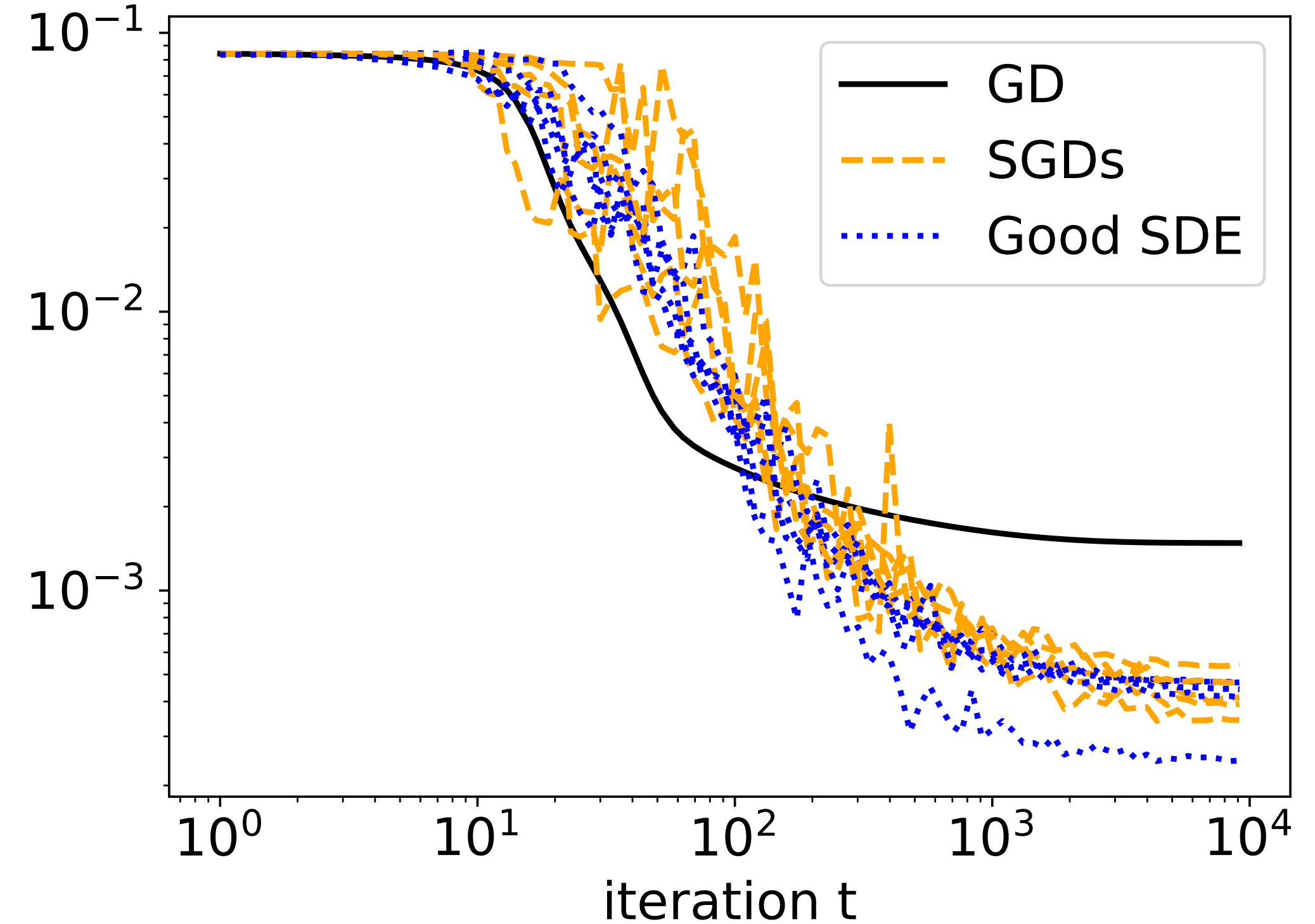
Bonus 0: other SDE modelisations

“Our SDE”:
$$du_t = -\nabla_u L(w_t)dt + 2\sqrt{\gamma n^{-1} L(w_t)}v_t \odot [X^\top dB_t]$$

Train losses



Test losses



Bonus 1: adding label noise

Perturb label at time t $\tilde{y}_{i_t} = y_{i_t} + \Delta_t$ where $\Delta_t \sim \text{unif}\{2\delta_t, -2\delta_t\}$

$$(\delta_t)_{t \in \mathbb{N}} \in \mathbb{R}_+^{\mathbb{N}}$$

“Slowed down loss”: $\tilde{L}(w_t) = L(w_t) + \delta_t^2$

Modified: $\tilde{\alpha}_\infty = \alpha \odot \exp\left(-2\gamma \text{diag}\left(\frac{X^\top X}{n}\right) \int_0^{+\infty} \tilde{L}(\beta_s) ds\right)$

Bonus 1: adding label noise

Perturb label at time t : $\tilde{y}_{i_t} = y_{i_t} + \Delta_t$ where $\Delta_t \sim \text{unif}\{2\delta_t, -2\delta_t\}$ $(\delta_t)_{t \in \mathbb{N}} \in \mathbb{R}_+^{\mathbb{N}}$

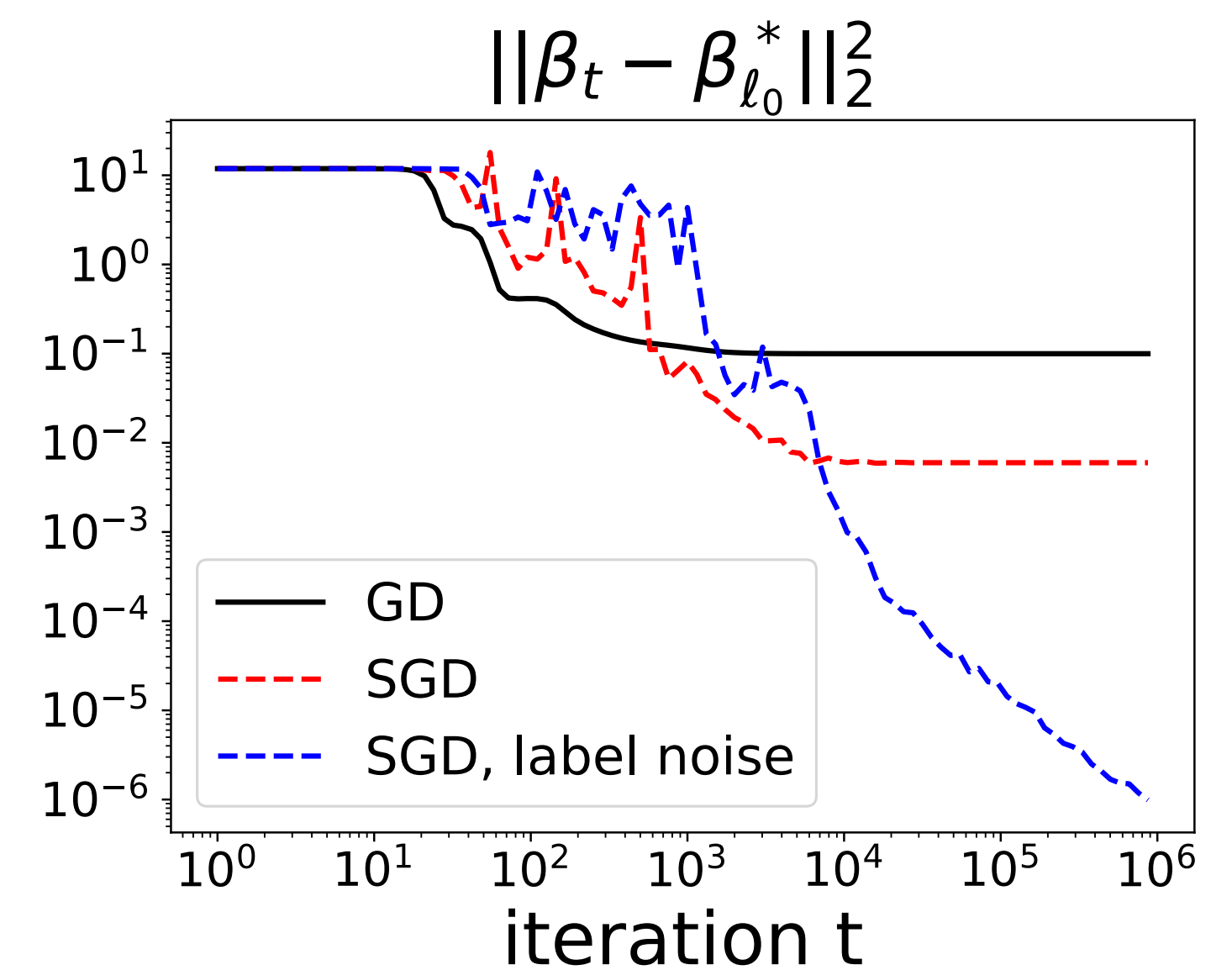
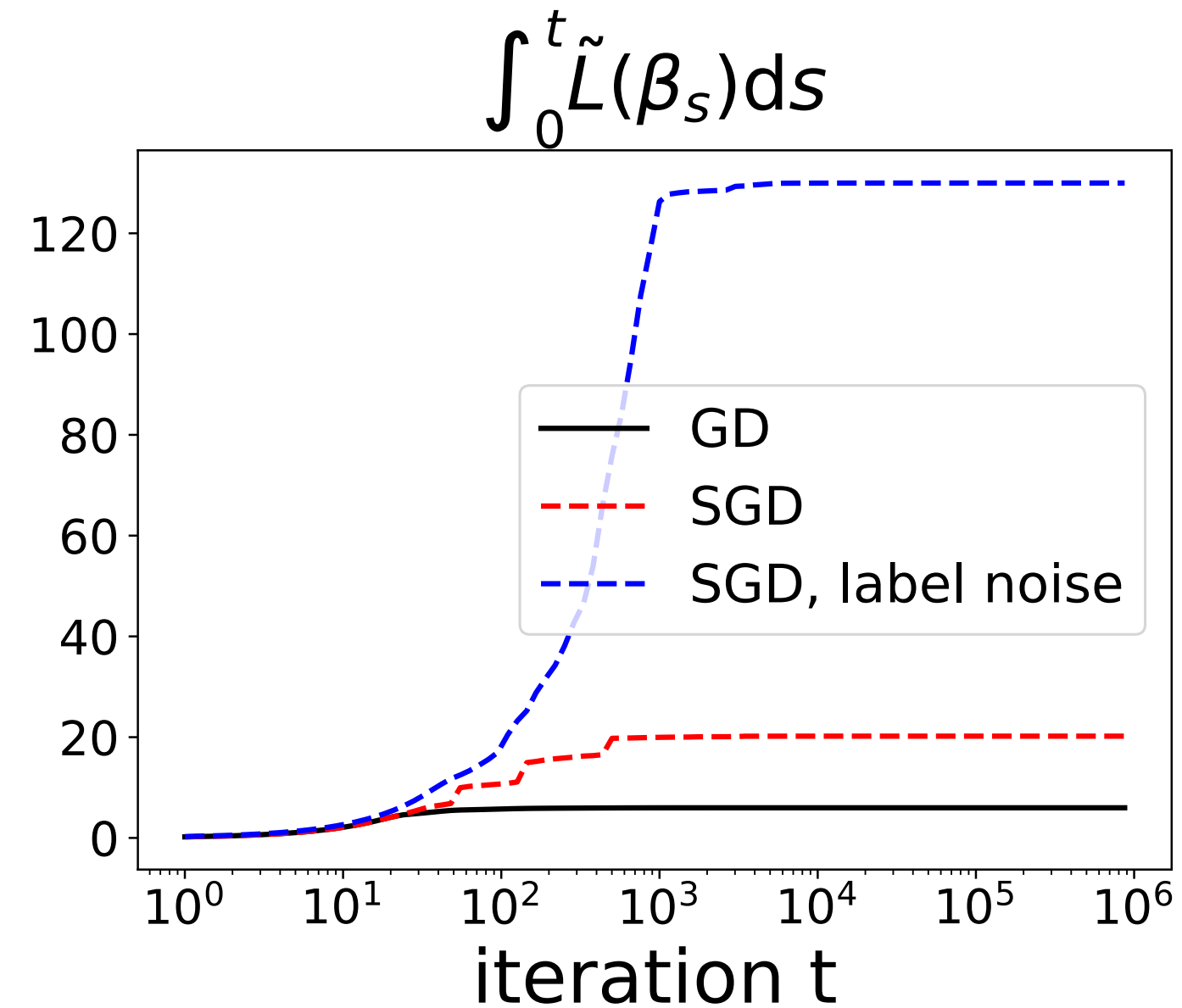
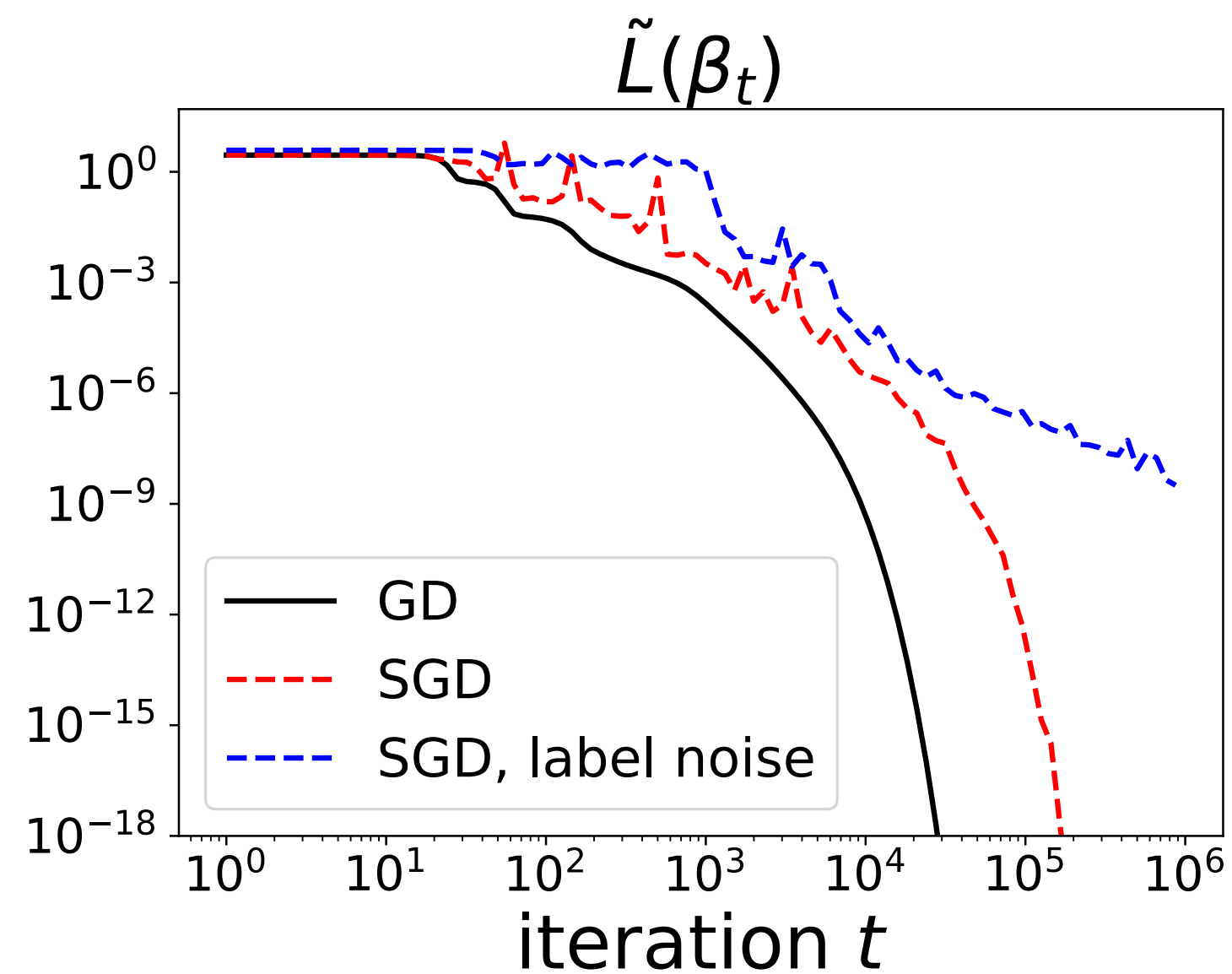
Experimental setup:

$$n = 40 \quad d = 100 \quad \|\beta_{\ell_0}^*\|_0 = 5$$

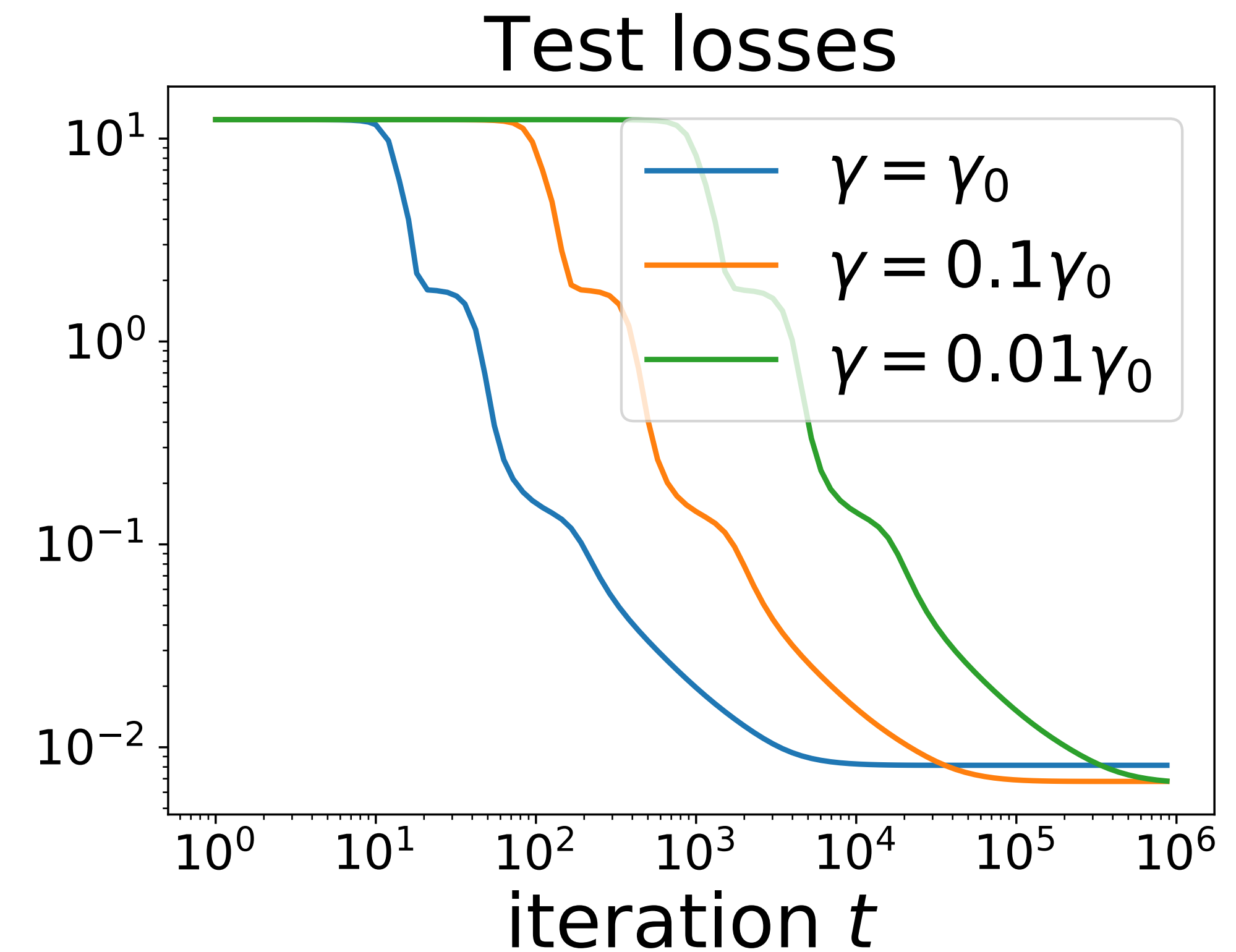
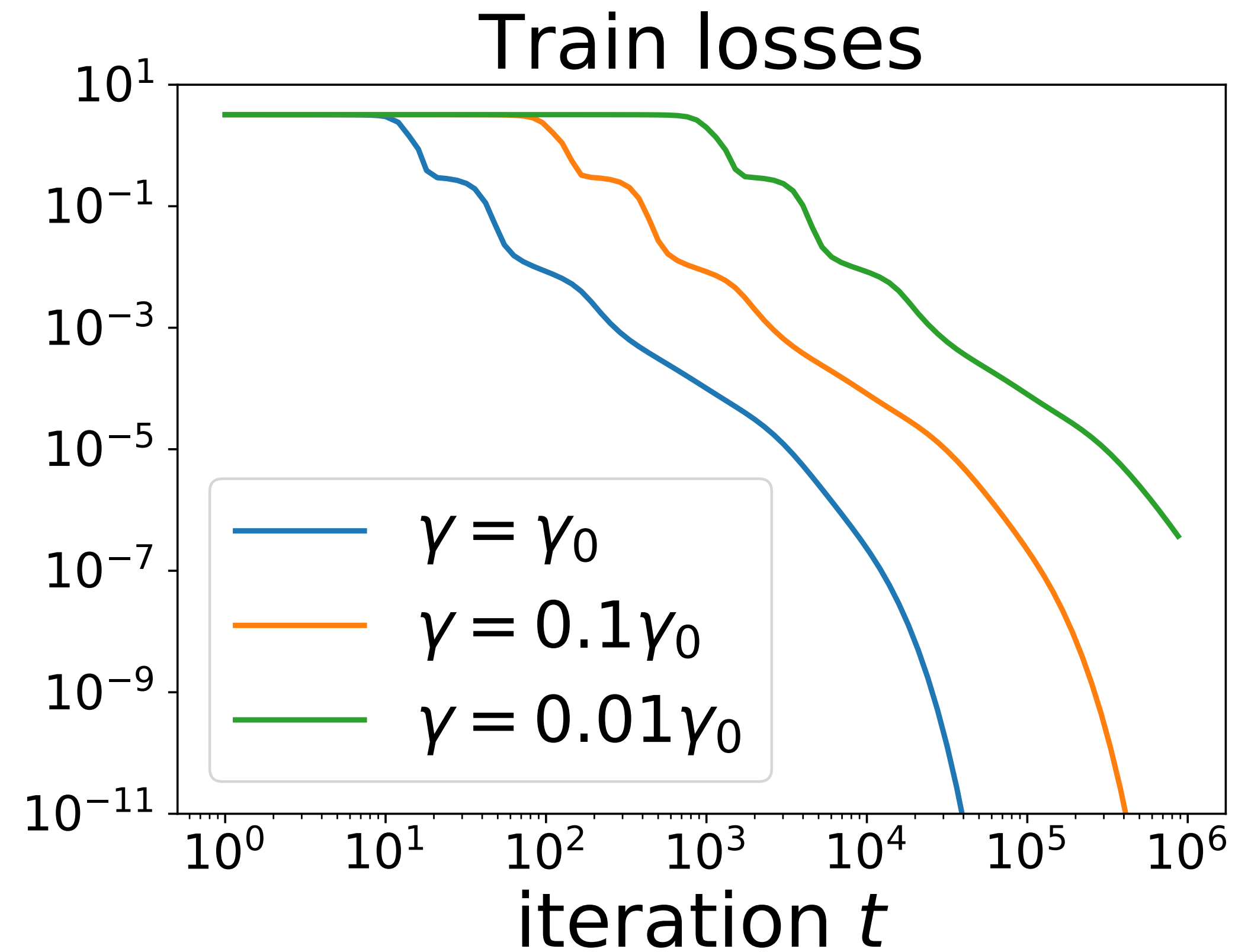
$$t \leq 10^3 : \delta_t = 1$$

$$x_i \sim \mathcal{N}(0, I) \quad y_i = \langle x_i, \beta_{\ell_0}^* \rangle$$

$$t > 10^3 : \delta_t = 0$$

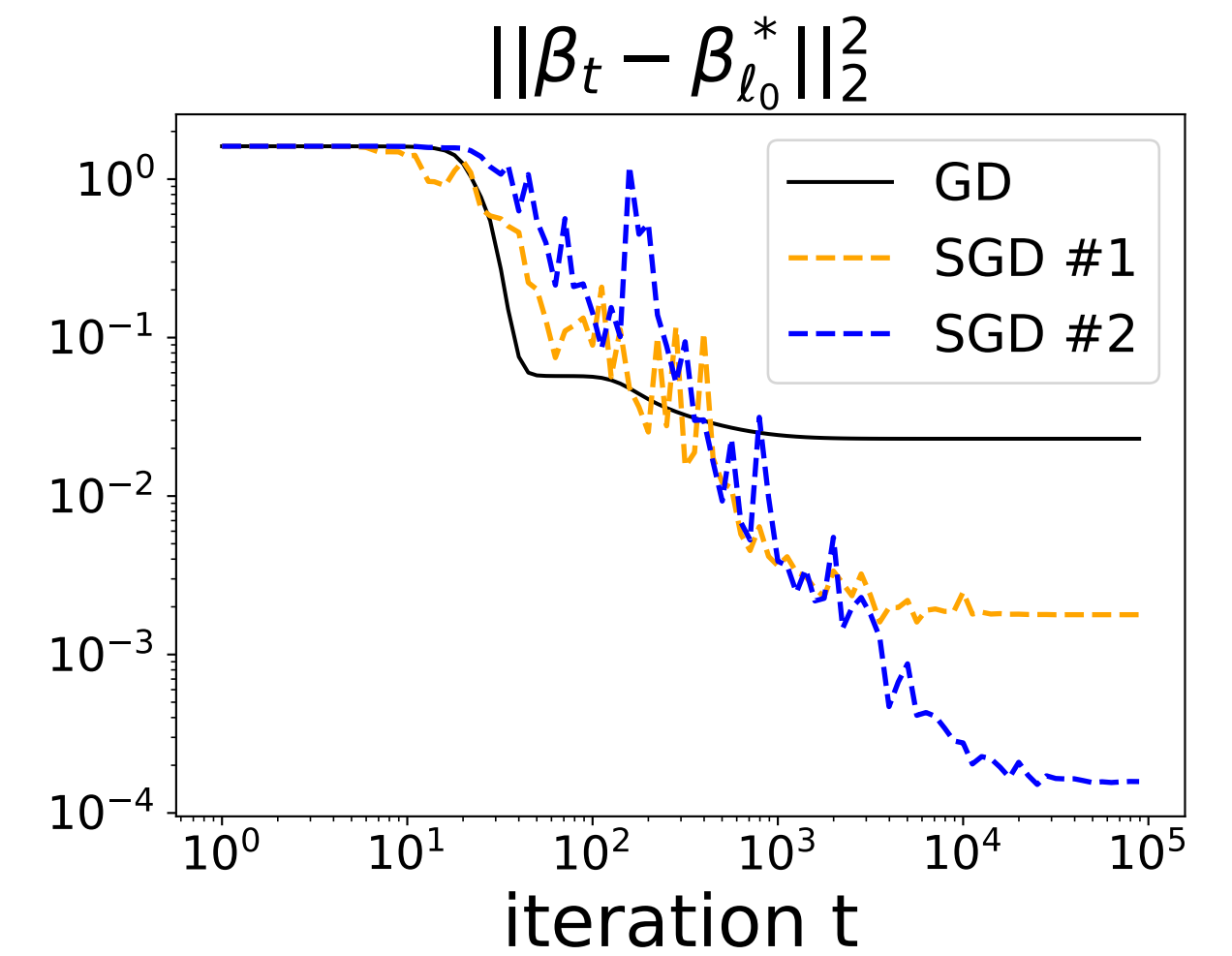
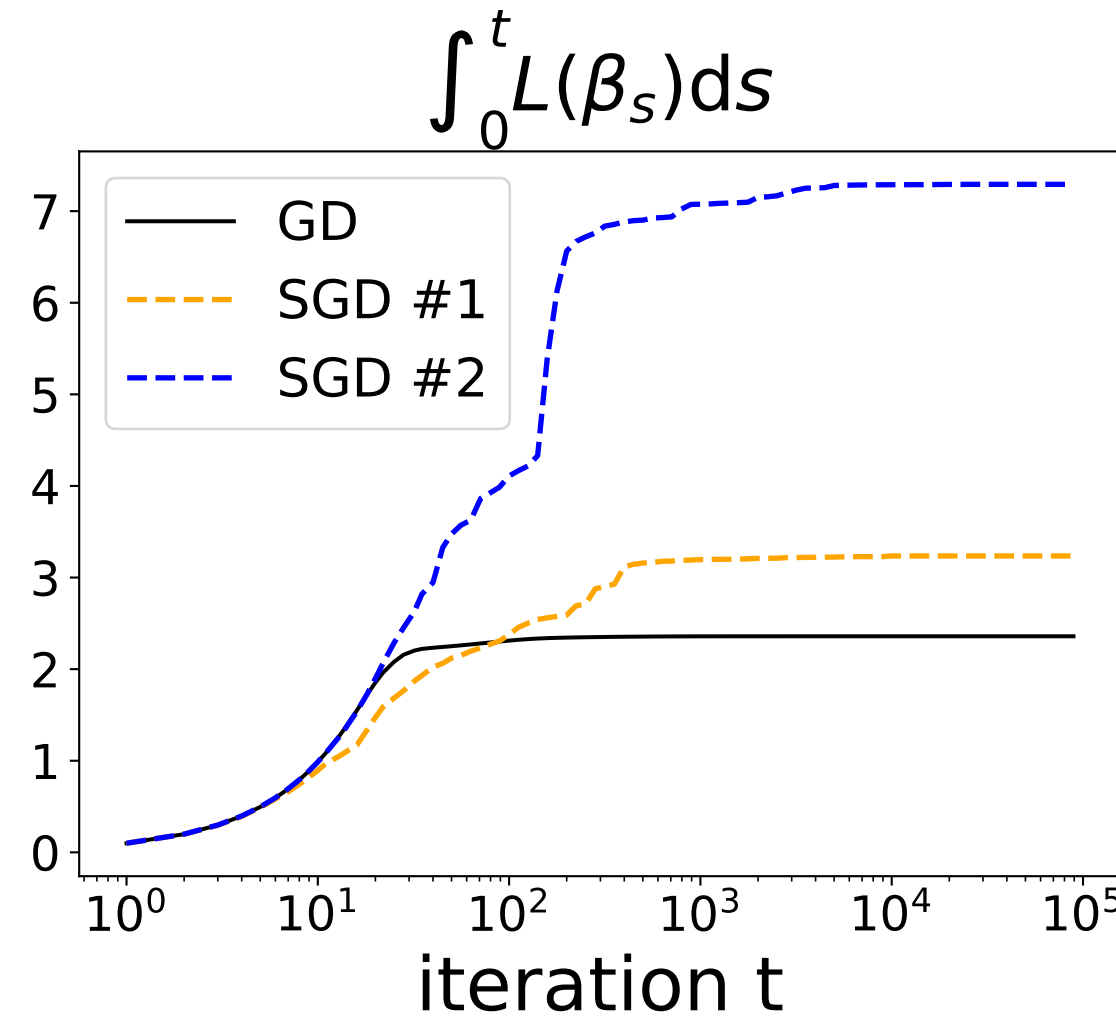
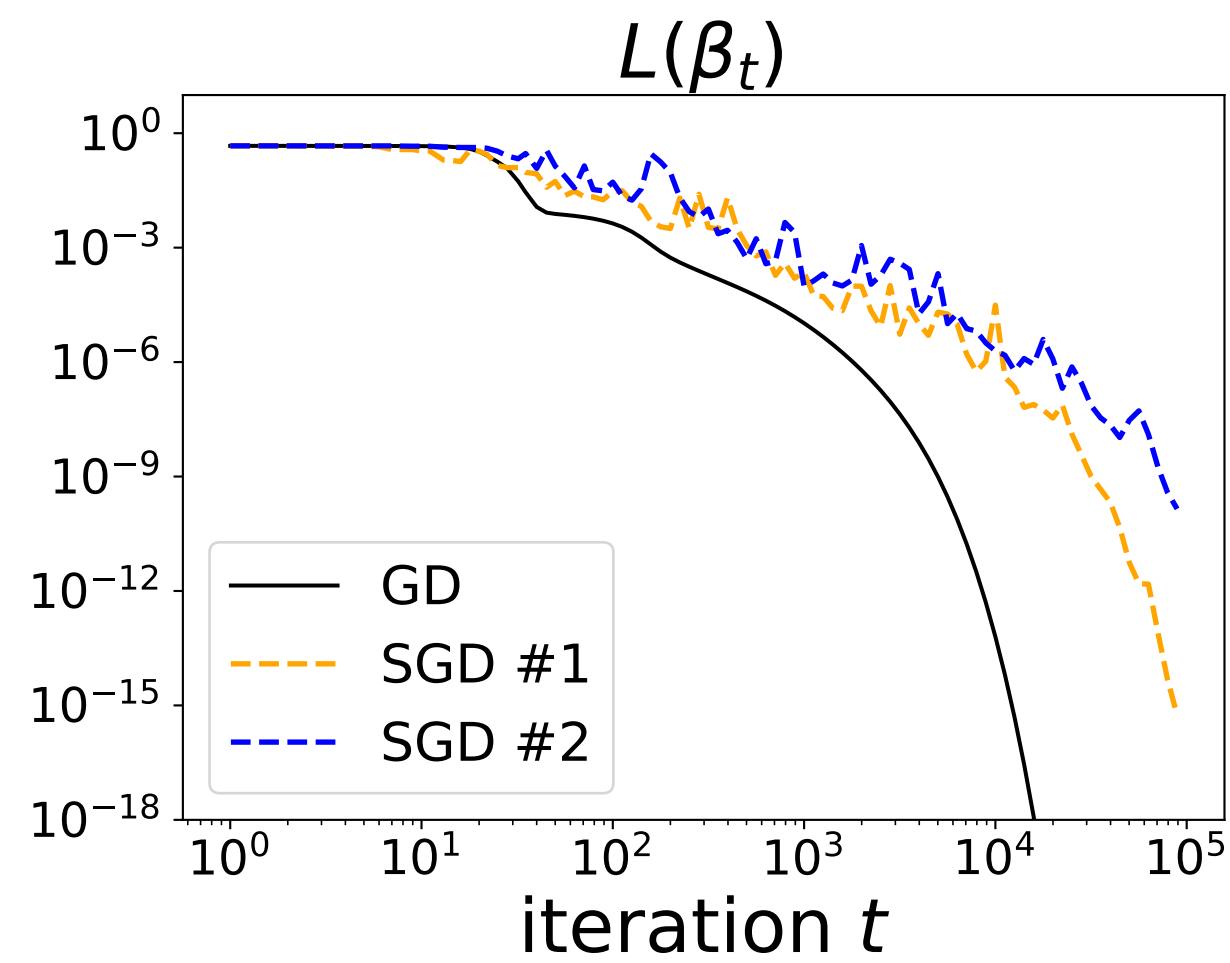


Bonus 2: step-size for GD has very little impact on the implicit bias



Bonus 3: the slower the training loss, the better the bias.

Setting: $n = 40$ $d = 100$ $\|\beta_{\ell_0}^*\|_0 = 5$
 $(\alpha = 0.1)$ $x_i \sim \mathcal{N}(0, I)$ $y_i = \langle x_i, \beta_{\ell_0}^* \rangle$



$$\underbrace{\alpha_\infty}_{\text{“effective” initialisation}} = \underbrace{\alpha}_{\text{initialisation scale}} \odot \exp\left(-2\gamma \text{diag}\left(\frac{X^\top X}{n}\right) \underbrace{\int_0^{+\infty} L(\beta_s) ds}_{\text{training loss}}\right) < \underbrace{\alpha}_{\text{initialisation scale}}$$

stochastic !