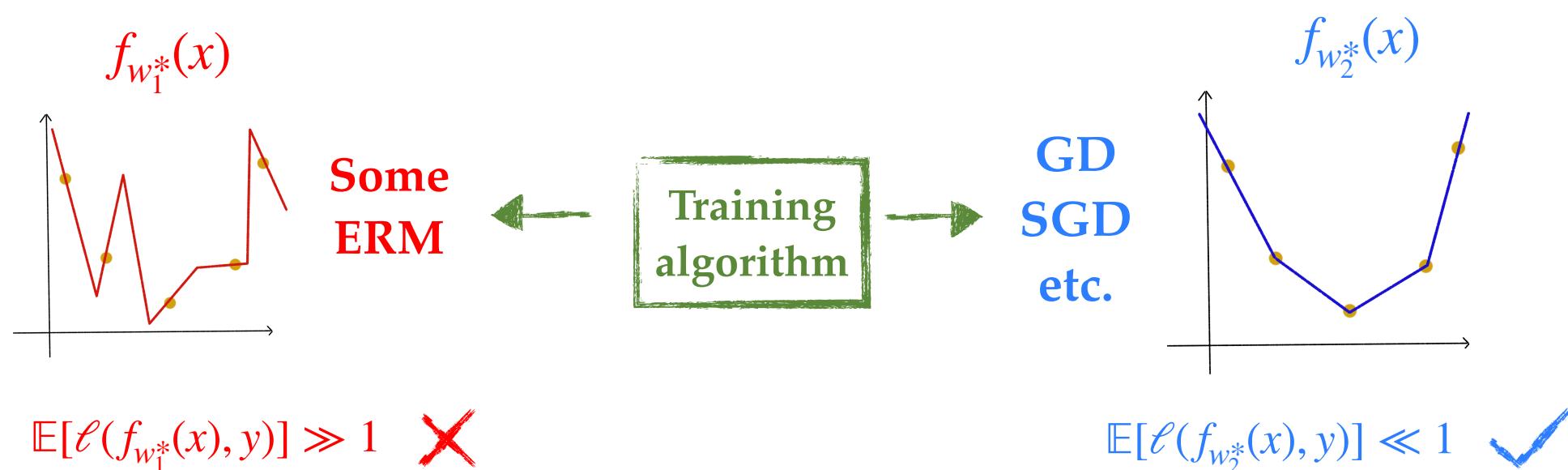


Implicit Bias of SGD for Diagonal Linear Networks: A Provable Benefit of Stochasticity

Scott Pesme, Lucas Pillaud-Vivien, Nicolas Flammarion
TML lab

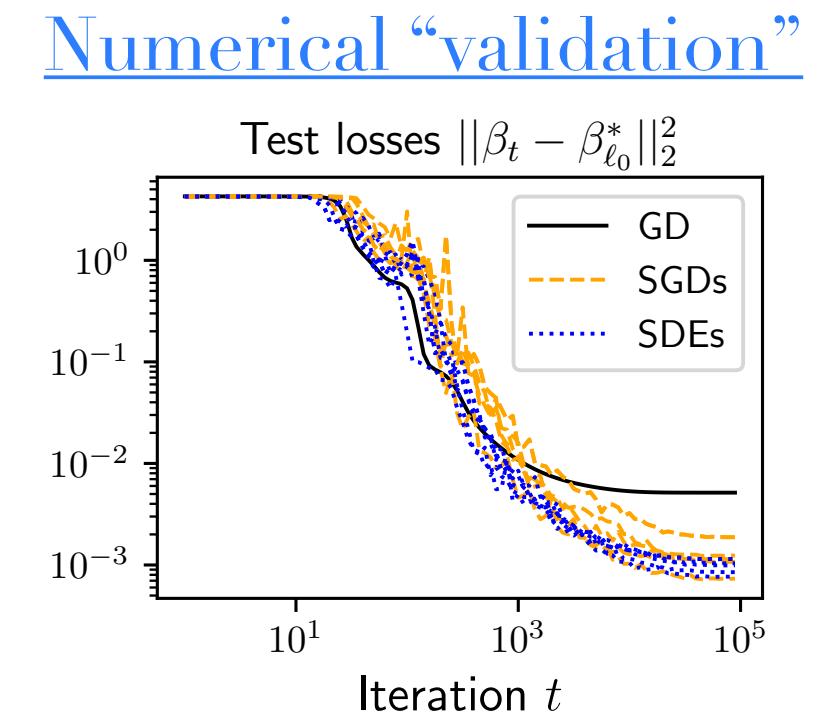
EPFL

Why the concept of implicit bias?

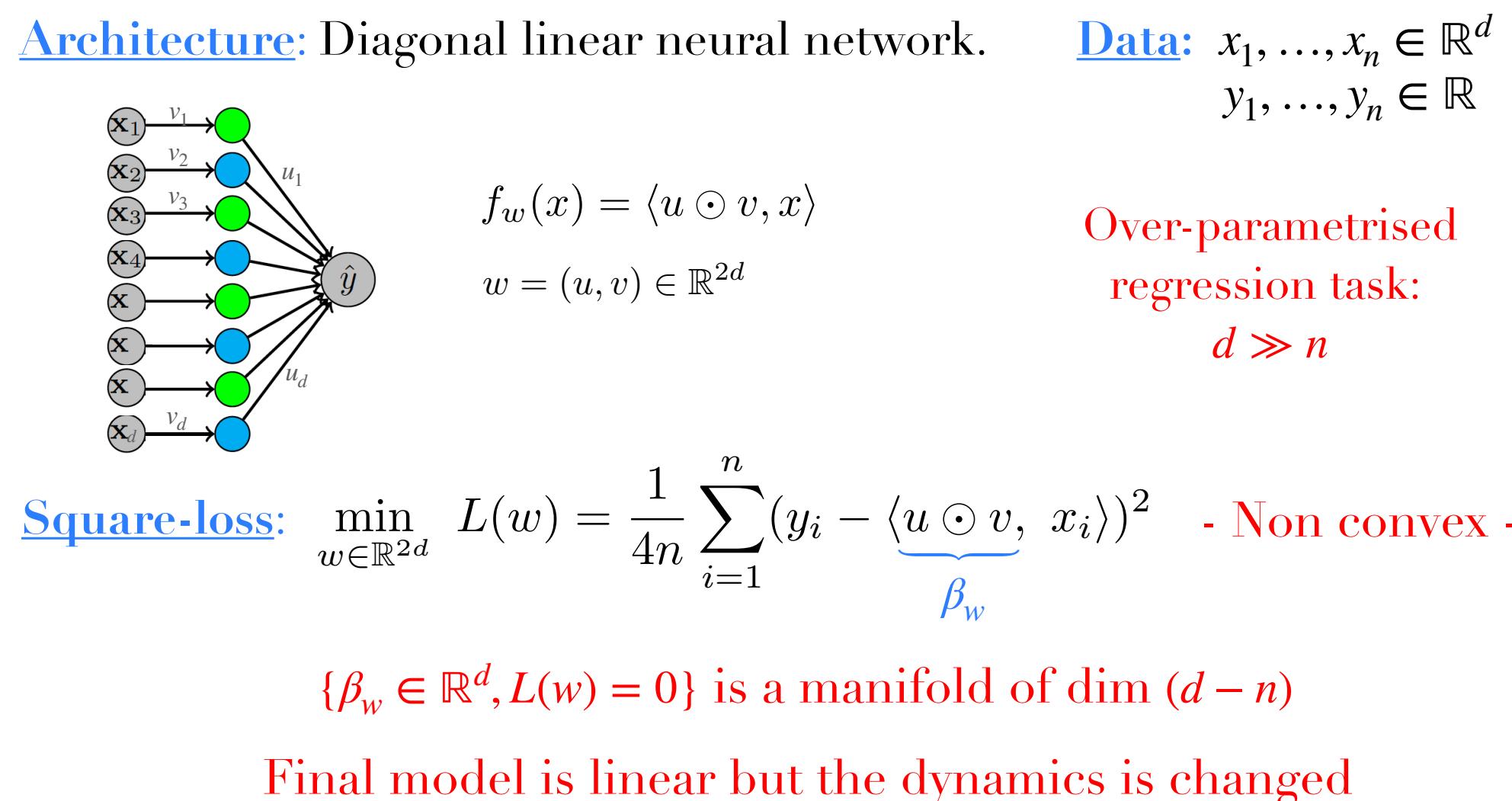


How can we accurately model SGD in a continuous way?

SGD $u_{t+1} = u_t - \gamma \langle \beta_{w_t} - \beta^*, x_{i_t} \rangle x_{i_t} \odot v_t$ $= u_t - \gamma \nabla_u L(w_t) + \underbrace{\gamma v_t \odot [X^\top \xi_{i_t}(w_t)]}_{\text{Zero mean noise}}$	Stochastic Gradient Flow (SGF) $du_t = -\nabla_u L(w_t) dt + 2\sqrt{\gamma n^{-1} L(w_t)} v_t \odot [X^\top dB_t]$	state dependent! (i) matching structure: belongs to $\text{span}(x_1 \odot v, \dots, x_n \odot v)$ (ii) matching covariance $\Sigma_{\text{SGD}}(w)$
---	--	--



2-layer diagonal linear network



Main result: convergence and implicit bias of the stochastic gradient flow

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 $\lambda_{\max} = \lambda_{\max}(X^\top X/n)$

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 β_∞^α

Implicit Bias \rightarrow • This solution β_∞^α satisfies

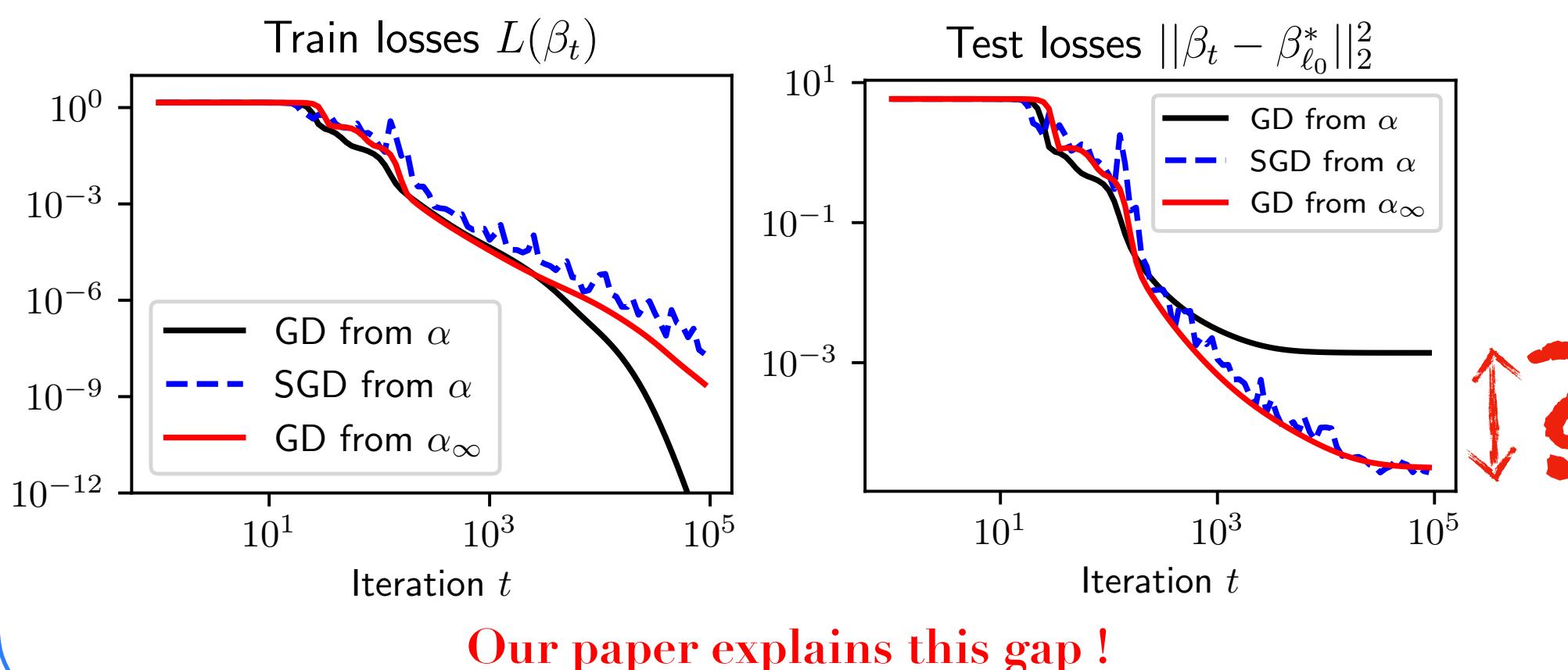
$$\beta_\infty^\alpha = \arg \min_{\beta \in \mathbb{R}^d, \langle \beta, x_i \rangle = y_i} \phi_{\alpha_\infty}(\beta) \quad \text{where } \underbrace{\alpha_\infty}_{\text{“effective”}} = \alpha \odot \exp\left(-2\gamma \text{diag}\left(\frac{X^\top X}{n}\right) \int_0^{+\infty} \underbrace{L(\beta_s)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

SGD empirically performs better than GD

Gradient flow: $dw_t = -\nabla_w L(w_t) dt, \quad \left| \begin{array}{l} u_{t=0} = \alpha \in \mathbb{R}^d \\ v_{t=0} = 0 \end{array} \right. \Rightarrow \beta_{w_{t=0}} = 0$

Implicit bias: $\beta_{w_t} \rightarrow \beta_\infty^\alpha = \arg \min_{\beta, \langle \beta, x_i \rangle = y_i} \phi_\alpha(\beta)$ where $\phi_\alpha(\beta) \underset{\alpha \rightarrow 0}{\sim} \|\beta\|_1$
(Woodworth et al. 2020)

What about SGD? Sparse gold model $\beta_{\ell_0}^*$ and $y_i = \langle x_i, \beta_{\ell_0}^* \rangle$



Interpretation and observations

GF vs SGF:

• Implicit bias of SGF is the same as GF but with an effective initialisation:

$$\alpha_\infty < \alpha \Rightarrow \beta_\infty^{\alpha, \text{SGF}} \text{ is “sparser” than } \beta_\infty^{\alpha, \text{GF}}$$

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

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

Under additional assumption (boundedness of the iterates):

$$\frac{\alpha_\infty}{\alpha} \underset{\alpha \rightarrow 0}{\leq} \left(\frac{\alpha^2}{\|\beta_{\ell_1}^*\|_1} \right)^\zeta \text{ for some } \zeta > 0$$

Convergence holds for a fixed step-size:

• This is due to the fact that the noise vanishes at the optimum

Sketch of proof

Stochastic mirror descent with time varying potential:

$$d\nabla\phi_{\alpha_t}(\beta_t) = -\nabla_\beta L(\beta_t) dt + \underbrace{\sqrt{\gamma n^{-1} L(\beta_t)}}_{\text{stochastic & time dependent}} X^\top dB_t \in \text{span}(x_1, \dots, x_n)$$

$$\alpha_t = \alpha \odot \exp\left(-2\gamma \text{diag}\left(\frac{X^\top X}{n}\right) \int_0^t L(\beta_s) ds\right)$$

Assuming convergence, the KKT conditions immediately give the result:

$$\left| \begin{array}{l} \nabla\phi_{\alpha_\infty}(\beta_\infty^\alpha) \in \text{span}(x_i) \\ L(\beta_\infty^\alpha) = 0 \end{array} \right. \xrightarrow{\text{(KKT)}} \beta_\infty^\alpha = \arg \min_{\beta \in \mathbb{R}^d, \langle \beta, x_i \rangle = y_i} \phi_{\alpha_\infty}(\beta)$$

Proving the convergence of the flow $(\beta_t)_{t \geq 0}$ is technically the hardest part:

- Use of appropriate stochastic Lyapunov functions
- Use of martingale concentration inequalities to control the stochastic terms