

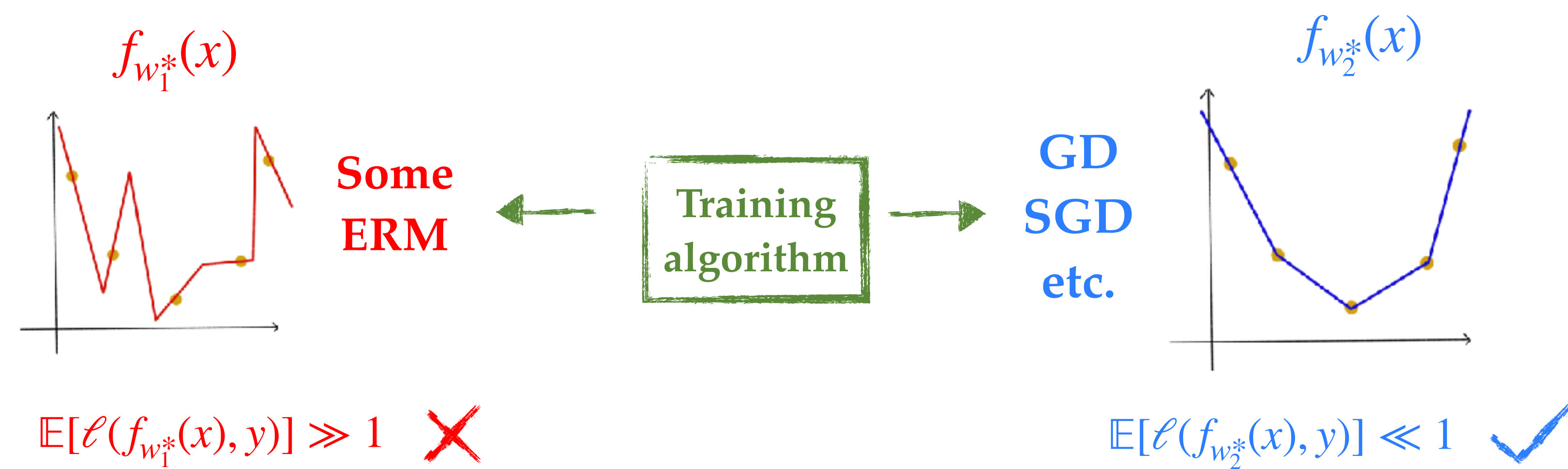
Implicit Bias of SGD for Diagonal Linear Networks:

A Provable Benefit of Stochasticity

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

EPFL

Why the concept of implicit bias?



How can we accurately model SGD in a continuous way?

SGD

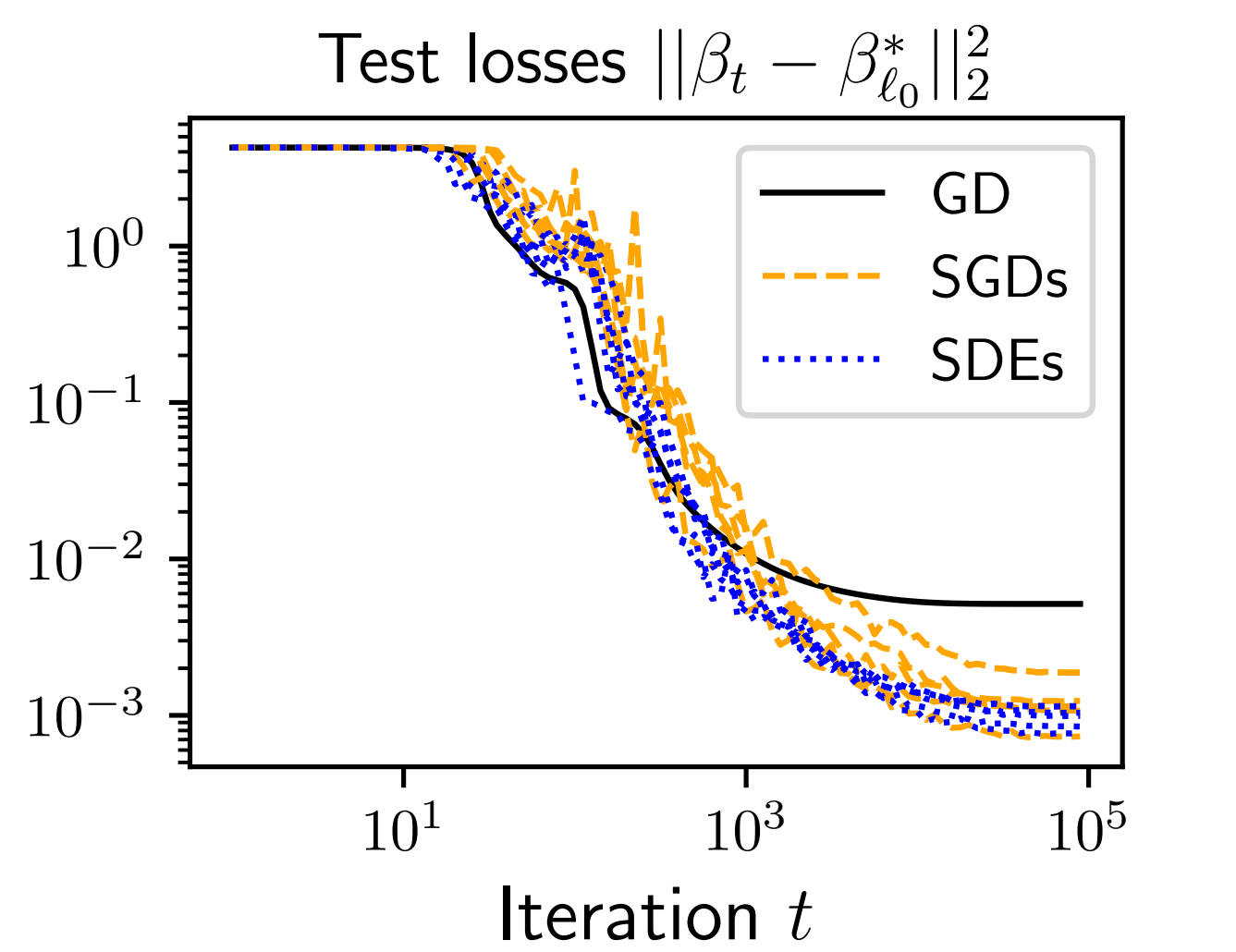
$$u_{t+1} = u_t - \gamma \langle \beta_{w_t} - \beta^*, x_{i_t} \rangle x_{i_t} \odot v_t \\ = u_t - \gamma \nabla_u L(w_t) + \underbrace{\gamma v_t \odot [X^\top \xi_{i_t}(w_t)]}_{\text{Zero mean noise}}$$

Stochastic Gradient Flow (SGF)

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

(i) matching structure: belongs to $\text{span}(x_1 \odot v, \dots, x_n \odot v)$ (ii) matching covariance $\Sigma_{\text{SGD}}(w)$

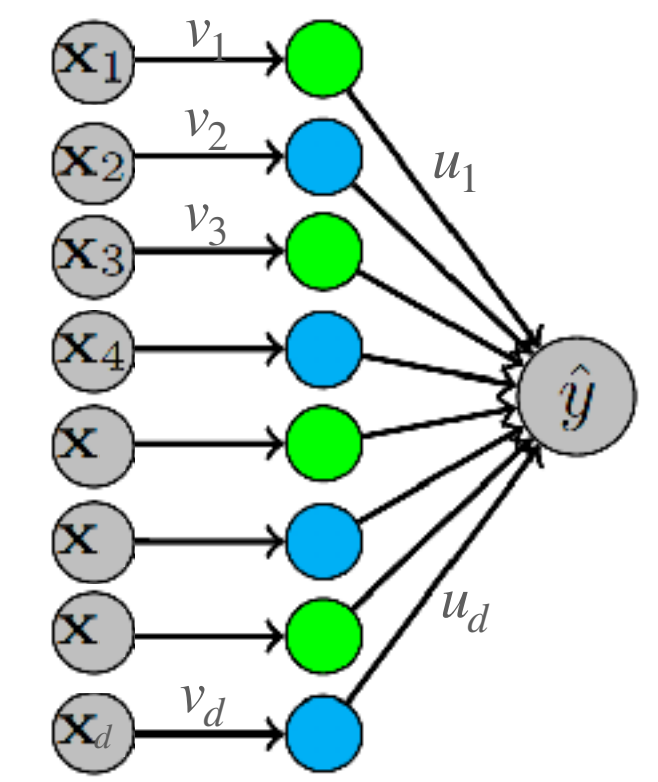
Numerical “validation”



2-layer diagonal linear network

Architecture: Diagonal linear neural network.

Data: $x_1, \dots, x_n \in \mathbb{R}^d$
 $y_1, \dots, y_n \in \mathbb{R}$



$$f_w(x) = \langle u \odot v, x \rangle \\ w = (u, v) \in \mathbb{R}^{2d}$$

Over-parametrised regression task:
 $d \gg n$

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 -

$\{\beta_w \in \mathbb{R}^d, L(w) = 0\}$ is a manifold of dim $(d - n)$

Final model is linear but the dynamics is changed

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

$$\lambda_{\max} = \lambda_{\max}(X^\top X/n) \\ \beta_{\ell_1}^* = \underset{\beta \text{ s.t. } X\beta=y}{\operatorname{argmin}} \|\beta\|_1$$

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

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

initialisation scale

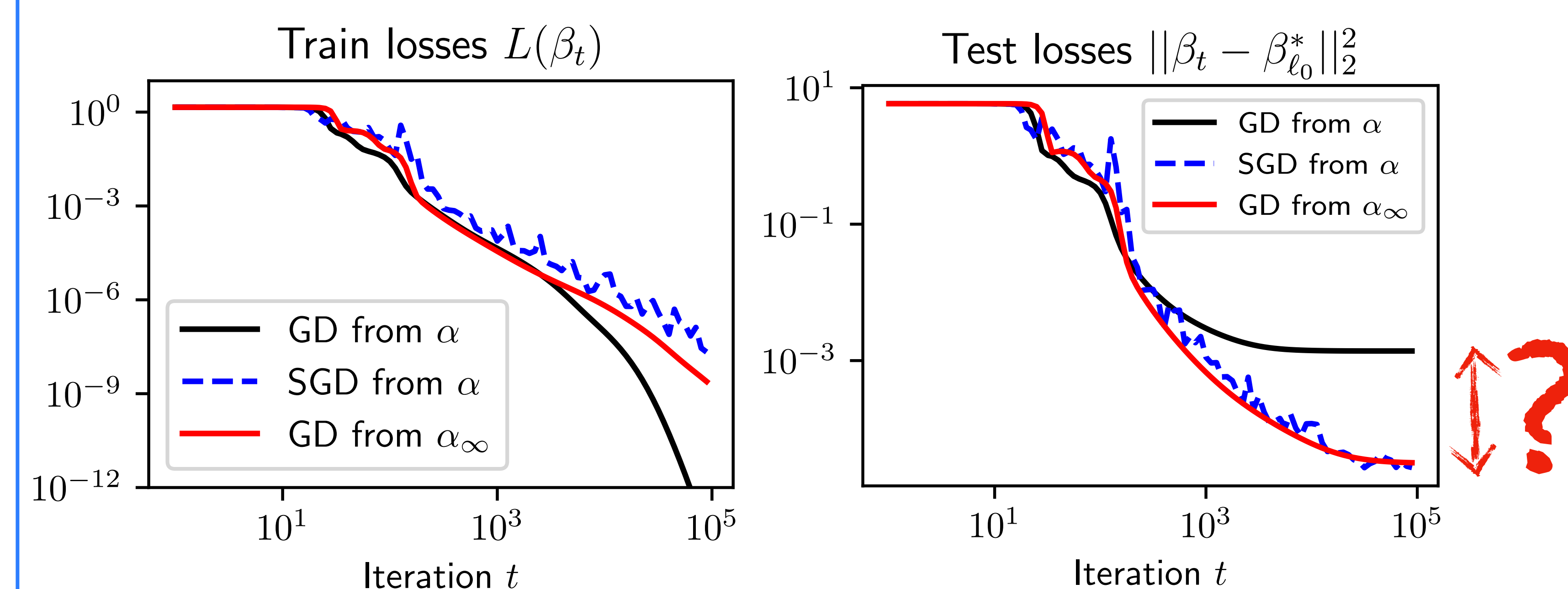
$$< \underbrace{\alpha}_{\text{initialisation scale}}$$

SGD empirically performs better than GD

Gradient flow: $dw_t = -\nabla_w L(w_t) dt$, $\begin{cases} u_{t=0} = \alpha \in \mathbb{R}^d \\ v_{t=0} = 0 \end{cases} \rightarrow \beta_{w_{t=0}} = 0$

Implicit bias: $\beta_{w_t} \rightarrow \beta_\infty^\alpha = \underset{\beta \in \mathbb{R}^d, \langle \beta, x_i \rangle = y_i}{\operatorname{argmin}} \phi_\alpha(\beta)$ where $\phi_\alpha(\beta) \underset{\alpha \rightarrow 0}{\sim} \|\beta\|_1$
(Woodworth et al. 2020) $\underset{\alpha \rightarrow \infty}{\sim} \|\beta\|_2$

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



Our paper explains this gap !

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:

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

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

Assuming convergence, the KKT conditions immediately give the result:

$$\begin{cases} \nabla \phi_{\alpha_\infty}(\beta_\infty) \in \text{span}(x_i) \\ L(\beta_\infty) = 0 \end{cases} \xrightarrow{\text{(KKT)}} \beta_\infty^\alpha = \underset{\beta \in \mathbb{R}^d, \langle \beta, x_i \rangle = y_i}{\operatorname{argmin}} \phi_{\alpha_\infty}(\beta)$$

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