

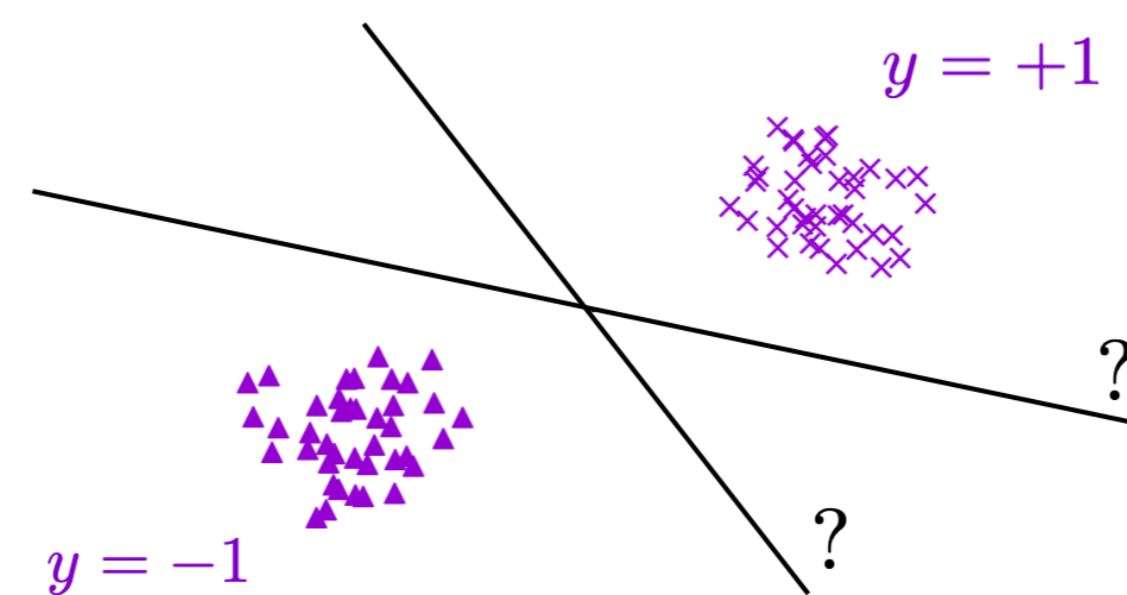
# Implicit Bias of Mirror Flow on Separable Data

Scott Pesme (Inria Grenoble)  
 Radu Dragomir (Télécom Paris)  
 Nicolas Flammarion (EPFL)



## Setup: logistic regression

$$\min_{\beta} L(\beta) = \sum_{i=1}^n \ln(1 + e^{-y_i \langle \beta, x_i \rangle})$$



**Assumption:** linearly separable data  $\rightarrow$  the loss is minimised 'at infinity':

$$\lim_{s \rightarrow \infty} L(s\beta^*) = 0 \quad \text{for } \beta^* \in \mathcal{S} := \{\beta^* \in \mathbb{R}^d, y_i \langle \beta^*, x_i \rangle \geq 1, \forall i\}$$

*set of vectors defining separating hyperplanes*

For mirror flow, what is the directional limit  $\lim_{t \rightarrow \infty} \frac{\beta_t}{\|\beta_t\|}$  of the iterates  $\beta_t$ ?

Many possible solutions in  $\mathcal{S}$ : which one is preferred by the method? (*implicit regularisation*)

This question is important to understand the generalisation properties!

## The gradient method: mirror flow

$$d\nabla \phi(\beta_t) = -\nabla L(\beta_t) dt$$

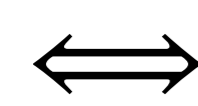
*strictly convex potential*

**Motivation:** reparametrisation  $\beta = F(\theta)$ , then under some (restrictive) conditions:

Gradient flow on  $\theta \mapsto L(F(\theta)) \implies$  Mirror flow on  $\beta \mapsto L(\beta)$

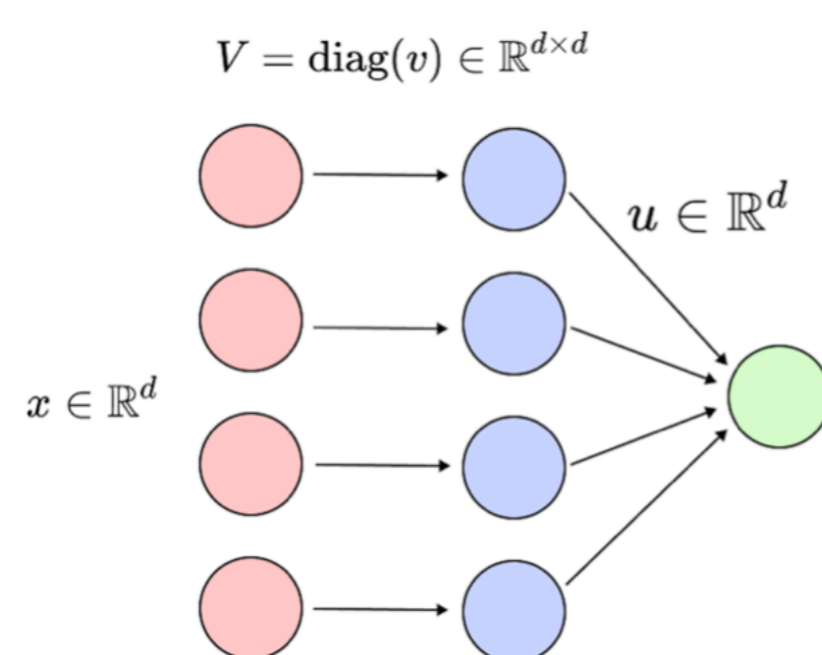
**Example:**  $\beta = F(u, v) = u \odot v$  "diagonal neural networks"

Gradient flow on  $L(u \odot v)$

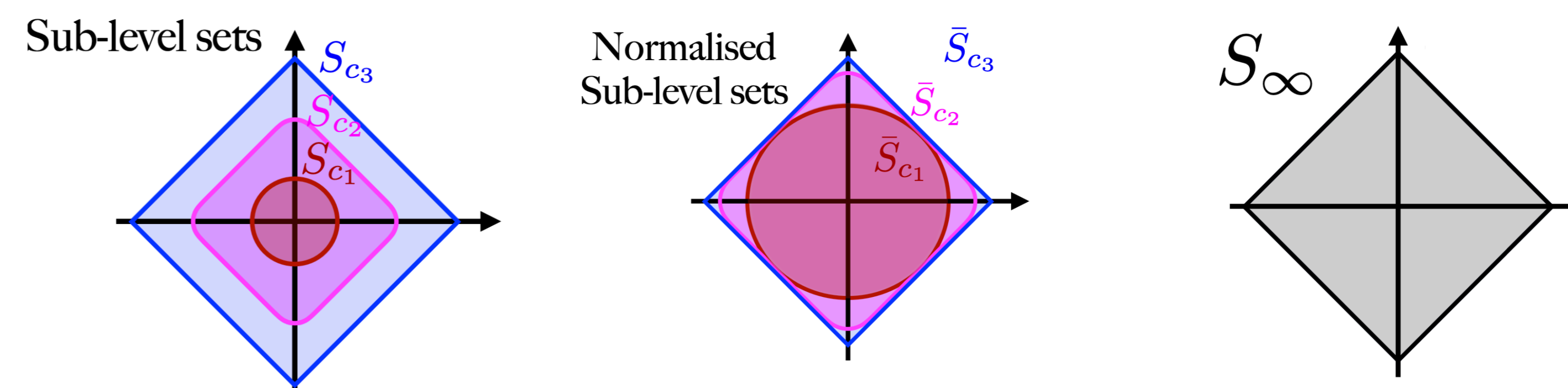


Mirror flow on  $L(\beta)$  with hyperbolic potential

$$\phi(\beta) = \sum_{i=1}^d \left( \beta_i \operatorname{arcsinh}(\beta_i) - \sqrt{\beta_i^2 + 1} \right)$$



## Horizon function : geometry of $\phi$ 'at infinity'



$$S_c = \{\beta : \phi(\beta) \leq c\}$$

$$\bar{S}_c = S_c / \max_{\beta \in S_c} \|\beta\|$$

$$S_\infty = \lim_{c \rightarrow \infty} \bar{S}_c$$

We say that  $\phi$  admits a horizon function if  $\lim_{c \rightarrow \infty} \bar{S}_c$  exists

**Horizon function:**  $\phi_\infty(\beta) = \inf\{r > 0 : \frac{\beta}{r} \in S_\infty\}$

*Minkowski gauge of  $S_\infty$   
 Asymmetric norm whose unit ball is  $S_\infty$*

## Main result: convergence and implicit bias

*e.g. polynomial, semialgebraic, subanalytic, log-exp...*

**Theorem:** if  $\phi$  is tame, it admits a horizon  $\phi_\infty$  and the mirror flow iterates  $\beta_t$  converge in direction towards the vector  $\bar{\beta}_\infty$  satisfying the  $\phi_\infty$ -max margin problem:

$$\lim_{t \rightarrow \infty} \frac{\beta_t}{\|\beta_t\|} =: \bar{\beta}_\infty \propto \operatorname{argmin}\{\phi_\infty(\beta^*) : \beta^* \in \mathcal{S}\}$$

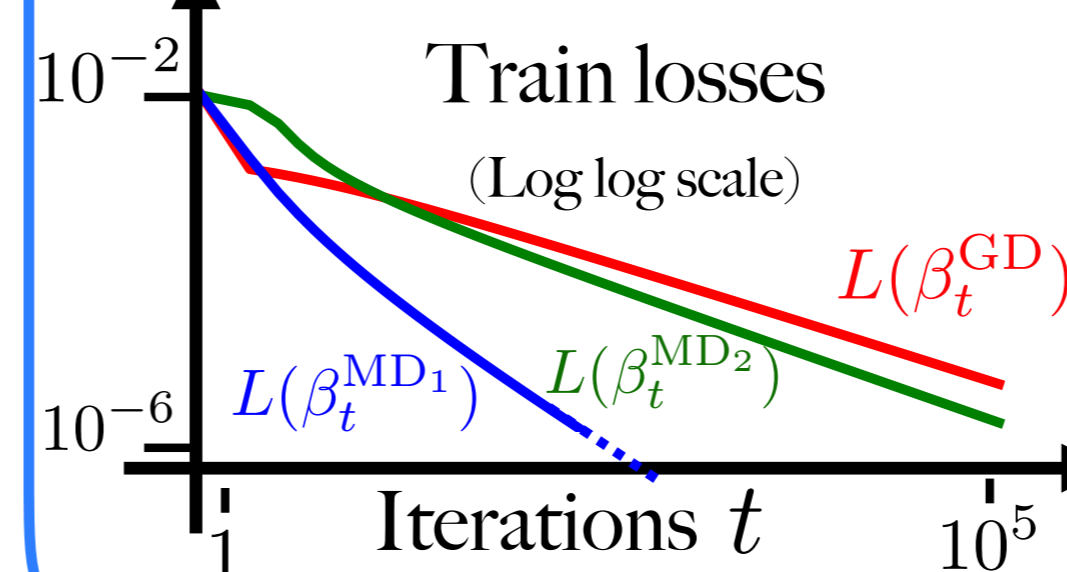
*set of separating hyperplanes*

## Experiments

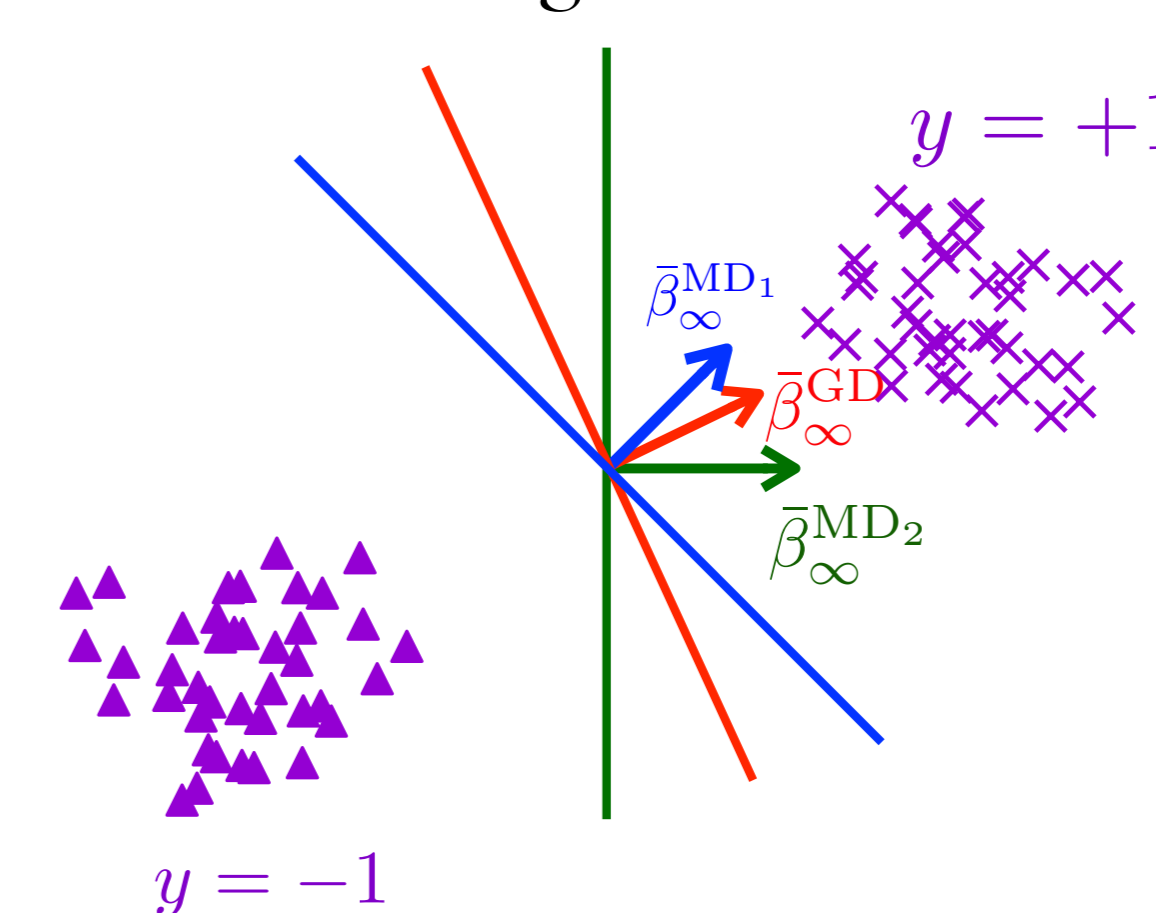
Mirror descent with:

$$\phi^{\text{GD}}(\beta) = \|\beta\|_2^2 \quad \phi^{\text{MD}_1}(\beta) = \sum_{i=1,2} \cosh(\beta_i)$$

$$\phi^{\text{MD}_2}(\beta) = \sum_{i=1,2} \beta_i \operatorname{arcsinh}(\beta_i) - \sqrt{\beta_i^2 + 1}$$



Limiting directions



$\phi_\infty$ -max-margins

