
Saddle-to-Saddle Dynamics in Diagonal Linear Networks

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In this paper we fully describe the trajectory of gradient flow over 2-layer diagonal
2 linear networks for the regression setting in the limit of vanishing initialisation. We
3 show that the limiting flow successively jumps from a saddle of the training loss to
4 another until reaching the minimum ℓ_1 -norm solution. We explicitly characterise
5 the visited saddles as well as the jump times through a recursive algorithm reminis-
6 cent of the Homotopy algorithm used for computing the Lasso path. Starting from
7 the zero vector, coordinates are successively activated until the minimum ℓ_1 -norm
8 solution is recovered, revealing an incremental learning. Our proof leverages a
9 convenient arc-length time-reparametrisation which enables to keep track of the
10 transitions between the jumps. Our analysis requires negligible assumptions on the
11 data, applies to both under and overparametrised settings and covers complex cases
12 where there is no monotonicity of the number of active coordinates. We provide
13 numerical experiments to support our findings.

14 1 Introduction

15 Strikingly simple algorithms such as gradient descent are driving forces for deep learning and have
16 led to remarkable empirical results. Nonetheless, understanding the performances of such methods
17 remains a challenging and exciting mystery: (i) their global convergence on highly non-convex losses
18 is far from being trivial and (ii) the fact that they lead to solutions which generalise well [40] is still
19 not fully understood.

20 To explain this second point, a major line of work has focused on the concept of implicit regularisation:
21 amongst the infinite space of zero-loss solutions, the optimisation process must be implicitly biased
22 towards solutions which have good generalisation properties for the considered real-world prediction
23 tasks. Many papers have therefore shown that gradient methods have the fortunate property of
24 asymptotically leading to solutions which have a well-behaving structure [31, 18, 11].

25 Aside from these results which mostly focus on characterising the asymptotic solution, a slightly
26 different point of view has been to try to describe the full trajectory. Indeed it has been experimentally
27 observed that gradient methods with small initialisations have the property of learning models of
28 increasing complexity across the training of neural networks [23]. This behaviour is usually referred
29 to as *incremental learning* or as a *saddle-to-saddle process* and describes learning curves which are
30 piecewise constant: the training process makes very little progress for some time, followed by a
31 sharp transition where a new “feature” is suddenly learned. In terms of optimisation trajectory, this
32 corresponds to the iterates “jumping” from a saddle of the training loss to another.

33 Several settings exhibiting such dynamics for small initialisation have been considered: matrix
34 and tensor factorisation [35, 21], simplified versions of diagonal linear networks [17, 6], linear
35 networks [16, 36, 20], 2-layer neural networks with orthogonal inputs [8] and matrix sensing [1, 27,

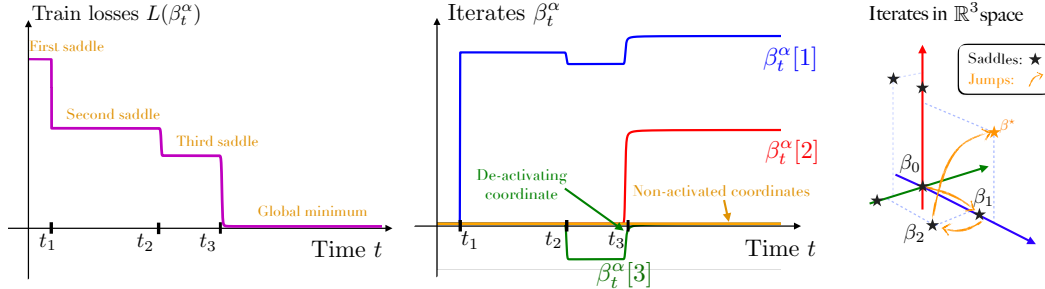


Figure 1: Gradient flow $(\beta_t^\alpha)_t$ with small initialisation scale α over a 2-layer diagonal linear network (for the precise experimental setting, see Appendix A). *Left*: Training loss across time, the learning is piecewise constant. *Middle*: The magnitudes of the coordinates are plotted across time: the process is piecewise constant. *Right*: In the \mathbb{R}^3 space in which the iterates evolve (the remaining coordinates stay at 0), the iterates jump from a saddle of the training loss to another. The jumping times t_i as well as the visited saddles β_i are entirely predicted by our theory.

36 22]. However, all these results require restrictive assumptions on the data or only characterise the
 37 first jump. Obtaining a complete picture of the saddle-to-saddle process by describing all the visited
 38 saddles and jump times is mathematically challenging and still missing. We intend to fill this gap
 39 by considering diagonal linear networks which are simplified neural networks that have received
 40 significant attention lately [39, 38, 19, 34, 14] as they are ideal proxy models for gaining a deeper
 41 understanding of complex phenomena such as saddle-to-saddle dynamics.

42 1.1 Informal statement of the main result

43 In this paper, we provide a full description of the trajectory of gradient flow over 2-layer diagonal
 44 linear networks in the limit of vanishing initialisation. The main result is informally presented here.

45 **Theorem 1** (Main result, informal). *In the regression setting and in the limit of vanishing initial-*
 46 *isation, the trajectory of gradient flow over a 2-layer diagonal linear network converges towards*
 47 *a limiting process which is piecewise constant: the iterates successively jump from a saddle of the*
 48 *training loss to another; each visited saddle and jump time can recursively be computed through an*
 49 *algorithm (Algorithm 1) reminiscent of the Homotopy algorithm for the Lasso.*

50 The incremental learning stems from the particular structure of the saddles as they correspond to min-
 51 imisers of the training loss with a constraint on the set of non-zero coordinates. The saddles therefore
 52 correspond to sparse vectors which partially fit the dataset. For simple datasets, a consequence of our
 53 main result is that **the limiting trajectory successively starts from the zero vector and successively**
 54 **learns the support of the sparse ground truth vector until reaching it. However, we make**
 55 **minimal assumptions on the data and our analysis also holds for complex datasets.** In that case,
 56 the successive active sets are not necessarily increasing in size and coordinates can deactivate as well
 57 as activate until reaching the minimum ℓ_1 -norm solution (see Figure 1 (middle) for an example of a
 58 deactivating coordinate). The regression setting and the diagonal network architecture are introduced
 59 in Section 2. Section 3 provides an intuitive construction of the limiting saddle-to-saddle dynamics
 60 and presents the algorithm that characterises it. Our main result regarding the convergence of the
 61 iterates towards this process is presented in Section 4 and further discussion is provided in Section 5.

62 2 Problem setup and leveraging the mirror structure

63 2.1 Setup

64 **Linear regression.** We study a linear regression problem with inputs $(x_1, \dots, x_n) \in (\mathbb{R}^d)^n$ and
 65 outputs $(y_1, \dots, y_n) \in \mathbb{R}^n$. We consider the typical quadratic loss:

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

66 We make no assumption on the number of samples n nor the dimension d . The only assumption we
 67 make on the data throughout the paper is that the inputs (x_1, \dots, x_n) are in *general position*. In order
 68 to state this assumption, let $X \in \mathbb{R}^{n \times d}$ be the feature matrix whose i^{th} row is x_i and let $\tilde{x}_j \in \mathbb{R}^n$ be
 69 its j^{th} column for $j \in [d]$.

70 **Assumption 1** (General position). *For any $k \leq \min(n, d)$ and arbitrary signs $\sigma_1, \dots, \sigma_k \in \{-1, 1\}$,
 71 the affine span of any k points $\sigma_1 \tilde{x}_{j_1}, \dots, \sigma_k \tilde{x}_{j_k}$ does not contain any element of the set $\{\pm \tilde{x}_j, j \neq$
 72 $j_1, \dots, j_k\}$.*

73 This assumption is slightly technical but is standard in the Lasso literature [37]. Note that it is
 74 not restrictive as it is almost surely satisfied when the data is drawn from a continuous probability
 75 distribution [37, Lemma 4].

76 **2-layer diagonal linear network.** In an effort to understand the training dynamics of neural networks,
 77 we consider a 2-layer diagonal linear network which corresponds to writing the regression vector β as

$$\beta_w = u \odot v \text{ where } w = (u, v) \in \mathbb{R}^{2d}. \quad (2)$$

78 This parametrisation can be interpreted as a simple neural network $x \mapsto \langle u, \sigma(\text{diag}(v)x) \rangle$ where u
 79 are the output weights, the diagonal matrix $\text{diag}(v)$ represents the inner weights, and the activation
 80 σ is the identity function. We refer to $w = (u, v) \in \mathbb{R}^{2d}$ as the *neurons* and to $\beta := u \odot v \in \mathbb{R}^d$
 81 as the *prediction parameter*. With the parametrisation (2), the loss function F over the parameters
 82 $w = (u, v) \in \mathbb{R}^{2d}$ is defined as:

$$F(w) := L(u \odot v) = \frac{1}{2n} \sum_{i=1}^n ((u \odot v, x_i) - y_i)^2. \quad (3)$$

83 Though this reparametrisation is simple, the associated optimisation problem is non-convex and
 84 highly non-trivial training dynamics already occur. The critical points of the function F exhibit a
 85 very particular structure, as highlighted in the following proposition proven in Appendix B.

86 **Proposition 1.** *All the critical points w_c of F which are not global minima, i.e., $\nabla F(w_c) = \mathbf{0}$ and
 87 $F(w_c) > \min_w F(w)$, are necessarily saddle points (i.e., not local extrema). They map to parameters
 88 $\beta_c = u_c \odot v_c$ which satisfy $|\beta_c| \odot \nabla L(\beta_c) = \mathbf{0}$ and:*

$$\beta_c \in \arg \min_{\beta[i]=0 \text{ for } i \notin \text{supp}(\beta_c)} L(\beta) \quad (4)$$

89 where $\text{supp}(\beta_c) = \{i \in [d], \beta_c[i] \neq 0\}$ corresponds to the support of β_c .

90 The optimisation problem in Eq. (4) states that the saddle points of the train loss F correspond to
 91 **sparse vectors that minimise the loss function L over its non-zero coordinates**. This property
 92 already shows that the saddle points possess interesting properties from a learning perspective. In the
 93 following we loosely use the term of ‘saddle’ to refer to points $\beta_c \in \mathbb{R}^d$ solution of Eq. (4) **that are**
 94 **not saddles of the convex loss function L** . We adopt this terminology because they correspond to
 95 points $w_c \in \mathbb{R}^{2d}$ that are indeed saddles of the non-convex loss F .

96 **Gradient Flow and necessity of “accelerating” time.** We minimise the loss F using gradient flow:

$$dw_t = -\nabla F(w_t) dt, \quad (5)$$

97 initialised at $u_0 = \sqrt{2}\alpha \mathbf{1} \in \mathbb{R}_{>0}^d$ with $\alpha > 0$, and $v_0 = \mathbf{0} \in \mathbb{R}^d$. This initialisation results in
 98 $\beta_0 = \mathbf{0} \in \mathbb{R}^d$ independently of the chosen neuron initialisation scale α . We denote $\beta_t^\alpha := u_t^\alpha \odot v_t^\alpha$
 99 the prediction iterates generated from the gradient flow to highlight its dependency on the initialisation
 100 scale α ¹. The origin $\mathbf{0} \in \mathbb{R}^{2d}$ is a critical point of the function F and taking the initialisation $\alpha \rightarrow 0$
 101 therefore arbitrarily slows down the dynamics. In fact, it can be easily shown for any fixed time t ,
 102 that $(u_t^\alpha, v_t^\alpha) \rightarrow \mathbf{0}$ as $\alpha \rightarrow 0$, indicating that the iterates are stuck at the origin. Therefore if we
 103 restrict ourselves to a finite time analysis, there is no hope of exhibiting the observed saddle-to-saddle
 104 behaviour. To do so, we must find an appropriate bijection \tilde{t}_α in $\mathbb{R}_{\geq 0}$ which “accelerates” time (i.e.
 105 $\tilde{t}_\alpha(t) \xrightarrow{\alpha \rightarrow 0} +\infty$ for all t) and consider the accelerated iterates $\beta_{\tilde{t}_\alpha(t)}^\alpha$ which can escape the saddles.
 106 Finding this bijection becomes very natural once the mirror structure is unveiled.

¹We point out that the trajectory of β_t^α exactly matches that of another common parametrisation $\beta_w := w_+^2 - w_-^2$, with initialisation $w_{+,0} = w_{-,0} = \alpha \mathbf{1}$.

107 **2.2 Leveraging the mirror flow structure**

108 While the iterates $(w_t^\alpha)_t$ follow a gradient flow on the non-convex loss F , it is shown in [4] that the
 109 iterates β_t^α follow a mirror flow on the convex loss L with potential ϕ_α and initialisation $\beta_{t=0}^\alpha = \mathbf{0}$:

$$d\nabla\phi_\alpha(\beta_t^\alpha) = -\nabla L(\beta_t^\alpha)dt, \quad (6)$$

110 where ϕ_α is the hyperbolic entropy function [15] defined as:

$$\phi_\alpha(\beta) = \frac{1}{2} \sum_{i=1}^d \left(\beta_i \operatorname{arcsinh}\left(\frac{\beta_i}{\alpha^2}\right) - \sqrt{\beta_i^2 + \alpha^4 + \alpha^2} \right). \quad (7)$$

111 Unveiling the mirror flow structure enables to leverage convex optimisation tools to prove convergence
 112 of the iterates to a global minimiser β_α^* as well as a simple proof of the implicit regularisation problem
 113 it solves. As shown by Woodworth et al. [39], in the overparametrised setting where $d > n$ and where
 114 there exists an infinite number of global minima, the limit β_α^* is the solution of the problem:

$$\beta_\alpha^* = \arg \min_{y_i = \langle x_i, \beta \rangle, \forall i} \phi_\alpha(\beta). \quad (8)$$

115 Furthermore, a simple function analysis shows that ϕ_α behaves as a rescaled ℓ_1 -norm as α goes
 116 to 0, meaning that the recovered solution β_α^* converges to the minimum ℓ_1 -norm solution $\beta_{\ell_1}^* :=$
 117 $\arg \min_{y_i = \langle x_i, \beta \rangle} \|\beta\|_1$ as α goes to 0. To bring to light the saddle-to-saddle dynamics which occurs
 118 as we take the initialisation to 0, we make substantial use of the nice mirror structure from Eq. (6).

119 **Appropriate time rescaling.** To understand the limiting dynamics of β_t^α , it is natural to consider
 120 the limit $\alpha \rightarrow 0$ in Eq. (6). However, the potential ϕ_α is such that $\phi_\alpha(\beta) \sim \ln(1/\alpha)\|\beta\|_1$ for small
 121 α and therefore degenerates as $\alpha \rightarrow 0$. Similarly, for $\beta \neq \mathbf{0}$, $\|\nabla\phi_\alpha(\beta)\| \rightarrow \infty$ as $\alpha \rightarrow 0$. The
 122 formulation from Eq. (6) is thus not appropriate to take the limit $\alpha \rightarrow 0$. We can nonetheless obtain a
 123 meaningful limit by considering the opportune time acceleration $\tilde{t}_\alpha(t) = \ln(1/\alpha) \cdot t$ and looking at
 124 the accelerated iterates

$$\tilde{\beta}_t^\alpha := \beta_{\tilde{t}_\alpha(t)}^\alpha = \beta_{\ln(1/\alpha)t}^\alpha. \quad (9)$$

125 Indeed, a simple chain rule leads to the ‘‘accelerated mirror flow’’: $d\nabla\phi_\alpha(\tilde{\beta}_t^\alpha) = -\ln(1/\alpha)\nabla L(\tilde{\beta}_t^\alpha)dt$.
 126 The accelerated iterates $(\tilde{\beta}_t^\alpha)_t$ follow a mirror descent with a rescaled potential:

$$d\nabla\tilde{\phi}_\alpha(\tilde{\beta}_t^\alpha) = -\nabla L(\tilde{\beta}_t^\alpha)dt, \quad \text{where} \quad \tilde{\phi}_\alpha := \frac{1}{\ln(1/\alpha)} \cdot \phi_\alpha, \quad (10)$$

127 with $\tilde{\beta}_{t=0} = \mathbf{0}$ and where ϕ_α is defined Eq. (7). Our choice of time acceleration ensures that the
 128 rescaled potential $\tilde{\phi}_\alpha$ is non-degenerate as the initialisation goes to 0 since $\tilde{\phi}_\alpha(\beta) \underset{\alpha \rightarrow 0}{\sim} \|\beta\|_1$.

129 **3 Intuitive construction of the limiting flow and saddle-to-saddle algorithm**

130 In this section, we aim to give a comprehensible construction of the limiting flow. We therefore
 131 choose to provide intuition over pure rigor, and defer the full and rigorous proof to the Appendix E.
 132 The technical crux of our analysis is to demonstrate the existence of a piecewise constant limiting
 133 process towards which the iterates $\tilde{\beta}^\alpha$ converge to. The convergence result is deferred to the following
 134 Section 4. **In this section we assume this convergence and refer to this piecewise constant**
 135 **limiting process as $(\tilde{\beta}_t^\circ)_t$.** Our goal is then to determine the jump times (t_1, \dots, t_p) as well as the
 136 saddles $(\beta_0, \dots, \beta_p)$ which fully define this process.

137 To do so, it is natural to examine the limiting equation obtained when taking the limit $\alpha \rightarrow 0$ in
 138 Eq. (10). We first turn to its integral form which writes:

$$-\int_0^t \nabla L(\tilde{\beta}_s^\alpha)ds = \nabla\tilde{\phi}_\alpha(\tilde{\beta}_t^\alpha). \quad (11)$$

139 Provided the convergence of the flow $\tilde{\beta}^\alpha$ towards $\tilde{\beta}^\circ$, the left hand side of the previous equation
 140 converges to $-\int_0^t \nabla L(\tilde{\beta}_s^\circ)ds$. For the right hand side, recall that $\tilde{\phi}_\alpha(\beta) \underset{\alpha \rightarrow 0}{\sim} \|\beta\|_1$, it is therefore

141 natural to expect the right hand side of Eq. (11) to converge towards an element of $\partial\|\tilde{\beta}_t^\circ\|_1$, where
 142 we recall the definition of the subderivative of the ℓ_1 -norm as:

$$\partial\|\tilde{\beta}\|_1 = \{1\} \text{ if } \tilde{\beta} > 0, \quad \{-1\} \text{ if } \tilde{\beta} < 0, \quad [-1, 1] \text{ if } \tilde{\beta} = 0.$$

143 The arising key equation which must satisfy the limiting process $\tilde{\beta}^\circ$ is then, for all $t \geq 0$:

$$-\int_0^t \nabla L(\tilde{\beta}_s^\circ) ds \in \partial\|\tilde{\beta}_t^\circ\|_1. \quad (12)$$

144 We show that **this equation uniquely determines the piecewise constant process** $\tilde{\beta}^\circ$ by imposing
 145 the number of jumps p , the jump times as well as the saddles which are visited between the jumps.
 146 Indeed the relation described in Eq. (12) provides 4 restrictive properties that enable to construct $\tilde{\beta}^\circ$.
 147 To state them, let $s_t = -\int_0^t \nabla L(\tilde{\beta}_s^\circ) ds$ and notice that it is continuous and piecewise linear since $\tilde{\beta}^\circ$
 148 is piecewise constant. For each coordinate $i \in [d]$, it holds that:

$$\begin{aligned} 149 \quad \text{(K1)} \quad s_t[i] \in [-1, 1] & \quad \text{(K2)} \quad s_t[i] = 1 \Rightarrow \tilde{\beta}_t^\circ[i] \geq 0 \text{ and } s_t[i] = -1 \Rightarrow \tilde{\beta}_t^\circ[i] \leq 0 \\ 150 \quad \text{(K3)} \quad s_t[i] \in (-1, 1) \Rightarrow \tilde{\beta}_t^\circ[i] = 0 & \quad \text{(K4)} \quad \tilde{\beta}_t^\circ[i] > 0 \Rightarrow s_t[i] = 1 \text{ and } \tilde{\beta}_t^\circ[i] < 0 \Rightarrow s_t[i] = -1 \end{aligned}$$

151 To understand how these conditions lead to the algorithm which determines the jump times and the
 152 visited saddles, we present a 2-dimensional example for which we can walk through each step. The
 153 general case then naturally follows from this simple example.

154 3.1 Construction of the saddle-to-saddle algorithm with an illustrative $2d$ example.

155 Let us consider $n = d = 2$ and data matrix $X \in \mathbb{R}^{2 \times 2}$ such that $X^\top X = ((1, 0.2), (0.2, -0.2))$.
 156 We consider $\beta^* = (-0.2, 2) \in \mathbb{R}^2$ and outputs $y = X\beta^*$. This setting is such that the loss L
 157 has β^* as its unique minimum and $L(\beta^*) = 0$. Furthermore the non-convex loss F has 3 saddles
 158 which map to: $\beta_{c,0} := (0, 0) = \arg \min_{\beta_i=0, \forall i} L(\beta)$, $\beta_{c,1} := (0.2, 0) = \arg \min_{\beta_{[2]=0}} L(\beta)$ and
 159 $\beta_{c,2} := (0, 1.6) = \arg \min_{\beta_{[1]=0}} L(\beta)$. The loss function L is sketched in Figure 2 (Left). Notice
 160 that by the definition of $\beta_{c,1}$ and $\beta_{c,2}$, the gradients of the loss at these points are orthogonal to the
 161 axis they belong to. When running gradient flow with a small initialisation over our diagonal linear
 162 network, we obtain the plots illustrated Figure 2 (Middle and Right). We observe three jumps: the
 163 iterates jump from the saddle at the origin to $\beta_{c,1}$ at time t_1 , then to $\beta_{c,2}$ at time t_2 and finally to the
 164 global minimum β^* at time t_3 .

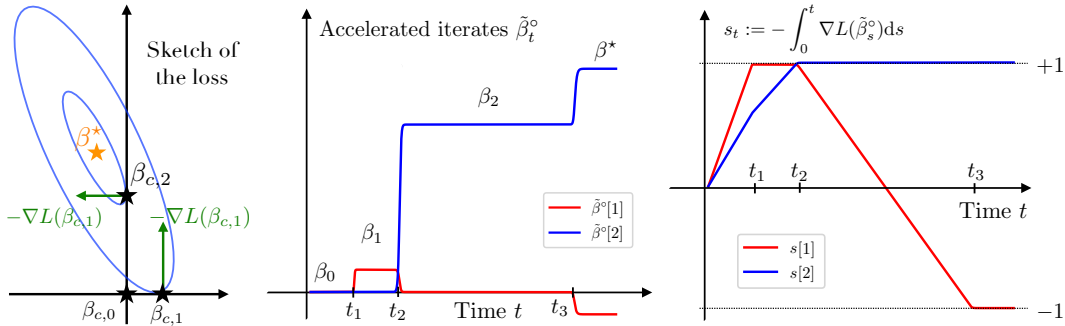


Figure 2: *Left*: Sketch of the $2d$ loss. *Middle and right*: Outputs of gradient flow with small initialisation scale: the iterates are piecewise constant and s_t is piecewise linear across time. We refer to the main text for further details.

165 Let us show how Eq. (12) enables us to theoretically recover this trajectory. A simple observation
 166 which we will use several times below is that for any $t' > t$ such that $\tilde{\beta}^\circ$ is constant equal to β over
 167 the time interval (t, t') , the definition of s enables to write that $s_{t'} = s_t - (t' - t) \cdot \nabla L(\beta)$.

168 **Zeroth saddle:** The iterates are at the saddle at the origin: $\tilde{\beta}_t^\circ = \beta_0 := \beta_{c,0}$ and therefore
 169 $s_t = -t \cdot \nabla L(\beta_0)$. Our key equation Eq. (12) is verified since $s_t = -t \cdot \nabla L(\beta_0) \in \partial\|\beta_0\|_1 = [-1, 1]^d$.
 170 However the iterates cannot stay at the origin after time $t_1 := 1/\|\nabla L(\beta_0)\|_\infty$ which corresponds to
 171 the time at which the first coordinate of s_t hits $+1$: $s_{t_1}[1] = 1$. If the iterates stayed at the origin after
 172 t_1 , (K1) for $i = 1$ would be violated. The iterates must hence jump.

173 **First saddle:** The iterates can only jump to a point different from the origin which maintains
 174 Eq. (12) valid. We denote this point as β_1 . Notice that:

- 175 • $s_{t_1}[2] = -t_1 \cdot \nabla L(\beta_0)[2] \in (-1, 1)$ and since s_t is continuous, we must have $\beta_1[2] = 0$ (K3)
- 176 • $s_{t_1}[1] = 1$ and hence for $t \geq t_1$, $s_t[1] = 1 - (t - t_1)\nabla L(\beta_1)[1]$. We cannot have $\nabla L(\beta_1)[1] <$
 177 0 (K1), and neither $\nabla L(\beta_1)[1] > 0$ since otherwise $s_t[1] \in (-1, 1)$ and $\beta_1 = \mathbf{0}$ (K3)

178 The two conditions $\beta_1[2] = 0$ and $\nabla L(\beta_1)[1] = 0$ **uniquely defines β_1 as equal to $\beta_{c,1}$** . We now
 179 want to know if and when the iterates jump again. We saw that $s_t[1]$ remains at the value +1. However
 180 since β_1 is not a global minimum, $\nabla L(\beta_1)[2] \neq 0$ and $s_t[2]$ hits +1 at time t_2 defined such that
 181 $-(t_1\nabla L(\beta_0) + (t_2 - t_1)\nabla L(\beta_1))[2] = 1$. The iterates must jump otherwise (K1) would break.

182 **The iterates cannot jump to β^* yet!** As the second coordinate of the iterates can activate, one
 183 could expect the iterates to be able to jump to the global minimum. However note that s_t is a
 184 continuous function and that s_{t_2} is equal to the vector (1, 1). If the iterates jumped to the global
 185 minimum, then the first coordinate of the iterates would change sign from +0.2 to -0.2. Due to (K4)
 186 this would lead s_t jumping from +1 to -1, violating its continuity.

187 **Second saddle:** We denote as β_2 the point to which the iterates jump. s_{t_2} is now equal to the vector
 188 (1, 1) and therefore (i) $\beta_2 \geq 0$ (coordinate-wise) from (K2 and K3) and the continuity of s . Since
 189 $s_t = s_{t_2} - (t - t_2)\nabla L(\beta_2)$, we must also have: (ii) $\nabla L(\beta_2) \geq 0$ from (K1) (iii) for $i \in \{1, 2\}$, if
 190 $\beta_2[i] \neq 0$ then $\nabla L(\beta_2)[i] = 0$ from (K4). The three conditions (i), (ii) and (iii) precisely correspond
 191 to the optimality conditions of the following problem:

$$\arg \min_{\beta[1] \geq 0, \beta[2] \geq 0} L(\beta).$$

192 The unique minimiser of this problem is $\beta_{c,2}$, hence $\beta_2 = \beta_{c,2}$, which means that the first coordinate
 193 deactivates. Similar to before, (K1) is valid until the time t_3 at which the first coordinate of $s_t =$
 194 $s_{t_2} - (t - t_2)\nabla L(\beta_2)$ reaches -1 due to the fact that $\nabla L(\beta_2)[1] > 0$.

195 **Global minimum:** We follow the exact same reasoning as for the second saddle. We now have
 196 s_{t_3} equal to the vector (-1, 1) and the iterates must jump to a point β_3 such that (i) $\beta_3[1] \leq 0$,
 197 $\beta_3[2] \geq 0$ (K2 and K3), (ii) $\nabla L(\beta_3)[1] \leq 0$, $\nabla L(\beta_3)[2] \geq 0$ (K1), (iii) for $i \in \{1, 2\}$, if $\beta_3[i] \neq 0$
 198 then $\nabla L(\beta_3)[i] = 0$ (K4). Again, these are the optimality conditions of the following problem:

$$\arg \min_{\beta[1] \leq 0, \beta[2] \geq 0} L(\beta).$$

199 β^* is the unique minimiser of this problem and $\beta_3 = \beta^*$. For $t \geq t_3$ we have $s_t = s_{t_3}$ and Eq. (12) is
 200 satisfied for all following times: the iterates do not have to move anymore.

201 3.2 Presentation of the full saddle-to-saddle algorithm

202 We can now provide the full algorithm (Algorithm 1) which computes the jump times (t_1, \dots, t_p) and
 203 saddles $(\beta_0 = \mathbf{0}, \beta_1, \dots, \beta_p)$ as the values and vectors such that the associated piecewise constant
 204 process satisfies Eq. (12) for all t . This algorithm therefore defines our limiting process $\tilde{\beta}^\circ$.

205 **Algorithm 1 in words.** The algorithm is a concise representation of the steps we followed in the
 206 previous section to construct $\tilde{\beta}^\circ$. We explain each step in words below. Starting from $k = 0$, assume
 207 we enter the loop number k at the saddle β_k computed in the previous loop:

- 208 • The set \mathcal{A}_k contains the set of coordinates "which are unstable": by having a non-zero
 209 derivative, the loss could be decreased by moving along each one of these coordinates and
 210 one of these coordinates will have to activate.
- 211 • The time gap Δ_k corresponds to the time spent at the saddle β_k . It is computed as being the
 212 elapsed time just before (K1) breaks if the coordinates do not jump.
- 213 • We update $t_{k+1} = t_k + \Delta_k$ and $s_{k+1} = s_k - \Delta_k \nabla L(\beta_k)$: t_{k+1} corresponds to the time at
 214 which the iterates leave the saddle β_k and s_{k+1} constrains the signs of the next saddle β_{k+1}
- 215 • The solution β_{k+1} of the constrained minimisation problem is the saddle to which the flow
 216 jumps to at time t_{k+1} . The optimality conditions of this problem are such that Eq. (12) is
 217 maintained for $t \geq t_{k+1}$.

Algorithm 1: Successive saddles and jump times of $\lim_{\alpha \rightarrow 0} \tilde{\beta}^\alpha$

Initialise: $(t, \beta, s) \leftarrow (0, \mathbf{0}, \mathbf{0})$;
while $\nabla L(\beta) \neq \mathbf{0}$ **do**
 $\mathcal{A} \leftarrow \{j \in [d], \nabla L(\beta)(j) \neq 0\}$
 $\Delta \leftarrow \inf \{\delta > 0 \text{ s.t. } \exists i \in \mathcal{A}, s(i) - \delta \nabla L(\beta)(i) = \pm 1\}$
 $(t, s) \leftarrow (t + \Delta, s - \Delta \cdot \nabla L(\beta))$
 $\beta \leftarrow \arg \min L(\beta)$ where $\beta \in \left\{ \beta \in \mathbb{R}^d \text{ s.t. } \begin{array}{l} \beta_i \geq 0 \text{ if } s(i) = +1 \\ \beta_i \leq 0 \text{ if } s(i) = -1 \\ \beta_i = 0 \text{ if } s(i) \in (-1, 1) \end{array} \right\}$
end
Output: Successive values of β and t

218 **Various comments on Algorithm 1.** First we point out that any solution β_c of the constrained min-
219 imisation problem which appears in Algorithm 1 also satisfies $\beta_c = \arg \min_{\beta[i]=0 \text{ for } i \notin \text{supp}(\beta_c)} L(\beta)$
220 as in Eq. (4): the algorithm hence indeed outputs saddles as expected. Up until now we have never
221 checked whether the algorithm's constrained minimisation problem has a unique minimum. This is
222 crucial otherwise the assignment step would be ill-defined. Showing the uniqueness is non-trivial
223 and is guaranteed thanks to the general position Assumption 1 on the data (see Proposition 6 in
224 Appendix D.1). In this same proposition, we also show that the algorithm terminates in at most
225 $\min(2^d, \sum_{k=0}^n \binom{d}{k})$ steps, that the loss strictly decreases at each step and that the final output β_p
226 is the minimum ℓ_1 -norm solution. These last two properties are expected given the fact that the
227 algorithm arises as being the limit process of $\tilde{\beta}^\alpha$ which follows the mirror flow Eq. (10).

228 **Links with the Homotopy algorithm for the Lasso.** Recall that the Lasso problem is formulated as:

$$\beta_\lambda^* = \arg \min_{\beta \in \mathbb{R}^d} L(\beta) + \lambda \|\beta\|_1. \quad (13)$$

229 The optimality condition of Eq. (13) writes $-\nabla L(\beta_\lambda^*) \in \lambda \partial \|\beta_\lambda^*\|_1$. Now notice the similarity
230 with Eq. (12): the two would be equivalent with $\lambda = 1/t$ if the integration on the left hand side
231 of Eq. (12) did not average over the whole trajectory but only on the final iterate, in which case
232 $-\int_0^t \nabla L(\tilde{\beta}_t^\circ) ds = -t \cdot \nabla L(\tilde{\beta}_t^\circ)$. Though the difference is small, the trajectories of our limiting
233 trajectory $\tilde{\beta}^\circ$ and the lasso path $(\beta_\lambda^*)_\lambda$ are quite different: one has jumps, whereas the other is
234 continuous. Nonetheless, the construction of Algorithm 1 shares many similarities with that of the
235 Homotopy algorithm (see, e.g., [37] and references therein) which is used to compute the Lasso
236 path. A notable difference however is the fact that each step of our algorithm depends on the whole
237 trajectory through the vector s , whereas the Homotopy algorithm can be started from any point on
238 the path.
239

240 **Outputs of the algorithm under a RIP assumption on the data.** Unlike previous results on
241 incremental learning, complex behaviours can occur when the feature matrix is ill designed: several
242 coordinates can activate and deactivate at the same time (see Appendix A for various cases). However,
243 if the feature matrix satisfies the $2r$ -restricted isometry property (RIP) [10] and there exists an r -sparse
244 solution β^* , the visited saddles can be easily approximated using Algorithm 1. Specifically, writing
245 $\beta^* = (\beta_1^*, \dots, \beta_r^*, 0, \dots, 0)$ such that w.l.o.g $|\beta_1^*| \geq \dots \geq |\beta_r^*| > 0$, then under properties depending
246 on the RIP constant and β^* , the algorithm terminates in exactly r loops, and outputs jump times and
247 saddles roughly equal to $t_i = 1/|\beta_i^*|$ and $\beta_i = (\beta_1^*, \dots, \beta_i^*, 0, \dots, 0)$ (we refer to Appendix D.2 for
248 the precise characterisation). Therefore, in simple settings, the support of the sparse vector is learnt a
249 coordinate at a time, without any deactivations.

250 4 Convergence of the iterates towards the process defined by Algorithm 1

251 We are now fully equipped to state our main result which formalises the convergence of the accelerated
252 iterates towards the limiting process $\tilde{\beta}^\circ$ which we built in the previous section.

253 **Theorem 2.** Let the saddles $(\beta_0 = \mathbf{0}, \beta_1, \dots, \beta_{p-1}, \beta_p = \beta_{\ell_1}^*)$ and jump times $(t_0 = 0, t_1, \dots, t_p)$
 254 be the outputs of Algorithm 1 and let $(\tilde{\beta}_t^\circ)_t$ be the piecewise constant process defined as follows:

$$(Saddles) \quad \tilde{\beta}_t^\circ = \beta_k \quad \text{for } t \in (t_k, t_{k+1}) \text{ and } 0 \leq k \leq p, \quad t_{p+1} = +\infty.$$

255 The accelerated flow $(\tilde{\beta}_t^\alpha)_t$ defined in Eq. (9) uniformly converges towards the limiting process $(\tilde{\beta}_t^\circ)_t$
 256 on any compact subset of $\mathbb{R}_{\geq 0} \setminus \{t_1, \dots, t_p\}$.

257 **Convergence result.** We recall that from a technical point of view, showing the existence of a
 258 limiting process $\lim_{\alpha \rightarrow 0} \tilde{\beta}^\alpha$ is the toughest part. Theorem 2 provides this existence as well as the
 259 uniform convergence of the accelerated iterates towards $\tilde{\beta}^\circ$ over all closed intervals of \mathbb{R} which do not
 260 contain the jump times. We highlight that this is the strongest type of convergence we could expect
 261 and a uniform convergence over all intervals of the form $[0, T]$ is impossible given that the limiting
 262 process $\tilde{\beta}^\circ$ is discontinuous. In Proposition 2, we give an even stronger result by showing a graph
 263 convergence of the iterates which takes into account the path followed between the jumps. We also
 264 point out that we can easily show the same type of convergence for the neurons $\tilde{w}_t^\alpha := w_{\tilde{t}^\alpha(t)}^\alpha$ using
 265 the bijective mapping which links the neurons w_t and the predictors β_t (see Lemma 1 in Appendix C).

266 **Estimates for the non-accelerated iterates β_t^α .** We point out that our result provides no speed
 267 of convergence of $\tilde{\beta}^\alpha$ towards $\tilde{\beta}^\circ$. We believe that a non-asymptotic result is challenging and
 268 leave it as future work. Note that we experimentally notice that the convergence rate quickly
 269 degrades after each saddle. Nonetheless, we can still write for the non-accelerated iterates that
 270 $\beta_t^\alpha = \tilde{\beta}_{t/\ln(1/\alpha)}^\alpha \sim \tilde{\beta}_{t/\ln(1/\alpha)}^\circ$ as $\alpha \rightarrow 0$. Hence, for α small enough the iterates β_t^α are roughly equal
 271 to 0 until time $t_1 \cdot \ln(1/\alpha)$ and the minimum ℓ_1 -norm interpolator is reached at time $t_p \cdot \ln(1/\alpha)$. **Such**
 272 **a precise estimate of the global convergence time is rather remarkable** and goes beyond classical
 273 Lyapunov analyses which only leads to $L(\beta_t^\alpha) \lesssim \ln(1/\alpha)/t$ (see Proposition 3 in Appendix C).

274 **Natural extensions of our setting.** More general initialisations can easily be dealt with. For instance,
 275 initialisations of the form $u_{t=0} = \alpha \mathbf{u}_0 \in \mathbb{R}^d$ lead to the exact same result as it is shown in [39]
 276 (Discussion after Theorem 1) that the associated mirror still converges to the ℓ_1 -norm. Initialisations
 277 of the form $[u_{t=0}]_i = \alpha^{k_i}$, where $k_i > 0$, lead to the associated potential converging towards a
 278 weighted ℓ_1 -norm and one should modify Algorithm 1 by accordingly weighting $\nabla L(\beta)$ in the
 279 algorithm. Also, deeper linear architectures of the form $\beta_w = w_+^D - w_-^D$ as in [39] do not change
 280 our result as the associated mirror still converges towards the ℓ_1 -norm. Though we only consider
 281 the square loss in the paper, we believe that all our results should hold for any loss of the type
 282 $L(\beta) = \sum_{i=1}^n \ell(y_i, \langle x_i, \beta \rangle)$ where for all $y \in \mathbb{R}$, $\ell(y, \cdot)$ is strictly convex with a unique minimiser at
 283 y . In fact, the only property which cannot directly be adapted from our results is showing the uniform
 284 boundedness of the iterates (see discussion before Proposition 4 in Appendix C).

285 4.1 High level sketch of proof of $\tilde{\beta}^\alpha \rightarrow \tilde{\beta}^\circ$ which leverages an arc-length parametrisation

286 In this section, we give the high level ideas concerning the proof of the convergence $\tilde{\beta}^\alpha \rightarrow \tilde{\beta}^\circ$ given
 287 in Theorem 2. A full and detailed proof can be found in Appendix E. The main difficulty stems
 288 from the non-continuity of the limit process $\tilde{\beta}^\circ$. To circumvent this difficulty, a clever trick which
 289 we borrow to [13, 29] is to “slow-down” time when the jumps occur by considering **an arc-length**
 290 **parametrisation of the path**. We consider the $\mathbb{R}_{\geq 0}$ arclength bijection τ^α and leverage it to define
 291 the ‘appropriately slowed down’ iterates $\hat{\beta}_\tau^\alpha$ as:

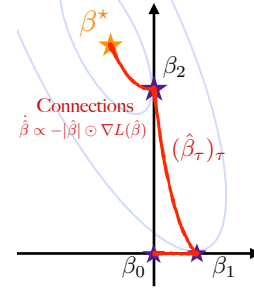
$$\hat{\beta}_\tau^\alpha = \tilde{\beta}_{\hat{t}^\alpha(\tau)}^\alpha \quad \text{where} \quad \hat{t}_\tau^\alpha = (\tau^\alpha)^{-1}(\tau) \quad \text{and} \quad \tau^\alpha(t) = t + \int_0^t \|\dot{\tilde{\beta}}_s^\alpha\| ds.$$

292 This time reparametrisation has the fortunate but crucial property of leading to $\dot{\hat{\beta}}^\alpha(\tau) + \|\dot{\hat{\beta}}_\tau^\alpha\| = 1$
 293 by a simple chain rule, which means that the speed of $(\hat{\beta}_\tau^\alpha)_\tau$ **is uniformly upperbounded by 1**
 294 **independently of α** . This behaviour is in stark contrast with the process $(\tilde{\beta}_t^\alpha)_t$ which has a speed
 295 which explodes at the jumps. This change of time now allows us to use Arzelà-Ascoli’s theorem

296 to extract a subsequence which uniformly converges to a limiting process which we denote $\hat{\beta}$.
 297 Importantly, $\hat{\beta}$ enables to keep track of the path followed between the jumps as we show that its
 298 trajectory has two regimes:

$$\text{Saddles: } \hat{\beta}_\tau = \beta_k \quad \text{Connections: } \dot{\hat{\beta}}_\tau = -\frac{|\hat{\beta}_\tau| \odot \nabla L(\hat{\beta}_\tau)}{\| |\hat{\beta}_\tau| \odot \nabla L(\hat{\beta}_\tau) \|}.$$

The process $\hat{\beta}$ is illustrated on the right: the red curves correspond to the paths which the iterates follow during the jumps. These paths are called *heteroclinic orbits* in the dynamical systems literature [25, 2]. To prove Theorem 2, we can map back the convergence of $\hat{\beta}^\alpha$ to show that of $\tilde{\beta}^\alpha$. Moreover from the convergence $\hat{\beta}^\alpha \rightarrow \hat{\beta}$ we get a more complete picture of the limiting dynamics of $\tilde{\beta}^\alpha$ as it naturally implies the convergence of the graph of the iterates $(\tilde{\beta}_t^\alpha)_t$ converges towards that of $(\hat{\beta}_\tau)_\tau$. The graph convergence result is formalised in this last proposition.



Proposition 2. For all $T > t_p$, the graph of the iterates $(\tilde{\beta}_t^\alpha)_{t \leq T}$ converges to that of $(\hat{\beta}_\tau)_\tau$:

$$\text{dist}(\{\tilde{\beta}_t^\alpha\}_{t \leq T}, \{\hat{\beta}_\tau\}_{\tau \geq 0}) \xrightarrow{\alpha \rightarrow 0} 0 \quad (\text{Hausdorff distance})$$

299 5 Further discussion and conclusion

300 **Link between incremental learning and saddle-to-saddle dynamics.** The incremental learning
 301 phenomenon and the saddle-to-saddle process are often complementary facets of the same idea and
 302 refer to the same phenomenon. Indeed for gradient flows $dw_t = -\nabla F(w_t)dt$, fixed points of the
 303 dynamics correspond to critical points of the loss. Stages with little progress in learning and minimal
 304 movement of the iterates necessarily correspond to the iterates being in the vicinity of a critical
 305 point of the loss. It turns out that in many settings (linear networks [24], matrix sensing [7, 32]),
 306 critical points are necessarily saddle points of the loss (if not global minima) and that they have a
 307 very particular structure (high sparsity, low rank, etc.). We finally note that an alternative approach to
 308 realising saddle-to-saddle dynamics is through the perturbation of the gradient flow by a vanishing
 309 noise as studied in [5].

310 **Characterisation of the visited saddles.** A common belief is that the saddle-to-saddle trajectory
 311 can be found by successively computing the direction of most negative curvature of the loss (i.e.
 312 the eigenvector corresponding to the most negative eigenvalue) and following this direction until
 313 reaching the next saddle [20]. However this statement cannot be accurate as it is inconsistent with
 314 our algorithm in our setting. In fact, it can be shown that this algorithm would match the orthogonal
 315 matching pursuit (OMP) algorithm [33, 12] which does not necessarily lead to the minimum ℓ_1 -norm
 316 interpolator. In the recent work [6], which is the closest to ours, the successive saddles are entirely
 317 characterised for the quadratic parametrisation $\beta = u^2$ but with restrictive assumptions on the data.

318 **Subdifferential equations and rate-independent systems.** As in Eq. (12), subdifferential inclusions
 319 of the form $\nabla L(\beta_t) \in \frac{d}{dt} \partial h(\beta_t)$ for non-differential functions h have been studied by Attouch
 320 et al. [3] but for strongly convex functions h . In this case, the solutions are continuous and do not
 321 exhibit jumps. On another hand, [13, 29, 30] consider so-called *rate-independent systems* of the form
 322 $\partial_q E(t, q_t) \in \partial h(\dot{q}_t)$ for 1-homogeneous *dissipation* potentials h . Examples of such systems are
 323 ubiquitous in mechanics and appear in problems related to friction, crack propagation, elastoplasticity
 324 and ferromagnetism to name a few [28, Ch. 6 for a survey]. As in our case, the main difficulty with
 325 such processes is the possible appearance of jumps when the energy E is non-convex.

326 **Conclusion.** Our study examines the behaviour of gradient flow with vanishing initialisation over
 327 diagonal linear networks. We prove that it leads to the flow jumping from a saddle point of the loss to
 328 another. Our analysis characterises each visited saddle point as well as the jumping times through an
 329 algorithm which is reminiscent of the Homotopy method used in the Lasso framework. There are
 330 several avenues for further exploration. The most compelling one is the extension of these techniques
 331 to broader contexts for which the implicit bias of gradient flow has not yet fully been understood.

References

- 332
- 333 [1] Sanjeev Arora, Nadav Cohen, Wei Hu, and Yuping Luo. Implicit regularization in deep matrix
334 factorization. *Advances in Neural Information Processing Systems*, 32, 2019.
- 335 [2] Peter Ashwin and Michael Field. Heteroclinic networks in coupled cell systems. *Arch. Ration.
336 Mech. Anal.*, 148(2):107–143, 1999.
- 337 [3] H. Attouch, J. Bolte, P. Redont, and M. Teboulle. Singular Riemannian barrier methods and
338 gradient-projection dynamical systems for constrained optimization. *Optimization*, 53(5-6):
339 435–454, 2004.
- 340 [4] Shahar Azulay, Edward Moroshko, Mor Shpigel Nacson, Blake E Woodworth, Nathan Srebro,
341 Amir Globerson, and Daniel Soudry. On the implicit bias of initialization shape: Beyond
342 infinitesimal mirror descent. In *Proceedings of the 38th International Conference on Machine
343 Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 468–477. PMLR,
344 18–24 Jul 2021.
- 345 [5] Yuri Bakhtin. Noisy heteroclinic networks. *Probab. Theory Related Fields*, 150(1-2):1–42,
346 2011.
- 347 [6] Raphaël Berthier. Incremental learning in diagonal linear networks. *arXiv preprint
348 arXiv:2208.14673*, 2022.
- 349 [7] Srinadh Bhojanapalli, Behnam Neyshabur, and Nati Srebro. Global optimality of local search for
350 low rank matrix recovery. In *Advances in Neural Information Processing Systems*, volume 29,
351 2016.
- 352 [8] Etienne Boursier, Loucas Pillaud-Vivien, and Nicolas Flammarion. Gradient flow dynamics of
353 shallow reLU networks for square loss and orthogonal inputs. In Alice H. Oh, Alekh Agarwal,
354 Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing
355 Systems*, 2022.
- 356 [9] Emmanuel J Candes. The restricted isometry property and its implications for compressed
357 sensing. *Comptes rendus mathématique*, 346(9-10):589–592, 2008.
- 358 [10] E. Candès, J. Romberg, and T. Tao. Stable signal recovery from incomplete and inaccurate
359 measurements. *Communications on Pure and Applied Mathematics*, 59(8):1207–1223, 2006.
- 360 [11] Lénaïc Chizat and Francis Bach. Implicit bias of gradient descent for wide two-layer neural
361 networks trained with the logistic loss. In *Proceedings of Thirty Third Conference on Learning
362 Theory*, volume 125 of *Proceedings of Machine Learning Research*, pages 1305–1338. PMLR,
363 09–12 Jul 2020.
- 364 [12] Geoff Davis, Stephane Mallat, and Marco Avellaneda. Adaptive greedy approximations.
365 *Constructive approximation*, 13:57–98, 1997.
- 366 [13] Messoud A. Efendiev and Alexander Mielke. On the rate-independent limit of systems with dry
367 friction and small viscosity. *J. Convex Anal.*, 13(1):151–167, 2006.
- 368 [14] Mathieu Even, Scott Pehme, Suriya Gunasekar, and Nicolas Flammarion. (s)gd over diagonal
369 linear networks: Implicit regularisation, large stepsizes and edge of stability. *arXiv preprint
370 arXiv:2302.08982*, 2023.
- 371 [15] Udaya Ghai, Elad Hazan, and Yoram Singer. Exponentiated gradient meets gradient descent. In
372 Aryeh Kontorovich and Gergely Neu, editors, *Proceedings of the 31st International Conference
373 on Algorithmic Learning Theory*, volume 117 of *Proceedings of Machine Learning Research*,
374 pages 386–407. PMLR, 08 Feb–11 Feb 2020.
- 375 [16] Gauthier Gidel, Francis Bach, and Simon Lacoste-Julien. Implicit regularization of discrete
376 gradient dynamics in linear neural networks. *Advances in Neural Information Processing
377 Systems*, 32, 2019.

- 378 [17] Daniel Gissin, Shai Shalev-Shwartz, and Amit Daniely. The implicit bias of depth: How incre-
379 mental learning drives generalization. In *International Conference on Learning Representations*,
380 2020.
- 381 [18] Suriya Gunasekar, Blake E Woodworth, Srinadh Bhojanapalli, Behnam Neyshabur, and Nati
382 Srebro. Implicit regularization in matrix factorization. In *Advances in Neural Information*
383 *Processing Systems*, volume 30, 2017.
- 384 [19] Jeff Z HaoChen, Colin Wei, Jason Lee, and Tengyu Ma. Shape matters: Understanding the
385 implicit bias of the noise covariance. In *Conference on Learning Theory*, pages 2315–2357.
386 PMLR, 2021.
- 387 [20] Arthur Jacot, François Ged, Berfin Şimşek, Clément Hongler, and Franck Gabriel. Saddle-to-
388 saddle dynamics in deep linear networks: Small initialization training, symmetry, and sparsity.
389 *arXiv preprint arXiv:2106.15933*, 2021.
- 390 [21] Liwei Jiang, Yudong Chen, and Lijun Ding. Algorithmic regularization in model-free over-
391 parametrized asymmetric matrix factorization. *arXiv preprint arXiv:2203.02839*, 2022.
- 392 [22] Jikai Jin, Zhiyuan Li, Kaifeng Lyu, Simon S Du, and Jason D Lee. Understanding incremen-
393 tal learning of gradient descent: A fine-grained analysis of matrix sensing. *arXiv preprint*
394 *arXiv:2301.11500*, 2023.
- 395 [23] Dimitris Kalimeris, Gal Kaplun, Preetum Nakkiran, Benjamin Edelman, Tristan Yang, Boaz
396 Barak, and Haofeng Zhang. Sgd on neural networks learns functions of increasing complexity.
397 *Advances in neural information processing systems*, 32, 2019.
- 398 [24] Kenji Kawaguchi. Deep learning without poor local minima. In *Advances in Neural Information*
399 *Processing Systems*, volume 29, 2016.
- 400 [25] M. Krupa. Robust heteroclinic cycles. *J. Nonlinear Sci.*, 7(2):129–176, 1997.
- 401 [26] Krzysztof Kurdyka. On gradients of functions definable in o-minimal structures. *Ann. Inst.*
402 *Fourier (Grenoble)*, 48(3):769–783, 1998.
- 403 [27] Zhiyuan Li, Yuping Luo, and Kaifeng Lyu. Towards resolving the implicit bias of gradient
404 descent for matrix factorization: Greedy low-rank learning. In *International Conference on*
405 *Learning Representations*, 2021.
- 406 [28] Alexander Mielke. Evolution of rate-independent systems. *Evolutionary equations*, 2:461–559,
407 2005.
- 408 [29] Alexander Mielke, Riccarda Rossi, and Giuseppe Savaré. Modeling solutions with jumps for
409 rate-independent systems on metric spaces. *Discrete Contin. Dyn. Syst.*, 25(2):585–615, 2009.
- 410 [30] Alexander Mielke, Riccarda Rossi, and Giuseppe Savaré. Variational convergence of gradient
411 flows and rate-independent evolutions in metric spaces. *Milan Journal of Mathematics*, 80:
412 381–410, 2012.
- 413 [31] Behnam Neyshabur. Implicit regularization in deep learning. *arXiv preprint arXiv:1709.01953*,
414 2017.
- 415 [32] Dohyung Park, Anastasios Kyrillidis, Constantine Carmanis, and Sujay Sanghavi. Non-square
416 matrix sensing without spurious local minima via the Burer-Monteiro approach. In *Proceedings*
417 *of the 20th International Conference on Artificial Intelligence and Statistics*, volume 54 of
418 *Proceedings of Machine Learning Research*, pages 65–74. PMLR, 20–22 Apr 2017.
- 419 [33] Yagyensh Chandra Pati, Ramin Rezaifar, and Perinkulam Sambamurthy Krishnaprasad. Or-
420 thogonal matching pursuit: Recursive function approximation with applications to wavelet
421 decomposition. In *Proceedings of 27th Asilomar conference on signals, systems and computers*,
422 pages 40–44. IEEE, 1993.
- 423 [34] Scott Pesme, Loucas Pillaud-Vivien, and Nicolas Flammarion. Implicit bias of sgd for diag-
424 onal linear networks: a provable benefit of stochasticity. In *Advances in Neural Information*
425 *Processing Systems*, 2021.

- 426 [35] Noam Razin, Asaf Maman, and Nadav Cohen. Implicit regularization in tensor factorization. In
427 *International Conference on Machine Learning*, pages 8913–8924. PMLR, 2021.
- 428 [36] Andrew M Saxe, James L McClelland, and Surya Ganguli. A mathematical theory of semantic
429 development in deep neural networks. *Proceedings of the National Academy of Sciences*, 116
430 (23):11537–11546, 2019.
- 431 [37] Ryan J. Tibshirani. The lasso problem and uniqueness. *Electron. J. Stat.*, 7:1456–1490, 2013.
- 432 [38] Tomas Vaskevicius, Varun Kanade, and Patrick Rebeschini. Implicit regularization for optimal
433 sparse recovery. *Advances in Neural Information Processing Systems*, 32, 2019.
- 434 [39] Blake Woodworth, Suriya Gunasekar, Jason D. Lee, Edward Moroshko, Pedro Savarese, Itay
435 Golan, Daniel Soudry, and Nathan Srebro. Kernel and rich regimes in overparametrized models.
436 In *Proceedings of Thirty Third Conference on Learning Theory*, volume 125 of *Proceedings of*
437 *Machine Learning Research*, pages 3635–3673. PMLR, 09–12 Jul 2020.
- 438 [40] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding
439 deep learning requires rethinking generalization. In *International Conference on Learning*
440 *Representations*, 2017.

441 **Organisation of the Appendix.**

- 442 1. In Appendix [A](#), we give the experimental setup and provide additional experiments.
- 443 2. In Appendix [B](#), we prove Proposition [1](#) and provide additional comments concerning the
444 unicity of the minimisation problem which appears in the proposition.
- 445 3. In Appendix [C](#), we provide some general results on the flow.
- 446 4. In Appendix [D](#), we give standalone properties of Algorithm [1](#).
- 447 5. In Appendix [E](#), we explain in more detail the arc-length parametrisation explained in the
448 main text as well as prove Theorem [2](#) and Proposition [2](#).
- 449 6. In Appendix [F](#), we provide technical lemmas which are useful to prove the main results.

450 **A Experimental setup and additional: experiments, extension, related works.**

451 **Experimental setup and additional experiments.** For each experiment we generate our dataset
 452 as $y_i = \langle x_i, \beta^* \rangle$ where $x_i = \mathcal{N}(\mathbf{0}, H)$ for a diagonal covariance matrix H and β^* is a vector of
 453 \mathbb{R}^d . Gradient descent is run with a small step size and from initialisation $u_{t=0} = \sqrt{2}\alpha \mathbf{1} \in \mathbb{R}^d$ and
 454 $v_{t=0} = \mathbf{0}$ for some initialisation scale $\alpha > 0$.

- 455 • Figure 1 and Figure 4 (Left): $(n, d, \alpha) = (5, 7, 10^{-120})$, $H = I_d$, $\beta^* =$
 456 $(10, 20, 0, 0, 0, 0, 0) \in \mathbb{R}^7$.
- 457 • Figure 4 (Right): $(n, d, \alpha) = (6, 6, 10^{-10})$, $H = \text{diag}(1, 10, 10, 10, 10, 10) \in \mathbb{R}^{6 \times 6}$,
 458 $\beta^* = (1, 0, 0, 0, 0, 0) \in \mathbb{R}^6$.
- 459 • Figure 3 (Left): $(n, d, \alpha_1, \alpha_2) = (7, 2, 10^{-100}, 10^{-10})$, $H = I_d$, $\beta^* = (10, 20) \in \mathbb{R}^7$.
- 460 • Figure 3 (Right): $(n, d, \alpha) = (3, 3, 10^{-100})$, X is the square root matrix of the matrix
 461 $((20, 6, -1.4), (6, 2, -0.4), (-1.4, -0.4, 0.12)) \in \mathbb{R}^{3 \times 3}$, $\beta^* = (1, 9, 10)$.

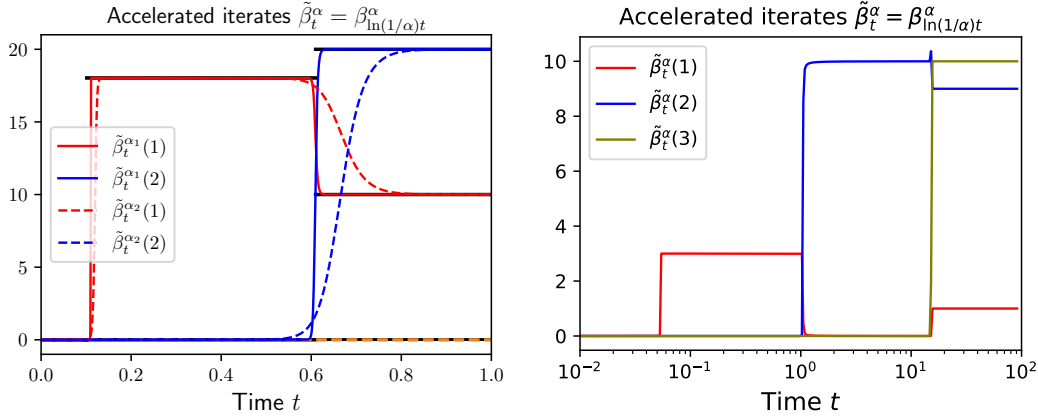


Figure 3: *Left:* Visualisation of the uniform convergence of $\tilde{\beta}^\alpha$ towards $\tilde{\beta}^\circ$ as $\alpha \rightarrow 0$. $\alpha_1 = 10^{-100} \ll \alpha_2 = 10^{-10}$ *Right:* In some cases, 2 coordinates can activate at the same time. Note that the time axis is in log-scale for better visualisation.

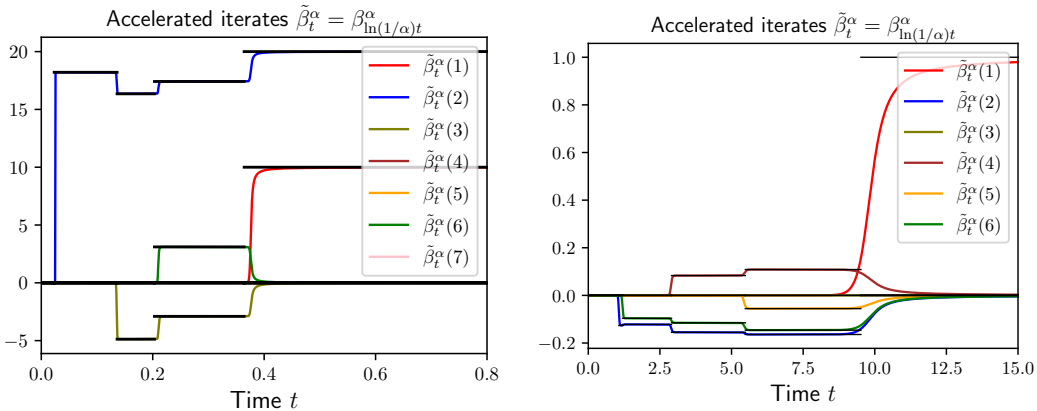


Figure 4: Complex dynamics can occur. *Left and right:* Coordinates are not monotonic and the number of active coordinates neither as several coordinates can deactivate at the same time. The piecewise constant process plotted in black is the limiting process $\tilde{\beta}^\circ$ predicted by our theory.

462 **B Proof of Proposition 1**

463 **Proposition 1.** *All the critical points w_c of F which are not global minima, i.e., $\nabla F(w_c) = \mathbf{0}$ and*
 464 *$F(w_c) > \min_w F(w)$, are necessarily saddle points (i.e., not local extrema). They map to parameters*
 465 *$\beta_c = u_c \odot v_c$ which satisfy $|\beta_c| \odot \nabla L(\beta_c) = \mathbf{0}$ and:*

$$\beta_c \in \arg \min_{\beta[i]=0 \text{ for } i \notin \text{supp}(\beta_c)} L(\beta) \quad (4)$$

466 where $\text{supp}(\beta_c) = \{i \in [d], \beta_c[i] \neq 0\}$ corresponds to the support of β_c .

467 *Proof. **Non-existence of maxima / non-global minima.** This is a simpler version of results which*
 468 *appear in [24], for the sake of completeness we provide here a simple proof adapted to our setting.*
 469 *The intuition follows the fact that if there existed a local maximum / non-global minimum for F then*
 470 *this would translate to the existence of a local maximum / non-global minimum for the convex loss L ,*
 471 *which is absurd.*

472 *Assume that there exists a local maximum $w^* = (u^*, v^*)$, i.e. assume that there exists $\varepsilon > 0$ such*
 473 *that for all $w = (u, v)$ such that $\|w - w^*\|_2^2 \leq \varepsilon$, $F(w) \leq F(w^*)$. We show that this would imply*
 474 *that $\beta^* = u^* \odot v^*$ is a local maximum of L , which is absurd.*

475 *The mapping $g : (u, v) \mapsto (u \odot v, \sqrt{(u^2 - v^2)/2})$ from $\mathbb{R}_{\geq 0}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^d \times \mathbb{R}_{\geq 0}^d$ is a bijection with*
 476 *inverse*

$$g^{-1} : (\beta, \alpha) \mapsto (\sqrt{\alpha^2 + \sqrt{\beta^2 + \alpha^4}}, \text{sign}(\beta) \odot \sqrt{-\alpha^2 + \sqrt{\beta^2 + \alpha^4}}). \quad (14)$$

477 *Also notice that $F(g^{-1}(\beta, \alpha)) = L(\beta)$ for all β and α . Now let $\tilde{\varepsilon} > 0$ and let $\beta \in \mathbb{R}^d$ such that*
 478 *$\|\beta - \beta^*\|_2^2 \leq \tilde{\varepsilon}$, then for $(u, v) = g^{-1}(\beta, \alpha_*)$ where $\alpha_* = \sqrt{((u^*)^2 - (v^*)^2)/2}$ we have that:*

$$\begin{aligned} \|(u, v) - (u^*, v^*)\|_2^2 &= 2 \left\| \left(\sqrt{\alpha_*^4 + \beta^2} - \sqrt{\alpha_*^4 + \beta^{*2}} \right)^2 \right\|_1 \\ &\leq 2 \|\beta^2 - \beta^{*2}\|_1 \\ &= 2 \|(\beta - \beta^*)^2 + 2(\beta - \beta^*)\beta^*\|_1 \\ &\leq 2 \|(\beta - \beta^*)^2\|_1 + 2 \|\beta^*\|_\infty \|\beta - \beta^*\|_1 \\ &\leq 2(1 + \sqrt{d}) \|\beta^*\|_\infty \tilde{\varepsilon} \\ &\leq \varepsilon \end{aligned}$$

479 *where the last inequality is for $\tilde{\varepsilon}$ small enough. This means that $L(\beta) = F(w) \leq F(w^*) = L(\beta^*)$*
 480 *and β^* is a local maximum of L , which is absurd.*

481 *The exact same proof holds to show that there are no local minima of F which are not global minima.*

482 **Critical points.** *The gradient of the loss function F writes:*

$$\nabla_w F(w) = \begin{pmatrix} \nabla_u F(w) \\ \nabla_v F(w) \end{pmatrix} = \begin{pmatrix} \nabla L(\beta) \odot v \\ \nabla L(\beta) \odot u \end{pmatrix} \in \mathbb{R}^{2d}.$$

483 *Therefore $\nabla F(w_c) = \mathbf{0} \in \mathbb{R}^{2d}$ implies that $\nabla L(\beta_c) \odot \beta_c = \mathbf{0} \in \mathbb{R}^d$. Now consider such a β_c and*
 484 *let $\text{supp}(\beta_c) = \{i \in [d] \text{ such that } \beta_c(i) \neq 0\}$ denote the support of β_c . Since $[\nabla L(\beta_c)]_i = 0$ for*
 485 *$i \notin \text{supp}(\beta_c)$, we can therefore write that*

$$\beta_c \in \arg \min_{\beta_i=0 \text{ for } i \notin \text{supp}(\beta_c)} L(\beta).$$

486 *Furthermore we point out that since $\text{supp}(\beta_c) \subset [d]$, there are at most 2^d distinct sets $\text{supp}(\beta_c)$, and*
 487 *therefore at most 2^d values $F(w_c) = L(\beta_c)$, where w_c is a critical point of F . \square*

488 **Additional comment concerning the uniqueness of $\arg \min_{\beta_i=0, i \notin \text{supp}(\beta_c)} L(\beta)$.**

489 We point out that the constrained minimisation problem (4) does not necessarily have a unique
490 solution, even when β_c is not a global solution. Though not required for any of our results, for the
491 sake of completeness, we show here that under an additional mild assumption on the data, we can
492 ensure that the minimisation problem (4) which appears in Proposition 1 has a unique minimum
493 when $L(\beta_c) > 0$. Under this additional assumption, there is therefore a finite number of saddles β_c .
494 Recall that we let $X \in \mathbb{R}^{n \times d}$ be the feature matrix and $(\tilde{x}_1, \dots, \tilde{x}_d)$ be its columns. Now assume
495 temporarily that the following assumption holds.

496 **Assumption 2** (Assumption used just in this short section). *Any subset of $(\tilde{x}_1, \dots, \tilde{x}_d)$ of size smaller
497 than $\min(n, d)$ is linearly independent.*

498 One can easily check that this assumption holds with probability 1 as soon as the data is drawn from
499 a continuous probability distribution, similarly to [37, Lemma 4]). In the following, for a subset
500 $\xi = \{i_1, \dots, i_k\} \subset [d]$, we write $X_\xi = (\tilde{x}_{i_1}, \dots, \tilde{x}_{i_k}) \in \mathbb{R}^{n \times k}$ (we extract the columns from X).
501 For a vector $\beta \in \mathbb{R}^d$ we write $\beta[\xi] = (\beta_{i_1}, \dots, \beta_{i_k})$ and $\beta[\xi^C] = (\beta_i)_{i \notin \xi}$. We distinguish two
502 different settings:

503 • Underparametrised setting ($n \geq d$): in this case, for any $\xi = \{i_1, \dots, i_k\} \subset [d]$, then
504 $\beta^* := \operatorname{argmin}_{\beta_i=0, i \notin \xi} L(\beta)$ is unique. Indeed we simply set the gradient to 0 and notice that
505 due to Assumption 2, there exists a unique solution, indeed it is β^* such that $\beta^*[\xi] =$
506 $(X_\xi^\top X_\xi)^{-1} X_\xi^\top y$ and $\beta^*[\xi^C] = 0$.

507 • Overparametrised setting ($d > n$): **Global solutions:** $\arg \min_{\beta \in \mathbb{R}^d} L(\beta)$ is an affine space
508 spanned by the orthogonal of (x_1, \dots, x_n) in \mathbb{R}^d . Since $\operatorname{span}(\tilde{x}_1, \dots, \tilde{x}_d) = \mathbb{R}^n$ from
509 Assumption 2, any $\beta^* \in \arg \min_{\beta \in \mathbb{R}^d} L(\beta)$ satisfies $X\beta^* = y$ and $L(\beta^*) = 0$. **"Saddle
510 points":** now let $\beta_c \in \mathbb{R}^d$ be such that we can write $\beta_c \in \arg \min_{\beta_i=0, i \notin \operatorname{supp}(\beta_c)} L(\beta)$ and
511 assume that $L(\beta_c) > 0$ (i.e., not a global solution), then: (1) β_c has at most n non-zero
512 entries, indeed if it were not the case, then y would necessarily belong to $\operatorname{span}(\tilde{x}_i)_{i \in \operatorname{supp}(\beta_c)}$
513 due to the assumption on the data, and this would lead to $L(\beta_c) = 0$, (2) therefore, similar
514 to the underparametrised case, $\arg \min_{\beta_i=0, i \notin \operatorname{supp}(\beta_c)} L(\beta)$ is unique, equal to β_c , and we
515 have that $\beta_c[\xi] = (X_\xi^\top X_\xi)^{-1} X_\xi^\top y$ and $\beta_c[\xi^C] = 0$ where $\xi = \operatorname{supp}(\beta_c)$.

516 Thus, in both the underparametrised and overparametrised settings, the minimisation problem (4)
517 appearing in Proposition 1 has a unique minimum when $L(\beta_c) > 0$ and Assumption 2 holds.

518 **C General results on the iterates**

519 In the following lemma we recall a few results concerning the gradient flow Eq. (5):

$$dw_t = -\nabla F(w_t)dt, \quad (15)$$

520 where F is defined in Eq. (3) as:

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

521

522 **Lemma 1.** *For an initialisation $u_0 = \sqrt{2}\alpha$, $v_0 = \mathbf{0}$, the flow $w_t^\alpha = (u_t^\alpha, v_t^\alpha)$ from Eq. (15) is such*
 523 *that the quantity $(u_t^\alpha)^2 - (v_t^\alpha)^2$ is constant and equal to $2\alpha^2\mathbf{1}$. Furthermore $u_t^\alpha > |v_t^\alpha| \geq 0$ and*
 524 *therefore from the bijection Eq. (14) we have that:*

$$u_t^\alpha = \sqrt{\alpha^2 + \sqrt{(\beta_t^\alpha)^2 + \alpha^4}}, \quad v_t^\alpha = \text{sign}(\beta_t^\alpha) \odot \sqrt{-\alpha^2 + \sqrt{(\beta_t^\alpha)^2 + \alpha^4}}.$$

525 *Proof.* From the expression of $\nabla F(w)$, notice that the derivative of $(u_t^\alpha)^2 - (v_t^\alpha)^2$ is equal to $\mathbf{0}$ and
 526 therefore equal to its initial value.

527 Since $(u_t^\alpha)^2 - (v_t^\alpha)^2 = (u_t^\alpha + v_t^\alpha)(u_t^\alpha - v_t^\alpha) > 0$, by continuity we get that $u_t^\alpha + v_t^\alpha > 0$ and
 528 $u_t^\alpha - v_t^\alpha > 0$ and therefore $u_t^\alpha > |v_t^\alpha|$. \square

529 In this section we consider the accelerated iterates Eq. (9) which follow:

$$d\nabla \tilde{\phi}_\alpha(\tilde{\beta}_t^\alpha) = -\nabla L(\tilde{\beta}_t^\alpha)dt, \quad \text{where} \quad \tilde{\phi}_\alpha := \frac{1}{\ln(1/\alpha)} \cdot \tilde{\phi}_\alpha \quad (16)$$

530 with $\tilde{\beta}_{t=0} = \mathbf{0}$ and where ϕ_α is defined Eq. (7).

531 **Proposition 3.** *For all $\alpha > 0$ and minimum $\beta^* \in \arg \min_\beta L(\beta)$, the loss values $L(\tilde{\beta}_t^\alpha)$ and the*
 532 *Bregman divergence $D_{\tilde{\phi}_\alpha}(\beta^*, \tilde{\beta}_t^\alpha)$ are decreasing. Moreover*

$$L(\tilde{\beta}_t^\alpha) - L(\beta^*) \leq \frac{\tilde{\phi}_\alpha(\beta^*)}{2t}, \quad (17)$$

$$L\left(\frac{1}{t} \int_0^t \tilde{\beta}_s^\alpha ds\right) - L(\beta^*) \leq \frac{\tilde{\phi}_\alpha(\beta^*)}{2t}. \quad (18)$$

533 *Proof.* The loss is decreasing since: $\frac{d}{dt} L(\tilde{\beta}_t^\alpha) = \nabla L(\tilde{\beta}_t^\alpha)^\top \dot{\tilde{\beta}}_t^\alpha = -\dot{\tilde{\beta}}_t^\alpha^\top \nabla^2 \tilde{\phi}_\alpha(\tilde{\beta}_t^\alpha) \dot{\tilde{\beta}}_t^\alpha \leq 0$.

534 $\frac{d}{dt} D_{\tilde{\phi}_\alpha}(\beta^*, \tilde{\beta}_t^\alpha) = -\nabla L(\tilde{\beta}_t^\alpha)^\top (\tilde{\beta}_t^\alpha - \beta^*) = -2(L(\tilde{\beta}_t^\alpha) - L(\beta^*))$ (since L is the quadratic loss),
 535 therefore the Bregman distance is decreasing. We can also integrate this last equality from 0 to t , and
 536 divide by $-2t$:

$$\begin{aligned} \frac{1}{t} \int_0^t L(\tilde{\beta}_s^\alpha) ds - L(\beta^*) &= \frac{D_{\tilde{\phi}_\alpha}(\beta^*, \beta_0^\alpha = \mathbf{0}) - D_{\tilde{\phi}_\alpha}(\beta^*, \tilde{\beta}_t^\alpha)}{2t} \\ &\leq \frac{\tilde{\phi}_\alpha(\beta^*)}{2t}. \end{aligned}$$

537 Since the loss is decreasing we get that $L(\tilde{\beta}_t^\alpha) - L(\beta^*) \leq \frac{\tilde{\phi}_\alpha(\beta^*)}{2t}$ and from the convexity of L we
 538 get that $L\left(\frac{1}{t} \int_0^t \tilde{\beta}_s^\alpha ds\right) - L(\beta^*) \leq \frac{\tilde{\phi}_\alpha(\beta^*)}{2t}$. \square

539 In the following proposition, we show that for α small enough, the iterates are bounded independently
 540 of α . Note that this result unfortunately only holds for the quadratic loss, we expect it to hold for
 541 other convex losses of the type $L(\beta) = \frac{1}{n} \sum_i \ell(y_i, \langle x_i, \beta \rangle)$ where $\ell(y, \cdot)$ is strictly convex has a
 542 unique root at y but we don't know how to show it. Also note that bounding the accelerated iterates
 543 $\tilde{\beta}_t^\alpha$ is equivalent to bounding the iterates β_t^α since $\tilde{\beta}_t^\alpha = \beta_{\ln(1/\alpha)t}^\alpha$.

544 **Proposition 4.** For $\alpha < \alpha_0$, where α_0 depends on $\beta_{\ell_1}^*$, the iterates $\tilde{\beta}_t^\alpha$ are bounded independently of
 545 α :

$$\|\tilde{\beta}_t^\alpha\|_\infty \leq 3\|\beta_{\ell_1}^*\|_1 + 1$$

546 *Proof.* From Eq. (16), integrating and using that L is the quadratic loss, we get:

$$\nabla \tilde{\phi}_\alpha(\tilde{\beta}_t^\alpha) = \frac{t}{n} X^\top (y - X \tilde{\beta}_t^\alpha) = -\frac{t}{n} X^\top X (\tilde{\beta}_t^\alpha - \beta^*),$$

547 where we recall that $X \in \mathbb{R}^{n \times d}$ is the input data represented as a matrix and where we denote the
 548 averaged iterate by $\bar{\beta}_t^\alpha = \frac{1}{t} \int_0^t \tilde{\beta}_s^\alpha ds$. Thus we get

$$\nabla \tilde{\phi}_\alpha(\tilde{\beta}_t^\alpha)^\top (\tilde{\beta}_t^\alpha - \beta^*) = -\frac{t}{n} (\tilde{\beta}_t^\alpha - \beta^*)^\top X^\top X (\tilde{\beta}_t^\alpha - \beta^*). \quad (19)$$

549 By convexity of $\tilde{\phi}_\alpha$ we have $\tilde{\phi}_\alpha(\beta_t^\alpha) - \tilde{\phi}_\alpha(\beta^*) \leq \nabla \tilde{\phi}_\alpha(\beta_t^\alpha)^\top (\beta_t^\alpha - \beta^*)$. By the Cauchy-Schwarz
 550 inequality, we also have $(\tilde{\beta}_t^\alpha - \beta^*)^\top X^\top X (\tilde{\beta}_t^\alpha - \beta^*) \leq \|X(\tilde{\beta}_t^\alpha - \beta^*)\| \|X(\tilde{\beta}_t^\alpha - \beta^*)\|$. Using
 551 Proposition 3: $\|X(\beta_t^\alpha - \beta^*)\|^2 \leq n\tilde{\phi}_\alpha(\beta^*)/t$ and $\|X(\tilde{\beta}_t^\alpha - \beta^*)\|^2 \leq n\tilde{\phi}_\alpha(\beta^*)/t$ we can further
 552 bound the right hand side of Eq. (19) as

$$-\frac{t}{n} (\tilde{\beta}_t^\alpha - \beta^*)^\top X^\top X (\tilde{\beta}_t^\alpha - \beta^*) \leq \tilde{\phi}_\alpha(\beta^*).$$

553 Thus it yields

$$\tilde{\phi}_\alpha(\beta_t^\alpha) - \tilde{\phi}_\alpha(\beta^*) \leq \tilde{\phi}_\alpha(\beta^*).$$

From [39] (proof of Lemma 1 in the appendix) we get that for

$$\alpha < \min \left\{ 1, \sqrt{\|\beta\|_1}, (2\|\beta\|_1)^{-1} \right\}$$

554 then:

$$\tilde{\phi}_\alpha(\beta) \leq \frac{3}{2} \|\beta\|_1,$$

555 and for all $\alpha < \exp(-d/2)$:

$$\begin{aligned} \tilde{\phi}_\alpha(\beta) &\geq \|\beta\|_1 - \frac{d}{\ln(1/\alpha^2)} \\ &\geq \|\beta\|_1 - 1, \end{aligned}$$

which finally leads for

$$\alpha < \alpha_0 := \min \left\{ 1, \sqrt{\|\beta_{\ell_1}^*\|_1}, (2\|\beta_{\ell_1}^*\|_1)^{-1}, \exp(-d/2) \right\}$$

556 to the result. □

557 The following proposition shows that we can bound the path length of the flow $\tilde{\beta}^\alpha$ independently of α .
 558 Keep in mind that the path length of $\tilde{\beta}^\alpha$ is equivalent to that of β^α as the first is just an acceleration
 559 of the second: $\tilde{\beta}_t^\alpha = \beta_{\ln(1/\alpha)t}^\alpha$.

Proposition 5. For $\alpha < \alpha_0$ where α_0 is the same as in Proposition 4, the path length of the iterates
 $(\beta_t^\alpha)_{t \geq 0}$ is bounded independently of $\alpha > 0$:

$$\int_0^{+\infty} \|\dot{\beta}_t^\alpha\| dt < C,$$

560 where C does not depend on α . Hence the path length of the accelerated flow $\tilde{\beta}^\alpha$ is also bounded
 561 independently of α .

562 *Proof.* Having shown that the iterates β_t^α are bounded independently of α , it also implies that the
 563 iterates $w_t = (u_t, v_t)$ are bounded following Lemma 1. Since the loss $w \mapsto F(w)$ is a multivari-
 564 ate polynomial function, it is a semialgebraic function and we can consequently apply the result
 565 of Kurdyka [26, Theorem 2] which grants that

$$\int_0^{+\infty} \|\dot{w}_t\| dt < C,$$

566 where the constant C only depends on the loss and on the bound on the iterates. We further use
 567 that $\dot{\beta} = \dot{u} \odot v + u \odot \dot{v}$ and $\|\dot{u} \odot v + u \odot \dot{v}\| \leq C_1(\|\dot{u}\| + \|\dot{v}\|)$ using that u and v are bounded and
 568 $\|\dot{u}\| + \|\dot{v}\| \leq C_2\|\dot{w}\|$ using the equivalence of norms. Therefore $\int_0^{+\infty} \|\dot{\beta}_t^\alpha\| dt < C$ for some C which
 569 is independent of the initialisation scale α . \square

570 D Standalone properties of Algorithm 1

571 D.1 “Well-definedness” of Algorithm 1 and upperbound on its number of loops

572 Notice that this proposition highlights the fact that Algorithm 1 is on its own an algorithm of interest
 573 for finding the minimum ℓ_1 -norm solution in an overparametrised regression setting. We point out that
 574 the provided upperbound on the number of iterations is very crude and could certainly be improved.

575 **Proposition 6.** *Algorithm 1 is well defined: at each iteration (i) the attribution of Δ is well defined*
 576 *as $\Delta < +\infty$, (ii) the constrained minimisation problem has a unique solution and the attribution*
 577 *of the value of β is therefore well-founded. Furthermore, along the loops: the iterates β have at*
 578 *most n non-zero coordinates, the loss is strictly decreasing and the algorithm terminates in at most*
 579 *$\min(2^d, \sum_{k=0}^n \binom{d}{k})$ steps by outputting the minimum ℓ_1 -norm solution $\beta_{\ell_1}^* := \arg \min_{\beta \in \arg \min L} \|\beta\|_1$.*

580 *Proof.* In the following, for the matrix X and for a subset $I = \{i_1, \dots, i_k\} \subset [d]$, we write
 581 $X_I = (\tilde{x}_{i_1}, \dots, \tilde{x}_{i_k}) \in \mathbb{R}^{n \times k}$ (we extract the columns from X). For a vector $\beta \in \mathbb{R}^d$ we write
 582 $\beta_I = (\beta_{i_1}, \dots, \beta_{i_k})$.

583 **(1) The constrained minimisation problem has a unique solution:** we follow the proof of [37,
 584 Lemma 2]. Following the notations in Algorithm 1, we define $I = \{i \in [d], |s_i| = 1\}$ and we
 585 point out that after k loops of the algorithm, the value of s is equal to $s = -(\Delta_1 \nabla L(\beta_0) + \dots +$
 586 $\Delta_k \nabla L(\beta_{k-1})) \in \text{span}(x_1, \dots, x_n)$. We can therefore write $s = X^\top r$ for some $r \in \mathbb{R}^n$.

587 Now assume that $\ker(X_I) \neq \{0\}$. Then, for some $i \in I$, we have $\tilde{x}_i = \sum_{j \in I \setminus \{i\}} c_j \tilde{x}_j$ where $c_j \in \mathbb{R}$.
 588 Without loss of generality, we can assume that $I \setminus \{i\}$ has at most n elements. Indeed, we can
 589 otherwise always find n elements $\tilde{I} \subset I \setminus \{i\}$ such that $\tilde{x}_i = \sum_{j \in \tilde{I}} c_j \tilde{x}_j$. Rewriting the previous
 590 equality, we get

$$s_i \tilde{x}_i = \sum_{j \in I \setminus \{i\}} (s_i s_j c_j) (s_j \tilde{x}_j). \quad (20)$$

591 Now by definitions of the set I and of r , we have that $\langle \tilde{x}_j, r \rangle = s_j \in \{+1, -1\}$ for any $j \in I$. Taking
 592 the inner product of Eq. (20) with r , we obtain that $1 = \sum_{j \in I \setminus \{i\}} (s_i s_j c_j)$. Consequently, we have
 593 shown that if $\ker(X_I) \neq \{0\}$, then we necessarily have for some $i \in I$,

$$s_i \tilde{x}_i = \sum_{j \in I \setminus \{i\}} a_j (s_j \tilde{x}_j),$$

594 with $\sum_{j \in I \setminus \{i\}} a_j = 1$, which means that $s_i \tilde{x}_i$ lies in the affine space generated by $(s_j \tilde{x}_j)_{j \in I \setminus \{i\}}$.
 595 This fact is however impossible due to Assumption 1 (recall that without loss of generality we
 596 have that $I \setminus \{i\}$ has at most n elements, and trivially less that d elements). **Therefore X_I is full**
 597 **rank**, and $\text{Card}(I) \leq n$. Now notice that the constrained minimisation problem corresponds to
 598 $\arg \min_{\substack{\beta_i \geq 0, i \in I_+ \\ \beta_i \leq 0, i \in I_-}} \|y - X_I \beta_I\|_2^2$. Since X_I is full rank, this restricted loss is strictly convex and the
 599 constrained minimisation problem **has a unique minimum**.

600 **(2) $\Delta < +\infty$:** Notice that the optimality conditions of

$$\beta = \arg \min_{\substack{\beta_i \geq 0, i \in I_+ \\ \beta_i \leq 0, i \in I_- \\ \beta_i = 0, i \notin I}} \|y - X_I \beta_I\|_2^2,$$

601 are (i) β satisfies the constraints, (ii) if $i \in I_+$ (resp $i \in I_-$) then $[-\nabla L(\beta)]_i \leq 0$ (resp $[-\nabla L(\beta)]_i \geq$
 602 0) and (iii) if $\beta_i \neq 0$ then $[\nabla L(\beta)]_i = 0$. One can notice that condition (ii) ensures that at each
 603 iteration, for $\delta \leq \Delta_k$, $s_{k-1} - \delta \nabla L(\beta_{k-1}) \in [-1, 1]$ coordinate wise. Also, if $L(\beta_{k-1}) \neq 0$, then a
 604 coordinate of the vector $|s_{k-1} - \delta \nabla L(\beta_{k-1})|$ must necessarily hit 1, this value of δ corresponds to
 605 Δ_k .

606 **(3) The loss is strictly decreasing:** Let $I_{k-1, \pm}$ and $I_{k, \pm}$ be the equicorrelation sets defined in the
 607 algorithm at step $k-1$ and k , and β_{k-1} and β_k the solutions of the constrained minimisation problems.
 608 Also, let i_k be the newly added coordinate which breaks the constraint at step k (which we assume
 609 to be unique for simplicity). Without loss of generality, assume that $s_k(i_k) = +1$. Since the sets

610 $I_{k-1,+} \setminus \{i_k\}$ and $I_{k-1,-} \setminus I_{k,-}$ are (if not empty) only composed of indexes of coordinates
611 of β_{k-1} which are equal to 0, one can notice that β_{k-1} also satisfies the new constraints at step
612 k . Therefore $L(\beta_k) \leq L(\beta_{k-1})$. Now since $[-\nabla L(\beta_{k-1})]_{i_k} > 0$, from the strict convexity of the
613 restricted loss on I_k , this means that $\beta_k(i_k) > 0$ (which also means that newly activated coordinate
614 i_k **must activate**), and therefore $\beta_{k-1} \neq \beta_k$ and $L(\beta_k) < L(\beta_{k-1})$.

(4) The algorithm terminates in at most $\min\left(2^d, \sum_{k=0}^n \binom{d}{k}\right)$ steps: Recall that we showed in part (1) of the proof that at each iteration k of the algorithm, I_k has at most $\min(n, d)$ elements. Since $\text{supp}(\beta_k) \subset I_k$, we have that β_k has at most $\min(n, d)$ non-zero elements, also recall that we always have $\beta_k = \arg \min_{\beta_i=0, i \notin \text{supp}(\beta_k)} L(\beta)$ (we here have unicity of this minimisation problem following part (1) of the proof). There are hence at most

$$\sum_{k=0}^{\min(n,d)} \binom{d}{k} = \min\left(2^d, \sum_{k=0}^n \binom{d}{k}\right)$$

615 such minimisation problems. The loss being strictly decreasing, the algorithm cannot output the same
616 solution β at two different loops, and the algorithm must terminate in at most $\min\left(2^d, \sum_{k=0}^n \binom{d}{k}\right)$
617 iterations by outputting a vector β^* such that $\nabla L(\beta^*) = 0$, i.e. $\beta^* \in \arg \min L(\beta)$.

618 **(5) The algorithm outputs the minimum ℓ_1 -norm solution.** Let β^* be the output of the algorithm
619 after p iterations. Notice that by the definition of the successive sets $I_{k,\pm}$ and of the constraints on the
620 minimisation problem, we have that at each iteration $s_k \in \partial \|\beta_k\|_1$. Therefore $s_p \in \partial \|\beta^*\|_1$. Also,
621 recall from part (1) of the proof that $s_p \in \text{span}(x_1, \dots, x_n)$ which means that there exists $r \in \mathbb{R}^n$
622 such that $s_p = X^\top r$. Putting the two together we get that $X^\top r \in \partial \|\beta^*\|_1$, this condition along with
623 the fact that $L(\beta^*) = \min L(\beta)$ are exactly the KKT conditions of $\arg \min_{\beta \in \arg \min L} \|\beta\|_1$. \square

624 D.2 Output of Algorithm 1 under the restricted isometry property (RIP)

625 As mentioned several times, for general feature matrices X complex behaviours can occur with
626 coordinates deactivating and changing sign several times. Here we show that for simple datasets
627 which have a feature matrix X that satisfy the restricted isometry property (RIP) [10], we can simply
628 determine the jump times and the saddles as a function of the sparse predictor which we seek to
629 recover.

630 The non-realistic but enlightening extreme case of the RIP assumption is to consider that the feature
631 matrix is such that $X^\top X/n = I_d$. In this case, by letting β^* be the unique vector such that
632 $y = \langle x, \beta^* \rangle$ and assuming that $\beta^* = (\beta_1^*, \dots, \beta_r^*, 0, \dots, 0)$ with $|\beta_1^*| > \dots > |\beta_r^*| > 0$, then the
633 loss writes $L(\beta) = \|\beta - \beta^*\|_2^2/2$ and one can easily check that Algorithm 1 would terminate in r
634 loops and output exactly $t_i = \frac{1}{|\beta_i^*|}$ and $\beta_i = (\beta_1^*, \dots, \beta_i^*, 0, \dots, 0)$ for $i \leq r$ (the case where several
635 coordinates of β^* are strictly equal can also be treated: for example if $\beta_1^* = \beta_2^*$ then the first output of
636 the algorithm is directly $\beta_1 = (\beta_1^*, \beta_2^*, 0, \dots, 0)$).

637 We now consider the more realistic RIP setting which is an adaptation of the previous observation.

638 **RIP setting and gap assumption.** We consider a sparse regression where there exists an r -sparse
639 vector β^* such that $y_i = \langle x_i, \beta^* \rangle$. Furthermore we assume that the feature matrix $X \in \mathbb{R}^{n,d}$ satisfies
640 the $2r$ -restricted isometry property with constant $\tilde{\varepsilon} < \sqrt{2} - 1 < 1/2$: for all submatrix X_s where we
641 extract any $s \leq 2r$ columns of X , the matrix $X_s^\top X_s/n$ of size $s \times s$ has all its eigenvalues in the
642 interval $[1 - \tilde{\varepsilon}, 1 + \tilde{\varepsilon}]$. Furthermore we assume that the r -sparse vector β^* has coordinates which
643 have a ‘‘sufficient gap’’. W.l.o.g we write $\beta^* = (\beta_1^*, \dots, \beta_r^*, 0, \dots, 0)$ with $|\beta_1^*| \geq \dots \geq |\beta_r^*| > 0$ and
644 we define $\lambda := \min_{i \in [r]} (\beta_i^* - \beta_{i+1}^*) \geq 0$ which corresponds to the smallest gap between the entries
645 of β^* . We assume that $5\tilde{\varepsilon} \|\beta^*\| < \lambda/2$ (**Gap assumption**) and we let $\varepsilon := 5\tilde{\varepsilon}$.

646 A classic result from compressed sensing (see Candes [9, Theorem 1.2]) is that the $2r$ -restricted
647 isometry property with constant $\sqrt{2} - 1$ ensures that the minimum ℓ_0 -minimisation problem has a
648 unique r -sparse solution which is β^* . Furthermore it ensures that the minimum ℓ_1 -norm solution is
649 unique and is equal to β^* . This means that Algorithm 1 will have β^* as a final output.

650 We can now characterise the outputs of Algorithm 1 when the data satisfies the previous assumptions.

651 **Proposition 7.** *Under the restricted isometry property and the gap assumption stated right above,*
 652 *Algorithm 1 terminates in r -loops and outputs:*

$$\begin{aligned} \beta_1 &= (\beta_1[1], 0, \dots, 0) && \text{with } \beta_1[1] \in [\beta_1^* - \varepsilon\|\beta^*\|, \beta_1^* + \varepsilon\|\beta^*\|] \\ \beta_2 &= (\beta_2[1], \beta_2[2], 0, \dots, 0) && \text{with } \begin{cases} \beta_2[1] \in [\beta_1^* - \varepsilon\|\beta^*\|, \beta_1^* + \varepsilon\|\beta^*\|] \\ \beta_2[2] \in [\beta_2^* - \varepsilon\|\beta^*\|, \beta_2^* + \varepsilon\|\beta^*\|] \end{cases} \\ \vdots & \\ \beta_{r-1} &= (\beta_{r-1}[1], \dots, \beta_{r-1}[r-1], 0, \dots, 0) && \text{with } \beta_{r-1}[i] \in [\beta_i^* - \varepsilon\|\beta^*\|, \beta_i^* + \varepsilon\|\beta^*\|] \\ \beta_r &= \beta^* = (\beta_1^*, \dots, \beta_r^*, 0, \dots, 0) \end{aligned}$$

653 *at times t_1, \dots, t_r such that*

$$t_i \in \left[\frac{1}{\beta_i^* + \varepsilon\|\beta^*\|}, \frac{1}{\beta_i^* - \varepsilon\|\beta^*\|} \right].$$

654 *Proof.* For simplicity we assume that $\beta_i^* > 0$ for all $i \in [r]$, the proof can easily be adapted to the
 655 general case. We first define $\xi := X^\top X/n - I_d$. By the restricted isometry property, for any $k \leq 2r$,
 656 we have that any $k \times k$ square matrix extracted from ξ which we denote ξ_{kk} has its eigenvalues in
 657 $[-\tilde{\varepsilon}, \tilde{\varepsilon}]$. It also means that the eigenvalues of $(I_k + \xi_{kk})^{-1} - I_k$ are in $[\frac{1}{1+\tilde{\varepsilon}} - 1, \frac{1}{1-\tilde{\varepsilon}} - 1] \subset [-2\tilde{\varepsilon}, 2\tilde{\varepsilon}]$.

658 We now proceed by induction with the following induction hypothesis:

- 659 • β_{k-1} has its support on its $(k-1)$ first coordinates with $|\beta_{k-1}[i] - \beta_i^*| \leq 5\tilde{\varepsilon}\|\beta^*\|$ for $i < k$
- 660 • $t_k \in \left[\frac{1}{\beta_k^* + 5\tilde{\varepsilon}\|\beta^*\|}, \frac{1}{\beta_k^* - 5\tilde{\varepsilon}\|\beta^*\|} \right]$ and $s_{t_k}[k] = 1$
- 661 • $s_{t_k}[i] \in [t_k(\beta_i^* - 5\tilde{\varepsilon}\|\beta^*\|), t_k(\beta_i^* + 5\tilde{\varepsilon}\|\beta^*\|)] \subset (-1, 1)$ for $i > k$

662 From the recurrence hypothesis, the output of the algorithm at step k is hence $\beta_k = \arg \min L(\beta)$
 663 under the constraint $\beta[i] \geq 0$ for $i \leq k$ and $\beta[i] = 0$ otherwise. We first search for the solution of the
 664 minimisation problem without the sign constraint and still (abusively) denote it β_k : we will show that
 665 it turns out to satisfy the sign constraint and that it is therefore indeed β_k .

666 In the following, for a vector v , we denote by $v[:k]$ its k first coordinates. Setting the k first
 667 coordinates of the gradient to 0, we get that $[X^\top X(\beta_k - \beta^*)][:k] = \mathbf{0}$, which leads to $(I_k + \xi_{kk})\beta_k[:k]$
 668 $= \beta^*[:k] + [\xi\beta^*][:k]$, which gives:

$$\begin{aligned} \beta_k[:k] &= (I_k + \xi_{kk})^{-1}(\beta^*[:k] + [\xi\beta^*][:k]) \\ &= \beta^*[:k] + [\xi\beta^*][:k] + v_1 \end{aligned}$$

669 where from the bound on the eigenvalues of $(I_k + \xi_{kk})^{-1} - I_k$ and $\|\xi\beta^*\| \leq \tilde{\varepsilon}\|\beta^*\|$:

$$\begin{aligned} \|v_1\| &\leq 2\tilde{\varepsilon}\|\beta^*[:k] + [\xi\beta^*][:k]\| \\ &\leq 2\tilde{\varepsilon}(\|\beta^*\| + \|\xi\beta^*\|) \\ &\leq 2\tilde{\varepsilon}(\|\beta^*\| + \tilde{\varepsilon}\|\beta^*\|) \\ &\leq 4\tilde{\varepsilon}\|\beta^*\|. \end{aligned}$$

670 Therefore

$$\beta_k[:k] = \beta^*[:k] + v_2$$

671 where $v_2 = [\xi\beta^*][:k] + v_1$ hence $\|v_2\|_\infty \leq \|v_2\| \leq 5\tilde{\varepsilon}\|\beta^*\|$. Notice that from the definition of λ
 672 and the fact that $5\tilde{\varepsilon}\|\beta^*\| < \lambda/2$ we have that $\beta_k[:k] \geq 0$ coordinate-wise, hence verifying the sign
 673 constraint. Also note that $\|\beta_k\| \leq \|\beta^*\| + 5\tilde{\varepsilon}\|\beta^*\| \leq 4\|\beta^*\|$.

674 For $t \geq t_k$, $s_t = s_{t_k} - (t - t_k)\nabla L(\beta_k)$, and $[\nabla L(\beta_k)][:k] = 0$ therefore $s_t[:k] = s_{t_k}[:k]$. Now
 675 for $i > k$, $[-\nabla L(\beta_k)]_i = n^{-1}[X^\top X(\beta^* - \beta_k)]_i = \beta_i^* + [\xi(\beta_k - \beta^*)]_i$. Now since $(\beta_k - \beta^*)$ is
 676 r -sparse we have that:

$$\begin{aligned} \|\xi(\beta_k - \beta^*)\|_\infty &\leq \|\xi(\beta_k - \beta^*)\| \\ &\leq \tilde{\varepsilon}\|\beta_k - \beta^*\| \\ &\leq \tilde{\varepsilon}(\|\beta_k\| + \|\beta^*\|) \\ &\leq 5\tilde{\varepsilon}\|\beta^*\| < \lambda/2, \end{aligned} \tag{21}$$

677 Now from the fact that $s_t[i] = s_{t_k}[i] + (t - t_k)\beta_i^* + (t - t_k)[\xi(\beta_k - \beta^*)]_i$ and using the recurrence
678 hypothesis: $s_{t_k}[i] \in [t_k(\beta_i^* - 5\tilde{\varepsilon}\|\beta^*\|), t_k(\beta_i^* + 5\tilde{\varepsilon}\|\beta^*\|)]$, we get (using the bound Eq. (21))
679 that $s_t[i] \in [t(\beta_i^* - 5\tilde{\varepsilon}\|\beta^*\|), t(\beta_i^* + 5\tilde{\varepsilon}\|\beta^*\|)]$. From the “separation assumption” we have that
680 $5\tilde{\varepsilon}\|\beta^*\| < \lambda/2$ and therefore the next coordinate to activate is necessarily the $(k + 1)^{th}$ at time t_{k+1}
681 with $s_{t_{k+1}}[k + 1] = 1$ and:

$$t_{k+1} \in \left[\frac{1}{\beta_{k+1}^* + 5\tilde{\varepsilon}\|\beta^*\|}, \frac{1}{\beta_{k+1}^* - 5\tilde{\varepsilon}\|\beta^*\|} \right].$$

682 This proves the recursion. The algorithm cannot stop before iteration r as β^* is the unique minimiser
683 of L that has at most r non-zero coordinates. But it stops at iteration r as β^* is the unique minimiser
684 of $L(\beta)$ under the constraints $\beta_i \geq 0$ for $i \leq r$ and $\beta_i = 0$ otherwise. \square

685 **E Proof of Theorem 2 and Proposition 2 through the arc-length**
686 **parametrisation**

687 In this section, we explain in more details the arc-length reparametrisation which circumvents the
688 apparition of discontinuous jumps and leads to the proof of Theorem 2. The main difficulty to show
689 the convergence stems from the non-continuity of the limit process $\tilde{\beta}^\circ$. Therefore we cannot expect
690 uniform convergence of $\tilde{\beta}^\alpha$ towards $\tilde{\beta}$ as $\alpha \rightarrow 0$. In addition, $\tilde{\beta}^\circ$ does not provide any insights into
691 the path followed between the jumps.

692 **Arc-length parametrisation.** The high-level idea is to “slow-down” time when the jumps occur. To
693 do so we follow the approach from [13, 29] and we consider an arc-length parametrisation of the
694 path, i.e., we consider τ^α equal to:

$$\tau^\alpha(t) = t + \int_0^t \|\dot{\tilde{\beta}}_s^\alpha\| ds.$$

695 In Proposition 5, we showed that the full path length $\int_0^{+\infty} \|\dot{\tilde{\beta}}_s^\alpha\| ds$ is finite and bounded independently
696 of α . Therefore τ^α is a bijection in $\mathbb{R}_{\geq 0}$. We can then define the following quantities:

$$\hat{t}_\tau^\alpha = (\tau^\alpha)^{-1}(\tau) \quad \text{and} \quad \hat{\beta}_\tau^\alpha = \tilde{\beta}_{\hat{t}_\tau^\alpha}^\alpha.$$

697 By construction, a simple chain rule leads to $\dot{\hat{t}}_\tau^\alpha + \|\dot{\hat{\beta}}_\tau^\alpha\| = 1$, which means that the speed of $(\hat{\beta}_\tau^\alpha)_\tau$
698 is always upperbounded by 1, independently of α . This behaviour is in stark contrast with the process
699 $(\tilde{\beta}_t^\alpha)_t$ which has a speed which explodes at the jumps. It presents a major advantage as we can now
700 use Arzelà-Ascoli’s theorem to extract a converging subsequence. A simple change of variable shows
701 that the new process satisfies the following equations:

$$- \int_0^\tau \dot{\hat{t}}_s^\alpha \nabla L(\hat{\beta}_s^\alpha) ds = \nabla \tilde{\phi}_\alpha(\hat{\beta}_\tau^\alpha) \quad \text{and} \quad \dot{\hat{t}}_\tau^\alpha + \|\dot{\hat{\beta}}_\tau^\alpha\| = 1 \quad (22)$$

702 started from $\hat{\beta}_\tau^\alpha = 0$ and $\hat{t}_0 = 0$. The next proposition states the convergence of the rescaled process,
703 up to a subsequence.

704 **Proposition 8.** *Let $T \geq 0$. For every $\alpha > 0$, let $(\hat{t}^\alpha, \hat{\beta}^\alpha)$ be the solution of Eq. (22). Then, there
705 exists a subsequence $(\hat{t}^{\alpha_k}, \hat{\beta}^{\alpha_k})_{k \in \mathbb{N}}$ and $(\hat{t}, \hat{\beta})$ such that as $\alpha_k \rightarrow 0$:*

$$(\hat{t}^{\alpha_k}, \hat{\beta}^{\alpha_k}) \rightarrow (\hat{t}, \hat{\beta}) \quad \text{in } (C^0([0, T], \mathbb{R} \times \mathbb{R}^d), \|\cdot\|_\infty) \quad (23)$$

$$(\dot{\hat{t}}^{\alpha_k}, \dot{\hat{\beta}}^{\alpha_k}) \rightharpoonup (\dot{\hat{t}}, \dot{\hat{\beta}}) \quad \text{in } L_1[0, T] \quad (24)$$

706 **Limiting dynamics.** *The limits $(\hat{t}, \hat{\beta})$ satisfy:*

$$- \int_0^T \dot{\hat{t}}_s \nabla L(\hat{\beta}_s) ds \in \partial \|\hat{\beta}_T\|_1 \quad \text{and} \quad \dot{\hat{t}}_T + \|\dot{\hat{\beta}}_T\| \leq 1 \quad (25)$$

707 **Heteroclinic orbit.** *In addition, when $\hat{\beta}_T$ is such that $|\hat{\beta}_T| \odot \nabla L(\hat{\beta}_T) \neq 0$, we have*

$$\dot{\hat{\beta}}_T = - \frac{|\hat{\beta}_T| \odot \nabla L(\hat{\beta}_T)}{\| |\hat{\beta}_T| \odot \nabla L(\hat{\beta}_T) \|} \quad \text{and} \quad \dot{\hat{t}}_T = 0. \quad (26)$$

708 *Furthermore, the loss strictly decreases along the heteroclinic orbits and the path length $\int_0^T \|\dot{\hat{\beta}}_\tau\| d\tau$
709 is upperbounded independently of T .*

710 *Proof.* Differentiating Eq. (22) and from the Hessian of $\tilde{\phi}_\alpha$ we get:

$$\begin{aligned} \dot{\hat{\beta}}_\tau^\alpha &= -\dot{\hat{t}}_\tau^\alpha (\nabla^2 \tilde{\phi}_\alpha(\hat{\beta}_\tau^\alpha))^{-1} \nabla L(\hat{\beta}_\tau^\alpha) \\ &= -(1 - \|\dot{\hat{\beta}}_\tau^\alpha\|) (\nabla^2 \tilde{\phi}_\alpha(\hat{\beta}_\tau^\alpha))^{-1} \nabla L(\hat{\beta}_\tau^\alpha). \end{aligned}$$

Therefore taking the norm on the right hand side we obtain that

$$\|\dot{\hat{\beta}}_\tau^\alpha\| = \frac{\|(\nabla^2 \tilde{\phi}_\alpha(\hat{\beta}_\tau^\alpha))^{-1} \nabla L(\hat{\beta}_\tau^\alpha)\|}{1 + \|(\nabla^2 \tilde{\phi}_\alpha(\hat{\beta}_\tau^\alpha))^{-1} \nabla L(\hat{\beta}_\tau^\alpha)\|},$$

711 and therefore

$$\dot{\hat{\beta}}_\tau^\alpha = - \frac{(\nabla^2 \tilde{\phi}_\alpha(\hat{\beta}_\tau^\alpha))^{-1} \nabla L(\hat{\beta}_\tau^\alpha)}{1 + \|(\nabla^2 \tilde{\phi}_\alpha(\hat{\beta}_\tau^\alpha))^{-1} \nabla L(\hat{\beta}_\tau^\alpha)\|}. \quad (27)$$

712 **Subsequence extraction.** By construction Eq. (22) we have $\dot{\hat{t}}_\tau^\alpha + \|\dot{\hat{\beta}}_\tau^\alpha\| = 1$, therefore the sequences
 713 $(\dot{\hat{t}}^\alpha)_\alpha$, $(\dot{\hat{\beta}}^\alpha)_\alpha$ as well as $(\hat{t}^\alpha)_\alpha$, $(\hat{\beta}^\alpha)_\alpha$ are uniformly bounded on $[0, T]$. The Arzelà-Ascoli theorem
 714 yields that, up to a subsequence, there exists $(\hat{t}, \hat{\beta})$ such that $(\hat{t}^{\alpha_k}, \hat{\beta}^{\alpha_k}) \rightarrow (\hat{t}, \hat{\beta})$ in $(C^0([0, T], \mathbb{R} \times$
 715 $\mathbb{R}^d), \|\cdot\|_\infty)$. Since $\|\dot{\hat{\beta}}_\tau^\alpha\|, \|\dot{\hat{t}}_\tau^\alpha\| \leq 1$ we have, applying the Banach–Alaoglu theorem, that up to a new
 716 subsequence

$$(\dot{\hat{t}}^{\alpha_k}, \dot{\hat{\beta}}^{\alpha_k}) \overset{*}{\rightharpoonup} (\dot{\hat{t}}, \dot{\hat{\beta}}) \text{ in } L_\infty(0, T) \quad (28)$$

and $\|\dot{\hat{\beta}}_\tau\| \leq \liminf_{\alpha_k} \|\dot{\hat{\beta}}_\tau^{\alpha_k}\| \leq 1$ and thus $\dot{\hat{t}}_\tau + \|\dot{\hat{\beta}}_\tau\| \leq 1$:

$$\int_0^T \|\dot{\hat{\beta}}_\tau\| d\tau \leq \int_0^T \liminf_{\alpha_k} \|\dot{\hat{\beta}}_\tau^{\alpha_k}\| d\tau \leq \int_0^{+\infty} \liminf_{\alpha_k} \|\dot{\hat{\beta}}_\tau^{\alpha_k}\| d\tau \leq \liminf_{\alpha_k} \int_0^{+\infty} \|\dot{\hat{\beta}}_\tau^{\alpha_k}\| d\tau < C,$$

where the third inequality is by Fatou's lemma. Note that since $[0, T]$ is bounded then it also implies the weak convergence in any $L_p(0, T)$, $1 \leq p < \infty$. Since $(\hat{\beta}^\alpha)$ converges uniformly on $[0, T]$, and ∇L is continuous, we have that $\nabla L(\hat{\beta}^\alpha)$ converges uniformly to $\nabla L(\hat{\beta})$. Since $\hat{t}^{\alpha_k} \rightarrow \hat{t}$ in $L_1(0, T)$, passing to the limit in the equation $\nabla \tilde{\phi}_\alpha(\hat{\beta}_\tau^\alpha) = - \int_0^\tau \dot{\hat{t}}_s^\alpha \nabla L(\hat{\beta}_s^\alpha) ds$ leads to

$$- \int_0^\tau \dot{\hat{t}}_s \nabla L(\hat{\beta}_s) ds \in \partial \|\hat{\beta}_\tau\|_1,$$

717 due to Lemma 2.

718 Recall from Eq. (27) and the definition of $\tilde{\phi}_\alpha$ that:

$$\dot{\hat{\beta}}_\tau^\alpha = - \frac{\sqrt{\hat{\beta}_\tau^\alpha + \alpha^4} \odot \nabla L(\hat{\beta}_\tau^\alpha)}{1/\ln(1/\alpha) + \|\sqrt{\hat{\beta}_\tau^\alpha + \alpha^4} \odot \nabla L(\hat{\beta}_\tau^\alpha)\|}. \quad (29)$$

719 Hence assuming that $\hat{\beta}_\tau$ is such that $\|\hat{\beta}_\tau \odot \nabla L(\hat{\beta}_\tau)\| \neq 0$, we can ensure that $\|\hat{\beta}_{\tau'} \odot \nabla L(\hat{\beta}_{\tau'})\| \neq 0$
 720 for $\tau' \in [\tau, \tau + \varepsilon]$ and ε small enough. We have then $\frac{\sqrt{\hat{\beta}_{\tau'}^\alpha + \alpha^4} \odot \nabla L(\hat{\beta}_{\tau'}^\alpha)}{1/\ln(1/\alpha) + \|\sqrt{\hat{\beta}_{\tau'}^\alpha + \alpha^4} \odot \nabla L(\hat{\beta}_{\tau'}^\alpha)\|}$ converges

721 uniformly toward $-\frac{|\hat{\beta}_{\tau'} \odot \nabla L(\hat{\beta}_{\tau'})|}{\|\hat{\beta}_{\tau'} \odot \nabla L(\hat{\beta}_{\tau'})\|}$ on $[\tau, \tau + \varepsilon]$. Using the dominated convergence theorem, we

722 have $\int_\tau^{\tau+\varepsilon} \frac{\sqrt{\hat{\beta}_{\tau'}^\alpha + \alpha^4} \odot \nabla L(\hat{\beta}_{\tau'}^\alpha)}{1/\ln(1/\alpha) + \|\sqrt{\hat{\beta}_{\tau'}^\alpha + \alpha^4} \odot \nabla L(\hat{\beta}_{\tau'}^\alpha)\|} d\tau' \rightarrow \int_\tau^{\tau+\varepsilon} \frac{|\hat{\beta}_{\tau'} \odot \nabla L(\hat{\beta}_{\tau'})|}{\|\hat{\beta}_{\tau'} \odot \nabla L(\hat{\beta}_{\tau'})\|} d\tau'$. We therefore obtain

723 $\dot{\hat{\beta}}_\tau = - \frac{|\hat{\beta}_\tau \odot \nabla L(\hat{\beta}_\tau)|}{\|\hat{\beta}_\tau \odot \nabla L(\hat{\beta}_\tau)\|}$ in $L_1[0, T]$. Consequently $\|\dot{\hat{\beta}}_\tau\| = 1$ and $\dot{\hat{t}}_\tau = 0$.

724 **Proof that the loss strictly decreases along the heteroclinic orbits.**

725 Assume $\hat{\beta}_\tau$ is such that $|\hat{\beta}_\tau \odot \nabla L(\hat{\beta}_\tau)| \neq 0$, then the flow follows

$$\dot{\hat{\beta}}_\tau = - \frac{|\hat{\beta}_\tau \odot \nabla L(\hat{\beta}_\tau)|}{\|\hat{\beta}_\tau \odot \nabla L(\hat{\beta}_\tau)\|}$$

726 Letting $\gamma(\tau) = \frac{1}{\|\hat{\beta}_\tau \odot \nabla L(\hat{\beta}_\tau)\|}$ we get:

$$dL(\hat{\beta}_\tau) = -\gamma(\tau) \sum_i |\hat{\beta}_\tau(i) \odot [\nabla L(\hat{\beta}_\tau)]_i|^2 d\tau < 0,$$

727 because $|\hat{\beta}_\tau \odot \nabla L(\hat{\beta}_\tau)|^2 \neq 0$. □

728 Borrowing terminologies from [13], we can distinguish two regimes: when $\dot{\hat{\beta}}_\tau = 0$, the system is
 729 *sticked* to the saddle point. When $\dot{\hat{t}}_\tau = 0$ and $\|\dot{\hat{\beta}}_\tau\| = 1$ the system switches to a *viscous slip* which
 730 follows the normalised flow Eq. (26). We use the term of *heteroclinic orbit* as in the dynamical
 731 systems literature since in the neuron space (u, v) it corresponds to a path with links two distinct
 732 critical points of the loss F . Since $\dot{\hat{t}}_\tau = 0$, this regime happens instantly for the original t time scale
 733 (*i.e.* a jump occurs).

734 From Proposition 8, following the same reasoning as in Section 3, we can show that the rescaled
 735 process converges uniformly to a continuous saddle-to-saddle process where the saddles are linked
 736 by normalized flows.

737 **Theorem 3.** *Let $T > 0$. For all subsequences defined in Proposition 8, there exist times $0 = \tau'_0 <$
 738 $\tau_1 < \tau'_1 < \dots < \tau_p < \tau'_p < \tau_{p+1} = +\infty$ such that the iterates $(\hat{\beta}_\tau^{\alpha_k})_\tau$ converge uniformly on
 739 $[0, T]$ to the following limit trajectory :*

$$\begin{aligned} \text{("Saddle")} \quad & \hat{\beta}_\tau = \beta_k && \text{for } \tau \in [\tau'_k, \tau_{k+1}] \text{ where } 0 \leq k \leq p \\ \text{("Orbit")} \quad & \dot{\hat{\beta}}_\tau = -\frac{|\hat{\beta}_\tau| \odot \nabla L(\hat{\beta}_\tau)}{\|\hat{\beta}_\tau| \odot \nabla L(\hat{\beta}_\tau)\|} && \text{for } \tau \in [\tau_{k+1}, \tau'_{k+1}] \text{ where } 0 \leq k \leq p-1 \end{aligned}$$

740 where the saddles $(\beta_0 = 0, \beta_1, \dots, \beta_p = \beta_{t_1}^*)$ are constructed in Algorithm 1. Also, the loss
 741 $(L(\hat{\beta}_\tau))_\tau$ is constant on the saddles and strictly decreasing on the orbits. Finally, independently of
 742 the chosen subsequence, for $k \in [p]$ we have $\hat{t}_{\tau_k} = \hat{t}_{\tau'_k} = t_k$ where the times $(t_k)_{k \in [p]}$ are defined
 743 through Algorithm 1.

744 *Proof.* Some parts of the proof are slightly technical. To simplify the understanding, we make use of
 745 auxiliary lemmas which are stated in Appendix F. The overall spirit follows the intuitive ideas given
 746 in Section 3 and relies on showing that Eq. (25) can only be satisfied if the iterates visit the saddles
 747 from Algorithm 1.

748 We let $\hat{s}_\tau := -\int_0^\tau \dot{\hat{t}}_s \nabla L(\hat{\beta}_s) ds$, which is continuous and satisfies $\hat{s}_\tau \in \partial \|\hat{\beta}_\tau\|_1$ from Eq. (25).
 749 Let $S = \{\beta \in \mathbb{R}^d, |\beta| \odot \nabla L(\beta) = \mathbf{0}\}$ denote the set of critical points and let (β_k, t_k, s_k) be the
 750 successive values of (β, t, s) which appear in the loops of Algorithm 1.

751 **We do a proof by induction:** we start by assuming that the iterates are stuck at the saddle β_{k-1} at
 752 time $\tau \geq \tau'_{k-1}$ where $\hat{t}_{\tau'_{k-1}} = t_{k-1}$ and $\hat{s}_{\tau'_{k-1}} = s_{k-1}$ (recurrence hypothesis), we then show that
 753 they can only move at a time τ_k and follow the normalised flow Eq. (26). We finally show that they
 754 must end up “stuck” at the new critical point β_k , validating the recurrence hypothesis.

755 *Proof of the jump time τ_k such that $\hat{t}_{\tau_k} = t_k$:* we set ourselves at time $\tau \geq \tau'_{k-1}$, stuck at the
 756 saddle β_{k-1} . Let $\tau_k := \sup\{\tau, \hat{t}_\tau \leq t_k\}$, we have that $\tau_k < \infty$ from Lemma 3. Note that by
 757 continuity of \hat{t}_τ it holds that $\hat{t}_{\tau_k} = t_k$. Now notice that $\hat{s}_\tau = \hat{s}_{\tau'_{k-1}} - (\hat{t}_\tau - \hat{t}_{\tau'_{k-1}}) \nabla L(\beta_{k-1}) =$
 758 $s_{k-1} - (\hat{t}_\tau - t_{k-1}) \nabla L(\beta_{k-1})$. We argue that for any $\varepsilon > 0$, we cannot have $\hat{\beta}_\tau = \beta_{k-1}$ on
 759 $(\tau_k, \tau_k + \varepsilon)$. Indeed by the definition of τ_k and from the algorithmic construction of time t_k , it would
 760 lead to $|\hat{s}_\tau(i)| > 1$ for some coordinate $i \in [d]$, which contradicts Eq. (25). Therefore the iterates
 761 must move at the time τ_k .

762 *Heterocline leaving β_{k-1} for $\tau \in [\tau_k, \tau'_k]$:* contrary to before, our time rescaling enables to capture
 763 what happens during the “jump”. We have shown that for any ε , there exists $\tau_\varepsilon \in (\tau_k, \tau_k + \varepsilon)$, such
 764 that $\hat{\beta}_{\tau_\varepsilon} \neq \beta_{k-1}$. From Lemma 4, since the saddles are distinct along the flow, we must have that
 765 $\hat{\beta}_{\tau_\varepsilon} \notin S$ for ε small enough. The iterates therefore follow a heterocline flow leaving β_{k-1} with a
 766 speed of 1 given by Eq. (26). We now define $\tau'_k := \inf\{\tau > \tau_k, \exists \varepsilon_0 > 0, \forall \varepsilon \in [0, \varepsilon_0], \hat{\beta}_{\tau+\varepsilon} \in S\}$
 767 which corresponds to the time at which the iterates reach a new critical point and stay there for at
 768 least a small time ε_0 . We have just shown that $\tau'_k > \tau_k$. Now from Proposition 8, the path length of $\hat{\beta}$
 769 is finite, and from Lemma 4 the flow visits a finite number of distinct saddles at a speed of 1. These
 770 two arguments put together, we get that $\tau'_k < +\infty$ and also $\hat{\beta}_{\tau'_k+\varepsilon} = \hat{\beta}_{\tau'_k}, \forall \varepsilon \in [0, \varepsilon_0]$. On another
 771 note, since $\dot{\hat{t}}_\tau = 0$ for $\tau \in [\tau_k, \tau'_k]$ we have $\hat{t}_{\tau'_k} = \hat{t}_{\tau_k} (= t_k)$ as well as $\hat{s}_{\tau_k} = \hat{s}_{\tau'_k} (= s_k)$.

772 *Proof of the landing point* β_k : we now want to find to which saddle $\hat{\beta}_{\tau'_k} \in S$ the iterates have moved
 773 to. To that end, we consider the following sets which also appear in Algorithm 1:

$$I_{\pm,k} := \{i \in \{1, \dots, d\}, \text{ s.t. } \hat{s}_{\tau'_k}(i) = \pm 1\} \quad \text{and} \quad I_k = I_{+,k} \cup I_{-,k}. \quad (30)$$

774 The set I_k corresponds to the coordinates of $\hat{\beta}_{\tau'_k}$ which “are allowed” (but not obliged) to be activated
 775 (i.e. non-zero). For $\tau \in [\tau'_k, \tau'_k + \varepsilon_0]$ we have that $\hat{s}_\tau = \hat{s}_{\tau'_k} - (\hat{t}_\tau - t_k) \nabla L(\hat{\beta}_{\tau'_k})$. By continuity of \hat{s}
 776 and the fact that $\hat{s}_\tau \in \partial \|\hat{\beta}_{\tau'_k}\|_1$, the equality translates into:

- 777 • if $i \notin I_k$, $\hat{\beta}_{\tau'_k}(i) = 0$
- 778 • if $i \in I_{+,k}$, then $[\nabla L(\hat{\beta}_{\tau'_k})]_i \geq 0$ and $\hat{\beta}_{\tau'_k}(i) \geq 0$
- 779 • if $i \in I_{-,k}$, then $[\nabla L(\hat{\beta}_{\tau'_k})]_i \leq 0$ and $\hat{\beta}_{\tau'_k}(i) \leq 0$
- 780 • for $i \in I_k$, if $\hat{\beta}_{\tau'_k}(i) \neq 0$, then $[\nabla L(\hat{\beta}_{\tau'_k})]_i = 0$

781 One can then notice that these conditions exactly correspond to the optimality conditions of the
 782 following constrained minimisation problem:

$$\begin{aligned} \arg \min \quad & L(\beta). \\ & \beta_i \geq 0, i \in I_{k,+} \\ & \beta_i \leq 0, i \in I_{k,-} \\ & \beta_i = 0, i \notin I_k \end{aligned} \quad (31)$$

783 We showed in Proposition 6 that the solution to this problem is unique and equal to β_k from Algorithm
 784 1. Therefore $\hat{\beta}_\tau = \beta_k$ for $\tau \in [\tau'_k, \tau'_k + \varepsilon_0]$. It finally remains to show that $\hat{\beta}_\tau = \beta_k$ while $\tau \leq \tau_{k+1}$,
 785 where $\tau_{k+1} := \sup\{\tau, \hat{t}_\tau = t_{k+1}\}$. For this let $\tau \in [\tau'_k, \tau_{k+1}]$, notice that for $i \notin I_k$, we necessarily
 786 have that $\hat{\beta}_\tau(i) = \beta_k(i) = 0$, otherwise we break the continuity of \hat{s}_τ . Similarly, for $i \in I_{k,+}$, we
 787 necessarily have that $\hat{\beta}_\tau(i) \geq 0$ and for $i \in I_{k,-}$, $\hat{\beta}_\tau(i) \leq 0$ for the same continuity reasons. Now
 788 assume that $\hat{\beta}_\tau(I_k) \neq \beta_k(I_k)$. Then from Lemma 4 and continuity of the flow, $\exists \tau' \in (\tau'_k, \tau)$ such
 789 that $\hat{\beta}_{\tau'} \notin S$ and there must exist a heterocline flow Eq. (26) starting from β_k which passes through
 790 $\hat{\beta}_{\tau'}$. This is absurd since along this flow the loss strictly decreases, which is in contradiction with the
 791 definition of β_k which minimises the problem Eq. (31). \square

792 E.1 Proof of Theorem 2

793 Theorem 3 enables to prove without difficulty Theorem 2 which we recall below. Indeed we can
 794 show that any extracted limit $\hat{\beta}$ maps back to the unique discontinuous process $\hat{\beta}^\circ$.

795 **Theorem 2.** *Let the saddles $(\beta_0 = \mathbf{0}, \beta_1, \dots, \beta_{p-1}, \beta_p = \beta_{\ell_1}^*)$ and jump times $(t_0 = 0, t_1, \dots, t_p)$
 796 be the outputs of Algorithm 1 and let $(\tilde{\beta}_t^\circ)_t$ be the piecewise constant process defined as follows:*

$$\text{(Saddles)} \quad \tilde{\beta}_t^\circ = \beta_k \quad \text{for } t \in (t_k, t_{k+1}) \text{ and } 0 \leq k \leq p, \quad t_{p+1} = +\infty.$$

797 *The accelerated flow $(\tilde{\beta}_t^\alpha)_t$ defined in Eq. (9) uniformly converges towards the limiting process $(\tilde{\beta}_t^\circ)_t$
 798 on any compact subset of $\mathbb{R}_{\geq 0} \setminus \{t_1, \dots, t_p\}$.*

799 *Proof.* We directly apply Theorem 3, let α_k be the subsequence from the theorem. Let $\varepsilon > 0$, for
 800 simplicity we prove the result on $[t_1 + \varepsilon, t_2 - \varepsilon]$, all the other compacts easily follow the same line
 801 of proof. Note that since $\hat{t}^{\alpha_k}(\tau'_1) \rightarrow t_1$ and $\hat{t}^{\alpha_k}(\tau_2) \rightarrow t_2$, for α_k small enough $\hat{t}^{\alpha_k}(\tau'_1) \leq t_1 + \varepsilon$ and
 802 $\hat{t}^{\alpha_k}(\tau_2) \geq t_2 - \varepsilon$, by the monotonicity of τ^{α_k} , this means that for α_k small enough, $\tau'_1 \leq \tau^{\alpha_k}(t_1 + \varepsilon)$
 803 and $\tau_2 \geq \tau^{\alpha_k}(t_2 - \varepsilon)$. Therefore

$$\begin{aligned} \sup_{t \in [t_1 + \varepsilon, t_2 - \varepsilon]} \|\tilde{\beta}_t^{\alpha_k} - \beta_1\| &= \sup_{t \in [t_1 + \varepsilon, t_2 - \varepsilon]} \|\hat{\beta}^{\alpha_k}(\tau_{\alpha_k}(t)) - \beta_1\| \\ &= \sup_{\tau \in [\tau^{\alpha_k}(t_1 + \varepsilon), \tau^{\alpha_k}(t_2 - \varepsilon)]} \|\hat{\beta}^{\alpha_k}(\tau) - \beta_1\| \\ &\leq \sup_{\tau \in [\tau'_1, \tau_2]} \|\hat{\beta}^{\alpha_k}(\tau) - \beta_1\|, \end{aligned}$$

804 which goes uniformly to 0 following Theorem 3. Since this result is independent of the subsequence
 805 α_k , we get the result of Theorem 2. \square

806 **E.2 Proof of Proposition 2**

807 We restate and prove Proposition 2 below.

Proposition 2. For all $T > t_p$, the graph of the iterates $(\tilde{\beta}_t^\alpha)_{t \leq T}$ converges to that of $(\hat{\beta}_\tau)_\tau$:

$$\text{dist}(\{\tilde{\beta}_t^\alpha\}_{t \leq T}, \{\hat{\beta}_\tau\}_{\tau \geq 0}) \xrightarrow{\alpha \rightarrow 0} 0 \quad (\text{Hausdorff distance})$$

808 *Proof.* For α small enough, we have that $\hat{t}_{\tau_p}^\alpha \leq t_p + \varepsilon \leq T$

$$\begin{aligned} \sup_{\tau \geq 0} d(\hat{\beta}_\tau, \{\tilde{\beta}_t^\alpha\}_{t \leq T}) &= \sup_{\tau \leq \tau_p'} d(\hat{\beta}_\tau, \{\tilde{\beta}_t^\alpha\}_{t \leq T}) \\ &\leq \sup_{\tau \leq \tau_p'} \|\hat{\beta}_\tau - \tilde{\beta}_{\hat{t}_\tau^\alpha}^\alpha\| \\ &= \sup_{\tau \leq \tau_p'} \|\hat{\beta}_\tau - \hat{\beta}_\tau^\alpha\| \xrightarrow{\alpha \rightarrow 0} 0, \end{aligned}$$

809 according to Theorem 3.

810 Similarly:

$$\begin{aligned} \sup_{t \leq T} d(\tilde{\beta}_t^\alpha, \{\hat{\beta}_{\tau'}\}_{\tau'}) &= \sup_{\tau \leq \tau_T^\alpha} d(\hat{\beta}_\tau^\alpha, \{\hat{\beta}_{\tau'}\}_{\tau'}) \\ &\leq \sup_{\tau \leq \tau_T^\alpha} \|\hat{\beta}_\tau^\alpha - \hat{\beta}_\tau\| \xrightarrow{\alpha \rightarrow 0} 0, \end{aligned}$$

811 according to Theorem 3, which concludes the proof. □

812 **F Technical lemmas**

813 The following lemma describes the behaviour of $\nabla\tilde{\phi}_\alpha(\beta^\alpha)$ as $\alpha \rightarrow 0$ in function of the subdifferen-
 814 tial $\partial\|\cdot\|_1$.

815 **Lemma 2.** *Let $(\beta^\alpha)_{\alpha>0}$ such that $\beta^\alpha \xrightarrow[\alpha\rightarrow 0]{} \beta \in \mathbb{R}^d$.*

816 • if $\beta_i > 0$ then $[\nabla\tilde{\phi}_\alpha(\beta^\alpha)]_i$ converges to 1

817 • if $\beta_i < 0$ then $[\nabla\tilde{\phi}_\alpha(\beta^\alpha)]_i$ converges to -1 .

818 Moreover if we assume that $\nabla\tilde{\phi}_\alpha(\beta^\alpha)$ converges to $\eta \in \mathbb{R}^d$, we have that:

819 • $\eta_i \in (-1, 1) \Rightarrow \beta_i = 0$

820 • $\beta_i = 0 \Rightarrow \eta_i \in [-1, 1]$.

821 Overall, assuming that $(\beta^\alpha, \nabla\tilde{\phi}_\alpha(\beta^\alpha)) \xrightarrow[\alpha\rightarrow 0]{} (\beta, \eta)$, we can write:

$$\eta \in \partial\|\beta\|_1.$$

822 *Proof.* We have that

$$\begin{aligned} [\nabla\tilde{\phi}_\alpha(\beta^\alpha)]_i &= \frac{1}{2\ln(1/\alpha)} \operatorname{arcsinh}\left(\frac{\beta_i^\alpha}{\alpha^2}\right) \\ &= \frac{1}{2\ln(1/\alpha)} \ln\left(\frac{\beta_i^\alpha}{\alpha^2} + \sqrt{\frac{(\beta_i^\alpha)^2}{\alpha^4} + 1}\right). \end{aligned}$$

823 Now assume that $\beta_i^\alpha \rightarrow \beta_i > 0$, then $[\nabla\tilde{\phi}_\alpha(\beta^\alpha)]_i \rightarrow 1$, if $\beta_i < 0$ we conclude using that $\operatorname{arcsinh}$ is
 824 an odd function. All the claims are simple consequences of this. \square

825 The following lemma shows that the extracted limits \hat{t} as defined in Proposition 8 diverge to ∞ . This
 826 divergence is crucial as it implies that the rescaled iterates $(\hat{\beta}_\tau)_\tau$ explore the whole trajectory..

827 **Lemma 3.** *For any extracted limit \hat{t} as defined in Proposition 8, we have that $\tau - C \leq \hat{t}_\tau$ where C
 828 is the upperbound on the length of the curves defined in proposition 5.*

829 *Proof.* Recall that

$$\tau^\alpha(t) = t + \int_0^t \|\dot{\hat{\beta}}_s^\alpha\| ds.$$

830 From Proposition 5, the full path length $\int_0^{+\infty} \|\dot{\hat{\beta}}_s^\alpha\| ds$ is finite and bounded by some constant C
 831 independently of α . Therefore τ^α is a bijection in $\mathbb{R}_{\geq 0}$ and we defined $\hat{t}_\tau^\alpha = (\tau^\alpha)^{-1}(\tau)$. Furthermore
 832 $\tau^\alpha(t) \leq t + C$ leads to $t \leq \hat{t}_\tau^\alpha(t + C)$ and therefore $\tau - C \leq \hat{t}_\tau^\alpha(\tau)$ for all $\tau \geq 0$. This inequality
 833 still holds for any converging subsequence, which proves the result. \square

834 Under a mild additional assumption on the data (see Assumption 2), we showed after the proof of
 835 Proposition 1 in Appendix B that the number of saddles of F is finite. Without this assumption, the
 836 number of saddles is *a priori* not finite. However the following lemma shows that along the flow of $\hat{\beta}$
 837 the number of saddles which can potentially be visited is indeed finite.

838 **Lemma 4.** *The limiting flow $\hat{\beta}$ as defined in Proposition 8 can only visit a finite number of critical
 839 points $\beta \in S := \{\beta \in \mathbb{R}^d, \beta \odot \nabla L(\beta) = 0\}$ and can visit each one of them at most once.*

840 *Proof.* Let $\tau \geq 0$, and assume that $\hat{\beta}_\tau \in S$, i.e., we are at a critical point at time τ . From Proposition 1,
 841 we have that

$$\hat{\beta}_\tau \in \arg \min_{\beta_i=0 \text{ for } i \notin \operatorname{supp}(\hat{\beta}_\tau)} L(\beta), \quad (32)$$

842 Let us define the sets

$$I_{\pm} := \{i \in \{1, \dots, d\}, \text{ s.t. } \hat{s}_{\tau}(i) = \pm 1\} \quad \text{and} \quad I = I_{+} \cup I_{-}.$$

843 The set I corresponds to the coordinate of $\hat{\beta}_{\tau}$ which “are allowed” (but not obliged) to be non-zero
844 since from Eq. (25), $\text{supp}(\hat{\beta}_{\tau}) \subset I$. Now given the fact that the sub-matrix $X_I = (\tilde{x}_i)_{i \in I} \in$
845 $\mathbb{R}^{n \times \text{card}(I)}$ is full rank (see part (1) of the proof of Proposition 6 for the explanation), the solution of
846 the minimisation problem (32) is unique and equal to $\hat{\beta}[\xi] = (X_{\xi}^{\top} X_{\xi})^{-1} X_{\xi}^{\top} y$ and $\beta[\xi^C] = 0$ where
847 $\xi = \text{supp}(\hat{\beta}_{\tau})$. There are $2^d = \text{Card}(P([d]))$ (where $P([d])$ contains all the subsets of $[d]$) number
848 of constraints of the form $\{\beta_i = 0, i \notin \mathcal{A}\}$, where $\mathcal{A} \subset [d]$, and $\hat{\beta}_{\tau}$ is the unique solution of one of
849 them. $\hat{\beta}_{\tau}$ can therefore take at most 2^d values (very crude upperbound). There is therefore a finite
850 number of critical points which can be reached by the flow $\hat{\beta}$. Furthermore, from Proposition 8, the
851 loss is strictly decreasing along the heteroclinic orbits, each of these critical points can therefore be
852 visited at most once. \square