

SGD vs GD: Rank Deficiency in Linear Networks

Aditya Varre

TML, EPFL

ADITYA.VARRE@EPFL.CH

Margarita Sagitova

TML, EPFL

MARGARITA.SAGITOVA@EPFL.CH

Nicolas Flammarion

TML, EPFL

NICOLAS.FLAMMARION@EPFL.CH

Abstract

In this article, we study the behaviour of continuous-time gradient methods on a two-layer linear network with square loss. A dichotomy between SGD and GD is revealed: GD preserves the rank at initialization while (label noise) SGD diminishes the rank regardless of the initialization. We demonstrate this rank deficiency by studying the time evolution of the *determinant* of a matrix of parameters. To further understand this phenomenon, we derive the stochastic differential equation (SDE) governing the eigenvalues of the parameter matrix. This SDE unveils a *repulsive force* between the eigenvalues: a key regularization mechanism which induces rank deficiency. Our results are well supported by experiments illustrating the phenomenon beyond linear networks and regression tasks.

1. Introduction

One of the remarkable features of deep learning models is their capacity to learn effective representations that generalize well across different tasks, even when they are heavily overparameterized [6]. This success in learning these representations is often attributed to the gradient-based algorithms used in training. These algorithms navigate complex non-convex landscapes to minimize the training objective and yield effective representations, while avoiding spurious features that arise from the models' vast number of parameters. Empirical studies have revealed that stochastic noise is a key factor in mitigating spurious features and enabling effective representation learning. In this paper, we investigate this overarching problem, specifically focusing on studying :

How does stochasticity facilitate the discovery of solutions with simplified structures?

We explore this question using a simplified model: a single hidden-layer linear network. Despite lacking non-linearity, such networks capture some intricate phenomena of real-world deep networks and have been extensively studied to understand convergence [3, 40], learning dynamics [45], and the implicit bias of optimization algorithms [23, 47]. Our work builds on this foundation by comparing stochastic algorithms with their deterministic counterparts, focusing on how these differences influence the learning of simpler structures. Specifically, we analyze vector regression on two-layer linear networks trained with both gradient flow and stochastic gradient flow methods and make the following contributions:

- In Section 3, we track the evolution of the determinant of the parameter matrix under gradient flow and stochastic gradient flow. We show that stochastic gradient flow drives the determinant towards zero, effectively removing irrelevant direction(s).
- In Section C, we derive a stochastic differential equation that describes the behavior of the eigenvalues of the parameter matrix. This analysis reveals a repulsive force between eigenvalues that pushes them apart and a geometric Brownian motion that pulls them toward zero.
- In Section D, we discuss the generalizability of our approach beyond square loss and various noise models, including discrete step sizes. Finally, we present experimental results in Section E that support our theoretical findings.

1.1. Related Work

Our work lies at the convergence of distinct research topics:

Effect of SGD on generalization. The relationship between the stochasticity of SGD and its generalization capabilities has been extensively examined [26, 28, 30, 33, 37]. Notably, SGD tends to yield models with superior generalization compared to gradient descent [26, 30, 32]. Various explorations into this phenomenon have been conducted through various approaches: hypothesizing that SGD favors flatter minima linked to better generalization, as opposed to sharp minima associated with poor generalization [1, 27, 32],

Stochastic dynamics and Label Noise. Recent literature has explored label noise-driven Gradient Descent as an effective method to probe the beneficial impact of stochasticity on generalization, with two distinct perspectives emerging. Firstly, an asymptotic view on general model parametrization is considered, where Blanc et al. [7], Damian et al. [14], Li et al. [36] suggest that stochastic dynamics preferentially optimize a hidden objective linked to the curvature of the loss.. Secondly, specifically for diagonal linear networks, HaoChen et al. [24], Pillaud-Vivien et al. [44] observe a similar collapsing effect due to label noise but with a finer characterization of the limiting process. Finally, in the absence of label noise, Even et al. [17], Pesme et al. [43] have characterized the solutions of stochastic GF and GD for diagonal linear networks. Recently, Ghosh et al. [19] further exhibit a similar sparser features effect for single-neuron autoencoder.

2. Linear networks and continuous-time gradient method

Vector regression. We study the vector regression problems with inputs x_1, \dots, x_n in $(\mathbb{R}^p)^n$ and outputs y_1, \dots, y_n in $(\mathbb{R}^k)^n$. We consider the minimization of the square loss over a class of parametric models $\mathcal{H} = \{f_\theta(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^k \mid \theta \in \mathbb{R}^p\}$ specified in the next paragraph. The train loss therefore can be written as $\mathcal{L}(\theta) = \frac{1}{2n} \sum_{i=1}^n (y_i - f_\theta(x_i))^2$.

Parameterization with a linear network. We focus on two-layer linear neural networks of width $l \in \mathbb{N}^*$. The model is described by the parameterization $\theta = (\mathbf{W}_1, \mathbf{W}_2)$, where $\mathbf{W}_1 \in \mathbb{R}^{p \times l}$ and $\mathbf{W}_2 \in \mathbb{R}^{l \times k}$, and the function $f_\theta(x) = \mathbf{W}_2^\top \mathbf{W}_1^\top x$. This model is linear with respect to the input x . In terms of expressivity, it is comparable to the linear class of predictors, represented as $f_\beta(x) = \beta^\top x$, where β equals $\mathbf{W}_1 \mathbf{W}_2$. Throughout our analysis, we denote the equivalent linear predictor of the network as β . It is important to note that this parameterization introduces some redundancy, a single linear predictor β can have multiple representations $\mathbf{W}_1, \mathbf{W}_2$ such that $\mathbf{W}_1 \mathbf{W}_2 = \beta$.

Some representations have a rich structure whereas other resemble random features. For example, consider the case of scalar regression ($k=1$), for a vector β there exists rich parameterizations where all the neurons, i.e., columns of \mathbf{W}_1 align with β and also some lazy structures where \mathbf{W}_1 resembles a random matrix [12, 51].

Train loss. By defining $X^\top = [x_1, \dots, x_n]$ and $Y^\top = [y_1, \dots, y_n]$, the loss function is given by:

$$\mathcal{L}(\mathbf{W}_1, \mathbf{W}_2) = \frac{1}{2n} \|X\mathbf{W}_1\mathbf{W}_2 - Y\|^2. \quad (2.1)$$

For simplicity, we adjust for the normalization factor n by rescaling the data to $(X, Y) \leftarrow (X/\sqrt{n}, Y/\sqrt{n})$, thereby implicitly considering it in the loss function without directly mentioning n in the formula. Note that the loss is non-convex in $\mathbf{W}_1, \mathbf{W}_2$.

Gradient flow. The dynamics induced in parameter space by running GF on Equation (2.1) is given by

$$d\mathbf{W}_1 = -\nabla_{\mathbf{W}_1} \mathcal{L}(\mathbf{W}_1, \mathbf{W}_2) dt = X^\top (Y - X\mathbf{W}_1\mathbf{W}_2) \mathbf{W}_2^\top dt, \quad (2.2)$$

$$d\mathbf{W}_2 = -\nabla_{\mathbf{W}_2} \mathcal{L}(\mathbf{W}_1, \mathbf{W}_2) dt = \mathbf{W}_1^\top X^\top (Y - X\mathbf{W}_1\mathbf{W}_2) dt. \quad (2.3)$$

Introducing the block matrix, $\Theta = [\mathbf{W}_1^\top \mid \mathbf{W}_2] \in \mathbb{R}^{l \times (p+k)}$ and denoting the residual matrix by $\mathbf{R} = X^\top (Y - X\mathbf{W}_1\mathbf{W}_2)$, the evolution of Θ can be written as

$$d\Theta = [d\mathbf{W}_1^\top \mid d\mathbf{W}_2] = [\mathbf{W}_2 \mathbf{R}^\top dt \mid \mathbf{W}_1^\top \mathbf{R} dt] = [\mathbf{W}_1^\top \mid \mathbf{W}_2] \begin{bmatrix} 0_{p \times p} & \mathbf{R} \\ \mathbf{R}^\top & 0_{k \times k} \end{bmatrix} dt.$$

The gradient flow can therefore be compactly written as

$$d\Theta = \Theta \mathbf{J} dt, \quad \text{where } \mathbf{J} = \begin{bmatrix} 0_{p \times p} & \mathbf{R} \\ \mathbf{R}^\top & 0_{k \times k} \end{bmatrix}. \quad (2.4)$$

The gradient flow (GF), when expressed in this form, reveals an inherent multiplicative structure with respect to Θ in the gradient of the loss. As we see in subsequent sections, this representation of the gradient flow with block matrices proves to be very convenient.

Label noise gradient descent. Label noise gradient descent (LNGD) is a theoretically studied alternative to SGD that mirrors its practical behavior by sharing the geometric properties of the noise [7, 14]. Let $\varepsilon_t \in \mathbb{R}^{n \times k}$, where each entry of ε_t is an independent Gaussian random variable. At iteration t , the labels are perturbed with this Gaussian noise at an intensity δ , i.e., $\tilde{Y} = Y + \sqrt{\delta} \varepsilon_t$. The LNGD algorithm updates the iterates with a step size η in the direction of the gradient computed after the labels have been perturbed, as follows:

$$\mathbf{W}_1^{t+1} = \mathbf{W}_1^t - \eta \nabla_{\mathbf{W}_1} \mathcal{L}(\tilde{Y}, X, \mathbf{W}_1^t, \mathbf{W}_2^t); \quad \mathbf{W}_2^{t+1} = \mathbf{W}_2^t - \eta \nabla_{\mathbf{W}_2} \mathcal{L}(\tilde{Y}, X, \mathbf{W}_1^t, \mathbf{W}_2^t),$$

where, by an abuse of notation, $\mathcal{L}(Y, X, \mathbf{W}_1, \mathbf{W}_2) = 1/2 \|X\mathbf{W}_1\mathbf{W}_2 - Y\|^2$. The iterates can then be restructured into a block matrix:

$$\Theta^{t+1} = \Theta^t - \eta \Theta^t \mathbf{J}_t - \eta \sqrt{\delta} \Theta^t \xi_t, \quad \text{where } \xi_t = \begin{bmatrix} 0_{p \times p} & X^\top \varepsilon_t \\ \varepsilon_t^\top X & 0_{k \times k} \end{bmatrix}, \quad (2.5)$$

and \mathbf{J}_t is defined as in Equation (2.4).

Stochastic gradient flow (SGF). We aim to model the aforementioned LNGD in continuous time using an appropriate SDE. Stochastic continuous-time counterparts of discrete stochastic gradient algorithms are favored for their enhanced amenability to theoretical analysis. We propose the following stochastic differential equation (SDE) to model LNGD in continuous time:

$$d\Theta = \Theta \left[\mathbf{J}dt + \sqrt{\eta\delta}d\xi \right], \text{ where } d\xi = \begin{bmatrix} 0_{p \times p} & X^\top d\mathbf{B}_t \\ d\mathbf{B}_t^\top X & 0_{k \times k}, \end{bmatrix} \quad (2.6)$$

where \mathbf{B}_t denotes a matrix Brownian motion in $\mathbb{R}^{n \times k}$. LNGD as defined in Equation (2.5), can be interpreted as the Euler-Maryama discretization of the above SGF with a stepsize η .

Initialization. The dynamics of gradient methods on homogeneous models are significantly influenced by initialization, which determines the regime they operate in—specifically, the lazy regime for large initializations and the rich regime for small ones [12, 54]. Thus, the scale of initialization has garnered significant interest, particularly its impact on the training of networks with GD [8, 54]. It is observed that stochastic methods eliminate the dependence on initialization [43].

Conserved quantities and balanceness. Gradient flows follow specific conservation laws along their trajectory [38], maintaining characteristics of the initial conditions. For linear networks, this conservation manifests as the *balanceness property* [15], described by:

$$\Delta = \mathbf{W}_1^\top \mathbf{W}_1 - \mathbf{W}_2 \mathbf{W}_2^\top = \mathbf{W}_1^\top(0) \mathbf{W}_1(0) - \mathbf{W}_2(0) \mathbf{W}_2^\top(0).$$

As a result, Arora et al. [2, 4], Saxe et al. [45] have adopted *balanced initialization*, where $\Delta(0) = 0$, to ensure that weight matrices remain low rank throughout the trajectory. However, unbalanced initialization do not preserve these simple low-rank structures, as aspects of the initial conditions persist. In contrast, stochastic methods do not adhere to these conservation laws [55] and the evolution of the imbalance Δ for SGF is $d\Delta = d(\mathbf{W}_1^\top \mathbf{W}_1 - \mathbf{W}_2 \mathbf{W}_2^\top) = \text{tr}(X X^\top) \mathbf{W}_2 \mathbf{W}_2^\top dt - k \mathbf{W}_1^\top X^\top X \mathbf{W}_1 dt$. While there is no diffusion term in the derivative, the matrices remain stochastic and no definitive conclusions can be drawn from this. However, in the case where $k = p$ and $X^\top X = I_p$, it can be shown that $\mathbf{W}_1^\top \mathbf{W}_1 - \mathbf{W}_2 \mathbf{W}_2^\top \rightarrow 0$, indicating that the stochastic noise eliminates initial imbalance.

Conclusion. Understanding how stochastic methods mitigate dependency on initialization requires exploring beyond the evolution of the imbalance Δ . To this end, we identify and discuss other conserved quantities, such as the determinant of the matrix $\Theta^\top \Theta$ in the following sections.

3. Separation between Gradient Flow through determinant

Here, we present our first separation result between GF and SGF. While the determinant of the parameters is preserved in GF, it is driven to zero by the stochasticity of SGF, leading to a simplistic low-rank structure.

Determinant evolution of the gradient flow The theorem below demonstrates that the determinant of the parameters is preserved in gradient flow.

Theorem 3.1 *For the gradient flow defined in Equation (2.4), the following property holds,*

$$d\left(\det\left(\Theta^\top \Theta\right)\right) = 0.$$

Hence, $\det\left(\Theta(t)^\top \Theta(t)\right) = \det\left(\Theta_0^\top \Theta_0\right)$, where $\Theta_0 = \Theta(0)$ is the initialisation at time $t = 0$.

The proof presented in the App. F.1, is based on straightforward computations of the derivative of the determinant and the fact that the matrix \mathbf{J} has zero trace. We note that the simplicity of the proof arises from the strategically chosen block structure of Θ . This result would have been less straightforward with different parametrizations, which likely explains why such a simple finding appears to be novel. The theorem implies that the determinant of \mathbf{M} along the trajectory remains equal to the determinant at initialization. If $\Theta_0^\top \Theta_0$ is full-rank initially, meaning the determinant is non-zero, the theorem ensures that the determinant of \mathbf{M} remains non-zero. Consequently, the rank of Θ does not diminish along the trajectory. When $l \geq p + k$, i.e., the hidden layer has a large width and $\mathbf{W}_1, \mathbf{W}_2$ are initialized randomly from a Gaussian distribution, $\Theta_0^\top \Theta_0$ has full rank almost surely. The theorem also reveals some implications regarding the impact of initialization scale. Note that $\lambda_{\min}(A) \leq \sqrt[p]{\det A}$, indicating that when the scale of initialization is very small, at least one singular value of Θ is small.

Determinant evolution of the stochastic gradient flow In contrast, the theorem presented below demonstrates that the determinant of the parameters converges to zero in stochastic gradient flow.

Theorem 3.2 *For the SDE, defined in the Equation (2.6), for $t \leq \tau_\infty$, the following property holds for the evolution of determinant*

$$d\left(\det\left(\Theta^\top \Theta\right)\right) = -2\eta\delta k \text{tr}\left(X^\top X\right)\det\left(\Theta^\top \Theta\right)dt.$$

Hence, $\det\left(\Theta(t)^\top \Theta(t)\right) = \det\left(\Theta_0^\top \Theta_0\right) \exp -2\eta\delta k \text{tr}\left(X^\top X\right)t$, where Θ_0 is the initialization.

Although the evolution of the parameters in SGF is random, the evolution of the determinant is deterministic. The theorem highlights a striking phenomenon: the noise in SGF diminishes the determinant along the trajectory, leading to a simplification of the network over time. The larger the noise and the stepsize, the faster the determinant vanishes. The vanishing of the determinant suggests that the rank of the parameters decreases by at least one, effectively eliminating some components. It holds for any initialization of Θ_0 and indicates how the SGF overrides some aspects of initialization. The proof uses the fact that stochastic Brownian term in the SDE, through Itô’s calculus, introduces a negative drift, driving the determinant to zero (refer to F.3 for the proof).

Limitations. Given the large width of the hidden layer, the determinant converging to zero does not fully reveal the complexity of the situation. It merely indicates that at least one singular value is approaching zero. Furthermore, the theorem provides limited insights when the determinant is already 0 at initialization, $\det \Theta_0 = 0$ which happens whenever $l < p + k$. Next, we explore the mechanisms behind this low-rank phenomenon, suggesting that the repulsive forces induced by stochasticity drive the spurious singular values to zero as seen in the right plot of Figure 1.

4. Conclusion

In this paper, we demonstrate a distinct separation between GF and SGF when trained on linear networks. This separation is obtained by tracking the evolution of the determinant of the parameter matrix. However, while the determinant is a significant factor, it does not fully capture the implicit regularization effects. We try to partially address this issue by studying the dynamics of singular values in Sec C. In section D, we extend our approach to shed some light on the training dynamics for losses other than square loss and discrete stepsize dynamics.

References

- [1] Maksym Andriushchenko, Francesco Croce, Maximilian Müller, Matthias Hein, and Nicolas Flammarion. A modern look at the relationship between sharpness and generalization. In *International Conference on Machine Learning*, pages 840–902. PMLR, 2023.
- [2] Sanjeev Arora, Nadav Cohen, and Elad Hazan. On the optimization of deep networks: Implicit acceleration by overparameterization. In *Proceedings of the 35th International Conference on Machine Learning*, 2018.
- [3] Sanjeev Arora, Nadav Cohen, Noah Golowich, and Wei Hu. A convergence analysis of gradient descent for deep linear neural networks. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=SkMQg3C5K7>.
- [4] Sanjeev Arora, Nadav Cohen, Wei Hu, and Yuping Luo. Implicit regularization in deep matrix factorization. *Advances in Neural Information Processing Systems*, 32, 2019.
- [5] Sanjeev Arora, Nadav Cohen, Wei Hu, and Yuping Luo. Implicit regularization in deep matrix factorization. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper_files/paper/2019/file/c0c783b5fc0d7d808f1d14a6e9c8280d-Paper.pdf.
- [6] Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE Trans. Pattern Anal. Mach. Intell.*, 2013.
- [7] Guy Blanc, Neha Gupta, Gregory Valiant, and Paul Valiant. Implicit regularization for deep neural networks driven by an ornstein-uhlenbeck like process. In *Conference on Learning Theory, COLT 2020*, Proceedings of Machine Learning Research. PMLR, 2020.
- [8] Etienne Boursier, Loucas Pillaud-Vivien, and Nicolas Flammarion. Gradient flow dynamics of shallow reLU networks for square loss and orthogonal inputs. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=L74c-iUxQ1I>.
- [9] Lukas Braun, Clémentine Carla Juliette Dominé, James E Fitzgerald, and Andrew M Saxe. Exact learning dynamics of deep linear networks with prior knowledge. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=1Jx2vng-KiC>.
- [10] Marie-France Bru. Diffusions of perturbed principal component analysis. *Journal of Multivariate Analysis*, 29(1):127–136, 1989. ISSN 0047-259X. doi: [https://doi.org/10.1016/0047-259X\(89\)90080-8](https://doi.org/10.1016/0047-259X(89)90080-8). URL <https://www.sciencedirect.com/science/article/pii/0047259X89900808>.
- [11] Marie-France Bru. Wishart processes. *Journal of Theoretical Probability*, 4:725–751, 1991.
- [12] Lénaïc Chizat, Edouard Oyallon, and Francis Bach. *On Lazy Training in Differentiable Programming*. Curran Associates Inc., Red Hook, NY, USA, 2019.

- [13] Jeremy Cohen, Simran Kaur, Yuanzhi Li, J Zico Kolter, and Ameet Talwalkar. Gradient descent on neural networks typically occurs at the edge of stability. In *International Conference on Learning Representations*, 2020.
- [14] Alex Damian, Tengyu Ma, and Jason D Lee. Label noise sgd provably prefers flat global minimizers. *Advances in Neural Information Processing Systems*, 34:27449–27461, 2021.
- [15] Simon S Du, Wei Hu, and Jason D Lee. Algorithmic regularization in learning deep homogeneous models: Layers are automatically balanced. In *Advances in Neural Information Processing Systems*, 2018.
- [16] Freeman J. Dyson. A Brownian-Motion Model for the Eigenvalues of a Random Matrix. *Journal of Mathematical Physics*, 3(6):1191–1198, 11 1962. ISSN 0022-2488. doi: 10.1063/1.1703862. URL <https://doi.org/10.1063/1.1703862>.
- [17] Mathieu Even, Scott Pesme, Suriya Gunasekar, and Nicolas Flammarion. (s)gd over diagonal linear networks: Implicit regularisation, large stepsizes and edge of stability. *Advances in Neural Information Processing Systems*, 2023.
- [18] Kenji Fukumizu. Effect of batch learning in multilayer neural networks. *Gen*, 1(04):1E–03, 1998.
- [19] Nikhil Ghosh, Spencer Frei, Wooseok Ha, and Bin Yu. The effect of sgd batch size on autoencoder learning: Sparsity, sharpness, and feature learning. *arXiv preprint arXiv:2308.03215*, 2023.
- [20] Gauthier Gidel, Francis Bach, and Simon Lacoste-Julien. Implicit regularization of discrete gradient dynamics in linear neural networks. *Advances in Neural Information Processing Systems*, 32, 2019.
- [21] Piotr Graczyk and Jacek Małecki. Multidimensional yamada-watanabe theorem and its applications to particle systems. *Journal of Mathematical Physics*, 54(2), 2013.
- [22] Piotr Graczyk and Jacek Małecki. Strong solutions of non-colliding particle systems. 2014.
- [23] Suriya Gunasekar, Blake E Woodworth, Srinadh Bhojanapalli, Behnam Neyshabur, and Nati Srebro. Implicit regularization in matrix factorization. *Advances in Neural Information Processing Systems*, 30, 2017.
- [24] Jeff Z. HaoChen, Colin Wei, Jason D. Lee, and Tengyu Ma. Shape matters: Understanding the implicit bias of the noise covariance. In Mikhail Belkin and Samory Kpotufe, editors, *Conference on Learning Theory, COLT 2021, 15-19 August 2021, Boulder, Colorado, USA*, 2021.
- [25] Charles R. Harris, K. Jarrod Millman, Stéfan J. van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern, Matti Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fernández del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-Marchant, Kevin Sheppard, Tyler Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and

- Travis E. Oliphant. Array programming with NumPy. *Nature*, 585(7825):357–362, September 2020. doi: 10.1038/s41586-020-2649-2. URL <https://doi.org/10.1038/s41586-020-2649-2>.
- [26] Fengxiang He, Tongliang Liu, and Dacheng Tao. Control batch size and learning rate to generalize well: Theoretical and empirical evidence. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- [27] Sepp Hochreiter and Jürgen Schmidhuber. Flat minima. *Neural Comput.*, 9(1):1–42, jan 1997.
- [28] Elad Hoffer, Itay Hubara, and Daniel Soudry. Train longer, generalize better: Closing the generalization gap in large batch training of neural networks. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS’17, page 1729–1739, 2017.
- [29] Nobuyuki Ikeda and Shinzo Watanabe. *Stochastic differential equations and diffusion processes*, volume 24 of *North-Holland Mathematical Library*. North-Holland Publishing Co., Amsterdam, 1981. ISBN 0-444-86172-6.
- [30] Stanislaw Jastrzebski, Zac Kenton, Devansh Arpit, Nicolas Ballas, Asja Fischer, Amos Storkey, and Yoshua Bengio. Three factors influencing minima in SGD. In *International Conference on Learning Representations*, 2018.
- [31] Ziwei Ji and Matus Telgarsky. Gradient descent aligns the layers of deep linear networks. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=HJflg30qKX>.
- [32] Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang. On large-batch training for deep learning: Generalization gap and sharp minima. In *International Conference on Learning Representations*, 2017.
- [33] Bobby Kleinberg, Yuanzhi Li, and Yang Yuan. An alternative view: When does SGD escape local minima? In *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 2698–2707. PMLR, 10–15 Jul 2018.
- [34] Aitor Lewkowycz, Yasaman Bahri, Ethan Dyer, Jascha Sohl-Dickstein, and Guy Gur-Ari. The large learning rate phase of deep learning: the catapult mechanism. *arXiv preprint arXiv:2003.02218*, 2020.
- [35] Qianxiao Li, Cheng Tai, and Weinan E. Stochastic modified equations and dynamics of stochastic gradient algorithms i: Mathematical foundations. *Journal of Machine Learning Research*, 20(40):1–47, 2019.
- [36] Zhiyuan Li, Tianhao Wang, and Sanjeev Arora. What happens after sgd reaches zero loss?—a mathematical framework. In *International Conference on Learning Representations*, 2021.
- [37] Stephan Mandt, Matthew D. Hoffman, and David M. Blei. A variational analysis of stochastic gradient algorithms. In *Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48*, ICML’16, page 354–363, 2016.

- [38] Sibylle Marcotte, Rémi Gribonval, and Gabriel Peyré. Abide by the law and follow the flow: conservation laws for gradient flows. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=kMueEV8Eyy>.
- [39] Eberhard Mayerhofer, Oliver Pfaffel, and Robert Stelzer. On strong solutions for positive definite jump diffusions. *Stochastic processes and their applications*, 121(9):2072–2086, 2011.
- [40] Hancheng Min, Salma Tarmoun, Rene Vidal, and Enrique Mallada. On the explicit role of initialization on the convergence and implicit bias of overparametrized linear networks. In *Proceedings of the 38th International Conference on Machine Learning*, 2021.
- [41] J. R. Norris, L. C. G. Rogers, and David Williams. Brownian motions of ellipsoids. *Transactions of the American Mathematical Society*, 294(2):757–765, 1986. ISSN 00029947. URL <http://www.jstor.org/stable/2000214>.
- [42] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019. URL <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>.
- [43] Scott Pesme, Loucas Pillaud-Vivien, and Nicolas Flammarion. Implicit bias of sgd for diagonal linear networks: a provable benefit of stochasticity. *Advances in Neural Information Processing Systems*, 34:29218–29230, 2021.
- [44] Loucas Pillaud-Vivien, Julien Reygner, and Nicolas Flammarion. Label noise (stochastic) gradient descent implicitly solves the lasso for quadratic parametrisation. In *Conference on Learning Theory*, pages 2127–2159. PMLR, 2022.
- [45] Andrew M. Saxe, James L. McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. In *2nd International Conference on Learning Representations, ICLR 2014*, 2014.
- [46] Andrew M Saxe, James L McClelland, and Surya Ganguli. A mathematical theory of semantic development in deep neural networks. *Proceedings of the National Academy of Sciences*, 116(23):11537–11546, 2019.
- [47] Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The implicit bias of gradient descent on separable data. *Journal of Machine Learning Research*, 19(70):1–57, 2018.
- [48] Salma Tarmoun, Guilherme Franca, Benjamin D Haeffele, and Rene Vidal. Understanding the dynamics of gradient flow in overparameterized linear models. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 10153–10161. PMLR, 18–24 Jul 2021.

- [49] James Townsend. Differentiating the singular value decomposition. Technical report, Technical Report 2016, <https://j-towns.github.io/papers/svd-derivative...>, 2016.
- [50] Guido Van Rossum and Fred L. Drake. *Python 3 Reference Manual*. CreateSpace, Scotts Valley, CA, 2009. ISBN 1441412697.
- [51] Aditya Vardhan Varre, Maria-Luiza Vladarean, Loucas Pillaud-Vivien, and Nicolas Flammarion. On the spectral bias of two-layer linear networks. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=FFdrXkm3Cz>.
- [52] Aditya Vardhan Varre, Maria-Luiza Vladarean, Loucas Pillaud-Vivien, and Nicolas Flammarion. On the spectral bias of two-layer linear networks. *Advances in Neural Information Processing Systems*, 36, 2024.
- [53] Zihan Wang and Arthur Jacot. Implicit bias of sgd in l_2 -regularized linear dnns: One-way jumps from high to low rank, 2023.
- [54] Blake Woodworth, Suriya Gunasekar, Jason D. Lee, Edward Moroshko, Pedro Savarese, Itay Golan, Daniel Soudry, and Nathan Srebro. Kernel and rich regimes in overparametrized models. In Jacob Abernethy and Shivani Agarwal, editors, *Proceedings of Thirty Third Conference on Learning Theory*, volume 125 of *Proceedings of Machine Learning Research*, pages 3635–3673. PMLR, 09–12 Jul 2020.
- [55] Liu Ziyin, Hongchao Li, and Masahito Ueda. Law of balance and stationary distribution of stochastic gradient descent. *arXiv preprint arXiv:2308.06671*, 2023.

Appendix A. Notations

Notation S_d, S_d^+, S_d^{++} denote the set of symmetric, positive semi-definite and positive definite matrices in $R^{d \times d}$. We use \odot to denote the Hadamard product. We use $\langle \cdot, \cdot \rangle$ to denote the inner product, i.e., $\langle u, v \rangle = u^\top v$ for vectors, and $\langle A, B \rangle = \text{Tr}(AB^\top)$ for matrices. I_d denotes the identity matrix of dimension d and $0_{p \times k}$ denote the matrix with all zero entries of dimension $p \times k$.

Appendix B. Further Related Work

Linear Networks. The study of two-layer linear networks has been explored extensively, particularly when optimized using gradient flow on the square loss, across various settings including zero-balance initialization and whitened data Braun et al. [9], Fukumizu [18], Saxe et al. [45, 46]. Early work by Saxe et al. [45, 46] elucidates the temporal changes in the singular values of the predictor, assuming decoupled dynamics and a specific data-dependent weight initialization. This condition is broadened by the analyses of Fukumizu [18] and Braun et al. [9], Tarmoun et al. [48], who apply solutions from a matrix Riccati equation to characterize the weights dynamics under full-rank network initialization. Furthermore, Gidel et al. [20] extends the existing framework by relaxing the whitened data assumption, conducting a perturbation analysis, and discussing the temporal evolution of the weight matrices' singular values. Additionally, Varre et al. [52] eliminates the need for zero-balanced and full-rank initializations. Their study provides detailed formulas for weight evolution as a function of the initial scale, also studies a simple version of a stochastic flow without the drift. Wang and Jacot [53] studied the implicit bias of SGD with ℓ_2 -regularization.

Matrix valued stochastic process and their eigenvalues. Stochastic process on the space of symmetric (or Hermitian) matrices and the evolution of their eigenvalues are well studied since Dyson [16]. These techniques were further developed by Bru [10, 11] to study perturbations of principal component analysis and the eigenvalues of Wishart processes. Graczyk and Małeckı [21], Norris et al. [41] applied SDE-based techniques to study the eigenvalues and eigenvectors of Brownian motion on ellipsoids.

Appendix C. Mechanism behind the low-rank phenomenon

In this section, we investigate the evolution of singular values under stochastic training to gain deeper insights into the low-rank phenomenon. To simplify the discussion, throughout the section we consider the case where $k = 1$ and for notational convenience, we let $\mathbf{W}_1 = \mathbf{W}$, $\mathbf{W}_2 = \mathbf{a}$. Additionally, we assume that $l \leq p$, however the results can be extended to any l .

Warm-up: Comparison with diagonal networks. Let $\mathbf{W} = \mathbf{U}\Sigma\mathbf{V}^\top$ be the singular value decomposition (assuming $l \leq p$). The predictor β can be expressed as

$$\mathbf{W}^\top \mathbf{a} = \mathbf{U}\Sigma\mathbf{V}^\top \mathbf{a} = \mathbf{U}[\boldsymbol{\sigma} \odot \mathbf{c}], \text{ where } \mathbf{c} = \mathbf{V}^\top \mathbf{a}.$$

This expression reveals a Hadamard product between $\boldsymbol{\sigma}$ and \mathbf{c} , reminiscent of diagonal networks which are widely studied to understand the nonconvex dynamics of gradient algorithms [43, 44, 54]. In the context of diagonal networks, SGD is known to provably induce sparsity in predictions. Similarly, for linear networks, SGF may induce sparsity in terms of the singular value σ . We next derive the SDE governing the evolution of the singular values Σ of the weight matrix to gain a clearer understanding of the low-rank phenomenon.

Scalar Regression. We assume that the data is isotropic, i.e., $X = I_p$. Under these conditions, the loss function for scalar regression can be written as

$$\mathcal{L}(\mathbf{W}, \mathbf{a}) = \frac{1}{2} \|y - \mathbf{W}\mathbf{a}\|^2. \quad (\text{C.1})$$

We train the above objective with SGF, formulated as follows,

$$d\mathbf{W} = (y - \mathbf{W}\mathbf{a})\mathbf{a}^\top dt + \sqrt{\eta\delta} d\mathbf{B}_t\mathbf{a}^\top; \quad d\mathbf{a} = \mathbf{W}^\top (y - \mathbf{W}\mathbf{a})dt + \sqrt{\eta\delta} \mathbf{W}^\top d\mathbf{B}_t. \quad (\text{C.2})$$

where \mathbf{B}_t is the standard Brownian motion in \mathbb{R}^p . For analytical convenience, we rescale the time $t \rightarrow t/\eta\delta$ and use the process $d\mathbf{X} = 1/\eta\delta(y - \mathbf{W}\mathbf{a})dt + d\mathbf{B}_t$. The SGF can then be rewritten as,

$$d\mathbf{W} = d\mathbf{X}\mathbf{a}^\top; \quad d\mathbf{a} = \mathbf{W}^\top d\mathbf{X}. \quad (\text{C.3})$$

Our focus is on understanding the evolution of the singular values of the matrix \mathbf{W} . This aim is facilitated by considering the symmetric matrix $\mathbf{M} = \mathbf{W}^\top \mathbf{W}$, whose eigenvalues are the squares of the singular values of \mathbf{W} . Taking the derivative of \mathbf{M} , we find

$$d\mathbf{M} = d\mathbf{W}^\top \mathbf{W} + \mathbf{W}^\top d\mathbf{W} + d\mathbf{W}^\top d\mathbf{W} = d\mathbf{X}^\top \mathbf{W} + \mathbf{W}^\top d\mathbf{X}\mathbf{a}^\top + p\mathbf{a}\mathbf{a}^\top dt. \quad (\text{C.4})$$

Note that $dxdy$ represents $d[x, y]$ for any continuous semi-martingales x, y [see, e.g., 29, chapter 3 for reference].

Eigenvalues of a matrix-valued stochastic process. We leverage tools from the study of eigenvalues of matrix-valued stochastic processes [10, 21] to derive the evolution of the eigenvalues of \mathbf{M} in the theorem that follows.

Theorem C.1 *Let $s_1 > \dots > s_l$ be the order of the eigenvalues of the matrix \mathbf{M} defined by Equation (C.4). Let the collision time for the eigenvalues be defined as*

$$\tau = \{\inf t : s_i(t) = s_j(t) \text{ for } 1 \leq i \neq j \leq l\}. \quad (\text{C.5})$$

For $t \leq \tau$, the eigenvalues are semi-martingales given by the solution of the following SDE

$$d(s_i) = p\mathbf{c}_i^2 dt + \sum_{\substack{j=1, \\ j \neq i}}^l \frac{s_i\mathbf{c}_j^2 + s_j\mathbf{c}_i^2}{s_i - s_j} dt + 2\sqrt{s_i\mathbf{c}_i^2} \left(d\tilde{\mathbf{X}} \right)_i \quad (\text{C.6})$$

where $\mathbf{c} = \mathbf{V}^\top \mathbf{a}$ and $\left(d\tilde{\mathbf{X}} \right)_i = 1/\eta\delta \left(\langle \mathbf{u}_i, y \rangle - \sqrt{s_i\mathbf{c}_i^2} \right) dt + d\varepsilon_i$ with \mathbf{u}_i being the i^{th} column of \mathbf{U} and $(\varepsilon_0, \dots, \varepsilon_{l-1})$ is the standard Brownian motion in \mathbb{R}^l . The evolution of \mathbf{c}_i and \mathbf{U} are presented in the appendix F.5.

This theorem can be interpreted as the stochastic counterpart to the evolution of eigenvalues previously described for linear networks by Arora et al. [5], Varre et al. [51]. The derivation of the eigenvalues is inspired by the work of [10].

The evolution of the eigenvalues features a key term highlighted in Equation (C.6) consisting of the sum of skew-symmetric elements $s_i\mathbf{c}_j^2 + s_j\mathbf{c}_i^2 / s_i - s_j$. For a pair of indices (i_0, j_0) with $i_0 < j_0$ and thus $s_{i_0} > s_{j_0}$, the term $s_{i_0}\mathbf{c}_{j_0}^2 + s_{j_0}\mathbf{c}_{i_0}^2 / s_{i_0} - s_{j_0}$ positively influences the evolution of the larger

eigenvalue ds_{i_0} and negatively affects the smaller eigenvalue ds_{j_0} . Therefore, this force is repulsive, driving the eigenvalues apart and increasing their gap. Another factor influencing the dynamics is the presence of Geometric Brownian motion, where the singular value σ_i multiplicatively influences the Brownian motion as $\sqrt{s_i c_i^2} (d\tilde{\mathbf{X}})_i$, similar to what is observed in diagonal linear networks (refer to the previous discussion for similarities). This effect tends to pull the singular values toward zero. Together with the fact that $(s_i, c_i) = (0, 0)$ represents a fixed point of the dynamics, these two forces collectively push redundant singular values toward zero.

To further understand the interplay of repulsive forces and geometric Brownian motion, we consider the evolution of the smaller singular value s_p for $l = p$. Using the Ito chain rule, we analyze the evolution of $\log s_p$, expressed as,

$$d(\log s_p) = p \frac{c_p^2}{s_p} dt + \frac{1}{s_p} \sum_{\substack{j=1, \\ j \neq p}}^p \frac{s_p c_j^2 + s_j c_p^2}{s_p - s_j} dt - 2 \frac{c_p^2}{s_p} + 2 \sqrt{\frac{c_p^2}{s_p}} (d\tilde{\mathbf{X}})_p.$$

Using that $s_p c_j^2 + s_j c_p^2 / s_p - s_j < -c_p^2$, for all indices j , the repulsive force accumulates to $-(p-1)(c_p^2/s_p)$ and the Ito correction term from the logarithm contributes an additional $-2(c_p^2/s_p)$ (the GBM component) thus offsetting the positive drift of $p(c_p^2/s_p)$. In the case of $l \neq p$, considering a polynomial x^α with an appropriate α would demonstrate similar behaviour. This discussion outlines the forces at play, yet a complete characterization of the solution of the SDE Equation (C.6) remains missing. Moreover, we have not established that the eigenvalues avoid a.s. collision, i.e., the explosion time $\tau_\infty = \infty$ which is in itself a significant challenge [10, 22].

A simplified two-vector problem. To enhance our understanding of the SDE governing the evolution of the eigenvalues detailed in Equation (C.6), we consider the large noise limit. In this scenario, the process described in Equation (C.3) simplifies to a purely noise-driven process without drift:

$$d\mathbf{W} = d\mathbf{B}_t \mathbf{a}^\top; \quad d\mathbf{a} = \mathbf{W}^\top d\mathbf{B}_t.$$

This SDE exhibits notable symmetry; allowing for an analysis using a matrix with sub-sampled columns. Let S be any subset of $1, \dots, l$, with $(\mathbf{w}_i)_{i=1}^l$ representing the columns of \mathbf{W} . We define $\mathbf{W}_S \in \mathbb{R}^{p \times |S|}$ as the subsampled matrix obtained by selecting columns \mathbf{w}_i where $i \in S$, and similarly, we define a subsampled vector \mathbf{a}_S by selecting the corresponding coordinates. The SDE restricted to the set S is structured as follows:

$$d\mathbf{W}_S = d\mathbf{B}_t \mathbf{a}_S^\top; \quad d\mathbf{a}_S = \mathbf{W}_S^\top d\mathbf{B}_t.$$

To demonstrate that the columns of \mathbf{W} align, we leverage the symmetry of the SDE by examining the restricted problem on every pair of rows $S = \{i, j\}$, and proving alignment within this subset. This approach leads us to consider the two vector problem ($l = 2$), where $\mathbf{W} = [\mathbf{w}_1 | \mathbf{w}_2]$ and $\mathbf{w}_1, \mathbf{w}_2 \in \mathbb{R}^p$, $\mathbf{a} \in \mathbb{R}^2$. We describe the behavior of the eigenvalues for this two-vector problem in the theorem below.

Theorem C.2 *In the large noise limit, let $s_0 > s_1$ be the eigenvalues of \mathbf{W} , the following properties hold, for $t \leq \tau$ defined by $\tau = \{\inf t : s_0(t) = s_1(t)\}$,*

- (a) s_0, s_1 are greater than zero almost surely,

(b) for $\alpha = (p - 3)/2$, $\mathbf{s}_0^{-\alpha}$ is a super-martingale while $\mathbf{s}_1^{-\alpha}$ is a sub-martingale.

This model for $l = 2$ mirrors the dynamics of the Wishart process studied by Bru [11], motivating the exploration of the evolution of an appropriately chosen exponent of $\mathbf{s}_0, \mathbf{s}_1$. The first part of the theorem arises from the fact that $\mathbf{s}_1^{-\alpha}\mathbf{s}_2^{-\alpha}$ is a local continuous martingale that cannot explode to infinity in finite time. The second part highlights a clear separation between the eigenvalues: one is a sub-martingale that consistently increases in expectation, while the other is a super-martingale that diminishes (note that the eigenvalues are raised to a negative power). This dynamic, coupled with the symmetry argument, suggests that for every pair of columns, there is a component that strengthens the alignment through its increases in expectation. Refer to App. F.6 for the proof.

Conclusion. In this section, we derive the SDE of eigenvalues for the matrix of parameters evolving under SGF. This derivation provides deeper insights into the mechanisms contributing to low-rank behavior. Specifically, repulsive forces drive the eigenvalues apart, while the geometric Brownian motion pulls them towards zero. These forces, unique to training with SGF, highlight the regularization effects of stochastic methods compared to gradient flow. However, fully characterizing the solution of this SDE remains a challenging open problem we let as future work.

Appendix D. Generalization to other settings

In this section, we generalize our results beyond the square loss and the label noise gradient flow. We consider the general framework of a loss function over the weight product $\mathbf{W}_1\mathbf{W}_2$ defined as

$$\mathcal{L}(\mathbf{W}_1, \mathbf{W}_2) = \widehat{\mathcal{L}}(\mathbf{W}_1\mathbf{W}_2) = \mathbb{E}_{(x,y)\sim\mathcal{D}}[\ell(\mathbf{W}_1\mathbf{W}_2; x, y)],$$

In this framework, the loss function ℓ combines the prediction loss directly with the parametrized model f_θ . This approach applies, for example, to classification problems using linear networks where ℓ might represent any classification loss and $f_\theta = \mathbf{W}_1\mathbf{W}_2$. It also directly extends to more complex architectures where $f_\theta = \sigma(\mathbf{W}_1\mathbf{W}_2)$ for an activation function σ , including settings like a self-attention layer with frozen value vectors. We denote the product by $\beta = \mathbf{W}_1\mathbf{W}_2$ noting it solely controls the loss. We investigate the evolution of the weight matrix determinant for a general loss across various algorithms, from gradient flow to gradient descent, and demonstrate that a similar separation occurs due to stochasticity.

Warm-up: Gradient flow. The gradient flow on the loss \mathcal{L} can be written as the following,

$$d\Theta = \Theta\mathbf{J}dt, \quad \text{where } \mathbf{J} = \begin{bmatrix} 0_{p\times p} & -\nabla\widehat{\mathcal{L}}(\beta) \\ -\nabla\widehat{\mathcal{L}}(\beta)^\top & 0_{k\times k} \end{bmatrix}. \quad (\text{D.1})$$

Following a similar proof as in Theorem 3.1, we obtain that $d(\Theta^\top\Theta) = 0$. For separable classification problem, the gradient flow converges to infinity [31, 47], hence, after appropriate rescaling, the layers are aligned, as shown by Ji and Telgarsky [31]. Next, we contrast this result with the outcomes observed in stochastic and discrete algorithms.

Continuous modelling of SGD. We consider the SGD algorithm with a batch size B . We denote the mini-batch version of the loss functions \mathcal{L} and $\widehat{\mathcal{L}}$ as \mathcal{L}_B and $\widehat{\mathcal{L}}_B$, respectively. The SGD update with stepsize η can be represented with the following block structure,

$$\Theta^{t+1} = \Theta^t - \eta \Theta^t \mathbf{J}^t - \eta \Theta^t \xi^t, \quad \text{where } \xi^t = \begin{bmatrix} 0_{p \times p} & -\left(\nabla \widehat{\mathcal{L}}(\beta) - \nabla \widehat{\mathcal{L}}_B(\beta)\right) \\ -\left(\nabla \widehat{\mathcal{L}}(\beta) - \nabla \widehat{\mathcal{L}}_B(\beta)\right)^\top & 0_{k \times k} \end{bmatrix}.$$

We denote the SGD noise as $g_t = \left(\nabla \widehat{\mathcal{L}}(\beta) - \nabla \widehat{\mathcal{L}}_B(\beta)\right)$ and the noise covariance as $\Sigma_t = \mathbb{E} \left[g^t (g^t)^\top \right]$ where the expectation is over all the minibatches. Following Li et al. [35], the SGD update can be modelled with the following SDE,

$$d\Theta = -\Theta \mathbf{J} dt - \sqrt{\eta} d\xi, \quad \text{where } d\xi = \begin{bmatrix} 0_{p \times p} & -\Sigma_t^{1/2} d\mathbf{B}_t \\ -\left(\Sigma_t^{1/2} d\mathbf{B}_t\right)^\top & 0_{k \times k} \end{bmatrix}. \quad (\text{D.2})$$

The main difference with SGF is that, in overparameterized problems, the noise covariance is time-varying and decreases to zero upon convergence. Using Theorem F.3, the evolution of the determinant of $\mathbf{M} = \Theta^\top \Theta$ is given by $d(\det(\mathbf{M})) = -\eta \det(\mathbf{M}) \text{Tr}(\Sigma(t)) dt$ and can be explicitly solved as

$$d(\det(\mathbf{M})(t)) = \det(\mathbf{M}(0)) \exp\left\{-\eta \int_0^t \text{Tr}(\Sigma(s)) ds\right\}.$$

Hence, the decay in the determinant is governed by the integral $\int_0^\infty \text{Tr}(\Sigma(t)) dt$ which is a stochastic quantity. $\text{Tr}(\Sigma(t))$ represents the strength of the stochastic noise, which, in over-parameterized regression, is proportional to the loss, i.e., $\text{Tr}(\Sigma(t)) \propto \mathcal{L}(\Theta)$ [43]. Therefore, the rate of decay in the determinant depends on $\int_0^\infty \mathcal{L}(\Theta(t)) dt$, with slower convergence leading to a simpler model at convergence, as observed in the case of diagonal networks by Pesme et al. [43]. The result above also holds for *non-separable* classification tasks where the noise of SGD drives the determinant to 0, a scenario not covered by the previous analysis of Ji and Telgarsky [31].

Discrete gradient algorithms. We can extend the previous results to discrete (possibly stochastic) gradient algorithm. Both algorithms can be written as

$$\Theta_{t+1} = \Theta_t (\mathbf{I}_{p+k} + \eta \mathbf{J}_t),$$

for stepsize η and \mathbf{J}_t the possibly stochastic block gradient matrix defined in Equation (D.1). In the context of discrete algorithms, the determinant is controlled by the following lemma.

Lemma D.1 *When $l = p + k$ and $\eta^2 \|\mathbf{J}_t\|_F^2 \leq 1$, the following property holds for the determinant,*

$$\|\det \Theta_{t+1}\| \leq \exp\left(-\frac{\eta^2}{2} \|\mathbf{J}_t\|_F^2\right) \|\det \Theta_t\|.$$

If the factor $\eta^2 \|\mathbf{J}_t\|_F^2 \leq 1$ at every iteration t , the determinant is reduced by the discrete step size. However, there is a tradeoff: the sum $S := \sum_{t=0}^\infty \eta^2 \|\mathbf{J}_t\|_F^2$ can be finite, indicating that

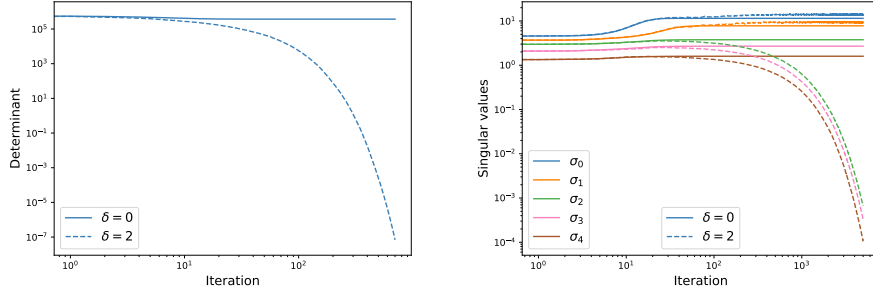


Figure 1: Evolution of the model characteristics for gradient flow ($\delta = 0$) and stochastic gradient flow ($\delta = 2$). Left: Determinant of \mathbf{M} . Right: Top-5 singular values of \mathbf{W}_1 .

it does not completely drive the determinant to zero. Increasing η to increase S might lead to instability and divergence. Furthermore, since $\|\mathbf{J}_t\|_F^2 \propto \mathcal{L}(\Theta_t)$, there is an additional tradeoff between convergence and the simplicity of the parameters. This illustrates how step sizes that produce non-convergent training loss patterns, such as the catapult effect [34] or the edge of stability mechanisms [13], can simplify the network’s parameters.

Appendix E. Experimental evidence

We consider a regression problem on synthetic data with $n = 1000$ samples of Gaussian data in \mathbb{R}^5 ($p = 5$) with labels in \mathbb{R}^2 ($k = 2$) generated by some ground truth $\beta \in \mathbb{R}^{5 \times 2}$, the width of the network is $l = 10$. We use Gaussian initialization of the network parameters with entries from $\mathcal{N}(0, 1)$. Experiments details can be found in the appendix G. In the left plot of Figure 1, we show the time evolution of the determinant of matrix \mathbf{M} . As suggested by theorems 3.1 and 3.2, in the case without label noise, $\det(\Theta^\top \Theta)$ stays constant, while with the Label Noise of intensity $\delta = 2$ it goes to zero with time. In the right plot of Figure 1, we demonstrate the time evolution of the top-5 singular values of the matrix \mathbf{W}_1 . Note that in the case of Gradient Flow all except the first k singular values (σ_0 and σ_1) stay at the same scale, while adding Label Noise forces smallest $d + l - k$ singular values (σ_2, σ_3 , and σ_4) to tend toward zero. Further experiments illustrate in Figure 2 the evolution of singular values of parameter matrix \mathbf{W}_1 when optimized with SGD, for classification tasks and with ReLU network. These results also confirm that the beneficial effects of stochasticity hold in these contexts.

Appendix F. Proofs

Theorem F.1 *For the gradient flow defined in Equation (2.4), the following property holds,*

$$d\left(\det\left(\Theta^\top \Theta\right)\right) = 0.$$

Hence, $\det(\Theta(t)^\top \Theta(t)) = \det(\Theta_0^\top \Theta_0)$, where $\Theta_0 = \Theta(0)$ is the initialisation at time $t = 0$.

First, we present a proof of this theorem, based on straightforward computations of the derivative of the determinant and the fact that the matrix \mathbf{J} has zero trace.

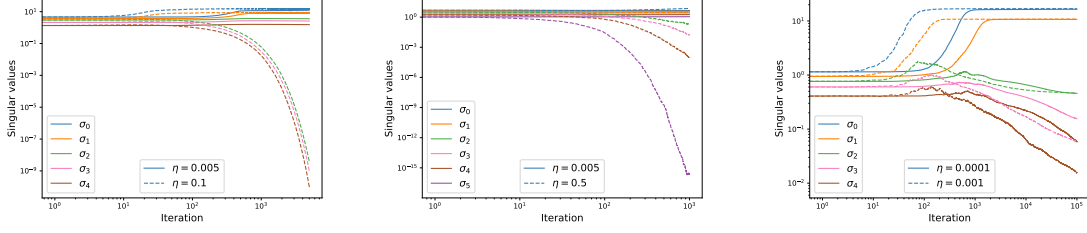


Figure 2: Evolution of the top-5 singular values of \mathbf{W}_1 for SGD with small and large stepsizes η . Left: Regression with MSE loss, linear network. Middle: Classification with logistic loss, linear network. Right: Regression with MSE loss, 2-layer ReLU network.

Proof Let $\mathbf{M} = \Theta^\top \Theta$. The dynamics of \mathbf{M} are governed by the ODE,

$$d\mathbf{M} = d\Theta^\top \Theta + \Theta^\top d\Theta = \Theta^\top \Theta \mathbf{J} dt + \mathbf{J} \Theta^\top \Theta dt = (\mathbf{M}\mathbf{J} + \mathbf{J}\mathbf{M})dt.$$

Using the gradient of the determinant given in Proposition F.2, the determinant of \mathbf{M} evolves as follows,

$$\begin{aligned} d(\det(\mathbf{M})) &= \langle \nabla \det(\mathbf{M}), d\mathbf{M} \rangle = \det(\mathbf{M}) \langle \mathbf{M}^{-1}, \mathbf{M}\mathbf{J} + \mathbf{J}\mathbf{M} \rangle dt, \\ &= \det(\mathbf{M}) \langle \mathbf{M}^{-1}, \mathbf{M}\mathbf{J} \rangle + \langle \mathbf{M}^{-1}, \mathbf{J}\mathbf{M} \rangle = 2\det(\mathbf{M}) \langle \mathbf{I}_{p+k}, \mathbf{J} \rangle = 2\det(\mathbf{M})\text{Tr}(\mathbf{J}). \end{aligned}$$

Given that $\text{Tr}(\mathbf{J}) = 0$, it follows that $d(\det(\mathbf{M})) = 0$. ■

Proposition F.2 For any matrix M in S_d , the first two derivatives of the determinant of M , denoted by $\det(M)$ are the following

- (i) $\nabla \det(M) = \det(M)M^{-1}$
- (ii) For $1 \leq a, b, k, l \leq d$, the second order partial derivative is given by

$$\frac{\partial^2 \det(M)}{\partial M_{ab} \partial M_{kl}} = \det(M) [(M^{-1})_{ba}(M^{-1})_{lk} - (M^{-1})_{bk}(M^{-1})_{la}] \quad (\text{F.1})$$

Theorem F.3 For a stochastic process given by the SDE,

$$d\Theta = \Theta [\mathbf{J}dt + d\xi] \quad (\text{F.2})$$

with $\text{Tr} \mathbf{J} = \text{Tr} \xi = 0$, the determinant of the $\mathbf{M} = \Theta^\top \Theta$ evolves as

$$d(\det(\mathbf{M})) = -\det(\mathbf{M})\text{Tr}[d\xi d\xi]. \quad (\text{F.3})$$

Proof First, we compute the evolution of $\mathbf{M} = \Theta^\top \Theta$ using the Ito's product rule,

$$d\mathbf{M} = d(\Theta^\top \Theta) = d\Theta^\top \Theta + \Theta^\top d\Theta + d\Theta^\top d\Theta$$

The last term is interpreted as a derivative of the finite variation and it should be computed using dt . $(d\mathbf{B}_t)_{ij} = 0$ and $(d\mathbf{B}_t)_{ij} \cdot (d\mathbf{B}_t)_{kl} = \delta_{i=k \wedge j=l} dt$. Using Eq. (2.6),

$$\begin{aligned} d\mathbf{M} &= [\mathbf{J}dt + d\xi] \Theta^\top \Theta + \Theta^\top \Theta [\mathbf{J}dt + d\xi] + d\xi \Theta^\top \Theta d\xi, \\ &= \mathbf{J}Mdt + \mathbf{M}Jdt + d\xi M d\xi + d\xi M + M d\xi. \end{aligned}$$

Using the Ito chain rule, we can compute the evolution of determinant as following,

$$d(\det(\mathbf{M})) = \langle \nabla \det(\mathbf{M}), d\mathbf{M} \rangle + \frac{1}{2} \sum_{a,b,k,l} \frac{\partial^2 \det(\mathbf{M})}{\partial M_{ab} \partial M_{kl}} dM_{ab} dM_{kl},$$

The first term is

$$\begin{aligned} \langle \nabla \det(\mathbf{M}), d\mathbf{M} \rangle &= \det(\mathbf{M}) \langle \mathbf{M}^{-1}, \mathbf{J}Mdt + \mathbf{M}Jdt + d\xi M d\xi + d\xi M + M d\xi \rangle, \\ &= 2 \det(\mathbf{M}) \langle \mathbf{I}_{p+k}, \mathbf{J} \rangle dt + 2 \det(\mathbf{M}) \langle \mathbf{I}_{p+k}, d\xi \rangle + \langle \mathbf{M}^{-1}, d\xi M \xi \rangle \end{aligned}$$

Using the property that $\text{Tr}(\mathbf{J}) = \text{Tr}(d\xi) = 0$. We get that $\langle \nabla \det(\mathbf{M}), d\mathbf{M} \rangle = \langle \mathbf{M}^{-1}, d\xi M \xi \rangle$.

For the second term

$$\begin{aligned} \frac{1}{2} \sum_{a,b,k,l} \frac{\partial^2 \det(\mathbf{M})}{\partial M_{ab} \partial M_{kl}} &= \frac{1}{2} \sum_{a,b,k,l} \det \mathbf{M} [(\mathbf{M}^{-1})_{ba} (\mathbf{M}^{-1})_{lk} - (\mathbf{M}^{-1})_{bk} (\mathbf{M}^{-1})_{la}] dM_{ab} dM_{kl}, \\ &= \frac{\det(\mathbf{M})}{2} \sum_{a,b,k,l} [(\mathbf{M}^{-1})_{ba} (\mathbf{M}^{-1})_{lk}] dM_{ab} dM_{kl} - \sum_{a,b,k,l} [(\mathbf{M}^{-1})_{bk} (\mathbf{M}^{-1})_{la}] dM_{ab} dM_{kl}, \end{aligned}$$

Rearranging the terms in the summation, we get,

$$\begin{aligned} \sum_{a,b,k,l} [(\mathbf{M}^{-1})_{ba} (\mathbf{M}^{-1})_{lk}] dM_{ab} dM_{kl} &= \sum_{a,b,k,l} [(\mathbf{M}^{-1})_{ba} dM_{ab}] [(\mathbf{M}^{-1})_{lk} dM_{kl}], \\ &= \sum_{b,l} \left[\sum_a (\mathbf{M}^{-1})_{ba} dM_{ab} \right] \left[\sum_k (\mathbf{M}^{-1})_{lk} dM_{kl} \right], \\ &= \sum_{b,l} (\mathbf{M}^{-1} d\mathbf{M})_{bb} (\mathbf{M}^{-1} d\mathbf{M})_{ll} = \sum_b (\mathbf{M}^{-1} d\mathbf{M})_{bb} \sum_l (\mathbf{M}^{-1} d\mathbf{M})_{ll}, \\ &= \text{Tr}(\mathbf{M}^{-1} d\mathbf{M}) \text{Tr}(\mathbf{M}^{-1} d\mathbf{M}). \end{aligned}$$

Similarly for the other term, we get,

$$\begin{aligned} \sum_{a,b,k,l} [(\mathbf{M}^{-1})_{bk} (\mathbf{M}^{-1})_{la}] dM_{ab} dM_{kl} &= \sum_{a,b,k,l} [(\mathbf{M}^{-1})_{bk} dM_{kl}] [(\mathbf{M}^{-1})_{la} dM_{ab}], \\ &= \sum_{b,l} \left[\sum_a (\mathbf{M}^{-1})_{ba} dM_{al} \right] \left[\sum_k (\mathbf{M}^{-1})_{bk} dM_{kl} \right], \\ &= \sum_b \left[\sum_l (\mathbf{M}^{-1} d\mathbf{M})_{bl} (\mathbf{M}^{-1} d\mathbf{M})_{lb} \right] = \sum_b (\mathbf{M}^{-1} d\mathbf{M} \mathbf{M}^{-1} d\mathbf{M})_{bb}, \\ &= \text{Tr}[(\mathbf{M}^{-1} d\mathbf{M}) (\mathbf{M}^{-1} d\mathbf{M})]. \end{aligned}$$

Note that the diffusion part of $\mathbf{M}^{-1}d\mathbf{M}$ is $d\xi + \mathbf{M}^{-1}d\xi\mathbf{M}$. Using this

$$\text{Tr}(\mathbf{M}^{-1}d\mathbf{M})\text{Tr}(\mathbf{M}^{-1}d\mathbf{M}) = \text{Tr}[d\xi + \mathbf{M}^{-1}d\xi\mathbf{M}]\text{Tr}[d\xi + \mathbf{M}^{-1}d\xi\mathbf{M}] = 0,$$

as $\text{Tr} d\xi = 0$. For the other term,

$$\begin{aligned} \text{Tr}[(\mathbf{M}^{-1}d\mathbf{M})(\mathbf{M}^{-1}d\mathbf{M})] &= \text{Tr}[(d\xi + \mathbf{M}^{-1}d\xi\mathbf{M})(d\xi + \mathbf{M}^{-1}d\xi\mathbf{M})], \\ &= 2\text{Tr}[d\xi d\xi] + 2\text{Tr}[\mathbf{M}^{-1}d\xi\mathbf{M}d\xi]. \end{aligned}$$

Putting everything together, we get,

$$\frac{1}{2} \sum_{a,b,k,l} \frac{\partial^2 \det(\mathbf{M})}{\partial \mathbf{M}_{ab} \partial \mathbf{M}_{kl}} = -\det \mathbf{M} (\text{Tr}[d\xi d\xi] + \text{Tr}[\mathbf{M}^{-1}d\xi\mathbf{M}d\xi])$$

which gives us

$$d(\det(\mathbf{M})) = -\det(\mathbf{M}) \text{Tr}[d\xi d\xi].$$

■

Lemma F.4 *When $l = p + k$ and $\eta^2 \|\mathbf{J}_t\|_F^2 \leq 1$, the following property holds for the determinant,*

$$\|\det \Theta_{t+1}\| \leq \exp\left(-\frac{\eta^2}{2} \|\mathbf{J}_t\|_F^2\right) \|\det \Theta_t\|.$$

Proof Note that because of the block structure of the matrix \mathbf{J}_t , its nonzero eigenvalues come in \pm pairs: $\pm\sigma_1, \dots, \pm\sigma_m$, moreover, since \mathbf{J}_t is symmetric, singular values of \mathbf{J}_t are the absolute values of eigenvalues, i.e. $\sigma_1, \dots, \sigma_m$. Then, the determinant of Θ_{t+1} can be written as the following,

$$\det \Theta_{t+1} = \det \Theta_t \det(\mathbf{I}_{p+k} + \eta \mathbf{J}_t) = \det \Theta_t \prod_{i=1}^m (1 - \eta^2 \sigma_i^2).$$

Using that $1 - x^2 \leq e^{-x^2}$ for all x , we can estimate

$$\prod_{i=1}^m (1 - \eta^2 \sigma_i^2) \leq \exp\left(-\eta^2 \sum_{i=1}^m \sigma_i^2\right) = \exp\left(-\frac{\eta^2}{2} \|\mathbf{J}_t\|_F^2\right).$$

We obtain the required inequality by observing that $\prod_{i=1}^m (1 - \eta^2 \sigma_i^2) = \left\| \prod_{i=1}^m (1 - \eta^2 \sigma_i^2) \right\|$ since each term $1 - \eta^2 \sigma_i^2 \geq 0$ when $\eta^2 \|\mathbf{J}_t\|_F^2 < 1$. ■

Theorem F.5 *Let $s_1 > \dots > s_l$ be the order of the eigenvalues of the matrix \mathbf{M} defined by Equation (C.4). Let the collision time for the eigenvalues be defined as*

$$\tau = \{\inf t : \mathbf{s}_i(t) = \mathbf{s}_j(t) \text{ for } 1 \leq i \neq j \leq l\}. \quad (\text{F.4})$$

For $t \leq \tau$, the eigenvalues are semi-martingales given by the solution of the following SDE

$$d(\mathbf{s}_i) = p\mathbf{c}_i^2 dt + \sum_{\substack{j=1, \\ j \neq i}}^l \frac{\mathbf{s}_i \mathbf{c}_j^2 + \mathbf{s}_j \mathbf{c}_i^2}{\mathbf{s}_i - \mathbf{s}_j} dt + 2\sqrt{\mathbf{s}_i \mathbf{c}_i^2} (d\tilde{\mathbf{X}})_i \quad (\text{F.5})$$

where $(d\tilde{\mathbf{X}})_i = 1/\eta\delta (\langle \mathbf{u}_i, y \rangle - \sqrt{\mathbf{s}_i \mathbf{c}_i^2}) dt + d\varepsilon_i$ with \mathbf{u}_i being the i^{th} column of \mathbf{U} and $(\varepsilon_0, \dots, \varepsilon_{l-1})$ is the standard Brownian motion in \mathbb{R}^l . The evolution of \mathbf{c}_i and \mathbf{U} are presented in the appendix.

Proof The proof follows the approach of Bru [10]. Let $\mathbf{W} = \mathbf{U}\Sigma\mathbf{V}^\top$ be the singularvalue decomposition (see Def.H.1 involved with $r = l$ and $l < p$ and it will be the rank). Our focus is on understanding the evolution of the singular values and singular vectors of the matrix \mathbf{W} . To derive the evolution of Σ, \mathbf{V} we can consider the eigen values and eigenvectors of the PSD matrix process \mathbf{M} . Note that $\mathbf{M} = \mathbf{V}\Sigma^\top\mathbf{V}^\top$, let $\mathbf{D} = \Sigma^2$.

Evolution of \mathbf{D} and \mathbf{V} Taking the derivative of \mathbf{M} , we find

$$d\mathbf{M} = d\mathbf{W}^\top\mathbf{W} + \mathbf{W}^\top d\mathbf{W} + d\mathbf{W}^\top d\mathbf{W} = a d\mathbf{X}^\top\mathbf{W} + \mathbf{W}^\top d\mathbf{X}a^\top + p\mathbf{a}\mathbf{a}^\top dt. \quad (\text{F.6})$$

We invoke the theorem H.2 we derived to give the eigenvalues of any matrix valued stochastic process. Note that $\mathbf{V}\mathbf{V}^\top = \mathbf{I}_l$, so some terms of the computation are not required.

$$d\mathbf{D} = \mathbf{I} \odot \tilde{\mathbf{N}} dt + \mathbf{I} \odot d\tilde{\mathbf{M}} dt + \mathbf{I} \odot (d\tilde{\mathbf{M}} (\mathbf{S} \odot d\tilde{\mathbf{M}})).$$

and the evolution of the eigenvectors,

$$d\mathbf{V} = \mathbf{V} (\mathbf{Q}_{\parallel} dt + \mathbf{S} \odot (\tilde{\mathbf{N}} dt + d\tilde{\mathbf{M}}))$$

where you define,

$$\mathbf{Q}_{\parallel} = \frac{\mathbf{I} \odot \left[(\mathbf{S} \odot d\tilde{\mathbf{M}}) (\mathbf{S} \odot d\tilde{\mathbf{M}}) \right]}{2} - \mathbf{S} \odot \left[(\mathbf{S} \odot d\tilde{\mathbf{M}}) [d\tilde{\mathbf{M}} \odot \mathbf{I}] \right] + \mathbf{S} \odot (d\tilde{\mathbf{M}} (\mathbf{S} \odot d\tilde{\mathbf{M}}))$$

where the matrix \mathbf{S} is given by

$$\begin{aligned} \mathbf{S}_{ij} &= \begin{cases} 0 & \text{if } i = j, \\ (\mathbf{s}_j - \mathbf{s}_i)^{-1} & \text{o.w.} \end{cases} \\ \tilde{\mathbf{N}} &= \mathbf{V}^\top (p\mathbf{a}\mathbf{a}^\top) \mathbf{V} = p\mathbf{c}\mathbf{c}^\top. \\ d\tilde{\mathbf{M}} &= \mathbf{V}^\top \left[a d\mathbf{X}^\top\mathbf{W} + \mathbf{W}^\top d\mathbf{X}a^\top \right] \mathbf{V}, \\ &= c d\mathbf{X}^\top \mathbf{U} \Sigma + \Sigma \mathbf{U}^\top d\mathbf{X} c^\top. \end{aligned}$$

Note that $\Sigma = \text{diag}((\sigma_0, \dots, \sigma_{l-1}))$ where $\sigma_0 > \sigma_1 \dots > \sigma_{l-1}$. Let $\mathbf{D} = \Sigma^2$ and denote the entires of \mathbf{D} as following, $\mathbf{D} = \text{diag}((\mathbf{s}_0, \dots, \mathbf{s}_{p-1}))$. Note that

$$\begin{aligned} \mathbf{U}^\top d\mathbf{X} &= \mathbf{U}^\top \left(\frac{1}{\eta\delta} (y - \mathbf{W}\mathbf{a}) dt + d\mathbf{B}_t \right), \\ &= \frac{1}{\eta\delta} \left[\mathbf{U}^\top y - \Sigma \mathbf{c} \right] dt + \mathbf{U}^\top d\mathbf{B}_t. \end{aligned}$$

Using Levy's characterization $\mathbf{U}^\top d\mathbf{B}_t$ is a Brownian motion in \mathbb{R}^l , lets call that $d\tilde{\mathbf{B}}_t$. The diffusion part of $d\tilde{\mathbf{M}}$ (say $d\mathbf{F}$)

$$\begin{aligned} d\mathbf{F} &= \Sigma \mathbf{V}^\top d\mathbf{B}_t \mathbf{c}^\top + \mathbf{c} d\mathbf{B}_t^\top \mathbf{V} \Sigma, \\ &= \left(\boldsymbol{\sigma} \odot d\tilde{\mathbf{B}}_t \right) \mathbf{c}^\top + \mathbf{c} \left(\boldsymbol{\sigma} \odot d\tilde{\mathbf{B}}_t \right)^\top \\ &= d\mathbf{m}_t \mathbf{c}^\top + \mathbf{c} d\mathbf{m}_t^\top \end{aligned}$$

where $d\mathbf{m}_t \stackrel{\text{def}}{=} (\boldsymbol{\sigma} \odot d\tilde{\mathbf{B}}_t)$. We are required to compute $d\mathbf{F}(\mathbf{S} \odot d\mathbf{F})$ to compute the evolution of eigenvalues. Using the lemma H.4, we get

$$\begin{aligned} d\mathbf{F}(\mathbf{S} \odot d\mathbf{F}) &= \mathbf{c} \mathbf{S}^\top \mathbf{S} \text{diag}(\mathbf{c}) dt - \mathbf{D} \text{diag}(\mathbf{S} \text{diag}(\mathbf{c}) \mathbf{c}) dt + \mathbf{D} \text{diag}(\mathbf{c}) \mathbf{S} \text{diag}(\mathbf{c}) dt, \\ \mathbf{I} \odot [d\mathbf{F}(\mathbf{S} \odot d\mathbf{F})] &= \mathbf{I} \odot \left[\mathbf{c} \mathbf{S}^\top \mathbf{S} \text{diag}(\mathbf{c}) dt - \mathbf{D} \text{diag}(\mathbf{S} \text{diag}(\mathbf{c}) \mathbf{c}) dt \right] \end{aligned}$$

The element wise computation of this term gives the required result for evolution of eigenvalues.

Evolution of \mathbf{c} . Note that $\mathbf{c} = \mathbf{V}^\top \mathbf{a}$. Computing the derivative using the Ito's product rule, we get,

$$\begin{aligned} d\mathbf{V}^\top \mathbf{a} &= \mathbf{V}^\top d\mathbf{a} + d\mathbf{V}^\top \mathbf{a} + d\mathbf{V}^\top d\mathbf{a}, \\ &= \mathbf{V}^\top d\mathbf{a} + d\mathbf{V}^\top \mathbf{V} \mathbf{V}^\top \mathbf{a} + d\mathbf{V}^\top \mathbf{V} \mathbf{V}^\top d\mathbf{a}, \\ d\mathbf{V}^\top \mathbf{V} &= \left[\left(\mathbf{Q}_{\parallel}^\top dt - \mathbf{S} \odot d\mathbf{X} \right) \right], \\ \mathbf{V}^\top d\mathbf{a} &= \mathbf{V}^\top \mathbf{W}^\top d\mathbf{B}_t + \frac{1}{\eta \delta} \left[\mathbf{U}^\top \mathbf{y} - \Sigma \mathbf{c} \right] dt = \Sigma d\tilde{\mathbf{B}}_t = d\mathbf{m}_t + \frac{1}{\eta \delta} \left[\mathbf{U}^\top \mathbf{y} - \Sigma \mathbf{c} \right] dt, \\ d\mathbf{V}^\top \mathbf{V} \mathbf{V}^\top d\mathbf{a} &= -(\mathbf{S} \odot d\mathbf{F}) d\mathbf{m}_t. \\ d\mathbf{V}^\top \mathbf{V} \mathbf{V}^\top d\mathbf{a} &= \left[\left(\mathbf{Q}_{\parallel}^\top dt - \mathbf{S} \odot (\tilde{\mathbf{N}} dt + d\tilde{\mathbf{M}}) \right) \right] \mathbf{c} \end{aligned}$$

Using the lemma H.6, H.5, H.4 and computing the element wise summation, we get the following evolution for $d\mathbf{c}$

$$\begin{aligned} d\mathbf{c}_i &= -\frac{1}{2} \sum_{j=1}^l \mathbf{S}_{ij} (\mathbf{s}_i \mathbf{c}_j^2 + \mathbf{s}_j \mathbf{c}_i^2) dt - \mathbf{c}_i \sum_{j=1}^l (\mathbf{S}_{ij} \mathbf{c}_j^2) \left(\sum_{k \neq i, j} \mathbf{s}_k \mathbf{S}_{ki} \right) \\ &\quad - (p-2) \mathbf{c}_i \sum_{j=1}^l \mathbf{S}_{ij} \mathbf{c}_i^2 dt - \sum_{j=1}^l \mathbf{S}_{ij} \mathbf{s}_j dt, \\ &\quad + \sigma_i (\mathbf{U}^\top d\mathbf{X})_i \left(1 - \sum_{j=1}^l \mathbf{S}_{ij} \mathbf{c}_j^2 \right) - \mathbf{c}_i \sum_j \mathbf{S}_{ij} \sigma_j \mathbf{c}_j (\mathbf{U}^\top d\mathbf{X})_j \end{aligned}$$

Evolution of \mathbf{U} . To compute the evolution of \mathbf{U} , we invoke the theorem H.2 on the evolution of $\mathbf{W} \mathbf{W}^\top = \mathbf{U} \mathbf{D} \mathbf{U}^\top$. We ignore it here as it does not have much consequence on our results. \blacksquare

Theorem F.6 *In the large noise limit, when $l = 2$, the following properties hold, for $t \leq \tau$,*

(a) s_0, s_1 are greater than zero almost surely.

(b) for $\alpha = (p-3)/2$, $s_0^{-\alpha}$ is a super-martingale while $s_1^{-\alpha}$ is a sub-martingale.

Proof First, note that in the large noise limit with $l = 2$, the evolution of the eigenvalues is expressed as

$$d(s_0) = pc_0^2 dt + \frac{s_0 c_1^2 + s_1 c_0^2}{s_0 - s_1} dt + 2\sqrt{s_0 c_0^2} (d\tilde{\mathbf{B}}_t)_0, \quad (\text{F.7})$$

$$d(s_1) = pc_1^2 dt - \frac{s_0 c_1^2 + s_1 c_0^2}{s_0 - s_1} dt + 2\sqrt{s_1 c_1^2} (d\tilde{\mathbf{B}}_t)_1. \quad (\text{F.8})$$

Using the Ito chain rule, for the evolution of $s_0^{-\alpha}$ we can write

$$\begin{aligned} d(s_0^{-\alpha}) &= \frac{\partial(s_0^{-\alpha})}{\partial s_0} \left(pc_0^2 dt + \frac{s_0 c_1^2 + s_1 c_0^2}{s_0 - s_1} dt + 2\sqrt{s_0 c_0^2} (d\tilde{\mathbf{B}}_t)_0 \right) + \frac{1}{2} \frac{\partial^2(s_0^{-\alpha})}{\partial^2 s_0} \left(2\sqrt{s_0 c_0^2} \right)^2 dt \\ &= -\alpha s_0^{-\alpha-1} \left(pc_0^2 dt + \frac{s_0 c_1^2 + s_1 c_0^2}{s_0 - s_1} dt - 2(\alpha+1)c_0^2 dt + 2\sqrt{s_0 c_0^2} (d\tilde{\mathbf{B}}_t)_0 \right) \\ &= -\alpha s_0^{-\alpha-1} \left(c_0^2 dt + \frac{s_0 c_1^2 + s_1 c_0^2}{s_0 - s_1} dt + 2\sqrt{s_0 c_0^2} (d\tilde{\mathbf{B}}_t)_0 \right), \end{aligned}$$

analogously

$$d(s_1^{-\alpha}) = -\alpha s_1^{-\alpha-1} \left(c_1^2 dt - \frac{s_0 c_1^2 + s_1 c_0^2}{s_0 - s_1} dt + 2\sqrt{s_1 c_1^2} (d\tilde{\mathbf{B}}_t)_1 \right),$$

and finally for $s_0^{-\alpha} s_1^{-\alpha}$

$$\begin{aligned} d(s_0^{-\alpha} s_1^{-\alpha}) &= d(s_0^{-\alpha}) s_1^{-\alpha} + s_0^{-\alpha} d(s_1^{-\alpha}) + d(s_0^{-\alpha}) d(s_1^{-\alpha}) \\ &= -\alpha s_0^{-\alpha-1} s_1^{-\alpha} \left(c_0^2 dt + \frac{s_0 c_1^2 + s_1 c_0^2}{s_0 - s_1} dt + 2\sqrt{s_0 c_0^2} (d\tilde{\mathbf{B}}_t)_0 \right) \\ &\quad -\alpha s_0^{-\alpha} s_1^{-\alpha-1} \left(c_1^2 dt - \frac{s_0 c_1^2 + s_1 c_0^2}{s_0 - s_1} dt + 2\sqrt{s_1 c_1^2} (d\tilde{\mathbf{B}}_t)_1 \right). \end{aligned}$$

Now, we can show that the drift term in the SDE that describes the dynamics of $s_0^{-\alpha} s_1^{-\alpha}$ is zero, which gives us the first part of the result by McKean's argument [39],

$$\begin{aligned} &-\alpha s_0^{-\alpha-1} s_1^{-\alpha-1} \left(s_1 c_0^2 + s_1 \frac{s_0 c_1^2 + s_1 c_0^2}{s_0 - s_1} + s_0 c_1^2 + s_0 \frac{s_0 c_1^2 + s_1 c_0^2}{s_0 - s_1} \right) \\ &= -\alpha s_0^{-\alpha-1} s_1^{-\alpha-1} \left(s_1 c_0^2 + s_0 c_1^2 + \frac{s_0 s_1 c_1^2 + s_1^2 c_0^2 - s_0^2 c_1^2 + s_0 s_1 c_0^2}{s_0 - s_1} \right) \\ &= -\alpha s_0^{-\alpha-1} s_1^{-\alpha-1} \left(s_1 c_0^2 + s_0 c_1^2 + \frac{(s_1 - s_0)(s_0 c_1^2 + s_1 c_0^2)}{s_0 - s_1} \right) = 0. \end{aligned}$$

The second part is obtained by noticing that

$$\begin{aligned} c_0^2 + \frac{s_0 c_1^2 + s_1 c_0^2}{s_0 - s_1} &= \frac{s_0 (c_1^2 + c_0^2)}{s_0 - s_1} \geq 0, \\ c_1^2 - \frac{s_0 c_1^2 + s_1 c_0^2}{s_0 - s_1} &= -\frac{s_1 (c_1^2 + c_0^2)}{s_0 - s_1} \leq 0, \end{aligned}$$

and hence the drift term of $d(\mathbf{s}_0^{-\alpha})$ is not positive, while the drift term of $d(\mathbf{s}_1^{-\alpha})$ is not negative. ■

Appendix G. Experiment details

In all the graphs we plot the values averaged on the 20 runs with different random seeds as well as the 95% confidence interval (lightly colored). To numerically emulate GF (Figure 1), we set a stepsize of $1e^{-6}$ in numerical simulation.

In the further experiments, we study the behaviour of the linear network for regression with the same synthetic data and same network initialization as in previous experiment. As seen in the left plot of the Figure 2, when the stepsize is large ($\eta = 0.1$), singular values exhibit behavior similar to the case of LNGF, while with the small stepsize ($\eta = 0.005$) the evolution of singular values is closer to GF case. Next, we examine the effect of SGD in the case of classification task with logistic loss, as illustrated in the middle plot of the Figure 2. We consider synthetic data with $n = 1000$ samples of Gaussian data in \mathbb{R}^5 ($d = 5$) constituting two clusters corresponding to two classes ($k = 1$). Note that larger stepsize ($\eta = 0.5$) in this case also forces the smallest singular value to tend to zero, however the effect is not so dramatic for the rest of singular values. Additionally, we study the 2-layer ReLU network optimized with SGD on the same regression task as before. As seen in the right plot of the Figure 2, the decrease of the last singular value σ_4 is much slower than in the case of the linear network, however, the larger stepsize still facilitates divergence of k largest (σ_0 and σ_1) and $p - k$ smallest (σ_2, σ_3 and σ_4) singular values.

All experiments are implemented with Python 3 [50] under PSF license, NumPy [25] under BSD license, and PyTorch [42] under BSD-3-Clause license.

The experiments were run on a Intel i5-8250U, 8-GB RAM, with OS Ubuntu 20.04.6.

Appendix H. Supplementary material

H.1. Notations and preliminary definitions

Definition H.1 (Eigen decomposition and Singular Value decomposition) *We discuss the eigen value decomposition for a symmetric square matrix, and the singular value decomposition for any matrix is defined as the following*

- (a) **Eigen decomposition.** *For any rank r matrix $\mathbf{R} \in S_p$, $\mathbf{R} = \mathbf{V}\mathbf{D}\mathbf{V}^\top$ is the eigen decomposition, where $\mathbf{V} \in \mathbb{R}^{p \times r}$, $\mathbf{D} \in \mathbb{R}^{r \times r}$, \mathbf{D} is a diagonal matrix and $\mathbf{V}^\top \mathbf{V} = \mathbf{I}_r$, however, $\mathbf{V}\mathbf{V}^\top$ is not necessarily an identity matrix unless $r = p$.*
- (b) **Singular Value Decomposition.** *For any rank r matrix $\mathbf{W} \in \mathbb{R}^{p \times l}$, $\mathbf{W} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$, where $\mathbf{U} \in \mathbb{R}^{p \times r}$, $\mathbf{V} \in \mathbb{R}^{l \times r}$, $\mathbf{\Sigma} \in \mathbb{R}^{r \times r}$, $\mathbf{\Sigma}$ is a diagonal matrix and $\mathbf{U}^\top \mathbf{U} = \mathbf{V}^\top \mathbf{V} = \mathbf{I}_r$, however the $\mathbf{U}\mathbf{U}^\top$ and $\mathbf{V}\mathbf{V}^\top$ are not necessarily identity unless $r = p$ or $r = l$ respectively.*

H.2. Eigenvalues of matrix valued stochastic process

Theorem H.2 *For a matrix-valued stochastic process on S_{p+k}^{++} ,*

$$d\mathbf{R} = \mathbf{N}dt + d\mathbf{M}$$

where $d\mathbf{M}$ is a local martingale process. Let $R = \mathbf{V}\mathbf{D}\mathbf{V}^\top$ is the eigenvalue decomposition of the process, the evolution of eigenvalues satisfy the SDE for time t less than the collision time,

$$d\mathbf{D} = \mathbf{I} \odot \tilde{\mathbf{N}} dt + \mathbf{I} \odot d\tilde{\mathbf{M}} dt + \mathbf{I} \odot \left(d\tilde{\mathbf{M}} \left(\mathbf{S} \odot d\tilde{\mathbf{M}} \right) \right) + \mathbf{D}^{-1} \odot \left(\mathbf{V}^\top d\mathbf{R} \left(\mathbf{I} - \mathbf{V}\mathbf{V}^\top \right) d\mathbf{R}\mathbf{V} \right).$$

where \mathbf{S} is defined as per Eq. H.1 and $d\tilde{\mathbf{M}} = \mathbf{V}^\top d\mathbf{M}\mathbf{V}$, $\tilde{\mathbf{N}} = \mathbf{V}^\top \mathbf{N}\mathbf{V}$. The evolution of the eigenvectors,

$$d\mathbf{V} = \mathbf{V} \left(\mathbf{Q}_\parallel dt + \mathbf{S} \odot d\mathbf{F} \right) + (\mathbf{I} - \mathbf{V}\mathbf{V}^\top) \left(\mathbf{Q}_\perp dt + d\mathbf{R} \mathbf{V}\mathbf{D}^{-1} \right).$$

where you define,

$$\begin{aligned} \mathbf{Q}_\parallel &= \frac{\mathbf{I} \odot \left[\left(\mathbf{S} \odot d\tilde{\mathbf{M}} \right) \left(\mathbf{S} \odot d\tilde{\mathbf{M}} \right) \right]}{2} - \frac{\mathbf{I} \odot \left[\mathbf{D}^{-1} \mathbf{V}^\top d\mathbf{R} \left(\mathbf{I} - \mathbf{V}\mathbf{V}^\top \right) d\mathbf{R}\mathbf{V}\mathbf{D}^{-1} \right]}{2} \\ &\quad - \mathbf{S} \odot \left[\left(\mathbf{S} \odot d\tilde{\mathbf{M}} \right) \left[d\tilde{\mathbf{M}} \odot \mathbf{I} \right] \right] + \mathbf{S} \odot \left(d\tilde{\mathbf{M}} \left(\mathbf{S} \odot d\tilde{\mathbf{M}} \right) \right) \\ &\quad + \mathbf{S} \odot \left(\mathbf{V}^\top d\mathbf{R} \left(\mathbf{I} - \mathbf{V}\mathbf{V}^\top \right) d\mathbf{R}\mathbf{V}\mathbf{D}^{-1} \right), \\ \mathbf{Q}_\perp &= \left[d\mathbf{R}\mathbf{V}\mathbf{D}^{-1} \right] \left[\left[\mathbf{S} \odot d\tilde{\mathbf{M}} \right] \mathbf{D} - d\tilde{\mathbf{M}} \right] \mathbf{D}^{-1}. \end{aligned}$$

Evolution of eigenvalues for general matrix SDE **Proof** Using the eigen decomposition, we have $\mathbf{R} = \mathbf{V}\mathbf{D}\mathbf{V}^\top$,

$$\begin{aligned} \mathbf{D} &= \mathbf{V}^\top \mathbf{R}\mathbf{V}, \\ d\mathbf{D} &= \mathbf{V}^\top d\mathbf{R}\mathbf{V} + \mathbf{V}^\top \mathbf{R}d\mathbf{V} + d\mathbf{V}^\top \mathbf{R}\mathbf{V} + \mathbf{V}^\top d\mathbf{R}d\mathbf{V} + d\mathbf{V}^\top d\mathbf{R}\mathbf{V} + d\mathbf{V}^\top \mathbf{R}d\mathbf{V}, \\ &= \mathbf{V}^\top d\mathbf{R}\mathbf{V} + \mathbf{D}\mathbf{V}^\top d\mathbf{V} + d\mathbf{V}^\top \mathbf{V}\mathbf{D} + \mathbf{V}^\top d\mathbf{R}d\mathbf{V} + d\mathbf{V}^\top d\mathbf{R}\mathbf{V} + \left(d\mathbf{V}^\top \mathbf{V} \right) \mathbf{D} \left(\mathbf{V}^\top d\mathbf{V} \right). \end{aligned}$$

The approach we follow is use the jacobian of the evolution of \mathbf{V} (see [49]) and solve the constrains equations to obtain the Ito correction term as done in [10]. Let (s_1, s_2, \dots, s_r) denote the diagonal entries of \mathbf{D} . Furthermore, we define the matrix \mathbf{S} , which plays a notable role in Jacobian w.r.t \mathbf{V} , as the following,

$$\mathbf{S}_{ij} = \begin{cases} 0 & \text{if } i = j, \\ (s_j - s_i)^{-1} & \text{o.w.} \end{cases} \quad (\text{H.1})$$

For the sake of brevity, we denote the evolution

$$\begin{aligned} d\mathbf{F} &\stackrel{\text{def}}{=} \mathbf{V}^\top d\mathbf{R}\mathbf{V} = \mathbf{V}^\top \mathbf{N}\mathbf{V} dt + \mathbf{V}^\top d\mathbf{M}\mathbf{V}, \\ &\stackrel{\text{def}}{=} \tilde{\mathbf{N}} dt + d\tilde{\mathbf{M}} \end{aligned}$$

The evolution of the eigenvectors,

$$d\mathbf{V} = \mathbf{V}d\Omega_{\mathbf{V}} + (\mathbf{I} - \mathbf{V}\mathbf{V}^\top)d\Xi_{\mathbf{V}}.$$

Using the Jacobian of the eigen vectors, we write,

$$\begin{aligned} d\Omega_{\mathbf{V}} &= \mathbf{Q}_{\parallel} dt + \mathbf{S} \odot d\mathbf{F}, \\ d\Xi_{\mathbf{V}} &= \mathbf{Q}_{\perp} dt + d\mathbf{R} \mathbf{V} \mathbf{D}^{-1}. \end{aligned}$$

Note that $\mathbf{V}^{\top} \mathbf{V} = \mathbf{I}_r$, using this we have,

$$\begin{aligned} 0 &= d(\mathbf{V}^{\top} \mathbf{V}) = d\mathbf{V}^{\top} \mathbf{V} + \mathbf{V}^{\top} d\mathbf{V} + d\mathbf{V}^{\top} d\mathbf{V}, \\ &= d\Omega_{\mathbf{V}}^{\top} + d\Omega_{\mathbf{V}} + d\mathbf{V}^{\top} \mathbf{V} \mathbf{V}^{\top} d\mathbf{V} + d\mathbf{V}^{\top} (\mathbf{I} - \mathbf{V} \mathbf{V}^{\top}) d\mathbf{V}, \\ &= d\Omega_{\mathbf{V}}^{\top} + d\Omega_{\mathbf{V}} + d\Omega_{\mathbf{V}}^{\top} d\Omega_{\mathbf{V}} + d\Xi_{\mathbf{V}}^{\top} (\mathbf{I} - \mathbf{V} \mathbf{V}^{\top}) d\Xi_{\mathbf{V}}, \\ &= d\Omega_{\mathbf{V}}^{\top} + d\Omega_{\mathbf{V}} - (\mathbf{S} \odot d\mathbf{F}) (\mathbf{S} \odot d\mathbf{F}) + \mathbf{D}^{-1} \mathbf{V}^{\top} d\mathbf{R} (\mathbf{I} - \mathbf{V} \mathbf{V}^{\top}) d\mathbf{R} \mathbf{V} \mathbf{D}^{-1}. \end{aligned}$$

Using $d\Omega_{\mathbf{V}}^{\top} = \mathbf{Q}_{\parallel}^{\top} dt - \mathbf{S} \odot d\mathbf{F}$, we have $d\Omega_{\mathbf{V}}^{\top} + d\Omega_{\mathbf{V}} = (\mathbf{Q}_{\parallel}^{\top} + \mathbf{Q}_{\parallel}) dt$.

$$(\mathbf{Q}_{\parallel} + \mathbf{Q}_{\parallel}^{\top}) dt = (\mathbf{S} \odot d\widetilde{\mathbf{M}}) (\mathbf{S} \odot d\widetilde{\mathbf{M}}) - \mathbf{D}^{-1} \mathbf{V}^{\top} d\mathbf{R} (\mathbf{I} - \mathbf{V} \mathbf{V}^{\top}) d\mathbf{R} \mathbf{V} \mathbf{D}^{-1}. \quad (\text{H.2})$$

Coming back to the evolution of singular values,

$$\begin{aligned} d\mathbf{D} &= \mathbf{V}^{\top} d\mathbf{R} \mathbf{V} + \mathbf{D} \mathbf{V}^{\top} d\mathbf{V} + d\mathbf{V}^{\top} \mathbf{V} \mathbf{D} + \mathbf{V}^{\top} d\mathbf{R} d\mathbf{V} + d\mathbf{V}^{\top} d\mathbf{R} \mathbf{V} + (d\mathbf{V}^{\top} \mathbf{V}) \mathbf{D} (\mathbf{V}^{\top} d\mathbf{V}). \\ &= d\mathbf{F} + (\mathbf{D} \mathbf{Q}_{\parallel} + \mathbf{Q}_{\parallel}^{\top} \mathbf{D}) dt + \mathbf{D} (\mathbf{S} \odot d\mathbf{F}) - (\mathbf{S} \odot d\mathbf{F}) \mathbf{D} + d\Omega_{\mathbf{V}}^{\top} \mathbf{D} d\Omega_{\mathbf{V}} \\ &\quad + \mathbf{V}^{\top} d\mathbf{R} [\mathbf{V} d\Omega_{\mathbf{V}} + (\mathbf{I} - \mathbf{V} \mathbf{V}^{\top}) d\Xi_{\mathbf{V}}] + [d\Omega_{\mathbf{V}}^{\top} \mathbf{V}^{\top} + d\Xi_{\mathbf{V}}^{\top} (\mathbf{I} - \mathbf{V} \mathbf{V}^{\top})] d\mathbf{R} \mathbf{V}, \end{aligned}$$

$$\begin{aligned} d\mathbf{D} &= \mathbf{I} \odot d\mathbf{F} + (\mathbf{D} \mathbf{Q}_{\parallel} + \mathbf{Q}_{\parallel}^{\top} \mathbf{D}) dt - (\mathbf{S} \odot d\widetilde{\mathbf{M}}) \mathbf{D} (\mathbf{S} \odot d\widetilde{\mathbf{M}}) + d\widetilde{\mathbf{M}} (\mathbf{S} \odot d\widetilde{\mathbf{M}}) \\ &\quad - (\mathbf{S} \odot d\widetilde{\mathbf{M}}) d\widetilde{\mathbf{M}} + \mathbf{V}^{\top} d\mathbf{R} (\mathbf{I} - \mathbf{V} \mathbf{V}^{\top}) d\mathbf{R} \mathbf{V} \mathbf{D}^{-1} + \mathbf{D}^{-1} \mathbf{V}^{\top} d\mathbf{R} (\mathbf{I} - \mathbf{V} \mathbf{V}^{\top}) d\mathbf{R}. \end{aligned} \quad (\text{H.3})$$

Note that $d\mathbf{D}$ is diagonal, hence, $\mathbf{I} \odot d\mathbf{D} = d\mathbf{D}$.

$$\begin{aligned} \mathbf{I} \odot d\mathbf{D} &= \mathbf{I} \odot d\mathbf{F} + \mathbf{I} \odot (\mathbf{D} \mathbf{Q}_{\parallel} + \mathbf{Q}_{\parallel}^{\top} \mathbf{D}) dt - \mathbf{I} \odot [(\mathbf{S} \odot d\widetilde{\mathbf{M}}) \mathbf{D} (\mathbf{S} \odot d\widetilde{\mathbf{M}})] \\ &\quad + 2\mathbf{I} \odot (d\widetilde{\mathbf{M}} (\mathbf{S} \odot d\widetilde{\mathbf{M}})) + 2\mathbf{I} \odot (\mathbf{D}^{-1} \mathbf{V}^{\top} d\mathbf{R} (\mathbf{I} - \mathbf{V} \mathbf{V}^{\top}) d\mathbf{R}) \end{aligned}$$

Note that $\mathbf{I} \odot (DM) = \mathbf{I} \odot (MD) = D \odot M$ for any matrix M and diagonal matrix D , using this property, we can simplify the above expression as,

$$\begin{aligned} d\mathbf{D} &= \mathbf{I} \odot d\mathbf{F} + \mathbf{D} \odot (\mathbf{Q}_{\parallel} + \mathbf{Q}_{\parallel}^{\top}) dt - \mathbf{I} \odot [(\mathbf{S} \odot d\widetilde{\mathbf{M}}) \mathbf{D} (\mathbf{S} \odot d\widetilde{\mathbf{M}})] \\ &\quad + 2\mathbf{I} \odot (d\widetilde{\mathbf{M}} (\mathbf{S} \odot d\widetilde{\mathbf{M}})) + 2\mathbf{D}^{-1} \odot (\mathbf{V}^{\top} d\mathbf{R} (\mathbf{I} - \mathbf{V} \mathbf{V}^{\top}) d\mathbf{R}) \end{aligned}$$

Using Eq. H.2, we have,

$$\begin{aligned} \mathbf{D} \odot (\mathbf{Q}_{\parallel} + \mathbf{Q}_{\parallel}^{\top}) dt &= \mathbf{D} \odot \left[(\mathbf{S} \odot d\widetilde{\mathbf{M}}) (\mathbf{S} \odot d\widetilde{\mathbf{M}}) - \mathbf{D}^{-1} \mathbf{V}^{\top} d\mathbf{R} (\mathbf{I} - \mathbf{V}\mathbf{V}^{\top}) d\mathbf{R}\mathbf{V}\mathbf{D}^{-1} \right], \\ &= \mathbf{I} \odot \left[(\mathbf{S} \odot d\widetilde{\mathbf{M}}) (\mathbf{S} \odot d\widetilde{\mathbf{M}}) \mathbf{D} \right] - \mathbf{D}^{-1} \odot \left(\mathbf{V}^{\top} d\mathbf{R} (\mathbf{I} - \mathbf{V}\mathbf{V}^{\top}) d\mathbf{R}\mathbf{V} \right). \end{aligned}$$

Using this,

$$\begin{aligned} d\mathbf{D} &= \mathbf{I} \odot d\mathbf{F} + \mathbf{I} \odot \left[(\mathbf{S} \odot d\widetilde{\mathbf{M}}) (\mathbf{S} \odot d\widetilde{\mathbf{M}}) \mathbf{D} \right] - \mathbf{I} \odot \left[(\mathbf{S} \odot d\widetilde{\mathbf{M}}) \mathbf{D} (\mathbf{S} \odot d\widetilde{\mathbf{M}}) \right] \\ &\quad + 2\mathbf{I} \odot \left(d\widetilde{\mathbf{M}} (\mathbf{S} \odot d\widetilde{\mathbf{M}}) \right) + \mathbf{D}^{-1} \odot \left(\mathbf{V}^{\top} d\mathbf{R} (\mathbf{I} - \mathbf{V}\mathbf{V}^{\top}) d\mathbf{R}\mathbf{V} \right), \\ &= \mathbf{I} \odot d\mathbf{F} + \mathbf{I} \odot \left[(\mathbf{S} \odot d\widetilde{\mathbf{M}}) \left[(\mathbf{S} \odot d\widetilde{\mathbf{M}}) \mathbf{D} - \mathbf{D} (\mathbf{S} \odot d\widetilde{\mathbf{M}}) \right] \right] \\ &\quad + 2\mathbf{I} \odot \left(d\widetilde{\mathbf{M}} (\mathbf{S} \odot d\widetilde{\mathbf{M}}) \right) + \mathbf{D}^{-1} \odot \left(\mathbf{V}^{\top} d\mathbf{R} (\mathbf{I} - \mathbf{V}\mathbf{V}^{\top}) d\mathbf{R}\mathbf{V} \right), \\ &= \mathbf{I} \odot d\mathbf{F} + \mathbf{I} \odot \left[(\mathbf{S} \odot d\widetilde{\mathbf{M}}) d\widetilde{\mathbf{M}} \right] \\ &\quad + 2\mathbf{I} \odot \left(d\widetilde{\mathbf{M}} (\mathbf{S} \odot d\widetilde{\mathbf{M}}) \right) + \mathbf{D}^{-1} \odot \left(\mathbf{V}^{\top} d\mathbf{R} (\mathbf{I} - \mathbf{V}\mathbf{V}^{\top}) d\mathbf{R}\mathbf{V} \right), \\ &= \mathbf{I} \odot d\mathbf{F} + \mathbf{I} \odot \left(d\widetilde{\mathbf{M}} (\mathbf{S} \odot d\widetilde{\mathbf{M}}) \right) + \mathbf{D}^{-1} \odot \left(\mathbf{V}^{\top} d\mathbf{R} (\mathbf{I} - \mathbf{V}\mathbf{V}^{\top}) d\mathbf{R}\mathbf{V} \right). \end{aligned}$$

Evolution of eigenvectors for general matrix SDE. Here, we derive the evolution of eigenvectors,

Using Eq. H.2, we have,

$$\left(\mathbf{Q}_{\parallel} \mathbf{D} + \mathbf{Q}_{\parallel}^{\top} \mathbf{D} \right) dt = \left(\mathbf{S} \odot d\widetilde{\mathbf{M}} \right) \left(\mathbf{S} \odot d\widetilde{\mathbf{M}} \right) \mathbf{D} - \mathbf{D}^{-1} \mathbf{V}^{\top} d\mathbf{R} (\mathbf{I} - \mathbf{V}\mathbf{V}^{\top}) d\mathbf{R}\mathbf{V}$$

Now further using the constrain that $d\mathbf{D}$ needs to be diagonal we get,

$$\begin{aligned} \left(\mathbf{D}\mathbf{Q}_{\parallel} + \mathbf{Q}_{\parallel}^{\top} \mathbf{D} \right) dt &= d\mathbf{D} - \mathbf{I} \odot d\mathbf{F} + \left(\mathbf{S} \odot d\widetilde{\mathbf{M}} \right) \mathbf{D} \left(\mathbf{S} \odot d\widetilde{\mathbf{M}} \right) - d\widetilde{\mathbf{M}} (\mathbf{S} \odot d\widetilde{\mathbf{M}}) + \left(\mathbf{S} \odot d\widetilde{\mathbf{M}} \right) d\widetilde{\mathbf{M}} \\ &\quad - \mathbf{V}^{\top} d\mathbf{R} (\mathbf{I} - \mathbf{V}\mathbf{V}^{\top}) d\mathbf{R}\mathbf{V}\mathbf{D}^{-1} - \mathbf{D}^{-1} \mathbf{V}^{\top} d\mathbf{R} (\mathbf{I} - \mathbf{V}\mathbf{V}^{\top}) d\mathbf{R}. \end{aligned}$$

$$\begin{aligned} \left(\mathbf{D}\mathbf{Q}_{\parallel} - \mathbf{Q}_{\parallel} \mathbf{D} \right) dt &= d\mathbf{D} - \mathbf{I} \odot d\mathbf{F} - \left(\mathbf{S} \odot d\widetilde{\mathbf{M}} \right) \left[\left(\mathbf{S} \odot d\widetilde{\mathbf{M}} \right) \mathbf{D} - \mathbf{D} \left(\mathbf{S} \odot d\widetilde{\mathbf{M}} \right) \right] - d\widetilde{\mathbf{M}} (\mathbf{S} \odot d\widetilde{\mathbf{M}}) \\ &\quad + \left(\mathbf{S} \odot d\widetilde{\mathbf{M}} \right) d\widetilde{\mathbf{M}} - \mathbf{V}^{\top} d\mathbf{R} (\mathbf{I} - \mathbf{V}\mathbf{V}^{\top}) d\mathbf{R}\mathbf{V}\mathbf{D}^{-1}, \\ &= d\mathbf{D} - \mathbf{I} \odot d\mathbf{F} - \left(\mathbf{S} \odot d\widetilde{\mathbf{M}} \right) \left[d\widetilde{\mathbf{M}} \odot \bar{\mathbf{I}} \right] - d\widetilde{\mathbf{M}} (\mathbf{S} \odot d\widetilde{\mathbf{M}}) \\ &\quad + \left(\mathbf{S} \odot d\widetilde{\mathbf{M}} \right) d\widetilde{\mathbf{M}} - \mathbf{V}^{\top} d\mathbf{R} (\mathbf{I} - \mathbf{V}\mathbf{V}^{\top}) d\mathbf{R}\mathbf{V}\mathbf{D}^{-1}, \\ &= d\mathbf{D} - \mathbf{I} \odot d\mathbf{F} + \left(\mathbf{S} \odot d\widetilde{\mathbf{M}} \right) \left[d\widetilde{\mathbf{M}} \odot \bar{\mathbf{I}} \right] - d\widetilde{\mathbf{M}} (\mathbf{S} \odot d\widetilde{\mathbf{M}}) \\ &\quad - \mathbf{V}^{\top} d\mathbf{R} (\mathbf{I} - \mathbf{V}\mathbf{V}^{\top}) d\mathbf{R}\mathbf{V}\mathbf{D}^{-1}. \end{aligned}$$

$$\begin{aligned} \bar{\mathbf{I}} \odot \left(\mathbf{D}\mathbf{Q}_{\parallel} - \mathbf{Q}_{\parallel} \mathbf{D} \right) dt &= \bar{\mathbf{I}} \odot \left(d\mathbf{D} - \mathbf{I} \odot d\mathbf{F} \right) + \bar{\mathbf{I}} \odot \left[\left(\mathbf{S} \odot d\widetilde{\mathbf{M}} \right) \left[d\widetilde{\mathbf{M}} \odot \bar{\mathbf{I}} \right] \right] - \bar{\mathbf{I}} \odot \left(d\widetilde{\mathbf{M}} (\mathbf{S} \odot d\widetilde{\mathbf{M}}) \right) \\ &\quad - \bar{\mathbf{I}} \odot \left(\mathbf{V}^{\top} d\mathbf{R} (\mathbf{I} - \mathbf{V}\mathbf{V}^{\top}) d\mathbf{R}\mathbf{V}\mathbf{D}^{-1} \right). \end{aligned}$$

$$(\bar{\mathbf{I}} \odot \mathbf{Q}_{\parallel}) dt = \mathbf{S} \odot \left[- \left(\mathbf{S} \odot d\tilde{\mathbf{M}} \right) \left[d\tilde{\mathbf{M}} \odot \mathbf{I} \right] + d\tilde{\mathbf{M}} \left(\mathbf{S} \odot d\tilde{\mathbf{M}} \right) + \mathbf{V}^{\top} d\mathbf{R} \left(\mathbf{I} - \mathbf{V}\mathbf{V}^{\top} \right) d\mathbf{R}\mathbf{V}\mathbf{D}^{-1} \right].$$

Combing these, we get the diagonal and off diagonal terms of \mathbf{Q}_{\parallel}

$$\begin{aligned} (\mathbf{I} \odot \mathbf{Q}_{\parallel}) dt &= \frac{1}{2} \mathbf{I} \odot \left(\mathbf{Q}_{\parallel} + \mathbf{Q}_{\parallel}^{\top} \right) dt, \\ &= \frac{\mathbf{I} \odot \left[\left(\mathbf{S} \odot d\tilde{\mathbf{M}} \right) \left(\mathbf{S} \odot d\tilde{\mathbf{M}} \right) \right]}{2} - \frac{\mathbf{I} \odot \left[\mathbf{D}^{-1} \mathbf{V}^{\top} d\mathbf{R} \left(\mathbf{I} - \mathbf{V}\mathbf{V}^{\top} \right) d\mathbf{R}\mathbf{V}\mathbf{D}^{-1} \right]}{2}. \end{aligned}$$

$$\begin{aligned} \mathbf{Q}_{\parallel} &= \frac{\mathbf{I} \odot \left[\left(\mathbf{S} \odot d\tilde{\mathbf{M}} \right) \left(\mathbf{S} \odot d\tilde{\mathbf{M}} \right) \right]}{2} - \frac{\mathbf{I} \odot \left[\mathbf{D}^{-1} \mathbf{V}^{\top} d\mathbf{R} \left(\mathbf{I} - \mathbf{V}\mathbf{V}^{\top} \right) d\mathbf{R}\mathbf{V}\mathbf{D}^{-1} \right]}{2} \\ &\quad - \mathbf{S} \odot \left[\left(\mathbf{S} \odot d\tilde{\mathbf{M}} \right) \left[d\tilde{\mathbf{M}} \odot \mathbf{I} \right] \right] + \mathbf{S} \odot \left(d\tilde{\mathbf{M}} \left(\mathbf{S} \odot d\tilde{\mathbf{M}} \right) \right) \\ &\quad + \mathbf{S} \odot \left(\mathbf{V}^{\top} d\mathbf{R} \left(\mathbf{I} - \mathbf{V}\mathbf{V}^{\top} \right) d\mathbf{R}\mathbf{V}\mathbf{D}^{-1} \right). \end{aligned}$$

Computing of \mathbf{Q}_{\perp} . Recalling the evolution of the eigenvectors,

$$d\mathbf{V} = \mathbf{V}d\Omega_{\mathbf{V}} + (\mathbf{I} - \mathbf{V}\mathbf{V}^{\top})d\Xi_{\mathbf{V}}.$$

Using the Jacobian of the eigen vectors, we write,

$$\begin{aligned} d\Omega_{\mathbf{V}} &= \mathbf{Q}_{\parallel} dt + \mathbf{S} \odot d\mathbf{F}, \\ d\Xi_{\mathbf{V}} &= \mathbf{Q}_{\perp} dt + d\mathbf{R} \mathbf{V}\mathbf{D}^{-1}, \\ d\mathbf{V} &= \mathbf{V} \left[\mathbf{Q}_{\parallel} dt + \mathbf{S} \odot d\mathbf{F} \right] + (\mathbf{I} - \mathbf{V}\mathbf{V}^{\top}) \left[\mathbf{Q}_{\perp} dt + d\mathbf{R} \mathbf{V}\mathbf{D}^{-1} \right], \\ d\mathbf{V}^{\top} &= \left[\mathbf{Q}_{\parallel}^{\top} dt - \mathbf{S} \odot d\mathbf{F} \right] \mathbf{V}^{\top} + \left[\mathbf{Q}_{\perp}^{\top} dt + \mathbf{D}^{-1} \mathbf{V}^{\top} d\mathbf{R} \right] (\mathbf{I} - \mathbf{V}\mathbf{V}^{\top}). \end{aligned}$$

Using the fact that $(\mathbf{I} - \mathbf{V}\mathbf{V}^{\top}) \mathbf{R} = 0$ and deriving it,

$$\begin{aligned} 0 &= (\mathbf{I} - \mathbf{V}\mathbf{V}^{\top}) \mathbf{R}, \\ 0 &= d \left[(\mathbf{I} - \mathbf{V}\mathbf{V}^{\top}) \mathbf{R} \right], \\ d\mathbf{R} &= d(\mathbf{V}\mathbf{V}^{\top} \mathbf{R}), \\ &= d\mathbf{V}\mathbf{V}^{\top} \mathbf{R} + \mathbf{V}d\mathbf{V}^{\top} \mathbf{R} + \mathbf{V}\mathbf{V}^{\top} d\mathbf{R} + d\mathbf{V}d\mathbf{V}^{\top} \mathbf{R} + d\mathbf{V}\mathbf{V}^{\top} d\mathbf{R} + \mathbf{V}d\mathbf{V}^{\top} d\mathbf{R}, \\ d\mathbf{R}\mathbf{V} &= d\mathbf{V}\mathbf{D} + \mathbf{V}d\mathbf{V}^{\top} \mathbf{V}\mathbf{D} + \mathbf{V}\mathbf{V}^{\top} d\mathbf{R}\mathbf{V} + d\mathbf{V}d\mathbf{V}^{\top} \mathbf{V}\mathbf{D} + d\mathbf{V}\mathbf{V}^{\top} d\mathbf{R}\mathbf{V} + \mathbf{V}d\mathbf{V}^{\top} d\mathbf{R}\mathbf{V}, \end{aligned}$$

$$\begin{aligned} d\mathbf{V}\mathbf{D} &= \mathbf{V} \left[\mathbf{Q}_{\parallel} \mathbf{D} dt + (\mathbf{S} \odot d\mathbf{F}) \mathbf{D} \right] + (\mathbf{I} - \mathbf{V}\mathbf{V}^{\top}) \left[\mathbf{Q}_{\perp} \mathbf{D} dt + d\mathbf{R} \mathbf{V} \right], \\ \mathbf{V}d\mathbf{V}^{\top} \mathbf{V}\mathbf{D} &= \mathbf{V} \left[\mathbf{Q}_{\parallel}^{\top} dt - \mathbf{S} \odot d\mathbf{F} \right] \mathbf{D}, \\ d\mathbf{V}d\mathbf{V}^{\top} \mathbf{V}\mathbf{D} &= -\mathbf{V} \left[\mathbf{S} \odot d\mathbf{F} \right] \left[\mathbf{S} \odot d\mathbf{F} \right] \mathbf{D} - \left[(\mathbf{I} - \mathbf{V}\mathbf{V}^{\top}) d\mathbf{R}\mathbf{V}\mathbf{D}^{-1} \right] \left[\mathbf{S} \odot d\mathbf{F} \right] \mathbf{D}, \\ d\mathbf{V}\mathbf{V}^{\top} d\mathbf{R}\mathbf{V} &= \mathbf{V} \left[\mathbf{S} \odot d\mathbf{F} \right] d\mathbf{F} + \left[(\mathbf{I} - \mathbf{V}\mathbf{V}^{\top}) d\mathbf{R}\mathbf{V}\mathbf{D}^{-1} \right] d\mathbf{F}, \\ \mathbf{V}d\mathbf{V}^{\top} d\mathbf{R}\mathbf{V} &= -\mathbf{V} \left[\mathbf{S} \odot d\mathbf{F} \right] d\mathbf{F} + \mathbf{V}\mathbf{D}^{-1} \mathbf{V}^{\top} d\mathbf{R} (\mathbf{I} - \mathbf{V}\mathbf{V}^{\top}) d\mathbf{R}\mathbf{V}. \end{aligned}$$

Adding the terms up we get,

$$\begin{aligned} & \mathbf{V} \left[\mathbf{Q}_{\parallel} + \mathbf{Q}_{\parallel}^{\top} \right] \mathbf{D} dt + \left(\mathbf{I} - \mathbf{V}\mathbf{V}^{\top} \right) \mathbf{Q}_{\perp} \mathbf{D} dt \\ & - \mathbf{V} [\mathbf{S} \odot d\mathbf{F}] [\mathbf{S} \odot d\mathbf{F}] \mathbf{D} - \left[\left(\mathbf{I} - \mathbf{V}\mathbf{V}^{\top} \right) d\mathbf{R}\mathbf{V}\mathbf{D}^{-1} \right] [\mathbf{S} \odot d\mathbf{F}] \mathbf{D} \\ & + \left[\left(\mathbf{I} - \mathbf{V}\mathbf{V}^{\top} \right) d\mathbf{R}\mathbf{V}\mathbf{D}^{-1} \right] d\mathbf{F} + \mathbf{V}\mathbf{D}^{-1}\mathbf{V}^{\top} d\mathbf{R}(\mathbf{I} - \mathbf{V}\mathbf{V}^{\top})d\mathbf{R}\mathbf{V} = 0. \end{aligned}$$

$$\left(\mathbf{I} - \mathbf{V}\mathbf{V}^{\top} \right) \mathbf{Q}_{\perp} \mathbf{D} dt - \left[\left(\mathbf{I} - \mathbf{V}\mathbf{V}^{\top} \right) d\mathbf{R}\mathbf{V}\mathbf{D}^{-1} \right] [\mathbf{S} \odot d\mathbf{F}] \mathbf{D} + \left[\left(\mathbf{I} - \mathbf{V}\mathbf{V}^{\top} \right) d\mathbf{R}\mathbf{V}\mathbf{D}^{-1} \right] d\mathbf{F} = 0.$$

$$\begin{aligned} \left(\mathbf{I} - \mathbf{V}\mathbf{V}^{\top} \right) \mathbf{Q}_{\perp} \mathbf{D} dt &= \left[\left(\mathbf{I} - \mathbf{V}\mathbf{V}^{\top} \right) d\mathbf{R}\mathbf{V}\mathbf{D}^{-1} \right] [\mathbf{S} \odot d\mathbf{F}] \mathbf{D} - \left[\left(\mathbf{I} - \mathbf{V}\mathbf{V}^{\top} \right) d\mathbf{R}\mathbf{V}\mathbf{D}^{-1} \right] d\mathbf{F}, \\ \left(\mathbf{I} - \mathbf{V}\mathbf{V}^{\top} \right) \mathbf{Q}_{\perp} &= \left[\left(\mathbf{I} - \mathbf{V}\mathbf{V}^{\top} \right) d\mathbf{R}\mathbf{V}\mathbf{D}^{-1} \right] [[\mathbf{S} \odot d\mathbf{F}] \mathbf{D} - d\mathbf{F}] \mathbf{D}^{-1} \end{aligned}$$

$$\mathbf{Q}_{\perp} = [d\mathbf{R}\mathbf{V}\mathbf{D}^{-1}] [[\mathbf{S} \odot d\mathbf{F}] \mathbf{D} - d\mathbf{F}] \mathbf{D}^{-1}$$

This gives the expression for \mathbf{Q}_{\perp} and this ends our computation. \blacksquare

Lemma H.3 For any matrix $A \in \mathbb{R}^{n \times m}$, $B \in \mathbb{R}^{n \times n}$, $m \times n$ -dimensional Brownian motion $d\mathbf{B}_t$, the following results hold on the covariance

$$d\mathbf{B}_t A d\mathbf{B}_t = A^{\top} dt, \quad (\text{H.4})$$

$$d\mathbf{B}_t B d\mathbf{B}_t^{\top} = \text{tr}(B) I_m dt. \quad (\text{H.5})$$

Lemma H.4 With \mathbf{S} defined in Equation (H.1), $d\mathbf{F} = d\mathbf{F} = \Sigma \mathbf{V}^{\top} d\mathbf{B}_t \mathbf{c}^{\top} + \mathbf{c} d\mathbf{B}_t^{\top} \mathbf{V} \Sigma$ and $d\mathbf{m}_t \stackrel{\text{def}}{=} (\boldsymbol{\sigma} \odot d\tilde{\mathbf{B}}_t)$.

$$d\mathbf{F}(\mathbf{S} \odot d\mathbf{F}) = \mathbf{c} \mathbf{s}^{\top} \mathbf{S} \text{diag}(\mathbf{c}) dt - \mathbf{D} \text{diag}(\mathbf{S} \text{diag}(\mathbf{c}) \mathbf{c}) dt + \mathbf{D} \text{diag}(\mathbf{c}) \mathbf{S} \text{diag}(\mathbf{c}) dt. \quad (\text{H.6})$$

Proof

$$\begin{aligned} \mathbf{S} \odot d\mathbf{F} &= [\text{diag}(\mathbf{c}) \mathbf{S} \text{diag}(d\mathbf{m}_t) + \text{diag}(d\mathbf{m}_t) \mathbf{S} \text{diag}(\mathbf{c})], \\ d\mathbf{F}(\mathbf{S} \odot d\mathbf{F}) &= \left(\mathbf{c} d\mathbf{m}_t^{\top} + d\mathbf{m}_t \mathbf{c}^{\top} \right) [\text{diag}(\mathbf{c}) \mathbf{S} \text{diag}(d\mathbf{m}_t) + \text{diag}(d\mathbf{m}_t) \mathbf{S} \text{diag}(\mathbf{c})], \\ &= \mathbf{c} \mathbf{s}^{\top} \mathbf{S} \text{diag}(\mathbf{c}) dt - \mathbf{D} \text{diag}(\mathbf{S} \text{diag}(\mathbf{c}) \mathbf{c}) dt + \mathbf{D} \text{diag}(\mathbf{c}) \mathbf{S} \text{diag}(\mathbf{c}) dt. \end{aligned}$$

\blacksquare

Lemma H.5 With \mathbf{S} defined in Equation (H.1), $d\mathbf{F} = d\mathbf{F} = \Sigma \mathbf{V}^{\top} d\mathbf{B}_t \mathbf{c}^{\top} + \mathbf{c} d\mathbf{B}_t^{\top} \mathbf{V} \Sigma$ and $d\mathbf{m}_t \stackrel{\text{def}}{=} (\boldsymbol{\sigma} \odot d\tilde{\mathbf{B}}_t)$.

$$(\mathbf{S} \odot d\mathbf{F})(\mathbf{S} \odot d\mathbf{F}) = \mathbf{D} \text{diag} \left(\mathbf{S} \text{diag}(\mathbf{c})^2 \mathbf{S} \right) dt + \text{diag}(\mathbf{c}) \mathbf{S} \mathbf{D} \mathbf{S} \text{diag}(\mathbf{c}) dt. \quad (\text{H.7})$$

Proof

$$\begin{aligned} (\mathbf{S} \odot d\mathbf{F}) &= \mathbf{S} \odot (\mathbf{d}\mathbf{m}_t \mathbf{c}^\top + \mathbf{c} \mathbf{d}\mathbf{m}_t^\top), \\ &= \text{diag}(\mathbf{c}) \mathbf{S} \text{diag}(\mathbf{d}\mathbf{m}_t) + \text{diag}(\mathbf{d}\mathbf{m}_t) \mathbf{S} \text{diag}(\mathbf{c}). \end{aligned}$$

Now, computing the product,

$$\begin{aligned} (\mathbf{S} \odot d\mathbf{F})(\mathbf{S} \odot d\mathbf{F}) &= [\text{diag}(\mathbf{c}) \mathbf{S} \text{diag}(\mathbf{d}\mathbf{m}_t) + \text{diag}(\mathbf{d}\mathbf{m}_t) \mathbf{S} \text{diag}(\mathbf{c})] [\text{diag}(\mathbf{c}) \mathbf{S} \text{diag}(\mathbf{d}\mathbf{m}_t) + \text{diag}(\mathbf{d}\mathbf{m}_t) \mathbf{S} \text{diag}(\mathbf{c})], \\ &= \mathbf{D} \text{diag}(\mathbf{S} \text{diag}(\mathbf{c})^2 \mathbf{S}) dt + \text{diag}(\mathbf{c}) \mathbf{S} \mathbf{D} \mathbf{S} \text{diag}(\mathbf{c}) dt. \end{aligned}$$

■

Lemma H.6

$$(\mathbf{S} \odot d\mathbf{F}) \mathbf{d}\mathbf{m}_t \mathbf{c}^\top d\mathbf{F} =$$

Proof

$$(\mathbf{S} \odot d\mathbf{F}) \mathbf{d}\mathbf{m}_t = [\text{diag}(\mathbf{c}) \mathbf{S} \text{diag}(\mathbf{d}\mathbf{m}_t) + \text{diag}(\mathbf{d}\mathbf{m}_t) \mathbf{S} \text{diag}(\mathbf{c})] \mathbf{d}\mathbf{m}_t = \text{diag}(\mathbf{c}) \mathbf{S} (\boldsymbol{\sigma} \odot \boldsymbol{\sigma})$$

■