
Bias in Motion: Theoretical Insights into the Dynamics of Bias in SGD Training

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Machine learning systems often acquire biases by leveraging undesired features
2 in the data, impacting accuracy variably across different sub-populations. This
3 paper explores the evolution of bias in a teacher-student setup modeling different
4 data sub-populations with a Gaussian-mixture model, by providing an analytical
5 description of the stochastic gradient descent dynamics of a linear classifier in this
6 setting. Our analysis reveals how different properties of sub-populations influence
7 bias at different timescales, showing a shifting preference of the classifier during
8 training. We empirically validate our results in more complex scenarios by training
9 deeper networks on real datasets including CIFAR10, MNIST, and CelebA.

10 1 Introduction

11 Machine learning (ML) systems not only reproduce existing biases in the data but also tend to amplify
12 them [19, 38, 11]. Given the complexity of the ML pipeline, isolating and characterising the key
13 drivers of this amplification is challenging. Theoretical results in this area (e.g., [35, 36]) are mostly
14 based on asymptotic analysis, leaving the transient learning regime poorly understood.

15 Our analysis addresses this gap by providing a precise characterisation of the transient dynamics
16 of online stochastic gradient descent (SGD) in a high dimensional prototypical model of linear
17 classification. We use the teacher-mixture (TM) framework [36], where different data sub-populations
18 are modeled with a mixture of Gaussians, each having its own linear rule (teacher) for determining the
19 labels. Adjusting the parameters of the data distribution in our framework connects models of fairness
20 and spurious correlations, providing a unifying framework and a general set of results applicable to
21 both domains. Remarkably, our study reveals a rich behaviour divided into three learning phases,
22 where different features of data bias the classifier and causing significant deviations from asymptotic
23 predictions. We reproduce our theoretical findings through numerical experiments in more complex
24 settings, demonstrating validity beyond the simplicity of our model.

25 2 Problem setup

26 We consider a standard supervised learning setup where the training data consists of pairs of a feature
27 vector $\mathbf{x} \in \mathbb{R}^d$ and a binary label $y = \pm 1$. To model subgroups within the data [33], we assume that
28 the feature vectors are structured as clusters c_1, \dots, c_m , respectively centered on some fixed attribute
29 vectors $\mathbf{v}_1, \dots, \mathbf{v}_m \in \mathbb{R}^d$. Specifically, \mathbf{x} is sampled from a mixture of m isotropic Gaussians:

$$\mathbf{x} \sim \sum_{j=1}^m \rho_j \mathcal{N}(\mathbf{v}_j / \sqrt{d}, \Delta_j \mathbb{I}_{d \times d}), \quad (1)$$

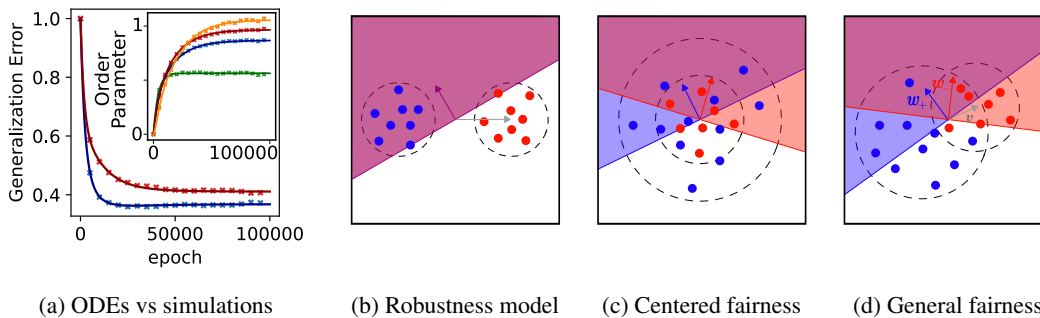


Figure 1: **Teacher-Mixture in fairness and robustness.** *Panel (a)* shows the generalisation errors—for the subpopulations + (blue) and – (red)—obtained through simulation (crosses) and predicted by the theory (solid lines) for a network with linear activation. The inset shows the same comparison for the *order parameters*: R_+ (blue), R_- (red), M (green), and Q (orange). *Panels (b-d)* exemplify the different scenarios achievable in the TM model investigated in Sec. 4. *Panel (b)* represent a model for robustness where a spurious feature—given by the shift vector—can mislead the classifier, see Sec. 4.1. *Panels (c,d)* are instead discussed in Sec. 4.2 and represent two models of fairness. First, *Panel (c)* has no shift, $v = 0$, allowing us to remove the confounding effects. Finally, *Panel (d)* shows the general fairness problem.

30 with mixing probabilities ρ_1, \dots, ρ_m and scalar variances $\Delta_1, \dots, \Delta_m$. Assuming the entries of
 31 \mathbf{v}_j are of order 1 as d gets large, the scaling factor $1/\sqrt{d}$ ensures that the Euclidean norm of the
 32 renormalised vector is of order 1. This prevents the problem from becoming either trivial or overly
 33 challenging in the high-dimensional limit [23, 22]. We adopt a teacher-mixture (TM) scenario [36]
 34 where each cluster has its own teacher rule:

$$\mathbf{x} \in c_j \implies y = \text{sign}(\bar{\mathbf{w}}_j^\top \mathbf{x} / \sqrt{d}). \quad (2)$$

35 This rule is characterised by the teacher vectors $\bar{\mathbf{w}}_j \in \mathbb{R}^d$, ensuring linear separability within each
 36 cluster. Fig. 1b-d illustrate the data distribution for two clusters with opposite mean vectors $\pm \mathbf{v}$,
 37 which will be the primary case study for our analysis.

38 **Model.** In this study we analyse a linear model applied to the above data distribution. We aim to
 39 learn a vector parameter \mathbf{w} , referred to as the ‘student’, such that predictions are given by

$$\hat{y}(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} / \sqrt{d}. \quad (3)$$

40 The training process involves applying online SGD on the squared loss $\hat{\epsilon} = (y - \hat{y})^2$ with learning rate
 41 $\eta/2 > 0$ (see Eq. C.17 in Appendix C). In our analysis, the model is evaluated by its generalisation
 42 error, or population loss, $\epsilon := \mathbb{E}[\hat{\epsilon}]$.

43 3 SGD analysis

44 We study the evolution of the generalisation error during training in the high dimensional setting (i.e.
 45 large d). Following a classical approach [32, 8], we streamline the problem by focusing on a small
 46 set of summary statistics, referred to as ‘order parameters’, which fully characterises the dynamics.
 47 As the dimension increases, it can be shown by concentration arguments that the evolution of these
 48 order parameters converges to the deterministic solution of a system of ODEs [14, 6, 3]. Notably, in
 49 our setting, we achieve an analytical solution of this ODE system.

50 3.1 Order parameters

51 In the setup described in Section 2, consider the following $2m + 1$ variables:

$$R_j = \frac{1}{d} \mathbf{w}^\top \bar{\mathbf{w}}_j, \quad M_j = \frac{1}{d} \mathbf{w}^\top \mathbf{v}_j, \quad Q = \frac{1}{d} \|\mathbf{w}\|^2, \quad (4)$$

52 for $1 \leq j \leq m$. These variables correspond to key statistics of the student, namely its alignment to
 53 the cluster teachers, its alignment to the cluster centers, and its magnitude, respectively. Lemma C.1

54 in Appendix C shows how the generalisation error depend on the model parameter \mathbf{w} only through
 55 these order parameters.

56 3.2 High dimensional dynamics

57 Let $\mathcal{S} := (S_i)_{1 \leq i \leq 2m+1}$ denote the collection of order parameters. Theorem C.3 in Appendix C
 58 states that as d gets large, the stochastic evolution \mathcal{S}^k of the order parameter gets uniformly close,
 59 with high probability, to the average continuous-time dynamics described by the ODE system:

$$\frac{d\bar{S}_i(t)}{dt} = f_i(\bar{\mathcal{S}}(t)), \quad 1 \leq i \leq 2m+1, \quad (5)$$

60 where the continuous *time* is given by the example number divided by the input dimension, $t = k/d$.

61 **Solving the ODEs.** We present the explicit solution of the ODEs in the case of two clusters ($m = 2$)
 62 with opposite mean vectors $\pm \mathbf{v}$, as in [36]. Henceforth, we refer to \mathbf{v} as the shift vector and to the
 63 two clusters as the ‘positive’ and ‘negative’ sub-populations, with mixing probabilities ρ and $(1 - \rho)$,
 64 variances Δ_{\pm} and teacher vectors $\bar{\mathbf{w}}_{\pm}$, respectively. The order parameters introduced in Eq. 4 are
 65 specifically denoted as $M = \mathbf{w}^{\top} \mathbf{v} / d$, $R_{+} = \mathbf{w}^{\top} \bar{\mathbf{w}}_{+} / d$, and $R_{-} = \mathbf{w}^{\top} \bar{\mathbf{w}}_{-} / d$ in this setting.

66 **Theorem 3.1.** *In the above setting, solutions to the order parameter evolution take the form*

$$M(t) = M_0 e^{-\eta(v + \Delta^{mix})t} + M^{\infty} (1 - e^{-\eta(v + \Delta^{mix})t}), \quad (6)$$

$$R_{\pm}(t) = R_{\pm}^0 e^{-\eta \Delta^{mix} t} + R_{\pm}^{\infty} (1 - e^{-\eta \Delta^{mix} t}) + k_1^{\pm} (e^{-\eta \Delta^{mix} t} - e^{-\eta(v + \Delta^{mix})t}), \quad (7)$$

$$\begin{aligned} Q(t) &= Q_0 e^{-\eta(2\Delta^{mix} - \eta \Delta^{2mix})t} + Q^{\infty} (1 - e^{-\eta(2\Delta^{mix} - \eta \Delta^{2mix})t}) \\ &+ k_2 (e^{-t(2\Delta^{mix} - \eta \Delta^{2mix})\eta} - e^{-t \Delta^{mix} \eta}) + k_3 (e^{-t(2\Delta^{mix} - \eta \Delta^{2mix})\eta} - e^{-t(v + \Delta^{mix})\eta}) \\ &+ k_4 (e^{-t(2\Delta^{mix} - \eta \Delta^{2mix})\eta} - e^{-t(2v + 2\Delta^{mix})\eta}), \end{aligned} \quad (8)$$

67 with $\Delta^{mix} = \rho \Delta_{+} + (1 - \rho) \Delta_{-}$, $\Delta^{2mix} = \rho \Delta_{+}^2 + (1 - \rho) \Delta_{-}^2$ and $v = \|\mathbf{v}\|^2 / d$.

68 The remaining constants are less significant and are reported in Appendix E.1 and discussed further
 69 in Appendix F. This solution allows us to describe important observables such as the generalisation
 70 error at any timestep. Fig. 1a plots the theoretical closed-form solutions along with values obtained
 71 through simulation when we set $d = 1000$. Note the remarkable agreement between the analytical
 72 ODE solution and simulations of the online SGD dynamics in this high dimensional data limit.

73 4 Insights

74 By examining the exponents in Eqs. 6-8, we can identify the relevant training timescales. Notably, M
 75 follows a straightforward behaviour dominated by a single timescale, whereas R_{\pm} and Q exhibit multiple timescales,
 76 leading to significant implications for the emergence and evolution of bias during training.
 77
 78
 79

80 Parameters specifying these different bias scenarios are the shift norm $v = \|\mathbf{v}\|^2 / d$ and relative representation ρ ,
 81 the subpopulation variances Δ_{\pm} , and the teacher overlap
 82 $T_{\pm} = \bar{\mathbf{w}}_{+}^{\top} \bar{\mathbf{w}}_{-} / d$. For simplicity we fix the teacher norm
 83 $\|\bar{\mathbf{w}}_{\pm}\|_2 = \sqrt{d}$, so that T_{\pm} is the cosine similarity between
 84 the two teachers.
 85

86 4.1 Spurious correlations

87 The emergence of spurious correlations during training
 88 illustrates a type of bias where a classifier favours a spu-
 89 rious feature over a core one. To isolate the impact
 90 of spurious correlation in our model while avoiding confounding effects, we consider perfectly
 91 overlapping teachers ($\bar{\mathbf{w}}_{+} = \bar{\mathbf{w}}_{-}$) and sub-populations with equal variance and representation

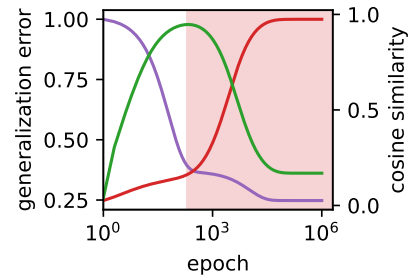


Figure 2: **Spurious correlations transient alignment.** Time-evolution of loss (purple), student-teacher (red) and student-shift (green) cosine similarities. The initial phase (green background) of learning aligns classifier and shift vector before aligning with the teacher (red background). Parameters: $v = 16$, $\rho = 0.5$, $\Delta_{-} = \Delta_{+} = 0.1$, $T_{\pm} = 1$, $\eta = 0.5$.

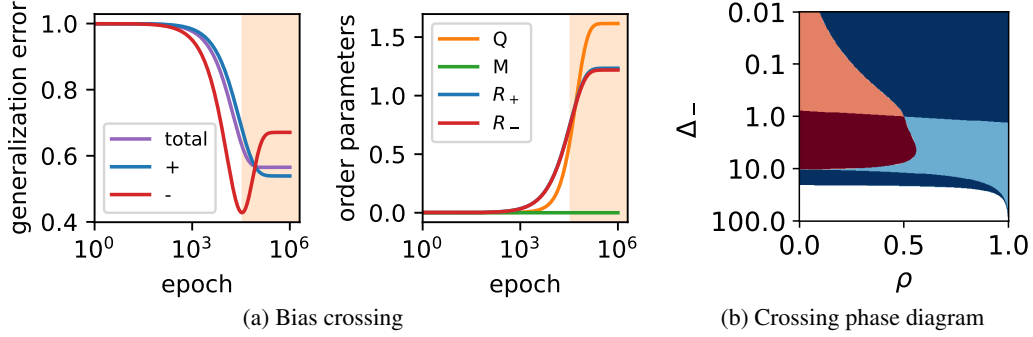


Figure 3: **The crossing phenomenon.** Panel (a) (left side) shows the loss curves of sub-population $-$ (in red) and sub-population $+$ in blue along with the overall loss (in purple). We observe a crossing cause by a higher variance but lower representation in sub-population $-$. The background colours represent the different phases of bias that are characterised by the evolution of the order parameters shown in Panel (a) (right side). Panel (b) shows the presence of the crossing phenomenon in a large portion of the parameter space using a phase diagram. Blue indicates an asymptotic preference for sub-population $+$ and red the opposite. Dark colours indicates regions where bias is consistent across training, while regions in light colours undergo a crossing phenomenon. White indicates that learning rate was too high and training diverged. Parameters: $v = 0, \Delta_+ = 1, T_{\pm} = 0.9, \eta = 0.1$.

92 ($\rho = 0.5, \Delta_+ = \Delta_-$). With non-perfectly overlapping clusters $v \neq 0$, we introduce a spurious
 93 correlation by adding a small cosine similarity between the shift vector and the teacher, creating a
 94 label imbalance within each sub-population (Fig. 1b).

95 From Eqs. 6-8, two relevant timescales for the problem are observed:

$$\tau_M = \frac{1}{\eta(v + \Delta^{mix})}, \quad \tau_R = 1/\eta\Delta^{mix}. \quad (9)$$

96 The shortest timescale, τ_M , indicates that the student first aligns with the spurious feature. By
 97 aligning with the shift vector, the student can predict most examples correctly, but not all. The effect
 98 of spurious correlations is transient; at $t \sim \tau_R$, the student starts disaligning from the spurious feature
 99 and aligns with the teacher vector, eventually achieving nearly perfect alignment (Fig. 2).

100 4.2 Fairness

101 In this section, we identify the properties of sub-populations that determine the bias during learning
 102 and show how bias evolves in three phases. To quantify bias, we use the *overall accuracy equality*
 103 metric [7], which measures the discrepancy in accuracy across groups. Intuitively, we aim for equal
 104 loss on both groups, considering any deviation from this condition as bias.

105 **Zero shift.** We first consider a simplified case where we assume that both clusters are centered at
 106 the origin $v = 0$ as shown in Fig. 1c. We will later reintroduce the shift and analyse the transient
 107 dynamics it introduces as per the discussion in section 4.1. This setting is particularly suited to
 108 analysing the effects of ‘group level’ features, such as group variance and relative representation, on
 109 the preference of the classifier.

110 In this simplified setting, $M(t)$ is always zero and the constants k_1^{\pm}, k_3, k_4 presented in equations 7
 111 and 8 are zero. Thus, the dynamics only involve two relevant timescales given by τ_R in Eq. 9 and

$$\tau_Q = 1/(\eta(2\Delta^{mix} - \eta\Delta^{2mix})). \quad (10)$$

112 Fig. 3a illustrates the changing preference of the classifier. Specifically, we observe that the variance
 113 of the sub-population is particularly relevant initially and the sub-population with higher variance
 114 (red) is *learnt* faster, i.e. its generalisation error drops faster. However, asymptotically we observe
 115 that the relative representation becomes more important wherein the student aligns itself with the
 116 teacher that has a higher product of representation and standard deviation (blue), i.e.

$$\rho\sqrt{\Delta_+} \geq (1 - \rho)\sqrt{\Delta_-} \iff R_+^{\infty} \geq R_-^{\infty}. \quad (11)$$

117 Thus, the network can advantage the cluster with higher variance initially but asymptotically advantage
 118 the other cluster if its representation is high enough. This leads to the ‘crossing’ of the losses on the
 119 two sub-populations shown in Fig. 3 (more in Appendix F.2).

120 *Initial dynamics.* The ratio between initial rate of change in generalisation errors is bounded by
 121 (derived in Appendix F.3):

$$T_{\pm} \sqrt{\frac{\Delta_+}{\Delta_-}} \leq \frac{d\epsilon_{g+}/dt|_{t=0}}{d\epsilon_{g-}/dt|_{t=0}} \leq \frac{1}{T_{\pm}} \sqrt{\frac{\Delta_+}{\Delta_-}}. \quad (12)$$

122 When the teachers are only slightly misaligned— $T_{\pm} \lesssim 1$ —the bound is tight and we can see that it is
 123 the ratio of the square roots of the variances that determines which cluster is learnt faster initially.
 124 Fig. 3b shows in a phase diagram the existence of ‘bias crossing’ across a wide range of variances
 125 and representations. The transition between the phases that represent an initial preference for the
 126 positive sub-population (light red and dark blue) and the phases that represent an initial preference for
 127 negative sub-population (dark red and light blue) is approximately given by the line $\Delta_- = \Delta_+ = 1$,
 128 independent of the representation as predicted by Eq. 12. The portion of the dark blue phase just
 129 above the white divergent phase marks a ‘quasi-divergent’ region wherein the generalisation error on
 130 the negative sub-population rises even at $t = 0$ because the learning rate is too large for such high
 131 variances and marks a region of impractical behaviour observed with poorly optimised learning rates.

132 *Asymptotic preference.* In the limit of small learning rates $\eta \rightarrow 0$, the student will asymptotically
 133 exhibit lower loss on whichever sub-population’s teacher it has better alignment with. Thus, Eq. 11
 134 provides a simple characterisation of asymptotic preference from representations and standard
 135 deviations in the small learning rate limit. However, the situation is more complex in the case of finite
 136 learning rate, which may disrupt learning in one or both clusters (more in Appendix F.4).

137 **General case.** We now consider the
 138 general case shown in Fig. 1d, where
 139 the shift is non zero and all three
 140 timescales identified so far play a role.

141 As observed in Sec. 4.1, when the shift
 142 norm v is large, the effect of spurious
 143 correlations becomes significant and
 144 the timescale associated with the spu-
 145 rious correlations is the fastest. In gen-
 146 eral, when $v \neq 0$ we observe an addi-
 147 tional phase due to the effect of spu-
 148 rious correlation. In this new first phase,
 149 the student advantages the cluster with
 150 higher representation and lower vari-
 151 ance since the salient information re-
 152 ceived from this cluster is more coher-
 153 ent and easier to access.

154 More precisely, in high dimensions
 155 the shift and the teachers are likely to
 156 exhibit a small cosine similarity leading to a class imbalance in the clusters and creating spurious
 157 correlation. The amount of label imbalance within a cluster is characterised by the value of α , as
 158 detailed in Appendix B. For smaller variances, α takes more extreme values leading to stronger
 159 spurious correlation of that cluster with the shift. If a cluster has more positive examples, we would
 160 observe a reduction in loss for that cluster if the student aligns with the mean of that cluster (and
 161 opposite to the mean if the cluster has mostly negative examples). When both clusters have different
 162 majority classes, the direction of spurious correlation for the two are same. However, when the
 163 majority classes are the same, we have competing directions for spurious correlation. The expression
 164 for M_{∞} in Appendix E.1 Eq. E.41 shows that in this case the relative representation comes into
 165 play and the mean of the cluster with greater representation and class imbalance will be chosen by
 166 the teacher to align with. Fig. 4 shows such a scenario with three phase bias evolution. First, the
 167 green phase is driven by spurious correlation where the positive cluster is advantaged since it has
 168 greater representation and class imbalance. Next, the red phase is driven by greater variance where

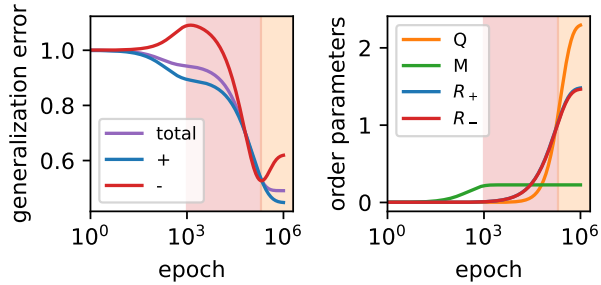


Figure 4: **Double crossing phenomenon.** (Left panel) shows the loss for the two sub-populations (blue and red lines) and the global one (in purple). (Right panel) shows the value of the order parameters across time. The behaviour of the order parameters across time provides a precise characterisation and understanding of the different phases. Parameters: $v = 100$, $\rho = 0.75$, $\Delta_+ = 0.1$, $\Delta_- = 0.5$, $\eta = 0.03$, $T_{\pm} = 0.9$, $\alpha_+ = 0.343$, $\alpha_- = 0.12$.

169 the negative cluster is learnt faster as discussed through Eq. 12. Finally, we observe the orange phase
 170 wherein the student starts aligning with the positive cluster as per the asymptotic rule in Eq. 11.

171 Our analysis thus shows that bias is a dynamical quantity that can vary non-monotonically during
 172 training and cannot be characterised by simply the initial and asymptotic values.

173 5 Ablations using numerical simulations

174 **Rotated MNIST.** We train a 2-
 175 layer neural network with 200 hid-
 176 den units, ReLU activation, and sig-
 177 moidal readout activation on a varia-
 178 tion of MNIST that mimics our
 179 model. Digits 0 to 4 and 5 to 9 are
 180 grouped to form the two subpopula-
 181 tions. With probability p_+ and p_- ,
 182 digits of both subpopulations are ro-
 183 tated with a subpopulation-specific
 184 angle—i.e. Fig. 5a uses angles of
 185 rotation $\theta_+ = 45^\circ$ and $\theta_- = -90^\circ$.
 186 The goal of the classifier is to detect
 187 rotations.

188 The experimental framework gives a
 189 correspondence between parameters
 190 of the generative model and proper-
 191 ties of a real dataset. We can con-
 192 trol relative representation by sub-
 193 sampling, teacher similarity by play-
 194 ing with angle difference, label im-
 195 balance by changing the probability of rotation, and saliency by increasing and decreasing the norm
 196 of the subpopulation using multiplicative factors Δ_\pm . The only parameter that we cannot control is
 197 the shift v which is a property of the data.

198 Therefore, in order to reproduce the zero-shift case of Sec. 4.2, we remove the label imbalance
 199 by setting the probability of rotation $p_+ = p_- = 0.5$. By properly calibrating the saliency Δ and
 200 the relative representation ρ , it is possible to bias the classifier towards one subpopulation at the
 201 beginning of training and the other in the end. This is shown in Fig. 5a where $\rho = 0.1$ and $\Delta_+ > \Delta_-$.
 202 The saliency difference favours subpopulation + initially while setting ρ small enough advantages
 203 subpopulation - later in training. This is precisely what we observe in the plot.

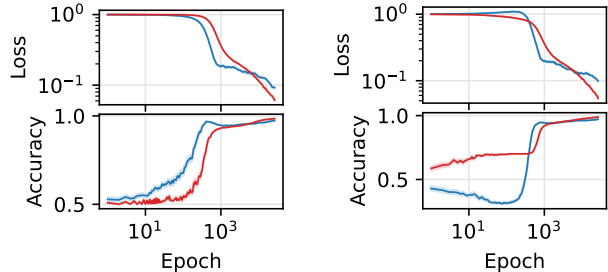
204 Finally, we consider the general fairness case. By creating label imbalance, i.e. setting $p_+ = 0.3$
 205 and $p_- = 0.7$, we observe an additional phase of bias evolution, wherein the classifier prefers dense
 206 regions with consistent labels. This advantages subpopulation - and indeed it is what we see in
 207 Fig. 5b. The result of the simulations matches the theory displaying a double crossing phenomenon.

208 **Additional numerical experiments.** In Appendix G, we provide additional experiments within our
 209 model and the CIFAR10 and CelebA, exploring different architectures and losses. We observe that
 210 bias presents different timescales and shows crossing behaviors.

211 6 Conclusion

212 This paper examined the dynamics of bias in a high dimensional synthetic framework, showing that it
 213 can be explicitly characterised to reveal transient behaviour. Our findings reveal that classifiers exhibit
 214 biases toward different data features during training, possibly alternating sub-population preference.
 215 Although our analysis is based on certain assumptions, numerical experiments that violate these
 216 assumptions still display the behaviour predicated by our theory.

217 We believe this line of research will have practical impacts in the medium term, aiding the design of
 218 mitigation strategies that account for transient dynamics. Future research will further explore this
 219 connection, proposing theory-based dynamical protocols for bias mitigation.



(a) Crossing phenomenon (b) Double crossing phenomenon

Figure 5: **Numerical simulations on MNIST.** The figure shows the average (solid lines) and standard deviation (shaded area) of 100 simulations run in this framework. In particular the upper plots show the test loss and lower plots the test accuracy for subpopulation + (blue) and - (red). *Panel (a)* an example of crossing phenomenon obtained by imposing $\sqrt{\Delta_+} = 1$, $\sqrt{\Delta_-} = 0.2$, and $\rho = 0.1$. *Panel (b)* shows the double crossing, obtained by introducing an additional timescale to the previous case by tuning label imbalance.

References

- 220
- 221 [1] Rangachari Anand, Kishan G Mehrotra, Chilukuri K Mohan, and Sanjay Ranka. An improved
222 algorithm for neural network classification of imbalanced training sets. *IEEE transactions on*
223 *neural networks*, 4(6):962–969, 1993.
- 224 [2] Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk mini-
225 mization. *arXiv preprint arXiv:1907.02893*, 2019.
- 226 [3] Luca Arnaboldi, Ludovic Stephan, Florent Krzakala, and Bruno Loureiro. From high-
227 dimensional & mean-field dynamics to dimensionless odes: A unifying approach to sgd in
228 two-layers networks. In Gergely Neu and Lorenzo Rosasco, editors, *Proceedings of Thirty Sixth*
229 *Conference on Learning Theory*, volume 195 of *Proceedings of Machine Learning Research*,
230 pages 1199–1227. PMLR, 12–15 Jul 2023.
- 231 [4] Devansh Arpit, Stanisław Jastrzębski, Nicolas Ballas, David Krueger, Emmanuel Bengio,
232 Maxinder S Kanwal, Tegan Maharaj, Asja Fischer, Aaron Courville, Yoshua Bengio, et al. A
233 closer look at memorization in deep networks. In *International conference on machine learning*,
234 pages 233–242. PMLR, 2017.
- 235 [5] Samuel James Bell and Levent Sagun. Simplicity bias leads to amplified performance disparities.
236 In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*,
237 pages 355–369, 2023.
- 238 [6] Gerard Ben Arous, Reza Gheissari, and Aukosh Jagannath. High-dimensional limit theorems
239 for sgd: Effective dynamics and critical scaling. *Advances in Neural Information Processing*
240 *Systems*, 35:25349–25362, 2022.
- 241 [7] Richard Berk, Hoda Heidari, Shahin Jabbari, Michael Kearns, and Aaron Roth. Fairness in
242 criminal justice risk assessments: The state of the art. *Sociological Methods & Research*,
243 50(1):3–44, 2021.
- 244 [8] Michael Biehl and Holm Schwarze. Learning by on-line gradient descent. *Journal of Physics A:*
245 *Mathematical and general*, 28(3):643, 1995.
- 246 [9] Joy Buolamwini and Timnit Gebru. Gender shades: Intersectional accuracy disparities in
247 commercial gender classification. In *Conference on fairness, accountability and transparency*,
248 pages 77–91. PMLR, 2018.
- 249 [10] Farzan Farnia, Jesse Zhang, and David Tse. A spectral approach to generalization and optimiza-
250 tion in neural networks. 2018.
- 251 [11] Yunhe Feng and Chirag Shah. Has ceo gender bias really been fixed? adversarial attacking and
252 improving gender fairness in image search. In *Proceedings of the AAAI Conference on Artificial*
253 *Intelligence*, volume 36, pages 11882–11890, 2022.
- 254 [12] Emanuele Francazi, Aurelien Lucchi, and Marco Baity-Jesi. Initial guessing bias: How untrained
255 networks favor some classes. *arXiv preprint arXiv:2306.00809*, 2023.
- 256 [13] Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel,
257 Matthias Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *Nature*
258 *Machine Intelligence*, 2(11):665–673, 2020.
- 259 [14] Sebastian Goldt, Madhu Advani, Andrew M Saxe, Florent Krzakala, and Lenka Zdeborová.
260 Dynamics of stochastic gradient descent for two-layer neural networks in the teacher-student
261 setup. *Advances in neural information processing systems*, 32, 2019.
- 262 [15] Suriya Gunasekar, Jason D Lee, Daniel Soudry, and Nati Srebro. Implicit bias of gradient
263 descent on linear convolutional networks. *Advances in neural information processing systems*,
264 31, 2018.
- 265 [16] Guo Haixiang, Li Yijing, Jennifer Shang, Gu Mingyun, Huang Yuanyue, and Gong Bing.
266 Learning from class-imbalanced data: Review of methods and applications. *Expert systems*
267 *with applications*, 73:220–239, 2017.

- 268 [17] Katherine L Hermann, Hossein Mobahi, Thomas Fel, and Michael C Mozer. On the foundations
269 of shortcut learning. In *The Twelfth International Conference on Learning Representations*,
270 2024.
- 271 [18] Kimmo Karkkainen and Jungseock Joo. Fairface: Face attribute dataset for balanced race,
272 gender, and age for bias measurement and mitigation. In *Proceedings of the IEEE/CVF winter*
273 *conference on applications of computer vision*, pages 1548–1558, 2021.
- 274 [19] Matthew Kay, Cynthia Matuszek, and Sean A Munson. Unequal representation and gender
275 stereotypes in image search results for occupations. In *Proceedings of the 33rd annual acm*
276 *conference on human factors in computing systems*, pages 3819–3828, 2015.
- 277 [20] Bartosz Krawczyk. Learning from imbalanced data: open challenges and future directions.
278 *Progress in Artificial Intelligence*, 5(4):221–232, 2016.
- 279 [21] Yann LeCun, Léon Bottou, Genevieve B Orr, and Klaus-Robert Müller. Efficient backprop. In
280 *Neural networks: Tricks of the trade*, pages 9–50. Springer, 2002.
- 281 [22] Marc Lelarge and Léo Miolane. Asymptotic bayes risk for gaussian mixture in a semi-supervised
282 setting. In *2019 IEEE 8th International Workshop on Computational Advances in Multi-Sensor*
283 *Adaptive Processing (CAMSAP)*, pages 639–643. IEEE, 2019.
- 284 [23] Thibault Lesieur, Caterina De Bacco, Jess Banks, Florent Krzakala, Cris Moore, and Lenka
285 Zdeborová. Phase transitions and optimal algorithms in high-dimensional gaussian mixture cluster-
286 ing. In *2016 54th Annual Allerton Conference on Communication, Control, and Computing*
287 *(Allerton)*, pages 601–608. IEEE, 2016.
- 288 [24] Evan Z Liu, Behzad Haghgoo, Annie S Chen, Aditi Raghunathan, Pang Wei Koh, Shiori Sagawa,
289 Percy Liang, and Chelsea Finn. Just train twice: Improving group robustness without training
290 group information. In *International Conference on Machine Learning*, pages 6781–6792.
291 PMLR, 2021.
- 292 [25] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Large-scale celebfaces attributes
293 (celeba) dataset. *Retrieved August, 15(2018):11*, 2018.
- 294 [26] Karttikeya Mangalam and Vinay Uday Prabhu. Do deep neural networks learn shallow learnable
295 examples first? In *ICML 2019 Workshop on Identifying and Understanding Deep Learning*
296 *Phenomena*, 2019.
- 297 [27] Vaishnavh Nagarajan, Anders Andreassen, and Behnam Neyshabur. Understanding the fail-
298 ure modes of out-of-distribution generalization. In *International Conference on Learning*
299 *Representations*, 2021.
- 300 [28] Preetum Nakkiran, Gal Kaplun, Dimitris Kalimeris, Tristan Yang, Benjamin L Edelman, Fred
301 Zhang, and Boaz Barak. Sgd on neural networks learns functions of increasing complexity.
302 *Advances in Neural Information Processing Systems*, 33, 2020.
- 303 [29] Behnam Neyshabur, Ryota Tomioka, and Nathan Srebro. In search of the real inductive bias:
304 On the role of implicit regularization in deep learning. *arXiv preprint arXiv:1412.6614*, 2014.
- 305 [30] Nasim Rahaman, Aristide Baratin, Devansh Arpit, Felix Draxler, Min Lin, Fred Hamprecht,
306 Yoshua Bengio, and Aaron Courville. On the spectral bias of neural networks. In *International*
307 *conference on machine learning*, pages 5301–5310. PMLR, 2019.
- 308 [31] Maria Refinetti, Alessandro Ingrassia, and Sebastian Goldt. Neural networks trained with sgd
309 learn distributions of increasing complexity. In *International Conference on Machine Learning*,
310 pages 28843–28863. PMLR, 2023.
- 311 [32] David Saad and Sara Solla. Dynamics of on-line gradient descent learning for multilayer neural
312 networks. *Advances in neural information processing systems*, 8, 1995.
- 313 [33] Shiori Sagawa, Pang Wei Koh, Tatsunori B. Hashimoto, and Percy Liang. Distributionally
314 robust neural networks. In *International Conference on Learning Representations*, 2020.

- 315 [34] Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally
316 robust neural networks for group shifts: On the importance of regularization for worst-case
317 generalization. In *International Conference on Learning Representations (ICLR)*, 2020.
- 318 [35] Shiori Sagawa, Aditi Raghunathan, Pang Wei Koh, and Percy Liang. An investigation of
319 why overparameterization exacerbates spurious correlations. In *International Conference on*
320 *Machine Learning*, pages 8346–8356. PMLR, 2020.
- 321 [36] Stefano Sarao Mannelli, Federica Gerace, Negar Rostamzadeh, and Luca Saglietti. Un-
322 fair geometries: exactly solvable data model with fairness implications. *arXiv preprint*
323 *arXiv:2205.15935*, 2022.
- 324 [37] Harshay Shah, Kaustav Tamuly, Aditi Raghunathan, Prateek Jain, and Praneeth Netrapalli.
325 The pitfalls of simplicity bias in neural networks. *Advances in Neural Information Processing*
326 *Systems*, 33:9573–9585, 2020.
- 327 [38] Harini Suresh and John Guttag. A framework for understanding sources of harm throughout the
328 machine learning life cycle. In *Equity and access in algorithms, mechanisms, and optimization*,
329 pages 1–9. 2021.
- 330 [39] Zhiqin John Xu. Understanding training and generalization in deep learning by fourier analysis.
331 *arXiv preprint arXiv:1808.04295*, 2018.
- 332 [40] Yu Yang, Eric Gan, Gintare Karolina Dziugaite, and Baharan Mirzasoleiman. Identifying
333 spurious biases early in training through the lens of simplicity bias. In *International Conference*
334 *on Artificial Intelligence and Statistics*, pages 2953–2961. PMLR, 2024.

335 Appendix

336 Contents

337	A Further related works	11
338	B Problem setup and notation	11
339	C Main theorems and proofs	12
340	C.1 Order parameters	13
341	C.2 High dimensional dynamics	13
342	D Derivation of the ODEs	15
343	D.1 Useful Averages	15
344	D.2 ODEs	16
345	Student-shift overlap M	16
346	Student-teacher + overlap R_+	16
347	Student-teacher – overlap R_-	16
348	Self-overlap Q	16
349	Continuous limit.	19
350	E ODE solutions	20
351	E.1 General case	20
352	M	20
353	R_+ :	20
354	R_- :	20
355	Q :	20
356	E.2 Spurious correlations setting	21
357	E.3 Fairness setting	22
358	F Deeper analysis of the learning dynamics equations	22
359	F.1 Single centered cluster	22
360	F.2 Analysis of teacher alignment (τ_R) and student magnitude (τ_Q) timescales	22
361	F.3 Initial Preference	23
362	F.4 Asymptotic preference	24
363	G Additional numerical simulations	24
364	G.1 CIFAR10	24
365	G.2 CelebA	25
366	G.3 Simulations on Synthetic Data and Deeper Networks	27

367 **A Further related works**

368 **Class imbalance and fairness.** A key element in our study is the presence of heterogeneous
 369 data distributions within the dataset. In the context of fairness, these distributions model different
 370 groups in a population. Sampling unbalance is particularly critical, as minority groups are often
 371 misclassified [9, 18]. However, theoretical studies on group imbalance have been limited to asymptotic
 372 analyses [36], which may not apply in practical settings. Related questions have been explored in
 373 the label imbalance literature [20], where it has long been known [1, 16] that underrepresented
 374 classes have slower convergence rate and may even experience increased errors early in training. Our
 375 work shows that pre-asymptotic analysis can reveal complex transient dynamics, which is practically
 376 relevant when learning slows down or training to convergence is not possible. Similar to our analysis,
 377 [12] has shown that supposedly neutral choices, like activation functions or pooling operations, can
 378 generate strong biases. In contrast to prior work, our focus on data properties identifies several
 379 timescales associated to different data features relevant to bias generation.

380 **Simplicity bias.** Several studies [29, 15, 39, 10, 30] have highlighted a bias of deep neural networks
 381 (DNNs) towards *simple* solutions, suggesting this bias is a key to their generalisation performance.
 382 Simplicity bias also influences learning dynamics: [4, 30, 26, 28, 31] have showed that DNNs learn
 383 progressively more complex functions during training, with a notion of complexity often defined
 384 implicitly by other DNNs or observations like the time to memorisation. Our results connect with
 385 simplicity bias by identifying interpretable properties of the data that make samples appear “simple”
 386 to a shallow network. Interestingly, our findings reveal that different phases of learning experience
 387 simplicity in different ways, leading to forgetting of previously learned features.

388 **Spurious correlations.** Simplicity bias can also lead to shortcomings [37] by excessively relying
 389 of spurious features in the data, possibly hurting generalisation, especially in out-of-distribution
 390 contexts [13]. Theoretical works [27, 35, 17] have identified statistical properties that cause a
 391 classifier to favour spurious features over potentially more complex but more predictive features.
 392 Various methods have been proposed to address this problem using explicit partitioning of the data
 393 [2, 34]; some approaches implicitly infer subgroups with various degrees of correlation as spurious
 394 features. Notably, [24, 40] rely on early stages of learning to detect bias and adjust sample importance
 395 accordingly. Our study provides a unifying view of learning in fairness and spurious correlation
 396 problems, highlighting the presence of ephemeral biases characterised by multiple timescales during
 397 training. This adds complexity to the understanding of learning dynamics and points out potential
 398 confounding effects in existing mitigation methods.

399 **B Problem setup and notation**

400 We begin by refreshing the problem description and notation introduced in the main body for the two
 401 cluster case (Sec. ??) as well as defining some new notation to make the presentations of the results
 402 more compact.

- 403 1. (\mathbf{x}, y) denotes a training example with $\mathbf{x} \in \mathbb{R}^d$ and $y \in \{-1, 1\}$.
- 404 2. \mathbf{x} is drawn from a mixture of two Gaussians with means \mathbf{v}/\sqrt{d} and $-\mathbf{v}/\sqrt{d}$ respectively,
 405 covariances $\Delta_+ I_{d \times d}$ and $\Delta_- I_{d \times d}$ respectively. These two Gaussians are henceforth referred
 406 to as the positive and negative Gaussians respectively.
- 407 3. ρ represents the probability of the data being drawn from the positive Gaussian.
- 408 4. $\langle \cdot \rangle$ denotes an average over x , $\langle \cdot \rangle_{\oplus}$ denotes an average over the positive Gaussian and $\langle \cdot \rangle_{\ominus}$
 409 denotes an average over the negative Gaussian.
- 410 5. $\bar{\mathbf{w}}_+$ and $\bar{\mathbf{w}}_-$ denote the teachers for the positive Gaussian and negative Gaussian respectively.
 411 \mathbf{w} is the learnt classifier (“the student”).
- 412 6. The true labels, y , are then given by:
 - 413 • $y = \text{sign}(\bar{\mathbf{w}}_+ \cdot \mathbf{x}/\sqrt{d})$ for the positive cluster;
 - 414 • $y = \text{sign}(\bar{\mathbf{w}}_- \cdot \mathbf{x}/\sqrt{d})$ for the negative cluster.
- 415 7. Our predictions are $\hat{y} = \mathbf{w} \cdot \mathbf{x}/\sqrt{d}$.
- 416 8. The student is trained to minimise L2 loss $= (y - \hat{y})^2$.

- 417 9. The student learns using online stochastic gradient descent.
418 10. $\eta/2$ is the learning rate.
419 11. ϵ denotes the generalisation error.
420 12. $a \cdot b$ denotes the dot product between vectors a and b .
421 13. We now define the following Order Parameters (where only the first 4 change with training):
422 • $Q = \mathbf{w} \cdot \mathbf{w}/d$;
423 • $R_+ = \mathbf{w} \cdot \bar{\mathbf{w}}_+/d$;
424 • $R_- = \mathbf{w} \cdot \bar{\mathbf{w}}_-/d$;
425 • $M = \mathbf{w} \cdot \mathbf{v}/d$;
426 • $T_{\pm} = \bar{\mathbf{w}}_+ \cdot \bar{\mathbf{w}}_-/d$;
427 • $M_+^* = \bar{\mathbf{w}}_+ \cdot \mathbf{v}/d$;
428 • $M_-^* = \bar{\mathbf{w}}_- \cdot \mathbf{v}/d$;
429 • $v = \mathbf{v} \cdot \mathbf{v}/d$.
430 14. For algebraic simplicity, we assume $\|\bar{\mathbf{w}}_+\|_2 = \|\bar{\mathbf{w}}_-\|_2 = \sqrt{d}$ (and thus, $\bar{\mathbf{w}}_+ \cdot \bar{\mathbf{w}}_+/d = 1$
431 and $\bar{\mathbf{w}}_- \cdot \bar{\mathbf{w}}_-/d = 1$). This has the consequence that T_{\pm} exactly equals the cosine similarity
432 between the two teachers.
433 15. We also define $\Delta^{mix} = \rho\Delta_+ + (1 - \rho)\Delta_-$ and $\Delta^{2mix} = \rho\Delta_+^2 + (1 - \rho)\Delta_-^2$.
434 16. For notational convenience we define:

$$\alpha_+ = \langle y \rangle_{\oplus} = 1 - 2\Phi\left(\frac{-M_+^*}{\sqrt{\Delta_+}}\right), \quad (\text{B.13})$$

$$\alpha_- = \langle y \rangle_{\ominus} = 1 - 2\Phi\left(\frac{-(-M_-^*)}{\sqrt{\Delta_-}}\right). \quad (\text{B.14})$$

435 Note, α_+ also has an intuitive meaning. It represents the difference between the probability
436 that an example drawn from the positive cluster has positive true label and the probability
437 that an example drawn from the positive cluster has negative true label. It is hence 0 when
438 the positive cluster has equal positive and negative examples, positive when the cluster has
439 more positive examples than negative, negative when the cluster has more negative examples
440 than positive. Similarly, α_- represents the difference in these probabilities for the negative
441 cluster.

- 442 17. Finally, we also define

$$\beta_+ = \sqrt{\frac{2\Delta_+}{\pi}} \exp\left(\frac{-M_+^{*2}}{2\Delta_+}\right), \quad (\text{B.15})$$

$$\beta_- = \sqrt{\frac{2\Delta_-}{\pi}} \exp\left(\frac{-M_-^{*2}}{2\Delta_-}\right). \quad (\text{B.16})$$

- 443 18. Lastly, we use t to denote continuous time given by (epoch number/ d).

444 C Main theorems and proofs

445 In our study we analyse the linear model in Eq. 3 trained with online SGD on the data distribution
446 Eq.1 with the square loss $\hat{\epsilon} = (y - \hat{y})^2$. At the k -th iteration, a feature vector \mathbf{x}^k is sampled from (1),
447 the ground truth label y^k and current model prediction \hat{y}^k are respectively given by (2) and (3), and
448 the parameter is updated as:

$$\Delta\mathbf{w}^k := \mathbf{w}^{k+1} - \mathbf{w}^k = -\frac{\eta}{2}\nabla\hat{\epsilon}^k(\mathbf{w}^k) = \frac{\eta}{\sqrt{d}}(y^k - \hat{y}^k)\mathbf{x}^k \quad (\text{C.17})$$

449 where $\eta/2 > 0$ denotes the learning rate. Note that in this online setting, the number of time steps is
450 equivalent to the number of training examples.

451 **C.1 Order parameters**

452 The following lemma shows how the generalisation error depend on the model parameter \mathbf{w} only
 453 through the order parameters defined in Eq. 4.

454 **Lemma C.1.** *The generalisation error can be written as an average $\epsilon = \sum_{j=1}^m \rho_j \epsilon_j$ over the clusters,*
 455 *where ϵ_j is a degree 2 polynomial in R_j, M_j and Q taking the form*

$$\epsilon_j = 1 - 2\alpha_j M_j + M_j^2 - \beta_j R_j + Q\Delta_j \quad (\text{C.18})$$

456 *where α_j, β_j are constants independent of the parameter \mathbf{w} .*

457 *Proof.* Denote with $\langle \cdot \rangle_j$ the expectation over samples from cluster j . The generalisation error reads
 458 $\epsilon = \sum_{j=1}^m \rho_j \epsilon_j$ with

$$\begin{aligned} \epsilon_j &:= \langle (y - \hat{y})^2 \rangle_j = \left\langle \left(y - \frac{\mathbf{w} \cdot \mathbf{x}}{\sqrt{d}} \right)^2 \right\rangle_j = \langle y^2 \rangle_j + \left\langle \left(\frac{\mathbf{w} \cdot \mathbf{x}}{\sqrt{d}} \right)^2 \right\rangle_j - 2 \left\langle y \frac{\mathbf{w} \cdot \mathbf{x}}{\sqrt{d}} \right\rangle_j \\ &= 1 + (Q\Delta_j + M_j^2) - 2(\alpha_j M_j + R_j \beta_j), \end{aligned}$$

459 where the second term comes from: isolating the mean and the definition of M_j , and the isotropy of
 460 x . The third term comes from the useful identity *Integral 1* Eq. D.30, derived in Appendix D.1, and
 461 the constants are given by

$$\alpha_j = 1 - 2\Phi\left(\frac{-M_j^*}{\sqrt{\Delta_j}}\right), \quad \beta_j = \sqrt{\frac{2\Delta_j}{\pi}} \exp\left(\frac{-(M_j^*)^2}{2\Delta_j}\right). \quad (\text{C.19})$$

462 where $M_j^* := \bar{\mathbf{w}}_j^\top \mathbf{v}_j / d$ and $\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-u^2/2} du$ is the cumulative distribution function of
 463 the standard normal.

464 The formula for the generalisation error specializes to the case of two clusters with opposite means as

$$\begin{aligned} \epsilon &= 1 + M^2 - (2\rho\alpha_+ - 2(1 - \rho)\alpha_-) M \\ &\quad - 2\rho\beta_+ R_+ - 2(1 - \rho)\beta_- R_- + \Delta^{mix} Q, \end{aligned} \quad (\text{C.20})$$

465 Notably, α_\pm has an intuitive meaning wherein it represents the difference between the fraction of
 466 positive and negatives in a cluster, i.e., $\alpha_+ = \langle y \rangle_{c=+}$ and $\alpha_- = \langle y \rangle_{c=-}$. \square

467 Our problem thus reduces to characterising the evolution of order parameters (4). Using the gradient
 468 update of the parameter in Eq. C.17 and the notation $\delta^k := y^k - \hat{y}^k$, we can write update equations
 469 for the order parameters as follows:

$$\Delta M_j^k = \frac{\eta}{d} \delta^k \frac{\mathbf{v}_j^\top \mathbf{x}^k}{\sqrt{d}}, \quad \Delta R_j^k = \frac{\eta}{d} \delta^k \frac{\bar{\mathbf{w}}_j^\top \mathbf{x}^k}{\sqrt{d}}, \quad \Delta Q^k = \frac{2\eta}{d} \delta^k \frac{\mathbf{w}_j^\top \mathbf{x}^k}{\sqrt{d}} + \frac{\eta^2}{d^2} (\delta^k)^2 \|\mathbf{x}^k\|^2. \quad (\text{C.21})$$

470 **C.2 High dimensional dynamics**

471 We build upon classic results [32, 8], recently put on rigorous grounds [14, 6, 3], leveraging the
 472 *self-averaging* property of the order parameters in the high dimensional limit $d \rightarrow \infty$. As a result,
 473 as the dimension gets large, the discrete, stochastic evolution (C.21) of the order parameters can be
 474 effectively described in terms of the deterministic solution of the average continuous-time dynamics.

475 Let $\mathcal{S} := (S_i)_{1 \leq i \leq 2m+1}$ denote the collection of order parameters. The following lemma shows that
 476 the average of the updates (C.21) over the sample \mathbf{x}^k can be expressed solely in terms of \mathcal{S}^k .

477 **Lemma C.2.** $\mathbb{E}[\Delta S_i^k] = \frac{1}{d} f_i(\mathcal{S}^k)$ for some functions $(f_i(\mathcal{S}))_{1 \leq i \leq 2m+1}$ in $O(1)$ as $d \rightarrow \infty$.

478 *Proof.* Explicit computations are carried out in Appendix D.2 below for the case of two clusters. \square

479 The theorem below states that as d gets large, the stochastic evolution \mathcal{S}^k of the order parameter gets
 480 uniformly close, with high probability, to the average continuous-time dynamics described by the
 481 ODE system:

$$\frac{d\bar{S}_i(t)}{dt} = f_i(\bar{\mathcal{S}}(t)), \quad 1 \leq i \leq 2m+1, \quad (\text{C.22})$$

482 where the continuous *time* is given by the example number divided by the input dimension, $t = k/d$.
 483 Formally,

484 **Theorem C.3.** Fix a time horizon $T > 0$. For $1 \leq i \leq 2m + 1$,

$$\max_{0 \leq k \leq dT} |\mathcal{S}_i^k - \bar{\mathcal{S}}_i(k/d)| \xrightarrow{P} 0 \quad \text{as } d \rightarrow \infty. \quad (\text{C.23})$$

485 where \xrightarrow{P} denotes convergence in probability. A proof is provided in Appendix C. We provide the
 486 explicit expression of the functions f_i in the ODEs (C.22) in Appendix D, focusing on $m = 2$ clusters
 487 for clarity.

488 *Proof.* Using the notation of Section C.2 and assuming Lemma C.2, we examine the update equations
 489 (C.21) written as a stochastic iterative process

$$\mathcal{S}^{k+1} = \mathcal{S}^k + \mathbb{E} \frac{1}{d} f(\mathcal{S}^k) + \frac{1}{\sqrt{d}} \xi_d^k, \quad \xi_d^k := \sqrt{d}(\Delta \mathcal{S}^k - \mathbb{E}[\Delta \mathcal{S}^k]) \quad (\text{C.24})$$

490 where the expectation is over the new sample \mathbf{x}^k and conditional on the past samples. The noise term
 491 ξ_d^k has zero mean $\mathbb{E}[\xi_d^k] = 0$ and conditional covariance $\Sigma_d := \mathbb{E}[\xi_d^k \xi_d^{k\top}]$.

492 Define the continuous-time rescaled process $S_d(t)$ as the linear interpolation of $S^{\lfloor td \rfloor}$:

$$S_d(t) = S^{\lfloor td \rfloor} + (td - \lfloor td \rfloor)(S^{\lfloor td \rfloor + 1} - S^{\lfloor td \rfloor}) \quad (\text{C.25})$$

493 Here we leverage existing stochastic process convergence results (e.g., [6], Theorem 2.3) showing
 494 that, if Σ_d converges to the matrix valued function $\Sigma(\mathcal{S})$ as $d \rightarrow \infty$ in some appropriate sense, then
 495 the sequence $S_d(t)$ converges weakly as $d \rightarrow \infty$ to the solution \tilde{S}_t of the stochastic differential
 496 equation:

$$d\tilde{S}_t = f(\tilde{S}_t)dt + \sqrt{\Sigma(\tilde{S}_t)}dB_t \quad (\text{C.26})$$

497 where B_t is a standard Brownian motion in \mathbb{R}^{2m+1} . In our case, we can show that $\Sigma_d \in \mathcal{O}(d^{-1})$ as
 498 $d \rightarrow \infty$, so that $\Sigma = 0$ and Eq. C.26 reduces to the ODE in Eq. C.22.

499 Let us sketch the scaling argument. Algebraic manipulations similar to those in Section D.2 show that

$$\Sigma_d = \nabla \mathcal{S}^{k\top} \mathbb{E}[\Phi^k \Phi^{k\top}] \nabla \mathcal{S}^k (1 + \mathcal{O}(d^{-1})), \quad \Phi^k := \eta(\delta^k \mathbf{x}^k - \mathbb{E}[\delta^k \mathbf{x}^k]) \quad (\text{C.27})$$

500 where ∇ denotes the gradient with respect to the student vector \mathbf{w} . Recall that S^k has $2m$ components
 501 that are linear in \mathbf{w} (corresponding to the order parameters R_j and M_j in Eq. 4) and one that is
 502 quadratic (corresponding to Q). By making the gradients $\nabla \mathcal{S}^k$ explicit using Eq. 4), we see that at
 503 leading order, the matrix entries Σ_d^{ij} , $1 \leq i, j \leq 2m + 1$ take the form

$$\Sigma_d^{ij} = \frac{1}{d} \mathbb{E}[\Phi_{\mathbf{a}_i}^k \Phi_{\mathbf{a}_j}^{k\top}], \quad \Phi_{\mathbf{a}_i}^k = \eta(\delta^k \frac{\mathbf{a}_i^\top \mathbf{x}^k}{\sqrt{d}} - \mathbb{E}[\delta^k \frac{\mathbf{a}_i^\top \mathbf{x}^k}{\sqrt{d}}]) \quad (\text{C.28})$$

504 where the vector \mathbf{a}_i is either one of the teacher vectors $\bar{\mathbf{w}}_j$, one of the shift vector \mathbf{v}_j , or the student
 505 vector \mathbf{w} , depending on the entry $i = 1, \dots, 2m + 1$. As can be shown explicitly as in Appendix D.1
 506 below, $\Phi_{\mathbf{a}_i}^k$ depend on \mathbf{x}^k only through auxiliary variables $\bar{\mathbf{w}}_j^\top \mathbf{x} / \sqrt{d}$, $\bar{\mathbf{v}}_j^\top \mathbf{x} / \sqrt{d}$, $\mathbf{w}^{k\top} \mathbf{x}^k / \sqrt{d}$, which
 507 jointly follow a multivariate distribution whose parameters depend on the student vector \mathbf{w}^k only
 508 through \mathcal{S}^k and are in $\mathcal{O}(1)$ as $d \rightarrow \infty$. As a result, $\Sigma_d^{ij} \in \mathcal{O}(d^{-1})$.

509 Finally, the weak convergence of $S_d(t)_t$ to \tilde{S}_t implies convergence in probability for the supremum
 510 norm on the interval $[0, T]$ for any $T > 0$. Specifically, for each $1 \leq i \leq 2m + 1$,

$$\sup_{0 \leq t \leq T} |S_{di}(t) - \bar{S}_i(t)| \xrightarrow{P} 0, \quad (\text{C.29})$$

511 where \xrightarrow{P} denotes convergence in probability. This result directly leads to Eq. C.23, thereby proving
 512 the theorem. \square

513 D Derivation of the ODEs

514 In this section we are going to explicitly derive the ODE describing the dynamics of the order
 515 parameters. Starting from the discrete updates of the order parameters, Eqs. C.21, we are going to
 516 consider the thermodynamic limit, $d \rightarrow \infty$. As proven in Thm. C.3, the updates concentrate to their
 517 typical value and the discrete evolution converges to differential equations. Therefore, the rest of the
 518 section is devoted to performing averages over the Gaussians in order to evaluate the typical values.
 519 Before proceeding with the evaluation of Eqs. C.21, it is useful to introduce two identities.

520 D.1 Useful Averages

521 Integral 1:

$$\langle a \cdot x \operatorname{sign}(b \cdot x + c) \rangle = (a \cdot \mu) \left(1 - 2\Phi \left(\frac{-(b \cdot \mu + c)}{\sqrt{\Delta b \cdot b}} \right) \right) + a \cdot b \sqrt{\frac{2\Delta}{b \cdot b \pi}} \exp \left(\frac{-(b \cdot \mu + c)^2}{2\Delta b \cdot b} \right) \quad (\text{D.30})$$

522 where x is multivariate normal distribution with mean μ and covariance ΔI , and the angular bracket
 523 notation indicates average with respect to x .

524 *Derivation.* Define the auxiliary random variables $z_1 = a \cdot x$ and $z_2 = b \cdot x + c$, that follow a
 525 multivariate normal distribution

$$\begin{bmatrix} z_1 \\ z_2 \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} a \cdot \mu \\ b \cdot \mu + c \end{bmatrix}, \Delta \begin{bmatrix} a \cdot a & a \cdot b \\ a \cdot b & b \cdot b \end{bmatrix} \right).$$

526 Using the law of iterated expectation, our average can be written as:

$$\begin{aligned} \langle a \cdot x \operatorname{sign}(b \cdot x + c) \rangle &= \mathbb{E}_{z_2} [\operatorname{sign}(z_2) \mathbb{E}_{z_1|z_2} [z_1]] \\ &= \mathbb{E}_{z_2} [\operatorname{sign}(z_2) (a \cdot \mu + \frac{a \cdot b}{b \cdot b} (z_2 - (b \cdot \mu + c)))] \\ &= (a \cdot \mu - \frac{a \cdot b}{b \cdot b} (b \cdot \mu + c)) \mathbb{E}_{z_2} [\operatorname{sign}(z_2)] + \frac{a \cdot b}{b \cdot b} \mathbb{E}_{z_2} [z_2 \operatorname{sign}(z_2)] \end{aligned}$$

527 The first expectation follows from the definition of the cumulative distribution function Φ

$$\mathbb{E}_{z_2} [\operatorname{sign}(z_2)] = \left(1 - 2\Phi \left(\frac{-(b \cdot \mu + c)}{\sqrt{\Delta b \cdot b}} \right) \right).$$

528 The second term is simply the mean of a folded normal distribution

$$\mathbb{E}_{z_2} [z_2 \operatorname{sign}(z_2)] = (\sqrt{\Delta b \cdot b}) \sqrt{\frac{2}{\pi}} \exp \left(\frac{-(b \cdot \mu + c)^2}{2\Delta b \cdot b} \right) + (b \cdot \mu + c) \left(1 - 2\Phi \left(\frac{-(b \cdot \mu + c)}{\sqrt{\Delta b \cdot b}} \right) \right).$$

529 Combining these three expressions we obtain the identity.

530 Integral 2:

$$\langle a \cdot x b \cdot x \rangle = (a \cdot \mu)(b \cdot \mu) + \Delta(a \cdot b) \quad (\text{D.31})$$

531 where x is defined as for the previous identity.

532 *Derivation.* We proceed as in the previous case. Define the auxiliary random variables $z_1 = a \cdot x$ and
 533 $z_2 = b \cdot x$. They follow a multivariate normal distribution

$$\begin{bmatrix} z_1 \\ z_2 \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} a \cdot \mu \\ b \cdot \mu \end{bmatrix}, \Delta \begin{bmatrix} a \cdot a & a \cdot b \\ a \cdot b & b \cdot b \end{bmatrix} \right).$$

534 Using the law of iterated expectation, our average may be written as:

$$\begin{aligned} \langle a \cdot x b \cdot x \rangle &= \mathbb{E}_{z_2} [z_2 \mathbb{E}_{z_1|z_2} [z_1]] \\ &= \mathbb{E}_{z_2} [z_2 (a \cdot \mu + \frac{a \cdot b}{b \cdot b} (z_2 - (b \cdot \mu)))] \\ &= (a \cdot \mu - \frac{a \cdot b}{b \cdot b} (b \cdot \mu)) \mathbb{E}_{z_2} [z_2] + \frac{a \cdot b}{b \cdot b} \mathbb{E}_{z_2} [z_2^2] \\ &= (a \cdot \mu - \frac{a \cdot b}{b \cdot b} (b \cdot \mu))(b \cdot \mu) + \frac{a \cdot b}{b \cdot b} (\Delta b \cdot b + (b \cdot \mu)^2) \\ &= (a \cdot \mu)(b \cdot \mu) + \Delta(a \cdot b). \end{aligned}$$

535 **D.2 ODEs**

536 We have now the building blocks to evaluate the expected values of Eqs. C.21. We refresh the
 537 notation that $\delta^\mu = y^\mu - \hat{y}^\mu$, $y^\mu = \text{sign}(\mathbf{x}^\mu \cdot \bar{\mathbf{w}}_\mu / \sqrt{d})$, and $\hat{y}^\mu = \mathbf{x}^\mu \cdot \mathbf{w} / \sqrt{d}$. Final step is to take
 538 the continuous limit. This is obtained by noticing that the RHS of the equations is factorised by $1/d$.
 539 Therefore by taking as time unit $1/d$ and defining time as $t = \mu/d$ the discrete updates converge to
 540 continuous increments as $d \rightarrow \infty$.

Student-shift overlap M .

$$\langle \Delta M \rangle = \frac{\eta}{d} (\rho v \alpha_+ + \rho M_+^* \beta_+ - (1 - \rho) v \alpha_- + (1 - \rho) M_-^* \beta_- - (M(v + \Delta^{mix}))) \quad (\text{D.32})$$

541 *Derivation.* Starting from the definition in Eq. C.21 for M

$$\langle \Delta M \rangle = \frac{\eta}{d} \left(\left\langle y \frac{\mathbf{x} \cdot \mathbf{v}}{\sqrt{d}} \right\rangle - \left\langle \frac{\mathbf{w} \cdot \mathbf{x}}{\sqrt{d}} \frac{\mathbf{x} \cdot \mathbf{v}}{\sqrt{d}} \right\rangle \right).$$

542 The first term can be evaluated using integral 1 and the second term using integral 2 yielding the
 543 result.

Student-teacher + overlap R_+ .

$$\begin{aligned} \langle \Delta R_+ \rangle = \frac{\eta}{d} & \left(\rho(M_+^* \alpha_+ + \beta_+) + (1 - \rho)(-M_+^* \alpha_- + T_\pm \beta_-) \right. \\ & \left. - \rho(M M_+^* + R_+ \Delta_+) - (1 - \rho)(M M_+^* + R_+ \Delta_-) \right) \end{aligned} \quad (\text{D.33})$$

544 *Derivation.*

$$\begin{aligned} \langle \Delta R_+ \rangle &= \frac{\eta}{d} \left\langle \left(y - \frac{\mathbf{w} \cdot \mathbf{x}}{\sqrt{d}} \right) \left(\frac{\mathbf{x} \cdot \bar{\mathbf{w}}_+}{\sqrt{d}} \right) \right\rangle \\ &= \frac{\eta}{d} \left(\rho \left\langle y \frac{\mathbf{x} \cdot \bar{\mathbf{w}}_+}{\sqrt{d}} \right\rangle_{\oplus} + (1 - \rho) \left\langle y \frac{\mathbf{x} \cdot \bar{\mathbf{w}}_+}{\sqrt{d}} \right\rangle_{\ominus} - \rho \left\langle \frac{\mathbf{w} \cdot \mathbf{x}}{\sqrt{d}} \frac{\mathbf{x} \cdot \bar{\mathbf{w}}_+}{\sqrt{d}} \right\rangle_{\oplus} - (1 - \rho) \left\langle \frac{\mathbf{w} \cdot \mathbf{x}}{\sqrt{d}} \frac{\mathbf{x} \cdot \bar{\mathbf{w}}_+}{\sqrt{d}} \right\rangle_{\ominus} \right). \end{aligned}$$

545 These 4 terms can be computed using integrals 1 and 2 yielding the result.

Student-teacher – overlap R_- .

$$\begin{aligned} \langle \Delta R_- \rangle &= \frac{\eta}{d} \left(\rho(M_-^* \alpha_+ + T_\pm \beta_+) + (1 - \rho)(-M_-^* \alpha_- + \beta_-) \right. \\ & \left. - \rho(M M_-^* + R_- \Delta_+) - (1 - \rho)(M M_-^* + R_- \Delta_-) \right) \end{aligned} \quad (\text{D.34})$$

546 *Derivation.* Same as for R_+ .

Self-overlap Q .

$$\begin{aligned} \langle \Delta Q \rangle &= \frac{2\eta}{d} (\rho(\alpha_+ M + \beta_+ R_+) + (1 - \rho)(-\alpha_- M + \beta_- R_+) - M^2 - Q \Delta^{mix}) \\ & \quad + \frac{\eta^2}{d} \left(\Delta^{mix} + Q \Delta^{2mix} + M^2 \Delta^{mix} \right. \\ & \quad \left. - 2(\rho \Delta_+(\alpha_+ M + \beta_+ R_+) + (1 - \rho) \Delta_-(-\alpha_- M + \beta_- R_+)) \right). \end{aligned} \quad (\text{D.35})$$

547 *Derivation.* This update requires additional steps with respect to the previous ones.

$$\langle \Delta Q \rangle = \frac{2\eta}{d} \left\langle \delta \frac{\mathbf{w}_j^\top \mathbf{x}}{\sqrt{d}} \right\rangle + \frac{\eta^2}{d} \left\langle (\delta^\mu)^2 \frac{\|\mathbf{x}^\mu\|^2}{d} \right\rangle.$$

548 The first term is

$$\begin{aligned} \frac{2\eta}{d} \left\langle \delta \frac{\mathbf{w}_j^\top \mathbf{x}}{\sqrt{d}} \right\rangle &= \frac{2\eta}{d} \left\langle y \frac{\mathbf{w} \cdot \mathbf{x}}{\sqrt{d}} - \left(\frac{\mathbf{w} \cdot \mathbf{x}}{\sqrt{d}} \right)^2 \right\rangle \\ &= \frac{2\eta}{d} (M(\rho\alpha_+ - (1-\rho)\alpha_-) + \rho\beta_+R_+ + (1-\rho)\beta_-R_- - M^2 - Q\Delta^{mix}). \end{aligned}$$

549 The second term

$$\begin{aligned} \frac{\eta^2}{d} \left\langle (\delta^\mu)^2 \frac{\|\mathbf{x}^\mu\|^2}{d} \right\rangle &= \frac{\eta^2}{d} \left\langle \left(y - \frac{\mathbf{w} \cdot \mathbf{x}}{\sqrt{d}} \right)^2 \frac{\mathbf{x} \cdot \mathbf{x}}{d} \right\rangle \\ &= \frac{\eta^2}{d} \left\langle y^2 \frac{\mathbf{x} \cdot \mathbf{x}}{d} + \left(\frac{\mathbf{w} \cdot \mathbf{x}}{\sqrt{d}} \right)^2 \frac{\mathbf{x} \cdot \mathbf{x}}{d} - 2y \frac{\mathbf{w} \cdot \mathbf{x} \mathbf{x} \cdot \mathbf{x}}{\sqrt{d} d} \right\rangle \end{aligned}$$

550 requires additional steps. We consider the three terms in the expression above, starting from the first
551 one

$$\begin{aligned} \left\langle y^2 \frac{\mathbf{x} \cdot \mathbf{x}}{d} \right\rangle &= \left\langle \frac{\mathbf{x} \cdot \mathbf{x}}{d} \right\rangle = \frac{1}{d} \left(\sum_{i=1}^d \langle x_i^2 \rangle \right) = \frac{1}{d} \left(\sum_{i=1}^d \rho \langle x_i^2 \rangle_{\oplus} + (1-\rho) \langle x_i^2 \rangle_{\ominus} \right) \\ &= \frac{1}{d} \left(\sum_{i=1}^d \rho(\Delta_+ + v_i^2/d) + (1-\rho)(\Delta_- + v_i^2/d) \right) = \Delta^{mix} + v/d \\ &= \Delta^{mix} + O(d^{-1}), \end{aligned}$$

552 Where we used the simplification $y^2 = 1$ independently of the cluster's teacher. However, the
553 remaining terms require us to split the expectation considering the probability of sampling from each
554 cluster. The second term

$$\left\langle \frac{\mathbf{x} \cdot \mathbf{x}}{d} \left(\frac{\mathbf{w} \cdot \mathbf{x}}{\sqrt{d}} \right)^2 \right\rangle = \rho \left\langle \frac{\mathbf{x} \cdot \mathbf{x}}{d} \left(\frac{\mathbf{w} \cdot \mathbf{x}}{\sqrt{d}} \right)^2 \right\rangle_{\oplus} + (1-\rho) \left\langle \frac{\mathbf{x} \cdot \mathbf{x}}{d} \left(\frac{\mathbf{w} \cdot \mathbf{x}}{\sqrt{d}} \right)^2 \right\rangle_{\ominus}.$$

555 We begin by analysing the average over the positive Gaussian and split \mathbf{x} as $\mathbf{x} = \mathbf{v}/\sqrt{d} + \tilde{\mathbf{x}}$ such that
556 $\tilde{\mathbf{x}}$ has zero mean. Then,

$$\left\langle \frac{\mathbf{x} \cdot \mathbf{x}}{d} \left(\frac{\mathbf{w} \cdot \mathbf{x}}{\sqrt{d}} \right)^2 \right\rangle_{\oplus} = \left\langle \left[\frac{\mathbf{v} \cdot \mathbf{v}}{d^2} + \frac{2\mathbf{v} \cdot \tilde{\mathbf{x}}}{d\sqrt{d}} + \frac{\tilde{\mathbf{x}} \cdot \tilde{\mathbf{x}}}{d} \right] \left[\left(\frac{\mathbf{w} \cdot \mathbf{v}}{d} \right)^2 + 2 \frac{\mathbf{w} \cdot \mathbf{v} \mathbf{w} \cdot \tilde{\mathbf{x}}}{d\sqrt{d}} + \left(\frac{\mathbf{w} \cdot \tilde{\mathbf{x}}}{\sqrt{d}} \right)^2 \right] \right\rangle_{\oplus}$$

557 Multiplying the terms in the brackets will give rise to 9 terms. We can see that the 3+3=6 terms
558 corresponding to $\mathbf{v} \cdot \mathbf{v}/d^2$ and $2\mathbf{v} \cdot \tilde{\mathbf{x}}/d\sqrt{d}$ will tend to 0 in the limit of infinite d due to their scaling.
559 We now analyse the other 3 terms:

560

561 Term 1:

$$\begin{aligned} \left\langle \frac{\tilde{\mathbf{x}} \cdot \tilde{\mathbf{x}}}{d} \left(\frac{\mathbf{w} \cdot \mathbf{v}}{d} \right)^2 \right\rangle_{\oplus} &= \left(\frac{\mathbf{w} \cdot \mathbf{v}}{d} \right)^2 \left\langle \frac{\tilde{\mathbf{x}} \cdot \tilde{\mathbf{x}}}{d} \right\rangle_{\oplus} + O(d^{-1}) \\ &= M^2\Delta_+ + O(d^{-1}). \end{aligned}$$

562 Term 2:

$$\begin{aligned} 2 \left\langle \frac{\tilde{\mathbf{x}} \cdot \tilde{\mathbf{x}} \mathbf{w} \cdot \mathbf{v} \mathbf{w} \cdot \tilde{\mathbf{x}}}{d d \sqrt{d}} \right\rangle_{\oplus} &= 2R \left\langle \frac{\tilde{\mathbf{x}} \cdot \tilde{\mathbf{x}} \mathbf{w} \cdot \tilde{\mathbf{x}}}{d \sqrt{d}} \right\rangle_{\oplus} \\ &= 2R \left\langle \frac{\tilde{\mathbf{x}} \cdot \tilde{\mathbf{x}}}{d} \right\rangle_{\oplus} \left\langle \frac{\mathbf{w} \cdot \tilde{\mathbf{x}}}{\sqrt{d}} \right\rangle_{\oplus} + O(d^{-1}) \\ &= 0 + O(d^{-1}). \end{aligned}$$

563 Term 3:

$$\left\langle \frac{\tilde{\mathbf{x}} \cdot \tilde{\mathbf{x}}}{d} \left(\frac{\mathbf{w} \cdot \tilde{\mathbf{x}}}{\sqrt{d}} \right)^2 \right\rangle_{\oplus} = \left\langle \frac{\tilde{\mathbf{x}} \cdot \tilde{\mathbf{x}}}{d} \right\rangle_{\oplus} \left\langle \left(\frac{\mathbf{w} \cdot \tilde{\mathbf{x}}}{\sqrt{d}} \right)^2 \right\rangle_{\oplus} + O(d^{-1})$$

$$= \Delta_+(\Delta_+Q) + O(d^{-1}) = Q\Delta_+^2 + O(d^{-1}).$$

564

565

566 Thus finally,

$$\begin{aligned} \left\langle \frac{\mathbf{x} \cdot \mathbf{x}}{d} \left(\frac{\mathbf{w} \cdot \mathbf{x}}{\sqrt{d}} \right)^2 \right\rangle &= \rho(M^2\Delta_+ + Q\Delta_+^2) + (1 - \rho)(M^2\Delta_- + Q\Delta_-^2) \\ &= M^2\Delta^{mix} + Q\Delta^{2mix}. \end{aligned}$$

567

568

569 For the the third term

$$\left\langle y \frac{\mathbf{w} \cdot \mathbf{x} \mathbf{x} \cdot \mathbf{x}}{\sqrt{d} d} \right\rangle = \rho \left\langle y \frac{\mathbf{w} \cdot \mathbf{x} \mathbf{x} \cdot \mathbf{x}}{\sqrt{d} d} \right\rangle_{\oplus} + (1 - \rho) \left\langle y \frac{\mathbf{w} \cdot \mathbf{x} \mathbf{x} \cdot \mathbf{x}}{\sqrt{d} d} \right\rangle_{\ominus}.$$

570 As before, we analyse the average over the positive Gaussian first and split \mathbf{x} into its mean and a zero

571 mean component:

$$\left\langle y \frac{\mathbf{x} \cdot \mathbf{x} \mathbf{w} \cdot \mathbf{x}}{d \sqrt{d}} \right\rangle_{\oplus} = \left\langle y \left[\frac{\mathbf{v} \cdot \mathbf{v}}{d^2} + \frac{2\mathbf{v} \cdot \tilde{\mathbf{x}}}{d\sqrt{d}} + \frac{\tilde{\mathbf{x}} \cdot \tilde{\mathbf{x}}}{d} \right] \left[\frac{\mathbf{w} \cdot \mathbf{v}}{d} + \frac{\mathbf{w} \cdot \tilde{\mathbf{x}}}{\sqrt{d}} \right] \right\rangle_{\oplus}.$$

572 This gives rise to 6 terms. We can see that the 2+2=4 terms corresponding to $\mathbf{v} \cdot \mathbf{v}/d^2$ and $2\mathbf{v} \cdot \tilde{\mathbf{x}}/d\sqrt{d}$
573 will tend to 0 in the limit of infinite d due to their scaling. We now analyse the other 2 terms:

574

575 Term 1:

$$\begin{aligned} \left\langle y \frac{\tilde{\mathbf{x}} \cdot \tilde{\mathbf{x}} \mathbf{w} \cdot \mathbf{v}}{d d} \right\rangle_{\oplus} &= M \left\langle y \frac{\tilde{\mathbf{x}} \cdot \tilde{\mathbf{x}}}{d} \right\rangle_{\oplus} \\ &= M \left\langle \text{sign} \left(\frac{\tilde{\mathbf{x}} \cdot \bar{\mathbf{w}}_+}{\sqrt{d}} + \frac{\bar{\mathbf{w}}_+ \cdot \mathbf{v}}{d} \right) \frac{\tilde{\mathbf{x}} \cdot \tilde{\mathbf{x}}}{d} \right\rangle_{\oplus} \\ &= M \left\langle \text{sign} \left(\frac{\tilde{\mathbf{x}} \cdot \bar{\mathbf{w}}_+}{\sqrt{d}} + \frac{\bar{\mathbf{w}}_+ \cdot \mathbf{v}}{d} \right) \right\rangle_{\oplus} \left\langle \frac{\tilde{\mathbf{x}} \cdot \tilde{\mathbf{x}}}{d} \right\rangle_{\oplus} + O(d^{-1}) \\ &= M \langle y \rangle_{\oplus} \Delta_+ + O(d^{-1}) \\ &= M\alpha_+ \Delta_+ + O(d^{-1}). \end{aligned}$$

576

577

578 Term 2:

$$\begin{aligned} \left\langle y \left(\frac{\tilde{\mathbf{x}} \cdot \tilde{\mathbf{x}}}{d} \right) \left(\frac{\mathbf{w} \cdot \tilde{\mathbf{x}}}{\sqrt{d}} \right) \right\rangle_{\oplus} &= \left\langle y \left(\frac{\mathbf{w} \cdot \tilde{\mathbf{x}}}{\sqrt{d}} \right) \right\rangle_{\oplus} \left\langle \left(\frac{\tilde{\mathbf{x}} \cdot \tilde{\mathbf{x}}}{d} \right) \right\rangle_{\oplus} + O(d^{-1}) \\ &= \Delta_+ \left\langle y \left(\frac{\mathbf{w} \cdot \tilde{\mathbf{x}}}{\sqrt{d}} \right) \right\rangle_{\oplus} + O(d^{-1}) \\ &= \Delta_+ R_+ \beta_+ + O(d^{-1}). \end{aligned}$$

579 Where the last equality follows using integral 1. Thus:

$$\left\langle y \frac{\mathbf{x} \cdot \mathbf{x} \mathbf{w} \cdot \mathbf{x}}{d \sqrt{d}} \right\rangle_{\oplus} = \Delta_+(\alpha_+M + \beta_+R_+) + O(d^{-1}).$$

580 We repeat the same analysis for the negative gaussian and get:

$$\left\langle y \frac{\mathbf{x} \cdot \mathbf{x} \mathbf{w} \cdot \mathbf{x}}{d \sqrt{d}} \right\rangle_{\ominus} = \rho\Delta_+(\alpha_+M + \beta_+R_+) + (1 - \rho)\Delta_-(-\alpha_-M + \beta_-R_+) + O(d^{-1}).$$

581 Collecting everything together and taking the infinite dimensional limit:

$$\langle \Delta \mathbf{w} \cdot \Delta \mathbf{w} / d \rangle = \frac{\eta^2}{d} (\Delta^{mix} + Q\Delta^{2mix} + M^2\Delta^{mix} - 2(\rho\Delta_+(\alpha_+M + \beta_+R_+) + (1 - \rho)\Delta_-(-\alpha_-M + \beta_-R_+)))$$

582 Thus,

$$\begin{aligned} \langle \Delta Q \rangle &= \frac{2\eta}{d} (\rho(\alpha_+M + \beta_+R_+) + (1 - \rho)(-\alpha_-M + \beta_-R_+) - M^2 - Q\Delta^{mix}) \\ &\quad + \frac{\eta^2}{d} (\Delta^{mix} + Q\Delta^{2mix} + M^2\Delta^{mix} - 2(\rho\Delta_+(\alpha_+M + \beta_+R_+) + (1 - \rho)\Delta_-(-\alpha_-M + \beta_-R_+))). \end{aligned}$$

583 **Continuous limit.** Final step of the derivation is taking the thermodynamics limit that leads to the
584 ODEs implicitly defined in Thm. C.3:

$$f_M(M, R_+, R_-, Q) = \eta \left(\rho v \alpha_+ + \rho M_+^* \beta_+ - (1 - \rho) v \alpha_- + (1 - \rho) M_-^* \beta_- - (M(v + \Delta^{mix})) \right), \quad (\text{D.36})$$

$$f_{R_+}(M, R_+, R_-, Q) = \eta \left(\rho(M_+^* \alpha_+ + \beta_+) + (1 - \rho)(-M_+^* \alpha_- + T_{\pm} \beta_-) - \rho(MM_+^* + R_+ \Delta_+) - (1 - \rho)(MM_+^* + R_+ \Delta_-) \right), \quad (\text{D.37})$$

$$f_{R_-}(M, R_+, R_-, Q) = \eta \left(\rho(M_-^* \alpha_+ + T_{\pm} \beta_+) + (1 - \rho)(-M_-^* \alpha_- + \beta_-) - \rho(MM_-^* + R_- \Delta_+) - (1 - \rho)(MM_-^* + R_- \Delta_-) \right), \quad (\text{D.38})$$

$$f_Q(M, R_+, R_-, Q) = 2\eta \left(\rho(\alpha_+ M + \beta_+ R_+) + (1 - \rho)(-\alpha_- M + \beta_- R_+) - M^2 - Q \Delta^{mix} \right) + \eta^2 \left(\Delta^{mix} + Q \Delta^{2mix} + M^2 \Delta^{mix} - 2(\rho \Delta_+ (\alpha_+ M + \beta_+ R_+) + (1 - \rho) \Delta_- (-\alpha_- M + \beta_- R_+)) \right). \quad (\text{D.39})$$

585 **E ODE solutions**

586 In this section we first present the general solutions of the ODEs sketched in Theorem 3.1, then we
587 specialise to the two scenarios discussed in the main text.

588 **E.1 General case**

589 From the previous section, we have a system of coupled ODEs for the order parameters of the form:

$$\begin{aligned}\frac{dM}{dt} &= c_1 + c_2M, \\ \frac{dR_-}{dt} &= c_{3-} + c_{4-}M + c_{5-}R_-, \\ \frac{dR_+}{dt} &= c_{3+} + c_{4+}M + c_{5+}R_+, \\ \frac{dQ}{dt} &= c_6 + c_7M + c_8M^2 + c_{9+}R_+ + c_{9-}R_- + c_{10}Q.\end{aligned}$$

590 This represent a linear system of ODEs which can be solved using standard methods like Laplace
591 transform, leading to Eqs. 6-8. We now report the equations including the exact expression of their
592 coefficients.

593 M :

$$M(t) = M_0 e^{-t\eta(v+\Delta^{mix})} + M_\infty (1 - e^{-t\eta(v+\Delta^{mix})}). \quad (\text{E.40})$$

594 Where,

$$M_\infty = \frac{(\rho M_+^* \beta_+ + (1-\rho)M_-^* \beta_-) + v(\rho\alpha_+ - (1-\rho)\alpha_-)}{v + \Delta^{mix}}. \quad (\text{E.41})$$

R_+ :

$$R_+(t) = R_+^0 e^{-t\eta\Delta^{mix}} + R_+^\infty (1 - e^{-t\eta\Delta^{mix}}) + k_{1+} (e^{-t\eta\Delta^{mix}} - e^{-t\eta(v+\Delta^{mix})}). \quad (\text{E.42})$$

595 Where,

$$R_+^\infty = \frac{(\rho\beta_+ + T_\pm(1-\rho)\beta_-) + M_+^*(\rho\alpha_+ - (1-\rho)\alpha_- - M_\infty)}{\Delta^{mix}}, \quad (\text{E.43})$$

$$k_{1+} = \frac{M_+^*(M_\infty - M_0)}{v}. \quad (\text{E.44})$$

R_- :

$$R_-(t) = R_-^0 e^{-t\eta\Delta^{mix}} + R_-^\infty (1 - e^{-t\eta\Delta^{mix}}) + k_{1-} (e^{-t\eta\Delta^{mix}} - e^{-t\eta(v+\Delta^{mix})}). \quad (\text{E.45})$$

596 Where,

$$R_-^\infty = \frac{(T_\pm\rho\beta_+ + (1-\rho)\beta_-) + M_-^*(\rho\alpha_+ - (1-\rho)\alpha_- - M_\infty)}{\Delta^{mix}}, \quad (\text{E.46})$$

$$k_{1-} = \frac{M_-^*(M_\infty - M_0)}{v}. \quad (\text{E.47})$$

Q :

$$\begin{aligned}Q(t) &= Q_0 e^{-t\eta(2\Delta^{mix} - \eta\Delta^{2mix})} + Q_\infty (1 - e^{-t\eta(2\Delta^{mix} - \eta\Delta^{2mix})}) \\ &\quad + k_2 (e^{-t\eta(2\Delta^{mix} - \eta\Delta^{2mix})} - e^{-t\eta\Delta^{mix}}) \\ &\quad + k_3 (e^{-t\eta(2\Delta^{mix} - \eta\Delta^{2mix})} - e^{-t\eta(v+\Delta^{mix})}) \\ &\quad + k_4 (e^{-t\eta(2\Delta^{mix} - \eta\Delta^{2mix})} - e^{-t\eta(2v+2\Delta^{mix})}).\end{aligned} \quad (\text{E.48})$$

597 Where,

$$Q_\infty = \frac{\eta\Delta^{mix} + 2\rho\beta_+ R_+^\infty (1 - \eta\Delta_+) + 2(1-\rho)\beta_- R_-^\infty (1 - \eta\Delta_-)}{2\Delta^{mix} - \eta\Delta^{2mix}},$$

$$+ \frac{M_\infty(M_\infty(\eta\Delta^{mix} - 2) + 2\rho\alpha_+(1 - \eta\Delta_+) - 2(1 - \rho)\alpha_-(1 - \eta\Delta_-))}{2\Delta^{mix} - \eta\Delta^{2mix}}, \quad (\text{E.49})$$

$$k_2 = \frac{2\rho\beta_+(1 - \eta\Delta_+)(R_+^\infty - R_+^0 - k_{1+}) + 2(1 - \rho)\beta_-(1 - \eta\Delta_-)(R_-^\infty - R_-^0 - k_{1-})}{\Delta^{mix} - \eta\Delta^{2mix}}, \quad (\text{E.50})$$

$$k_3 = \frac{2\rho\beta_+(1 - \eta\Delta_+)k_{1+} + 2(1 - \rho)\beta_-(1 - \eta\Delta_-)k_{1-}}{\Delta^{mix} - \eta\Delta^{2mix} + v},$$

$$+ \frac{(M_\infty - M_0)(M_\infty(\eta\Delta^{mix} - 2) + 2\rho\alpha_+(1 - \eta\Delta_+) - 2(1 - \rho)\alpha_-(1 - \eta\Delta_-))}{\Delta^{mix} - \eta\Delta^{2mix} + v}, \quad (\text{E.51})$$

$$k_4 = \frac{(\eta\Delta^{mix} - 2)(M_\infty - M_0)^2}{\eta\Delta^{2mix} + 2v}. \quad (\text{E.52})$$

598 E.2 Spurious correlations setting

599 Under the setting discussed in the Sec. 4.1 ($\rho = 0.5, \Delta_+ = \Delta_- = \Delta, T_\pm = 1$), we can make the
600 following simplifications:

- 601 1. $\Delta^{mix} = \Delta,$
- 602 2. $\Delta^{2mix} = \Delta^2,$
- 603 3. $\alpha_+ = -\alpha_- = \alpha,$
- 604 4. $\beta_+ = \beta_- = \beta,$
- 605 5. $M_+^* = M_-^* = M^*,$
- 606 6. $R_+ = R_- = R.$

607 The equations then take the form:

$$M(t) = M_0 e^{-t\eta(v+\Delta)} + M_\infty(1 - e^{-t\eta(v+\Delta)}),$$

$$R(t) = R^0 e^{-t\eta\Delta} + R^\infty(1 - e^{-t\eta\Delta}) + k_1(e^{-t\eta\Delta} - e^{-t\eta(v+\Delta)}),$$

$$Q(t) = Q_0 e^{-t\eta(2\Delta - \eta\Delta^2)} + Q_\infty(1 - e^{-t\eta(2\Delta - \eta\Delta^2)})$$

$$+ k_2(e^{-t\eta(2\Delta - \eta\Delta^2)} - e^{-t\eta\Delta})$$

$$+ k_3(e^{-t\eta(2\Delta - \eta\Delta^2)} - e^{-t\eta(v+\Delta)})$$

$$+ k_4(e^{-t\eta(2\Delta - \eta\Delta^2)} - e^{-t\eta(2v+2\Delta)}).$$

608 Where,

$$M_\infty = \frac{M^*\beta + v\alpha}{v + \Delta},$$

$$R_\infty = \frac{\beta + M^*(\alpha - M_\infty)}{\Delta},$$

$$k_1 = \frac{(M_\infty - M_0)}{v},$$

$$Q_\infty = \frac{\eta\Delta + 2\beta R_\infty(1 - \eta\Delta)}{2\Delta - \eta\Delta^2} + \frac{M_\infty(M_\infty(\eta\Delta - 2) + 2\alpha(1 - \eta\Delta))}{2\Delta - \eta\Delta^2},$$

$$k_2 = \frac{2\beta(1 - \eta\Delta)(R_\infty - R_0 - k_1)}{\Delta - \eta\Delta^2},$$

$$k_3 = \frac{2\beta(1 - \eta\Delta)k_1}{\Delta - \eta\Delta^2 + v} + \frac{(M_\infty - M_0)(M_\infty(\eta\Delta - 2) + 2\alpha(1 - \eta\Delta))}{\Delta - \eta\Delta^2 + v},$$

$$k_4 = \frac{(\eta\Delta - 2)(M_\infty - M_0)^2}{\eta\Delta^2 + 2v}.$$

609 **E.3 Fairness setting**

610 The general fairness case coincides with the general case discussed above (E.1), therefore we limit
611 our discussion to the simplified case with centered clusters.

612 Under the zero shift $v = 0$, the equations take the simplified form wherein M, v, M_{\pm}^* are 0, the
613 transient term in R_{\pm} vanishes and Q only has one transient term. Specifically:

$$\begin{aligned} R_+(t) &= R_+^0 e^{-t\eta\Delta^{mix}} + R_+^{\infty}(1 - e^{-t\eta\Delta^{mix}}), \\ R_-(t) &= R_-^0 e^{-t\eta\Delta^{mix}} + R_-^{\infty}(1 - e^{-t\eta\Delta^{mix}}), \\ Q(t) &= Q_0 e^{-t\eta(2\Delta^{mix} - \eta\Delta^{2mix})} + Q_{\infty}(1 - e^{-t\eta(2\Delta^{mix} - \eta\Delta^{2mix})}) + Q_{trans}(e^{-t\eta(2\Delta^{mix} - \eta\Delta^{2mix})} - e^{-t\eta\Delta^{mix}}). \end{aligned}$$

614 Where

$$\begin{aligned} R_+^{\infty} &= \sqrt{\frac{2}{\pi}} \frac{\rho\sqrt{\Delta_+} + T_{\pm}(1 - \rho)\sqrt{\Delta_-}}{\Delta^{mix}}, \\ R_-^{\infty} &= \sqrt{\frac{2}{\pi}} \frac{T_{\pm}\rho\sqrt{\Delta_+} + (1 - \rho)\sqrt{\Delta_-}}{\Delta^{mix}}, \\ Q_{\infty} &= \frac{\eta\Delta^{mix} + 2\sqrt{\frac{2}{\pi}}\rho\sqrt{\Delta_+}R_+^{\infty}(1 - \eta\Delta_+) + 2\sqrt{\frac{2}{\pi}}(1 - \rho)\sqrt{\Delta_-}R_-^{\infty}(1 - \eta\Delta_-)}{2\Delta^{mix} - \eta\Delta^{2mix}}, \\ Q_{trans} &= \sqrt{\frac{2}{\pi}} \frac{2\rho\sqrt{\Delta_+}(1 - \eta\Delta_+)(R_+^{\infty} - R_+^0) + 2(1 - \rho)\sqrt{\Delta_-}(1 - \eta\Delta_-)(R_-^{\infty} - R_-^0)}{\Delta^{mix} - \eta\Delta^{2mix}}. \end{aligned}$$

615 **F Deeper analysis of the learning dynamics equations**

616 This section provides insights into the learning dynamics — particularly those relevant to bias
617 evolution — that arise out of the expressions for order parameter evolution. We shall provide intuitive
618 explanations behind the various mathematical terms that appear.

619 **F.1 Single centered cluster**

620 Consider first a single cluster centered at the origin—i.e. $\rho = 1, v = 0$ with variance Δ . In this setting,
621 the minimum generalisation error is achieved when the student perfectly aligns with the teacher and
622 optimises its norm such that $Q_{opt} = \frac{2}{\pi\Delta}$, achieving the generalisation error $\epsilon_{\min} = 1 - \frac{2}{\pi}$.

623 Importantly, this is not 0 since the student and the teacher are mismatched—i.e. the student is linear
624 whereas the teacher has a $sign(\cdot)$ activation function. From the equations, we observe that the
625 asymptotic generalisation error when training using online stochastic gradient descent in this setting
626 is

$$\epsilon_{\infty} = \frac{1 - 2/\pi}{1 - \eta\Delta/2} = \left(1 - \frac{2}{\pi}\right) \left(1 + \frac{\eta\Delta}{2} + O(\eta^2\Delta^2)\right). \quad (\text{F.53})$$

627 Thus, as the learning rate increases, the generalisation error increases until it reaches the critical
628 learning rate beyond which training is unstable and the loss grows unboundedly. In the single cluster
629 case, Eq. F.53 this is $2/\Delta$ which matches the classical result from convex optimisation [21]. We can
630 similarly find the critical learning rate for two clusters to be $2\Delta^{mix}/\Delta^{2mix}$ by ensuring exponential
631 terms decay to zero in equation 8.

632 **F.2 Analysis of teacher alignment (τ_R) and student magnitude (τ_Q) timescales**

633 We now consider the fairness setting with zero shift as illustrated in Fig. 1c. As discussed in section
634 4.2, the relevant timescales in this setting are

$$\tau_R = \frac{1}{\eta\Delta^{mix}}, \quad \tau_Q = \frac{1}{\eta(2\Delta^{mix} - \eta\Delta^{2mix})},$$

635 since $M(t)$ is always zero. Fig. 6 shows the crossing phenomena of the loss curves along with the
636 order parameter evolution and other insightful terms. The alignment of the student is governed by the

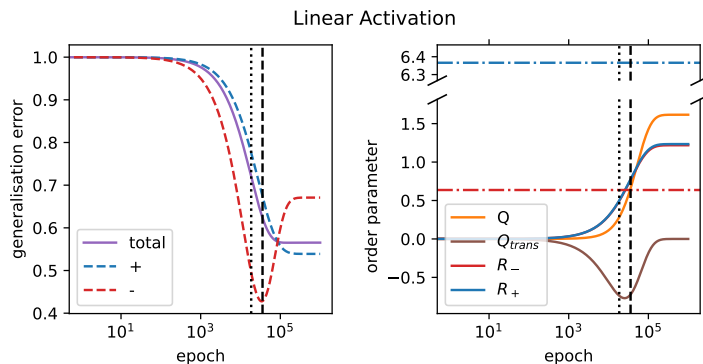


Figure 6: **The Crossing Phenomenon** The left shows the ‘crossing’ of the loss curves on the negative sub-population in red (higher variance and lower representation) and positive sub-population in blue (lower variance but greater representation) along with the overall loss in purple obtained as a weighted average of the two. It also marks τ_R as the dashed vertical line and τ_Q as the dotted vertical line. The right side shows the evolution of the order parameters and a transient term. The horizontal blue and red dash-dotted line mark the optimal value of Q for the positive-subpopulation and negative sub-populations respectively. The parameters are $v = 0, \rho = 0.8, \Delta_+ = 0.1, \Delta_- = 1, T_{\pm} = 0.9, \eta = 0.1$.

637 timescale τ_R and the change in its magnitude is governed by the timescale τ_Q . Initially, the classifier
638 has a small magnitude and its alignment roughly matches the two teachers which are themselves quite
639 similar ($T_{\pm} = 0.9$). Indeed, we see that the R_+ and R_- have very similar trajectories. However,
640 smaller magnitudes advantage higher variances as discussed in Appendix F.1 (Q_{opt} is inversely
641 proportional to the cluster variance).

642 We mark the optimal values of Q using horizontal lines in Fig.6 on the left side with blue for the
643 positive sub-population (lower variance) and red for the negative sub-population (higher variance). As
644 the magnitude of the student grows, we observe a sharp drop in the generalisation error on the higher
645 variance sub-population till Q crosses the horizontal red line. Beyond this point, the generalisation
646 error on the higher variance sub-population rises since the magnitude of the student has exceeded the
647 optimal value (horizontal red line) and the generalisation error on the lower variance sub-population
648 continues to fall as the magnitude of the student approaches the horizontal blue line. Finally, an
649 inspection of the timescales reveals that τ_Q (vertical dotted line) is less than t_R (vertical dashed
650 line) and hence we may expect the student magnitude to saturate before its alignment. However,
651 Q_{trans} , the transient term associated with Q (third line of equation 8), is always negative and hence
652 suppresses the growth of Q initially.

653 In summary, we observe a two phase behaviour. First the student shifts its alignment and increases
654 magnitude leading to a sharper drop in the higher variance generalisation error. Second, we observe
655 that as the student continues increasing magnitude while keeping its alignment fixed, it advantages
656 the lower variance cluster.

657 F.3 Initial Preference

658 Starting from a small initialisation, the initial rate of change of the generalisation error for sub-
659 population + is

$$\frac{d\epsilon_{g+}}{dt} \Big|_{t=0} = -\eta^2 \Delta^{mix} \Delta_+ \left(\sqrt{\frac{2}{\pi} \frac{R_+^{\infty}}{\Delta_+}} \frac{1}{\eta} - 1 \right) \quad (\text{F.54})$$

660 and analogously for $-$. The learning rate η must be chosen to be small enough such that the
661 generalisation errors decrease and hence the first term in the brackets must dominate over the 1.
662 Since $R_{\pm}^{\infty}/R_{\pm}^{\infty} \in [T_{\pm}; 1/T_{\pm}]$ (for $T_{\pm} > 0$), the ratio between generalisation error rates is therefore
663 bounded by

$$T_{\pm} \sqrt{\frac{\Delta_+}{\Delta_-}} \leq \frac{d\epsilon_{g+}/dt|_{t=0}}{d\epsilon_{g-}/dt|_{t=0}} \leq \frac{1}{T_{\pm}} \sqrt{\frac{\Delta_+}{\Delta_-}}. \quad (\text{F.55})$$

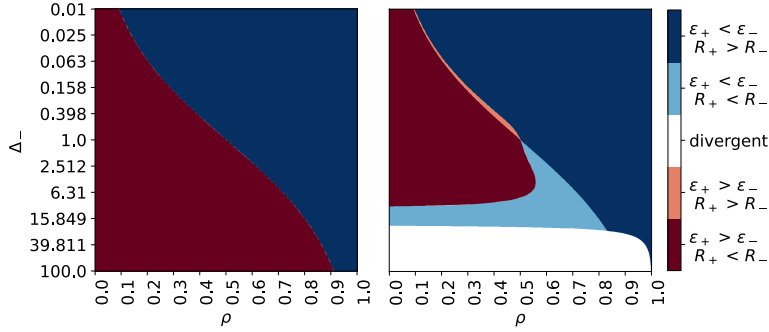


Figure 7: **Initial and Asymptotic student preferences** We set $v = 0, \Delta_+ = 1, T_{\pm} = 0.9, \eta = 0.1$ and study the values of ρ, Δ_- . The figure studies only asymptotic preferences under $v = 0, \Delta_+ = 1, T_{\pm} = 0.9$. When the learning rate is small ($\eta \rightarrow 0^+$ on *left side*), the cluster which has better alignment with the teacher must also have lower generalisation error. However, for non-zero learning rates ($\eta = 0.1$ on *right side*), behaviour is more complicated leading to the light colored phases where despite better asymptotic alignment with the teacher, the generalisation error is higher. Parameters: $\eta \rightarrow 0^+$ (left) vs $\eta = 0.1$ (right).

664 F.4 Asymptotic preference

665 This section discusses the asymptotic generalisation errors of our classifier when $v = 0$ as a function
 666 of representation and variances. Firstly, as discussed in section 4.2,

$$R_+^{\infty} > R_-^{\infty} \iff \rho\sqrt{\Delta_+} > (1 - \rho)\sqrt{\Delta_-}.$$

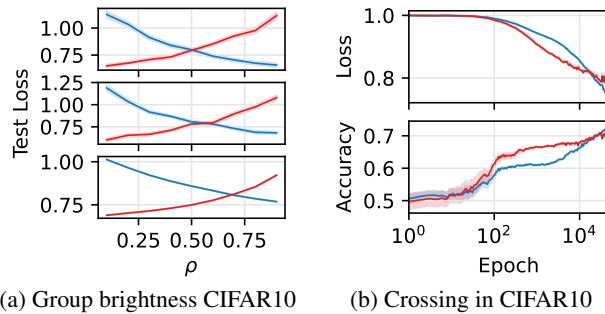
667 Intuitively, one might expect that the asymptotically lower generalisation error is achieved on the
 668 population whose teacher has better asymptotic alignment with the student. Indeed, when the learning
 669 rate tends to 0, we observe exactly this as illustrated by the two dark phases in Fig. 7 on the left
 670 side. However, when the learning rate is greater than zero, we observe more complex behaviour.
 671 Fig. 7 (*right*) shows the emergence two new phases (light red and light blue) wherein the classifier
 672 exhibits higher generalisation error on a sub-population despite having better alignment with its
 673 corresponding teacher. This behaviour can be traced back to equation F.53 wherein the increase in
 674 asymptotic generalisation error due to non-zero learning rates is amplified by the cluster variance.
 675 Thus, our analysis shows how a large learning rate can also become a source of bias in our classifier
 676 by advantaging the sub-population with smaller variance.

677 G Additional numerical simulations

678 G.1 CIFAR10

679 We consider the same architecture
 680 and pre-processing described for
 681 MNIST in Sec. 5 on a CIFAR10 clas-
 682 sification task. We select 8 classes
 683 and assign 4 of them to the pos-
 684 itive group and 4 to the negative
 685 group. Inside each group, 2 classes
 686 are labelled as negative and 2 as
 687 positive. This simulation frame-
 688 work is similar to the one considered
 689 by [5] where the authors used sub-
 690 populations with only 2 classes each.

691 The average brightness of the sam-
 692 ples in each cluster plays the same
 693 role as the parameter Δ in the syn-
 694 thetic model. Our theory predicts



(a) Group brightness CIFAR10

(b) Crossing in CIFAR10

Figure 8: **Numerical simulations on CIFAR10.** The figure shows experiments of a 2L neural network on CIFAR10 where classes were grouped together to form the subpopulations. The plots show the average performance—measure by loss or accuracy—achieved over 100 simulations (for *Panel (a)*) and 10 simulations (for *Panel (b)*, respectively) using the shaded area to quantify the standard deviation. *Panel (a)* shows the result at the end of training changing relative representation ρ , while *Panel (b)* shows the training trajectories in a particular instance, see text for more details.

695 that the classifier will advantage the
696 group with highest average bright-
697 ness, see Eq. 11. In order to achieve
698 the same generalisation error on both
699 subpopulations, the less bright group
700 needs more samples (larger ρ). This
701 is shown in Fig. 8a, where the three panels correspond to different assignment of the classes: in the
702 top panel classes are randomly assigned to the two groups; in the middle panel classes are randomly
703 partitioned in two groups and the brighter one is assigned to group $-$; finally the last panel assigns
704 the brightest classes to group $-$ and least bright to group $+$. As predicted, we need increasingly high
705 relative representation ρ to achieve a balance in losses at the end of training.
706 When labels are balanced, our theory predicts that the classifier is initially attracted by the larger
707 Δ and eventually—if the relative representation of the group with smaller Δ is large enough—it
708 switches and favours the other group. This effect is indeed verified in the CIFAR10 experiments.
709 Starting from the partitioning in Fig. 8a (bottom) with $\rho = 0.8$, the dynamics is initially attracted by
710 group $-$ before advantaging the other group, giving rise to a crossing as shown in Fig. 8b.

711 G.2 CelebA

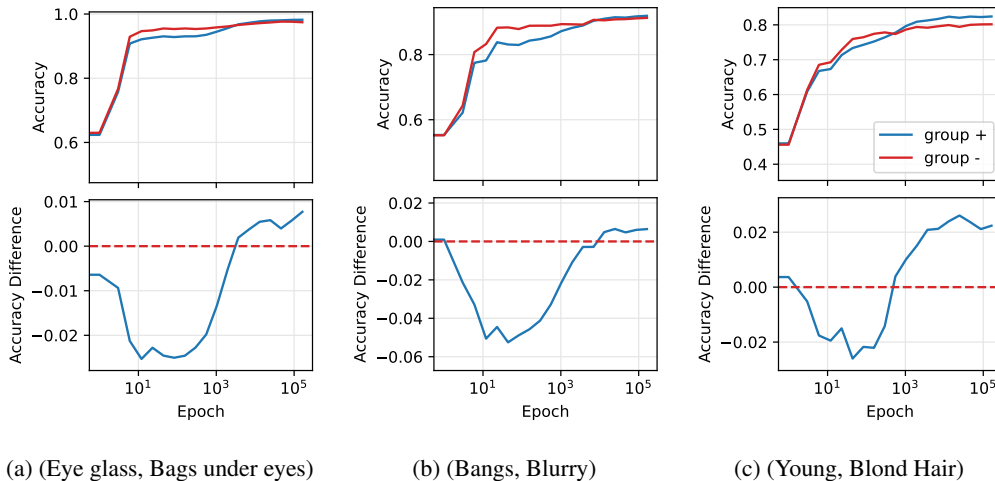


Figure 9: **Numerical simulations in the CelebA dataset.** Figure shows the average accuracy (solid lines) and standard deviation (shaded area) of 4 different runs in this framework. The top row depicts the test accuracy over the course of training for different pairs of target and group attributes. The bottom row illustrates the difference in test accuracies between the $+$ and $-$ subpopulations, highlighting the crossing phenomenon observed during training. *Panels (a), (b), and (c)* depict this for the pairs of target and group attributes of (Eye glass, Bags under eyes), (Bangs, Blurry), and (Young, Blond Hair), respectively.

712 The goal of this experiment is to show the emergence of different timescales in realist scenarios of
713 relevance for the fairness literature.

714 The CelebA dataset [25] contains over 200k celebrity images annotated with 40 attribute labels,
715 covering a wide range of facial attributes such as gender, age, and expressions. For this experiment,
716 we consider different pairs as the target and group attributes. The task is to predict the target attribute
717 while the group attribute defines the $+$ and $-$ subpopulations.

718 For the model, we select a pretrained ResNet-18 model on ImageNet and add an additional fully
719 connected layer, with only the latter being optimised during training. We use cross-entropy as the
720 loss objective and train via online SGD.

721 We randomly selected target-label pairs, making sure to avoid attributes that are pathologically
722 underrepresented in the dataset and would hinder the significance of the result. In the plots shown

723 in Fig. 9 we show some of the pairs that show a crossing phenomenon. Each panel in Fig. 9
724 show the accuracy and accuracy gap over the course of training. Notice how the classifier favours
725 sub-population – in the initial phase of training before changing preference.

726 This result shows that bias can change over the course of training even in standard setting. This
727 does not imply that it will always occur and indeed several of the pairs in the dataset do not show a
728 crossing phenomenon. However, understanding when and why this phenomenon occurs can affect
729 the algorithmic choices that we make in our ML pipeline.

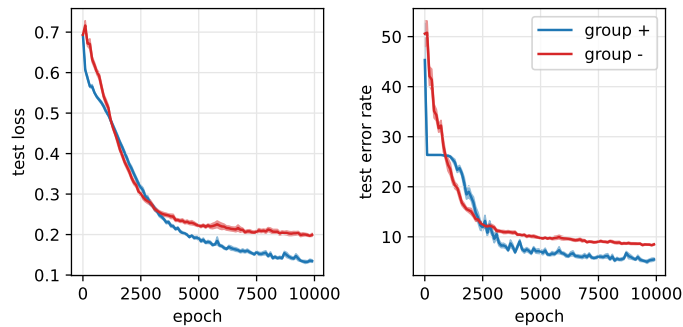


Figure 10: **Simulations on Synthetic Data and Deeper Networks** We observe the ‘double-crossing’ phenomena in not only the loss curves, but also the error curves for the positive sub-population (blue) and the negative sub-population (red). The shaded areas quantify the standard deviation obtained across 10 seeds. The data distribution parameters are $d = 100, v = 4, \rho = 0.75, \Delta_+ = 0.1, \Delta_- = 1, T_{\pm} = 0.9, \eta = 0.01, \alpha_+ = 0.473, \alpha_- = -0.200$

731 In this section we test the validity of the prediction of our model in more realistic settings. Specifically,
 732 assuming the same data distribution, we now train a multilayer perceptron (MLP) having one hidden
 733 layer of 200 units. We use ReLU activation and a sigmoid activation on the output. We train using
 734 online stochastic gradient descent and use binary cross entropy as our loss function. We sample
 735 training and test data from the data distribution and use the test data to obtain estimates of the loss as
 736 well as error rates (percentage of test examples misclassified).

737 For the general fairness case (sec. 4.2), we observe the three phase behaviour predicted by our
 738 model. The positive sub-population is initially advantaged more since it exhibits stronger spurious
 739 correlation. Then, the negative sub-population is advantaged since it has a higher variance. Finally,
 740 as per Eq. 11, the positive-sub-population is advantaged once more since it has sufficiently high
 741 representation. We not only observe the ‘double-crossing’ phenomena in the losses, but also in the test
 742 errors demonstrating the robustness of our model beyond the linearity and MSE loss assumptions.