1  We thank the reviewers for their reviews and comments. We will incorporate your suggestions in the final version.

2  We analyze an extremely simple model (as explicitly discussed and motivated, e.g. on lines 59-61 and the broader
3  impact section), but have reason to believe insights, behaviours and analysis methods will carry over to more complex
4  and realistic models. The paper Woodworth et al. 2020 [30] already shows how phenomena in this simple model
5  predicted by theory can be observed also for complex realistic networks. Insights about the extreme training accuracy
6  needed to reach the limiting behaviour are already helping us and others understand apparent discrepancies regarding
7  matrix factorization models. In addition, as we comment in lines 36-37, the effect of width is equivalent to that of
8  initialization scale (more precisely: increasing width corresponds to effectively scaling initialization as sqrt(width) in
9  matrix factorization / wide diagonal nets, see [30]). Therefore, our analysis is also relevant to more-overparameterized
10  models, with width$\to \infty$, and the number of weights entering whenever we discuss init scale.

11  **Normalization of the loss function (Reviewers #1 and #2):** The loss is defined to be the average over all examples.
12  Thus the interval of $\epsilon$ should be $(0,1]$, this is a typo that does not cause difficulty with analysis and we will fix it.

13  **For reviewer #1:** - There is a typo in the inequality in line 146, we will fix it. Also, in eq (2),(3),(4) it should be $\forall n$,
14  we will add it. All other minor typos will be fixed, thank you.
15  - Note that the algebraic form of the model appears in line 147.
16  - On the role of the number of observations: We analyze the linearly separable case. Overparameterization is a generic
17  way of achieving separability, but once the data is separable, the number of observations does not play a direct role.
18  - Comparison to assumptions/setting of Chizat and Bach 2020 [6] / Lyu and Li 2020 [19]: These papers analyze more
19  general models, as we discuss in Section 2. Our analytic results are only for the linear diagonal nets - which is a special
20  case of the models in those papers. For this homogeneous model all the assumptions of previous works hold.

21  **For reviewer #2:** - The discrepancy between gradient flow and discrete gradient descent is indeed very interesting. We
22  actually think that for large step-sizes one would exit the kernel regime much faster. Characterizing this is an important
23  challenge, but seems difficult (both we, and others, have tried, and are still trying to study this). In this paper we focus
24  on the training accuracy ($\epsilon$), but we certainly are interested also in the effects of other hyperparameters, including, and
25  perhaps most importantly, the step-size.
26  - Note that $Q_\mu$ is defined for general $\mu$, and $\mu = \alpha^D$ is only for the square loss setting studied in Woodworth et al. [30]
27  and takes different values in our Theorem 6. We will make it more clear.
28  - In line 285 it should be $\epsilon = \exp(-10000)$. We will fix that, thank you.

29  **For reviewer #3:** - Comparison with Ji and Telgarsky 2019 [16] and Lyu and Li 2020 [19]: References [16] and [19]
30  analyze the asymptotic max-margin solution only for *infinite* training accuracy (ie infinite training time or zero training
31  loss) and fixed initialization, i.e. the "rich regime", as discussed in the introduction, and in lines 125-140 in Section 2
32  (including Theorem 2, and footnote 3 where we discuss how [16] fits in). Section 2 is dedicated to discussing how this
33  is only one extreme endpoint (the other, corresponding to the "kernel regime", occurs when the training accuracy is
34  finite and the initialization grows to infinity). This is also visualized in Figure 1 and in the table above the figure. As
35  discussed explicitly in lines 55-58, the submission fills the gaps between the two extremes and studies the transition
36  between them. Moreover, [16]'s results are for fully connected linear networks where the rich regime and kernel regime
37  both lead to $L_2$ max-margin, while for the diagonal network, the rich regime limit ($L_1$ or $L_{2/D}$) is very different from
38  the kernel regime ($L_2$).
39  - Comparison with Woodworth et al. 2020 [30]: As discussed in the intro in lines 30-33 and in lines 119-125, [30]
40  analyzes the transition for *regression with the squared loss*, and does not study training accuracy as a controlling hyper-
41  parameter. The effect of the training accuracy is more important for classification problems since unlike regression,
42  with exponential loss, no finite parameter can ever achieve zero training loss. Thus, as discussed in lines 40-45, the
43  classification setting we study is substantially different, and in order to understand it we need to study the effect of
44  training accuracy, which was not studied in [30].

45  **For reviewer #4:** - In eq (4) we take $\alpha$ to infinity only since we want to emphasize the importance of the order of limits
46  when both $\alpha$ and time go to infinity (in contrast with eq. 2).
47  - Indeed, our theoretical results are asymptotic and obtained when *both* initialization and accuracy go to infinity at some
48  relative rates. The case of finite initialization and infinite accuracy is already known in Lyu and Li 2020 [19], Nacson et
49  al. 2019 [20], and so is the case of infinite initialization and finite accuracy Chizat et al. 2019 [7]. We fill in the gap by
50  letting both go to infinity while maintaining some finite relationship between them. Non-asymptotic results for finite
51  initialization and finite accuracy seem very difficult to obtain, but we see in figure 5(b) that the relationship predicted by
52  our asymptotic results is already seen also for moderate finite $\alpha$.
53  - The technical difficulty compared to [Woodworth et al. 2020] stems from the difference in KKT conditions. [Chizat et
54  al 2019] derived only the kernel regime (large initialization and finite time).