

A Effect of signal-to-noise ratio and nonlinearity

A.1 RF model

In the RF model, varying r can easily be achieved analytically and yields interesting results, as shown in Fig. 10⁶.

In the top panel, we see that the parameter-wise profile exhibits double descent for all degrees of linearity r and signal-to-noise ratio SNR, except in the linear case $r = 1$ which is monotonously decreasing. Increasing the degree of nonlinearity (decreasing r) and the noise (decreasing the SNR) simply makes the nonlinear peak stronger.

In the bottom panel, we see that the sample-wise profile is more complex. In the linear case $r = 1$, only the linear peak appears (except in the noiseless case). In the nonlinear case $r < 1$, the nonlinear peak appears is always visible; as for the linear peak, it is regularized away, except in the strong noise regime $\text{SNR} > 1$ when the degree of nonlinearity is small ($r > 0.8$), where we observe the triple descent.

Notice that both in the parameter-wise and sample-wise profiles, the test loss profiles change smoothly with r , except near $r = 1$ where the behavior abruptly changes, particularly at low SNR.

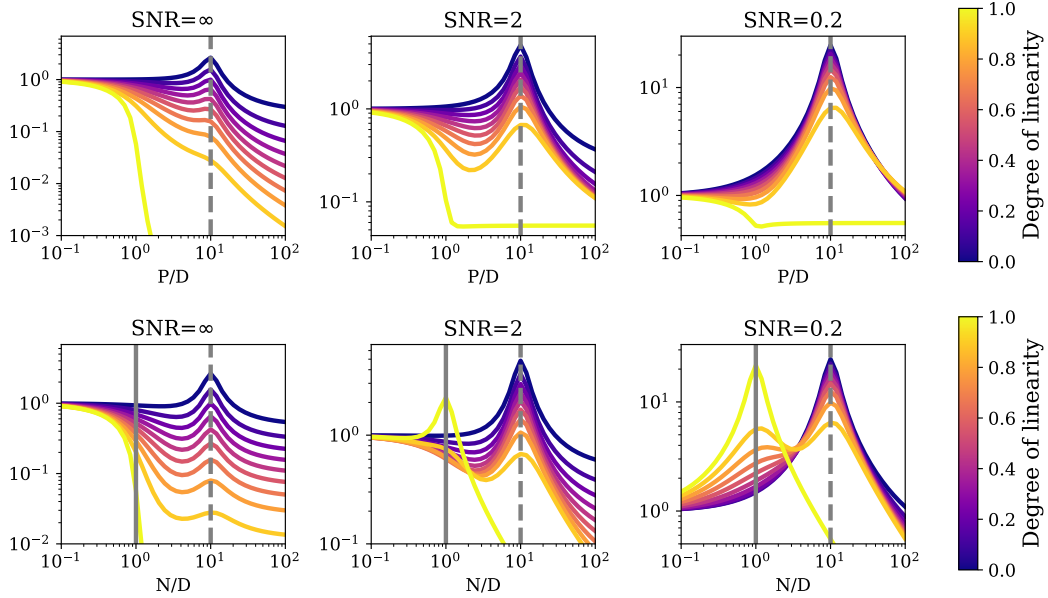


Figure 10: Analytical parameter-wise (**top**, $N/D = 10$) and sample-wise (**bottom**, $P/D = 10$) test loss profiles of the RF model. **Left**: noiseless case, $\text{SNR} = \infty$. **Center**: low noise, $\text{SNR} = 2$. **Right**: high noise, $\text{SNR} = 0.2$. We set $\gamma = 10^{-1}$.

One can also mimick these results numerically by considering, as in [30], the following family of piecewise linear functions:

$$\sigma_{\alpha}(x) = \frac{[x]_{+} + \alpha[-x]_{+} - \frac{1+\alpha}{\sqrt{2\pi}}}{\sqrt{\frac{1}{2}(1+\alpha^2) - \frac{1}{2\pi}(1+\alpha)^2}}, \quad (8)$$

for which

$$r_{\alpha} = \frac{(1-\alpha)^2}{2(1+\alpha^2) - \frac{2}{\pi}(1+\alpha)^2}. \quad (9)$$

⁶We focus here on the practically relevant setup $N/D \gg 1$. Note from the (P, N) phase-space that things can be more complex at $N/D \lesssim 1$.

Here, α parametrizes the ratio of the slope of the negative part to the positive part and allows to adjust the value of r continuously. $\alpha = -1$ ($r = 1$) will correspond to a (shifted) absolute value, $\alpha = 1$ ($r = 0$) will correspond to a linear function, $\alpha = 0$ will correspond to a (shifted) ReLU. In Fig. 11 we show the effect of sweeping α uniformly from 1 to -1 (which causes r to range from 0 to 1). As expected, we see the linear peak become stronger and the nonlinear peak become weaker.

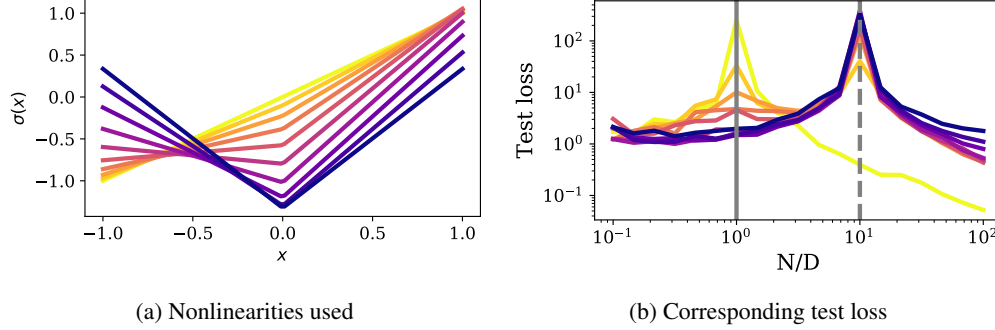


Figure 11: Moving from a purely nonlinear function to a purely linear function (dark to light colors) strengthens the linear peak and weakens the nonlinear peak.

363 A.2 NN model

We show in the top row of Fig. 12 the effect of varying the SNR on the (P, N) phase space for $\sigma = \text{Tanh}$ in the NN model. Just like in the RF model, triple descent only appears at $\text{SNR} < 1$ (right panel).

In the bottom row of the same figure, we show the effect of replacing Tanh ($r \sim 0.92$) by ReLU ($r = 0.5$). In the low SNR setup, we still distinguish the two peaks of triple descent, but the linear peak is much weaker, as expected from the stronger degree of nonlinearity.

Notice that in the intermediate signal-to-noise scenario, $1 < \text{SNR} < \infty$, results are different from the RF model where we only observed the nonlinear peak. For Tanh , we observe only the linear peak, whereas for ReLU , we observe something intermediate between the linear peak and the nonlinear peak.

374 B Origin of the linear peak

In this section, we follow the lines of [28], where the test loss is decomposed in the following way (Eq. D.6):

$$\mathcal{L}_g = \rho + Q - 2M \quad (10)$$

$$\rho = \frac{1}{D} \|\beta\|^2, \quad M = \frac{\sqrt{\zeta}}{D} \mathbf{b} \cdot \beta, \quad Q = \frac{\zeta}{D} \|\mathbf{b}\|^2 + \frac{\eta - \zeta}{P} \|\mathbf{a}\|^2, \quad \mathbf{b} = \Theta \mathbf{a} \quad (11)$$

As before, β denotes the linear teacher vector and Θ, \mathbf{a} respectively denote the (fixed) first and (learnt) second layer of the student. This insightful expression shows that the loss only depends on the norm of the second layer $\|\mathbf{a}\|$, the norm of the linearized network $\|\mathbf{b}\|$, and its overlap with the teacher $\mathbf{b} \cdot \beta$.

We plot these three terms in Fig. 13, focusing on the triple descent scenario $\text{SNR} < 1$. In the left panel, we see that the overlap of the student with the teacher is monotonically increasing, and reaches its maximal value at a certain point which increases from D to P as we decrease r from 1 to 0. In the central panel, we see that $\|\mathbf{a}\|$ peaks at $N = P$, causing the nonlinear peak as expected, but nothing special happens at $N = D$ (except for $r = 1$). However, in the right panel, we see that the norm of the linearized network peaks at $N = D$, where we know from the spectral analysis that the gap of the linear part of the spectrum is minimal. This is the origin of the linear peak.

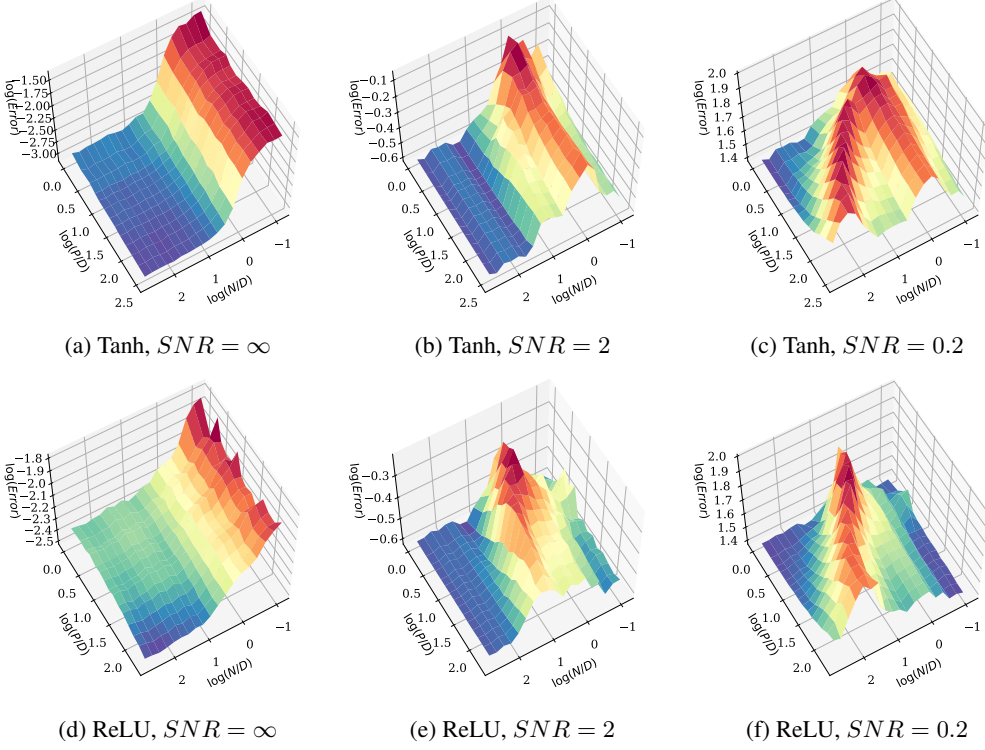


Figure 12: Logarithmic plot of the test loss in the phase space defined by number of parameters. **Left:** Single descent at low SNR. **Center:** Double descent at intermediate SNR. **Right:** Triple descent at low SNR.

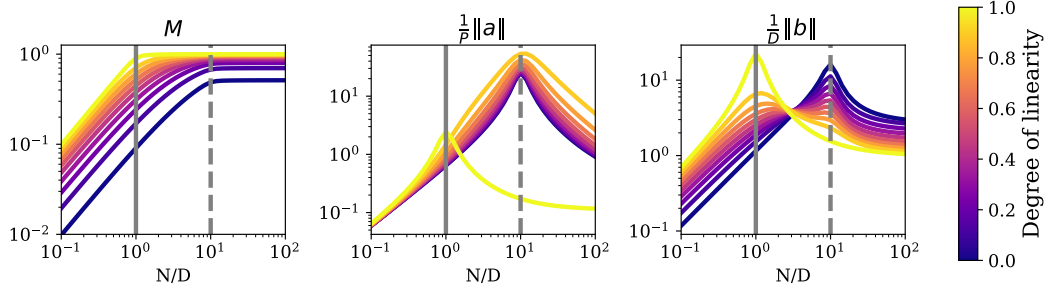


Figure 13: Terms entering Eq. [11], plotted at $SNR = 0.2$, $\gamma = 10^{-1}$.

388 C Structured datasets

389 In this section, we examine how our results are affected by considering the realistic case of correlated
 390 data. To do so, we replace the Gaussian i.i.d. data by MNIST data, downsampled to 10×10 images
 391 for the RF model ($D = 100$) and 14×14 images for the NN model ($D = 196$).

392 C.1 RF model

393 We refer to the results in Fig [14]. Interestingly, the triple descent profile is weakly affected by the
 394 correlated structure of this realistic dataset. However, the spectral properties of $\Sigma = \frac{1}{N} Z^\top Z$ are
 395 changed in an interesting manner: the two parts of the spectrum are now contiguous, there is no gap
 396 between the linear part and the nonlinear part.

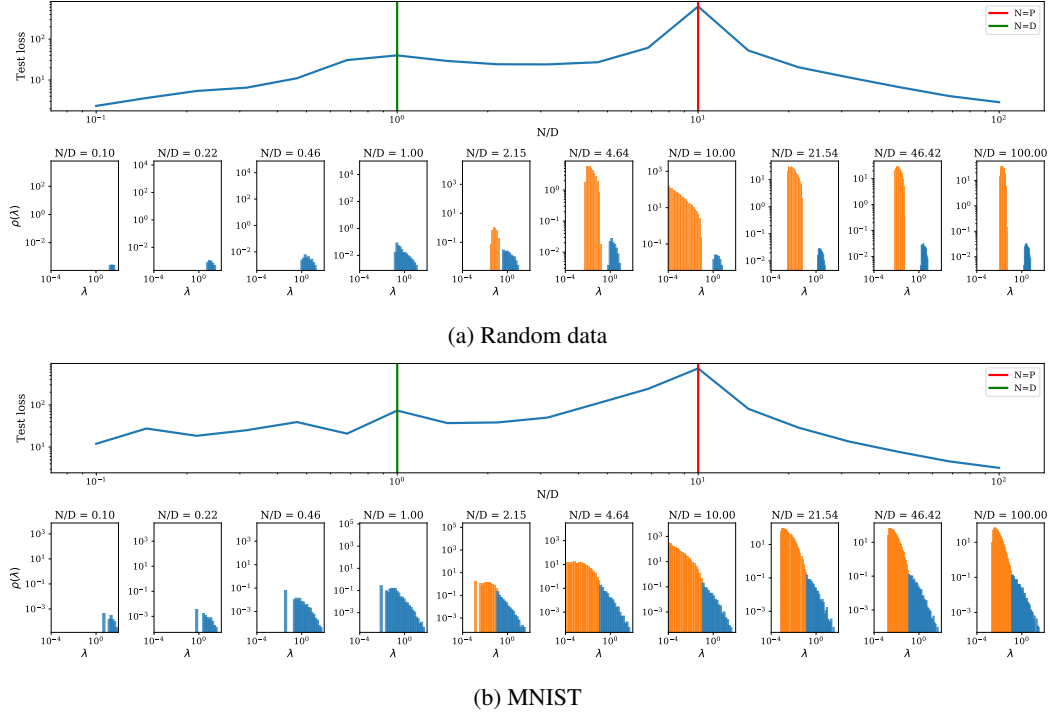


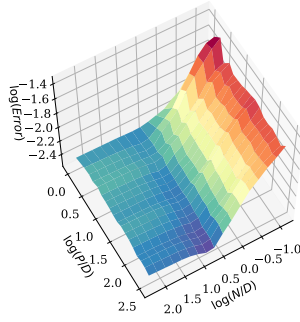
Figure 14: Spectrum of the covariance of the projected features $\Sigma = \frac{1}{N} \mathbf{Z}^\top \mathbf{Z}$ at various values of N/D , with the corresponding loss curve shown above. We set $\sigma = \text{Tanh}$, $\gamma = 10^{-5}$.

397 C.2 NN model

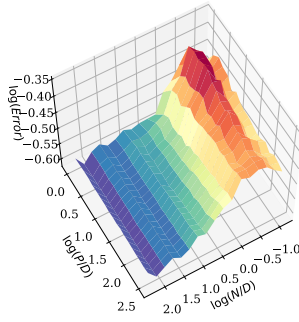
398 As shown in the top row of Fig. 15, the NN model is qualitatively different on the structured dataset:
 399 the two peaks at $N = D$ and $N = P$ are not well separated at $\text{SNR} < 1$ anymore. The single peak
 400 which appears is somewhat intermediate between the $N = D$ and $N = P$. However, by considering
 401 the time evolution in the bottom row of the same figure, we see that this peak shifts across the phase
 402 space during training, just like in the case of random data (Fig. 9).

403 At early times, it is located along a line of constant N , which makes it akin to a linear peak. At late
 404 times, it is rather reminiscent of a nonlinear peak, though it does not seem to be located at $P \sim N$ as
 405 before, but rather at $N \sim P^\alpha$ with $\alpha < 1$. This sublinear scaling is a consequence of the fact that
 406 structured data is easier to memorize than random data [17], and may blur the distinction between the
 407 two peaks.

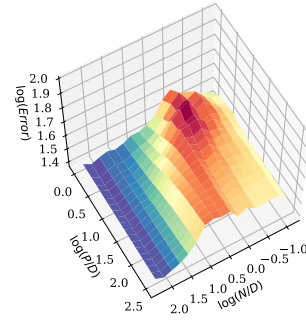
408 Interestingly, at early times, the peak does not occur at $N = D$ as expected, but rather at $N = D_{\text{eff}} \sim$
 409 $D/10 \sim 20$. We hypothesize that D_{eff} may be related to the intrinsic dimension of the input data
 410 [39, 40, 41]. Although the linear peak still occurs at $N = D$ for MNIST data in the RF model, in the
 411 NN setup feature learning occurs. When the dataset is highly correlated like MNIST, feature learning
 412 compresses the dataset down to a more compact representation, likely causing the $N = D$ peak to
 413 shift to lower values. A study of this crucial question is deferred to future work.



(a) MNIST, $SNR = \infty$



(b) MNIST, $SNR = 2$

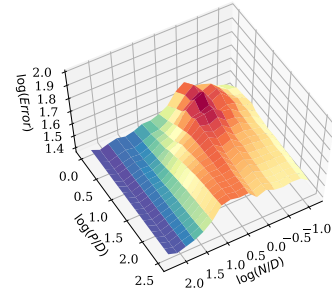
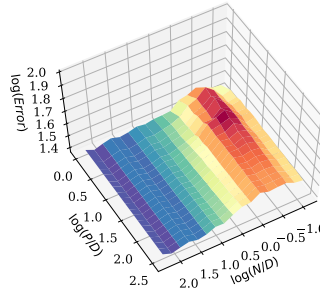
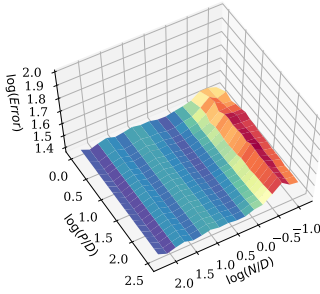


(c) MNIST, $SNR = 0.2$

t=37 epochs

t=162 epochs

t=695 epochs



(d) Dynamics on MNIST at $SNR = 0.2$

Figure 15: Test loss phase space on MNIST with $\sigma = \text{ReLU}$. **Top:** After 1000 epochs, for various values of the SNR. **Bottom:** at three different times during training in the low SNR case.