MULTIPLE DESCENTS IN UNSUPERVISED AUTO-ENCODERS: THE ROLE OF NOISE, DOMAIN SHIFT AND ANOMALIES

Anonymous authors

006

007

008 009 010

011

013

014

015

016

017

018

019

021

023

025

026

027

Paper under double-blind review

Abstract

The phenomenon of *double descent* has recently gained attention in supervised learning. It challenges the conventional wisdom of the bias-variance trade-off by showcasing a surprising behavior. As the complexity of the model increases, the test error initially decreases until reaching a certain point where the model starts to overfit the train set, causing the test error to rise. However, deviating from classical theory, the error exhibits another decline when exceeding a certain degree of over-parameterization. We study the presence of double descent in unsupervised learning, an area that has received little attention and is not yet fully understood. We conduct extensive experiments using under-complete auto-encoders (AEs) for various applications, such as dealing with noisy data, domain shifts, and anomalies. We use synthetic and real data and identify model-wise, epoch-wise, and sample-wise double descent for all the aforementioned applications. Finally, we assessed the usability of the AEs for detecting anomalies and mitigating the domain shift between datasets. Our findings indicate that over-parameterized models can improve performance not only in terms of reconstruction, but also in enhancing capabilities for the downstream task.

028 029 1 INTRODUCTION

In recent years, there has been a surge in the use of extremely large models for both supervised 031 and unsupervised tasks. This trend is driven by a desire to solve challenging machine-learning 032 tasks. However, this pursuit contradicts the well-known bias-variance trade-off, which suggests 033 that larger models tend to overfit the training data and perform poorly on the test set (Hastie et al., 034 2009). Despite this, many over-parameterized models have been able to generalize well (Krizhevsky 035 et al., 2012; He et al., 2016). This challenges common assumptions regarding the generalization capabilities of models (Zhang et al., 2021; Advani et al., 2020; Neyshabur et al., 2018), as over-037 parameterized models often exhibit significantly superior performance compared to smaller models, 038 even when interpolating the training data (Belkin et al., 2019b; 2018).

Recently, Belkin et al. (2019a) conducted a study on the bias-variance trade-off for large, complex deep neural network models. They discovered an interesting phenomenon called double descent. Initially, as the complexity of the model increases, the test error decreases. Specifically, as the complexity continues to increase, the variance term starts to dominate the test loss, resulting in an increase, which is known as the classical bias-variance trade-off. However, at a certain point, termed the "interpolation threshold" (Nakkiran et al., 2021), the test loss stops increasing and begins to decline again in the over-parameterized regime, yielding a curve with two decent regimes.

The phenomenon of double descent has been observed in many frameworks in supervised learning (see a survey in (Dar et al., 2021)). Model-wise double descent was demonstrated in (Spigler et al., 2018), while (Nakkiran et al., 2021; Gamba et al., 2022; Hastie et al., 2022) explore the impact of label noise and Signal-to-Noise ratio (SNR) on the double descent curve respectively. Nakkiran et al. (2021); Gamba et al. (2022) also demonstrated the phenomenon to epoch-wise and sample-wise double descent. Multiple descents were discussed in (Adlam & Pennington, 2020; Liang et al., 2020; Chen et al., 2021), and d'Ascoli et al. (2020) revealed that the interpolation threshold depends on the linearity and non-linearity of the model. However, the existence of double descent in core tasks in unsupervised learning is not yet fully understood.

Bas-Variance Modern Regime 1.0 0.8 0.6 0.4 0.2 0.100 200 300 400 500 Hidden Layer Size

Figure 1: Demonstration of the double descent phenomenon in unsupervised learning. We present the test loss for varying epochs and hidden layer sizes for an under-complete AE.

062 In this study, we analyze the double descent phenomenon and its implications for crucial unsuper-063 vised tasks such as domain adaptation, anomaly detection, and robustness to noisy data, utilizing 064 under-complete AEs. AEs are general architectures that have been used in unsupervised learning for 065 numerous tasks, including denoising (Vincent et al., 2008; 2010), manifold learning (Duque et al., 066 2020; Wang et al., 2014), clustering (Song et al., 2013; Yang et al., 2019), anomaly detection (Zhou 067 & Paffenroth, 2017; Sakurada & Yairi, 2014), feature selection (Han et al., 2018; Gong et al., 2022), 068 domain adaptation (Deng et al., 2014; Yang et al., 2021), segmentation (Baur et al., 2021; Myronenko, 2019), and generative models (Kingma, 2013; Doersch, 2016), making them a prominent use 069 case when studying double descent in the field of unsupervised learning. We utilized under-complete AEs to ensure the model learns meaningful representations in the latent space, as over-complete 071 AEs might simply learn the identity function and fail to capture the underlying data structure. 072

073 We present extensive empirical evidence showing that double and triple descent phenomena occur 074 in unsupervised AEs when the data is contaminated with noise. Our findings reveal that "memoriza-075 tion" plays a critical role, as models can overfit the noise rather than capture the underlying signal, leading to poor test performance. However, we also discovered that sufficiently large models can still 076 achieve superior test performance on clean data, even when fitting noisy training samples. This sug-077 gests that these models successfully extract the true signal despite the presence of noise. Through experiments on both synthetic and real-world data, we show that double and triple descent curves man-079 ifest in under-complete AEs exposed to various types of data contamination. Specifically, we show that different levels of sample noise, feature noise, domain shift, and the proportion of outliers all 081 significantly influence the shape of the double descent curve. We observe model-wise, sample-wise, 082 and epoch-wise double descent patterns across these settings. For instance, in Figure 1, we illustrate 083 a double descent curve obtained by training an AE on data generated using the "sample noise" model 084 described in Subsection 3.1. We follow (Nakkiran et al., 2021) in categorizing models as under-085 parameterized or over-parameterized, referring to those on the left or right of the critical regime, respectively. In these regimes, increasing the model size leads to a reduction in test loss. The critical regime is characterized by models for which changes in size can either decrease or increase test loss. 087

Our results have important implications for real-world applications, particularly in unsupervised learning tasks. We demonstrate that over-parameterized models trained on source domain data can adapt more effectively to target domains, even under distributional shifts. Additionally, we uncover non-monotonic behavior in anomaly detection performance as model complexity increases, further underscoring the practical relevance of our findings when noise and domain shifts are prevalent.

093 094

054

060

061

2 RELATED WORK

095

096 The discovery of double descent for neural networks (NNs) has led to extensive research aimed at understanding the behavior of generalization errors. It has also provided insights into why larger 098 models perform better than smaller or intermediate ones. Most studies have been conducted in a supervised learning setting, as detailed in (Belkin et al., 2019a; Nakkiran et al., 2021; Dar et al., 099 2021; Spigler et al., 2018; Gamba et al., 2022; Adlam & Pennington, 2020; Liang et al., 2020; Chen 100 et al., 2021; Xia et al., 2022; Kausik et al., 2023). Recent studies (Nakkiran et al., 2021; Hastie 101 et al., 2022; Bartlett et al., 2020; Li et al., 2020) have introduced label noise, feature noise, and 102 different levels of SNRs and demonstrated that large over-parameterized NNs can "memorize" the 103 noise while still generalizing better than smaller models. 104

The phenomenon of double descent has not been extensively studied in the context of unsupervised learning, and there are some contradictions in the literature regarding its presence. Principal Component Analysis (PCA) (Shlens, 2014b) and Principal Component Regression (PCR) (Massy, 1965), which are special types of linear AEs, are widely used unsupervised and supervised learning

models respectively and can serve as an interesting case study for exploring double descent. Gedon et al. (2022) argued that there is no double descent in PCA while (Xu & Hsu, 2019; Teresa et al., 2022) showed evidence for double descent in PCR. Lupidi et al. (2023) used a specific subspace data model and argued that there is no sign of model-wise double descent in both linear and non-linear AEs. Sonthalia & Nadakuditi (2023) and Dubova (2022) demonstrated sample-wise double descent for denoising AEs with different SNR values. Zhang et al. (2023) used a self-supervised learning framework for signal processing and found epoch-wise double descent for different levels of noise.

115 Our analysis of double descent differs from previously published studies in three significant ways. 116 Firstly, when trained on noisy data, we demonstrate that standard under-complete AEs experience 117 double and even triple descent at the model-wise, sample-wise, and epoch-wise levels. We have also 118 partitioned the model's size into bottleneck and other hidden layer dimensions to understand the phenomenon better. Secondly, we show that the noise magnitude and the number of noisy samples affect 119 the double descent curve. Thirdly, we show that double descent also occurs in common realistic con-120 tamination settings in unsupervised learning, such as source-to-target domain shift, anomalous data, 121 and additive sample and feature noise. Finally, we demonstrate the implications of multiple descents 122 in unsupervised learning tasks using real-world data, extending beyond reconstruction. 123

124 125

126 127

128 129

130

146 147

3 DATA MODEL

This section outlines the data and contamination models used to study double descent.

3.1 LINEAR SUBSPACE DATA

We utilized the synthetic dataset of Lupidi et al. (2023) to challenge their assertion that "double descent does not occur in self-supervised settings". First, we sample N random i.i.d. Gaussian vectors, each of size d, representing random features in a latent space, $z_i \sim \mathcal{N}(0, I_d)$. Next, we embed the vectors $\{z_i\}_{i=1}^N$ into a higher dimensional space of size n by multiplying each z_i by D of size $n \times d$, Dz_i , where $D_{ij} \sim \mathcal{N}(0, 1)$. This setting can be thought of as measuring $\{z_i\}_{i=1}^N$ with a measurement tool D, resulting in higher-dimensional data. Our dataset differs from (Lupidi et al., 2023) in several ways, and we will investigate four scenarios as part of our study:

Sample Noise. We aim to investigate the impact of the number of noisy training samples on the test 138 loss curve. In contrast to Lupidi et al. (2023), which adds noise to all samples, we vary the number 139 of noisy training samples to identify memorization. We do this by introducing a new variable, p, 140 representing the probability of a sample being noisy. Thus, $p \cdot 100\%$ represents the percentage of 141 noisy samples in the data. As noise is added, we control the SNR by defining the parameter θ . 142 Another significant change from Lupidi et al. (2023) is the chosen values of SNR, which can be 143 found in Appendix A, table 2, along with its calculation to derive θ , in Appendix B. This leads to 144 the following equation, which describes our model for the sample noise scenario: 145

$$x_i = \begin{cases} \theta D z_i + \epsilon_i, & \text{with probability } p, \\ \theta D z_i, & \text{with probability } 1 - p \end{cases}$$

where $\epsilon_i \sim \mathcal{N}(0, I_n)$ is an additive white Gaussian noise (AWGN), representing the noise added to samples with probability p. This setting can be likened to using a noisy measurement device. To illustrate this generation we present in Appendix A, Figure 14 a visualization of the data model.

Feature Noise. We further study the impact of the number of noisy training features on the test loss curve. In this scenario, each sample $\{\theta Dz_i\}_{i=1}^N$ is affected by noise in certain features. We denote the probability of a feature being noisy by p, controlling each sample's noisy features. We simulate a scenario where we have n measuring tools, each measuring a different feature. To introduce noise, we select the same set of features to be noisy across all samples. This mimics a situation where $[n \cdot p]$ of the measuring tools are unreliable or noisy. The SNR calculation is explained in Appendix B and Appendix A, Figure 14 depicts the data generation for this setting.

Domain Shift. We aim to explore how the test loss curve behaves when there is a domain shift between the train and test datasets. First, we partition the vectors in the latent space $\{z_i\}_{i=1}^N$ to train and test vectors, denoted as z_{train}^i and z_{test}^i respectively. Then, the train vectors are projected to a higher dimensional space with the matrix D, and the test vectors are projected with a different matrix D'', modeling a domain shift. To control the shift, we define $D'' = D + s \cdot D'$, where D is the matrix multiplying the train vectors and $D'_{ij} \sim \mathcal{N}(0,1)$ is a new random matrix added to cause perturbations at each entry of D and the parameter s > 0 controls the shift between D and D''.

$$x_i = \begin{cases} Dz_{train}^i, & \text{if } train, \\ D'' z_{test}^i, & \text{if } test. \end{cases}$$

Since D and D' are i.i.d., D'' follows a normal distribution $\mathcal{N}(0, (1 + s^2)I)$. To obtain the same norm in the test data, we divide D'' by $\sqrt{1 + s^2}$. This scenario is similar to the case where two different measuring instruments (i.e., D, D'') are measuring the same phenomenon. This data model is illustrated in Appendix A, Figure 15, and the definition of the SNR is detailed in Appendix B.

Anomalies. We conduct an experiment to investigate the impact of anomalies in the training set on the test loss curve. To represent clean samples, we utilize $\{\theta D z_i\}_{i=1}^N$. For generating anomalies, we sample from a normal distribution $\mathcal{N}(0, I_n)$. We introduce a metric termed Signal-to-Anomaly ratio (SAR), which regulates the magnitude ratio between the clean and anomaly samples through the parameter θ . Subsequently, we substitute $p \cdot 100\%$ of the normal samples with anomalies. This generation is illustrated in Appendix A, Figure 16.

178

180

165 166 167

179 3.2 SINGLE-CELL RNA DATA

We utilized single-cell RNA sequencing data from (Tran et al., 2020) to illustrate our findings using
real-world data. The data exhibits diverse domain shifts across different laboratory environments
and measurement technologies. This dataset is crucial for assessing the impact of domain shifts on
the test loss curve. Since this data is from a real-world setting, we are unable to control the shifts
between the training (source) and testing (target) datasets, as explained in Subsection 3.1. We also
use this dataset to show double descent when noise is injected to the samples and features (i.e.,
sample and feature noise) manually. We refer to Appendix A for more details.

188 3.3 CELEBA DATA

We incorporate real-world data to investigate anomaly detection across various model sizes. Specifically, we leverage the CelebA attributes dataset used in (Han et al., 2022), comprising over 200K samples and 4,547 anomalies, each characterized by 40 binary attributes. We sub-sample 3000 clean samples and replace $p \cdot 100\%$ of them with anomalies to create the training set.

194 195

196

4 Results

197 Our results presented in the main text are primarily based on a multi-layer perceptron (MLP) undercomplete AEs. However, we have also conducted additional evaluations using convolutional NNs 199 (CNNs). These findings are presented in Appendix E.2. Each of the reported results is based on 5 to 200 15 random seeds. Complete implementation details and discussion on the high computational load that each figure requires can be found in Appendix A. All models are trained using contaminated 201 datasets and tested on clean data. Consequently, the test loss serves as an indicator of whether the 202 model has memorized the noise (high test loss) or learned the signal (low test loss). Over-complete 203 AEs are beyond the scope of this discussion, as they can learn the identity function, leading to trivial 204 and uninformative data learning. Our emphasis is on standard (unsupervised) AEs, where in the 205 training process, we minimize the mean squared error (MSE) between the input and the model's 206 output. Train loss figures corresponding to all test losses depicted in this section are provided in 207 Appendix C and more results with a non-linear synthetic dataset are presented in Appendix E.3.

208 209

210

4.1 MODEL-WISE DOUBLE DESCENT

This section analyzes the test loss with increasing model sizes. For AEs, we break down the well-known "double descent" phenomenon into two interconnected variations: "hidden-wise" and "bottleneck-wise" and show how both contribute to the double descent behavior in the test loss. We also study the influence of several contaminations described in Section 3 and conclude that the interpolation threshold location and value can be manipulated by these factors. We also found that double descent typically occurs with high levels of sample noise and low SNR values. In these settings, the



Figure 2: Test and train losses as a function of model size. The test loss demonstrates clear double descent when varying the bottleneck or hidden layer size. The AEs were trained on the linear subspace model with sample noise = 90% and SNR = -15 [dB] (see details in Subsection 3.1).

noise predominates the training set, leading models in the critical regime to focus on interpolating the noise rather than learning the underlying signal to reduce the training loss. This, in turn leads to higher test loss. This explains why Lupidi et al. (2023) did not observe the phenomenon, as they used only high SNR levels (10 dB) in the model-wise sample noise scenario.

In Figure 2, we provide visual evidence of the bottleneck-wise and hidden-wise double descent. This not only helps to distinguish between various model sizes but also underscores the significance of our different architectural choices. The training loss consistently decreases as the dimensions of the model increase. In contrast, both the bottleneck and hidden layers exhibit the characteristic double descent curve, as seen in the decrease in test loss, followed by an increase and then another decline. This demonstrates that AEs trained on contaminated data can exhibit double descent.

Sample and feature noise. Interestingly, Figure 3a shows that the height of the test loss increases and the interpolation threshold location shifts towards larger models as the level of sample noise increases. This can be clarified by the observation that increased noise adversely affects model learning. Moreover, we need a bigger model to overfit the noisy samples. The absence of double descent for 0-20% sample noise can be attributed to the insufficient number of noisy samples in the training data. In Figure 3b, we demonstrate triple descent using single-cell RNA data, where we notice a similar behavior for the test loss, specifically for each of the two peaks. Evidence for double descent using the feature noise data model introduced in Subsection 3.1 is presented in Appendix D.

In addition, we observed the phenomenon of final ascent, characterized by double descent following
a final increase in the test loss, initially discussed in the case of supervised learning (Xue et al., 2022). We present the results of final ascent for unsupervised AEs with the single-cell RNA data
for 0-20% sample noise in Appendix E.4, Figure 55b. We also show how double and triple descent patterns emerge under different types of noise, such as Laplacian noise, and find double descent for sparse AEs, as detailed in Appendix E.5.

SNR. We observed that the SNR plays a crucial role in the test loss, which in turn affects its height. A higher SNR value reduces the impact of noise, allowing the model to learn the underlying signal from the training set, resulting in a lower test loss. Conversely, a lower SNR value amplifies the influence of noise, disrupting the model's ability to learn the signal leading to inferior results in the test loss. In Figure 4a, for SNR = 0 [dB], the double descent curve is absent because it prevents models in the critical regime from memorizing the noise, as the noise is not sufficiently dominant. Figure 4 and Appendix D, Figure 28 present results for the sample and feature noise settings respectively.



(a) Linear subspace data. SNR = -15 [dB].

(b) **Single-cell RNA data**. SNR = -17 [dB].





316 317

318

319

320

321

322

Figure 4: The effect of SNR for the case of **noisy samples** on the test loss curve. Train losses are illustrated in Appendix C, Figure 20.

281 **Domain shift.** We study the existence of double descent when the distribution of the training 282 (source) data, differs from that of the testing (target) data. We investigate the impact of the model 283 size on learning shared representations for both source and target datasets and reducing the shift 284 between them. By training the model on the source data and testing it on different targets, we un-285 veil non-monotonic behavior for the linear subspace dataset, shown in Figure 5. Additionally, the test loss rises as the shift is more dominant, and for lower levels of domain shift (shift = 0.1, 0.5), 287 non-monotonic behavior does not occur since the source and target domains are closely aligned. In 288 these cases, even models that interpolate the source data and learn domain-specific representations 289 perform well on the target data, preventing an increase in the test loss. Furthermore, we identify that over-parameterized models result in lower test loss, leading to improved target data reconstruction. 290 Subsection 5.1, Figure 12 presents double and triple descent results for the single-cell RNA data and 291 further insights about the connection of model size and domain adaptation utilizing real-world data. 292

293 **Anomaly detection.** We also identify double descent occurring when anomalies, deviating from the expected behavior of the data are introduced into the training set. Following the unsupervised setting, we consider the scenario where there is no anomaly-free dataset available for training, 295 making it more challenging to differentiate between normal and anomalous data (Cheng et al., 296 2021). In particular, we use the anomaly dataset mentioned in Subsection 3.1 with high number 297 of anomalies in the train set, akin to (Lindenbaum et al., 2024; Lerman & Maunu, 2018b;a), which 298 included anomalies with up to 99.5% for the subspace recovery setting and (Han et al., 2022), which 299 includes up to 40% anomalies in the case of unsupervised anomaly detection. We study the test loss 300 curves by varying the amounts of anomalous training samples. We used a common method where 301 data points with reconstruction loss surpassing a defined threshold are identified as anomalies as 302 discussed in (Lindenbaum et al., 2024; Malhotra et al., 2016; Borghesi et al., 2019). 303

We evaluate the anomaly detection capabilities using the receiver operating characteristic area 304 under the curve (ROC-AUC) metric. This metric employs the reconstruction error to measure 305 the model's ability to distinguish between clean and anomalous data (anomalies are identified as 306 data points with errors crossing a defined threshold). A higher ROC-AUC value signifies superior 307 performance. The models in the critical regime depicted in Figure 6a result in higher test loss of the 308 clean samples, complicating the differentiation between clean and anomaly data, leading to worse 309 ROC-AUC results. Scaling up the model size results in a secondary descent in the test loss of the 310 clean data. This double descent curve is particularly evident under conditions of low SAR and a 311 high number of anomalies in the training set, similar to the results of the sample noise scenario (the anomalies play the role of the noise). This secondary descent facilitates the model's ability 312 to differentiate between clean and anomalous data, resulting in performance comparable to that 313 of the under-parameterized models in terms of ROC-AUC, while learning meaningful embedding 314 for both clean data and outliers, resulting in lower test losses. Figure 6b demonstrates the absence 315



Figure 5: Linear subspace data exhibits model-wise nonmonotonic behavior for varying domain shifts. Appendix C, Figure 21 shows the train loss behavior.



(b) Synthetic anomaly data with SAR = 0 [dB].

338 Figure 6: Left: test loss of the clean samples. A double descent pattern emerges for low SARs and 339 high anomaly presence in the training data. Middle: test loss of the anomaly data. Right: Non-340 monotonic behavior of the ROC-AUC.

341 of double descent due to high SAR. However, similar to Figure 6a, intermediate models exhibit 342 poorer ROC-AUC performance compared to under and over-parameterized models. We present 343 more insights on anomaly detection utilizing real-world data in Subsection 5.2.

344 In conclusion, as contamination setups become more severe, such as higher noise levels, significant 345 domain shifts, many anomalies, or low SNR, the double descent phenomenon becomes more 346 pronounced, and the test loss increases. In some instances, these noise levels also cause the critical 347 regime to shift to the right. In Appendix A, Table 1, we compare the existence of double descent 348 between unsupervised and supervised learning for varying contamination setups and conclude that 349 they result in similar behaviors.

4.2 EPOCH-WISE DOUBLE DESCENT

350 351

352

363

364 365

366

367 368 369

370

371

372

373

374

375

353 In this section, we explore the presence of double descent versus the number of epochs. This study represents the first unsupervised investigation of this kind, expanding on similar research carried out 354 by (Nakkiran et al., 2021) for supervised learning. Figure 7 and Appendix D, Figure 29 show the 355 impact of the number of noisy samples and features in the train set on the test loss curve respectively. 356 Increasing noise makes it harder for the model to learn the signal, leading to a higher test loss. A 357 similar effect is obtained when varying the SNR, where as it decreases, the noise becomes more 358 dominant, resulting in an increase in the test loss. This is illustrated in Figure 8 and in Appendix D, 359 Figure 30 for the case of sample and feature noise respectively. Epoch-wise double descent is also 360 present when there is a domain shift between the train and test sets, as illustrated in Figure 9. 9a 361 shows that the stronger the shift, the higher the test loss. 362

4.3 SAMPLE-WISE DOUBLE DESCENT

In this section, we study the impact of the number of training samples on the test loss curve. The complexity of a model and the number of samples it is trained on both play a crucial role in determining whether the model is over (small sample size) or under-parameterized (large sample size).







Figure 7: Epoch-wise double descent influenced by the number of noisy samples. Train losses are depicted in Appendix C, Figure 22.



Figure 8: Epoch-wise double descent for the case of sample noise influenced by the SNR. Train losses are exhibited in Appendix C, Figure 23.

This causes the interpolation threshold's location to change, as shown in Figure 10. This adjustment 390 can sometimes result in a model that performs worse than a model trained on a smaller set of training 391 samples. A similar phenomenon was demonstrated in (Nakkiran et al., 2021) in a supervised setting. 392

393 We also investigate the impact of gradually increasing the number of training samples on the test loss 394 curve while keeping the model's size fixed. Remarkably, we identify a non-monotonic trend in the test loss curve at Figures 11b, 11c, and Appendix D, Figures 31a and 31b, which sometimes results 395 in double descent as noticed in 11a. The emergence of non-monotonic behavior is defined by a phase 396 where an increased number of samples negatively impacts performance, resulting in higher test loss. 397 Appendix C, Figure 26 depicts the training loss plotted against the number of samples in the scenario 398 of sample noise. Figure 11 showcases only the results from the linear subspace dataset due to the 399 insufficient amount of samples in the single-cell RNA dataset. More results regarding sample-wise 400 double descent can be found in Appendix E.2, E.3. The impact of the noise level, SNR, and the 401 domain shift on the test loss is consistent with the analyses conducted in Subsections 4.1 and 4.2. 402

403 404

405 406

407

408

387

388 389

5 **REAL WORLD APPLICATIONS**

In this section, we demonstrate how our findings can be applied to important tasks in machine learning, such as domain adaptation and anomaly detection. Our objective is to emphasize the significance of model size selection rather than to compete with state-of-the-art techniques.

409 410 411

412

413

414

415

416

417

418

419

420

5.1 DOMAIN ADAPTATION

Many frameworks in machine learning are exposed to domain shifts. The difference in distribution between the training and testing data can lead to inferior results when the model is employed on new, unseen data. Numerous domain adaptation methods have been proposed for both supervised and unsupervised settings (Zhou et al., 2022; Peng et al., 2019; Chang et al., 2019; Rozner et al., 2023; Yampolsky et al., 2023) to minimize the shift between the source and target domains. This is an ongoing challenge in biology, where researchers attempt to integrate datasets collected under different environmental conditions that cause distribution shifts. Many studies have been conducted to develop strategies to mitigate this shift, known in biology as "batch effect" (Tran et al., 2020).



428 429 430

Figure 9: Epoch-wise double descent influenced by the amount of **domain shift**. Left: we introduce some noise to emphasize the double descent curve. Appendix C, Figure 24 shows the train losses.



Figure 10: Model-wise double descent for the linear subspace data with different number of training samples. In the yellow interval, models trained with 10000 samples perform worse compared to those trained with 5000 samples. Sample noise = 70% and SNR = -15 [dB]. Train loss results are shown in Appendix C, Figure 25.

440 In this section, we study the relation between model size and its ability to alleviate distribution 441 shifts in real world single-cell RNA data 3.2. Our work is the first to show the advantage of over-442 parameterized models in unsupervised tasks under real domain shifts through the emergence of 443 double descent. Tripuraneni et al. (2021) and Kausik et al. (2023) focused on supervised learning 444 and showed related results, each under their own set of assumptions. We visualize the source and target datasets using UMAP embeddings (McInnes et al., 2018) in Appendix E.1, Figure 32b. The 445 top two sub-figures in Figure 12 present the test and train losses respectively for models trained on 446 source and tested on target datasets. We observed that testing models on the 'Wang' dataset results 447 in triple descent, while all other targets result in double descent. 448

449 We used the KL-divergence (KLD) (Shlens, 2014a) metric to calculate the distribution shift between 450 our source data ('Baron') and the target datasets ('Segerstolpe', 'Xin', 'Mutaro', and 'Wang'). Our findings show that as the shift between the source and target increases (higher KLD), the test loss 451 rises, inline with the results of the simulated experiment yielding Figure 5. To evaluate how different 452 models perform in terms of domain adaptation, we measure how much of the shift was removed by 453 analyzing the bottleneck representations of the AEs. Precisely, we compute the k = 10 nearest 454 neighbors of each bottleneck vector and determine the proportion belonging to the same biological 455 batch as mentioned in (Schilling, 1986), Section 3. We call this metric "k-nearest neighbors domain 456 adaptation test" (KNN-DAT), indicating the extent of mixing between different domains. KNN-DAT 457 of 1 implies complete separation, while a lower value indicates better mixing of different domains. 458 That is, lower values of KNN-DAT imply that the embedding of samples from the target domain is 459 more similar to the embedding of samples from the source domain.

460 The bottom row in Figure 12 presents the UMAP representations based on the embeddings of the 461 learned AEs. For under-parameterized models, KNN-DAT results are better compared to the mod-462 els in the critical regime. However, they achieve lower KNN-DAT at the expense of learning the 463 source data inadequately (high train loss). The model in the critical regime learns domain-specific 464 features, resulting in high KNN-DAT. We find that over-parameterized models yield the best KNN-465 DAT results, achieving a score of 0.75. They also lead to reduced test loss, resulting in improved 466 reconstruction of the target data. This suggests that over-parameterized models facilitate the tran-467 sition between source and target datasets, serving as a viable domain adaptation strategy. We also display results for the linear subspace dataset in Appendix E.1. 468

5.2 ANOMALY DETECTION

Unsupervised anomaly detection is a crucial task in machine learning. It has various applications across scientific fields, and many studies have utilized AEs for anomaly detection (Chandola et al., 2009; Lindenbaum et al., 2024; Chen et al., 2018; Rozner et al., 2024). We train our AEs on both



483 Figure 11: Sample-wise non-monotonicity and double descent for the linear subspace data. Re-484 sults for the feature noise scenario are presented in Appendix D, Figure 31 and the training losses in 485 Appendix C, Figure 26.

9

469



Figure 12: Top: test and train losses of different model sizes utilizing the single-cell RNA dataset. Training is done on the source data while testing on the target data. Bottom: UMAP of bottleneck vectors extracted from the encoder's output and KNN-DAT results for different model sizes.



Figure 13: Left, middle: test loss of clean and anomaly data respectively. Right: non-monotonic behavior of the ROC-AUC for the celebA dataset.

normal and anomaly data and detect anomalies based on the reconstruction loss as detailed in Sub-513 section 4.1. We used the CelebA attributes dataset and conducted an experiment similar to the one in 514 Subsection 4.1 to investigate how the model's size affect its ability to detect anomalies. As expected, 515 in line with the findings in (Han et al., 2022), small models outperform larger models in anomaly 516 detection (the highest ROC-AUC is achieved by the smallest models in Figure 13). Since we do 517 not control the SAR value in this data, which is positive, we do not observe a double descent in the 518 test loss curves. Nonetheless, we identify a non-monotonic behavior of the ROC-AUC curve. Ini-519 tially, it decreases for intermediate models, followed by an increase for over-parameterized models. 520 In conclusion, when employing a model for unsupervised anomaly detection, it is recommended to avoid selecting intermediate models, as their anomaly detection performance is inferior to under and over-parameterized models. 522

CONCLUSIONS 6

524 525

521

523

502

504 505

506

507

508

509

510

511

512

In our study, we identified various instances of multiple descents and non-monotonic behaviors 526 in unsupervised learning. These phenomena occur at the model-wise, epoch-wise, and sample-527 wise levels. We used under-complete AEs to investigate these phenomena and found compelling 528 evidence for their robustness across diverse datasets, model types, and experimental scenarios. We 529 examined four distinct use cases: sample noise, feature noise, domain shift, and anomalies. Our 530 experiments revealed multiple instances of consecutive descents, with most of them resulting in 531 improved (lower) test loss. Additionally, we found a connection between the model's size and 532 its real-world performance. Specifically, over-parameterized models can serve as effective domain 533 adaptation strategies when there is a distribution shift between the source and target data. In the 534 realm of anomaly detection, we find that it is important to avoid selecting intermediate models that yield lower ROC-AUC outcomes.Our work was limited by computational resources, preventing us 536 from using larger datasets for training. However, we hope that our findings will benefit research 537 groups with greater computational capabilities, enabling them also to explore other frameworks in unsupervised learning, such as generative models. Another exciting direction for future research is 538 developing theoretical frameworks that explain our findings, using similar ideas such as in (Curth et al., 2024; Curth, 2024).

540 REPRODUCIBILITY

541 542

546 547

548

552

553

554

560

567

568

569

570

576

580

581

Please refer to Appendix A for all the necessary information for reproducing the results. This includes a detailed explanation of the datasets, which is also discussed in Section 3. Additionally, Appendix A covers more datasets which are used in Appendix E, model types, hyperparameters, and the loss function used during training. The SNR calculations are provided in Appendix B.

- References
- Ben Adlam and Jeffrey Pennington. The neural tangent kernel in high dimensions: Triple descent and a multi-scale theory of generalization. In *International Conference on Machine Learning*, pp. 74–84. PMLR, 2020.
 - Madhu S Advani, Andrew M Saxe, and Haim Sompolinsky. High-dimensional dynamics of generalization error in neural networks. *Neural Networks*, 132:428–446, 2020.
- Peter L Bartlett, Philip M Long, Gábor Lugosi, and Alexander Tsigler. Benign overfitting in linear
 regression. *Proceedings of the National Academy of Sciences*, 117(48):30063–30070, 2020.
- ⁵⁵⁷ Christoph Baur, Stefan Denner, Benedikt Wiestler, Nassir Navab, and Shadi Albarqouni. Autoen ⁵⁵⁸ coders for unsupervised anomaly segmentation in brain mr images: a comparative study. *Medical* ⁵⁵⁹ *Image Analysis*, 69:101952, 2021.
- Mikhail Belkin, Daniel J Hsu, and Partha Mitra. Overfitting or perfect fitting? risk bounds for classification and regression rules that interpolate. *Advances in neural information processing systems*, 31, 2018.
- Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine learning practice and the classical bias-variance trade-off. *Proceedings of the National Academy* of Sciences, 116(32):15849–15854, 2019a.
 - Mikhail Belkin, Alexander Rakhlin, and Alexandre B Tsybakov. Does data interpolation contradict statistical optimality? In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 1611–1619. PMLR, 2019b.
- Andrea Borghesi, Andrea Bartolini, Michele Lombardi, Michela Milano, and Luca Benini. Anomaly
 detection using autoencoders in high performance computing systems. In *Proceedings of the AAAI Conference on artificial intelligence*, volume 33, pp. 9428–9433, 2019.
- Varun Chandola, Arindam Banerjee, and Vipin Kumar. Anomaly detection: A survey. ACM computing surveys (CSUR), 41(3):1–58, 2009.
- Woong-Gi Chang, Tackgeun You, Seonguk Seo, Suha Kwak, and Bohyung Han. Domain-specific
 batch normalization for unsupervised domain adaptation. In *Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition*, pp. 7354–7362, 2019.
 - Lin Chen, Yifei Min, Mikhail Belkin, and Amin Karbasi. Multiple descent: Design your own generalization curve. *Advances in Neural Information Processing Systems*, 34:8898–8912, 2021.
- Zhaomin Chen, Chai Kiat Yeo, Bu Sung Lee, and Chiew Tong Lau. Autoencoder-based network anomaly detection. In 2018 Wireless telecommunications symposium (WTS), pp. 1–5. IEEE, 2018.
- Zhen Cheng, Siwei Wang, Pei Zhang, Siqi Wang, Xinwang Liu, and En Zhu. Improved autoencoder
 for unsupervised anomaly detection. *International Journal of Intelligent Systems*, 36(12):7103–
 7125, 2021.
- Alicia Curth. Classical statistical (in-sample) intuitions don't generalize well: A note on bias-variance tradeoffs, overfitting and moving from fixed to random designs. *arXiv preprint arXiv:2409.18842*, 2024.
- Alicia Curth, Alan Jeffares, and Mihaela van der Schaar. A u-turn on double descent: Rethinking
 parameter counting in statistical learning. *Advances in Neural Information Processing Systems*, 36, 2024.

622

394	Yehuda Dar, Vidya Muthukumar, and Richard G Baraniuk. A farewell to the bias-variance
595	tradeoff? an overview of the theory of overparameterized machine learning arXiv preprint
596	arXiv:2109.02355, 2021.
597	

- Stéphane d'Ascoli, Levent Sagun, and Giulio Biroli. Triple descent and the two kinds of overfitting:
 Where & why do they appear? Advances in Neural Information Processing Systems, 33:3058–3069, 2020.
- Jun Deng, Zixing Zhang, Florian Eyben, and Björn Schuller. Autoencoder-based unsupervised
 domain adaptation for speech emotion recognition. *IEEE Signal Processing Letters*, 21(9):1068–
 1072, 2014.
- Carl Doersch. Tutorial on variational autoencoders. *arXiv preprint arXiv:1606.05908*, 2016.
- Marina Dubova. Generalizing with overly complex representations. In *NeurIPS 2022 Workshop on* Information-Theoretic Principles in Cognitive Systems, 2022.
- Andrés F Duque, Sacha Morin, Guy Wolf, and Kevin Moon. Extendable and invertible manifold
 learning with geometry regularized autoencoders. In 2020 IEEE International Conference on Big Data (Big Data), pp. 5027–5036. IEEE, 2020.
- Matteo Gamba, Erik Englesson, Mårten Björkman, and Hossein Azizpour. Deep double descent via
 smooth interpolation. *arXiv preprint arXiv:2209.10080*, 2022.
- Leo Gao, Tom Dupré la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya Sutskever, Jan Leike, and Jeffrey Wu. Scaling and evaluating sparse autoencoders. *arXiv preprint arXiv:2406.04093*, 2024.
- Daniel Gedon, Antônio H Ribeiro, and Thomas B Schön. No double descent in pca: Training and pre-training in high dimensions. 2022.
- Kiaoling Gong, Ling Yu, Jian Wang, Kai Zhang, Xiao Bai, and Nikhil R Pal. Unsupervised feature
 selection via adaptive autoencoder with redundancy control. *Neural Networks*, 150:87–101, 2022.
- Kai Han, Yunhe Wang, Chao Zhang, Chao Li, and Chao Xu. Autoencoder inspired unsupervised feature selection. In 2018 IEEE international conference on acoustics, speech and signal processing (ICASSP), pp. 2941–2945. IEEE, 2018.
- Songqiao Han, Xiyang Hu, Hailiang Huang, Minqi Jiang, and Yue Zhao. Adbench: Anomaly
 detection benchmark. *Advances in Neural Information Processing Systems*, 35:32142–32159, 2022.
- Trevor Hastie, Robert Tibshirani, Jerome H Friedman, and Jerome H Friedman. *The elements of statistical learning: data mining, inference, and prediction*, volume 2. Springer, 2009.
- Trevor Hastie, Andrea Montanari, Saharon Rosset, and Ryan J Tibshirani. Surprises in high dimensional ridgeless least squares interpolation. *Annals of statistics*, 50(2):949, 2022.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Chinmaya Kausik, Kashvi Srivastava, and Rishi Sonthalia. Double descent and overfitting under noisy inputs and distribution shift for linear denoisers. *arXiv preprint arXiv:2305.17297*, 2023.
- 640 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980, 2014.
 642
- 643 DP Kingma. Auto-encoding variational bayes. arXiv preprint arXiv:1312.6114, 2013.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25, 2012.
- 647 Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.

648 Gilad Lerman and Tyler Maunu. Fast, robust and non-convex subspace recovery. Information and 649 Inference: A Journal of the IMA, 7(2):277–336, 2018a. 650 Gilad Lerman and Tyler Maunu. An overview of robust subspace recovery. Proceedings of the 651 *IEEE*, 106(8):1380–1410, 2018b. 652 653 Zhu Li, Weijie Su, and Dino Sejdinovic. Benign overfitting and noisy features. arXiv preprint 654 arXiv:2008.02901, 2020. 655 Tengyuan Liang, Alexander Rakhlin, and Xiyu Zhai. On the multiple descent of minimum-norm 656 interpolants and restricted lower isometry of kernels. In Conference on Learning Theory, pp. 657 2683-2711. PMLR, 2020. 658 659 Ofir Lindenbaum, Yariv Aizenbud, and Yuval Kluger. Probabilistic robust autoencoders for outlier 660 detection. The Conference on Uncertainty in Artificial Intelligence (UAI), 2024. 661 Alisia Lupidi, Yonatan Gideoni, and Dulhan Jayalath. Does double descent occur in self-supervised 662 learning? arXiv preprint arXiv:2307.07872, 2023. 663 664 Pankaj Malhotra, Anusha Ramakrishnan, Gaurangi Anand, Lovekesh Vig, Puneet Agarwal, and Gautam Shroff. Lstm-based encoder-decoder for multi-sensor anomaly detection. arXiv preprint 665 arXiv:1607.00148, 2016. 666 667 William F Massy. Principal components regression in exploratory statistical research. Journal of 668 the American Statistical Association, 60(309):234–256, 1965. 669 Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation and 670 projection for dimension reduction. arXiv preprint arXiv:1802.03426, 2018. 671 672 Andriy Myronenko. 3d mri brain tumor segmentation using autoencoder regularization. In Brain-673 lesion: Glioma, Multiple Sclerosis, Stroke and Traumatic Brain Injuries: 4th International Work-674 shop, BrainLes 2018, Held in Conjunction with MICCAI 2018, Granada, Spain, September 16, 675 2018, Revised Selected Papers, Part II 4, pp. 311-320. Springer, 2019. 676 Preetum Nakkiran, Gal Kaplun, Yamini Bansal, Tristan Yang, Boaz Barak, and Ilya Sutskever. Deep 677 double descent: Where bigger models and more data hurt. Journal of Statistical Mechanics: 678 Theory and Experiment, 2021(12):124003, 2021. 679 Behnam Neyshabur, Zhiyuan Li, Srinadh Bhojanapalli, Yann LeCun, and Nathan Srebro. To-680 wards understanding the role of over-parametrization in generalization of neural networks. arXiv 681 preprint arXiv:1805.12076, 2018. 682 683 Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. Moment matching 684 for multi-source domain adaptation. In Proceedings of the IEEE/CVF international conference 685 on computer vision, pp. 1406–1415, 2019. 686 Amit Rozner, Barak Battash, Lior Wolf, and Ofir Lindenbaum. Domain-generalizable multiple-687 domain clustering. Transactions on Machine Learning Research, 2023. 688 689 Amit Rozner, Barak Battash, Henry Li, Lior Wolf, and Ofir Lindenbaum. Anomaly detection with 690 variance stabilized density estimation. The Conference on Uncertainty in Artificial Intelligence (UAI), 2024. 691 692 Mayu Sakurada and Takehisa Yairi. Anomaly detection using autoencoders with nonlinear dimen-693 sionality reduction. In Proceedings of the MLSDA 2014 2nd workshop on machine learning for 694 sensory data analysis, pp. 4-11, 2014. 695 Mark F Schilling. Multivariate two-sample tests based on nearest neighbors. Journal of the American 696 Statistical Association, 81(395):799-806, 1986. 697 Jonathon Shlens. Notes on kullback-leibler divergence and likelihood. arXiv preprint 699 arXiv:1404.2000, 2014a. 700 Jonathon Shlens. A tutorial on principal component analysis. arXiv preprint arXiv:1404.1100, 701

2014b.

702 703	Chunfeng Song, Feng Liu, Yongzhen Huang, Liang Wang, and Tieniu Tan. Auto-encoder based data clustering. In Progress in Pattern Recognition. Image Analysis, Computer Vision and An							
704	plications: 18th Iberoamerican Converses CIARP 2013 Havana Cuba November 20-23 2013							
705	Proceedings, Part I 18, pp. 117–124. Springer, 2013.							
706								
707	Rishi Sonthalia and Raj Rao Nadakuditi. Training data size induced double descent for denoising							
708	feedforward neural networks and the role of training noise. <i>Transactions on Machine Learning</i>							
709	Research, 2025.							
710	Stefano Spigler, Mario Geiger, Stéphane d'Ascoli, Levent Sagun, Giulio Biroli, and Matthieu Wyart.							
711	A jamming transition from under-to over-parametrization affects loss landscape and generation							
712	tion. arXiv preprint arXiv:1810.09665, 2018.							
71/	Ningyuan Teresa David W Hogg and Soledad Villar Dimensionality reduction regularization and							
715	generalization in overparameterized regressions SIAM Journal on Mathematics of Data Scien							
716	4(1):126–152, 2022.							
717								
718	Hoa Thi Nhu Tran, Kok Siong Ang, Marion Chevrier, Xiaomeng Zhang, Nicole Yee Shin Lee,							
719	Michelle Goh, and Jinmiao Chen. A benchmark of batch-effect correction methods for single-cell							
720	Tha sequencing data. Genome biology, 21.1–32, 2020.							
721	Nilesh Tripuraneni, Ben Adlam, and Jeffrey Pennington. Overparameterization improves robustness							
722	to covariate shift in high dimensions. Advances in Neural Information Processing Systems, 34:							
723	13883–13897, 2021.							
724	Descel Vincent Hugo Larochelle, Voshua Rengio, and Dierro Antoine Manzagol. Extracting and							
725	composing robust features with denoising autoencoders. In <i>Proceedings of the 25th international</i>							
726 727	conference on Machine learning, pp. 1096–1103, 2008.							
728	Pascal Vincent Hugo Larochelle Isabelle Lajoie, Yoshua Bengio, Pierre-Antoine Manzagol, and							
729	Léon Bottou. Stacked denoising autoencoders: Learning useful representations in a deep network							
730 731	with a local denoising criterion. Journal of machine learning research, 11(12), 2010.							
732	Wei Wang, Yan Huang, Yizhou Wang, and Liang Wang. Generalized autoencoder: A neural network							
733	framework for dimensionality reduction. In <i>Proceedings of the IEEE conference on computer</i> vision and pattern recognition workshops, pp. 490–497, 2014.							
734								
735 736	Mengzhou Xia, Mikel Artetxe, Chunting Zhou, Xi Victoria Lin, Ramakanth Pasunuru, Danqi Chen, Luke Zettlemoyer, and Ves Stoyanov. Training trajectories of language models across scales.							
737 738	arXiv preprint arXiv:2212.09803, 2022.							
739	Ji Xu and Daniel J Hsu. On the number of variables to use in principal component regression.							
740	Advances in neural information processing systems, 32, 2019.							
741	Yihao Xue, Kyle Whitecross, and Baharan Mirzasoleiman. Investigating the impact of model width							
742 743	and density on generalization in presence of label noise. arXiv preprint arXiv:2208.08003, 2022.							
744	Tamir Baruch Yampolsky, Ronen Talmon, and Ofir Lindenbaum. Domain and modality adaptation							
745	using multi-kernel matching. In 2023 31st European Signal Processing Conference (EUSIPCO),							
746	pp. 1285–1289. IEEE, 2023.							
747	Shuai Yang, Kui Yu, Fuyuan Cao, Hao Wang, and Xindong Wu. Dual-representation-based autoen-							
748 749	coder for domain adaptation. <i>IEEE Transactions on Cybernetics</i> , 52(8):7464–7477, 2021.							
750	Xu Yang, Cheng Deng, Feng Zheng, Junchi Yan, and Wei Liu. Deep spectral clustering using							
751	dual autoencoder network. In <i>Proceedings of the IEEE/CVF conference on computer vision and</i>							
752	pattern recognition, pp. 4066–4075, 2019.							
753								
754 755	Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning (still) requires rethinking generalization. <i>Communications of the ACM</i> , 64(3):107–115, 2021.							

756 757 758	Zhengming Zhang, Taotao Ji, Haoqing Shi, Chunguo Li, Yongming Huang, and Luxi Yang. A self- supervised learning-based channel estimation for irs-aided communication without ground truth. <i>IEEE Transactions on Wireless Communications</i> , 2023.
759	Chang They and Dandy C. Dofferenth Anomaly detection with reduct dam outcomedars. In Dre
760	choing Zhou and Randy C Partenrouti. Anomary detection with robust deep autoencoders. In Pro-
761	mining pp 665 674 2017
762	mining, pp. 005–074, 2017.
763	Kaiyang Zhou, Ziwei Liu, Yu Qiao, Tao Xiang, and Chen Change Loy. Domain generalization:
764	A survey. IEEE Transactions on Pattern Analysis and Machine Intelligence, 45(4):4396–4415,
765	2022.
766	
767	
768	
769	
770	
771	
772	
773	
774	
775	
776	
777	
778	
779	
780	
781	
782	
783	
784	
785	
786	
787	
788	
789	
790	
791	
792	
793	
794	
796	
797	
798	
799	
800	
801	
802	
803	
804	
805	
806	
807	
808	
809	

А IMPLEMENTATION DETAILS

In this section, we provide complete implementation details for all experiments conducted in the paper. Illustrations of the linear subspace dataset generation introduced in Subsection 3.1 for the scenarios of sample and feature noise, domain shift, and anomalies are displayed in Figures 14, 15, and 16 respectively. We also provide Table 1 for comparison between the double descent results of supervised and unsupervised regimes.



Figure 14: Data generation for the scenarios of sample and feature noise with p = 0.5. The first (leftmost) matrix depicts the latent vectors Z. The second matrix illustrates the latent vectors being projected into a higher dimensional space, and the rightmost matrices contain clean (blue) and noisy (red) samples / features respectively.



Figure 15: Data generation for the scenario of domain shift. The matrix on the left depicts the latent vectors Z and the two middle matrices represent the separation to source (Z_{train}) and target (Z_{test}) . The two rightmost matrices illustrate the latent vectors of the train and test data being projected into a higher dimensional space with different matrices (D, D''), resulting in a domain shift.



Figure 16: Data generation for the case of anomalies and p = 0.5. The matrix on the left depicts the anomalous data, the middle matrix represents the clean data, and the rightmost matrix contains both clean (blue) and outlier (red) samples.

Parameters. Table 2 details the hyper-parameters and other variables for the training process with the linear subspace, non-linear subspace (Appendix E.3), single-cell RNA, CelebA, and MNIST (Appendix E.2) datasets. The training optimizer utilized was Adam (Kingma & Ba, 2014), and the loss function for reconstruction is the mean squared error, which is mentioned in this Section.

Data. For the linear subspace, non-linear subspace, and MNIST datasets, we generate 5000 samples for training and 10000 for testing across all scenarios (sample noise, feature noise, domain shift, and anomaly detection). Regarding the single-cell RNA data, we have focused on dataset number 4 from (Tran et al., 2020), which includes 5 distinct domains (biological batches) named 'Baron', 'Mutaro,' 'Segerstolpe,' 'Wang,' and 'Xin', each representing 15 different cell types. Each cell (sample) in this dataset contains over 15000 genes (features). To facilitate the training of deep

865

Setup	Unsupervis	rvised Supervised			or Notes	Notes (unsup.)	
Sample noise Feature noise Domain shift Anomalies	exists (our pa exists (our pa exists (our pa exists (our pa	• paper)exists (see Sections 1, 2)• paper)exists (see Section 2)• paper)exists (see Subsection 5.1)• paper)not explored		2) similar) similar 5.1) similar -	low Sl	low SNR needed low SNR needed	
	Tal	ble 2: 1	Parameters and hyper-pa	arameters			
Parame	ters L	inear/	non-linear Subspace	RNA	CelebA	MNIST	
Model	М	ILP		MLP	MLP	CNN	
Learning rate	0.	001		0.001	0.001	0.001	
Optimizer	А	dam		Adam	Adam	Adam	
Epochs	20	00		1000	200	1000	
Batch size	10	C		128	10	128	
Data's latent si	ze (d) 20)		-	-	-	
Number of feat	tures (n) 50	C		1000	40	784	
Train dataset s	ize 50	000		5000	3000	5000	
SNR/ SAR [dE	3]	-20, -15, -10, -7, -5, -2, 0, 2					
Sample/ featur	e noise (p)		0, 0.1	, 0.2,,1			
Domain shift s	cale (s) 1,	2, 3, 4	ŀ	-	-	-	
Bottleneck lay	er size 25	5, 30, 4	15	20, 100, 300	25 10, 1	30, 50, 500	
Hidden layer s	ize 4	- 500		10 - 3000	4-400	-	
Channels	-			-	-	1-64	

Table 1: The existence of double descent in unsupervised and supervised learning.

models while preserving the domain shift, we have retained the top 1000 prominent features. We 891 utilize the 'Baron' biological batch as our source data for the scenario of domain shift, comprising 892 5000 training samples, while the target batches are 'Mutaro' (2122 samples), 'Segerstople' (2127 893 samples), 'Wang' (457 samples), and 'Xin' (1492 samples). As for the sample and feature noise 894 scenarios, we use the 'Baron' domain for both sample and feature noise scenarios due to its largest 895 sample size (8569). We allocate 5000 samples for training and introduce additive white Gaussian 896 noise (AWGN) to specific samples and features, as described in subsection 3.1. The calculations 897 of the SNR for both sample and feature noise cases are provided in Section B. The reserved 3569 898 samples are for testing. Please be aware that all the domains in this dataset are inherently noisy, 899 reflecting their real-world nature. Therefore, even when no additional noise is applied (p = 0), the data remains noisy. This may account for why the test loss does not decrease monotonically as the 900 901 model size increases for cases with low noise levels, as shown in Appendix E.4, Figure 55b.

For the **celebA** dataset, we sub-sample 3000 clean samples and replace $\lfloor 3000 \cdot p \rfloor$ of them with anomalies to ensure that $\sim p \cdot 100\%$ of the data is contaminated with anomalies. Due to the limited availability of anomaly data (4547 samples), the test set includes $\lfloor (1-p) \cdot 4547 \rfloor$ anomalies along with an equal number of clean samples.

907 Models. All experiments, including the linear subspace, non-linear subspace, single-cell RNA, 908 and celebA datasets are conducted using the same MLP AE architecture. To facilitate the exploration of double descent in both bottleneck layer size and hidden layer size, we employ a simplified 909 model mentioned in (Lupidi et al., 2023) consisting of a single hidden layer for both the encoder and 910 decoder, as depicted in Figure 17. We also utilize a CNN AE architecture consisting of three convo-911 lution layers in the encoder part, followed by a bottleneck layer, and then a decoder part consisting 912 of three deconvolution layers trained on the MNIST dataset as illustrated in Figure 18 (results for 913 the CNN AE are reported in Appendix E.2). 914

We work with under-complete AEs to encourage the acquisition of a meaningful embedding in the
latent space and prevent the model from learning the identity function. The size of these models is
determined by the sizes of the hidden layers (for MLP), the number of channels (for CNN), and the
bottleneck layer, while the width of these models remains constant.



Figure 17: Left: Demonstration of the MLP-based AE model structure. Right: model's number of parameters for single-cell RNA settings (bottleneck layer size = 300 and input size = 1000 features).



Figure 18: Upper: Demonstration of CNN AE model structure. Lower: number of parameters.

Loss function. All AEs are trained with the mean squared error (MSE) loss function:

$$MSE = \frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2.$$

Where n is the number of data samples, y_i is the true value, and \hat{y}_i is the predicted value. Due to contamination in the training dataset, the norm of train samples tends to be higher than that of the clean test samples. As the MSE loss is not scale-invariant, we opt to normalize both train and test losses only after the training process is complete, using $\frac{1}{n}\sum_{i=1}^{n}(y_i - \bar{y})^2$, Where \bar{y} is the mean of $\{y_i\}_{i=1}^n$. This strategy enables us to continue utilizing the MSE loss function while facilitating a fair and meaningful comparison between train and test losses.

Results. Ensuring the robustness of the findings across various model initializations and enhancing their reliability, all figures combine several results of different random seeds. The bolded curves in each figure represent the average across the results of different seeds, and the transparent curve around each bolded curve represents the ± 1 standard error from the mean.

Environments and Computational Time. All experiments were conducted on NVIDIA RTX 6000 Ada Generation with 47988 MiB, NVIDIA GeForce RTX 3080 with 10000 MiB, Tesla V100-SXM2-32GB with 34400 MiB, and NVIDIA GeForce GTX 1080 Ti with 11000 MiB.

Each result in Figure 3 represents an average over 10 seeds. The hidden layer sizes for the linear subspace data range from 4 to 500 with a step size of 4, and for the single-cell RNA data, they range from 10 to 500 with a step size of 10, and from 500 to 3000 with a step size of 50. This

972 results in 125 and 110 models trained for each dataset, respectively. Figure 3a illustrates 10 different 973 sample noise levels, requiring the training of $125 \times 10 \times 10 = 12,500$ models. Similarly, Figure 974 3b depicts 4 different sample noise levels, corresponding to $110 \times 10 \times 4 = 4,400$ trained models. 975 In total, 16,900 models, each with up to 8 million parameters trained on 5,000 data points were 976 needed to obtain the results. In Appendix E.2, Figure 34, we present results for CNNs including between 1 to 64 channels trained on 5,000 images from the MNIST dataset including 5 levels of 977 sample noise and 6 levels of feature noise for 10 different seeds. These experiments require training 978 $64 \times 11 \times 10 = 7,040$ models with up to 13 million parameters. Each evaluation of a specific 979 experiment takes several days if trained on the NVIDIA RTX 6000 and weeks if trained on the other 980 mentioned GPUs to obtain the results. 981

982 983

984

988 989

В **SNR** CALCULATIONS

In this section, we will outline our approach for calculating the signal-to-noise ratio (SNR) for all 985 experiments involving the addition of noise. Initially, we convert the SNR from decibels to linear 986 SNR using the formula: 987

$$SNR = 10^{\left(\frac{SNR[dB]}{20}\right)}.$$
(1)

We have a closed-form equation for the linear subspace dataset to determine the scalar θ required to 990 multiply the train samples and achieve the desired linear SNR value. We use the fact that both train 991 and noise are sampled from an i.i.d. normal distribution and calculate θ for the sample noise, feature 992 noise, domain shift, and anomalies. 993

994 Notations:

 $z - d \times 1$ vector. Represents a vector in a latent space of size d. 995 $D - n \times d$ matrix. Represents a random matrix to project z from a d dimensional space into a 996

higher-dimensional space (n > d). 997

 $\epsilon - n \times 1$ vector. Represents the noise added to a vector with n dimensions.

1000 For the scenario of sample noise, where a particular sample is affected by noise across all its features:

1001 1002

1003 1004

1007

1 1 1

998 999

$$SNR^{2} = \frac{E[\|\theta Dz\|_{2}^{2}]}{E[\|\epsilon\|_{2}^{2}]} = \frac{E[\theta^{2}z^{T}D^{T}Dz]}{E[\epsilon^{T}\epsilon]} = \frac{\theta^{2}E_{z}[E_{D|z}[z^{T}D^{T}Dz|z]]}{E[\sum_{i=1}^{n}\epsilon_{i}^{2}]} \underbrace{=}_{(a)}$$
(2)

$$\frac{\theta^2 E_z[z^T E_{D|z}[D^T D]z]}{n} \underbrace{=}_{\text{(b)}} \frac{\theta^2 E_z[z^T n \cdot I_{d \times d}z]}{n} = \frac{\theta^2 \cdot n \cdot E_z[z^T z]}{n} = \theta^2 E\left[\sum_{i=1}^d z_i^2\right] \underbrace{=}_{\text{(a)}} \theta^2 \cdot d.$$

Isolating θ , we get that $\theta = \frac{\text{SNR}}{\sqrt{d}}$. 1008

1009 (a) Given a vector $a \sim \mathcal{N}(0, I_n)$ of *n* i.i.d. samples, $E\left[\sum_{i=1}^n a_i^2\right] = \sum_{i=1}^n E[a_i^2] = \sum_{i=1}^n 1 = n$. 1010 (b) Given a matrix $M \sim \mathcal{N}(0, I_n)$ of size $n \times n$ where all entries are i.i.d., then 1011

1019 For the scenario of feature noise, each train sample has only $n \cdot p$ noisy features, meaning the noise 1020 vector contains values for only $n \cdot p$ entries. Consequently, θ is determined by $\sqrt{\frac{p}{d}} \cdot \text{SNR}$. For 1021 practitioners who want to explore the scenario involving domain shift, where the source and target 1022 are noisy, note that the matrix responsible for projecting z_{test} into a higher-dimensional space is 1023 denoted as $D^{''} = D + s \cdot D^{'}$ where $D^{'}$ is sampled from a standard normal distribution $\mathcal{N}(0, I)$ and 1024 both D and D' are i.i.d. Consequently, $D_{ij}^{''} \sim \mathcal{N}(0, 1 + s^2)$. Substituting D with $D^{''}$ in equation 1025 equation 2, we find that $E_{D''|z}[D''^T D''] = n \cdot (1+s^2) \cdot I_d$, leading to $\text{SNR}^2 = (1+s^2) \cdot \theta^2 \cdot d$,

therefore $\theta = \frac{\text{SNR}}{\sqrt{(s^2+1)d}}$. In other words, since the covariance matrix of $D^{''}$ is $(1+s^2)I$, we need to make sure we first normalize the matrix by $\sqrt{1+s^2}$ to maintain the identity covariance matrix. For other datasets, such as the single-cell RNA dataset, we normalize each sample x by its norm ||x||, and similarly normalize each noise vector n, yielding: $\hat{x} = \frac{x}{||x||}$ and $\hat{n} = \frac{n}{||n||}$. This ensures that the ratio $\frac{\hat{x}}{\hat{\alpha}}$ equals 1. By employing equation equation 1, we attain the intended linear SNR factor θ , and then scale down \hat{n} by θ , yielding $\hat{n}_{scaled} = \frac{\hat{n}}{\theta}$. This guarantees that the linear SNR is $\frac{\hat{x}}{\hat{n_{scaled}}} = \theta.$

C TRAIN LOSS RESULTS



In this section, we provide the train loss figures corresponding to each of the test losses mentioned in the main paper.



1134 D RESULTS FOR FEATURE NOISE

In this section, we present the results for the feature noise scenario. Feature noise adds complexity since each sample contains noise in some of its features. As a result, the model never encounters samples with entirely clean features, making it unable to isolate and focus on clean data. Consequently, the model experiences difficulty in learning the correct data structure. Surprisingly, increasing feature noise actually leads to a decrease in the test loss for the single-cell RNA dataset (Figure 27b). This can also be observed in Appendix E.2, Figure 34b and Appendix E.3, Figure 44b. Moreover, the peak shifts left as the number of noisy features rise in Figure 27a.



(a) Double descent for the linear subspace data
trained with SNR = -13 [dB]. The model also exhibits the final ascent phenomenon (Xue et al., 2022).



(b) Non-monotonic behavior for the **Single-cell RNA data** trained with SNR = -12 [dB]. Beyond a hidden layer size of 2000, the test loss continues to decrease, while the train loss increases.





(a) Linear subspace data with 40% noisy features. Beyond hidden layer of size 300, the test loss rises.



(b) **Single-cell RNA data** with 10% noisy features. Beyond a hidden layer size of 2600, the test loss continues to decrease, and the train loss increases.

1187

Figure 28: The effect of SNR for the case of **noisy features** on the test loss curve.



1242 E ADDITIONAL EXPERIMENTS

1249

1250

1251

1257

1259 1260

1261

1244 E.1 More results for domain adaptation

This section presents the UMAP visualizations of the different domains for both the linear subspace and single-cell RNA data in Figure 32. Results for different model sizes trained on the linear subspace dataset are also reported in Figure 33.





(a) The UMAP representation shows a clear domain shift between the source and target datasets.

(b) Clusters represent different cell types. Different domains are represented by different colors.

Figure 32: UMAP representations of source and target datasets for the linear subspace dataset (left) and single-cell RNA dataset (right).

1265 Figure 33 illustrates the results based on a similar experiment conducted in Section 5.1 for the lin-1266 ear subspace data. As expected, the interpolating models exhibit the poorest KNN-DAT outcomes. 1267 Over-parameterized models introduce a decrease in the test loss indicating an improved reconstruc-1268 tion of the target data. In this scenario, we noticed that smaller models perform better than over-1269 parameterized models based on KNN-DAT results. We think that the small size of the hidden layer 1270 (4) and the high dimensionality of the dataset (50 features) result in significant information loss in these layers. This could lead to closely clustered vectors in the embedding space, ultimately causing 1271 low KNN-DAT results. However, a hidden layer of size 4 indicates insufficient capacity to represent 1272 the signal, as shown by the high values of test and train losses in Figure 33. 1273



Figure 33: UMAP of the latent (bottleneck) vectors with a size of 45 and KNN-DAT results for different model sizes trained on the **linear subspace** dataset for a shift of 3.



1296 E.2 DOUBLE DESCENT RESULTS FOR CNNS TRAINED ON MNIST



1317 noise (left) and feature noise (right).

1315

1318 In this section, we demonstrate that the double descent phenomenon can be reproduced in other 1319 unsupervised AE architectures. We employed the MNIST dataset (LeCun et al., 1998) and trained 1320 under-complete CNNs as detailed in Figure 18. For the case of sample noise, the noise is added to 1321 $p \cdot 100\%$ of the images, and for the feature noise scenario, noise is introduced to $p \cdot 100\%$ of the 1322 pixels of each image. To demonstrate the phenomenon with the presence of domain shift, the model is trained on the MNIST-M and MNIST datasets and tested on MNIST and MNIST-M, respectively. 1323 Results for model-wise double descent for varying amounts of sample and feature noise cases are 1324 presented in Figure 34. 1325

1326 In Figure 35, we show the test and train loss results (top two sub-figures) for three different mod-1327 els trained on MNIST with 50% sample noise and an SNR of -15 dB and find out that overparameterized models can reduce the noise levels in an image. The smallest model, with 3 channels, 1328 is under-parameterized. The second model, within the critical regime, with 5 channels, performs 1329 poorly, while the third is over-parameterized, containing 60 channels. Interestingly, We noticed that 1330 even though our AE was not trained to remove noise (as in denoising AEs (Vincent et al., 2008; 1331 2010)), over-parameterized models were able to reduce noise to some extent. In contrast, models 1332 within the critical regime performed significantly worse. 1333

1334 After training, we evaluated each model by feeding it images with varying SNR values and examining the reconstructed outputs (bottom sub-figure in Figure 35). The over-parameterized model 1335 produced the best-quality reconstructed images. Following that, the under-parameterized model per-1336 formed moderately well, and the model in the critical regime generated the noisiest images. This is 1337 because the critical model focused on memorizing the noise during training instead of learning the 1338 underlying signal, resulting in consistently noisy outputs. In contrast, the over-parameterized model 1339 had enough capacity to memorize the noise and learn the signal. While the under-parameterized 1340 model cleans the images better than the critical model, it still distorts some details compared to the 1341 over-parameterized model due to its limited capacity. 1342

To quantify noise reduction, we used the Peak Signal-to-Noise Ratio (PSNR), a metric that assesses signal quality by comparing the original image to its noisy version. PSNR measures the ratio between the maximum possible value of a signal (R^2) and the power of the noise (MSE). Higher PSNR values indicate better quality, meaning less noise. The formula for PSNR is

$$PSNR = 10 \cdot \log\left(\frac{R^2}{MSE(x, f(x+n))}\right)$$
1348

where x + n represents the noisy image (*n* is the noise), and *x* is the clean version. This metric, expressed in decibels, allows us to evaluate how well each model cleans the images. As shown,



Figure 35: Models trained on 50% noisy MNIST images with SNR = -15 [dB] and tested on MNIST images with different values of SNR.

the over-parameterized model consistently achieves the highest PSNR values (highlighted in bold green), while the poorly interpolating model, which primarily memorized noise, produces the lowest PSNR values (in red). In conclusion, *over-parameterized models are capable of reducing noise when trained on noisy data, even without being explicitly tasked to do so.*

We proceed by illustrating the impact of SNR on the test loss curve for both sample and feature noise scenarios in Figure 36. As expected, the test loss increases for low SNR values. We then investigate the effect of domain shifts between the training and testing datasets in two cases. First, models are trained on the MNIST dataset and tested on the MNIST-M dataset, as shown in Figure 37a. Second, models are trained on MNIST-M and tested on MNIST, as seen in Figure 37b. In both cases, the model-wise double descent curve is observed.

1403 We further illustrate this phenomenon along the epochs axis, displaying non-monotonic behavior and double descent under different levels of sample and feature noise (Figure 38) and showing the



Figure 36: Model-wise double descent for CNN trained on MNIST with varying levels of SNRs.Left: sample noise, right: feature noise.



(b) Source data: MNIST-M, target data: MNIST.

Figure 37: Model-wise double descent for CNN trained and tested on different domains. impact of SNR variation (Figure 39). Additionally, we provide similar results under domain shift conditions between the train and test datasets (Figure 40). Sample-wise double descent and nonmonotonic behavior is observed as well in all contamination setups. The cases of varying levels of sample noise and feature noise are displayed in Figure 41 and for varying SNRs for both scenarios in Figure 42. Sample-wise double descent is also illustrated in Figure 43 for when a domain shift is present between the training data (MNIST) and the testing data (MNIST-M).

1449 1450

1451

E.3 DOUBLE DESCENT RESULTS FOR THE NON-LINEAR SUBSPACE DATASET

(a) Source data: MNIST, target data: MNIST-M.

Building on the linear subspace dataset discussed in Subsection 3.1, we have developed a new dataset
with non-linear characteristics to investigate the double descent phenomenon in more complex scenarios. Although the single-cell RNA dataset is already non-linear, we have created this dataset to
demonstrate the reproducibility of the double descent phenomenon across various datasets.

As in the linear subspace model discussed in Subsection 3.1, we sample N latent vectors $\{z_i\}_{i=1}^N$ from a normal distribution and project them to a higher dimension using a random matrix D_1 . The key difference is the inclusion of non-linear components z_i^2 and z_i^3 , each projected to a higher



Figure 38: Epoch-wise non-monotonic behavior for varying levels of sample noise (left) and feature noise (right).



Figure 39: Epoch-wise double descent and non-monotonic behavior for varying SNRs. Left: sample noise, right: feature noise.

1496 1497

1498

1499 1500 1501

1506 1507 dimensional space with different random matrices D_2 and D_3 . To create contaminated setups of sample and feature noise, noise is added to $p \cdot 100\%$ of the data where θ controls the SNR:

$$x_{i} = \begin{cases} \theta(D_{1}z_{i} + D_{2}z_{i}^{2} + D_{3}z_{i}^{3}) + \epsilon_{i}, & \text{with probability } p, \\ \theta(D_{1}z_{i} + D_{2}z_{i}^{2} + D_{3}z_{i}^{3}), & \text{with probability } 1 - p, \end{cases}$$

For the domain shift scenario, we divide the latent vectors into training and testing sets and use the same parameter 's' as described in Subsection 3.1 to control the shift between the train and test sets in the following manner: $D''_i = D_i + s \cdot D'$ for $1 \le i \le 3$ and get:

$$x_i = \begin{cases} D_1 z_{train}^i + D_2 (z_{train}^i)^2 + D_3 (z_{train}^i)^3, & \text{if } train, \\ D_1^{''} z_{test}^i + D_2^{''} (z_{test}^i)^2 + D_3^{''} (z_{test}^i)^3, & \text{if } test. \end{cases}$$

For anomaly detection, clean samples are represented by $\theta(D_1z_i + D_2z_i^2 + D_3z_i^3)$, with $p \cdot 100\%$ of them replaced by anomalies sampled from a normal distribution, as detailed in Subsection 3.1.

1511 We start by presenting results for the sample and feature noise scenarios as depicted in Figure 44. As shown, the test loss results for the case of sample noise (Figure 44a) resemble those of the linear



Figure 41: Sample-wise double descent and non-monotonic behavior for varying levels of sample noise (left) and feature noise (right).

subspace data model presented in Figure 3a. Figure 44b demonstrates the model-wise final ascent phenomenon for the case of feature noise as elaborated in Appendix E.4. Figure 45 shows how the SNR affects the test loss curve for both sample and feature noise cases. As observed, the test loss increases with decreasing SNR. Additionally, the final ascent in the test loss is depicted in 45b for the feature noise scenario, where the slope becomes steeper as the SNR decreases. We also demonstrate the double descent and final ascent results regarding the domain shift scenario in Figure 46 and the anomaly detection capabilities in Figure 47.

We also observed epoch-wise double descent and non-monotonic behavior for this dataset, as shown in Figure 48 for different percentages of sample and feature noise and in Figure 49 for varying SNR levels under the same noise conditions. Additionally, epoch-wise double descent is also observed when a domain shift is present between the train and test sets, as depicted in Figure 50. Instances of sample-wise double descent and non-monotonic curves are also reported and displayed in Figure 51 for varying levels of sample and feature noise, Figure 52 for varying levels of SNR, and in Figure 53 for domain shift.



Figure 42: Sample-wise double descent and non-monotonic behavior for varying levels of SNR. Left: sample noise, right: feature noise.



Figure 43: Sample-wise double descent for models trained on the MNIST dataset and tested on the MNIST-M dataset.

1596 E.4 FINAL ASCENT PHENOMENON

While training various models on different datasets contaminated with sample and feature noise at 1598 different SNR levels and domain shifts between train and test sets, we observed a final ascent phe-1599 nomenon characterized by a pattern of decreasing-increasing-decreasing-increasing test loss. The phenomenon was first observed in Xue et al. (2022) in supervised learning with label noise. We sus-1601 pect a potential connection to this phenomenon in unsupervised learning, which we have yet to fully analyze. We refer to Figure 55a, which illustrates the final ascent results for the linear subspace 1603 dataset under extreme conditions of 100% sample noise, as a continuation of Figure 3a. We also 1604 present the final ascent results for the single-cell RNA dataset in Figure 55b. Another instance of final ascent with the presence of varying feature noise is illustrated in Figure 27a for the linear subspace dataset and in Figures 44b, 45b for the non-linear subspace dataset. Results are also replicated 1606 using the non-linear subspace dataset under various domain shifts, as observed in Figure 46.

1608

1597

1584

1609 E.5 MULTIPLE DESCENTS UNDER DIFFERENT NOISE TYPES AND SPARSE AES

This section explores the emergence of double and triple descent for noise distributions beyond
Gaussian noise and sparse AEs. Figure 55 illustrates the phenomenon for the linear subspace and
single-cell RNA datasets when subjected to Laplacian noise. The experimental setup mirrors that of
Figure 3. As shown, both datasets exhibit similar results under these conditions.

We extend our research to recent applications of AEs, including sparse AEs, which are increasingly utilized in explainable AI (XAI) (Gao et al., 2024) and have been adopted by Google in their Gemini project. Using sparse CNN AEs, we trained models on the MNIST dataset containing 80% noisy samples and observed the emergence of double descent. The models were configured with a bottleneck layer of size 550, and the parameter k, determining the top k highest bottleneck values to retain, was set to 500. The results are illustrated in Figure 56.



(a) Sample noise scenario with SNR = -15 [dB].

1638

1639

1640

1661

(Appendix E.4) with SNR = -20 [dB].

Figure 44: Model-wise double descent for the non-linear subspace data with varying levels of sample noise (left) and final ascent with varying levels of feature noise (right).



Figure 45: Effect of SNR on the test loss curve as a function of model size. Left: sample noise scenario. Right: feature noise scenario. 1662



Figure 46: Model-wise double descent and final ascent for the scenario of domain shift.

1682

1683

1684



Figure 47: Non-linear anomaly data with SAR = -15 [dB]. Left: test loss of the clean samples. A double descent pattern emerges for low SARs and high anomaly presence in the training data. Middle: test loss of the anomaly data. Right: Non-monotonic behavior of the ROC-AUC.



Figure 48: Epoch-wise double descent and non-monotonic behavior for varying levels of sample noise (left) and feature noise (right). For the scenario of feature noise, we mostly noticed the non-monotonic curve at 10%.



Figure 49: Epoch-wise double descent and non-monotonic behavior for varying levels of SNR. Left:sample noise, right: feature noise.



Figure 51: Sample-wise double descent and non-monotonic behavior for varying levels of sample noise (left) and feature noise (right).



Figure 52: Sample-wise double descent and non-monotonic behavior for varying levels of SNR.
Left: sample noise, right: feature noise.





Figure 56: Test loss exhibits model-wise double descent for sparse CNN AEs trained on MNIST with bottleneck layer size of 550 and k = 500.