
Absolute Variation Distance: an Inversion Attack Evaluation Metric for Federated Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Federated Learning (FL) has emerged as a pivotal approach for training models on
2 decentralized data sources by sharing only model gradients. However, the shared
3 gradients in FL are susceptible to inversion attacks which can expose sensitive
4 information. While several defense and attack strategies have been proposed,
5 their effectiveness is often evaluated using metrics that may not necessarily reflect
6 the success rate of an attack or information retrieval, especially in the context
7 of multidimensional data such as images. Traditional metrics like the Structural
8 Similarity Index (SSIM), Peak Signal-to-Noise Ratio (PSNR), and Mean Squared
9 Error (MSE) are typically used as lightweight metrics, assume only pixel-wise
10 comparison, but fail to consider the semantic context of the recovered data. This
11 paper introduces the Absolute Variation Distance (AVD), a lightweight metric
12 derived from total variation, to assess data recovery and information leakage in
13 FL. Unlike traditional metrics, AVD offers a continuous measure for extracting
14 information in noisy images and aligns closely with human perception. Our results
15 are combined with a user experience survey demonstrate that AVD provides a more
16 accurate and consistent measure of data recovery. It also matches the accuracy of
17 the more costly and complex Neural Network based metric, the Learned Perceptual
18 Image Patch Similarity (LPIPS). Hence it offers an effective tool for automatic
19 evaluation of data security in Federation and a reliable way of studying defence
20 and inversion attacks strategies in FL.

21 1 Introduction

22 In the age of Large Models (LM), with size of billions of parameters, data play a crucial role for
23 their continuous development. Therefore, the availability of large amounts of data for their training
24 and fine-tuning is critical. Traditionally, data was concentrated in centralized repositories, but the
25 increased awareness of privacy and the decentralized nature of information generation (mobile phones,
26 multiple regional data centres) has necessitated a more nuanced approach. In this regard, Federated
27 Learning (FL) enables models to learn from a multitude of decentralized edge devices or servers
28 holding local data samples, obviating the need to exchange raw data. The standard FL configuration
29 is achieved with a central aggregator node that exchanges gradients to train a centralised model
30 (McMahan et al., 2017). Particularly, at each training step t , a client node receives neural network
31 model weights, W_t , from an aggregator server and calculates loss l with local data (x_t, y_t) for a batch,
32 B , which generates gradients with respect to the model weights:

$$\Delta W_t = -\frac{\gamma}{B} \sum_{b < B} \frac{\partial l(F_{W_t}(x_{t,b}, y_{t,b}))}{\partial W_t}, \quad (1)$$

33 where F_{W_t} is a neural network parameterized by W_t . The gradients are typically averaged in the
34 server with a rate, γ . Because of their flexibility and the client anonymity that they offer, FL models

35 have been deployed in a variety of real-world applications (Yang et al., 2019; Rieke et al., 2020;
 36 Nguyen et al., 2022)

37 However, the gradients, ΔW_t , shared by the client are vulnerable to inversion attacks instigated by a
 38 malicious eavesdropper that can expose the original sensitive data. Existing literature on inversion
 39 attacks (Zhu et al., 2019b; Zhao et al., 2020a; Geiping et al., 2020; Yin et al., 2021; Balunović et al.,
 40 2022) have shown that these inversion attacks can be highly successful, with potentially recovering
 41 large batch of data after several rounds of learning, and at a pixel resolution Geiping et al. (2020).
 42 In general, these attacks more or less follow the same paradigm, generate a dummy dataset (usually
 43 images) and then use a loss function with priors distributions (Balunović et al., 2022) (as regulators
 44 or with additional generative models) to minimise the loss between the FL model and the dummy
 45 gradients. The success of the inversion attacks prevents FL from becoming a fully trustful framework
 46 for distributed training.

47 To mitigate inversion attacks in FL, several defence strategies were proposed to reduce the leakage of
 48 information (Sikandar et al., 2023; Huang et al., 2021a; Chen et al., 2022; Wainakh et al., 2022). These
 49 include, data transformation from the client side (Huang et al., 2021b), homomorphic encryption
 50 techniques (Phong et al., 2018), data sanitation methods (Zhu et al., 2019b), and defense strategies
 51 originated from Differential Privacy approaches (DP) (Dwork, 2006). In FL, DP techniques can be
 52 achieved by adding noise to the gradients shared or input data, with inevitably, a potential loss in
 53 model training performance (Zhu et al., 2019a; Zhao et al., 2020b; Eloul et al., 2022).

54 One important step towards the development of robust defence strategies and the prevention and
 55 understanding of such attacks is by assessing their success. Therefore, an inversion attack is studied
 56 experimentally by comparing the information revealed from the recovered input data to that of the
 57 original dataset (Huang et al., 2021a). There exist a few metrics that are commonly used to measure
 58 the reconstruction quality, with the most popular being the structural similarity index (SSIM) (Wang
 59 et al., 2004), the peak signal-to-noise ratio (PSNR) (Cahn, 1961), the learned perceptual image patch
 60 similarity (LPIPS) (Zhang et al., 2018), and the mean squared error (MSE).



Figure 1: Random recovered vectors from the MNIST dataset. Sub-figure (a) shows the images sorted by the AVD metric. Sub-figure (b) shows the images sorted by the MSE metric. The smaller figures represent the original images associated with image above. The figures include complex images, such as Sub-figure (a), row 1 & column 3, that is a combination of digits 4 and 3. Sub-figure (b), row 3 & column 4, that is a combination of digits 2 and 8.

61 In this paper, we show how the aforementioned metrics are insufficient to properly assess the success
 62 of an inversion attack on the gradients of an FL model. Especially, for cases that the FL model
 63 generates multidimensional outputs that contain contextual information (like images). The reason is
 64 that these metrics assume spatial independence when comparing the recovered image to the original
 65 one. Therefore, in many cases they lack to produce accurate results of the semantic context, and
 66 fail to reveal minimal information out of noisy image. For example, a common defense mechanism

67 for an FL model is to add noise to an image that depicts a number from the MNIST dataset. An
 68 attack on the gradients of this model may recover an approximate image with the same number
 69 but considerably different background colour (see Fig. 1a, row 4 & column 2). Metrics, such as
 70 SSIM, PSNR, and MSE fail to provide a consistent and accurate result and indicate this attack as
 71 unsuccessful (in the image example $MSE = 0.52$, which is considered high). That is, because of the
 72 use of Euclidean distance-per-pixel measure, they miss the fact that the attacker has recovered the
 73 most important element, the actual number; regardless of how noisy or changed the background of
 74 the image is. Therefore, these metrics discard the contextual information in the image that a human’s
 75 vision would have otherwise recognized, e.g. edges and points of interests (another example is 1a,
 76 row 2 & column 3, where $MSE = 1.06$ but a human would have read the number). It becomes
 77 even a larger challenge to use these metrics as a mechanism to approve data for FL in real-time.
 78 Since the development of attacks models are mostly empirical and data dependant, it is plausible
 79 to have an automatic verification (e.g. as a smart contract/client service) to assess the security of
 80 data by applying brute-force attacks before submission of gradients. For that purpose, a reliable and
 81 lightweight metric is needed.

82 This problem has not gone unnoticed and efforts to address these challenges have resulted in proposal
 83 of specialized metrics that may be computationally expensive, for example, Learned Perceptual
 84 Image Patch Similarity (LPIPS) from Zhang et al. (2018). The authors use the power of deep neural
 85 networks (DNN) to create LPIPS that is aligned with human perception metric. A major downside is
 86 that LPIPS is a complex and a computationally costly metric that is difficult to interpret due to the
 87 underlying DNN themselves that may require training for new data.

88 This paper introduces a new distance metric to assess data recovery, the *Absolute Variation Distance*
 89 (AVD). It is derived from total variation and in contrast to standard methods (MSE, SSIM), it offers
 90 a continuous metric for extracting information in noisy images. Furthermore, we show via a user
 91 study that AVD is highly correlated with human perception, but at the same time it is computationally
 92 more efficient and interpretable compared to LPIPS. Our results show that recovery of data is more
 93 visible as AVD decreases in a continuous manner. In contrast the MSE metric for MNIST fluctuates
 94 drastically when the image is not completely clear or a blend, and can obtain various values similar
 95 or higher than the MSE for the pure noise input.

Table 1: Types of gradient inversion attacks employed in our study.

Attack Name	Main Objective Function	Description
2-norm	g^2 (3)	Euclidean distance and initial label determination.
Angle & var	$g^{ang} + TV$ (4)	Geiping et al. (2020) proposed to leverage cosine similarity, total variation (TV) and initial label determination.
Angle & var & Orth_regulators	$g^{ang} + TV + Orth$	Cosine distance with orthogonal regulator for the input + initial label determination. (Qian et al., 2021)



Figure 2: Random recovered vectors from the LFW face dataset. Sub-figure (a) shows the images sorted by the AVD metric. Sub-figure (b) shows the images sorted by the MSE metric. The smaller figures represent the original images associated with image above.

96 **2 Absolute Variation Distance**

97 In this paper we developed AVD, a variant of total variation metric (Rudin et al., 1992), which is a
 98 more suitable indicator to compare the spatial gradient of the recovered image and source image.
 99 Given two images v^{source} and v^{target} , we define AVD between them as following:

$$\begin{aligned} \text{AVD}(v^{source}, v^{target}) = & \\ & \|(|\nabla v^{source}| - |\nabla v^{target}|)\| + \\ & \|(|\nabla^2 v^{source}| - |\nabla^2 v^{target}|)\| \end{aligned} \quad (2)$$

100 where $\nabla v = \frac{dv}{di} + \frac{dv}{dj}$ is the spatial gradient and $\nabla^2 v = \frac{d^2v}{di^2} + \frac{d^2v}{dj^2}$ is the second order gradient.
 101 Here we treat the image as a 2-D array with values $v(i, j)$. Therefore, because AVD measures
 102 distance in gradient space it allows to consider boundaries and edges in images which are a common
 103 discriminator in visual recognition, whilst the gradient of noise remains as noise.

104 **2.1 Inversion Attack Algorithms**

105 In our setting (and typically) the gradient inversion attack is carried out by choosing x'_t, y'_t on a proxy
 106 model, $F'(x'_t, y'_t)$, and then minimizing an objective function that measures the distance between
 107 gradients computed the proxy model $\Delta W'_t$ and the original gradients. A typical objective can be the
 108 norm of the gradients' difference:

$$g^{l2}(x'_t, y'_t) = \min \|\Delta W'_t - \Delta W_t\| \quad (3)$$

109 This solution searches for a model $F'(x'_t, y'_t)$ that matches the size of the gradient vector observed
 110 by the client. Although further empirical studies have found the cosine distance to provide better
 111 convergence results (Geiping et al., 2020):

$$g^{ang}(x'_t, y'_t) = \min 1 - \frac{\langle \Delta W'_t, \Delta W_t \rangle}{\|\Delta W'_t\| \cdot \|\Delta W_t\|} \quad (4)$$

112 Various regularisation terms were shown to improve convergence. For example, regularisation that
 113 penalises high variations in the input images and constrains the search to high-fidelity images and
 114 de-noised solutions (Geiping et al., 2020; Yin et al., 2021). In mini-batches the orthogonality (Qian
 115 et al., 2021) between input vectors in the batch has been shown to bias the search towards different
 116 vectors in the batch. Additionally it has been found that determining the label from the gradients is
 117 important for initialisation of the numerical optimisation (Yin et al., 2021).

118 In our study in section 4 we apply various types of attacks and regularisation terms to provide
 119 a comprehensive analysis without any prior assumption on the performance of the attack. As
 120 summarised in Table 1, we utilise both the Euclidean distance and cosine similarity objective
 121 functions proposed by recent prior work (Zhu et al., 2019b; Geiping et al., 2020) including a selection
 122 of popular regularisation functions.

123 **3 Experiments**

124 We conduct gradient inversion attack experiments on two benchmark datasets, MNIST Handwritten
 125 Digit (LeCun et al., 2010) and Labelled Faces in the Wild (LFW) (Huang et al., 2007), to illustrate
 126 how our proposed metric successfully evaluates information leakage that aligns to human perception.
 127 These two dataset are commonly used among researchers to study attacks (Zhu et al., 2019b; Zhao
 128 et al., 2020a; Melis et al., 2019; Shokri et al., 2017). We explore the privacy of the input data with
 129 the standard LeNET convolutional neural network (LeCun et al., 1990). We analyse the impact of
 130 different loss functions (MSE, LPIPS, SSIM, PSNR). For the attacks, in terms of the optimisation
 131 scheme, we utilized the standard LFBGS, with learning rate (lr) of 0.05, batch size of 4, and 300
 132 iterations for running a proxy model to attack.

133 We also carried out a complementary user study by asking 10 individuals for their feedback, to rank a
 134 series of inverse attack images from 0 – 5. With 0 being that they can very clearly observe underlying
 135 information (e.g. they can see the number 9 in a mnist recovered image, see example Fig. 1), to 5 that
 136 they cannot extract any useful information (e.g. the image is pure noise). We randomly generated 6
 137 groups, with 100 images each (total 600 images); specifically, three 10×10 frames of LFW images
 138 and three 10×10 frames of MNIST images (see Appendix, Fig. 17 and 18 for an example of LFW
 139 and MNIST). The images were not ranked by noise level/clarity, but they were randomly allocated
 140 into the 10×10 frame. After the individuals ranked them then we averaged the scores they gave for
 141 each of the 100 images and we compared it to the score each image achieved from MSE, AVD, and
 142 LPIPS metrics. We used Pearson correlation (ρ) as a measure of similarity. The results can be studied
 143 in the heatmap in Fig. 3.

	Noisy	Reference	Noisy	Reference	Noisy	Reference	Noisy	Reference	Reference	Noisy	Reference	
AVD		0.00		0.62		0.50		0.48		0.40		0.80
MSE		0.00		0.51		1.32		1.05		0.92		0.60
SSIM		0.99		0.09		0.13		0.12		0.18		0.12
PSNR		50.77		15.45		14.88		13.63		15.81		17.76
LPIPS		0.00		0.26		0.49		0.26		0.34		0.82

Figure 3: Comparison table of the most widely used metrics in FL for evaluation of inversion attacks (MSE, SSIM, PSNR, LPIPS), plus our own novel metric, the AVD. Each column compares the metrics between a noisy (attack generated) image and its reference (original) image. We included a wide array of examples, from no noise (column 1), to a mixture with two references (column 4), and complete noise (column 5).

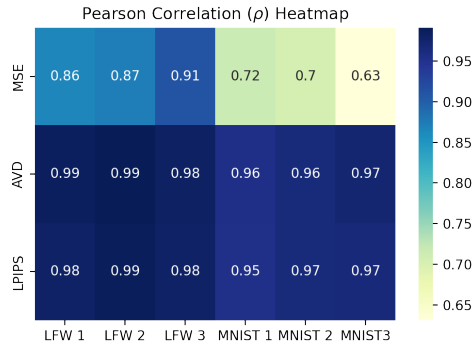


Figure 4: Pearson correlation ρ heatmap between the average ranking score (0-5) of 10 people for each image, and the MSE, LPIPS, AVD scores of these images. A ranking score of 0 means that the human can perceive very clear information in the image and a rank of 5 means the image is pure noise. The x-axis shows 6 groups (LFW1, LFW2, LFW3, MNIST1, MNIST2, MNIST3) of images. Each group has a 10×10 frame, 100 images in each frame. The y-axis shows the the three inversion attack metrics. The correlation is between the metrics (y-axis) and the average ranking score by the users. For more details refer to Section 3.

144 4 Results

145 For Fig. 1 we ran the experiment for the MNIST dataset. The first sub-figure (a) sorts the images
 146 by AVD value. For comparison, the second sub-figure (b) sorts the images by the MSE. The

147 smaller images represent the original MNIST data. Low score ranking represents that the attack
148 was successful and the image matches the original. High score ranking shows that the attack was
149 unsuccessful and the generated attack images are very noisy with no observable pattern. From Fig. 1b,
150 row one and columns four and five, the $MSE_{1,4} = 1.67$ and $MSE_{1,5} = 3.22$. These images clearly
151 have a pattern, a user might be able to infer that the number of the last image relates to the number
152 5. So the attacker can extract private information from an FL model. But, according to the MSE
153 these images are more private (noisy) when compared against images two and three from the same
154 row ($MSE_{1,2} = 0.97$ and $MSE_{1,3} = 1.01$). On the other hand, our metric AVD captures these
155 irregularities, with $AVD_{1,4} = 0.69$ and $AVD_{1,5} = 0.43$ being lower than $AVD_{1,2} = 0.88$ and
156 $AVD_{1,3} = 0.85$. The AVD scores also agree with the LPIPS benchmark in these images, which
157 indicates the AVD follows the human perception to evaluate the success of an inversion attack. In
158 the LFW dataset, Fig. 2, we can observe the same phenomenon when using the MSE as a score to
159 evaluate FL inversion attacks. In Fig. 3 we compare different inversion attack metrics, including
160 PSNR, SSIM, MSE, LPIPS, and AVD. The LPIPS and AVD results are consistent and agree very well
161 with human consensus; they attribute the lowest value ($LPIPS = 0$ & $AVD = 0$) to column one that
162 the two images are identical, and the largest value at column five ($LPIPS = 0.82$ & $AVD = 0.80$),
163 where the generated image is just noise.

164 In our final experiment, we contacted a qualitative survey amongst 10 people, Fig 3. When we
165 evaluate inversion attacks in multidimensional data that exhibit strong intercorrelation amongst the
166 data-points, such as images, then a contextual interpretation of the image is imperative for an accurate
167 evaluation. Therefore, the similarity metric should be able to showcase human-like perception. Our
168 survey results further support the quantitative analysis that we conducted in Fig. 3 and show that
169 the AVD is highly correlated ($\rho \geq 0.96$) with how a human would have recognised information
170 from an attack generated image. For the LFW group of images, the MSE had a correlation between
171 $0.86 \leq \rho \leq 0.91$ with the average human score. On the other hand, for the same images, the LPIPS
172 and the AVD achieved very high levels of correlation, $0.98 \leq \rho \leq 0.99$. For the MNIST group, the
173 MSE showed correlation between $0.63 \leq \rho \leq 0.72$. The AVD though retained consistently high
174 levels of correlation with the human score, $0.96 \leq \rho \leq 0.97$. It seems that the mixing of numbers and
175 the change of the background that we observed in the MNIST examples (Fig. 1 and Fig. 3) "confuse"
176 the MSE score and drive the results further away from human perception, reducing its accuracy.

177 5 Conclusion

178 In this paper, we have addressed a significant challenge in the field of FL - the evaluation of the
179 success of inversion attacks and the effectiveness of defense strategies. Traditional metrics such as
180 the SSIM, PSNR, and MSE have been shown to be insufficient for accurately assessing the success
181 of these attacks, particularly in the context of multidimensional outputs like images. These metrics,
182 which assume spatial independence, fail to consider the semantic context of the recovered data,
183 leading to potentially misleading evaluations.

184 To overcome these limitations, we introduced the AVD, a metric for assessing data recovery and
185 information leakage in FL. Derived from total variation, AVD offers a continuous measure for
186 extracting information in noisy images, aligning closely with human perception. It is computationally
187 more efficient and mathematically more interpretable than the LPIPS, a deep learning-based metric.
188 The quantitative experiments demonstrated that AVD provides a more accurate and consistent
189 measure of data recovery, thereby offering a more reliable tool for evaluating defense strategies
190 against inversion attacks in FL. Also, the survey that we contacted amongst 10 people, asking them
191 to rank random generated images by scoring the success of the recovered image, showed that human
192 perception had high correlation with the AVD scores.

193 By providing a more accurate measure of data recovery, AVD allows researchers to better understand
194 the effectiveness of their defense strategies and to develop robust FL evaluation of data security.
195 We hope that our work will inspire further advancements in the field of FL and contribute to the
196 development of more secure and reliable distributed learning systems.

197 References

198 Balunović, M., Dimitrov, D. I., Staab, R., and Vechev, M. Bayesian framework for gradient leakage,
199 2022.

- 200 Cahn, C. A note on signal-to-noise ratio in band-pass limiters. *IRE Transactions on Information*
201 *Theory*, 7(1):39–43, 1961. doi: 10.1109/TIT.1961.1057616.
- 202 Chen, Y., Gui, Y., Lin, H., Gan, W., and Wu, Y. Federated learning attacks and defenses: A survey,
203 2022.
- 204 Dwork, C. Differential privacy. In *33rd International Colloquium on Automata, Languages and*
205 *Programming, part II (ICALP 2006)*, volume 4052 of *Lecture Notes in Computer Science*, pp.
206 1–12. Springer Verlag, July 2006. ISBN 3-540-35907-9. URL <https://www.microsoft.com/en-us/research/publication/differential-privacy/>.
- 208 Eloul, S., Silavong, F., Kamthe, S., Georgiadis, A., and Moran, S. J. Enhancing privacy against
209 inversion attacks in federated learning by using mixing gradients strategies, 2022.
- 210 Geiping, J., Bauermeister, H., Dröge, H., and Moeller, M. Inverting Gradients - How easy is it to
211 break privacy in federated learning? In *Advances in Neural Information Processing Systems 33:*
212 *Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December*
213 *6-12, 2020, virtual*, 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/c4ede56bbd98819ae6112b20ac6bf145-Abstract.html>.
- 215 Huang, G. B., Ramesh, M., Berg, T., and Learned-Miller, E. Labeled faces in the wild: A database
216 for studying face recognition in unconstrained environments. Technical Report 07-49, University
217 of Massachusetts, Amherst, October 2007.
- 218 Huang, Y., Gupta, S., Song, Z., Li, K., and Arora, S. Evaluating gradient inversion attacks and
219 defenses in federated learning. In Beygelzimer, A., Dauphin, Y., Liang, P., and Vaughan, J. W.
220 (eds.), *Advances in Neural Information Processing Systems*, 2021a. URL <https://openreview.net/forum?id=0CDKgyYaxC8>.
- 222 Huang, Y., Song, Z., Li, K., and Arora, S. Instahide: Instance-hiding schemes for private distributed
223 learning, 2021b.
- 224 LeCun, Y., Boser, B., Denker, J., Henderson, D., Howard, R., Hubbard, W., and Jackel, L. Handwritten
225 Digit Recognition with a Back-Propagation Network. In *Advances in Neural Information Process-*
226 *ing Systems*, volume 2. Morgan-Kaufmann, 1990. URL <https://proceedings.neurips.cc/paper/1989/file/53c3bce66e43be4f209556518c2fcb54-Paper.pdf>.
- 228 LeCun, Y., Cortes, C., and Burges, C. Mnist handwritten digit database. *ATT Labs [Online]*. Available:
229 <http://yann.lecun.com/exdb/mnist>, 2, 2010.
- 230 McMahan, B., Moore, E., Ramage, D., Hampson, S., and Arcas, B. A. y. Communication-Efficient
231 Learning of Deep Networks from Decentralized Data. In *Proceedings of the 20th International*
232 *Conference on Artificial Intelligence and Statistics*, volume 54 of *Proceedings of Machine Learning*
233 *Research*, pp. 1273–1282. PMLR, 20–22 Apr 2017. URL <https://proceedings.mlr.press/v54/mcmahan17a.html>.
- 235 Melis, L., Song, C., Cristofaro, E. D., and Shmatikov, V. Exploiting unintended feature leakage
236 in collaborative learning. In *2019 IEEE Symposium on Security and Privacy, SP 2019, San*
237 *Francisco, CA, USA, May 19-23, 2019*, pp. 691–706. IEEE, 2019. doi: 10.1109/SP.2019.00029.
238 URL <https://doi.org/10.1109/SP.2019.00029>.
- 239 Nguyen, D. C., Pham, Q.-V., Pathirana, P. N., Ding, M., Seneviratne, A., Lin, Z., Dobre, O., and
240 Hwang, W.-J. Federated learning for smart healthcare: A survey. *ACM Computing Surveys (CSUR)*,
241 55(3):1–37, 2022.
- 242 Phong, L. T., Aono, Y., Hayashi, T., Wang, L., and Moriai, S. Privacy-preserving deep learning via
243 additively homomorphic encryption. *IEEE Transactions on Information Forensics and Security*,
244 13(5):1333–1345, 2018. doi: 10.1109/TIFS.2017.2787987.
- 245 Qian, J., Nassar, H., and Hansen, L. K. Minimal model structure analysis for input reconstruction in
246 federated learning, 2021.

- 247 Rieke, N., Hancox, J., Li, W., Milletari, F., Roth, H., Albarqouni, S., Bakas, S., Galtier, M., Landman,
248 B., Maier-Hein, K., Ourselin, S., Sheller, M., Summers, R., Trask, A., Xu, D., Baust, M., and
249 Cardoso, M. The future of Digital Health with Federated Learning. *npj Digital Medicine*, 3(1),
250 December 2020. ISSN 2398-6352. doi: 10.1038/s41746-020-00323-1.
- 251 Rudin, L. I., Osher, S., and Fatemi, E. Nonlinear total variation based noise removal algorithms. In
252 *Proceedings of the Eleventh Annual International Conference of the Center for Nonlinear Studies*
253 *on Experimental Mathematics: Computational Issues in Nonlinear Science: Computational Issues*
254 *in Nonlinear Science*, pp. 259–268, USA, 1992. Elsevier North-Holland, Inc.
- 255 Shokri, R., Stronati, M., Song, C., and Shmatikov, V. Membership inference attacks against machine
256 learning models. In *2017 IEEE Symposium on Security and Privacy, SP 2017, San Jose, CA,*
257 *USA, May 22-26, 2017*, pp. 3–18. IEEE Computer Society, 2017. doi: 10.1109/SP.2017.41. URL
258 <https://doi.org/10.1109/SP.2017.41>.
- 259 Sikandar, H. S., Waheed, H., Tahir, S., Malik, S. U. R., and Rafique, W. A detailed survey on
260 federated learning attacks and defenses. *Electronics*, 12(2), 2023. ISSN 2079-9292. doi: 10.3390/
261 electronics12020260. URL <https://www.mdpi.com/2079-9292/12/2/260>.
- 262 Wainakh, A., Zimmer, E., Subedi, S., Keim, J., Grube, T., Karuppayah, S., Guinea, A. S., and
263 Mühlhäuser, M. Federated learning attacks revisited: A critical discussion of gaps, assumptions,
264 and evaluation setups, 2022.
- 265 Wang, Z., Bovik, A., Sheikh, H., and Simoncelli, E. Image quality assessment: from error visibility
266 to structural similarity. *IEEE Transactions on Image Processing*, 13(4):600–612, 2004. doi:
267 10.1109/TIP.2003.819861.
- 268 Yang, Q., Liu, Y., Chen, T., and Tong, Y. Federated Machine Learning: Concept and Applications.
269 *ACM Trans. Intell. Syst. Technol.*, 10(2), jan 2019. ISSN 2157-6904. doi: 10.1145/3298981. URL
270 <https://doi.org/10.1145/3298981>.
- 271 Yin, H., Mallya, A., Vahdat, A., Alvarez, J., Kautz, J., and Molchanov, P. See through Gradients:
272 Image Batch Recovery via GradInversion. In *2021 IEEE/CVF Conference on Computer Vision*
273 *and Pattern Recognition (CVPR)*, pp. 16332–16341, 2021. doi: 10.1109/CVPR46437.2021.01607.
- 274 Zhang, R., Isola, P., Efros, A. A., Shechtman, E., and Wang, O. The unreasonable effectiveness of
275 deep features as a perceptual metric, 2018.
- 276 Zhao, B., Mopuri, K. R., and Bilen, H. idlg: Improved deep leakage from gradients, 2020a.
- 277 Zhao, Y., Zhao, J., Yang, M., Wang, T., Wang, N., Lyu, L., Niyato, D., and Lam, K.-Y. Local
278 differential privacy based federated learning for internet of things, 2020b.
- 279 Zhu, L., Liu, Z., and Han, S. Deep leakage from gradients. In Wallach, H., Larochelle, H., Beygelz-
280 imer, A., d'Alché-Buc, F., Fox, E., and Garnett, R. (eds.), *Advances in Neural Information Process-*
281 *ing Systems*, volume 32. Curran Associates, Inc., 2019a. URL https://proceedings.neurips.cc/paper_files/paper/2019/file/60a6c4002cc7b29142def8871531281a-Paper.pdf.
- 282
283 Zhu, L., Liu, Z., and Han, S. Deep leakage from gradients. In *Advances in Neural Information*
284 *Processing Systems*, volume 32. Curran Associates, Inc., 2019b. URL <https://proceedings.neurips.cc/paper/2019/file/60a6c4002cc7b29142def8871531281a-Paper.pdf>.
- 285

286 Appendix

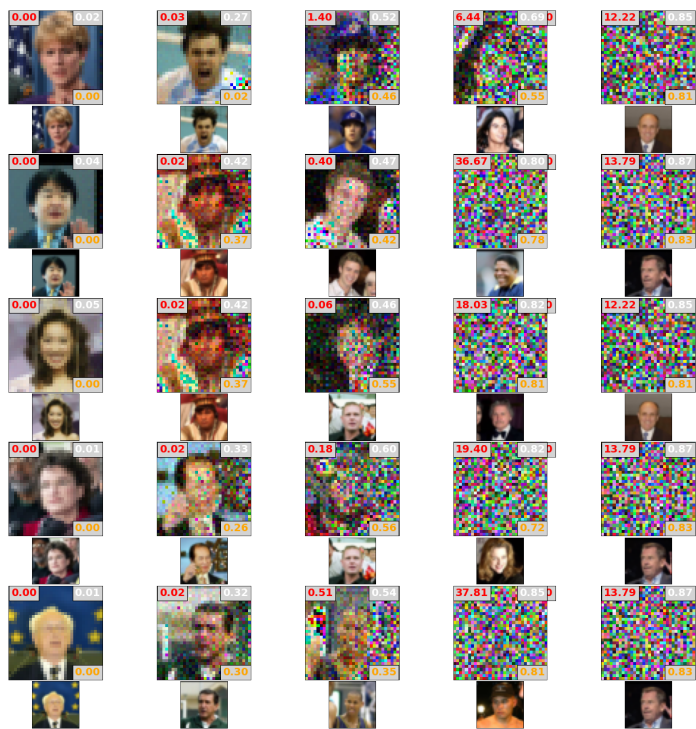


Figure 5: Random recovered vectors from LFW datasets, column-wise sorted via the AVD.



Figure 6: Random recovered vectors from LFW datasets, column-wise sorted via the AVD.

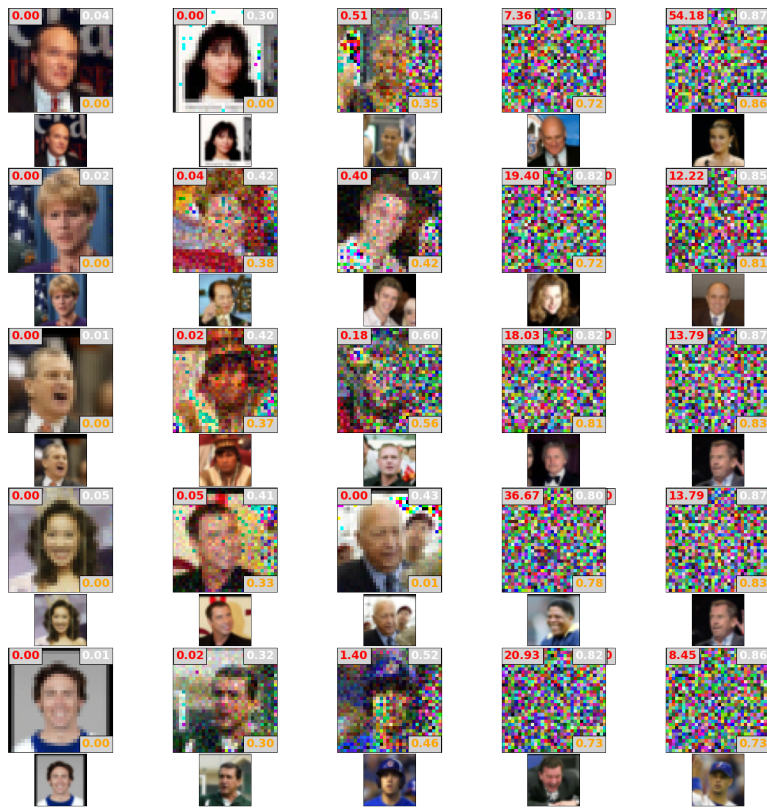


Figure 7: Random recovered vectors from LFW datasets, column-wise sorted via the AVD.



Figure 8: Random recovered vectors from LFW datasets, column-wise sorted via the MSE.

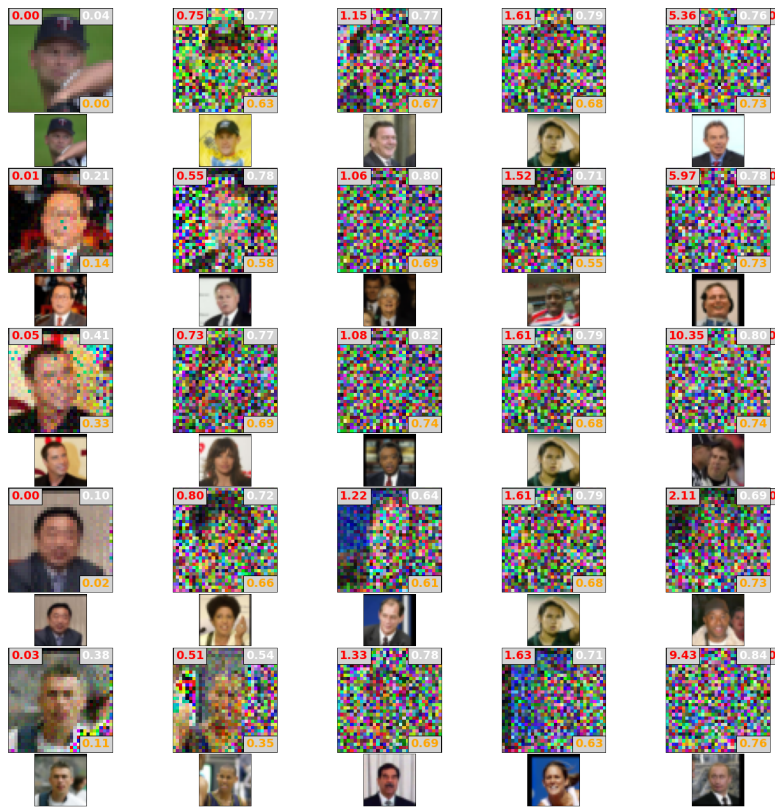


Figure 9: Random recovered vectors from LFW datasets, column-wise sorted via the MSE.

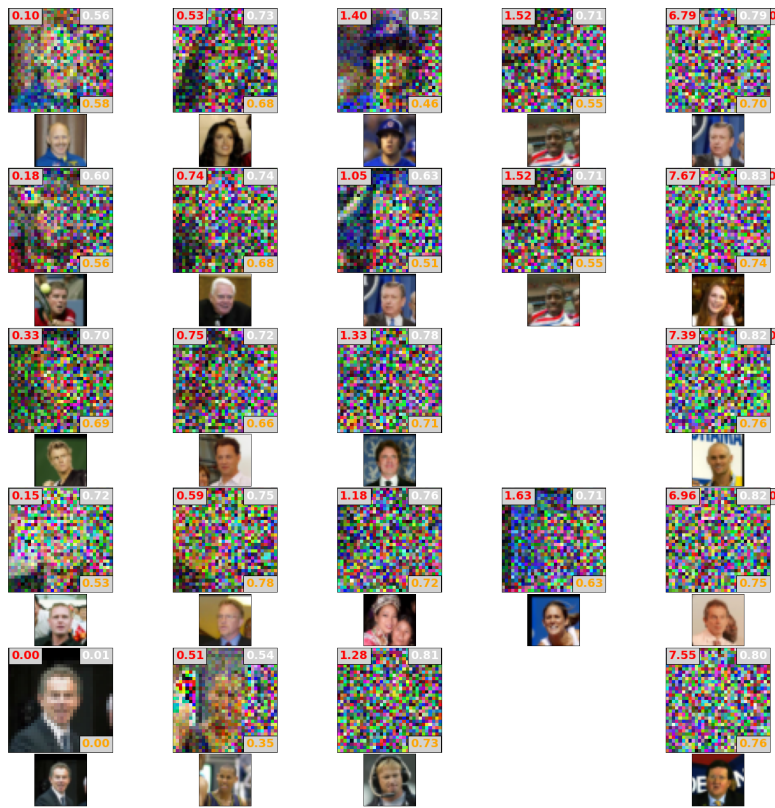


Figure 10: Random recovered vectors from LFW datasets, column-wise sorted via the MSE.

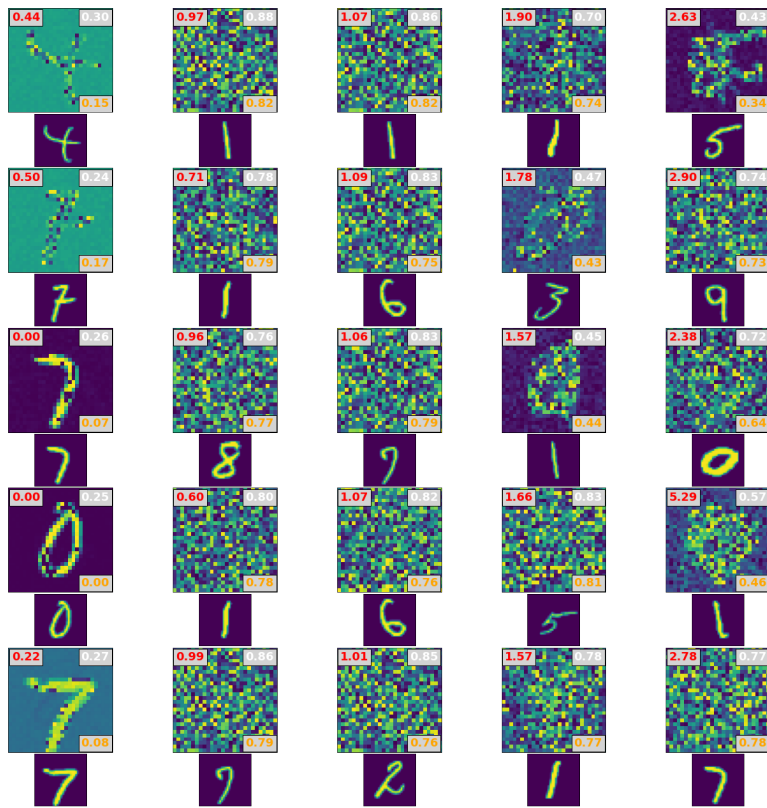


Figure 11: Random recovered vectors from MNIST datasets, column-wise sorted via the MSE.

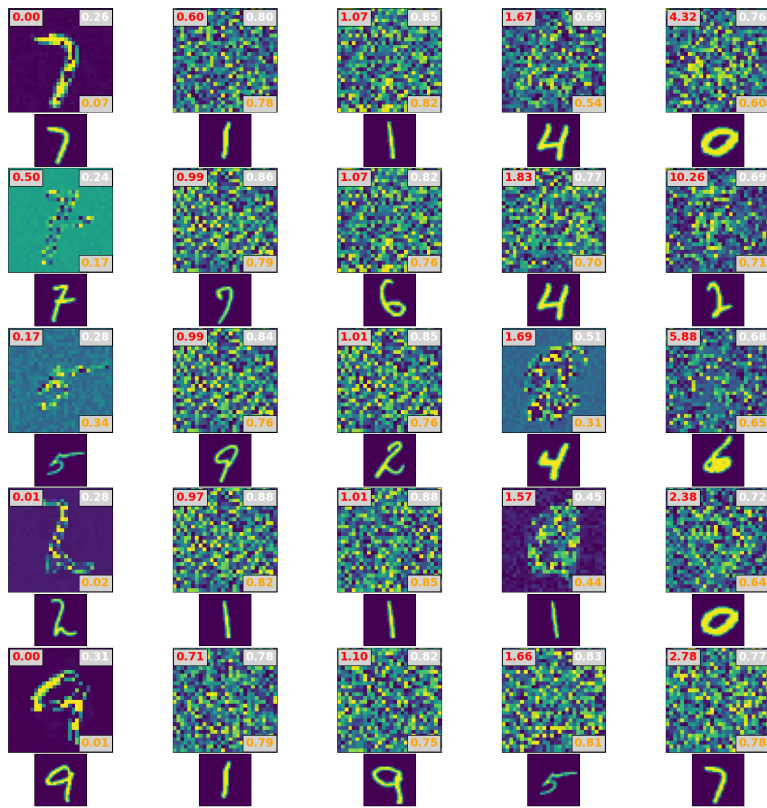


Figure 12: Random recovered vectors from MNIST datasets, column-wise sorted via the MSE.

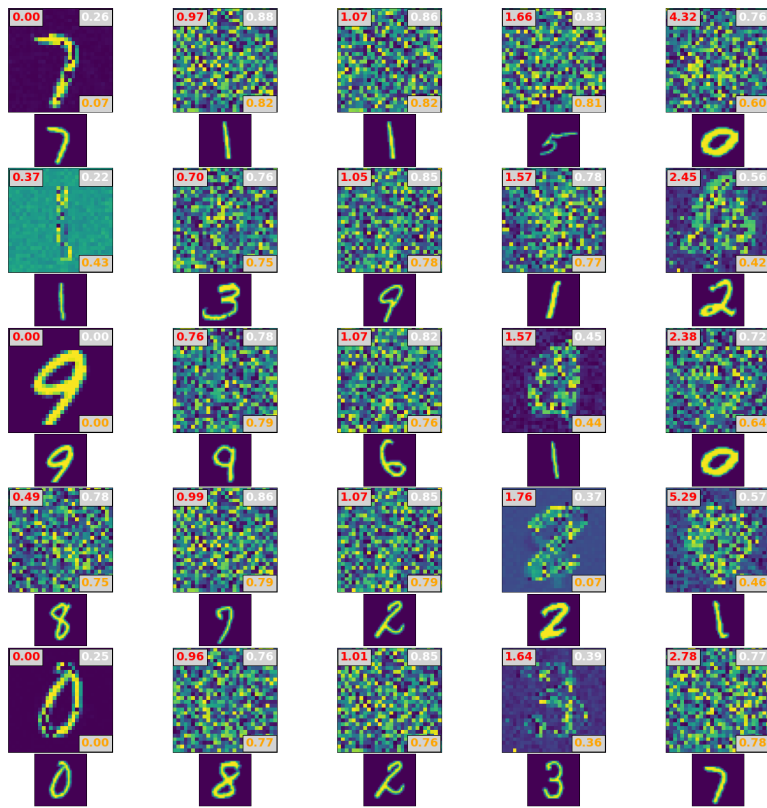


Figure 13: Random recovered vectors from MNIST datasets, column-wise sorted via the MSE.

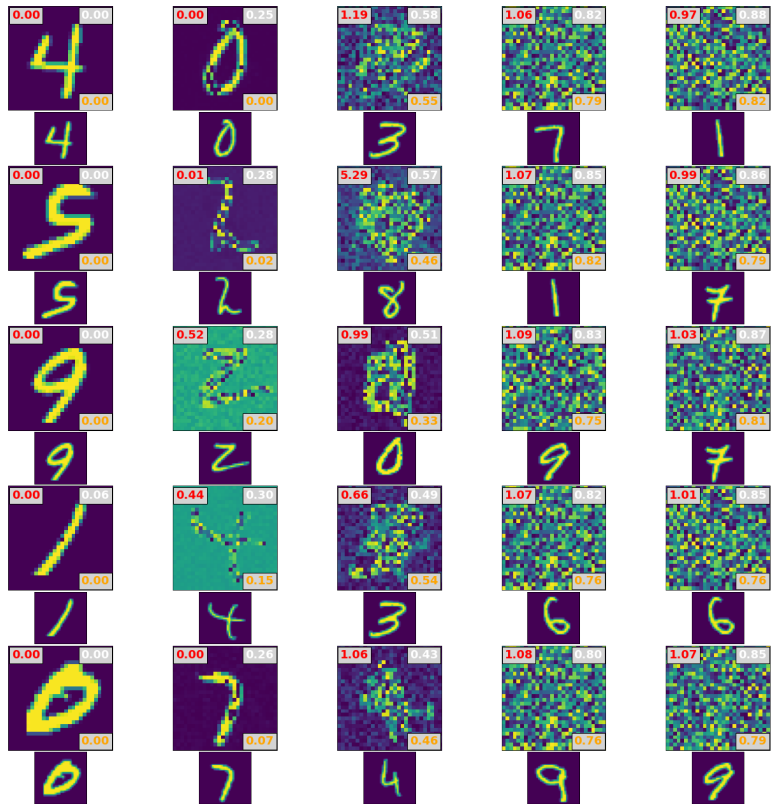


Figure 14: Random recovered vectors from MNIST datasets, column-wise sorted via the AVD.

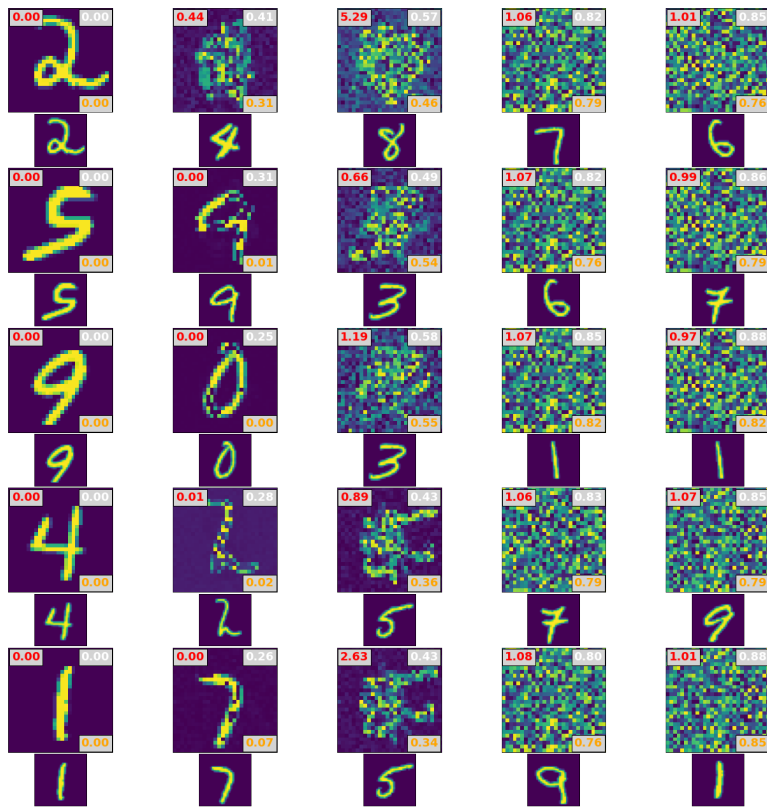


Figure 15: Random recovered vectors from MNIST datasets, column-wise sorted via the AVD.

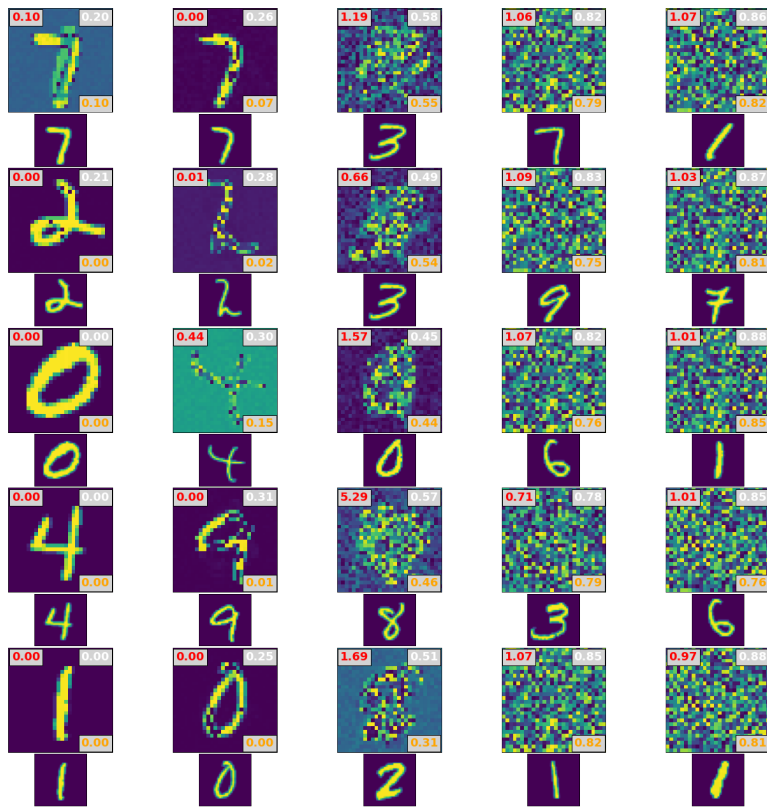


Figure 16: Random recovered vectors from MNIST datasets, column-wise sorted via the AVD.

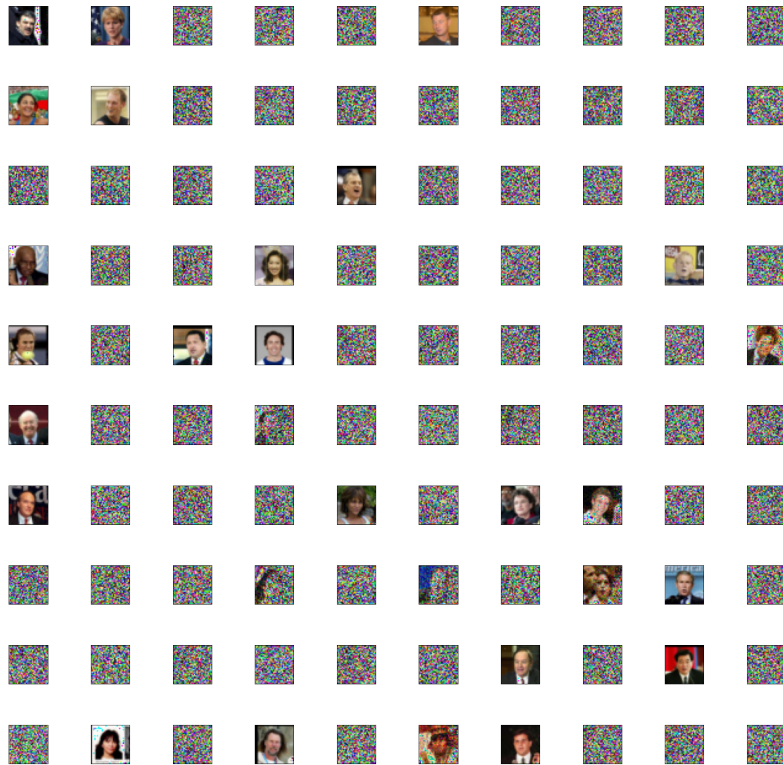


Figure 17: Random recovered vectors from MNIST datasets, column-wise sorted via the AVD.

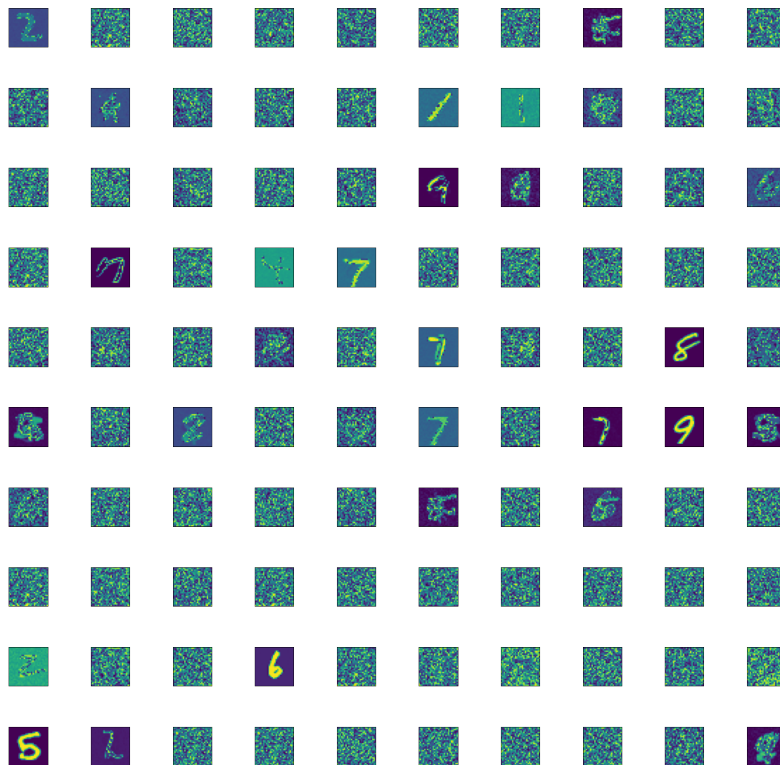


Figure 18: Random recovered vectors from MNIST datasets, column-wise sorted via the AVD.