



HAL
open science

Random Matrix Theory Proves that Deep Learning Representations of GAN-data Behave as Gaussian Mixtures

Mohamed El Amine Seddik, Cosme Louart, Mohamed Tamaazousti, Romain Couillet

► **To cite this version:**

Mohamed El Amine Seddik, Cosme Louart, Mohamed Tamaazousti, Romain Couillet. Random Matrix Theory Proves that Deep Learning Representations of GAN-data Behave as Gaussian Mixtures. ICML 2020: Thirty-seventh International Conference on Machine Learning, 2020, Online, France. hal-02971185

HAL Id: hal-02971185

<https://hal.science/hal-02971185>

Submitted on 19 Oct 2020

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Random Matrix Theory Proves that Deep Learning Representations of GAN-data Behave as Gaussian Mixtures

Mohamed El Amine Seddik^{1 2} Cosme Louart^{1 3} Mohamed Tamaazousti¹ Romain Couillet^{2 3}

Abstract

This paper shows that deep learning (DL) representations of data produced by generative adversarial nets (GANs) are random vectors which fall within the class of so-called *concentrated* random vectors. Further exploiting the fact that Gram matrices, of the type $\mathbf{G} = \mathbf{X}^\top \mathbf{X}$ with $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{p \times n}$ and \mathbf{x}_i independent concentrated random vectors from a mixture model, behave asymptotically (as $n, p \rightarrow \infty$) as if the \mathbf{x}_i were drawn from a Gaussian mixture, suggests that DL representations of GAN-data can be fully described by their first two statistical moments for a wide range of standard classifiers. Our theoretical findings are validated by generating images with the BigGAN model and across different popular deep representation networks.

1. Introduction

The performance of machine learning methods depends strongly on the choice of the data representation (or features) on which they are applied. This data representation should ideally contain *relevant information* about the learning task in order to achieve learning with *simple* models and *small* amount of samples. Deep neural networks (Rumelhart et al., 1988) have particularly shown impressive results by automatically learning representations from raw data (e.g., images). However, due to the complex structure of deep learning models, the characterization of their hidden representations is still an open problem (Bengio et al., 2009).

Specifically, quantifying what makes a given deep learning representation better than another is a fundamental question in the field of *Representation Learning* (Bengio et al., 2013). Relying on (Montavon et al., 2011) a data representation is said to be *good* when it is possible to build *simple* models

on top of it that are *accurate* for the given learning problem. (Montavon et al., 2011) have notably quantified the layer-wise evolution of the representation in deep networks by computing the principal components of the Gram matrix $\mathbf{G}_\ell = \{\phi_\ell(\mathbf{x}_i)^\top \phi_\ell(\mathbf{x}_j)\}_{i,j=1}^n$ at each layer for n input data $\mathbf{x}_1, \dots, \mathbf{x}_n$, where $\phi_\ell(\mathbf{x})$ is the representation of \mathbf{x} at layer ℓ of the given DL model, and the number of components controls the model simplicity. In their study, the impact of the representation at each layer is quantified through the prediction error of a linear predictor trained on the principal subspace of \mathbf{G}_ℓ .

Pursuing on this idea, given a certain representation model $\mathbf{x} \mapsto \phi(\mathbf{x})$, we aim in this article at theoretically studying the large dimensional behavior, and in particular the spectral information (i.e., eigenvalues and dominant eigenvectors), of the corresponding Gram matrix $\mathbf{G} = \{\phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)\}_{i,j=1}^n$ in order to determine the information encoded (i.e., the sufficient statistics) by the representation model on a set of real data $\mathbf{x}_1, \dots, \mathbf{x}_n$. Indeed, standard classification and regression algorithms –along with the last layer of a neural network (Yeh et al., 2018)– retrieve the data information directly from functionals or the eigenspectrum of \mathbf{G}^1 . To this end, though, one needs a statistical model for the representations given the distribution of the raw data (e.g., images) which is generally unknown. Yet, due to recent advances in generative models since the advent of Generative Adversarial Nets (Goodfellow et al., 2014), it is now possible to generate complex data structures by applying successive *Lipschitz* operations to Gaussian random vectors. In particular, GAN-data *are* used in practice as substitutes of real data for data augmentation (Antoniou et al., 2017). On the other hand, the fundamental concentration of measure phenomenon (Ledoux, 2005) tells us that Lipschitz-ally transformed Gaussian vectors satisfy a concentration property. Precisely, defining the class of *concentrated* vectors $\mathbf{x} \in E$ through concentration inequalities of $f(\mathbf{x})$, for any real Lipschitz observation $f : E \rightarrow \mathbb{R}$, implies that deep learning representations of GAN-data fall within this class of random vectors, since the mapping $\mathbf{x} \mapsto \phi(\mathbf{x})$ is Lipschitz. Thus, GAN-data are concentrated random vec-

¹CEA List, Palaiseau, France ²Centralesupélec, Gif-sur-Yvette, France ³University Grenoble-Alpes, Grenoble, France. Correspondence to: MEA.Seddik <mohamedelamine.seddik@cea.fr>.

¹For instance, spectral clustering uses the dominant eigenvectors of \mathbf{G} , while support vector machines use functionals (quadratic forms) involving \mathbf{G} .

tors and thus an appropriate statistical model of realistic data.

Targeting classification applications by assuming a mixture of concentrated random vectors model, this article studies the spectral behavior of Gram matrices \mathbf{G} in the large n, p regime. Precisely, we show that these matrices have asymptotically (as $n, p \rightarrow \infty$ with $p/n \rightarrow c < \infty$) the same first-order behavior as for a Gaussian Mixture Model (GMM). As a result, by generating images using the BigGAN model (Brock et al., 2018) and considering different commonly used deep representation models, we show that the spectral behavior of the Gram matrix computed on these representations is the same as on a GMM model with the same p -dimensional means and covariances. A surprising consequence is that, for GAN data, the aforementioned *sufficient statistics* to characterize the quality of a given representation network are only the *first* and *second* order statistics of the representations. This behavior is shown by simulations to extend beyond random GAN-data to real images from the Imagenet dataset (Deng et al., 2009).

The rest of the paper is organized as follows. In Section 2, we introduce the notion of concentrated vectors and their main properties. Our main theoretical results are then provided in Section 3. In Section 4 we present experimental results. Section 5 concludes the article.

Notation: In the following, we use the notation from (Goodfellow et al., 2016). $[n]$ denotes the set $\{1, \dots, n\}$. Given a vector $\mathbf{x} \in \mathbb{R}^n$, the ℓ_2 -norm of \mathbf{x} is given as $\|\mathbf{x}\|^2 = \sum_{i=1}^n x_i^2$. Given a $p \times n$ matrix \mathbf{M} , its Frobenius norm is defined as $\|\mathbf{M}\|_F^2 = \sum_{i=1}^p \sum_{j=1}^n M_{ij}^2$ and its spectral norm as $\|\mathbf{M}\| = \sup_{\|\mathbf{x}\|=1} \|\mathbf{M}\mathbf{x}\|$. \odot for the Hadamard product. An application $\mathcal{F} : E \rightarrow F$ is said to be $\|\mathcal{F}\|_{lip}$ -Lipschitz, if $\forall (\mathbf{x}, \mathbf{y}) \in E^2$, $\|\mathcal{F}(\mathbf{x}) - \mathcal{F}(\mathbf{y})\|_F \leq \|\mathcal{F}\|_{lip} \cdot \|\mathbf{x} - \mathbf{y}\|_E$ and $\|\mathcal{F}\|_{lip}$ is finite.

2. Basic notions of concentrated vectors

Being the central tool of our study, we start by introducing the notion of concentrated vectors. While advanced concentration notions have been recently developed in (Louart & Couillet, 2019) in order to specifically analyze the behavior of large dimensional sample covariance matrices, for simplicity, we restrict ourselves here to the sufficient so-called q -exponentially concentrated random vectors.

Definition 2.1 (q -exponential concentration). *Given a normed space $(E, \|\cdot\|_E)$ and a real q , a random vector $\mathbf{x} \in E$ is said to be q -exponentially concentrated if for any 1-Lipschitz real function $f : E \rightarrow \mathbb{R}$, there exists $C \geq 0$ independent of $\dim(E)$ and $\sigma > 0$ such that for all $t \geq 0$*

$$\mathbb{P}\{|f(\mathbf{x}) - \mathbb{E}f(\mathbf{x})| > t\} \leq C e^{-(t/\sigma)^q} \quad (1)$$

which we denote $\mathbf{x} \in \mathcal{E}_q(\sigma | E, \|\cdot\|_E)$. We simply write

$\mathbf{x} \in \mathcal{E}_q(1 | E, \|\cdot\|_E)$ if the tail parameter σ does not depend on $\dim(E)$, and $x \in \mathcal{E}_q(1)$ for x a scalar real random variable.

Therefore, concentrated vectors are defined through the concentration of any 1-Lipschitz real scalar ‘‘observation’’. One of the most important examples of concentrated vectors are standard Gaussian vectors. Precisely, we have the following proposition. See (Ledoux, 2005) for more examples such as uniform and Gamma distribution.

Proposition 2.2 (Gaussian vectors (Ledoux, 2005)). *Let $d \in \mathbb{N}$ and $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_d)$. Then \mathbf{x} is a 2-exponentially concentrated vector independently on the dimension d , i.e. $\mathbf{x} \in \mathcal{E}_2(1 | \mathbb{R}^d, \|\cdot\|)$.*

Concentrated vectors have the interesting property of being stable by application of $\mathbb{R}^d \rightarrow \mathbb{R}^p$ vector-Lipschitz transformations. Indeed, Lipschitz-ally transformed concentrated vectors remain concentrated according to the following proposition.

Proposition 2.3 (Lipschitz stability). *Let $\mathbf{x} \in \mathcal{E}_q(1 | E, \|\cdot\|_E)$ and $\mathcal{G} : E \rightarrow F$ a Lipschitz application with Lipschitz constant $\|\mathcal{G}\|_{lip}$ which may depend on $\dim(F)$. Then the concentration property on \mathbf{x} is transferred to $\mathcal{G}(\mathbf{x})$, if $\mathbf{x} \in \mathcal{E}_q(1 | E, \|\cdot\|_E)$ then*

$$\mathcal{G}(\mathbf{x}) \in \mathcal{E}_q(\|\mathcal{G}\|_{lip} | F, \|\cdot\|_F). \quad (2)$$

Note importantly for the following that the Lipschitz constant of the transformation \mathcal{G} must be controlled, in order to constrain the tail parameter of the obtained concentration.

In particular, we have the coming corollary to Proposition 2.3 of central importance in the following.

Corollary 2.4. *Let $\mathcal{G}_1, \dots, \mathcal{G}_n : \mathbb{R}^d \rightarrow \mathbb{R}^p$ a set of n Lipschitz applications with Lipschitz constants $\|\mathcal{G}_i\|_{lip}$. Let $\mathcal{G} : \mathbb{R}^{d \times n} \rightarrow \mathbb{R}^{p \times n}$ be defined for each $\mathbf{X} \in \mathbb{R}^{d \times n}$ as $\mathcal{G}(\mathbf{X}) = [\mathcal{G}_1(\mathbf{X}_{:,1}), \dots, \mathcal{G}_n(\mathbf{X}_{:,n})]$. Then,*

$$\begin{aligned} \mathbf{Z} &\in \mathcal{E}_q(1 | \mathbb{R}^{d \times n}, \|\cdot\|_F) \\ \Rightarrow \mathcal{G}(\mathbf{Z}) &\in \mathcal{E}_q\left(\sup_i \|\mathcal{G}_i\|_{lip} | \mathbb{R}^{p \times n}, \|\cdot\|_F\right). \end{aligned} \quad (3)$$

Proof. This is a consequence of Proposition 2.3 since the map \mathcal{G} is $\sup_i \|\mathcal{G}_i\|_{lip}$ -Lipschitz with respect to (w.r.t.) the Frobenius norm. Indeed, for $\mathbf{X}, \mathbf{H} \in \mathbb{R}^{d \times n} : \|\mathcal{G}(\mathbf{X} + \mathbf{H}) - \mathcal{G}(\mathbf{X})\|_F^2 \leq \sum_{i=1}^n \|\mathcal{G}_i\|_{lip}^2 \cdot \|\mathbf{H}_{:,i}\|^2 \leq \sup_i \|\mathcal{G}_i\|_{lip}^2 \cdot \|\mathbf{H}\|_F^2$. \square

3. Main Results

3.1. GAN data: An Example of Concentrated Vectors

Concentrated random vectors are particularly interesting from a practical standpoint for real data modeling. In fact,

unlike simple Gaussian vectors, the former do not suffer from the constraint of having independent entries which is quite a restrictive assumption when modeling real data such as images or their non-linear features (*e.g.*, DL representations). The other modeling interest of concentrated vectors lies in their being already present in practice as alternatives to real data. Indeed, adversarial neural networks (GANs) have the ability nowadays to generate random *realistic* data (for instance realistic images) by applying successive Lipschitz operations to standard Gaussian vectors (Goodfellow et al., 2014).

A GAN architecture involves two networks, a generator model which maps random Gaussian noise to new plausible synthetic data and a discriminator model which classifies real data as real (from the dataset) or fake (for the generated data). The discriminator is updated directly through a binary classification problem, whereas the generator is updated through the discriminator. As such, the two models are trained alternatively in an adversarial manner, where the generator seeks to better deceive the discriminator and the former seeks to better identify the fake data (Goodfellow et al., 2014).

In particular, once both models are trained (when they reach a Nash equilibrium), DL representations of GAN-data –and GAN-data themselves– are schematically constructed in practice as follows:

$$\text{Real Data} \approx \text{GAN Data} = \mathcal{F}_N \circ \dots \circ \mathcal{F}_1(z), \quad (4)$$

where $z \sim \mathcal{N}(0, \mathbf{I}_d)$, d stands for the input dimension of the generator model, N the number of layers, and the \mathcal{F}_i 's either Fully Connected Layers, Convolutional Layers, Pooling Layers, Up-sampling Layers and Activation Functions,

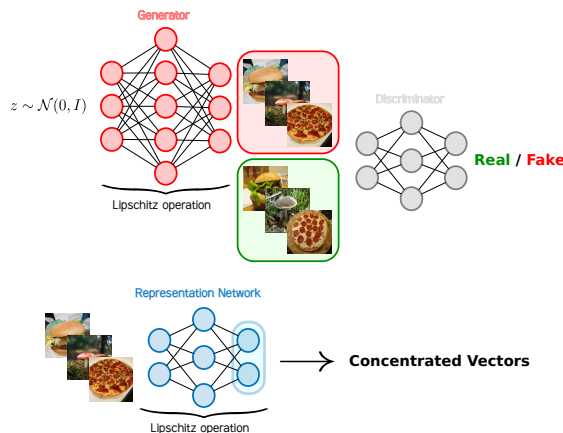


Figure 1. Deep learning representations of GAN-data are constructed by applying successive Lipschitz operations to Gaussian vectors, therefore they are *concentrated* vectors by design, since Gaussian vectors are concentrated and thanks to the Lipschitz stability in Proposition 2.3.

Residual Layers or Batch Normalizations. All these operations happen to be *Lipschitz* applications. Precisely,

- **Fully Connected Layers and Convolutional Layers:** These are affine operations which can be expressed as

$$\mathcal{F}_i(\mathbf{x}) = \mathbf{W}_i \mathbf{x} + \mathbf{b}_i,$$

for \mathbf{W}_i the weight matrix and \mathbf{b}_i the bias vector. Here the Lipschitz constant is the operator norm (the largest singular value) of the weight matrix \mathbf{W}_i , that is $\|\mathcal{F}_i\|_{lip} = \sup_{\mathbf{u} \neq 0} \frac{\|\mathbf{W}_i \mathbf{u}\|_2}{\|\mathbf{u}\|_2}$.

- **Pooling Layers and Activation Functions:** Most commonly used activation functions and pooling operations are

$$\text{ReLU}(\mathbf{x}) = \max(0, \mathbf{x}),$$

$$\text{MaxPooling}(\mathbf{x}) = [\max(\mathbf{x}_{S_1}), \dots, \max(\mathbf{x}_{S_q})]^\top,$$

where S_i 's are patches (*i.e.*, subsets of $[\dim(\mathbf{x})]$). These are at most 1-Lipschitz operations with respect to the Frobenius norm. Specifically, the maximum absolute sub-gradient of the ReLU activation function is 1, thus the ReLU operation has a Lipschitz constant of 1. Similarly, we can show that the Lipschitz constant of MaxPooling layers is also 1.

- **Residual Connections:** Residual layers act the following way

$$\mathcal{F}_i(\mathbf{x}) = \mathbf{x} + \mathcal{F}_i^{(1)} \circ \dots \circ \mathcal{F}_i^{(\ell)}(\mathbf{x}),$$

where the $\mathcal{F}_i^{(j)}$'s are Fully Connected Layers or Convolutional Layers with Activation Functions, and which are Lipschitz operations. Thus \mathcal{F}_i is a Lipschitz operation with Lipschitz constant bounded by $1 + \prod_{j=1}^{\ell} \|\mathcal{F}_i^{(j)}\|_{lip}$.

- **Batch Normalization (BN) Layers:** They consist in statistically standardizing (Ioffe & Szegedy, 2015) the vectors of a small batch $\mathcal{B} = \{\mathbf{x}_i\}_{i=1}^b \subset \mathbb{R}^d$ as follows: for each $\mathbf{x}_k \in \mathcal{B}$

$$\mathcal{F}_i(\mathbf{x}_k) = \text{diag} \left(\frac{\mathbf{a}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \right) (\mathbf{x}_k - \mu_{\mathcal{B}} \mathbf{1}_d) + \mathbf{b}$$

where $\mu_{\mathcal{B}} = \frac{1}{db} \sum_{k=1}^b \sum_{i=1}^d [\mathbf{x}_k]_i$, $\sigma_{\mathcal{B}}^2 = \frac{1}{db} \sum_{k=1}^b \sum_{i=1}^d ([\mathbf{x}_k]_i - \mu_{\mathcal{B}})^2$, $\mathbf{a}, \mathbf{b} \in \mathbb{R}^d$ are parameters to be learned and $\text{diag}(\mathbf{v})$ transforms a vector \mathbf{v} to a diagonal matrix with its diagonal entries being those of \mathbf{v} . Thus BN is a Lipschitz transformation with Lipschitz constant $\|\mathcal{F}_i\|_{lip} = \sup_i \left| \frac{\mathbf{a}_i}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \right|$.

Therefore, as illustrated in Figure 1, since standard Gaussian vectors are concentrated vectors as mentioned in Proposition 2.2 and since the notion of concentrated vectors is stable by Lipschitz transformations thanks to Proposition 2.3, GAN-data (and their DL representations) are concentrated vectors by design given the construction in Equation (4).

Moreover, in order to generate data belonging to a specific class, Conditional GANs have been introduced (Mirza & Osindero, 2014); once again data generated by these models are concentrated vectors as a consequence of Corollary 2.4. Indeed, a generator of a Conditional GAN model can be seen as a set of multiple generators where each generates data of a specific class conditionally on the class label (e.g., BigGAN model (Brock et al., 2018)).

Yet, in order to ensure that the resulting Lipschitz constant of the combination of the above operations does not scale with the network or data size, so to maintain good concentration behaviors, a careful control of the learned network parameters is needed. This control happens to be already considered in practice in order to ensure the stability of GANs during the learning phase, notably to generate realistic and high-resolution images (Roth et al., 2017; Brock et al., 2018). The control of the Lipschitz constant of representation networks is also needed in practice in order to make them robust against adversarial examples (Szegedy et al., 2013; Gulrajani et al., 2017). This control is particularly ensured through spectral normalization of the affine layers (Brock et al., 2018), such as Fully Connected Layers, Convolutional Layers and Batch Normalization. Indeed, spectral normalization (Miyato et al., 2018) consists in applying the operation $\mathbf{W} \leftarrow \mathbf{W}/\sigma_1(\mathbf{W})$ to the affine layers at each backward iteration of the back-propagation algorithm, where $\sigma_1(\mathbf{W})$ stands for the largest singular value of the weight matrix \mathbf{W} . (Brock et al., 2018), have notably observed that, without spectral constraints, a subset of the generator layers grow throughout their GAN training and explode at collapse. They thus suggested the following spectral normalization—which happens to be less restrictive than the standard spectral normalization $\mathbf{W} \leftarrow \mathbf{W}/\sigma_1(\mathbf{W})$ (Miyato et al., 2018)—to the affine layers:

$$\mathbf{W} \leftarrow \mathbf{W} - (\sigma_1(\mathbf{W}) - \sigma_*) \mathbf{u}_1(\mathbf{W})\mathbf{v}_1(\mathbf{W})^\top \quad (5)$$

where $\mathbf{u}_1(\mathbf{W})$ and $\mathbf{v}_1(\mathbf{W})$ denote respectively the left and right largest singular vectors of \mathbf{W} , and σ_* is an hyperparameter fixed during training.

To get an insight about the influence of this operation and to ensure that it controls the Lipschitz constant of the generator, the following proposition provides the dynamics of a random walk in the space of parameters along with the spectral normalization in Equation (5). Indeed, since stochastic gradient descent (SGD) consists in estimating the gradient of the loss function on randomly selected batches of data, it can be assimilated to a random walk in the space of parameters (Antognini & Sohl-Dickstein, 2018).

Proposition 3.1 (Lipschitz constant control). *Let $\sigma_* > 0$ and \mathcal{G} be a neural network composed of N affine layers, each one of input dimension d_{i-1} and output dimension d_i for $i \in [N]$, with 1-Lipschitz activation functions. Assume that the weights of \mathcal{G} at layer $i + 1$ are initialized as*

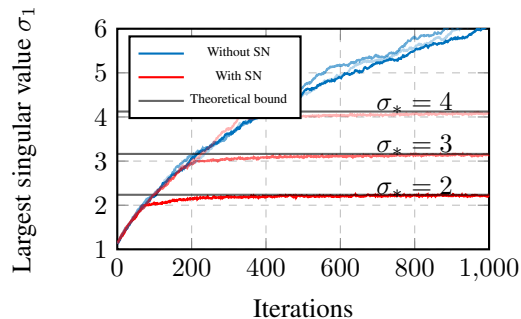


Figure 2. Behavior of the largest singular value of a weight matrix in terms of the iterations of a random walk (see proposition 3.1), without spectral normalization in (blue) and with spectral normalization in (red). The (black) lines correspond to the theoretical bound $\sqrt{\sigma_*^2 + \eta^2 d_1 d_0}$ for different σ_* 's. We took $d_0 = d_1 = 100$ and $\eta = 1/d_0$.

$\mathcal{U}([-1/\sqrt{d_i}, 1/\sqrt{d_i}])$, and consider the following dynamics with learning rate η :

$$\begin{aligned} \mathbf{W} &\leftarrow \mathbf{W} - \eta \mathbf{E}, \text{ with } \mathbf{E}_{i,j} \sim \mathcal{N}(0, 1) \\ \mathbf{W} &\leftarrow \mathbf{W} - \max(0, \sigma_1(\mathbf{W}) - \sigma_*) \mathbf{u}_1(\mathbf{W})\mathbf{v}_1(\mathbf{W})^\top. \end{aligned} \quad (6)$$

Then, $\forall \varepsilon > 0$, the Lipschitz constant of \mathcal{G} is bounded at convergence with high probability as:

$$\|\mathcal{G}\|_{lip} \leq \prod_{i=1}^N \left(\varepsilon + \sqrt{\sigma_*^2 + \eta^2 d_i d_{i-1}} \right). \quad (7)$$

Proof. The proof is provided in the supplementary material. \square

Proposition 3.1 shows that the Lipschitz constant of a neural network is controlled when trained with the spectral normalization in Equation (5). In particular, recalling the notations in Proposition 3.1, in the limit where $d_i \rightarrow \infty$ with $\frac{d_i}{d_{i-1}} \rightarrow \gamma_i \in (0, \infty)$ for all $i \in [N]$ and choosing the learning rate $\eta = \mathcal{O}(d_0^{-1})$, the Lipschitz constant of \mathcal{G} is of order $\mathcal{O}(1)$ if it has finitely many layers N and σ_* is constant. Therefore, with this spectral normalization, it can be assumed that $\|\mathcal{G}\|_{lip} = \mathcal{O}(1)$ when dimensions grow. Figure 2 depicts the behavior of the Lipschitz constant of a linear layer with and without spectral normalization in the setting of Proposition 3.1, which confirms the obtained bound.

3.2. Mixture of Concentrated Vectors

In this section, we assume data to be a mixture of concentrated random vectors with controlled $\mathcal{O}(1)$ Lipschitz constant (e.g., DL representations of GAN-data as we discussed in the previous section). Precisely, let $\mathbf{x}_1, \dots, \mathbf{x}_n$

be a set of mutually independent random vectors in \mathbb{R}^p . We suppose that these vectors are distributed as one of k classes of distribution laws μ_1, \dots, μ_k with distinct means $\{\mathbf{m}_\ell\}_{\ell=1}^k$ and ‘‘covariances’’ $\{\mathbf{C}_\ell\}_{\ell=1}^k$ defined respectively as

$$\mathbf{m}_\ell = \mathbb{E}_{\mathbf{x}_i \sim \mu_\ell}[\mathbf{x}_i], \quad \mathbf{C}_\ell = \mathbb{E}_{\mathbf{x}_i \sim \mu_\ell}[\mathbf{x}_i \mathbf{x}_i^\top]. \quad (8)$$

For some $q > 0$, we consider a q -exponential concentration property on the laws μ_ℓ , in the sense that for any family of independent vectors $\mathbf{y}_1, \dots, \mathbf{y}_s$ sampled from μ_ℓ , $[\mathbf{y}_1, \dots, \mathbf{y}_s] \in \mathcal{E}_q(1 | \mathbb{R}^{p \times s}, \|\cdot\|_F)$. Without loss of generality, we arrange the \mathbf{x}_i 's in a data matrix $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ such that, for each $\ell \in [k]$, $\mathbf{x}_{1+\sum_{j=1}^{\ell-1} n_j}, \dots, \mathbf{x}_{\sum_{j=1}^{\ell} n_j} \sim \mu_\ell$, where n_ℓ stands for the number of \mathbf{x}_i 's sampled from μ_ℓ . In particular, we have the concentration of \mathbf{X} as

$$\mathbf{X} \in \mathcal{E}_q(1 | \mathbb{R}^{p \times n}, \|\cdot\|_F). \quad (9)$$

Such a data matrix \mathbf{X} can be constructed through Lipschitz-ally transformed Gaussian vectors ($q = 2$), with controlled Lipschitz constant, thanks to Corollary 2.4. In particular, DL representations of GAN-data are constructed as such, as shown in Section 3.1. We further introduce the following notations that will be used subsequently.

$$\mathbf{M} = [\mathbf{m}_1, \dots, \mathbf{m}_k] \in \mathbb{R}^{p \times k}, \quad \mathbf{J} = [\mathbf{j}_1, \dots, \mathbf{j}_k] \in \mathbb{R}^{n \times k}$$

and $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_n] \in \mathbb{R}^{p \times n}$,

where $\mathbf{j}_\ell \in \mathbb{R}^n$ stands for the canonical vector selecting the \mathbf{x}_i 's of distribution μ_ℓ , defined by $(\mathbf{j}_\ell)_i = \mathbf{1}_{\mathbf{x}_i \sim \mu_\ell}$, and the \mathbf{z}_i 's are the centered versions of the \mathbf{x}_i 's, i.e. $\mathbf{z}_i = \mathbf{x}_i - \mathbf{m}_\ell$ for $\mathbf{x}_i \sim \mu_\ell$.

3.3. Gram Matrices of Concentrated Vectors

Now we study the behavior of the Gram matrix $\mathbf{G} = \frac{1}{p} \mathbf{X}^\top \mathbf{X}$ in the large n, p limit and under the model of the previous section. Indeed, \mathbf{G} appears as a central component in many classification, regression and clustering methods. Precisely, a finer description of the behavior of \mathbf{G} provides access to the internal functioning and performance evaluation of a wide range of machine learning methods such as Least Squares SVMs (AK et al., 2002), Semi-supervised Learning (Chapelle et al., 2009) and Spectral Clustering (Ng et al., 2002). Indeed, the performance evaluation of these methods has already been studied under GMM models in (Liao & Couillet, 2017; Mai & Couillet, 2017; Couillet & Benaych-Georges, 2016) through RMT. On the other hand, analyzing the spectral behavior of \mathbf{G} for DL representations quantifies their quality –through its principal subspace (Montavon et al., 2011)– as we have discussed in the introduction. In particular, the Gram matrix decomposes as

$$\mathbf{G} = \frac{1}{p} \mathbf{J} \mathbf{M}^\top \mathbf{M} \mathbf{J}^\top + \frac{1}{p} \mathbf{Z}^\top \mathbf{Z} + \frac{1}{p} (\mathbf{J} \mathbf{M}^\top \mathbf{Z} + \mathbf{Z}^\top \mathbf{M} \mathbf{J}^\top).$$

Intuitively \mathbf{G} decomposes as a low-rank informative matrix containing the class canonical vectors through \mathbf{J} and a noise term represented by the other matrices and essentially $\mathbf{Z}^\top \mathbf{Z}$. Given the form of this decomposition, RMT predicts –through an analysis of the spectrum of \mathbf{G} and under a GMM model (Benaych-Georges & Couillet, 2016)– the existence of a threshold ξ function of the ratio p/n and the data statistics for which the dominant eigenvectors of \mathbf{G} contain information about the classes only when $\|\mathbf{M}^\top \mathbf{M}\| \geq \xi$ asymptotically (i.e., only when the means of the different classes are sufficiently distinct).

In order to characterize the spectral behavior (i.e., eigenvalues and leading eigenvectors) of \mathbf{G} under the concentration assumption in Equation (9) on \mathbf{X} , we will be interested in determining the spectral distribution $L = \frac{1}{n} \sum_{i=1}^n \delta_{\lambda_i}$ of \mathbf{G} , with $\lambda_1, \dots, \lambda_n$ the eigenvalues of \mathbf{G} , where δ_x stands for the Dirac measure at point x . Essentially, to determine the limiting eigenvalue distribution as $p, n \rightarrow \infty$ and $p/n \rightarrow c \in (0, \infty)$, a conventional approach in RMT consists in determining an estimate of the *Stieltjes transform* (Silverstein & Choi, 1995) m_L of L , which is defined for some $z \in \mathbb{C} \setminus \text{Supp}(L)$

$$m_L(z) = \int_{\lambda} \frac{dL(\lambda)}{\lambda - z} = \frac{1}{n} \text{tr}((\mathbf{G} - z\mathbf{I}_n)^{-1}). \quad (10)$$

Hence, quantifying the behavior of the *resolvent* of \mathbf{G} defined as $\mathbf{R}(z) = (\mathbf{G} + z\mathbf{I}_n)^{-1}$ determines the limiting measure of L through $m_L(z)$. Furthermore, since $\mathbf{R}(z)$ and \mathbf{G} share the same eigenvectors with associated eigenvalues $\frac{1}{\lambda_i + z}$, the projector matrix corresponding to the top m eigenvectors $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_m]$ of \mathbf{G} can be calculated through a Cauchy integral $\mathbf{U} \mathbf{U}^\top = \frac{-1}{2\pi i} \oint_{\gamma} \mathbf{R}(-z) dz$ where γ is an oriented complex contour surrounding the top m eigenvalues of \mathbf{G} .

To study the behavior of $\mathbf{R}(z)$, we look for a so-called *deterministic equivalent* (Hachem et al., 2007) $\tilde{\mathbf{R}}(z)$ for $\mathbf{R}(z)$, which is a deterministic matrix that satisfies for all $\mathbf{A} \in \mathbb{R}^{n \times n}$ and all $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$ of respectively bounded spectral and Euclidean norms, $\frac{1}{n} \text{tr}(\mathbf{A} \mathbf{R}(z)) - \frac{1}{n} \text{tr}(\mathbf{A} \tilde{\mathbf{R}}(z)) \rightarrow 0$ and $\mathbf{u}^\top (\mathbf{R}(z) - \tilde{\mathbf{R}}(z)) \mathbf{v} \rightarrow 0$ almost surely as $n \rightarrow \infty$. In the following, we present our main result which gives such a deterministic equivalent under the concentration assumption on \mathbf{X} in Equation (9) and under the following assumptions.

Assumption 3.2. As $p \rightarrow \infty$,

1. $p/n \rightarrow c \in (0, \infty)$,
2. The number of classes k is bounded,
3. $\|\mathbf{m}_\ell\| = \mathcal{O}(\sqrt{p})$.

Theorem 3.3 (Deterministic Equivalent for $\mathbf{R}(z)$). *Under the model described in Section 3.2 and Assumptions 3.2, we*

have $\mathbf{R}(z) \in \mathcal{E}_q(p^{-1/2} | \mathbb{R}^{n \times n}, \|\cdot\|_F)$. Furthermore,

$$\begin{aligned} \|\mathbb{E}\mathbf{R}(z) - \tilde{\mathbf{R}}(z)\| &= \mathcal{O}\left(\sqrt{\frac{\log(p)}{p}}\right) \\ \tilde{\mathbf{R}}(z) &= \frac{1}{z} \text{diag} \left\{ \frac{\mathbf{I}_{n_\ell}}{1 + \delta_\ell^*(z)} \right\}_{\ell=1}^k + \frac{1}{pz} \mathbf{J} \Omega_z \mathbf{J}^\top \end{aligned} \quad (11)$$

with $\Omega_z = \mathbf{M}^\top \tilde{\mathbf{Q}}(z) \mathbf{M} \odot \text{diag} \left\{ \frac{\delta_\ell^*(z)-1}{\delta_\ell^*(z)+1} \right\}_{\ell=1}^k$ and $\tilde{\mathbf{Q}}(z) = \left(\frac{1}{ck} \sum_{\ell=1}^k \frac{\mathbf{C}_\ell}{1 + \delta_\ell^*(z)} + z \mathbf{I}_p \right)^{-1}$,

where $\delta^*(z) = [\delta_1^*(z), \dots, \delta_k^*(z)]^\top$ is the unique fixed point of the system of equations for each $\ell \in [k]$

$$\delta_\ell(z) = \frac{1}{p} \text{tr} \left(\mathbf{C}_\ell \left(\frac{1}{ck} \sum_{j=1}^k \frac{\mathbf{C}_j}{1 + \delta_j(z)} + z \mathbf{I}_p \right)^{-1} \right).$$

Sketch of proof. The first step of the proof is to show the concentration of $\mathbf{R}(z)$. This comes from the fact that the application $\mathbf{X} \mapsto \mathbf{R}(z)$ is $2z^{-3/2}p^{-1/2}$ -Lipschitz w.r.t. the Frobenius norm, thus we have by Proposition 2.3 that $\mathbf{R}(z) \in \mathcal{E}_q(p^{-1/2} | \mathbb{R}^{n \times n}, \|\cdot\|_F)$. The second step consists in estimating $\mathbb{E}\mathbf{R}(z)$ through a deterministic matrix $\tilde{\mathbf{R}}(z)$. Indeed, $\mathbf{R}(z)$ can be expressed as a function of $\mathbf{Q}(z) = (\mathbf{X}\mathbf{X}^\top/p + z\mathbf{I}_p)^{-1}$ as $\mathbf{R}(z) = z^{-1}(\mathbf{I}_n - \mathbf{X}^\top \mathbf{Q}(z) \mathbf{X}/p)$, where the statistical dependency between \mathbf{X} and $\mathbf{Q}(z)$ is handled through Lemmas 1.1 and 1.2 (see supplementary material) and finally exploiting the result of (Louart & Couillet, 2019) which shows that $\mathbb{E}\mathbf{Q}(z)$ can be estimated through $\tilde{\mathbf{Q}}(z)$, we obtain the estimator $\tilde{\mathbf{R}}(z)$ for $\mathbb{E}\mathbf{R}(z)$. A more detailed proof is provided in the supplementary material. \square

This result allows specifically to (i) describe the limiting eigenvalues distribution of \mathbf{G} , (ii) determine the spectral detectability threshold mentioned above, (iii) evaluate the asymptotic ‘‘content’’ of the leading eigenvectors of \mathbf{G} and, much more fundamentally, (iv) infer the asymptotic performances of machine learning algorithms that are based on simple functionals of \mathbf{G} (e.g., LS-SVM, spectral clustering etc.). Looking carefully at Theorem 3.3 we see that the spectral behavior of the Gram matrix \mathbf{G} computed on concentrated vectors only depends on the *first* and *second* order statistics of the laws μ_ℓ (their means m_ℓ and ‘‘covariances’’ \mathbf{C}_ℓ). This suggests the surprising result that \mathbf{G} has the same behavior as when the data follow a GMM model with the same means and covariances. The asymptotic spectral behavior of \mathbf{G} is therefore *universal* with respect to the data distribution laws which satisfy the aforementioned concentration properties (for instance DL representations of GAN-data). We illustrate this universality result in the next section by considering data as CNN representations of GAN generated images.



Figure 3. (Top) GAN generated images using the BigGAN model (Brock et al., 2018). (Bottom) Real images selected from the Imagenet dataset (Deng et al., 2009). We considered $n = 1500$ images from $k = 3$ classes which are {mushroom, pizza, hamburger}.

4. Application to CNN Representations of GAN-generated Images

This section presents experiments that confirm the result of Theorem 3.3. In particular, we compare, in the first part, the eigenvalues distribution and the largest eigenvectors of the Gram matrix computed on deep learning representations with those of the Gram matrix computed on Gaussian data with the same first and second order moments. In the second part of this section, we evaluate the performance of a linear SVM model on the principal subspace of the Gram matrix (computed on the representations or on the corresponding Gaussian data) by varying the number of components in the same vein as the work of (Montavon et al., 2011). In the following, all representation networks are standard convolutional neural networks pre-trained on the Imagenet dataset (Deng et al., 2009), in particular, we used pre-trained models of the Pytorch deep learning framework.

4.1. Spectrum and Dominant Eigenspace of the Gram Matrix

In this section, we consider $n = 1500$ data $x_1, \dots, x_n \in \mathbb{R}^p$ as CNN representations –across popular CNN architectures of different sizes p – of GAN-generated images using the generator of the Big-GAN model (Brock et al., 2018). We further use real images from the Imagenet dataset (Deng et al., 2009) for comparison. In particular, we empirically compare the spectrum of the Gram matrix of this data with the Gram matrix of a GMM model with the same means and covariances. We also consider the leading 2-dimensional eigenspace of the Gram matrix which contains clustering

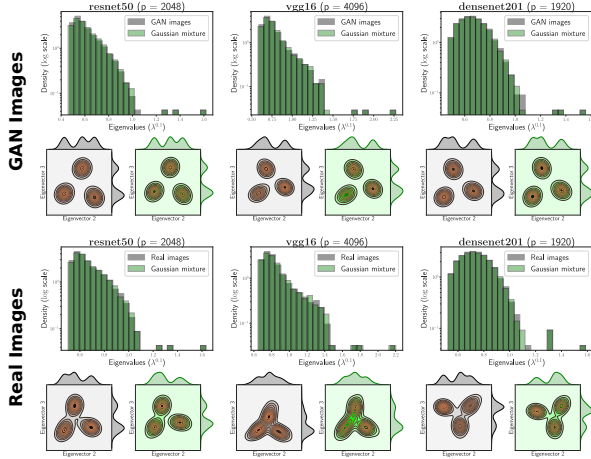


Figure 4. **(Top)** Spectrum and leading eigenspace of the Gram matrix for CNN representations of GAN generated images using the BigGAN model (Brock et al., 2018). **(Bottom)** Spectrum and leading eigenspace of the Gram matrix for CNN representations of real images selected from the Imagenet dataset (Deng et al., 2009). Columns correspond to the three representation networks which are *resnet50*, *vgg16* and *densenet201*. We used $n = 1500$ images and considered $k = 3$ classes as depicted in Figure 3.

information as detailed in the previous section. Figure 3 depicts some images generated using the Big-GAN model (Top) and the corresponding real class images from the Imagenet dataset (Bottom). The Big-GAN model is visually able to generate highly realistic images which are by construction concentrated vectors, as discussed in Section 3.1 and therefore satisfy the assumptions of Theorem 3.3.

Figure 4 depicts the spectrum and leading 2D eigenspace of the Gram matrix computed on CNN representations of GAN generated and real images (in gray), and the corresponding GMM model with same first and second order statistics (in green). The Gram matrix is seen to follow the same spectral behavior for GAN-data as for the GMM model which is a natural consequence of the universality result of Theorem 3.3 with respect to the data distribution. Besides, and perhaps no longer surprisingly, we further observe that the spectral properties of \mathbf{G} for real data (here CNN representations of real images) are conclusively matched by their Gaussian counterpart. Figure 5 shows more results about the Gram matrix spectrum of the representations (in black) and the corresponding Gaussian data (in green), by considering more representation networks and using $k = 6$ classes for both GAN images and real images, which confirms the result of Theorem 3.3. This both theoretically and empirically confirms that the proposed random matrix framework is fully compliant with the theoretical analysis of real machine learning datasets. As a consequence, recalling the work of (Montavon et al., 2011), the *quality* of a given representation is quantified through the prediction

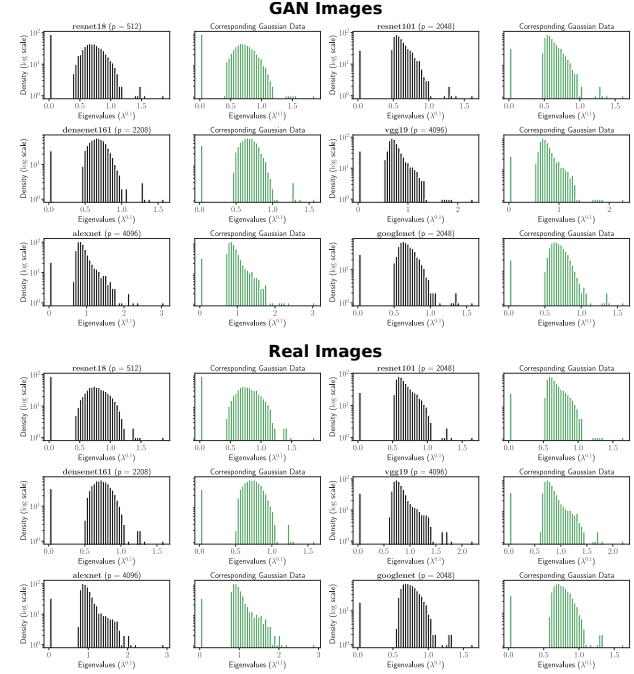


Figure 5. Spectrum of the Gram matrix for CNN representations in **(black)** and the corresponding Gaussian data **(in green)** for GAN generated images using the BigGAN model (Brock et al., 2018) **(Top)** and for real images randomly selected from the Imagenet dataset (Deng et al., 2009) **(Bottom)**. The considered representation network are *resnet18*, *resnet101*, *densenet161*, *vgg19*, *alexnet* and *googlenet*. We used $n = 600$ images selected among $k = 6$ classes {hamburger, mushroom, pizza, strawberry, coffee, daisy} (100 images per class).

accuracy of a linear classifier trained on the principal Gram matrix eigenvectors of the representations computed on a set of samples. Given our result in Theorem 3.3, and the fact that the top m eigenvectors $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_m]$ of \mathbf{G} are related to the resolvent matrix $\mathbf{R}(z)$ through the Cauchy integral $\mathbf{U}\mathbf{U}^\top = \frac{-1}{2\pi i} \oint_\gamma \mathbf{R}(-z) dz$ where γ is an oriented complex contour surrounding the top m eigenvalues of \mathbf{G} , we should expect that the prediction accuracy of a linear classifier trained on the principal eigenvectors of \mathbf{G} be the same for the representations themselves as for the corresponding Gaussian data with the same first and second moments. Therefore, the purpose of the following section is to show simulations which confirm this result.

4.2. Linear SVM Performance on the Dominant k -dimensional Eigenspace of \mathbf{G}

In this section, we compare the performance of a linear SVM model trained on the dominant Gram matrix's k -dimensional eigenspace of the representations versus the corresponding Gaussian data with the same first and second

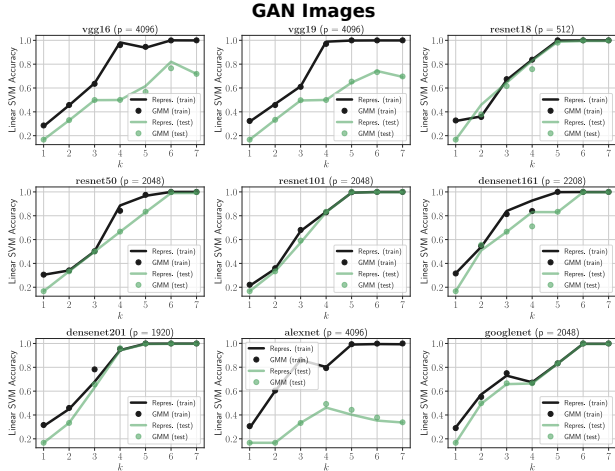


Figure 6. Train and test accuracy of a linear SVM model trained on the top k eigenvectors of the Gram matrix, computed on the representations of GAN generated images. We generated with the BigGAN model (Brock et al., 2018) $n = 6000$ images belonging to the 6 classes {hamburger, mushroom, pizza, strawberry, coffee, daisy} (1000 images per class). We considered 9 representation networks which are *vgg16*, *vgg19*, *resnet18*, *resnet50*, *resnet101*, *densenet161*, *densenet201*, *alexnet* and *googlenet*. Lines represent the performance of the SVM model on the representations themselves whereas dots represent the performance of the SVM model on Gaussian data with the same first and second order moments. We used a train vs test split of 2/3 and 1/3 respectively.

order moments. Experiments were made in the following settings:

- **Data types:** We do the experiments for both GAN generated images using the BigGAN model (Brock et al., 2018) and for real images randomly selected for the Imagenet dataset (Deng et al., 2009). In both cases we consider $n = 6000$ images.
- **Classes:** We consider $k = 6$ classes which are: hamburger, mushroom, pizza, strawberry, coffee and daisy.
- **Representation networks:** We consider 9 representation networks pre-trained on the Imagenet dataset (Deng et al., 2009) which are: *vgg16* ($p = 4096$), *vgg19* ($p = 4096$), *resnet18* ($p = 512$), *resnet50* ($p = 2048$), *resnet101* ($p = 2048$), *densenet161* ($p = 2208$), *densenet201* ($p = 1920$), *alexnet* ($p = 4096$) and *googlenet* ($p = 2048$).

Figure 6 depicts the train and test accuracy of a linear SVM trained on the top k eigenvectors of \mathbf{G} , for the representations (of GAN generated images) and the corresponding Gaussian data, for different values of k . As we can notice, the performance of the SVM model on the representations matches its performance on the corresponding Gaussian data

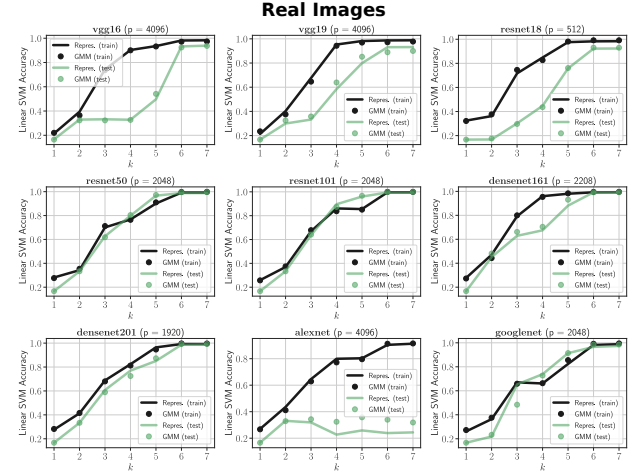


Figure 7. Train and test accuracy of a linear SVM model trained on the top k eigenvectors of the Gram matrix, computed on the representations of Real images. We randomly sampled $n = 6000$ images belonging to the 6 classes {hamburger, mushroom, pizza, strawberry, coffee, daisy} (1000 images per class) of the Imagenet dataset (Deng et al., 2009). We considered 9 representation networks *vgg16*, *vgg19*, *resnet18*, *resnet50*, *resnet101*, *densenet161*, *densenet201*, *alexnet* and *googlenet*. Lines represent the performance of the SVM model on the representations themselves whereas dots represent the performance of the SVM model on Gaussian data with the same first and second order moments. We used a train vs test split of 2/3 and 1/3 respectively.

with the same first and second order statistics as predicted by Theorem 3.3. This matching seems to extend beyond GAN images (which are concentrated vectors) to real images as depicted in Figure 7. As a consequence, our results suggest that the *quality* of a given representation network can be quantified through their first two statistical moments.

5. Conclusion

Leveraging on random matrix theory (RMT) and the concentration of measure phenomenon, we have shown through this paper that DL representations of GAN-data behave as Gaussian mixtures for linear classifiers, a fundamental *universal* property which is only valid in high-dimension of data. To the best of our knowledge, this result constitutes a new approach towards the theoretical understanding of complex objects such as DL representations, as well as the understanding of the behavior of more elaborate machine learning algorithms for complex data structures. In addition, the article explicitly demonstrated our ability, through RMT, to anticipate the behavior of a certain range of widely used standard classifiers for data as complex as DL representations of the realistic and surprising images generated by GANs. This opens the way to a more systematic analysis and improvement of machine learning algorithms on real datasets by means of large dimensional statistics.

References

- AK, S. J. et al. *Least squares support vector machines*. World Scientific, 2002.
- Antognini, J. and Sohl-Dickstein, J. Pca of high dimensional random walks with comparison to neural network training. In *Advances in Neural Information Processing Systems*, pp. 10307–10316, 2018.
- Antoniou, A., Storkey, A., and Edwards, H. Data augmentation generative adversarial networks. *arXiv preprint arXiv:1711.04340*, 2017.
- Benaych-Georges, F. and Couillet, R. Spectral analysis of the gram matrix of mixture models. *ESAIM: Probability and Statistics*, 20:217–237, 2016.
- Bengio, Y., Courville, A., and Vincent, P. Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8):1798–1828, Aug 2013. ISSN 0162-8828. doi: 10.1109/TPAMI.2013.50.
- Bengio, Y. et al. Learning deep architectures for ai. *Foundations and trends® in Machine Learning*, 2(1):1–127, 2009.
- Brock, A., Donahue, J., and Simonyan, K. Large scale gan training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2018.
- Chapelle, O., Scholkopf, B., and Zien, A. Semi-supervised learning (chapelle, o. et al., eds.; 2006)[book reviews]. *IEEE Transactions on Neural Networks*, 20(3):542–542, 2009.
- Couillet, R. and Benaych-Georges, F. Kernel spectral clustering of large dimensional data. *Electronic Journal of Statistics*, 10(1):1393–1454, 2016.
- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. Generative adversarial nets. In *NIPS*, pp. 2672–2680, 2014.
- Goodfellow, I., Bengio, Y., Courville, A., and Bengio, Y. *Deep learning*, volume 1. MIT Press, 2016.
- Gulrajani, I., Ahmed, F., Arjovsky, M., Dumoulin, V., and Courville, A. C. Improved training of wasserstein gans. In *Advances in Neural Information Processing Systems 30*, pp. 5767–5777. 2017.
- Hachem, W., Loubaton, P., Najim, J., et al. Deterministic equivalents for certain functionals of large random matrices. *The Annals of Applied Probability*, 17(3):875–930, 2007.
- Ioffe, S. and Szegedy, C. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- Ledoux, M. *The concentration of measure phenomenon*. Number 89. American Mathematical Soc., 2005.
- Liao, Z. and Couillet, R. Random matrices meet machine learning: A large dimensional analysis of ls-svm. In *ICASSP*, pp. 2397–2401. IEEE, 2017.
- Louart, C. and Couillet, R. Concentration of measure and large random matrices with an application to sample covariance matrices. *submitted*, 2019.
- Mai, X. and Couillet, R. A random matrix analysis and improvement of semi-supervised learning for large dimensional data. *arXiv preprint arXiv:1711.03404*, 2017.
- Mirza, M. and Osindero, S. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014.
- Miyato, T., Kataoka, T., Koyama, M., and Yoshida, Y. Spectral normalization for generative adversarial networks. *arXiv preprint arXiv:1802.05957*, 2018.
- Montavon, G., Braun, M. L., and Miller, K.-R. Kernel analysis of deep networks. *Journal of Machine Learning Research*, 12(Sep):2563–2581, 2011.
- Ng, A. Y., Jordan, M. I., and Weiss, Y. On spectral clustering: Analysis and an algorithm. In *Advances in neural information processing systems*, pp. 849–856, 2002.
- Roth, K., Lucchi, A., Nowozin, S., and Hofmann, T. Stabilizing training of generative adversarial networks through regularization. In *Advances in Neural Information Processing Systems 30*, pp. 2018–2028. 2017.
- Rumelhart, D. E., Hinton, G. E., Williams, R. J., et al. Learning representations by back-propagating errors. *Cognitive modeling*, 5(3):1, 1988.
- Silverstein, J. W. and Choi, S.-I. Analysis of the limiting spectral distribution of large dimensional random matrices. *Journal of Multivariate Analysis*, 54(2):295–309, 1995.
- Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., and Fergus, R. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.
- Yeh, C.-K., Kim, J., Yen, I. E.-H., and Ravikumar, P. K. Representer point selection for explaining deep neural networks. In *Advances in Neural Information Processing Systems*, pp. 9291–9301, 2018.