

# Neural Networks Classify through the Class-wise Means of their Representations

Mohamed El Amine Seddik<sup>1</sup>, Mohamed Tamaazousti<sup>2</sup>

<sup>1</sup>Huawei Paris Research Center, 92012 Boulogne-Billancourt, France

<sup>2</sup>CEA List, 91120 Palaiseau, France

mohamed.seddik@huawei.com, mohamed.tamaazousti@cea.fr

## Abstract

In this paper, based on an asymptotic analysis of the Softmax layer, we show that when training neural networks for classification tasks, the weight vectors corresponding to each class of the Softmax layer tend to converge to the class-wise means computed at the representation layer (for specific choices of the representation activation). We further show some consequences of our findings to the context of transfer learning, essentially by proposing a simple yet effective initialization procedure that significantly accelerates the learning of the Softmax layer weights as the target domain gets closer to the source one. Experiments are notably performed on the datasets: MNIST, Fashion MNIST, Cifar10, and Cifar100 and using a standard CNN architecture.

## Introduction

One of the big trends of the moment is opening up access to the Deep Learning technology to non-experts. Therefore, there is an emergent need to move from manual design of neural networks models by experts to an automatic design for non-experts. This methodological shift is known as AutoML (for Automatic Machine Learning). Technically, this consists in designing automatically the architecture of a neural network as well as its hyperparameters, e.g. the Neural Architecture Search (NAS) approach (Zoph and Le 2017; Zoph et al. 2018; Cai, Zhu, and Han 2019; Liu, Simonyan, and Yang 2019; Real et al. 2019).

Basically, this trend implies to handle more and more specific data domains with generally few data provided by the user for each domain. In this context, the NAS approach, which needs a lot of training data, is no longer relevant. Consequently, the practical solution generally adopted to overcome this problem is to “recycle”, as much as possible, the networks already learned, with a minimum of adaptation. This is facilitated by the fact that there are more and more pretrained neural networks available online (open source pretrained models<sup>1</sup>).

To this end, one of the most simple and widespread approach is the transfer of pre-trained representations referred to as Transfer Learning (TL) (Kuzborskij, Orabona, and Cauto 2013; Sharif Razavian et al. 2014; Yosinski et al. 2014;

Ginsca et al. 2015; Azizpour et al. 2015; Kornblith, Shlens, and Le 2019; Tamaazousti et al. 2019) in the rest of the document. In principle, TL consists in learning a neural network on a source domain for which we have a large amount of data and then using this network on a target domain (generally with few available data) by updating only the last layer (very often with Softmax activation) of weight associated with the task to be solved (e.g., classification task). Therefore, TL plays a key role in this context by reducing drastically the hyper-parameters space and replacing the neural architecture search by the best available source model/representation search.

Unfortunately, a scale up of Deep Learning based on this TL methodology will be with a huge computational consequences. Indeed, the TL approach implies to duplicate the search for the best hyper-parameters, generally obtained by a very costly grid search, as many times as there are available source models. To the best of our knowledge, there is no clear solution to an *a priori* (i.e., without training the last layer) choice of the best source model from the data associated to the target domain.

To tackle this issue, our (1) first contribution consists in analyzing theoretically the behavior of the Softmax layer in neural networks classifiers (see section *Asymptotic analysis of the Softmax classifier weights*). Our (2) second contribution consists in showing that Softmax classifies through the class-wise means of the data representations (see section *Expression of the class-weights vectors for near-optimal representations*) for specific choices of the representation activation. This second result is obtained under additional hypotheses on the representations statistics of the data, based on commonly accepted interpretations of the properties expected from the representations of DNNs (Bengio, Courville, and Vincent 2013; Tamaazousti et al. 2019). Our results are experimentally confirmed by showing that when training neural networks for classification tasks, the weight vectors corresponding to each class of the Softmax layer tend to be similar to the class-wise means computed at the representation layer.

Based on this finding, our (3) third contribution then consists in addressing the source model selection in the TL context for classification. In addition, this result provides a simple yet effective initialization procedure (through the class-wise means of the representations computed on the target

domain) which significantly accelerates TL as the target domain gets closer to the source one. Experiments on real data confirm empirically the relevance of these technical choices (see section *Experiments*). Note that, our second and main contribution might be of independent interest even outside the context of TL. Indeed, our theoretical results also support the remarkable observed performances of the Nearest Class Mean Classifier (Guerriero, Caputo, and Mensink 2018).

*Notation:* Vectors are represented by boldface lowercase letters and matrices are represented by boldface capital letters. The  $i$ -th entry of a vector  $\mathbf{v}$  is denoted by  $v_i$ , and element  $(i, j)$  of a matrix  $\mathbf{M}$  is denoted by  $M_{i,j}$ . The  $i$ -th row and  $j$ -th column of  $\mathbf{M}$  are denoted respectively by  $\mathbf{M}_{i,:}$  and  $\mathbf{M}_{:,j}$ . Let  $[p]$  denote the set of integers  $\{1, \dots, p\}$ .  $\delta_{i,j}$  denotes the Kronecker delta. The notation  $\mathbf{v}_1 = \mathbf{v}_2 + \mathcal{O}_{\|\cdot\|}(u_p)$  means that for any unitary vector  $\mathbf{u}$ ,  $\mathbf{u}^\top(\mathbf{v}_1 - \mathbf{v}_2) = \mathcal{O}(u_p)$ .  $\mathbf{1}_p$  denotes the vector of size  $p$  full of ones.  $\|\cdot\|$  denotes the  $\ell_2$  norm for vectors and the operator norm for matrices. Denote by  $\text{Cos}(\mathbf{u}, \mathbf{v}) \equiv (\mathbf{u}/\|\mathbf{u}\|)^\top (\mathbf{v}/\|\mathbf{v}\|)$  the cosine similarity between two vectors  $\mathbf{u}$  and  $\mathbf{v}$ .

## Settings

In all the paper we consider the following settings. Let  $\mathcal{E}$  denote some data space and  $\mathcal{T} = \{(\mathbf{d}_i, \ell_i)\}_{i=1}^n$  be a training set containing  $n$  data points  $\mathbf{d}_i$  belonging to one of  $k$  different classes  $\mathcal{C}_1, \dots, \mathcal{C}_k \subset \mathcal{E}$ , where  $\ell_i \in [k]$  is the label for data point  $\mathbf{d}_i$ . Let  $n_j$  denote the number of training points in class  $j$  so that  $n = \sum_{j=1}^k n_j$ . Let  $c : [n] \rightarrow [k]$  be the function which returns the class-index of a datum  $\mathbf{x}_i$  so is defined such that  $c(i) = \ell$  if  $\mathbf{x}_i \in \mathcal{C}_\ell$ . We further denote

$$\mathbf{x} \equiv \varphi \circ \phi(\mathbf{d}; \Theta) \in \mathbb{R}^p, \quad (1)$$

the representation of dimension  $p$  for  $\mathbf{d} \in \mathcal{E}$ , where  $\phi : \mathcal{E} \rightarrow \mathbb{R}^p$  is typically implemented by a deep CNN model (with the final layer being a dense layer with a linear activation function) parameterized by  $\Theta$ , and  $\varphi : \mathbb{R} \rightarrow \mathbb{R}$  stands for the activation function at the representation layer applied element-wise to  $\phi(\mathbf{d}; \Theta)$ . The final class prediction  $\tilde{\ell} \in [k]$  is given by a classifier function  $\psi : \mathbb{R}^p \rightarrow \mathbb{R}^k$ , such that  $\tilde{\ell} = \arg \max \psi(\mathbf{x})$ . Generally,  $\psi$  is a Softmax classifier with weight matrix  $\mathbf{W} \in \mathbb{R}^{p \times k}$  and bias vector  $\mathbf{b} \in \mathbb{R}^k$  so that  $\psi(\mathbf{x}) = \text{softmax}(\mathbf{W}^\top \mathbf{x} + \mathbf{b})$ . And the Softmax function is defined for a vector  $\mathbf{v} \in \mathbb{R}^k$  such that  $\text{softmax}(\mathbf{v})_i = e^{v_i} / \sum_{j=1}^k e^{v_j}$ . We consider in all what follows that the bias vector  $\mathbf{b}$  is null since its influence is negligible (Kang et al. 2019).

Define the following statistics of the representation layer, for each  $\ell \in [k]$

$$\mathbf{m}_\ell \equiv \mathbb{E}_{\mathbf{d} \in \mathcal{C}_\ell} [\mathbf{x}], \quad \mathbf{C}_\ell \equiv \mathbb{E}_{\mathbf{d} \in \mathcal{C}_\ell} [(\mathbf{x} - \mathbf{m}_\ell)(\mathbf{x} - \mathbf{m}_\ell)^\top], \quad (2)$$

which are the expected class-wise means and covariances of the representations as defined in equation 1. Authors in (Sedik et al. 2020) have shown that the performances of a linear classifier applied on top of the representations  $\mathbf{x}_i$ 's can be quantified through their class-wise means and covariances, namely the quantities  $\mathbf{m}_\ell$  and  $\mathbf{C}_\ell$ , when the data  $\mathbf{d}_i$ 's are assumed to be generated by a GAN. A result which suggests

that the representations behave as Gaussian mixtures for linear classifiers. Let  $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_k] \in \mathbb{R}^{p \times k}$  where  $\mathbf{w}_\ell$  are the Softmax classifier weights corresponding to class  $\ell$ . In following, we aim to explicit the relationship between the class-weights  $\mathbf{w}_\ell$  and the quantities  $\mathbf{m}_\ell$  and  $\mathbf{C}_\ell$ , when the Softmax classifier is trained on the representations. Namely, assuming a Gaussian Mixture Model on the data representations  $\mathbf{x}_i$  we give subsequently an implicit expression between  $\mathbf{w}_\ell$  and the representations statistics (cf. equation 8). And under specific assumptions on  $\mathbf{m}_\ell$  and  $\mathbf{C}_\ell$  which we refer to as *the near-optimal representations* assumption, we give a more explicit relationship between  $\mathbf{w}_\ell$  and the class-wise means  $\mathbf{m}_\ell$  (cf. equation 10).

## Main Results

### Preliminaries

Before introducing our results, we give the following technical tool which will be of central interest.

**Proposition 0.1.** *Let  $\mathbf{m} \in \mathbb{R}^p$  a deterministic (mean) vector and  $\mathbf{C} \in \mathbb{R}^{p \times p}$  a positive semi-definite (covariance) matrix, consider  $\mathbf{x} \sim \mathcal{N}(\mathbf{m}, \mathbf{C})$  and let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be some differentiable function and  $\mathbf{w} \in \mathbb{R}^p$  a deterministic vector. Then,*

$$\mathbb{E}[f(\mathbf{w}^\top \mathbf{x})] = \mathbb{E}[f(\mathbf{w}^\top \mathbf{x})]\mathbf{m} + \mathbb{E}[f'(\mathbf{w}^\top \mathbf{x})]\mathbf{C}\mathbf{w}. \quad (3)$$

*Proof.* See supplementary material.  $\square$

**Remark 0.2.** *Proposition 0.1 can be extended beyond the Gaussian assumption when the dimension  $p$  is getting large. Indeed, if  $\mathbf{x} = \mathbf{m} + \mathbf{C}^{\frac{1}{2}}\mathbf{z}$  for some random vector  $\mathbf{z}$  with random i.i.d. entries with zero mean, unit variance and bounded fourth order moment. As  $p \rightarrow \infty$ , assuming further that  $\|\mathbf{C}\| = \mathcal{O}(1)$ , the result of Proposition 0.1 can be extended thanks to the central limit theorem under this more general setting and we get  $\mathbb{E}[f(\mathbf{w}^\top \mathbf{x})] = \mathbb{E}[f(\mathbf{w}^\top \mathbf{x})]\mathbf{m} + \mathbb{E}[f'(\mathbf{w}^\top \mathbf{x})]\mathbf{C}\mathbf{w} + \mathcal{O}_{\|\cdot\|}(\|\mathbf{w}\| p^{-\frac{1}{2}})$ .*

### Asymptotic analysis of the Softmax classifier weights

Given a set of  $n$  labeled representations  $(\mathbf{x}_1, \mathbf{y}^{(1)}), \dots, (\mathbf{x}_n, \mathbf{y}^{(n)})$  belonging to  $k$  classes  $\mathcal{C}_1, \dots, \mathcal{C}_k$ , where  $\mathbf{y}^{(i)} \in \mathbb{R}^k$  are classically one-hot encoded vectors such that  $y_\ell^{(i)} = 1$  and  $y_j^{(i)} = 0$  for  $j \neq \ell$  if  $\mathbf{x}_i \in \mathcal{C}_\ell$ . In the following we will consider a generalized expression for the labels given by:

$$y_\ell^{(i)} = \alpha_{c(i)} \frac{|\delta_{\ell, c(i)} - \varepsilon|}{1 + (k-2)\varepsilon}, \quad (4)$$

where  $\alpha_{c(i)}$  are hyper-parameters and  $\varepsilon > 0$  which will be converging to zero in our analysis,  $c(i) \in [k]$  returns the class index of the  $i$ -th datum and  $\delta_{i,j}$  stands for the Kronecker delta. We will see in the following that a careful choice of these parameters  $\alpha_{c(i)}$  can make the Softmax classifier weights independent of the class proportions, which is a desired property for datasets with unbalanced classes also known as the long-tail recognition problem (Kang et al. 2019). In particular, the classical labels are recovered by setting  $\alpha_{c(i)} = 1$  and  $\varepsilon \rightarrow 0$ .

The Softmax classification task consists in optimizing the categorical-cross entropy loss function

$$\mathcal{L}(\mathbf{w}_1, \dots, \mathbf{w}_k) = -\frac{1}{n} \sum_{i=1}^n \sum_{\ell=1}^k y_{\ell}^{(i)} \log p_{\ell}^{(i)} \quad (5)$$

where

$$p_{\ell}^{(i)} = \frac{e^{\mathbf{w}_{\ell}^{\top} \mathbf{x}_i}}{\sum_{j=1}^k e^{\mathbf{w}_j^{\top} \mathbf{x}_i}} \quad (6)$$

The gradient of which with respect to each class-weight vector  $\mathbf{w}_{\ell}$  is given by

$$\begin{aligned} \nabla_{\mathbf{w}_{\ell}} \mathcal{L} &= \frac{1}{n} \sum_{i=1}^n \alpha_{c(i)} \left( \frac{e^{\mathbf{w}_{\ell}^{\top} \mathbf{x}_i}}{\sum_{j=1}^k e^{\mathbf{w}_j^{\top} \mathbf{x}_i}} - \frac{|\delta_{\ell, c(i)} - \varepsilon|}{1 + (k-2)\varepsilon} \right) \mathbf{x}_i \\ &\equiv \frac{1}{n} \sum_{i=1}^n \alpha_{c(i)} f_{\ell, i}(\mathbf{w}_{\ell}^{\top} \mathbf{x}_i) \mathbf{x}_i. \end{aligned} \quad (7)$$

**Expression of the class-weights vectors under a GMM model** In order to characterize the behavior of the weight vectors  $\mathbf{w}_{\ell}$ , we need to assume a statistical prior on the data representations. Motivated by the result of (Seddik et al. 2020) which shows that deep learning representations tend to behave as Gaussian mixtures when their dimension grow, we make the assumption that the  $\mathbf{x}_i$  follow a Gaussian Mixture Model (GMM) as follows.

**Assumption 0.3** (Statistical model on the representations). *We assume that the data representations are independent random vectors which follow a GMM with means  $\mathbf{m}_{\ell}$  and covariances  $\mathbf{C}_{\ell}$  for each class  $\mathcal{C}_{\ell}$ , further we denote by  $\pi_{\ell} \equiv \lim_n \frac{n_{\ell}}{n} \in (0, 1)$  where  $n_{\ell}$  stands for the cardinality of the class  $\mathcal{C}_{\ell}$ . Formally, if  $\mathbf{x}_i \in \mathcal{C}_{\ell}$  then  $\mathbf{x}_i \sim \mathcal{N}(\mathbf{m}_{\ell}, \mathbf{C}_{\ell})$ . We further assume that both  $p, n \rightarrow \infty$ .*

Under assumption 0.3, the loss function in equation 5 is a sum of i.i.d. random variables and therefore converges to its expectation for sufficiently large  $n$  thanks to the central limit theorem. Therefore, the class-weights vectors  $\mathbf{w}_{\ell}$  concentrate around some deterministic vectors  $\bar{\mathbf{w}}_{\ell}$  which are defined as  $\mathbb{E}[\nabla_{\bar{\mathbf{w}}_{\ell}} \mathcal{L}] = 0$  since the expectation and the gradient are linear operators. The following proposition provides an implicit expression for the asymptotic class-weights vectors  $\bar{\mathbf{w}}_{\ell}$ .

**Proposition 0.4** (Asymptotic class-weight vectors). *Let  $\bar{\mathbf{w}}_1, \dots, \bar{\mathbf{w}}_k$  be the deterministic vectors satisfying  $\mathbb{E}[\nabla_{\bar{\mathbf{w}}_{\ell}} \mathcal{L}] = 0$ . Thus, under Assumption 0.3, each  $\bar{\mathbf{w}}_{\ell}$  satisfies the implicit equation*

$$\bar{\mathbf{w}}_{\ell} = \left( \sum_{j=1}^k \alpha_j \pi_j \mathbb{E}_j[f'_{\ell, i}(\bar{\mathbf{w}}_{\ell}^{\top} \mathbf{x}_i)] \mathbf{C}_j \right)^{-1} \left( \sum_{j=1}^k \alpha_j \pi_j \mathbb{E}_j[f_{\ell, i}(\bar{\mathbf{w}}_{\ell}^{\top} \mathbf{x}_i)] \mathbf{m}_j \right) \quad (8)$$

where the notation  $\mathbb{E}_j[g(\mathbf{x}_i)] \equiv \mathbb{E}[g(\mathbf{x}_i) | \mathbf{x}_i \in \mathcal{C}_j]$  for some  $g : \mathbb{R} \rightarrow \mathbb{R}$ .

*Proof.* The result is straightforward thanks to Proposition 0.1. Indeed,

$$\begin{aligned} \mathbb{E}[\nabla_{\mathbf{w}_{\ell}} \mathcal{L}] &= \sum_{j=1}^k \frac{n_j}{n} \alpha_j \mathbb{E}_j[f_{\ell, i}(\mathbf{w}_{\ell}^{\top} \mathbf{x}_i)] \mathbf{m}_j \\ &\quad + \sum_{j=1}^k \frac{n_j}{n} \alpha_j \mathbb{E}_j[f'_{\ell, i}(\mathbf{w}_{\ell}^{\top} \mathbf{x}_i)] \mathbf{C}_j \mathbf{w}_{\ell}. \end{aligned}$$

Finally, letting  $\mathbb{E}[\nabla_{\mathbf{w}_{\ell}} \mathcal{L}] = 0$  completes the proof.  $\square$

**Remark 0.5** (Relaxation of the Gaussian assumption). *Note, thanks to Remark 0.2, that the Gaussian assumption can be relaxed and one may assume that  $\mathbf{x}_i = \mathbf{m}_{\ell} + \mathbf{C}_{\ell}^{\frac{1}{2}} \mathbf{z}_i$  with  $\mathbf{z}_i$  having i.i.d. entries with zero mean, unit variance and bounded fourth order moment. In which case, we have*

$$\begin{aligned} \bar{\mathbf{w}}_{\ell} &= \left( \sum_{j=1}^k \alpha_j \pi_j \mathbb{E}_j[f'_{\ell, i}(\bar{\mathbf{w}}_{\ell}^{\top} \mathbf{x}_i)] \mathbf{C}_j \right)^{-1} \\ &\quad \left( \sum_{j=1}^k \alpha_j \pi_j \mathbb{E}_j[f_{\ell, i}(\bar{\mathbf{w}}_{\ell}^{\top} \mathbf{x}_i)] \mathbf{m}_j \right) + \mathcal{O}_{\|\cdot\|} \left( \|\bar{\mathbf{w}}_{\ell}\| p^{-\frac{1}{2}} \right). \end{aligned}$$

Looking carefully at the expression in equation 8, we see that in order to make the weight vectors  $\bar{\mathbf{w}}_{\ell}$  independent of the class proportions  $\pi_j$  which is a desirable property for the long-tail recognition problem (Kang et al. 2019), a natural choice of the parameters  $\alpha_j$ 's in the expression of the generalized labels in equation 4 is such that  $\alpha_j = (k \pi_j)^{-1}$ . We hence have the following corollary to Proposition 0.4.

**Corollary 0.6.** *Recalling Proposition 0.4, setting the weight labels  $\alpha_j$ 's such that  $\alpha_j = (k \pi_j)^{-1}$ . The class-weight vectors become asymptotically independent of the class proportions  $\pi_j$ . Specifically,*

$$\begin{aligned} \bar{\mathbf{w}}_{\ell} &= \left( \sum_{j=1}^k \mathbb{E}_j[f'_{\ell, i}(\bar{\mathbf{w}}_{\ell}^{\top} \mathbf{x}_i)] \mathbf{C}_j \right)^{-1} \\ &\quad \left( \sum_{j=1}^k \mathbb{E}_j[f_{\ell, i}(\bar{\mathbf{w}}_{\ell}^{\top} \mathbf{x}_i)] \mathbf{m}_j \right). \end{aligned} \quad (9)$$

Corollary 0.6 provides a general implicit expression for the class-weight vectors in function of the representations statistics namely the means  $\mathbf{m}_j$  and the covariances  $\mathbf{C}_j$ . In the following, we will provide an explicit expression for the class-weights  $\mathbf{w}_{\ell}$  under specific assumptions on the representations statistics which we refer to as the case of *near-optimal representations*.

**Expression of the class-weights vectors for near-optimal representations**

**Assumption 0.7** (Near-optimal representations). *Let  $\epsilon > 0$ , as  $p \rightarrow \infty$*

(A1)  $\mathbf{m}_i^{\top} \mathbf{m}_j = \delta_{i,j} \mu_1 + (1 - \delta_{i,j}) \mu_2$  with  $\mu_1 = \mathcal{O}(1)$  and  $\mu_2 = \mathcal{O}(p^{-\epsilon})$ .

$$(A2) \ C_j = \sigma_{1,j}^2 \mathbf{I}_p + \sigma_{2,j}^2 (\mathbf{1}_p \mathbf{1}_p^\top - \mathbf{I}_p) \text{ with } \sigma_{1,j}^2 = \mathcal{O}(p^{-\epsilon}) \\ \text{and } \sigma_{2,j}^2 = \mathcal{O}(p^{-1-3\epsilon}).$$

An optimal representation for a given dataset is typically a representation which perfectly separates the different classes by maximizing the between-class variance and minimizing the within-class variance while having a large number of degrees of freedom (Bengio, Courville, and Vincent 2013; Tamaazousti et al. 2019). Assumption 0.7 makes a formal definition for what we call a *near-optimal representation* which satisfies for a sufficiently large number of degrees of freedom (dimension  $p$  of the representation) two main properties: (A1) which ensures that the between-class means are asymptotically orthogonal;  $\mathbf{m}_i^\top \mathbf{m}_j \rightarrow \delta_{i,j} \mu_1$  (maximize the between-class variance); and (A2) which makes the within-class covariances asymptotically isotropic;  $\mathbf{C}_j \cong \sigma_{1,j}^2 \mathbf{I}_p$  at the first order (which models the fact that an optimal representation should have independent features) and with low variance;  $\sigma_{1,j}^2 \rightarrow 0$  (minimize the within-class variance). The term *near-optimal* is used in the sense that these properties hold as the dimension  $p$  grows. We will see in the experiments part that when training jointly the representation and the classification parts of a neural network classifier (for a specific choice of the representation activation function), its representations tend to satisfy the near-optimal representations assumption. Under this assumption, we have the following proposition which gives an explicit relationship between the class-weights vectors  $\bar{\mathbf{w}}_\ell$  and the class-wise means.

**Proposition 0.8** (Expression of the weights for near-optimal representations). *Under Assumption 0.7, for sufficiently large  $p$  and letting  $\epsilon \rightarrow 0$  in the expression of the generalized labels in equation 4, the class-weight vectors are asymptotically proportional to the centred class-wise means as*

$$\bar{\mathbf{w}}_\ell = \frac{\gamma_\ell k e^{-\kappa \mu_1}}{1 + (k-1) e^{-\kappa \mu_1}} \left( \mathbf{m}_\ell - \frac{1}{k} \sum_{j=1}^k \mathbf{m}_j \right) + \mathcal{O}_{\|\cdot\|} (1) \quad (10)$$

$$\text{for some constant } \kappa > 0 \text{ and } \gamma_\ell = \left( \sum_{j=1}^k \sigma_{1,j}^2 \mathbb{E}_j[-f'_{\ell,i}(\bar{\mathbf{w}}_\ell^\top \mathbf{x}_i)] \right)^{-1} \geq \frac{4}{k \sigma_{1,\max}^2}.$$

*Proof.* See supplementary material.  $\square$

Proposition 0.8 provides an explicit link between  $\bar{\mathbf{w}}_\ell$  and the class-wise means of the representations when the latter satisfy the near-optimal assumption. Note that since  $\gamma_\ell \geq \frac{4}{k \sigma_{1,\max}^2} = \mathcal{O}(p^\epsilon)$ , the  $\bar{\mathbf{w}}_\ell$ 's are equal at the first order to the (scaled) centred class-wise means, as a consequence, the cosine similarity between them is asymptotically high as in the following corollary.

**Corollary 0.9.** *Denote the centred class-wise means by  $\bar{\mathbf{m}}_\ell \equiv \mathbf{m}_\ell - k^{-1} \sum_{j=1}^k \mathbf{m}_j$  and the scaling factor  $\varrho_\ell = \frac{\gamma_\ell k e^{-\kappa \mu_1}}{1 + (k-1) e^{-\kappa \mu_1}}$ , therefore  $\text{Cos}(\varrho_\ell^{-1} \bar{\mathbf{w}}_\ell, \bar{\mathbf{m}}_\ell) = 1 - \mathcal{O}(p^{-\epsilon})$ .*

We now turn to the experiments to show the application of this result to actual neural network architecture applied to the classification of image datasets and present the consequences to TL.

## Experiments

In all the experiments we consider the following settings.

**Datasets.** All the experiments performed on four datasets which are: MNIST (LeCun 1998), Fashion MNIST (Xiao, Rasul, and Vollgraf 2017), Cifar10 (Krizhevsky and Hinton 2010) and Cifar100 (Krizhevsky, Nair, and Hinton 2009). Each contains about  $n = 60000$  training images and 10000 test images belonging to  $k = 10$  classes, except Cifar100 which contains  $n = 50000$  training images belonging to  $k = 100$  classes. All the images are normalized in  $[0, 1]$ .

**Architecture.** We consider a CNN architecture for the representation  $\phi(\mathbf{d}; \Theta)$ . The dimension of the representation is taken to be  $p = 512$  and the activation function  $\varphi$  is either ReLU or tanh, and other choices of the activation are considered in the experiments part of the supplementary material. The classification layer is a dense layer with Softmax activation and without biases since their role is negligible (Kang et al. 2019). We use the Keras (Géron 2019) deep learning framework for the implementation. More details about the architecture are provided in the supplementary material.

**Training parameters.** Four CNN networks having the above architecture are trained on the four considered datasets for 50 epochs (except Cifar100 with 100 epochs) and using a batch-size of 1000 images, using the Adam optimizer (Kingma and Ba 2014).  $\mathbf{m}_\ell$  and  $\mathbf{C}_\ell$  are estimated by their empirical estimates.

### Near-optimal representations in practice

Under a specific choice of the activation function  $\varphi$  at the representation layer, it turns out that Assumption 0.7 can hold once the network has been trained. Figure 1 depicts the empirical values of the parameters  $\mu_1$ ,  $\mu_2$ ,  $\sigma_{1,j}^2$  and  $\sigma_{2,j}^2$  for two choices of the activation function at the representation layer. Here we considered the MNIST dataset and for each class, the means are estimated through their empirical estimates and the parameters  $\sigma_{1,j}^2$ ,  $\sigma_{2,j}^2$  are estimated respectively by the average of the diagonal and off-diagonal entries of the sample covariance matrix corresponding to class  $j$ .

As we can see in Figure 1, when considering the ReLU activation function at the representation layer, the parameters  $\mu_1 \approx \|\mathbf{m}_i\|^2$  and  $\mu_2 \approx \mathbf{m}_i^\top \mathbf{m}_j$  for  $i \neq j$  are of the same order which is not consistent with Assumption 0.7. Besides, when considering the tanh activation function,  $\mu_1 \gg \mu_2$  and  $\sigma_{1,j}^2 \gg \sigma_{2,j}^2$  which agrees with the near-optimal representation assumption. Therefore, when using the ReLU activation function the classes-separability information is shared between the class-wise means and covariances, while in the case of tanh all the information is encoded in the class-wise means. As a consequence, we will see in the next section that, when considering the tanh activation, the weights vectors of the Softmax classifier get aligned with the centred class-wise means as suggested by Corollary 0.9 under the near-optimal representations assumption.



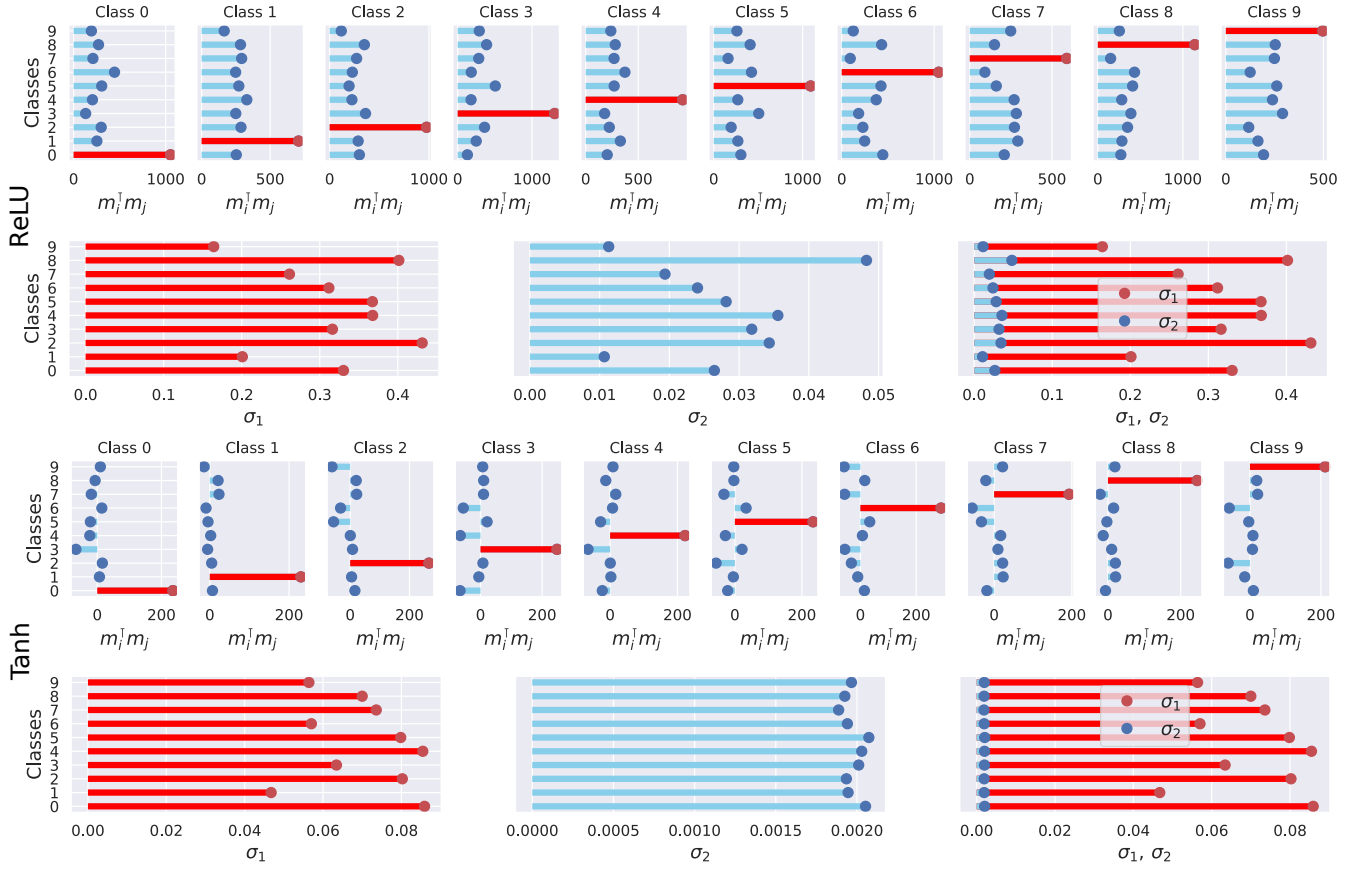


Figure 1: Empirical values of the parameters  $\mu_1$ ,  $\mu_2$ ,  $\sigma_{1,j}^2$  and  $\sigma_{2,j}^2$  as in Assumption 0.7. The top figures correspond to a ReLU activation in the representation layer and the bottom figures for a tanh activation. As we can see, Assumption 0.7 is verified when using the tanh activation. Indeed, the means satisfy  $\mu_1 \approx \|\mathbf{m}_i\|^2 \gg \mathbf{m}_i^T \mathbf{m}_j \approx \mu_2$  for  $i \neq j$  and the covariance scalars satisfy  $\sigma_{1,j}^2 \gg \sigma_{2,j}^2$ . While these conditions are not satisfied when using the ReLU activation (for which  $\mu_1 \approx \mu_2$  and  $\sigma_1^{\text{ReLU}} \approx 10 \times \sigma_1^{\text{tanh}}$ ). The CNN model has been trained on the MNIST dataset, and when using both activations (tanh or ReLU) we get the same test accuracy 99.6%.

The fact that we observe a difference in the behavior between different choices of the activation function at the representation layer may be associated to the symmetry properties of the considered activation function. Indeed, tanh is an odd function while ReLU does not have any symmetries, a deeper analysis of the interplay of the non-linearity and the data statistics has been done in (Liao and Couillet 2018) which provides insights in the case of random features maps. More examples as in Figure 1 using different activation functions and datasets are provided in the supplementary material.

Note however that the choice of the activation function at the representation layer does not affect much the performances of the model. Indeed, in the case of the MNIST dataset in Figure 1, the model has the same test accuracy of 99.6% for both choices of the activation function (ReLU or tanh). The same property hold for the datasets Fashion MNIST, Cifar10 and Cifar100 as provided in the supplementary material.

### Similarity between $\bar{\mathbf{w}}_\ell$ and $\bar{\mathbf{m}}_\ell$

Corollary 0.9 states that the cosine similarity between the learned  $\mathbf{w}_\ell$  and  $\bar{\mathbf{m}}_\ell$  is, at a scaling factor, asymptotically high when the representations statistics  $\mathbf{m}_\ell$  and  $\mathbf{C}_\ell$  satisfy the near-optimal representations assumption 0.7. As we discussed in the previous section, this assumption can be verified when the CNN model is trained for a specific choice of the representation activation  $\varphi$ . Figure 2 shows the empirically computed cosine similarities  $\text{Cos}(\mathbf{w}_\ell, \bar{\mathbf{m}}_\ell)$  during training for both activations ReLU and tanh. As we can notice,  $\text{Cos}(\mathbf{w}_\ell, \bar{\mathbf{m}}_\ell)$  is relatively high when considering the tanh activation since it yields to near-optimal representations as assumed in 0.7.

We also provide the curves of accuracy during training of the CNN model (in blue) along with the accuracy when replacing the weight matrix  $\mathbf{W}$  of the Softmax layer by the normalized matrix  $\bar{\mathbf{M}}$  containing the centred class-wise means  $\bar{\mathbf{m}}_\ell$  (in red). Precisely,  $\bar{\mathbf{M}} = \mathbf{M}\mathbf{P}/\|\mathbf{M}\mathbf{P}\|$  with  $\mathbf{M} = [\mathbf{m}_1, \dots, \mathbf{m}_k]$  and  $\mathbf{P}$  stands for the centring matrix

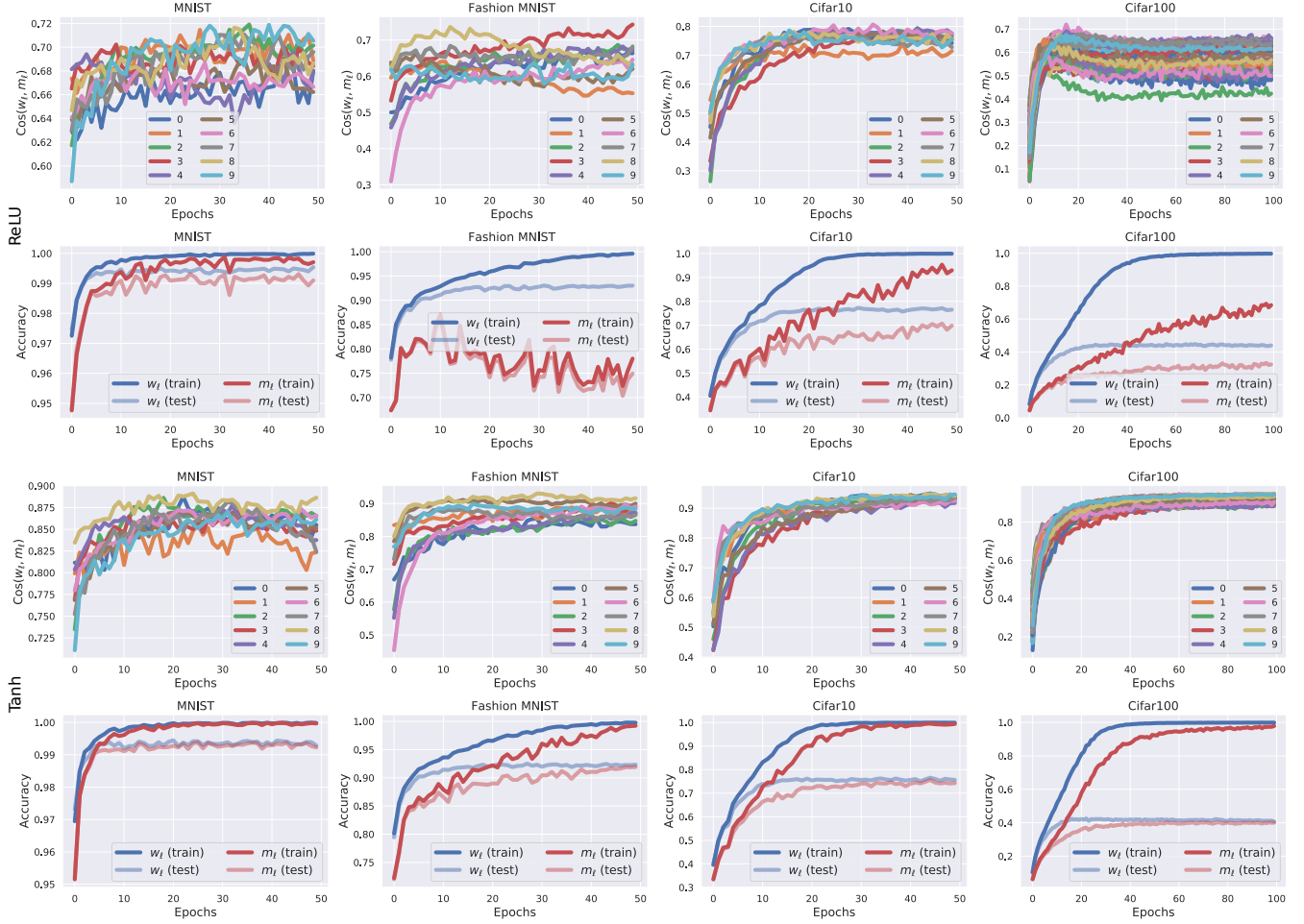


Figure 2: The first and third rows show, respectively for ReLU and tanh activation at the representation, the cosine similarity during the training between the learned weights  $w_\ell$  and the centred class-wise means  $\bar{m}_\ell$  for different classes and across different datasets (columns). The second and fourth rows show the train and test accuracy during training, when considering the learned weights (in blue) and the centred class-wise means (in red). As we can see, the cosine similarity between  $w_\ell$  and  $\bar{m}_\ell$  for tanh is higher than ReLU since Assumption 0.7 holds when using tanh (see Figure 1).

of dimension  $k$  and is given by  $P = I_k - \frac{1}{k} \mathbf{1}_k \mathbf{1}_k^\top$ . The normalization by  $\|MP\|$  is considered since the scaling factor  $\varrho_\ell$  depends implicitly on  $w_\ell$ , but still one can use spectral normalization (Miyato et al. 2018) as we considered here by the term  $\|MP\|$ .

We can see in Figure 2 that, for  $\varphi = \tanh$ , there is a matching between the train and test accuracy when using the default weights  $W$  and when replacing them by  $\bar{M}$  at the end of the training phase. Therefore, this result suggests that the weight matrix  $W$  encode the class-wise means of the representations, and hence the neural network classifier make decisions through the  $\bar{m}_\ell$ 's. The similarity between  $\bar{w}_\ell$  and  $\bar{m}_\ell$  notably supports the remarkable observed performances of the Nearest Class Mean Classifier (Guerriero, Caputo, and Mensink 2018), which relies on classifying through the euclidean distances to the class-wise means. Note that the performance of the CNN model for  $\varphi = \text{ReLU}$  is the same as when considering  $\varphi = \tanh$ , but the latter

has the advantage to be more interpretable through Corollary 0.9.

### Consequences for transfer learning

As a consequence of our finding in the previous section, in a deep transfer learning scheme (Bengio, Courville, and Vincent 2013; Tamaazousti et al. 2019), one can initialize the weights of the Softmax layer with the scaled centred class-wise means  $\bar{M}$  of the target data representations. To stress this point, we first train a network on the used datasets considered as source domain and the target domain is considered as noisy versions of the input data ( $n = 5000$  randomly samples) by adding Gaussian noise with variance  $\sigma^2$ . We particularly compare three methods which are: CM (Centred Means) the weights are replaced with the scaled centred means  $\bar{M}$  without further training; LWM (Learned Weights of Means) a learned weighted sum of means  $MW$  with  $W$  a  $k \times k$  weight matrix to be trained; WFS (Weights Learned

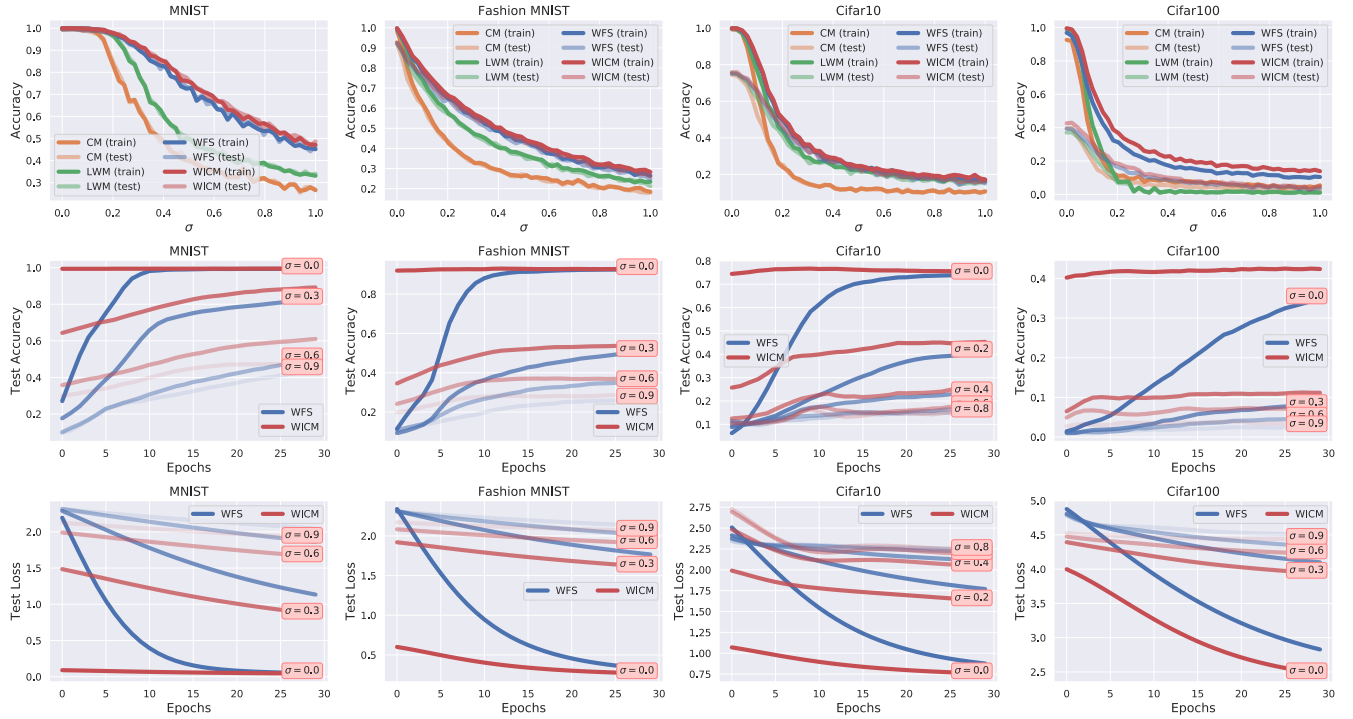


Figure 3: Training of the Softmax layer on target data. The target data ( $n = 5000$  sampled at random) being noisy versions of the considered datasets with different Gaussian noise variance  $\sigma^2$ . Accuracy (**Top**) and dynamics (**Bottom**) of different baselines in terms of the added noise variance  $\sigma^2$ . The different methods are: CM (Centred Means) the weights are replaced with the scaled centred means; LWM (Learned Weights of Means) a learned weighted sum of means; WFS (Weights Learned From Scratch); WICM (Weights Initialized with Centred Means).

From Scratch); WICM (Weights Initialized with Centred Means).

Figure 3 depicts the accuracy on the target domain (along the different datasets) in terms of the noise variance  $\sigma^2$ . We particularly notice from this figure that when the weights of the Softmax layer are initialized with the class-wise means of the representations (WICM), we get a slightly better accuracy than initializing them randomly (WFS). Moreover, internalizing with the class-wise means yield to faster convergence as the target domain gets closer to the source one, i.e., as  $\sigma \rightarrow 0$ . Besides, the source model selection can be performed by comparing the accuracies, of different representations, of the Softmax layer with the weights being the centred class-wise means of the representations (CM method in Figure 3), thereby not requiring the training of the Softmax layer multiple times.

## Conclusion

We have analyzed throughout the paper the behavior of the Softmax layer in neural networks classifiers. We have notably shown that the weight vectors corresponding to each class of the Softmax layer tend to converge to the class-wise means computed at the representation layer, a result which unfolds from the near-optimal representations assumption. Our findings suggested three main procedures for efficient transfer learning: (i) use of symmetric representation activa-

tions to ensure the near-optimal representations assumption; (ii) we provided a systematic approach to perform source model selection without training the Softmax layer; (iii) and finally we provided an initialization procedure which accelerates the training of the Softmax layer as the target domain gets closer to the source domain.

## References

- Azizpour, H.; Razavian, A. S.; Sullivan, J.; Maki, A.; and Carlsson, S. 2015. Factors of transferability for a generic convnet representation. *IEEE transactions on pattern analysis and machine intelligence*, 38(9): 1790–1802.
- Bengio, Y.; Courville, A.; and Vincent, P. 2013. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8): 1798–1828.
- Cai, H.; Zhu, L.; and Han, S. 2019. Proxylessnas: Direct neural architecture search on target task and hardware. *ICLR*.
- Géron, A. 2019. *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems*. O’Reilly Media.
- Ginsca, A. L.; Popescu, A.; Le Borgne, H.; Ballas, N.; Vo, P.; and Kanellos, I. 2015. Large-scale image mining with

- flickr groups. In *International Conference on Multimedia Modeling*, 318–334. Springer.
- Guerriero, S.; Caputo, B.; and Mensink, T. 2018. Deep nearest class mean classifiers. In *International Conference on Learning Representations, Workshop Track*.
- Kang, B.; Xie, S.; Rohrbach, M.; Yan, Z.; Gordo, A.; Feng, J.; and Kalantidis, Y. 2019. Decoupling Representation and Classifier for Long-Tailed Recognition. *arXiv:1910.09217*.
- Kingma, D. P.; and Ba, J. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Kornblith, S.; Shlens, J.; and Le, Q. V. 2019. Do better imagenet models transfer better? In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2661–2671.
- Krizhevsky, A.; and Hinton, G. 2010. Convolutional deep belief networks on cifar-10. *Unpublished manuscript*, 40(7): 1–9.
- Krizhevsky, A.; Nair, V.; and Hinton, G. 2009. Cifar-10 and cifar-100 datasets. URL: <https://www.cs.toronto.edu/kriz-cifar.html>, 6.
- Kuzborskij, I.; Orabona, F.; and Caputo, B. 2013. From  $n$  to  $n+1$ : Multiclass transfer incremental learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 3358–3365.
- LeCun, Y. 1998. The MNIST database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>.
- Liao, Z.; and Couillet, R. 2018. On the spectrum of random features maps of high dimensional data. *arXiv preprint arXiv:1805.11916*.
- Liu, H.; Simonyan, K.; and Yang, Y. 2019. Darts: Differentiable architecture search. *ICLR*.
- Miyato, T.; Kataoka, T.; Koyama, M.; and Yoshida, Y. 2018. Spectral normalization for generative adversarial networks. *arXiv preprint arXiv:1802.05957*.
- Real, E.; Aggarwal, A.; Huang, Y.; and Le, Q. V. 2019. Regularized evolution for image classifier architecture search. In *Proceedings of the aaai conference on artificial intelligence*, volume 33, 4780–4789.
- Seddik, M. E. A.; Louart, C.; Tamaazousti, M.; and Couillet, R. 2020. Random Matrix Theory Proves that Deep Learning Representations of GAN-data Behave as Gaussian Mixtures. *arXiv preprint arXiv:2001.08370*.
- Sharif Razavian, A.; Azizpour, H.; Sullivan, J.; and Carlsson, S. 2014. CNN features off-the-shelf: an astounding baseline for recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 806–813.
- Tamaazousti, Y.; Le Borgne, H.; Hudelot, C.; Seddik, M. E. A.; and Tamaazousti, M. 2019. Learning more universal representations for transfer-learning. *IEEE transactions on pattern analysis and machine intelligence*.
- Xiao, H.; Rasul, K.; and Vollgraf, R. 2017. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*.
- Yosinski, J.; Clune, J.; Bengio, Y.; and Lipson, H. 2014. How transferable are features in deep neural networks? In *Advances in neural information processing systems*, 3320–3328.
- Zoph, B.; and Le, Q. V. 2017. Neural architecture search with reinforcement learning. *ICLR*.
- Zoph, B.; Vasudevan, V.; Shlens, J.; and Le, Q. V. 2018. Learning transferable architectures for scalable image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 8697–8710.