

SEMI-SUPERVISED DEEP LEARNING BY METRIC EMBEDDING

Elad Hoffer

Technion - Israel Institute of Technology
Haifa, Israel
ehoffer@tx.technion.ac.il

Nir Ailon

Technion - Israel Institute of Technology
Haifa, Israel
nailon@cs.technion.ac.il

ABSTRACT

Deep networks are successfully used as classification models yielding state-of-the-art results when trained on a large number of labeled samples. These models, however, are usually much less suited for semi-supervised problems because of their tendency to overfit easily when trained on small amounts of data. In this work we will explore a new training objective that is targeting a semi-supervised regime with only a small subset of labeled data. This criterion is based on a deep metric embedding over distance relations within the set of labeled samples, together with constraints over the embeddings of the unlabeled set. The final learned representations are discriminative in euclidean space, and hence can be used with subsequent nearest-neighbor classification using the labeled samples.

1 INTRODUCTION

Deep neural networks have been shown to perform very well on various classification problems, often yielding state-of-the-art results. Key motivation for the use of these models, is the assumption of hierarchical nature of the underlying problem. This assumption is reflected in the structure of NNs, composed of multiple stacked layers of linear transformations followed by non-linear activation functions. The NN final layer is usually a softmax activated linear transformation indicating the likelihood of each class, which can be trained by cross-entropy using the known target of each sample, and back-propagated to previous layers. The hierarchical property of NNs has been observed to yield high-quality, discriminative representations of the input in intermediate layers. These representative features, although not explicitly part of the training objective, were shown to be useful in subsequent tasks in the same domain as demonstrated by Razavian et al. (2014). One serious problem occurring in neural network is their susceptibility to overfit over the training data. Due to this fact, a considerable part of modern neural network research is devoted to regularization techniques and heuristics such as Srivastava et al. (2014); Ioffe & Szegedy (2015); Wan et al. (2013); Szegedy et al. (2015), to allow the networks to generalize to unseen data samples. The tendency to overfit is most apparent with problems having a very small number of training examples per class, and these are considered ill-suited to solve with neural network models. Because of this property, semi-supervised regimes in which most data is unlabeled, are considered hard to learn and generalize with NNs.

In this work we will consider a new training criterion designed to be used with deep neural networks in semi-supervised regimes over datasets with a small subset of labeled samples. Instead of a usual cross-entropy between the labeled samples and the ground truth class indicators, we will use the labeled examples as targets for a metric embedding. Under this embedding, which is the mapping of a parameterized deep network, the features of labeled examples will be grouped together in euclidean space. In addition, we will use these learned embeddings to separate the unlabeled examples to belong each to a distinct cluster formed by the labeled samples. We will show this constraint translates to a minimum entropy criterion over the embedded distances. Finally, because of the use of euclidean space interpretation of the learned features, we are able to use a subsequent nearest-neighbor classifier to achieve state-of-the-art results on problems with small number of labeled examples.

2 RELATED WORK

2.1 LEARNING METRIC EMBEDDING

Previous works have shown the possible use of neural networks to learn useful metric embedding. One kind of such metric embedding is the “Siamese network” framework introduced by Bromley et al. (1993) and later used in the works of Chopra et al. (2005). One use for this methods is when the number of classes is too large or expected to vary over time, as in the case of face verification, where a face contained in an image has to be compared against another image of a face. This problem was recently tackled by Schroff et al. (2015) for training a convolutional network model on triplets of examples. Learning features by metric embedding was also shown by Hoffer & Ailon (2015) to provide competitive classification accuracy compare to conventional cross-entropy regression. This work is also related to Rippel et al. (2015), who introduced Magnet loss - a metric embedding approach for fine-grained classification. The Magnet loss is based on learning the distribution of distances for each sample, from K clusters assigned for each classified class. It then uses an intermediate k-means clustering, to reposition the different assigned clusters. This proved to allow better accuracy than both margin-based Triplet loss, and softmax regression. Using metric embedding with neural network was also specifically shown to provide good results in the semi-supervised learning setting as seen in Weston et al. (2012).

2.2 SEMI-SUPERVISED LEARNING BY ADVERSARIAL REGULARIZATION

As stated before, a key approach to generalize from a small training set, is by regularizing the learned model. Regularization techniques can often be interpreted as prior over model parameters or structure, such as L_p regularization over the network weights or activations. More recently, neural network specific regularizations that induce noise within the training process such as Srivastava et al. (2014); Wan et al. (2013); Szegedy et al. (2015) proved to be highly beneficial to avoid overfitting. Another recent observation by Goodfellow et al. (2015) is that training on adversarial examples, inputs that were found to be misclassified under small perturbation, can improve generalization. This fact was explored by Feng et al. (2016) and found to provide notable improvements to the semi supervised regime by Miyato et al. (2015).

2.3 SEMI-SUPERVISED LEARNING BY AUXILIARY RECONSTRUCTION LOSS

Recently, a stacked set of denoising auto-encoders architectures showed promising results in both semi-supervised and unsupervised tasks. A stacked what-where autoencoder by Zhao et al. (2015) computes a set of complementary variables that enable reconstruction whenever a layer implements a many-to-one mapping. Ladder networks by Rasmus et al. (2015) - use lateral connections to allow higher levels of an auto-encoder to focus on invariant abstract features by applying a layer-wise cost function.

Generative adversarial network (GAN) is a recently introduced model that can be used in an unsupervised fashion Goodfellow et al. (2014). Adversarial Generative Models use a set of networks, one trained to discriminate between data sampled from the true underlying distribution (e.g., a set of images), and a separate generative network trained to be an adversary trying to confuse the first network. By propagating the gradient through the paired networks, the model learns to generate samples that are distributed similarly to the source data. As shown by Radford et al. (2015), this model can create useful latent representations for subsequent classification tasks. The usage for these models for semi-supervised learning was further developed by Springenberg (2016) and Salimans et al. (2016), by adding a $N + 1$ way classifier (number of classes + and additional “fake” class) to the discriminator. This proved to allow excellent accuracy with only a small subset of labeled examples.

2.4 SEMI-SUPERVISED LEARNING BY ENTROPY MINIMIZATION

Another technique for semi-supervised learning introduced by Grandvalet & Bengio (2004) is concerned with minimizing the entropy over expected class distribution for unlabeled examples. Regularizing for minimum entropy can be seen as a prior which prefers minimum overlap between observed classes. This can also be seen as a generalization of the “self-training” wrapper method

described by Triguero et al. (2015), in which unlabeled examples are re-introduced after being labeled with the previous classification of the model. This is also related to the “Transductive support vector machines” (TSVM) Vapnik & Vapnik (1998) which introduces a maximum margin objective over both labeled and unlabeled examples.

Comparing with previous works such as Chopra et al. (2005) and Weston et al. (2012), this work uses a novel objective which is composed of a distance ratio measure (unlike the contrastive, hinge based loss used before), and an entropy minimization on the distance measure to labeled samples. Although distance ratio loss (Hoffer & Ailon (2015)) and entropy minimization (Grandvalet & Bengio (2004)) are not new, this is the first attempt to our knowledge, of combining these ideas for semi-supervised metric learning.

3 OUR CONTRIBUTION: NEIGHBOR EMBEDDING FOR SEMI-SUPERVISED LEARNING

In this work we are concerned with a semi-supervised setting, in which learning is done on data of which only a small subset is labeled. Given observed sets of labeled data $X_L = \{(x, y)\}_{i=1}^l$ and unlabeled data $X_U = \{x\}_{i=l+1}^n$ where $x \in \mathcal{X}$, $y \in \mathcal{C}$, we wish to learn a classifier $f : \mathcal{X} \rightarrow \mathcal{C}$ to have a minimum expected error on some unseen test data X_{test} .

We will make a couple of assumptions regarding the given data:

- The number of labeled examples is small compared to the whole observed set $l \ll n$.
- Structure assumption - samples within the same structure (such as a cluster or manifold) are more likely to share the same label. This assumption is shared with many other semi-supervised approaches as discussed in Chapelle et al. (2009), Weston et al. (2012).

Using these assumptions, we are motivated to learn a metric embedding that forms clusters such that samples can be classified by their L_2 distance to the labeled examples in a nearest-neighbor procedure.

We will now define our learning setting on the semi-labeled data, using a neural network model denoted as $F(x; \theta)$ where x is the input fed into the network, and θ are the optimized parameters (dropped henceforward for convenience). The output of the network for each sample is a vector of features of D dimensions $F(x) \in \mathbb{R}^D$ which will be used to represent the input.

Our two training objectives which we aim to train our embedding networks by are:

- Create feature representation that form clusters from the labeled examples $\{(x, y)\} \in X_L$ such that two examples x_1, x_2 sharing the same label $y_1 = y_2$ will have a smaller embedded distance than any third example x_3 with a different label $y_1 \neq y_3$

$$\|F(x_1) - F(x_2)\|_2 < \|F(x_1) - F(x_3)\|_2$$

- For each unlabeled example, its feature embedding will be close to the embeddings of one specific label occurring in L :

For all $x \in X_U$, $z \in X_L$, there exists a specific class $l \in \mathcal{C}$ such that

$$\|F(x) - F(z_l)\|_2 \ll \|F(x) - F(z_k)\|_2$$

where z_l is any labeled example of class l and z_k is any example from class $k \in \mathcal{C} \setminus \{l\}$.

As the defined objectives create embeddings that target a nearest-neighbor classification with regard to the labeled set, we will refer to it as “Neighbor embedding”.

4 LEARNING BY DISTANCE COMPARISONS

We will define a discrete distribution for the embedded distance between a sample $x \in \mathcal{X}$, and c labeled examples $z_1, \dots, z_c \in X_L$ each belonging to a different class:

$$P(x; z_1, \dots, z_c)_i = \frac{e^{-\|F(x) - F(z_i)\|^2}}{\sum_{j=1}^c e^{-\|F(x) - F(z_j)\|^2}}, i \in \{1 \dots c\} \quad (1)$$

This definition assigns a probability $P(x; z_1, \dots, z_c)_i$ for sample x to be classified into class i , under a 1-nn classification rule, when z_1, \dots, z_c neighbors are given. It is similar to the stochastic-nearest-neighbors formulation of Goldberger et al. (2004), and will allow us to state the two underlying objectives as measures over this distribution.

4.1 DISTANCE RATIO CRITERION

Addressing objective (i), we will use a sample $x_l \in X_L$ from the labeled set belonging to class $k \in \mathcal{C}$, and another set of sampled labeled examples $z_1, \dots, z_c \in X_L$. In this work we will sample in uniform over all available samples for each class.

Defining the class-indicator $I(x)$ as

$$I(x_l)_i = \begin{cases} 1 & \text{if } i = k \\ 0 & \text{otherwise} \end{cases}$$

we will minimize the cross-entropy between $I(x_l)$ and the distance-distribution of x with respect to z_1, \dots, z_c :

$$L(x_l, z_1, \dots, z_c)_L = H(I(x_l), P(x_l; z_1, \dots, z_c)) \quad (2)$$

This is in fact a slightly modified version of distance ratio loss introduced in Hoffer & Ailon (2015).

$$L(x_l, z_1, \dots, z_c)_L = -\log \frac{e^{-\|F(x_l) - F(z_k)\|^2}}{\sum_{i=1}^c e^{-\|F(x_l) - F(z_i)\|^2}} \quad (3)$$

This loss is aimed to ensure that samples belonging to the same class will be mapped to have a small embedded distance compared to samples from different classes.

4.2 MINIMUM ENTROPY CRITERION

Another part of the optimized criterion, inspired by Grandvalet & Bengio (2004), is designed to reduce the overlap between the different classes of the unlabeled samples.

We will promote this objective by minimizing the entropy of the underlying distance distribution of x , again with respect to labeled samples z_1, \dots, z_c :

$$L(x, z_1, \dots, z_c)_U = H(P(x; z_1, \dots, z_c)) \quad (4)$$

which is defined as

$$L(x, z_1, \dots, z_c)_U = -\sum_{i=1}^c \frac{e^{-\|F(x) - F(z_i)\|^2}}{\sum_{j=1}^c e^{-\|F(x) - F(z_j)\|^2}} \cdot \log \frac{e^{-\|F(x) - F(z_i)\|^2}}{\sum_{j=1}^c e^{-\|F(x) - F(z_j)\|^2}} \quad (5)$$

We note that entropy is lower if the distribution is sparse, and higher if the distribution is dense, and this intuition is compatible with our objectives.

Our final objective will use a sampled set of labeled examples, where each class is represented $\{z_1, \dots, z_c\}$ and additional labeled x_l and unlabeled x_u examples, combining a weighted sum of both 3 and 5 to form:

$$L(x_l, x_u, \{z_1, \dots, z_c\}) = \lambda_L L(x_l, z_1, \dots, z_c)_L + \lambda_U L(x_u, z_1, \dots, z_c)_U \quad (6)$$

Where $\lambda_L, \lambda_U \in [0, 1]$ are used to determine the weight assigned to each criterion.

This loss is differentiable and hence can be used for gradient-based training of deep models by existing optimization approaches and back-propagation (Rumelhart et al.) through the embedding neural network. The optimization can further be accelerated computationally by using mini-batches of both labeled and unlabeled examples.

5 QUALITIES OF NEIGHBOR EMBEDDING

We will now discuss some observed properties of neighbor embeddings, and their usefulness to semi-supervised regimes using neural network models.

5.1 REDUCING OVERFIT

Usually, when using NNs for classification, a cross-entropy loss minimization is employed by using a fixed one-hot indicator (similar to 2) as target for each labeled example, thus maximizing a log-likelihood of the correct label. This form of optimization over a fixed target tend to cause an overfitting of the neural-network, especially on small labeled sets. This was lately discussed and addressed by Szegedy et al. (2015) using added random noise to the targets by sampling uniformly from the set of classes, effectively smoothing the cross-entropy target distribution. This regularization technique was shown empirically to yield better generalization by reducing the overfitting over the training set.

Training on distance ratio comparisons, as shown in our work, provides a natural alternative to this problem. By setting the optimization target to be the embeddings of labeled examples, we create a continuously moving target that is dependent on the current model parameters. We speculate that this reduces the model’s ability to overfit easily on the training data, allowing very small labeled datasets to be exploited.

5.2 EMBEDDING INTO EUCLIDEAN SPACE

By training the model to create feature embedding that are discriminative with respect to their distance in euclidean space, we can achieve good classification accuracy using a simple nearest-neighbor classifier. This embedding allows an interpretation of semantic relation in euclidean space, which can be useful for various tasks such as information retrieval, or transfer learning.

5.3 COMBINING SUPERVISED AND UNSUPERVISED OBJECTIVES

Neighbor embedding is composed of both supervised 3 and unsupervised 5 objectives that are weighted by λ_L, λ_U coefficients. This can be used to balance or possibly anneal over time (Zamora-Martínez et al. (2016)) to adapt for the availability of labeled samples. This form of balancing was found previously to allow for better representation learning using unlabeled data.

5.4 INCORPORATING PRIOR KNOWLEDGE

We also note that prior knowledge about a problem at hand can be incorporated into the expected measures with respect to the distance distribution 1. E.g, knowledge of relative distance between classes can be used to replace $I(x)$ as target distribution in eq. 3 and knowledge concerning overlap between classes can be used to relax the constraint in eq. 5.

6 EXPERIMENTS

All experiments were conducted using the Torch7 framework by Collobert et al. (2011). Code reproducing these results will be available at <https://github.com/eladhoffer/SemiSupContrast>. For every experiment we chose a small random subset of examples, with a balanced number from each class and denoted by X_L . The remaining training images are used without their labels to form X_U . Finally, we test our final accuracy with a disjoint set of examples X_{test} . No data augmentation was applied to the training sets.

In each iteration we sampled uniformly a set of labeled examples $z_1, \dots, z_{|C|} \in X_L$. In addition, batches of uniformly sampled examples were also sampled again from the labeled set X_L , and the unlabeled set X_U .

A batch-size of $b = 32$ was used for all experiments, totaling a sampled set of $2 \cdot b + |C|$ examples for each iteration, where $|C| = 10$ for both datasets. We used 6 as optimization criterion, where $\lambda_L = \lambda_U = 1$. Optimization was done using the Accelerated-gradient method by Nesterov (1983)

Table 1: Results for MNIST. Using 100 labeled examples, no data-augmentation.

Model	Test error %
EmbedCNN Weston et al. (2012)	7.75
SWWAE Zhao et al. (2015)	9.17
Ladder network Rasmus et al. (2015)	0.89 (\pm 0.50)
Conv-CatGAN Springenberg (2016)	1.39 (\pm 0.28)
Ours	0.78 (\pm 0.3)

with an initial learning rate of $lr_0 = 0.1$ which was decreased by a factor of 10 after every 30 epochs. Both datasets were trained on for a total of 90 epochs. Final test accuracy results was achieved by using a k-NN classifier with best results out of $k = \{1, 3, 5\}$. These results were average over 10 random subsets of labeled data. The choices for λ and k-NN parameters were made using a validation set. We did not found any substantial difference between the values we explored, so they were usually left as the default value for simplicity.

As the embedding model was chosen to be a convolutional network, the spatial properties of input space are crucial. We thus omit results on permutation-invariant versions of these problems, noting they usually tend to achieve worse classification accuracies. We also note that the neural network models themselves are very simple to ensure performance achieved is due to the proposed objective, and not the network architecture.

6.1 RESULTS ON MNIST

The MNIST database of handwritten digits introduced by LeCun et al. (1998) is one of the most studied dataset benchmark for image classification. The dataset contains 60,000 examples of handwritten digits from 0 to 9 for training and 10,000 additional examples for testing, where each sample is a 28 x 28 pixel gray level image.

We followed previous works ((Weston et al., 2012),(Zhao et al., 2015),Rasmus et al. (2015)) and used semi-supervised regime in which only 100 samples (10 for each class) were used as X_L along with their labels. For the embedding network, we used a convolutional network with 5-convolutional layers, where each layer is followed by a ReLU non-linearity and batch-normalization layer Ioffe & Szegedy (2015). The full network structure is described in Appendix table 3. Results are displayed in table 1 and reflect that our approach yields state-of-the-art results in this regime.

We also attempted to visualize the outcome of using this method, by training an additional model with a final 2-dimensional embedding. Figure 1b shows the final embeddings, where labeled examples are marked in color with their respective class, and unlabeled examples are marked in gray. We can see that, in accordance with our objectives, the labeled examples formed clusters in euclidean space separate by their labels, while unlabeled examples were largely grouped to belong each to one of these clusters.

6.2 RESULTS ON CIFAR-10

Cifar-10 introduced by Krizhevsky & Hinton (2009) is an image classification benchmark dataset containing 50,000 training images and 10,000 test images. The image sizes 32×32 pixels, with color. The classes are airplanes, automobiles, birds, cats, deer, dogs, frogs, horses, ships and trucks.

Following a commonly used regime, we trained on 4000 randomly picked samples (400 for each class). As the convolutional embedding network, we used a network similar to that of Lin et al. (2013) which is described in table 3. The test error results are brought in table 2.

As can be observed, we achieve competitive results with state-of-the-art in this regime. We also note that current best results are from generative models such as Springenberg (2016) and Salimans et al. (2016) that follow an elaborate and computationally heavy training procedure compared with our approach.

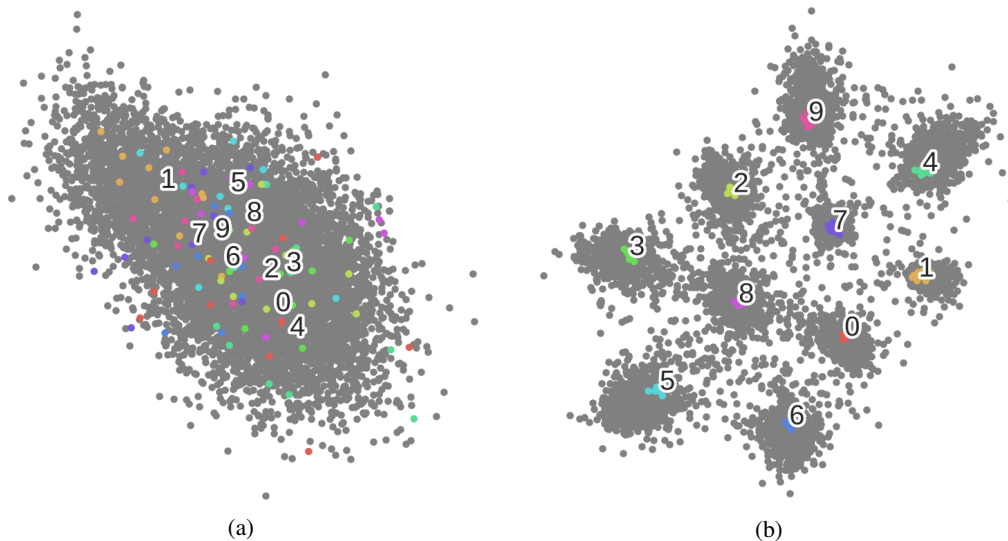


Figure 1: MNIST 2d visualization before (a) and after (b) training. 100 colored labeled samples, unlabeled samples marked in gray

Table 2: Results for Cifar-10. Using 4000 labeled samples, no data-augmentation.

Model	Test error %
Spike-and-Slab Sparse Coding Goodfellow et al. (2012)	31.9
View-Invariant k-means Hui (2013)	27.4 (± 0.7)
Exemplar-CNN Dosovitskiy et al. (2014)	23.4 (± 0.2)
Ladder network Rasmus et al. (2015)	20.04 (± 0.47)
Conv-CatGan Springenberg (2016)	19.58 (± 0.58)
ImprovedGan Salimans et al. (2016)	18.63 (± 2.32)
Ours	20.3 (± 0.5)

7 CONCLUSIONS

In this work we have shown how neural networks can be used to learn in a semi-supervised setting using small sets of labeled data, by replacing the classification objective with a metric embedding one. We introduced an objective for semi-supervised learning formulated as minimization of entropy over a distance encoding distribution. This objective is compliant with standard techniques of training deep neural network and requires no modification of the embedding model. Using the method in this work, we were able to achieve state-of-the-art results on MNIST with only 100 labeled examples and competitive results on Cifar10 dataset. We speculate that this form of learning is beneficial to neural network models by decreasing their tendency to overfit over small sets of training data. The objectives formulated here can potentially leverage prior knowledge on the distribution of classes or samples, as well as incorporating this knowledge in the training process. For example, utilizing the learned embedded distance, we speculate that a better sampling can be done instead of a uniform one over the entire set.

Further exploration is needed to apply this method to large scale problems, spanning a large number of available classes, which we leave to future work.

REFERENCES

Jane Bromley, James W Bentz, Léon Bottou, Isabelle Guyon, Yann LeCun, Cliff Moore, Eduard Säckinger, and Roopak Shah. Signature verification using a siamese time delay neural network. *International Journal of Pattern Recognition and Artificial Intelligence*, 7(04):669–688, 1993.

- Olivier Chapelle, Bernhard Scholkopf, and Alexander Zien. Semi-supervised learning (chapelle, o. et al., eds.; 2006)[book reviews]. *IEEE Transactions on Neural Networks*, 20(3):542–542, 2009.
- Sumit Chopra, Raia Hadsell, and Yann LeCun. Learning a similarity metric discriminatively, with application to face verification. In *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*, volume 1, pp. 539–546. IEEE, 2005.
- Ronan Collobert, Koray Kavukcuoglu, and Clément Farabet. Torch7: A matlab-like environment for machine learning. In *BigLearn, NIPS Workshop*, number EPFL-CONF-192376, 2011.
- Alexey Dosovitskiy, Jost Tobias Springenberg, Martin Riedmiller, and Thomas Brox. Discriminative unsupervised feature learning with convolutional neural networks. In *Advances in Neural Information Processing Systems*, pp. 766–774, 2014.
- Jiashi Feng, Tom Zahavy, Bingyi Kang, Huan Xu, and Shie Mannor. Ensemble robustness of deep learning algorithms. *arXiv preprint arXiv:1602.02389*, 2016.
- Jacob Goldberger, Geoffrey E Hinton, Sam T Roweis, and Ruslan Salakhutdinov. Neighbourhood components analysis. In *Advances in neural information processing systems*, pp. 513–520, 2004.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems*, pp. 2672–2680, 2014.
- Ian J. Goodfellow, Aaron Courville, and Yoshua Bengio. Large-scale feature learning with spike-and-slab sparse coding. 2012.
- Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. 2015.
- Yves Grandvalet and Yoshua Bengio. Semi-supervised learning by entropy minimization. In *Advances in neural information processing systems*, pp. 529–536, 2004.
- Elad Hoffer and Nir Ailon. Deep metric learning using triplet network. In *Similarity-Based Pattern Recognition*, pp. 84–92. Springer, 2015.
- Ka Y Hui. Direct modeling of complex invariances for visual object features. In *Proceedings of the 30th International Conference on Machine Learning (ICML-13)*, pp. 352–360, 2013.
- Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *Proceedings of The 32nd International Conference on Machine Learning*, pp. 448–456, 2015.
- Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. *Computer Science Department, University of Toronto, Tech. Rep*, 2009.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Min Lin, Qiang Chen, and Shuicheng Yan. Network in network. *ICLR2014*, 2013.
- Takeru Miyato, Shin-ichi Maeda, Masanori Koyama, Ken Nakae, and Shin Ishii. Distributional smoothing by virtual adversarial examples. *stat*, 1050:2, 2015.
- Yurii Nesterov. A method of solving a convex programming problem with convergence rate $o(1/k^2)$. 1983.
- Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.
- Antti Rasmus, Mathias Berglund, Mikko Honkala, Harri Valpola, and Tapani Raiko. Semi-supervised learning with ladder networks. In *Advances in Neural Information Processing Systems*, pp. 3532–3540, 2015.

- Ali Razavian, Hossein Azizpour, Josephine Sullivan, and Stefan Carlsson. Cnn features off-the-shelf: an astounding baseline for recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 806–813, 2014.
- Oren Rippel, Manohar Paluri, Piotr Dollar, and Lubomir Bourdev. Metric learning with adaptive density discrimination. *stat*, 1050:18, 2015.
- David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back-propagating errors. *Cognitive modeling*, 5(3):1.
- Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. *arXiv preprint arXiv:1606.03498*, 2016.
- Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 815–823, 2015.
- Jost Tobias Springenberg. Unsupervised and semi-supervised learning with categorical generative adversarial networks. In *International Conference on Learning Representations (ICLR)*. 2016. URL <https://arxiv.org/abs/1511.06390>.
- Nitish Srivastava, Geoffrey E Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15(1):1929–1958, 2014.
- Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. *arXiv preprint arXiv:1512.00567*, 2015.
- Isaac Triguero, Salvador García, and Francisco Herrera. Self-labeled techniques for semi-supervised learning: taxonomy, software and empirical study. *Knowledge and Information Systems*, 42(2): 245–284, 2015.
- Vladimir Naumovich Vapnik and Vlamimir Vapnik. *Statistical learning theory*, volume 1. Wiley New York, 1998.
- Li Wan, Matthew Zeiler, Sixin Zhang, Yann L Cun, and Rob Fergus. Regularization of neural networks using dropconnect. In *Proceedings of the 30th International Conference on Machine Learning (ICML-13)*, pp. 1058–1066, 2013.
- Jason Weston, Frédéric Ratle, Hossein Mobahi, and Ronan Collobert. Deep learning via semi-supervised embedding. In *Neural Networks: Tricks of the Trade*, pp. 639–655. Springer, 2012.
- Francisco Zamora-Martínez, Javier Muñoz-Almaraz, and Juan Pardo. Integration of unsupervised and supervised criteria for deep neural networks training. In *International Conference on Artificial Neural Networks*, pp. 55–62. Springer, 2016.
- Junbo Zhao, Michael Mathieu, Ross Goroshin, and Yann Lecun. Stacked what-where auto-encoders. *arXiv preprint arXiv:1506.02351*, 2015.

8 APPENDIX

Table 3: Convolutional models - (feature-maps, kernel, stride, padding) for each layer. Convolutional layers are each followed by ReLU and Batch-norm.

Model	
MNIST	Cifar-10
Input: 28×28 monochrome	Input: 32×32 RGB
Conv-ReLU-BN (16, 5x5, 1x1, 1x1)	Conv-ReLU-BN (192, 5x5, 1x1, 2x2)
Max-Pooling (2x2, 2x2)	Conv-ReLU-BN (160, 1x1, 1x1)
Conv-ReLU-BN (32, 3x3, 1x1, 1x1)	Conv-ReLU-BN (96, 1x1, 1x1)
Conv-ReLU-BN (64, 3x3, 1x1, 1x1)	Max-Pooling (3x3, 2x2)
Conv-ReLU-BN (64, 3x3, 1x1, 1x1)	Conv-ReLU-BN (96, 5x5, 1x1, 2x2)
Max-Pooling (2x2, 2x2)	Conv-ReLU-BN (192, 1x1, 1x1)
Conv-ReLU-BN (128, 3x3, 1x1, 1x1)	Conv-ReLU-BN (192, 1x1, 1x1)
Avg-Pooling (6x6, 1x1)	Max-Pooling (3x3, 2x2)
	Conv-ReLU-BN (192, 3x3, 1x1, 1x1)
	Conv-ReLU-BN (192, 1x1, 1x1)
	Avg-Pooling (7x7, 1x1)