

# Learning about an exponential amount of conditional distributions

Mohamed Ishmael Belghazi<sup>1,2</sup> Maxime Oquab<sup>1</sup> Yann LeCun<sup>1</sup> David Lopez-Paz<sup>1</sup>

## Abstract

We introduce the Neural Conditioner (NC), a self-supervised machine able to learn about all the conditional distributions of a random vector  $X$ . The NC is a function  $\text{NC}(x \cdot a, a, r)$  that leverages adversarial training to match each conditional distribution  $P(X_r | X_a = x_a)$ . After training, the NC generalizes to sample conditional distributions never seen, including the joint distribution. The NC is also able to auto-encode examples, providing data representations useful for downstream classification tasks. In sum, the NC integrates different self-supervised tasks (each being the estimation of a conditional distribution) and levels of supervision (partially observed data) seamlessly into a single learning experience.

## 1. Introduction

Supervised learning estimates the conditional distribution of a target variable given values for a feature variable (Vapnik, 1998). Supervised learning is the backbone to build state-of-the-art prediction models using large amounts of labeled data, with unprecedented success in domains spanning image classification, speech recognition, and language translation (LeCun et al., 2015). Unfortunately, collecting large amounts of labeled data is an expensive task painstakingly performed by humans (for instance, consider labeling the objects appearing in millions of images). If our ambition to transition from machine learning to artificial intelligence is to be met, we must build algorithms capable of learning effectively from inexpensive unlabeled data without human supervision (for instance, millions of unlabeled images). Furthermore, we are interested in the case where the available unlabeled data is partially observed. Thus, the goal of this paper is unsupervised learning, defined as understanding the underlying process generating some partially observed unlabeled data.

<sup>\*</sup>Equal contribution <sup>1</sup>Facebook AI Research, Paris, France  
<sup>2</sup>Montréal Institute for Learning Algorithms (MILA), University of Montréal. Correspondence to: Mohamed Ishmael Belghazi <ishmael.belghazi@gmail.com>.

Currently, unsupervised learning strategies come in many flavors, including component analysis, clustering, energy modeling, and density estimation (Hastie et al., 2009). Each of these strategies targets the estimation of a particular statistic from high-dimensional data. For example, principal component analysis extracts a set of directions under which the data exhibits maximum variance (Jolliffe, 2011). However, powerful unsupervised learning should not commit to the estimation of a particular statistic from data, but extract general-purpose features useful for downstream tasks.

An emerging, more general strategy to unsupervised learning is the one of self-supervised learning (Hinton & Salakhutdinov, 2006, for instance). The guiding principle behind self-supervised learning is to set up a supervised learning problem based on unlabeled data, such that solving that supervised learning problem leads to partial understanding about the data generating process (Kolesnikov et al., 2019). More specifically, self-supervised learning algorithms transform the unlabeled data into one set of input features and one set of output features. Then, a supervised learning model is trained to predict the output features from the input features. Finally, the trained model is later leveraged to solve subsequent learning tasks efficiently. As such, self-supervision turns unsupervised learning into the supervised learning problem of estimating the conditional expectation of the output features given the input features. A common example of a self-supervised problem is image in-painting. Here, the central patch of an image (output feature) is predicted from its surrounding pixel values (input feature), with the hope that learning to in-paint leads to the learning of non-trivial image features (Pathak et al., 2016; Liu et al., 2018). Another example of a self-supervised learning problem extracts a pair of patches from one image as the input feature, and requests their relative position as the target output feature (Doersch et al., 2015). These examples hint one potential pitfall of “specialized” self-supervised learning algorithms: in order to learn a single conditional distribution from the many describing the data, it may be acceptable to throw away most of the information about the sought generative process, which in fact we would like to keep for subsequent learning tasks.

Thus, a general-purpose unsupervised learning machine should not commit to the estimation of a particular conditional distribution from data, but attempt to learn as much

structure (i.e., interactions between variables) as possible. This is a daunting task, since joint distributions can be described in terms of an exponential amount of conditional distributions. This means that, unsupervised learning, when attacked in its most general form, is analogous to an exponential amount of supervised learning problems.

Our challenges do not end here. Being realistic, learning agents never observe the entire world. For instance, occlusions and camera movements hide portions of the world that we would otherwise observe. Therefore, we are interested in unsupervised learning algorithms able to learn about the structure of unlabeled data from partial observations.

In this paper, we address the task of unsupervised learning from partial data by introducing the Neural Conditioner (NC). In a nutshell, the NC is a function  $\text{NC}(x \cdot a, a, r)$  that leverages adversarial training to match each conditional distribution  $P(X_r | X_a = x_a)$ . The set of available variables  $a$ , the set of requested variables  $r$ , and the set of available values  $x \cdot a$  can be either determined by the pattern of missing values in data, or randomly by the self-supervised learning process. The set of available variables  $a$  and the set of requested variables  $r$  are not necessarily complementary, and index an exponential amount of conditional distributions (each associated to a single self-supervised learning problem). After trained, the NC generalizes to sample from conditional distributions never seen during training, including the joint distribution. Furthermore, trained NC's are also able to auto-encode examples, providing data representations useful for downstream classification tasks. Since the NC does not commit to a particular conditional distribution but attempts to learn a large amount of them, we argue that our model is a small step towards general-purpose unsupervised learning. Our contributions are as follows:

- We introduce the Neural Conditioner (NC) (Section 2), a method to perform unsupervised learning from partially observed data.
- We explain the multiple uses of NCs (Section 3), including the generation of conditional samples, unconditional samples, and feature extraction from partially observed data.
- We provide insights on how NCs work and should be regularized (Section 4).
- Throughout a variety of experiments on synthetic and image data, we show the efficacy of NCs in generation and prediction tasks (Section 5).

## 2. The Neural Conditioner (NC)

Consider the dataset  $(x_1, \dots, x_n)$ , where each  $x_i \in \mathbb{R}^d$  is an identically and independently distributed (iid) example drawn from some joint probability distribution  $P(X)$ .

Without any further information, we could consider  $O(3^d)$  different prediction problems about the random vector  $X$ , where each prediction problem partitions the coordinates  $x_i$  into features, targets, or unobserved variables. We may index this exponential amount of supervised learning problems using binary vectors of *available features*  $a \in \{0, 1\}^d$  and *requested features*  $r \in \{0, 1\}^d$ . In statistical terms, a pair of available and requested vectors  $(r, a)$  instantiates the supervised learning problem of estimating the conditional distribution  $P(X_r | X_a = x_a)$ , where  $x_r = (x_i : r_i = 1)$ , and  $x_a = (x_i : a_i = 1)$ .

By making use of the notations above, we can design a single supervised learning problem to estimate all the conditional distributions contained in the random vector  $X$ . Since learning algorithms are often designed to deal with inputs and outputs with a fixed number of dimensions, we will consider the augmented supervised learning problem of mapping the feature vector  $(x \cdot a, a, r)$  into the target vector  $x \cdot r$ , where the operation “ $\cdot$ ” denotes entry-wise multiplication. In short, our goal is to learn a Neural Conditioner (NC) producing samples:

$$\hat{x} \sim \text{NC}(x \cdot a, a, r) : \hat{x}_r \sim P(X_r | X_a = x_a) \quad \forall (x, a, r). \quad (1)$$

The previous equation manifests the ambition of NC to model the entire conditional distribution  $P(X_r | X_a = x_a)$  when given a triplet  $(x, a, r)$ . Therefore, given the dataset  $(x_1, \dots, x_n)$ , learning a NC translates into minimizing the distance between the estimated conditional distributions  $\text{NC}(x \cdot a, a, r)$  and the true conditional distributions  $P(X_r | X_a = x_a)$ , based on their samples. In particular, we will follow recent advances in implicit generative modeling, and implement NC training using tools from generative adversarial networks (Goodfellow et al., 2014). Other alternatives to train NCs would include maximum mean discrepancy metrics (Gretton et al., 2012), energy distances (Székely et al., 2007), or variational inference (Kingma & Welling, 2013). If the practitioner is only interested in recovering a particular statistic from the exponentially many conditional distributions (e.g. the conditional means), training a NC with a scoring rule  $D$  for such statistic (e.g. the mean squared error loss) would suffice.

Training a NC is an iterative process involving six steps, illustrated in Figures 1 and 2:

1. A data sample  $x$  is drawn from  $P(X)$ .
2. Available and requested masks  $(r, a)$  are drawn from some data-defined or user-defined distribution  $P(R, A)$ . These masks are not necessarily complementary, enabling the existence of unobserved (neither requested or observed) variables. If a coordinate equals to one in both  $r$  and  $a$ , we zero it at the requested mask.

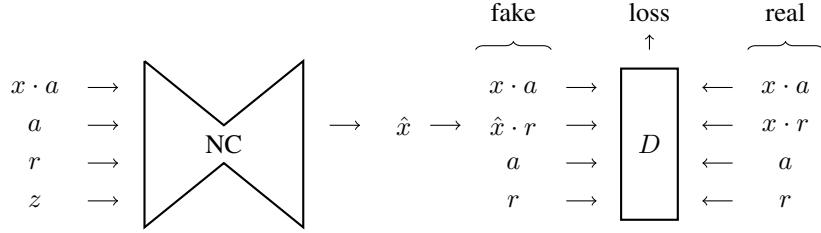


Figure 1. The proposed NC, where data  $x \sim P(X)$ , available/requested masks  $a, r \sim P(a, r)$ , and noise  $z \sim \mathcal{N}(0, I)$ .



Figure 2. Example of masks and masked images. At each iteration, the NC learns to predict  $x \cdot r$  from  $x \cdot a$ .

3. A noise vector  $z$  is sampled from an external source of noise, following some user-defined distribution  $P(Z)$ .
4. A sample is generated as  $\hat{x} = \text{NC}(x \cdot a, a, r, z)$ .
5. A discriminator  $D$  provides the final scalar objective function by distinguishing between data samples (scored as  $D(x \cdot r, x \cdot a, a, r)$ ) and generated samples (scored as  $D(\hat{x} \cdot r, x \cdot a, a, r)$ ).
6. The NC parameters are updated to minimize the scalar objective function, while the parameters of the discriminator are updated to maximize it, in what becomes an adversarial training game (Goodfellow et al., 2014).

Mathematically, our general objective function is:

$$\begin{aligned} \min_{\text{NC}} \max_D & \mathbb{E}_{x, a, r} \log D(x \cdot r, x \cdot a, a, r) + \\ & \mathbb{E}_{x, a, r, z} \log(1 - D(\text{NC}(x \cdot a, a, r, z) \cdot r, x \cdot a, a, r)). \end{aligned} \quad (2)$$

### 3. Using NCs

Once trained, one NC serves many purposes.

The most direct use is perhaps the *multimodal prediction of any subset of variables given any subset of variables*. More specifically, a NC is able to leverage any partially observed vector  $x_a$  to predict about any partially requested vector  $x_r$ . Importantly, the combination of test values, available, and requested masks  $(x, a, r)$  could be novel and never seen during training. Since NCs leverage an external source of noise  $z$  to make their predictions, NCs provide a conditional distribution for each triplet  $(x, a, r)$ .

Two special cases of masks deserve special attention. First, properly regularized NCs are able to *compress and reconstruct samples* when provided with the full requested mask  $r = 0$  and the full available mask  $a = 1$ . This turns NCs into autoencoders able to *extract feature representations of data*, as well as allowing *latent interpolations between pairs of examples*. Second, when provided with the full requested mask  $r = 1$  and the empty available mask  $a = 0$ , NCs are able to *generate full samples from the data joint distribution*  $P(X)$ , even in the case when the training never provided the NC with this mask combination, as our experiments verify.

NCs are able to seamlessly *deal with missing features and/or labels during both training and testing time*. Such “missingness” of features and labels can be real (as given by incomplete or unlabeled examples) or simulated by designing an appropriate distribution of masks  $P(A, R)$ . This blurs the lines that often separate unsupervised, semi-supervised, and supervised learning, integrating all types of data and supervision into a new learning paradigm.

Finally, a trained NC can be used to understand relations between variables, for instance by using a complete test vector  $x$  and querying different available and requested masks. The strongest relations between variables can also be analyzed in terms of gradients with respect to  $(a, r)$ .

### 4. Understanding NCs

To better understand how NCs work, this section describes i) how NCs look like in the Gaussian case, ii) what the optimal discriminator minimizes, iii) the relationship between NC training and the usual reconstruction error minimized by auto-encoders, and iv) some regularization techniques.

#### 4.1. The Gaussian case

Let us consider the case where the data joint distribution is a Gaussian  $P(X) = \mathcal{N}(\mu, \Sigma)$ . Then, the closed-form expression of the conditional distribution implied by any triplet  $(x, a, r)$  is  $P(X_r | X_a = x_a) = \mathcal{N}(\mu_{r|a}, \Sigma_{r|a})$ , where

$$\mu_{r|a} = \mu_r + \Sigma_{ra} \Sigma_{aa}^{-1} (x_a - \mu_a), \quad (3)$$

$$\Sigma_{r|a} = \Sigma_{rr} - \Sigma_{ra} \Sigma_{aa}^{-1} \Sigma_{ar}. \quad (4)$$

The previous expressions highlight an interesting fact: even in the case of Gaussian distributions, computing the conditional moments implied by  $(x, a, r)$  is a non-linear operation. When fixing  $(a, r) = (a_0, r_0)$ , learning the conditional distribution implied by triplets  $(x, a_0, r_0)$  can be understood as linear heteroencoding (Roweis & Brody, 1999).

The motivation behind self-supervised learning is that learning about a conditional distribution is an effective way to learn about the joint distribution. In part, this is because learning conditional distributions allows to deploy the powerful machinery of supervised learning. To formalize this, we consider the amount of information contained in a probability distribution in terms of its differential entropy. Then, we show that learning conditional distributions is easier than learning joint distributions, where “difficult” is measured in terms of how much information is to be learned. This argument can be made by considering the chain rule of the differential entropy (Cover & Thomas, 2012):

$$h(X) = \sum_{i=1}^d h(X_i | X_1, \dots, X_{i-1}), \quad (5)$$

where, in the case of partitioning  $X = (X_a, X_r)$ , we have:

$$h(X) = h(X_r | X_a) + h(X_a). \quad (6)$$

The previous shows that  $h(X_r | X_a) \leq h(X_r)$ , where equality is achieved if and only if  $X_a$  and  $X_r$  are independent. This reveals a “blessing of structure” of sorts: to reduce the difficulty of learning about a joint distribution, we should construct self-supervised learning problems associated to conditional distributions between highly coupled blocks of input and output features. Indeed, if all of our variables are independent, self-supervised learning is hopeless. For the case of a  $d$ -dimensional Gaussian with covariance matrix  $\Sigma$ , the differential entropy can be stated in terms of the covariance function:

$$h(\Sigma) = \frac{d}{2}(1 + \log(2\pi)) + \frac{1}{2} \log(|\Sigma|), \quad (7)$$

which allows to choose good self-supervised learning problems based on the log-determinant of empirical covariances.

A successful evolution from single self-supervised learning problems to NCs rests on the existence of relationships between different conditional distributions. More formally, the

success of NCs relies on assuming a smooth landscape of conditionals. If smoothness across conditional distributions is satisfied, learning about some conditional distribution should inform us about other, perhaps never seen, conditionals. This is akin to supervised learning algorithms relying on smoothness properties of the function to be learned. For NCs we do not consider the smoothness of a single function, but the smoothness of the “conditioning operator”  $C_x(a, r) = \text{NC}(x \cdot a, a, r)$ . The smoothness of this conditioning operator is related to the smoothness of the covariance operator studied in kernel embeddings of distributions (Muandet et al., 2017). In one extreme case, product distributions will lead to non-smooth conditional operators, since all conditionals are independent. In the other extreme case, where all the components of the random vector  $X$  are copies, the smoothness operator is constant, so learning about one conditional distribution informs us about all other conditional distributions. This is another instantiation of the “blessing of structure”, present in data such as images.

#### 4.2. Training objective, discriminator’s point of view

Next, we would like to understand about the problem that the discriminator  $D$ , involved in NC training, is trying to solve. To this end, consider the mapping

$$D \mapsto V[D] = \mathbb{E}[D(X_A, X_R, A, R)] - \log(\mathbb{E}[\exp(D(X_A, \hat{X}_R, A, R))]), \quad (8)$$

where  $\hat{X}_R$  is the generated sample from  $\text{NC}(X_A, A, R)$ . Use  $h \in \mathcal{C}_b$  and  $\alpha \in [0, \infty)$  to form the Gateaux derivative

$$dV(D; h) = \mathbb{E}[h] - \mathbb{E}\left[\frac{e^D}{Z} h\right], \quad (9)$$

where differentiation under the integral is justified by dominated convergence, and  $Z$  is the *Gibbs* partition function  $Z = \mathbb{E}[e^D]$ . Since the map  $D \mapsto V[D]$  is concave, the maximum is reached at the critical point of the Gateaux derivative. This derivative is zero if and only if the optimal discriminator  $D^*$  satisfies:

$$\frac{e^{D^*}}{Z} = \frac{P(X_R | X_A, A, R)}{P(\hat{X}_R | X_A, A, R)}. \quad (10)$$

The previous expression shows that i) NCs need to estimate all the conditional when fighting powerful discriminators, and ii) we need to provide the discriminator with full tuples  $((X_R, X_A, A, R), (\hat{X}_R, X_A, A, R))$  versus  $((X_R), (\hat{X}_R))$ .

#### 4.3. Training objective, NC’s point of view

The previous section shed some light to the objective function of the discriminator  $D$  during training. This section considers the opposite question: what is the objective function minimized by NC? In particular we are interested in the

intriguing fact of how NCs is able to complete and reconstruct samples, when the discriminator is never presented with pairs of real and generated requested variables. First, consider the “augmented” data joint distribution

$$P(X_a, X_r, A, R) = q(X_r | X_a, A, R)p(X_a, A, R), \quad (11)$$

and the augmented model joint distribution

$$P_\theta(X_a, X_r, A, R) = q_\theta(X_r | X_a, A, R)p(X_a, A, R). \quad (12)$$

Next, consider the negative log-likelihood  $L(x_a, a, r) = -\mathbb{E}_q \log q_\theta$  and its expectation  $L = -\mathbb{E}_P \log q_\theta$ . Recall that the latter expectation is the objective function minimized by generators in the usual non-saturating GAN objective (Goodfellow et al., 2014), such as it happens in NC. Then, we can see

$$L(X_a, A, R) = -\mathbb{E}_a \log \left( q_\theta \cdot \frac{q}{q} \right) \quad (13)$$

$$= -\mathbb{E}_a \left\{ \log \frac{q_\theta}{q} + \log q \right\} \quad (14)$$

$$= \int q \log q - \int q \log \frac{q_\theta}{q}. \quad (15)$$

Integrating wrt  $p(X_a, A, R)$ , we see that NCs minimize:

$$L = D_{KL}(P \| P_\theta) + H(X_R | X_A) \quad (16)$$

$$= D_{KL}(P \| P_\theta) - I(X_A, X_R) + H(X_R) \quad (17)$$

$$= D_{KL}(P \| P_\theta) - I(X_A, X_R) + H(X_A). \quad (18)$$

Where  $H$  stands for (conditional) entropy and  $I$  for mutual information. Lastly, by the positivity of the KL divergence, we have that  $L \geq H(X_R | X_A)$ .

We summarize the previous results as follows. If a NC is able to match the distributions  $(P, P_\theta)$ , there will be a residual reconstruction error of  $H(X_R, X_A)$ . Thus, if  $X_A$  and  $X_R$  are independent, such residual reconstruction error reduces to  $H(X_R)$ . This can happen if  $A = 0$ , or if  $X_A$  holds no information about  $X_R$ . Moreover the reconstruction error is a decreasing function of the amount of information that  $X_A$  holds about  $X_R$ . Since the entropic term  $H(X_A)$  does not depend on  $\theta$ , it will have no effect on the learning of NC. Therefore, the learning signal about the reconstruction error is bounded by  $I(X_A, X_R)$ .

#### 4.4. Regularization

We close this section with a few words on how to regularize NC training. We found, during our experiments, gradient based regularization on the discriminator to be crucial. Following (Roth et al., 2017) we augment the discriminator’s loss with the expected gradient with respect to the inputs for both the positive and negative examples; Less succinctly,

we add  $\frac{1}{2}(\mathbb{E}[D(X_A, X_R, A, R)] + \mathbb{E}[D(X_A, \hat{X}_R, A, R)])$  to the discriminator’s loss.

For NC to generalize to unobserved conditional distributions and prevent memorizing the observed ones, we have found that regularization of the latent space to be essential. In information theoretic terms, we would like to control the mutual information between  $X_A$  and  $Z := enc(X_A, \epsilon)$ . One could use a variational approximation of the conditional entropy (Alemi et al., 2016) or an adversarial approach (Belghazi et al., 2018a). The former requires an encoder with tractable conditional density (e.g. Gaussian), the latter, while allowing general encoders, introduces an additional training loop in the algorithm. We opt for another approach by controlling the encoder’s Lipschitz constant using a one-sided variation of spectral normalization (Miyato et al., 2018).

## 5. Experiments

In this section we conduct experiments on Gaussian and image data to showcase the uses and performance of NC. We defer implementation details to the Supplementary Material.

### 5.1. Gaussian data

We train a single NC to model all the conditionals of a three-dimensional Gaussian distribution. Given that in this example we know that the data generating process is fully determined by the first two moments, we train two versions of NCs: one that uses moment-matching, and one that uses our full adversarial training pipeline. Both strategies train NC given minibatches of triplets  $(x, a, r)$  observed from the same Gaussian distribution. This allows us to better understand the impact of adversarial training when dealing with NCs. For these experiments, both the discriminator and the NC have 2 hidden layers of 64 units each, and ReLU nonlinearities. We regularize the latent space of the NC using one-sided spectral normalization (Miyato & Koyama, 2018). We train the networks for 10,000 updates, with a batch-size of 512, and the Adam optimizer with a learning rate of  $10^{-4}$ ,  $\beta_1 = 0.5$ , and  $\beta_2 = 0.999$ . The training set contains  $10^4$  fixed samples sampled from a Gaussian with mean  $(2, 4, 6)$  and covariance  $((1, 0.5, 0.25), (0.5, 1, 0), (0.25, 0, 1))$ .

Figure 3 illustrates the capabilities of NC to perform one-dimensional and two-dimensional conditional distribution estimation. We also show the embeddings of the conditional distributions as given by the bottleneck of NC. These show a higher dependence for variables that are more tightly coupled. Table 1 shows the error on the conditional parameter estimation for the NC (both using moment matching and adversarial training) as well as the VAEAC (Ivanov et al., 2019a), a VAE-based analog to the NC. Finally, Table 2 shows the importance of conditioning both the discriminator and NC on both available and requested masks.

Figure 3. Illustration of the NC on a three-dimensional Gaussian dataset. We show a) one-dimensional conditional estimation, b) two-dimensional conditional estimation, and c,d) the representation of the conditional distributions in the hidden space.

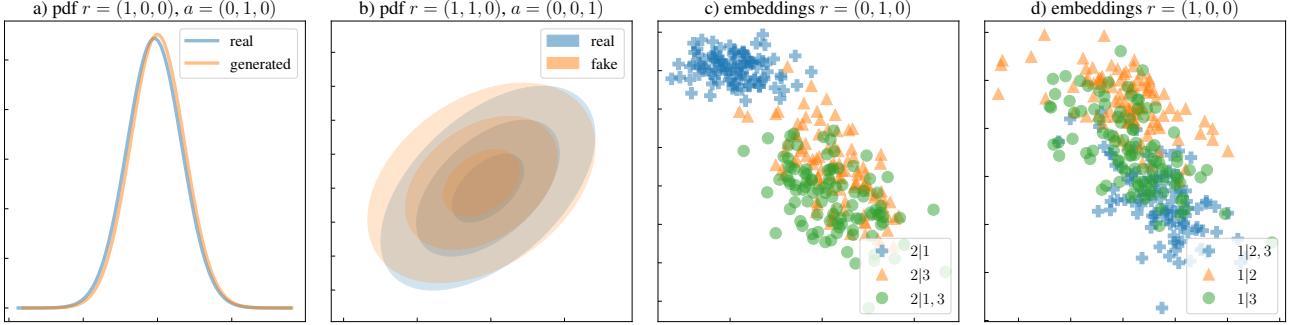


Table 1. Error norms  $\|\theta_{r|a} - \hat{\theta}_{r|a}\|$  (averaged over ten runs) in the task of estimating the conditional moments  $\theta_{r|a} = (\mu_{r|a}, \Sigma_{r|a})$  of Gaussian data. We show results for NC trained with Moment-Matching (MM) or the full Adversarial Training (AT). VAEAC only supports complementary masks, therefore some results are unavailable.

$a$	$r$	NC (MM)	NC (AT)	VAEAC
(1, 0, 0)	(0, 0, 1)	0.09	0.11	NA
	(0, 1, 0)	0.10	0.08	NA
	(0, 1, 1)	0.67	0.13	0.68
(0, 1, 0)	(0, 0, 1)	0.16	0.08	NA
	(1, 0, 0)	0.20	0.05	NA
	(1, 0, 1)	0.28	0.14	0.73
(0, 0, 1)	(0, 1, 0)	0.13	0.11	NA
	(1, 0, 0)	0.08	0.09	NA
	(1, 1, 0)	0.29	0.17	0.71
(1, 0, 1)	(0, 1, 0)	0.22	0.13	0.50
	(1, 1, 0)	0.15	0.08	0.43
	(0, 1, 0)	0.27	0.07	0.35

Table 2. Euclidean error between the true and estimated Gaussian parameters as a function of masks conditioning in the discriminator and NC (averages over ten runs).

	NC conditioning	
	$\emptyset$	$(a, r)$
discriminator	$\emptyset$	0.12
conditioning	$(a, r)$	0.15
		0.07

## 5.2. Image data

We train NCs on SVHN and CelebA. We use rectangular  $a, r$  masks spanning between 10% and 50% of the images.

We evaluate our setup in several ways. First qualitatively: generating full samples (using the never seen mask configuration  $a = 0, r = 1$ , Fig 4) and reconstructing samples (Figures 5 for denoising and 6 for inpainting). These experiments share the goal of showing that our model is able to generalize to conditional distributions not observed during training. Second, we evaluate our models quantitatively: that is, their ability to provide useful features for downstream classification tasks (see Tables 3 and 4). Our results show that NC-based figures systematically outperform state-of-art hand-crafted features, while being competitive with deep unsupervised features.

Figures 5 and 6 show samples and in-paintings using masks configurations unobserved during training to illustrate that our model is able to generalize to conditional distributions and construct representation of the data solely through partial observation. Figure 4 shows samples from the joint distribution ( $a = 0, r = 1$ ), even though these masks were never observed during training.

### 5.2.1. FEATURE EXTRACTION

**SVHN** As a feature extraction procedure, we retrieve the latent code created by the PAE while feeding an image in *compress and reconstruct* mode ( $a = r = 1$ ). Then, we use a linear SVM to assess the quality of the extracted encoding, and show in Table 3 that our approach is competitive with deep unsupervised feature extractor.

**CelebA** The multimodality presented by the CelebA attributes provides an ideal test mode to quantify our model ability to construct a global understanding out of local and partial observations.



Figure 4. SVHN and CelebA samples from the joint distribution. The model never observed a complete example during training.

Table 3. Test errors on SVHN classification experiments.

model	test error
VAE (M1 + M2) (Kingma et al., 2014)	36.02
SWWAE with dropout (Zhao et al., 2015)	23.56
DCGAN + L2-SVM (Radford et al., 2015)	22.18
SDGM (Maaløe et al., 2016)	16.61
ALI + L2-SVM (Dumoulin et al., 2016)	19.14
NC (L2-SVM) (ours)	<b>17.12</b>

Following Berg & Belhumeur (2013); Liu et al. (2015), we train 40 linear SVMs on learned representation representations extracted from the encoder using full available and requested masks ( $a = r = 1$ ) on the CelebA validation set. We measure the performance on the test set. As in (Berg & Belhumeur, 2013; Huang et al., 2016; Kalayeh et al., 2017), we report the *balanced accuracy* in order to evaluate the attribute prediction performance. Please note that our model was trained on entirely unsupervised data and masking configurations unobserved during training. Attribute labels were only used to train the linear SVM classifiers.

## 6. Related work

Self-supervised learning is an emerging technique for unsupervised learning. Perhaps the earliest example of self-supervised learning is auto-encoding (Baldi & Hornik, 1989; Hinton & Salakhutdinov, 2006), which in the language of NCs amounts to full available and requested masks. Auto-encoders evolved into more sophisticated variants such as

Table 4. Test balanced accuracies on CelebA classification experiments.

model	mean	stdv
Triplet-kNN (Schroff et al., 2015)	71.55	12.61
PANDA (Zhang et al., 2014)	76.95	13.33
Anet (Liu et al., 2015)	79.56	12.17
LMLE-kNN (Huang et al., 2016)	83.83	12.33
VAE (Kingma & Welling, 2013)	73.30	9.65
ALI (Dumoulin et al., 2016)	73.88	10.16
HALI (Belghazi et al., 2018b)	83.75	8.96
VAEAC (Ivanov et al., 2019b)	66.06	6.98
NC (Ours)	82.21	7.63

denoising auto-encoders (Vincent et al., 2010), a family of models including NC. Recent trends in generative adversarial networks (Goodfellow et al., 2014) are yet another example of self-supervised training. The connection between auto-encoders and generative adversarial training was first instantiated by Larsen et al. (2015). Auto-regressive models (Bengio & Bengio, 2000) such as the masked autoencoder (Germain et al., 2015), neural autoregressive distribution estimators (Larochelle & Murray, 2011; Uria et al., 2014), and Pixel RNNs (Oord et al., 2016) are other examples of casting unsupervised learning using a simple self-supervision strategy: order the variables, and then predict each of them using the previous in the ordering.

Moving further, the task of unsupervised learning with partially observed data was also considered by others, often in



Figure 5. Denoising SVHN images corrupted with 50% missing pixels using a model *trained on square masks*.



Figure 6. In-painting SVHN images using masks of size and shapes *not seen during training*.



Figure 7. Predicting partially-observed CelebA images. From left to right:  $x \cdot a$ ,  $x \cdot r$ ,  $\hat{x} \cdot r$ ,  $\hat{x}$ ,  $(x \cdot a + \hat{x} \cdot r)$ ,  $x$ . Saturation patterns happen only for pixels where  $a = 1$ .

terms of estimating transition operators (Goyal et al., 2017; Bordes et al., 2017; Sohl-Dickstein et al., 2015). Generative adversarial imputation nets (Yoon et al., 2018) considered the case of learning missing feature predictions using adversarial training. In a different thread of research, the literature in kernel mean embeddings (Song et al., 2009; Lever et al., 2012; Muandet et al., 2017) is an early consideration of the problem of learning distributions.

Moving to applications, successful incarnations of self-supervised are pioneered by word embeddings (Mikolov et al., 2013). In the image domain, self-supervised setups include image in-painting (Pathak et al., 2016), colorization (Zhang et al., 2016), clustering (Caron et al., 2018), de-rotation (Gidaris et al., 2018), and patch reordering (Doersch

et al., 2015; Noroozi & Favaro, 2016). In the video domain, common self-supervised strategies include enforcing similar feature representations for nearby frames (Mobahi et al., 2009; Goroshin et al., 2015; Wang & Gupta, 2015), or predicting ambient sound statistics from video frames (Owens et al., 2016). These applications yield representations useful for downstream tasks, including classification (Caron et al., 2018), multi-task learning (Doersch & Zisserman, 2017), and RL (Pathak et al., 2017).

Finally, the most similar piece of literature to our research is the concurrent work on VAE with Arbitrary Conditioning, or VAEAC (Ivanov et al., 2019a). The VAEAC is proposed as a fast alternative to the also related universal marginalizer (Douglas et al., 2017). Similarly to our setup, the VAEAC augments a VAE with a mask of requested variables; the complimentary set of variables is provided as the available information for prediction. Our work extends VAEAC by employing adversarial training to obtain better sample quality and features for downstream tasks. To sustain these claims, a comparison between NC and VAEAC was performed in Section 5. As commonly assumed in VAE-like architectures, the conditional encoding and decoding distributions are assumed Gaussian, which may not be a good fit for complex multimodal data such as natural images. The VAEAC work was mainly applied to the problem of feature imputation. Here we hope to provide a more holistic perspective on the uses of NCs, including feature extraction and semi-supervised learning. Furthermore, training VAEACs involves learning two encoders, and their use of complementary masks  $r = 1 - a$  complicates the definition of missing (neither requested or available) variables.

## 7. Conclusion

We presented the Neural Conditioner (NC), an adversarially-learned neural network able to learn about the exponentially many conditional distributions describing some partially observed unlabeled data. Once trained, one NC serves many purposes: sampling from (unseen) conditional distributions to perform multimodal prediction, sampling from the (unseen) joint distribution, and auto-encode (partially observed) data to extract data representations useful for (semi-supervised) downstream tasks. Our biggest ambition when we built the NC is to make one small step towards holistic machine learning. That is, can we shift from the pure supervised learning (where features and labels are prescribed) into a fully self-supervised learning paradigm where everything can be predicted from anything? We wish that the NC serves as inspiration to researchers who, like us, hope to unleash the full power of unlabeled data to build better intelligent systems.

## References

- Alemi, A. A., Fischer, I., Dillon, J. V., and Murphy, K. Deep variational information bottleneck. *arXiv preprint arXiv:1612.00410*, 2016.
- Baldi, P. and Hornik, K. Neural networks and principal component analysis: Learning from examples without local minima. *Neural networks*, 2(1):53–58, 1989.
- Belghazi, M. I., Baratin, A., Rajeshwar, S., Ozair, S., Bengio, Y., Courville, A., and Hjelm, D. Mutual information neural estimation. In Dy, J. and Krause, A. (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 531–540, Stockholm, Sweden, 10–15 Jul 2018a. PMLR. URL <http://proceedings.mlr.press/v80/belghazi18a.html>.
- Belghazi, M. I., Rajeswar, S., Mastropietro, O., Rosamzadeh, N., Mitrovic, J., and Courville, A. Hierarchical adversarially learned inference. *International Conference on Machine Learning Workshop on Theoretical Foundations and Applications of Deep Generative Models*, 2018b.
- Bengio, S. and Bengio, Y. Taking on the curse of dimensionality in joint distributions using neural networks. *IEEE Transactions on Neural Networks*, 11(3):550–557, 2000.
- Berg, T. and Belhumeur, P. N. Poof: Part-based one-vs.-one features for fine-grained categorization, face verification, and attribute estimation. In *CVPR*, pp. 955–962, 2013.
- Bordes, F., Honari, S., and Vincent, P. Learning to generate samples from noise through infusion training. *arXiv*, 2017.
- Caron, M., Bojanowski, P., Joulin, A., and Douze, M. Deep clustering for unsupervised learning of visual features. *ECCV*, 2018.
- Cover, T. M. and Thomas, J. A. *Elements of information theory*. John Wiley & Sons, 2012.
- Doersch, C. and Zisserman, A. Multi-task self-supervised visual learning. *ICCV*, 2017.
- Doersch, C., Gupta, A., and Efros, A. A. Unsupervised visual representation learning by context prediction. *CVPR*, 2015.
- Douglas, L., Zarov, I., Gourgoulias, K., Lucas, C., Hart, C., Baker, A., Sahani, M., Perov, Y., and Johri, S. A universal marginalizer for amortized inference in generative models. *arXiv*, 2017.
- Dumoulin, V., Belghazi, I., Poole, B., Mastropietro, O., Lamb, A., Arjovsky, M., and Courville, A. Adversarially learned inference. *arXiv*, 2016.
- Germain, M., Gregor, K., Murray, I., and Larochelle, H. Made: Masked autoencoder for distribution estimation. In *ICML*, pp. 881–889, 2015.
- Gidaris, S., Singh, P., and Komodakis, N. Unsupervised representation learning by predicting image rotations. *ICLR*, 2018.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. Generative adversarial nets. In *Advances in neural information processing systems*, pp. 2672–2680, 2014.
- Goroshin, R., Bruna, J., Tompson, J., Eigen, D., and LeCun, Y. Unsupervised learning of spatiotemporally coherent metrics. In *Proceedings of the IEEE international conference on computer vision*, pp. 4086–4093, 2015.
- Goyal, A. G. A. P., Ke, N. R., Ganguli, S., and Bengio, Y. Variational walkback: Learning a transition operator as a stochastic recurrent net. In *NeurIPS*, pp. 4392–4402, 2017.
- Gretton, A., Borgwardt, K. M., Rasch, M. J., Schölkopf, B., and Smola, A. A kernel two-sample test. *JMLR*, 13(1):723–773, 2012.
- Hastie, T., Tibshirani, R., and Friedman, J. Unsupervised learning. In *The elements of statistical learning*, pp. 485–585. Springer, 2009.
- Hinton, G. E. and Salakhutdinov, R. R. Reducing the dimensionality of data with neural networks. *science*, 313 (5786):504–507, 2006.
- Huang, C., Li, Y., Change Loy, C., and Tang, X. Learning deep representation for imbalanced classification. In *CVPR*, pp. 5375–5384, 2016.
- Ivanov, O., Figurnov, M., and Dmitry, V. Variational autoencoder with arbitrary conditioning. *ICLR*, 2019a.
- Ivanov, O., Figurnov, M., and Vetrov, D. Variational autoencoder with arbitrary conditioning. In *ICLR*, 2019b. URL <https://openreview.net/forum?id=SyxtJh0qYm>.
- Jolliffe, I. Principal component analysis. In *International encyclopedia of statistical science*, pp. 1094–1096. Springer, 2011.
- Kalayeh, M. M., Gong, B., and Shah, M. Improving facial attribute prediction using semantic segmentation. *arXiv*, 2017.

- Kingma, D. P. and Welling, M. Auto-encoding variational bayes. *arXiv*, 2013.
- Kingma, D. P., Mohamed, S., Rezende, D. J., and Welling, M. Semi-supervised learning with deep generative models. In *NeurIPS*, pp. 3581–3589, 2014.
- Kolesnikov, A., Zhai, X., and Beyer, L. Revisiting self-supervised visual representation learning. *arXiv*, 2019.
- Larochelle, H. and Murray, I. The neural autoregressive distribution estimator. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, pp. 29–37, 2011.
- Larsen, A. B. L., Sønderby, S. K., Larochelle, H., and Winther, O. Autoencoding beyond pixels using a learned similarity metric. *arXiv*, 2015.
- LeCun, Y., Bengio, Y., and Hinton, G. Deep learning. *nature*, 521(7553):436, 2015.
- Lever, G., Baldassarre, L., Patterson, S., Gretton, A., Pontil, M., and Grünewälder, S. Conditional mean embeddings as regressors. In *International Conference on Machine Learning (ICML)*, volume 5, 2012.
- Liu, G., Reda, F. A., Shih, K. J., Wang, T.-C., Tao, A., and Catanzaro, B. Image inpainting for irregular holes using partial convolutions. *ECCV*, 2018.
- Liu, Z., Luo, P., Wang, X., and Tang, X. Deep learning face attributes in the wild. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 3730–3738, 2015.
- Maaløe, L., Sønderby, C. K., Sønderby, S. K., and Winther, O. Auxiliary deep generative models. *arXiv*, 2016.
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., and Dean, J. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pp. 3111–3119, 2013.
- Miyato, T. and Koyama, M. cGANs with projection discriminator. *arXiv*, 2018.
- Miyato, T., Kataoka, T., Koyama, M., and Yoshida, Y. Spectral normalization for generative adversarial networks. In *ICLR*, 2018. URL <https://openreview.net/forum?id=B1QRgziT->.
- Mobahi, H., Collobert, R., and Weston, J. Deep learning from temporal coherence in video. In *ICML*, pp. 737–744. ACM, 2009.
- Muandet, K., Fukumizu, K., Sriperumbudur, B., Schölkopf, B., et al. Kernel mean embedding of distributions: A review and beyond. *Foundations and Trends® in Machine Learning*, 10(1-2):1–141, 2017.
- Noroozi, M. and Favaro, P. Unsupervised learning of visual representations by solving jigsaw puzzles. In *ECCV*, pp. 69–84. Springer, 2016.
- Oord, A. v. d., Kalchbrenner, N., and Kavukcuoglu, K. Pixel recurrent neural networks. *arXiv*, 2016.
- Owens, A., Wu, J., McDermott, J. H., Freeman, W. T., and Torralba, A. Ambient sound provides supervision for visual learning. In *ECCV*, pp. 801–816. Springer, 2016.
- Pathak, D., Krahenbuhl, P., Donahue, J., Darrell, T., and Efros, A. A. Context encoders: Feature learning by inpainting. pp. 2536–2544, 2016.
- Pathak, D., Agrawal, P., Efros, A. A., and Darrell, T. Curiosity-driven exploration by self-supervised prediction. *ICML*, 2017.
- Radford, A., Metz, L., and Chintala, S. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv*, 2015.
- Roth, K., Lucchi, A., Nowozin, S., and Hofmann, T. Stabilizing training of generative adversarial networks through regularization. In *Advances in Neural Information Processing Systems*, pp. 2018–2028, 2017.
- Roweis, S. and Brody, C. Linear heteroencoders, 1999.
- Schroff, F., Kalenichenko, D., and Philbin, J. Facenet: A unified embedding for face recognition and clustering. 2015.
- Sohl-Dickstein, J., Weiss, E. A., Maheswaranathan, N., and Ganguli, S. Deep unsupervised learning using nonequilibrium thermodynamics. *arXiv*, 2015.
- Song, L., Huang, J., Smola, A., and Fukumizu, K. Hilbert space embeddings of conditional distributions with applications to dynamical systems. In *ICML*, pp. 961–968. ACM, 2009.
- Székely, G. J., Rizzo, M. L., Bakirov, N. K., et al. Measuring and testing dependence by correlation of distances. *The Annals of Statistics*, 35(6):2769–2794, 2007.
- Uria, B., Murray, I., and Larochelle, H. A deep and tractable density estimator. In *ICML*, pp. 467–475, 2014.
- Vapnik, V. *Statistical learning theory*. Wiley, New York, 1998.
- Vincent, P., Larochelle, H., Lajoie, I., Bengio, Y., and Manzagol, P.-A. Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion. *JMLR*, 11(Dec):3371–3408, 2010.

Wang, X. and Gupta, A. Unsupervised learning of visual representations using videos. In *ICCV*, pp. 2794–2802, 2015.

Yoon, J., Jordon, J., and van der Schaar, M. Gain: Missing data imputation using generative adversarial nets. *arXiv*, 2018.

Zhang, N., Paluri, M., Ranzato, M., Darrell, T., and Bourdev, L. Panda: Pose aligned networks for deep attribute modeling. In *CVPR*, pp. 1637–1644, 2014.

Zhang, R., Isola, P., and Efros, A. A. Colorful image colorization. In *ECCV*, pp. 649–666. Springer, 2016.

Zhao, J., Mathieu, M., Goroshin, R., and Lecun, Y. Stacked what-where auto-encoders. *arXiv*, 2015.