

Data Generation as Sequential Decision Making

Authored by:

Doina Precup
Philip Bachman

Abstract

We connect a broad class of generative models through their shared reliance on sequential decision making. Motivated by this view, we develop extensions to an existing model, and then explore the idea further in the context of data imputation – perhaps the simplest setting in which to investigate the relation between unconditional and conditional generative modelling. We formulate data imputation as an MDP and develop models capable of representing effective policies for it. We construct the models using neural networks and train them using a form of guided policy search. Our models generate predictions through an iterative process of feedback and refinement. We show that this approach can learn effective policies for imputation problems of varying difficulty and across multiple datasets.

1 Paper Body

Directed generative models are naturally interpreted as specifying sequential procedures for generating data. We traditionally think of this process as sampling, but one could also view it as making sequences of decisions for how to set the variables at each node in a model, conditioned on the settings of its parents, thereby generating data from the model. The large body of existing work on reinforcement learning provides powerful tools for addressing such sequential decision making problems. We encourage the use of these tools to understand and improve the extended processes currently driving advances in generative modelling. We show how sequential decision making can be applied to general prediction tasks by developing models which construct predictions by iteratively refining a working hypothesis under guidance from exogenous input and endogenous feedback. We begin this paper by reinterpreting several recent generative models as sequential decision making processes, and then show how changes inspired by this point of view can improve the performance of the LSTM-based model introduced in [3]. Next, we explore the connections between directed generative models and reinforcement learning more fully by developing an approach to training policies for sequential data imputation. We base our approach on formulating imputation as a finitehorizon Markov Decision Process

which one can also interpret as a deep, directed graphical model. We propose two policy representations for the imputation MDP. One extends the model in [3] by inserting an explicit feedback loop into the generative process, and the other addresses the MDP more directly. We train our models/policies using techniques motivated by guided policy search [9, 10, 11, 8]. We examine their qualitative and quantitative performance across imputation problems covering a range of difficulties (i.e. different amounts of data to impute and different missingness mechanisms?), and across multiple datasets. Given the relative paucity of existing approaches to the general imputation problem, we compare our models to each other and to two simple baselines. We also test how our policies perform when they use fewer/more steps to refine their predictions. As imputation encompasses both classification and standard (i.e. unconditional) generative modelling, our work suggests that further study of models for the general imputation problem is worthwhile. The performance of our models suggests that sequential stochastic construction of predictions, guided by both input and feedback, should prove useful for a wide range of problems. Training these models can be challenging, but lessons from reinforcement learning may bring some relief. 1

2

Directed Generative Models as Sequential Decision Processes

Directed generative models have grown in popularity relative to their undirected counter-parts [6, 14, 12, 4, 5, 16, 15] (etc.). Reasons include: the development of efficient methods for training them, the ease of sampling from them, and the tractability of bounds on their log-likelihoods. Growth in available computing power compounds these benefits. One can interpret the (ancestral) sampling process in a directed model as repeatedly setting subsets of the latent variables to particular values, in a sequence of decisions conditioned on preceding decisions. Each subsequent decision restricts the set of potential outcomes for the overall sequence. Intuitively, these models encode stochastic procedures for constructing plausible observations. This section formally explores this perspective. 2.1

Deep AutoRegressive Networks

The deep autoregressive networks investigated in [4] define distributions of the following form: $T \times Y \ p(x) = p(x-z)p(z)$, with $p(z) = p_0(z_0) \prod_{t=1}^T p_t(z_t - z_0, \dots, z_{t-1})$ (1)

$t=1$

in which x indicates a generated observation and z_0, \dots, z_T represent latent variables in the model. The distribution $p(x-z)$ may be factored similarly to $p(z)$. The form of $p(z)$ in Eqn. 1 can represent arbitrary distributions over the latent variables, and the work in [4] mainly concerned approaches to parameterizing the conditionals $p_t(z_t - z_0, \dots, z_{t-1})$ that restricted representational power in exchange for computational tractability. To appreciate the generality of Eqn. 1, consider using z_t that are univariate, multivariate, structured, etc. One can interpret any model based on this sequential factorization of $p(z)$ as a non-stationary policy $p_t(z_t - s_t)$ for selecting each action z_t in a state s_t , with each s_t determined by all z_{t_0} for $t_0 \leq t$, and train it using some

form of policy search. 2.2

Generalized Guided Policy Search

We adopt a broader interpretation of guided policy search than one might initially take from, e.g., [9, 10, 11, 8]. We provide a review of guided policy search in the supplementary material. Our expanded definition of guided policy search includes any optimization of the general form:

$$\text{minimize } \mathbb{E} \mathbb{E} [\ell(\tau, i_q, i_p)] + \lambda \text{div}(q(\tau \sim i_q, i_p), p(\tau \sim i_p)) \quad (2) \quad p, q$$

$$i_q \sim I_q, i_p \sim I_p(\tau \sim i_q) \quad \tau \sim q(\tau \sim i_q, i_p)$$

in which p indicates the primary policy, q indicates the guide policy, I_q indicates a distribution over information available only to q , I_p indicates a distribution over information available to both p and q , $\ell(\tau, i_q, i_p)$ computes the cost of trajectory τ in the context of i_q/i_p , and $\text{div}(q(\tau \sim i_q, i_p), p(\tau \sim i_p))$ measures dissimilarity between the trajectory distributions generated by p/q . As $\lambda \rightarrow 0$ goes to infinity, Eqn. 2 enforces the constraint $p(\tau \sim i_p) = q(\tau \sim i_q, i_p)$, $\lambda \rightarrow \infty$, i_q, i_p . Terms for controlling, e.g., the entropy of p/q can also be added. The power of the objective in Eq. 2 stems from two main points: the guide policy q can use information i_q that is unavailable to the primary policy p , and the primary policy need only be trained to minimize the dissimilarity term $\text{div}(q(\tau \sim i_q, i_p), p(\tau \sim i_p))$. For example, a directed model structured as in Eqn. 1 can be interpreted as specifying a policy for a finite-horizon MDP whose terminal state distribution encodes $p(x)$. In this MDP, the state at time $1 \leq t \leq T+1$ is determined by $\{z_0, \dots, z_{t-1}\}$. The policy picks an action $z_t \in \mathcal{Z}_t$ at time $1 \leq t \leq T$, and picks an action $x \in \mathcal{X}$ at time $t = T+1$. I.e., the policy can be written as $p_t(z_t \sim z_0, \dots, z_{t-1})$ for $1 \leq t \leq T$, and as $p(x \sim z_0, \dots, z_T)$ for $t = T+1$. The initial state $z_0 \sim Z_0$ is drawn from $p_0(z_0)$. Executing the policy for a single trial produces a trajectory $\tau = \{z_0, \dots, z_T, x\}$, and the distribution over x s from these trajectories is just $p(x)$ in the corresponding directed generative model. The authors of [4] train deep autoregressive networks by maximizing a variational lower bound on the training set log-likelihood. To do this, they introduce a variational distribution q which provides $q_0(z_0 \sim x)$ and $q_t(z_t \sim z_0, \dots, z_{t-1}, x)$ for $1 \leq t \leq T$, with the final step $q(x \sim z_0, \dots, z_T, x)$ given by a Dirac-delta at x . Given these definitions, the training in [4] can be interpreted as guided policy search for the MDP described in the previous paragraph. Specifically, the variational distribution q provides a guide policy $q(\tau \sim x)$ over trajectories $\tau = \{z_0, \dots, z_T, x\}$: $T \prod_{t=1}^T q_t(\tau \sim x)$, $q(x \sim z_0, \dots, z_T, x) q_0(z_0 \sim x) \prod_{t=1}^T q_t(z_t \sim z_0, \dots, z_{t-1}, x)$ (3) $t=1$

2

The primary policy p generates trajectories distributed according to: $T \prod_{t=1}^T p_t(\tau)$, $p(x \sim z_0, \dots, z_T) p_0(z_0) \prod_{t=1}^T p_t(z_t \sim z_0, \dots, z_{t-1})$

(4)

$t=1$

which does not depend on x . In this case, x corresponds to the guide-only information $i_q \sim I_q$ in Eqn. 2. We now rewrite the variational optimization as:

$$\text{minimize } \mathbb{E} \mathbb{E} [\ell(\tau, x)] + \text{KL}(q(\tau \sim x) \parallel p(\tau)) \quad (5) \quad p, q$$

$$x \sim \mathcal{D}\mathcal{X}$$

$$\int q(z \rightarrow x) dz$$

where $q(z, x)$, 0 and DX indicates the target distribution for the terminal state of the primary policy p . When expanded, the KL term in Eqn. 5 becomes: $KL(q(z \rightarrow x) \parallel p(z)) = \int \int q_0(z_0 \rightarrow x) \prod_{t=1}^T q_t(z_t \rightarrow z_{t-1}, x) \log p(z_0) \prod_{t=1}^T p(z_t \rightarrow z_{t-1}, x) dz_0 \dots dz_T$. Thus, the variational approach used in [4] for training directed generative models can be interpreted as a form of generalized guided policy search. As the form in Eqn. 1 can represent any finite directed generative model, the preceding derivation extends to all models we discuss in this paper.

Time-reversible Stochastic Processes

One can simplify Eqn. 1 by assuming suitable forms for X and Z_0, \dots, Z_T . E.g., the authors of [16] proposed a model in which $Z_t = X$ for all t and $p_0(x_0)$ was Gaussian. We can write their model as: $p(x_T) =$

$$\int \prod_{t=1}^T p(x_t \rightarrow x_{t-1}) p_0(x_0) dx_0, \dots, x_{T-1} \quad (7)$$

where $p(x_T)$ indicates the terminal state distribution of the non-stationary, finite-horizon Markov process determined by $\{p_0(x_0), p_1(x_1 \rightarrow x_0), \dots, p_T(x_T \rightarrow x_{T-1})\}$. Note that, throughout this paper, we (ab)use sums over latent variables and trajectories which could/should be written as integrals. The authors of [16] observed that, for any reasonably smooth target distribution DX and sufficiently large T , one can define a "reverse-time" stochastic process $q_t(x_{t-1} \rightarrow x_t)$ with simple, time-invariant dynamics that transforms $q(x_T)$, DX into the Gaussian distribution $p_0(x_0)$. This q is given by: $q_0(x_0) =$

$$\int \prod_{t=1}^T q_1(x_0 \rightarrow x_1) DX(x_T) dx_1, \dots, x_{T-1} \quad (8)$$

Next, we define $q(\cdot)$ as the distribution over trajectories $\{x_0, \dots, x_T\}$ generated by the reverse-time process determined by $\{q_1(x_0 \rightarrow x_1), \dots, q_T(x_T \rightarrow x_{T-1}), DX(x_T)\}$. We define $p(\cdot)$ as the distribution over trajectories generated by the "forward-time" process in Eqn. 7. The training in [16] is equivalent to guided policy search using guide trajectories sampled from q , i.e. it uses the objective: $\int \prod_{t=1}^T q_1(x_0 \rightarrow x_1) \prod_{t=1}^T p_t(x_t \rightarrow x_{t-1}) \log p_0(x_0) \log p(x_T \rightarrow x_{T-1})$ minimize $E \log p + \log p_0(x_0) \log p(x_T \rightarrow x_{T-1})$ (9) which corresponds to minimizing $KL(q \parallel p)$. If the log-densities in Eqn. 9 are tractable, then this minimization can be done using

basic Monte-Carlo. If, as in [16], the reverse-time process q is not PT trained, then Eqn. 9 simplifies to: minimize $\mathbb{E}_q \left[\sum_{t=1}^T \log p_0(x_0) + \log p_t(x_t - x_{t-1}) \right]$. This trick for generating guide trajectories exhibiting a particular distribution over terminal states x_T i.e. running dynamics backwards in time starting from x_T may prove useful in settings other than those considered in [16]. E.g., the LapGAN model in [1] learns to approximately invert a fixed (and information destroying) reverse-time process. The supplementary material expands on the content of this subsection, including a derivation of Eqn. 9 as a bound on $\mathbb{E}_q \left[\sum_{t=1}^T \log p(x_t) \right]$. We could pull the $\sum_{t=1}^T \log p(x_t - z_0, \dots, z_T)$ term from the KL and put it in the cost $\ell(\theta, x)$, but we prefer the path-wise KL formulation for its elegance. We abuse notation using $\text{KL}(q(x = x^*) \parallel p(x))$, $\sum_{t=1}^T \log p(x_t)$. This also includes all generative models implemented and executed on an actual computer. 1

3

2.4

Learning Generative Stochastic Processes with LSTMs

The authors of [3] introduced a model for sequentially-deep generative processes. We interpret their model as a primary policy p which generates trajectories $\theta, \{z_0, \dots, z_T, x\}$ with distribution: $p(\theta), p(x - s^* | \theta_{1:T}) p_0(z_0)$

T Y

$p_t(z_t)$, with $\theta_{1:T}, \{z_0, \dots, z_T\}$

(10)

$t=1$

in which $\theta_{1:T}$ indicates a latent trajectory and $s^* | \theta_{1:T}$ indicates a state trajectory $\{s_0, \dots, s_T\}$ computed recursively from $\theta_{1:T}$ using the update $s_t = f(s_{t-1}, z_t)$ for $t = 1$. The initial state s_0 is given by a trainable constant. Each state s_t , $[h_t; v_t]$ represents the joint hidden/visible state h_t/v_t of an LSTM and f (state, input) computes a standard LSTM update. The authors of [3] defined all $p_t(z_t)$ as isotropic Gaussians and defined the output distribution $p(x - s^* | \theta_{1:T})$ as $p(x - c_T)$, where $P_T c_T, c_0 + \sum_{t=1}^T v_t$. Here, c_0 is a trainable constant and $\sum_{t=1}^T v_t$ is, e.g., an affine transform of v_t . Intuitively, $\sum_{t=1}^T v_t$ transforms v_t into a refinement of the working hypothesis c_{t-1} , which gets updated to $c_t = c_{t-1} + \sum_{t=1}^T v_t$. p is governed by parameters θ which affect $f, \sum_{t=1}^T v_t, s_0$, and c_0 . The supplementary material provides pseudo-code and an illustration for this model. To train p , the authors of [3] introduced a guide policy q with trajectory distribution: $q(\theta - x^*), q(x - s^* | \theta_{1:T}, x^*) q_0(z_0 - x^*)$

T Y

$q_t(z_t - s_t, x^*)$, with $\theta_{1:T}, \{z_0, \dots, z_T\}$

(11)

$t=1$

in which $s^* | \theta_{1:T}$ indicates a state trajectory $\{s_0, \dots, s_T\}$ computed recursively from $\theta_{1:T}$ using the guide policy's state update $s_t = f(s_{t-1}, g(s_{t-1}, x^*))$. In this update s_{t-1} is the previous guide state and $g(s_{t-1}, x^*)$ is a deterministic function of x^* and the partial (primary) state

was 85.1. When the model used the alternate $cT, L? (hT)$, the raw/finetuned test scores were 85.9/85.3. Fig. 1 shows samples from the model. Model/test code is available at [http://github.com/Philip-Bachman/ Sequential-Generation](http://github.com/Philip-Bachman/Sequential-Generation).

3

Figure 1: The left block shows $? (ct)$ for $t \in \{1, 3, 5, 9, 16\}$, for a policy p with $ct, c0 + Pt 0 t0 = 1 L? (vt)$. The right block is analogous, for a model using $ct, L? (ht)$.

Developing Models for Sequential Imputation

The goal of imputation is to estimate $p(x_u \text{---} x_k)$, where $x, [x_u ; x_k]$ indicates a complete observation with known values x_k and missing values x_u . We define a mask $m \in M$ as a (disjoint) partition of x into x_u / x_k . By expanding x_u to include all of x , one recovers standard generative modelling. By shrinking x_u to include a single element of x , one recovers standard classification/regression. Given distribution DM over $m \in M$ and distribution DX over $x \in X$, the objective for imputation is:

$$\text{minimize } E_{x \sim DX} E_{m \sim DM} \log p(x_u \text{---} x_k) \quad (14)$$

We now describe a finite-horizon MDP for which guided policy search minimizes a bound on the objective in Eqn. 14. The MDP is defined by mask distribution DM , complete observation distribution DX , and the state spaces $\{Z_0, \dots, Z_T\}$ associated with each of T steps. Together, DM and DX define a joint distribution over initial states and rewards in the MDP. For the trial determined by $x \sim DX$ and $m \sim DM$, the initial state $z_0 \sim p(z_0 \text{---} x_k)$ is selected by the policy p based on the known values x_k . The cost $c(z, x_u, x_k)$ suffered by trajectory $z, \{z_0, \dots, z_T\}$ in the context (x, m) is given by $-\log p(x_u \text{---} z, x_k)$, i.e. the negative log-likelihood of p guessing the missing values x_u after following trajectory z , while seeing the known values x_k . QT We consider a policy p with trajectory distribution $p(z \text{---} x_k), p(z_0 \text{---} x_k) t=1 p(z_t \text{---} z_0, \dots, z_{t-1}, x_k)$, where x_k is determined by x/m for the current trial and p can't observe the missing values x_u . With these definitions, we can find an approximately optimal imputation policy by solving:

$$\text{minimize } E_{x \sim DX} E_{m \sim DM} E_{z \sim p(z \text{---} x_k)} [-\log p(x_u \text{---} z, x_k)] \quad (15)$$

I.e. the expected negative log-likelihood of making a correct imputation on any given trial. This is a valid, but loose, upper bound on the imputation objective in Eq. 14 (from Jensen's inequality). We can tighten the bound by introducing a guide policy (i.e. a variational distribution). As with the unconditional generative models in Sec. 2, we train p to imitate a guide policy q shaped by additional information (here it's x_u). This q generates trajectories with distribution $q(z \text{---} x_u, x_k), QT q(z_0 \text{---} x_u, x_k) t=1 q(z_t \text{---} z_0, \dots, z_{t-1}, x_u, x_k)$. Given this p and q , guided policy search solves:

$$\text{minimize } E_{x \sim DX} E_{m \sim DM} E_{z \sim p(z \text{---} x_k)} [-\log q(x \text{---} z, i_q, i_p)] + KL(q(z \text{---} i_q, i_p) \text{---} p(z \text{---} i_p)) \quad (16)$$

$$p, q$$

$$x \sim DX, m \sim DM$$

$$z \sim p(z \text{---} i_q, i_p)$$

$$\text{where we define } i_q, x_u, i_p, x_k, \text{ and } q(x_u \text{---} z, i_q, i_p), p(x_u \text{---} z, i_p). \quad 56$$

Data splits from: http://www.cs.toronto.edu/~larocheh/public/datasets/binarized_mnist
The model in [3] significantly improves its score to 80.97 when using an image-specific architecture.

5

3.1

A Direct Representation for Sequential Imputation Policies

We define an imputation trajectory as $c^? , \{c_0, \dots, c_T\}$, where each partial imputation $c_t^? \in X$ is computed from a partial step trajectory $z_t^? , \{z_1, \dots, z_t\}$. A partial imputation $c_{t-1}^?$ encodes the policy's guess for the missing values x_u immediately prior to selecting step z_t , and c_T gives the policy's final guess. At each step of iterative refinement, the policy selects a z_t based on $c_{t-1}^?$ and the known values x_k , and then updates its guesses to c_t based on $c_{t-1}^?$ and z_t . By iteratively refining its guesses based on feedback from earlier guesses and the known values, the policy can construct complexly structured distributions over its final guess c_T after just a few steps. This happens naturally, without any post-hoc MRFs/CRFs (as in many approaches to structured prediction), and without sampling values in c_T one at a time (as required by existing NADE-type models [7]). This property of our approach should prove useful for many tasks. We consider two ways of updating the guesses in c_t , mirroring those described in Sec. 2. The first way sets $c_t^? = c_{t-1}^? + \theta(z_t)$, where $\theta(z_t)$ is a trainable function. We set $c_0, [c_{u0}; c_{k0}]$ using a trainable bias. The second way sets $c_t^? = \theta(z_t)$. We indicate models using the first type of update with the suffix -add, and models using the second type of update with -jump. Our primary policy $p^?$ selects z_t at each step $1 \leq t \leq T$ using $p^?(z_t - c_{t-1}^?, x_k)$, which we restrict to be a diagonal Gaussian. This is a simple, stationary policy. Together, the step selector $p^?(z_t - c_{t-1}^?, x_k)$ and the imputation constructor $\theta^?(z_t)$ fully determine the behaviour of the primary policy. The supplementary material provides pseudo-code and an illustration for this model. We construct a guide policy q similarly to p . The guide policy shares the imputation constructor $\theta^?(z_t)$ with the primary policy. The guide policy incorporates additional information $x, [x_u; x_k]$, i.e. the complete observation for which the primary policy must reconstruct some missing values. The guide policy chooses steps using $q^?(z_t - c_{t-1}^?, x)$, which we restrict to be a diagonal Gaussian. We train the primary/guide policy components $\theta^?, p^?$, and $q^?$ simultaneously on the objective:

$$\min_{\theta, p, q} E_{x \sim p} [-\log q(x - c_T)] + KL(q(\theta - x, x) - p(\theta - x)) \quad (17)$$

$\theta^? \in \mathbb{R}^{D_X \times D_M}$

$\theta^? \in \mathbb{R}^{D_X \times D_M}$

where $q(x_u - c_{uT}), p(x_u - c_{uT})$. We train our models using Monte-Carlo roll-outs of q , and stochastic backpropagation as in [6, 14]. Full implementations and test code are available from <http://github.com/Philip-Bachman/Sequential-Generation>.

3.2 Representing Sequential Imputation Policies using LSTMs

To make it useful for imputation, which requires conditioning on the exogenous information x_k , we modify the LSTM-based model from Sec. 2.5 to include

a `read` operation in its primary policy p . We incorporate a read operation by spreading p over two LSTMs, p_r and p_w , which respectively `read` and `write` an imputation trajectory c , $\{c_0, \dots, c_T\}$. Conveniently, the guide policy q for this model takes the same form as the primary policy's reader p_r . This model also includes an `infinite mixture` initialization step, as used in Sec. 2.5, but modified to incorporate conditioning on x and m . The supplementary material provides pseudo-code and an illustration for this model. Following the infinite mixture initialization step, a single full step of execution for p involves several k substeps: first p updates the reader state using $srt \leftarrow f_r(srt_{t-1}, p_r(ct_{t-1}, sw_{t-1}, x))$, then p selects a r w w step $zt \leftarrow p(zt \leftarrow vt)$, then p updates the writer state using $st \leftarrow f_w(st_{t-1}, zt)$, and finally p updates r, w its guesses by setting $ct \leftarrow ct_{t-1} + \eta_w(vtw)$ (or $ct \leftarrow \eta_w(hw, [hr, w_t; vt]_t)$). In these updates, st_r, w refer to the states of the (r)reader and (w)writer LSTMs. The LSTM updates f and the read/write operations $\eta_{r,w}$ are governed by the policy parameters θ . We train p to imitate trajectories sampled from a guide policy q . The guide policy shares the primary policy's writer updates f_w and write operation η_w , but has its own reader updates f_q and read operation η_q . At each step, the guide policy: updates the guide state $sqt \leftarrow f_q(sqt_{t-1}, \eta_q(ct_{t-1}, sw_{t-1}, x))$, w w then selects $zt \leftarrow q(zt \leftarrow vt_q)$, then updates the writer state $sw \leftarrow f_w(s, z)$, and finally updates $t_{t-1} \leftarrow$ its guesses $ct \leftarrow ct_{t-1} + \eta_w(vtw)$ (or $ct \leftarrow \eta_w(hw)$). As in Sec. 3.1, the guide policy's read operation η_q gets to see the complete observation x , while the primary policy only gets to see the known values x_k . We restrict the step distributions p_{θ}/q_{θ} to be diagonal Gaussians whose means and log-variances are affine functions of vtr/vtq . The training objective has the same form as Eq. 17. 6

350 300 250 200

Imputation NLL vs. Available Information

88

TM-orc TM-hon VAE-imp GPSI-add GPSI-jump LSTM-add LSTM-jump

86 84 82

Imputation NLL vs. Available Information

98

GPSI-add GPSI-jump LSTM-add LSTM-jump

94 92 90

78 88

76

86

74 100

84

72 0.60

0.65

GPSI-add GPSI-jump

80

150

50 0.55

The Effect of Increased Refinement Steps

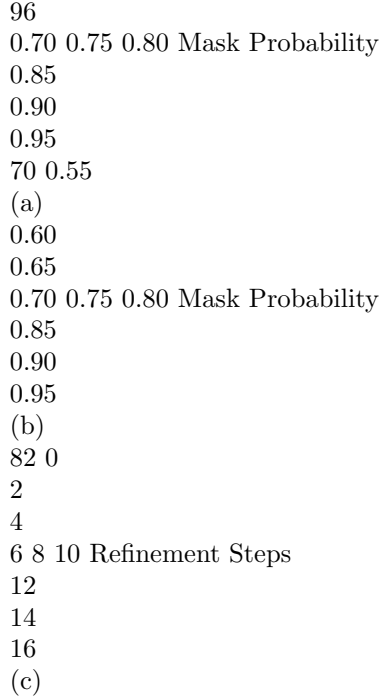


Figure 2: (a) Comparing the performance of our imputation models against several baselines, using MNIST digits. The x-axis indicates the % of pixels which were dropped completely at random, and the scores are normalized by the number of imputed pixels. (b) A closer view of results from (a), just for our models. (c) The effect of increased iterative refinement steps for our GPSI models.

Experiments

We tested the performance of our sequential imputation models on three datasets: MNIST (28x28), SVHN (cropped, 32x32) [13], and TFD (48x48) [17]. We converted images to grayscale and shift/scaled them to be in the range [0...1] prior to training/testing. We measured the imputation log-likelihood $\log q(x_u | c_u)$ using the true missing values x_u and the models' guesses given by c_u . We report negative log-likelihoods, so lower scores are better in all of our tests. We refer to variants of the model from Sec. 3.1 as GPSI-add and GPSI-jump, and to variants of the model from Sec. 3.2 as LSTM-add and LSTM-jump. Except where noted, the GPSI models used 6 refinement steps and the LSTM models used 16.7 We tested imputation under two types of data masking: missing completely at random (MCAR) and missing at random (MAR). In MCAR, we masked pixels uniformly at random from the source images, and indicate removal of $d\%$ of the pixels by MCAR- d . In MAR, we masked square regions, with the occlusions located uniformly at random within the borders of the source image. We indicate occlusion of a $d \times d$ square by

MAR-d. On MNIST, we tested MCAR-d for $d \in \{50, 60, 70, 80, 90\}$. MCAR-100 corresponds to unconditional generation. On TFD and SVHN we tested MCAR-80. On MNIST, we tested MAR-d for $d \in \{14, 16\}$. On TFD we tested MAR-25 and on SVHN we tested MAR-17. For test trials we sampled masks from the same distribution used in training, and we sampled complete observations from a held-out test set. Fig. 2 and Tab. 1 present quantitative results from these tests. Fig. 2(c) shows the behavior of our GPSI models when we allowed them fewer/more refinement steps.

LSTM-add LSTM-jump GPSI-add GPSI-jump VAE-imp

MNIST MAR-14 MAR-16 170 167 172 169 177 175 183 177 374 394

TFD MCAR-80 MAR-25 1381 1377 ? ? 1390 1380 1394 1384 1416 1399

SVHN MCAR-80 MAR-17 525 568 ? ? 531 569 540 572 567 624

Table 1: Imputation performance in various settings. Details of the tests are provided in the main text. Lower scores are better. Due to time constraints, we did not test LSTM-jump on TFD or SVHN. These scores are normalized for the number of imputed pixels. We tested our models against three baselines. The baselines were ?variational auto-encoder imputation?, honest template matching, and oracular template matching. VAE imputation ran multiple steps of VAE reconstruction, with the known values held fixed and the missing values re-estimated with each reconstruction step.⁸ After 16 refinement steps, we scored the VAE based on its best 7 GPSI stands for ?Guided Policy Search Imputer?. The tag ?-add? refers to additive guess updates, and ?-jump? refers to updates that fully replace the guesses. ⁸ We discuss some deficiencies of VAE imputation in the supplementary material.

7

(a)

(b)

(c)

Figure 3: This figure illustrates the policies learned by our models. (a): models trained for (MNIST, MAR-16). From top?bottom the models are: GPSI-add, GPSI-jump, LSTM-add, LSTM-jump. (b): models trained for (TFD, MAR-25), with models in the same order as (a) ? but without LSTMjump. (c): models trained for (SVHN, MAR-17), with models arranged as for (b). guesses. Honest template matching guessed the missing values based on the training image which best matched the test image?s known values. Oracular template matching was like honest template matching, but matched directly on the missing values. Our models significantly outperformed the baselines. In general, the LSTM-based models outperformed the more direct GPSI models. We evaluated the log-likelihood of imputations produced by our models using the lower bounds provided by the variational objectives with respect to which they were trained. Evaluating the template-based imputations was straightforward. For VAE imputation, we used the expected log-likelihood of the imputations sampled from multiple runs of the 16-step imputation process. This provides a valid, but loose, lower bound on their log-likelihood. As shown in Fig. 3, the imputations produced by our models appear promising. The imputations are generally of high quality, and the models are capable of capturing strongly

multi-modal reconstruction distributions (see subfigure (a)). The behavior of GPSI models changed intriguingly when we swapped the imputation constructor. Using the -jump imputation constructor, the imputation policy learned by the direct model was rather inscrutable. Fig. 2(c) shows that additive guess updates extracted more value from using more refinement steps. When trained on the binarized MNIST benchmark discussed in Sec. 2.5, i.e. with binarized images and subject to MCAR-100, the LSTMadd model produced raw/fine-tuned scores of 86.2/85.7. The LSTM-jump model scored 87.1/86.3. Anecdotally, on this task, these ‘closed-loop’ models seemed more prone to overfitting than the ‘open-loop’ models in Sec. 2.5. The supplementary material provides further qualitative results.

5

Discussion

We presented a point of view which links methods for training directed generative models with policy search in reinforcement learning. We showed how our perspective can guide improvements to existing models. The importance of these connections will only grow as generative models rapidly increase in structural complexity and effective decision depth. We introduced the notion of imputation as a natural generalization of standard, unconditional generative modelling. Depending on the relation between the data-to-generate and the available information, imputation spans from full unconditional generative modelling to classification/regression. We showed how to successfully train sequential imputation policies comprising millions of parameters using an approach based on guided policy search [9]. Our approach outperforms the baselines quantitatively and appears qualitatively promising. Incorporating, e.g., the local read/write mechanisms from [3] should provide further improvements. 8

2 References

- [1] Emily L Denton, Soumith Chintala, Arthur Szlam, and Robert Fergus. Deep generative models using a laplacian pyramid of adversarial networks. arXiv:1506.05751 [cs.CV], 2015.
- [2] Alex Graves. Generating sequences with recurrent neural networks. arXiv:1308.0850 [cs.NE], 2013.
- [3] Karol Gregor, Ivo Danihelka, Alex Graves, and Daan Wierstra. Draw: A recurrent neural network for image generation. In International Conference on Machine Learning (ICML), 2015.
- [4] Karol Gregor, Ivo Danihelka, Andriy Mnih, Charles Blundell, and Daan Wierstra. Deep autoregressive networks. In International Conference on Machine Learning (ICML), 2014.
- [5] Diederik P Kingma, Danilo J Rezende, Shakir Mohamed, and Max Welling. Semi-supervised learning with deep generative models. In Advances in Neural Information Processing Systems (NIPS), 2014.
- [6] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In International Conference on Learning Representations (ICLR), 2014.
- [7] Hugo Larochelle and Iain Murray. The neural autoregressive distribution estimator. In International Conference on Machine Learning (ICML), 2011.
- [8] Sergey Levine and Pieter Abbeel. Learning neural network policies

with guided policy search under unknown dynamics. In Advances in Neural Information Processing Systems (NIPS), 2014. [9] Sergey Levine and Vladlen Koltun. Guided policy search. In International Conference on Machine Learning (ICML), 2013. [10] Sergey Levine and Vladlen Koltun. Variational policy search via trajectory optimization. In Advances in Neural Information Processing Systems (NIPS), 2013. [11] Sergey Levine and Vladlen Koltun. Learning complex neural network policies with trajectory optimization. In International Conference on Machine Learning (ICML), 2014. [12] Andriy Mnih and Karol Gregor. Neural variational inference and learning in belief networks. In International Conference on Machine Learning (ICML), 2014. [13] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. NIPS Workshop on Deep Learning and Unsupervised Feature Learning, 2011. [14] Danilo Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In International Conference on Machine Learning (ICML), 2014. [15] Danilo J Rezende and Shakir Mohamed. Variational inference with normalizing flows. In International Conference on Machine Learning (ICML), 2015. [16] Jascha Sohl-Dickstein, Eric A. Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In International Conference on Machine Learning (ICML), 2015. [17] Joshua Susskind, Adam Anderson, and Geoffrey E Hinton. The toronto face database. 2010.