

Sparse-Coding Variational Auto-Encoders

Gabriel Barello¹, Adam S. Charles², Jonathan W. Pillow²

¹Institute of Neuroscience, University of Oregon.

²Princeton Neuroscience Institute & Dept. of Psychology, Princeton University.

Keywords: sparse coding, variational auto-encoders

Abstract

The sparse coding model posits that the visual system has evolved to efficiently code natural stimuli using a sparse set of features from an overcomplete dictionary. The classic sparse coding model suffers from two key limitations, however: (1) computing the neural response to an image patch requires minimizing a nonlinear objective function, which was initially not easily mapped onto a neurally plausible feedforward mechanism and (2) fitting the model to data relied on an approximate inference method that ignores uncertainty. Here we address these two shortcomings by formulating a variational inference method for the sparse coding model inspired by the variational auto-encoder (VAE) framework. The sparse-coding variational auto-encoder (SVAE) augments the classic sparse coding model with a probabilistic recognition model, parametrized by a deep neural network. This recognition model provides a neurally plausible implementation for the mapping from image patches to neural activities, and enables a principled method for fitting the sparse coding model to data via maximization of the evidence lower bound (ELBO). The SVAE differs from the traditional VAE in three important ways: the generative model is the sparse coding model instead of a deep network; the latent representation is overcomplete, with more latent dimensions than image pixels; and the prior over latent variables is a sparse or heavy-tailed instead of Gaussian. We fit the SVAE to natural image data under different assumed prior distributions, and show that it obtains higher test performance than previous fitting methods. Finally, we examine the response properties of the recognition network and show that it captures important nonlinear properties of neurons in the early visual pathway.

1 Introduction

Generative models have played an important role in computational neuroscience by offering normative explanations of observed neural response properties (Olshausen & Field, 1996a,b, 1997; Lewicki & Olshausen, 1999; Dayan et al., 2003; Berkes & Wiskott, 2005; Coen-Cagli et al., 2012). These models seek to model the distribution of stimuli in the world $P(\mathbf{x})$ in terms of a conditional probability distribution $P(\mathbf{x}|\mathbf{z})$, the probability of a stimulus \mathbf{x} given a set of latent variables \mathbf{z} , and a prior over the latent variables $P(\mathbf{z})$. The advantage of this approach is that it mimics the causal structure of the world: an image falling on the retina is generated by “sources” in the world (e.g., the identity and pose of a face, and the light source illuminating it), which are typically latent or hidden from the observer. Perception is naturally formulated as the statistical inference problem of identifying the latent sources that generated a particular sensory stimulus D. Knill & Richards (1996); Weiss et al. (2002); D. C. Knill & Pouget (2004); Moreno-Bote et al. (2011). Mathematically, this

corresponds to applying Bayes' rule to obtain the posterior over latent sources given sensory data: $P(\mathbf{z}|\mathbf{x}) \propto P(\mathbf{x}|\mathbf{z})P(\mathbf{z})$, where the terms on the right-hand-side are the likelihood $P(\mathbf{x}|\mathbf{z})$ and prior $P(\mathbf{z})$, which come from the generative model.

Perhaps the famous generative model in neuroscience is the *sparse coding model*, introduced by Olshausen and Field (Olshausen & Field, 1996a,b) to account for the response properties of neurons in visual cortex. The sparse coding model posits that neural activity represents an estimate of the latent features underlying a natural image patch under a linear generative model. The model's key feature is sparsity: a heavy-tailed prior over the latent variables ensures that neurons are rarely active, so each image patch must be explained by a small number of active features. Remarkably, the feature vectors obtained by fitting this model to natural images resemble the localized, oriented receptive fields found in the early visual cortex (Olshausen & Field, 1996a). Subsequent work showed the model could account for a variety of properties of neural activity in the visual pathway (e.g. classical and non-classical receptive field effects (Rozell et al., 2008; Karklin & Lewicki, 2009; Lee et al., 2007)).

Although the sparse coding model is a linear generative model, recognition (inferring the latent variables from an image) and learning (optimizing the dictionary of features) are both computationally difficult problems. Early work on the sparse coding model did not provide a neurally plausible mechanism for neurons to compute their responses to an image patch. Moreover, fitting the model to data relied on approximate optimization methods.

In this paper, we propose a solution to these two important problems using ideas from the variational auto-encoder (VAE) (Kingma & Welling, 2014; Rezende et al., 2014). The VAE is a framework for training a complex generative model by coupling it to a recognition model parametrized by a deep neural network. This deep network offers tractable inference for latent variables from data and allows for gradient-based learning of the generative model parameters using a variational objective. Here we adapt the VAE methodology to the sparse coding model by adjusting its structure and prior assumptions. We compare the resulting *sparse-coding VAE* (SVAE) to fits using the original methodology, and show that our model achieves higher log-likelihood on test data. Furthermore, we show that the recognition model of the trained SVAE performs accurate inference under the sparse coding model, and captures important response properties of neurons in visual cortex, including orientation tuning, surround suppression, and frequency tuning.

2 Background

2.1 The Sparse Coding Model

The sparse coding model (Olshausen & Field, 1996a,b) posits that the spike response \mathbf{z} from a population of neurons can be interpreted as a set of sparse latent variables underlying an image \mathbf{x} presented to the visual system. This can be formalized by the following generative model:

$$\mathbf{z} \sim P(\mathbf{z}) \quad (\text{prior over neural activity}) \quad (1)$$

$$\mathbf{x} = \Phi\mathbf{z} + \epsilon, \quad (\text{noisy linear mapping to images}) \quad (2)$$

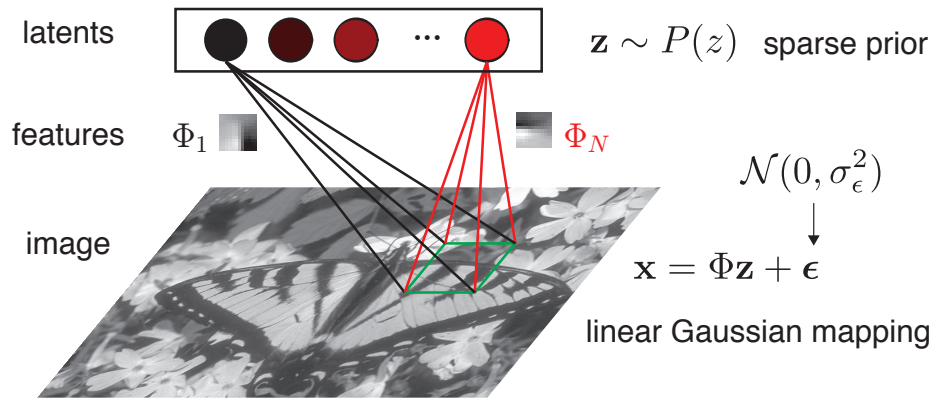


Figure 1: The sparse coding model. A sample from the model is generated by first sampling the latent variables \mathbf{z} from the sparse prior distribution $P(\mathbf{z})$, which provide linear weights on feature vectors Φ_i , and finally adding Gaussian noise. The image shown is taken from the BSDS300 dataset Martin et al. (2001).

where $P(\mathbf{z})$ is a sparse or heavy-tailed prior over neural activities, Φ is matrix of dictionary elements or “features”, and $\epsilon \sim \mathcal{N}(0, \sigma_\epsilon^2 \mathbf{I})$ is isotropic Gaussian noise with variance σ_ϵ^2 . Critically, the sparse coding model is *overcomplete*, meaning that the dimension of the latent variable is larger than that of the data: $\text{Dim}(\mathbf{z}) > \text{Dim}(\mathbf{x})$. As a result, multiple settings of \mathbf{z} can reproduce any given image patch \mathbf{x} .

Although this linear generative model is easy to sample, the sparse prior over \mathbf{z} make both inference and model fitting difficult. Full Bayesian inference involves computing the posterior distribution over latent variables given an image, which is (according to Bayes’ rule):

$$P(\mathbf{z} | \mathbf{x}) = \frac{\mathcal{N}(\mathbf{x} | \Phi \mathbf{z}, \sigma_\epsilon^2 \mathbf{I}) P(\mathbf{z})}{\int \mathcal{N}(\mathbf{x} | \Phi \mathbf{z}, \sigma_\epsilon^2 \mathbf{I}) P(\mathbf{z}) d\mathbf{z}}, \quad (3)$$

where $\mathcal{N}(\mathbf{x} | \Phi \mathbf{z}, \sigma_\epsilon^2 \mathbf{I})$ denotes a multivariate normal density in \mathbf{x} with mean $\Phi \mathbf{z}$ and covariance $\sigma_\epsilon^2 \mathbf{I}$. Unfortunately, the denominator cannot be computed in closed form for most sparsity-promoting priors, such as the Cauchy and Laplace distribution.

Olshausen & Field (1996a) considered a simpler inference problem by setting neural activity equal to the MAP estimate of \mathbf{z} given the image:

$$\hat{\mathbf{z}}_{map} | \mathbf{x} = \arg \max_{\mathbf{z}} P(\mathbf{z} | \mathbf{x}) = \arg \max_{\mathbf{z}} \frac{1}{2\sigma_\epsilon^2} \|\mathbf{x} - \Phi \mathbf{z}\|_2^2 + \log P(\mathbf{z}), \quad (4)$$

which does not depend on the denominator in (eq. 3). Here $\|\mathbf{x} - \Phi \mathbf{z}\|_2^2$ is the sum of squared errors between the image \mathbf{x} and its reconstruction in the basis defined by Φ , and $\log P(\mathbf{z})$ serves as a “sparsity penalty” given by $-\sum_i \log(1 + z_i^2)$ for a Cauchy prior or $-\sum_i |z_i|$ for a Laplace prior. The noise variance σ_ϵ^2 serves to trade off the importance of minimizing the image reconstruction error against the importance of sparsity imposed by the prior. (In the limit $\sigma_\epsilon \rightarrow 0$, the effect of the sparsity penalty $\log P(\mathbf{z})$ vanishes, and $\hat{\mathbf{z}}_{map}$ becomes a least-squares estimate). However, finding $\hat{\mathbf{z}}_{map}$ requires numerical optimization of the right-hand-side of (eq. 4), which does not easily map onto a model of feed-forward processing in the visual pathway.

2.2 Fitting the sparse coding model

The problem of fitting the sparse coding model to data is even more difficult than the inference problem. The maximum likelihood estimate of the model parameters Φ for a dataset of images $X = \{\mathbf{x}_1, \dots, \mathbf{x}_L\}$ is given by

$$\arg \max_{\Phi} P(X|\Phi) = \arg \max_{\Phi} \prod_{l=1}^L \left(\int P(\mathbf{x}_l|\mathbf{z}_l, \Phi) P(\mathbf{z}_l) d\mathbf{z}_l \right), \quad (5)$$

where $Z = \{\mathbf{z}_1, \dots, \mathbf{z}_L\}$ are the latent variables corresponding to the images in X . Once again, the integrals over \mathbf{z}_l have no closed-form solution for the relevant sparsity-promoting priors of interest.

To circumvent this intractable integral, (Olshausen & Field, 1996a) employed an approximate iterative method for optimizing the dictionary Φ . After the initializing the dictionary randomly, iterate:

1. Take a group of training images $X^{(i)}$ and compute the MAP estimate of the latent variables $\hat{Z}^{(i)}$ for each image using the current dictionary $\Phi^{(i)}$:

$$\hat{Z}^{(i)} = \arg \max_Z \log P(X^{(i)}|Z, \Phi) + \log P(Z) \quad (6)$$

2. Update the dictionary using the gradient of the log-likelihood, conditioned on $Z^{(i)}$:

$$\Phi^{(i+1)} = \Phi^{(i)} + \eta \left(\nabla_{\Phi} \log P(X^{(i)}|Z^{(i)}, \Phi) \right) \quad (7)$$

where η is a learning rate.

2.3 Fitting as Variational EM

The iterative fitting algorithm from Olshausen & Field (1996a,b) can be seen to relate closely to a form variational expectation-maximization (EM). Variational EM employs a surrogate distribution $Q(Z)$ to optimize a tractable lower-bound on the log-likelihood known as the *evidence lower bound* (ELBO):

$$\log P(X | \Phi) = \log \int P(X, Z | \Phi) dZ = \log \int Q(Z) \frac{P(X, Z | \Phi)}{Q(Z)} dZ \quad (8)$$

$$\geq \int Q(Z) \log \left(\frac{P(X, Z | \Phi)}{Q(Z)} \right) dZ \quad (9)$$

$$= \log P(X) - D_{KL}[Q(Z) || P(Z|X, \Phi)] \triangleq \mathcal{L}(Q, \Phi), \quad (10)$$

where the inequality in (eq. 9) follows from Jensen's inequality, and $D_{KL}[Q(Z)||P(Z)] = \int Q(Z) \log[Q(Z)/P(Z)]dZ$ denotes the Kullback-Leibler (KL) divergence between two distributions $Q(Z)$ and $P(Z)$ (Bishop, 2005; Blei et al., 2016).

The expectation or ‘‘E’’ step of variational EM involves setting $Q(Z)$ to minimize the KL-divergence between $Q(Z)$ and $P(Z|X, \Phi)$, the posterior over the latent given the data and model parameters. Note that when $Q(Z) = P(Z|X, \Phi)$, the KL term (eq. 10) is zero and

the lower bound is tight, so \mathcal{L} achieves equality with the log-likelihood. The maximization or “M” step involves maximizing $\mathcal{L}(Q, \Phi)$ for Φ with Q held fixed, which is tractable for appropriate choices of surrogate distribution $Q(Z)$. Variational EM can therefore be viewed as alternating coordinate ascent of the ELBO:

$$\begin{aligned} \text{E-step: } \quad Q &= \arg \max_Q \mathcal{L}(Q, \Phi) = \arg \min_Q D_{KL} \left[Q(Z) \parallel P(Z|X, \Phi) \right] \\ \text{M-step: } \quad \Phi &= \arg \max_{\Phi} \mathcal{L}(Q, \Phi) = \arg \min_{\Phi} \mathbb{E}_Q \left[\log P(X, Z|\Phi) \right]. \end{aligned} \quad (11)$$

It is now possible to see that the fitting algorithm of Olshausen & Field (1996a) is closely related to variational EM in which the surrogate distribution $Q(Z)$ is a product Dirac delta functions:

$$Q(Z) = \prod_{l=1}^L \delta(\mathbf{z}_l - \gamma_l), \quad (12)$$

where $\{\gamma_1, \dots, \gamma_L\}$ are the variational parameters governing Q . With a Dirac delta posterior, the EM algorithm is formally undefined. Examining Equation 11, the KL divergence in the E-step between the Dirac delta distribution and the latent variable posterior is infinite, and the ELBO is formally $-\infty$, always. We can see the source of this catastrophe by writing out the KL divergence of the E-step above as

$$\begin{aligned} D_{KL} \left[Q(Z) \parallel P(Z|X, \Phi) \right] &= \int Q(Z) [\log(Q(Z)) - \log(P(Z|X, \Phi))] \\ &= -P(\gamma|X, \Phi) + \int Q(Z) \log(Q(Z)), \end{aligned} \quad (13)$$

where in the second line we have used our delta-function variational distribution $Q(X)$ (eq.12) to evaluate the first term, which is simply the posterior evaluated at γ . The second term is infinite; note however that this infinite term is independent of Φ , and so the gradient of the KL divergence with respect to Φ can be computed sensibly, and optimized to perform the E-step.

If we neglect the infinite term (which is constant in Φ), the ELBO reduces to:

$$\mathcal{L}(Q, \Phi) = \sum_l \int \delta(\mathbf{z}_l - \gamma_l) \log P(\mathbf{x}_l, \mathbf{z}_l|\Phi) d\mathbf{z}_l = \sum_l \log P(\mathbf{x}_l|\gamma_l, \Phi) + \log P(\gamma_l) \quad (14)$$

The E-step involves setting γ_l to the MAP estimate of the latent vector \mathbf{z}_l for each image \mathbf{x}_l , while the M-step involves updating the dictionary Φ to minimize the reconstruction error of all images given the inferred latent variables. The original sparse coding algorithm differed in that it operated on randomly selected sub-sets of the data (stochastic mini-batches), and took a single gradient step in Φ in place of a full M-step at each iteration. The algorithm can therefore be seen as a stochastic gradient descent version of variational EM (R. Neal & Hinton, 1998).

Employing a Dirac delta variational posterior introduces bias whereby Φ can grow without bound to increase the ELBO¹. In the case that a variational distribution with finite variance is used, the norm of Φ is regularized by the KL divergence. To compensate for this

¹Note that the ELBO isn't actually increasing because, as noted above, it is always formally $-\infty$; we could however justify this approach with a careful appeal a finite-variance $Q(x)$ that approaches delta function only in the limit.

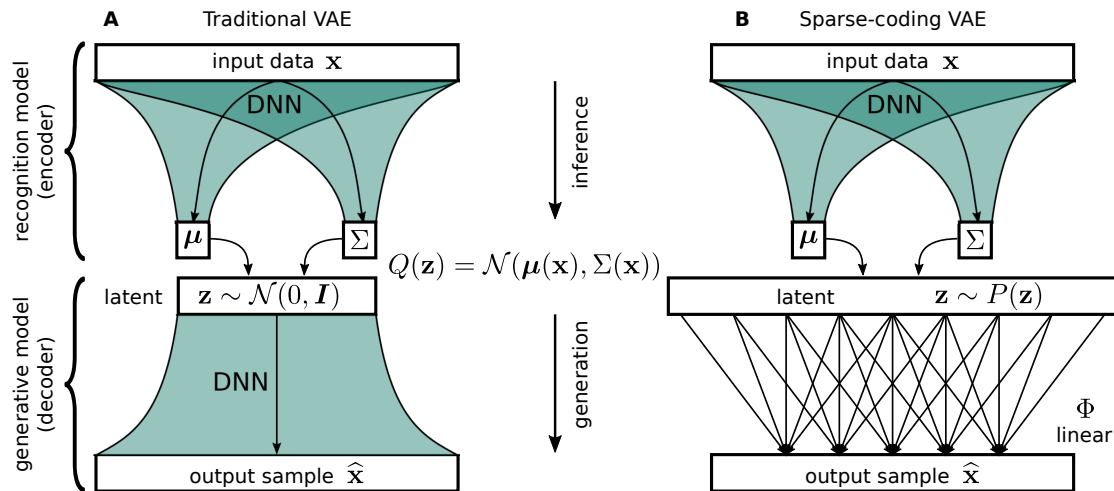


Figure 2: Schematics of traditional VAE and the sparse-coding VAE. **(A)** The traditional VAE consists of a generative model (decoder) with a recognition model (encoder) stacked on top to subserve variational inference. Both models are parametrized by deep neural networks with Gaussian noise. **(B)** The sparse-coding VAE results from replacing the generative model by the sparse coding model. It differs from the original VAE in three key respects: the latent variable z has higher dimension than the data x ; the prior over z is sparse (e.g., Cauchy or Laplace) instead of Gaussian; and the generative model is linear instead of a deep neural network.

bias, Olshausen and Field introduced a post-hoc normalization procedure which effectively matches the z variance to the data variance, regularizing Φ in the process. It is not clear that this post-hoc normalization can be understood in the variational EM framework, though it bears some resemblance to a proximal gradient step. Other normalization methods have been proposed, including placing priors on Φ , which could more easily be mapped onto the variational EM framework (see Discussion).

2.4 Variational Auto-Encoders (VAEs)

The variational auto-encoder is a powerful framework for training complex generative models (Kingma & Welling, 2014; Rezende et al., 2014; Doersch, 2016). In the first papers describing the VAE, the generative model was defined by a standard Gaussian latent variable mapped through a deep neural network with additive Gaussian noise:

$$z \sim \mathcal{N}(0, I) \quad (15)$$

$$x = G_{\theta}(z) + \epsilon, \quad (16)$$

where $G_{\theta}(z)$ denotes the output of a deep neural network with parameters θ , and $\epsilon \sim \mathcal{N}(0, \sigma_{\epsilon}^2 I)$ represents Gaussian noise added to the output of this network. By using a sufficiently large, deep network, this model should be able to generate data with arbitrarily complicated continuous distributions.

The key idea of the VAE framework is to perform variational inference for this model using a variational distribution parametrized by another pair of deep neural networks:

$$Q(z) = \mathcal{N}(\mu_{\gamma}(x), \Sigma_{\gamma}(x)), \quad (17)$$

where $\mu_\gamma(\mathbf{x})$ and $\Sigma_\gamma(\mathbf{x})$ denote neural networks, with parameters $\gamma = \{\gamma_\mu, \gamma_\sigma\}$, that map data samples \mathbf{x} to the mean and covariance (respectively) of a Gaussian variational distribution over \mathbf{z} given \mathbf{x} . This pair of networks provide an approximate recognition model for inferring \mathbf{z} from \mathbf{x} under the generative model (eqs. 15-16). In VAE parlance, $\mu_\gamma(\cdot)$ and $\Sigma_\gamma(\cdot)$ comprise the *encoder*, mapping data samples to distributions over the latent variables, while the generative network $G_\theta(\cdot)$ is the *decoder*, mapping latent variables to the data space. This terminology is inspired by the model's structural similarity to classic auto-encoders, and is consistent with the fact that \mathbf{z} is typically lower-dimensional than \mathbf{x} , providing a compressed representation of the structure contained in the data when using the encoder. Note that the analogy is imprecise, however: the output of the encoder in the VAE is not a point \mathbf{z} in latent space, but an entire distribution over latent variables $Q(\mathbf{z})$.

Fitting the VAE to data involves stochastic gradient ascent of the ELBO (eq. 10) for model parameters θ and variational parameters γ . This can be made practical using several clever tricks. The first trick is to evaluate the expectation over $Q(\mathbf{z})$ as a Monte Carlo integral. The contribution to the ELBO from a single data sample \mathbf{x} is therefore:

$$\mathcal{L}(Q, \theta|\mathbf{x}) = \int Q(\mathbf{z}) \log \frac{P(\mathbf{x}, \mathbf{z}|\theta)}{Q(\mathbf{z})} d\mathbf{z} \approx \frac{1}{m} \sum_{j=1}^m \log \frac{P(\mathbf{x}, \mathbf{z}_j|\theta)}{Q(\mathbf{z}_j)}, \quad (18)$$

for $\mathbf{z}_1, \dots, \mathbf{z}_m \sim Q(\mathbf{z})$, where m is the number of samples used to evaluate the Monte Carlo integral.

The second trick is the *reparametrization trick*, which facilitates taking derivatives of the above expression with respect to the variational parameters γ governing $Q(\mathbf{z})$. Instead of sampling $\mathbf{z}_j \sim Q(\mathbf{z})$, the idea is to map samples from a standard normal distribution into samples from the desired variational distribution $Q(\mathbf{z}) = \mathcal{N}(\mu_\gamma(\mathbf{x}), \Sigma_\gamma(\mathbf{x}))$ via a differentiable transformation, namely:

$$\mathbf{z}_j = \Sigma_\gamma^{\frac{1}{2}}(\mathbf{x}) \mathbf{n}_j + \mu_\gamma(\mathbf{x}) \quad (19)$$

for $\mathbf{n}_1, \dots, \mathbf{n}_m \sim N(0, I)$, where $\Sigma_\gamma^{\frac{1}{2}}(\mathbf{x})$ denotes the matrix square root of the covariance matrix $\Sigma_\gamma(\mathbf{x})$.

Combining these two tricks, and plugging in the VAE generative and recognition model components for $P(\mathbf{x}, \mathbf{z}|\theta)$ and $Q(\mathbf{z})$, a Monte Carlo evaluation of the per-datum ELBO can be written:

$$\mathcal{L}_{MC}(Q, \theta|\mathbf{x}) = \frac{1}{m} \sum_{j=1}^m \left(\frac{1}{2\sigma_\epsilon^2} \left\| \mathbf{x} - G_\theta \left(\Sigma_\gamma^{\frac{1}{2}}(\mathbf{x}) \mathbf{n}_j + \mu_\gamma(\mathbf{x}) \right) \right\|^2 + \log \frac{P(\Sigma_\gamma^{\frac{1}{2}}(\mathbf{x}) \mathbf{n}_j + \mu_\gamma(\mathbf{x}))}{Q(\Sigma_\gamma^{\frac{1}{2}}(\mathbf{x}) \mathbf{n}_j + \mu_\gamma(\mathbf{x}))} \right), \quad (20)$$

where the first term is the the mean squared error between the original image patch \mathbf{x} and its reconstruction after passing \mathbf{x} through the noisy encoder and decoder, and the second term is the log-ratio of the latent prior $P(\mathbf{z}) = \mathcal{N}(0, I)$ and the variational distribution $Q(\mathbf{z}) = \mathcal{N}(\mu_\gamma(\mathbf{x}), \Sigma_\gamma(\mathbf{x}))$, evaluated at the sample point $\mathbf{z}_j = \Sigma_\gamma^{\frac{1}{2}}(\mathbf{x}) \mathbf{n}_j + \mu_\gamma(\mathbf{x})$.

It is worth noting that the first term in (eq. 20) relies on a stochastic pass through the encoder, given by a noisy sample of the latent from $Q(\mathbf{z})$ (which is really an approximation

to $P(\mathbf{z}|\mathbf{x})$, the conditional distribution of the latent given \mathbf{x}). This sample is then passed deterministically through the decoder network G_θ . The generative model noise variance σ_ϵ^2 serves as an inverse weight that determines how much to penalize reconstruction error relative to the second term in the ELBO. This second term, in turn, can be seen as a Monte Carlo estimate for $-D_{KL}(Q(\mathbf{z})||P(\mathbf{z}))$, the negative KL divergence between the variational posterior and the prior over \mathbf{z} . Because both these distributions are Gaussian, the standard approach is to replace the Monte Carlo evaluation of this term with its true expectation, using the fact that the KL divergence between two Gaussians can be computed analytically (Kingma & Welling, 2014; Rezende et al., 2014).

In contrast to the iterative variational EM algorithm used for the classic sparse coding model, optimization of the VAE is carried out by simultaneous gradient ascent of the ELBO with respect to generative network parameters θ and variational parameters γ . During training, the per-datum ELBO in (eq. 20) is summed over a mini-batch of data for each stochastic gradient ascent step.

3 Sparse Coding VAEs

In this paper, we adapt the VAE framework to the sparse coding model, a sparse generative model motivated by theoretical coding principles. This involves three changes to the standard VAE: (1) We replace the deep neural network from the VAE generative model with a linear feed-forward network; (2) We change the under-complete latent variable representation to an overcomplete one, so that the dimension of the latent variables is larger than the dimension of the data; and (3) We replace the standard normal prior over the latent variables with heavy-tailed, sparsity-promoting priors (e.g., Laplace and Cauchy).

We leave the remainder of the VAE framework intact, including the conditionally Gaussian variational distribution $Q(\mathbf{z})$, parametrized by a pair of neural networks that output the mean and covariance as a function of \mathbf{x} . The parameters of the SVAE are therefore given by $\{\Phi, \sigma_\epsilon^2, \gamma\}$ and a prior $P(\mathbf{z})$, where $P(\mathbf{z})$, Φ and σ_ϵ^2 specify the elements of sparse coding model (a sparse prior, generative weight matrix, and Gaussian noise variance, respectively), and variational parameters γ are the weights of the recognition networks $\mu_\gamma(\cdot)$ and $\Sigma_\gamma(\cdot)$ governing the variational distribution $Q(\mathbf{z}) = \mathcal{N}(\mu(\mathbf{x}), \Sigma_\gamma(\mathbf{x}))$. Fig. 2 shows a schematic comparing the two models.

4 Methods

4.1 Data Preprocessing

We fit the SVAE to 12×12 pixel image patches sampled from the BSDS300 dataset (Martin et al., 2001). Before fitting, we preprocessed the images and split the data into train, test, and validation sets. In the original sparse coding model, the images were whitened in the frequency domain, followed by a low-pass filtering stage. In (Olshausen & Field, 1996a) the whitening step was taken to expedite learning, accentuating high frequency features that would be far less prominent for natural image data, which is dominated by the low-frequency features. This is due to the fact that the Fourier (frequency) components are

approximately the principle components of natural images, however the overall variance of each component scales with the inverse frequency squared (the well known $1/f$ spectral properties of natural images), producing large differences in variance between the high- and low-frequency features. This poses a problem for gradient-based approaches to fitting the sparse coding model since low variance directions dominate the gradients. The low-pass filtering stage then served to reduce noise and artifacts from the rectangular sampling grid.

We perform a slight variation of these original preprocessing steps, but with the same overall effect. We whiten the data by performing PCA on the natural images and normalizing each component by its associated eigenvalue. Our low-pass filtering is achieved by retaining only the $100\pi/4\%$ most significant components, which correspond to the $100\pi/4\%$ lowest frequency modes, which also roughly corresponds to a circumscribed circle in the square Fourier space of the image, removing the noisy, high frequency corners of Fourier space.

4.2 SVAE Parameters

We implemented an SVAE with $\dim(\mathbf{z}) = 169$ latent dimensions, resulting in a latent code that is 1.5 times overcomplete, and set the Gaussian output noise variance to $\sigma_\epsilon^2 = \exp(-2)$. We implemented the SVAE with three choices of prior which included Laplace (Eq. 21), Cauchy (Eq. 22), and Gaussian (Eq. 23).

$$P_L(\mathbf{z}) = \prod_{l=1}^N \frac{1}{2} \exp(-|z_l|) \quad (21)$$

$$P_C(\mathbf{z}) = \prod_{l=1}^N \frac{1}{\pi} \left(1 + z_l^2\right)^{-1} \quad (22)$$

$$P_G(\mathbf{z}) = \prod_{l=1}^N \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}z_l^2\right) \quad (23)$$

We parametrized the recognition models $\mu_\gamma(\cdot)$ and $\Sigma_\gamma(\cdot)$ using feed-forward deep neural networks (Figure 2). The two networks took as input a vector of preprocessed data and had a shared initial hidden layer of 128 rectified linear units. Each network had two additional hidden layers of 256 and 512 rectified linear units respectively, which were not shared. These networks, parameterizing $\mu_\gamma(\cdot)$ and $\Sigma_\gamma(\cdot)$, each output a $\dim(\mathbf{z}) = 169$ dimensional vector which encode the mean and main diagonal of the covariance of the posterior, respectively. We set the off-diagonal elements of the posterior covariance matrix to zero. The final hidden layer in each network is densely connected to the output layer, with no nonlinearity for the mean output layer and a sigmoid nonlinearity for the variance output layer. In principle the variance values should be encoded by a non-saturating positive-definite nonlinearity; however we found that this led to instability during the fitting process and the sigmoid nonlinearity resulted in more stable behavior. Intuitively, given that our priors have scales of one, the posteriors will generally have variances less than 1, and can be expressed sufficiently well with the sigmoid nonlinearity.

4.3 Optimization

We optimized the SVAE using the `lasagne` (Dieleman et al., 2015) python package with `theano` (Theano Development Team, 2016) back-end. Gradient descent was performed for 10^6 steps with the Adam optimizer (Kingma & Ba, 2014) with default parameters, a batch size of 32, and a learning rate of $\eta = .001$. The networks always converged well within this number of gradient descent steps. We took the number of Monte-Carlo integration samples to be $m = 1$ and tested higher values, up to $m = 5$, but found that this parameter did not influence the results. We used the same learning hyperparameters for all three priors.

For comparison we fit our own version of the sparse coding model using the methods of Olshausen and Field (Olshausen & Field, 1996a). We utilized the same learning hyperparameters as in the original work, except during the normalization of Φ (see Section 2.3). In the original work the norm of Φ was driven to match the variance of the inferred \mathbf{z} to the data variance, however the data variance is arbitrary. Instead we chose to normalize Φ so that the inter-quartile range of the inferred \mathbf{z} matched the inter-quartile range of the prior $P(\mathbf{Z})$. This puts a constraint on the inferred \mathbf{z} to have a distribution which more closely matches that of the true generative model. Explicitly, at the end of each mini-batch we update each dictionary element Φ_i according to

$$\Phi_i \rightarrow \Phi_i \left(\frac{\text{IQR}(\gamma_i)}{\text{IQR}(P(\mathbf{Z}))} \right)^\alpha \quad (24)$$

where $\text{IQR}(\cdot)$ denotes the inter-quartile range, γ are the inferred latent variables from the most recent minibatch, and $\alpha = .05$.

5 Results

5.1 Quality of Fit

We first assess our model by calculating goodness-of-fit measures. We compare here the various choices in prior distributions using both the SVAE optimized with VAE-based inference and the original sparse coding model fit with the method of Olshausen & Field (1996a). To perform this assessment, we computed the log-likelihood of test data under the fit parameters. To calculate log-likelihoods (which can be a difficult quantity to obtain for large models), we used annealed importance sampling (AIS) (R. M. Neal, 2001; Wu et al., 2016) with 200 annealing steps and an initial distribution of either a standard normal distribution (for the Olshausen and Field implementations) or the variational posterior (for VAE implementations). Note that AIS only gives a lower bound on the log-likelihood (Wu et al., 2016). Even though an exact expression is available for the log likelihood in the Gaussian case we use AIS for the Gaussian prior as well so that the resulting values are directly comparable.

We first compare our chosen priors within the sparse coding VAE framework. Figure 3A depicts how the log-likelihood is almost universally monotonically increasing over training for all priors. Small dips are due to the randomness of the stochastic gradient descent (SGD) learning rule. We observe that, of the priors tested, sparser distributions result in higher log-

likelihoods, with the Cauchy prior providing the best fit. To explore the utility of the sparse coding VAE over the approximate EM method, we similarly calculated the log-likelihoods for the method in (Olshausen & Field, 1996a), again for each of the three prior distributions. We observe that the log-likelihood goodness-of-fit for the VAE is higher than the equivalent fits obtained with the method of (Olshausen & Field, 1996a) (Table 1). This is due to the fact that the sparse coding VAE uses a more robust approximation of the log-likelihood, and the variational posterior of VAEs are more informative than simply using the posterior mode. Finally, Figure 3B depicts how the inferred latent values of the sparse coding VAEs provide a better approximation of the prior than those obtained with the Dirac-delta approximation.

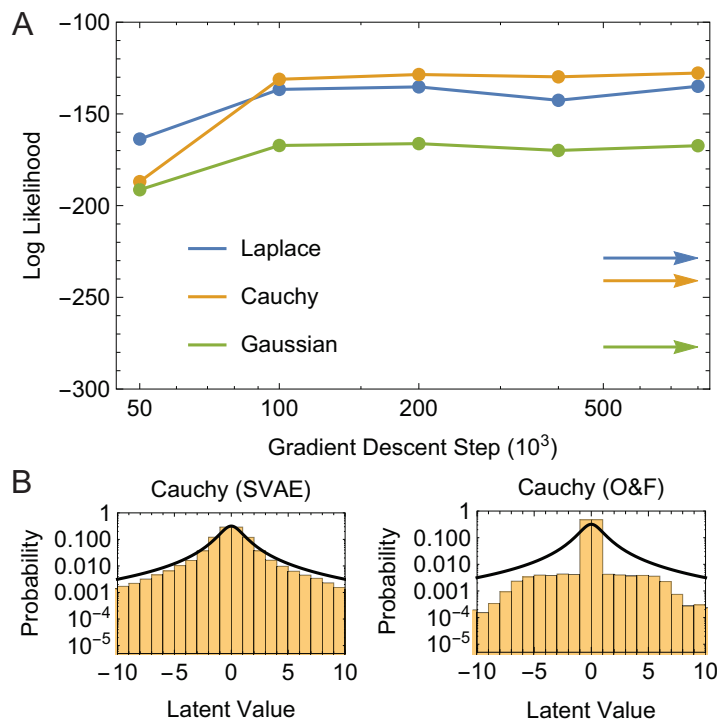


Figure 3: Quantitative analysis of quality of fit. **Panel A:** Log-Likelihood of VAE fit throughout training for Laplace, Cauchy and Gaussian priors. The initial log-likelihoods for the sparse VAE models are translated to a horizontal position of 50 for visualization purposes. Arrows indicate the final log-likelihood for the corresponding fit using the method of (Olshausen & Field, 1996a). **Panel B:** Table of final Log-Likelihoods for VAE trained models (our result) and models fit using the method of (Olshausen & Field, 1996a). Values are rounded to the nearest integer and reported in nats. **Panel C:** Histogram of inferred latent variables sampled from the variational posterior under our fit (left) and the fit method of (Olshausen & Field, 1996a) (right) along with the true prior (black line). Note that the variational posterior in the method of (Olshausen & Field, 1996a) is a Dirac-delta distribution at the MAP estimate of z .

5.2 Feed Forward Inference Model

As a consequence of training a VAE, we obtain a neural network which performs approximate Bayesian inference. Previous mechanistic implementations of sparse coding use the

Implementation	Prior		
	Laplace	Cauchy	Gaussian
VAE	-135	-128	-167
Traditional Method	-229	-241	-277

Table 1: Log-likelihood calculations for the sparse coding VAE (top) versus the traditional Dirac-delta approximation. Values were calculated using the AIS and show improvement across all prior choices.

MAP estimates under the true posterior to model trial-averaged neural responses (Rozell et al., 2008; Boerlin & Denève, 2011; Gregor & LeCun, 2010; Martins et al., 2011). In our case, the recognition model performs approximate inference using a more expressive approximation, suggesting that it may serve as an effective model of visual cortical responses. To study the response properties of the feed-forward inference network we simulated neural activity as the mean of the recognition distribution $Q(\mathbf{z}) \approx P(\mathbf{z}|\mathbf{x})$ with stimulus \mathbf{x} taken to be sinusoidal gratings of various sizes, contrasts, angles, frequencies, and phases. For each set of grating parameters, we measured responses to both a cosine and sine phase grating. To enforce non-negative responses and approximate phase invariance, the responses shown in Figure 4 are the root sum-of-squares of responses to both phases.

Fig. 4 shows the performance of the recognition network in response to sinusoidal gratings of various sizes, contrasts, angles, and phases. We found that the responses of the recognition model exhibited frequency and orientation tuning (Fig. 4B,C), reproducing important characteristics of cortical neurons (Hubel & Wiesel, 1962) and reflecting the Gabor-like structure of the dictionary elements. Additionally, Fig. 4D demonstrates that the orientation tuning of these responses is invariant to grating contrast, which is observed in cortical neurons (Troyer et al., 1998). Fig. 4A shows that the recognition model exhibits the surround suppression, in which the response to an optimally-tuned grating first increases with grating size and then decreases as the scale of the grating increases beyond the classical receptive field (Sceniak et al., 1999). Lastly, Figure 4E shows the receptive fields of recognition model neurons, as measured by a reverse-correlation experiment (Ringach & Shapley, 2004), verifying that the linear receptive fields exhibit the same Gabor-like properties as the dictionary elements.

6 Discussion

We have introduced a variational inference framework for the sparse coding model based on the VAE. The resulting SVAE model offers a more principled and accurate method for fitting the sparse coding model, and comes equipped with a neural implementation of feed-forward inference under the model. We showed, first of all, that the classic fitting method of Olshausen & Field is equivalent to variational inference under a delta-function variational posterior. We then extended the VAE framework to incorporate the sparse coding model as generative model. In particular, we replaced the standard deep network of the VAE with an overcomplete latent variable governed by a sparse prior, and showed that variational inference using a conditionally Gaussian recognition distribution provided accurate, neurally plausible inference of latent variables from images. Additionally, the SVAE provided

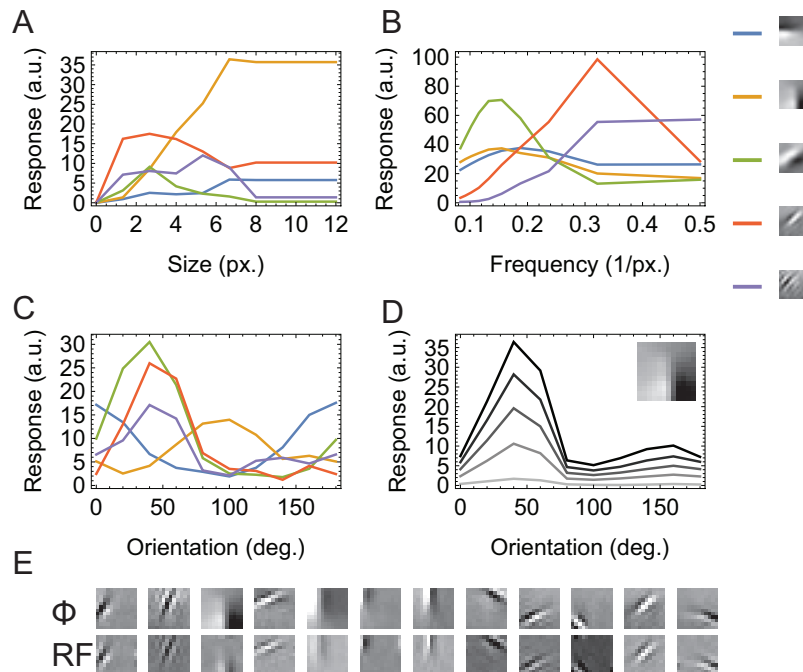


Figure 4: Visualization of feed-forward inference model behavior. Curves in Panels A-C are generated using full-contrast gratings. Dictionary elements associated with each neuron are visualized on the right. **Panel A:** size tuning curves of five sample neurons in response to a grating with optimal orientation and frequency. **Panel B:** Frequency response function of five model neurons in response to a full-field, full-contrast grating with optimal orientation. **Panel C:** Orientation preference curves of five model neurons in response to full-field, full-contrast gratings. **Panel D:** Example orientation preference curves of a single model neuron in response to a full-field grating at varying contrasts (denoted by line darkness) demonstrating contrast invariance. **Panel E:** Comparison of dictionary elements (top) with receptive fields (bottom) measured as in a reverse-correlation experiment for measuring receptive field (Ringach & Shapley, 2004). The stimulus used for reconstruction was Gaussian (low-pass) filtered white noise with Gaussian filter width $s = 0.6$ px. and receptive fields were averaged over 10000 random stimulus presentations.

improved fitting of the sparse coding model to natural images, as measured by the test log-likelihood. Moreover, we showed that the associated recognition model recapitulates important response properties of neurons in the early mammalian visual pathway.

Given this demonstration of VAEs for fitting tailor-made generative models, it is important to ask whether VAEs have additional applications in theoretical neuroscience. Specifically, many models are constrained by their ability to be fit. Our technique may allow more powerful, yet still highly structured, generative models to be practically applied by making learning tractable. The particular property that VAEs model the entire posterior means that connections can also now be drawn between generative models, (e.g., sparse coding) and models that depend explicitly on neural variability — a property often tied to the confidence levels in encoding. Current models of over-dispersion (e.g. (Goris et al., 2014; Charles et al., 2018))

shy from proposing mechanistic explanations of the explanatory statistical models. Another important line of future work then is to explore whether the posterior predictions made by SVAEs, or VAEs tuned to other neuroscience models, can account for observed spiking behaviors.

Relationship to previous work

The success of the sparse coding model as an unsupervised learning method for the statistics of the natural world has prompted an entire field of study into models for sparse representation learning and implementations of such models in artificial and biological neural systems. In terms of the basic mathematical model, many refinements and expansions have been proposed to better capture the statistics of natural images, including methods to more strictly induce sparsity (Garrigues & Olshausen, 2008; Girolami, 2001), hierarchical models to capture higher order statistics (Karklin & Lewicki, 2003, 2005, 2009; Garrigues & Olshausen, 2010), constrained dictionary learning for non-negative data (Charles et al., 2011), and applications to other modalities, such as depth (Töscic et al., 2011), motion (Cadieu & Olshausen, 2009) and auditory coding (Smith & Lewicki, 2006).

Given the success of such models in statistically describing visual responses, the mechanistic question as to how a neural substrate could implement sparse coding also became an important research area. Neural implementations of sparse coding have branched into two main directions: recurrent (Rozell et al., 2008; Boerlin & Denève, 2011) and feed-forward neural networks (Gregor & LeCun, 2010; Martins et al., 2011). The recurrent network models have been shown to provably solve the sparse-coding problem (Rozell et al., 2008; Shapero et al., 2014; Schwemmer et al., 2015). Furthermore, recurrent models can implement hierarchical extensions (Charles et al., 2012) as well as replicate key properties of visual cortical processing, such as non-classical receptive fields (Zhu & Rozell, 2013). The feed-forward models are typically based off of either mimicking the iterative processing of recursive algorithms, such as the iterative soft-thresholding algorithm (ISTA) (Daubechies et al., 2004) or by leveraging unsupervised techniques for learning deep neural networks, such as optimizing auto-encoders (Makhzani & Frey, 2013). The resulting methods, such as the learned ISTA (LISTA) (Gregor & LeCun, 2010; Borgerding et al., 2017), provide faster feed-forward inference than their RNN counterparts², at the cost of losing theoretical guarantees on the estimates.

Although the details of these implementations vary, all essentially retain the Dirac-delta posterior approximation and are thus constructed to calculate MAP estimates of the coefficients for use in a gradient-based feedback to update the dictionary. None of these methods reassess this basic assumption, and so are limited in the overall accuracy of the marginal-log-likelihood estimation of the model Φ , as well as their ability to generalize beyond inference-based networks to other theories of neural processing, such as probabilistic coding (Fiser et al., 2010; Orbán et al., 2016). In this work we have taken advantage of the refinement of VAEs in the machine learning literature to revisit this initial assumption from the influential early work and create just such a posterior-seeking neural network. Specifically, VAEs can provide a nontrivial approximation of the posterior via a fully Bayesian learning procedure in a feed-forward neural network model of inference under the sparse coding model.

²We note that this is true for digital processing only, and that analog recurrent systems can be faster (Shapero et al., 2014).

Obtaining the posterior distribution is especially important given that neural variability and spike-rate over-dispersion can be related to the uncertainty in the generative coefficients' posterior (e.g., via probabilistic coding (Fiser et al., 2010; Orbán et al., 2016)). Finding a neural implementation of sparse coding that also estimates the full posterior would be an important step towards bridging the efficient and probabilistic coding theories. Our work complements other recent efforts to connect the sparse coding model with tractable variational inference networks. These works have focused on either the non-linear sparsity model (Salimans, 2016) or the linear generative model of sparse coding (Aitchison et al., 2018). Other related work uses traditional VAEs and then performs sparse coding in the latent space (Sun et al., 2018). To date, however, no variational method has been designed to capture the three fundamental characteristics of sparse coding: overcomplete codes, sparse priors, and a linear-generative model.

Limitations and future directions

The state-of-the-art results of this work are primarily the result of orienting sparse coding in a variational framework where more expressive variational posterior distributions can be used for model fitting. Nevertheless this work represents only the first steps in this direction. One area for improvement is our selection of a Gaussian variational posterior with diagonal covariance matrix $\Sigma_\gamma(\cdot)$. This choice restricts the latent variables to be uncorrelated under the posterior distribution and can limit the variational inference method's performance (Mnih & Gregor, 2014; Turner & Sahani, 2011). For example, SVAE posterior employed here cannot directly account for the “explaining away” effect which occurs between the activations of overlapping dictionary elements. An important next step is thus to expand this work to the case of variational posteriors with nontrivial correlations.

In addition to improved learning, a more complex posterior would give the recognition model the potential to exhibit interesting phenomena associated with correlations between dictionary element activations. While some computational models aim to account for such correlations (e.g., (Cadieu & Olshausen, 2009; Karklin & Lewicki, 2009; Averbek et al., 2006)), the SVAE framework would allow for the systematic analysis of the many assumptions possible in the population coding layer within the sparse coding framework. These various assumptions can thus be validated against the population correlations observed in biological networks (e.g., (Ecker et al., 2011; Cohen & Kohn, 2011)). For example, if V1 responses are interpreted as arising via the sampling hypothesis (Fiser et al., 2010) then “explaining away” may account for correlations in neural variability observed in supra- and infra-granular layers of V1 (Hansen et al., 2012). Additionally, such correlations can be related to probabilistic population coding (e.g. (Fiser et al., 2010; Orbán et al., 2016)) where correlated variability represents correlated uncertainty in the neural code.

In this work we restricted the sparse coding model by choosing the magnitude of the output noise variance σ_ϵ *a priori*. This was done in order to make this work comparable to the original sparse coding implementations of Olshausen & Field (1996a). Nevertheless, this parameter can be fit in a data-driven way as well, providing additional performance beyond the current work. In future explorations this constraint may be relaxed.

One of the favorable properties of sparse coding using MAP inference with a Laplace prior is that truly sparse representations are produced in the sense that a finite fraction of latent variables are inferred to be exactly zero. As a consequence of the more robust variational

inference we perform here, the SVAE no longer has this property. Sparse representations could be regained within the SVAE framework if truly sparse priors, which have a finite fraction of their probability mass at exactly zero, were used such as “spike and slab” priors Garrigues & Olshausen (2008); Ziniel & Schniter (2013). Such sparse priors are not differentiable and thus cannot be directly incorporated into the framework presented here, however continuous approximations of sparse priors do exist and during this work we implemented an approximate spike and slab prior using a sum of Gaussians with a small and large variance respectively. We were unable in our setting, however, to replicate the superior performance of such priors seen elsewhere, and recommend additional explorations into incorporating hard-sparse priors such as the spike-and-slab or concrete Maddison et al. (2016) distributions.

We have restricted ourselves in this work to a relatively simple generative model of natural images. Hierarchical variants of the sparse coding model (Karklin & Lewicki, 2003, 2005, 2009; Garrigues & Olshausen, 2010) provide superior generative models of natural images. These more complex generative models can be implemented and fit using the same methods we present here by explicitly constructing the VAE generative model in their image.

7 Conclusion

In summary, we have cast the sparse coding model in the framework of variational inference, and demonstrated how modern tools such as the VAE can be used to develop neurally plausible algorithms for inference under generative models. We feel that this work strengthens the connection between machine learning methods for unsupervised learning of natural image statistics and efforts to understand neural computation in the brain, and hope it will inspire future studies along these lines.

Acknowledgments

The authors would like to acknowledge Ryan Pyle for involvement in early stages of this work. GB acknowledges the Marine Biological Laboratory in Woods Hole, NIMH funding for the Methods in Computational Neuroscience course (R25MH062204), and support from the Simons Foundation. JWP was supported by grants from the Simons Foundation (SCGB AWD1004351 and AWD543027), the NIH (R01EY017366), and the CAREER award (IIS-1150186).

References

- Aitchison, L., Hennequin, G., & Lengyel, M. (2018). Sampling-based probabilistic inference emerges from learning in neural circuits with a cost on reliability. *arXiv preprint arXiv:1807.08952*.
- Averbeck, B. B., Latham, P. E., & Pouget, A. (2006). Neural correlations, population coding and computation. *Nature reviews neuroscience*, 7(5), 358.

- Berkes, P., & Wiskott, L. (2005). Slow feature analysis yields a rich repertoire of complex cell properties. *Journal of Vision*, 5(6), 579–602. Retrieved from <http://dx.doi.org/10.1167/5.6.9>
- Bishop, C. M. (2005). *Neural networks for pattern recognition*. Oxford University Press.
- Blei, D. M., Kucukelbir, A., & McAuliffe, J. D. (2016). Variational inference: A review for statisticians. *arXiv preprint arXiv:1601.00670*.
- Boerlin, M., & Denève, S. (2011). Spike-based population coding and working memory. *PLoS computational biology*, 7(2), e1001080.
- Borgerding, M., Schniter, P., & Rangan, S. (2017). Amp-inspired deep networks for sparse linear inverse problems. *IEEE Transactions on Signal Processing*, 65(16), 4293–4308.
- Cadiou, C., & Olshausen, B. A. (2009). Learning transformational invariants from natural movies. In *Advances in neural information processing systems* (pp. 209–216).
- Charles, A. S., Garrigues, P., & Rozell, C. J. (2012). A common network architecture efficiently implements a variety of sparsity-based inference problems. *Neural computation*, 24(12), 3317–3339.
- Charles, A. S., Olshausen, B. A., & Rozell, C. J. (2011). Learning sparse codes for hyper-spectral imagery. *IEEE Journal of Selected Topics in Signal Processing*, 5(5), 963–978.
- Charles, A. S., Park, M., Weller, J. P., Horwitz, G. D., & Pillow, J. W. (2018). Dethroning the fano factor: A flexible, model-based approach to partitioning neural variability. *Neural computation*, 30(4), 1012–1045.
- Coen-Cagli, R., Dayan, P., & Schwartz, O. (2012). Cortical surround interactions and perceptual salience via natural scene statistics. *PLoS Computational Biology*, 8(3), e1002405.
- Cohen, M. R., & Kohn, A. (2011, 6). Measuring and interpreting neuronal correlations. *Nat. Neurosci.*, 14(7), 811–819.
- Daubechies, I., Defrise, M., & De Mol, C. (2004). An iterative thresholding algorithm for linear inverse problems with a sparsity constraint. *Communications on pure and applied mathematics*, 57(11), 1413–1457.
- Dayan, P., Sahani, M., & Deback, G. (2003). Adaptation and unsupervised learning. In *Advances in neural information processing systems 15*. MIT Press.
- Dieleman, S., Schläijter, J., Raffel, C., Olson, E., SÅynderby, S. K., Nouri, D., . . . Degraeve, J. (2015, August). *Lasagne: First release*. Retrieved from <http://dx.doi.org/10.5281/zenodo.27878> doi: 10.5281/zenodo.27878
- Doersch, C. (2016). Tutorial on variational autoencoders. *arXiv preprint arXiv:1606.05908*.
- Ecker, A. S., Berens, P., Tolias, A. S., & Bethge, M. (2011, 10). The effect of noise correlations in populations of diversely tuned neurons. *J. Neurosci.*, 31(40), 14272–14283.
- Fiser, J., Berkes, P., Orban, G., & Lengyel, M. (2010, March). Statistically optimal perception and learning: from behavior to neural representations. *Trends Cogn. Sci. (Regul. Ed.)*, 14(3), 119–130.

- Garrigues, P., & Olshausen, B. A. (2008). Learning horizontal connections in a sparse coding model of natural images. In *Advances in neural information processing systems* (pp. 505–512).
- Garrigues, P., & Olshausen, B. A. (2010). Group sparse coding with a laplacian scale mixture prior. In *Advances in neural information processing systems* (pp. 676–684).
- Girolami, M. (2001). A variational method for learning sparse and overcomplete representations. *Neural Computation*, *13*(11), 2517–2532.
- Goris, R. L. T., Movshon, J. A., & Simoncelli, E. P. (2014, April 28). Partitioning neuronal variability. *Nature Neuroscience*, *17*, 858. Retrieved from <http://dx.doi.org/10.1038/nn.3711> (Article)
- Gregor, K., & LeCun, Y. (2010). Learning fast approximations of sparse coding. In *Proceedings of the 27th international conference on international conference on machine learning* (pp. 399–406).
- Hansen, B. J., Chelaru, M. I., & Dragoi, V. (2012, November 08). Correlated variability in laminar cortical circuits. *Neuron*, *76*(3), 590-602. Retrieved from <http://dx.doi.org/10.1016/j.neuron.2012.08.029> doi: 10.1016/j.neuron.2012.08.029
- Hubel, D. H., & Wiesel, T. N. (1962). Receptive fields, binocular interaction and functional architecture in the cat's visual cortex. *J. Physiol. (Lond.)*, *160*, 106–154.
- Karklin, Y., & Lewicki, M. S. (2003, Aug). Learning higher-order structures in natural images. *Network*, *14*(3), 483–499.
- Karklin, Y., & Lewicki, M. S. (2005). A hierarchical bayesian model for learning nonlinear statistical regularities in nonstationary natural signals. *Neural computation*, *17*(2), 397–423.
- Karklin, Y., & Lewicki, M. S. (2009, Jan). Emergence of complex cell properties by learning to generalize in natural scenes. *Nature*, *457*(7225), 83–86. Retrieved from <http://dx.doi.org/10.1038/nature07481>
- Kingma, D. P., & Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Kingma, D. P., & Welling, M. (2014). Auto-encoding variational bayes. *arXiv:1312.6114*. Retrieved from <http://arxiv.org/abs/1312.6114>
- Knill, D., & Richards, W. (1996). *Perception as Bayesian inference*. Cambridge University Press.
- Knill, D. C., & Pouget, A. (2004). The bayesian brain: the role of uncertainty in neural coding and computation. *TRENDS in Neurosciences*, *27*(12), 712–719.
- Lee, H., Battle, A., Raina, R., & Ng, A. Y. (2007). Efficient sparse coding algorithms. In *Advances in neural information processing systems* (pp. 801–808).
- Lewicki, M., & Olshausen, B. (1999). Probabilistic framework for the adaptation and comparison of image codes. *JOSA A*, *16*(7), 1587–1601.

- Maddison, C. J., Mnih, A., & Teh, Y. W. (2016). The concrete distribution: A continuous relaxation of discrete random variables. *CoRR*, *abs/1611.00712*. Retrieved from <http://arxiv.org/abs/1611.00712>
- Makhzani, A., & Frey, B. (2013). K-sparse autoencoders. *arXiv preprint arXiv:1312.5663*.
- Martin, D., Fowlkes, C., Tal, D., & Malik, J. (2001). A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics. In *Proc. 8th int'l conf. computer vision* (Vol. 2, pp. 416–423).
- Martins, A. F., Smith, N. A., Aguiar, P. M., & Figueiredo, M. A. (2011). Structured sparsity in structured prediction. In *Proceedings of the conference on empirical methods in natural language processing* (pp. 1500–1511).
- Mnih, A., & Gregor, K. (2014). Neural variational inference and learning in belief networks. In *ICML*.
- Moreno-Bote, R., Knill, D. C., & Pouget, A. (2011, Jul). Bayesian sampling in visual perception. *Proc Natl Acad Sci U S A*, *108*(30), 12491–12496. Retrieved from <http://dx.doi.org/10.1073/pnas.1101430108>
- Neal, R., & Hinton, G. E. (1998). A view of the em algorithm that justifies incremental, sparse, and other variants. In *Learning in graphical models* (pp. 355–368). Kluwer Academic Publishers.
- Neal, R. M. (2001, April). Annealed importance sampling. *Statistics and Computing*, *11*(2), 125–139. Retrieved from <https://doi.org/10.1023/A:1008923215028> doi: 10.1023/A:1008923215028
- Olshausen, B. A., & Field, D. J. (1996a, June 13). Emergence of simple-cell receptive field properties by learning a sparse code for natural images. *Nature*, *381*, 607. Retrieved from <http://dx.doi.org/10.1038/381607a0>
- Olshausen, B. A., & Field, D. J. (1996b, May). Natural image statistics and efficient coding. *Network*, *7*(2), 333–339. Retrieved from <http://dx.doi.org/10.1088/0954-898X/7/2/014>
- Olshausen, B. A., & Field, D. J. (1997). Sparse coding with an overcomplete basis set: A strategy employed by v1? *Vision research*, *37*(23), 3311–3325.
- Orbán, G., Berkes, P., Fiser, J., & Lengyel, M. (2016). Neural variability and sampling-based probabilistic representations in the visual cortex. *Neuron*, *92*(2), 530–543.
- Rezende, D. J., Mohamed, S., & Wierstra, D. (2014). Stochastic backpropagation and approximate inference in deep generative models. In *Proceedings of the 31st international conference on machine learning (icml-14)* (pp. 1278–1286).
- Ringach, D., & Shapley, R. (2004, 03). Reverse correlation in neurophysiology. *Cognitive Science*, *28*, 147-166.
- Rozell, C. J., Johnson, D. H., Baraniuk, R. G., & Olshausen, B. A. (2008). Sparse coding via thresholding and local competition in neural circuits. *Neural computation*, *20*(10), 2526–2563.

- Salimans, T. (2016). A structured variational auto-encoder for learning deep hierarchies of sparse features. *arXiv preprint arXiv:1602.08734*.
- Sceniak, M. P., Ringach, D. L., Hawken, M. J., & Shapley, R. (1999). Contrast's effect on spatial summation by macaque v1 neurons. *Nature neuroscience*, 2(8), 733.
- Schwemmer, M. A., Fairhall, A. L., Denève, S., & Shea-Brown, E. T. (2015). Constructing precisely computing networks with biophysical spiking neurons. *Journal of Neuroscience*, 35(28), 10112–10134.
- Shapero, S., Zhu, M., Hasler, J., & Rozell, C. (2014). Optimal sparse approximation with integrate and fire neurons. *International journal of neural systems*, 24(05), 1440001.
- Smith, E. C., & Lewicki, M. S. (2006). Efficient auditory coding. *Nature*, 439(7079), 978.
- Sun, J., Wang, X., Xiong, N., & Shao, J. (2018). Learning sparse representation with variational auto-encoder for anomaly detection. *IEEE Access*, 6, 33353–33361.
- Theano Development Team. (2016, May). Theano: A Python framework for fast computation of mathematical expressions. *arXiv e-prints*, [abs/1605.02688](https://arxiv.org/abs/1605.02688). Retrieved from <http://arxiv.org/abs/1605.02688>
- Troyer, T. W., Krukowski, A. E., Priebe, N. J., & Miller, K. D. (1998). Contrast-invariant orientation tuning in cat visual cortex: Thalamocortical input tuning and correlation-based intracortical connectivity. *Journal of Neuroscience*, 18(15), 5908–5927. Retrieved from <http://www.jneurosci.org/content/18/15/5908> doi: 10.1523/JNEUROSCI.18-15-05908.1998
- Turner, R. E., & Sahani, M. (2011). Two problems with variational expectation maximisation for time-series models. In D. Barber, T. Cemgil, & S. Chiappa (Eds.), *Bayesian time series models* (pp. 109–130). Cambridge University Press.
- Tóšić, I., Olshausen, B. A., & Culpepper, B. J. (2011). Learning sparse representations of depth. *IEEE journal of selected topics in signal processing*, 5(5), 941–952.
- Weiss, Y., Simoncelli, E. P., & Adelson, E. (2002). Motion illusions as optimal percepts. *Nature Neuroscience*, 5, 598–604.
- Wu, Y., Burda, Y., Salakhutdinov, R., & Grosse, R. B. (2016). On the quantitative analysis of decoder-based generative models. *CoRR*, [abs/1611.04273](https://arxiv.org/abs/1611.04273).
- Zhu, M., & Rozell, C. J. (2013). Visual nonclassical receptive field effects emerge from sparse coding in a dynamical system. *PLoS computational biology*, 9(8), e1003191.
- Ziniel, J., & Schniter, P. (2013). Dynamic compressive sensing of time-varying signals via approximate message passing. *IEEE transactions on signal processing*, 61(21), 5270–5284.