

Autoencoders, GANs, and Diffusion

Theories of Deep Learning: C6.5, Lecture / Video 13 Prof. Jared Tanner Mathematical Institute University of Oxford



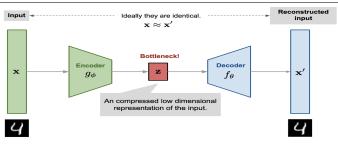


(Variational) Auto Encoders; VAEs

Autoencoder (AE) Illustration

Restricting the number of data parameters





The parameters, (θ, ϕ) , of the autoencoder are then learned:

$$\mathcal{L}(\theta,\phi) = m^{-1} \sum_{\mu=1}^{m} I(x_{\mu}, f_{\theta}(g_{\phi}(x_{\mu})))$$

https://lilianweng.github.io/lil-log/2018/08/12/from-autoencoder-to-beta-vae.html

Principal component analysis (PCA) as an AE

Learning with one layer hard thresholding



The parameters, (θ, ϕ) , of the autoencoder are then learned:

$$\mathcal{L}(\theta,\phi) = n^{-1} \sum_{\mu=1}^{n} I(x_{\mu}, f_{\theta}(g_{\phi}(x_{\mu})))$$

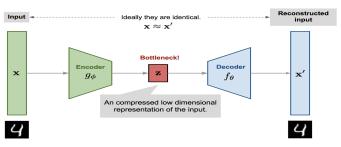
Consider a simple model where the encoder and decoder are linear, that is $g_{\phi}(x) = \Phi x$ where $\Phi \in \mathbb{R}^{r \times p}$ with r < p, and the linear decoder $f_{\theta}(z) = \Theta z$ with $\Theta \in \mathbb{R}^{p \times r}$.

Moreover, consider an entrywise ℓ_2^2 error for $I(x_\mu, f_\theta(g_\phi(x_\mu)))$, then

$$\mathcal{L}(\theta,\phi) = n^{-1} \|X - \Theta \Phi X\|_F^2$$

where $\Theta\Phi$ is a learned rank r matrix, whose optimal solution is the projector of X to its leading r singular space.





The autoencoder framework allows $g_{\phi}(\cdot)$ and $f_{\theta}(\cdot)$ to be more general than linear, and in particular to benefit from the expressively of depth and introduce variation.

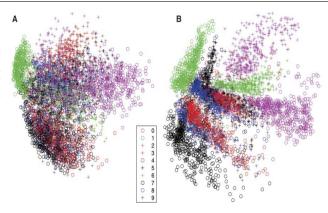
https://lilianweng.github.io/lil-log/2018/08/12/from-autoencoder-to-beta-vae.html

PCA vs 3 layer Autoencoder: MNIST (Hinton et al. 06')

Improved separation of data classes



Fig. 3. (A) The twodimensional codes for 500 digits of each class produced by taking the first two principal components of all 60,000 training images. (B) The two-dimensional codes found by a 784-1000-500-250-2 autoenoder. For an alternative visualization, see (B).



http://science.sciencemag.org/content/313/5786/504

k-sparse autoencoders (Makhzani et al. 13')

Low dimensionality through sparsity



k-Sparse Autoencoders:

Training:

- 1) Perform the feedforward phase and compute $z = W^{T}x + b$
- 2) Find the k largest activations of \boldsymbol{z} and set the rest to zero.
- $z_{(\Gamma)^c} = 0$ where $\Gamma = \text{supp}_k(z)$
- 3) Compute the output and the error using the sparsified z.

$$\hat{\boldsymbol{x}} = W\boldsymbol{z} + \boldsymbol{b'}$$

$$E = \|\boldsymbol{x} - \hat{\boldsymbol{x}}\|_2^2$$

3) Backpropagate the error through the k largest activations defined by Γ and iterate.

Sparse Encoding:

Compute the features $\boldsymbol{h} = W^{\mathsf{T}} \boldsymbol{x} + \boldsymbol{b}$. Find its αk largest activations and set the rest to zero.

$$h_{(\Gamma)^c} = 0$$
 where $\Gamma = \operatorname{supp}_{\alpha k}(h)$

This framework includes nonlinearity and can be rigorously analysed using techniques from sparse approximation, but it lacks depth. https://arxiv.org/pdf/1312.5663.pdf

k-sparse autoencoders (Makhzani et al. 13')

Learned elements: MNIST



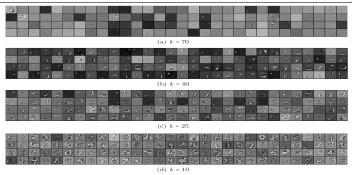


Figure 1. Filters of the k-sparse autoencoder for different sparsity levels k, learnt from MNIST with 1000 hidden units.

Elements learned depend on number of components, sparsity, allowed; k small are class elements, k large are basis elements. https://arxiv.org/pdf/1312.5663.pdf

k-sparse autoencoders (Makhzani et al. 13')

Performance vs other autoencoders



	Error Rate
Raw Pixels	7.20%
RBM	1.81%
Dropout Autoencoder (50% hidden)	1.80%
Denoising Autoencoder	1.95%
(20% input dropout)	
Dropout + Denoising Autoencoder	1.60%
(20% input and 50% hidden)	
k-Sparse Autoencoder, $k = 40$	1.54%
k-Sparse Autoencoder, $k = 25$	1.35%
k-Sparse Autoencoder, $k = 10$	2.10%

Table 1. Performance of	unsupervised	learning	${\it methods}$
(without fine-tuning) with	$1000\ \mathrm{hidden}$ u	nits on M	NIST.

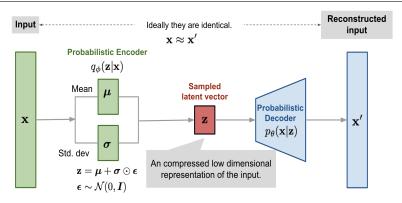
	Error
Without Pre-Training	1.60%
RBM + F.T.	1.24%
Shallow Dropout AE + F.T.	1.05%
(%50 hidden)	
Denoising $AE + F.T.$	1.20%
(%20 input dropout)	
Deep Dropout AE + F.T.	0.85%
(Layer-wise pre-training, %50 hidden)	
k-Sparse AE + F.T.	1.08%
(k=25)	
Deep k -Sparse AE + F.T.	0.97%
(Layer-wise pre-training)	

Table 3. Performance of supervised learning methods on MNIST. Pre-training was performed using the corresponding unsupervised learning algorithm with 1000 hidden units. and then the model was fine-tuned

Variational Autoencoders (VAE) (Kingma et al. 13')

Introduction of noise as a generative model





https://lilianweng.github.io/lil-log/2018/08/12/from-autoencoder-to-beta-vae.html

Gaussian MPL as encoder-decoder



Examples of encoder or decoder structure is a multivariate Gaussian with a diagonal covariance structure; e.g. for the decoder:

$$h = \tanh(W_3z + b_3)$$
with $\mu = W_4h + b_4$

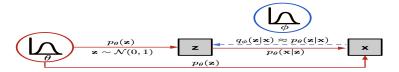
$$\log(\sigma^2) = W_5h + b_5$$

$$\log p(x|z) = \log \mathcal{N}(x; \mu, \sigma^2 I)$$

Variational Autoencoders (VAE) (Kingma et al. 13')

Distribution for generative model





 $p_{\theta}(x|z)$ acts as the generators, analogous to the decoder $f_{\theta}(x|z)$, and is called a probabilistic decoder

 $q_{\phi}(z|x)$ acts as the encoder, analogous to $g_{\phi}(z|x)$, and is used to approximate $p_{\theta}(z|x)$.

The parameters ϕ , θ for a model are then learned so minimize a distance, or divergence, between $q_{\phi}(z|x)$ and $p_{\theta}(z|x)$; Kingma proposed minimising the Kullback-Leibler divergence.

https://lilianweng.github.io/lil-log/2018/08/12/from-autoencoder-to-beta-vae.html

12



Loss function for VAE

The formulae for $\mathcal{L}_{ELBO}(\phi, \theta; X)$, the evidence lower bound (ELBO), follows from minimising a lower bound of $\sum_{\mu=1}^{n} \log p_{\theta}(x_{\mu})$:

$$\begin{split} \log p_{\theta}(x) &= \log \left(\int p_{\theta}(x|z) p_{\theta}(z) dz \right) \\ &= \log \left(\int p_{\theta}(x|z) \frac{p_{\theta}(z)}{q_{\phi}(z)} q_{\phi}(z) dz \right) \\ &\geq \int \log \left(p_{\theta}(x|z) \frac{p_{\theta}(z)}{q_{\phi}(z)} \right) q_{\phi}(z) dz \\ &= \mathbb{E}_{q_{\phi}(z|x)} \log(p_{\theta}(x|z)) - D_{KL}(q_{\phi}(z|x)|p_{\theta}(z|x)) \\ &=: \mathcal{L}_{ELBO}(\phi, \theta; x) \end{split}$$

Variational Autoencoder: manifold (Kingma et al. 13')

Learned two dimensional space: faces and MNIST



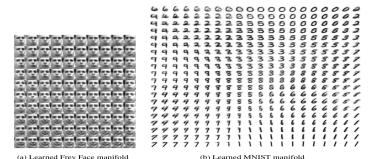


Figure 4: Visualisations of learned data manifold for generative models with two-dimensional latent space, learned with AEVB. Since the prior of the latent space is Gaussian, linearly spaced coordinates on the unit square were transformed through the inverse CDF of the Gaussian to produce values of the latent variables z. For each of these values z, we plotted the corresponding generative $p_0(\mathbf{x}|\mathbf{z})$ with the learned parameters θ .

Variational Autoencoder: manifold (Kingma et al. 13')

Dependence on manifold dimension: MNIST



```
$165764674 1831385738 8208933408
             8591682162 8387798538 7519
             6103288638 3599439513 8962882924
               68912041 1418833492 1986317961
             5191018350 1136470283 5449199910
             6861491788 5970583848
             1343913470 6943618552 15821
             4582970189 8490800066 7939299396
0401232088 6194272345 7476303601 4524390184
             2 6 4 5 6 0 9 9 9 8
                          2120131850
                                        8872516231
9 9 5 4 9 3 4 8 8 1
(a) 2-D latent space
              (b) 5-D latent space
                           (c) 10-D latent space
                                        (d) 20-D latent space
```

Figure 5: Random samples from learned generative models of MNIST for different dimensionalities of latent space.

Inference Variational Autoencoders (Zhao et al. 17')



 $p_{\theta}(x|z)$ acts as the generators, analogous to the decoder $f_{\theta}(x|z)$, and is called a probabilistic decoder; $q_{\phi}(z|x)$ acts as the encoder, analogous to $g_{\phi}(z|x)$, and is used to approximate $p_{\theta}(z|x)$. The parameters ϕ , θ for a model are then learned to minimize a distance, or divergence, between $q_{\phi}(z|x)$ and $p_{\theta}(z|x)$; Kingma proposed minimising the Kullback-Leibler divergence, giving the evidence lower bound (ELBO)

$$\mathcal{L}_{\textit{ELBO}} := \mathbb{E}_{q_{\phi(z|x)}} \log(p_{ heta}(x|z)) - \beta D_{\textit{KL}}(q_{\phi}(x,z)||p_{ heta}(x,z))$$

VAEs originally use $\beta=1$, with larger $\beta>1$ called β -VAEs. Zhao et al. propose including a mutual information term to avoid mode separation and collapse.

Inference Variational Autoencoders (Zhao et al. 17')

Impact of VAE objective ($\lambda = \beta$, in β -VAE



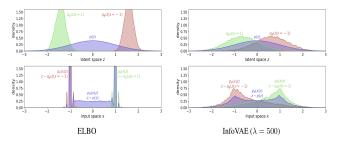


Figure 1: Verification of Proposition $\overline{\mathbf{I}}$ where the dataset only contains two examples $\{-1,1\}$. **Top:** density of the distributions $q_{\phi}(z|x)$ when x=1 (red) and x=-1 (green) compared with the true prior p(z) (purple). **Bottom:** The "reconstruction" $p_{\theta}(x|z)$ when z is sampled from $q_{\phi}(z|x=1)$ (green) and $q_{\phi}(z|x=-1)$ (red). Also plotted is $p_{\theta}(x|z)$ when z is sampled from the true prior p(z) (purple). When the dataset consists of only two data points, ELBO (**left**) will push the density in latent space Z away from 0, while InfoVAE (**right**) does not suffer from this problem.

Inference Variational Autoencoders (Zhao et al. 17')

Impact of VAE objective



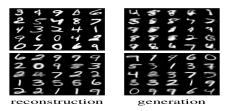


Figure 5: Samples generated by ELBO vs. MMD InfoVAE $(\lambda = 1000)$ after training on 500 samples (plotting mean of $p_{\theta}(x|z)$). **Top:** Samples generated by ELBO. Even though ELBO generates very sharp reconstruction for samples on the training set, model samples $p(z)p_{\theta}(x|z)$ is very poor, and differ significantly from the reconstruction samples, indicating over-fitting, and mismatch between $q_{\phi}(z)$ and p(z). **Bottom:** Samples generated by InfoVAE. The reconstructed samples and model samples look similar in quality and appearance, suggesting better generalization in the latent space.

β-VAEs disentangling features pt. 1 (Higgins et al. 17')





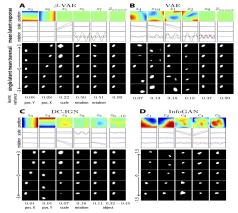


Figure 7: A: Representations learnt by a β -VAE ($\beta = 4$). Each column represents a latent z_i , ordered according to the learnt Gaussian variance (last row). Row 1 (position) shows the mean activation (red represents high values) of each latent z, as a function of all 32x32 locations averaged across objects, rotations and scales. Row 2 and 3 show the mean activation of each unit z_i as a function of scale (respectively rotation), averaged across rotations and positions (respectively scales and positions). Square is red, oval is green and heart is blue. Rows 4-8 (second group) show reconstructions resulting from the traversal of each latent z, over three standard deviations around the unit Gaussian prior mean while keeping the remaining 9/10 latent units fixed to the values obtained by running inference on an image from the dataset. B: Similar analysis for VAE $(\beta = 1)$. C: Similar analysis for DC-IGN, clamping a single latent each for scale, positions, orientation and 5 for shape. D: Similar analysis for InfoGAN, using 5 continuous latents regularized using the mutual information cost, and 5 additional unconstrained noise latents (not shown).

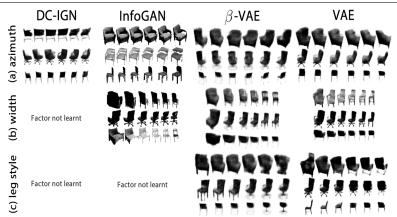
https://arxiv.org/pdf/1706.02262.pdf

19

β -VAEs disentangling features pt. 2 (Higgins et al. 17')

Explainable latent space: chairs





β-VAEs disentangling features pt. 3 (Higgins et al. 17')



Architectures for encoder-decoders

Dataset	Optimiser		Architecture	
2D shapes	Adagrad	Input	4096 (flattened 64x64x1).	
(VAE)	1e-2	Encoder	FC 1200, 1200. ReLU activation.	
		Latents	10	
		Decoder	FC 1200, 1200, 1200, 4096. Tanh activation. Bernoulli.	
2D shapes	rmsprop	Input	64x64x1.	
(DC-IGN)	(as in Kulkarni et al., 2015)	Encoder	Conv 96x3x3, 48x3x3, 48x3x3 (padding 1).	
			ReLU activation and Max pooling 2x2.	
		Latents	10	
		Decoder	Unpooling, Conv 48x3x3, 96x3x3, 1x3x3.	
			ReLU activation, Sigmoid.	
2D shapes	Adam	Generator	FC 256, 256, Decony 128x4x4, 64x4x4 (stride 2). Tanh.	
(InfoGAN)	1e-3 (gen) 2e-4 (dis)	Discriminator	Conv and FC reverse of generator. Leaky ReLU activation. FC 1. Sigmoid activation.	
20-4 (uis)	(,	Recognition	Conv and FC shared with discriminator. FC 128, 5. Gaussian	
	Latents	10: $z_{15} \sim Unif(-1,1), c_{15} \sim Unif(-1,1)$		
Chairs	Adam	Input	64x64x1.	
(VAE) le-4	Encoder	Conv 32x4x4 (stride 2), 32x4x4 (stride 2), 64x4x4 (stride 2), 64x4x4 (stride 2), FC 256. ReLU activation.		
		Latents	32	
		Decoder	Decony reverse of encoder, ReLU activation, Bernoulli,	

https://arxiv.org/pdf/1706.02262.pdf

4 日 N 4 個 N 4 国 N 4 国 N



- Principal component analysis (PCA) reveals the low dimensional latent space within a data matrix by projecting to the space of low-rank matrices.
- ▶ Autoencoders (AE) extend this notion allowing more general maps to and from a low dimensional parameter space.
- Variational AEs (VAEs) give a probabilistic notion that gives a natural generative model.
- Inference VAEs and β-VAEs are further extensions to improve VAEs and for interpretability respectively.
- ► (V)AEs are a flexible structure allowing general maps; an area open for great further analysis.



Generative Adversarial Networks (GANs)

Generative deep nets (Goodfellow et al. 14')

Generative model from 100 latent variables



Example of a deep convolutional generator:

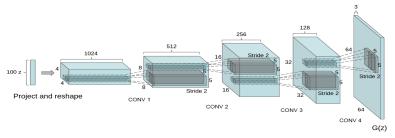


Figure 1: DCGAN generator used for LSUN scene modeling. A 100 dimensional uniform distribution Z is projected to a small spatial extent convolutional representation with many feature maps.

https://arxiv.org/pdf/1511.06434.pdf https://arxiv.org/pdf/1406.2661.pdf

Generative deep nets (Goodfellow et al. 14')

Generative model from 100 latent variables



Train the two network parameters using the objective

$$\min_{G} \max_{D} n^{-1} \sum_{\mu=1}^{n} \log(D(x_{\mu}, y_{\mu})) + p^{-1} \sum_{p} \log(1 - D(G(z_{p}), y_{p}))$$

Algorithm 1 Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator, k, is a hyperparameter. We used k=1, the least expensive option, in our experiments.

for number of training iterations do

- for k steps do
 - Sample minibatch of m noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_q(z)$.
 - Sample minibatch of m examples $\{x^{(1)}, \dots, x^{(m)}\}$ from data generating distribution
 - Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^{m} \left[\log D\left(\boldsymbol{x}^{(i)}\right) + \log\left(1 - D\left(G\left(\boldsymbol{z}^{(i)}\right)\right)\right) \right].$$

end for

- Sample minibatch of m noise samples {z⁽¹⁾,...,z^(m)} from noise prior p_q(z).
- · Update the generator by descending its stochastic gradient:

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^{m} \log \left(1 - D\left(G\left(\boldsymbol{z}^{(i)}\right)\right)\right).$$

The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

https://arxiv.org/pdf/1406.2661.pdf

4 D > 4 A > 4 B > 4 B > B 900

Generative deep nets (Radford et al. 16')

Early training examples





Figure 2: Generated bedrooms after one training pass through the dataset. Theoretically, the model could learn to memorize training examples, but this is experimentally unlikely as we train with a small learning rate and minibatch SGD. We are aware of no prior empirical evidence demonstrating memorization with SGD and a small learning rate.

https://arxiv.org/pdf/1511.06434.pdf

4 D > 4 P > 4 B > 4 B >

Generative deep nets (Radford et al. 16')

Later training examples





Figure 3: Generated bedrooms after five epochs of training. There appears to be evidence of visual under-fitting via repeated noise textures across multiple samples such as the base boards of some of the beds.

https://arxiv.org/pdf/1511.06434.pdf

Wasserstein GAN (Arjovsky et al. 17')





One of the central challenges with GANs is the ability to train the parameters. Improvements have been made through choice of generative architecture (DC-GAN of Radford) and through different training objective functions (W-GAN)

```
Algorithm 1 WGAN with gradient penalty. We use default values of \lambda = 10, n_{\text{critic}} = 5, \alpha =
0.0001, \beta_1 = 0, \beta_2 = 0.9.
```

Require: The gradient penalty coefficient λ , the number of critic iterations per generator iteration n_{critic} , the batch size m, Adam hyperparameters α, β_1, β_2 . **Require:** initial critic parameters w_0 , initial generator parameters θ_0 .

```
1: while \theta has not converged do
 2:
            for t = 1, ..., n_{critic} do
 3:
                   for i = 1, ..., m do
 4:
                          Sample real data x \sim \mathbb{P}_r, latent variable z \sim p(z), a random number \epsilon \sim U[0,1].
 5:
                          \tilde{x} \leftarrow G_{\theta}(z)
                          \hat{x} \leftarrow \epsilon \hat{x} + (1 - \epsilon)\tilde{x}
 6:
                          L^{(i)} \leftarrow D_w(\tilde{x}) - D_w(x) + \lambda(\|\nabla_{\tilde{x}}D_w(\hat{x})\|_2 - 1)^2
 7:
                   w \leftarrow \operatorname{Adam}(\nabla_w \frac{1}{m} \sum_{i=1}^m L^{(i)}, w, \alpha, \beta_1, \beta_2)
 9:
10:
            Sample a batch of latent variables \{z^{(i)}\}_{i=1}^m \sim p(z).
11:
            \theta \leftarrow \operatorname{Adam}(\nabla_{\theta} \frac{1}{m} \sum_{i=1}^{m} -D_{w}(G_{\theta}(z)), \theta, \alpha, \beta_{1}, \beta_{2})
12:
13: end while
```

```
https://arxiv.org/pdf/1704.00028.pdf
https://arxiv.org/pdf/1701.07875.pdf
```

4日 × 4周 × 4 国 × 4 国 × 国

Wasserstein GAN (Arjovsky et al. 17')

Examples of output from GAN architectures

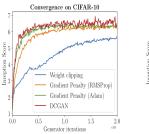




Figure 2: Different GAN architectures trained with different methods. We only succeeded in training every architecture with a shared set of hyperparameters using WGAN-GP.

https://arxiv.org/pdf/1704.00028.pdf





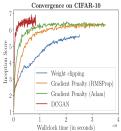


Figure 3: CIFAR-10 Inception score over generator iterations (left) or wall-clock time (right) for four models: WGAN with weight clipping, WGAN-GP with RMSProp and Adam (to control for the optimizer), and DCGAN. WGAN-GP significantly outperforms weight clipping and performs comparably to DCGAN.

https://arxiv.org/pdf/1704.00028.pdf

Large scale WGAN (Karras et al. 18')

Growing the encoder/decoder complexity



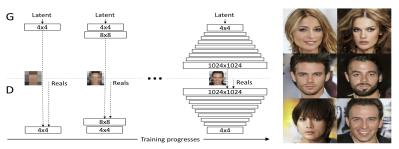


Figure 1: Our training starts with both the generator (G) and discriminator (D) having a low spatial resolution of 4×4 pixels. As the training advances, we incrementally add layers to G and D, thus increasing the spatial resolution of the generated images. All existing layers remain trainable throughout the process. Here $N\times N$ refers to convolutional layers operating on $N\times N$ spatial resolution. This allows stable synthesis in high resolutions and also speeds up training considerably. One the right we show six example images generated using progressive growing at 1024×1024 .

https://arxiv.org/abs/1710.10196

Large scale WGAN (Karras et al. 18')

Examples of synthetic faces





Figure 10: Top: Our CELEBA-HQ results. Next five rows: Nearest neighbors found from the training data, based on feature-space distance. We used activations from five VGG layers, as suggested by Chen & Koltun (2017). Only the crop highlighted in bottom right image was used for comparison in order to exclude image background and focus the search on matching facial features.

https://arxiv.org/abs/1710.10196

Understanding individual units in a DNN (Bau et al. 20')

Single units which reliably detect object classes



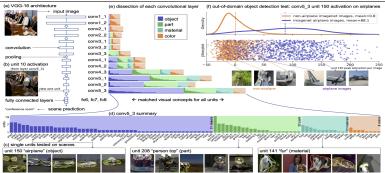


Fig. 1. The emergence of single-unit object detectors within a VGG-16 some classifier. (a) VGG-16 consists of 13 convolutional layers, conv1_1 through conv5_3, followed by trace bully connected layers, cof_7, for_6. (b) The activation of a single filter on an input image can be visualized as the region where the filter activates beyond its 16 pt 16, quantities that the convolution of the property of the control of the convolution of the property of the control of the convolutional layer are summarized, showing a broad diversity of electors for objects, object parts, materials, and colors. Many concepts are associated with multiple units, (a) Company all the layers of the most revenue that convolutional layers are summarized. The convolutional layers are summarized through a control of the convolutional layers are summarized through a control of the convolutional layers. (b) Although the training set contains no object detector for object detector in the convolutional layers. (b) Although the training set contains no object labels, unit 150 emerges as an 'airplane' object detector that activates much more strongly on airplane objects, as tested applies at a dataset of labels object displaced object may be reviewed, because the reviews. The filter plot shows past of the convolution of t

https://arxiv.org/abs/2009.05041

Understanding individual units in a DNN (Bau et al. 20')

Single units which reliably generate object classes



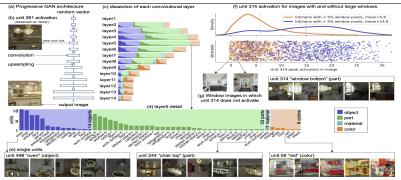


Fig. 3. The emergence of object- and part-specific units within a Progressive GAN generator(*19), (a) The analyzed Progressive GAN consists of 15 convolutional layers that transform a random input vector into a synthesized image of a kitchen, o) A nignle filter is visualized as the region of the object and extractional power that filter activates beyond its top remarks that the progressive GAN consists of the p

https://arxiv.org/abs/2009.05041

Editing individual units in a DNN (Bau et al. 20')

Removing units of classes decreases their generation



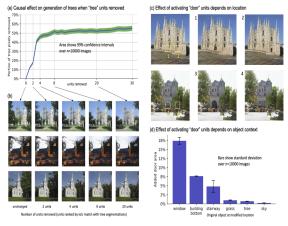


Fig. 4. The causal effect of altering units within a GAN generator. (a) When successively larger sets of units are removed from a GAN trained to generate outdoor church scenes, the tree area of the generated images is reduced. Removing 20 tree units removes more than half the generated tree pixels from the output. (b) Qualitative results: removing tree units affects trees while leaving other objects intact. Building parts that were previously occluded by trees are rendered as if revealing the objects that were behind the trees. (c) Doors can be added to buildings by activating 20 door units. The location, shape, size, and style of the rendered door depends on the location of the activated units. The same activation levels produce different doors. or no door at all (case 4) depending on locations. (d) Similar context dependence can be seen quantitatively: doors can be added in reasonable locations such as at the location of a window, but not in abnormal locations such as on a tree or in the sky.

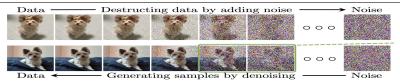
https://arxiv.org/abs/2009.05041



Diffusion models

36





Given an input x_0 , repeatedly scale and add additive noise

$$x_t = (1 - \beta_t)^{1/2} x_{t-1} + \mathcal{N}(0, \beta_t)$$

which is repeated from initial data x_0 to $x_t = \bar{\alpha}_t^{1/2} x_0 + (1 - \bar{\alpha}_t)^{1/2} \epsilon$ where $\epsilon \sim \mathcal{N}(0, (1 - \bar{\alpha}_t)I)$ with $\alpha_t = 1 - \beta_t$ and $\bar{\alpha}_t := \Pi_{s=0}^t \alpha_s$. The deep network is trained to denoise x_t to recover x_0 , unlike in a normal GAN x_t is not purely noise. New data generated by drawing from $\mathcal{N}(0, (1 - \bar{\alpha}_t)I)$ as if the x_0 is lost in the noise.

https://arxiv.org/abs/2209.00796