

Neural Image Processing: Style Transfer, Image Transformation, Neural Compression

Bolei Zhou IERG4190/IEMS5707

image credit to neuralstyle.art

Review: Generative Adversarial Networks

A two-player minimax game between a generator and a discriminator



Goodfellow et al. NIPS'14

Training objective for discriminator:

$$\max_{D} V(G, D) = E_{\mathbf{x} \sim p_{\text{data}}}[\log D(\mathbf{x})] + E_{\mathbf{x} \sim p_{G}}[\log(1 - D(\mathbf{x}))]$$

Training objective for generator:

$$\min_{G} V(G, D) = E_{\mathbf{x} \sim p_{\text{data}}}[\log D(\mathbf{x})] + E_{\mathbf{x} \sim p_{G}}[\log(1 - D(\mathbf{x}))]$$

$$\min_{\theta} \max_{\phi} V(G_{\theta}, D_{\phi}) = E_{\mathbf{x} \sim p_{\text{data}}}[\log D_{\phi}(\mathbf{x})] + E_{\mathbf{z} \sim p(\mathbf{z})}[\log(1 - D_{\phi}(G_{\theta}(\mathbf{z})))]$$

https://github.com/pytorch/examples/blob/master/dcgan/main.py

Example Code of Adversarial Training

for i, data in enumerate(dataloader, 0):

netD.zero_grad()

```
real_cpu = data[0].to(device)
batch_size = real_cpu.size(0)
label = torch.full((batch_size,), real_label, device=device)
```

```
output = netD(real cpu)
```

```
errD_real = criterion(output, label)
errD_real.backward()
D_x = output.mean().item()
```

```
# train with fake
noise = torch.randn(batch_size, nz, 1, 1, device=device)
fake = netG(noise)
label.fill_(fake_label)
output = netD(fake.detach())
errD_fake = criterion(output, label)
errD_fake.backward()
D_G_z1 = output.mean().item()
errD = errD_real + errD_fake
optimizerD.step()
```


https://github.com/pytorch/examples/blob/master/dcgan/main.py

Image Transformation through Latent Manipulation



InterFaceGAN, Shen, Gu, Tang, Zhou, CVPR'20

Image to Image Transformation

noise in image out





image in image out







Neural Image Processing

Super-resolution









Style transfer





Colorization

Image compression













Basic model Image Reconstruction through Auto-Encoder



Reconstruction Loss d(y', y)



Reconstruction Loss d(y', y)

- per-pixel loss does not capture perceptual differences between output and ground-truth images.
- Consider two identical images offset from each other by a few pixels; despite their perceptual similarity they would be very different as measured by per-pixel losses.









per-pixel loss will be huge

Image Reconstruction using CNN Feature Loss



Mahendran, A., Vedaldi, A.: Understanding deep image representations by inverting them. In: Proceedings of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR). (2015)

Image Reconstruction using CNN Feature Loss







feature loss could be at different layers



It means different layer's feature activation captures knowledge of the data

Mahendran, A., Vedaldi, A.: Understanding deep image representations by inverting them. In: Proceedings of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR). (2015)

Neural Style Transfer



Style Image

https://towardsdatascience.com/neural-style-transfer-tutorial-part-1-f5cd3315fa7f

Neural Style Transfer

Style *The Starry Night*, Vincent van Gogh, 1889





Style *The Muse*, Pablo Picasso, 1935





Neural Style Transfer: Live Demo

• https://reiinakano.com/arbitrary-image-stylization-tfjs/

How to Represent Content and Style?



Style Image

https://towardsdatascience.com/neural-style-transfer-tutorial-part-1-f5cd3315fa7f

How to Represent Content and Style?

Feature Loss $\|\phi_j(\hat{y}) - \phi_j(y)\|_2^2$



Style Image

https://towardsdatascience.com/neural-style-transfer-tutorial-part-1-f5cd3315fa7f

How to Represent Content and Style?

Feature Loss $\|\phi_j(\hat{y}) - \phi_j(y)\|_2^2$



Content Image





Style Image



Generated image

Representing Style as Texture

• Texture information is preserved as the correlation of features in CNNs





L. A. Gatys, A. S. Ecker, and M. Bethge. Texture Synthesis Using Convolutional Neural Networks. NeurIPS, 2015.

Representing Style with Gram Matrix



Transferring Style through back-propagation on both content and style losses



Gatys et al. Image Style Transfer Using Convolutional Neural Networks. CVPR'16



Drawback: back-propagation is slow and computationally expensive



Train an Image Transform Network



$$V(y)=\sum_n |y_{n+1}-y_n|.$$

Result of Style Transfer



Result of Super-Resolution



Ground TruthBicubicOurs (ℓ_{pixel}) SRCNN [11]Ours (ℓ_{feat}) This Image21.69 / 0.584021.66 / 0.588122.53 / 0.652421.04 / 0.6116Set14 mean25.99 / 0.730125.75 / 0.699427.49 / 0.750324.99 / 0.6731BSD100 mean25.96 / 0.68225.91 / 0.668026.90 / 0.710124.95 / 63.17

What are the metrics PSNR/SSIM?

Metrics for Image Reconstruction

- Pairwise similarity:
 - PSNR
 - MS-SSIM
- Distribution similarity (commonly used in GANs):
 - FID

Metric: PSNR (Peak Signal-to-Noise Ratio)

- **Peak signal-to-noise ratio (PSNR)** is an engineering term for the ratio between the maximum possible power of a signal and the power of corrupting noise that affects the fidelity of its representation.
- Commonly used to measure the quality of lossy compression reconstruction, an approximation to human perception (30dB to 50 dB)

$$\mathit{MSE} = rac{1}{m\,n}\sum_{i=0}^{m-1}\sum_{j=0}^{n-1}[I(i,j)-K(i,j)]^2$$

The PSNR (in dB) is defined as:

$$egin{aligned} PSNR &= 10 \cdot \log_{10} \left(rac{MAX_I^2}{MSE}
ight) \ &= 20 \cdot \log_{10} \left(rac{MAX_I}{\sqrt{MSE}}
ight) \ &= 20 \cdot \log_{10} (MAX_I) - 10 \cdot \log_{10} (MSI) \end{aligned}$$



Q=30, PSNR 36.81dB

Q=10, PSNR 31.45dB

Metric: MS-SSIM (Multi-scale structural similarity index measure)

(1)

(2)

(3)

• Perceptual image quality assessment

luminance:
$$l(\mathbf{x}, \mathbf{y}) = \frac{2 \mu_x \mu_y + C_1}{\mu_x^2 + \mu_y^2 + C_1}$$
,

SSIM:

contrast:
$$c(\mathbf{x}, \mathbf{y}) = \frac{2\sigma_x \sigma_y + C_2}{\sigma_x^2 + \sigma_y^2 + C_2},$$

structure: $s(\mathbf{x}, \mathbf{y}) = \frac{\sigma_{xy} + C_3}{\sigma_x \sigma_y + C_3},$

$$\mathrm{SSIM}(\mathbf{x}, \mathbf{y}) = \left[l(\mathbf{x}, \mathbf{y})\right]^{\alpha} \cdot \left[c(\mathbf{x}, \mathbf{y})\right]^{\beta} \cdot \left[s(\mathbf{x}, \mathbf{y})\right]^{\gamma}, \quad (5)$$

SSIM
$$(\mathbf{x}, \mathbf{y}) = \frac{(2 \,\mu_x \,\mu_y + C_1) \left(2 \,\sigma_{xy} + C_2\right)}{(\mu_x^2 + \mu_y^2 + C_1) \left(\sigma_x^2 + \sigma_y^2 + C_2\right)},$$
 (6)

where C_1 , C_2 and C_3 are small constants given by

$$C_1 = (K_1 L)^2, C_2 = (K_2 L)^2 \text{ and } C_3 = C_2/2,$$
 (4)



Metric: Frechet Inception Distance (FID)

- FID is the Wasserstein distance between two multidimensional Gaussian distributions $N(u, \Sigma)$ and $N(u_w, \Sigma_w)$
- Rather than modeling the image pixels, it models the deep features from CNN called Inception v3
- Commonly used in GAN evaluation

$$\mathrm{FID} = \left| \mu - \mu_w
ight|^2 + \mathrm{tr}(\Sigma + \Sigma_w - 2(\Sigma \Sigma_w)^{1/2}).$$



https://machinelearningmastery.com/how-to-implement-the-frechet-inception-distance-fid-from-scratch/ https://en.wikipedia.org/wiki/Fr%C3%A9chet_inception_distance

Result of Super-Resolution



BSD100 mean 25.96 / 0.682 25.91 / 0.6680 26.90 / 0.7101 24.95 / 63.17

PSNR/SSIM

Example: In-Domain GAN inversion



$$\begin{split} \min_{\Theta_E} \mathcal{L}_E &= ||\mathbf{x}^{real} - G(E(\mathbf{x}^{real}))||_2 + \lambda_{vgg} ||F(\mathbf{x}^{real}) - F(G(E(\mathbf{x}^{real})))||_2 \\ &- \lambda_{adv} \mathop{\mathbb{E}}_{\mathbf{x}^{real} \sim P_{data}} [D(G(E(\mathbf{x}^{real})))], \end{split}$$

Zhu et al. In-Domain GAN Inversion for Real Image Editing. ECCV'21

Example: In-Domain GAN inversion

GAN inversion: optimize z such that x = G(z)





Zhu et al. In-Domain GAN Inversion for Real Image Editing. ECCV'21

Image to Image Translation with Conditional GAN

Image to Image Translation



Isola et al. Image-to-Image Translation with Conditional Adversarial Nets. CVPR'17



loss function:

 $\min_{G} \max_{D} \mathbb{E}_{\mathbf{x},\mathbf{y}} [\log D(G(\mathbf{x})) + \log(1 - D(\mathbf{y}))]$
Input

Output

Ground-truth





data downloaded from the Google Map

Isola et al. Image-to-Image Translation with Conditional Adversarial Nets. CVPR'17

Input



L1 loss

GAN loss



More structured loss

Amazing applications based on Pix2Pix



#edges2cats [Chris Hesse]



https://phillipi.github.io/pix2pix/

Amazing applications based on Pix2Pix

#edges2cats



Christopher Hesse trained our model on converting edge maps to photos of cats, and included this in his interactive demo. Apparently, this is what the Internet wanted most, and #edges2cats briefly went viral. The above cats were designed by Vitaly Vidmirov (@vvid).

Alternative Face



Mario Klingemann used our code to translate the appearance of French singer Francoise Hardy onto Kellyanne Conway's infamous "alternative facts" interview. Interesting articles about it can be read here and here.

Person-to-Person



Brannon Dorsey recorded himself mimicking frames from a video of Ray Kurzweil giving a talk. He then used this data to train a Dorsey—Kurzweil translator, allowing him to become a kind of puppeter in control of Kurzweil's appearance.

Color palette completion



Colormind adapted our code to predict a complete 5-color palette given a subset of the palette as input. This application stretches the definition of what counts as "image-to-image translation" in an exciting way: if you can visualize your input/output data as images, then imageto-image methods are applicable! (not that this is necessarily the best choice of representation, just one to think about.)

Interactive Anime



Bertrand Gondouin trained our method to translate sketches—Pokemon, resulting in an interactive drawing tool.

Background masking



Kaihu Chen performed a number of interesting experiments using our method, including getting it to mask out the background of a portrait as shown above.

https://phillipi.github.io/pix2pix/



Unpaired data











Y

slides from Phillip Isola

Paired translation

Unpaired translation



["pix2pix", Isola, Zhu, Zhou, Efros, 2017]

["CycleGAN", Zhu*, Park*, Isola, Efros, 2017]



slides from Phillip Isola

Cycle Consistency





Cycle Consistency



$$\begin{split} \mathcal{L}_{\text{GAN}}(G, D_Y, X, Y) &= \mathbb{E}_{y \sim p_{\text{data}}(y)} [\log D_Y(y)] \\ &+ \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log(1 - D_Y(G(x)))], \\ \text{Our full objective is:} \\ \mathcal{L}_{\text{cyc}}(G, F) &= \mathbb{E}_{x \sim p_{\text{data}}(x)} [\|F(G(x)) - x\|_1] \\ &+ \mathbb{E}_{y \sim p_{\text{data}}(y)} [\|G(F(y)) - y\|_1]. \\ \end{split}$$

Amazing applications based on CycleGAN



Check this out: https://junyanz.github.io/CycleGAN/

Wide Applications of Adversarial Loss



 $\min_{G} \max_{D} \mathbb{E}_{\mathbf{x},\mathbf{y}} [\log D(G(\mathbf{x})) + \log(1 - D(\mathbf{y}))]$

Application: Compression Artifact Removal



Leonardo Galteri, Lorenzo Seidenari, Marco Bertini, and Alberto Del Bimbo. Deep generative adversarial compression artifact removal. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 4826–4835, 2017.

Application: Denoising





(a) Original





(c) CDnCNN-B (d) Our method Jingwen Chen, Jiawei Chen, Hongyang Chao, and Ming Yang. Image blind denoising with generative adversarial network based noise modeling. In The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), June 2018.

Application: Super-resolution

bicubic (21.59dB/0.6423)



4 x Upsampling SRResNet (23.53dB/0.7832)



SRGAN (21.15dB/0.6868)



original



Ledig et al. Photo-realistic single image superresolution using a generative adversarial network. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 4681–4690, 2017.

Application: Super-resolution



SRGAN



 $\min_{ heta_G} \max_{ heta_D} \mathbb{E}_{I^{HR} \sim p_{\text{train}}(I^{HR})}[\log D_{ heta_D}(I^{HR})] +$

(2)



Ledig et al. Photo-realistic single image superresolution using a generative adversarial network. CVPR, 2017.

ESRGAN (Enhanced SRGAN)





ESRGAN (Enhanced SRGAN)



https://arxiv.org/pdf/1809.00219.pdf https://github.com/xinntao/BasicSR



ESRGAN on Resident Evil 2 demo video: <u>https://www.youtube.com/watch?v=DZZ2HiN5_zc</u>

Face Hallucination versus Face Super-resolution 32 x upsampling









Face Hallucination as an Inversion Problem

$$\|DS(G(z)) - I_{LR}\|_p^p \le \epsilon.$$





PULSE: Self-Supervised Photo Upsampling via Latent Space Exploration of Generative Models. CVPR'20

- 🔶 🗙 🏠 🗈 theverge.com/21298762/face-depixelizer-ai-machine-learning-tool-pulse-stylegan-obama... 🛧 📀 関 🥴

REPORT \ TECH \ ARTIFICIAL INTELLIGENCE

What a machine learning tool that turns Obama white can (and can't) tell us about AI bias

A striking image that only hints at a much bigger problem By James Vincent | Jun 23, 2020, 3:45pm EDT

🛉 🔰 📝 share











Quantifying the Al Biases



Quantifying the Al Biases





Improving the Fairness of Deep Generative Models without Retraining. Tan, Shen, Zhou. <u>https://arxiv.org/pdf/2012.04842.pdf</u>

Model

Input

Neural Compression

Neural Image Compression



Input

r(y) = -log(P(y))

Reconstruction

Objective:
$$\mathcal{L}_{EG} = \mathbb{E}_{x \sim p_X} [\lambda \ r(y) + d(x, x')].$$

* P is the distribution model of y, using P and an entropy coding algorithm (e.g., arithmetic coding [30]), we can store y losslessly using bitrate r(y)

Neural Image Compression with Adversarial Loss



Input

Reconstruction

Objective:

$$\mathcal{L}_{EGP} = \mathbb{E}_{x \sim p_X} [\lambda r(y) + d(x, x') - \beta \log(D(x', y))],$$
$$\mathcal{L}_D = \mathbb{E}_{x \sim p_X} [-\log(1 - D(x', y))] + \mathbb{E}_{x \sim p_X} [-\log(D(x, y))].$$

* P is the distribution of code y, then an entropy coding algorithm on y

High-Fidelity Generative Image Compression (NeurIPS'20)



* Q is quantization, P is hyper-prior model

$$\mathcal{L}_{EGP} = \mathbb{E}_{x \sim p_X} [\lambda r(y) + d(x, x') - \beta \log(D(x', y))],$$

$$\mathcal{L}_D = \mathbb{E}_{x \sim p_X} [-\log(1 - D(x', y))] + \mathbb{E}_{x \sim p_X} [-\log(D(x, y))].$$

.https://hific.github.io/

Experiments of High Fidelity Compression (HiFiC)

- Training set: It consists of a large set of high-resolution images collected from the Internet
- Testing set:
 - Kodak [23] dataset (24 images),
 - CLIC2020 [46] testset (428 images)
 - DIV2K [2] validation set (100 images)

bpp: bit per pixel



Calculating bpp for an image

- Say we have a jpg image with 420x920, its size is 20635 bytes. What is the bpp in this image?
- We know 1 byte = 8 bits, bpp (bits per pixel) = 20635 * 8 bits / (420*920)
- Other knowledge about units:
 - 1 Mb = 1024 kb
 - 1 kb = 1024 bytes

Experiments of BPG





VARIATIONAL IMAGE COMPRESSION WITH A SCALE HYPERPRIOR. ICLR'18, https://arxiv.org/pdf/1802.01436.pdf



Interactive demo images: https://hific.github.io/

Learned Lossless Image Compression





[bpsp]	ImageNet32	Learned
L3C (ours)	4.76	\checkmark
PixelCNN [46]	3.83	\checkmark
MS-PixelCNN [32]	3.95	\checkmark
PNG	6.42	
JPEG2000	6.35	
WebP	5.28	
FLIF	5.08	

Table 3: Comparing bits per sub-pixel (bpsp) on the 32×32 im-

Practical Full Resolution Learned Lossless Image Compression. CVPR'20. https://arxiv.org/pdf/1811.12817.pdf

Learned Lossless Image Compression



Practical Full Resolution Learned Lossless Image Compression. CVPR'20. https://arxiv.org/pdf/1811.12817.pdf

Real-Time Adaptive Image Compression



Rippel and Bourdev. Real-Time Adaptive Image Compression. ICML'17. https://arxiv.org/pdt/1705.05823.pdt

Learned Video Compression



Rippel et al. Learned Video Compression. ICCV'19. <u>https://arxiv.org/pdf/1811.06981.pdf</u>


Building Al-native video.

WaveOne is leveraging the latest advancements in machine learning and deep learning to reinvent the way video is compressed, transported, stored, and analyzed.

https://www.wave.one/

Domain-Specific

Al-native compression can be optimized for specific types of video, dramatically improving coding efficiency.

















Semantics-Aware

Al-native compression can be trained to associate each pixel

with its object type, prioritizing the pixels that matter most.



Video Conferencing

Street Views

Gaming

Satellite

Baseline Bit Allocation

Semantics-Aware Bit Allocation

Semantics-Aware Result

Source

Other Resources on Neural Compression

- Neural Compression Workshop: https://neuralcompression.github.io/
- Workshop and Challenge on Learned Image Compression: <u>http://clic.compression.cc/</u>

Big tech companies really care about this technology!

Sponsored by



Other Application: Action Recognition on Compressed Video

- I-frames (intra-coded frames)
- P/B-frames (predictive frames)
 - Both B- and P- frames capture only what changes in the video





Other Application: Action Recognition on Compressed Video

- I-frames (intra-coded frames)
- P/B-frames (predictive frames)
 - Both B- and P- frames capture only what changes in the video





$$\mathcal{J}_i^{(t,k)} := \mu_{\mathcal{T}^{(k+1)}} \circ \cdots \circ \mu_{\mathcal{T}^{(t)}}(i).$$

$$I_i^{(t)} = I_{i-\mathcal{D}_i^{(t)}}^{(0)} + \mathcal{R}_i^{(t)}, \quad t = 1, 2, \dots,$$

$$\begin{split} \mathcal{D}_{i}^{(t)} &:= i - \mathcal{J}_{i}^{(t,k)}, \text{ and} \\ \mathcal{R}_{i}^{(t)} &:= \Delta_{\mathcal{J}_{i}^{(t,k+1)}}^{(k+1)} + \dots + \Delta_{\mathcal{J}_{i}^{(t,t-1)}}^{(t-1)} + \Delta_{i}^{(t)}, \end{split}$$

Wu et al. Compressed Video Action Recognition. CVPR'18 <u>https://arxiv.org/pdf/1712.00636.pdf</u>

Other Application: Action Recognition on Compressed Video

 Use ResNet-152 to model I-frames and ResNet-18 to model the motion vectors and residual



Wu et al. Compressed Video Action Recognition. CVPR'18 <u>https://arxiv.org/pdf/1712.00636.pdf</u>

New Frontier: Neural Image Processing

Super-resolution







Style transfer





Colorization

Image compression











