



ТЕХНОСФЕРА

Generative adversarial networks

Part 1

Галков Михаил

17 апреля 2017 г.

Understanding data

When we're learning to see, nobody's telling us what the right answers are — we just look. Every so often, your mother says "that's a dog", but that's very little information. You'd be lucky if you got a few bits of information — even one bit per second — that way. The brain's visual system has 10^{14} neural connections. And you only live for 10^9 seconds. So it's no use learning one bit per second. You need more like 10^5 bits per second. And there's only one place you can get that much information: from the input itself. — Geoffrey Hinton, 1996 (quoted in (Gorder 2006)).

Problems in supervised learning

Проблемы с размеченными датасетами

- ▶ Данных много, но меток обычно мало/очень шумные/дорогие
- ▶ Transfer learning - нетривиальная проблема
- ▶ Распределение на обучающем множестве не совпадает с распределением 'in the wild'
- ▶ Распределение данных меняется со временем/адаптируется к вашей системе

Подходы к решению

- ▶ Semi-supervised learning (в широком смысле)
- ▶ Online learning (например: stock prediction)

Generative modeling

Наша задача - научиться генерировать реалистичные изображения. Под generative modeling понимают решение двух задач.

1. Понимание того, как устроено распределение $p_{data}(x)$.
2. Построение функции G , которая по заданному вектору будет генерировать изображение.

Learning to generate Chairs, Tables and Cars pt.1

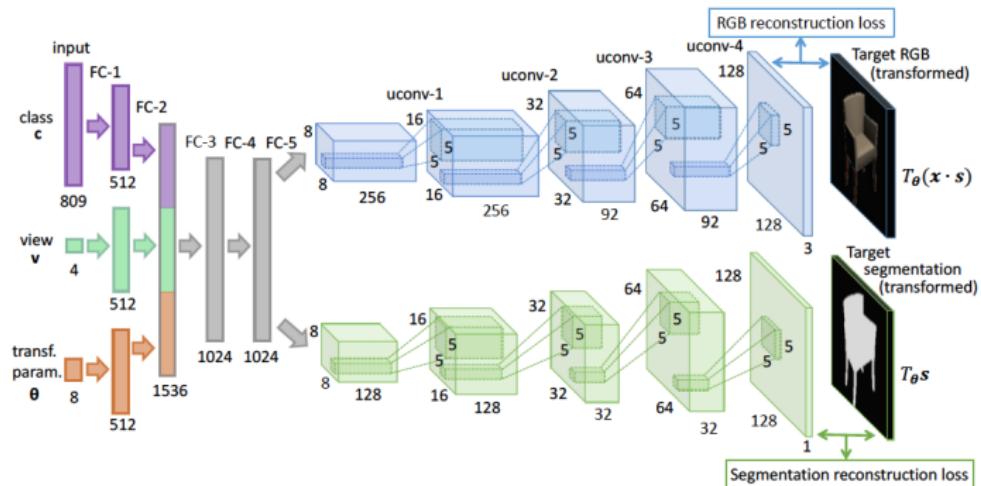


Fig. 1. Architecture of a 2-stream network that generates 128×128 pixel images. Layer names are shown above: FC - fully connected, uconv - unpooling+convolution.

$$\min_{\mathbf{W}} \sum_{i=1}^N L_{RGB}(T_{\theta^i}(x^i \cdot s^i), u_{RGB}(h(c^i, v^i, \theta^i))) + \lambda \cdot L_{segm}(T_{\theta^i}s^i, u_{segm}(h(c^i, v^i, \theta^i))),$$

Learning to generate Chairs, Tables and Cars pt.2

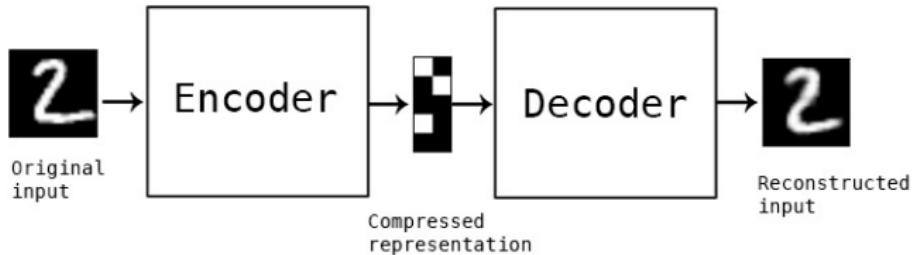


Почему получилось так удачно? - **Огромное** количество информации на входе

1. категория семпла
2. углы и повороты
3. заданные параметры трансформаций
4. сегментация

Можно ли автоматизировать извлечение этих фичей?

Autoencoder



Идея очень похожая на PCA: найти сжимающее отображение исходных данных (Encoder) в пространство меньшей размерности, такое, что из него возможно восстановить исходное изображение (Decoder).

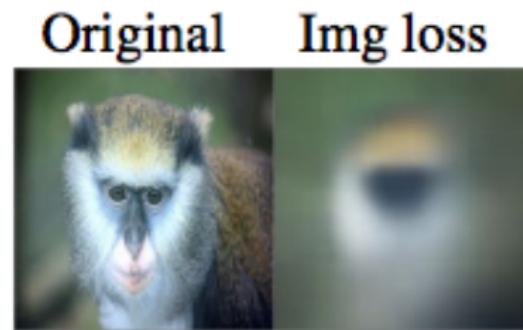
$$X \rightarrow Z \rightarrow X' \tag{1}$$

$$X \approx X' \tag{2}$$

Наивный подход

Естественно, хочется взять MSE как лосс-функцию и попробовать оптимизировать.

Но, к сожалению, результаты не впечатляют.



Возможно, мы взяли слишком простой лосс, можно придумать что-то лучше?

Deep reconstructions pt.1

layer	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
name	conv1	relu1	mpool1	norm1	conv2	relu2	mpool2	norm2	conv3	relu3	conv4	relu4	conv5	relu5	mpool5	fc6	relu6	fc7	relu7	fc8
type	cnv	relu	mpool	nrm	cnv	relu	mpool	nrm	cnv	relu	cnv	relu	cnv	relu	mpool	cnv	relu	cnv	relu	cnv
channels	96	96	96	96	256	256	256	256	384	384	384	384	256	256	256	4096	4096	4096	4096	1000
rec. field	11	11	19	19	51	51	67	67	99	99	131	131	163	163	195	355	355	355	355	355

Table 2. **CNN-A structure.** The table specifies the structure of CNN-A along with receptive field size of each neuron. The filters in layers from 16 to 20 operate as “fully connected”: given the standard image input size of 227×227 pixels, their support covers the whole image. Note also that their receptive field is larger than 227 pixels, but can be contained in the image domain due to padding.

$\Phi : \mathbb{R}^{H \times W \times C} \rightarrow \mathbb{R}^d$ - representation function

$\Phi_0 = \Phi(x_0)$ - representation of specific image

L - loss function

R - regularization term

$$\mathbf{x}^* = \underset{\mathbf{x} \in \mathbb{R}^{H \times W \times C}}{\operatorname{argmin}} \ell(\Phi(\mathbf{x}), \Phi_0) + \lambda \mathcal{R}(\mathbf{x})$$

Deep reconstructions pt.2

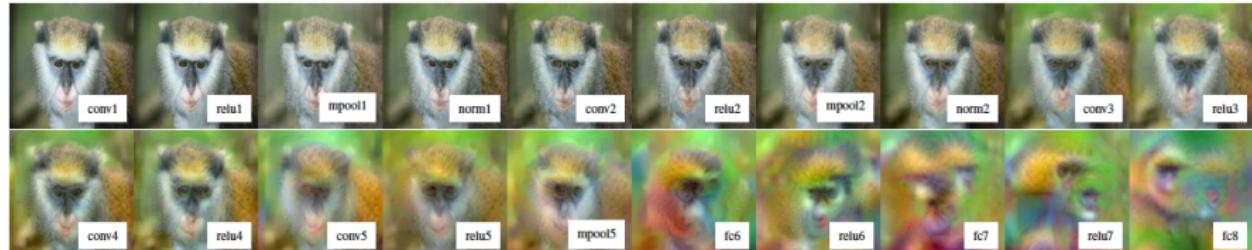


Figure 6. **CNN reconstruction.** Reconstruction of the image of Fig. 5.a from each layer of CNN-A. To generate these results, the regularization coefficient for each layer is chosen to match the highlighted rows in table 3. This figure is best viewed in color/screen.

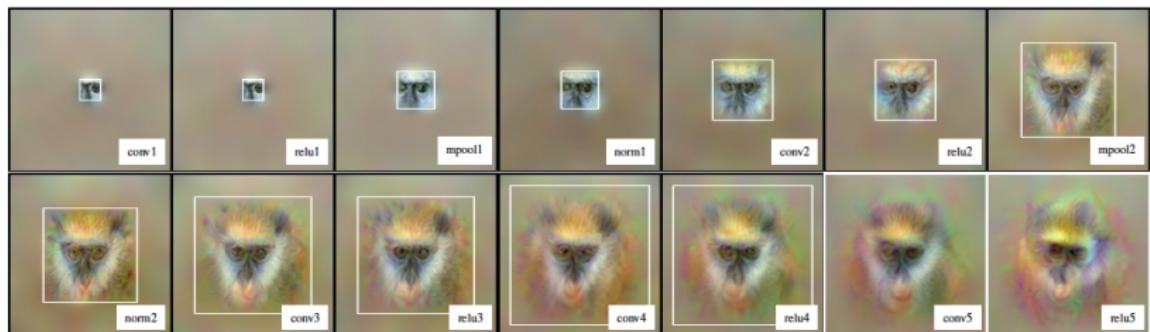
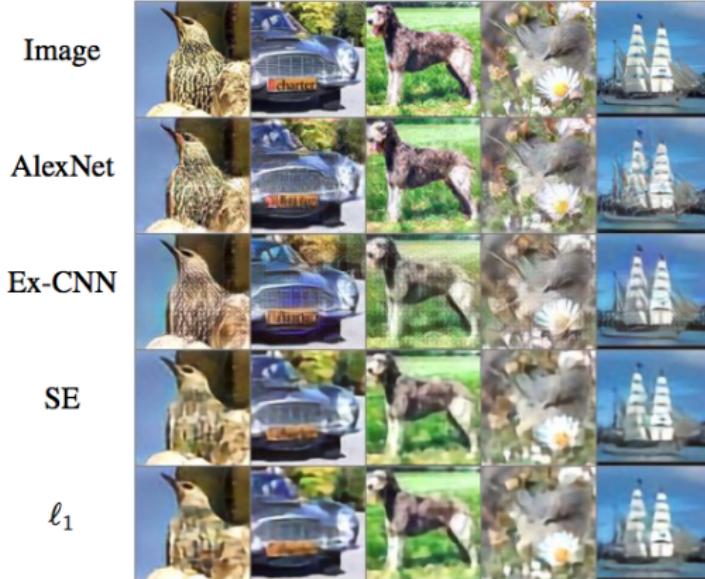


Figure 9. **CNN receptive field.** Reconstructions of the image of Fig. 5.a from the central 5×5 neuron fields at different depths of CNN-A. The white box marks the field of view of the 5×5 neuron field. The field of view is the entire image for conv5 and relu5.

Основные выводы

- ▶ Как глубина влияет на степень реконструкции? (текстура, цвет, положение)
- ▶ Каков эффект max-pooling?
- ▶ Как можно использовать готовые классификаторы?

Эмпирические результаты



На практике

Действительно работает, настроить параметры несложно.

- ▶ Лучше всего использовать комбинацию loss-функций
- ▶ Image-to-image loss: L1
- ▶ Deep features loss: VGG 16/19
- ▶ Регуляризация: total variation

Но можно еще лучше!

Generative adversarial networks pt.1

Введем основные определения:

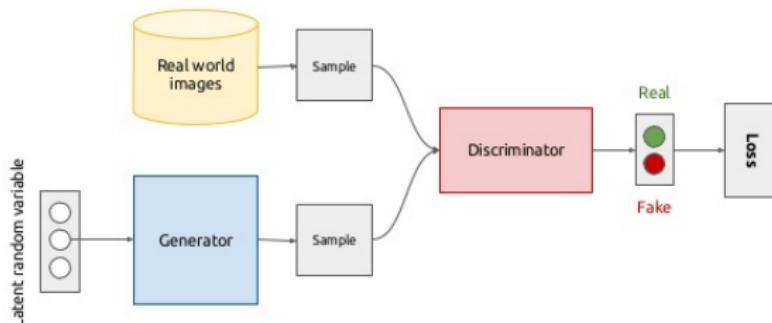
$z \sim p_z(z)$ - noise vector

$p_g(z)$ - распределение сгенерированных картинок из noise

$p_{data}(x)$ - распределение настоящих картинок

$G(z)$ - генератор (генерирует картинку из z)

$D(x)$ - дискриминатор (отличает реальные от сгенерированных)



Generative adversarial networks pt.2

Задача генератора - сгенерировать картинку, которую дискриминатор сочтет реалистичной.

Задача дискриминатора - отличить сгенерированную от реальной.

Математически - это игра двух игроков:

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))].$$

Generative adversarial networks pt.3

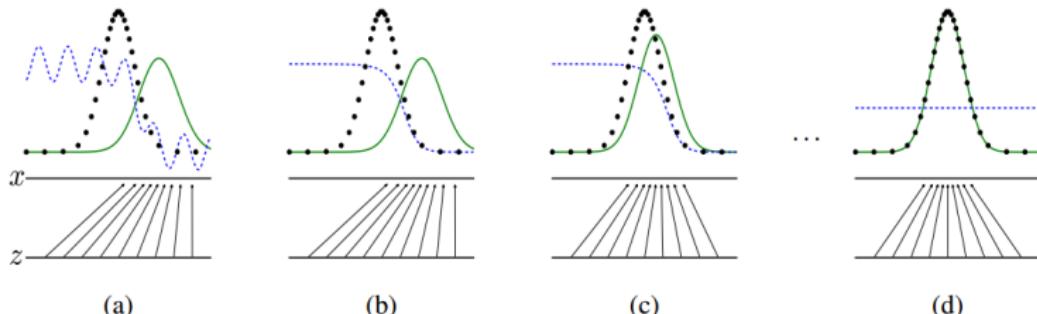


Figure 1: Generative adversarial nets are trained by simultaneously updating the discriminative distribution (D , blue, dashed line) so that it discriminates between samples from the data generating distribution (black, dotted line) p_x from those of the generative distribution p_g (G) (green, solid line). The lower horizontal line is the domain from which z is sampled, in this case uniformly. The horizontal line above is part of the domain of x . The upward arrows show how the mapping $x = G(z)$ imposes the non-uniform distribution p_g on transformed samples. G contracts in regions of high density and expands in regions of low density of p_g . (a) Consider an adversarial pair near convergence: p_g is similar to p_{data} and D is a partially accurate classifier. (b) In the inner loop of the algorithm D is trained to discriminate samples from data, converging to $D^*(x) = \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)}$. (c) After an update to G , gradient of D has guided $G(z)$ to flow to regions that are more likely to be classified as data. (d) After several steps of training, if G and D have enough capacity, they will reach a point at which both cannot improve because $p_g = p_{\text{data}}$. The discriminator is unable to differentiate between the two distributions, i.e. $D(x) = \frac{1}{2}$.

Theoretical results pt. 1

Theorem

Пусть G - фиксирован, тогда $D_G^*(x) = \frac{p_{data}(x)}{p_{data}(x) + p_g(x)}$

Доказательство.

$$\begin{aligned} V(G, D) &= \mathbb{E}_{x \sim p_{data}(x)} \log D(x) + \mathbb{E}_{z \sim p_z(z)} \log(1 - D(G(z))) \\ &= \int_x p_{data}(x) \log D(x) + p_g(x) \log(1 - D(x)) \end{aligned} \tag{3}$$

Теперь воспользуемся тем, что функционал $f(y) = a \log(y) + b \log(1 - y)$ на $[0, 1]$ достигает максимума в точке $\frac{a}{a+b}$. □

Theoretical results pt. 2

Подставим оптимальный дискриминатор в value-function

$$\begin{aligned} C(G) &= \int_x p_{data}(x) \log D^*(x) + p_g(x) \log(1 - D^*(x)) \\ &= \int_x p_{data}(x) \log \frac{p_{data}(x)}{p_{data}(x) + p_g(x)} \\ &\quad + p_g(x) \log \frac{p_g(x)}{p_{data}(x) + p_g(x)} \end{aligned}$$

Если подставить $p_{data}(x)$ вместо $p_g(x)$, получим:

$C(G) = -\log 4$, и если вынести $-\log 4$, то

$$C(G) = -\log 4 + KL(p_{data} \parallel \frac{p_{data} + p_g}{2}) + KL(p_g \parallel \frac{p_{data} + p_g}{2}) \quad (4)$$

Training algorithm

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 $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ from noise prior $p_g(\mathbf{z})$.
- Sample minibatch of m examples $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ from data generating distribution $p_{\text{data}}(\mathbf{x})$.
- Update the discriminator by ascending its stochastic gradient:

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

end for

- Sample minibatch of m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ from noise prior $p_g(\mathbf{z})$.
- Update the generator by descending its stochastic gradient:

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

end for

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

Takeaway

Преимущества GAN

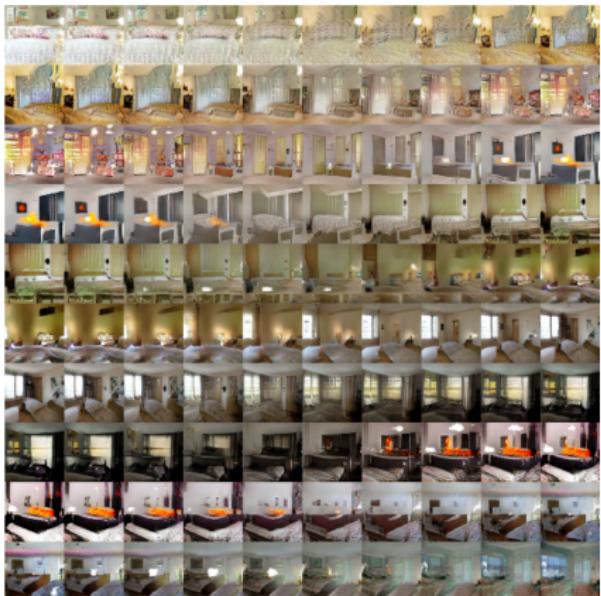
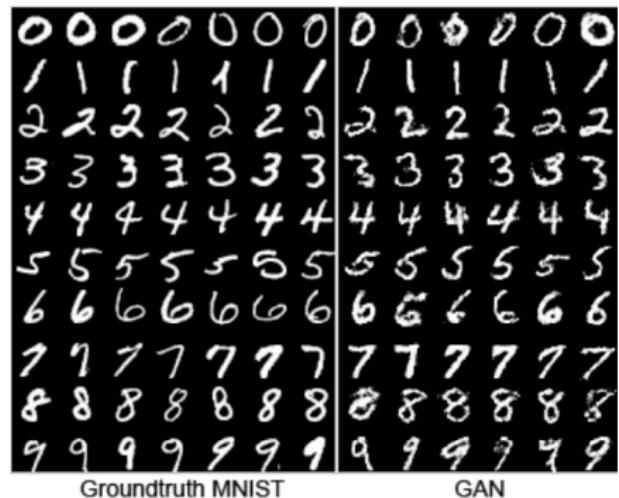
- ▶ Теоретические гарантии сходимости
- ▶ Можно обучать обычным SGD
- ▶ Решает в явном виде задачу generative modeling
- ▶ Но неявным образом (нейросети)

Недостатки GAN

- ▶ Нестабильное обучение
- ▶ Очень долгая сходимость
- ▶ Mode-collapsing
- ▶ Generator/Discriminator starvation
- ▶ Поиск оптимальных параметров - pure luck

Тем не менее, на практике GAN почти всегда дает заметные улучшения. Следующая лекция о том, как улучшить GAN.

Картинки для привлечения внимания



GAN gif



Вопросы

