

Generative Models, part I

Machine Learning in High Energy Physics

Maxim Borisyak

National Research University Higher School of Economics

June 30, 2019

Generative models

-

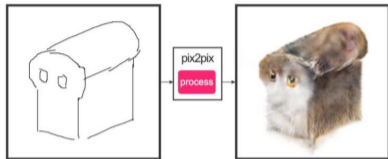
Generative models

Given samples of a random variable X , find X' , such that:

$$P(X) \approx P(X')$$

Applications

- auxiliary tasks:
 - pivoted models;
- data manipulation:
 - realistic image-to-image translation;
- approximation of existing generators:
 - fast Monte-Carlo;
- data compression.



Images are from <https://arxiv.org/abs/1611.07004>

Types of generative models

Density estimation:

- usually, density is known up to a constant (e.g. RBM):

$$f(x) = C \cdot p(x)$$

- sampling is via MCMC;
- challenging in high dimensional spaces.

Types of generative models

Sampling procedure:

- learning transformation from a simple random variable to the target one:

$$Z \sim \mathcal{N}^n(0, 1);$$

$$X' = f(Z).$$

- density is often intractable:

$$p(x) = \sum_{z|f(z)=x} p(z) \left| \frac{\partial}{\partial z} f(z) \right|^{-1}$$

Kernel Density Estimation

Kernel Density Estimation

Main idea: place a small gaussian-like function around each sample:

- sample $\{x_i\}_{i=1}^N$
- kernel $k(x)$:

$$\int_{\mathcal{X}} k(x) dx = 1;$$

- kernel width h ;

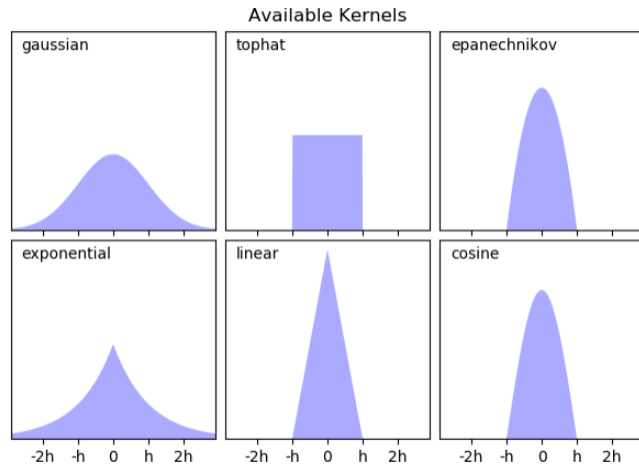
$$f_{KDE}(x) = \frac{1}{hN} \sum_i k\left(\frac{x - x_i}{h}\right);$$

(almost) theorem

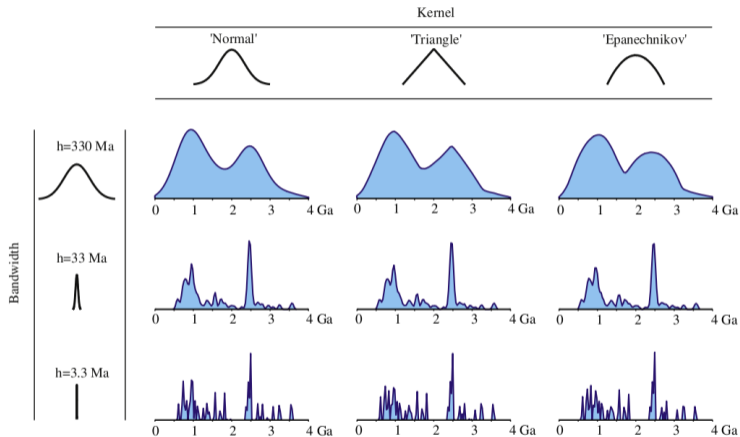
When $N \rightarrow \infty$ and $h \rightarrow 0$:

$$f_{KDE}(x) \rightarrow p(x).$$

Kernel Density Estimation



Kernel Density Estimation



Kernel Density Estimation

- kernel selection:
 - Gaussian, Epanechnikov kernels produce smooth densities;
 - Epanechnikov, linear kernels are faster to compute:
 - depends on nearest neighbour algorithm;
- kernel width:
 - a hyper-parameter;
 - selected by e.g. cross-validation.

Kernel Density Estimation

Advantages:

- (almost) simple;
- (almost) trivial sampling;
- (almost) no training;
- explicit density estimation;
- great for 1-2 dimensional problems.

Disadvantages:

- complexity $O(\log N)$ for evaluation in one point;
- performs poorly on high-dimensional problems;
- need to memorize the whole training set.

Gaussian Mixtures

Gaussian Mixtures

A similar to KDE idea — describe probability density as a mixture of Gaussians:

- unlike KDE, centers and widths of Gaussians are learnt;
- usually, number of components n is much smaller than the number of samples N .

$$f(x) = \frac{1}{n} \sum_j \phi(x | \mu_j, \Sigma_j)$$

where:

- $\phi(x | \mu, \Sigma)$ — density of the Gaussian distribution with mean μ and covariance matrix Σ .

Gaussian Mixtures

Training is done by maximizing likelihood:

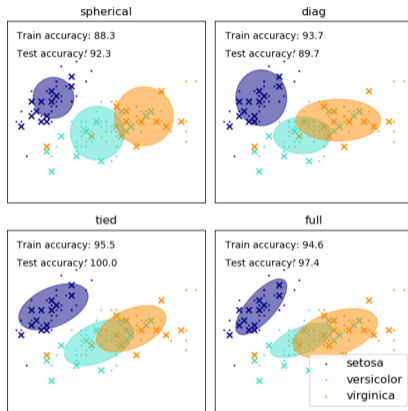
$$L = \sum_i \log \left(\frac{1}{N} \sum_j \phi(x_i | \mu_j, \Sigma_j) \right) \rightarrow \max$$

- Expectation-Maximization algorithm;
- priors are possible (ML \rightarrow MAP).

Gaussian Mixtures

Types of GMM:

- spherical: $\Sigma = \sigma \cdot \mathbb{I}$;
- diagonal: $\Sigma = \text{diag}(\sigma_1, \sigma_2, \dots)$;
- tied: $\Sigma_1 = \Sigma_2 = \dots$;
- ...
- full;



<https://scikit-learn.org/stable/modules/mixture.html>

Gaussian Mixtures

Advantages:

- simple;
- trivial sampling;
- explicit density estimation;
- great for 1-2 dimensional problems;
- tends to have fewer parameters than KDE.

Disadvantages:

- complexity $O(n)$ for evaluation in one point;
- performs poorly on high-dimensional problems;

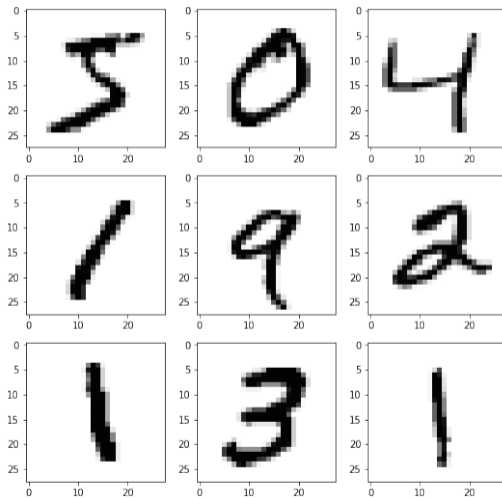
Variational AutoEncoder

Latent variables revisited

Before generating a sample, model should first decide what it should generate:

- which digit to generate: 0, 1, ..., 9
- width of stokes;
- 'speed';
- etc.

Such decision can be represented as **latent variables**.



Variational AutoEncoder

VAE **non-deterministically** transforms latent variables Z into samples X :

1. given latent variables z , VAE computes $f(z)$;
2. $f(z)$ represents parameters of some distribution;
3. examples are sampled from $P(x | f(z))$.

This section is largely based on <https://arxiv.org/abs/1606.05908>.

Variational AutoEncoder

Common choices:

- continuous data: $P(x | f(z)) = \mathcal{N}(x | f(z), \sigma^2 \mathbb{I})$;
 - σ – hyper-parameter;
 - \mathbb{I} – identity matrix.
- discrete data:
 - $P(x | f(z)) = \text{Bi}(x | f(z))$;
 - $P(x | f(z)) = \text{Multi}(x | f(z))$.

How to choose latent variables?

How to choose latent variables?

- let \hat{P} be a magical optimal choice of latent variables;
- let $z \sim \mathcal{N}^m(0, 1)$;
- if model \mathcal{G} has enough capacity, then

$$\exists g \in \mathcal{G} : g(z) \sim \hat{P}.$$

Let the network assign the meaning of the latent variables.

VAE training

Maximum Likelihood:

$$\sum_i \log P(x_i) \rightarrow \max$$

where $\{x_i\}_{i=1}^N$ – observed data.

$$P(x) = \int P(x | z)P(z)dz = \mathbb{E}_Z P(x | Z)$$

- for the most of z . $P(x | z) \approx 0$.

How to deal with the integral?

Variational bound

$$P(x) = \int P(x | z)P(z)dz = \mathbb{E}_Z P(x | Z)$$

In order to make sampling tractable, $P(z)$ can be replaced by some $Q(z | x)$:

$$P(x) = \mathbb{E}_Z P(x | Z) \rightarrow \mathbb{E}_{Z \sim Q(z|x)} P(x | Z)$$

Let's consider KL divergence:

$$\text{KL} (Q(z | x) \| P(z | x)) = \mathbb{E}_{Z \sim Q(z|x)} [\log Q(Z | x) - \log P(Z | x)]$$

Variational bound

$$\text{KL} (Q(z | x) | P(z | x)) =$$

$$\mathbb{E}_{Z \sim Q(z|x)} [\log Q(Z | x) - \log P(Z | x)] =$$

$$\mathbb{E}_{Z \sim Q(z|x)} [\log Q(Z | x) - \log P(x | Z) - \log P(Z)] + \log P(x)$$

$$\log P(x) - \text{KL} (Q(z | x) || P(z | x)) = \mathbb{E}_{Z \sim Q(z|x)} \log P(x | Z) - \text{KL} (Q(z | x) || P(z))$$

Variational bound

$$l(x) = \underbrace{\log P(x)}_{\text{MLE objective}} - \underbrace{\text{KL}(Q(z|x) \| P(z|x))}_{\text{inference penalty, } \geq 0} =$$
$$\underbrace{\mathbb{E}_{Z \sim Q(z|x)} \log P(x|Z)}_{\text{reconstruction error}} - \underbrace{\text{KL}(Q(z|x) \| P(z))}_{\text{regularization}}$$

$$\boxed{\log P(x) \geq l(x) \rightarrow \max}$$

$$\mathcal{L} = \sum_i \left[\mathbb{E}_{Z \sim Q(z|x_i)} \log P(x_i | Z) - \text{KL} (Q(z | x_i) \parallel P(z)) \right]$$

- reconstruction error can be estimated by sampling z from $Q(z | x_i)$:

$$\mathbb{E}_{Z \sim Q(z|x_i)} \log P(x_i | Z) \rightarrow \text{RE}(x_i, z)$$

- regularization term is, usually, computed analytically.

$$\text{RE}(x, z) = \log P(x | z)$$

- for Gaussian posterior i.e. $P(x | z) = \mathcal{N}(x | f(z), \sigma^2 I)$:

$$\text{RE}(x, z) \propto (f(z) - x)^2$$

- for Bernoulli posterior (e.g. for discrete output) $P(X = 1 | z) = f(z)$:

$$\text{RE}(x, z) = x \log f(z) + (1 - x) \log(1 - f(z))$$

Limitations



(a)



(b)



(c)

Image (b) — slightly altered image (a), image (c) — image (a) shifted by several pixels.
Under MSE metric, image (b) is much closer to (a), than (c) to (a).

Regularization

Consider:

- $Q(z | x) = \mathcal{N}(z | \mu(x), \Sigma(x));$
- $P(z) = \mathcal{N}(0, I);$

$$\text{KL}(\mathcal{N}(x | \mu(z), \Sigma(z)) \| \mathcal{N}(x | \mu(z), \Sigma(z))) =$$

$$\frac{1}{2} (\text{tr}(\Sigma(x)) + \|\mu(x)\|^2 - k - \log \det \Sigma(x)) =$$

$$\frac{1}{2} \left(\|\mu(x)\|^2 + \sum_i \Sigma_{ii}(x) - \log \Sigma_{ii}(x) \right) - \frac{k}{2}$$

Training time

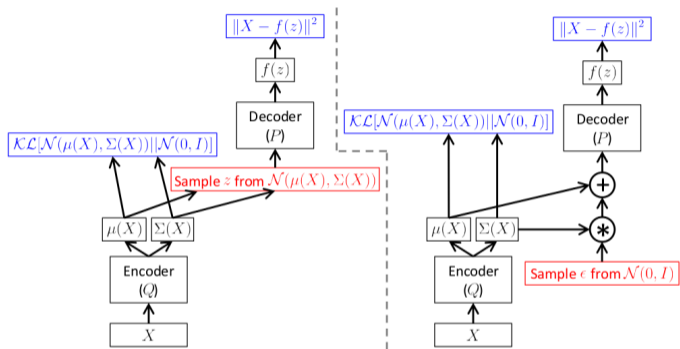


Figure 4: A training-time variational autoencoder implemented as a feed-forward neural network, where $P(X|z)$ is Gaussian. Left is without the “reparameterization trick”, and right is with it. Red shows sampling operations that are non-differentiable. Blue shows loss layers. The feedforward behavior of these networks is identical, but backpropagation can be applied only to the right network.

Testing time

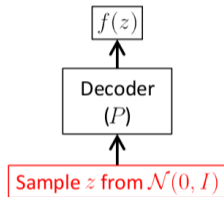


Figure 5: The testing-time variational “autoencoder,” which allows us to generate new samples. The “encoder” pathway is simply discarded.

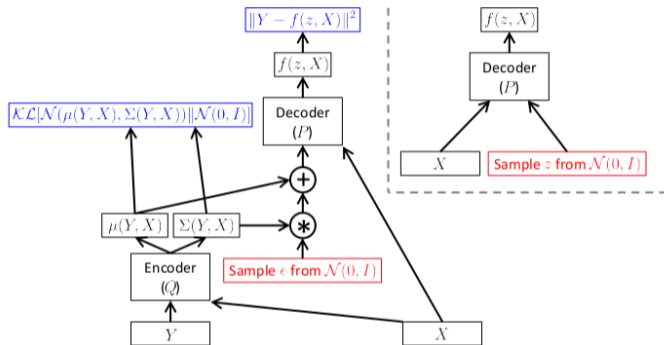


Figure 6: Left: a training-time conditional variational autoencoder implemented as a feedforward neural network, following the same notation as Figure 4. Right: the same model at test time, when we want to sample from $P(Y|X)$.

Generative Adversarial Networks

Fitting Distributions

Notation: Q - ground truth distribution, P - model distribution.

Maximum Likelihood:

$$\mathcal{L} = \sum_i \log P(x_i) \approx \mathbb{E}_{X \sim Q} \log P(X) \rightarrow_P \min;$$

$$\text{KL}(Q \parallel P) = \mathbb{E}_{X \sim Q} \log Q(X) - \mathbb{E}_{X \sim Q} \log P(X) \rightarrow_P \min.$$

Jensen-Shannon distance:

$$\text{JS}(P, Q) = \frac{1}{2} [\text{KL}(P \parallel M) + \text{KL}(Q \parallel M)] \rightarrow_P \min;$$

$$M = \frac{1}{2}(P + Q).$$

Approximating JS distance

$$\begin{aligned} \text{JS}(P, Q) &= \frac{1}{2} \left[\mathbb{E}_{X \sim P} \log \frac{P(X)}{M(X)} + \mathbb{E}_{X \sim Q} \log \frac{Q(X)}{M(X)} \right] = \\ & \frac{1}{2} \left[\mathbb{E}_{X \sim P} \log \frac{P(X)}{P(X) + Q(X)} + \mathbb{E}_{X \sim Q} \log \frac{Q(X)}{P(X) + Q(X)} \right] + \log 2 = \\ & \mathbb{E}_{X \sim M} \frac{P(X)}{P(X) + Q(X)} \log \frac{P(X)}{P(X) + Q(X)} + \mathbb{E}_{X \sim M} \frac{Q(X)}{P(X) + Q(X)} \log \frac{Q(X)}{P(X) + Q(X)} + \log 2 \end{aligned}$$

Approximating JS distance

Let's introduce binary indicator y : $y = 1$ if x is sampled from P and $y = 0$ for Q :

$$\text{JS}(P, Q) - \log 2 =$$

$$\mathbb{E}_{X \sim M} \frac{P(X)}{P(X) + Q(X)} \log \frac{P(X)}{P(X) + Q(X)} + \mathbb{E}_{X \sim M} \frac{Q(X)}{P(X) + Q(X)} \log \frac{Q(X)}{P(X) + Q(X)} =$$

$$\mathbb{E}_{X \sim M, Y} P(Y = 1 | X) \log P(Y = 1 | X) + P(Y = 0 | X) \log P(Y = 0 | X) =$$

$$\max_f \mathbb{E}_{X, Y} Y \log f(X) + (1 - Y) \log(1 - f(X))$$

Approximating JS distance

$$\text{JS}(P, Q) =$$

$$\log 2 + \max_f \mathbb{E}_{X,Y} Y \log f(X) + (1 - Y) \log(1 - f(X)) =$$

$$\log 2 - \min_f \mathcal{L}(f | P, Q)$$

where \mathcal{L} — cross-entropy loss.

Approximating JS distance

$$\arg \min_P \text{JS}(P, Q) = \arg \max_P \left[\min_f \mathcal{L}(f | P, Q) \right]$$

Generative Adversarial Networks

GAN makes no assumptions about nature of P :

- the most popular choice is via a **generator** g .

$$Z \sim \mathcal{N}^m(0, 1);$$

$$X = g(Z).$$

Discriminator

Minimization of $\mathcal{L}(f|P, Q)$ is a classical classification problem:

- f is often defined by a neural network — **discriminator**;
- Q is defined by given dataset;
- P is defined by the generator.

Algorithm 1 Discriminator Training

```
while not enough do
    sample  $x$  from the dataset;
    sample latent variables  $z$  from  $\mathcal{N}^m(0, 1)$ ;

     $\theta \leftarrow \theta + \lambda_\theta \nabla_\theta [\log f_\theta(x) + \log (1 - f_\theta(g_\psi(z)))]$ 
end while
```

- θ – parameters of the discriminator f_θ ;
- ψ – parameters of the generator g_ψ ;
- λ_θ – SGD learning rate.

Generator training

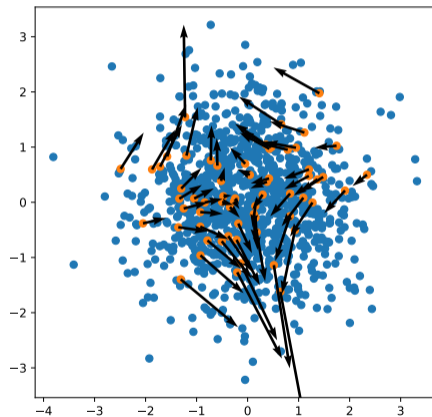
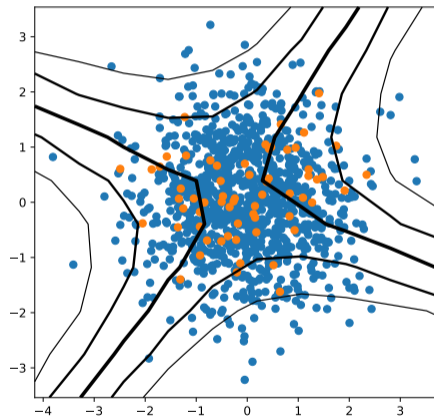
Generator is often trained by gradient methods, using:

$$\Delta\psi \propto \nabla_{\psi} \mathbb{E}_Z \log(1 - f(g_{\psi}(Z))) \Big|_{f=f^*}$$

as subderivative, where:

- $f^* = \arg \min_f \mathcal{L}(f^* | P_{\psi}, Q)$.

Generator training

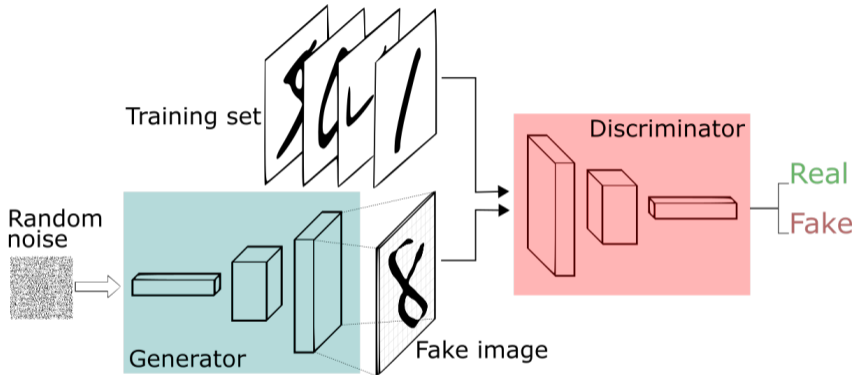


Algorithm 2 Generative Adversarial Training

```
while not enough do
  for  $i := 1, \dots, n$  do
    sample  $x$  from the dataset;
    sample latent variables  $z$  from  $\mathcal{N}^m(0, 1)$ ;
     $\theta \leftarrow \theta + \lambda_\theta \nabla_\theta [\log f_\theta(x) + \log (1 - f_\theta(g_\psi(z)))]$ 
  end for

  sample latent variables  $z$  from  $\mathcal{N}^m(0, 1)$ ;
   $\psi \leftarrow \psi - \lambda_\psi \nabla_\psi [\log (1 - f_\theta(g_\psi(z)))]$ 
end while
```

Generative Adversarial Networks



Source: <https://sthalles.github.io/assets/dcgan/GANs.png>

Game interpretation

$$\mathcal{L}(\theta, \psi) = -\frac{1}{2} \left[\mathbb{E}_{X \sim Q} \log f_{\theta}(X) + \mathbb{E}_{Z \sim Z} \log (1 - f_{\theta}(g_{\psi}(Z))) \right]$$

Min-max game:

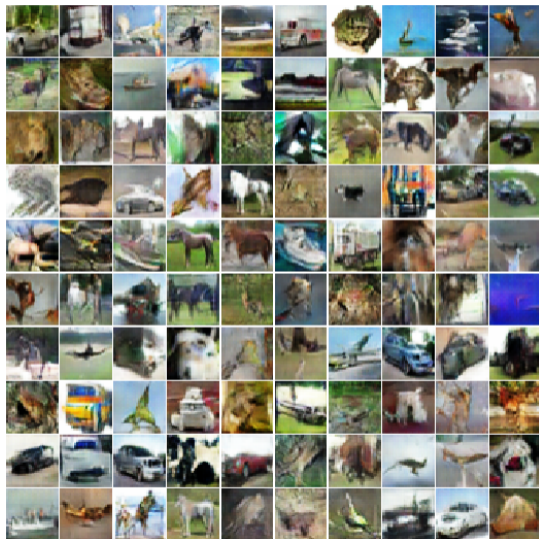
- goal of discriminator: distinguish between real and generated samples:

$$\mathcal{L}(\theta, \psi) \rightarrow_{\theta} \min$$

- goal of generator: 'fool' discriminator:

$$\mathcal{L}(\theta, \psi) \rightarrow_{\psi} \max$$

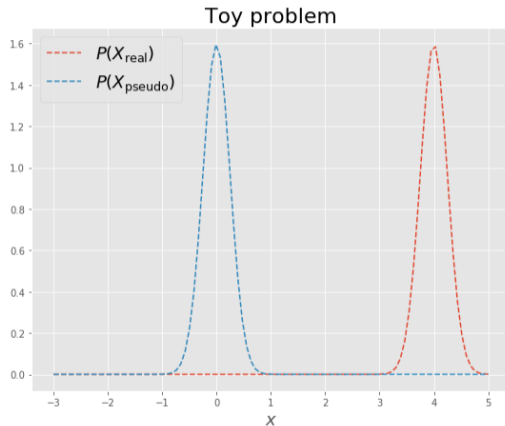
CIFAR examples



Vanishing gradients

Consider toy problem:

- powerfull discriminator;
- (almost) disjoint supports:
 - unlucky initial guess;
 - target data is on low-dimensional manifold;



Vanishing gradients

After training discriminator:

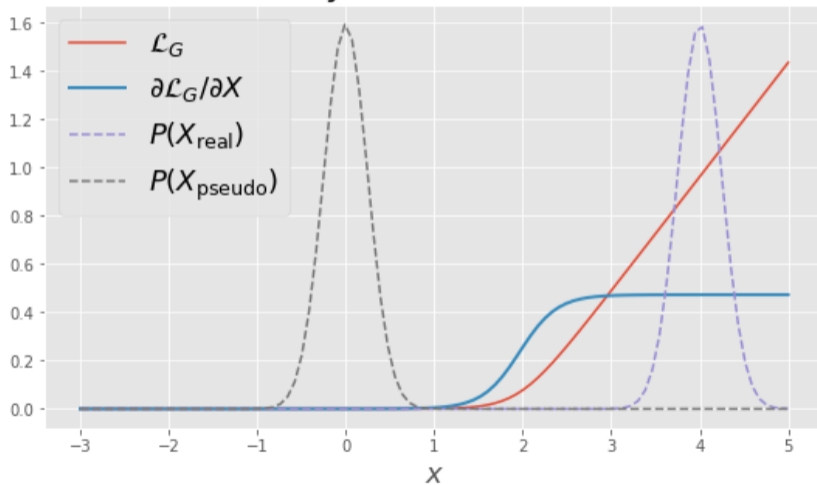
$$\frac{\partial \mathcal{L}(\theta, \psi)}{\partial \psi} = -\frac{1}{1 - f(g(z))} \cdot \frac{\partial f}{\partial g} \cdot \frac{\partial g}{\partial \psi};$$

$$f(g(z)) \approx 0;$$

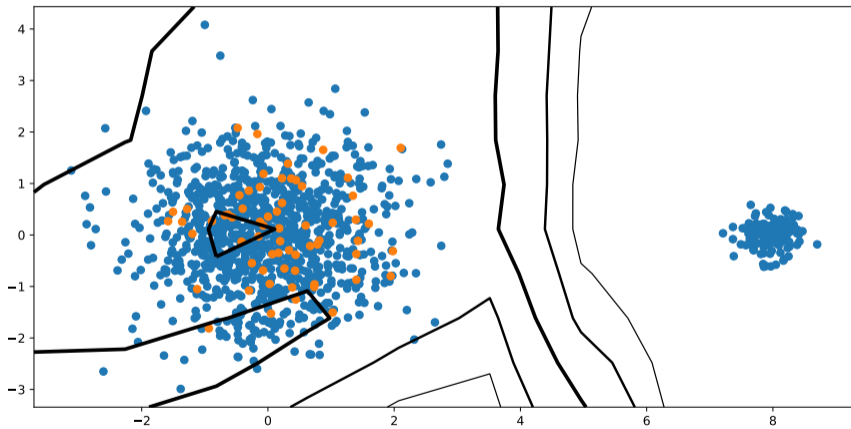
$$\frac{\partial f}{\partial g} \approx 0.$$

⇒ gradients tend to vanish on early stages.

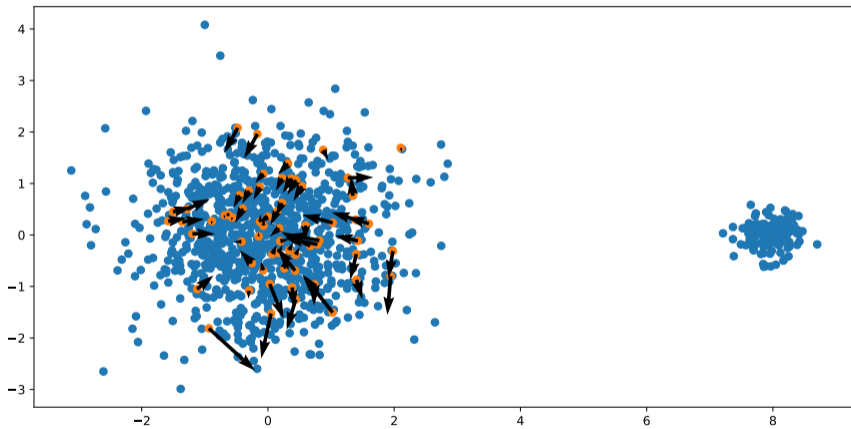
Toy classical GAN



Mode collapse



Mode collapse



GAN training tricks

Fight for the gradients

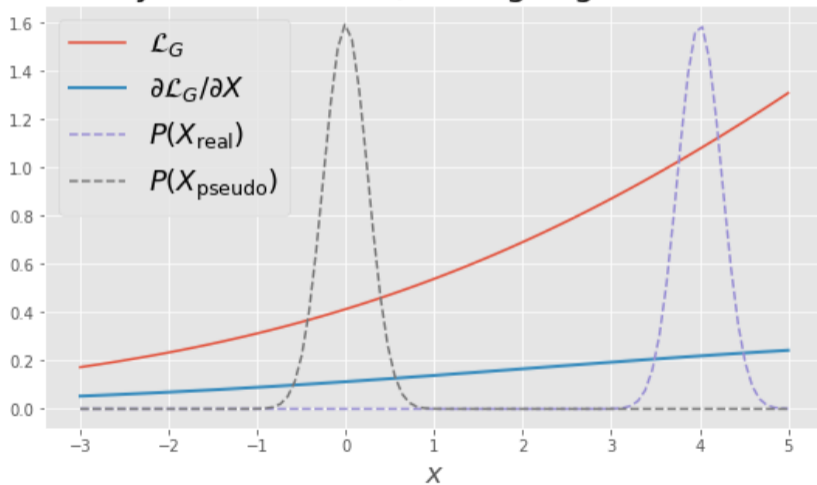
Start with heavily restricted discriminator:

- don't train discriminator fully:
 - poor-man solution;
- add noise to the samples:
 - nicely works for target on low-dimensional manifolds;
 - easy to control.
- heavy regularization:
 - might interfere with the convergence.

As learning progresses gradually relax restrictions.

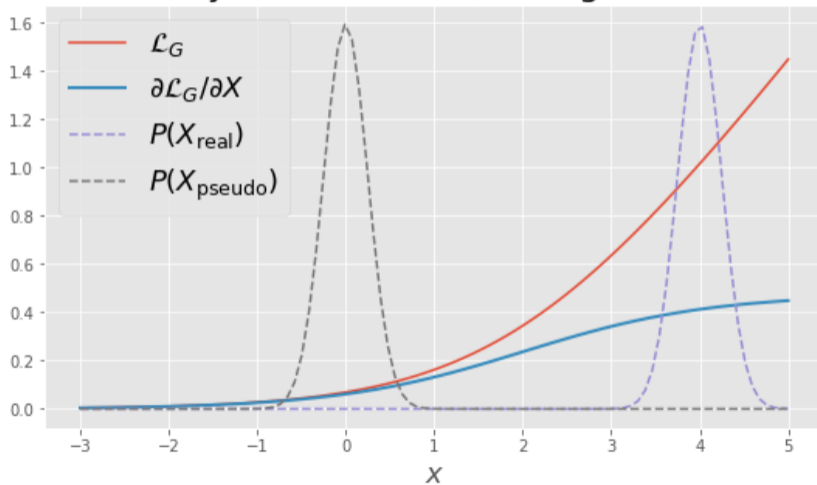
Fight for the gradients

Toy classical GAN, strong regularization



Fight for the gradients

Toy classical GAN, strong noise



Generator collapse

Often generator learns to output constant or just few values. This is a syndrome of poorly trained discriminator:

- generator aims to maximize discriminator loss;
- discriminator does not adapt quickly enough;
- generator collapses into a current maxima of discriminator.

Feature matching

Let h be some **feature**, then feature matching is an auxiliary objective:

$$\mathcal{L}_{\text{fm}} = \left\| \mathbb{E}_{X \sim \text{data}} h(X) - \mathbb{E}_Z h(g(Z)) \right\|^2$$

Alternatively, adversarial objective might be used as well:

- just add feature h to the discriminator input;
- use a separate (simple) discriminator.

Summary

Kernel Density Estimation:

- good for 1-2 dimensional problems;
- might be computationally expensive;
- explicit probability density;

Gaussian Mixture Models:

- similar to KDE;
- explicit probability density;

Variational Auto-Encoder:

- a powerful generative model;
- easy to train;

Generative Adversarial:

- a powerful generative model;
- hard to train;
- a huge number of modifications:
 - Wasserstein-GAN solves problem of vanishing gradients;
 - BiGAN, ALI add inference;
 - CycleGAN allows to learn transformation between two unpaired sets;
 - and many more.

References I

- Bengio Y. Learning deep architectures for AI. Foundations and trends® in Machine Learning. 2009 Nov 15;2(1):1-27.
- Hinton GE, Osindero S, Teh YW. A fast learning algorithm for deep belief nets. Neural computation. 2006 Jul;18(7):1527-54.
- Nair V, Hinton GE. Rectified linear units improve restricted boltzmann machines. In Proceedings of the 27th international conference on machine learning (ICML-10) 2010 (pp. 807-814).

References II

- Hinton G. A practical guide to training restricted Boltzmann machines. *Momentum*. 2010 Aug 2;9(1):926.
- Tieleman T. Training restricted Boltzmann machines using approximations to the likelihood gradient. In *Proceedings of the 25th international conference on Machine learning* 2008 Jul 5 (pp. 1064-1071). ACM.