

# 1 Deep Learning Basics

## Deep Learning: write in the front

Deep Learning is a class of machine learning methods that use **neural networks** to learn **representations** from raw data. A DL algorithm *always* consists of neural network (model) structure, training procedure (loss function) and inference procedure.

*Always* use **backpropagation**.

*Always* use **mini-batch** of data to compute gradient, and then do **gradient descent** on each dimension; the optimizer can vary, such as SGD, RMSprop or Adam.

Convention: if a function  $f$  is a neural network and its parameters are  $\theta$ , we write  $f_\theta(\cdot)$  or  $f(\cdot; \theta)$  and say  $f$  is **parameterized** by  $\theta$ . The loss function is  $\mathcal{L}$ , the dataset is  $\mathcal{D}$ .

## Deep Learning: about loss function

The loss function should *always* be **differentiable**.

We *always prefer* **tractable** loss. If the loss is not tractable, we have to estimate the gradient (using Monte-Carlo), which introduces instability and error.

What is tractable?  
 $\mathbb{E}_{x \sim \mathcal{D}}$  and  $\mathbb{E}_{z \sim \mathcal{N}(0, I)}$  are both tractable, since we can sample a mini-batch from dataset  $\mathcal{D}$  and Gaussian  $\mathcal{N}$  (or other known distribution) during training.

The **direct output** of the model.

**CAVEAT:** If the model inputs  $z$  and outputs  $x$ , then  $p_\theta(x|z)$  is tractable, but  $p_\theta(z|x)$  is often not! This is what “direct output” means.

## Optimization

Gradient Descent (GD):  $\theta \leftarrow \theta - \eta \cdot \nabla_\theta \mathcal{L}$ .

Stochastic GD (SGD) (in DL scenario): sample a small batch of data, then compute average loss as the estimation of the true  $\mathcal{L}$ .

**Momentum:**  $\theta_{t+1} \leftarrow \theta_t - \eta \cdot \nabla_\theta \mathcal{L}(\theta_t) + \beta \cdot (\theta_t - \theta_{t-1})$ , where  $\theta_t$  is the parameters at step  $t$ . **No convergence guarantee**.

**Nesterov’s Method:**  $\theta_{t+1} \leftarrow \theta_t - \eta \cdot \nabla_\theta \mathcal{L}(\theta_t + \beta \cdot (\theta_t - \theta_{t-1})) + \beta \cdot (\theta_t - \theta_{t-1})$ . **Faster convergence guarantee in smooth functions**.

**AdaGrad:**  $\theta_{t+1}[i] = \theta_t[i] - \frac{\eta}{\sqrt{G_t[i] + \epsilon}} (\nabla_\theta \mathcal{L})[i]$ ,  $G_t[i] =$

$$\sum_{s \leq t} |(\nabla_\theta \mathcal{L})[i]|^2.$$

**RMSProp:** Change  $G_t[i] = \gamma \cdot G_{t-1}[i] + (1 - \gamma)|(\nabla_\theta \mathcal{L})[i]|^2$  (“moving average”)

**Adam:** Change the gradient term  $\nabla_\theta \mathcal{L}$  in RMSProp into the momentum version

$$M_t = \delta M_{t-1} + (1 - \delta) \nabla_\theta \mathcal{L}$$

compute bias-corrected first moment and second raw moment estimate

$$\hat{G}_t = \frac{G_t}{1 - \delta^t}, \quad \hat{M}_t = \frac{M_t}{1 - \gamma^t}$$

update rule

$$\theta_{t+1}[i] = \theta_t[i] - \frac{\eta}{\sqrt{\hat{G}_t[i] + \epsilon}} \hat{M}_t[i]$$

RMSProp and Adam are common choices in modern DL.

Essentially, *AdaGrad and RMSProp are making the learning rate in each dimension different (instead of a unified constant  $\eta$ )*. For example, in AdaGrad, the learning rate of dimension  $i$  is  $\frac{\eta}{\sqrt{G_t[i] + \epsilon}}$ .

## Optimization: Practices

Why SGD but not GD? Computational affordable; avoid saddle points.

What is the benefit of RMSProp compared to AdaGrad? It can avoid vanishing learning rate.

T/F: Adam, AdaGrad are making the learning rate different and disentangled in different dimension. True.

T/F: Although Nesterov Momentum has no convergence guarantee, it is commonly used in real training. False. None of them is true.

T/F: SGD and GD have the same mean and variance in training. False.

T/F: Second-order optimization is more stable than gradient descent. False. It can diverge for non-convex functions due to negative eigenvalues.

## 2 Model Architecture

• Global pooling, getting a  $C$ -dim feature.

• A MLP layer to output the final logits.

## Important Tricks

• Dropout & early-stopping & regularization → alleviate overfitting.

• Initialization tricks (Xavier / Kaiming init for MLP; zero-init for residual blocks).

• Learning rate schedule (lr decay; lr warmup).

• Data augmentation.

• Gradient Clipping.

## Model Architecture: Practices

• Skip connection is a common technique in neural network design, e.g. ResNet, DenseNet, Transformers. True.

• When training a CNN with BatchNorm, it’s equivalent to use batch size 32 for one step and use batch size 2 for 16 steps while accumulating the gradient. False. Intuitively, all data in a batch “communicates” through the batch statistics  $\mu, \sigma^2$ .

• In evaluation, BatchNorm normalizes the input by the mean and variance of a batch. False. In evaluation (inference, sampling), we use the moving average of  $\mu, \sigma^2$ .

• When the batch size is 1, BatchNorm and LayerNorm are equivalent. False. The group separation is different.

• The input of DenseNet contains more information than ResNet. Why people still use ResNet? DenseNet makes the layer wider, memory consumption will be large. ResNet, however, does not increase the layer width.

## 3 Generative Models in Vision

### Generative Models: Ideas

**Universal Principle:** learn the data distribution  $p(X)$ . But we can construct this distribution in different ways:

• Directly model  $p(X)$ , then do sampling → **EBM**

• Model the score field  $\nabla_X \log p(X)$ , then do Langevin Dynamics sampling → **SBM**

• Model the score field  $s(X, \sigma_i)$  for each noise level, then gradually denoise in sampling → **Diffusion & SDE/ODE Models**

• Model  $p(X)$  by creating an explicit bijection from  $Z \sim \mathcal{N}(0, I)$

fact, EBM is indeed hard to scale up since calculating the loss requires inefficient Monte-Carlo sampling.

• T/F: The training of EBM involves estimating the value of normalizing constant  $Z$ . Thus, we often use Monte-Carlo method to estimate the value of  $Z$ . False. We are estimating  $\nabla_\theta \log Z$ , since we **only care about** the gradient.

• T/F: Gibbs Sampling is a special case of the MH algorithm, where the acceptance rate is always 1. True.

### 3.2 Score-Based Models (SBM) & Diffusion Models & SDE/ODE Models

#### SBM: Langevin Dynamics

We want to sample from a distribution  $p(x)$ . If we have the **score field**  $s(x) = \nabla_x \log p(x)$ , starting from a random point  $x_0$ , we iteratively sample:

$$x_{t+1} \leftarrow x_t + \epsilon \cdot s(x_t) + \sqrt{2\epsilon} \cdot z, \quad z \sim \mathcal{N}(0, I).$$

To use this sampling method, only the score field  $s(x)$  is required.

#### SBM: Basic Idea

Besides explicitly modeling the distribution  $p(X)$ , we can represent the distribution  $p(X)$  by its score function  $s(x) = \nabla \log p(x)$ . In EBM, we maximize  $\mathbb{E}_{x \sim p_{\text{data}}} \log p_\theta(x) - \log Z(\theta)$ . In SBM, we want to force  $\nabla \log p_{\text{data}}(x) = \nabla \log p_\theta(x)$ .

#### SBM: Training

To match  $\nabla \log p_{\text{data}}(x) = \nabla \log p_\theta(x)$ , the most intuitive way is minimizing the MSE:

$$\mathcal{L}(\theta) = \frac{1}{2} \mathbb{E}_{x \sim p_{\text{data}}} [|\nabla_x \log p_{\text{data}}(x) - \nabla_x \log p_\theta(x)|^2]$$

which is the **Fisher divergence** between  $p_{\text{data}}$  and  $p_\theta$ . With some math,

$$\mathcal{L}(\theta) = \mathbb{E}_{x \sim p_{\text{data}}} \left[ \frac{1}{2} |\nabla_x \log p_\theta(x)|^2 + \text{tr}(\nabla_x^2 \log p_\theta(x)) \right] + \text{const}$$

Notice: No partition function  $Z(\theta)$  anymore!

#### NF: Math

**Theorem:** Given a bijection  $f(z) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ , probability of  $x = f(z)$ :

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

If  $f = f_K \circ \dots \circ f_2 \circ f_1$ , each  $f_i$  is invertible, then

$$\log p(x) = \log p(z_0) - \sum_i \log \left| \det \frac{\partial f_i(z_{i-1})}{\partial z_i} \right|$$

Notice:  $\det \frac{\partial f(z)}{\partial z}$  requires  $O(d^3)$  computation

#### Ways to compute Jacobian

Goal: find  $f$  such that computing  $\det \frac{\partial f(z)}{\partial z}$  is cheap — **f** with triangular Jacobian

• NICE (Nonlinear Independent Components Estimation):

$$x_{1:m} = z_{1:m}, \quad x_{m+1:d} = z_{m+1:d} - \mu_\theta(z_{1:m}), \quad \det J = 1$$

• Real-NVP

$$x_{1:m} = z_{1:m}, \quad x_{m+1:d} = (z_{m+1:d} - \mu_\theta(z_{1:m})) \exp(-\alpha_\theta(z_{1:m}))$$

$$\det J = \prod_{i=m+1}^d \exp(-\alpha_\theta(z_{1:m}))$$

• GLOW (Generative Flow with Invertible 1x1 Convolutions):

$$x_{ij} = W_{ij}, \quad \log \det J = hw \log \det W$$

#### AF & IAF

Autoregressive Flow (AF):

• Forward: fast train  $x \rightarrow z$ ,

$$z_i = (x_i - \mu(x_{<i})) \odot \exp(-\alpha(x_{<i}))$$

• Reverse: slow sample  $z \rightarrow x$ ,

$$x_i = z_i \odot \exp(\alpha(x_{<i})) + \mu(x_{<i})$$

$$\sum_{s \leq t} |(\nabla_\theta \mathcal{L})[i]|^2.$$

**RMSProp:** Change  $G_t[i] = \gamma \cdot G_{t-1}[i] + (1 - \gamma)|(\nabla_\theta \mathcal{L})[i]|^2$  (“moving average”)

**Adam:** Change the gradient term  $\nabla_\theta \mathcal{L}$  in RMSProp into the momentum version

$$M_t = \delta M_{t-1} + (1 - \delta) \nabla_\theta \mathcal{L}$$

compute bias-corrected first moment and second raw moment estimate

$$\hat{G}_t = \frac{G_t}{1 - \delta^t}, \quad \hat{M}_t = \frac{M_t}{1 - \gamma^t}$$

update rule

$$\theta_{t+1}[i] = \theta_t[i] - \frac{\eta}{\sqrt{\hat{G}_t[i] + \epsilon}} \hat{M}_t[i]$$

RMSProp and Adam are common choices in modern DL.

Essentially, *AdaGrad and RMSProp are making the learning rate in each dimension different (instead of a unified constant  $\eta$ )*. For example, in AdaGrad, the learning rate of dimension  $i$  is  $\frac{\eta}{\sqrt{G_t[i] + \epsilon}}$ .

## Optimization: Practices

Why SGD but not GD? Computational affordable; avoid saddle points.

What is the benefit of RMSProp compared to AdaGrad? It can avoid vanishing learning rate.

T/F: Adam, AdaGrad are making the learning rate different and disentangled in different dimension. True.

T/F: Although Nesterov Momentum has no convergence guarantee, it is commonly used in real training. False. None of them is true.

T/F: SGD and GD have the same mean and variance in training. False.

T/F: Second-order optimization is more stable than gradient descent. False. It can diverge for non-convex functions due to negative eigenvalues.

## 2 Model Architecture

## Multi-layer Perceptron (MLP)

Definition: alternating between **linear layer** and **activation layer**.

Linear layer:  $x_{\text{output}} = Wx_{\text{input}} + b$ , where  $W$  and  $b$  are learnable.

Activation layer: a non-linear layer applied to each neuron.

• **ReLU:**  $f(x) = \max\{x, 0\}$

• **LeakyReLU:**  $f(x) = \max\{x, kx\}$  ( $0 < k < 1$ )

• **Sigmoid:**  $f(x) = \frac{1}{1 + e^{-x}}$ ,  $f'(x) = f(x)(1 - f(x))$ .

• **Tanh:**  $f(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$ ,  $f'(x) = 1 - f(x)^2$ .

## Convolutional Layer

Key insight: **invariance** (In Vision, the image is *spatial invariant*; in sequential data, there is *temporal invariance*).

The most common usage is in CV

<p>• Reconstruction loss:</p> $\mathcal{L}_{\text{recon}} = -\mathbb{E}_{\epsilon \sim \mathcal{N}(0, I), y \sim q_\phi(y x)} [\log p_\theta(x \mu(x) + \epsilon \cdot \sigma_\phi(x), y)]$ <p>Backpropagation through <math>\mathbb{E}_{y \sim q_\phi(y x)}</math> via expanding expectation</p> <p><b>VAE: Practices</b></p> <ul style="list-style-type: none"> <li>• T/F: Optimizing ELBO w.r.t. <math>p, q</math> will implicitly maximize NLL. Thus, assuming infinite power of neural networks, if we train a VAE, the margin of NLL and ELBO will finally converge to 0. False. <math>q_\phi(z x) = p_\theta(z x)</math> will never hold since <math>q_\phi(z x)</math> is restricted to be a Gaussian, but <math>p_\theta(z x)</math> is not.</li> <li>• T/F: We can still directly use VAE to generate when training with <math>\beta = 0</math>. False. The latent distribution is unknown: we don't know how to sample <math>z</math> (<math>\beta = 0</math> is standard AE).</li> <li>• T/F: If we train a VAE to reconstruct a randomly masked image, this model can do impairing. True.</li> <li>• T/F: Intuitively, large <math>\beta</math> in <math>\beta</math>-VAE makes the latent space disentangled. True.</li> <li>• T/F: The reconstruction loss of VAE is commonly the <math>\ell_2</math> loss between the input image and the predicted image. True, using the formulation we have introduced.</li> </ul> <p><b>Discrete Latent VAE: Gumbel-softmax trick</b></p> <p>Method for discrete latent, not widely used today. Want the latent distribution <math>p(z) = \text{Uniform}(1, 2, \dots, K)^L</math>. The regularization loss is naturally <math>D_{\text{KL}}(q(z x) \  p(z))</math>:</p> $\mathcal{L}_{\text{recon}} = -\mathbb{E}_{z \sim q_\phi(z x)} \log p_\theta(x z),$ <p>where <math>q_\phi(z x) = \text{softmax}(l_\phi(x))</math>. For the backpropagation of reconstruction loss, we need to reparameterize the "sampling" from a discrete distribution <math>q(z x)</math>.</p> <p><b>Discrete Latent VAE: Gumbel-softmax trick</b></p> <p><b>Theorem:</b> Define <math>Q(0, 1)</math> as the probability distribution <math>p(x) = \exp\{-x - \exp\{-x\}\}</math>. Sample i.i.d. <math>\epsilon_1, \dots, \epsilon_n \sim Q(0, 1)</math>. Given <math>x_1, \dots, x_n \in \mathbb{R}</math>, define a random variable</p> $k = \arg \max_{1 \leq i \leq n} \{x_i + \epsilon_i\}.$	<p>Then, <math>p(k) = \text{softmax}(\{x_i\}_{1 \leq i \leq n})[k], 1 \leq k \leq n</math>. Sample <math>\epsilon_1, \dots, \epsilon_K \sim Q(0, 1)</math>, using the Theorem above</p> $v = \lim_{\tau \rightarrow 0} \text{softmax} \left\{ \frac{x_i + \epsilon_i}{\tau} \right\}_{1 \leq i \leq K} = \text{one-hot}(\arg \max_i \{x_i + \epsilon_i\})$ <p>equivalent to sample from one-hot vector from softmax(<math>x_1, \dots, x_K</math>), hence</p> $\mathcal{L}_{\text{recon}} = -\mathbb{E}_{z \sim q_\phi(z x)} \log p_\theta(x z) \approx -\mathbb{E}_{\epsilon \sim Q(0, 1)} \log p_\theta \left( x   \text{softmax} \left( \frac{l_\phi(x) + \epsilon}{\tau} \right) \right) (\tau \ll 1)$ <p><b>Vector Quantized VAE (VQVAE): Motivation</b></p> <ul style="list-style-type: none"> <li>• Idea: make the prior <math>p(z)</math> discrete.</li> <li>• Motivation: <math>p(z) \sim \mathcal{N}(0, I)</math> is "unimodal"; but a categorical distribution <math>\{1, \dots, K\}^L</math> (<math>L</math> is the latent size) is naturally "multimodal".</li> <li>• Specifically, given an image <math>x</math>, the encoder should output a discrete vector <math>v = (v_1, \dots, v_L) \in \{1, 2, \dots, K\}^L</math>. Then, using this discrete vector, we want a decoder to reconstruct <math>x</math>.</li> </ul> <p><b>VQVAE: Encoder &amp; Decoder</b></p> <p>Using a normal encoder, we can encode <math>x</math> into a latent <math>z = (z_1, \dots, z_L)</math>. However, <math>z</math> is not discrete; so we need to <b>quantize</b> it.</p> <ul style="list-style-type: none"> <li>• KEY IDEA: learn a <b>codebook</b> <math>\mathbf{Z} = \{e_1, e_2, \dots, e_K\}</math>.</li> <li>• Each <math>z_i</math> is encoded into a discrete value <math>v_i \in \{1, 2, \dots, N\}</math> by finding the nearest neighbor:</li> </ul> $v_i = \arg \min_{1 \leq j \leq K} \ e_j - z_i\ .$ <p><b>Encoder:</b> <math>x \rightarrow z</math> (<math>z_e \rightarrow v</math>; <math>v \rightarrow z</math> uses a neural network encoder; <math>z \rightarrow v</math> uses <b>quantization</b> by a learned codebook. <b>Decoder:</b> in reverse order, <math>v \rightarrow z_q \rightarrow x</math>. <math>z \rightarrow x</math> part can be implemented by another NN; for <math>v \rightarrow z</math>, we can <b>look-up</b> the codebook, <math>z_i = e_{v_i}</math>.</p> <p><b>VQVAE: Loss</b></p> <p>We need to define the <b>reconstruction loss</b> and the <b>regularization loss</b>.</p> <ul style="list-style-type: none"> <li>• <b>Reconstruction:</b> just the <math>\ell_2</math> between input and output.</li> </ul>	<p>• <b>Regularization:</b> minimize the quantization error, i.e., the <math>\ell_2</math> distance between <math>z_e</math> and <math>z_q</math>. In the original paper:</p> $\mathcal{L}_{\text{reg}} = \ z_e - \text{sg}(z_q)\ ^2 + \beta \cdot \ z_q - \text{sg}(z_e)\ ^2.$ <p>Essentially, this <math>\beta</math> is to <b>separate the learning rate</b> for <math>z_e</math> and <math>z_q</math>.</p> <p><b>VQVAE: "Reparameterization"</b></p> <p><b>NOTICE:</b> taking arg min and look-up are not differentiable! We should make the whole process differentiable, so that the gradient of the reconstruction loss can backpropagate. "<b>Straight-through</b>" trick: <math>z_q \leftarrow z_e + \text{sg}(z_q - z_e)</math>. Then, the gradient will flow back to <math>z_e</math>, bypassing the quantization step.</p> <p><b>VQVAE: Sampling</b></p> <p>To sample images, we first need to sample <math>v</math>, then use the decoder to turn <math>v</math> into <math>z_q</math>, and then <math>x</math>. <b>NOTICE:</b> We do not know the latent distribution <math>p(v)</math>! However, we can learn it. For instance, regard <math>p(v) = p(v_1, \dots, v_L)</math> in an autoregressive manner, then we can use a Transformer decoder to learn</p> $p(v_t   v_{1:t-1}), 1 \leq t \leq L.$ <p><b>VQVAE: Practices</b></p> <ul style="list-style-type: none"> <li>• After training a VQ-VAE, can we sample the latent <math>z</math> from a uniform distribution and do generation? If not, what should we do? We cannot sample directly from a uniform distribution. We need the "second stage": train the distribution of <math>(z, v)</math>.</li> <li>• T/F: If the latent has high dimension (e.g. <math>10^3</math>), dictionary embedding is better than the Gumbel softmax trick. True.</li> </ul>	<p>is a good sample". The idea: use a neural classifier <math>D_\phi</math> to measure whether "an image looks like an image".</p> <p><b>GAN: Training</b></p> <p>Two core questions:</p> <ul style="list-style-type: none"> <li>• Given <math>D_\phi</math>, how to train <math>G_\theta</math>? Maximize the score given by the classifier:</li> </ul> $\theta = \arg \max_{\theta} \mathbb{E}_{z \sim \mathcal{N}(0, I)} D_\phi(G_\theta(z))$ <ul style="list-style-type: none"> <li>• Given <math>G_\theta</math>, how to train <math>D_\phi</math>? Try to classify images from dataset and generated by <math>G_\theta</math>:</li> </ul> $\phi = \arg \max_{\phi} \mathbb{E}_{x \sim p_{\text{data}}} \log D_\phi(x) + \mathbb{E}_{z \sim \mathcal{N}(0, I)} \log (1 - D_\phi(G_\theta(z)))$ <p><b>GAN: Objective</b></p> <p>The goal:</p> $\min_{\theta} \max_{\phi} \mathbb{E}_{x \sim p_{\text{data}}} \log D_\phi(x) + \mathbb{E}_{z \sim \mathcal{N}(0, I)} \log (1 - D_\phi(G_\theta(z)))$ <p>The optimal discriminator:</p> $D_\phi(x) = \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_G(x)}$ <p>Under the optimal discriminator, the goal becomes:</p> $\min_{\theta} 2 \text{JSD}(p_G, p_{\text{data}})$ <p>where JSD is the Jensen-Shannon divergence.</p> <p><b>GAN: Evaluation</b></p> <p>Since <math>p_G(x)</math> is not tractable, likelihood-based evaluation is impossible. Common metrics are based on a classifier (e.g., Inception v3 pretrained on ImageNet):</p> <ul style="list-style-type: none"> <li>• <b>IS (Inception Score):</b> <math>IS = \exp(\mathbb{E}_{x \sim G} [KL(f(y x)    p_f(y))])</math>, higher is better.</li> <li>• <b>FID (Fréchet Inception Distance):</b> <math>FID = \ \mu_p - \mu_G\ ^2 + \text{tr}(\Sigma_p + \Sigma_G - 2\sqrt{\Sigma_p \Sigma_G})</math>, lower is better.</li> </ul> <p><b>GAN: Problems</b></p> <p>Two main issues:</p>
<p>• Mode collapse: generating a few samples can cheat the loss.</p> <p>• Training instability: discriminator and generator may keep oscillating.</p> <p>Solutions include DCGAN, improved training techniques, WGAN, and more extensions.</p> <p><b>DCGAN</b></p> <p>Tricks for stability:</p> <ul style="list-style-type: none"> <li>• Fully convolutional network (no pooling / MLP).</li> <li>• Batch normalization.</li> <li>• Leaky ReLU activation.</li> <li>• Small learning rate and momentum.</li> </ul> <p><b>Improved GAN</b></p> <p>Additional tricks:</p> <ul style="list-style-type: none"> <li>• Match features of the final activation layer in <math>D_\phi</math>.</li> <li>• <math>D_\phi</math> classifies a sample via batch information.</li> <li>• Historical averaging: <math>\mathcal{L}(\theta) = \ \theta - \frac{1}{T} \sum_t \theta_t\ ^2</math>.</li> <li>• Change positive label to 0.9.</li> <li>• Use a fixed, pre-selected batch for normalization.</li> </ul> <p><b>WGAN</b></p> <p>Change the distance metric from JSD to Wasserstein distance (more smooth):</p> $W(p_G, p_{\text{data}}) = \sup_{\ f\ _L \leq 1} \mathbb{E}_{x \sim p_{\text{data}}} f(x) - \mathbb{E}_{x \sim p_G} f(x)$ <p>The goal:</p> $\min_{\theta} \max_{\phi, D_\phi} \mathbb{E}_{x \sim p_{\text{data}}} D_\phi(x) - \mathbb{E}_{z \sim \mathcal{N}(0, I)} D_\phi(G_\theta(z))$ <p>Improved WGAN adds a gradient penalty.</p> <p><b>GAN: A Bit More</b></p> <p>Improvements:</p> <ul style="list-style-type: none"> <li>• BigGAN: Large model + tricks.</li> <li>• GigaGAN: Larger model + tricks.</li> <li>• R3GAN: New loss and architectures.</li> </ul> <p>Extensions:</p> <ul style="list-style-type: none"> <li>• Semi-Supervised Learning: assign fake label for data without label.</li> </ul>	<p>• Representation Learning: similar distribution for <math>(z, G(z))</math> and <math>(E(x), x)</math>.</p> <p>• Style transfer: sample <math>z</math> from another dataset (see CycleGAN for unpaired data).</p> <p><b>GAN: Practices</b></p> <ul style="list-style-type: none"> <li>• T/F: In GAN training, we should first train the discriminator to make it converge, then train the generator. False. We should train generator and discriminator jointly.</li> <li>• T/F: We want to fit a multimodal distribution by normal distribution and using KL loss function, but this results in collapsing to one mode. Change the loss function to JSD can fix this problem. False. JSD also suffers from mode collapse issue, that's why we need WGAN.</li> </ul> <p><b>4 Sequential Models</b></p> <p><b>4.1 Recurrent Networks (RNN &amp; LSTM) and Language Models</b></p> <p><b>Language Models</b></p> <p><b>seq2seq model:</b> input: sequence <math>x = (x_1, x_2, \dots)</math>; output: sequence <math>y = (y_1, y_2, \dots)</math>.</p> <p>In a language model, sequence <math>x, y</math> are natural language sentences. Each <math>x_i</math> or <math>y_i</math> is a <b>token</b> but not a word.</p> <p>There is often a finite set of possible tokens (of size <math>N</math>). Then, each <math>x_i, y_i \in \{1, 2, \dots, N\}</math>.</p> <p>We want to learn a <math>N</math>-dimensional probability distribution for each position <math>y_i</math>, say <math>p_\theta(y_i)</math>. The loss function is the NLL loss</p> $\mathcal{L} = \sum_i -\log p_\theta(y_i).$ <p>A sequential <b>generative model</b> is commonly modeled in a <b>autoregressive</b> manner:</p> $p(y_{1:L} x) = p_\theta(y_1 x) \cdot p_\theta(y_2 x, y_1) \cdots p_\theta(y_L x, y_{1:L-1}).$ <p>This is called "<b>next-token prediction</b>". In this way, the NLL loss is equivalent to <math>-\log p(y x)</math>.</p> <p>Model structures:</p> <ul style="list-style-type: none"> <li>• RNN;</li> </ul>	<p>• LSTM;</p> <p>• Transformer.</p> <p>Other important concepts:</p> <ul style="list-style-type: none"> <li>• Word embedding;</li> <li>• Beam search.</li> </ul> <p><b>Recurrent NN (RNN): Idea</b></p> <p>Idea: maintain a "latent (hidden) state" <math>h_i</math>, combining information of <math>x_1, \dots, x_i</math>. <i>Intuitively, <math>h</math> represents what your brain is thinking when you hear the first <math>i</math> words.</i></p> <ul style="list-style-type: none"> <li>• <math>h_{i+1}</math> is inherited from <math>h_i</math>, can be influenced by <math>x_{i+1}</math>;</li> <li>• <math>y_i</math> is directly computed using <math>h_i</math>.</li> </ul> <p><b>RNN: Formulation</b></p> <ul style="list-style-type: none"> <li>• Hidden state transition: <math>h_{i+1} = f_1(W_{hx}x_{i+1} + W_{hh}h_i + b_h)</math>;</li> <li>• Output: <math>y_i = f_2(W_{yh}h_i + b_y)</math>.</li> <li>• <math>f_1, f_2</math> are activations. <math>W, b</math> are learnable parameters.</li> </ul> <p><b>RNN: BPTT</b></p> <p>Note that the parameters are used for <math>L</math> times in a sequence of length <math>L</math>. So, when we backpropagate the gradient, the gradient will accumulate <math>L</math> times. This is called "<b>backpropagation through time (BPTT)</b>".</p> <p>PROBLEM: gradient vanishing &amp; exploding (toy example: all layers are</p>	