## Bidirectional RNN

## Long short-term memory

- Motivation
  - vanishing or exploding gradient for larger time step
  - need to "preserve" the earlier hidden node activations for prediction at current time $t$ (Hochreiter and Schmidhuber, 1997)
    * remove and add information to the cell
    * gated mechanism

## Gated recurrent unit

Two gates and one combined state

Reset gate
$$r_t = \sigma(W_{xr}x_t + W_{hr}h_{t-1} + b_r)$$

Update gate
$$z_t = \sigma(W_{xz}x_t + W_{hz}h_{t-1} + b_z)$$

Output state
$$\tilde{h}_t = \tanh(W_x x_t + W_h(r_t \odot h_{t-1}) + b)$$

$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \tilde{h}_t$$

## Stacked recurrent neural network (Joulin and Mikolov, 2015)

- prediction model for next symbol in a stream of discrete data
- a structured memory inspired by pushdown automaton using a stack
- learn how to operate this memory through optimization tools

(a) neural network extended with pushdown stack and a controlling mechanism that learns what action (among PUSH, POP and NO-OP)
(b) the same model extended with a doubly-linked list with actions INSERT, LEFT, RIGHT and NO-OP

## Recurrent neural network

- Unfolding in RNN
- Gradient vanishing and gradient exploding

Memory cell consists of input gate $i_t$, forget gate $f_t$, cell $c_t$, output gate $o_t$ and hidden state $h_t$ which are operated as

$$i_t = \sigma(W_{xi}x_t + W_{hi}h_{t-1} + W_{ci}c_{t-1} + b_i)$$
$$f_t = \sigma(W_{xf}x_t + W_{hf}h_{t-1} + W_{cf}c_{t-1} + b_f)$$
$$c_t = f_t c_{t-1} + i_t \tanh(W_{xc}x_t + W_{hc}h_{t-1} + b_c)$$
$$o_t = \sigma(W_{xo}x_t + W_{ho}h_{t-1} + W_{co}c_t + b_o)$$
$$h_t = o_t \tanh(c_t)$$

## LSTM versus GRU

- Long short-term memory
  - input gate, forget gate and output gate
  - possess an internal memory
  - apply a nonlinearity transformation
- Gated recurrent unit
  - update gate and reset gate
  - no internal memory
  - no nonlinearity on the output
- RNNs with tanh, LSTM and GRU units (Chung et al., 2)
  - GRU and LSTM outperformed tanh
  - GRU and LSTM converged faster than tanh
  - choosing GRU or LSTM depends on the dataset

## Neural Turing machine versus memory network

- Most machine learning models lack an easy way to
  - read and write to part of a long-term memory component
  - combine this seamlessly with inference
- Neural Turing machine (Graves et al., 2014)
  - learns to read from and write to memory cells without explicit supervision
  - allows end-to-end training via content-based soft attention
  - emulates algorithmic mechanism in a way that allows gradient-based optimization
- End-to-end memory network (Sukhbaatar et al., 2015)
  - includes memory cells that can be accessed via an addressing mechanism
  - combines learning strategies for inference with a memory component that can be read and written to

## Sequence to sequence learning

- Traditional DNN was sensibly encoded with vectors with a fixed dimensionality
- Many important problems are best expressed with sequences whose lengths are unknown a priori
- An input sequence "ABC" is encoded and decoded to produce "WXYZ" as the output sequence (Sutskever et al., 2014)
- LSTM architecture is applied to deal with this problem

## Sequence learning

- RNN can not deal with sequential learning with input and output sequences in different lengths
- Sequence to sequence learning is performed by
  - first, map the input sequence to a fixed-sized vector using on RNN
  - second, map the vector to the target sequence using another RNN
- LSTM is used to estimate $p(y_1, \dots, y_{T'}|x_1, \dots, x_T)$ where $\{x_1, \dots, x_T\}$ is an input sequence and $\{y_1, \dots, y_{T'}\}$ is its output sequence whose length $T'$ may differ from $T$
- LSTM language model is calculated by
$$p(y_1, \dots, y_{T'}|x_1, \dots, x_T) = \prod_{t=1}^{T'} p(y_t|v, y_1, \dots, y_{t-1})$$
- LSTM computes this probability by obtaining the fixed dimensional $v$ of $\{x_1, \dots, x_T\}$ given by the last hidden state of LSTM

## CTC decoding

- Expected transcription loss is minimized
$$\mathcal{L}(x) = \sum_y \Pr(y|x)\mathcal{L}(x, y)$$
$$= \sum_y \sum_{a \in B^{-1}(y)} \Pr(a|x)\mathcal{L}(x, y)$$
$$= \sum_a \Pr(a|x)\mathcal{L}(x, B(a))$$
- RNN transducers perform decoding with beam search to yield an $n$-best list of candidate transcriptions
- CTC fits naturally into the existing neural network classifiers
- CTC is similar to HMM, but differs slightly in interpretation of network outputs
- CTC outperformed HMM and HMM-RNN (Graves et al., 2006)

## Machine translation

(Sutskever et al., 2014) propose a sequence to sequence model
- compress all the information into a fixed length vector
- performance decreases as the input sentence increases

LSTM: decoder
LSTM: encoder

## Evaluation for text generation

BLEU (bilingual evaluation understudy) (Papineni et al., 2002)
- evaluate the quality of the generated text
- calculate the corresponding scores
$$BLEU\text{-}N(C, S) = b(C, S)\exp\left(\sum_{n=1}^{N} w_n \log CP_n(C, S)\right)$$

n-gram precision compares a candidate translation against multiple reference translations
$$CP_n(C, S) = \frac{\sum_i \sum_k \min(h_k(c_i), \max_{j \in m}(h_k(s_{ij})))}{\sum_i \sum_k h_k(c_i)}$$

Brevity penalty (BP) forces a high-scoring candidate translation must match the reference translations in length
$$b(C, S) = \begin{cases} 1, & \text{if } l_c > l_r \\ e^{(1-l_r/l_c)}, & \text{if } l_c \le l_r \end{cases}$$

## Image caption generator

- Describe the content of an image is a challenging task
  - capture the objects in images
  - express relations between objects
- (Vinyals et al., 2015) presents an end-to-end system for the problem
  - CNN encoder
  - LSTM decoder

## Attention network

$$\alpha_i = \text{Attend}(s_{i-1}, \alpha_{i-1}, h) \qquad g_i = \sum_{j=1}^{L} \alpha_{i,j} h_j$$
$$y_i \sim \text{Generate}(s_{i-1}, g_i) \qquad s_i = \text{Recurrency}(s_{i-1}, g_i, y_i)$$

## Content-based attention

- $\alpha_i = \text{Attend}(s_{i-1}, \alpha_{i-1}, h) \to \alpha_i = \text{Attend}(s_{i-1}, h)$
- estimation of attention weight
$$e_{i,j} = \text{Score}(s_{i-1}, h_j), \qquad \alpha_{i,j} = \frac{\exp(e_{i,j})}{\sum_{j=1}^{L}\exp(e_{i,j})}$$
- limitation: identical or very similar elements of $h$ are scored equally regardless of their position in the sequence

## Location-based attention

- $\alpha_i = \text{Attend}(s_{i-1}, \alpha_{i-1}, h) \to \alpha_i = \text{Attend}(s_{i-1}, \alpha_{i-1})$
- in speech recognition, it aims to predict the distance between consequent phonemes using $s_{i-1}$ only
- hybrid attention mechanism is fitted to speech recognition
- Score function in content-based attention was calculated by
$$e_{i,j} = w^\top \tanh(Ws_{i-1} + Vh_j + b)$$
- Content-based attention can be location-aware through alignment produced at the previous step
  - extract $k$ vectors $f_{i,j} \in \mathbb{R}^k$ for every position $j$ of the previous alignment $\alpha_{i-1}$ by convolving it with a matrix $F \in \mathbb{R}^{k \times r}$
$$f_i = F * \alpha_{i-1}$$
  - scoring mechanism is updated by using
$$e_{i,j} = w^\top \tanh(Ws_{i-1} + Vh_j + Uf_{i,j} + b)$$

## Machine translation with attention

- (Bahdanau et al., 2015) introduces attention mechanism into sequence to sequence model
  - alignment model
  - translation model
$$c_i = \sum_{j=1}^{T_x} \alpha_{ij}h_j$$
- Compute attention weight
$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x}\exp(e_{ik})}$$
where $e_{ij} = g(s_{i-1}, h_j)$

LSTM: decoder
LSTM: encoder

## Image caption with attention

CNN: encoder
LSTM: decoder

## WaveNet

- Causal convolution
  - cannot depend on any of the future timesteps
  - shifting the output of a normal convolution by a few timesteps
  - CNN is faster than RNN
- Dilated convolution
  - filter is applied over an area larger than its length by skipping input values with a certain step
  - similar to pooling or strided convolutions, but the output has the same size as the input
  - dilation 1 yields the standard convolution
  - receptive field to grow exponentially with depth

## Convolutional encoder

- Encoder consists of two stacked convolutional networks (Gehring et al., 2017)
  - $CNN_a$ produces the key vector $z_j$
  $$z_j = CNN_a(e_j)$$
  - $CNN_c$ produces the value vector $v_j$
  $$v_j = CNN_c(e_j)$$
- Conditional input $c_i$ to the decoder is
$$a_i = \text{Attention}(z_j, s_j)$$
$$c_i = \sum_{j=1}^{T} a_{ij}v_j$$
- Gated linear unit (Dauphin et al., 2017) is calculated via convolution operation $*$ for hidden layers $h_0, \dots, h_L$ as
$$h_l(E) = (E * W + b) \otimes \sigma(E * V + c)$$
  - LSTM style with no forget and input gates required
  - only possess output gate in which information to be propagated
- Gated tanh units (GTU)
  - mimic gated mechanism
$$h_l(E) = \tanh(E * W + b) \otimes \sigma(E * V + c)$$
- Residual and skip connections
  - residual connection - vertical
  - skip connection - horizontal

## Autoencoder vs variational autoencoder

- Autoencoder
  - learn a compressed representation of input automatically
  - compress the input and decompress it back to match the original input
- Variational autoencoder (VAE)
  - learn the parameters of a probability distribution representing the data
  - sample from the distribution and generate new input data samples

## Mean field variational inference

- Variational inference aims to find a variational distribution $q(z|x)$ that optimally close to the original true posterior $p(z|x)$
- $q(z|x) = \prod_{n=1}^{N} q(z_n|x_n)$ is assumed
- $p(z)$ is standard normal distribution $\mathcal{N}(0, I)$
- Model parameters $\Theta = \{\theta, \phi\}$ are learned by maximizing $\mathcal{L}_\Theta$
$$\log p(x) \ge \mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x|z)] - KL(q_\phi(z|x)||p(z))$$
$$= \mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x, z) - \log q_\phi(z|x)]$$
$$= \mathbb{E}_{q_\phi(z|x)}[f_\Theta(x, z)] \triangleq \mathcal{L}_\Theta$$
- Encoder $q_\phi(z|x)$ and decoder $p_\theta(x|z)$ are represented by Gaussians where means and variances are calculated through multilayer

## Stochastic gradient variational Bayes

Objective:
Gradient:
$$\mathcal{L}_\Theta = \mathbb{E}_{q_\phi(z|x)}[f_\Theta(x, z)]$$

Step1: sample $\epsilon^{(l)}$ from $\mathcal{N}(0, I)$
Step2: $z^{(l)} = \mu_x + \sigma_x \odot \epsilon^{(l)}$
Step3: $\mathcal{L}_\Theta \simeq f_\Theta(x, z^{(l)})$
Step4: $\nabla_\Theta \mathcal{L}_\Theta = \nabla_\Theta f_\Theta(x, z^{(l)})$

- Reduce the variance caused by directly sampling z

## Lower bound & KL divergence

- Introduce a variational distribution $q(z)$ and adopt the Jensen's inequality for convex function $-\log(\cdot)$ to obtain
$$\log p(x|\theta) = \log \sum_z \frac{p(x, z|\theta)}{q(z)}q(z) = \log \mathbb{E}_q\left[\frac{p(x, z|\theta)}{q(z)}\right]$$
$$\ge \mathbb{E}_q\left[\log \frac{p(x, z|\theta)}{q(z)}\right] \triangleq \mathcal{L}(q, \theta)$$
$$\sum_z q(z)\log p(x|\theta) \quad \mathcal{L}(q, \theta) - \sum_z q(z)\log\left\{\frac{p(x|\theta)}{q(z)}\right\} = KL(q||p)$$

Evidence Decomposition
$$\log p(x|\theta) = KL(q||p) + \mathcal{L}(q, \theta)$$

## Auto-encoding variational Bayes algorithm

Initialize parameters $\theta, \phi$
repeat
$X^M \leftarrow$ random minibatch of $M$ data points
$\epsilon \leftarrow$ random samples from noise distribution $p(\epsilon)$
$g \leftarrow \nabla_{\theta,\phi}\mathcal{L}_{\theta,\phi}(X^M, \epsilon)$ gradients of minibatch estimator
$\theta, \phi \leftarrow$ updating parameters using gradients g
until convergence of parameters $(\theta, \phi)$
return $\theta, \phi$

## Maximum likelihood

$$KL(q||p) = -\mathbb{E}_q[\log p(z|x, \theta)] - \mathbb{H}_q[z]$$
$$\mathcal{L}(q, \theta) = \mathbb{E}_q[\log p(x, z|\theta)] + \mathbb{H}_q[z]$$
- Maximizing $p(x|\theta)$ is equivalent to first setting $KL(q||p) = 0$ or approximating (E-step)
$$q(z) = p(z|x, \theta^{old})$$
then maximizing the resulting lower bound (M-step)
$$\mathcal{L}(q, \theta) \triangleq Q(\theta, \theta^{old}) + const$$
where $Q(\theta, \theta^{old}) \triangleq \mathbb{E}_q[\log p(x, z|\theta)|x, \theta^{old}]$ is concave

## Expectation-maximization algorithm

- Likelihood function for observations $x$ in latent variable model with latent variable $z$
$$p(x|\theta) = \sum_z p(x, z|\theta)$$
- Expectation (E) step: calculate an auxiliary function
$$Q(\theta, \theta^{old}) = \mathbb{E}_z[\log p(x, z|\theta)|x, \theta^{old}]$$
- Maximization (M) step: find a new estimate $\theta^{new}$ via
$$\theta^{new} = \arg\max_\theta Q(\theta, \theta^{old})$$
- EM algorithm (Dempster et al., 1977) for ML can be extended for MAP

## Why approximate inference?

- There are a number of latent variables in model-based methods
  - semantic topics in topic model
  - hidden units in layered network
- Latent variables are coupled in their posteriors
- Posterior distribution of multiple latent variables should be factorizable to find analytical solution to inference algorithm
- Exact inference does not exist in a complicated system
- Evolution of inference algorithms
  - maximum likelihood
  - maximum a posteriori
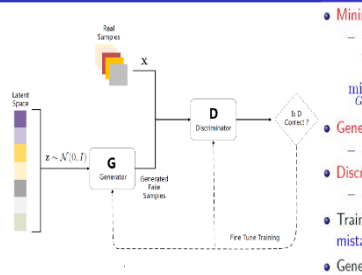  - variational inference
  - Gibbs sampling

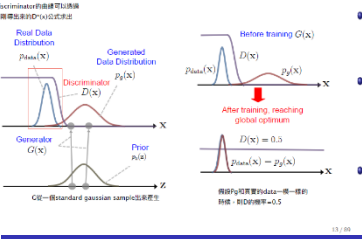## VAE versus CVAE

VAE只能生成一個training data的像而已，無法控制生成什麼

- **Variational** autoencoder
  - generation of data is **not controllable**
  - **labels** of data are **unnecessary**
  - **lower bound** and objective function $\mathcal{L}(\mathbf{x}, \theta, \phi)$

$$\log p_\theta(\mathbf{x}) \geq -KL[q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z})] + \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})] \triangleq \mathcal{L}(\mathbf{x}, \theta, \phi)$$

- **Conditional** variational autoencoder    conditional的label
  - **generation of data is controllable**    透過conditional likelihood控制生成data
  - **labels** of data are **required**
  - **lower bound** and objective function $\mathcal{L}(\mathbf{x}, \mathbf{y}, \theta, \phi)$

$$\log p_\theta(\mathbf{y}|\mathbf{x}) \geq -KL[q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{y})||p(\mathbf{z})] + \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{y}|\mathbf{x}, \mathbf{z})]$$
$$\triangleq \mathcal{L}(\mathbf{x}, \mathbf{y}, \theta, \phi)$$



CVAE    VAE    AAE

## Adversarial autoencoder

- Discriminator is trained to **distinguish** samples of **latent distribution** $q(\mathbf{z})$ from samples of **prior distribution** $p(\mathbf{z})$



- **Objective function**
  - $\mathcal{L}_{\theta_{dis}} = -\mathbb{E}_{\mathbf{z}\sim p(\mathbf{z})}[\log(D(\mathbf{z}))] - \mathbb{E}_{\mathbf{x}\sim p_{enc}}[\log(1 - D(\text{Enc}(\mathbf{x})))]$
  - $\mathcal{L}_{\theta_{enc}} = \mathbb{E}_{\mathbf{x}\sim p_{enc}}[\log(1 - D(\text{Enc}(\mathbf{x})))]$
  - $\mathcal{L}_{enc,dec} = \mathbb{E}_{\mathbf{x}\sim p_{enc}}[||\mathbf{x} - \text{Dec}(\text{Enc}(\mathbf{x}))||^2]$

**Algorithm**

Initialize parameters $\theta_{enc}, \theta_{dis}, \theta_{dec}$
repeat
  $\mathbf{X} \leftarrow$ random minibatch of $M$ data points
  $\mathbf{z} \leftarrow$ sample from prior $p(\mathbf{z})$
  $g_{\theta_{enc,dec}} \leftarrow -\nabla_{\theta_{enc,dec}}\mathcal{L}_{\theta_{enc,dec}}(\mathbf{X})$ gradients of updating encoder and decoder to reconstruct data
  $g_{\theta_{dis}} \leftarrow -\nabla_{\theta_{dis}}\mathcal{L}_{\theta_{dis}}(\mathbf{X},\mathbf{z})$ gradient of updating discriminator to distinguishes real data from fake data
  $g_{\theta_{enc}} \leftarrow -\nabla_{\theta_{enc}}\mathcal{L}_{\theta_{enc}}(\mathbf{X})$ gradients of updating encoder to confuse discriminator
  $\theta_{enc}, \theta_{dec}, \theta_{dis} \leftarrow$ updating parameters using $g_{\theta_{enc,dec}}, g_{\theta_{dis}}, g_{\theta_{enc}}$
until convergence

## Generative model

- **Conventional generative model**
  - supposes training examples are obtained by $\mathbf{x} \sim p_{data}(\mathbf{x})$
  - builds a model that can draw samples $\mathbf{x} \sim p_{model}(\mathbf{x})$
  - assumes $p_{model}(\mathbf{x}) = p_{data}(\mathbf{x})$
  - **maximum likelihood** estimate of model parameters $\theta$ is obtained by

$$\theta_{ML} = \arg\max_\theta \mathbb{E}_{\mathbf{x}\sim p_{data}}[\log p_{model}(\mathbf{x}|\theta)]$$

- **Generative adversarial network** (GAN) (Goodfellow et al., 2014)
  - general idea to learn a **sampling mechanism**
  - use **latent code**
  - **asymptotically consistent**
  - often regarded as producing the **best samples**
  - no Markov chains involved so no issue of **mixing**

## Taxonomy of generative models



## Step 1: finding optimal discriminator

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

把期望值寫成積分
$$= \int_\mathbf{x} p_{data}(\mathbf{x})\log D(\mathbf{x})d\mathbf{x} + \int_\mathbf{z} p_\mathbf{z}(\mathbf{z})\log(1 - D(G(\mathbf{z})))d\mathbf{z}$$

把z用x換掉
$$= \int_\mathbf{x} \{p_{data}(\mathbf{x})\log(D(\mathbf{x})) + p_g(\mathbf{x})\log(1 - D(\mathbf{x}))\}d\mathbf{x}$$
Pg(x代表generator的分佈)
第二項的D(x的)x是由P_g(x)產生

- For the fixed $G$, the optimal discriminator $D$, maximizing $V(D,G)$, is obtained by

做微分，令等於0
可得到D(x)的最佳解
$$D_G^*(\mathbf{x}) = \frac{p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})}$$

## Step 2: finding optimal generator

$$C(G) = \max_D V(D,G)$$
$$= \max_D \int_\mathbf{x} p_{data}(\mathbf{x})\log(D(\mathbf{x})) + p_g(\mathbf{x})\log(1 - D(\mathbf{x}))d\mathbf{x}$$
$$= \int_\mathbf{x} p_{data}(\mathbf{x})\log(D_G^*(\mathbf{x})) + p_g(\mathbf{x})\log(1 - D_G^*(\mathbf{x}))d\mathbf{x}$$
$$= \int_\mathbf{x} p_{data}(\mathbf{x})\log\left(\frac{p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})}\right) + p_g(\mathbf{x})\log\left(\frac{p_g(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})}\right)d\mathbf{x}$$   trick 讓分子多些log2 最後在消回去
$$= \int_\mathbf{x} p_{data}(\mathbf{x})\log\left(\frac{2p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})}\right) + p_g(\mathbf{x})\log\left(\frac{2p_g(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})}\right)d\mathbf{x} - \log 4$$
$$= KL\left(p_{data}(\mathbf{x})\left\|\frac{p_{data}(\mathbf{x}) + p_g(\mathbf{x})}{2}\right.\right) + KL\left(p_g(\mathbf{x})\left\|\frac{p_{data}(\mathbf{x}) + p_g(\mathbf{x})}{2}\right.\right) - \log 4$$   套用對稱KL divergence的公式
$$= 2\,JS(p_{data}(\mathbf{x})||p_g(\mathbf{x})) - \log 4$$   當兩個都是最佳解就是Pdata=Pg 也就是那一個距離的0 此時C的最小值=-log4

- **Global minimum** of $C(G)$ is achieved if and only if $p_{data}(\mathbf{x}) = p_g(\mathbf{x})$
- $C(G)$ achieves the value $-\log 4$

## Generative adversarial nets

### Minimax game



- **Minimax two-player game** over $G$ and $D$
  - two models trainable simultaneously through back-propagation over value function $V(D,G)$ using neural network

$$\min_G \max_D \left\{\mathbb{E}_{\mathbf{x}\sim p_{data}(\mathbf{x})}[\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z}\sim p_\mathbf{z}(\mathbf{z})}[\log(1 - D(G(\mathbf{z})))]\right\}$$

- **Generative** model $G$
  - captures the **data distribution**
- **Discriminative** model $D$
  - classifies a sample coming from the **training data** rather than $G$
- Train a model $G$ via **maximizing** the probability of $D$ making a mistake
- Generator $G$ **minimizes** the **classification accuracy** of $D$

### Algorithm 1 SGD training for generative adversarial net

**for** number of training iterations **do**
  **for** $k$ steps **do**
    • sample minibatch of $m$ noise samples $\{\mathbf{z}^{(1)}, \ldots, \mathbf{z}^{(m)}\}$ from prior $p_g(\mathbf{z})$
    • sample minibatch of $m$ examples $\{\mathbf{x}^{(1)}, \ldots, \mathbf{x}^{(m)}\}$ from data generating distribution $p_{data}(\mathbf{x})$
    • 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 $\{\mathbf{z}^{(1)}, \ldots, \mathbf{z}^{(m)}\}$ from prior $p_\mathbf{z}(\mathbf{z})$
  • 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**
Gradient-based updates with momentum can be used

## Global optimum at $p_{data} = p_g$



## Gradient vanish

- GAN is difficult to optimize and training is **unstable**
- Generator and discriminator must be carefully maintained for converging
  - balance the capacity between $G$ and $D$
- **Discriminator easily takes samples of generator as fake**
  - $D(G(\mathbf{z})) = 0$
  - $\nabla_{\theta_g}\frac{1}{m}\sum_{i=1}^m \log(1 - D(G(\mathbf{z}^{(i)}))) = 0$
  - gradient vanish in $G$
- To deal with it, **reformulate the optimization objective for $G$ as a minimization of**

$$\mathbb{E}_{\mathbf{z}\sim p_\mathbf{z}(\mathbf{z})}[-\log D(G(\mathbf{z}))]$$

## Mode collapse

- Lack of **diversity** in generated samples
  - $G$ learns to map different inputs $\mathbf{z}$ to the same output point $G(\mathbf{z})$
  - $G$ is likely to generate sample that $D$ believes to be real rather than fake



## Conditional GAN

- **Original GAN objective function**

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

- **Conditional GAN objective function**

$$V(D,G) = \mathbb{E}_{\mathbf{x}\sim p_{data}(\mathbf{x})}[\log D(\mathbf{x}|\mathbf{y})] + \mathbb{E}_{\mathbf{z}\sim p(\mathbf{z})}[\log(1 - D(G(\mathbf{z}|\mathbf{y})))]$$

- $G$ and $D$ are conditioned on some **extra information** $y$
- Combine prior input noise $p(z)$ and $y$ in joint hidden representa

## Information GAN

- **InfoGAN objective function** (Chen et al., 2016)

$$\min_G \max_D V(D,G) - \lambda I(\mathbf{c}; G(\mathbf{z}, \mathbf{c}))$$

- $I(\mathbf{c}; G(\mathbf{z}, \mathbf{c}))$ is the **mutual information** between **generated samples** and the **latent code**

$$I(\mathbf{c}, G(\mathbf{z}, \mathbf{c}))$$
$$= H(\mathbf{c} - H(\mathbf{c}|G(\mathbf{z}, \mathbf{c})))$$
$$= \mathbb{E}_{\mathbf{x}\sim G(\mathbf{z},\mathbf{c})}[\mathbb{E}_{\mathbf{c}'\sim p(\mathbf{c}|\mathbf{x})}[\log p(\mathbf{c}'|\mathbf{x})]] + H(\mathbf{c})$$
$$= \mathbb{E}_{\mathbf{x}\sim G(\mathbf{z},\mathbf{c})}[D_{KL}(p(.|\mathbf{x})||q(.|\mathbf{x})) + \mathbb{E}_{\mathbf{c}'\sim p(\mathbf{c}|\mathbf{x})}[\log q(\mathbf{c}'|\mathbf{x})]] + H(\mathbf{c})$$
KL-divergence是非負的  $\geq 0$
$$\geq \mathbb{E}_{\mathbf{x}\sim G(\mathbf{z},\mathbf{c})}[\mathbb{E}_{\mathbf{c}'\sim p(\mathbf{c}|\mathbf{x})}[\log q(\mathbf{c}'|\mathbf{x})]] + H(\mathbf{c})$$

## f-divergence GAN

- **f-GAN** (Nowozin et al., 2016)
  - generalize Jensen-Shannon divergence to $f$-divergence
- **$f$-divergence is defined**

$$\mathcal{D}_f(P||Q) = \int_\mathbf{x} q(\mathbf{x})f\left(\frac{p(\mathbf{x})}{q(\mathbf{x})}\right)d\mathbf{x}$$

  - two distributions $P$ and $Q$
  - absolutely continuous density functions $p(\mathbf{x})$ and $q(\mathbf{x})$
  - $f$ is **convex and lower-semicontinuous**
  - $p(\mathbf{x}) = q(\mathbf{x})$ and $f(1) = 0 \Rightarrow \mathcal{D}_f(P||Q) = 0$

| Name | $\mathcal{D}_f(P\|Q)$ |
|---|---|
| Kullback-Leibler | $\int p(\mathbf{x})\log\frac{p(\mathbf{x})}{q(\mathbf{x})}d\mathbf{x}$ |
| Reverse KL | $\int q(\mathbf{x})\log\frac{q(\mathbf{x})}{p(\mathbf{x})}d\mathbf{x}$ |
| Pearson $\chi^2$ | $\int\frac{(q(\mathbf{x})-p(\mathbf{x}))^2}{p(\mathbf{x})}$ |
| Squared Hellinger | $\int(\sqrt{p(\mathbf{x})} - \sqrt{q(\mathbf{x})})^2 d\mathbf{x}$ |
| Jensen-Shannon | $\frac{1}{2}\int p(\mathbf{x})\log\frac{2p(\mathbf{x})}{p(\mathbf{x})+q(\mathbf{x})} + q(\mathbf{x})\log\frac{2q(\mathbf{x})}{p(\mathbf{x})+q(\mathbf{x})}d\mathbf{x}$ |
| GAN | $\int p(\mathbf{x})\log\frac{2p(\mathbf{x})}{p(\mathbf{x})+q(\mathbf{x})} + q(\mathbf{x})\log\frac{2q(\mathbf{x})}{p(\mathbf{x})+q(\mathbf{x})}d\mathbf{x} - \log(4)$ |

- **$f$-GAN objective function**

$$V(G,D) = \mathbb{E}_{\mathbf{x}\sim p_{data}(\mathbf{x})}[D(\mathbf{x})] - \mathbb{E}_{\mathbf{z}\sim p(\mathbf{z})}[f^*(D(G(\mathbf{z})))]$$

- **GAN objective function**

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

- **$f$-GAN is reduced to GAN, when**
  - $f^*(t) = -\log(1 - e^t)$    GAN是一個f-GAN的special case
  - output activation function of $D$
    - $f$-GAN $= -\log(1 + e^{-v(\mathbf{x})})$
    - GAN $\to \frac{1}{1+e^{-v(\mathbf{x})}}$
    - $v(\mathbf{x})$ is the output before activation function
- Consider measuring with difference divergences
  - $W(p_0||p_\theta) = |\theta|$
  - $JS(p_0||p_\theta) = \begin{cases} \log 2 & \text{if } \theta \neq 0 \\ 0 & \text{if } \theta = 0 \end{cases}$
  - $KL(p_\theta||p_0) = \begin{cases} +\infty & \text{if } \theta \neq 0 \\ 0 & \text{if } \theta = 0 \end{cases}$
  - $\delta(p_0||p_\theta) = \begin{cases} 1 & \text{if } \theta \neq 0 \\ 0 & \text{if } \theta = 0 \end{cases}$

## Deep convolutional GAN

- **Deep convolutional GAN (DCGAN)** (Radford et al., 2015)
  - **convolutional neural networks + adversarial learning**
- Architecture guidelines for stable DCGAN
  - replace any pooling layers with convolutions in $D$
  - **batchnormlization** in both $G$ and $D$
  - **remove** fully connected hidden layers
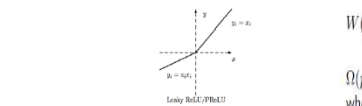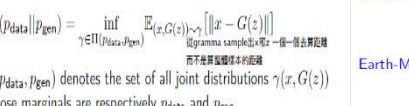  - use **LeakyReLU** for all layers in $D$



## Wasserstein GAN

**Wasserstein GAN (WGAN)** (Arjovsky et al., 2017) is optimized by the Wasserstein distance to avoid the problem in original GAN

- gradient vanish
- mode collapse

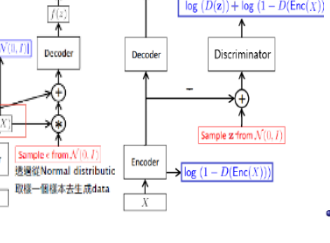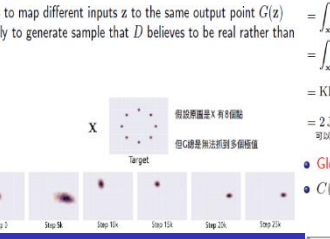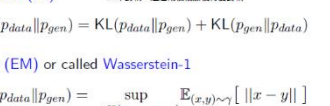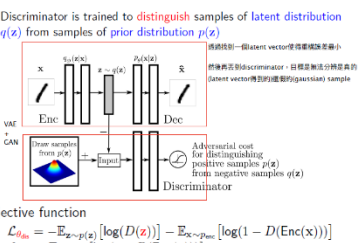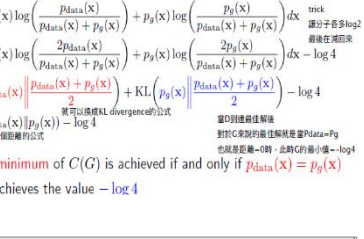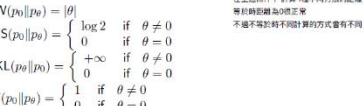**Use the Earth-Mover (EM) distance or Wasserstein-1 distance**

$$W(p_{data}||p_{gen}) = \inf_{\gamma\in\Pi(p_{data}, p_{gen})} \mathbb{E}_{(x,G(z))\sim\gamma}[\|x - G(z)\|]$$

$\Pi(p_{data}, p_{gen})$ denotes the set of all joint distributions $\gamma(x, G(z))$ whose marginals are respectively $p_{data}$ and $p_{gen}$

## Total Variance (TV)

直接算2個distribution的距離
$$\delta(p_{data}||p_{gen}) = \sup_{A\in\Sigma}|p_{data}(A) - p_{gen}(A)|$$

## Kullback-Leibler (KL)

直接算2個distribution的差異
$$KL(p_{data}||p_{gen}) = \int\log\left(\frac{p_{data}(x)}{p_{gen}(x)}\right)P_r(x)d\mu(x)$$

## Jensen-Shannon (JS)

KL不對稱，透過兩個相加使得變對稱
$$JS(p_{data}||p_{gen}) = KL(p_{data}||p_{gen}) + KL(p_{gen}||p_{data})$$

## Earth-Mover (EM) or called Wasserstein-1

$$W(p_{data}||p_{gen}) = \sup_{\gamma\in\Pi(p_{data}, p_{gen})} \mathbb{E}_{(x,y)\sim\gamma}[\|x - y\|]$$

## Image generation

- **Super resolution GAN (SRGAN)** (Ledig et al., 2016)

**SRGAN objective function**
$$E_{\mathbf{x}^{HR}\sim p_{data}(\mathbf{x})}[\log D(\mathbf{x}^{HR})] + E_{\mathbf{x}^{LR}\sim p_g(\mathbf{x})}[\log(1 - D(G(\mathbf{x}^{LR})))]$$

- input a low-resolution image $\mathbf{x}^{LR}$ to $G$ instead of latent code $z$
- $\mathcal{L}^{SR}$ is weighted combination of content loss and adversarial loss
  - $\mathcal{L}^{SR} = \mathcal{L}_{content} + \lambda\mathcal{L}_{adv}$

**Objective function**
$$\mathcal{L}^{SR}(G(\mathbf{x}^{LR}), \mathbf{x}^{HR})$$
$$\mathcal{L}_{content/i,j} = \frac{1}{W_{i,j}H_{i,j}}\sum_{n=1}^{W_{i,j}}\sum_{m=1}^{H_{i,j}}\left([\phi_{i,j}(\mathbf{x}^{HR})]_{n,m} - [\phi_{i,j}(G(\mathbf{x}^{LR}))]_{n,m}\right)^2$$

$\phi_{i,j}(\cdot)$: **feature map** of $j^{th}$ convolution before $i^{th}$ maxpooling
$W_{i,j}$ and $H_{i,j}$: **dimensions** of feature map

**Adversarial loss**
- encourages network to favour images that reside in manifold of natural images
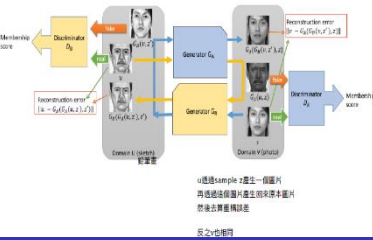
$$\mathcal{L}_{adv} = -\log D(G(\mathbf{x}^{LR}))$$

## Architecture of CycleGAN

輸入是一張圖片，而不再是一個gaussian sample



再看生成的回去generator A 希望重構誤差越小越好

## CycleGAN



- **Adversarial Loss**
$$\mathcal{L}_{GAN}(G, D_Y, X, Y) = \mathbb{E}_{y\sim p_{data}}[\log D_Y(y)] + \mathbb{E}_{x\sim p_{data}(x)}[\log(1 - D_Y(G(x)))]$$

- **Cycle Consistency Loss**
$$\mathcal{L}_{cyc}(G, F) = \mathbb{E}_{x\sim p_{data}(x)}[||F(G(x)) - x||_1] + \mathbb{E}_{y\sim p_{data}(y)}[||G(F(y)) - y||_1]$$
  - $F$ is another generator, which can be viewed as an inverse function of $G$

- **Objective function**
$$\mathcal{L}(G, D_X, D_Y) = \mathcal{L}_{GAN}(G, D_Y, X, Y) + \mathcal{L}_{GAN}(F, D_X, Y, X) + \lambda\mathcal{L}_{cyc}(G, F)$$
$$G^*, F^* = \arg\min_{G,F}\max_{D_X,D_Y}\mathcal{L}(G, F, D_X, D_Y)$$

## Architecture of DiscoGAN



$$\mathbf{x}_{AB} = G_{AB}(\mathbf{x}_A)$$
$$\mathbf{x}_{ABA} = G_{BA}(\mathbf{x}_{AB}) = G_{BA}\circ G_{AB}(\mathbf{x}_A)$$   A轉到B再轉回A
$$\mathcal{L}_{rec_A} = d(G_{BA}\circ G_{AB}(\mathbf{x}_A), \mathbf{x}_A)$$
$$\mathcal{L}_{GAN_B} = -\mathbb{E}_{\mathbf{x}_A\sim P_A}[\log D_B(G_{AB}(\mathbf{x}_A))]$$

- $G_{AB}$ receives two types of losses

$$\mathcal{L}_{G_{AB}} = \mathcal{L}_{\mathsf{GAN}_B} + \mathcal{L}_{\mathsf{rec}_A}$$
$$\mathcal{L}_{D_B} = -\mathbb{E}_{\mathbf{x}_B \sim P_B}[\log D_B(\mathbf{x}_B)] - \mathbb{E}_{\mathbf{x}_A \sim P_A}[\log(1 - D_B(G_{AB}(\mathbf{x}_A)))]$$

  – reconstruction loss $\mathcal{L}_{\mathsf{rec}_A}$ that measures how well the original input is reconstructed
  – standard GAN loss $\mathcal{L}_{\mathsf{GAN}_B}$ that measures how realistic the generated image is in domain $B$
- DiscoGAN objective function

$$\mathcal{L}_G = \mathcal{L}_{G_{AB}} + \mathcal{L}_{G_{BA}} = \mathcal{L}_{\mathsf{GAN}_B} + \mathcal{L}_{\mathsf{rec}_A} + \mathcal{L}_{\mathsf{GAN}_A} + \mathcal{L}_{\mathsf{rec}_B}$$
$$\boxed{\mathcal{L}_D = \mathcal{L}_{D_A} + \mathcal{L}_{D_B}}$$ 跟道是cyclegan 只是搬化D的loss

## Architecture of DualGAN



u透過sample z產生一個圖片
再透過這個圖片產生回來原本圖片
然後去算量構誤差

反之y也相同

## Transfer learning

- Transferring knowledge from source domain to enhance learning capability in target domain



- The distribution of source domain and target domain are different, we can't directly use source domain data to train the model for target domain
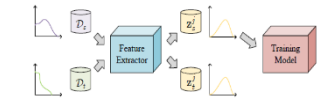
- Labels of target domain data are often not observed

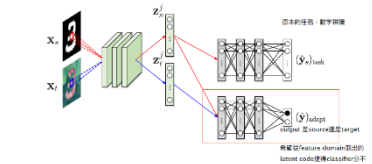- We can solve these problems by applying distribution matching

## Distribution matching

Domain adaptation aims to find a feature extractor where the domain of output features coming from source domain and target domain are same or similar

The model trained by features of source domain can work properly with the features of target domain



不同的domain再feature上 要取出相似的資訊

- Maximum mean discrepancy (Gretton et al., 2012) is a kernel method that measures the discrepancy between distributions

## Adversarial domain adaptation



標本的任務：數字辨識

output 是source還是target
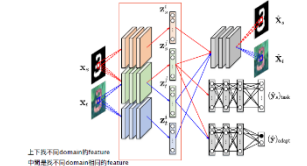
希望從feature domain取出的
latent code使得classifier分不
出來是source還是target

- DAN (Ganin et al., 2016) incorporates the adversarial learning into the training procedure of domain adaptation
- Task classification loss is measured by cross entropy

$$\mathcal{L}_{\mathsf{task}}(\mathsf{Enc}^j, C_{\mathsf{task}}) = -\frac{1}{N}\sum_{n=1}^{N_s}\sum_{k=1}^{K} y_{snk}\log(\hat{y}_{snk})_{\mathsf{task}}$$

## Domain separation network

跟上面差不多 只是多了一型feature

保證編的私分(min 最小化重檔距量)



上下兩不同domain的feature
中間是共不同domain相同的feature

- DSN (Bousmalis et al., 2016) introduces two additional private encoder to capture individual features