



## Generative Model

classifier: Given  $(x, y)$ , find the probability  $P_{x,y}(x, y)$   
 Generator: Given  $y$ , generate sample from  $x \sim P_{\cdot, y}$

### ① Auto Encoder

- Given a sample  $x$ , transform  $x$  into a "code" variable  $z$   
 then reconstruct  $x'$  which is close to original  $x$
- $P(z)$  training is fast, but  
 issue: - sampling (test time) is slow  
 - inverse may not be unique
- soln: VAEs

### ② VAEs

- build decoder as a probabilistic model  $P_\theta(x|z)$
- Invert using variational inference  
 to get approximation  $P_\theta(z|x)$
- construct an approximate encoder from decoder

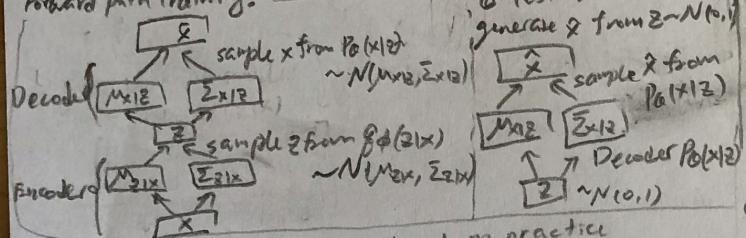
#### Encoder

- approximate inverse density  
 $g_\phi(z|x) \propto 1/N \propto 1/\text{param}$

#### Decoder

- Generator ( $\theta$ ) generate  $N$  samples  
 from  $x \sim P_\theta(x|z)$

#### Forward path training:



- Issue: computing  $g_\phi(z|x)$  is hard in practice  
 Soln: Reparameterization Trick  
 Instead, compute  $z = g_\phi(x, \epsilon)$  s.t.  $\epsilon \sim N(0, I)$

#### Loss func for VAE

$$\begin{aligned} \text{max likelihood } P(x) &= \max_x \log P(x) \\ &\geq \max_x \mathbb{E}_{z \sim g_\phi(z|x)} \left[ \log \left( \frac{P(x, z)}{g_\phi(z|x)} \right) \right] \end{aligned}$$

variational lowerbound (convex)

- Issue:  $L(\theta, \phi|x)$  is not accurate  
 we sampling estimate

Soln: KL-Divergence

$$L(\theta, \phi|x) \approx -D_{KL}(g_\phi(z|x) \parallel P_\theta(z)) + \frac{1}{N} \sum_i \log P_\theta(x_i|z_i)$$

no sampling needed by sampling

#### Issue w/ VAE

- ① Density estimate for  $g_\phi(z|x)$   
 $z = g_\phi(x, \epsilon) = g_\phi(x) + g_\phi(1, \epsilon)$  assumes  $\epsilon$  is indep  
 → image becomes blurry

Soln: Auto regressive Model

- ② approximation w/  $D_{KL}(g_\phi(z|x) \parallel P_\theta(z))$

Soln: GAN

### ③ GAN

want: avoid approximations

- ① compare real/synthetic images in image space
- ② use  $x = g(z)$  that is general

Generator

discriminator

- learns a transform  $P_\theta(z)$  from  $P(z) \sim N(0, I)$
- it doesn't estimate density  $P(x)$  from image  $x$

#### Discriminator

use JS divergence b/w real  $P_\theta(x)$  & synthetic  $P_\phi(x)$

$$JSD(P_\theta(x) \parallel P_\phi(x)) = \frac{1}{2} \mathbb{E}_{x \sim P_\theta(x)} [\log \frac{P_\theta(x)}{m(x)}] + \frac{1}{2} \mathbb{E}_{x \sim P_\phi(x)} [\log \frac{P_\phi(x)}{m(x)}]$$

$$\text{where } m(x) = \frac{1}{2} (P_\theta(x) + P_\phi(x)) = \text{avg image}$$

#### Adversarial Training

$$\min_{\theta} \max_{\phi} \mathbb{E}_{x \sim P_\theta(x)} [\log \frac{P_\phi(x)}{m(x)}] + \mathbb{E}_{x \sim P_\phi(x)} [\log \frac{P_\theta(x)}{m(x)}]$$

↑ train discriminator,  
 train generator to minimize JSD to fool discriminator

#### ① DC GAN

→ adopted from CNN

#### ② Conditional GAN (CGAN)

→ use image + label

#### ③ AC GAN

→ conditioned on class label  
 → discriminator produces class label as output

#### Issue w/ GAN

- ① Mode collapse → soln: unrolled GAN

- ② Difficulty in training → soln: Wasserstein GAN (WGAN)  
 → replace JSD w/  
 Wasserstein distance

#### ③ Difficulty in understanding global structure

→ soln: Progressive GAN

→ incrementally add layers

④ Difficulty in evaluating GAN

→ soln: Inception score,

Frechet inception score

⑤ Unpaired conditional image generation → soln: X-GAN; Cycle GAN

- Autoregressive Model ex) Image Transformer

→ pixels are highly correlated locally

→ predict new pixel conditioned on already predicted pixel

→ Fully generative model (can compute density)

Issue: condition on all pixels already predicted is expensive

Soln: local Autoregressive model

→ condition on a few nearby pixels

① Pixel CNN  $\Rightarrow O(n)$

② Pixel RNN (Row LSTM)  $\Rightarrow$  @train:  $O(n \text{ rows}^2 \text{ image})$

③ Pixel RNN (diagonal LSTM)  $\Rightarrow$  @test:  $O(n \text{ rows}^2)$

→ use all prev pixels

**CNN**

- Assumes locality: important info is contained locally
- Locality invariance: can reuse same param for whole image

Output layer Dim:  $\frac{W \times H}{\text{stride}} \times \text{filtersize}$  (Ex: input = 32x32x3, stride = 2, filtersize = 5x5)  $D_o = D_f = \left\lfloor \frac{(W+2P-F)}{\text{stride}} \right\rfloor + 1$

Depth = # filters

# Parameters:  $\left[ (f_w \times f_h) \times (f_d \cdot \text{depth}) + 1 \right] \times (\text{filters})$  (filters depth bias)

ALEX Net

uses ReLU, conv → ReLU → Max pool → Batch Normal

Need big data to train → solve! Transfer learning

VGG Net

If small data, fix all weights in classifier  
If medium, increase # layers & smaller filters, use weights for initialization  
If large, small # params

GoogleNet

Inception module: good for human-like behavior  
Run diff kinds of CNN in parallel  
Issue: computationally expensive  $\Rightarrow$  soln: Bottlenecking  
As depth grows, reduce # params

ResNet

Add many layers, but maybe they are close to Identity

Residual Connection:  $x \rightarrow \text{Conv} \rightarrow \text{ReLU} \rightarrow F(x) + x$

Dense Network: Connect all layers to each other

Computer Vision

- Classification: classifies single obj
- Classification + Localization: classifies single object & bounding box

Intersection over Union (IoU) =  $\frac{\text{Area of Intersec}}{\text{Area of Union}} \in [0, 1]$

ImageNet CNN: affine classif  $\rightarrow$  softmax loss  
affine Box pred  $\rightarrow$  L2 loss  $\rightarrow$  sum loss

Object Detection:

- foreach object, category label
- Sliding window regression classif (bounding box)
- need to check many position  $\Rightarrow$  soln: region proposal
- R-CNN (region proposal + CNN)
- finding region prop takes long
- classifying each region takes long

Ideas:

- Fast R-CNN (swap conv layer in region proposal)
- ROI pooling (efficient cropping of resizing features)
- Faster R-CNN
- Add NN for region proposal
- use anchor box

use mAP (mean avg precision)

Instance Segmentation: pixel labels for each thing object

Semantic Segmentation: label every pixel for both things

CNN w/ downsampling & upsampling

Pooling/Strided conv

U-Net (FCNN + Residual)

downscale  $\xrightarrow{\text{batch} \times \text{height} \times \text{width}}$  map  $\xrightarrow{\text{res}}$  upscale by only one dimension (bottleneck)  $\xrightarrow{\text{res}}$

**Training NM**

Good for non-linear control

Sigmoid:  $f(x) = \frac{1}{1 + e^{-x}}$   $\hat{f}(x, i)$

Not zero-centered  $\Rightarrow$  if input  $> 0$ ,  $\frac{dx}{d\theta} > 0 \Rightarrow$  in high dim, grad becomes very small

Soln: zero-mean data! If input  $< 0$ ,  $\frac{dx}{d\theta} < 0$

Saturation as  $x \rightarrow \pm \infty$ ,  $(\frac{dx}{d\theta} = 0)$

Exp is expensive

Softplus  $\hat{f}(-1, i)$

Good for RNN

Zero-centered

But saturation

rand weight init:  $W \sim N(0, 0.01)$

Want to want to break symmetry

Xavier Initialization allows  $G_{in} = G_{out}$  Every layer  $F_{in} = \# \text{input neurons that contribute to output}$

For CNN,  $F_{in} \times F_{out} \times C$

Forget  $F_{in} = F_{out}$ , need  $\text{Var}(W) = F_{in}$

Dropout

Randomly set neuron to 0 w/p p

Avoid redundant info

train, norm / (1-p)

test, do nothing

Ensemble

Train k indep models

At test, take avg

Bagging (parallel): reduce var

train base model on bootstrap samples & take majority vote/avg

good when model doesn't have bias & variance is high

good for deep Net

use when models are same

Takes long time to train

Boosting (sequence): misclassified sample get higher weight

Each learned tries to reduce error on hard examples

good when we want to reduce bias & possibly variance

not for NN

How to implement?

- True ensemble: train a few models indep @ tail layers
  - avg the prediction probability (almost always work)
  - avg the parameters (Never works b/c too many possible space)
- Model Snapshot: keep training a single base model & take periodic snapshot of parameters
- Prediction averaging: always work
- Parameter averaging: often works b/c parameter space is narrow

Only need to keep a single model for future prediction

**Visualization**

t-SNE (Stochastic Neighbor Embedding)

Preserve locally pairwise distance

KL divergence (measure of how diff 2 dist are)

$D_{KL}(P||Q) = \sum_i P(i) \log \left( \frac{P(i)}{Q(i)} \right)$  s.t.  $P(i) = P(X_i | \text{class})$ ;  $Q(i) = P(Y_i | \text{class})$ ; anchor

$D_{KL}(P||Q) \geq 0$  iff  $P$  is almost everywhere

$D_{KL}(P||Q) \geq 0$   $\forall p, q$

because cluster is not sphere, t-means doesn't work

use DBSCAN, HDBSCAN

**Visualize filters** as only work in low layers

Saliency Approach: generate image that maximize class score

Gradient: picel layer for backprop

Perceptron: Guided Backprop

Neural Style Transfer

Idiot 20x20 don't remove 0 or negative  $\Rightarrow (f(x) > 0) \wedge (\text{don't } > 0) \wedge (\text{out})$

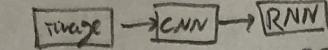
Attention

- learns saliency
- make NN interpretable
- reduce computation
- ↳ weighted comb/attention of some inputs
- ↳ differentiable  $\Rightarrow$  can back-prop

### Hard Attention

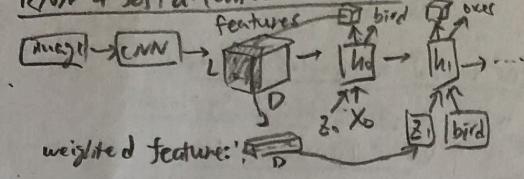
- ↳ attend to a single input location
- ↳ not differentiable  $\Rightarrow$  need RL

### RNN-captions Model



Issue: don't remember parts of image

### RNN + soft attention Model



weighted feature:  $D = \sum_i z_i h_i$   
 $D = \# \text{ features}$   
 $L = \text{width} \times \text{height}$

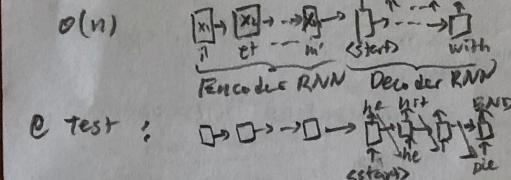
### Neural Machine Translation (NMT)

#### - Seq2Seq RNN

↳ Conditional lang Model

↳ directly calculates  $P(y|x) = P(y|x_1) \cdot P(y|x_2) \cdot \dots \cdot P(y|x_n)$

① Train :



Issue of: - forgets the beginning of source  $\Rightarrow$  Bi-directional RNN  
 - bottleneck: single path b/w Encoder & Decoder

#### - BLEU Score for Translation

↳ compare machine translation vs human translation

n-gram precision =  $\frac{\text{correct n-gram in ref sentence}}{\text{unique n-gram in test sentence}}$

+ clipping:  $\min\left(\frac{n}{2}, \log(P_n)\right)$   $\rightarrow$  N-gram precision

BLEU =  $\sum_{n=1}^N p_n \cdot e^{-\lambda \sum_{m=1}^{n-1} |c_m|}$  mean of N-gram precision  
 where  $p_n = \begin{cases} 1 & \text{if } c_n \in m \\ \frac{1}{2} & \text{if } c_n \in m \\ 0 & \text{else} \end{cases}$

#### - English to English (Summarization)

In : length N passage  
 length M summary

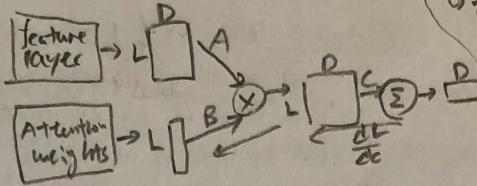
$O(N^2) + O(M^2) + O(NM)$   
 Encoder Decoder cross-attention  
 self-att self-attention

Soln: remove input/encoder self-attention  
 $O(N^2) + O(NM)$

### Reinforcement learning (RL)

- $x_t \rightarrow S_t$  action env
- reward  $r_t$
- maximize reward @ each step

#### Soft Attention Mechanism



$$\frac{\partial L}{\partial B} = \sum_i A_i \frac{\partial L}{\partial C_i} = \text{saliency}$$

$$= \begin{cases} > 0 & \text{if increase attention} \\ < 0 & \text{if decrease} \end{cases}$$

### Statistical Machine Translation (SMT)

$$y^* = \arg \max_y P(y|x) = \arg \max_y P(x|y) P(y) \quad \begin{matrix} \text{English} \\ \text{French} \end{matrix} \quad \begin{matrix} \text{Language Model} \\ \text{Translation Model} \end{matrix}$$

### NLP Challenge

↳ Tokenization challenge

- ① UNK token
- ② Increase vocab size
- ③ word piece tokenization
- ↳ break down unknown words

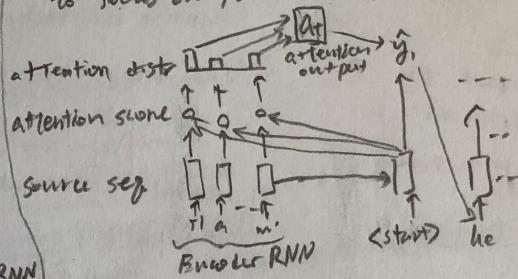
↳ Tokenization is not unique & not reversible

#### - Unresolved issue w/ translation

- out-of-vocab words
- domain mismatch
- maintain context over long text
- what if we have little text available?
- machine don't have common sense
- learns bias (race, gender, etc)

#### - Soft Attention for Translation

↳ at each step of the decoder, use direct connection to the encoder to focus on particular part of input



#### Adv:

- solves bottleneck
- helps vanishing gradient
- finds alignment automatic
- Issue: hard to generalize to NN
- Soln: Global Attention

Issue w/RNN Model: - lack long range dependency  
 & soln: - @ training time,  $O(n)$  time

#### - The Transformer

↳ QKV attention  
 ↳  $O(1)$  training time

#### - Self Attention

↳ Q, K, V come from Encoder/Decoder

#### - Cross Attention

↳ Q come from Decoder  
 K, V come from Encoder

#### - Bidirectional attention

↳ each word attend to its left & right

↳ used for transformer Encoder

↳ can't be used for generation

#### - Masked Attention

↳ each word attend to only to its left

↳ used for transformer Decoder

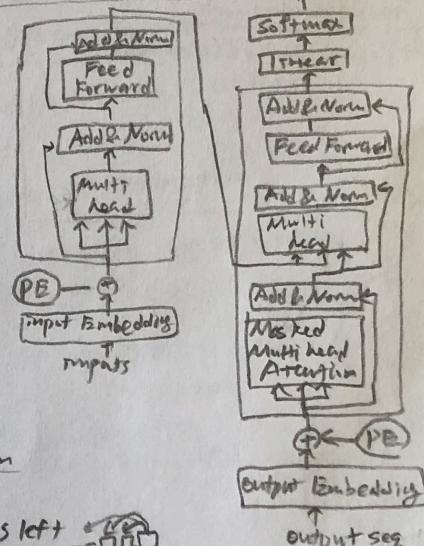
↳ can be used for generating next word

#### - Position Encoding

↳ learns relative displacement  $t-X$  b/w 2 language

(@ train,  $O(1)$  time!!)

(@ test,  $O(\text{length of target})$ )



## Text Semantics

### - Propositional / Formal Semantic

- ↳ convert text into a formula
- ↳ goal for inference in well-defined domain
- ex) dog bites man  $\Rightarrow$  bites(dog, man)

### - vector representation

- ↳ embed text into high-dim space

### ① Bag of Words (BoW)

$$\text{ex) } \text{vect}(\text{"dog"}) = (0.2, -0.3, 1.5, \dots) \in \mathbb{R}^n$$

$$\text{vect}(\text{"bites"}) = (\quad) \in \mathbb{R}^n$$

$$\text{vect}(\text{"man"}) = (\quad) \in \mathbb{R}^n$$

" "      vect("dog bites man")

issue: order of words not considered

### - Lexical semantics $\rightarrow$ sum of words

- ↳ focus on the meaning of individual words

### ① t-SNE

- ↳ preserve the distance for close pts

### ② Word2Vec

- ↳ minimize global distance

### ③ Co-occurrence Matrix

### - Compositional Semantics

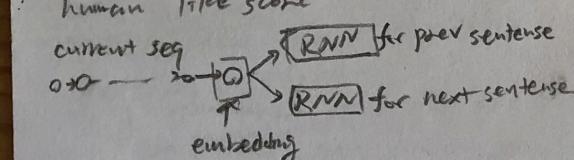
- ↳ meaning depends on how words are combined

### ① Skip-thought vectors

- ↳ seq-to-seq RNNs to predict the next & previous sentences

### ↳ Semantic Related Evaluation Scores

human-like score



### ② Siamese Network

- ↳ semantic relatedness

## vector Embeddings

### ⑤ Paralinguistic word similarity

- ↳ similar words occur in similar context
- $\Rightarrow$  can exchange words

ex) My dog barked, ...

Hec dog barked, ...

### ⑥ Dimensionality Reduction (PCA)

- ↳ issue w/ vector encoding: len of word vector is vocab size

### ⑦ Latent Semantic Analysis (LSA)

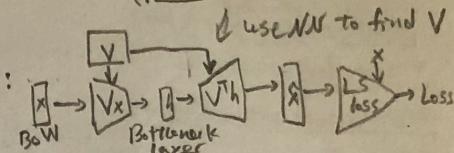
- ↳ encode a set of documents in a matrix  $T$

$T_{ij}$  = count of word  $j$  in doc  $i$

- ↳  $T$  is sparse  $\Rightarrow$  use SVD

- ↳ auto-encoding method  $\Rightarrow$  use NN:

$$N_{\text{docs}} \times T = U \Sigma V^T$$



### ⑧ Word2Vec

- ↳ focus on local context

skip-gram  $\equiv$  a pair of center word & context word

ex) It was a great day yesterday

### - skip-gram model

- ↳ predict context word from center word

$$P(j|i) = \frac{e^{u_i^T v_j}}{\sum e^{u_i^T v_i}}$$

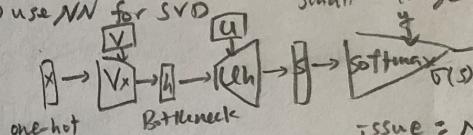
$v_j$  = emb. of context word  $j$   
 $v_i$  = emb. of center word  $i$

- ↳ SVD minimize dist b/w original & embedding pts

- ↳ capture relationship b/w similar words

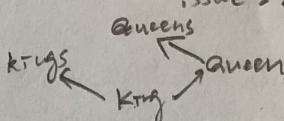
ex) Japan - Tokyo  
small - large

- ↳ use NN for SVD



issue: Normalization is expensive

- ↳ composition



### ③ Co-occurrence Matrix

#### - Window-Based Co-Occur Matrix

- ↳ count each pair of words

- ↳ similar to word2vec

issue: - increase in size w/ vocab

- classification model is sparse

solution: use SVD

#### Glove Model

$$\text{min} \sum_{i,j} f(C_{ij})(u_i^T v_j + b_i + b_j - \log(C_{ij}))^2$$

$$\text{s.t. } f(x) = \begin{cases} \left(\frac{x}{x_{\max}}\right)^a & \text{if } x < x_{\max} \\ 1 & \text{otherwise} \end{cases}$$

- ↳ scalable w/ huge corpus

- ↳ efficient

## natural lang understanding (NLU) vs

Given a text, build a representation of the inputs that can be used for many tasks

### ex) Extractive Summarization

Extract K sentences that summarizes texts

### - Bert (Bidirectional Encoder Representation from Transformer)

↳ masked lang model + Is Next pretrained

↳ select 15% random words ↳ predict next sentence

↳ replace w/ <Mask> token

↳ predict what the original words are

### - Fine tune BERT

#### ① Sentence pair classification Task.

↳ Given 2 sentences, are they saying same thing? T/F

#### ② Question Answering Task. e.g. SQuAD dataset

#### ③ natural lang inference (NLI)

↳ Given a premise P, classify whether hypothesis H is entailment/contradiction/neutral

#### ④ Corpus of Language Acceptability (CoLA)

↳ Given a sentence, classify if grammatically correct

## Security to Adversarial Attack

↳ Since Traditional ML assumes Training data & Test data, add adversarial noise to fool model

### White Box Attack

↳ optimization based attack

$$x^* = x + \beta \left( \frac{\nabla_x L(f(x), y)}{\|\nabla_x L(f(x), y)\|} \right)$$

adv. perturbation      adv. target label

#### ① Fast Gradient method

$$x^* = x + \beta \operatorname{sign}(\nabla_x L(f(x), y))$$

#### ② Fast Gradient Sign method

$$x^* = x + \beta \operatorname{sign}(\nabla_x L(f(x), y))$$

ISSUE: need to know lot of info  
(weight, learning rate, etc)

- GAN w/ Adv Network

- Defense to Adv Image

↳ ensembling, but GAN can still defeat

- Attack on RL

↳ proxy Network Attack.

## Natural Lang Generation (NLG)

↳ Given texts, produce a sequence of words

ex) Abstractive summarization

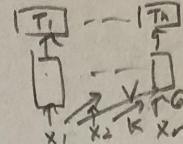
↳ produce a new summary

perplexity = evaluation score for comparing 2 lang models

### - GPT

↳ generative model is computationally efficient & scalable

↳ transformer based generator, but w/o encoder



### NLG tasks

#### ① Conversational Q&A (CoQA)

↳ Given a dialogue (conversation) b/w teacher & student, teacher must generate answer

↳ ISSUE: coreference = need to be asked previously

#### ② Dialog System (DSTC)

mistaken to classify specific class

### Black Box Attack

#### ① Zero-Query Attack

↳ randomly estimate gradient by perturbation

↳ transferability-Based attack ↳ Targeted attack

↳ Train a model on some chosen model (SVM, kNN, etc)  
↳ produce adv. example

#### ② Ensemble Targeted Black Box attack

ISSUE: estimating gradient, weights, etc is hard

### ② Query-Based Attack

↳ Finite differentiating gradient Estimation

Given  $x \in \mathbb{R}^d$ , we can estimate gradient w/ 2d queries

$$\text{FD}_x(g(x), \delta) = \begin{bmatrix} \frac{1}{\delta}(g(x+\delta e_i) - g(x-\delta e_i)) \\ \vdots \\ \frac{1}{\delta}(g(x+\delta e_d) - g(x-\delta e_d)) \end{bmatrix}$$

ISSUE: need to estimate in all d dimension

SOLN: Query Reduced gradient Estimation  
↳ random feature grouping, PCA, etc