

②③④ Basic Principles

Setup: x, y : input, labels
 $f_\theta(\cdot)$ model
 $l(y, \hat{y})$ loss function optimal $\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n l(y_i, f_\theta(x_i))$
 Problem: $\hat{y} = f_\theta(x)$

Neural Nets: differentiable operations allow non-linearities
 1) Expressivity: express patterns we want to learn $f_\theta(\cdot)$
 2) Reliably learnable

Represent piecewise linear func that are differentiable
 ReLU using elbows
 weight w : slope = w
 located at bias b : $-\frac{b}{w}$

Initialization: He initialization: Gaussian $N(0, \frac{2}{n})$
 Regularization: to prevent overfitting

1) Explicit Regularization: add regularizer term
 $\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \left[\frac{1}{n} \sum_{i=1}^n l_{\text{train}}(y_i, f_\theta(x_i)) + R(\theta) \right]$

ridge regularization $R(\theta) = \lambda \|\theta\|^2$ prevents θ too large
 OLS Ridge Regression: $\hat{\beta} = (X^T X + \lambda I)^{-1} X^T Y$
 GD Ridge Regression: gradient: $\eta(2X^T(\hat{y} - X\hat{\beta}) - 2\lambda\hat{\beta})$
 update: $\hat{\beta}_{(t+1)} = (1 - 2\eta\lambda)\hat{\beta}_{(t)} + \eta 2X^T(\hat{y} - X\hat{\beta}_{(t)})$

2) Data Augmentation: fake observations of features
 3) Implicit Regularization: optimizer as implicit regularizer

Maximum A Posteriori (MAP): RLO responds to a prior

Gradient Descent (GD): iterative approach to find local optima

Errors: 1) Irreducible error: noise or randomness from $Y|X$
 2) Approximation error: not flexible enough for true signal
 3) Estimation error: bias: systematic error of learning process
 variance: randomness of training process

Features: representation of data-driven, allow generalized linear model to work well

⑤⑥ Survey of Architectures & Problems

Network Architectures:

- 1) Multi-layer Perceptron (MLP): fully connected
- 2) Convolution Neural Nets (CNNs): spatial regularity embedded, images
- 3) Recurrent Neural Nets (RNN): CNNs w/ internal state over time
- 4) Graph Neural Nets (GNN): nearby items more related
- 5) Transformers: access input data elsewhere and weight share

Types of Problems:

- 1) Regression
- 2) Classification
- 3) Generation
- 4) Recommendation

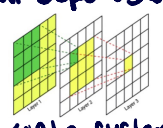
Optimization:

- 1) Gradient Descent: descend using learning rate
 To ensure grad descent stable, $\eta < \frac{1}{\sigma^2}$
- 2) Momentum based methods: low pass filter to make learning rate bigger w/o trouble for large singular values
 LPF the directions that oscillate, update as any grad direction
 "Vanilla" Momentum $\nabla L(w_t)$
 "Nesterov" Momentum $\nabla L(w_t - \eta(1-\beta)\nabla L(w_t))$
- 3) Adaptive approaches: change learning rates according to singular val
 $\hat{\alpha}_{k+1} = \frac{\hat{\alpha}_{k1}}{1 - \beta^k}$ $\hat{\beta}_{k+1} = \frac{\hat{\beta}_{k1}}{1 - (\beta')^k}$ $u_{k+1}[i] = u_k[i] - \eta \frac{\hat{\alpha}_{k+1}[i]}{\sqrt{\hat{\beta}_{k+1}[i]} + \epsilon}$

⑦⑧⑨ CNN

class of neural nets used to analyze images
 - respect "locality", "invariances",
 - support hierarchical structure, multi-resolution understanding
 - weight sharing: neighborhood of pixels is processed by kernel
 pooling 2×2 elements w/ stride to desired pixel
 Output: $\frac{N - K + 2P}{S} + 1$

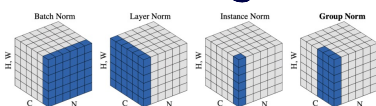
Receptive field: pixels in original image the output depends on
 receptive fields grow linearly but exponentially w/ stride/padding



Data Augmentation: transfer domain knowledge into system
 ex) auto-contrast, rotation, translation

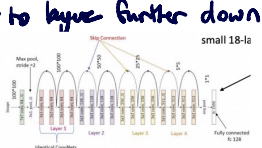
Standardization & Normalization: convert data to zero mean & unit variance, so raw magnitude doesn't have effect on gradient

- 1) Batch Normalization: normalizes layer for mini batch
- 2) Layer normalization: normalize across all channels
- 3) Instance normalization: each channel each training image
- 4) Group normalization: over group of channels



⑩ CNN Architectures, Dropout

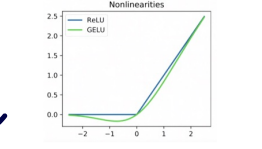
Residual Network: skip connections, output to layer further down
 ex) ResNet, allowed more expressive networks, learn salient features



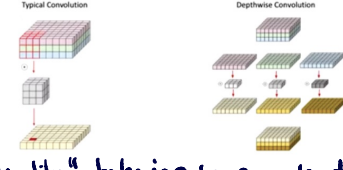
ConvNext: layer normalization + Gaussian ReLU

Gaussian ReLU (GReLU): $x \cdot \Phi(x)$
 $\Phi(x)$: Gaussian CDF

- Smoother, non-convex, non-monotonic



Depthwise Convolution: broken into channels, convolve, concat



Dropout: "ensemble-like" behavior to promote diversity and redundancy, learning rate can compensate for training slower
 - kill certain units setting to 0

Stochastic Depth Regularization: drop entire residual blocks

Label Smoothing: probabilities in goal array
 $y = \left[\frac{1 - \alpha \cdot k - 1}{k}, \frac{\alpha}{k}, \frac{\alpha}{k}, \dots \right]^T$ $k = \text{classes}$

⑪ GNN

Graph w/ information in nodes, "generalized" CNN
 GNN vs. CNN: diff number of neighbors in GNN, no ordering

Weight Sharing: consider itself and neighbor nodes, regardless of ordering and have learnable func $f_w[m_e, \sum_{\text{neighbors}} s_{i2}(m_e, \text{them})] g_{\text{out}}(k_{\text{sum}})$

Pooling: can cluster and pool or not as useful
 Can run if we have lots of different graphs as data in training

12) RNN

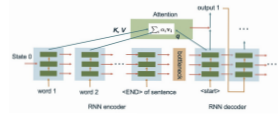
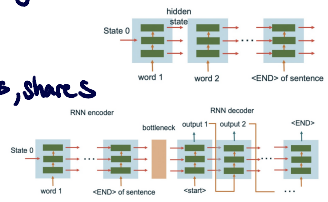
Residual Connection blocks issue of vanishing gradients
 Speed Up Training by:
 1) larger batch size, linearly scale learning rate
 2) distributed data-parallel training
 3) regularize by aggressive data augmentation

14) Attention/self-supervision

RNN encoder: good for sequential words, shares same weight

Bottleneck: include input statement
 Attention: allow to look back at words originally embedded, allow info to flow along layers
 use hash table to store key, values,
 To get values for query, scan for closest matches of queries to keys

$e_i = \frac{q \cdot k_i}{d}$

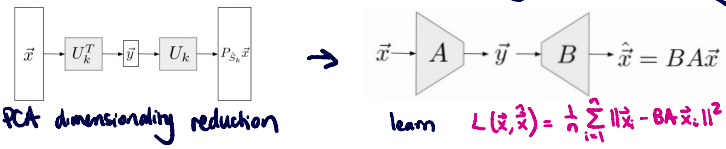


To track order information, we can use vector expression (complex expression)

15) Self-Supervision and Autoencoders

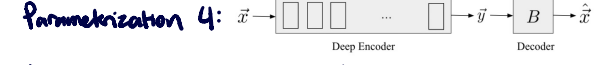
Self-supervision: for unsupervised learning (PCA, clustering), design to use loss func, grad to supervised

Easier to obtain unlabeled data, use dimensionality reduction clustering



Autoencoders: want grad descent to learn x-hat to P_k x-hat
 doesn't necessarily need to be linear, k-dimensional structure
 can replace A, B w nonlinear encoder/decoder

A Weight Sharing: $A = (B^T B)^{-1} B^T$ $\hat{x} = BA\bar{x} = B(B^T B)^{-1} B^T \bar{x}$



Data Augmentation can help prevent learning identity

Other

Beam Search Runtime: $O(TkM \log M)$ & $O(T^2 k M \log M)$

LSTM: $C_t \rightarrow \oplus \rightarrow \oplus \rightarrow C_{t+1}$

f_t : sigmoid $(Wx_t + Wf_t + W_3(t + bias))$
 input: $\tanh(Wx_t + Wf_t + W_3(t + bias))$

Practice Exams

- Newton's Method can converge to global optimum when loss func optimized is convex
- SGD does not find same empirical gradient
- batch normalization introduces dependence between data pts in one-mini batch
- convolutional block reduces # of channels to speed up forward/backward

Misc

Linear Algebra

- L2 Norm (Euclidean): $\|x\|_2 := \sqrt{\sum_{i=1}^n x_i^2} = \sqrt{x^T x}$
- L1 Norm: $\|x\|_1 = \sum_{i=1}^n |x_i|$
- L-infinity Norm: $\|x\|_\infty := \max_{1 \leq i \leq n} |x_i|$
- Frobenius Norm: $\|A\|_F := \sqrt{\sum_{i,j} A_{ij}^2} = \sqrt{\text{Tr}(A^T A)}$
- Orthogonal: $U^T U_j = 0, i \neq j$ and $\|u_i\|_2 = 1$
- Trace: $\text{Tr} A = \sum_{i=1}^n A_{ii}$: sum of diagonal elements
- Cauchy-Schwarz Inequality: $|z^T y| \leq \|y\|_2 \cdot \|z\|_2$

Fundamental Theorem of Linear Algebra
 Range of a matrix is the orthogonal complement of the nullspace of its transpose

$R(A^T) = N(A)$

Spectral Theorem (Symmetric eigenvalue decomposition (SED))

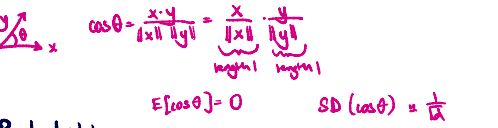
$A = \sum \lambda_i u_i u_i^T = U \Lambda U^T, \Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$
 $\text{Tr} \Sigma = \frac{1}{m} \sum X X^T = \frac{1}{m} \|X\|_F^2$

- Total Variance: $\text{Tr} \Sigma = \text{Tr}(U \Lambda U^T) = \text{Tr}(U^T U \Lambda) = \text{Tr} \Lambda = \lambda_1 + \dots + \lambda_n$
- rank: linearly independent columns $X: \mathbb{R}^{m \times n}$
- Rank-Nullity Theorem: $n - \dim(\text{nullspace}(X)) = \text{rank}(X)$
- $\dim(\text{rowspace}(X)) = \dim(\text{columnspace}(X^T)) = \text{rank}(X^T) = \text{rank}(X)$
- symmetric matrix has real eigenvalues
- eigenvalues neg if concave
- Axis scaled by square roots of eigenvalues of Σ
- Cover = $\frac{X X^T}{n}$ $X \in \mathbb{R}^{n \times d}$

Symmetric matrix M

- Positive Definite: if $w^T M w > 0$ all $w \neq 0 \Rightarrow$ pos eigenvalues
- Positive Semidefinite: if $w^T M w \geq 0$ all $w \Rightarrow$ nonnegative eigenvalues
- Indefinite: if pos & neg eigenvalue
- Invertible: no zero eigenvalue

convex: x^2
 concave: \sqrt{x}



Probability

- Bayes Rule: $P(Y=1|X) = \frac{P(X|Y=1)P(Y=1)}{P(X)}$
- $P(A, B) = \sum_c P(A, B, C)$
- Chain Rule: $P(A, B, C) = P(A, B|C) P(C) = P(A|B, C) P(B|C) P(C)$
- A conditionally independent given C: $P(A, B|C) = P(A|C) P(B|C)$
- A independent of B given C: $P(A|B, C) = P(A|C)$
- A, B independent: $P(A, B) = P(A)P(B)$

$P(A|B, C) = \frac{P(A, B, C)}{P(B, C)} = \frac{P(A, B|C)}{P(B|C)}$
 $P(X) = P(X|Y=1)P(Y=1) + P(X|Y=-1)P(Y=-1)$

Matrix Derivatives

- $\frac{\partial x^T a}{\partial x} = \frac{\partial a^T x}{\partial x} = a$
- $\frac{\partial a^T X b}{\partial X} = a b^T$
- $\frac{\partial a^T X^T b}{\partial X} = b a^T$
- $\frac{\partial a^T X a}{\partial X} = \frac{\partial a^T X^T a}{\partial X} = a a^T$
- $\frac{\partial X}{\partial X_{ij}} = J^{ij}$
- $\frac{\partial (X A)_{ij}}{\partial X_{mn}} = \delta_{im} (A)_{nj} = (J^{mn} A)_{ij}$
- $\frac{\partial (X^T A)_{ij}}{\partial X_{mn}} = \delta_{in} (A)_{mj} = (J^{mn} A)_{ij}$
- $\frac{\partial}{\partial X_{ij}} \sum_{klmn} X_{kl} X_{mn} = 2 \sum_{kl} X_{kl}$
- $\frac{\partial b^T X^T X c}{\partial X} = X (b c^T + c b^T)$
- $\frac{\partial (B x + b)^T C (D x + d)}{\partial x} = B^T C (D x + d) + D^T C^T (B x + b)$
- $\frac{\partial (X^T B X)_{kl}}{\partial X_{ij}} = \delta_{ij} (X^T B)_{ki} + \delta_{ij} (B X)_{il}$
- $\frac{\partial (X^T B X)}{\partial X_{ij}} = X^T B J^{ij} + J^{ij} B X \quad (J^{ij})_{kl} = \delta_{ik} \delta_{jl}$
- $\frac{\partial x^T B x}{\partial x} = (B + B^T) x$
- $\frac{\partial b^T X^T D X c}{\partial X} = D^T X b c^T + D X c b^T$
- $\frac{\partial}{\partial X} (X b + c)^T D (X b + c) = (D + D^T) (X b + c) b^T$

- $\text{Tr}(AB) = \text{Tr}(BA)$
- $(A^T)^T = A$
- $(A+B)^T = A^T + B^T$
- $(AB)^T = B^T A^T$
- PSD $Q = P D P^T$
 $P D = D^{1/2} P^T$
 $= U U^T$
 $A = A^{1/2} A^{1/2}$
 $= (P D^{1/2} P^T) (P D^{1/2} P^T)$
- Eigenvalues 2x2
 $\begin{bmatrix} a & b \\ c & d \end{bmatrix} \rightarrow m \pm \sqrt{m^2 - p}$
 $\frac{a+d}{2} \rightarrow \mu$
 $\sqrt{(\frac{a-d}{2})^2 + bc} \rightarrow \sigma$

16) Self Supervision

Autoencoders: $\vec{x} \rightarrow \text{encoder} \rightarrow \vec{z} \rightarrow \text{decoder} \rightarrow \hat{\vec{x}} \rightarrow \text{surrogate obj}$
 learn $\vec{z} \rightarrow \text{neural network} \rightarrow \hat{y}$
 ↑ use larger set of unlabeled

Surrogate tasks:

vanilla autoencoding: rate approach, $L(\vec{x}, D(E(\vec{x})))$

denoising autoencoding: $L(\vec{x}, D(E(\vec{x} + \vec{n})))$, $\vec{n} \sim \mathcal{N}(0, \sigma^2)$

masked autoencoding: $L(\vec{x}, D(E(\vec{x}')))$ $\vec{x}' = [x_1, \dots, x_3, x_4?]^T$

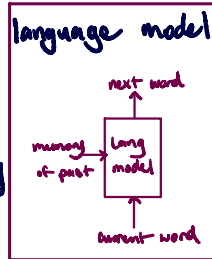
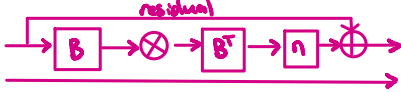
1) First Parameterization

$$\vec{x} \rightarrow [A] \rightarrow \vec{z} \rightarrow [B] \rightarrow \hat{\vec{x}}$$

2) Second Parameterization

$$\vec{x} \rightarrow [B^T B] \rightarrow \vec{z} \rightarrow [B] \rightarrow \hat{\vec{x}}$$

3) Third Parameterization: RNN w/ weight sharing



Beam Search: keep bag of current best possibilities, choose k best, least negative, add them together

17) 18) 19) Transformer Models

attention

1) drop recurrence from RNN, keep weight sharing

- 1) Soft approx hash table
- 2) Queryable softmax pooling

hash table pairs k_i and v_i , apply a query \vec{q}

$$e_{i,t} = \frac{\langle \vec{q}, \vec{k}_i \rangle}{\Delta}, \quad \alpha_{i,t} = \frac{e^{e_{i,t}}}{\sum_j e^{e_{j,t}}}, \quad \text{output} = \sum \alpha_{i,t} \vec{v}_i$$

2) can also use multiple channels in transformers, multi head attention query through multiple heads & concatenate outputs

- for a query, output is approx corresponds to nearest key,

inner product $\text{sim}(q, k) = \frac{\langle q, k \rangle}{\Delta}$

- attention mechanism to learn dependencies across input sequence, agg info, no learnable parameters

Transformers

- 1) inputs $x_{t,1}$
- 2) create key vector $k_{t,1}$, query vector $q_{t,1}$, value vector $v_{t,1}$
- 3) store key vector $k_{t,1}$, $v_{t,1}$ into table
- 4) pass query vec $q_{t,1}$ as input to attention block, softmax
- 5) attention to linear layer W , combined with $x_{t,1}$ with skip connect
- 6) layer Norm \rightarrow MLP \rightarrow LN

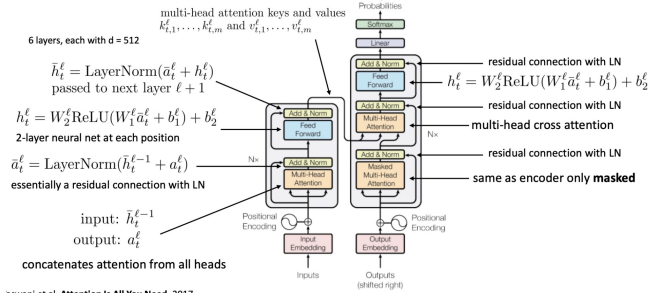
Masking: should not perform inner products w/ key-value pairs from future, set future values to 0

Cross-Attention: queries from decoder, key-value from encoder

Position Encoding: add positions to transformer input tokens concatenating position to vector or adding it

Transformers: self attention layer, relevant to next, cross-attention w key-value from encoder, weight matrix has learnable weights
 Represent words w/ embedding to show similarity

The Transformer



20) Pretraining and fine-tuning

input \rightarrow tokenizer \rightarrow lookup table (id \rightarrow vector) \rightarrow transformer (tokens \rightarrow id)

Word2vec: find vector representations of each word where similar words have similar embeddings

$$\arg \max \sum_{v_0} \log \frac{\exp(u_0^T v_0)}{\sum_{v \in V} \exp(u_0^T v)}$$

$$\arg \max \sum (\log \sigma(u_0^T v_0)) + \sum \log \sigma(-u_0^T v_0)$$

Pretrained models:

- 1. train language model on surrogate task
- 2. Run it on a sentence
- 3. Take hidden state from model and treat it as embedding

masked self-attention: make sure model doesn't just look ahead

BERT: bidirectional Transformer Language Models, mask percentage of inputs, don't need masked self-attention, trained on predicting mask and next sentence prediction, to predict if two sentences swapped

Finetuning can either freeze BERT and retrain classifier at end or retrain end-to-end (given pretrained transform)

GPT: autoregressive model, one directional transformer

21) 22) 23) Fine tuning

- Finetuning:**
- 1) pretrain large model w/ self supervision & lots data
 - 2) finetune on task specific data or objective

Feature Extraction: decapitate training head and train new head
 ↑ less params to retrain, scalability, ↓ ok performance

Finetuning: replace head, retrain entire model

↑: increased capacity, ↓ overfitting harder scale divergent backbone grad

Prompt Engineering: prompts to make model do what we want

↑: no training, ↓ bad performance, limited training data

Prompt Tuning: create prompts in vector language

Catastrophic Forgetting: forgetting how to perform old tasks

when trained to do new ones,

Continual learning: learns series of tasks sequentially

- earlier layers are somewhat task specific

- skip connections allow early layers to learn task specific

↳ Naive approach is to batch learn randomly

↳ **Replay dummy training:** examples of old task when training on new, (preferred solution)

↳ **Learning w/o forgetting:** create pseudo labels after running through old heads, knowledge distillation, generate analog labels

T5/BART: BERT: encoder only, masked autoencoder
GPT: decoder only, predict next token

encoder-decoder transformers, masked autoencoder
 mask spans tokens

Soft prompts: population of attention tables in between layers

- fine-tuning can improve performance, but if model params too small, prompt cannot work

24) Meta Learning

Meta-learning: find system to quickly, reliably learn new tasks

1) Feature extract 2) Fine tuning

1) randomly initialize task specific head

2) either fine tune entire model or just head

3) use sgd to update params w/ differentiable loss func

NAML ↓ exploding grad while training, memory

Alternatus: train on union of tasks

ANIL/Meta Opt Net/RZD2: freeze "feature extractor"

optimize task head, differentiate w/ params

Variational Autoencoders

VAEs fall in the class of likelihood-based models. The goal is to learn a model $p_\theta(x)$ that is close to the true distribution p_{data} . This is done by maximizing the likelihood of the observed data under the model, i.e. $\max_\theta \mathbb{E}_{x \sim p_{data}} [\log p_\theta(x)]$.

Evidence Lower Bound

Given a latent-variable model, $z \rightarrow x$, we have an alternative objective to maximize the likelihood of the data, i.e. $\max_\theta \mathbb{E}_{x \sim p_{data}} [\log p_\theta(x)]$. In particular, consider the expression

$$\begin{aligned} \log p_\theta(x) &= \log \int_z p_\theta(x, z) dz \\ &= \log \int_z p_\theta(x|z) p(z) dz \end{aligned}$$

- latent var z between encoder and decoder is noisy
 - find probability $X_t | X_{t+1}$ as cond distribution, procedure for denoising

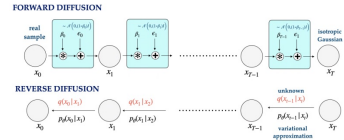


Figure 4: Diffusion Model iteratively adds noise to the input skewing the distribution towards well-understood distributions.

Diffusion models are a class of generative models that are *multi-step*, meaning that they perform multiple steps (say T) of inference. In particular, diffusion models are composed of two phases:

• **Forward:** In this phase, we start from samples from the real distribution and iteratively add noise to the samples. In particular, starting from $x_0 \sim p_{data}$ we have:

$$q(x_t|x_{t-1}) = \mathcal{N}(\sqrt{\beta_t}x_{t-1}, (1 - \beta_t)\mathbb{I})$$

Notably, the variance of the distribution is a function of the step t , and the variance increases as we move forward in time. A typical example³ of this would be linearly increasing from $\beta_0 = 10^{-4}$ to $\beta_T = 0.02$.

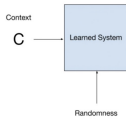
• **Reverse:** Sampling from the diffusion model is done by starting from a sample from our target distribution (e.g. isotropic Gaussian) and iteratively denoising the noisy inputs. In the Markov chain illustrated in fig. 4 this corresponds to starting at x_T and moving backwards in time, such that x_0 is the denoised sample.

Performing this denoising step, however requires us to estimate $q(x_{t-1}|x_t)$ which is not tractable. Instead, we perform a variational

³Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33: 6840-6851, 2020.

25) 26) Generative Tasks

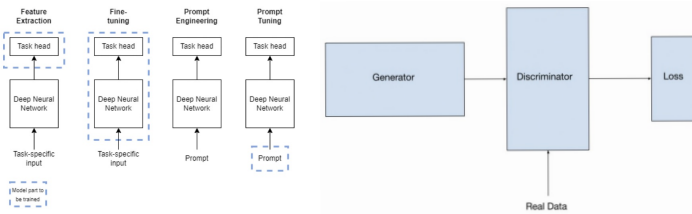
Generation Model: generate unseen example of data



typically, provide context (prefix) which is passed through layers

Generative Adversarial Networks: generator to create images and classifier (discriminator) to tell if image real or fake, loss trains both at same time.

↓ delicate, mode collapse if same output, always rejects

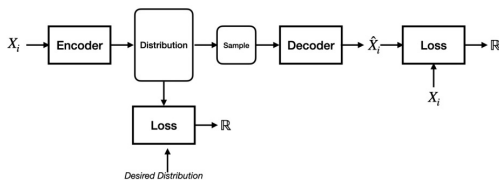


- 1) Train discriminator given frozen generator
- 2) Train generator given frozen discriminator, given real or fake either minimizer or max sgd steps

Mode collapse: if generator lacks diversity

27) Diffusion

Variational Autoencoders: input to decoder random during training, add loss term on distribution of z , parameterize sgd, estimate density



Disc

- robotic navigation task
- encoder/decoder (T5), encode target, inputs history of trajectory in decoder, auto regressively masking, output time logits, cross-entropy loss
- subword tokenizers preferred, word tokenizers map words not in vocab to <UNK>

175B param, 100k data → soft prompting
 90B param, no data classification → hard prompting
 1B param, 100 data image classify → feature extraction
 1B param, 100M data ImageNet → full finetuning