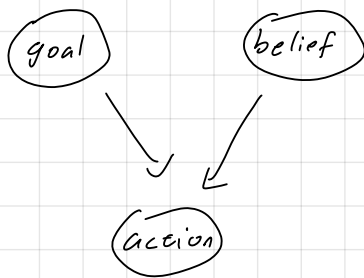


limitation in LLMs.

- language reasoning.
 - Embodied reasoning.
 - social reasoning.
(understand human)
- Human: internal world model
- Human: strategic planning
- internal model to predict states.
- simulation of alternative plans.
 - access outcomes to refine / pick the best.
- ↓
- Theory of mind: reasoning about hidden mental variables



Human conduct model-based reasoning based on world models and agents.

World Model in Humans

- Perceiving physical properties
 - Predicting dynamics
 - Model-based control/planning.
- human tool use: can learn tool through a few trials.

world model as state transition probability.

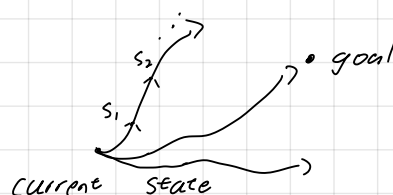
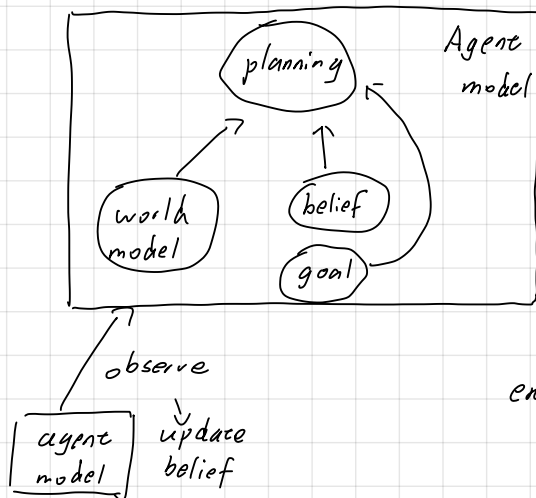
causal relationship between action and state change.

$$P(s' | s, a)$$

↑ ↑ ↑
next current action.
state state

Agent model in human:

- strength
- goal
- relationship
- moral judgment
- beliefs.



(simulate plan via world model).

ex) Alpha Go

MCTS with a known world model.

level - 1 agent model:

- agent as a observer for observing level - 0 agent

model-based theory of mind:

$$P(\underbrace{\text{mind}}_{\substack{\text{understand} \\ \text{the mind through state, action}}} \mid \text{state, actions}) \propto \underbrace{P(\text{actions} \mid \text{state, mind})}_{\text{level-0 agent}} P(\text{mind})$$

Behavior Prediction

$$P(\text{future action} \mid \text{state, mind}).$$

Human - Ai interaction:

$$P(\text{action}_{AI} \mid \text{state, mind}_{AI}, \text{mind}_{human}).$$

Recursive social reasoning:

level - 2 agent

↓ observe

level - 1 agent

⋮

Experience:

- Data example

1. language reasoning

- math, logic, ...

2. space/time

- vision, audio, sensor2, ...

3. multi-agency

- adversary, collaboration, ...

4. law of nature

- biology, ...

- Rules / Constraints

- Knowledge graphs

- Rewards

- Auxiliary agents

- Adversaries

- Master classes.

...

problems with few label:

- Privacy, security

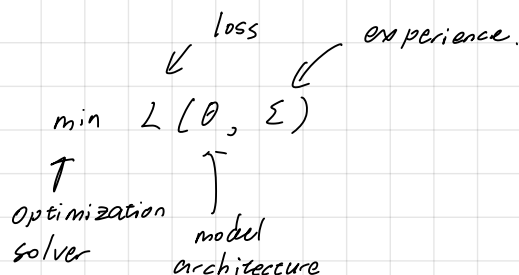
- Expensive to collect/annotate

- Controllable content generation.

- Difficult / expertise-demanding to annotate.

- Specific domain

ML Solution:



Algorithm marketplace

Designs driven by: experience, task, loss function, training procedure ...



Algorithm marketplace listing various ML techniques:

- maximum likelihood estimation
- reinforcement learning as inference
- data re-weighting
- policy optimization
- inverse RL
- active learning
- data augmentation
- reward-augmented maximum likelihood
- label smoothing
- imitation learning
- softmax policy gradient
- actor-critic
- adversarial domain adaptation
- posterior regularization
- constraint-driven learning
- knowledge distillation
- GANs
- reward
- generalized expectation
- prediction minimization
- regularized Bayes
- learning from measurements
- energy-based GANs
- weak/distant supervision



unified, standardized ML principles

A "standardized formalism" of ML

$$\min_{q, \theta} -H + D - \bar{E}$$

Annotations:

- $-H$ is labeled "uncertainty".
- D is labeled "Divergence".
- \bar{E} is labeled "Experience".

LLM

NLP before 2017:

Automated understanding and generation of natural language

- Name Entity Recognition.

ex) Adam Driver was born in San Diego.
person city.

- Sentiment Analysis.

Since 2017.

Every year, model size increases by 10x.

- google Bert → understanding
- OpenAI GPT → generative
- GPT-4 → mixture of experts model.

Language Model:

next word prediction given context.

$$P(w_i | w_1, \dots, w_{i-1}).$$

↓
probability distribution

↓
sample from distribution

↓
generate.

→ (sampling decoding).
→ (greedy decoding).
→ (top-k decoding).

Implementation → Transformer

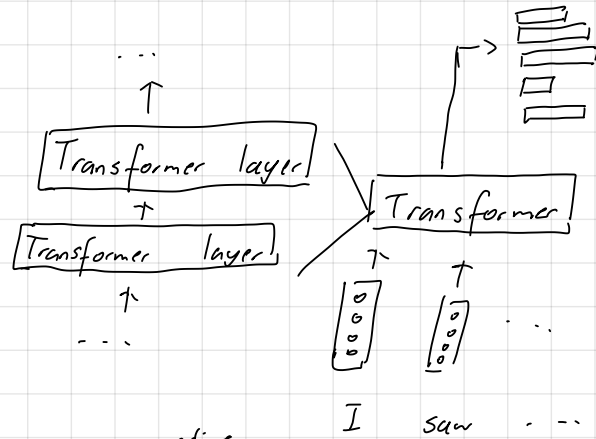


Attention mechanisms.

en)

$$p(* | \text{I saw a cat on a})$$

Training of ML:



Self-supervised learning:

- predict the future from the past $||||| \rightarrow \text{predict}$
- predict the future from the recent past $||| ||| \rightarrow \text{predict}$
- predict the past from the present. $\leftarrow ||| |||$
- predict the top from the bottom.

(predict any part of the past, present or future percepts from whatever information is available).

Motivation: • successfully learning to predict everything from everything else would result in the accumulation of lots of background knowledge about how the world works.

- Mass amount of data.

SSL in Language Model.

- sentence:

$$y = (y_1, y_2, \dots, y_T)$$

$$p_{\theta}(y) = \prod_{t=1}^T p_{\theta}(y_t | y_{1:t-1}).$$

Training:

given data example y^* ,

minimizes negative log-likelihood of the data.

$$\min_{\theta} \mathcal{L}_{MLE} = -\log p_{\theta}(y^*) = -\prod_{t=1}^T p_{\theta}(y_t^* | y_{1:t-1}^*).$$

- A transformer-based LM with 125M to 175B parameters

- Trained on massive text data.

Word Embedding:

- Conventional word embedding
 - Word2vec, Glove
 - A pre-trained matrix, each row is an embedding vector

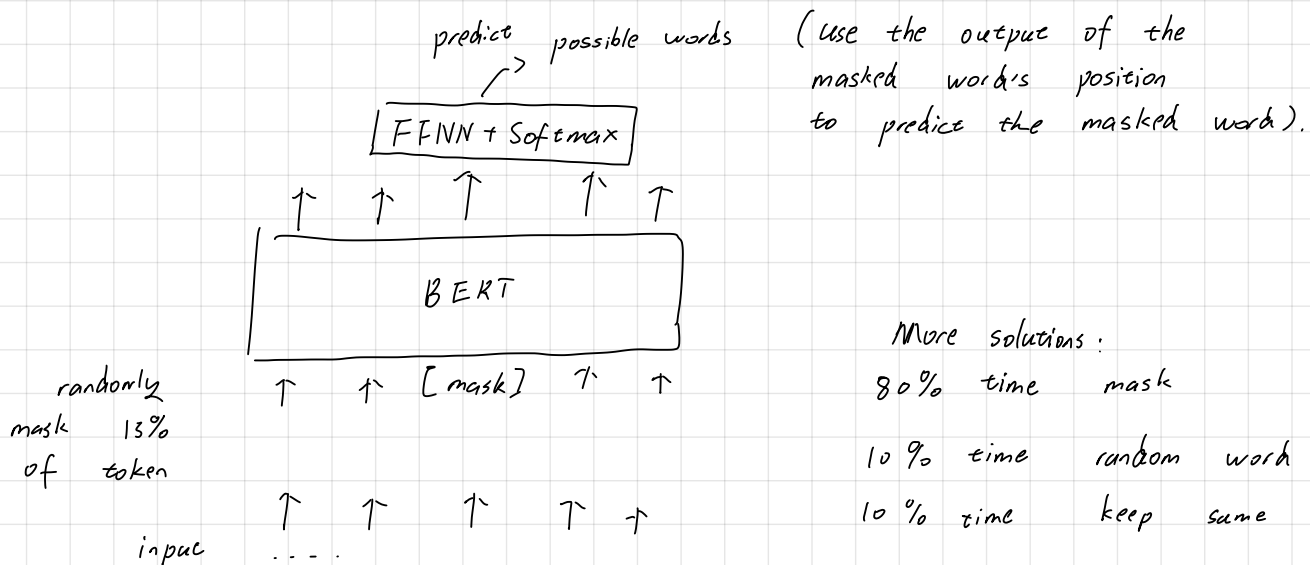
corpus \rightarrow Word2vec \rightarrow embedding matrix

- Problem: word embedding are applied in a context free manner
- Solution: train contextual representations on text corpus.

↓

BERT: a bidirectional model to extract contextual word embedding

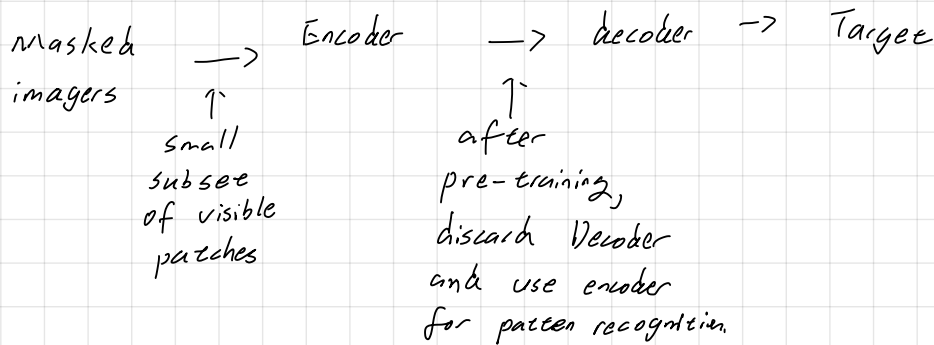
- Training procedure
 - Masked language model (masked LM).
 - masks some percent of words from the input and has to reconstruct those words from context.



- Two sentence task:
 - To understand relationship between sentences:

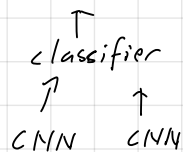
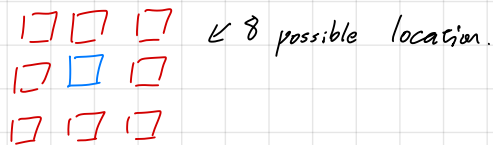
- Predict likelihood that sentence B belongs after sentence A.
- BERT: Downstream Fine-tuning:
 - BERT for sentence classification.

Masked Autoencoder (MAE)



Relative Positioning:

Train network to predict relative position of two regions in the same image.



pre-train CNN using self-supervision

1. Input image
2. Extract region proposals
3. Compute CNN features.

• Colorization:

Train network to predict pixel colour from a monochrome input

SSL from videos:

- Video Sequence Order
 - Sequential Verification
- Video Direction
 - Predict if video playing forwards or backwards.

• Video Tracking:

- Given a color video, colorize all frame of a gray scale version using a reference frame.

Enhancing LLM

- Limitation:
- Lack World and Agent knowledge → Need richer learning mechanism
 - Inefficiency of the language modality → Need multi-modal capabilities

• Richer learning mechanisms

- learning with embodied experience

- Where to go
- How to go
- How to learn

• Embodied simulators

- Household activities
- Touchdown
- Minecraft
- OS
- Simulated websites

• How to get:

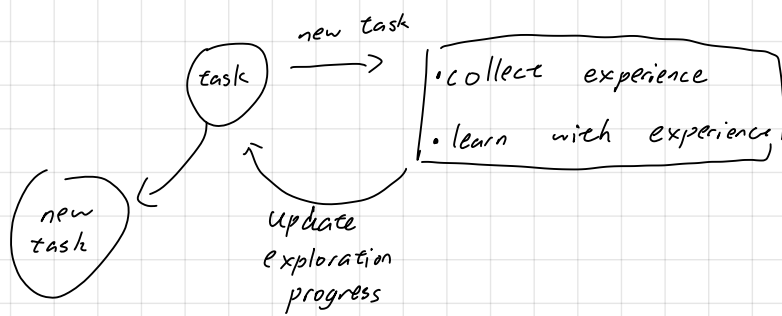
- Goal-oriented

- Collecting experiences by completing a given task



• Auto curriculum

- proposing new tasks automatically by prompting GPT-4 to generate new tasks.



- Random Exploration:
- Finetuning LMs with the experiences
- But also want to preserve the original language capabilities of LMs.
 - instead of overfitting
 - Continual learning with EWC (Elastic Weight Consolidation).

$$\Downarrow$$

$$\bar{F}_{i,i} = \frac{1}{n} \sum_{i=1}^N \left(\frac{\partial L_u^{(i)}}{\partial \theta_{u,i}^*} \right)^2$$

$$L(\theta) = L_v(\theta) + \lambda \sum_i \bar{F}_{i,i} (\theta_i - \theta_{u,i}^*)^2$$

↑
conventional
finetuning
objective

↑
regularizer to
preserve important weights.

- Updating external memory
 - instead of changing LM parameters

Multi-modal Capabilities

existing limitation

- can understand images / cannot generate images
- can do interleaved generation of image and text.
 - however generated images are not consistent → Lack internal world model.

- Video diffusion model \rightarrow generates video given actions
 - predict future frame given previous frame and action.

$$P(s_t | s_{t-1}, a_{t-1}).$$

- Text-to-Video generation.
 - ex) Sora

\Rightarrow Need a more general world model.

- integrate different spaces
- generalise capability
- real time control.

Latent-space Reasoning

\downarrow
which fuses information of different observed modalities.

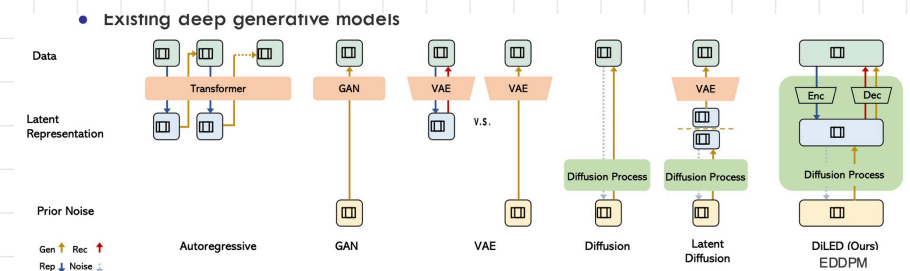
- Multi-level latent space
 - immediate next moves
 - mid-term and long-term planning and thought experiments.

\downarrow
How to learn a good latent space?

- Compact and well-structured representation of the world, realistic generation.
consistent reconstruction.

Examples of latent-space

- Autoregressive
- GAN
- VAE
- Diffusion
- Latent diffusion.
- DILEED (EDDPM).



• No Free Lunch Theorem:

- No single model is universally best-performing algorithm for all problems.
- All algorithms perform equally well when their performance is averaged across all problems.

Unsupervised Learning

probability:

sample space: space of all possible outcomes

- $p(x, y)$ joint probability

- $p(y|x) = \frac{p(x, y)}{p(x)}$

- $E[f(x)] = \sum_x f(x) p(x) \quad \text{or} \quad = \int_x f(x) p(x) dx$

- $p(x) = \sum_y p(x, y)$

- $p(x_1) = \sum_{x_2} \dots \sum_{x_N} p(x_1, \dots, x_N)$

- $p(x, y) = p(y|x) p(x)$

- $p(x_1, \dots, x_N) = p(x_1) p(x_2|x_1) \dots p(x_N|x_1, \dots, x_{N-1})$

- Bayes rule:

$$p(\theta|D) = \frac{p(D|\theta) p(\theta)}{p(D)}$$

↑
posterior belief

↖ likelihood
↘ prior belief

- $p(x, y) = p(x) p(y)$

- $p(x, y|z) = p(x|z) p(y|z)$

- Gaussian Distribution

- Multinomial Distribution

Entropy:

Shannon entropy $H(p) = - \sum_x p(x) \log p(x)$

average level of "information"

KL Divergence: measure closeness of two distributions $p(x)$ and $q(x)$

$$KL(q(x) || p(x)) = \sum_x q(x) \log \frac{q(x)}{p(x)}$$

also the relative entropy.

q high, p high \Rightarrow low KL

• $KL \geq 0$

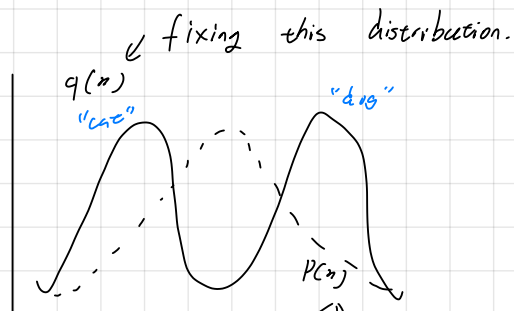
q high, p low \Rightarrow high KL.

• Not a true "distance"

If q is low, low KL regardless of p .

• Not commutative $KL(p || q) \neq KL(q || p)$.

• Not satisfy triangle inequality



see how close $p(x)$ is to $q(x)$

min KL is to approach $p(x)$ to $q(x)$.

• Supervised learning:

learn $p_\theta(x)$

observe full data.

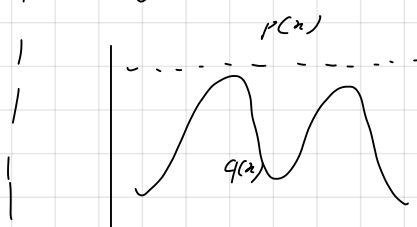
$$MLE: \min_{\theta} -E_{x \sim \tilde{p}(x)} [\log p_\theta(x)]$$



MLE is minimizing KL between data and model distribution.

$$\begin{aligned} KL(\tilde{p}(x) || p_\theta(x)) &= -E_{\tilde{p}(x)} [\log p_\theta(x)] + H(\tilde{p}(x)) \\ &= \sum_x \tilde{p}(x) \log \frac{\tilde{p}(x)}{p_\theta(x)} \\ &= \sum_x \tilde{p}(x) \log \tilde{p}(x) - \sum_x \tilde{p}(x) \log p_\theta(x) \\ &= H(\tilde{p}(x)) - \underbrace{E_{\tilde{p}(x)} [\log p_\theta(x)]}_{\text{cross entropy}} \end{aligned}$$

$$KL(p(x) || q(x))$$



VAE mode seeking

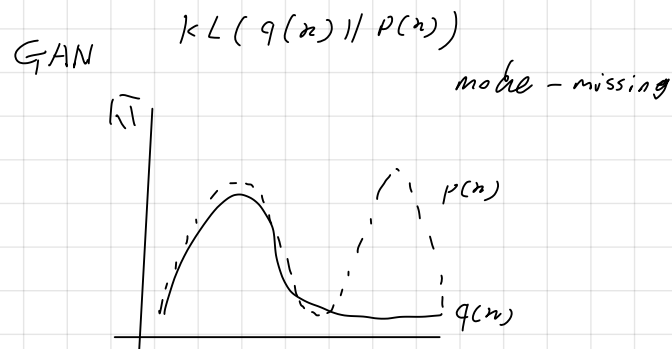
Unsupervised Learning

- observed variables x
- latent variables z .

want to learn $p_\theta(x, z)$.

latent variable:

- Discrete latent variable can be used for partition/cluster data into subgroups.
- Continuous latent variable can be used for



en) Gaussian Mixture Models.

a mixture of K Gaussian components

$$p(x_n | \mu, \Sigma) = \sum_k \pi_k N(x_n | \mu_k, \Sigma_k).$$

↑ mixture proportion
↑ mixture component

- z as latent class indicator

$$p(z_n) = \text{multi}(z_n; \pi) = \prod_k (\pi_k)^{z_n^k}$$

- x as conditional gaussian variable with class specific mean/covariance

$$p(x_n | z_n^k = 1, \mu, \Sigma)$$

$$= \frac{1}{(2\pi)^{m/2} |\Sigma_k|^{1/2}} \exp \left\{ \dots \right\}$$

- likelihood

taking the expectation.

$$p(x_n | \mu, \Sigma) = \sum_k \pi_k p(z^k = 1 | \pi) p(x_n | z^k = 1, \mu, \Sigma)$$

$$= \sum_k \pi_k N(x | \mu_k, \Sigma_k).$$

↑ mixture proportion, weight of gaussian distribution
↑ same different gaussian distribution how data fits in

MLE: $\ell(\theta; D)$ observed data (supervised).

$$\begin{aligned}
 &= \log \prod_n p(z_n, x_n) \\
 &= \log \prod_n p(z_n | \pi) p(x_n | z_n, \mu, \sigma) \\
 &= \sum_n \log \prod_k \pi_k^{z_n^k} + \sum_n \log \prod_k N(x_n; \mu_k, \sigma) ^{z_n^k} \\
 &= \sum_n \sum_k z_n^k \log \pi_k - \sum_n \sum_k z_n^k \frac{1}{2\sigma^2} (x_n - \mu_k)^2 + C
 \end{aligned}$$

to find unknown parameter

$$\hat{\pi}_k$$

$$\hat{\mu}_k$$

$$\hat{\sigma}_k$$

In unobserved data case:

Incomplete (or marginal) log likelihood:

with z unobserved, our objective becomes log of a marginal probability:

$$\ell(\theta, x) = \log p(x | \theta) = \log \sum_z p(x, z | \theta)$$

hard to learn \uparrow \uparrow marginalize over z

$$\begin{aligned}
 &\frac{\partial p(x_n | \mu, \Sigma, \pi)}{\partial \pi} \\
 &= \frac{\partial \frac{p(x_n | \mu, \Sigma)}{\partial \pi}}{p(x_n | \mu, \Sigma, \pi)}
 \end{aligned}$$

$$GMM: \log p(x_n | \mu, \Sigma) = \log \sum_k p(z^k = 1 | \pi) p(x_n | z^k = 1, \mu, \Sigma) \quad \theta_z$$

Complete

log likelihood of "known latent" model

$$p(A, B) = p(A) p(B | A)$$

$$\begin{aligned}
 \ell_c(\theta; x, z) &= \log p(x, z | \theta) = \log p(z | \theta) p(x | z, \theta) \\
 &= \log p(z | \theta_z) + \log p(x | z, \theta_x)
 \end{aligned}$$

$$GMM: \log p(x_n, z_n | \mu, \Sigma) = \log p(z_n | \pi) p(x_n | z_n, \mu, \Sigma)$$

Expectation Maximization (EM)

Intuition:

supervised MLE is easy : $\max_{\theta} \ell_c(\theta; x, z) = \log p(z, x | \theta)$

unsupervised MLE is hard : $\max_{\theta} \ell(\theta; x) = \log p(x | \theta) = \log \sum_z p(x, z | \theta)$

EM:

supervised MLE



M-step: $\max_{\theta} E_{q(z|x)} [\log p(x, z | \theta)]$ pretend we also observe z .

↑
distribution
that we "imagine"

E-step: $q(z|x) = p(z|x, \theta)$ can't observe q , estimate it.

↑
posterior
distribution.
using the current
parameter, estimate the
 q distribution.

EM:

(optimize q distribution)

E-step: $q^{t+1}(z|x) = p(z|x, \theta^t)$

iterative process

M-step: $\max_{\theta} E_{q^{t+1}(z|x)} [\log p(x, z | \theta)]$
(optimize model parameters)

similar idea GAN { generator
discriminator.

(coordinated descent)

Formally,

For any distribution $q(z|x)$, define expected complete log likelihood.

$$E_q[\ell_c(\theta; x, z)] = \sum_z q(z|x) \log p(x, z | \theta) = E_{q(z|x)} [\log p(x, z | \theta)]$$

• a deterministic function of θ supervised

• inherits the factorizability of $\ell_c(\theta; x, z)$

use as the surrogate objective



Maximizing this yields a maximizer of likelihood?

$$\ell(\theta; x) \geq E_q[\ell_c(\theta; x, z)] + H(q)$$

marginal
unsupervised

↑
EM optimize this lower bound
which in turn optimize this

proof of inequality:

$$\ell(\theta; x) \geq \mathbb{E}_q[\ell_c(\theta; x, z)] + H(q).$$

$$\ell(\theta; x) = \log p(x|\theta)$$

$$= \log \sum_z p(x|z, \theta)$$

$$= \log \sum_z q(z|x) \frac{p(x, z|\theta)}{q(z|x)}$$

$$\rightarrow = \log \mathbb{E}_{q(z|x)} \left[\frac{p(x, z|\theta)}{q(z|x)} \right]$$

$$\geq \mathbb{E}_{q(z|x)} \left[\log \frac{p(x, z|\theta)}{q(z|x)} \right] \quad (\text{Evidence Lower Bound})$$

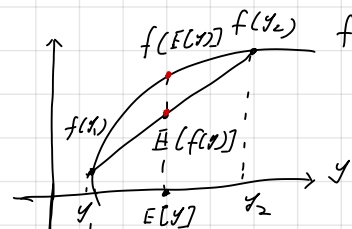
$$= \mathbb{E}_{q(z|x)} [\log p(x, z|\theta) - \log q(z|x)]$$

$$= \mathbb{E}_{q(z|x)} [\log p(x, z|\theta)] - \mathbb{E}_{q(z|x)} [\log q(z|x)] \quad \square$$

log
function
concave

Note:
convex: Jensen's inequality
 $\mathbb{E}_{p(y)}[f(y)] \geq f(\mathbb{E}_{p(y)}[y])$

concave:
 $\mathbb{E}_{p(y)}[f(y)] \leq f(\mathbb{E}_{p(y)}[y])$



ELBO

Thus, conclusion,

$$\ell(\theta; x) \geq \mathbb{E}_q[\ell_c(\theta; x, z)] + H(q) = \mathbb{E}_{q(z|x)} \left[\log \frac{p(x, z|\theta)}{q(z|x)} \right] \quad (\text{ELBO})$$

We can show that (HW).



$\uparrow \uparrow$ b/c $KL \geq 0$.

estimated distribution
given parameters.

$$\ell(\theta; x) = \mathbb{E}_{q(z|x)} \left[\log \frac{p(x, z|\theta)}{q(z|x)} \right] + KL(\underbrace{q(z|x)}_{\text{ideal distribution}} \parallel \underbrace{p(z|x, \theta)}_{\text{given parameters}}).$$

(ELBO)

Back to EM

M-step

E-step.

For fixed data x , define a functional called the (variational) free energy

$$F(q, \theta) = -\mathbb{E}_q[\ell_c(\theta; x, z)] - H(q) \geq -\ell(\theta; x)$$

function
of function

$$\uparrow$$

$$\min F(q, \theta)$$

Thus, EM is coordinated descent on F :

$$\text{E-step: } q^{t+1} = \argmin_q F(q, \theta^t)$$

$$\text{M-step: } \theta^{t+1} = \argmin_{\theta} F(q^{t+1}, \theta)$$

逐生逐
逐生逐
algorithm

(guest lecture)

Pre-training through language modeling.

- $P_{\theta}(w_t | w_{1:t-1})$, prob. discription of next word given previous context.

high quality pre-training data

⌞

Internet.

- however, web-data is noisy, dirty, biased.
- copyright and usage constraint.
- data is contaminated with auto-generated text.

⌞

training on these data will cause model collapse.

Diversity of data matter.

① Pre-training

↓

← limited data
adapt to task.

② Fine-tuning

SGD intuition:

- approach general pre-training loss
- approach local min fine-tuning loss

Full fine-tuning → update all parameters.

Parameter-efficient fine-tuning → update a few existing / new parameters.

(less overfitting).

↓

by injecting a adapter layers into original network, keep other parameters frozen.
(randomly initialized).

prefix finetuning \rightarrow learning a small continuous task-specific vector to each transformer block, while keeping the pre-trained LM frozen.

prompt fine-tuning \rightarrow A single "soft prompt" prepended to the embedding input.

low-rank Adaptation (LoRA) \rightarrow update the ^{of full weights.} low-rank representation parameters.
learn a low-rank "diff" between pre-trained and fine-tuned

Pre-training:

Encoder \rightarrow bidirectional, condition on the future context.
Autoencoder model, masked language model.

Encode & decoder \rightarrow seq2seq model

Decoder only \rightarrow Autoregressive model, Left-to-right language model.

Encoder: text reconstruction \rightarrow encode info. from both bidirection

Prompting:

- complete a sentence

- prompting

\downarrow

inference

- ask model in format.

- Few-shot prompting

- provide a few example together with the tasks in prompting

- Add "role", "name", "content" as JSON.

- In-context learning

- order of example matters

- label coverage matters
- label balance matters.
- Replacing correct label with random label will barely hurt the accuracy in in-context learning.
- Chain-of-Thought
 - get the model to explain its reasoning in output
- Zero-shot } - CoT → add a prompt to encourage model to reason
 Few-shot ("let's think step by step").
- Structuring output as program code
 (reasoning through coding)

Prompt engineering :

- Manual engineering : format
 - ↓
 - format of prompt matters in task accuracy.
- Paraphrase a existing prompt
- Reason + Action (ReAct)
- Persona-based prompting
 - role-playing ——— adapt to user
 - personalization ——— adapt to environment.
- self-refinement prompting :

EM continued:

E-step : minimization of $F(q, \theta)$ w.r.t. q .

Optimal solution: $q^{t+1} = \underset{q}{\operatorname{argmin}} F(q, \theta^t) = p(z|x, \theta^t)$

(the posterior distribution over the latent variables given the data and current parameters).

hint: $\ell(\theta; x) = E_{q(z|x)} \left[\log \frac{p(x, z|\theta)}{q(z|x)} \right] + KL(q(z|x) \parallel p(z|x, \theta)).$

$$= -F(q, \theta) + \underbrace{KL(q(z|x) \parallel p(z|x, \theta))}_{\text{approaches 0, } q(z|x) = p(z|x, \theta) \text{ as possible.}}$$

$\min \ell(\theta, x)$ approaches 0, $q(z|x) = p(z|x, \theta)$ as possible.

VI step: minimization of $F(q, \theta)$ w.r.t. θ .

$$F(q, \theta) = - \underbrace{E_q[\ell_c(\theta; x, z)]}_{\text{expected complete log likelihood}} - H(q) \geq -\ell(\theta; x).$$

\Downarrow
so, $\theta^{t+1} = \underset{\theta}{\operatorname{argmax}} E_q[\ell_c(\theta; x, z)]$
given as estimate

$$= \underset{\theta}{\operatorname{argmax}} \sum_z q^{t+1}(z|x) \log p(x, z|\theta)$$

solve as supervised.

Gaussian Mixture Model (GMM)

E-step: posterior of z_n given the current estimate of the parameters

$$\begin{aligned} p(z^k=1 | x, \theta^t) &= \frac{p(z^k=1) p(x|z^k=1)}{p(x)} \\ &= \frac{p(z^k=1) p(x|z^k=1)}{\sum_{j=1}^K p(z^j=1) p(x|z^j=1)} \\ &= \frac{\pi_k N(x|\mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x|\mu_j, \Sigma_j)} \end{aligned}$$

mixture proportion given parameter.

$$= \gamma_k$$

M-step: computing the parameter given current estimate of z

$$q^{t+1}(z^k | x) = p(z^k | x, \theta^t) = \gamma^k$$

$$\theta^{t+1} = \underset{\theta}{\operatorname{argmax}} \sum_k q^{t+1}(z^k = 1 | x) \log p(x, z^k = 1 | \theta)$$

$$E_{q^{t+1}} [\log p(x, z | \theta)]$$

$$= \sum_k \gamma_k (\log p(z^k = 1 | \theta) + \log p(x | z^k = 1; \theta))$$

$$= \sum_k \gamma_k \log \pi_k + \sum_k \gamma_k \log N(x; \mu_k, \Sigma_k)$$

we can maximize by solving parameter directly

EM for GMM:

- initialize μ_k, Σ_k, π_k
- iterate until convergence:

E-step.

M-step.

EM on LLM reasoning

$$E_{q(y|x)} [\log p(o=1 | x, y)] - KL[q(y|x) || p_\theta(y|x)]$$

$$= E_{q(y|x)} \log p(o=1 | x, y) - \log q(y|x) + \log p_\theta(y|x)$$

$$\leftarrow -E_q \log q(y|x) - q^*(y|x)$$

$$= -KL(\underline{q(y|x)} || \underline{p(o=1|x, y) p_\theta(y|x)})$$

$$\uparrow q^*(y|x) = p(o=1|x, y) p_\theta(y|x)$$

reweighting the answer

M-step.

E-step

- Generate data

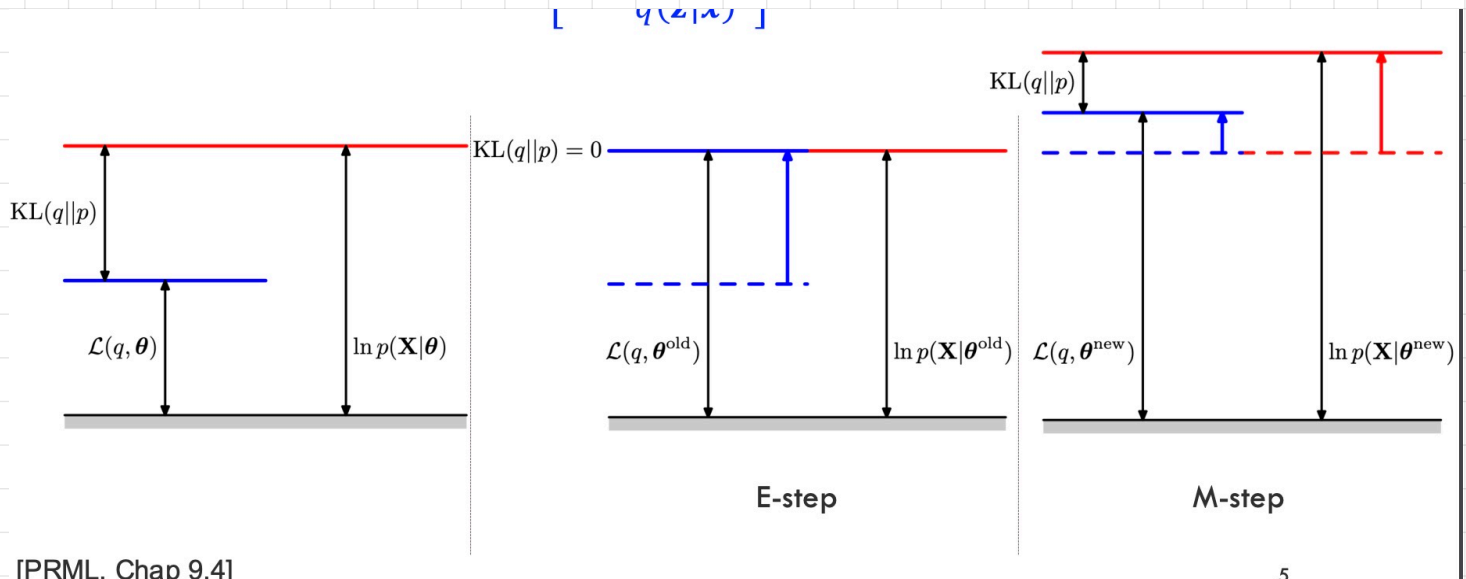
M-step

- optimise θ with positive data.

Replace PPO with EM.

Each EM iteration guarantees to improve the likelihood.

[4(212)]



5

limitation: need to compute $p(z|x, \theta)$ in E-step.

$$\frac{p(z, x | \theta^c)}{\sum_z p(z, x | \theta^c)}$$

might be hard to compute.
especially if it is continuous.

\Rightarrow Variational inference.

• Observed variable x

latent variable z .

• Variational Bayes:

used to approximate the posterior distribution over the latent variables.

$$p(z|x, \theta) = \frac{p(z, x|\theta)}{\sum_z p(z, x|\theta)}$$

$$p(A, B) = p(A) p(B|A)$$

$$p(z, x|\theta) = p(x|\theta) p(z|x, \theta)$$

$$\downarrow$$
$$\sum_z p(z, x|\theta)$$

in E-M, we assume $q(z|x)$ can be any distribution.

& E-step shows the optimal $q(z|x)$ is the posterior distribution

idea:

1. choose a family of distribution over latent variable $z_{1:m}$ with its own set of variational parameters v .

$$q(z_{1:m}|v)$$

2. Then we find the setting of the parameters that makes our approximation q close to the posterior distribution.

\Rightarrow optimization problem.

3. use q with the fitted parameter in place of posterior.

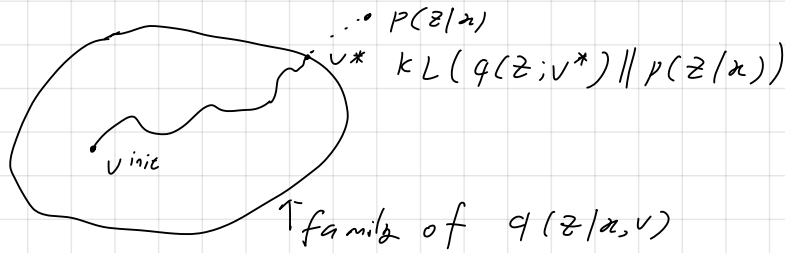
\Rightarrow

$$\min_q KL(\underbrace{q(z|x, v)}_{\text{approximation}} \parallel \underbrace{p(z|x)}_{\text{posterior}}).$$

$$\text{how } \Rightarrow \underbrace{Q(\theta, x)}_{\text{constant}} + E_{q(z|x, v)} \left[\log \frac{p(x, z|\theta)}{q(z|x, \theta)} \right] = -KL(q(z|x) \parallel p(z|x, \theta)).$$
$$\downarrow$$
$$\arg\max E_{q(z|x, v)} \left[\frac{\log p(x, z|\theta)}{q(z|x, \theta)} \right]$$

$$= \operatorname{argmax}_v \mathbb{E}_{q(z|x, v)} [\log p(x, z|\theta)] - \mathbb{E}_{q(z|x, v)} [\log q(z|x, v)].$$

idea:



using EM:

$$q(z|x, v^*)$$

$$= \min_v \text{KL}(q(z|x, v) || p(z|x, \theta^*))$$

$$= \min_v F(q(z|x, v), \theta^*) + \text{const.}$$

Q: How to choose variational family of $q(z|x, v)$?

- factorized distribution \rightarrow mean field VI
- mixture of gaussian distribution \rightarrow black-box VI
- neural-based distribution \rightarrow variational autoencoder (VAEs)

Mean field VI:

Assume variational distribution over latent variable factorizes as:

$$q(z) = q(z_1, \dots, z_m) = \prod_{i=1}^m q(z_i).$$

\uparrow
independent z_i

"local variational approximation".

ex) Bayesian mixture of gaussians.

Assume mean-field $q(u_{1:k}, z_{1:n}) = \prod_k q(u_k) \prod_i q(z_i).$

For each data example $i \in D.$

update local variational distribution \rightarrow E-step

update global distribution & parameters. \rightarrow M-step.

Until ELBO converges.

Or Stochastic VI

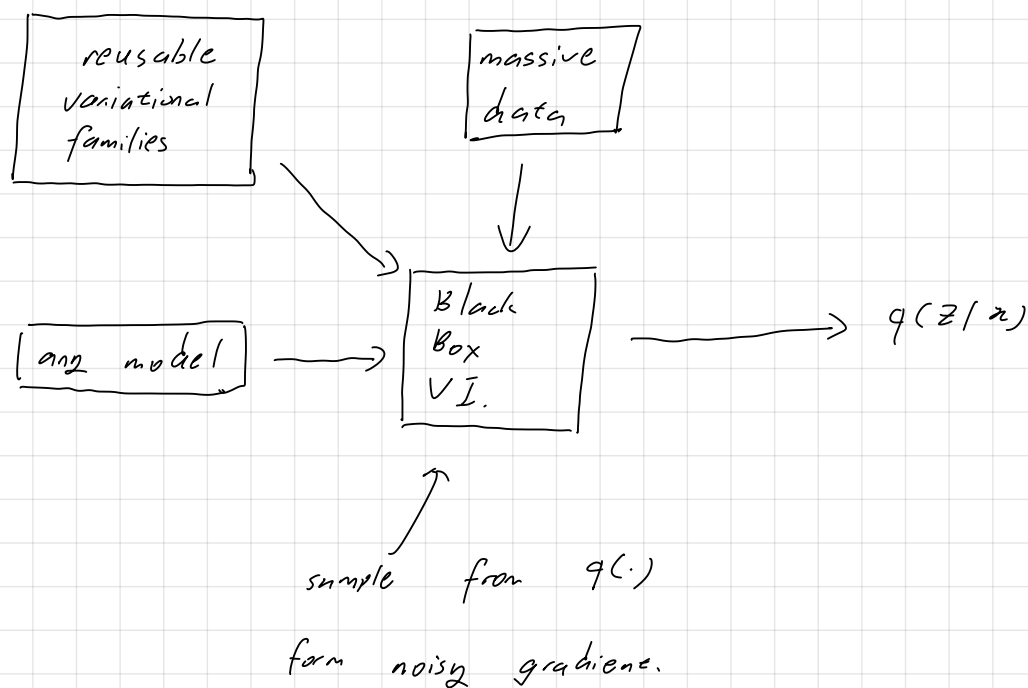
where we sample data example randomly
and update $q(z_i)$.

update global $q(\theta_k)$ with natural gradient ascent.

Black-box VI

limitation: we have to define update rules

How can we use VI with any model:

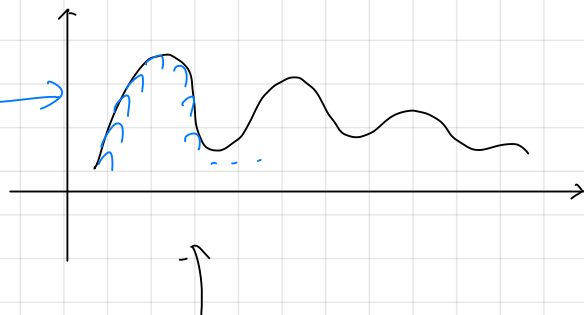


• Variational distribution $q_\lambda(z|x)$ with parameter λ .

Gaussian mixture distribution as universal approximator

ELBO:

$$\mathcal{L}(\lambda) = \mathbb{E}_{q(z|\lambda)} [\log p(x, z)] - \mathbb{E}_{q(z|\lambda)} [\log q(z|\lambda)].$$



Why deep neural network good?

\Rightarrow inductive bias

Bottleneck: approximate well only
with enough components.
impractical.

- compute exact gradient is not feasible.

$$L = E_{q_{\lambda}(z)} [f_{\lambda}(z)].$$



- Score function
- Reparametrization trick

$$\Rightarrow \text{Loss: } L = E_{q_{\lambda}(z)} [f_{\lambda}(z)] = \int q_{\lambda}(z) f_{\lambda}(z)$$

assume we can express the distribution $q_{\lambda}(z)$ with a transformation.

$$\begin{aligned} \epsilon &\sim s(\epsilon) \\ z &= t(\epsilon, \lambda) \end{aligned} \Leftrightarrow z \sim q(z|\lambda).$$

$$\text{ex) } \begin{aligned} \epsilon &\sim \text{Normal}(0, 1). \\ z &= \epsilon \sigma + \mu. \end{aligned} \Leftrightarrow z \sim \text{Normal}(\mu, \sigma^2).$$

\Rightarrow After reparameterization:

$$L = E_{\epsilon \sim s(\epsilon)} [f_{\lambda}(z(\epsilon, \lambda))].$$

\Downarrow chain rule.

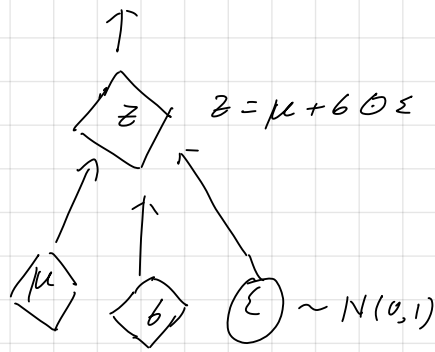
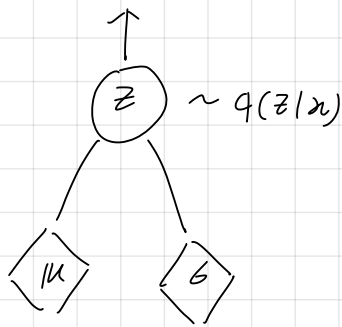
$$\text{so, gradient } \nabla_{\lambda} L = E_{\epsilon \sim s(\epsilon)} [\nabla_z f_{\lambda}(z) \nabla_{\lambda} t(\epsilon, \lambda)].$$

$$= \frac{1}{N} \sum_i \nabla_{z_i} f_{\lambda}(z_i) \nabla_{\lambda} z_i \quad \text{MC Simulation.}$$

Pro: low variance of gradient estimate

Con: Not all distribution can be reparameterized.

$$\nabla_{\lambda} L = E_{\epsilon \sim s(\epsilon)} [\nabla_z [\log p(x, z) - \log q(z)] \nabla_{\lambda} t(\epsilon, \lambda)].$$



VAE

- Variational inference
- Variational distribution parameterized as neural networks.
- reparameterization

Model: $p_{\theta}(x, z) = p_{\theta}(x|z)p(z)$.

↑ generative model
↑ prior (gaussian).



Assume variational distribution. $q_{\phi}(z|x)$.

a gaussian distribution parameterized
as deep neural networks.
(probabilistic encoder).

ELBO:

$$\mathcal{L}(\theta, \phi; x) = \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x, z)] + H(q_{\phi}(z|x)).$$

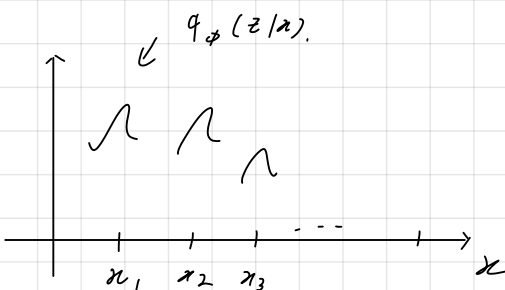
$$= \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x|z)] - \underbrace{KL(q_{\phi}(z|x) || p(z))}_{\text{divergence from prior.}}$$

both
gaussian
distribution.

↓
analytical
form

↑
generator

↓
reconstruction



• reparameterization:

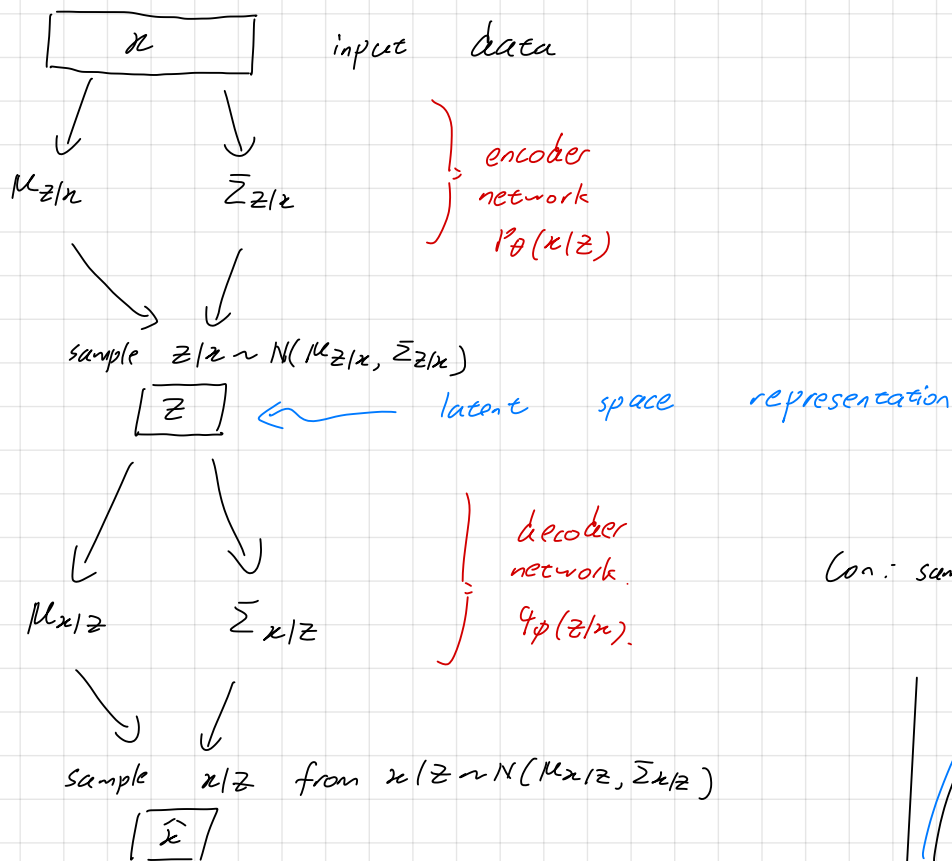
$$[\mu; \sigma] = f_{\phi}(x) \quad (\text{a neural network}).$$

$$z = \mu + \sigma \epsilon, \quad \epsilon \sim N(0, 1).$$

$$\nabla_{\phi} L = \mathbb{E}_{\epsilon \sim N(0, 1)} [\nabla_z [\log p_{\theta}(x, z) - \log q_{\phi}(z|x)] \nabla_{\phi} z(\epsilon, \phi)].$$

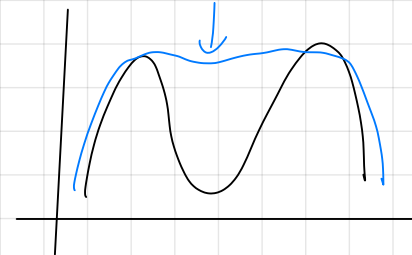
$$\nabla_{\theta} L = \mathbb{E}_{q_{\phi}(z|x)} [\nabla_{\theta} \log p_{\theta}(x, z)].$$

Networks:

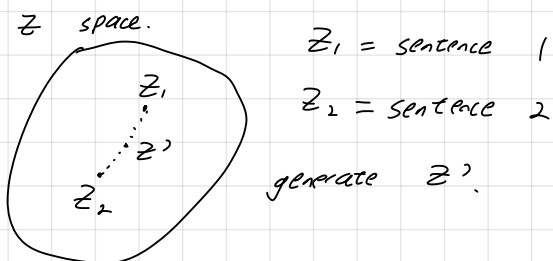


Con.: sample blurrier and lower quality

$$KL(p_{\text{data}} || p_{\theta})$$



VAE for text: multinomial gaussian



Amortized Variational Inference:

Variational distribution as an inference model $q_\phi(z|x)$.



- Amortized the cost of inference by learning a single data-dependent inference model.
↓
for quick inference.

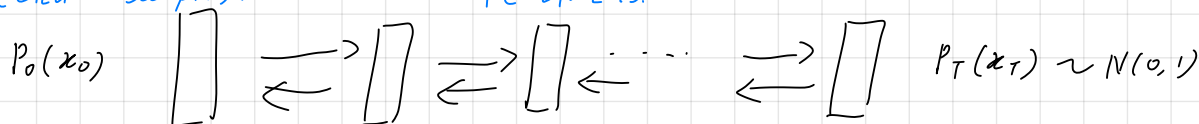
Diffusion Model

sample data $p(x_0) \rightarrow$ turn to noise

(Clean sample).

$q(x_t|x_{t-1})$

(Pure noise).



reverse / denoising process

turn into data \leftarrow sample noise $p_T(x_T)$

$p_\phi(x_{t-1}|x_t)$

- Better than VAE b/c
- decompose into many steps (2000 steps)
 - denoising as supervised learning (ground true data exists).
 - ...

PSC 190 PI.

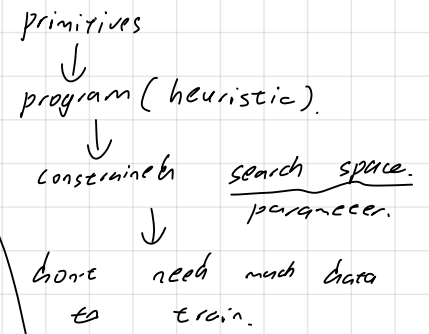
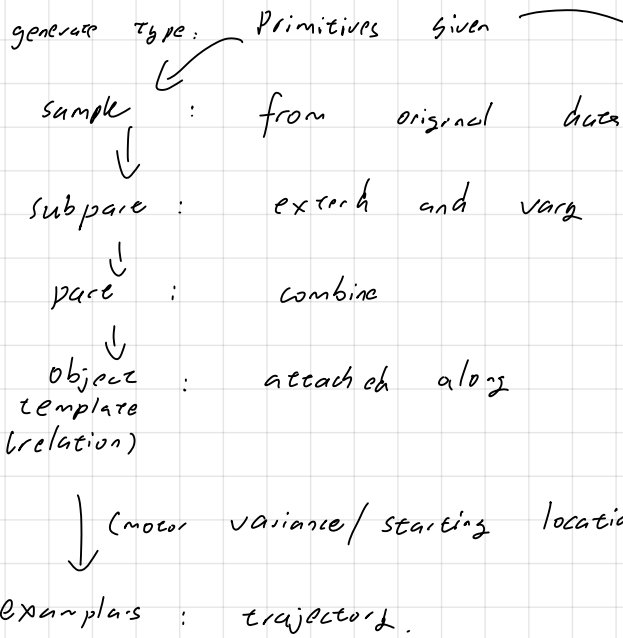
Concept learning through probabilistic program induction.

experience: massive data examples.

Human-level concept learning require much less data.
richer representation.

Bayesian Program Induction \rightarrow to mimic human-level learning.

(\hookrightarrow) The program: a generative model of handwritten characters.



Bayesian inference: $\text{program}^* = \underset{\text{program}}{\text{argmax}} P(\text{program} | \text{data}).$

posterior

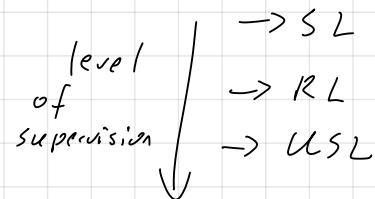
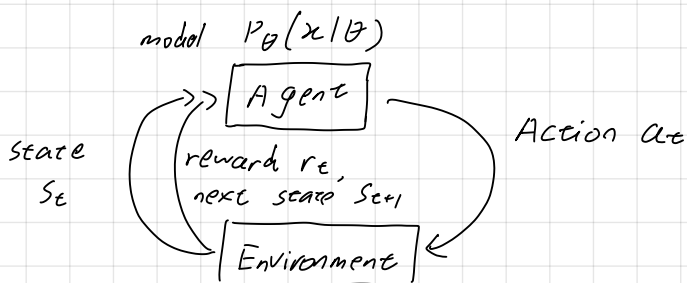
constrained searching:

$$= \underset{\text{program}}{\text{argmax}} \frac{P(\text{data} | \text{program}) \cdot P(\text{program})}{P(\text{data})}$$

Reinforcement Learning → Post-training for LLM

RL HF

take actions that maximize rewards.



ex)

state: angle & position

Action: torque applied on joints

Reward: 1 at each time step upright & forward movement.

sparse reward

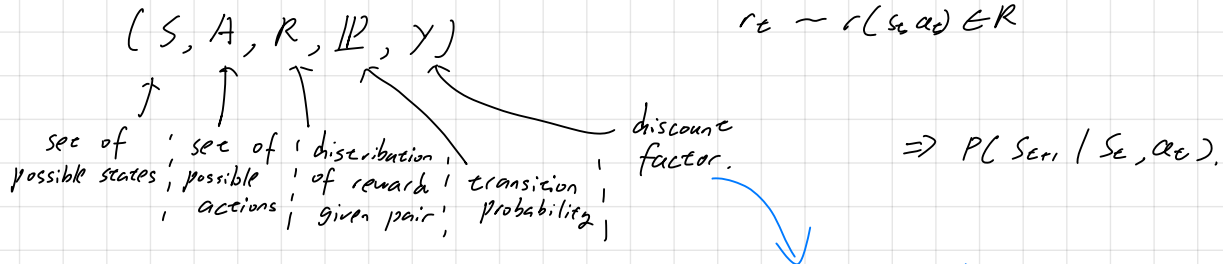
↓

Reward shaping

MDP

$$S_t \rightarrow a_t$$

markov property: current state completely characterises the state of the world.



At $t=0$, initial state $S_0 \sim P(S_0).$

to simulate uncertainty in the future for reward.
(exponential decay).

- action a_t
- sample reward $r_t \sim R(\cdot | s_t, a_t)$
- sample next state $S_{t+1} \sim P(\cdot | s_t, a_t)$
- receive reward r_t and next state S_{t+1}

Policy π is a function from S to A that specifies what action to take in each state.

Objective: find policy π^* that maximizes cumulative discounted reward

example)

Grid-world: action {right, left, up, down}
rewards from one grid to the other
negative "reward" for each transition.

$$\sum_{t=0}^{\infty} \gamma^t r_t.$$

Want to find optimal policy π^* , how to handle randomness

Maximizes expected sum of rewards

$$\pi^* = \underset{\pi}{\operatorname{argmax}} E \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid \pi \right] \quad \text{with} \quad \begin{aligned} S_0 &\sim P(S_0) \\ a_t &\sim \pi(\cdot | s_t) \\ S_{t+1} &\sim P(\cdot | s_t, a_t). \end{aligned}$$

Following a policy results in sample trajectory: $S_0, a_0, r_0, S_1, a_1, r_1, \dots$

Value Function

at state s , the expected cumulative reward from following policy from state s .

$$V^\pi(s) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s, \pi \right].$$

↑
assume following a policy.

Q-value function: (value func for state-action pair).

$$Q^\pi(s, a) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s, a_0 = a, \pi \right]. \quad (\text{future expected cumulative reward})$$

expanded form:

$$Q^\pi(s, a) = \mathbb{E} \left[\underbrace{r_0}_{\substack{\pi(a|s) \\ p(s'|a,s)}} + \gamma \underbrace{r_1}_{\substack{s_{t+1} \\ a_{t+1}}} + \gamma^2 r_2 + \dots \right]$$

↑ ↑
 another Q-value

taking action a in state s ,
the following the policy

$$= \mathbb{E} \left[r_t + \gamma Q^\pi(s_{t+1}, a_{t+1}) \right]$$

$\pi(a|s)$
 $p(s'|a,s)$

Bellman Equation → also follows recursive fashion

optimal Q-value function Q^* maximum expected cumulative reward

$$Q^*(s, a) = \max_{\pi} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s, a_0 = a, \pi \right]$$

satisfies

$$Q^*(s, a) = \mathbb{E}_{s' \sim \mathcal{E}} \left[r + \gamma \max_{a'} Q^*(s', a') \mid s, a \right].$$

↙ future value

intuition: if optimal state-action values for the next time-step $Q^*(s', a')$ are known,

then the optimal strategy is to take action that max expected value of

$$r + \gamma Q^*(s', a').$$

$$a^* = \arg \max_a Q^*(s, a)$$

← $\pi^*(s|a^*)$
take best action in any state as specified by Q^* ↑ optimal policy

Value iteration algorithm: use bellman equation as an iterative update:

$$Q_{i+1}(s, a) = E [r + \gamma \max_{a'} Q_i(s', a') | s, a]$$

Q_i converge to Q^* as $i \rightarrow \infty$.

Problem:

Not scalable. must compute all state-action pair.

small scale \rightarrow Table-based value iteration.

\downarrow

	a_1	a_2
s_1
s_2
s_3

Solution: a function approximator to estimate $Q(s, a)$
(neural network)

$$Q^\theta(s, a) \rightarrow Q$$

Deep Q network (Alpha Go).

Q-learning

$$Q(s, a; \theta) \approx Q^*(s, a)$$

function approximator to estimate action-value function.

Forward pass:

experience database

surprised task.

$$\text{loss function: } L(\theta_i) = E_{s, a \sim p(s, a)} [(y_i - Q(s, a; \theta_i))^2]$$

$$\text{where } y_i = E_{s', r \sim \mathcal{L}} [r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a]$$

bellman eq as ground true

Backward Pass:

learns a function to satisfies
bellman equation.

Gradient update:

$$\nabla_{\theta_i} L_i(\theta_i) = \mathbb{E}_{s, a \sim p(\cdot); s' \sim \pi} [r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) - Q(s, a; \theta_i)] \nabla_{\theta_i} Q(s, a; \theta_i)$$

ex) Atari Game:

$Q(s, a, \theta)$.

FC-4 (Q-val) \rightarrow 4 actions

\uparrow

FC-256

\uparrow

32 4×4 conv

\uparrow

16 8×8 conv

\uparrow

game screen

(all states)

(where if actions is continuous)

ex) control

\downarrow

pass in action as well

\downarrow

enumerate all possible a for all state
(huge space to search).

Training Q-network trick

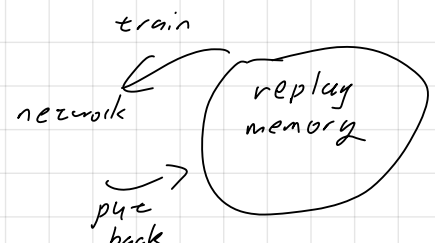
• Experience Replay

learning from batches of consecutive samples is problematic.

- sample are correlated
- bad feedback loop.

\Rightarrow - continually update a replay memory table of transitions

- train Q network on random minibatches of transitions from replay memory



Algorithm 1 Deep Q-learning with Experience Replay

Initialize replay memory \mathcal{D} to capacity N
Initialize action-value function Q with random weights *replay memory + Q network*
for episode = 1, M **do**
 Initialize sequence $s_1 = \{x_1\}$ and preprocessed sequence $\phi_1 = \phi(s_1) \leftarrow$ initialize state *initialize state*
 for $t = 1, T$ **do**
 With probability ϵ select a random action $a_t \leftarrow$ exploration/exploitation.
 otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$ *greedy that maximize Q value.*
 execute action { Execute action a_t in emulator and observe reward r_t and image x_{t+1}
 Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$
 Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in $\mathcal{D} \rightarrow$ store data
 Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from $\mathcal{D} \rightarrow$ sample from experience replay and do G.D.
 Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$
 Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3
 end for
end for

Derive policy : $a^* = \underset{a}{\operatorname{argmax}} Q_{\theta}(s, a).$

Policy Gradient :

- on-policy
- policy-based

define a class of parameterized policies: $\Pi = \{\pi_{\theta}, \theta \in \mathbb{R}^m\}.$

for each policy, define its value.

$$J(\theta) = E \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid \pi_{\theta} \right]$$

policy based

Want to find optimal policy $\theta^* = \underset{\theta}{\operatorname{argmax}} J(\theta)$

Gradient descent on policy parameter.

$$\Rightarrow J(\theta) = E_{\tau \sim p(\tau; \theta)} [r(\tau)]$$

$$= \int_{\tau} \underbrace{r(\tau)}_{\text{reward of trajectory}} p(\tau; \theta) d\tau$$

$$\nabla_{\theta} J(\theta) = \int_{\tau} r(\tau) \nabla_{\theta} p(\tau; \theta) d\tau.$$

use trick:

$$\nabla_{\theta} p(\tau; \theta) = p(\tau; \theta) \frac{\nabla_{\theta} p(\tau; \theta)}{p(\tau; \theta)} = p(\tau; \theta) \nabla_{\theta} \log p(\tau; \theta).$$

$$\nabla_{\theta} J(\theta) = \int_{\tau} (r(\tau) \nabla_{\theta} \log p(\tau; \theta)) p(\tau; \theta) d\tau.$$

$$= \underbrace{E_{\tau \sim p(\tau; \theta)}}_{\text{policy based}} \left[\underbrace{r(\tau) \nabla_{\theta} \log p(\tau; \theta)}_{\text{estimate with Monte Carlo sampling}} \right]$$

$$p(\tau; \theta) = \prod_{t \geq 0} p(s_{t+1} | s_t, a_t) \pi_{\theta}(a_t | s_t)$$

$$\log p(\tau; \theta) = \sum_{t \geq 0} \log p(s_{t+1} | s_t, a_t) + \log \pi_{\theta}(a_t | s_t).$$

$$\nabla_{\theta} \log p(\tau; \theta) = \sum_{t \geq 0} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \quad (\text{does depend on transition probabilities}).$$

Intuition:

if $r(\tau)$ is high, push up probabilities of actions seen.

if $r(\tau)$ is low, push down probabilities of actions seen.

In expectation, if a trajectory is good, then all its actions were good; then it averages out.

RL for LLM

(Autoregressive) text generation model:

Sequence $y = (y_0 \dots y_T)$

↑
trajectory, $\tau = s_0, a_0, r_0, s_1, a_1, \dots$

← logits

$$\Pi_{\theta}(y_t | y_{<t}) = \text{softmax}(f_{\theta}(y_t | y_{<t})).$$

↑
action

↑
state

policy $\Pi_{\theta}(a_t | s_t)$

• Reward $r_t = r(s_t, a_t)$

↓
sparse, $r_t = 0$ for $t < T$. (until sequence ends).

• General RL objective: maximum cumulative reward $J(\pi) = \mathbb{E}_{\tau \sim \pi} \left[\sum_{t=0}^T \gamma^t r_t \right]$.

• Q function: expected future reward of taking action a_t in state s_t .

$$Q^{\pi}(s_t, a_t) = \mathbb{E}_{\pi} \left[\sum_{t'=t}^T \gamma^{t'} r_{t'} \mid s_t, a_t \right]$$

GPT3.5 → ChatGPT:

(supervised fine-tuning & RLHF).

1. collect demonstration data & train a supervised policy

① prompt samples from prompt dataset.

② labeler demonstrate desired output

③ data used to fine-tune LLM (supervised learning).

2. collect comparison data & train a reward model.

↓
LLM output

labeler ranks output → used to train our reward model.

3. Optimize policy against reward model.

new prompt sampled from dataset.

↓
policy generates output

↓
reward model calculates a reward for the output.

update policy using PPO/bpo...