

CNN idea: treat convolutional matrix as params to be learnt.

$$W_{in} = \text{input width} \quad W_{out} = \left\lfloor \frac{W_{in} - K_w + 2 \times P}{S} \right\rfloor + 1$$

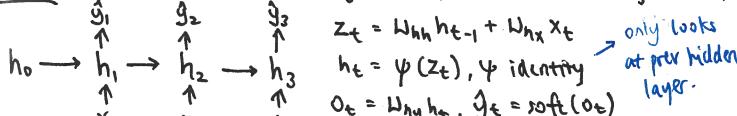
K_w = kernel width

P = padding

S = stride

hyperparameters. detect same features at diff positions.
parameter sharing (using same set of weights)

RNN $x_t \in \mathbb{R}^3$, $h_t \in \mathbb{R}^4$, $y_t \in \mathbb{R}^2$, $W_{hx} \in \mathbb{R}^{4 \times 3}$, $W_{hy} \in \mathbb{R}^{2 \times 4}$, $W_{hh} \in \mathbb{R}^{4 \times 4}$



$$\frac{\partial J}{\partial a_t} = y_t - \hat{y}_t \quad \left(\frac{\partial y_t}{\partial a_t} = \hat{y}_t (\mathbb{I}[t=t] - \hat{y}_t) \right), \quad \frac{\partial J}{\partial g_i} = -\frac{\partial y_t}{\partial g_i} \Rightarrow \frac{\partial J}{\partial a_t} = \sum_i \frac{\partial y_t}{\partial g_i} \frac{\partial g_i}{\partial a_t}$$

$$\frac{\partial J}{\partial h_i} = \sum_{t=i}^T (W_{hh})^{t-i} \quad W_{hy} \frac{\partial J}{\partial a_t}, \quad \frac{\partial J}{\partial W_{hh}} = \sum_{i=1}^T \frac{\partial J}{\partial h_i} \frac{\partial h_i}{\partial W_{hh}} \rightarrow h_{i-1} + \sum_{j=1}^{i-1} W_{hh} h_{i-1-j}$$

greedy decoding: linear in max pattern select edge is lower neg. log prob.

seg info, incor info from prev step to current & prev steps of var length.

RNN LM: convert all prev words to fixed len vector & def dist. $p(W_t | f_{W_{t-1}, \dots, W_1})$

$\hookrightarrow p(W_1, \dots, W_T) = p(W_1|h_1) \dots p(W_T|h_T)$ choose next word according to structures designed to

changes in dist of return activation due to update in return neurons during training.

problem: ↑ network depth, ↑ internal cov. shift (small change in layer layers → higher)

fix: layer normalization. $a \rightarrow [b = f(a)] \rightarrow b$ to $a \rightarrow [b' = f(a)] \rightarrow b = b' + a$

problem: perf. degradation → reduced connections. making e.g. $(0 \dots 0)$ to maintain causality

$$S_{i,j} = \frac{K_T \cdot q_i}{N \cdot k} \quad \text{attention: query} = XW_q, \text{ key} = XW_k, \text{ value} = XW_v$$

transformer layer score $S = \text{softmax}(A^T / \sqrt{d_k})$, attention $A = \text{softmax}(S)$

output $X' = A V$. applied row-wise

looks back at hidden vector of curr & prev timestep. $\hookrightarrow X'_4 = \sum_{j=1}^4 a_{4j} v_j$ & if x_i, x_j swap, a_{ij}, a_{ji} swap too

+ tokenization * cache keys & values.

batching: padding and truncation. discard query/similarity score/attn weights

complexity: $O(N^2 d_k)$ → each head pay attention to individual entries.

training transformer LM: 1 sentence = 1 training ex. $J = \log p(W_i) = \sum_i \log p(W_i | h_i)$

Kernel: a matrix/tensor of weights that is slid over image tensor. at each position (dictated by hyperparams such as padding/kernel size), perform element-wise mult. of kernel with underlying part of img. tensor, sum resulting

stride defines how many pixels the kernel moves at each step as it passes over the image. stride, sum resulting

mult. of image tensor. ↑ stride reduces computational power but lose more info (livelock updrd of model's accuracy)

↑ stride preserves size of output img. & can id more fine grained features, but ↑ computational cost.

padding surrounds image tensor to make all of zeros, by adding appropriate amt, we can ensure output shape = input shape

and allow every pixel to be included in convolution. also, padding helps focus on corner just as much as middle

downsampling reduce output dimensionsality (e.g. maxpool), appropriate for (pair our corner mult. then not just one).

larger images will be used to reduce the/marginally overfitting by making input to future convolutions less complex.

upsampling increase output dimensionsality when we want output = curr to be larger (e.g. to match input img in dim)

leg. assign label to each pixel of output & match it to input).

RL $\rightarrow p(s'|s, a)$ assumes s' only determined by s & a

$MDP = (S, A, T, R, \gamma, S_0)$, $T: S \times A \times S \rightarrow [0, 1]$, $R: S \times A \times S \rightarrow \mathbb{R}$.

$\hookrightarrow s_0 \xrightarrow{a_1} s_1 \xrightarrow{a_2} \dots, r_t = R(s_t, a_t, s_{t+1})$, then discounted payoff = $r_0 + \gamma r_1 + \gamma^2 r_2 + \dots$

policy $\Pi: S \rightarrow A$, optimal value function $V^*(s) = \max_{a \in A} \sum_{s' \in S} p(s'|s, a) (R(s, a, s') + \gamma V^*(s'))$

expected discounted future value of taking action a in state s s.t. $V^*(s) = \max_{a \in A} Q^*(s, a)$.

- exploration: take action to visit (state, action) pairs never unseen. \rightarrow no better than random.

- exploitation: take action to visit (s, a) pairs known to have high reward. \rightarrow never choose optimal.

learning $V^*(s)$: fixed pt. iterations on Bellman eq. $V^*(s) = E[r + \gamma V^*(s')]$. $V_{\text{opt}}^*(s) \approx r + \gamma V_{\text{opt}}^*(s')$.

similarly, $Q^*(s_t, a_t) \approx r_t + \gamma Q^*(s_{t+1}, a_{t+1})$, then switch to $\Pi_{\text{opt}}(s) = \arg \max_a Q^*(s, a)$ SARSA.

↳ good if Q-tables in lookup table (policy improvement guaranteed to Δ number). \hookrightarrow O($|A|^{|S|} \cdot |S|^2$) \rightarrow value iteration (instead)

↳ too aggressive if NN used to estimate Q-val: find policy params θ that maximizes $E_{\pi_{\theta}}[r_t + \gamma r_{t+1}]$

Policy Gradient Theorem: observe trajectory $(s_1^n, a_1^n, r_1^n, \dots, s_T^n, a_T^n, r_T^n)$ by following $\pi_{\theta}(a_t^n | s_t^n)$. depends on future traj.

$Q_{\theta}^n = \sum_{t=1}^T r_t^n$ (empirical total reward starting from step t) $d = \# \text{params in policy, so } \theta \in \mathbb{R}^d$ (per iteration)

$U_{\theta}^n = \frac{d}{dt} \ln \pi_{\theta}(a_t^n | s_t^n) \in \mathbb{R}^d$ (depends on action) (action score vector, from autodiff) even if env/state is not/partially observable.

$g^n = \sum_{t=1}^T Q_{\theta}^n U_{\theta}^n \in \mathbb{R}^d$ gradient estimate \rightarrow unbiased estimate of $\frac{d}{d\theta} J(\theta)$.

Monte Carlo algo: $\theta^{n+1} \leftarrow \theta^n + \eta g^n$, g^n from PGT. (rather than stochastic estimate $\frac{1}{m} \sum_{i=1}^m \frac{d}{d\theta} J(\theta)|_{\theta=\theta_i}$).

↳ failure modes: (1) cancellation (low SNR): total return θ can scale to horizon, and vary a lot due to randomness in policy/env, overall grad can be small (terms w opp signs).

(2) getting stuck: when may get close to boundary of simplex, store vectors for unlikely actions

$\Pi_{\theta} \rightarrow Q_{\theta}^n$ get large, prob of x very likely then get small, in the limit, 0. \rightarrow this is generic (using θ^n). usually not implementable.

Actor Critic: $g_{AC} = \sum_{t=1}^T \frac{1}{\pi_{\theta}(a_t | s_t)} U_{\theta}^n$ $\rightarrow E[\pi_{\theta}(a_t | s_t, a_t)]$: conditional expectation to dec. variance.

$\hookrightarrow E[U_t | s_t] = 0$, $g = \sum_t (\pi_{\theta} - V_{\phi}(s_t)) U_t$: baseline: no bias, but higher variance (compared to Q_{θ}^n of AC).

Advantage Actor Critic: $A^*(s, a) = Q^*(s, a) - V^*(s)$, $g = \sum_t A_t \nabla_{\theta} \ln \pi_{\theta}(a_t | s_t) \rightarrow s = \delta(s, a)$

* Q-learning, online learning/deterministic transition $[Q(s, a) \leftarrow r + \gamma \max_{a'} Q(s', a')] \rightarrow s' \sim p(s' | s, a)$

[non-deterministic/temporal difference error update] $Q(s, a) \leftarrow (1-\alpha) Q(s, a) + \alpha (r + \gamma \max_{a'} Q(s', a'))$ \rightarrow gradient param.

\hookrightarrow take greedy action $\pi_{\theta} \rightarrow$ random action o.w. learning rate α (trust param)

N-step returns $G_t^{(n)} = \sum_{k=t}^{t+n-1} \gamma^k r_{t+k} + \gamma^n V_{\phi}(s_{t+n})$, $A(s_t, a_t) \approx G_t^{(n)} - V_{\phi}(s_t)$ \rightarrow else $\pi_{\theta} = r_t + \gamma V(S_{t+1})$ (adv. based on actual obs. reward)

↳ (1) ↑ n to reduce bias, increase variance, (2) fast backprop (n steps instead of 1), (3) better grad quality

few shot learning: $D = \{(x_i, y_i)\}_{i=1}^n$ with small N & k . large pre-trained model.

A. supervised fine-tuning: use standard supervised obj, backprop to compute grad + adam optimizer, req. model weights

B. feed train ex as prompt, allow LM to infer patterns in many existing utterances // transformer LM req. $O(N^2)$ time/space, $N = \# \text{context}$

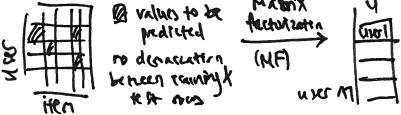
Q-learning: model-free/rrq. $\pi(s, a) \rightarrow Q(s, a)$ or $\pi(s, a)$, but only guaranteed to conv & erg. state visited no times/

Value-iteration: iteratively update $V^*(s)$ (M) then derive optimal policy $\Pi^* = \arg \max_{\pi} \sum_{s' \in S} p(s'|s, a) (R(s, a, s') + \gamma V^*(s'))$ \rightarrow slow conv, req. exp. strategy.

Policy-iteration: update $V^*(s) \leftarrow \text{Thres}^*(s)$, conv. faster than value iteration. \rightarrow guarantees conv to opt. policy, req. full world

actor-critic: actor update Π based on critic, critic estimates V , update rules use policy gradient. more sample efficient.

Recommender Systems

- Content filtering: analyzes item attributes (features)/metadata/keywords/categories) & match them to user known preference, ex. Movie (Genre proj. assumes access to ride info) ↗ doesn't work on new items w/o ratings.
- Collaborative filtering: analyzes user-item interactions patterns (rating)/(clicks) to find similar users/items, e.g. Netflix (incl. latent factor & neighborhood methods).
 - Low-rank factorization seek $U, V \in \mathbb{R}^{n \times k}$ s.t. $\text{rank}(UV) \leq k$.
 - if R has rank $> k$, cannot be represented exactly: minimize residual $E = R - UV^T$
 - 
 - ↳ values to be predicted → matrix factorization (MF)
 - ↳ no de-correlation between rows & cols → user M
 - ↳ each col is 1 latent factor (rank k) → $r_{ij} = \theta_i + \beta_j + u_i^T v_j$
 - ↳ U, V same #cols → $U^T = U^{-1} = U^{-1} = \frac{1}{\lambda} I_n$, $V = (\frac{\lambda}{\lambda - \beta})^{-1} V$
 - ↳ $J_{\lambda}(UV) = \frac{1}{2} \|Uv_j - U^T v_j\|^2 + \lambda \|U\|_F^2 + \|V\|_F^2$
 - ↳ $\nabla_{U,V} J_{\lambda}(UV) = -\epsilon_{ij} v_j + \lambda U_i$ (from v_j)

unconstrained MF: [SGD] sample $(i, j) \sim Z$, compute $e_{ij} = r_{ij} - U_i^T v_j$, $U_i \leftarrow U_i - \eta \nabla_{U_i} J_{\lambda}(UV)$, sim for \tilde{v}_j , [Alternating Least Squares] $U = \arg \min_{U \in \mathbb{R}^{n \times k}} J(U, V) \rightarrow \sum_{i,j} e_{ij}^2 (U_i^T v_j)^2$ is non-convex, so use iterative for each of U, V , keeping other fixed, Singular Value Decomposition: $R = U \Sigma V^T$, cols of U, V are orthonormal ($U^T U = I$), $R = \sum_{i=1}^k \sigma_i u_i v_i^T$ → rank k , $V = P_k$ constrained optimisation [nonneg/no dislike] $UV = \arg \min_{U, V} \frac{1}{2} \|R - UV^T\|_F^2$, train fit on D , test on D' , matrix to R under L_2

- Ensemble Methods** ① bootstrap: resample D from D to represent D' thus, in expectation $\frac{1}{k}$ excluded features norm ↗ returns accuracy $\hat{a}(t) = \arg \max_{a \in \text{obs}} \text{accuracy}_a$, model ensemble. The $\hat{a}_{i,t}$ to predict for obs x_i , return $\frac{1}{k} \sum_{t=1}^k h_t(x_i)$ ↗ T: hyperparam (t , more ex., finer quantifications of uncertainty) ↗ diversity of set of learned classifiers ↗ bagging = bootstrapping. ↗ diversity helps to cover posterior ↗ change hypothesis class for each bag: random forest: ↗ uses unweighted for every leaf, or each split — split feature randomisation ↗ bagging + col resampling + small DTs. ↗ sum.

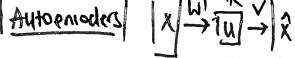
- ② column subsampling: use only subset of features for each classifier ↗ if base classifier H_t , can resample in beginning of training, ↗ classify above chance perf.

Boosting: if weak learner can get error = $\frac{1}{2} - \epsilon$, final boosted classifier can get 0 train & test on test (ϵ assumption) ↗ weighted rule $\sum_{t=1}^T \alpha_t h_t(x_i) \geq 0$ ↗ A: col corresponds to h_t , row to training ex $x^{(i)}$, $\bar{a}^{(i)} = (h_t(x^{(i)}))_{t=1}^T$, $\bar{x} \cdot \bar{a} \rightarrow$ try to optimize loss that favors tree margins, $y^{(i)} \bar{a}^{(i)} \alpha > 0$ ↗ H1 ↗ perceptron/SGD: $\alpha < 0$, $|H|$ big.

AdaBoost: initialize datapoint weights $w_0^{(i)} = \frac{1}{N}$, for $t=1, \dots, T$, train weak learner H_t by minimizing weighted (E_t)

$$E_t = \sum_{i=1}^N w_0^{(i)} \mathbb{I}(y^{(i)} \neq h_t(\bar{x}^{(i)})), \text{ rate weight of } h_t = \alpha_t = \frac{1 - \epsilon_t}{\epsilon_t}, w_t^{(i)} = \frac{w_0^{(i)}}{Z_t} e^{\alpha_t h_t(\bar{x}^{(i)})}, \bar{y} = \text{sgn} \left(\sum_t \alpha_t h_t(\bar{x}^{(i)}) \right)$$

SPS to fix \bar{y} , train till all margins ≥ 0 , up to $1 - \delta$, true error of final rule \leq boosting loss + $O \left(\frac{A + \sqrt{N} \gamma \delta}{\sqrt{N}} \right)$.

Autoencoders  $\min_{W, V} \sum_{i=1}^N \|W V \bar{x}^{(i)} - \bar{x}^{(i)}\|^2$. ↗ solve for best \bar{a}_i , keeps V fixed where $\sum_t \alpha_t \leq A$.

one training ex for each dim of $x^{(i)}$ ↗ $V_j = \text{feature vec for ex. } j$, $x_j^{(i)}$: target output for ex. j , $a^{(i)}$: linear regression weights ↗ minimize $\sum_{j=1}^d (x_j^{(i)} - V_j^T a^{(i)})^2$ min norm sol: $\bar{a}^{(i)} = W^T \bar{x}^{(i)}$, no to optimiz linear autoencoders, find W, V , then find \bar{a} as $\min \|W^T - x\|_F^2$ or $\sum_j (V_j^T a^{(i)} - x_j^{(i)})^2$ (collaborative filtering)

if $k=1$ hidden dim, sol is $V, V \in \mathbb{R}^{d \times d}$ and $W = \frac{V}{\sqrt{d}} V^T$, to avoid scaling ambiguity, $\|V\|_F = 1$, $W = V$,

Principal Component Analysis

$V_1 = \arg \min_{V \in \mathbb{R}^{n \times k}} \|V\|_F^2 \text{ s.t. } \|V\|_F = 1$ construct residuals $e_1^{(i)} = x^{(i)} - (V_1^T x^{(i)}) V_1$ ↗ in long dictinct equiv: the variance of $x^{(i)}$ maximises of V_1

1. the 1st principal component is the unit vector that minimizes mean-sq reconstruction error ↗ 2. 2nd PC is unit vector that minimizes mean-sq error of residuals while remaining orthogonal to V_1 , max $\min_{V \in \mathbb{R}^{n \times k}} \|V^T x^{(i)}\|^2$ ↗ k-th PC..

$V_2 = \arg \min_{V \in \mathbb{R}^{n \times k}} \sum_{i=1}^N \|e_1^{(i)} L(V e_1^{(i)}) V\|^2 \text{ s.t. } \|V\|_F = 1, V^T V_1 = 0, e_2^{(i)} = e_1^{(i)} - (V_2^T e_1^{(i)}) V_2$ ↗ ortho. V_1, \dots, V_{k-1}

j-th principal value = variance of proj. onto j-th PC, bc PCs orthogonal sum of $V = \text{var. of } X$ ↗ diagonal: var. of feature $x_j^{(i)}$ off-diag: cov. of pair of feature. $\lambda_j = \arg \min_{V \in \mathbb{R}^{n \times k}} \sum_{i=1}^N \|X - V\|_F^2 \rightarrow 1 - \text{tot. var.} \sum_{i=1}^N \|e_i^{(i)}\|^2$

→ discrete autoencoder $X \xrightarrow{VQ} \text{discrete hidden } Z \xrightarrow{V} \text{pred. = center w.r.t. } Z \rightarrow$ implements clustering

train to minimize $\text{MSE} = \frac{1}{N} \sum_{i=1}^N \|X^{(i)} - \bar{X}^{(i)}\|^2$, $\bar{X} = V_2 = V \times V Q(x, V)$, $V = Z V^T$, $V \in \mathbb{R}^{n \times k}$, $Z \in \mathbb{R}^{k \times 1}$ (1-hot)

[blocks word decoder] given feature vectors $D = \{x^{(1)}, \dots, x^{(N)}\}$, initialize matrix of centers V (each col V_j)

repeat: minimize wrt $Z = V_i$, s.t. $Z^{(i)}$ to map $x^{(i)}$ to closest center, then minimize wrt V_i , V_j : min MSE from V_j ($\arg \min_{V, Z} \sum_{i=1}^N \|X^{(i)} - V^{(i)}\|^2$) ↗ to its assigned pts. $V_j = \text{mean}\{x^{(i)} | Z^{(i)} \leq \delta\}$ ↗ equal sized Gaussian

Initialization: (1) random, (2) furthest pt. heuristic, (3) k-means++ ↗ Pr[each initial center in diff Gaussian]

1. may get stuck in poor local minima (if we didn't pick one cent per Gaussian = k-means output bad only) $\propto \frac{k!}{k^n} \propto \frac{1}{k^n}$

2. pick first cluster pt C_1 randomly, pick each subseq C_j s.t. it is as far as possible from C_1, \dots, C_{j-1} ok if data Gaussian

3. let $D(X)$ be distance between pt x & nearest center, choose next center b.p. proto $P(x)^2$. (if n. of random).

↳ THM: attains $O(\log k)$ approx to optimal k-means sol. in expectation. ↗ kNN supervised, k-means unsup

PCA: factorize $X \in \mathbb{R}^{n \times n}$ as $Z^T U \approx X$, $Z \in \mathbb{R}^{n \times m}$, $U \in \mathbb{R}^{m \times m}$, rows of X contain data, rows of U the PC prototypes

— removing rows from U/Z takes variance associated w/ PC. \Rightarrow inc. recon. error (var removed) $\bar{X} = Z^T U$.

— rows of U eigenvectors $\Rightarrow U^T U = I$. (bc. PCs unit vectors) ↗ can produce representation that has same # dimensions as input.

— goal of PCA to produce underlying structure to data that preserves largest amount of variance (output variable never provided)