

Linear Algebra

1.1.1. Least Squares

1.1.2. Matrix Factorization

1.1.3. Eigenvalues and Eigenvectors

1.1.4. Singular Value Decomposition

1.1.5. Positive Semidefiniteness

1.1.6. Fundamentals of Linear Algebra

1.2. Probability

1.3. PCA

1.4. Bias-Variance Tradeoff

1.5. Neighborhood Embedding (NE) / Isomap

1.6. t-SNE

1.7. Diagonalizing an Ellipse

Bayes classifier: attains minimum risk if a classifier can achieve it.

Bayes risk/error rate: lowest risk any classifier can attain.

$f^*(x) = \operatorname{argmin}_i E_{Y|X} [L(i, y) | X] = \operatorname{argmax}_i P(Y=i | X)$

equating the risks of 2 classes should help solve for Bayes classifier decision boundary eg:

$$[\lambda_1 \cdot \lambda_2] - \lambda_1 \cdot \text{loss for classifying } j \text{ as } i \\ \lambda_1 \cdot P(C_i | X) + \lambda_2 \cdot P(C_j | X) = \lambda_2 \cdot P(C_i | X) + \lambda_1 \cdot P(C_j | X)$$

$$\frac{P(C_i | X)}{P(C_j | X)} = \frac{\lambda_2}{\lambda_1}$$

Decision Theory

classifier/decision rule $f: X \rightarrow \{1, \dots, K\}$ where X is space of inputs & we have K classes

loss function $L(i, j)$: penalty classifier receives for predicting class i when true class is j

risk $R(f) = E_{X,Y} [L(f(x), y)] = E_x [E_{y|x} [L(f(x), y) | X]] = E_x [\sum_j L(f(x), j) P(Y=j | X)]$

$$= \int_X \sum_{j=1}^K L(f(x), j) P(Y=j | X) P(x) dx = \int_X \sum_{j=1}^K L(f(x), j) P(X|Y=j) P(Y=j) dx$$

Bayes classifier: attains minimum risk if a classifier can achieve it.

Bayes risk/error rate: lowest risk any classifier can attain.

$f^*(x) = \operatorname{argmin}_i E_{Y|X} [L(i, y) | X] = \operatorname{argmax}_i P(Y=i | X)$

equating the risks of 2 classes should help solve for Bayes classifier decision boundary eg:

$$[\lambda_1 \cdot \lambda_2] - \lambda_1 \cdot \text{loss for classifying } j \text{ as } i \\ \lambda_1 \cdot P(C_i | X) + \lambda_2 \cdot P(C_j | X) = \lambda_2 \cdot P(C_i | X) + \lambda_1 \cdot P(C_j | X)$$

$$\frac{P(C_i | X)}{P(C_j | X)} = \frac{\lambda_2}{\lambda_1}$$

FC NNs - more expressive & have greater model capacity due to their incr. # of weights (rel. to CNNs)

CNNs less weights & less data needed but need ↑ architecture params to work well

capacity of a model is a hyperparam that we choose based on the error on a validation set

If capacity < optimal \Rightarrow underfitting

convolution " " " " over global pooling: translation invariant self-attn: permutation equivariant GNN: aggregation + learning of the representation containing the aggregation set as input

representational containment: conv & attn Eqn pass be general msg passing allows arbitrarily learnable msg passing but conv & attn pass extra constraints translational equiv - useful for pixel/node-level tasks rotational invariance for graph & img-level tasks

Markov Models

$Q = q_1, q_2, \dots, q_n$: set of N states

$A = a_{ij}, a_{ij} \dots a_{nn}$: transition prob matrix A , each a_{ij} repr. $P(q_i \rightarrow q_j)$ s.t. $\sum_j a_{ij} = 1$

$\pi = \pi_1, \dots, \pi_N$: initial prob dist over states

HMMs have the above but also have $O = o_1, \dots, o_T$: sequence of T observations, each drawn from a vocab $V = v_1, v_2, \dots, v_n$

$B = b_{(c,o)}$: sequence of observation likelihoods/ emission probabilities $P(o_t \text{ generated from state } c)$

HMM probs:

- 1) Likelihood: given specified HMM, compute likelihood of observation sequence O
- 2) Decoding: given HMM, find best seq. of hidden states
- 3) Learning: given observation seq. O & set of states, learn HMM params $A \& B$

probabilistic graphical models

each node repr. random var. & edges repr. dependence DAGs help us achieve tractability thru cond. independence

$P(R=T | G=T) = P[G=T, R=T]$

$$P[G=T] = \frac{\sum_{x \in \mathcal{X}, f \in \mathcal{F}} P[G=T, S=x, R=T]}{\sum_{x \in \mathcal{X}, f \in \mathcal{F}} P[G=T, S=x]}$$

$P[G=T, S=T, R=T] = P[G=T | S=T, R=T] P[S=T | R=T] P[R=T]$

Joint Factorization P_S, P_R, P_G

$$P_S = P[X=x] P[Z=z] P[Y=y | X=x, Z=z] P[S=z | X=x, Y=y]$$

each node is conditionally indep. of all of its ancestor nodes, given all of its parents

$S \perp\!\!\!\perp Y \text{? No. } S \perp\!\!\!\perp X \perp\!\!\!\perp Y \text{? No. }$

Markov Decision Processes

$s_t \in S$: state at timestep t . $a_t \in A$ action at time t . r_t : discount factor

$R_{t+1} \in R$: reward generated from A_t ; S_{t+1} : state generated from A_t

one step dynamics: $P(S', r | s, a) = P[S_{t+1}=s' | S_t=s, A_t=a]$

return following time t : $G_t = R_{t+1} + \gamma^2 R_{t+2} + \gamma^4 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} = R_{t+1} + \gamma G_{t+1}$

$\Pi(a|s) = P(A=a | S_t=s)$

state-value func $V_\Pi(s)$ of state s under policy Π :

$$V_\Pi(s) = E_{\Pi}[G_t | S_t=s] = E_{\Pi}[R_{t+1} + \gamma G_{t+1} | S_t=s]$$

action-value func $q_\Pi(s, a)$ of taking action a in state s under policy Π starting at state s

$$q_\Pi(s, a) = E_{\Pi}[G_t + \gamma R_{t+1} | S_t=s, A_t=a] = E_{\Pi}[R_{t+1} + \gamma G_{t+1} | S_t=s, A_t=a]$$

$E_{\Pi}[R_{t+1} | S_t=s] = \sum_r r P(s', r | s, a) \Pi(a|s)$

$$q_\Pi(s, a) = \sum_r r V_\Pi(s')$$

$V_\Pi(s) = E_{\Pi}[R_{t+1} + \gamma V_\Pi(S_{t+1}) | S_t=s, A_t=a] = \sum_{s', r} P(s', r | s, a) [r + \gamma V_\Pi(s')]$

Bellman expectation eqn:

$$V_\Pi(s) = \sum_a \Pi(a|s) \sum_{s', r} P(s', r | s, a) [r + \gamma V_\Pi(s')]$$

optimal state-value func (Bellman optimality eqn)

$$V_\star(s) = \max_{\Pi} V_\Pi(s) = \max_{\Pi} \sum_{s', r} P(s', r | s, a) [r + \gamma V_\star(s')] = \max_a q_\star(s, a)$$

idea: optimal policy yields max expected total reward

Value iteration: compute optimal values of states by iterative updates until convergence [ie $V_{K+1}(s) = V_K(s)$]

- ① Vs, initialize $V_0(s) = 0$
- ② Vs Es, until convergence: $V_{K+1}(s) = \max_{a, s', r} P(s', r | s, a) [r + \gamma V_K(s')]$

Policy iteration: use policy evaluation & extraction to iteratively converge to optimal policy; usually outperforms value iteration bc policies usually converge faster than state's values

- ① Define initial policy (can be arbitrary)
- ② Until convergence: (ie $\Pi_{t+1} = \Pi_t$)
- A) Policy evaluation: $V_{\Pi_t}(s) = \sum_{s', r} P(s', r | s, a) [r + \gamma V_{\Pi_t}(s')]$
- B) Policy improvement: $\Pi'(s) = \operatorname{argmax}_a \sum_{s', r} P(s', r | s, a) [r + \gamma V_{\Pi_t}(s')]$

Decision Theory

classifier/decision rule $f: X \rightarrow \{1, \dots, K\}$ where X is space of inputs & we have K classes

loss function $L(i, j)$: penalty classifier receives for predicting class i when true class is j

risk $R(f) = E_{X,Y} [L(f(x), y)] = E_x [E_{y|x} [L(f(x), y) | X]] = E_x [\sum_j L(f(x), j) P(Y=j | X)]$

$$= \int_X \sum_{j=1}^K L(f(x), j) P(Y=j | X) P(x) dx = \int_X \sum_{j=1}^K L(f(x), j) P(X|Y=j) P(Y=j) dx$$

Bayes classifier: attains minimum risk if a classifier can achieve it.

Bayes risk/error rate: lowest risk any classifier can attain.

$f^*(x) = \operatorname{argmin}_i E_{Y|X} [L(i, y) | X] = \operatorname{argmax}_i P(Y=i | X)$

equating the risks of 2 classes should help solve for Bayes classifier decision boundary eg:

$$[\lambda_1 \cdot \lambda_2] - \lambda_1 \cdot \text{loss for classifying } j \text{ as } i \\ \lambda_1 \cdot P(C_i | X) + \lambda_2 \cdot P(C_j | X) = \lambda_2 \cdot P(C_i | X) + \lambda_1 \cdot P(C_j | X)$$

$$\frac{P(C_i | X)}{P(C_j | X)} = \frac{\lambda_2}{\lambda_1}$$

Perf. classifier

FP = false positive rate

TP = sensitivity

AUC: area under curve \rightarrow larger area = better model

ROC curve: T -> varying threshold of our classifier

GNNs

compute VV msgs across edges as hidden memory

might have more power dep. on complexity of parameterization of χ

gradient data ∇C can be used to regularize neighborhood structure

Graph Neural Networks

graph that can operate on arbitrary relational structures

structure \mathcal{G} : able to learn repr. of nodes that depend on structure of the graph

message passing used to learn these representations

$G = (V, E)$: node set V , edge set E , set of node feats χ

want to generate learned node embeddings (χ_u)

$\chi_u, u \in V$

during each message passing operation a hidden embedding $\chi_u^{(k)}$ is generated, representing the updated embedding of node $u \in V$ in the k iteration, based on the info aggregated from u 's graph neighborhood $N(u)$ (which can incl. u)

Message Passing Update:

$$\chi_u^{(k)} = \phi^{(k)} [\chi_u^{(k-1)}, \bigoplus \{\psi^{(k)}(\chi_v^{(k-1)}, \chi_u^{(k-1)}) \mid v \in N(u)\}]$$

update func $\phi^{(k)}$: aggregation (addit. non-parametric) message func (arbitrarily diff. able) (parameterized by NNs)

$$= \phi^{(k)} [\chi_u^{(k-1)}, \bigoplus \{m_{vu}^{(k)} \mid v \in N(u)\}]$$

message from sender node v to receiver u

$$m_{vu}^{(k)} = \phi^{(k)} (\chi_v^{(k-1)}, \chi_u^{(k-1)})$$

$$= \phi^{(k)} [\chi_u^{(k-1)}, m_u^{(k)}]$$

aggregated msg from neighbors of u

single GNN layer can have 1 or more rounds of msg passing

3 main blocks of msg passing layer: construction, aggregation, update

initial embeddings @ $k=0$ are set to input fts of each node $\chi_u^{(0)} = \chi_u, u \in V$

After K iterations of msg passing, $\chi_u^{(K)}$ encodes info abt u nodes in u 's K -hop neighborhood that is relevant for the training objective

if we run K iterations of msg passing in total, we can use output of final layer to define embedding for each node: $\chi_u = \chi_u^{(K)}, u \in V$

3 flavors of GNN layer

- ① Message passing
- ex: $\chi_u^{(k)} = \phi^{(k)} (W_{self}^{(k)} \chi_u^{(k-1)} + W_{neig} \sum_{v \in N(u)} \chi_v^{(k-1)} + b^{(k)})$

$\phi^{(k)}(h_{v,u}) = \phi(h_{v,u}^{(k-1)} + m_{v,u}^{(k)})$

$$m_{v,u}^{(k)} = \phi^{(k)} (m_{v,u}^{(k-1)} + \sum_{v \in N(u)} \chi_v^{(k-1)})$$

Note: if $W_i^{(k)}$ $\in \mathbb{R}^{d \times d}$ the # of learnable params in the MP update func above is $2 \cdot d \cdot m$

Convolutional

$$\chi_u^{(k)} = \phi(\chi_u^{(k-1)}, \bigoplus \{ (W_i^{(k)}, h_{i,u}^{(k-1)}, \forall v \in N(u)) \})$$

e.g. input $I \in \mathbb{R}^{L \times H \times W \times C}$ (eg image (L, H)) \rightarrow input $\chi \in \mathbb{R}^{L \times H \times W \times C \times d}$ (eg filter (L, H)) \rightarrow output $\chi^{(k)} \in \mathbb{R}^{L \times H \times W \times C \times d}$

msg only dep on sender node v but not u & weights don't dep on either u or v

Attentional

$$\chi_u^{(k)} = \phi(\chi_u^{(k-1)}, \bigoplus \{ (\alpha_i^{(k)}, h_{i,u}^{(k-1)}, \forall v \in N(u)) \})$$

$\alpha_i^{(k)}, h_{i,u}^{(k-1)}$ are attn. weights

suppose attn. weights are computed by single-head scaled dot-prod (self-attn): $\alpha_{i,j}^{(k)} = \frac{1}{\sqrt{d}} \text{softmax}(\text{dotprod}(h_{i,u}^{(k-1)}, h_{j,v}^{(k-1)}))$

weights as

$$\alpha_{i,j}^{(k)} = \exp \left\{ \frac{1}{\sqrt{d}} (W_{key} h_{i,u}^{(k-1)})(W_{query} h_{j,v}^{(k-1)})^\top \right\} / \sum_{i=1}^n \exp \left\{ \frac{1}{\sqrt{d}} (W_{key} h_{i,u}^{(k-1)})(W_{query} h_{j,v}^{(k-1)})^\top \right\}$$

to emulate self-attn more, could also multiply $h_{i,u}^{(k-1)}$ by weights that are ft-dep. & messages dep. on sender & receiver

