

## Standard search problem

**State space:** Number tiles in each cell position

**Initial state:** [8, -6, 5, 4, 7, 2, 3, 1]

**Actions:** Move tile {Left, Right, Up, Down}

**Transition Model:** Update tiles in current & target cell pos

**Path cost:** Number of moves

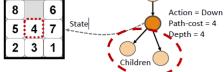
**Goal test:** Compare to positions in goal state

**State:** rep. of a physical configuration

**Node:** data struct, part of a search tree

- Comprising state, parent-node, child-node(s), action, path-cost, depth unlike a state.

**Solution:** Seq of act from init to goal state



**Expand function creates new nodes** (& their various fields) using the Act & Transition Model to create the corresponding states

**Search strategy:** picking the order of node expansion

**b:** Max branching factor of the search tree

**d:** depth of the least-cost solution

**m:** max depth of the state-space (inf possible)

## Tree Search:

- Offline, simulated exploration of state-space by expanding states via generation of successors of already-explored states

- However repeated states results in redundant paths that can cause a tractable problem to become intractable (inf loop, m = inf)

## Graph Search:

- Modified tree search to keep track of previously visited states (to not revisit) but a potentially large number of states to track

## Types of Search:

### a) Uninformed Search

No additional information about states beyond that in the problem definition (blind search)

### b) Informed Search

Uses problem-specific knowledge beyond the definition of the problem itself

### c) Adversarial Search

Used in a multi-agent environment where the agent needs to consider the actions of other agents and how they affect its own performance.

## BFS (FIFO, Graph):

Init: A

Expand & pop A: AB, AC

Expand & pop B: AC, ABD, ABE

Expand & pop C: ABD, ABE, ACF

Expand & pop D: ABE, ACF, ABDG, E visited

Expand & pop E: ACF, ABDG, ABEH, F visited

Expand & pop F: ABDG, ABEH, H visited

Expand & pop G: ABEH, ABDGX, H visited

Expand & pop X: X is goal state

Frontier empty, ret no solution

AC3: V1 -> V2: D1 = RB, D2 = RG

V2 -> V1: D2 = RG, D1 = RB

V2 -> V3: D2 = RG, D3 = G

V3 -> V2: D3 = G, D2 = R

V1 -> V3: D1 = RB, D3 = G

V3 -> V1: D3 = G, D1 = RB

(+) V1 -> V2: D1 = RB, D2 = R

(+) V3 -> V1: D3 = G, D1 = B

(+) V2 -> V1: D2 = R, D1 = B

Queue empty, all arc consistent, ret ...

Pure Backtracking: V1: R

V2: R (Inconsist. assign, bt)

V2: G

V3: G (Inconsist. assign., backtrack)

V1: B (V2 lv none left)

V2: R

V3: G, all var assigned, ret

V1: B, V2: R, V3: G

## BT with FC:

V1: R, remove any inconsist in

remaining var of neighbours

V2: G (FC)

V3: Empty, terminate

V1: B

V2: R

V3: G all var assigned, ret

V1: B, V2: R, V3: G

## Informed Searches

$g(n)$  ignores direction of the search, instead we use  $h(n)$  heuristic fn: est dist from node n to goal state (how close a state n is to the goal state)

$h(n)$  admissible if  $h(n) \leq h^*(n)$

$h^*(n)$  = true cost from node n to goal state (e.g. straight line distance  $\leq$  the actual road distance)

$h(n)$  consistent if  $h(n) \leq h(n') + c(n, a, n')$

$c(n, a, n')$  = cost of getting from node n to n'

$h_1(n)$  dominates  $h_2(n)$  if  $h_1(n) \geq h_2(n)$

Provided that both  $h(n)$  are admissible More dominant  $h(n)$  potentially explores lesser branches, hence better

Admissible heuristics can be derived from the exact solution cost of a relaxed version of the problem

e.g. Relaxed: Tiles can be directly moved to any square.  $h_1(n)$ : # of misplaced tiles OR Relaxed: Tiles can be moved to any adjacent square.  $h_2(n)$ : total Manhattan distance

Admissible heuristics can be derived from the solution cost of a sub-problem of the problem (e.g. Solve 4 tiles first)

## Greedy-BFS (Graph):

Init: A(7)

Expand & pop A: AB(4), AC(5)

Expand & pop B: ABD(3), AC(5)

Expand & pop X: X is goal state, ret ABX

A\*: Init: A(0+7=7)

Expand & pop A: AB(2+4=6), AC(2+5=7)

Expand & pop B: ABD(2+2+3=7), AC(7), ABX(2+6=8)

Expand & pop D: ABDX(2+2+1+0=5), ABDF(2+2+1+2=7), AC(7), ABX(8)

Expand & pop X: X is goal state, ret ABDX

## Early Detection:

### 1) Fwd Checking Algo

- Keep track of remaining lv for unassigned var, terminate search when any var has no lv  
- Propagates info from assigned to unassigned var only, only provide early detection for some failures. Need to repeatedly enforce constraints locally using constraint propagation.

### 2) AC-3 Algo $O(n^2d^3)$ order arcs no matter

- Start with queue of all directed arcs

### RemoveInconsistent(X, Y) :

Check for consistency ( $X > Y$  is consistent iff for every val of X there is some allowed y)

- If removed any, add the neighbours (into X) into queue

- Continue till queue is empty

$O(n^2)$ : Need to check all (potentially)  $n^2$  edges (directional)

$O(d^2)$ : For every edge, compare both domains

$O(d)$ : Each var change, repropagate to its neighbour at most max d times

## Variable Assignment Order:

### 1) Minimum remaining values

- Choose var with fewest legal values left (most constraint var, most likely to fail)

### If same MRV, Degree Heuristic:

- Choose var with the most constraint on remaining var

### 2) Least constraining value:

- Choose one that rules out the fewest values in the remaining var

## Constraint Satisfaction Problem:

**State:** Variables,  $X_i$  with values from domain,  $D_i$

**Goal Test:** Set of constraint,  $C_i$  specifying allowable combinations of values for subsets of var

(Nodes) Finite set of variables,  $X = \{X_1, \dots, X_n\}$

Non-empty domain D of k possible values for each variable,  $D_i = \{v_1, \dots, v_k\}$

(Edges) Finite set of constraints,  $C = \{C_1, \dots, C_m\}$  where they limit the values variables can take

**Complete assignment:** every var assigned a val

**Consistent assignment** meets all constraint

**CSP solution** = complete & consistent for all var

**Formal representation language** that can be used to formalize many problems types

Able to use **general-purpose solver**, which are generally more efficient than standard search.

- Constraints allow us to focus the search to valid branches, invalid branches removed

- Non-trivial to do this for standard search (need manual selection of actions)

## Varieties of CSP:

### o Discrete variables:

- **Finite domains:**  $O(d^n)$  complete assignments for n variables, domain size d (e.g. map-colouring, scheduling with time limits)

- **Infinite domains:** Int, str, etc. (e.g. job

scheduling (e.g. StartJob1 + 5 < StartJob3))

**o Conti. variables:** (e.g. start/end times of an obs)

### o Unary/Binary/Higher-order Constraints:

# of variables involved in constraint

### o Preference (soft) Constraints:

(e.g. red is better than green represented by some cost of each var, aka Constr. Opti Problem)

## Uninformed Searches

### Breadth First Search:

Expand shallowest unexpanded node, implement using (FIFO) Queue

- **Completeness:** Yes (if b is finite)

- **Optimality:** Yes (Not optimal in general)

- **Time complexity:**  $O(b^d)$

- **Space:** Same (keeps every node in mem)

Recall BFS checks level by level, if the step cost != 1/ not uniform, then we might miss out certain nodes with a cheaper step cost

### Uniform Cost Search: BFS w g(n)=1

Expand unexpanded node with lowest path cost g(n) implement using Queue ordered by g(n). If visited node, only replace if cost is lower

- **Completeness:** Yes (if step cost some positive constant else lead to loop pattern)

- **Optimality:** Yes

- **Time/Space complexity:**  $O(b^{C^*})$

$C^*$  is the cost of opt solution

Goal test after popping only, since there might be some shorter path to that node. Only do goal test when popping that node.

### Depth First Search:

Expand deepest unexpanded node, implement using (LIFO) Queue

- **Completeness:** No (if m inf, keep going down inf path)

- **Optimality:** No (might be searching down non-optimal path)

- **Time complexity:**  $O(b^m)$  terrible if m>>d, as need to back track all the way back up

- **Space:**  $O(bm)$  linear , store m depth, and at each level, keep track of b child nodes



### DFS (LIFO, Graph, insert lowest alphabetical order 1st):

Init: A

Expand & pop A: AB, AC

Expand & pop C: AB, ACF

Expand & pop F: AB, ACFE, ACFH

Expand & pop H: AB, ACFE, ACFHG, ACFHX

Expand & pop X: X is goal state, ret ACFHX

### DLS (with I = 2):

Init: A

Expand & pop A: AB, AC

Expand & pop C: AB, ACF

Expand & pop F: AB, ACFE, ACFH

Expand & pop H: AB, ACFE, ACFHG, ACFHX

Expand & pop X: X is goal state, ret ACFHX

DLS (with I = 2): Init: A

Expand & pop A: AB, AC

Expand & pop C: AB, ACF

Expand & pop F: AB, ACFE, ACFH

Expand & pop H: AB, ACFE, ACFHG, ACFHX

Expand & pop X: X is goal state, ret ACFHX

### AC3: V1 -> V2: D1 = RB, D2 = RG

V2 -> V1: D2 = RG, D1 = RB

V2 -> V3: D2 = RG, D3 = G

V3 -> V1: D1 = RB, D3 = G

V3 -> V2: D3 = G, D1 = RB

(+) V1 -> V2: D1 = RB, D2 = R

(+) V3 -> V1: D3 = G, D1 = B

Queue empty, all arc consistent, ret ...

### Pure Backtracking:

V1: R

V2: R (Inconsist. assign, bt)

V2: G

V3: G (Inconsist. assign., backtrack)

V1: B (V2 lv none left)

V2: R

V3: G, all var assigned, ret

V1: B, V2: R, V3: G

### BT with FC:

V1: R, remove any inconsist in

remaining var of neighbours

V2: G (FC)

V3: Empty, terminate

V1: B

V2: R

V3: G all var assigned, ret

V1: B, V2: R, V3: G

### Backtracking Search:

(DFS with single-var assignment)

CSP var assignments are commutative (order not imp)

Hence only consider assignment to 1 var at each level/depth, branch factor at all levels d,  $d^n$  leaves

CSP:  $O(d^n)$  If subproblem have c var out of n total,

worst case is  $O(\frac{n}{c} \cdot d^n)$

### Tree Structured CSP:

If constraint graph has no loop:  $O(nd^2)$  time since tree with n nodes has at most n edges, compare 2 nodes which each side at most d values

Domain = {1,2,3}

Constraints = A<B, B<C

Order var from root (choose one) to leaves st every node's parent precedes it in the ordering

From j = n to 2, apply RemoveInconsistent(parent(X<sub>j</sub>, X<sub>j</sub>))

RI(B,C): B<1,2,3, C<1,2,3

RI(A,B): A<1,2,3, B<1,2

From j = 1 ton, assign X<sub>j</sub> consistently with Parent(X<sub>j</sub>)

d = 3: {R, G, B}

Depth = 1 # var = n # leaves = nd

Branching factor: nd

Depth = 2 # var = n-1 # leaves = n!d^n

Branching factor: (n-1)d

Depth = n, # var = 1, # leaves = n!d^n

### Nearly Tree Structured CSP:

Conditioning: Instantiate a var, prune its neighbours' domain

Cutset conditioning: instantiate (in all ways) a set of variables at the remaining const. graph is a tree

Cutset size c implies  $O(d^c \cdot (n - c)d^2)$ , very fast for small c

Rest all pruned to G, B, solve as tree-CSP

Inst: R

WA NT Q NSW V SA T

Pruned solved as normal CSP c = 1 here

From j = 1 ton, assign X<sub>j</sub</sub>

**Supervised (Inductive) Learning:** Train + labels (desired output) e.g. classification, regression, obj detection, semantic segmentation

**Unsupervised:** Train (no desired output) e.g. clust, dim reduction

**Semi-supervised:** Train + few desired output

**Self-supervised:** Pretext task (e.g. predict if rotated image), learn the latent features (output labels) intrinsically from data

**RL:** Rewards from seq actions

**Evaluation:** Acc (must be class specific in case of imbalanced data)

E.g. If only 30% 1 & predict all 0, we have 70% acc! totally biased

Predicted

		+ve	-ve	Recall/ Sensitivity: $\frac{TP}{TP+FN}$
Actual	+ve	TP	FN (Type II Error)	
	-ve	FP (Type I Error)	TN	Specificity: $\frac{TN}{TN+FP}$
	Precision: $\frac{TP}{TP+FP}$	-ve Pred Val: $\frac{TN}{TN+FN}$	Acc: $\frac{(TP+TN)}{(All)}$	
	F1: (Precision * Recall)/(Precision + Recall)			

**MSE (L2 Loss) for regression:**

$$MSE = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

**Receiver Operating Curve (ROC):**  
TPR vs FPR (plot against all threshold, AUC = 1 best)

**Information:** (Generally) Measured in bits ( $\log_2$ ): Amt of "surprise" arising from a given (stochastic) event E:  $I(E) = -\log(P(E))$

E.g. Three boys: Tom (20%), Bob (1%), Sam (79%)  
Bob =  $-\log_2(0.01) = 6.6$ , high surprise vs  
Sam =  $-\log_2(0.79) = 0.34$ , low surprise

**Entropy:** Avg lvl of information inherent in possible outcomes:

Higher Entropy signifies more uncertainty in what action an agent would choose:  $H(X) = -\sum_i P_X(x_i) \log_b P_X(x_i)$

$$Entropy = -(0.20 \log_2(0.20) + \dots + 0.79 \log_2(0.79)) = 0.26$$

**KL-Divergence (relative entropy):**

Expression of Surprise, diff btw 2 prob distribution P & Q:  
Under the assumption P & Q are not close, KL divergence high (surprise). If close then KL divergence low

**Likelihood ratio (n indep samples):**  $LR = \prod_{i=0}^n \frac{p(x_i)}{q(x_i)}$

(how much more likely samples came from P dist than Q dist)

**Log likelihood ratio: LLR** =  $\sum_{i=0}^n \log\left(\frac{p(x_i)}{q(x_i)}\right)$

**Expected LLR:**

(Avg, how much each sample gives evidence of P(x) over Q(x))

$$D_{KL}(P \parallel Q) = \sum_{i=0}^n p(x_i) \log\left(\frac{p(x_i)}{q(x_i)}\right)$$

$$D_{KL}(P \parallel Q) = \sum_{i=0}^n p(x_i) \log(p(x_i)) - \sum_{i=0}^n p(x_i) \log(q(x_i))$$

Entropy (Info content of P) - IC of Q weighted by P

If P is the "true" distribution, then the KL divergence is the amount of information "lost" when expressing it via Q.

**Cross-Entropy (min unlike max likelihood):**

Cross entropy is, at its core, a way of measuring the "distance" between two probability distributions P and Q. (Entropy on its own is just a measure of a single probability distribution)

CE = Combination of the entropy of the "true" distribution P and the KL divergence between P and Q:

$$H(p, q) = H(p) + D_{KL}(p \parallel q) = -\sum_{i=0}^n p(x_i) \log(q(x_i))$$

E.g. Predict one-hot vector: True {0, 0, 1, 0, 0}

Output (of softmax layer): Predicted {0.01, 0.02, 0.75, 0.05}

Predicted class (argmax): 2 (zero index)

$$-(0 \cdot \log(0.01) + 0 \cdot \log(0.02) + 1 \cdot \log(0.75) + 0 \cdot \log(0.05))$$

**Cross-entropy:**

Avg # of total bits needed to encode data from a source with dist p when we use model q

**KL-Divergence:**

Avg # of extra bits needed to encode the data, when using dist q instead of the true dist p.

**Logistic Regression:**

Given an input space X, the goal is to predict for every sample  $x \in X$  to which class it belongs. For 2 class classification, the goal is:

predict output space  $Y = \{0, 1\}$  or  $\{-1, +1\}$

Linear:  $f(x) = w \cdot x + b$ , unbounded space

Need to map  $(-\infty, \infty)$  to (0, 1) st 0 map to 0.5

Use logistic sigmoid fn  $s(a)$ :

$$s(a) = \frac{\exp(a)}{1 + \exp(a)} = \frac{1}{\exp(-a) + 1} \exp(a) = \frac{1}{\exp(-a) + 1}$$

Note if multi class (k) class, use Softmax:

$$\sigma(z)_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}} \quad \text{for } i = 1, \dots, K \quad \text{and } z = (z_1, \dots, z_K) \in \mathbb{R}^K$$

$$h(x) = s(f_w, b(x)) \in (0, 1) \text{ aka prob(sample } x \text{ class label } y = +1)$$

**Non-deterministic Game:**

Adversarial search where the actions/transitions are non-deterministic

**ExpectiMax**, similar to **MiniMax** but accounts for chance nodes

**Cross-Entropy loss (bin class):**

**CE Loss:** (1 sample:  $(x_i, y_i)$ )

$-\bar{y}_i \log(p_1(x_i)) - (1 - \bar{y}_i) \log(1 - p_1(x_i))$  where  $p_1(x_i) = 1 - p_2(x_i)$

**Avg CE Loss:** (dataset: n indep samples)

$$\frac{1}{n} \sum_i -\bar{y}_i \log(p_1(x_i)) - (1 - \bar{y}_i) \log(1 - p_1(x_i))$$

Opti problem (convex) to be solved is (finding best parameter w min CE loss):

$$(w^*, b^*) = \operatorname{argmin}_{w,b} (\text{avg CE Loss})$$

$$\text{where } p_1(x_i) = h(x_i) = \frac{1}{\exp(-w \cdot x_i - b + 1)} = s(w \cdot x_i)$$

$$\nabla_w L = \nabla_w \left( (-1) \cdot \sum_{i=1}^n y_i \log(s(w \cdot x_i)) + (1 - y_i) \log(1 - s(w \cdot x_i)) \right) = \sum_{i=1}^n x_i (s(w \cdot x_i) - y_i) = \sum_{i=1}^n x_i (h(x_i) - y_i)$$

**Data Set:** (No best split) **Training:** run your learning algo

**Development:** tune parameters **Test:** evaluate perf

**Common Error:**

1) Data mismatch (e.g. diff. Resolution img)

2) Bias and variance

**Bias:** error rate on the training set

**Variance:** how much worse test set error than train (e.g. 15% error (85% acc) on train, bias = 15%). But 16% error on test, variance = 1%)

**Set Distribution:**

If all same dist, then optimising for web imgs mostly here, which is diff from target dist.

100k imgs (from Web) 8k imgs (from Target dist)

Train (104k imgs) Dev (2k imgs) Test (2k imgs)

Instead make test & dev from the same dist such that optimizing to perform well on the target dist.

100k imgs (from Web) 4k imgs 4k imgs

Train (104k imgs) Dev (2k imgs) Test (2k imgs)

However, the training dist is now diff from the dev/test distribution: it will take longer & more effort to optimize the model. More importantly, cannot easily tell if the classifier error on the dev set relative to the error on the train set is a variance error, a data mismatch error, or a combination of both.

You can use a **bridge set** (training data not used for training) to trouble shoot where the error comes from.

100k imgs (from Web) 4k imgs 4k imgs

Train (102k) Bridge (2k) Dev (2k) Test (2k)

Error 2% 3% 10%

Error 8% (btw dev & train): 1% var + 7% data mismatch Hence our classifier performed well on data from the same distribution but poorly when from a diff distribution, (data mismatch problem). Else if higher var & lower data mismatch error, implies overfitted to train set.

**Rules of thumb:**

**High avoidable bias: (underfitting)**

- Increase the size of your model
- Modify the input features based on error analysis
- Reduce or eliminate regularisation (reduces bias but increases variance)
- Modify model architecture



**High variance: (overfitting)**

- Add data to your training set (e.g. by data augmentation)
- Add regularisation, dropout
- Add early stopping
- Feature selection to decrease # of features (may increase bias)
- Decrease the model size (May increase bias)
- Modify input features based on error analysis
- Modify model architecture



**Low bias, low variance**



**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_{1..m}\}$ ; Parameters to be learned:  $\gamma, \beta$

**Output:**  $y_i = BN_{\gamma, \beta}(x_i)$

$$\mu_B \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{mini-batch mean}$$

$$\sigma_B^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2 \quad // \text{mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} \quad // \text{normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv BN_{\gamma, \beta}(x_i) \quad // \text{scale and shift}$$

**Minimax-Algo:** (Deterministic + Fully obs game). Idea: choose move to position with highest minimax val

**Completeness:** Yes if tree is finite (exist terminal states)

**Optimality:** Yes against a rational (opti) opponent

**Time Complexity:**  $O(b^m)$

**Space Complexity:**  $O(bm)$  similar to DFS

Yet, for chess e.g. branching factor (possible moves),  $b \approx 35$  & depth (number of plies)  $m \approx 100$ , exact solution is infeasible

**Optimizing Param:**

1. Init weights randomly  $\sim N(0, \sigma^2)$

2. Loop until convergence:

$$\text{a. Compute gradient: } \frac{\partial J(W_n)}{\partial W_n}$$

b. Update coeff:

$$W_{n+1} = W_n - \alpha \cdot \frac{\partial J(W_n)}{\partial W_n}$$

3. Return weights, grad to 0

Hence gradient wrt w: (sum of points  $x_i$  weighted by the diff btw the predicted value  $h(x_i) = s(w \cdot x_i)$  and true val  $y_i$ )

**MSE (of unseen test pt):** bias & var tradeoff + inherent noise,  $\varepsilon$ :  $y = f(x) + \varepsilon$  where  $\varepsilon \sim rv(0, \sigma_\varepsilon^2)$

$$\text{Total Error} = E_x [bias[\hat{f}(x)]^2] + E_x [\text{var}(\hat{f}(x))] + \sigma_\varepsilon^2$$

**MSE:** Avg squared diff of a pred  $\hat{f}(x)$  from its true val y:

$$MSE = E[(y - \hat{f}(x))^2]$$

**Bias:** Diff of avg val of pred (over diff realization of train data) to the true underlying function  $f(x)$  for a given unseen test point x:  $bias[\hat{f}(x)] = E[(\hat{f}(x))] - f(x)$

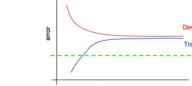
Note realization implies access to only a portion of the underlying data as training data (unrepresentative of the underlying population)

**Variance:** mean squared deviation of  $\hat{f}(x)$  from its expected value  $E[\hat{f}(x)]$  over diff realizations of training data:

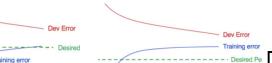
$$var[\hat{f}(x)] = E[(\hat{f}(x) - E[\hat{f}(x)])^2]$$

**Learning Curves:**

1) **High Bias:** Error plateau, dev error higher than train error, very high train error



2) **High Var, Low Bias:** Train error relatively low, dev error much higher than train error, bias is small but variance is large. Perhaps add more training data will close the gap



3) **High Bias, High Var:** Train error large and much higher than desired perf, dev error also larger than train.

**Regularisation:**

Used to discourage the complexity of the model by penalizing the loss function (help solve overfitting)

1) **L1 (Lasso/L1 Norm)**

Feature selection by assigning insig input features with 0 weight & useful features with a non-0 weight

$\lambda$  is the penalty (reg) param: If 0, back to original loss function no penalty, if too large, weights become close to 0, hence high bias (underfitting)

2) **L2 (Ridge)**

Forces the weights to be small but  $\neq 0$  (non-sparse)

Not robust to outliers as square terms blows up the error differences of the outliers &  $\lambda$  tries to fix it by penalizing the weights.

Performs better when all the input features influence the output and all with weights are of roughly equal size

$$L(x, y) \equiv \sum_{i=1}^n (y_i - h_\theta(x_i))^2 + \lambda \sum_{i=1}^p |\theta_i| \text{ vs } \sum_{i=1}^p \theta_i^2$$

$\theta_i$  is the coeff of the i-th term

Put the L1/L2 into the optimizer or use torch.norm(), total\_loss = loss + L1/L2, then backward prop

**Resource Limitations: Standard Approaches:**

a) **Cutoff test:** On the depth limit

b) **Evaluation fn:** Heuristic scoring on the est. desirability of current state

Note for the eval fn values, behaviour is preserved under any monotonic transformation of itself, only the order matters

$\alpha - \beta$  prunning

$\alpha$ : best value (to MAX) found so far off the current path

$\beta$ : best value (to MIN) found so far off the current path

Prune if  $\alpha \geq \beta$

- Pruning only affect the size of the search tree, not result

- Good move ordering improves effectiveness of pruning

- 'Perfect ordering':  $O\left(\frac{m}{b^2}\right)$ , hence we can double the depth of search

$\hat{y} = g(w_0 + x^T w)$  linear decis<sup>o</sup> no matter network size

Nonlinear  $g(a) \Rightarrow$  approx. of arbit complex f<sup>o</sup>

- i) Threshold:  $[a > 0]$
  - ii) Sigmoid:  $\sigma = \frac{1}{1+e^{-a}}$
  - iii) ReLU:  $\max(0, a)$
- non 0 centered, non diff @ 0, dying problem

a) Reg.: Linear  $(-\infty, +\infty)$  / ReLU  $(0, +\infty)$   $\nparallel$  MSE loss

b) Bin class.:  $\sigma(0, 1)$   $\nparallel$  BCE loss:  $\sigma + c\epsilon$  loss =  $\sigma$  CE loss

c) Categorical: Softmax  $(0, 1)$   $\forall$  class,  $\sum_i$  to 1  $\nparallel$  CE loss: Softmax g+ CE loss = Softmax loss

d) Multiclass cat:  $\rightarrow$   $\nparallel$  BCE loss  $\nparallel$  class

Empirical loss: Total loss over entire data

$$J(W) = \frac{1}{N} \sum_{i=1}^N \text{MSE}(f(x^{(i)}, w), y^{(i)})$$

$$\text{BCE} = \frac{1}{N} \sum_{i=1}^N [y^{(i)} \log(f(x^{(i)}, w)) + (1-y^{(i)}) \log(1-f(x^{(i)}, w))]$$

Training: Find  $w^* = \underset{w}{\operatorname{argmin}} J(W)$ , lowest loss

- i) Init weights rand  $\sim N(0, \sigma^2)$
- ii) Loop until converge:  $\nabla W = \frac{\partial J(W)}{\partial W}$ ,  $W \leftarrow W - h\nabla W$
- iii) Return  $W$
- Small  $h$ : converge quickly (local optima)  $\nparallel$  Use adapt!
- Large  $h$ : unstable & diverges

- Easy compute, implement & understand

- My trap at local min (need convex f<sup>o</sup> too)

- Update weights only after calculating  $\nabla$  on entire dataset (too large  $\Rightarrow$  years to compute + large mem)

SGD: Update model param more freq.

- Converges in less time  $\nparallel$  may get new minima
- High var in model param (may shoot up even after global min)
- To get some convergence as SGD, need slowly to h

Pick Batch B of data  $\left\{ \begin{array}{l} \text{Fast compute } \nabla \text{ allow } // \text{compute} \\ \nabla_W = \frac{1}{B} \sum_{i=1}^B \frac{\partial J_W(w)}{\partial W} \end{array} \right\}$

$\nabla_W$  better est of true  $\nabla$  than 1 datapoint, smoother converge, allow larger  $h$

Adagrad: Changes  $h$  if param  $l$  t (need 2<sup>nd</sup> order der)

- No need tune  $h$  manually & trainable on sparse data

- It always  $\downarrow$  slow training - Comp. exp  $\xrightarrow{\text{Adadelta}}$

ADAM: momentum of 1st+2nd order - Comp. exp. RMS Prob

- Very fast converge + rectifies decaying  $\nabla$  & high var

$$\frac{\partial E_{\text{Total}}}{\partial w_5} = \frac{\partial E_{\text{Total}}}{\partial w_5} \left( \frac{\partial \text{act}_1}{\partial w_5} \right) \left( \frac{\partial \text{act}_2}{\partial w_5} \right)$$

where  $\frac{\partial \text{act}_1}{\partial w_5} = \text{act}_1 \cdot \frac{\partial \text{act}_1}{\partial w_5} = \text{act}_1(1 - \text{act}_1)$

$$\frac{\partial E_{\text{Total}}}{\partial w_5} = -(\text{target}_1 - \text{act}_1) \cdot w_5 = w_5 - h \frac{\partial E_{\text{Total}}}{\partial w_5}$$

For  $b_1$ :  $\text{net}_{b_1} = i_1 \cdot w_1 + i_2 \cdot w_2 + b_1$ ,  $\text{out}_{b_1} = 1/(1+e^{-\text{net}_{b_1}})$

For  $b_2$ :  $\text{net}_{b_2} = i_1 \cdot w_3 + i_2 \cdot w_4 + b_2$ ,  $\text{out}_{b_2} = 1/(1+e^{-\text{net}_{b_2}})$

For  $w_1$ :  $\text{net}_{w_1} = a_1 \cdot w_1 + a_2 \cdot w_2 + a_3 \cdot w_3 + b_1$ ,  $\text{out}_{w_1} = 1/(1+e^{-\text{net}_{w_1}})$

For  $w_2$ :  $\text{net}_{w_2} = a_1 \cdot w_4 + a_2 \cdot w_3 + a_3 \cdot w_1 + b_2$ ,  $\text{out}_{w_2} = 1/(1+e^{-\text{net}_{w_2}})$

$E_{\text{Total}} = E_{b_1} + E_{b_2}$

$$= \frac{1}{2} (\text{target}_1 - \text{act}_1)^2 + \frac{1}{2} (\text{target}_2 - \text{act}_2)^2$$

Generative Modeling: Take input train samples from some dist and learn a model that represents that dist (Debiasing)

- Capable of uncovering underlying features in a dataset

- Detect outliers in the dist for training

AE: Encoder: Unsupervised learning of a lower dim feature space (under complete),  $z$  (salient features + easier to work with) representation from unlabelled train data (Note no labels used.)

Decoder: Use  $z$  to reconstruct the obs  $\hat{x}$

$$L(x, \hat{x}) = \|x - \hat{x}\|^2$$

- Form of compression, smaller  $z \rightarrow$  larger training bottleneck
- Bottleneck hidden layers forces network to learn compressed latent rep. while reconstruct loss forces  $z$  to capture as much information of the data

Regularizer AE: Overcomplete case with a loss function that encourages the model to have other properties beside copy-paste

- sparsity rep, smallness of rep derivatives (similar img, sim lv), robust to noise/missing inputs

Sparse AE:  $L(x, g(f(x)) + \Omega(z))$

decoder, encoder output + sparse penalty

Denoising AE: Add noise to input image, AE learns to remove noise



Bagword (BoW): Columns: Sentences, Rows: all words  
Uses Binary one-hot encoding (very sparse)

Term-Freq-Inv doc. freq. (TF-IDF): Rare words carry more meaning (also very sparse)

TF = # appearances for a term / Total # words in document

IDF = log (# docs / # docs term appear in)

TF-IDF = TF  $\times$  IDF

Less-sparse Rep: Co-occurrence models (Singular-value-decomp)  
Vector embedding models (word2Vec)

Word Embeddings: Each word = a point in a  $n$  dimen space  
Represented by vector of len  $n$  ( $\sim 100-300$ ), 1 layer

i) Cont.-BoW: Pred center word from surrounding context

ii) Skip-grams: Pred surr context words from given center word

- Maximum likelihood of context, given focus word  
e.g. I love pizza,  $P(1|pizza)$  &  $P(\text{love}|pizza)$

Train word2vec: (weights of layer) but no preserve order!

- Decide window size & embedding size & vocab size

- Training Obj: Negative sampling (binary LR class)

- Fast accurate model that captures semantic content

"Gensim" library, pre-trained word embedding

- Interested in the 1 hidden (latent rep) layer, not the output

CNN: Preprocessing calls e.g. Normalize to increase perf

- Decrease # params to learn, preserves locality (no need flatten image)

- For all conv. layer, # activation map = # filters/feature maps/kernels

- Depth of input (channels) = Depth of filters too (Note depth of activation map output is 1 nevertheless)

- For 2D image H with 2D kernel F, Conv operator:

$$G = H * F, G[i, j] = \sum_{u=-k}^k \sum_{v=-k}^k H[u, v] F[i-u, j-v]$$

However in CNN, do not flip kernel; instead of more complex above equation, use:

$\boxed{[x]}$  symmetric, kernels transposed

$$G[i, j] = \sum_{u=-k}^k \sum_{v=-k}^k H[u, v] F[i+u, j+v]$$

pos in resulting act map, centred at 0

pos in filter (given dim  $2k+1$ )

- Still no difference in results, same output, just convention
- Stride: Output size = (N, input size - F, feature size / Stride) + 1
- Padding: Add padding to output, decrease problem size

To preserve spatial size: Stride = 1, F by F filter size, with 0-padding of  $(F-1)/2$  size

- Pooling: Non-linear down sampling (common max/ avg), no learning here, just the stride size and pooling window size, non overlapping best, if too large lose information

L to R, Top -down

Input  $\rightarrow$  Conv+Relu, Pooling  $\rightarrow \dots \rightarrow$  Flatten+FC+Softmax

Feature learn Classification

- Can be used for any vector/matrix representation  
e.g. Audio: Spectrogram, Word: Embeddings

- Waveforms: Bunch of numbers over time

- More informative are spectrograms: a visual representation of the spectrum of frequencies of sound or another signal as they vary with time:

- Typical spectrogram of a few spoken words.

Y-axis: frequencies, X-axis: time, lower frequencies can be seen as denser (brighter), like a male voice

- The power spectrum of a time signal describes the power distribution into freq components composing that signal. Fourier analysis can decompose a signal into discrete freq (a spectrum of freq) over a continuous range.

- Effective: Mel Freq Cepstrum (non-linear) represents human hearing better. TA library nnAudio

Variational AE (VAE): Sample from mean & std dev to compute  $z$  (probabilistic, no longer deterministic).

E.g.  $z = \begin{cases} \text{smile: } \xrightarrow{\text{act}} \\ \text{samp 1} \\ \vdots \\ \text{samp 2} \end{cases}$  latent dist.

$$L(\Theta, \theta, x) = \text{reconst loss} + \text{reg term} = \|x - \hat{x}\|^2 + \sum_j KL(P_j(z|x) || p(z))$$

- Choice of prior on id: common  $p(z) \sim N(0, 1)$ , encourages encodings to distribute evenly around center of  $z$ , penalizes if tries to cluster points tgt (memorise data)

- Generation (only sampling for  $z + \text{decoder}$ ) vs
- Training (Encoder: produces  $z(\mu + \sigma)$ , calc reg loss, then ...., calc reconstruct loss, minimize both losses tgt)

- Cannot backprop grad through sampling layers,  $\therefore$  reparametrize sampling layer (sum fixed  $\mu$  &  $\sigma$  scaled by random const drawn from prior dist.  $z = \mu + \sigma \cdot \epsilon \sim N(0, 1)$ )

- Goal is to find disentangled uncorrelated meaningful lv  
Latent Perturbation: Slowly increase/dec 1 lv, keep rest same for interpretation

Conditional VAE: Feed label  $Y$  into model, since VAE no control over the generated data

RNN: Input  $x_t$  vector, Output  $\hat{y} = W_{hy}^T h_t$

Update hidden state:  $h_t = \tanh(W_{hh}^T h_{t-1} + W_{xh}^T x_t)$

- Many values >1: Exploding gradients, need grad clipping

- Many values <1: Small numbers multiplication, gradient increasingly smaller, bias, param unable to capture short term dependencies

- Weights matrices updated after entire timesteps

a) ReLU, prevent  $f'$  from shrinking grad when  $x > 0$

b) Weight Init: Weights use Identity matrix, Bias use 0, prevent weights shrinking to 0

c) Using gates to track and control information through time

LSTM: - Forget irrelevant info of prev states

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

$\circ$ : pointwise multiplication,

$\overline{C}_t = \tanh(W_c \cdot [h_{t-1}, x_t] + b_c)$  Hadamard mult

- Update cell state selectively

$$C_t = f_t \circ C_{t-1} + i_t \circ \overline{C}_t$$

- Output gate: Control what to send to next t

$$O_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

$$h_t = O_t \circ \tanh(C_t)$$

Cellstate highway: Allows uninterrupted grad flow

input gate output gate

Transfer Learning (Opti): Higher start, slope & asymptote

- Identify related task with abundant data, reuse pretrained model

E.g. LeNet 5, AlexNet, VGG-16,

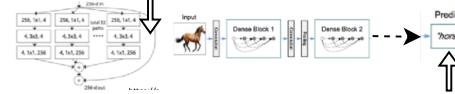
Inception (GoogleLeNet)

- More layers  $\rightarrow$  effect, degradation pro

$f(x) \rightarrow f(x) \circ \text{relu}$

ResNet: residual blocks at intermediate layers (winner)

ResNext: "Split-transform-merge" + residual block connection



DenseNet: Each layer's feature map concatenated to the input of every successive layer within dense block

Generative Adversarial Networks (GAN):

Generator: Noise  $\rightarrow$  imitation of data to try trick D

Discriminator: Tries identify data from fakes generated by G

$\text{noise } z \xrightarrow{G} \text{fake } x \xrightarrow{D} \text{real } y$

2 nn competing

Data: Real  $\rightarrow$  D use as +ve e.g., Fake  $\rightarrow$  D use as -ve e.g.

Alternative training:

D Training: D ignores G loss & uses its own loss

- Sample random noise to produce G output

- classify & calculate D loss

BP: through D & G weight, only update G weights

G Training: Opposite

- G loss penalizes G for failing to fool D

If G perfect, D acc = 50%, D feedback gets less meaningful overtime, may collapse, sin ce GAN not stable, but fleeting convergence

D output real

Minimax loss:  $\min_{\theta_g} \max_{\theta_d} [E_{X \sim P_{\text{data}}} \log D_{\theta_d}(c) + E_{Z \sim P(z)} \log (1 - D_{\theta_d}(G_{\theta_g}(z)))]$

G output | $x$

D output fake

(CE loss btw real & fake dist)

- D max probability real & output prob fake low

- G no influence on 1st term, tries min 2nd term (1-...)  $\rightarrow$  get high output for fake to fool D

- Can get stuck at early stages of training,

$\therefore$  use Modified minimax loss: max log  $D(G(z))$  instead of min  $(1 - D(G(z)))$

Deep Conv GANs: Uses leaky ReLu (-ve slope before  $x = 0$ ), sin ce ReLu is non-0 centered, non diff at 0 & dying problem

Conditional GANs: same feed label  $Y$ (paired data)

Used for object transfiguration

a) Pixel2pix(cGAN with original img as condition)

- Paired data, D compares  $G(x)$  and  $y$  given  $x$  aka generate img in different style from original img

- To make img less blurry use L1/L2

$$G^* = \arg \min_{G} \max_{D} L_{cGAN}(G, D) + \lambda L_{L1}(G)$$

Cycle GAN (Unpaired data): Domain transformation

- Uses adversarial loss + inverse mapping (sin ce mapping highly constrained). Introduces a cycle consistency loss to push  $F(G(x)) \approx X$

- Not just img, e.g. waveform, but not spectrogram

- Used to convert accented speech

