- [ T: Theorem ] Search > Uninformed
- [Fringe:, Expanded:,] [N.=Node I=Implementation]

## BES

- Expands Shallowest unexpanded N. (1: put children of the Expanded N. at the end of the fringe ) Complete: Yes (if b is finite) Optimal: No in General (Yes if Step Cost is the same)

- Time: 1+b+b^2+...b^d = O(b^d), expontnl Space: O(b^d), (keeps every node in memory)
- Uniform Cost (UCS)
- Expands Least-Cost Unexpanded N, g(n) (1: insert nodes in the fringe in order of increasing path cost From the root)
- Complete: Yes, if step cost > 0
- Optimal: Yes Time & Space:
- HIN with gscost of Optimal Soln O(b^d), (depends on path costs, not depths, difficult to caterize in terms of b, d) UCS = BFS, when g(n)=depth(n)
- Expands Deepest Unexpanded N
- ( I: insert successors at front of fringe )
- Complete: No. fails in infinite-depth spaces (i.e. m= ∞) Ontimal: No
- Time: 1+b+b^2+...b^m = O(b^m), (higher than BFS, as M>>d (m=max
- depth, d=least cost soln path) ) Space: O(bm), linear, excellent.

## Depth Limited Search

- DFS w/ depth limit I Complete: No in general, Yes in finite spc.
- Optimal: No Time: 1+b^2+...b^l= O(b^l) (as BFS)
- Space: O(bl) (as DFS) Iterative Deepening (IDS)
- Expands deepest unexpanded node within level I. Complete: Yes (as BFS)
- Optimal: Yes, if step cost=1 (as BFS)
- Time: O(b^d) (as BFS)
- Space: O(bd), linear (as DFS) ( can be modified to explore uniform-c

## Search > Informed (heuristic)

(+ve over uninformed, Knows if a non-goal node > another, typically more efficient )

- [Best First Search Algorithms : Expands the most desirable unexpanded node]
- [ N=Node, g(n)=pathCost, h(n)=heuristic ] Greedy
- Expands N. with smallest h(n)
- Complete: Yes, in finite space. Fails in ∞ space (+ can get stuck in loops) Conjuete. res, in mine space. rais in ~ space ( Optimal: No. Time: O(b<sup>+</sup>m), but good heurstic can Improve lots Space: O(b<sup>+</sup>m), keeps every node in memory Greedy=BFS, h(n) = depth(n), ties+L}R

- Greedy=DFS, h(n)=-depth(n), ties+ deepest 1st. Greedy=UCS, h(n)=g(n)

- Comparison (Company)
   Company N, with smallest (n)=g(n)+h(n)
   Expands N, with smallest (n)=g(n)+h(n)
   Tif h is an admissible heuristic, than A\* is Complete And Optimal. (only with the second state of the theory of the second state of the theory of the second state of the secon w/ Tree S.) T: If h is consistent, than A\* is Optimally Efficient, among all optimal search algorithms. (always true ^v) (It will not revisit states (as in graph
- search))
- Complete: Yes, unless there are ∞ //v many nodes with f≤f(G), G–Optimal Goal State Goal State Optimal: Ves, with admissible heuristic Time: (Qb<sup>\*4</sup>), exponential, b<sup>\*</sup>=effective branching factor Space: Exponential, keeps all nodes in memory. [Both Time&Space are probs for A\*, typically running out of SPACE b4

- time]
- May be better to settle for a non-admissible h that works well ever May be better to settle nor a non-acmission in that works well even though completeness and optimality are no longer guaranteed. Simpler, faster h may be better even though more N.s expan+. - A "=UCS, if h(n)=0 for all N. A" (foraph Search) - If h is admissible, Complete & Optimal, if revising repeated (:||) states ellewed (measuring closed N.b. close art entities)

- allowed (reopening closed N.s), else not optimal
- If h is Consistent, we avoid : || states. Enforced Consistency, using 'PathMax': set child's f value to parent's f
- (If done as we search, may not solve problem of reopening N.s from Closed. Better to ensure before search starts that h(n) is consistent ).
- Heuristics (h(n))
- Admissible h(n) are optimistic (smaller than true cost), e.g. SLD. Dominance, when an Admissible Heuristic is better than another admissible one. A better estimate of true \$ to G, and expanding fewer
- nodes in A\* Inventing, h(n)'s by creating for a relaxed version of prolem (1 w/ less restrictions on the actions)
- [ T: \$ of an optimal soln to a Relaxed problem is an Admisible h(n) for the
- [1:3 of an optimal solin to a Kelaxed problem is an Admistible h(n) for the Original problem.] (composite heuristics h(n)=max(h1(n),h2(n)...) are admissible, good to use when there's a bunch of h(n), none dominating one another, composite will Dominate.) **Consistent (monotonic) Heurisic**, if all such pairs in the search graph which the thereis the source the search graph.

- satisfy the triangle inequality: [ h(n), par. 5 cost(n,n)+h(n), chi for all n ] + (f(n)), child  $\geq$  (f(n), parent : f is non-decreasing along any path (Admissible (Consistent))

# Local Search Algorithms

- (Optimisation Problems) Hill-Climbing Finds closest local min.imum / max.imum. (may not be global)

- Soln: Evaluation f() should be applied only to positions that are

Search, extending search to make sure there's no hidden pitfall

(c) or distance/instance-based learning) An eg. of Lazy Learning, stores all training eg.s + doesn't build a classifier until new eg. needs to be classified. Opposite to Eager learning (constructing classifier before recieving new eg.s, like 1R,DT,NB.)

Lazy classifiers are Faster at training(=memorizing), but Slower at

Minkowski distance – generalization of Eucledian & Manhattan

 $D(A,B) = (|a_1 - b_1|^q + |a_2 - b_2|^q + ... + |a_n - b_n|^q)^{1/q}$  q - positive int

Need for **Normalization**, as when calculating distances between 2 examples, the effect of the attributes with smaller scale will be less significant than those larger. i.e. Normalize (between 0 and 1)

Classification>Lookup: O(mn), m training examples w/ dimensionality n

Classification 2 Jookup: U(m), m training examples w amenisonaim Memory: U(m), need remember each eg BAD for large datasets, **slow**. CNearest Neighbour k majority volting Very Sensitive to value of k general nule: [k ≤ sqrt(#training\_es)]] Very Sensitive to value of uverside and data (sensitive) by uverside addition.

Also usable w/ numeric prediction (regression) by averaging values Distance for Nominal Attributes:

Ustance for Nominal Authoutes:
0 - 2 the same, 1 otherwise
for Missing Values:
0 - both same & NON-MISSING, else 1
(if numeric, d=max(v, 1.0-v))
d((red, new, 7, 7, 7),(blue, new, 7, 0.3, 0.8)) = 1 + 0 + 1 + 0.7 + 0.8
Variation: Weighted Nearest Neighbor
Closer neighbors count more than distant neighbours. Instead of k, all
training ans

 $w_i = \frac{1}{D(X_{new},X_i)}_{-ve: \text{ slower algorithm}}$  - Curse of Dimensionality: NN great in low dimensions (up to 6), but

become ineffective as dimensionality increases. As more examples are

far from one another, and close to the boundaries. (Notion of nearness

ter non not articular, and odde of matedouillante (retained inductor) becomes ineffective in high-dim space) Sohr. Feature selection (attrs) to reduce dimensionality. Produces articliarially shaped decision boundary defined by a subset of the Voronoi edges.

Standard algorithm makes predictions based on LOCAL info. 1R, DT,

For each attribute value makes rule by majority class. Calculates error rate of rules. Chooses rule w/ smallest error rate Missing Values: Treated as another attribute Value

Treated as another autonue value Nominal Attivitues are discretized to nominal. - May lead to overfitting due to noise in data. - Soln: impose min num of egs of majority class in each partition, merge. Simple, Computationally Cheap.

- Assumes attributes are equally important and independent of one

 $P(sunny \mid yes) = \underbrace{\bigoplus_{9+3}}_{9+3} = \frac{1}{12}$ 

 $P(overcast | yes) = \frac{4+1}{9+3} = \frac{5}{12}$ 

 $P(rainy | yes) = \frac{3+1}{9+3} = \frac{4}{12}$ 

Missing Nominal Values: Ommit Value from P(yes|E) and P(no|E)

If Numeric: Calc. Probability Distribution using the Probability Densitiy Function (assuming normal distribution). μ=mean, Ô=sd.

+ves: simple,
 Excellent Computational Complexity: Requires 1 scan of the training data to calculate statistics (for both nominal & continuous attributes assuming normal distribution). O(pk), p=#training\_egs,

- rousust to isolate noise points (avgd out)
 -ves: - Correlated attributes reduce power (violation of independence assumption) - Solin: feature selection b4hand.
 - many numeric features not normally distrubted - Soln: other types of distributions /transform attribute to normally distributed one /discretize data first.

 $f(temperature = 66 | yes) = \frac{1}{6.2\sqrt{2\pi}}e^{-\frac{(66/7)^2}{2^*6.2^2}}$ 

(note, in tut examples, only the value is 2d, not cousins )

 $(x-\mu)^2$  $\frac{1}{\sigma\sqrt{2\pi}}e^{-\frac{(1-\mu)^2}{2\sigma^2}}$ 

· Robust to isolate noise points (avgd out)

lace Correction to handle Zero-Numerators. (add 1/Num Attrs to all

NNs, try to find a GLOBAL model that fits training set

Weight Contribution based off distance to new example:

ExpectMinMiMax

MAX

CHANCE

(Non-deterministic Games)

3=2\*0.5+4\*0.5

Supervised Learning

Classification - categorica Regression - Numeric

NEAREST NEIGHBOUR

- Nearest determined by distance Eucledian distance – most frequently used

Manhattan (or city block) distance

 $D(A,B) = \sqrt{(a_1 - b_1)^2 + (a_2 - b_2)^2 + \dots + (a_n - b_n)^2}$ 

 $D(A,B) = |a_1 - b_1| + |a_2 - b_2| + \dots + |a_n - b_n|$ 

Training=Fast, no model built, just storing.

classification.

training eqs.

Sensitive to Noise

NAIVE BAYES

anothe

attrs.)

f(x) = -

ves: • simple

k=valued attributes

Max

Min

(Statistical-Based Classification) • P(H|E) = P(E|H).P(H) / P(E)

P(yes|E) = P(E1|yes).P(E2|yes).P(yes)/P(E)

· (all P(Ex)ves) have same denominator)

36

Time: O(b<sup>m</sup> . n<sup>m</sup>) n=# of distinct dice rolls

0.5

quiescent, unlikely to change extremely in near future. + Secondary

Evaluating Classifier

(~2/3.1/3)

eval

3.

Suppose that

Confusion Matrix

q = 2

q = 1

in hoth sets

-ve: high computation cost i.e. more useful for small data sets

C1 C2 Fold 1 95% 91% d1=|95-91|=4

Fold 2 82% 85% d2=|82-85|=3

mean 91.3% 89.4% d\_mean=3.5

4. Interval contains 0 -difference not signific

 $Z = 3.5 \pm 2.26 \times 0.5 = 3.5 \pm 1.13$ 

d\_mean=3.5 and the standard deviation is 0.5

Calculate the differences di

- We use 10 fold CV => k=10

tp

 $P = \frac{\cdot_r}{tp + fp}$ 

examples

Precision = (R+R)/Retrieved

call = (R+R)/Relevant

Inductive Learning

correctly).

DECISION TREES

0or1=0hit

=H(S)-

# from class yes

# from class no

Where retrieved=#retrievedDocs\_relevant=#relDocs

Choice of H important (eg sin.f() vs polyn).

Empirical Evaluation: - best thing we can do. (Test data).

Comparing Classifiers

Holdout Procedure - split data into 2 independent sets; Training & Test

(~25,10) Accuracy (Success Rate) = 1.0 - Error Rate Validation Set (for DTs, NNs) - Classifier built from Training Set, - Tuned wi Validation Set, Evaluated wi Test Set. • DTs - training set used to build tree, validation set used to prune, test to

Prob: egs in training set may not be representative of all classes. Soln: Stratification, ensures each class is represented w/ ~= proportions

In oom sets. Holdout more reliable by repeating (Repeated Holdout Method) - which can be improved by ensuring Test sets don't overlap - Cross Validation. Leave-One-Out Cross Validation - n-fold cross-validation, where

 $\hat{\sigma}^2 = \frac{\sum_{j=1}^{k} (di - d_mean)^2}{2}$ 

2PR

P + R

# assigned<br/>to class yes# assigned to<br/>class notrue positivesfalse negatives

 (tp)
 (fn)

 false positives
 true negatives

 (fp)
 (tn)

k(k-1)

n=#Egs in data set. • +ve: greates possible amount of data used, deterministic procedure.

Calculate the variance of the difference (an estimate of the true v If k is sufficiently large, di is normally distributed

If it is submetric as  $p_{i}$ ,  $d_{i}$  is  $q_{i}$ ,  $d_{i}$  is  $q_{i}$ . Obtained from a prob  $Z = d \_mean \underbrace{(1-\alpha)(k)}_{(1-\alpha)(k-1)} \underbrace{\sigma}_{1-\alpha}$  - confidence level  $k_{1}$  - degree of freedom Interval contains 0 - difference not significant, else significant

- We are interested in significance at 95% confidence level

Interval does not contain 9 => difference is statistically significant

Conusion matrix (remember accuracy=(tp+tn)/(tp+tn+fp+fn)) Information retrieval (IR) uses *recall (R), precision (P)* and their combination *FI measure (FI)* as performance measures

tp

Mucure Learning Supervised Learning is Inductive Learning. Induction: inducing the universal from the particular. We can generate many hypothese h, the set of all possible h form the hypothesis space H, a good h will generalize well (i.e. predict new egs

JECISION IREES (top-down recursive divide-and-conquer) Growing the tree: hill climbing search guided by information gain. Can represent any boolean function. Compact Decision tree with most pure (high entropy) attributes. Entropy is measured in bits. For binary classification; value:0.5=1bit, value:

Entropy H(Y) measures: • the disorder of a set of training egs w/ respect

Linutgy in () ineasures, the usature to a set of training edge in respect to class Y, shows the amount of surprise of the receiver by the answer Y based on probability of answers • the smallest # of bits per symbol (on any) needed to transmit a stream of symbols drawn from Ys distribution. Information Gain is the expected reduction in entropy caused by the

 $Gain(S \mid A) = H(S) - \sum_{i \in V \in V} P(A = v_i) H(S \mid A = v) =$ 

 $\frac{|S_v|}{|S|}H(S \mid A = v_j)$ 

(u). The use autubule has ingress (an) Overfitting, due to D's growing each branch deeply to perfectly classify, coupled with noise in training data, +/or small training set. Soln: Pre/Post Pruning.
Stop Pruning: Estimate accuracy w/ validation set (acts as safety net).

partitioning of the set of examples using that attribute

 $-\sum_{j \in values(A)} \frac{|S_v|}{|S|}$ 

(DT: The best attribute has Highest Gain.)

however tree is built on less data. Post-**Pruning** > Tree Pruning: • Sub-Tree Replacement:

start from leaves+work to root. For each candidate node

with that node

accuracy and prune it no nodes are pruned if the new tree is worse than the old

remove the sub-tree rooted at it

make it a leaf and assign the most common label of the

training examples affiliated

Calculate accuracy of the new tree on validation set and compare with the old tree

Choose the node whose removal

results in the highest increase of

- ontinue iteratively till further pruning=harmful
 Sub-Tree Raising
 potentially time consuming operation,

restricted to raising the sub-tree of most popular branch.

NNs - validation set used to stop training, prevent overtaining

- Finds doeses local minimum (maximum, (may not be global) Soln found depends on Initial State: Can run several times starting from  $\partial$  rand, points. Plateus: (random walk no change in v, wander endlessly, revisiting prev.
- N.s): Can keep track of # of times v is the same and don't allow revisitng of
- nodes w/ same v. Ridges: (cur. local max. not good 'nuff): Can combine 2/+ moves in a macro, or allow limited # of look-ahead

search Keeps track of k BEST states (not 1) (\* 2 vrsns: 1. start wil 1 given state OR 2. k randomly generated states At each iteration (lvl): gen.erate all successors of all k states
 If any one is a goal state, stop; else select k best successors and continue

# Beam Search and Hill-Climbing Search

- Compare beam search with 1 initial state and hill climbing m-1 start node, at each step keeps k best nodes
- Hill climbing 1 start node, at each step keeps 1 best node Compare beam search with 1 initial state and hill-climbing with k random initial states
- Beam 1 start node, k search threads are run in parallel and useful information is passed among them
- Hill climbing k starting positions, k threads run individually, no passing of information among them
- Can be used w/ A\*, +ve: memory efficieny, -ve: !complete, !optimal Variations: Keep only nodes that are at most €(beam width) worse than
- best N. Simulated Annealing
- (similar to hill climbing, but selects random successor)
- (similar to finil cimolog, but selects random (select finitistate s. set cur. N. to s Randomly select m, one of N.'s succssrs if v(m) > v(n), n=m //accept m else n=m w/ small probability //accept m w/ small prob Anneal T, Repeat ktimes/goodnuff) Parbehöltt p. = 0.6/ v(m).V(T)

- Probability:  $P = e^{(v(m)-v(n))/T}$ i.e. bad move v(n)>v(m) asuming looking for min., P decreases expo. w
- i.e. bad move v(n)>v(m) asuming looking for min., P decreases expo. badness of move. T decreases, anneals, w/ time, e.g. T\*=.8 Thm:: fischedule lowers T slowly enough, algorithm will find global optimum. Complete & Optimal, given a long enough cooling schedule. Difficult to set 'slowly enough' (T).
- Genetic Algorithms (\* Select best individuals, from fitness f()
- CrossOver, about an init random point
- Mutate, random change of bits) Success depends on represental *lcomplete, loptimal* ntation (encoding)

# Games

# [ Deterministic vs Chance 1

1. 2. 3.

4

[Zero-Sum vs Nor-0-2]: (1 P's gain is another's loss) ] [Processes forward, blackward from goal, cause often too many goal states. + if goal state too far, will not provide any useful info on termin

# MinMax Algorithm

## fect, Deterministic, Assumes both P's (Max, Min) play optimally Minimax Algorithm - Typical Implementation xample from http://nages.cs.wisc.edu/~sl

if MIN's move, minimum value of children backed up if MAX's move, maximum value of children backed up choose move with the maximum of minimax values of children

minimax values gradually propagate upwards as DFS proceeds: i.e minimax values propagate up in "idrito-oright" fashion minimax values for sub-tree backed up "as we go", so only *O(bd)* nodes need to be kept in memory at any time

<=2

200

each move by MAX (first player): Perform depth-first search to a terminal s Evaluate each terminal state Propagate upwards the minimax values

Only O(bd) nodes need be kept in memory at a time If Min doesn't play optimally, Max will do even better

Optimal: res.
 Time: O(b^m) as in DFS – main Problem.
 Space: O(bm) as in DFS
 Alpha Beta Pruning

>=3>=4

<=10 <=8 4

Assumptions; branching factor b, all terminal Nodes ad depth d.

10 8 4

At each non-leaf node store a value indicating the best

MAX nodes: alpha value = best (max) value found so far MIN nodes: beta value = best (min) value found so far Given a node n, cutoff the search below n if:

*n* is a MAX node and  $abla(n) \ge beta (i)$  for some MIN node *i* that is ancestor of *n* (beta cutoff) *n* is a MIN node and beta(*n*)  $\le alpha(i)$  for some MAX node *i* that is ancestor of *n* ( $abpha \ cut$ )

Traverse the search tree in depth-first order

Apply utility function at each leaf node generated ompute values backwards

Implemented as DFS

Ontimal: Yes

 $\bigcirc$ 

8

cked up value so far

Pruning doesn't effect final result

Worst Case: No Pruning: O(b^d)

Best Case: Perfect Ordering O(b^(d/2))

[ Both MinMax & AlphaBeta require too much time ]

ation at leaf node Probs: Horizon Effect (hidden pitfalls)

3

Imperfect

e. Heuristic Evalu

# Post-Pruning > Rule Pruning:

- Grow Tree 
   Convert tree to equivalent set of rules by creating 1 rule Grow Iree - Convert tree to equivalent set of rules by creating 1 rule per path (fristement) + Prune each rule by removing any pre-conditions that result in improving estimated accuracy - Sort pruned rules by estimated accuracy, consider in this sequence when classifying subsequent ristances.
   Converting DT to Rules Be4 Pruning: Provides bigger flexibility, when
- trees are pruned; can only remove node completely or retain. when rules are pruned there are less restrictions; pre-conditions, not nodes, are removed, each branch in tree (i.e. each rule) treated separately, removes distinction between attribute tests that occur near the root of the tree and distinction between authorize tests that occur hear the root of the tree and those near the leeves. +ve > trees, easier to read. Mumerical attributes need discretization, in DTs we're restricted to binary

- split. Problem: if an attribute is **highly-branchng** then likely selected by
- Information Gain, can lead to **Overfitting**. Soln: **Gain Ration**. a modification of the Gain that reduces its bias
- towards highly branching attributes. it penalizes highly-branching attributes by incorporating

SplitInformation  $splitInformation(S \mid A) = -\sum_{i=1}^{c} \frac{|S_i|}{|S|} \log_2 \frac{|S_i|}{|S|}$ 

# SplitInformation is the entropy of S w.r.t. the values of A $GainRatio(S \mid A) = \frac{Gain(S \mid A)}{SplitInformation(S \mid A)}$

### Gain ratio:

- Missing values: handled by unique value A(x)=most common value among training egs at n w/ dask() · sample fractioning strategy - assign probability to each possible value of A, calc prob.s, use frequencies of values of A among examples at n Altributes w/ Different Costs, Incorporating cost in Gain (penalizing attrs
- $\begin{array}{c} \text{Automatical with high cost})\\ \text{mmer (90)} \\ Gain<sup>2</sup>(S \mid A) \end{array}$ 
  - Cost(A)
  - $2^{Gain(S|A)} 1$  where  $w \in [0,1]$  determines cost importance
- $\overline{(Cost(A)+1)^w}$

- Efficient: Cost of Building tree O(mnlogn), n-instances and m attributes
- · Cost of pruning tree w/ sub-tree replacement O(n)
- Cost of pruning by subtree lifting
- O(n(loan)^2)
- ∑ (build+prune): O(mnlogn)+O(n(logn)^2)
   Resulting Hypthoese Easy to interpret by humans if DT not too big
- NEURAL NETWORKS

# Single-Neuron Perceptron -Investigation of the Decision Boundaries

• 4. The decision boundary is:  $n = \mathbf{w}\mathbf{p} + b = w_1 p_1 + w_2 p_2 + b = p_1 + p_2 - I = 0$ i.e. a line in the input space 5. Draw the decision boundary  $p_1 = 0 \Rightarrow p_2 = I; (p_2 \text{ intersect})$  $p_2 = 0 \Rightarrow p_1 = I; (p_1 \text{ intersect})$ • 1. Output of the net:  $a = hardlim(\mathbf{wp}+b) =$ = hardlim(w<sub>1</sub> p<sub>1</sub> + w<sub>2</sub> p<sub>2</sub> + b) • 2. Decision boundary: w p+b=0  $n = \mathbf{w}\mathbf{p} + b = w_1 p_1 + w_1 p_2 + b = 0$ · 3. Let's assign values to the parameters of the perceptron (w and b): a = 0 $w_1 = 1; w_2 = 1; b = -1;$ 



- Initialize weights (including biases) to small random values, set *epoch-I*.
   Choose a random example (input-output pair (p,t))from the training set
   Calculate the response of the network for this example a (also called
   network activation) asbaching webb.setII
   Compute the output error e-ta
- Update weights: Add a matrix  $\Delta W$  to the weight matrix W, which is proportional to the product ep<sup>T</sup> between the error vector and the input:
- W<sup>new</sup> = W<sup>old</sup> + ep<sup>T</sup> from example
- Add a vector Ab to the bias vector b, which is proportional to the error

 $\mathbf{b}^{new} = \mathbf{b}^{old} + \mathbf{e}$ 

excluse 6. Repeat steps 2-5 (by choosing another example from the training data) 7. At the end of each *epoch* check if the stopping criteria is satisfied: all examples are correctly classified or a maximum number of epochs is reached; if yes-stop, eise *epoch++* and repeat from 2.

- Limitations: Binary Output if training examples are linearly separable guarantees a soln in finite # of steps. Doesn't try to find 'optimal' line.
- (The weight vector is orgthoganal to the decision boundary.)

- BACKPROPAGATION ALGORITHM (Powerful, Can learn Non-Linear Decision Boundaries, but hard to Tune) (Uses Steepest descent algorithm for minimizing the mean square error)
- Multi-Layer NNs trained w/ backpropagation popular. To learn input-output (I/O) mapping Error F() formulated (eg.  $\Sigma$  of squared errors between target & actual output) & use a learning rule that Minimizes this error.
- Sum of Squared Errors (E) is a classical measure of error
- E for a single training example over all output neurons d<sub>i</sub> desired, a<sub>i</sub> actual network output for <u>output neuron</u> *i*

# $E = \frac{1}{2}\sum_{i}e_{i}^{2} = \frac{1}{2}\sum_{i}(d_{i} - a_{i})^{2}$ a Feedforward netowrk, a Fully conected network (typically), weights

initialized to small rand. values. 5) Each neuton (except the input neurons) computes the weighed sum of its inputs, and then applies a differentiable transfer function  $a = f(\mathbf{w}\mathbf{p} + b)$ -248

⋛⋝→∕/→ 

any differentiable transfer function f can be used, i.e. the derivative should exist for every point;
 most frequently used: sigmoid and tan-sigmoid (hyperbolic tangent sigmoid)



 Numerical data requires no Encoding, unlike nominal - typically binary encoded. Output Encoding:

- •1 output neuron different output values represent different classes, e.g. <0.2 class 1, > <0.8 class 2, in between ambiguous class (class 3) istributed (binary, 1-of-n) encoding typically used in multi class problet Number of outputs = number of classes Example 3 classes, 3 output neurons; class 1 is represented as 1 0 0, class 2 as 0 1 0 and class 3 as 0 0 1 Another representation of the targets; use 0.1 instead of 0 and 0.9 instead of 1 ativation for absention bio.
- v Provides more degree of freedom to represent the target function (n
- There are a set of the set of t
- Heuristic to start w/ : 1 hidden layer w/ n hidden neurons, n=(inputs
- +ouput\_nurons)/2 BackPropagation (BP) adjusts weights backwards by propagatong the
- weight a
- training examples are applied and a total error calculated. Backpropagation Pula - Summary

Dackpropagation Rule - Summary		$o_1 \bigcirc o_k$
$w_{_{N}}(t)$ - weight from node $p$ to node $q$ at time $t$ $w_{_{N}}(t+1)=w_{_{N}}(t)+\Delta w_{_{N}}$	δ <sub>q</sub> ⊕ q	<sup>w<sub>q1</sub></sup> δ <sub>q</sub> δ <sub>q</sub> w <sub>t</sub>
$\Delta w_{pq} = \eta \cdot \delta_q \cdot o_p  \text{- weight change}$	o <sup>b</sup> Ob	0 <sub>p</sub> (
. The weight shange is prepartional to the output activation of neuron n and		

- The weight change  $\delta$  is calculated in 2 different ways
- q is an output neuron  $\delta_q = (d_q o_q) \cdot (ret_q)$ • q is a hidden neuron  $\delta_q = f'(net_q) \sum_i w_{qi} \delta_i$  (*i* is over the nodes in the layer above q

# Backpropagation - Example (cont. 1)

Forward pass for ex. 1 - calculate the outputs  $o_6$  and o  $o_1=0.6, o_2=0.1$ , target output 10, i.e. class 1 • Activations of the hidden neurons:

- $net_3 = o_1 * w_{13} + o_2 * w_{23} + b_3 = 0.6 * 0.1 + 0.1 * (-0.2) + 0.1 = 0.14$   $o_3 = 1/(1 + e^{-net3}) = 0.53$
- $\begin{array}{l} net_{4} = o_{1} \ ^{*}w_{14} + \ o_{2} \ ^{*}w_{24} + b_{4} = 0.6 \ ^{*}0 + 0.1 \ ^{*}0.2 + 0.2 = 0.22 \\ o_{4} = 1/(1 + e^{\cdot net4}) = 0.55 \end{array}$
- $\begin{array}{l}t_{5} = o_{1} \ ^{*}w_{15} + \ o_{2} \ ^{*}w_{25} + b_{5} = 0.6 \ ^{*}0.3 + 0.1 \ ^{*}(-0.4) + 0.5 = 0.64 \\ = 1/(1 + e^{-nt5}) = 0.65\end{array}$
- $\begin{array}{c} net_{7}=o_{3}*w_{37}+o_{4}*w_{47}+o_{5}*w_{57}+b_{7}=0.53*0.2+0.55*(-0.1)+0.65*(-0.2)+0.6=0.52\\ o_{7}=1/(1+e^{-net_{7}})=0.63\\ \end{array}$

## **Backpropagation** – Example (cont. 2) Backward pass for ex. 1

• Calculate the output errors  $\delta_6$  and  $\delta_7$  (note that  $d_6=1$ ,  $d_7=0$  for class 1)  $\begin{array}{c} \bullet \ \delta_6 = (d_6 - o_6) \ast \ o_6 \ast (1 - o_6) = (1 - 0.53) \ast 0.53 \ast (1 - 0.53) = 0.12 \\ \bullet \ \delta_7 = (d_7 - o_7) \ast \ o_7 \ast (1 - o_7) = (0 - 0.63) \ast 0.63 \ast (1 - 0.63) = -0.15 \end{array}$ 

1s (η=0.1) 

 $\Delta w_{37} = \eta * \delta_7 * o_3 = 0.1 * - 0.15 * 0.53 = -0.008$ 

a = 1

×

 $w_{37}^{new} = w_{37}^{old} + \Delta w_{37} = 0.2-0.008=-0.19$ Similarly for  $w_{46}^{new}$ ,  $w_{47}^{new}$ ,  $w_{56}^{new}$  and  $w_{57}^{new}$ 

apper 1 tenerone For the biases  $b_g$  and  $b_\gamma$  (remember: biases are weights with inp  $\Delta b_g = \eta + \delta_g + 1 = 0.1 + 0.12 = 0.012$   $b_\mu^{aee} = b_\mu^{aee} + \Delta b_g = -0.1 + 0.012 = 0.012$ Similarly for  $B_\gamma$  Koprinka, irrangit myd.edu.at COMP310874608 AL, week 10, 2011 with input 1):

- Every Boolean f() of incputs can be represented by network with a single en laver
- hidden layer. Any continuous (f) can be approximated w/ arbitrary small error by a network w/ 1 hidden layer. Any (f) (inc. discontinuous) able to arbitral small error by a network w/ 2 hidden layers. **Overfitting**, occurs w/ **Noise**, OR when # of free (trainable) params is bigger than # of training examples. OR network been trained too long. Soln: Use network that's 'just large enuff', network shouldn't have more free nerres than there are training one.
- Som: Use network thats 'just large enum', network shound it have more free params than there are training egs. Validation Set can be used to STOP training if error increases for a pre-specified # of iterations, the weights,bias' at the min. are returned. Prob w Viaidation Sets: Small tast ests. Son 2x. KFold Cross Validation, get mean number of optimum epochs. Final Run: train
- network on ep\_mean.

## With Gradient Descent:

- Small Learning Rate = slow Convergence. Large Learning Rate = oscillation, overshooting of minimum. Momentum used to stabilize the algorithm.
  - First Order Linear Filter

# Oscillation in the filter output y(k) is less than the

oscillation in the filter input w(k)

- As the momentum coefficient increases, the oscillation in the output is reduced
- The average filter output is the same as the average filter input
  - Although as the momentum increases the filter output is slow to respond

# The filter tends to reduce the amount of oscillation, while still tracking the average value

- (All layers have bias' linput) (BP iterates till it minimizes the sum of the squared errors of the output
- values over all training examples)

# Support Vector Machines

- Maximize margin.
   Transform data into a higher dimensional space where it's more likely to be linearly separable
- ( not to high tho, overfitting danger ) 3. Kernel Trick - Do calculations in the original, not the new

Unsupervised > Clustering

– number of clusters

Issues: • data should be normalized

don't∂

SSF

cluster

K-medoids

Repeat

(don't know class labels, may not know #classes, want to group similar egs) [Single link (MIN), Complete link (MAX), Average link (avg distance between each element in one cluster to each ele. in other) Centroid,

between each element in one custor x between each element in the formula of the f

D(xi) -mean-squared distance from the points in the cluster *i* to the center

 D(x,y) - astace between the centers of cluster rand / [Types: + Partitional - creates one set of cluster + ilerarchical • Density-based • Model-Based (generative) • Fuzzy Clustring]

 K.Means Clustering Algorithm

 Requires K, Houlsters, to be specified.

 • elect K points as initial centrols • (Repeat: • form K clusters based on elect K points as initial centrols • (Repeat: • form K clusters based on

closest centroid · recompute centroid of each cluster) Until: Centroids

Issues: • data should be normalized • momial data needs to 2 h o humeric • #epochs for convergence typically much smaller than #points, converges quickly -) • Stopping Criterion usually e.g. <1%, not ==, • SENSITIVE to choice of Initial Seeds (could run several times w/ ∂ initial centroids)

Not sensitive to order of input egs :) Doesn't work well for clusters w/ non-spherical / non-convex shapes.

Doesn't work well w/ data containing Outliners (Soln: preprocess, cut Deent Work wen w von some outliers) Space: Modest; O( (m+k)n ), m=#egs, n=#attributes, k=#clusters Time: Expensive; O(Hmn), t=#iterations Not Practical for Large datasets.

Can be viewed as an optimization problem; find k clusters that minimize

May find local minimum instead of global. Variations. • improving chances of finding global min. ; ∂ ways to initialize., allow splitting&merging of clutsers. • can be used for hierarchical clustering. w/ k=2 & recursively ;|| w/ each

- Use medoid instead of cluster means.
 - Reduces sensitivity to outliers
 Select K points (items) as the initial medoids

other non-medoid

Form K clusters by assigning all items to the closest medoid

Re-compute the medoids for each cluster (search for better medoids): • Initial (parent) state = current set of medoids

Evaluate the children states – are they better than the paren

evaluation function: cost=sum of absolute distances (item, medoid). • Choose the best state (set of medoids with the smallest cost)

Until medoids don't change or cost doesn't decrease

Computationally expensive, not suitable for large databases: Time: O(n(n-k)), Space: $(O(n^2)$ -needs proximity matrix). Doesn't depend on order of examples.

tearest weighoor Uustering Algorithm 1 pass algorithm, Partitional algorithm. Sensitive to Input Order. Puts ite, s in cluster of itself, single-link distance between item and new cluster, threshold t determines merging/creation. Space & Time: O(n\*2), n=#items

Suitable for domains w/ natural nesting relationships between clusters

- Computationally Expensive, limiting applicability to high dimensional data

dendrogram. Time: O(n<sup>2</sup>), n levels, at each of them, n<sup>2</sup> proximity matrix must be searched & updated (reducable to O(n<sup>2</sup>2 .logn) if distances store in

Not incrementa, assumes an usia saw.
 Aglomerative Clustering
 bottom up, merges clusters iteratively (wisingle-link (min) clustering))
 - d=0 • compute proximity matrix + let each data pt be a cluster • (Increment d, merge clusters w/ dist ≤ d, update Proximity Matrix [DRAW]

using Complete Link, is less sensitive to Noise and Outliers than Single

Divisive Clustering - top-down, splits cluster iteratively. reverse of agglomerative, less

Space: O(n^2), n=#items, to store proximity distance matrix &

Can be used when oly Distances are given and not raw data.

Nearest Neighbor Clustering Algorithm

(Hierarchical Clustering) <sup>0</sup> A B C D

Not Incremental, assumes all data static.

Link. Generating more compact clusters.

sorted list.

:||)

Generate the children states by swapping each medoid with each

pose the best state (set of medoids with the smallest cost)

Hill-climbing search for the set of medoids minimizing the distance from an example to a medoid

 $DB = \frac{1}{c} \sum_{i=1}^{c} \max_{j,j \neq i} = \left[ \frac{D(x_i) + D(x_j)}{D(x_i, x_j)} \right]^{-1}$ 

 $D(x_i, x_j)$  – distance between the centers of cluster *i* and *j* 

# higher di. space.

# $\overline{(\mathbf{X}(\mathbf{u},\mathbf{v}) = \Phi(\mathbf{u}) \cdot \Phi(\mathbf{v}) \neq (\mathbf{u} \cdot \mathbf{v})^2}$ "Kernel function of the new vector and the support vectors", instead of dot

Compared to Perceptrons ('linear only'), and Backpopragation NNs ('tuning,localMinima') (which SVM can

reduce too), SVM's: are relatively efficient training algorithms

SVM: an optimization problem using Lagrange multipliers (/\)

The optimization problem can be transformed into its dual

 $\max \mathbf{w}(\lambda) = \sum_{i=1}^{N} \lambda_i - \frac{1}{2} \sum_{i,j=1}^{N} \lambda_i \lambda_j \sqrt{j} \sqrt{\mathbf{x}_i \cdot \mathbf{x}_j}$  Dot product of pairs of raining vectors

 $f = \mathbf{w} \cdot \mathbf{z} + b = \sum_{i=1}^{N} \lambda_i y_i \mathbf{x} + b$  Dot product of the new vector and the support

sign(f) i.e. the new example belongs to class 1, if f>0 or class -1 if f<0

 $\mathbf{w} = \sum_{i=1}^{N} \lambda_{i} y_{i} \mathbf{x}_{i}$ 

Soft Margin: allows some misclassifications. Tradeoff

between margin width & #misclassifications. Soln is same as

w/ hard margins but there is an upper bound C on values of /

Ensemble of Classifiers [Works when individual classifiers are highly accurate and diverse (uncorrelated, don't make same mistake. Generated by manipulating ...,) & when base classifiers are good nuff i.e. better than random guessing]

Creates M Bootstrap samples • each sample used to build a classifier •

Creates M Bootstrap samples - each sample used to build a classifie classify a new gb y getting majority vote
 Effective for Unstable classifiers, (small *∂*s in the training set results in large *∂*s in predictions, e.g. DTs, Neural Networks).
 May slightly decay performance of stable classifiers (e.g. k-nn).
 Applicable to regression (votes avg'd)

Combo of weighted votes - each training weight has associated weight - higher the weight, more imporant the eg during training. - training eg weights init = 1 • generate classifier • correctly classified egs decrease in weight + vice versa • repeat • finally combine the M hypothese, each weighed according to performance on training set Adaboost - typically performs better than individual classes. If beca learning endertism is a weak learning decitime than 40-Boost

- If base learning algorithm is a weak learning algorithm, then AdaBoost will return a hypothesis that classifies the Training data Perfectly for

Win return a hyboresia via classifies de retening data renecity tor Large enough M. classifiers 2 'complex' Boosting allow building a powerful combined classifier from Very simple ones, eg. Simple DTs generated by 1R.

Bagging vs Boosting
- Similarities: • Use voting (for classification) and averaging (for prediction)

to combine the outputs of the individual learners. . Combines models of

the Same Type Differences: •Bagging builds individual models separately, Boosting builds them iteratively. • Bagging weighs opinions equally, Boosting weights by performance. Boosting, typically more Accurate. but Boosting more sensitive to Noise.

(this generates diversity, reducing correlation) Proven that RF does not overfit. RF Faster than Adaboost, gives comparable accuracy results.

wanipudating ceaning algorithms applied to same dataset but w/ *∂* params (e.g. NNs w/ *∂* architecture / params ). Train them on same training data to create M classifiers, which can output a Majority Vote.

Instead of voting, uses a (level-1) metalearner to learn which (level-0)

Separates training data into training, validation sets, trains level-0 classifiers, applies validation set to classifiers and use the predictions to build training data for level-1 model above - could use cross validation instead of training & validation sets. Slow,

Can be applied to numeric predictions, instead of a class-value a numeric target value is attached to level-1 training eg

level-0 base learners do most work. level-1 can be just a simple

Training vectors class

vectors

that can learn non-linear decision boundaries.

to maximize the margin of the decision boundary.

Linear Constraint:  $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \ge 1, \forall i(1)$ 

- product) ass of usable f()s K) cer's T:
- Scales well in high dimensions.
   Multi-Class problems need to be Transformed to 2 class problems. (Slows them down).

Margin Width: d=2/||w||

w: the optimal decision boundary:

[Enlarges Hypothesis Space]

(Manipulating Training Data)

Combo of weighted votes

the Same Type

(Manipulating Attributes)

(Using a Meta-Learner)

base classifiers are reliable.

but allows level-0 full-use of data.

Stacking

classifier.

Random Forest
 Bagging & Random selection of features

(Manipulating Learning Algorithm)

Majority Vote

Boosting

Bagging "Bootstrap Aggregation

) input meneore G