

CSE 592 Applications of Artificial Intelligence

Henry Kautz
Winter 2003

Today's Agenda

- Inductive learning
- Decision trees
break
- Bayesian learning
- Neural nets

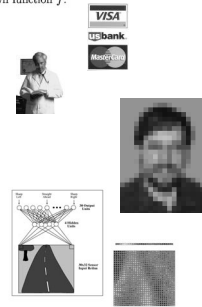
Inductive Learning

Supervised Learning

- **Given:** Training examples $(x, f(x))$ for some unknown function f .
- **Find:** A good approximation to f .

Example Applications

- **Credit risk assessment**
 x : Properties of customer and proposed purchase.
 $f(x)$: Approve purchase or not.
- **Disease diagnosis**
 x : Properties of patient (symptoms, lab tests)
 $f(x)$: Disease (or maybe, recommended therapy)
- **Face recognition**
 x : Bitmap picture of person's face
 $f(x)$: Name of the person.
- **Automatic Steering**
 x : Bitmap picture of road surface in front of car.
 $f(x)$: Degrees to turn the steering wheel.



Appropriate Applications for Supervised Learning

- **Situations where there is no human expert**
 x : Bond graph for a new molecule.
 $f(x)$: Predicted binding strength to AIDS protease molecule.
- **Situations where humans can perform the task but can't describe how they do it.**
 x : Bitmap picture of hand-written character
 $f(x)$: Ascii code of the character
- **Situations where the desired function is changing frequently**
 x : Description of stock prices and trades for last 10 days.
 $f(x)$: Recommended stock transactions
- **Situations where each user needs a customized function f**
 x : Incoming email message.
 $f(x)$: Importance score for presenting to user (or deleting without presenting).

A Learning Problem



Example	x_1	x_2	x_3	x_4	y
1	0	0	1	0	0
2	0	1	0	0	0
3	0	0	1	1	1
4	1	0	0	1	1
5	0	1	1	0	0
6	1	1	0	0	0
7	0	1	0	1	0

Hypothesis Spaces

- **Complete Ignorance.** There are $2^{16} = 65536$ possible boolean functions over four input features. We can't figure out which one is correct until we've seen every possible input-output pair. After 7 examples, we still have 2^9 possibilities.

x_1	x_2	x_3	x_4	y
0	0	0	0	?
0	0	0	1	?
0	0	1	0	0
0	0	1	1	1
0	1	0	0	0
0	1	0	1	0
0	1	1	0	?
0	1	1	1	?
1	0	0	0	?
1	0	0	1	?
1	0	1	0	?
1	0	1	1	?
1	1	0	0	0
1	1	0	1	?
1	1	1	0	?
1	1	1	1	?

Hypothesis Spaces (2)

- **Simple Rules.** There are only 16 simple conjunctive rules.

Rule	Counterexample
$\Rightarrow y$	1
$x_1 \Rightarrow y$	3
$x_2 \Rightarrow y$	2
$x_3 \Rightarrow y$	1
$x_4 \Rightarrow y$	7
$x_1 \wedge x_2 \Rightarrow y$	3
$x_1 \wedge x_3 \Rightarrow y$	3
$x_2 \wedge x_3 \Rightarrow y$	3
$x_2 \wedge x_4 \Rightarrow y$	3
$x_3 \wedge x_4 \Rightarrow y$	4
$x_1 \wedge x_2 \wedge x_3 \Rightarrow y$	3
$x_1 \wedge x_2 \wedge x_4 \Rightarrow y$	3
$x_1 \wedge x_3 \wedge x_4 \Rightarrow y$	3
$x_2 \wedge x_3 \wedge x_4 \Rightarrow y$	3
$x_1 \wedge x_2 \wedge x_3 \wedge x_4 \Rightarrow y$	3

No simple rule explains the data. The same is true for simple clauses.

Hypothesis Space (3)

- **m-of-n rules.** There are 32 possible rules (includes simple conjunctions and clauses).

variables	Counterexample			
	1-of	2-of	3-of	4-of
{ x_1 }	3	-	-	-
{ x_2 }	2	-	-	-
{ x_3 }	1	-	-	-
{ x_4 }	7	-	-	-
{ x_1, x_2 }	3	3	-	-
{ x_1, x_3 }	4	3	-	-
{ x_1, x_4 }	6	3	-	-
{ x_2, x_3 }	2	3	-	-
{ x_2, x_4 }	2	3	-	-
{ x_3, x_4 }	4	4	-	-
{ x_1, x_2, x_3 }	1	3	3	-
{ x_1, x_2, x_4 }	2	3	3	-
{ x_1, x_3, x_4 }	1	3	3	-
{ x_2, x_3, x_4 }	1	5	3	-
{ x_1, x_2, x_3, x_4 }	1	5	3	3

Bingo!
Is it necessarily the answer?

Hypothesis Space (3)

- **m-of-n rules.** There are 32 possible rules (includes simple conjunctions and clauses).

variables	Counterexample			
	1-of	2-of	3-of	4-of
{ x_1 }	3	-	-	-
{ x_2 }	2	-	-	-
{ x_3 }	1	-	-	-
{ x_4 }	7	-	-	-
{ x_1, x_2 }	3	3	-	-
{ x_1, x_3 }	4	3	-	-
{ x_1, x_4 }	6	3	-	-
{ x_2, x_3 }	2	3	-	-
{ x_2, x_4 }	2	3	-	-
{ x_3, x_4 }	4	4	-	-
{ x_1, x_2, x_3 }	1	3	3	-
{ x_1, x_2, x_4 }	2	3	3	-
{ x_1, x_3, x_4 }	1	3	3	-
{ x_2, x_3, x_4 }	1	5	3	-
{ x_1, x_2, x_3, x_4 }	1	5	3	3

Bingo!
Is it necessarily the answer?

$$x_4 \wedge \neg x_3$$

Bias in Learning

- Hypothesis space
- Preferences over hypothesis
- Other prior knowledge

Without bias learning is impossible!

Terminology

- **Training example.** An example of the form $(x, f(x))$.
- **Target function (target concept).** The true function f .
- **Hypothesis.** A proposed function h believed to be similar to f .
- **Concept.** A boolean function. Examples for which $f(x) = 1$ are called **positive examples** or **positive instances** of the concept. Examples for which $f(x) = 0$ are called **negative examples** or **negative instances**.
- **Classifier.** A discrete-valued function. The possible values $f(x) \in \{1, \dots, K\}$ are called the **classes** or **class labels**.
- **Hypothesis Space.** The space of all hypotheses that can, in principle, be output by a learning algorithm.
- **Version Space.** The space of all hypotheses in the hypothesis space that have not yet been ruled out by a training example.

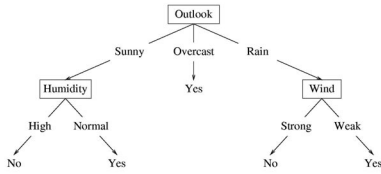
Decision Trees

Training Examples

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Decision Tree Hypothesis Space

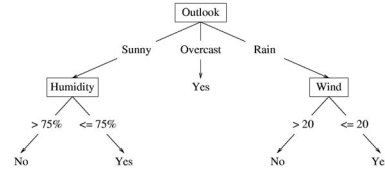
- **Internal nodes** test the value of particular features x_j and branch according to the results of the test.
- **Leaf nodes** specify the class $h(x)$.



Suppose the features are **Outlook** (x_1), **Temperature** (x_2), **Humidity** (x_3), and **Wind** (x_4). Then the feature vector $\mathbf{x} = (\text{Sunny}, \text{Hot}, \text{High}, \text{Strong})$ will be classified as **No**. The **Temperature** feature is irrelevant.

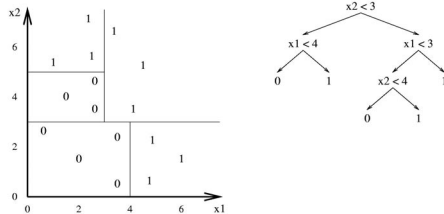
Decision Tree Hypothesis Space

If the features are continuous, internal nodes may test the value of a feature against a threshold.

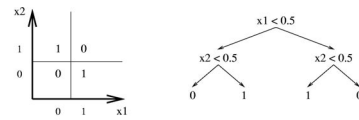


Decision Tree Decision Boundaries

Decision trees divide the feature space into axis-parallel rectangles, and label each rectangle with one of the K classes.



Decision Trees Can Represent Any Boolean Function



The tree will in the worst case require exponentially many nodes, however.

Variable-sized hypothesis space

Number of possible hypotheses grows with depth of tree

Top-Down Induction of Decision Trees

Main loop:

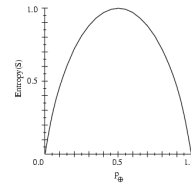
1. $A \leftarrow$ the "best" decision attribute for next *node*
2. Assign A as decision attribute for *node*
3. For each value of A , create new descendant of *node*
4. Sort training examples to leaf nodes
5. If training examples perfectly classified, Then STOP, Else iterate over new leaf nodes

Which attribute is best?



How can this algorithm be viewed as a state-space search problem?

Entropy



- S is a sample of training examples
- p_{\oplus} is the proportion of positive examples in S
- p_{\ominus} is the proportion of negative examples in S
- Entropy measures the impurity of S

$$Entropy(S) \equiv -p_{\oplus} \log_2 p_{\oplus} - p_{\ominus} \log_2 p_{\ominus}$$

Entropy

$Entropy(S)$ = expected number of bits needed to encode class (\oplus or \ominus) of randomly drawn member of S (under the optimal, shortest-length code)

Why?

Information theory: optimal length code assigns $-\log_2 p$ bits to message having probability p .

So, expected number of bits to encode \oplus or \ominus of random member of S :

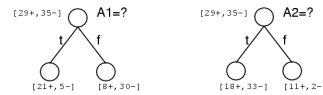
$$p_{\oplus}(-\log_2 p_{\oplus}) + p_{\ominus}(-\log_2 p_{\ominus})$$

$$Entropy(S) \equiv -p_{\oplus} \log_2 p_{\oplus} - p_{\ominus} \log_2 p_{\ominus}$$

Information Gain

$Gain(S, A)$ = expected reduction in entropy due to sorting on A

$$Gain(S, A) \equiv Entropy(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} Entropy(S_v)$$

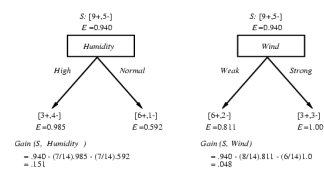


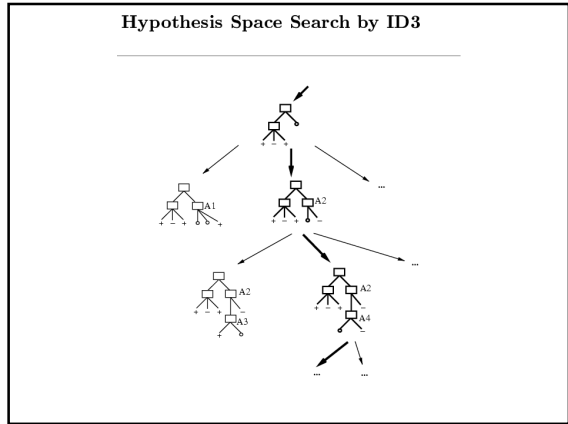
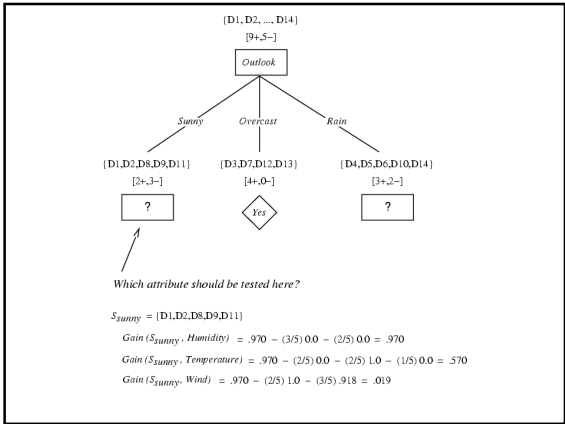
Training Examples

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Selecting the Next Attribute

Which attribute is the best classifier?





Hypothesis Space Search by ID3

- Hypothesis space is complete!
 - Target function surely in there...
- Outputs a single hypothesis (which one?)
 - Can't play 20 questions...
- No back tracking
 - Local minima...
- Statistically-based search choices
 - Robust to noisy data...
- Inductive bias: approx "prefer shortest tree"

Occam's Razor

Why prefer short hypotheses?

Argument in favor:

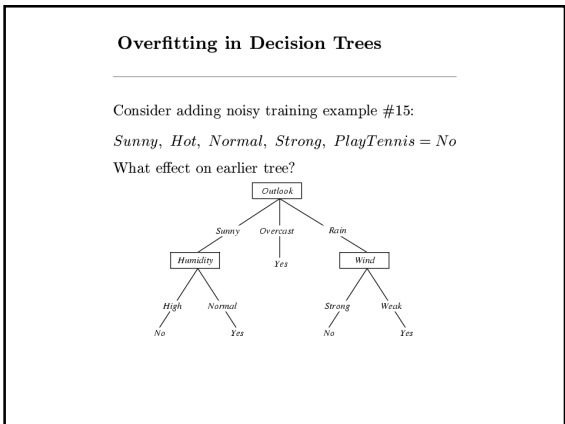
- Fewer short hyps. than long hyps.

→ a short hyp that fits data unlikely to be coincidence

→ a long hyp that fits data might be coincidence

Argument opposed:

- There are many ways to define small sets of hyps
- e.g., all trees with a prime number of nodes that use attributes beginning with "Z"
- What's so special about small sets based on size of hypothesis??



Overfitting

Consider error of hypothesis h over

- training data: $error_{train}(h)$
- entire distribution \mathcal{D} of data: $error_{\mathcal{D}}(h)$

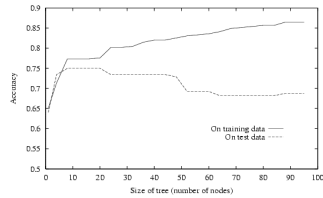
Hypothesis $h \in H$ **overfits** training data if there is an alternative hypothesis $h' \in H$ such that

$$error_{train}(h) < error_{train}(h')$$

and

$$error_{\mathcal{D}}(h) > error_{\mathcal{D}}(h')$$

Overfitting in Decision Tree Learning



Avoiding Overfitting

How can we avoid overfitting?

- stop growing when data split not statistically significant
- grow full tree, then post-prune

How to select “best” tree:

- Measure performance over training data
- Measure performance over separate validation data set
- MDL: minimize $size(tree) + size(misclassifications(tree))$

Reduced-Error Pruning

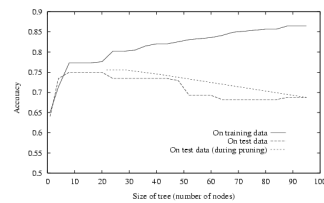
Split data into *training* and *validation* set

Do until further pruning is harmful:

1. Evaluate impact on *validation* set of pruning each possible node (plus those below it)
2. Greedily remove the one that most improves *validation* set accuracy

- produces smallest version of most accurate subtree
- What if data is limited?

Effect of Reduced-Error Pruning



Attributes with Costs

Consider

- medical diagnosis, *BloodTest* has cost \$150
- robotics, *Width_from_1ft* has cost 23 sec.

How to learn a consistent tree with low expected cost?

One approach: replace gain by

- Tan and Schlimmer (1990) $\frac{Gain^2(S, A)}{Cost(A)}$
- Nunez (1988) $\frac{2^{Gain(S, A)} - 1}{(Cost(A) + 1)^w}$

where $w \in [0, 1]$ determines importance of cost

Scaling Up

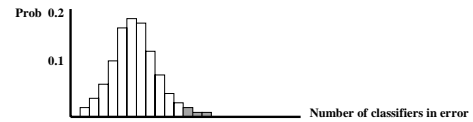
- ID3, C4.5, etc. assume data fits in main memory (OK for up to hundreds of thousands of examples)
- SPRINT, SLIQ: multiple sequential scans of data (OK for up to millions of examples)
- VFDT: at most one sequential scan (OK for up to billions of examples)

Ensembles of Classifiers

- Idea: instead of training one classifier (decision tree)
- Train k classifiers and let them vote
 - Only helps if classifiers disagree with each other
 - Trained on different data
 - Use different learning methods
- Amazing fact: can help a lot!

How voting helps

- Assume errors are independent
- Assume majority vote
- Probability majority is wrong = area under binomial dist



- If individual area is 0.3
- **Area under curve for ≥ 11 wrong is 0.026**
- Order of magnitude improvement!

Constructing Ensembles

- Bagging
 - Run classifier k times on m examples drawn randomly with replacement from the original set of n examples
- Cross-validated committees
 - Divide examples into k disjoint sets
 - Train on k sets corresponding to original *minus* $1/k$ -th
- Boosting (Shapire)
 - Maintain a probability distribution over set of training examples
 - On each iteration, use distribution to sample
 - Use error rate to modify distribution
 - Create harder and harder learning problems

Summary

- Inductive learning
- Decision trees
 - Representation
 - Tree growth
 - Heuristics
 - Overfitting and pruning
 - Scaling up
- Ensembles

Break!

Bayesian Learning

Preview

- Bayes' theorem
- MAP learners
- Bayes optimal classifier
- Naive Bayes learner
- Example: text classification
- Bayesian networks
- EM algorithm

Two Roles for Bayesian Methods

Practical learning algorithms:

- Naive Bayes learning
- Bayesian network learning
- Combine prior knowledge with observed data
- Require prior probabilities

Useful conceptual framework:

- "Gold standard" for evaluating other learners
- Tools for analysis

Bayes' Theorem

$$P(h|D) = \frac{P(D|h)P(h)}{P(D)}$$

- $P(h)$ = prior probability of hypothesis h
- $P(D)$ = prior probability of training data D
- $P(h|D)$ = probability of h given D
- $P(D|h)$ = probability of D given h

Choosing Hypotheses

Find most probable hypothesis given training data

Maximum a posteriori hypothesis h_{MAP} :

$$\begin{aligned} h_{MAP} &= \arg \max_{h \in H} P(h|D) \\ &= \arg \max_{h \in H} \frac{P(D|h)P(h)}{P(D)} \\ &= \arg \max_{h \in H} P(D|h)P(h) \end{aligned}$$

Assuming $P(h_i) = P(h_j)$ we can further simplify, and choose the *Maximum likelihood* (ML) hypothesis

$$h_{ML} = \arg \max_{h_i \in H} P(D|h_i)$$

Example

Does patient have cancer or not?

A patient takes a lab test and the result comes back positive. The test returns a correct positive result in only 98% of the cases in which the disease is actually present, and a correct negative result in only 97% of the cases in which the disease is not present. Furthermore, 0.008 of the entire population have this cancer.

$$\begin{aligned} P(\text{cancer}) &= \\ P(\neg\text{cancer}) &= \\ P(+|\text{cancer}) &= \\ P(-|\text{cancer}) &= \\ P(+|\neg\text{cancer}) &= \\ P(-|\neg\text{cancer}) &= \\ P(\text{cancer}|+) &= \end{aligned}$$

Brute-Force MAP Hypothesis Learner

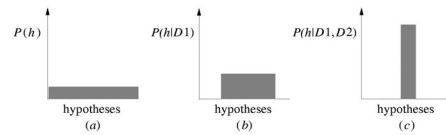
1. For each hypothesis h in H , calculate the posterior probability

$$P(h|D) = \frac{P(D|h)P(h)}{P(D)}$$

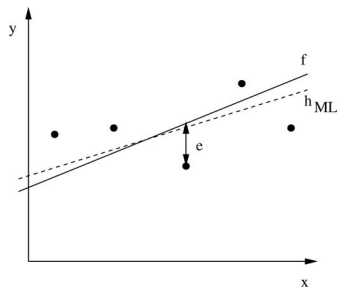
2. Output the hypothesis h_{MAP} with the highest posterior probability

$$h_{MAP} = \underset{h \in H}{\operatorname{argmax}} P(h|D)$$

Evolution of Posterior Probabilities



Learning a Real-Valued Function



Consider any real-valued target function f

Training examples $\langle x_i, d_i \rangle$, where d_i is noisy training value

- $d_i = f(x_i) + e_i$
- e_i is random variable (noise) drawn independently for each x_i according to some Gaussian distribution with mean=0

Then the maximum likelihood hypothesis h_{ML} is the one that minimizes the sum of squared errors:

$$h_{ML} = \underset{h \in H}{\operatorname{argmin}} \sum_{i=1}^m (d_i - h(x_i))^2$$

Most Probable Classification of New Instances

So far we've sought the most probable *hypothesis* given the data D (i.e., h_{MAP})

Given new instance x , what is its most probable *classification*? Not $h_{MAP}(x)$!

Consider:

- Three possible hypotheses:
 $P(h_1|D) = .4$, $P(h_2|D) = .3$, $P(h_3|D) = .3$
- Given new instance x ,
 $h_1(x) = +$, $h_2(x) = -$, $h_3(x) = -$
- What's most probable classification of x ?

Bayes Optimal Classifier

Bayes optimal classification:

$$\underset{v_j \in V}{\operatorname{argmax}} \sum_{h_i \in H} P(v_j|h_i)P(h_i|D)$$

Example:

$$\begin{aligned} P(h_1|D) &= .4, & P(-|h_1) &= 0, & P(+|h_1) &= 1 \\ P(h_2|D) &= .3, & P(-|h_2) &= 1, & P(+|h_2) &= 0 \\ P(h_3|D) &= .3, & P(-|h_3) &= 1, & P(+|h_3) &= 0 \end{aligned}$$

Classify instance D as:

Naive Bayes Classifier

Assume target function $f : X \rightarrow V$, where each instance x described by attributes $\langle a_1, a_2 \dots a_n \rangle$.

Most probable value of $f(x)$ is:

$$\begin{aligned} v_{MAP} &= \operatorname{argmax}_{v_j \in V} P(v_j | a_1, a_2 \dots a_n) \\ v_{MAP} &= \operatorname{argmax}_{v_j \in V} \frac{P(a_1, a_2 \dots a_n | v_j) P(v_j)}{P(a_1, a_2 \dots a_n)} \\ &= \operatorname{argmax}_{v_j \in V} P(a_1, a_2 \dots a_n | v_j) P(v_j) \end{aligned}$$

Naive Bayes assumption:

$$P(a_1, a_2 \dots a_n | v_j) = \prod_i P(a_i | v_j)$$

which gives

Naive Bayes classifier:

$$v_{NB} = \operatorname{argmax}_{v_j \in V} P(v_j) \prod_i P(a_i | v_j)$$

Naive Bayes Algorithm

Naive_Bayes_Learn(*examples*)

For each target value v_j

- $\hat{P}(v_j) \leftarrow$ estimate $P(v_j)$
- For each attribute value a_i of each attribute a
 $\hat{P}(a_i | v_j) \leftarrow$ estimate $P(a_i | v_j)$

Classify_New_Instance(x)

$$v_{NB} = \operatorname{argmax}_{v_j \in V} \hat{P}(v_j) \prod_{a_i \in x} \hat{P}(a_i | v_j)$$

Naive Bayes: Example

Consider *PlayTennis* again, and new instance

$\langle \text{Outlk} = \text{sun}, \text{Temp} = \text{cool}, \text{Humid} = \text{high}, \text{Wind} = \text{strong} \rangle$

Want to compute:

$$v_{NB} = \operatorname{argmax}_{v_j \in V} P(v_j) \prod_i P(a_i | v_j)$$

$$P(y) P(\text{sun}|y) P(\text{cool}|y) P(\text{high}|y) P(\text{strong}|y) = .005$$

$$P(n) P(\text{sun}|n) P(\text{cool}|n) P(\text{high}|n) P(\text{strong}|n) = .021$$

$$\rightarrow v_{NB} = n$$

Learning to Classify Text

Why?

- Learn which news articles are of interest
- Learn to classify web pages by topic

Naive Bayes is among most effective algorithms

What attributes shall we use to represent text documents?

Learning to Classify Text

Target concept *Interesting?* : Document $\rightarrow \{+, -\}$

1. Represent each document by vector of words:
one attribute per word position in document
2. Learning: Use training examples to estimate

- $P(+)$
- $P(-)$
- $P(\text{doc}|+)$
- $P(\text{doc}|-)$

Naive Bayes conditional independence assumption

$$P(doc|v_j) = \prod_{i=1}^{length(doc)} P(a_i = w_k | v_j)$$

where $P(a_i = w_k | v_j)$ is probability that word in position i is w_k , given v_j

One more assumption:

$$P(a_i = w_k | v_j) = P(a_m = w_k | v_j), \forall i, m$$

LEARN_NAIVE_BAYES_TEXT(*Examples*, V)

1. Collect all words & tokens that occur in *Examples*
 - *Vocabulary* \leftarrow all distinct words & tokens in *Examples*
2. Compute all probabilities $P(v_j)$ and $P(w_k | v_j)$
 - For each target value v_j in V do
 - *docs_j* \leftarrow *Examples* for which the target value is v_j
 - $P(v_j) \leftarrow \frac{|docs_j|}{|Examples|}$
 - *Text_j* \leftarrow concatenate all members of *docs_j*
 - $n \leftarrow$ total number of words in *Text_j* (counting duplicate words multiple times)
 - for each word w_k in *Vocabulary*
 - * $n_k \leftarrow$ number of times word w_k occurs in *Text_j*
 - * $P(w_k | v_j) \leftarrow \frac{n_k + 1}{n + |Vocabulary|}$

CLASSIFY_NAIVE_BAYES_TEXT(*Doc*)

- *positions* \leftarrow all word positions in *Doc* that contain tokens found in *Vocabulary*
- Return v_{NB} , where

$$v_{NB} = \underset{v_j \in V}{\operatorname{argmax}} P(v_j) \prod_{i \in positions} P(a_i | v_j)$$

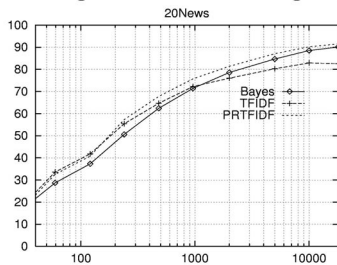
Example: 20 Newsgroups

Given 1000 training documents from each group
Learn to classify new documents according to which newsgroup it came from

comp.graphics	misc.forsale
comp.os.ms-windows.misc	rec.autos
comp.sys.ibm.pc.hardware	rec.motorcycles
comp.sys.mac.hardware	rec.sport.baseball
comp.windows.x	rec.sport.hockey
alt.atheism	sci.space
soc.religion.christian	sci.crypt
talk.religion.misc	sci.electronics
talk.politics.mideast	sci.med
talk.politics.misc	talk.politics.guns

Naive Bayes: 89% classification accuracy

Learning Curve for 20 Newsgroups



Accuracy vs. Training set size (1/3 withheld for test)

Learning Bayesian Networks

Several variants of this learning task

- Network structure might be *known* or *unknown*
- Training examples might provide values of *all* network variables, or just *some*

If structure known and no missing values,
it's as easy as training a Naive Bayes classifier

The EM Algorithm

Suppose structure known, variables partially observable

E.g., observe *ForestFire*, *Storm*, *BusTourGroup*, *Thunder*, but not *Lightning*, *Campfire* ...

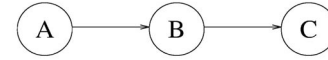
Initialize parameters ignoring missing information

Repeat until convergence:

E step: Calculate expected vals of unobserved variables, assuming current parameter values

M step: Calculate new parameter values to maximize probability of data (observed & estimated)

Example



Examples:	0	1	1
	1	0	0
	1	1	1
	1	?	0

Initialization: $P(B|A) = P(C|B) =$
 $P(A) = P(B|\neg A) = P(C|\neg B) =$

E-step: $P(? = 1) = P(B|A, \neg C) = \frac{P(A, B, \neg C)}{P(A, \neg C)} = \dots = 0$

M-step: $P(B|A) = P(C|B) =$
 $P(A) = P(B|\neg A) = P(C|\neg B) =$

E-step: $P(? = 1) = 0$ (converged)

Unknown Structure

Search:

- Initial state: empty network, prior network
- Operators: Add arc, delete arc, reverse arc
- Evaluation: Posterior probability

Bayesian Learning: Summary

- Optimal prediction
- Naive Bayes learner
- Text classification
- Bayesian networks
- EM algorithm

Neural Networks

Preview

- Perceptrons
- Gradient descent
- Multilayer networks
- Backpropagation

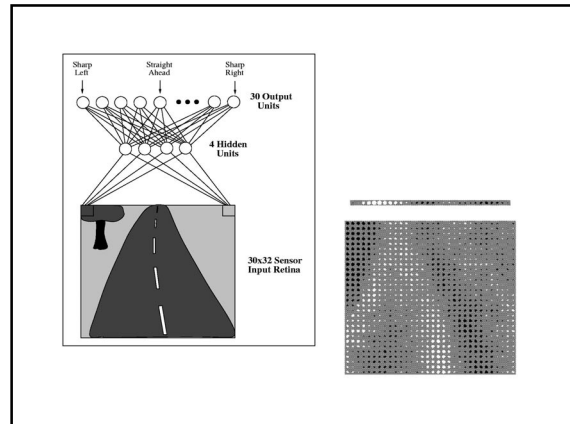
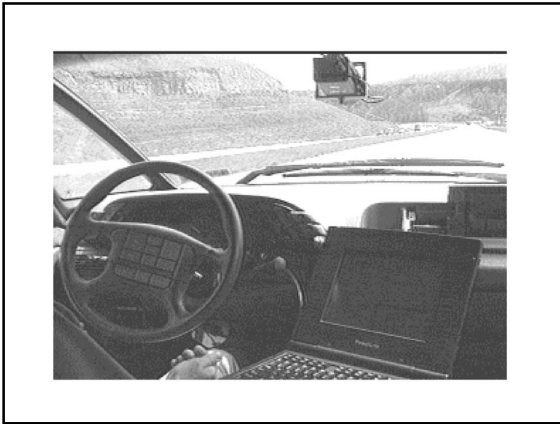
Connectionist Models

Consider humans:

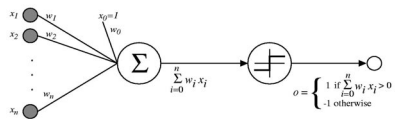
- Neuron switching time $\sim .001$ second
 - Number of neurons $\sim 10^{10}$
 - Connections per neuron $\sim 10^4-5$
 - Scene recognition time $\sim .1$ second
 - 100 inference steps doesn't seem like enough
- \Rightarrow Much parallel computation

Properties of neural nets:

- Many neuron-like threshold switching units
- Many weighted interconnections among units
- Highly parallel, distributed process
- Emphasis on tuning weights automatically



Perceptron

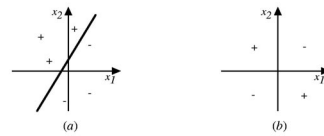


$$o(x_1, \dots, x_n) = \begin{cases} 1 & \text{if } w_0 + w_1 x_1 + \dots + w_n x_n > 0 \\ -1 & \text{otherwise.} \end{cases}$$

Sometimes we'll use simpler vector notation:

$$o(\vec{x}) = \begin{cases} 1 & \text{if } \vec{w} \cdot \vec{x} > 0 \\ -1 & \text{otherwise.} \end{cases}$$

Decision Surface of a Perceptron



Represents some useful functions

- What weights represent $g(x_1, x_2) = AND(x_1, x_2)$?

But some functions not representable

- All not linearly separable
- Therefore, we'll want networks of these...

Perceptron Training Rule

$$w_i \leftarrow w_i + \Delta w_i$$

where

$$\Delta w_i = \eta(t - o)x_i$$

Where:

- $t = c(\vec{x})$ is target value
- o is perceptron output
- η is small constant (e.g., 0.1) called *learning rate*

Perceptron Training Rule

Can prove it will converge if

- Training data is linearly separable
- η sufficiently small

Gradient Descent

To understand, consider simpler *linear unit*, where

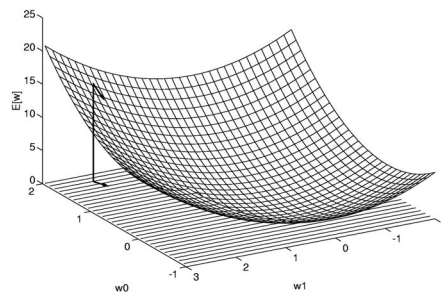
$$o = w_0 + w_1x_1 + \dots + w_nx_n$$

Let's learn w_i 's that minimize the squared error

$$E[\vec{w}] \equiv \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

Where D is set of training examples

Gradient Descent



Gradient:

$$\nabla E[\vec{w}] \equiv \left[\frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right]$$

Training rule:

$$\Delta \vec{w} = -\eta \nabla E[\vec{w}]$$

I.e.:

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i}$$

Gradient Descent

$$\begin{aligned} \frac{\partial E}{\partial w_i} &= \frac{\partial}{\partial w_i} \frac{1}{2} \sum_d (t_d - o_d)^2 \\ &= \frac{1}{2} \sum_d \frac{\partial}{\partial w_i} (t_d - o_d)^2 \\ &= \frac{1}{2} \sum_d 2(t_d - o_d) \frac{\partial}{\partial w_i} (t_d - o_d) \\ &= \sum_d (t_d - o_d) \frac{\partial}{\partial w_i} (t_d - \vec{w} \cdot \vec{x}_d) \\ \frac{\partial E}{\partial w_i} &= \sum_d (t_d - o_d)(-x_{i,d}) \end{aligned}$$

Gradient Descent

GRADIENT-DESCENT(*training_examples*, η)

Initialize each w_i to some small random value

Until the termination condition is met, Do

- Initialize each Δw_i to zero.
- For each $\langle \vec{x}, t \rangle$ in *training_examples*, Do
 - Input instance \vec{x} to unit and compute output o
 - For each linear unit weight w_i , Do

$$\Delta w_i \leftarrow \Delta w_i + \eta(t - o)x_i$$

- For each linear unit weight w_i , Do

$$w_i \leftarrow w_i + \Delta w_i$$

Summary

Perceptron training rule guaranteed to succeed if

- Training examples are linearly separable
- Sufficiently small learning rate η

Linear unit training rule uses gradient descent

- Guaranteed to converge to hypothesis with minimum squared error
- Given sufficiently small learning rate η
- Even when training data contains noise
- Even when training data not separable by H

Batch vs. Incremental Gradient Descent

Batch Mode Gradient Descent:

Do until convergence

1. Compute the gradient $\nabla E_D[\vec{w}]$
2. $\vec{w} \leftarrow \vec{w} - \eta \nabla E_D[\vec{w}]$

Incremental Mode Gradient Descent:

Do until convergence

For each training example d in D

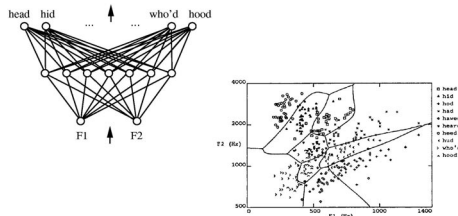
1. Compute the gradient $\nabla E_d[\vec{w}]$
2. $\vec{w} \leftarrow \vec{w} - \eta \nabla E_d[\vec{w}]$

$$E_D[\vec{w}] \equiv \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

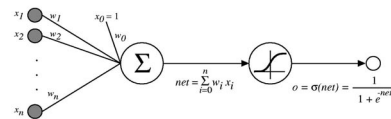
$$E_d[\vec{w}] \equiv \frac{1}{2} (t_d - o_d)^2$$

Incremental Gradient Descent can approximate *Batch Gradient Descent* arbitrarily closely if η made small enough

Multilayer Networks of Sigmoid Units



Sigmoid Unit



$\sigma(x)$ is the sigmoid function

$$\frac{1}{1 + e^{-x}}$$

Nice property: $\frac{d\sigma(x)}{dx} = \sigma(x)(1 - \sigma(x))$

We can derive gradient descent rules to train

- One sigmoid unit
- *Multilayer networks* of sigmoid units → Backpropagation

Error Gradient for a Sigmoid Unit

$$\begin{aligned} \frac{\partial E}{\partial w_i} &= \frac{\partial}{\partial w_i} \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2 \\ &= \frac{1}{2} \sum_d \frac{\partial}{\partial w_i} (t_d - o_d)^2 \\ &= \frac{1}{2} \sum_d 2(t_d - o_d) \frac{\partial}{\partial w_i} (t_d - o_d) \\ &= \sum_d (t_d - o_d) \left(-\frac{\partial o_d}{\partial w_i} \right) \\ &= -\sum_d (t_d - o_d) \frac{\partial o_d}{\partial net_d} \frac{\partial net_d}{\partial w_i} \end{aligned}$$

But we know:

$$\frac{\partial o_d}{\partial net_d} = \frac{\partial \sigma(net_d)}{\partial net_d} = o_d(1 - o_d)$$

$$\frac{\partial net_d}{\partial w_i} = \frac{\partial (\vec{w} \cdot \vec{x}_d)}{\partial w_i} = x_{i,d}$$

So:

$$\frac{\partial E}{\partial w_i} = -\sum_{d \in D} (t_d - o_d) o_d (1 - o_d) x_{i,d}$$

Let: $\delta_k = -\frac{\partial E}{\partial net_k}$

$$\begin{aligned} \frac{\partial E}{\partial net_j} &= \sum_{k \in \text{Outs}(j)} \frac{\partial E}{\partial net_k} \frac{\partial net_k}{\partial net_j} \\ &= \sum_{k \in \text{Outs}(j)} -\delta_k \frac{\partial net_k}{\partial net_j} \\ &= \sum_{k \in \text{Outs}(j)} -\delta_k \frac{\partial net_k}{\partial o_j} \frac{\partial o_j}{\partial net_j} \\ &= \sum_{k \in \text{Outs}(j)} -\delta_k w_{kj} \frac{\partial o_k}{\partial net_j} \\ &= \sum_{k \in \text{Outs}(j)} -\delta_k w_{kj} o_j (1 - o_j) \\ \delta_j &= -\frac{\partial E}{\partial net_j} = o_j (1 - o_j) \sum_{k \in \text{Outs}(j)} \delta_k w_{kj} \end{aligned}$$

Backpropagation Algorithm

Initialize all weights to small random numbers
Until convergence, Do

For each training example, Do

1. Input it to network and compute network outputs
2. For each output unit k

$$\delta_k \leftarrow o_k(1 - o_k)(t_k - o_k)$$

3. For each hidden unit h

$$\delta_h \leftarrow o_h(1 - o_h) \sum_{k \in \text{outputs}} w_{h,k} \delta_k$$

4. Update each network weight $w_{i,j}$

$$w_{i,j} \leftarrow w_{i,j} + \Delta w_{i,j}$$

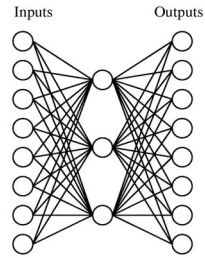
where $\Delta w_{i,j} = \eta \delta_j x_{i,j}$

More on Backpropagation

- Gradient descent over entire *network* weight vector
- Easily generalized to arbitrary directed graphs
- Will find a local, not necessarily global error minimum
 - In practice, often works well (can run multiple times)
- Often include weight *momentum* α

$$\Delta w_{i,j}(n) = \eta \delta_j x_{i,j} + \alpha \Delta w_{i,j}(n-1)$$
- Minimizes error over *training* examples
 - Will it generalize well to subsequent examples?
- Training can take thousands of iterations → slow!
- Using network after training is very fast

Learning Hidden Layer Representations



A target function:

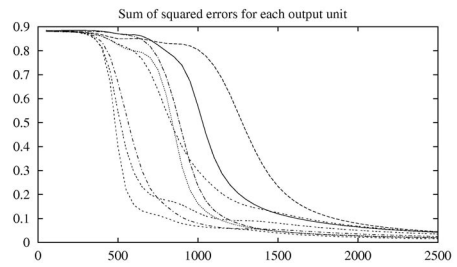
Input	Output
10000000	→ 10000000
01000000	→ 01000000
00100000	→ 00100000
00010000	→ 00010000
00001000	→ 00001000
00000100	→ 00000100
00000010	→ 00000010
00000001	→ 00000001

Can this be learned?

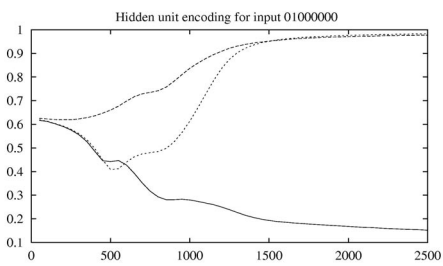
Learned hidden layer representation:

Input	Hidden Values	Output
10000000	→ .89 .04 .08	→ 10000000
01000000	→ .01 .11 .88	→ 01000000
00100000	→ .01 .97 .27	→ 00100000
00010000	→ .99 .97 .71	→ 00010000
00001000	→ .03 .05 .02	→ 00001000
00000100	→ .22 .99 .99	→ 00000100
00000010	→ .80 .01 .98	→ 00000010
00000001	→ .60 .94 .01	→ 00000001

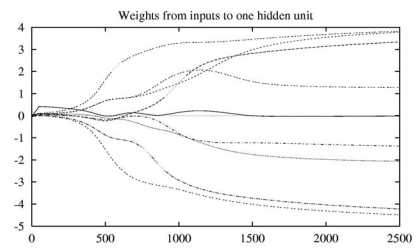
Training



Training



Training



Convergence of Backpropagation

Gradient descent to some local minimum

- Perhaps not global minimum...
- Add momentum
- Stochastic gradient descent
- Train multiple nets with different initial weights

Nature of convergence

- Initialize weights near zero
- Therefore, initial networks near-linear
- Increasingly non-linear functions possible as training progresses

Expressiveness of Neural Nets

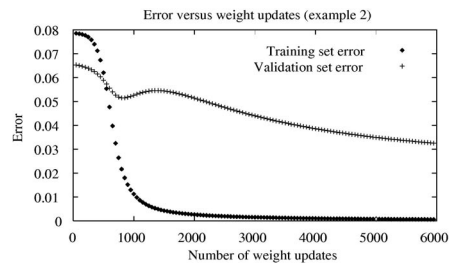
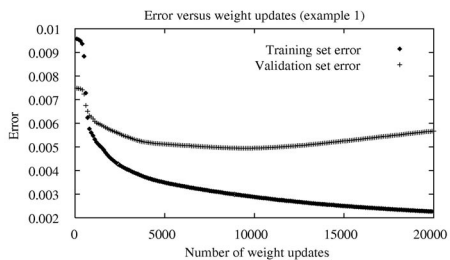
Boolean functions:

- Every Boolean function can be represented by network with single hidden layer
- But might require exponential (in number of inputs) hidden units

Continuous functions:

- Every bounded continuous function can be approximated with arbitrarily small error, by network with one hidden layer
- Any function can be approximated to arbitrary accuracy by a network with two hidden layers

Overfitting in Neural Nets



Overfitting Avoidance

Penalize large weights:

$$E(\vec{w}) \equiv \frac{1}{2} \sum_{d \in D} \sum_{k \in \text{outputs}} (t_{kd} - o_{kd})^2 + \gamma \sum_{i,j} w_{ji}^2$$

Train on target slopes as well as values:

$$E(\vec{w}) \equiv \frac{1}{2} \sum_{d \in D} \sum_{k \in \text{outputs}} \left[(t_{kd} - o_{kd})^2 + \mu \sum_{j \in \text{inputs}} \left(\frac{\partial t_{kd}}{\partial x_d^j} - \frac{\partial o_{kd}}{\partial x_d^j} \right)^2 \right]$$

Weight sharing

Early stopping

Neural Networks: Summary

- Perceptrons
- Gradient descent
- Multilayer networks
- Backpropagation