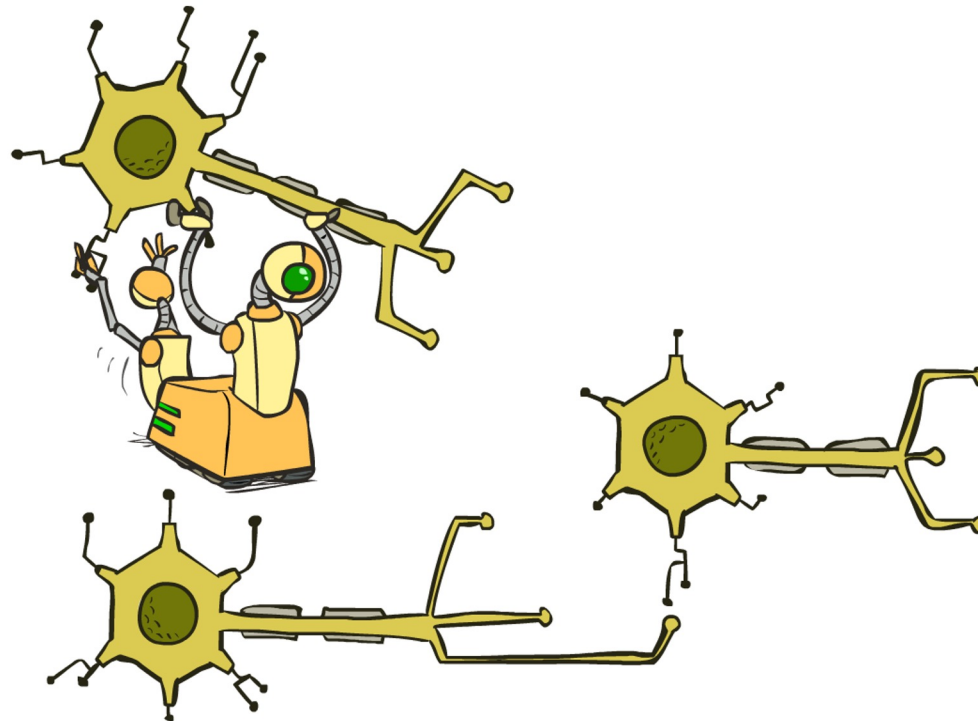# CS 188: Artificial Intelligence
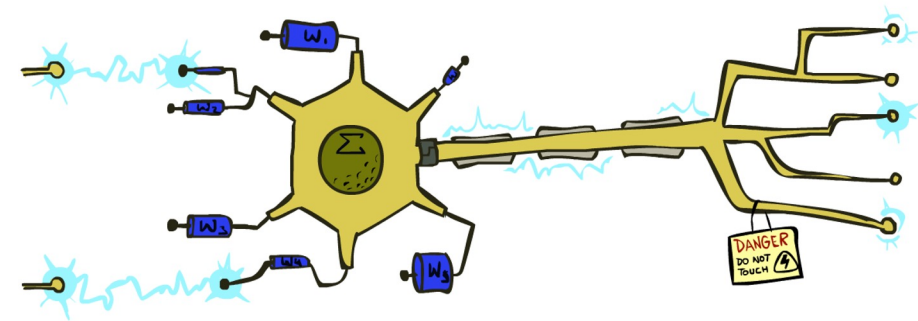
## Neural Networks

Summer 2024: Eve Fleisig & Evgeny Pobachienko
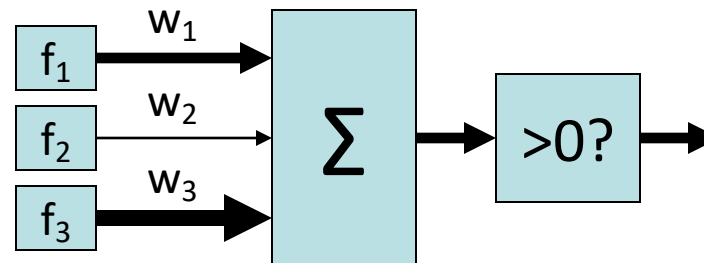
# Reminder: Linear Classifiers

- Inputs are feature values
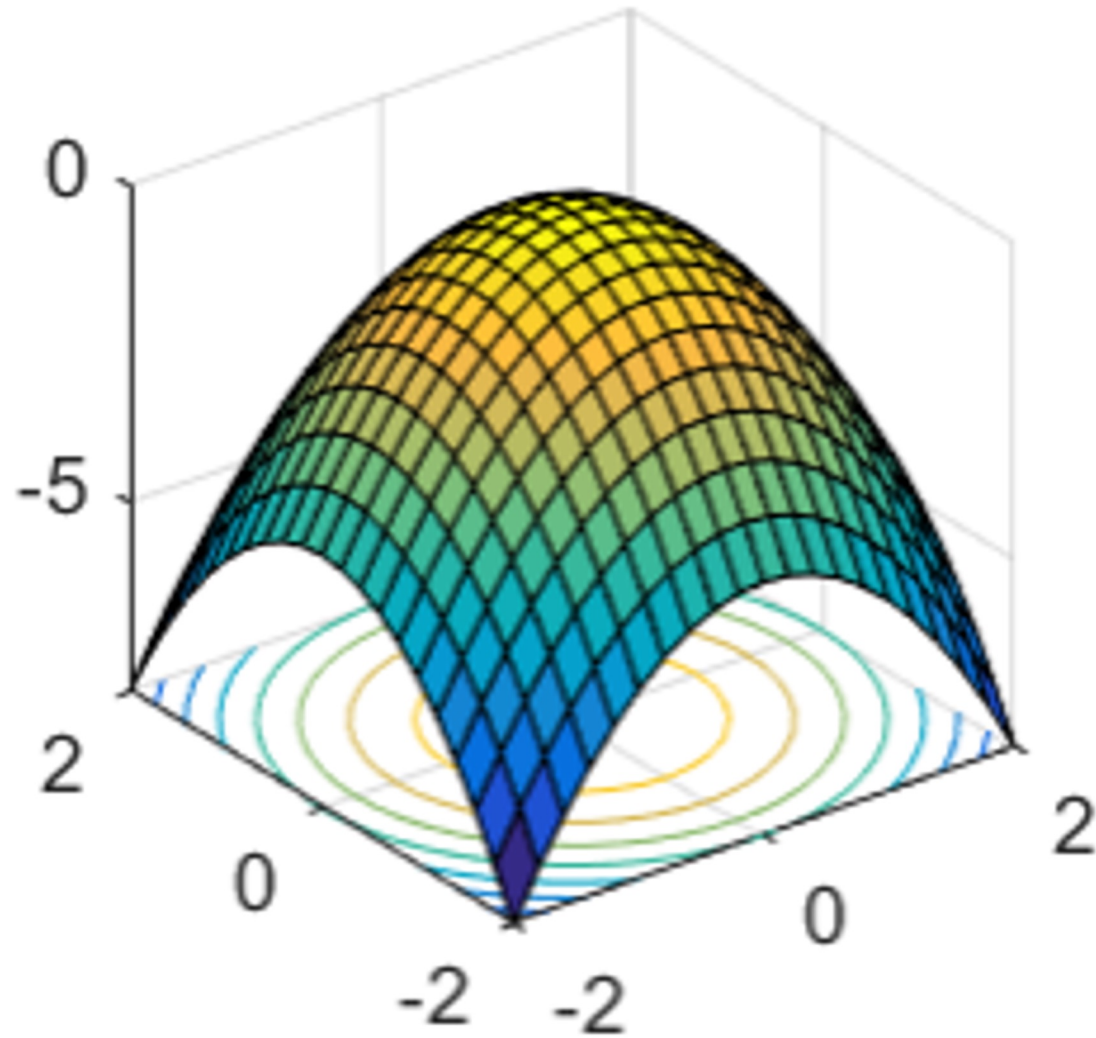- Each feature has a weight
- Sum is the activation

$$\text{activation}_w(x) = \sum_i w_i \cdot f_i(x) = w \cdot f(x)$$

- If the activation is:
  - Positive, output +1
  - Negative, output -1

# 2-D Optimization



Source: offconvex.org

# Neural Networks

# Multi-class Logistic Regression

- = special case of neural network



$$P(y_1|x;w) = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$P(y_2|x;w) = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$P(y_3|x;w) = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

# Deep Neural Network = Also learn the features!



$$P(y_1|x;w) = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$P(y_2|x;w) = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$P(y_3|x;w) = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

# Deep Neural Network = Also learn the features!



$$z_i^{(k)} = g(\sum_j W_{i,j}^{(k-1,k)} z_j^{(k-1)})$$

# Deep Neural Network = Also learn the features!



$$z_i^{(k)} = g(\sum_j W_{i,j}^{(k-1,k)} z_j^{(k-1)})$$

# Importance of Nonlinear Activation Functions

- What happens if we add more layers?
- $z_2 = W_2(W_1x + b_1) + b_2$
- $z_2 = W_2(W_1x + b_1) + b_2 = W_2W_1x + W_2b_1 + b_2 = W_{new}x + b_{new}$
- No gain to adding more linear layers!
- Idea: add nonlinearities to capture more complex relationships

# Deep Neural Network = Also learn the features!



$$z_i^{(k)} = g(\sum_j W_{i,j}^{(k-1,k)} z_j^{(k-1)})$$

**g = nonlinear activation function**

# Common Activation Functions

## Sigmoid Function



$$g(z) = \frac{1}{1 + e^{-z}}$$

$$g'(z) = g(z)(1 - g(z))$$

## Hyperbolic Tangent



$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$g'(z) = 1 - g(z)^2$$

## Rectified Linear Unit (ReLU)



$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$

# Deep Neural Network: Also Learn the Features!

- Training the deep neural network is just like logistic regression:

$$\max_w \ ll(w) = \max_w \sum_i \log P(y^{(i)}|x^{(i)}; w)$$

just w tends to be a much, much larger vector ☺

just run gradient ascent

+ stop when log likelihood of hold-out data starts to decrease

# Neural Networks Properties

- Theorem (Universal Function Approximators).  A two-layer neural network with a sufficient number of neurons can approximate any continuous function to any desired accuracy.

- Practical considerations
  - Can be seen as learning the features

  - Large number of neurons
    - Danger for overfitting
    - (hence early stopping!)

# Universal Function Approximation Theorem*

**Hornik theorem 1:** Whenever the activation function is *bounded and nonconstant*, then, for any finite measure $\mu$, standard multilayer feedforward networks can approximate any function in $L^p(\mu)$ (the space of all functions on $R^k$ such that $\int_{R^k} |f(x)|^p d\mu(x) < \infty$) arbitrarily well, provided that sufficiently many hidden units are available.

**Hornik theorem 2:** Whenever the activation function is *continuous, bounded and nonconstant*, then, for arbitrary compact subsets $X \subseteq R^k$, standard multilayer feedforward networks can approximate any continuous function on $X$ arbitrarily well with respect to uniform distance, provided that sufficiently many hidden units are available.

- <u>In words:</u> Given any continuous function f(x), if a 2-layer neural network has enough hidden units, then there is a choice of weights that allow it to closely approximate f(x).

Cybenko (1989) "Approximations by superpositions of sigmoidal functions"
Hornik (1991) "Approximation Capabilities of Multilayer Feedforward Networks"
Leshno and Schocken (1991) "Multilayer Feedforward Networks with Non-Polynomial Activation Functions Can Approximate Any Function"

# Universal Function Approximation Theorem*

Mathematics of Control,
Signals, and Systems
© 1989 Springer-Verlag New York Inc.

## Approximation by Superpositions of a Sigmoidal Function*

G. Cybenko†

**Abstract.** In this paper we demonstrate that finite linear combinations of compositions of a fixed, univariate function and a set of affine functionals can uniformly approximate any continuous function of $n$ real variables with support in the unit hypercube; only mild conditions are imposed on the univariate function. Our results settle an open question about representability in the class of single hidden layer neural networks. In particular, we show that arbitrary decision regions can be arbitrarily well approximated by continuous feedforward neural networks with only a single internal, hidden layer and any continuous sigmoidal nonlinearity. The paper discusses approximation properties of other possible types of nonlinearities that might be implemented by artificial neural networks.

**Key words.** Neural networks, Approximation, Completeness.

### 1. Introduction

A number of diverse application areas are concerned with the representation of general functions of an $n$-dimensional real variable, $x \in \mathbb{R}^n$, by finite linear combinations of the form

$$\sum_{j=1}^{N} \alpha_j \sigma(y_j^T x + \theta_j), \qquad (1)$$

where $y_j \in \mathbb{R}^n$ and $\alpha_j, \theta \in \mathbb{R}$ are fixed. ($y^T$ is the transpose of $y$ so that $y^T x$ is the inner product of $y$ and $x$.) Here the univariate function $\sigma$ depends heavily on the context of the application. Our major concern is with so-called sigmoidal $\sigma$'s:

$$\sigma(t) \to \begin{cases} 1 & \text{as} \quad t \to +\infty, \\ 0 & \text{as} \quad t \to -\infty. \end{cases}$$

Such functions arise naturally in neural network theory as the activation function of a neural node (or *unit* as is becoming the preferred term) [L1], [RHM]. The main result of this paper is a demonstration of the fact that sums of the form (1) are dense in the space of continuous functions on the unit cube if $\sigma$ is any continuous sigmoidal

*Date received: October 21, 1988. Date revised: February 17, 1989. This research was supported in part by NSF Grant DCR-8619103, ONR Contract N000-86-G-0202 and DOE Grant DE-FG02-85ER25001.
† Center for Supercomputing Research and Development and Department of Electrical and Computer Engineering, University of Illinois, Urbana, Illinois 61801, U.S.A.

303

---

*ORIGINAL CONTRIBUTION*

## Approximation Capabilities of Multilayer Feedforward Networks

Kurt Hornik

Technische Universität Wien, Vienna, Austria

(Received 30 January 1990; revised and accepted 25 October 1990)

**Abstract**—We show that standard multilayer feedforward networks with as few as a single hidden layer and arbitrary bounded and nonconstant activation function are universal approximators with respect to $L^p(\mu)$ performance criteria, for arbitrary finite input environment measures $\mu$, provided only that sufficiently many hidden units are available. If the activation function is continuous, bounded and nonconstant, then continuous mappings can be learned uniformly over compact input sets. We also give very general conditions ensuring that networks with sufficiently smooth activation functions are capable of arbitrarily accurate approximation to a function and its derivatives.

**Keywords**—Multilayer feedforward networks, Activation function, Universal approximation capabilities, Input environment measure, $L^p(\mu)$ approximation, Uniform approximation, Sobolev spaces, Smooth approximation.

### 1. INTRODUCTION

The approximation capabilities of neural network architectures have recently been investigated by many authors, including Carroll and Dickinson (1989), Cybenko (1989), Funahashi (1989), Gallant and White (1988), Hecht-Nielsen (1989), Hornik, Stinchcombe, and White (1989, 1990), Irie and Miyake (1988), Lapedes and Farber (1988), Stinchcombe and White (1989, 1990). (This list is by no means complete.)

If we think of the network architecture as a rule for computing values at $l$ output units given values at $k$ input units, hence implementing a class of mappings from $R^k$ to $R^l$, we can ask how well arbitrary mappings from $R^k$ to $R^l$ can be approximated by the network, in particular, if as many hidden units are required for internal representation and computation may be employed.

How to measure the accuracy of approximation depends on how we measure closeness between functions, which in turn varies significantly with the specific problem to be dealt with. In many applications, it is necessary to have the network perform *simultaneously* well on all input samples taken from some compact input set $X$ in $R^k$. In this case, closeness is

measured by the uniform distance between functions on $X$, that is,

$$\rho_{u,X}(f, g) = \sup_{x \in X} |f(x) - g(x)|.$$

In other applications, we think of the inputs as random variables and are interested in the *average performance* where the average is taken with respect to the input environment measure $\mu$, where $\mu(R^k) < \infty$. In this case, closeness is measured by the $L^p(\mu)$ distances

$$\rho_{p,\mu}(f, g) = \left[ \int_{R^k} |f(x) - g(x)|^p \, d\mu(x) \right]^{1/p},$$

$1 \le p < \infty$, the most popular choice being $p = 2$, corresponding to mean square error.

Of course, there are many more ways of measuring closeness of functions. In particular, in many applications, it is also necessary that the *derivatives* of the approximating function implemented by the network closely resemble those of the function to be approximated, up to some order. This issue was first taken up in Hornik et al. (1990), who discuss the sources of need of smooth functional approximation in more detail. Typical examples arise in robotics (learning of smooth movements) and signal processing (analysis of chaotic time series); for a recent application to problems of nonparametric inference in statistics and econometrics, see Gallant and White (1989).

All papers establishing certain approximation ca-

Requests for reprints should be sent to Kurt Hornik, Institut für Statistik und Wahrscheinlichkeitstheorie, Technische Universität Wien, Wiedner Hauptstraße 8-10/107, A-1040 Wien, Austria.

251

---

## MULTILAYER FEEDFORWARD NETWORKS WITH NON-POLYNOMIAL ACTIVATION FUNCTIONS CAN APPROXIMATE ANY FUNCTION

by

Moshe Leshno
Faculty of Management
Tel Aviv University
Tel Aviv, Israel 69978

and

Shimon Schocken
Leonard N. Stern School of Business
New York University
New York, NY 10003

September 1991

Center for Research on Information Systems
Information Systems Department
Leonard N. Stern School of Business
New York University

Working Paper Series

STERN IS-91-26

Appeared previously as *Working Paper No. 21/91* at The Israel Institute Of Business Research

---

Cybenko (1989) "Approximations by superpositions of sigmoidal functions"
Hornik (1991) "Approximation Capabilities of Multilayer Feedforward Networks"
Leshno and Schocken (1991) "Multilayer Feedforward Networks with Non-Polynomial Activation
Functions Can Approximate Any Function"

# Fun Neural Net Demo Site

- Demo-site:
  - http://playground.tensorflow.org/

# How about computing all the derivatives?

- Derivatives tables:

$$\frac{d}{dx}(a) = 0$$

$$\frac{d}{dx}(x) = 1$$

$$\frac{d}{dx}(au) = a\frac{du}{dx}$$

$$\frac{d}{dx}(u + v - w) = \frac{du}{dx} + \frac{dv}{dx} - \frac{dw}{dx}$$

$$\frac{d}{dx}(uv) = u\frac{dv}{dx} + v\frac{du}{dx}$$

$$\frac{d}{dx}\left(\frac{u}{v}\right) = \frac{1}{v}\frac{du}{dx} - \frac{u}{v^2}\frac{dv}{dx}$$

$$\frac{d}{dx}(u^n) = nu^{n-1}\frac{du}{dx}$$

$$\frac{d}{dx}(\sqrt{u}) = \frac{1}{2\sqrt{u}}\frac{du}{dx}$$

$$\frac{d}{dx}\left(\frac{1}{u}\right) = -\frac{1}{u^2}\frac{du}{dx}$$

$$\frac{d}{dx}\left(\frac{1}{u^n}\right) = -\frac{n}{u^{n+1}}\frac{du}{dx}$$

$$\frac{d}{dx}[f(u)] = \frac{d}{du}[f(u)]\frac{du}{dx}$$

$$\frac{d}{dx}[\ln u] = \frac{d}{dx}[\log_e u] = \frac{1}{u}\frac{du}{dx}$$

$$\frac{d}{dx}[\log_a u] = \log_a e\frac{1}{u}\frac{du}{dx}$$

$$\frac{d}{dx}e^u = e^u\frac{du}{dx}$$

$$\frac{d}{dx}a^u = a^u \ln a\frac{du}{dx}$$

$$\frac{d}{dx}(u^v) = vu^{v-1}\frac{du}{dx} + \ln u \; u^v\frac{dv}{dx}$$

$$\frac{d}{dx}\sin u = \cos u\frac{du}{dx}$$

$$\frac{d}{dx}\cos u = -\sin u\frac{du}{dx}$$

$$\frac{d}{dx}\tan u = \sec^2 u\frac{du}{dx}$$

$$\frac{d}{dx}\cot u = -\csc^2 u\frac{du}{dx}$$

$$\frac{d}{dx}\sec u = \sec u \tan u\frac{du}{dx}$$

$$\frac{d}{dx}\csc u = -\csc u \cot u\frac{du}{dx}$$

[source: http://hyperphysics.phy-astr.gsu.edu/hbase/Math/derfunc.html

# How about computing all the derivatives?

- But neural net f is never one of those?
  - No problem: CHAIN RULE:

If
$$f(x) = g(h(x))$$

Then
$$f'(x) = g'(h(x))h'(x)$$

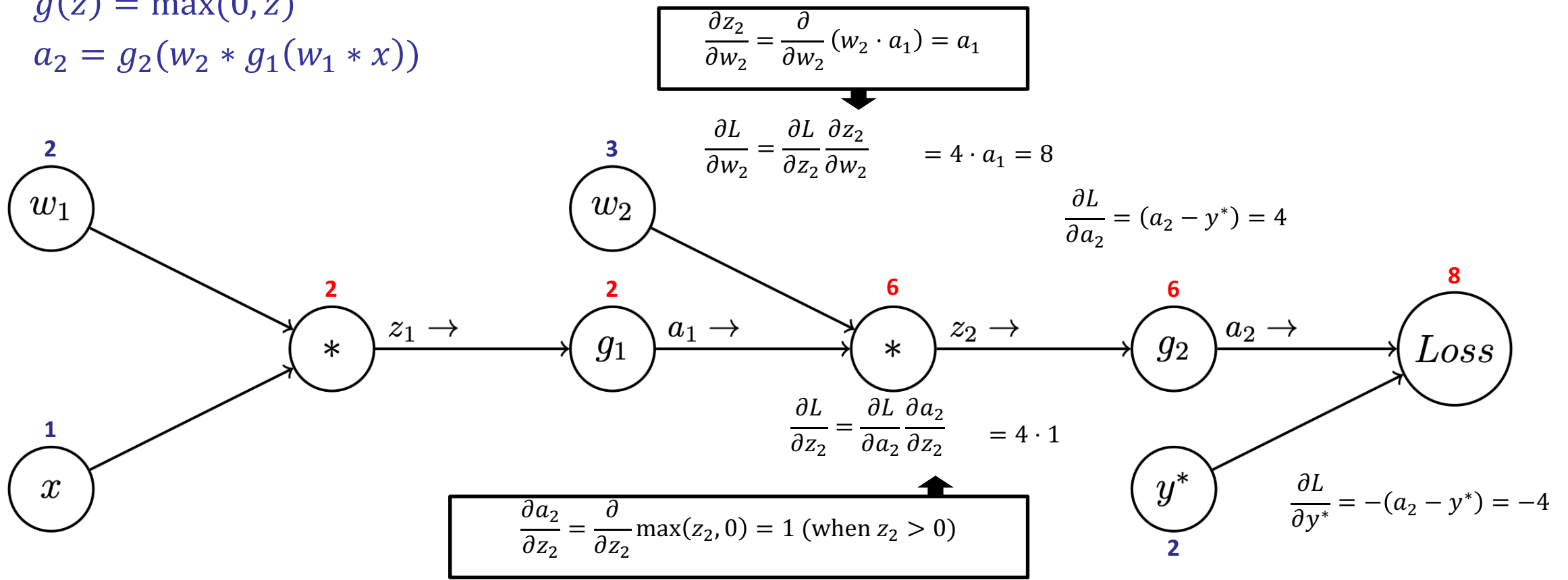**Derivatives can be computed by following well-defined procedures**

# Automatic Differentiation

- **Automatic differentiation software**
    - e.g., PyTorch, TensorFlow, Jax
    - Only need to program the function g(x,y,w)
    - Can automatically compute all derivatives w.r.t. all entries in w
    - This is typically done by caching info during forward computation pass of f, and then doing a backward pass = "backpropagation"
    - Autodiff / Backpropagation can often be done at computational cost comparable to the forward pass
- **Need to know this exists**
- **How this is done?  --  outside of scope of CS188**

# Example: Automatic Differentiation

- Build a computation graph and apply chain rule: $f(x) = g(h(x))$    $f'(x) = h'(x) \cdot g'(h(x))$

- Example: neural network with quadratic loss: $L(a_2, y^*) = \frac{1}{2}(a_2 - y^*)^2$ and ReLU activations $g(z) = \max(0, z)$

- $a_2 = g_2(w_2 * g_1(w_1 * x))$

$$\frac{\partial z_2}{\partial w_2} = \frac{\partial}{\partial w_2}(w_2 \cdot a_1) = a_1$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial z_2}\frac{\partial z_2}{\partial w_2} \qquad = 4 \cdot a_1 = 8$$

$$\frac{\partial L}{\partial a_2} = (a_2 - y^*) = 4$$



$$\frac{\partial L}{\partial z_2} = \frac{\partial L}{\partial a_2}\frac{\partial a_2}{\partial z_2} \qquad = 4 \cdot 1$$

$$\frac{\partial a_2}{\partial z_2} = \frac{\partial}{\partial z_2}\max(z_2, 0) = 1 \text{ (when } z_2 > 0)$$

$$\frac{\partial L}{\partial y^*} = -(a_2 - y^*) = -4$$

# Preventing Overfitting in Neural Networks

- Early stopping:



- Weight regularization: $\max\limits_{w} \sum_i \log P\left( y^{(i)} \mid x^{(i)}; w \right) - \frac{\lambda}{2} \sum_j w_j^2$

- Dropout:

# Dropout

"Damage" the network during training to increase redundancy

At each training step, with probability (1-p) set an activation to zero (i.e., drop it)

When making predictions, don't apply dropout, but multiply weights by p (rescaling)



(a) Standard Neural Net

(b) After applying dropout.

# Preventing Overfitting in Neural Networks

- Early stopping:



- Weight regularization: $\max\limits_{w} \sum_i \log P\left( y^{(i)} \mid x^{(i)}; w \right) - \frac{\lambda}{2} \sum_j w_j^2$

- Dropout:



(a) Standard Neural Net          (b) After applying dropout.

# Summary of Key Ideas

- Optimize probability of label given input $\quad \max\limits_{w} \quad ll(w) = \max\limits_{w} \sum\limits_{i} \log P(y^{(i)}|x^{(i)}; w)$

- Continuous optimization
  - Gradient ascent:
    - Compute steepest uphill direction = gradient (= just vector of partial derivatives)
    - Take step in the gradient direction
    - Repeat (until held-out data accuracy starts to drop = "early stopping")

- Deep neural nets
  - Last layer = still logistic regression
  - Now also many more layers before this last layer
    - = computing the features
    - → the features are learned rather than hand-designed
  - Universal function approximation theorem
    - `If` neural net is large enough
    - `Then` neural net can represent any continuous mapping from input to output with arbitrary accuracy
    - But remember: need to avoid overfitting / memorizing the training data → early stopping!
  - Automatic differentiation gives the derivatives efficiently

# Inductive Learning

# Inductive Learning (Science)

- Simplest form: learn a function from examples
  - A target function: $g$
  - Examples: input-output pairs $(x, g(x))$
  - E.g. $x$ is an email and $g(x)$ is spam / ham
  - E.g. $x$ is a house and $g(x)$ is its selling price

- Problem:
  - Given a hypothesis space $H$
  - Given a training set of examples $x_i$
  - Find a hypothesis $h(x)$ such that $h \sim g$

- Includes:
  - Classification (outputs = class labels)
  - Regression (outputs = real numbers)

- How do perceptron and naïve Bayes fit in?  ($H, h, g$, etc.)

$g$

$h$

$H$

# Inductive Learning

- Curve fitting (regression, function approximation):



- Consistency vs. simplicity
- Ockham's razor

# Consistency vs. Simplicity

- Fundamental tradeoff: bias vs. variance

- Usually algorithms prefer consistency by default (why?)

- Several ways to operationalize "simplicity"
  - Reduce the hypothesis space
    - Assume more: e.g. independence assumptions, as in naïve Bayes
    - Have fewer, better features / attributes: feature selection
    - Other structural limitations (decision lists vs trees)
  - Regularization
    - Smoothing: cautious use of small counts
    - Many other generalization parameters (pruning cutoffs today)
    - Hypothesis space stays big, but harder to get to the outskirts

# Decision Trees

# Reminder: Features

- Features, aka attributes
  - Sometimes: TYPE=French
  - Sometimes: $f_{\text{TYPE=French}}(x) = 1$

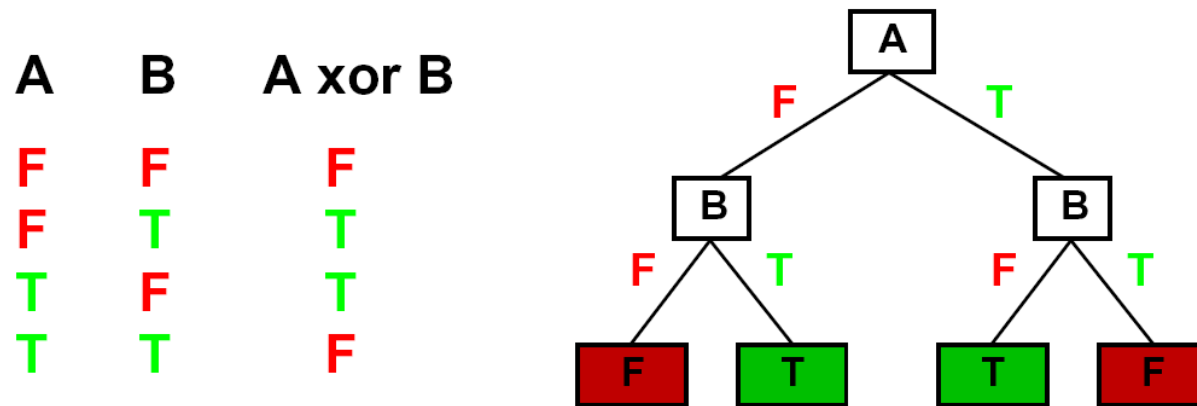| Example | Attributes | | | | | | | | | | Target |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | *Alt* | *Bar* | *Fri* | *Hun* | *Pat* | *Price* | *Rain* | *Res* | *Type* | *Est* | *WillWait* |
| $X_1$ | T | F | F | T | Some | $$$ | F | T | French | 0–10 | T |
| $X_2$ | T | F | F | T | Full | $ | F | F | Thai | 30–60 | F |
| $X_3$ | F | T | F | F | Some | $ | F | F | Burger | 0–10 | T |
| $X_4$ | T | F | T | T | Full | $ | F | F | Thai | 10–30 | T |
| $X_5$ | T | F | T | F | Full | $$$ | F | T | French | >60 | F |
| $X_6$ | F | T | F | T | Some | $$ | T | T | Italian | 0–10 | T |
| $X_7$ | F | T | F | F | None | $ | T | F | Burger | 0–10 | F |
| $X_8$ | F | F | F | T | Some | $$ | T | T | Thai | 0–10 | T |
| $X_9$ | F | T | T | F | Full | $ | T | F | Burger | >60 | F |
| $X_{10}$ | T | T | T | T | Full | $$$ | F | T | Italian | 10–30 | F |
| $X_{11}$ | F | F | F | F | None | $ | F | F | Thai | 0–10 | F |
| $X_{12}$ | T | T | T | T | Full | $ | F | F | Burger | 30–60 | T |

# Decision Trees

- Compact representation of a function:
  - Truth table
  - Conditional probability table
  - Regression values

- True function
  - Realizable: in $H$

# Expressiveness of DTs

- Can express any function of the features

| A | B | A xor B |
|---|---|---------|
| F | F | F |
| F | T | T |
| T | F | T |
| T | T | F |



$$P(C|A,B)$$

- However, we hope for compact trees

# Comparison: Perceptrons

- What is the expressiveness of a perceptron over these features?

| Example | Attributes | | | | | | | | | | Target |
|---------|-----|-----|-----|-----|------|-------|------|-----|--------|-------|----------|
|  | Alt | Bar | Fri | Hun | Pat | Price | Rain | Res | Type | Est | WillWait |
| $X_1$ | T | F | F | T | Some | $$$ | F | T | French | 0–10 | T |
| $X_2$ | T | F | F | T | Full | $ | F | F | Thai | 30–60 | F |

- For a perceptron, a feature's contribution is either positive or negative
  - If you want one feature's effect to depend on another, you have to add a new conjunction feature
  - E.g. adding "PATRONS=full $\wedge$ WAIT = 60" allows a perceptron to model the interaction between the two atomic features

- DTs automatically conjoin features / attributes
  - Features can have different effects in different branches of the tree!

- Difference between modeling relative evidence weighting (NB) and complex evidence interaction (DTs)
  - Though if the interactions are too complex, may not find the DT greedily

# Hypothesis Spaces

- **How many distinct decision trees with n Boolean attributes?**
  = number of Boolean functions over n attributes
  = number of distinct truth tables with $2^n$ rows
  = 2^($2^n$)
  - E.g., with 6 Boolean attributes, there are
    18,446,744,073,709,551,616 trees

- **How many trees of depth 1 (decision stumps)?**
  = number of Boolean functions over 1 attribute
  = number of truth tables with 2 rows, times n
  = 4n
  - E.g. with 6 Boolean attributes, there are 24 decision stumps

- **More expressive hypothesis space:**
  - Increases chance that target function can be expressed (good)
  - Increases number of hypotheses consistent with training set (bad, why?)
  - Means we can get better predictions (lower bias)
  - But we may get worse predictions (higher variance)

# Decision Tree Learning

- Aim: find a small tree consistent with the training examples
- Idea: (recursively) choose "most significant" attribute as root of (sub)tree

---

function DTL(*examples, attributes, default*) returns a decision tree

    if *examples* is empty then return *default*
    else if all *examples* have the same classification then return the classification
    else if *attributes* is empty then return MODE(*examples*)
    else
        $best \leftarrow$ CHOOSE-ATTRIBUTE(*attributes, examples*)
        $tree \leftarrow$ a new decision tree with root test *best*
        for each value $v_i$ of *best* do
            $examples_i \leftarrow \{$elements of *examples* with $best = v_i\}$
            $subtree \leftarrow$ DTL($examples_i, attributes - best$, MODE(*examples*))
            add a branch to *tree* with label $v_i$ and subtree *subtree*
    return *tree*

# Choosing an Attribute

- Idea: a good attribute splits the examples into subsets that are (ideally) "all positive" or "all negative"



- So: we need a measure of how "good" a split is, even if the results aren't perfectly separated out

# Entropy and Information

- Information answers questions
  - The more uncertain about the answer initially, the more information in the answer
  - Scale: bits
    - Answer to Boolean question with prior <1/2, 1/2>?
    - Answer to 4-way question with prior <1/4, 1/4, 1/4, 1/4>?
    - Answer to 4-way question with prior <0, 0, 0, 1>?
    - Answer to 3-way question with prior <1/2, 1/4, 1/4>?

- A probability p is typical of:
  - A uniform distribution of size 1/p
  - A code of length log 1/p

# Entropy

- General answer: if prior is $\langle p_1, \ldots, p_n \rangle$:
  - Information is the expected code length

$$H(\langle p_1, \ldots, p_n \rangle) = E_p \log_2 1/p_i$$

$$= \sum_{i=1}^{n} -p_i \log_2 p_i$$

- Also called the entropy of the distribution
  - More uniform = higher entropy
  - More values = higher entropy
  - More peaked = lower entropy
  - Rare values almost "don't count"

1 bit

0 bits

0.5 bit

# Information Gain

- Back to decision trees!
- For each split, compare entropy before and after
  - Difference is the information gain
  - Problem: there's more than one distribution after split!



  - Solution: use expected entropy, weighted by the number of examples

# Next Step: Recurse

- Now we need to keep growing the tree!
- Two branches are done (why?)
- What to do under "full"?
  - See what examples are there…



| Example | Attributes | | | | | | | | | | Target |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | Alt | Bar | Fri | Hun | Pat | Price | Rain | Res | Type | Est | WillWait |
| $X_1$ | T | F | F | T | Some | $$$ | F | T | French | 0–10 | T |
| $X_2$ | T | F | F | T | Full | $ | F | F | Thai | 30–60 | F |
| $X_3$ | F | T | F | F | Some | $ | F | F | Burger | 0–10 | T |
| $X_4$ | T | F | T | T | Full | $ | F | F | Thai | 10–30 | T |
| $X_5$ | T | F | T | F | Full | $$$ | F | T | French | >60 | F |
| $X_6$ | F | T | F | T | Some | $$ | T | T | Italian | 0–10 | T |
| $X_7$ | F | T | F | F | None | $ | T | F | Burger | 0–10 | F |
| $X_8$ | F | F | F | T | Some | $$ | T | T | Thai | 0–10 | T |
| $X_9$ | F | T | T | F | Full | $ | T | F | Burger | >60 | F |
| $X_{10}$ | T | T | T | T | Full | $$$ | F | T | Italian | 10–30 | F |
| $X_{11}$ | F | F | F | F | None | $ | F | F | Thai | 0–10 | F |
| $X_{12}$ | T | T | T | T | Full | $ | F | F | Burger | 30–60 | T |

# Example: Learned Tree

- Decision tree learned from these 12 examples:



- Substantially simpler than "true" tree
  - A more complex hypothesis isn't justified by data
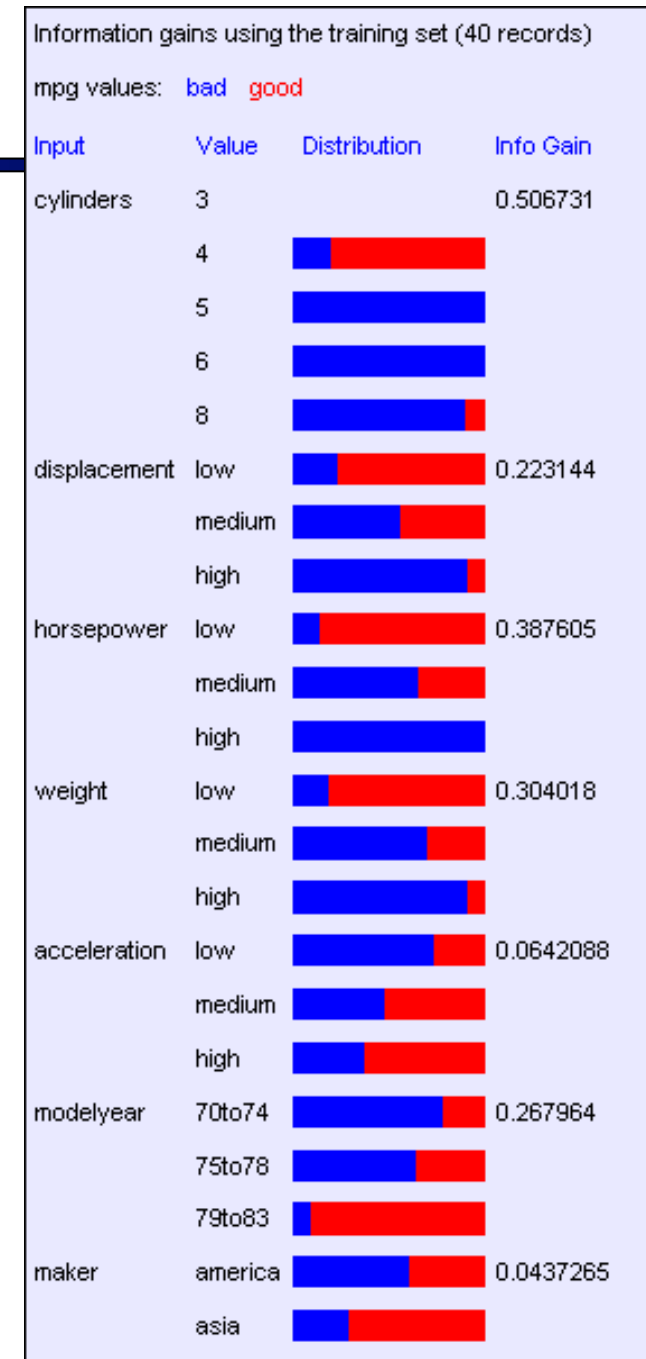- Also: it's reasonable, but wrong

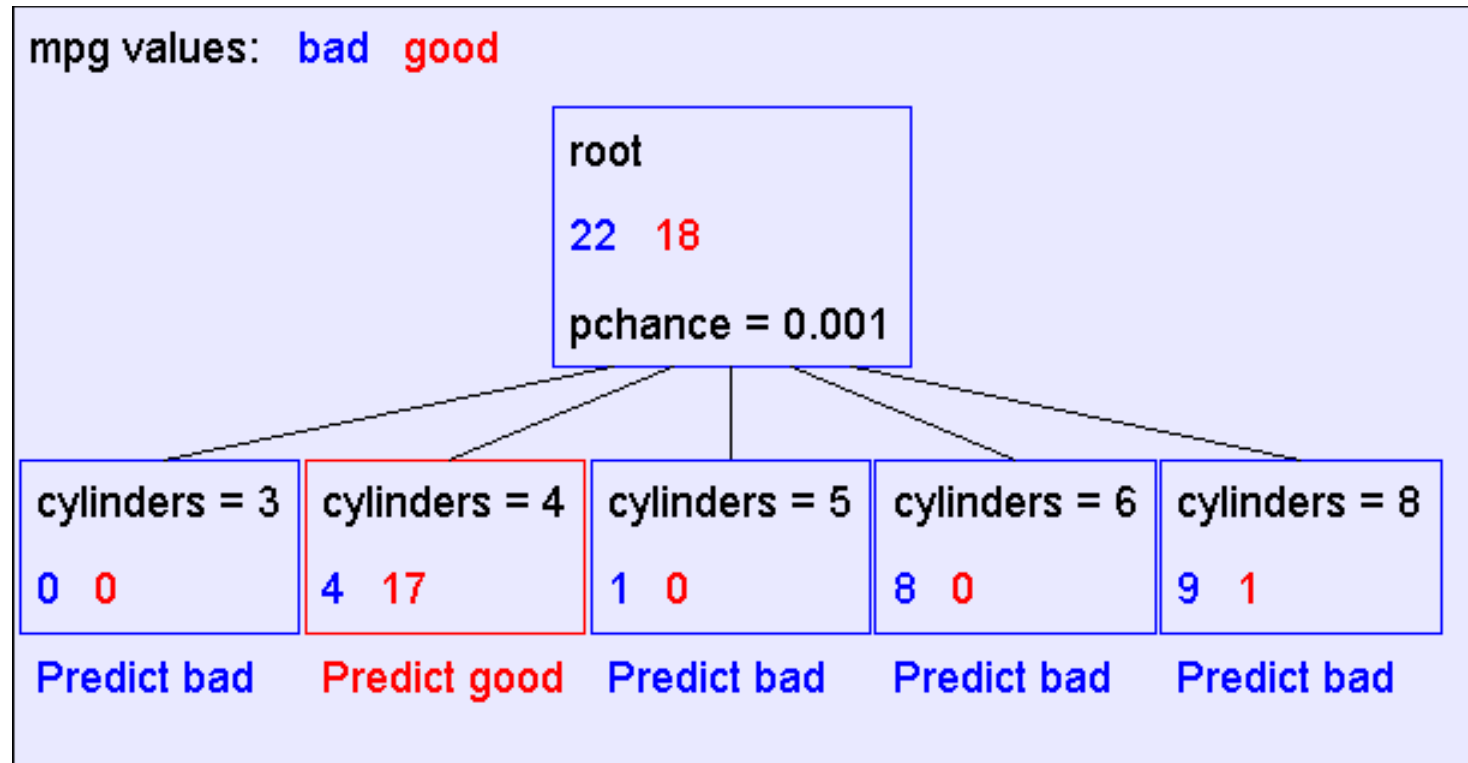# Example: Miles Per Gallon

**40 Examples**

| mpg | cylinders | displacement | horsepower | weight | acceleration | modelyear | maker |
|-----|-----------|--------------|------------|--------|--------------|-----------|-------|
| | | | | | | | |
| good | 4 | low | low | low | high | 75to78 | asia |
| bad | 6 | medium | medium | medium | medium | 70to74 | america |
| bad | 4 | medium | medium | medium | low | 75to78 | europe |
| bad | 8 | high | high | high | low | 70to74 | america |
| bad | 6 | medium | medium | medium | medium | 70to74 | america |
| bad | 4 | low | medium | low | medium | 70to74 | asia |
| bad | 4 | low | medium | low | low | 70to74 | asia |
| bad | 8 | high | high | high | low | 75to78 | america |
| : | : | : | : | : | : | : | : |
| : | : | : | : | : | : | : | : |
| : | : | : | : | : | : | : | : |
| bad | 8 | high | high | high | low | 70to74 | america |
| good | 8 | high | medium | high | high | 79to83 | america |
| bad | 8 | high | high | high | low | 75to78 | america |
| good | 4 | low | low | low | low | 79to83 | america |
| bad | 6 | medium | medium | medium | high | 75to78 | america |
| good | 4 | medium | low | low | low | 79to83 | america |
| good | 4 | low | low | medium | high | 79to83 | america |
| bad | 8 | high | high | high | low | 70to74 | america |
| good | 4 | low | medium | low | medium | 75to78 | europe |
| bad | 5 | medium | medium | medium | medium | 75to78 | europe |

# Find the First Split
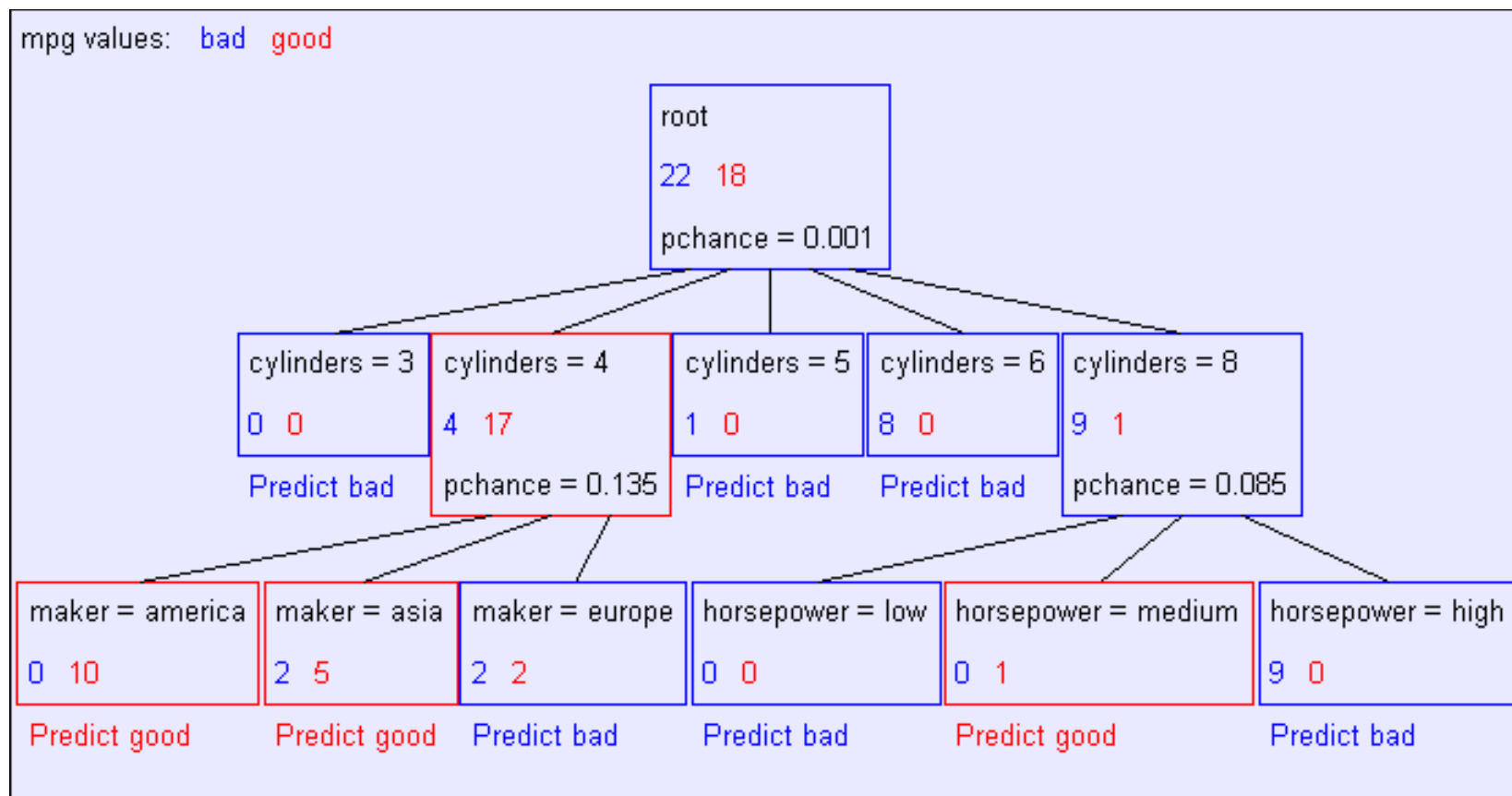
- Look at information gain for each attribute

- Note that each attribute is correlated with the target!
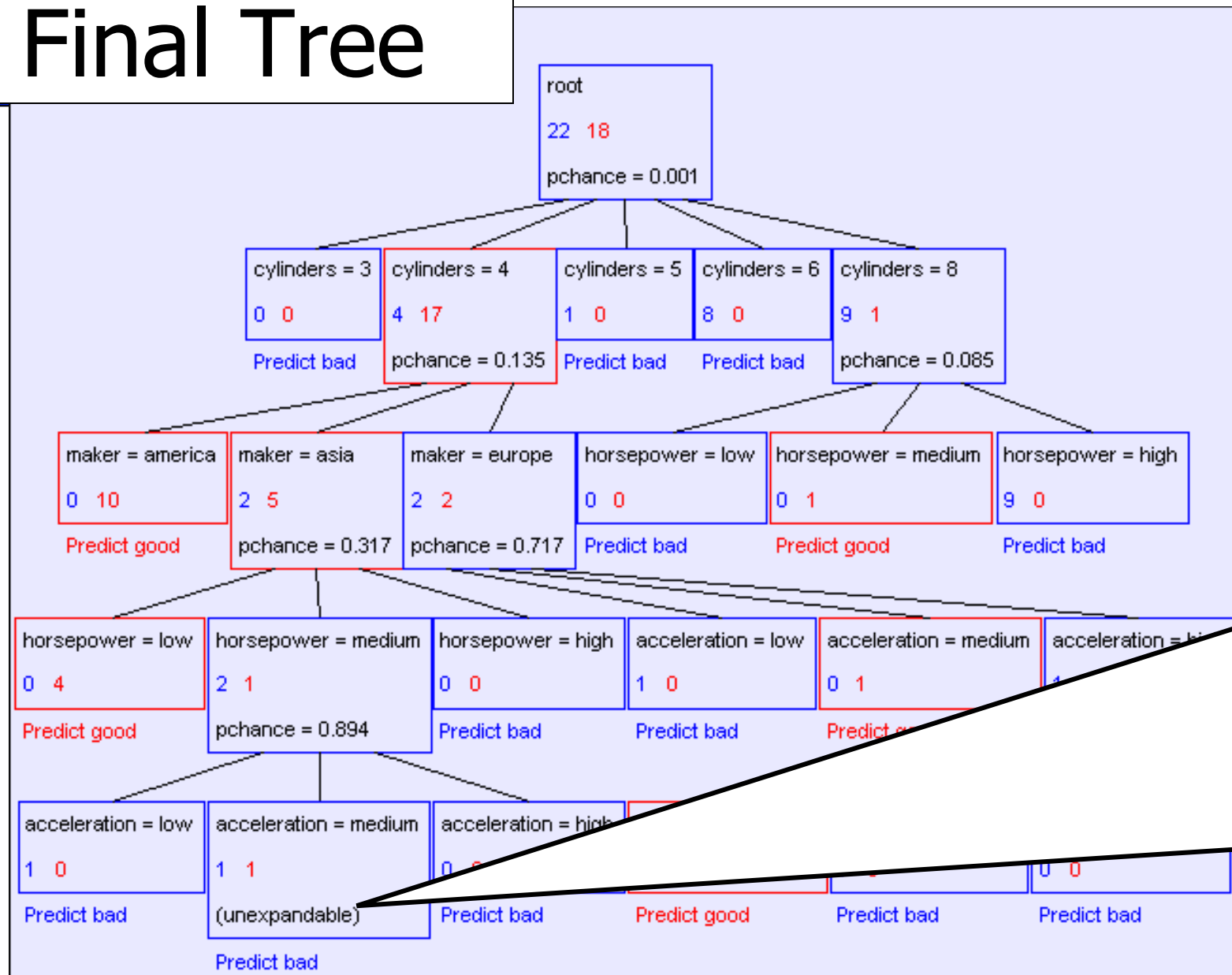
- What do we split on?



Information gains using the training set (40 records)

mpg values: bad good

| Input | Value | Distribution | Info Gain |
|---|---|---|---|
| cylinders | 3 | | 0.506731 |
| | 4 | | |
| | 5 | | |
| | 6 | | |
| | 8 | | |
| displacement | low | | 0.223144 |
| | medium | | |
| | high | | |
| horsepower | low | | 0.387605 |
| | medium | | |
| | high | | |
| weight | low | | 0.304018 |
| | medium | | |
| | high | | |
| acceleration | low | | 0.0642088 |
| | medium | | |
| | high | | |
| modelyear | 70to74 | | 0.267964 |
| | 75to78 | | |
| | 79to83 | | |
| maker | america | | 0.0437265 |
| | asia | | |

# Result: Decision Stump

# Second Level

# Final Tree

root

22 18

pchance = 0.001

| cylinders = 3 | cylinders = 4 | cylinders = 5 | cylinders = 6 | cylinders = 8 |
|---|---|---|---|---|
| 0  0 | 4  17 | 1  0 | 8  0 | 9  1 |
| Predict bad | pchance = 0.135 | Predict bad | Predict bad | pchance = 0.085 |

| maker = america | maker = asia | maker = europe | horsepower = low | horsepower = medium | horsepower = high |
|---|---|---|---|---|---|
| 0  10 | 2  5 | 2  2 | 0  0 | 0  1 | 9  0 |
| Predict good | pchance = 0.317 | pchance = 0.717 | Predict bad | Predict good | Predict bad |

| horsepower = low | horsepower = medium | horsepower = high | acceleration = low | acceleration = medium | acceleration = high |
|---|---|---|---|---|---|
| 0  4 | 2  1 | 0  0 | 1  0 | 0  1 | 1 |
| Predict good | pchance = 0.894 | Predict bad | Predict bad | Predict good | |

| acceleration = low | acceleration = medium | acceleration = high | | |
|---|---|---|---|---|
| 1  0 | 1  1 | 0 | | 0  0 |
| Predict bad | (unexpandable) | Predict bad | Predict good | Predict bad |
| | Predict bad | | | |

# Reminder: Overfitting

- Overfitting:
  - When you stop modeling the patterns in the training data (which generalize)
  - And start modeling the noise (which doesn't)

- We had this before:
  - Naïve Bayes: needed to smooth
  - Perceptron: early stopping

mpg values:  bad   good

root

22  18

pchance = 0.001

| | Num Errors | Set Size | Percent Wrong |
|---|---|---|---|
| Training Set | 1 | 40 | 2.50 |
| Test Set | 74 | 352 | 21.02 |

...epower = high

...ict bad

horsepower = low    horsepower = medium    horsepower = high    acceleration = low    acceleration = medium    acceleration = high

= 0.717

= 79to83

The test set error is much worse than the training set error...

...why?

Predict bad    (unexpandable)    Predict bad    Predict good    Predict bad    Predict bad

Predict bad

# Significance of a Split

- Starting with:
    - Three cars with 4 cylinders, from Asia, with medium HP
    - 2 bad MPG
    - 1 good MPG

- What do we expect from a three-way split?
    - Maybe each example in its own subset?
    - Maybe just what we saw in the last slide?

- Probably shouldn't split if the counts are so small they could be due to chance

- A chi-squared test can tell us how likely it is that deviations from a perfect split are due to chance*

- Each split will have a significance value, $p_{CHANCE}$

# Keeping it General

- Pruning:
  - Build the full decision tree
  - Begin at the bottom of the tree
  - Delete splits in which

    $p_{CHANCE} > MaxP_{CHANCE}$

  - Continue working upward until there are no more prunable nodes
  - Note: some chance nodes may not get pruned because they were "redeemed" later

y = a XOR b

| a | b | y |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

# Pruning example

- With MaxP$_{CHANGE}$ = 0.1:



mpg values:  bad  good

root
22  18
pchance = 0.001

| cylinders = 3 | cylinders = 4 | cylinders = 5 | cylinders = 6 | cylinders = 8 |
|---|---|---|---|---|
| 0  0 | 4  17 | 1  0 | 8  0 | 9  1 |
| Predict bad | Predict good | Predict bad | Predict bad | Predict bad |

Note the improved test set accuracy compared with the unpruned tree

|  | Num Errors | Set Size | Percent Wrong |
|---|---|---|---|
| Training Set | 5 | 40 | 12.50 |
| Test Set | 56 | 352 | 15.91 |

# Regularization

- MaxP$_{CHANCE}$ is a regularization parameter
- Generally, set it using held-out data (as usual)

# Two Ways of Controlling Overfitting

- Limit the hypothesis space
  - E.g. limit the max depth of trees
  - Easier to analyze

- Regularize the hypothesis selection
  - E.g. chance cutoff
  - Disprefer most of the hypotheses unless data is clear
  - Usually done in practice

# Summary of Key Ideas

- Optimize probability of label given input    $\max\limits_{w} \quad ll(w) = \max\limits_{w} \sum\limits_{i} \log P(y^{(i)}|x^{(i)}; w)$

- Continuous optimization
  - Gradient ascent:
    - Compute steepest uphill direction = gradient (= just vector of partial derivatives)
    - Take step in the gradient direction
    - Repeat (until held-out data accuracy starts to drop = "early stopping")

- Deep neural nets
  - Last layer = still logistic regression
  - Now also many more layers before this last layer
    - = computing the features
    - the features are learned rather than hand-designed
  - Universal function approximation theorem
    - `If`      neural net is large enough
    - `Then`  neural net can represent any continuous mapping from input to output with arbitrary accuracy
    - But remember: need to avoid overfitting / memorizing the training data    early stopping!
  - Automatic differentiation gives the derivatives efficiently (how? = outside of scope of 188)

# How well does it work?

# Computer Vision

# Manual Feature Design

# Features and Generalization



[HoG: Dalal and Triggs, 2005]

# Features and Generalization



Image



HoG

# Performance



*graph credit Matt Zeiler, Clarifai*

# Performance



*graph credit Matt Zeiler, Clarifai*

# Performance

## ImageNet Error Rate 2010-2014

● Traditional CV    ● Deep Learning

Error Rate

79%
60%
40%
20%
7%

AlexNet

2010    2011    2012    2013    2014

*graph credit Matt Zeiler, Clarifai*

# Performance



*graph credit Matt Zeiler, Clarifai*

# Performance



graph credit Matt
Zeiler, Clarifai

# MS COCO Image Captioning Challenge



"man in black shirt is playing guitar."

"construction worker in orange safety vest is working on road."

"two young girls are playing with lego toy."

"boy is doing backflip on wakeboard."

"girl in pink dress is jumping in air."

"black and white dog jumps over bar."

"young girl in pink shirt is swinging on swing."

"man in blue wetsuit is surfing on wave."

Karpathy & Fei-Fei, 2015; Donahue et al., 2015; Xu et al, 2015; many more

# Visual QA Challenge

Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C. Lawrence Zitnick, Devi Parikh

# Speech Recognition



graph credit Matt Zeiler, Clarifai

# Machine Translation

# Next: More Neural Net Applications!