#### CS 188: Artificial Intelligence

#### Neural Nets (wrap-up) and Decision Trees



Instructors: John Canny and Oliver Grillmeyer --- University of California, Berkeley

[These slides were created by Dan Klein and Pieter Abbeel for CS188 Intro to AI at UC Berkeley. All CS188 materials are available at http://ai.berkeley.edu.]

#### Announcements

- Project 5 (last project)
  - Due Friday 4/25 at 11:59pm
- HW9
  - Due Wednesday 4/16 at 11:59pm
- HW10 (last homework)
  - Due Wednesday 4/23 at 11:59pm
- Final Exam
  - Thursday 5/15 from 3:00-6:00pm
  - See Exam Logistics on CS 188 website

## Today

- Neural Nets -- wrap
- Enhanced Training
- Formalizing Learning
  - Consistency
  - Simplicity
- Decision Trees
  - Expressiveness
  - Information Gain
  - Overfitting

#### Refresh: Deep Neural Network



#### **Computer Vision: Object Detection**



#### Traditional CV: Features and Generalization





Image



#### Performance

## ImageNet Error Rate 2010-2014



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



#### **Machine Translation**



## YOLO object detection



#### YOLO v3 network Architecture

https://towardsdatascience.com/yolo-v3-object-detection-53fb7d3bfe6b



## **Transfer Learning**

- Problem: how do we efficiently build machine learning models
- Data Labeling is a very time consuming operation that requires human input
- Can we leverage off of existing, similar models?
- Transfer Learning entails using the weights of a similar network as a starting point in training a new model
- Domain Adaptation is a simple form of Transfer Learning in which an existing model is further trained using a new smaller data set
- Transfer Learning can involve freezing the feature detection part of the network to just learn to discriminate and classify

## Formalizing Learning: 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<sub>i</sub>
- 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.)



## 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**



#### Features

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

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

#### **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



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
<b></b>	Alt	Bar	Fri	Hun	Pat	Price	Rain	Res	Type	Est	WillWait
$X_1$	T	F	F	Т	Some	\$\$\$	F	Т	French	0–10	Т
$X_2$	T	F	F	Т	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 
     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<sup>n</sup> rows
- = 2^(2<sup>n</sup>)
- 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  $\langle p_1, ..., p_n \rangle$ :
  - Information is the expected code length

$$H(\langle p_1, \dots, p_n \rangle) = E_p \log_2 1/p_i$$
$$= \sum_{i=1}^n -p_i \log_2 p_i$$

- 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"



## **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!





## 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										
<b>L</b> iterinpie	Alt	Bar	Fri	Hun	Pat	Price	Rain	Res	Type	Est	WillWait	
$X_1$	T	F	F	Т	Some	\$\$\$	F	T	French	0–10	Т	
$X_2$	T	F	F	Т	Full	\$	F	F	Thai	30–60	F	
$X_3$	F	T	F	F	Some	\$	F	F	Burger	0–10	Т	
$X_4$	T	F	Т	Т	Full	\$	F	F	Thai	10–30	Т	
$X_5$	T	F	Т	F	Full	\$\$\$	F	Т	French	>60	F	
$X_6$	F	T	F	Т	Some	\$\$	T	T	Italian	0–10	Т	
$X_7$	F	T	F	F	None	\$	T	F	Burger	0–10	F	
$X_8$	F	F	F	Т	Some	\$\$	Т	Т	Thai	0–10	Т	
$X_9$	F	T	Т	F	Full	\$	Т	F	Burger	>60	F	
$X_{10}$	T	T	Т	Т	Full	\$\$\$	F	Т	Italian	10–30	F	
<i>X</i> <sub>11</sub>	F	F	F	F	None	\$	F	F	Thai	0–10	F	
$X_{12}$	T	T	Т	Т	Full	\$	F	F	Burger	30–60	Т	

#### 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

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

40 Examples

## 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?



#### **Result: Decision Stump**



#### Next Lecture: Large Language Models & Transformers

- Wrap up Decision Trees
  - Pruning
- Large Langauge Models
  - Transformers