# Lecture 6

prəˌnʌnsiˈeɪʃən ˈmɑdəlɪŋ

Michael Picheny, Bhuvana Ramabhadran, Stanley F. Chen,
Markus Nussbaum-Thom

Watson Group
IBM T.J. Watson Research Center
Yorktown Heights, New York, USA
{picheny,bhuvana,stanchen,nussbaum}@us.ibm.com

24 February 2016 and Mar 2 2016

# Administrivia

- Lab 2
  - Not graded yet; handed back next lecture.
- Lab 3
  - Due nine days from now (Friday, Mar. 11) at 6pm.

# Lab Grading

- How things work:
  - Overall scale: -1 to +1 (-2 if don't hand in).
  - Programming part: max score +0.5.
  - Short answers: default score 0, with small bonus/pens.
  - +0.5 bonus: if total pens of at most 0.2 (think $\sqrt{+}$).
- Lab 1: +0.5 bonus wasn't applied in original grading.
  - If your score changed, should have recv'd E-mail.
  - Contact Stan if you still have questions.

# Feedback

- Clear (4), mostly clear (3), unclear (1).
- Pace: fast (2), OK (1).
- Muddiest: pronunciation modeling (1), Laplace smoothing (1).
- Comments (2+ votes)
    - Handing out grades distracting, inefficient (2).

# Review to date

- Learned about features (MFCCs, etc.)
- Learned about Gaussian Mixture Models
- Learned about HMMs and basic operations (finding best path, training models)
- Learned about basic Language modeling.

# Where Are We?
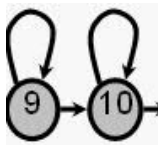
# Where Are We?

# In the beginning...

- ... . was the whole word model.
- For each word in the vocabulary, decide on an HMM structure.
- Often the number of states in the model is chosen to be proportional to the number of phonemes in the word.
- Train the HMM parameters for a given word using examples of that word in the training data.
- Good domain for this approach: digits.
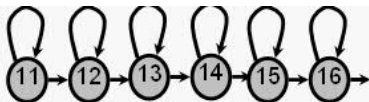
# Example topologies: Digits

- Vocabulary consists of ("zero", "oh", "one", "two", "three", "four", "five", "six", "seven", "eight", "nine").
- Assume we assign two states per phoneme.
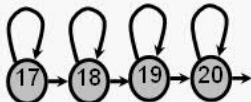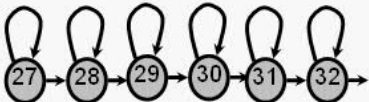- Models look like:
- "zero".



- "oh".

"one"  W AA N

"two"  T UW

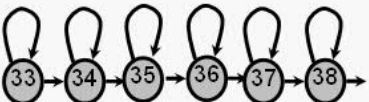"three"  TH R IY

"four"  F AO R

"five"  F AY V

# Whole-word model limitations

- The whole-word model suffers from two main problems.
    - Cannot model unseen words. In fact, we need several samples of each word to train the models properly.
    - Cannot share data among models – data sparseness problem.
    - The number of parameters in the system is proportional to the vocabulary size.
- Thus, whole-word models are best on small vocabulary tasks with lots of data per word.
    - n.b. as the amount of public speech data continues to increase this wisdom may be thrown into question.

# Where Are We?

# Subword Units

- To reduce the number of parameters, we can compose word models from sub-word units.
- These units can be shared among words. Examples include

| Units | Approximate number |
|-------|--------------------|
| Phones | 50. |
| Diphones | 2000. |
| Syllables | 5,000. |

- Each unit is small in terms of amount of speech modeled.
- The number of parameters is proportional to the number of units (not the number of words in the vocabulary as in whole-word models.).

# Phonetic Models

- We represent each word as a sequence of phonemes. This representation is the "baseform" for the word.

  BANDS   ->   B AE N D Z

- Some words need more than one baseform.

  THE   ->   DH UH
        ->   DH IY

# Baseform Dictionary

- To determine the pronunciation of each word, we look it up in a dictionary.
- Each word may have several possible pronunciations.
- Every word in our training script and test vocabulary must be in the dictionary.
- The dictionary is generally written by hand.
- Prone to errors and inconsistencies.

| | | |
|---|---|---|
| 2nd looks wrong | acapulco | \| AE K AX P AH L K OW |
| AA K .. is missing | acapulco | \| AE K AX P UH K OW |
| | accelerator | \| AX K S EH L AX R EY DX ER |
| Shouldn't the choices | accelerator | \| IX K S EH L AX R EY DX ER |
| For the 1st phoneme | acceleration | \| AX K S EH L AX R EY SH IX N |
| Same for these 2 words | acceleration | \| AE K S EH L AX R EY SH IX N |
| | accent | \| AE K S EH N T |
| | accept | \| AX K S EH P T |
| | acceptable | \| AX K S EH P T AX B AX L |
| Don't these words | access | \| AE K S EH S |
| All start the same? | accessory | \| AX K S EH S AX R IY |
| | accessory | \| EH K S EH S AX R IY |

# Phonetic Models, cont'd

- We can allow for a wide variety of phonological variation by representing baseforms as graphs.

acapulco   AE K AX P AH L K OW
acapulco   AA K AX P UH K OW

# Phonetic Models, cont'd

- Now, construct a Markov model for each phone.
- Examples:

# Embedding

- Replace each phone by its Markov model to get a model for the entire word

acapulco    AE K AX P AH L K OW
acapulco    AA K AX P UH K OW

# Reducing Parameters by Tying

- Consider the three-state model.



- Note that.
  - $t_1$ and $t_2$ correspond to the beginning of the phone.
  - $t_3$ and $t_4$ correspond to the middle of the phone.
  - $t_5$ and $t_6$ correspond to the end of the phone.
- If we force the output distributions for each member of those pairs to be the same, then the training data requirements are reduced.

# Tying

- A set of arcs in a Markov model are tied to one another if they are constrained to have identical output distributions.
- Similarly, states are tied if they have identical transition probabilities.
- Tying can be explicit or implicit.

# Implicit Tying

- Occurs when we build up models for larger units from models of smaller units.
- Example: when word models are made from phone models.
- First, consider an example without any tying.
  - Let the vocabulary consist of digits 0,1,2,... 9.
- We can make a separate model for each word.
- To estimate parameters for each word model, we need several samples for each word.
- Samples of "0" affect only parameters for the "0" model.

# Implicit Tying, cont'd

- Now consider phone-based models for this vocabulary.



```
0       Z  IY (R) OW
1       W  AA  N
2       T  UW
3        TH (R) IY
4       F  AO (R)
5       F  AY  V
6       S  IH  K  S
7       S  EH  V  AX  N
8        EY  T
9        N  AY  N
```

- Training samples of "0" will also affect models for "3" and "4".
- Useful in large vocabulary systems where the number of words is much greater than the number of phones.

# Explicit Tying

- Example:



- 6 non-null arcs, but only 3 different output distributions because of tying.
- Number of model parameters is reduced.
- Tying saves storage because only one copy of each distribution is saved.
- Fewer parameters mean less training data needed.

# Where Are We?

# Variations in realizations of phonemes

- The broad units, phonemes, have variants known as allophones
    - Example: *p* and *p$^h$* (un-aspirated and aspirated *p*).
    - Exercise: Put your hand in front of your mouth and pronounce *spin* and then *pin* Note that the *p* in *pin* has a puff of air,. while the *p* in *spin* does not.

# Variations in realizations of phonemes

- Articulators have inertia, thus the pronunciation of a phoneme is influenced by surrounding phonemes. This is known as <span style="color:red">co-articulation</span>
  - Example: Consider k in different contexts.
    - In *keep* the whole body of the tongue has to be pulled up to make the vowel.
    - Closure of the k moves forward compared to *coop*

keep

coop

# Phoneme Targets

Phonemes have idealized articulator target positions that may or may not be reached in a particular utterance.

- Speaking rate
- Clarity of articulation



Inherent variability of Speech biggest challenge

How do we model all this variation?

# Triphone models

- Model each phoneme in the context of its left and right neighbor.
- E.g. $_K IY_P$ is a model for *IY* when *K* is its left context phoneme and *P* is its right context phoneme.
  - "keep" → K IY P → $_{wb}K_{IY}$ $_K IY_P$ $_{IY}P_{wb}$
- If we have 50 phonemes in a language, we could have as many as $50^3$ triphones to model.
- Not all of these occur, or only occur a few times. Why is this bad?
- Suggestion: Combine similar triphones together
  - For example, map $_K IY_P$ and $_K IY_F$ to common model

# "Bottom-up" (Agglomerative) Clustering

- Start with each item in a cluster by itself.
- Find "closest" pair of items.
- Merge them into a single cluster.
- Iterate.

# Triphone Clustering

- Helps with data sparsity issue
- BUT still have an issue with unseen data
- To model unseen events, we can "back-off" to lower order models such as bi-phones and uni-phones. But this is still sort of ugly.
- So instead, we use **Decision Trees** to deal with the sparse/unknown data problem.

# Where Are We?

# Decision Trees



**TREES**
by Joyce Kilmer

I think that I shall never see
A poem lovely as a tree.

A tree whose hungry mouth is prest
Against the earth's sweet flowing breast;

A tree that looks at God all day,
And lifts her leafy arms to pray;

A tree that may in summer wear
A nest of robins in her hair;

Upon whose bosom snow has lain;
Who intimately lives with rain.

Poems are made by fools like me,
But only God can make a tree.

Features (Questions)

Appendicitis 100
Flu 100
Strep 100
Heart Attack 100
Cold 100
None 100

Pain?

Abdomen    Throat

Chest

None

Appendicitis 90
Flu 5
Strep 0
Heart Attack 5
Cold 10
None 5

Fever?

Yes

No

Appendicitis 0
Flu 0
Strep 0
Heart Attack 90
Cold 0
None 5

Cough?

Yes

No

Appendicitis 0
Flu 50
Strep 0
Heart Attack 0
Cold 20
None 5

Appendicitis 0
Flu 0
Strep 90
Heart Attack 0
Cold 10
None 5

Fever?

Yes

No

Appendicitis 10
Flu 5
Strep 10
Heart Attack 5
Cold 10
None 70

Appendicitis 0
Flu 40
Strep 0
Heart Attack 0
Cold 10
None 5

Appendicitis 0
Flu 0
Strep 0
Heart Attack 0
Cold 40
None 5

Outputs

# Types of Features

- Nominal or categorical: Finite set without any natural ordering (e.g., occupation, marital status, race).
- Ordinal: Ordered, but absolute differences between values is unknown (e.g., preference scale, severity of an injury).
- Numerical: Domain is numerically ordered (e.g., age, income).

# Types of Outputs

- Categorical: Output is one of *N* classes
  - Diagnosis: Predict disease from symptoms
  - Language Modeling: Predict next word from previous words in the sentence
  - Spelling to sound rules: Predict phone from spelling
- Continuous: Output is a continuous vector
  - Allophonic variation: Predict spectral characteristics from phone context

# Where Are We?

# Decision Trees: Letter-to-Sound Example

- Let's say we want to build a tree to decide how the letter "p" will sound in various words.
- Training examples:

|   |   |
|---|---|
| p | loophole peanuts pay apple |
| f | physics telephone graph photo |
| $\phi$ | apple psycho pterodactyl pneumonia |

- The pronunciation of "p" depends on its letter context.
- Task: Using the above training data, devise a series of questions about the letters to partition the letter contexts into equivalence classes to minimize the uncertainty of the pronunciation.

# Decision Trees: Letter-to-Sound Example, cont'd

- Denote the context as $\ldots L_2\ L_1\ p\ R_1\ R_2 \ldots$
- Ask potentially useful question: $R_1 =$ "h"?
- At this point we have two equivalence classes: 1. $R_1 =$ "h" and 2. $R_1 \neq$ "h".



- The pronunciation of class 1 is either "p" or "f", with "f" much more likely than "p".
- The pronunciation of class 2 is either "p" or "$\phi$"

Four equivalence classes. Uncertainty only remains in class 3.

Five equivalence classes, which is much less than enumerating each of the possibilities. No uncertainty left in the classes.

A node without children is called a leaf. Otherwise it is called an internal node

Although effective on the training data, this tree does not generalize well. It was constructed from too little data.

# Where Are We?

# Decision Tree Construction

## How to Grow a Tree

1. Find the best question for partitioning the data at a given node into 2 equivalence classes.
2. Repeat step 1 recursively on each child node.
3. Stop when there is insufficient data to continue or when the best question is not sufficiently helpful.

- Previous example - picked questions "out of the air"
- Need more principled way to chose questions

# Basic Issues to Solve

- How do we determine the best question at a node?
  - Nature of questions to be asked (next 10-15 slides or so)
  - Criterion for deciding between questions (the next set of slides after that)
- When to declare a node terminal or to continue splitting (the final part of the lecture)

# Decision Tree Construction – Fundamental Operation

- There is only 1 fundamental operation in tree construction:
  - Find the best question for partitioning a subset of the data into two smaller subsets.
  - i.e. Take a node of the tree and split it (and the data at the node) into 2 more-specific classes.

# Decision Tree Greediness

- Tree construction proceeds from the top down – from root to leaf.
- Each split is locally optimal.
- Constructing a tree in this "greedy" fashion usually leads to a good tree, but probably not globally optimal.
- Finding the globally optimal tree is an NP-complete problem: it is not practical.
- n.b.: nor does it probably matter.....

# Splitting

- At each internal node, ask a question.
  - Goal is to split data into two "purer" pieces.
- Example questions:
  - Age <= 20 (numeric).
  - Profession in (student, teacher) (categorical).
  - 5000*Age + 3*Salary – 10000 > 0 (function of raw features).

# Dynamic Questions

- The best question to ask at a node about some discrete variable x consists of the subset of the values taken by x that best splits the data.
- Search over all subsets of values taken by x. (This means generating questions on the fly during tree construction.).

$$x \in \{A, B, C\}$$

Q1:$x \in \{A\}$?    Q2:$x \in \{B\}$?    Q3:$x \in \{C\}$?

Q4:$x \in \{A, B\}$?  Q5:$x \in \{A, C\}$?  Q6:$x \in \{B, C\}$?

- Use the best question found.
- Potential problems:
  - Requires a lot of CPU. For alphabet size A there are $\sum_j \binom{A}{j}$ questions.
  - Allows a lot of freedom, making it easy to overtrain.

# Pre-determined Questions

- The easiest way to construct a decision tree is to create in advance a list of possible questions for each variable.
- Finding the best question at any given node consists of subjecting all relevant variables to each of the questions, and picking the best combination of variable and question.
- In acoustic modeling, we typically ask about 2-4 variables: the 1-2 phones to the left of the current phone and the 1-2 phones to the right of the current phone. Since these variables all span the same alphabet (phone alphabet) only one list of questions is needed.
  - Each question on this list consists of a subset of the phonetic phone alphabet.

# Sample Questions

| Phones | Letters |
|---|---|
| {P} | {A} |
| {T} | {E} |
| {K} | {I} |
| {B} | {O} |
| {D} | {U} |
| {G} | {Y} |
| {P,T,K} | {A,E,I,O,U} |
| {B,D,G} | {A,E,I,O,U,Y} |
| {P,T,K,B,D,G} | |

# More Formally - Discrete Questions

- A decision tree has a question associated with every non-terminal node.
- If x is a discrete variable which takes on values in some finite alphabet A, then a question about x has the form: $x \in S$? where S is a subset of A.
  - Let L denote the preceding letter in building a spelling-to-sound tree. Let S=(A,E,I,O,U). Then $L \in S$? denotes the question: Is the preceding letter a vowel?
  - Let R denote the following phone in building an acoustic context tree. Let S=(P,T,K). Then $R \in S$? denotes the question: Is the following phone an unvoiced stop?

# Continuous Questions

- If x is a continuous variable which takes on real values, a question about x has the form x<q? where q is some real value.
- In order to find the threshold q, we must try values which separate all training samples.



- We do not currently use continuous questions for speech recognition.

# Types of Questions

- In principle, a question asked in a decision tree can have any number (greater than 1) of possible outcomes.
- Examples:
  - Binary: Yes No.
  - 3 Outcomes: Yes No Don't_Know.
  - 26 Outcomes A B C ... Z

- In the case of determining speech recognition allophonic variation, only binary questions are used to build decision trees.

# Simple Binary Question

- A simple binary question consists of a single Boolean condition, and no Boolean operators.
- $X_1 \in S_1$? Is a simple question.
- $((X_1 \in S_1)\&\&(X_2 \in S_2))$? is not a simple question.
- Topologically, a simple question looks like:

# Complex Binary Question

- A complex binary question has precisely 2 outcomes (yes, no) but has more than 1 Boolean condition and at least 1 Boolean operator.
- $((X_1 \in S_1)\&\&(X_2 \in S_2))$? Is a complex question.
- Topologically this question can be shown as:



- All complex binary questions can be represented as binary trees with terminal nodes tied to produce 2 outcomes.

# Where Are We?

# Configurations Currently Used

- All decision trees currently used for determining allophonic variation in speech recognition use:
  - a pre-determined set
  - of simple,
  - binary questions.
  - on discrete variables.

# Tree Construction - Detailed Recap

- Let $x_1 \ldots x_n$ denote $n$ discrete variables whose values may be asked about. Let $Q_{ij}$ denote the $j$th pre-determined question for $x_i$.

- Starting at the root, try splitting each node into 2 sub-nodes:

  1. For each $x_i$ evaluate questions $Q_{i1}, Q_{i2}, \ldots$ and let $Q_i'$ denote the best.
  2. Find the best pair $x_i, Q_i'$ and denote it $x', Q'$
  3. If $Q'$ is not sufficiently helpful, make the current node a leaf.
  4. Otherwise, split the current node into 2 new sub-nodes according to the answer of question $Q'$ on variable $x'$.

- Stop when all nodes are either too small to split further or have been marked as leaves.

# Question Evaluation

- The best question at a node is the question which maximizes the likelihood of the training data at that node after applying the question.

# Question Evaluation, cont'd

- For simplicity, assume the output is a single discrete variable $x$ with $M$ outcomes (e.g., illnesses, pronunciations, etc.)
- Let $x^1, x^2, \ldots, x^N$ be the data samples
- Let each of the $M$ outcomes occur $c_j$ times in the overall sample, $j = 1 \ldots M$
- Let $Q_i$ be a question which partitions this sample into left and right sub-samples of sizes $N = n^l + n^r$.
- Let $c_j^l, c_j^r$ denote the frequency of the $j$th outcome in the left and right sub-samples, $n^l = \sum_j c_j^l$, $n^r = \sum_j c_j^r$
- The best question $Q'$ for is defined to be the one which maximizes the conditional (log) likelihood of the combined sub-samples.

# log likelihood computation

- The likelihood of the data, given that we ask question Q

$$L(x^1, \ldots, x^N | Q) = \prod_{j=1}^{M} (p_j^l)^{c_j^l} \prod_{j=1}^{M} (p_j^r)^{c_j^r}$$

$$\log L(x^1, \ldots, x^N | Q) = \sum_{j=1}^{M} c_j^l \log p_j^l + \sum_{j=1}^{M} c_j^r \log p_j^r$$

- The above assumes we know the "true" probabilities $p_j^l, p_j^r$

# log likelihood computation (continued)

- Using the maximum likelihood estimates of $p_j^l, p_j^r$ gives:

$$
\begin{aligned}
\log L(x^1, \ldots, x^N | Q) &= \sum_{j=1}^{M} c_j^l \log \frac{c_j^l}{n^l} + \sum_{j=1}^{M} c_j^r \log \frac{c_j^r}{n^r} \\
&= \sum_{j=1}^{M} c_j^l \log c_j^l - \log n^l \sum_{j=1}^{M} c_j^l + \sum_{j=1}^{M} c_j^r \log c_j^r - \log n^r \sum_{j=1}^{M} c_j^r \\
&= \sum_{j=1}^{M} \{c_j^l \log c_j^l + c_j^r \log c_j^r\} - n^l \log n^l - n^r \log n^r
\end{aligned}
$$

- The best question is the one which maximizes this simple expression. $c_j^l, c_j^r, n^l, n^r$ are all non-negative integers.
- The above expression can be computed very efficiently using a precomputed table of $n \log n$ for non-nonegative integers $n$

# Ballad of 5.60

Free energy and entropy were swirling in his brain,
With partial differentials and Greek letters in their train,
For Delta, Sigma, Gamma, Theta, Epsilon, and Pi's,
Were driving him distracted as they danced before his eyes.

**Chorus:** Glory, Glory, dear old Thermo,
Glory, Glory, dear old Thermo,
Glory, Glory, dear old Thermo,
It'll get you by and by.

# Entropy

- Let $x$ be a discrete random variable taking values $a_1, \ldots, a_M$ with probabilities $p_1, \ldots, p_M$ respectively.

- Define the entropy of the probability distribution $p = (p_1 p_2 \ldots p_M)$

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

$$H = 0 \Leftrightarrow p_j = 1 \text{ for some } j \text{ and } p_i = 0 \text{ for } i \neq j$$

- $H >= 0$

- Entropy is maximized when $p_i = 1/M$ for all $i$. Then $H = \log_2 M$

- Thus $H$ tells us something about the sharpness of the distribution $p$.

# What does entropy look like for a binary variable?

# Entropy and Likelihood

- Let $x$ be a discrete random variable taking values $a_1, \ldots a_M$ with probabilities $p_1, \ldots, p_M$ respectively.
- Let $x^1, \ldots, x^M$ be a sample of $x$ in which $a_i$ occurs $c_i$ times
- The sample log likelihood is: $\log L = \sum\limits_{i=1}^{M} c_i \log p_i$
- The maximum likelihood estimate of $p_i$ is $\hat{p}_i = c_i/N$
- Thus, an estimate of the sample log likelihood is
  $\log \hat{L} = \sum\limits_{i=1}^{M} N\hat{p}_i \log_2 \hat{p}_i \propto -\hat{H}$
- Therefore, maximizing likelihood $\Leftrightarrow$ minimizing entropy.

# "p" tree, revisited

| | | |
|---|---|---|
| p | loophole peanuts pay apple | $c_p = 4$ |
| f | physics telephone graph photo | $c_f = 4$ |
| $\phi$ | apple psycho pterodactyl pneumonia | $c_\phi = 4, N = 12$ |

- Log likelihood of the data at the root node is
  - $\log_2 L(x^1, \ldots, x^{12}) = \sum\limits_{i=1}^{3} c_i \log_2 c_i - N \log_2 N$
  - $= 4 \log_2 4 + 4 \log_2 4 + 4 \log_2 4 - 12 \log_2 12 = -19.02$
- Average entropy at the *root node* is
  - $H(x^1, \ldots, x^{12}) = -\log_2 L(x^1, \ldots, x^{12})/N$
  - $= 19.02/12 = 1.58$ bits
- Let's now apply the above formula to compare three different questions.

# "p" tree revisited: Question A

Remember formulae for Log likelihood of data:

$$\sum_{i=1}^{M} \{ c_i^l \log c_i^l + c_i^r \log c_i^r \} - n^l \log n^l - n^r \log n^r$$

Log likelihood of data after applying question A is:

$$\log_2 L(x^1, \ldots, x^{12} | Q_A) = \overbrace{1 \log_2 1}^{c_p^l} + \overbrace{4 \log_2 4}^{c_t^l} + \overbrace{3 \log_2 3}^{c_p^r} + \overbrace{4 \log_2 4}^{c_\phi^r} - \overbrace{5 \log_2 5}^{n^l} - \overbrace{7 \log_2 7}^{n^r} = -10.51$$

Average entropy of data after applying question A is

$$H(x^1, \ldots, x^{12} | Q_A) = -\log_2 L(x^1, \ldots, x^{12} | Q_A) / N = 10.51/12 = .87 \text{ bits}$$

Increase in log likelihood due to question A is -10.51 - (-19.02) = 8.51
Decrease in entropy due to question A is 1.58-.87 = .71 bits

Knowing the answer to question A provides 0.71 bits of information about the pronunciation of p. A further 0.87 bits of information is still required to remove all the uncertainty about the pronunciation of p.

# "p" tree revisited: Question B

Log likelihood of data after applying question B is:

$\log_2 L(x^1, \ldots, x^{12} | Q_B) =$
$2 \log_2 2 + 2 \log_2 2 + 3 \log_2 3 + 2 \log_2 2 + 2 \log_2 2 - 7 \log_2 7 - 5 \log_2 5 = -18.51$

Average entropy of data after applying question B is

$$H(x^1, \ldots, x^{12} | Q_B) = -\log_2 L(x^1, \ldots, x^{12} | Q_B)/N = 18.51/12 = .87 \text{ bits}$$

Increase in log likelihood due to question B is -18.51 - (-19.02) = .51
Decrease in entropy due to question B is 1.58-1.54 = .04 bits

Knowing the answer to question B provides 0.04 bits of information (very little) about the pronunciation of p.

# "p" tree revisited: Question C

Log likelihood of data after applying question C is:

$\log_2 L(x^1, \ldots, x^{12} | Q_C) =$
$2 \log_2 2 + 2 \log_2 2 + 2 \log_2 2 + 2 \log_2 2 + 4 \log_2 4 - 4 \log_2 4 - 8 \log_2 8 = -16.00$
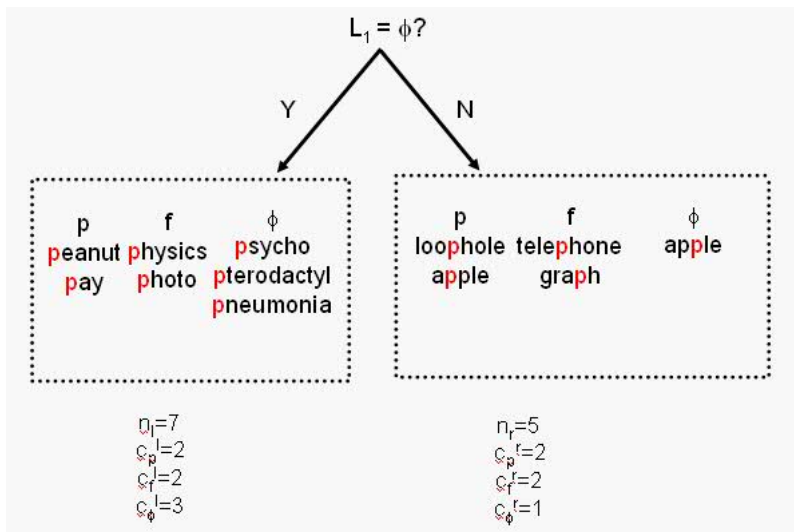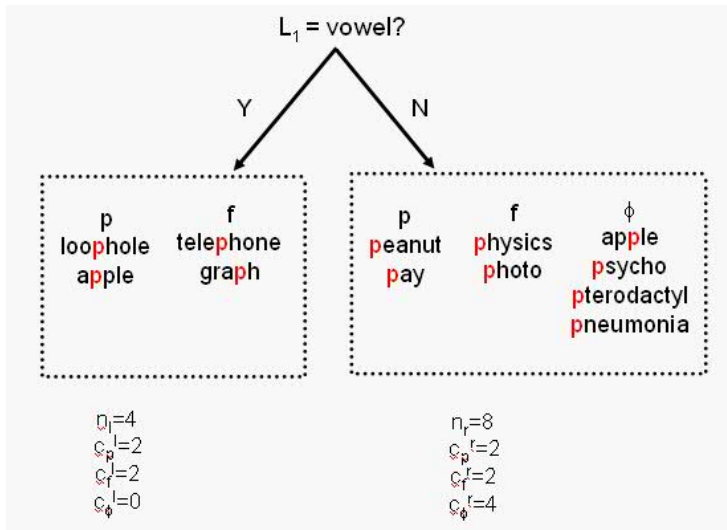
Average entropy of data after applying question C is

$$H(x^1, \ldots, x^{12} | Q_C) = - \log_2 L(x^1, \ldots, x^{12} | Q_C) / N = 16/12 = 1.33 \text{ bits}$$

Increase in log likelihood due to question C is -16 + 19.02 = 3.02
Decrease in entropy due to question C is 1.58-1.33 = .25 bits

Knowing the answer to question C provides 0.25 bits of information about the pronunciation of p.

# Comparison of Questions A, B, C

- Log likelihood of data given question:
  - A -10.51.
  - B -18.51.
  - C -16.00.
- Average entropy (bits) of data given question:
  - A 0.87.
  - B 1.54.
  - C 1.33.
- Gain in information (in bits) due to question:
  - A 0.71.
  - B 0.04.
  - C 0.25.
- These measures all say the same thing:
  - Question A is best. Question C is 2nd best. Question B is worst.

# Where Are We?

# Using Decision Trees to Model Context Dependence in HMMs

**Listen closely, this is the whole point of this lecture!**

- Remember that the pronunciation of a phone depends on its context.
- Enumeration of all triphones is one option but has problems
- Idea is to use decision trees to group triphones in a top-down manner.

# Using Decision Trees to Model Context Dependence in HMMs

- Align training data (feature vectors) against set of phonetic-based HMMs
- For each feature vector, tag it with ID of current phone and the phones to left and right.



|     | $L_1$ | C  | $R_1$ |
|-----|-------|----|-------|
| 07: | \|    | K  | AE    |
| 08: | \|    | K  | AE    |
| 09: | K     | AE | T     |
| 10: | K     | AE | T     |
| 18: | K     | AE | T     |
| 19: | AE    | T  | \|    |

# Using Decision Trees to Model Context Dependence in HMMs

- For each phone, create a decision tree by asking questions about the phones on left and right to maximize likelihood of data.
- Leaves of tree represent context dependent models for that phone.
- During training and recognition, you know the phone and its context (why?) so no problem in identifying the context-dependent models on the fly.

# New Problem: dealing with real-valued data

- We grow the tree so as to maximize the likelihood of the training data (as always), but now the training data are real-valued vectors.
- Can't use the discrete distribution we used for the spelling-to-sound example (why?)
- instead, estimate the likelihood of the acoustic vectors during tree construction using a diagonal Gaussian model.

# Diagonal Gaussian Likelihood

Let $Y = y_1, y_2 \ldots, y_n$ be a sample of independent p-dimensional acoustic vectors arising from a diagonal Gaussian distribution with mean $\vec{\mu}$ and variances $\vec{\sigma^2}$. Then

$$\log L(Y|DG(\vec{\mu}, \vec{\sigma_2})) = \frac{1}{2} \sum_{i=1}^{n} \{p \log 2\pi + \sum_{j=1}^{p} \log \sigma_j^2 + \sum_{j=1}^{p} (y_{ij} - \mu_j)^2/\sigma_j^2\}$$

The maximum likelihood estimates of $\vec{\mu}$ and $\vec{\sigma^2}$ are

$$\hat{\mu}_j = 1/n \sum_{i=1}^{n} y_{ij}, j = 1, \ldots, p$$

$$\hat{\sigma_j^2} = 1/n \sum_{i=1}^{n} y_{ij}^2 - \mu_j^2, j = 1, \ldots p$$

Hence, an estimate of log L(Y) is:

$$\log L(Y|DG(\vec{\mu}, \vec{\sigma_2})) = 1/2 \sum_{i=1}^{n} \{p \log 2\pi + \sum_{j=1}^{p} \log \hat{\sigma_j^2} + \sum_{j=1}^{p} (y_{ij} - \hat{\mu}_j)^2/\hat{\sigma_j^2}\}$$

# Diagonal Gaussian Likelihood

Now

$$\sum_{i=1}^{n}\sum_{j=1}^{p}(y_{ij}-\hat{\mu}_j)^2/\hat{\sigma}_j^2 = \sum_{j=1}^{p}\frac{1}{\hat{\sigma}_j^2}\sum_{i=1}^{n}(y_{ij}^2) - 2\hat{\mu}_j\sum_{i=1}^{n}y_{ij} + n\hat{\mu}_j^2$$

$$= \sum_{j=1}^{p}\frac{1}{\hat{\sigma}_j^2}\left\{(\sum_{i=1}^{n}y_{ij}^2) - n\hat{\mu}_j^2\right\}$$

$$= \sum_{j=1}^{p}\frac{1}{\hat{\sigma}_j^2}n\hat{\sigma}_j^2 = \sum_{j=1}^{p}n$$

Hence

$$\log L(Y|DG(\hat{\mu},\hat{\sigma}^2)) = -1/2\{\sum_{i=1}^{n}p\log 2\pi + \sum_{i=1}^{n}\sum_{j=1}^{p}\hat{\sigma}_j^2 + \sum_{j=1}^{p}n\}$$

$$= -1/2\{np\log 2\pi + n\sum_{j=1}^{p}\hat{\sigma}_j^2 + np\}$$

## Diagonal Gaussian Splits

- Let Q be a question which partitions $Y$ into left and right sub-samples $Y_l$ and $Y_r$, of size $n_l$ and $n_r$.
- The best question is the one which maximizes $\log L(Y_l) + logL(Y_r)$
- Using a diagonal Gaussian model.

$$\log L(Y_l \mid DG(\hat{\mu}_l, \hat{\sigma}_l^2)) + \log L(Y_r \mid DG(\hat{\mu}_r, \hat{\sigma}_r^2))$$

$$= -\frac{1}{2}\{n_l p \log(2\pi) + n_l \sum_{j=1}^{p} \log \hat{\sigma}_{lj}^2 + n_l p$$

Common to all splits

$$+ n_r p \log(2\pi) + n_r \sum_{j=1}^{p} \log \hat{\sigma}_{rj}^2 + n_r p\}$$

$$= -\frac{1}{2}\{np \log(2\pi) + np\} - \frac{1}{2}\{n_l \sum_{j=1}^{p} \log \hat{\sigma}_{lj}^2 + n_r \sum_{j=1}^{p} \log \hat{\sigma}_{rj}^2\}$$

# Diagonal Gaussian Splits, cont'd

Thus, the best question $Q$ minimizes:

$$D_Q = n_l \sum_{j=1}^{p} \log \hat{\sigma}_{lj}^2 + n_r \sum_{j=1}^{p} \log \hat{\sigma}_{rj}^2$$

Where

$$\hat{\sigma}_{lj}^2 = 1/n^l \sum_{y \in Y_l} y_j^2 - 1/{n^l}^2 (\sum_{y \in Y_l} y_j)^2$$

$$\hat{\sigma}_{rj}^2 = 1/n^r \sum_{y \in Y_r} y_j^2 - 1/{n^r}^2 (\sum_{y \in Y_r} y_j)^2$$

$D_Q$ involves little more than summing vector elements and their squares.

# How Big a Tree?

- Cross-validation.
  - Measure performance on a held-out data set.
  - Choose the tree size that maximizes the likelihood of the held-out data.
- In practice, simple heuristics seem to work well.
- A decision tree is fully grown when no terminal node can be split.
- Reasons for not splitting a node include:
  - Insufficient data for accurate question evaluation.
  - Best question was not very helpful / did not improve the likelihood significantly.
  - Cannot cope with any more nodes due to CPU/memory limitations.

# Recap

- Given a word sequence, we can construct the corresponding Markov model by:
  - Re-writing word string as a sequence of phonemes.
  - Concatenating phonetic models.
  - Using the appropriate tree for each phone to determine which allophone (leaf) is to be used in that context.
- In actuality, we make models for the HMM arcs themselves
  - Follow same process as with phones - align data against the arcs
  - Tag each feature vector with its arc id and phonetic context
  - Create decision tree for each arc.

# Example

The rain in Spain falls ....

Look these words up in the dictionary to get:

DH AX | R EY N | IX N | S P EY N | F AA L Z | ...

Rewrite phones as states according to phonetic model

$DH_1 DH_2 DH_3 AX_1 AX_2 AX_3 R_1 R_2 R_3 EY_1 EY_2 EY_3$ ...

Using phonetic context, descend decision tree to find leaf sequences

$DH_{1\_5} DH_{2\_27} DH_{3\_14} AX_{1\_53} AX_{2\_37} AX_{3\_11} R_{1\_42} R_{2\_46}$ ....

Use the Gaussian mixture model for the appropriate leaf as the observation probabilities for each state in the Hidden Markov Model.

# Some Results

| System | T1 | T2 | T3 | T4 |
|--------|-----|-----|-----|-----|
| Monophone | 5.7 | 7.3 | 6.0 | 9.7 |
| Triphone | 3.7 | 4.6 | 4.2 | 7.0 |
| Arc-Based DT | 3.1 | 3.8 | 3.4 | 6.3 |

- From Julian Odell's PhD Thesis (Cambridge U., 1995)
- Word error rates on 4 test sets associated with 1000 word vocabulary (Resource Management) task

# Strengths & Weaknesses of Decision Trees

- **Strengths**.
    - Easy to generate; simple algorithm.
    - Relatively fast to construct.
    - Classification is very fast.
    - Can achieve good performance on many tasks.
- **Weaknesses**.
    - Not always sufficient to learn complex concepts.
    - Can be hard to interpret. Real problems can produce large trees...
    - Some problems with continuously valued attributes may not be easily discretized.
    - Data fragmentation.

# Course Feedback

- Was this lecture mostly clear or unclear?
- What was the muddiest topic?
- Other feedback (pace, content, atmosphere, etc.).