

MP1 – Audio Features/Nearest Neighbor

ECE 417 – Multimedia Signal Processing
University of Illinois at Urbana-Champaign
Spring 2014

Instructor: Professor Hasegawa-Johnson

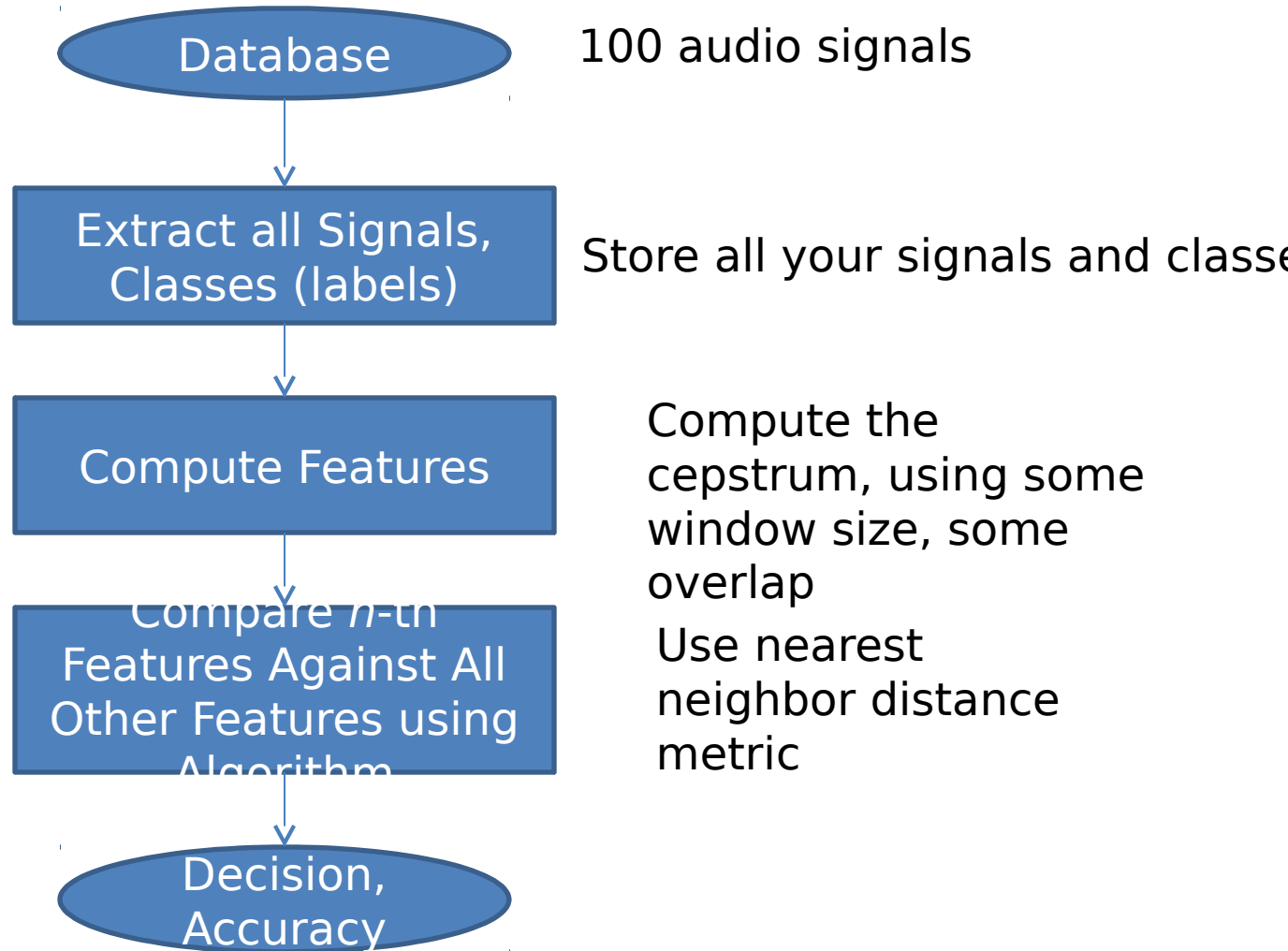
TA's: Daniel Soberal, Po-Sen Huang

Prepared by Daniel Soberal

Goals

- Given a dataset of N different audio samples of people speaking, be able to:
 - Extract cepstral features
 - Use nearest neighbor to perform speaker recognition
 - Use nearest neighbor to perform speech recognition

The System: Highly Generalized



The Data

- 100 different audio samples
 - 4 different speakers, labelled A,B,C,D.
 - 5 different spoken digits, labelled 1,2,3,4,5.
 - Various instances, labelled a,b,c,d,e
- File format:
 - Name is [Speaker][Digit][Instance].wav
- Each audio sample is called an *observation*.

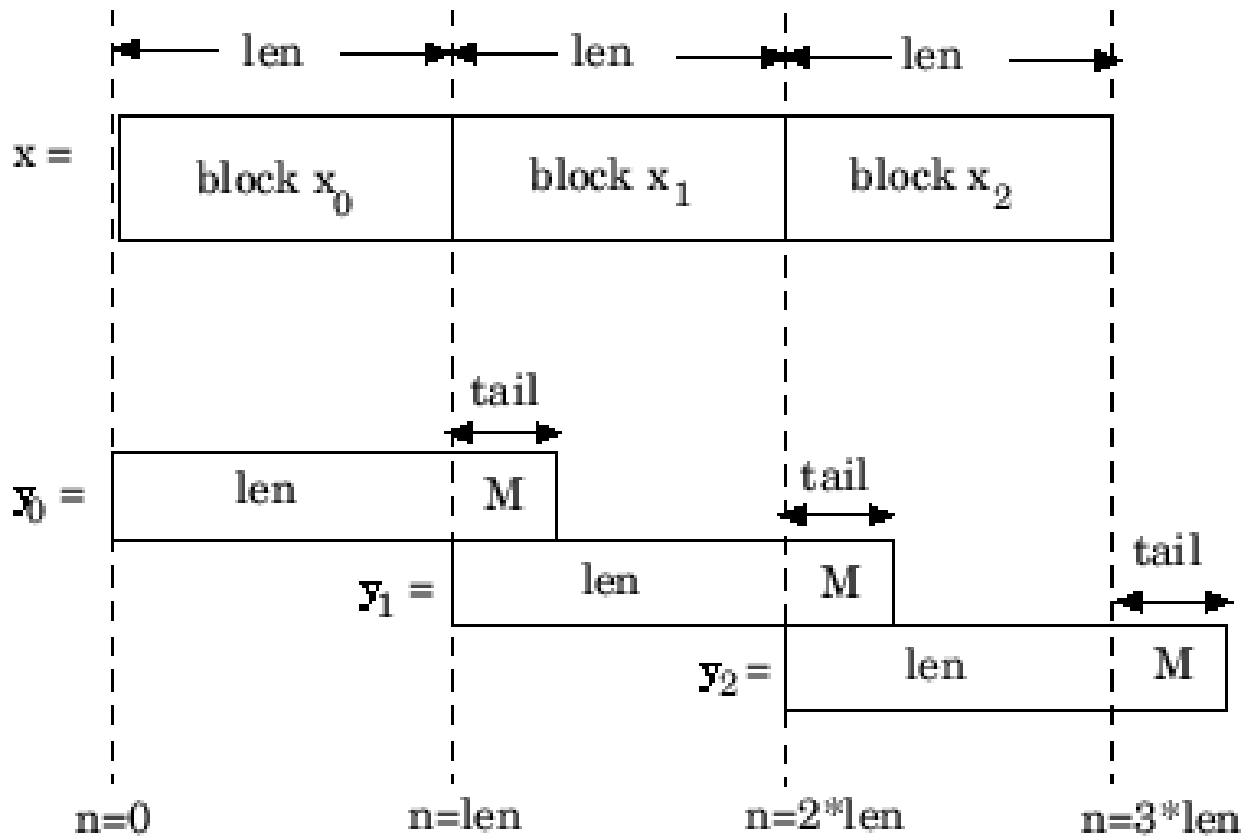
Audio Extraction

- Each file is a .wav file.
- Use audioread (or wavread) to extract the data.
- Result is stereo vector (2 channels).
Take the right channel.

Signal Pre-conditioning

- Speech is a non-stationary, time-varying signal.
- For most phonemes, the properties of the speech remain invariant for a short period of time (5-100 ms)
- We want to separate the signal into overlapping frames.
 - 10% Overlap
- Since we are inherently windowing the signal anyway by chopping it up into frames, let's use a window with small sidelobes.
 - Hamming window

Signal Pre-Conditioning (Continued)



The Cepstrum

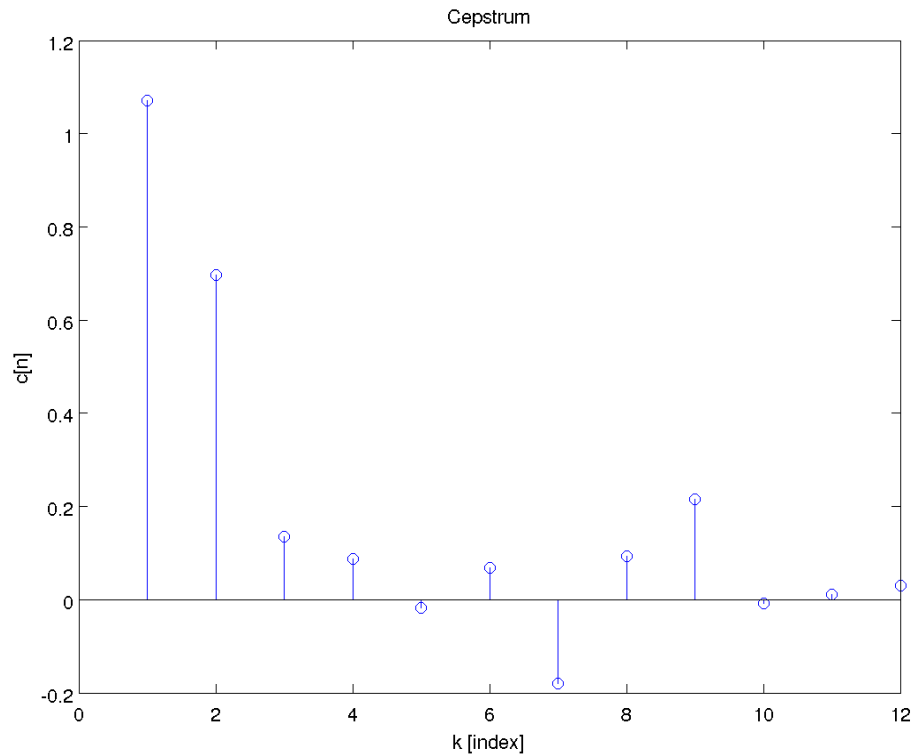
- General formula:

$$c[n] = \mathcal{F}^{-1}\{\log_{10}|\mathcal{F}\{x[n]\}|\}$$

- The FFT and IFFT function in Matlab's signal processing toolbox will be useful here.
 - Implementation of the FFT itself is best left to courses on the computational aspects of DSP...
- Apply the cepstrum calculation to each frame of windowed data. If there are M frames of length L, you get an LxM matrix.
 - We are only concerned with the first 12 coefficients. This reduces your matrix to 12xM.
- Unroll this into a single column vector that is (12M)x1

Cepstrum (Continued)

- Cepstrum example with first coefficient removed.



K-Nearest Neighbor

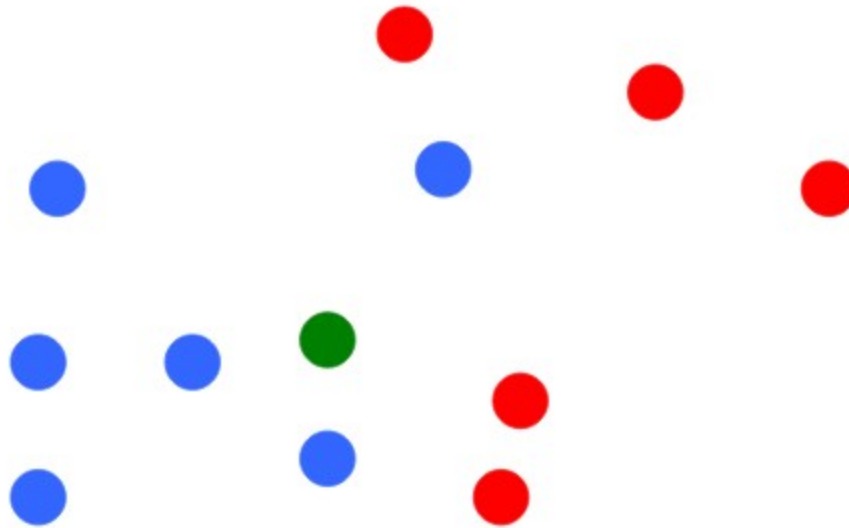
- Non-parametric approach (uses data, not model)
- Steps:
 - Given a data point, find the distance between that point and all other points in the dataset.

$$d_{mn} = \sum_{k=0}^{K-1} |c_m[k] - c_n[k]|^2$$

- Find the K closest distances and corresponding classes to the given point.
- Of these, find the class that occurs the most. This class is the one that matches the input data.

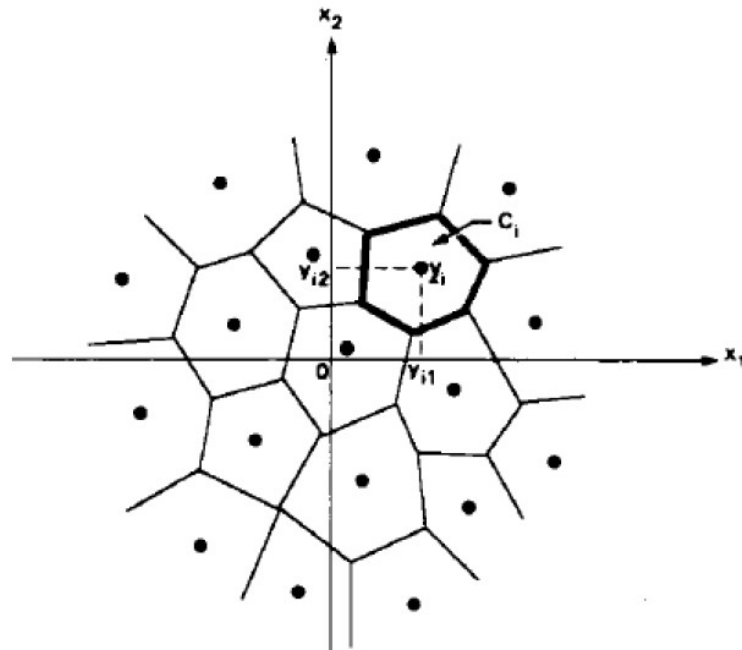
Illustration

- The green dot should be classified as what class (blue or red)?



K-Nearest Neighbor (Continued)

- This results in a voronoi tessellation.
- Essentially, breaks up space into polygons, where each polygon is associated with a point.



Overall Procedure

- Get all audio samples, resize them so that they are all the same length.
- Compute the cepstrum using 10% overlap, a Hamming window, and various window sizes (100, 500, 10000).
- Perform nearest neighbor:
 - Remove all data corresponding to the same speaker for every sample when doing speech recognition.
 - Remove all data corresponding to the same digit for every sample when doing speaker recognition.

A Few Tricks for Removing Data and Indexing

- Logical Indexing will be your friend in this MP
- Separate labels into arrays of strings
- Store data in cells or arrays
- Use commands like 'find' or just use logical indexing to find indices of like classes (labels)
- Some Example Code
 - `CurrentData = data{1};`
 - `CurrentSpeaker = SpeakerLabel(1,:);`
 - `CurrentDigit = SpeechLabel(1,:);`
 - `IndToKeep = find(SpeechLabel~=CurrentDigit);`
 - `Phi = CurrentData(IndToKeep);`

Results (Tables)

- Tables for the 1NN and 5NN results for each speaker, digit, and overall accuracy in both cases.
- This must be done for the raw data case and the 3 window-length cases.

Results (Graphs)

- Graphs for
 - Recognition Accuracy vs Window Length for 12-Coefficient Cepstrum Using 1-NN (individual digits and overall results overlaid)
 - Recognition Accuracy vs Window Length for 12-Coefficient Cepstrum Using 5 NN (individual digits and overall results overlaid)
 - Speaker Recognition vs Window Length using Nearest Neighbor (individual speakers and overall results overlaid)
 - Speaker Recognition vs Window Length using 5-Nearest Neighbor (individual speakers and overall results overlaid)