

## Lecture 19

---

### Input Processing

- ◆ Mel Frequency Scale
- ◆ Input Signal Preprocessing
  - Discrete Cosine Transform
  - Time Derivative estimates
  - Feature Decorrelation
  - Feature Vector length reduction
- ◆ Gaussian Mixtures
- ◆ Speech Recognition Summary

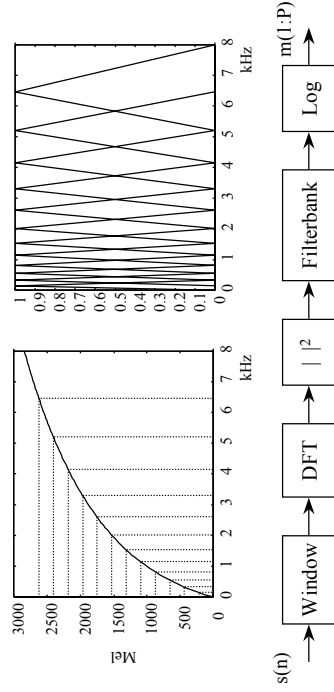
## Feature Vector Requirements

---

- ◆ When different people say the same phoneme, the feature vectors should have similar values.
- ◆ Different phonemes from the same or different speakers should give dissimilar values.
- ◆ For different examples of the same phoneme, the features should be independent and uncorrelated: this allows us to multiply their probabilities.
- ◆ For different examples of the same phoneme, each feature should preferably follow a probability distribution that is well described as a sum of gaussians.
- ◆ The features should not be affected by the amplitude of the speech signal otherwise recognition performance would vary with your distance from the microphone.

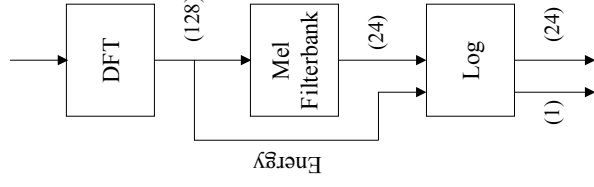
## Mel Frequency Scale

- ◆ The feature vector must discriminate between speech sounds using as few components as possible to reduce computation.
- ◆ The human ear has better frequency resolution at low frequencies. The *mel* scale relates perceived pitch to frequency: linear at low  $f$ , logarithmic at high  $f$ :
  - $mel(f) = 2595 \log_{10}(1 + f/700)$  where  $f$  is in Hz
  - Form a mel-spaced filterbank by setting the centre frequencies to equally spaced mel values.



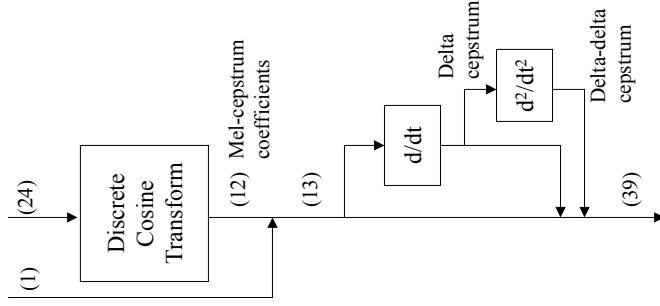
## Preprocessing: Stage 1

- ◆ Divide signal into overlapping 25 ms segments at 10 ms intervals
- ◆ Apply Hamming window and take FFT
- ◆ Smooth the spectrum with a mel filterbank
  - mel filterbank concentrates data values in the more significant part of the spectrum
- ◆ Take the log of the mel spectrum
  - variations in signal level just cause a DC shift in the log spectrum
  - gaussian approximation is more nearly true for log spectrum than for the power spectrum directly



## Preprocessing: Stage 2

- ◆ Discrete Cosine Transform (DCT)
  - reduces correlation between coefficients
  - compresses information into fewer low-order coefficients
  - output is the *mel-cepstrum*
  - DC component is ignored to make it independent of signal level
- ◆ First and Second time derivatives
  - provide additional information about how the spectrum is changing with time
- ◆ Result is a 39 element feature vector (or 38 if you drop the log energy).



## Discrete Cosine Transform

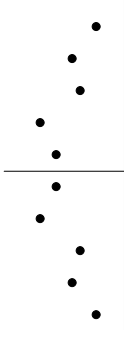
- ◆ The discrete cosine transform (DCT) of  $m_1, \dots, m_p$  is defined by

$$c_k = \sum_{p=1}^p m_p \cos(k(p-1/2)\pi / P)$$

- ◆ The DCT of these points



is equal to the DFT of these points



with a phase shift to centre the time origin.

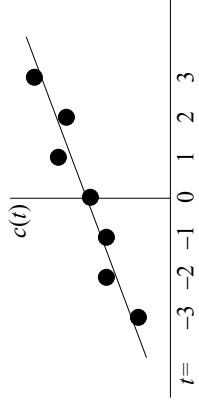
- ◆ Taking the DCT of the +ve frequency spectrum is essentially the same as taking the DFT of the symmetrical ±ve frequency spectrum.
- ◆ There are efficient algorithms for calculating the DCT

## Polynomial Fitting

- ◆ To fit a polynomial to a set of points  $x_i, y_i$   
for  $i=1, 2, \dots, N$ : 
$$y_i = \sum_{k=0}^P a_k x_i^k$$
- ◆ Error  $E = \sum_{i=1}^N e_i^2$  where  $e_i = y_i - \sum_{k=0}^P a_k x_i^k$
- ◆ Minimize  $E$  by differentiating w.r.t.  $a_m, m=0:P$   
$$\frac{\partial E}{\partial a_m} = \sum_{i=1}^N 2e_i \frac{\partial e_i}{\partial a_m} = -2 \sum_{i=1}^N \left( y_i - \sum_{k=0}^P a_k x_i^k \right) \times x_i^m$$
- ◆ Hence we get  $P+1$  equations (same as LPC)  
$$\sum_{k=0}^P \left( a_k \sum_{i=1}^N x_i^{k+m} \right) = \sum_{i=1}^N y_i x_i^m \quad \text{for } m = 0, 1, \dots, P$$
- ◆ In matrix form (each value of  $m$  gives one row):

$$\begin{pmatrix} \sum x^0 & \sum x^1 & \sum x^2 & \dots & \sum x^P \\ \sum x^1 & \sum x^2 & \sum x^3 & \dots & \sum x^{P+1} \\ \sum x^2 & \sum x^3 & \sum x^4 & \dots & \sum x^{P+2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \\ a_2 \\ \vdots \\ a_P \end{pmatrix} = \begin{pmatrix} \sum yx^0 \\ \sum yx^1 \\ \sum yx^2 \\ \vdots \end{pmatrix}$$

## Cepstral Time-Derivatives

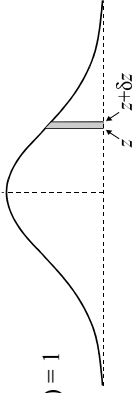


- ◆ Want to estimate  $dc/dt$  by fitting a line
  - Few points  $\Rightarrow$  noisy estimate
  - Many points  $\Rightarrow$  can't follow time variations
- ◆ Fit a 1<sup>st</sup>-order polynomial to  $2T+1$  points:
 
$$\begin{pmatrix} \sum_{t=-T}^T t^0 \\ \sum_{t=-T}^T t^1 \\ \sum_{t=-T}^T t^2 \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \end{pmatrix} = \begin{pmatrix} \sum_{t=-T}^T c(t)t^0 \\ \sum_{t=-T}^T c(t)t^1 \end{pmatrix}$$
- ◆ this simplifies to
 
$$\begin{pmatrix} 2T+1 & 0 \\ 0 & \sum_{t=-T}^T t^2 \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \end{pmatrix} = \begin{pmatrix} \sum_{t=-T}^T c(t) \\ \sum_{t=-T}^T t c(t) \end{pmatrix} \Rightarrow a_1 = \frac{\sum_{t=-T}^T t c(t)}{\sum_{t=-T}^T t^2}$$
- ◆ Typically  $T=5$  for 1<sup>st</sup> derivative and  $T=1$  for 2<sup>nd</sup>

## Multivariate Gaussian Distributions

- ◆  $Z$  is a random variable with a standard Gaussian (or Normal) probability density func:  

$$\text{pr}(Z \in [z, z + \delta z]) = (2\pi)^{-1/2} \exp(-1/2z^2) \delta z$$



- ◆ Mean:  $E(Z) = 0$   
 Variance:  $E(Z^2) = 1$

- ◆ A linear sum of multiples of Gaussian random variables gives another Gaussian random variable. This property is unique to Gaussians.
- ◆ If we have a column random vector  $Z$  with  $P$  elements each of which is an *independent* standard Gaussian random variable then

$$\begin{aligned} \text{pr}(Z \in [z, z + dz]) &= \prod_{i=1}^P (2\pi)^{-1/2} \exp(-1/2z_i^2) dz_i \\ &= (2\pi)^{-P/2} \exp\left(-1/2 \sum_{i=1}^P z_i^2\right) dz = (2\pi)^{-P/2} \exp(-1/2z^T z) dz \end{aligned}$$

- ◆ Note too that because the  $z_i$  are independent.  
 $E(z_i z_j) = 0$  whenever  $i \neq j \Rightarrow E(z z^T) = \mathbf{I}$

## Correlated Gaussian Distributions

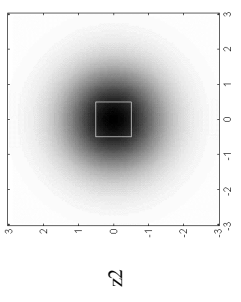
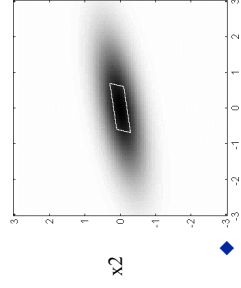
- ◆ Now suppose that  $x = A z$  where  $A$  is a non-singular matrix, then  $dx = |A| dz$  and  $z = A^{-1} x$ . Note that  $X$  is gaussian and  $E(x)$  is 0.

- ◆ The covariance matrix of  $x$  is  $C = E(xx^T)$  and is symmetric and positive definite

$$C = E(xx^T) = E(Az z^T A^T) = AE(zz^T)A^T = AA^T$$

- ◆ We can work out the pdf of  $x$

$$\begin{aligned} \text{pr}(X \in [x, x + dx]) / dx &= (2\pi)^{-P/2} \exp\left(-1/2(A^{-1}x)^T (A^{-1}x)\right) |A|^{-1} \\ &= (2\pi)^{-P/2} |A|^{-1} \exp\left(-1/2x^T A^{-T} A^{-1}x\right) \\ &= (2\pi)^{-P/2} |C|^{-1/2} \exp\left(-1/2x^T C^{-1}x\right) \end{aligned}$$



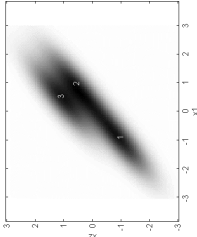
$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 1.3 & 0.1 \\ 0.2 & 0.4 \end{pmatrix} \begin{pmatrix} z_1 \\ z_2 \end{pmatrix} \Rightarrow C = \begin{pmatrix} 1.7 & 0.3 \\ 0.3 & 0.2 \end{pmatrix} \text{ and } |C| = |A|^2 = 0.25$$

## Computational Costs

- ◆ The log prob density of correlated gaussians:
 
$$\log(\text{pd}(x)) = \log\left((2\pi)^{-1/2} |C|^{-1/2} \exp\left(-\frac{1}{2} x^T C^{-1} x\right)\right) \\ = -\frac{1}{2} (P \log(2\pi) + \log(|C|) + x^T C^{-1} x)$$
- ◆ The first two terms are independent of  $x$  and can be precalculated for each state.
- ◆ For  $F$  features, the final term involves  $F^2 + F$  multiplications and  $F^2 - 1$  additions:  $39^2 = 1521$
- ◆ If the features are (or are assumed to be) independent,  $C$  is diagonal and we need  $2F$  multiplications and  $F - 1$  additions
- ◆ Probability calculations consume most of the computation in a recogniser: almost all recognisers assume feature independence
  - DCT on log spectrum improves independence
  - We can do even better by applying a linear transformation to the feature vector.

## Feature Decorrelation

- ◆ We can apply a linear transformation to our feature vectors,  $x$ , to reduce correlations.
 

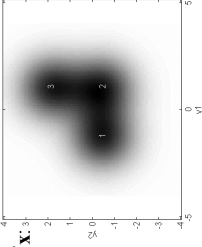


$W_s$  is the covariance matrix of state  $s$ .

$W$  is the average of the  $W_s$ : the average within-state covariance matrix.
- ◆ If we multiply the feature vectors by a matrix  $F^T$ ,  $y = F^T x$ , then the covariance matrix of  $y$  within state  $s$  is given by:
 
$$E((y - \bar{y}_s)(y - \bar{y}_s)^T) = E(F^T (x - \bar{x}_s)(x - \bar{x}_s)^T F) = F^T W_s F$$

where  $\bar{x}_s$  is the mean value of  $x$  in state  $s$ .
- ◆ We transform our data with an  $F^T$  satisfying  $F^T W F = I$ .
 

$y = F^T x$



## Eigenvectors

- ◆  $d$  is an eigenvalue of  $\mathbf{W}$  and  $\mathbf{y}$  is an associated eigenvector if
- ◆ Since  $\mathbf{W}$  is symmetric and positive definite, we can find  $F$  orthonormal eigenvectors and make them the columns of a matrix:

$$\mathbf{W}\mathbf{y}=\mathbf{y}d$$

$$\mathbf{W}\mathbf{Y}=\mathbf{Y}\mathbf{D}$$

- where  $\mathbf{D}$  is a diagonal matrix of eigenvalues
- ◆ The orthonormality of the eigenvectors means that

$$\mathbf{Y}^T\mathbf{Y}=\mathbf{I}$$

- ◆ Now we define
- ◆ This gives
- ◆ In MATLAB:

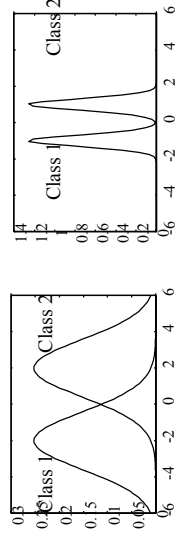
$$\mathbf{F}=\mathbf{Y}\mathbf{D}^{-1/2}$$

$$\mathbf{F}^T\mathbf{W}\mathbf{F}=\mathbf{D}^{-1/2}\mathbf{Y}^T\mathbf{W}\mathbf{Y}\mathbf{D}^{-1/2}=\mathbf{D}^{-1/2}\mathbf{Y}^T\mathbf{Y}\mathbf{D}\mathbf{D}^{-1/2}=\mathbf{I}$$

```
[Y,D]=eig(W);
F=Y * sqrt(inv(D));
```

## Class Discrimination

- ◆ We would like to make our feature vector as short as possible while preserving its ability to discriminate.



- ◆ The graphs show two possible distributions of a parameter for two different speech sounds (or classes).
- ◆ For a single parameter, Fisher's  $F$  Ratio is a measure of discriminability (the bigger the better):

$$F = \frac{\text{Variance of the class means}}{\text{Average variance within a class}}$$

- ◆ For a parameter vector, this generalises to:

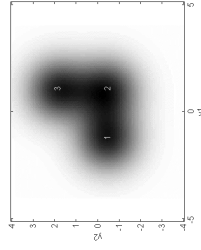
$$F = \text{trace}(\mathbf{W}^{-1}\mathbf{B})$$

where  $\mathbf{W}$  and  $\mathbf{B}$  are "average within-class" and "between-class" covariance matrices

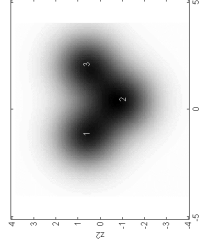
## Dimensionality Reduction

- ◆ We define **B** to be the between-state covariance matrix:
 
$$\mathbf{B} = \frac{1}{S} \sum_{s=1}^S (\bar{\mathbf{y}}_s - \bar{\bar{\mathbf{y}}})(\bar{\mathbf{y}}_s - \bar{\bar{\mathbf{y}}})^T$$
- ◆ As before we can find the eigenvalues of **B**

$$\mathbf{B}\mathbf{G}=\mathbf{G}\mathbf{L}$$
- ◆ where **G** is orthogonal and **L** diagonal.
- ◆ Set  $\mathbf{z}=\mathbf{G}^T\mathbf{y}=\mathbf{G}^T\mathbf{F}^T\mathbf{x}$
- ◆ The between-state covariance matrix is now  $\mathbf{G}^T\mathbf{B}\mathbf{G}=\mathbf{L}$
- ◆ We can discard any elements of **z** for which the corresponding element of **L** is very small. Gives reduced feature set with equal (or even better) discrimination.



$\mathbf{y}=\mathbf{F}^T\mathbf{x}$



$\mathbf{z}=\mathbf{G}^T\mathbf{y}=\mathbf{G}^T\mathbf{F}^T\mathbf{x}$

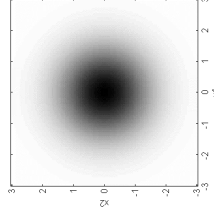
## Gaussian Mixtures

- ◆ For large vocabularies, independent gaussian model is too simple. Use instead a mixture of gaussians:

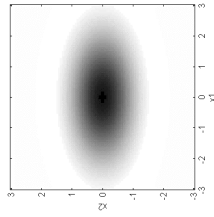
$$pd(\mathbf{x}) = \sum_{i=1}^K w_i \mathcal{N}(\mathbf{m}_i, \mathbf{C}_i) \quad \text{with} \quad \sum_{i=1}^K w_i = 1$$

$$\text{where } \mathcal{N}(\mathbf{m}_i, \mathbf{C}_i) = (2\pi)^{-1/2 P} |\mathbf{C}_i|^{-1/2} \exp(-1/2 \mathbf{x}^T \mathbf{C}_i^{-1} \mathbf{x})$$

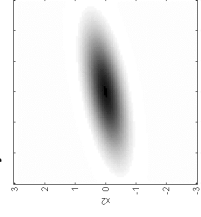
For simplicity we restrict **C** to be diagonal.



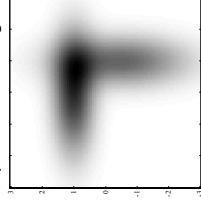
Symmetric: **W=I**



Independent: **W** diagonal



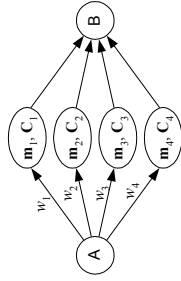
Correlated



Diagonal Gaussian Mixture



## Mixtures = Alternate HMM states



- ◆ The total probability of all paths from A to B is the sum of the individual path probabilities

$$pd(\mathbf{x}) = \sum_{i=1}^K w_i N(\mathbf{m}_i, C_i)$$

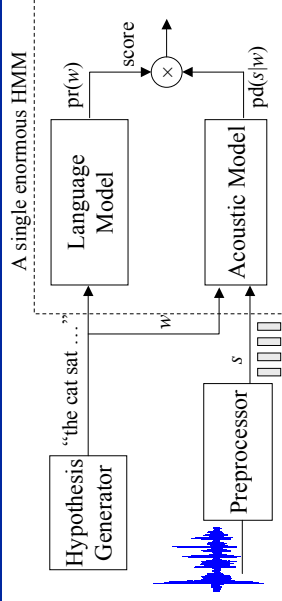
this is identical to the gaussian mixture expression.

- ◆ Once we have initial values for the model parameters we can use *Viterbi* and *Baum-Welch* procedures to train them.
- ◆ We can view gaussian mixtures as describing alternative pronunciations of a particular speech sound

## K-means Algorithm

- ◆ We need to form an initial estimate for the K mixture means,  $\mathbf{m}_i$ , and covariances,  $C_i$ .
- ◆ First create and train models with only one mixture using Viterbi training.
- ◆ Use Viterbi alignment to determine which training frames correspond to each state.
- ◆ For each state
  - Set the  $\mathbf{m}_i$  to K randomly chosen training frames
- Repeat until convergence occurs:
  - Allocate each training frame to whichever  $\mathbf{m}_i$  it is nearest to.
  - Update each  $\mathbf{m}_i$  to the mean of all the frames that were allocated to it
  - If no frames were allocated to  $\mathbf{m}_i$ , set it to a randomly chosen point from one of the other distributions.
- Set  $C_i$  to the covariance of the frames allocated to  $\mathbf{m}_i$ .

## Speech Recognition



- ◆ Preprocessor
  - Mel Cepstrum + Velocity + Acceleration
  - Linear Transform to decorrelate & reduce  $F$
- ◆ Acoustic Model
  - 60,000 triphones  $\times$  3 states  $\times$  20 features  $\times$  10 mixtures = 72,000,000 parameters to train.
- ◆ Language Model
  - Phonetic description of each word in vocabulary + trigram or quadram transition probabilities
- ◆ Dynamic model creation
  - Create storage only for models when needed: use pruning to delete models with a hopelessly low probability.
  - Trade-off memory/computation versus accuracy