

Speech Features and Speaker Classification

CSC401/2511 – Natural Language Computing – Spring 2026

Lecture 11

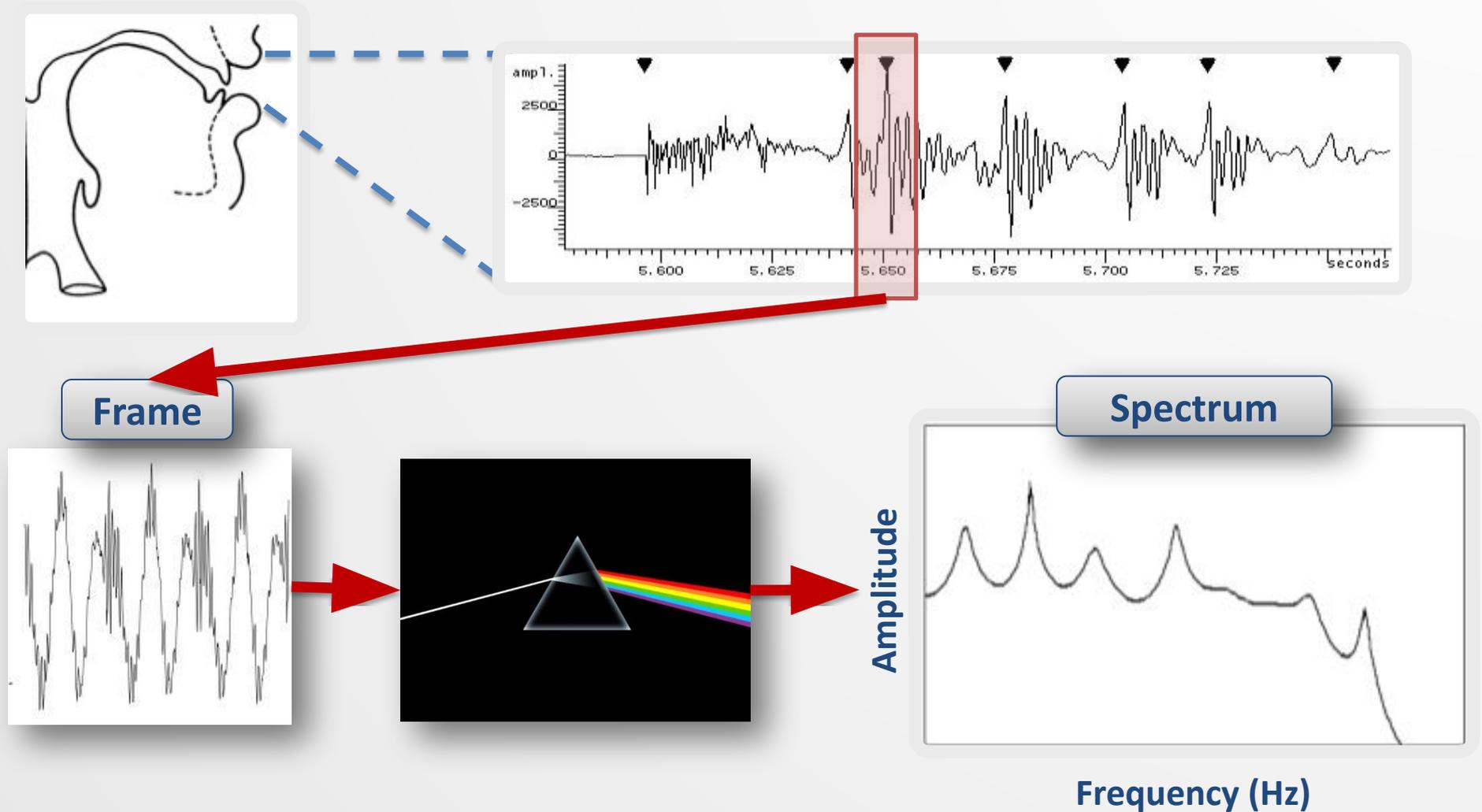
University of Toronto

Contents

- Define some common feature vectors for speech processing
- Use them as input to a GMM-based speaker classification system
- All of this is part of A3

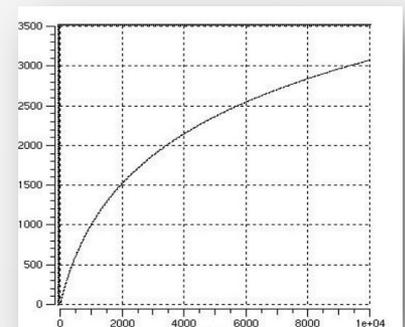
SPEECH FEATURES

Recall the spectrogram pipeline



Problems with spectrograms

- As input to speech systems, spectrograms are...
- **Too big**
 - The discrete signal is usually 16,000 samps/sec
 - 100 frames/sec x 400 samps/frame = 40,000 samps/sec!
- **Too linear**
 - Pitch perception is log-linear (recall Mels)
 - Lots of coefficients wasted on high frequencies
- **Too entangled**
 - Speaker and phoneme info is correlated

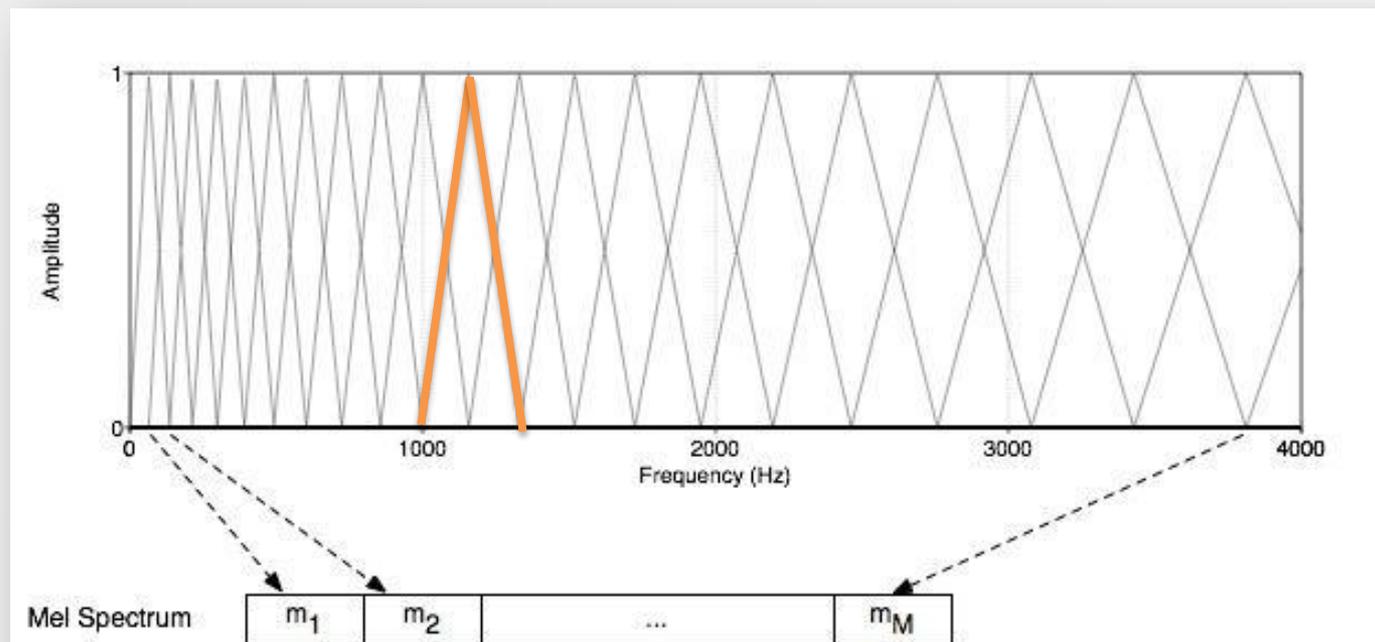


Filtering

- To reduce the size of the spectra, we **filter** it with **filters** from a **filter bank**
- Each filter is a signal whose spectrum $F_m \in \mathbb{R}^N$ picks out small a range (or **band**) of frequencies
- The bands of the M filters are overlapping and span the spectrum
- A **filter coefficient** is computed as the **log** of the dot product of the **magnitude** of the frame X_t and filter F_m spectra:
$$c_{t,m} = \log \sum_{n=1}^N |X_t|[n]|F_m|[n]$$
- If there are T frames, this gives us a real-valued feature matrix of size $T \times M$
 - $M = 40$ is a lot smaller than 400!

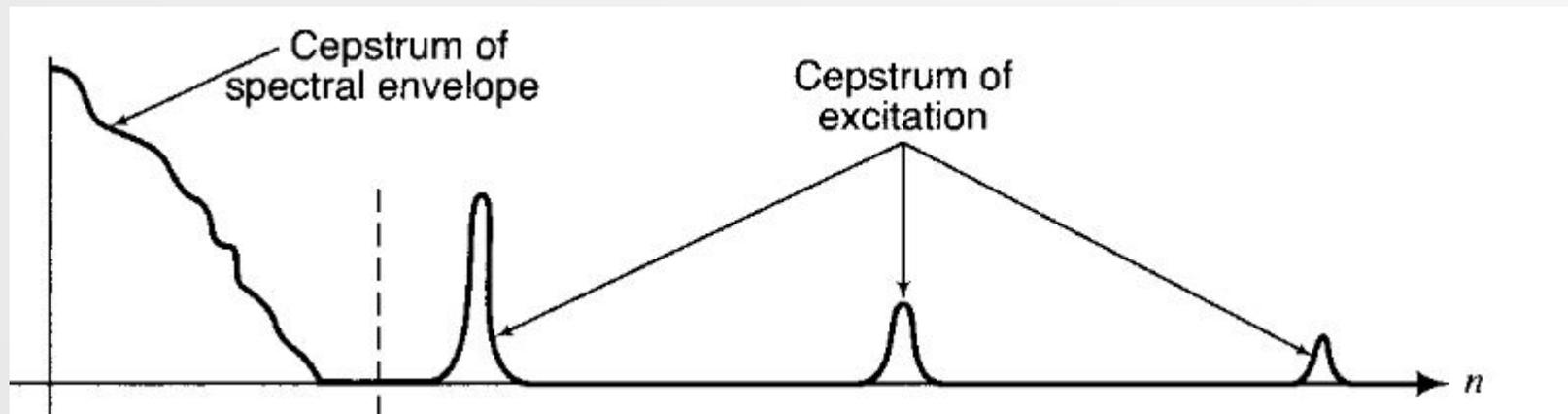
The mel-scale filter bank

- The mel-scale triangular overlapping filter bank, or **f-bank**, is a popular choice
- The filter's vertices are arranged along the mel-scale
 - Ascending frequency = wider bands



The source-filter model

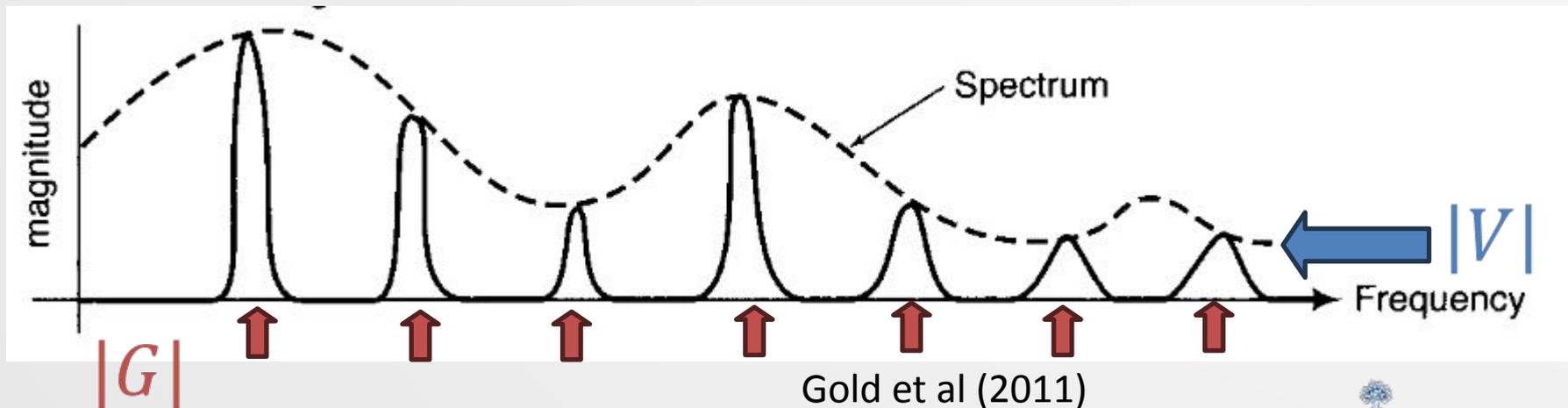
- In vowels, the sound signal emitted from the glottis g is filtered by the vocal tract v
- The **source-filter model** of speech assumes
$$|X[n]| = |G[n]||V[n]|$$
- $|V|$ is responsible for the smooth shape (envelope)
- $|G|$ is responsible for all the bumps (F0 harmonics)



Gold et al (2011)

The cepstrum

- We can get at $|V|$ by computing the **cepstrum** \hat{x}
- The cepstrum is $\log|X|$ transformed by the inverse DFT
- Because $\log|X| = \log|G| + \log|V|$, and DFT^{-1} is linear
$$\hat{x}[n] = \hat{g}[n] + \hat{v}[n]$$
- $DFT^{-1} \approx DFT$, so \hat{x} is like the spectrum of $\log|X|$
- $|V|$ is slower-moving than $|G|$, so $\hat{v}[n]$ is higher for lower n (lower frequency of frequency)



Mel-Frequency Cepstral Coefficients

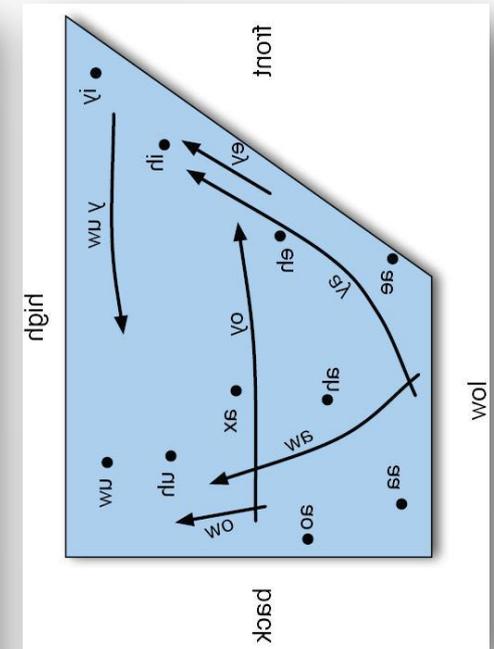
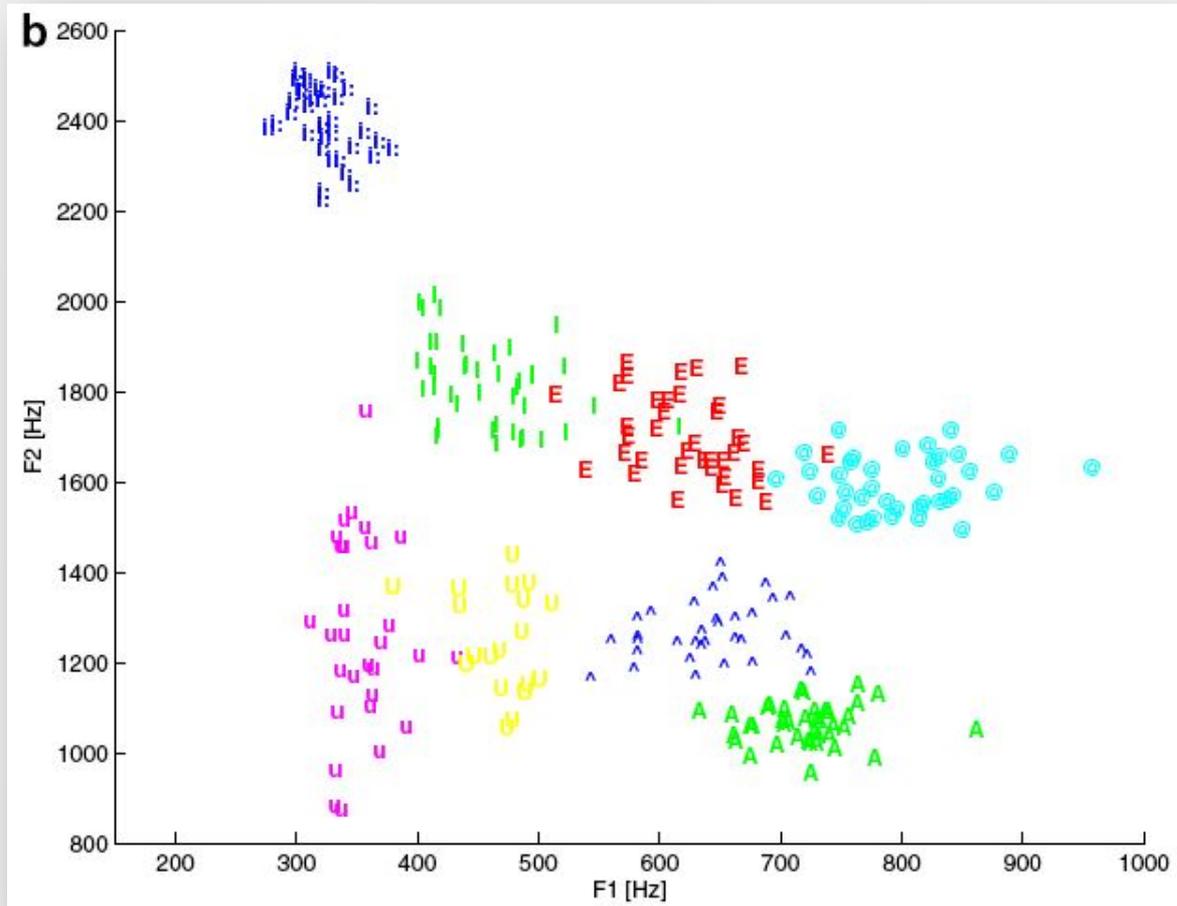
- **MFCCs** are the coefficients of the cepstrum of F-bank coefficients
- Altogether



- MFCCs are useful for models which can't handle speaker correlations themselves, like (diagonal) GMMs
- F-banks are better for those which can, like NNs

GAUSSIAN MIXTURES

Classifying speech sounds



Note: The vowel trapezoid's dimensions were physical

- Speech sounds can cluster. This graph shows vowels, each in their own colour, according to the second two formants.

Classify speakers by cluster attributes

- Similarly, all of the speech produced by one **speaker** will cluster differently in the **Mel space** than speech from another speaker.
 - We can \therefore decide if a given observation comes from one speaker or another.

		0	1	...	T
M F C C	1			...	
	2			...	
	3			...	

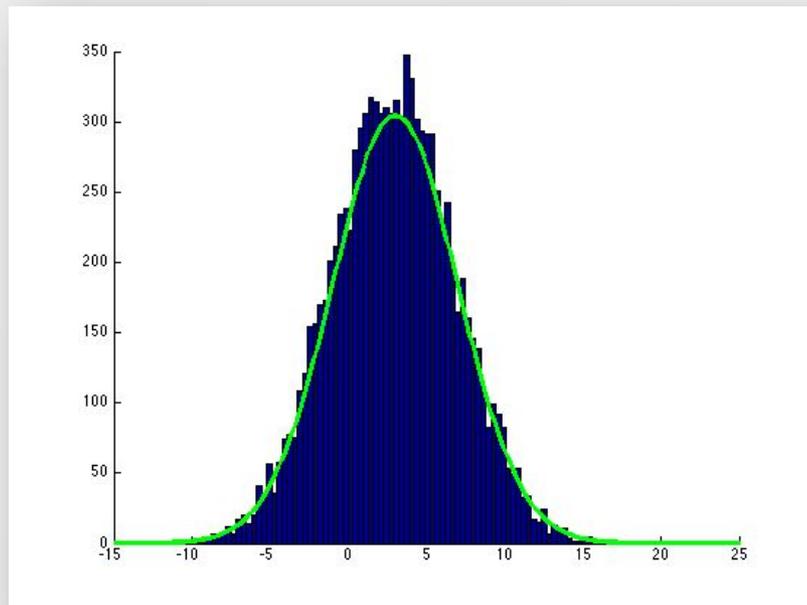
	42			...	

Observation matrix

$$P(\text{orange bar} \mid \text{woman on phone}) > P(\text{orange bar} \mid \text{man in uniform on phone})$$

Fitting continuous distributions

- Since we are operating with **continuous** variables, we need to **fit continuous probability** functions to a **discrete number** of observations.

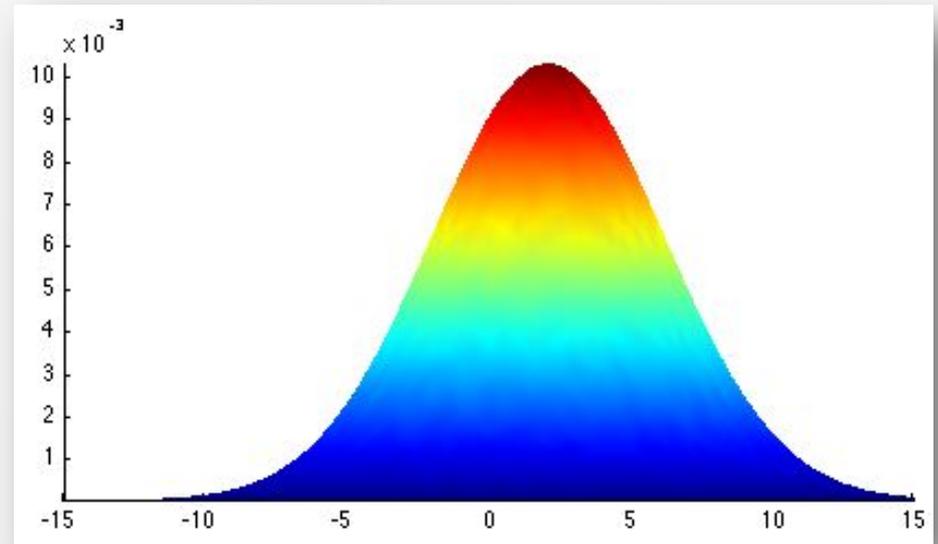


- If we *assume* the 1-dimensional data in **this histogram** are normally distributed, we can fit a continuous Gaussian function simply in terms of the mean μ and variance σ^2 .

(Aside) Univariate (1D) Gaussians

- Also known as **Normal** distributions, $N(\mu, \sigma)$

- $$P(x; \mu, \sigma) = \frac{\exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)}{\sqrt{2\pi}\sigma}$$



- The parameters we can modify are $\theta = \langle \mu, \sigma^2 \rangle$
 - $\mu = E(x) = \int x \cdot P(x)dx$ (**mean**)
 - $\sigma^2 = E((x - \mu)^2) = \int (x - \mu)^2 P(x)dx$ (**variance**)

But we don't have samples for all x ...

Maximum likelihood estimation

- Given data $X = \{x_1, x_2, \dots, x_n\}$, MLE produces an estimate of the parameters $\hat{\theta}$ by maximizing the **likelihood**, $L(X, \theta)$:

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} L(X, \theta)$$

where $L(X, \theta) = P(X; \theta) = \prod_{i=1}^n P(x_i; \theta)$.

- Since $L(X, \theta)$ provides a **surface** over all θ , in order to find the **highest likelihood**, we look at the derivative

$$\frac{\delta}{\delta\theta} L(X, \theta) = 0$$

to see **at which point** the likelihood **stops growing**.

MLE with univariate Gaussians

- Estimate μ :

$$L(X, \mu) = P(X; \mu) = \prod_{i=1}^n P(x_i; \theta) = \prod_{i=1}^n \frac{\exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right)}{\sqrt{2\pi}\sigma}$$

$$\log L(X, \mu) = -\frac{\sum_i (x_i - \mu)^2}{2\sigma^2} - n \log(\sqrt{2\pi}\sigma)$$

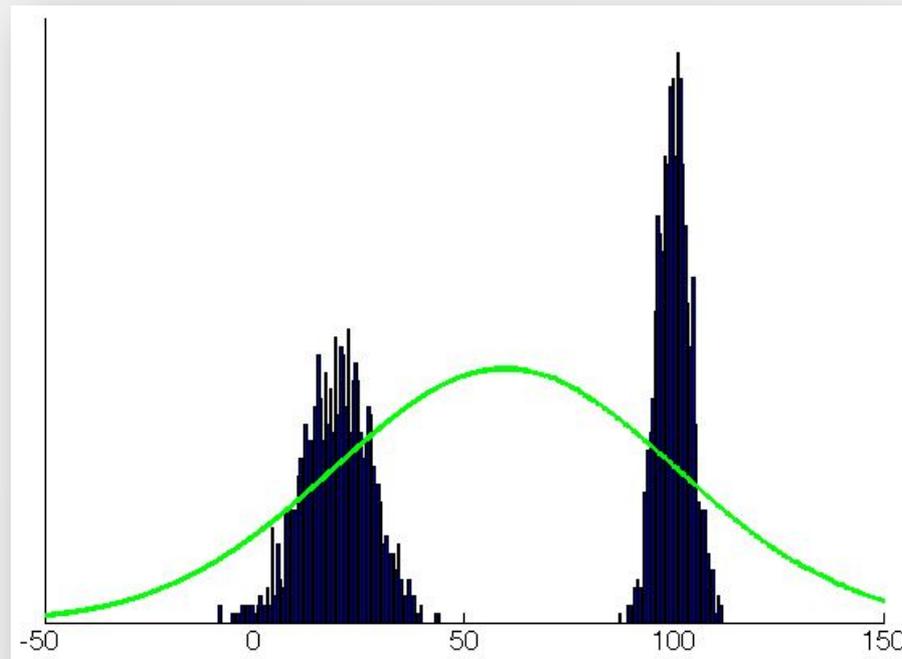
$$\frac{\delta}{\delta\mu} \log L(X, \mu) = \frac{\sum_i (x_i - \mu)}{\sigma^2} = 0$$

$$\mu = \frac{\sum_i x_i}{n}$$

- Similarly, $\sigma^2 = \frac{\sum_i (x_i - \mu)^2}{n}$

Non-Gaussian observations

- Speech data are generally *not* unimodal.
- The observations below are **bimodal**, so fitting one Gaussian would not be representative.



Multivariate Gaussians

- When data is **d -dimensional**, the input variable is

$$\vec{x} = \langle x[1], x[2], \dots, x[d] \rangle$$

the **mean** is

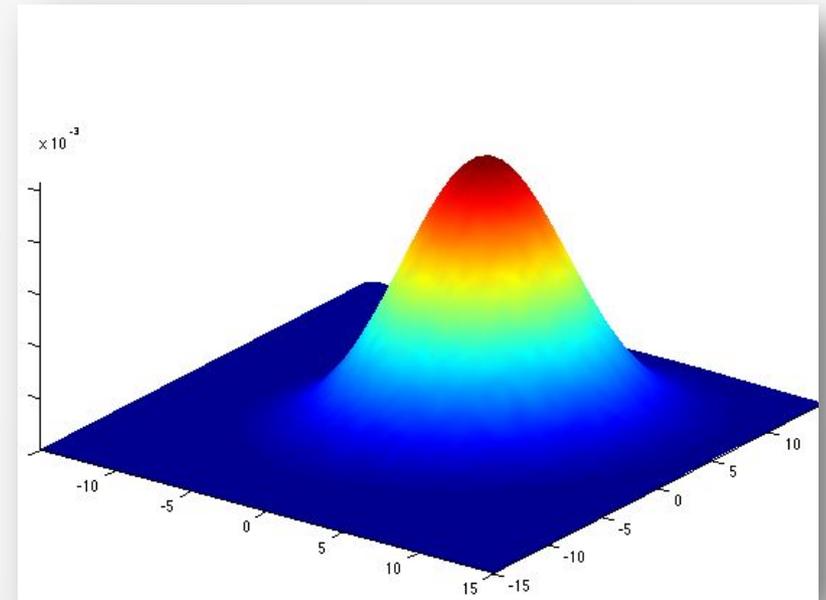
$$\vec{\mu} = E(\vec{x}) = \langle \mu[1], \mu[2], \dots, \mu[d] \rangle$$

the **covariance matrix** is

$$\Sigma[i, j] = E(x[i]x[j]) - \mu[i]\mu[j]$$

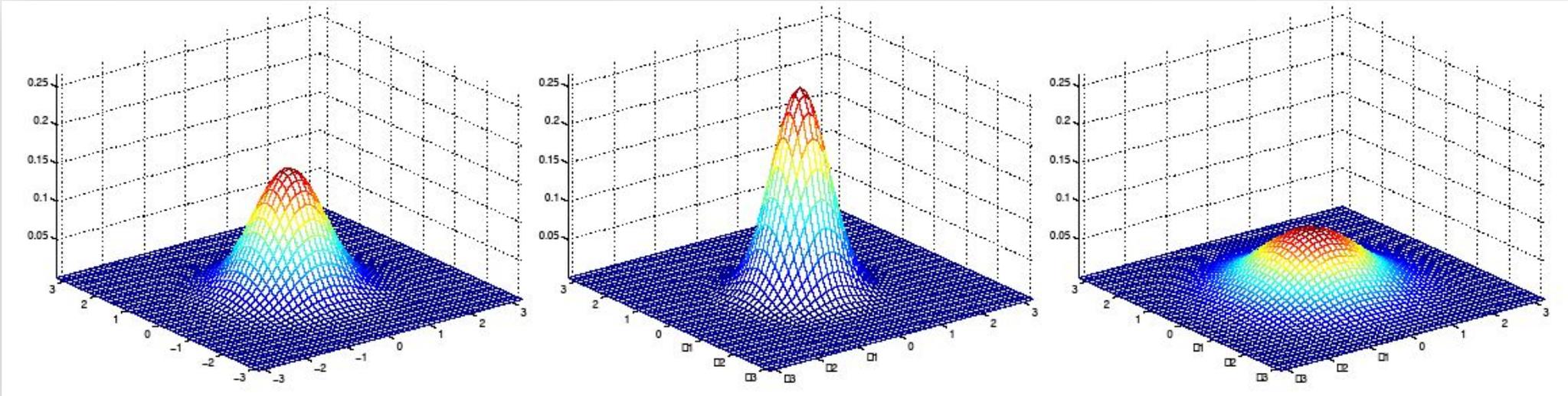
and

$$P(\vec{x}) = \frac{\exp\left(-\frac{(\vec{x} - \vec{\mu})^T \Sigma^{-1} (\vec{x} - \vec{\mu})}{2}\right)}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}}$$



A^T is the **transpose** of A
 A^{-1} is the **inverse** of A
 $|A|$ is the **determinant** of A

Intuitions of covariance



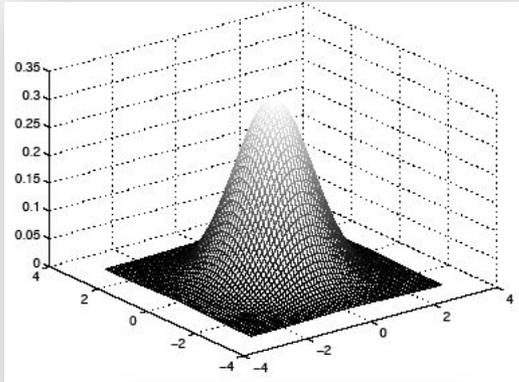
$$\mu = [0 \ 0]$$
$$\Sigma = I$$

$$\mu = [0 \ 0]$$
$$\Sigma = 0.6I$$

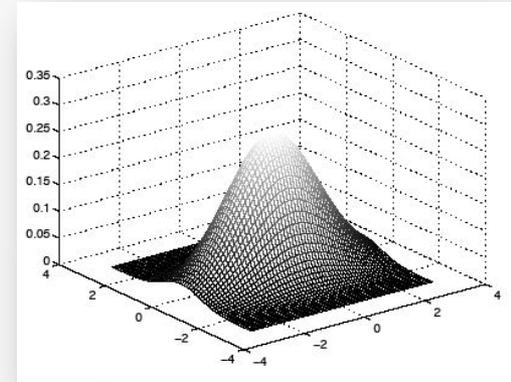
$$\mu = [0 \ 0]$$
$$\Sigma = 2.0I$$

- As values in Σ become larger, the Gaussian spreads out.
- (I is the identity matrix)

Intuitions of covariance



$$\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

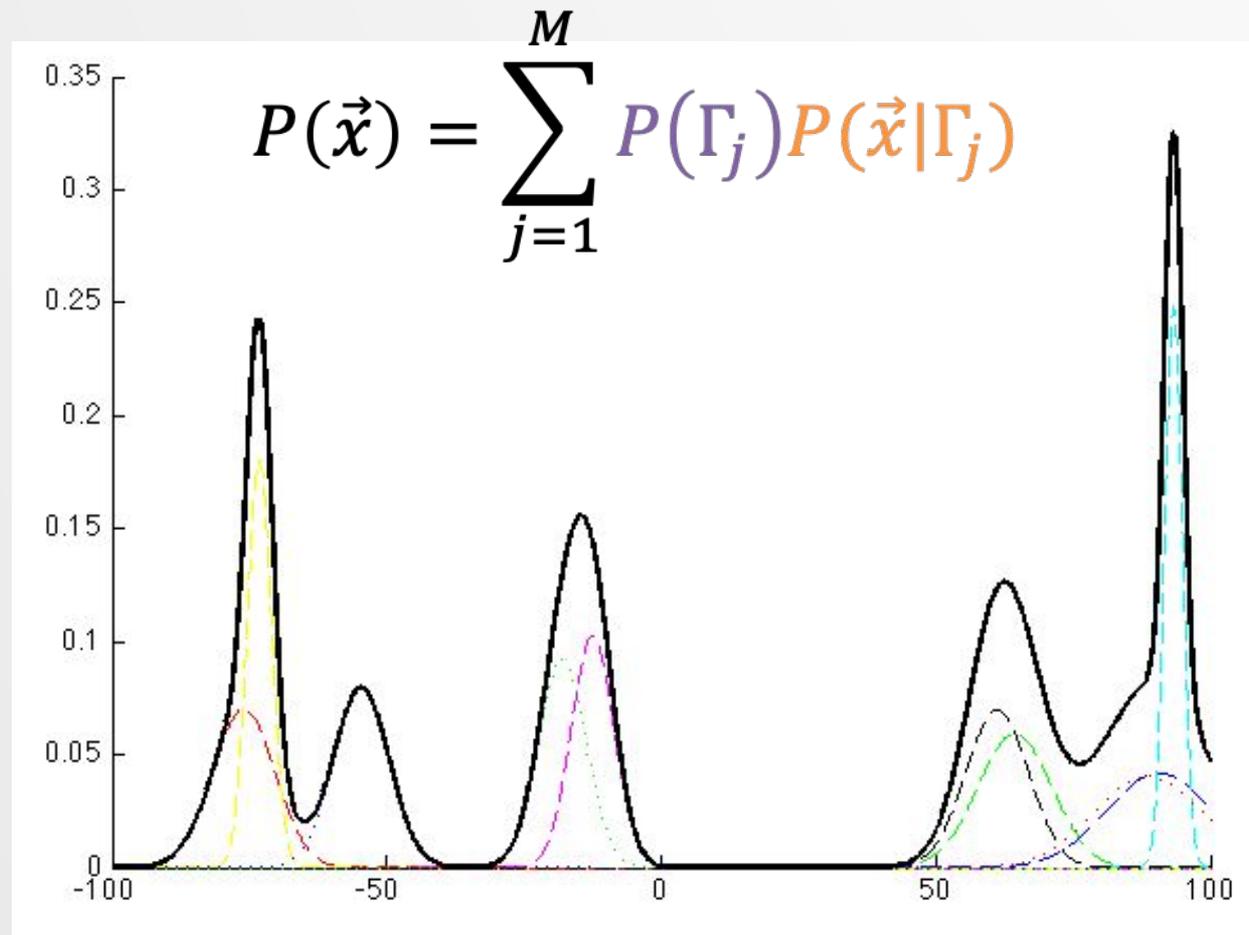


$$\Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 0.6 \end{bmatrix}$$

- Different values on the diagonal result in different variances in their respective dimensions

Mixtures of Gaussians

- **Gaussian mixture models (GMMs)** are a **weighted** linear combination of M component Gaussians, $\langle \Gamma_1, \Gamma_2, \dots, \Gamma_M \rangle$:



Observation likelihoods

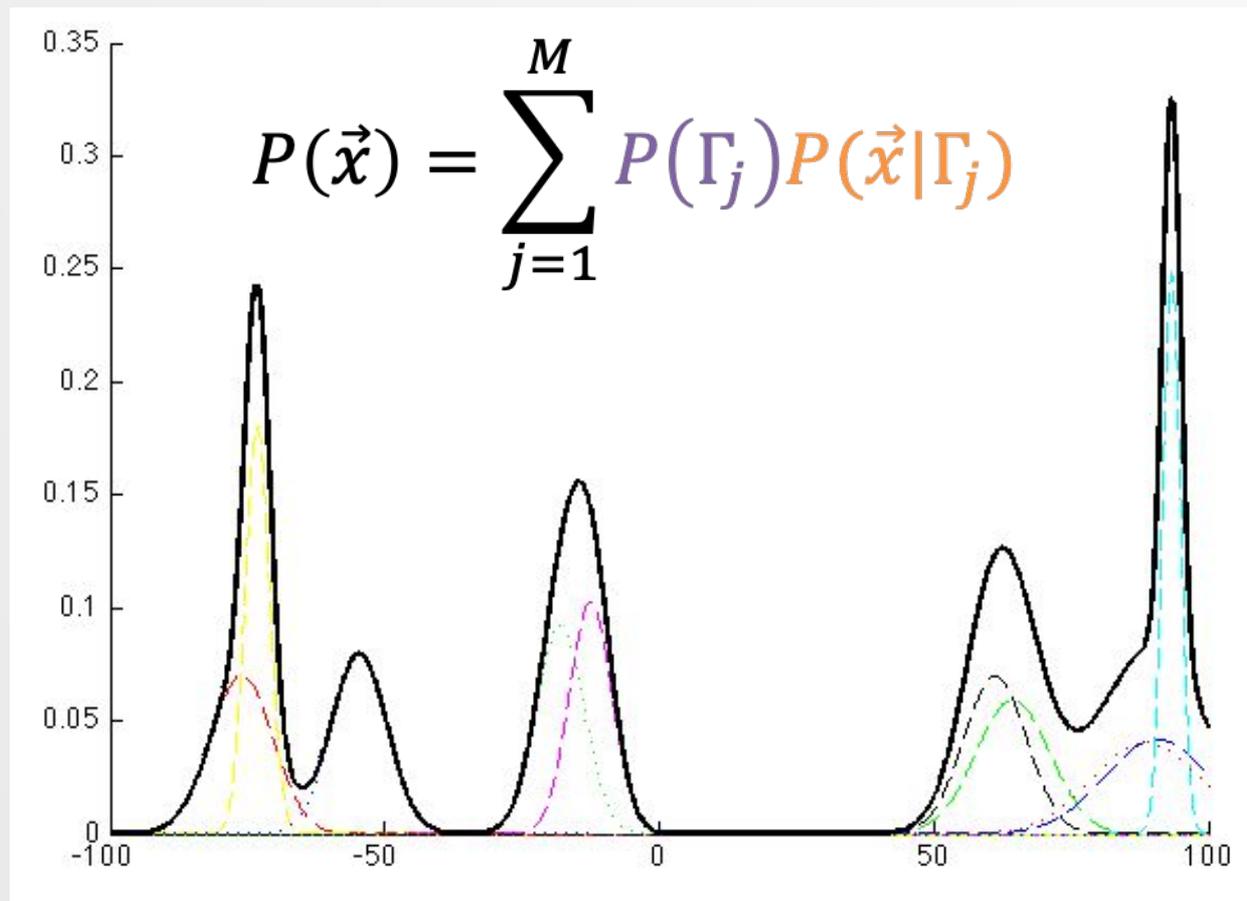
- Assuming MFCC dimensions are independent of one another, the **covariance matrix is diagonal** – i.e., 0 off the diagonal.
- Therefore, the probability of an observation vector given a Gaussian becomes

$$P(\vec{x}|\Gamma_m) = \frac{\exp\left(-\frac{1}{2}\sum_{i=1}^d \frac{(x[i] - \mu_m[i])^2}{\Sigma_m[i]}\right)}{(2\pi)^{\frac{d}{2}} \left(\prod_{i=1}^d \Sigma_m[i]\right)^{\frac{1}{2}}}$$

- *Imagine* that a GMM first *chooses a Gaussian*, then *emits an observation* from that Gaussian.

Mixtures of Gaussians

- If we knew *which* Gaussian generated each sample (which we don't), then $\langle \vec{\mu}_m, \Sigma_m \rangle$ could be learned by MLE.
- We must learn $P(\Gamma_j)$ as well.



Expectation-Maximization for GMMs

- Overall idea:
 - First, initialize a set of model parameters.
 - “Expectation”: Compute the expected probabilities of observation, given these parameters.
 - “Maximization”: Update the parameters to maximize the aforementioned probabilities.
 - Repeat.
- Let’s look at the detailed steps in the next a few slides...

Expectation-Maximization for GMMs

- Let $\omega_m = P(\Gamma_m)$ and $b_m(\vec{x}_t) = P(\vec{x}_t | \Gamma_m)$,

'weight'

'component observation likelihood'

$$P_{\theta}(\vec{x}_t) = \sum_{m=1}^M \omega_m b_m(\vec{x}_t)$$

where $\theta = \langle \omega_m, \vec{\mu}_m, \Sigma_m \rangle$ for $m = 1..M$

- To estimate θ , we solve $\nabla_{\theta} \log L(X, \theta) = 0$ where

$$\log L(X, \theta) = \sum_{t=1}^T \log P_{\theta}(\vec{x}_t) = \sum_{t=1}^T \log \sum_{m=1}^M \omega_m b_m(\vec{x}_t)$$

Expectation-Maximization for GMMs

- We **differentiate** the log likelihood function w.r.t . $\mu_m[n]$ and set this to 0 to find the value of $\mu_m[n]$ at which the likelihood stops growing.

$$\frac{\delta \log L(X, \theta)}{\delta \mu_m[n]} = \sum_{t=1}^T \frac{1}{P_{\theta}(\vec{x}_t)} \left[\frac{\delta}{\delta \mu_m[n]} \omega_m b_m(\vec{x}_t) \right] = 0$$

Expectation-Maximization for GMMs

- The **expectation step** gives us:

$$b_m(\vec{x}_t) = P(\vec{x}_t | \Gamma_m)$$

$$P(\Gamma_m | \vec{x}_t; \theta) = \frac{\omega_m b_m(\vec{x}_t)}{P_\theta(\vec{x}_t)}$$

Proportion of overall probability contributed by m

- The **maximization step** gives us:

~ “number of points explained by m ”

$$\widehat{\mu}_m = \frac{\sum_t P(\Gamma_m | \vec{x}_t; \theta) \vec{x}_t}{\sum_t P(\Gamma_m | \vec{x}_t; \theta)}$$

$$\widehat{\Sigma}_m = \frac{\sum_t P(\Gamma_m | \vec{x}_t; \theta) \vec{x}_t^2}{\sum_t P(\Gamma_m | \vec{x}_t; \theta)} - \widehat{\mu}_m^2$$

$$\widehat{\omega}_m = \frac{1}{T} \sum_{t=1}^T P(\Gamma_m | \vec{x}_t; \theta)$$

Recall from slide 13, MLE wants:

$$\mu = \frac{\sum_i x_i}{n}$$

$$\sigma^2 = \frac{\sum_i (x_i - \mu)^2}{n}$$

Some notes...

- In the previous slide, the square of a vector, \vec{a}^2 , is elementwise (i.e., `numpy.multiply`)
 - E.g., $[2, 3, 4]^2 = [4, 9, 16]$
- Since Σ is diagonal, it can be represented as a vector.
- Can $\widehat{\sigma}_m^2 = \frac{\sum_t P(\Gamma_m | \vec{x}_t; \theta) \vec{x}_t^2}{\sum_t P(\Gamma_m | \vec{x}_t; \theta)} - \widehat{\mu}_m^2$ become negative?
 - No.
 - This is left as an exercise, but only if you're interested.

Speaker recognition

- **Speaker recognition:** n . the identification of a speaker among several speakers given only acoustics.
- Each **speaker** will produce speech according to **different** probability distributions.
 - We train a **Gaussian mixture model for each speaker**, given annotated data (mapping utterances to speakers).
 - We choose the speaker whose model gives the highest probability for an observation.



Recipe for GMM EM

- For each speaker, we learn a GMM given all T frames of their training data.

- 1. Initialize:** Guess $\theta = \langle \omega_m, \vec{\mu}_m, \Sigma_m \rangle$ for $m = 1..M$ either uniformly, randomly, or by k -means clustering.
- 2. E-step:** Compute $b_m(\vec{x}_t)$ and $P(\Gamma_m | \vec{x}_t; \theta)$.
- 3. M-step:** Update parameters for $\langle \omega_m, \vec{\mu}_m, \Sigma_m \rangle$ as described on slide 28.