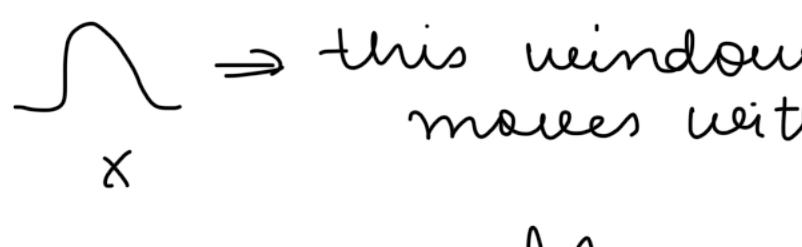


Speech Tech

Week-2

L1 Speech Product", Percep" & Freq. Analysis

- speech product model
 - ↳ terminal analogous model
 - ↳ it is similar to the LTI systems, it interacts with input.)
 - ↳ linear time invariant filter.
- STFT (using the window of small size)
 - in small window, the frequency is quasi-stationary
 - ↳ this is an assumption.

 \Rightarrow this window moves with some overlap.

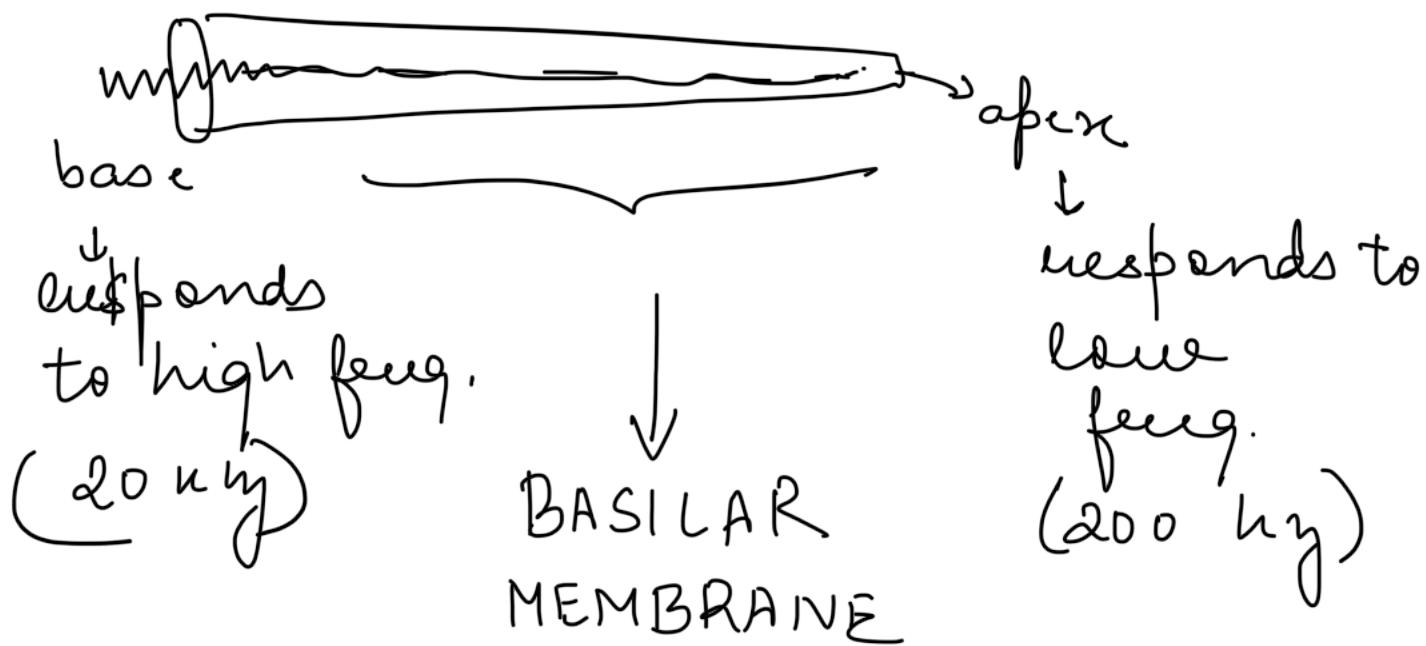


\Downarrow \Rightarrow it will be we get the Fourier Transform $\xrightarrow{\text{DFT}}$.

- wideband spectrogram
 - small window
 - good time resolvⁿ
 - resonant freq.
envelope of the spectrum
- narrowband spectrogram
 - large
 - poor
 - pitch harmonics
- In speech, window len = 25 ms
overlap = 10 ms
- Speech perceptⁿ → actual & perceived speech attributes, mel scale, masking
- Physical attribute of a acoustic signal
 - intensity
 - freq.
- Psycho physical Observaⁿ
 - auditory sys.
 - loudness
 - pitch.
- pitch & freq.
 - perceptual counterpart of freq.
 - units is mel
 - pitch of a 1000 hz is taken as ¹⁰⁰⁰ mels.

→ obtained by the freq. as perceived by the listener as the signal freq. keep varying.

- the basilar membrane for over body acts as the frequency analyzer

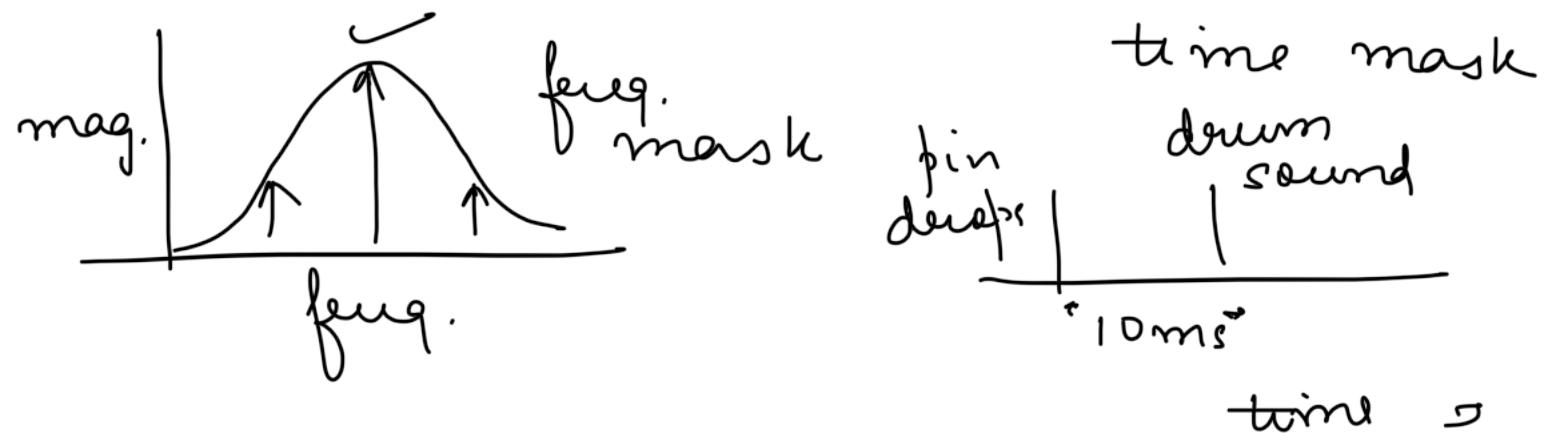


2 Perceptual Masking, Cepstrum & Filtering

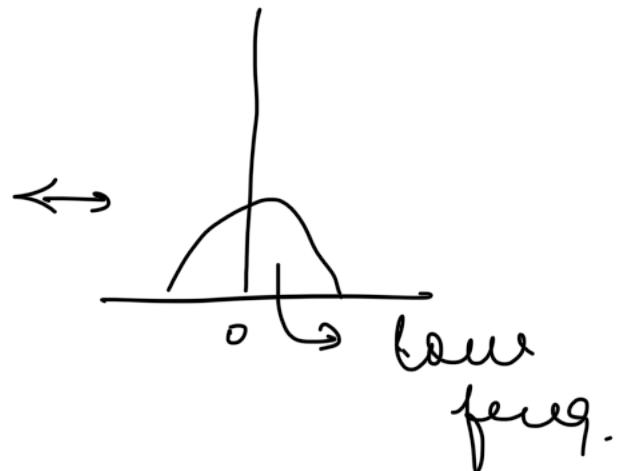
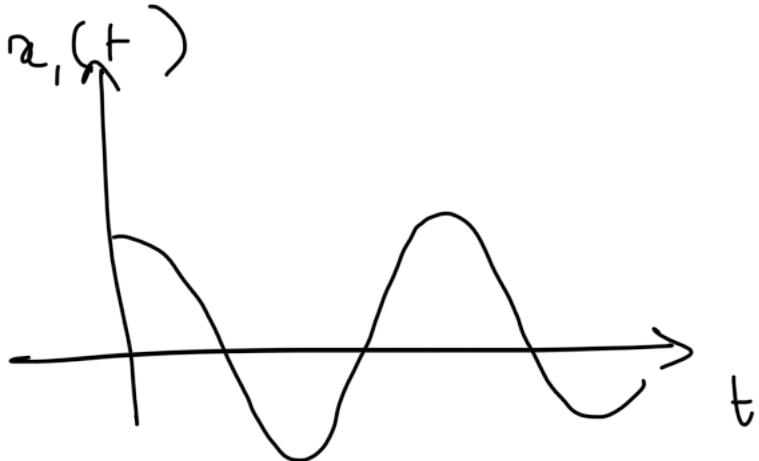
- MEL scale - quasi-log funcⁿ of acoustic frequency s.t. perceptually similar pitch intervals appear equal in width while over the full hearing range.

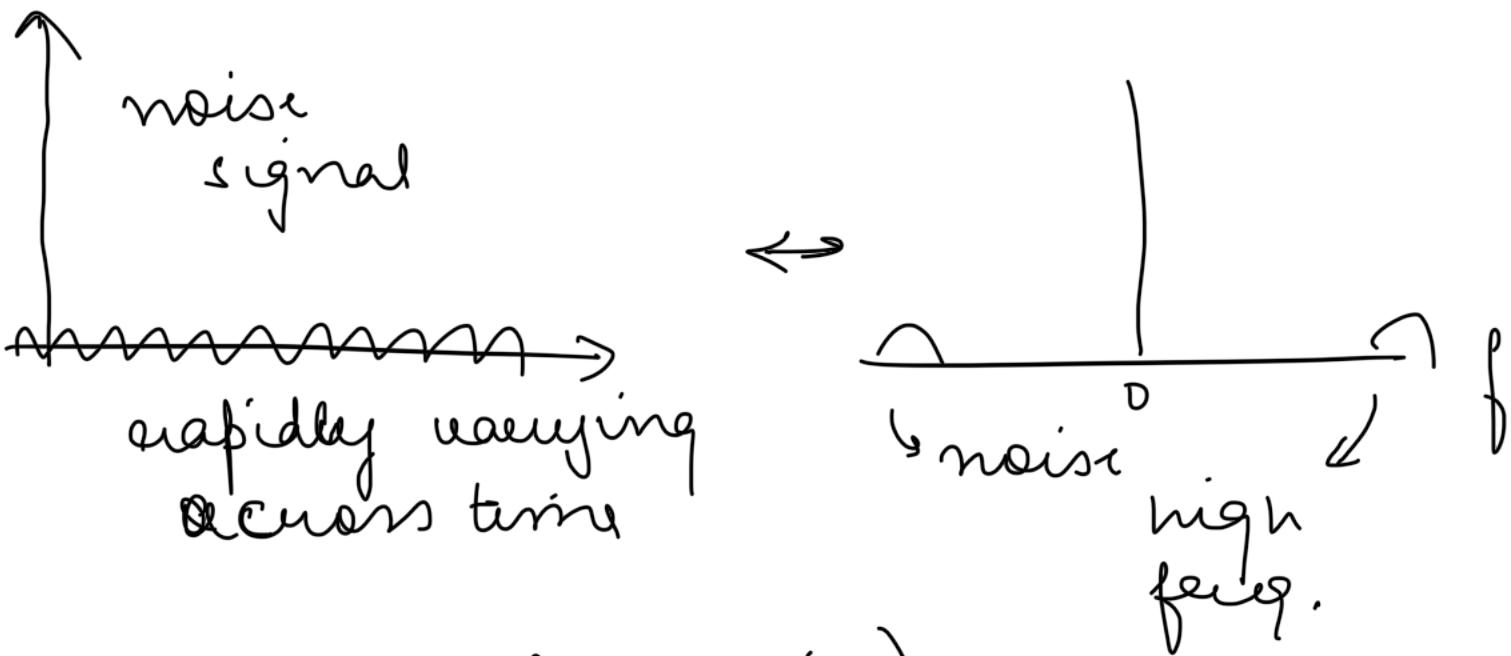
$$MEL = 1127 \log_e \left(1 + \frac{F}{700} \right)$$

- masking of sounds in v of other sounds.

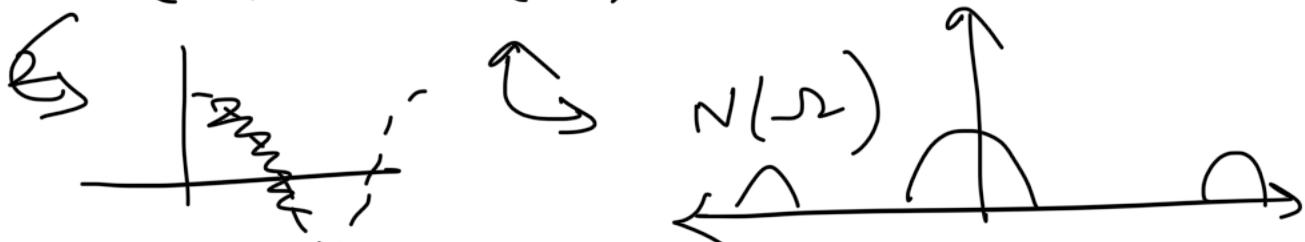


- cepstrum
 - recall speech producer model
 - pitch are the vibrations
 - speech comes as a result of convolution of vocal tract response with glottis pulses / turbulent air.
- We have to sep. out the vocal tract response -
 - yes, can be done.
 - no, only filtering doesn't work
 - if 2 signals are added, then filtering needs.
 - what abt. convolved signals?





$$z(t) = x(t) + n(t)$$



from a given $z(t)$, how to get back $x(t)$?



$$\hat{x}(t) = \frac{1}{2T_0} \int_{-T_0}^{T_0} z(t-\tau) d\tau$$

Practical implementation $X(\omega) = z(\omega) \cdot H(\omega)$

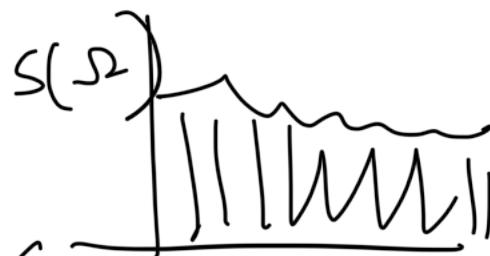
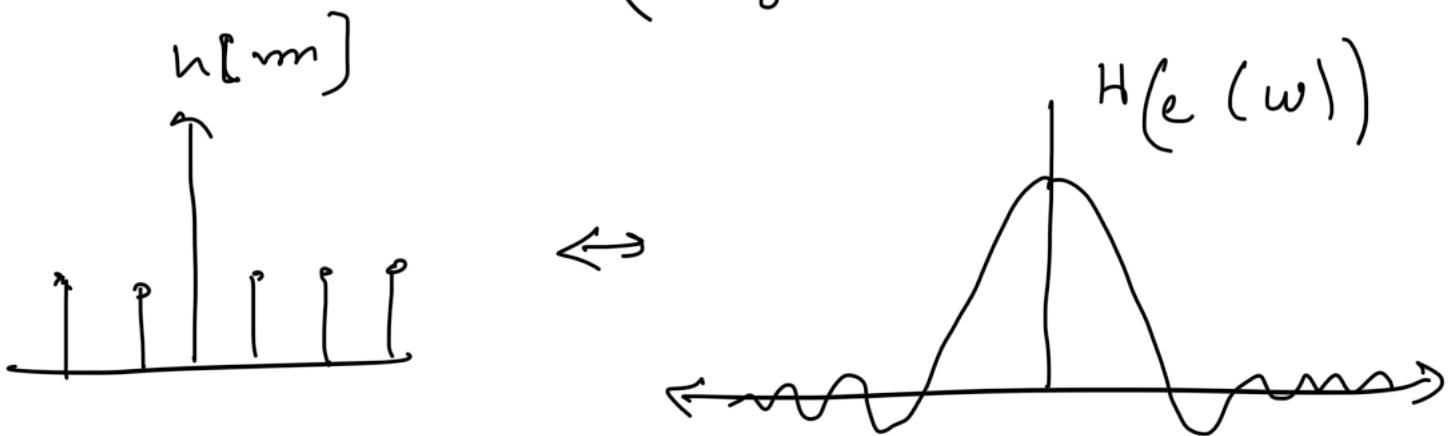
low pass filter

The received $\hat{x}(t) \approx x(t)$.

$$\hat{x}[n] = \frac{1}{2M+1} \sum_{m=-M}^M x[n-m]$$

$$= \sum_n h[m] x(n-m)$$

$$h[m] = \begin{cases} \frac{1}{2M+1}, & -M \leq m \leq M \\ 0, & \text{otherwise} \end{cases}$$



Cepstrum
Analysis

\downarrow Fourier transform of speech signal
we are recovering $x(n)$ from here.

$$\log |S(j\omega)| = \log |V(\omega)| + \log |\tilde{E}(j\omega)|$$

→ take out $\log |V(\omega)|$

yes, by avg. or low-pass filter (LIFT-ing)

L3 Feature Extract" for Speech Processing - I

$$|s(\omega)| = V(\omega) E(\omega)$$

analogous

$$y[n] = h[n] * x[n]$$

↳ low pass filter

- "convolve" is "multiplication" in freq. dom.
- Applying lag in freq. domain makes the multiplication as addition.
- slowly varying vocal-response can then be separated from a rapidly varying pitch harmonics (instead of filtering, it becomes lifting)
- now low pass filtering can extract the spectrum related vocal tract response alone. We use a rectangular moving avg. filter.
- feature extract" → filter bank, MFCC
- to obtain MFCC - a triangular moving avg. filters are used on DFT. They are of not same size & width. Inspired by mel scale

Speech → DFT → Take lag → low pass lifting.

- Mel Filter bank - features
- A DCT-2 is performed to decorrelate the features to obtain MFCC.
- basis functions are type of cosine.

$$\Delta \text{MFCC}_m[n] = \text{MFCC}_m[n] - \text{MFCC}_{m-1}[n]$$

↳ Mel Filter

Speech



Pre-Emphasis



Windowing

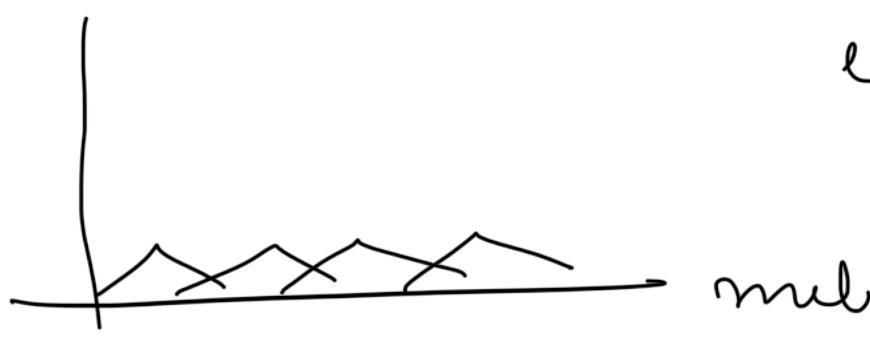


FFT



Mel Filter Bank

MFCC + Delta +
Delta Delta
↑
DCT (MFCC)
↑
deg



equal spaced
equal B.W.



$$f_{\text{mel}} = 1125 \ln\left(1 + \frac{f}{f_{00}}\right)$$

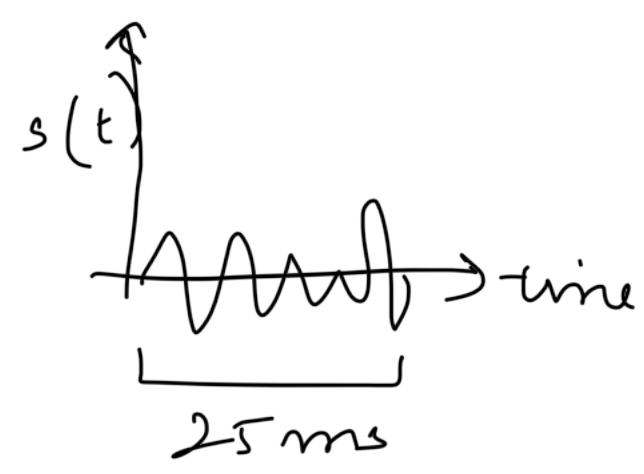
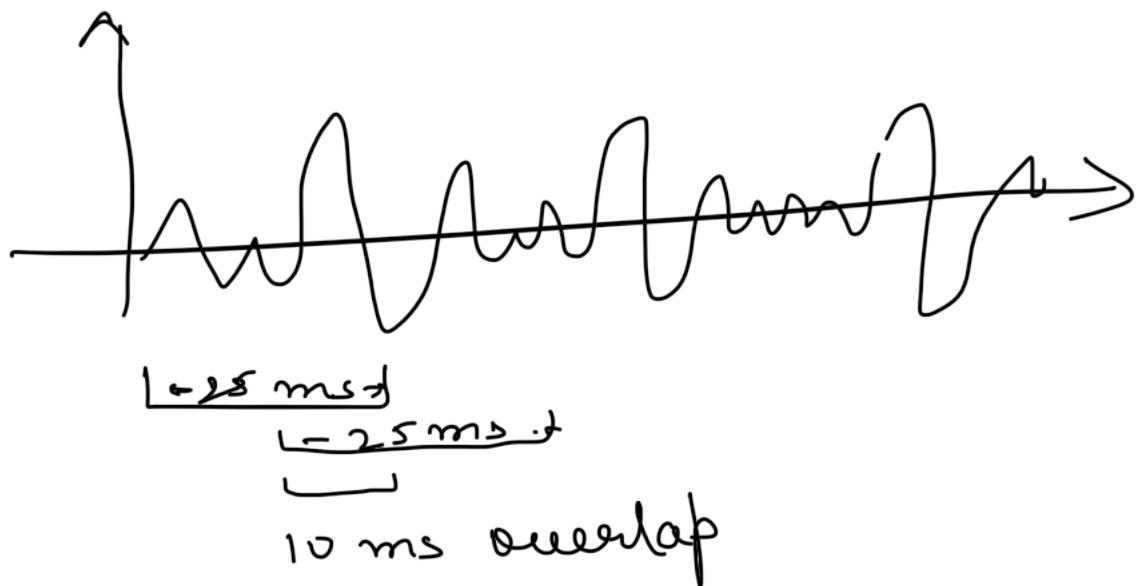
↓

leg \Rightarrow in the magnitude direct leg
mel filter bank coeff.

40 mel filters \Rightarrow telephone (8 kHz)

80 mel filters \Rightarrow wide band speech (16 kHz)

↪ in DL generally this is used.



→ freq.
domain
DFT
(FFT)



↓ Triangular
filter aug.



Log



it will be the feature vector []
for the first 25 ms. 80x1

⇒ seq. of features, they follow
time ordering

- GMM ⇒ often we assume diagonal
covar. matrices ↓ uncorrelated
 - for speech, log mel filter bank
↓ DCT-2 (discrete cosine transform - 2)
mel freq, cepstral coefficient
- we retain only first 13 cepstral coeff.

Δ \Rightarrow diff. btw. adjacent MFCC

$\Delta\Delta$ \Rightarrow diff. of diff. btw. adjacent MFCC.

↓
concatenate to form the
MFCC vectors.

L5 MFCC

1. Impt. of speech \rightarrow speaker & listener
2. phonemens \Rightarrow spoken words
3. freq. vs perceived mel

$$f_{\text{mel}} = 1125 \ln \left(1 + \frac{f_{\text{hz}}}{700} \right)$$

intensity vs loudness

masking

4. processing the signal

features: mel filter bank

mel frequency cepstral coeff.
(MFCC)

excitation



$$|s(f)|$$

$$= v(f) \times E(f)$$

- we use 25ms, so small period of time,

that we assume that the speech signal will be stationary during this time.

- the seq. of mel filter bank features. Almost we have 100 such $[]_{80 \times 1}$ frames per second.

L2.6 Gaussian Review

$$f_x(x; \mu, \sigma) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

\downarrow mean \downarrow SD

↳ Gaussian distribution (bell-shaped curve)

- $(x - \mu)$ some kind of shift operation is happening.
- we use all data to calculate μ & σ of gaussian distri. Generally, some estimation techniques are used.
use MLE (maximum likelihood estimation)

$$\mu = \frac{1}{1000} \sum_{i=1}^{1000} x_i \quad \sigma = \sqrt{\frac{1}{1000} \sum_{i=1}^{1000} (x_i - \mu)^2}$$

- likelihood

$$L(\mu, \sigma) = f_x(x_1, \dots, x_{1000}; \mu, \sigma)$$

$$\begin{matrix} \mu^*, \sigma^* \\ \hookrightarrow MLE^* \end{matrix} = \underset{\mu, \sigma}{\operatorname{argmax}} (L(\mu, \sigma))$$

$$= \underset{\mu, \sigma}{\operatorname{argmax}} \prod_{i=1}^{1000} f_x(x_i; \mu, \sigma)$$

$$= \underset{\mu, \sigma}{\operatorname{argmax}} \sum_{i=1}^{1000} \log(f_x(x_i; \mu, \sigma))$$

- we can also use the gaussian model
for classification problem.