

"Statistical models" for ASR



CDEEP
IIT Bombay

So far, we've seen "min. dist" classific.

Instead, we consider now minimizing $P(\text{error})$

EE 679 L **22** / Slide **1**

\Rightarrow min. classific. error

For a test vector \bar{x} belonging to one of the classes $\{w_i\}$. Bayes' decision rule \Rightarrow

choose the class that maximizes $P(w_i | \bar{x})$

$$P(\text{error}) = \int_{-\infty}^{\infty} P(\text{error} | x) \cdot P(x) dx$$





CDEEP
IIT Bombay

EE 679 L 22 / Slide 2

By Bayes Theorem :

$$P(w_i | \bar{x}) = \frac{P(\bar{x} | w_i) P(w_i)}{P(\bar{x})}$$

$$\left. \begin{array}{l} P(\bar{x}) = \sum_{j=1}^N P(\bar{x} | w_j) P(w_j) \end{array} \right\} N = \# \text{ classes}$$

$P(w_i)$ = Prior probab (of class w_i)

$P(w_i | \bar{x})$ = Posterior probab (of class w_i given obs. \bar{x})

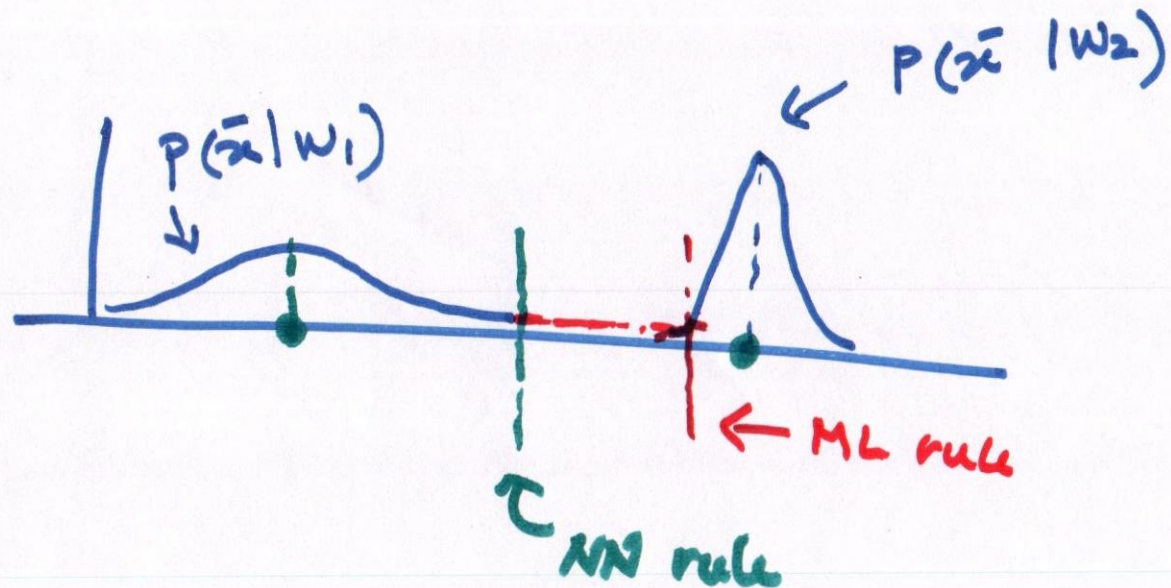
$P(\bar{x} | w_i)$ = Likelihood (cond. probab. of \bar{x} given w_i)

$P(\bar{x}) \equiv$ normalisⁿ const., does not affect decision.

Bayes Decision Rule is therefore MAP (max. a posteriori)

MAP reduces ML when $P(w_i) = \text{const.}$

Ex. using ML rule \rightarrow maximize $P(\bar{x} | w_i)$



CDEEP
IIT Bombay

EE 679 L 22 / Slide 3

Classification of a single spectral vector

One observⁿ frame at time "t"

Assume that the possible phoneme hypotheses are

$$\lambda_i \quad \lambda_1 = /i/ \quad , \quad \lambda_2 = /a/$$

Then, "best" hypothesis according to the MAP rule:

$$\hat{\lambda} = \arg \max_{\lambda_i} p(\lambda_i | O_t)$$

$$\text{We have } p(\lambda_i | O_t) = \frac{p(O_t | \lambda_i) \cdot p(\lambda_i)}{p(O_t)}$$

$$\hat{\lambda} = \arg \max_{\lambda_i} p(O_t | \lambda_i) \cdot \underbrace{p(\lambda_i)}_{\text{Lang. model}}$$



CDEEP
IIT Bombay

EE 679 L 22 / Slide 4

Gaussian probability models



CDEEP
IIT Bombay

EE 679 L 22 / Slide 5

$$p(o|\mu_i) = N(o; \mu_i; \Sigma_i)$$

\nearrow
L-dim spectral vector

$$= \frac{1}{\sqrt{(2\pi)^L |\Sigma_i|}} \exp \left[-\frac{1}{2} (o - \mu_i)^T \Sigma_i^{-1} (o - \mu_i) \right]$$

where

$$\mu_i = \frac{1}{N} \sum_{n=1}^N o_n = \mathbb{E} \{ o_n \}$$

$N = \#$ training vectors

$$\begin{aligned} \Sigma_i &= \mathbb{E} \{ (o_n - \mu_i)(o_n - \mu_i)^T \} \\ &= \frac{1}{N} \sum_{n=1}^N (o_n - \mu_i)(o_n - \mu_i)^T \end{aligned}$$



CDEEP
IIT Bombay

EE 679 L 22 / Slide 6

In computations, we use "negative log likelihood"

$$-\log(p(o | \lambda_i))$$

$$= \frac{1}{2} \log((2\pi)^L |\Sigma_i|) + 0.5 (o - \mu_i)^T \Sigma_i^{-1} (o - \mu_i)$$

Minimizing the negative log-likelihood \equiv MLE

If all class covariances ^{(Σ_i)} are assumed to be equal,

we have :

minimise $(o - \mu_i)^T \Sigma^{-1} (o - \mu_i)$

Mahalanobis
dist

Further if $\Sigma = \mathbf{I}$,

\Rightarrow minimize $(o - \mu_i)^T (o - \mu_i) \leftarrow \text{E.D.}!$

Gaussian Mixture models

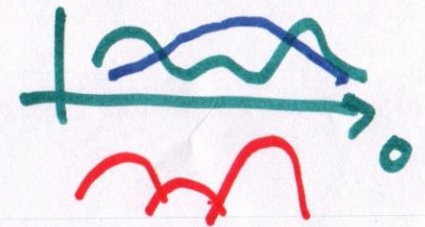


CDEEP
IIT Bombay

EE 679 L 22 / Slide 7

$$\begin{aligned} p(o | \lambda_j) &= \sum_{k=1}^K p(o | \theta_{jk}) \cdot p(\theta_{jk} | \lambda_j) \\ &= \sum_{k=1}^K c_{jk} \mathcal{N}(o; \mu_{jk}; \Sigma_{jk}) \end{aligned}$$

Handwritten annotations:
- A green arrow points from " j^{th} class" to the subscript j in $p(o | \lambda_j)$.
- A green arrow points from " k^{th} Gaussian of j^{th} class" to the subscript k in the summation.
- A green arrow points from the term $c_{jk} \mathcal{N}(o; \mu_{jk}; \Sigma_{jk})$ to the summation index k .



For large enough K , a GMM can represent any cont. prob. distribⁿ with arbitrarily good precision.

The mixture distribⁿ parameters for class j
 $\equiv \{c_{jk}; \mu_{jk}; \Sigma_{jk}\}$ are "trained" on labeled data of j .

Training a GMM via EM algo



CDEEP
IIT Bombay

EE 679 L 22 Slide 9

EM attempts to estimate the parameters of the distribⁿ (model) that maximize the $\log(P(x|\theta))$
 $\underbrace{\hspace{10em}}$ parameters of the model
 \uparrow all the training data

where $P(x|\theta) = \prod_{n=1}^N p(x_n|\theta)$, $x = \{x_1, x_2, \dots, x_N\}$

Steps:

Initialization: Set k cluster parameters $\{\mu_k, \sigma_k^2, c_k\}$ by guessing.



CDEEP
IIT Bombay

EE 679 L 22 / Slide 9

Assignment step: take x_n

$$p_{kn} = \frac{p(x_n | \mu) \cdot P(k)}{p(x_n)}$$

\equiv "degree of belonging" to cluster k

Update step: Adjust params $\{\mu_k, \sigma_k^2, c_k\}$
based on the data assignments

$$\hat{c}_k = \frac{1}{N} \sum_{n=1}^N p_{kn} ; \quad \hat{\mu}_k = \frac{\sum_N p_{kn} x_n}{\sum_N p_{kn}}, \text{ etc.}$$

Repeat