

PART IV

Automatic Speech Recognition (ASR)

ASR: What for?

- Office/desktop:
 - voice control of PC/Workstations, of programs, dictation systems
- Manufacturing/Business:
 - aid in manufacturing process, quality control, stock control and management
- Medical/Legal:
 - creation of medical/legal reports, briefs, diagnostics...
- Others:
 - games, aid to handicapped, interactive kiosk information systems

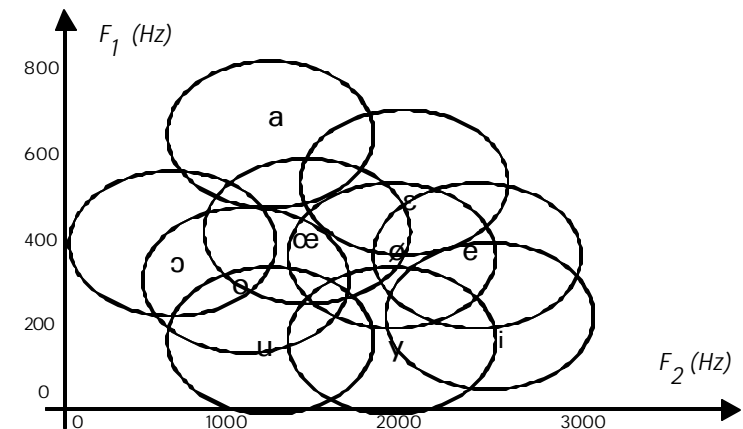
ASR: What for?

- Telecommunications:
 - access to data or services over the telephone (e.g. AT&T's Maxwell personal telephone attendant)



Challenges

- Coarticulation!



Challenges

- **Inter-speaker** variability
 - Vocal tract, gender, dialects
- **Language** variability
 - From isolated words to continuous speech
 - Out-of-vocabulary words
- **Noise**
 - Convolutional: recording/transmission conditions, reverberation
 - Additive: recording environment, transmission SNR
 - Intra-speaker variability: stress, age, humor, Lombard effect, ...

Typology of ASR systems

- **Speaker**-dependent vs. -independent
- **Language** constraints:
 - isolated word recognition
 - connected word recognition
 - keyword spotting
 - continuous speech recognition
- **Robustness** constraints
 - laboratory (office) conditions: imposed microphone, no ambient noise
 - (quiet) telephone system (Δ mic., Δ noise in a given range)
 - real-life (human-like) ASR...

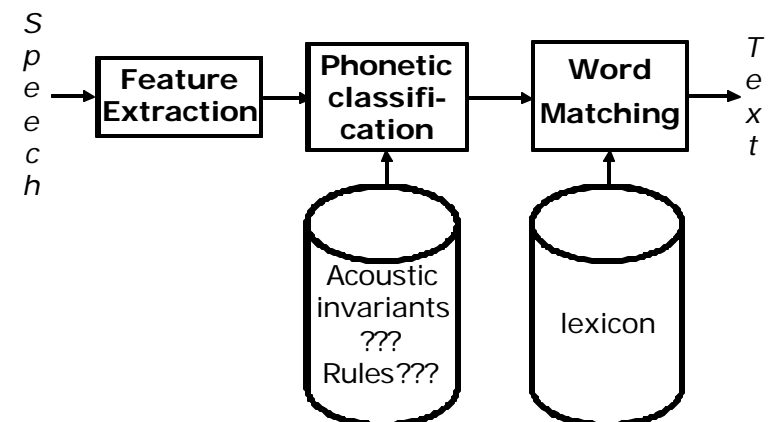
+ **vocabulary** :
small (100),
medium (5000),
large (50000)

+ **perplexity**

Levels of complexity

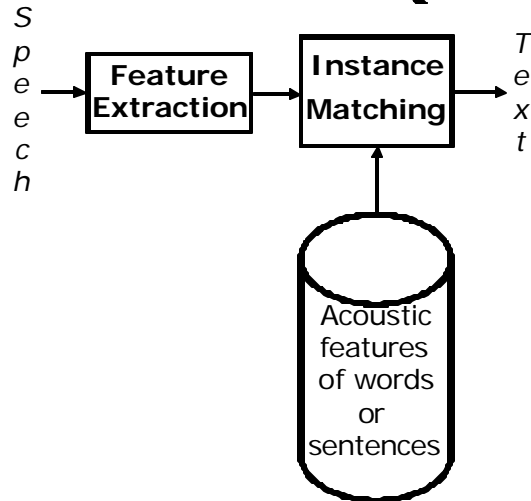
	<i>Isolated</i>		<i>Connected</i>		<i>Continuous</i>	
<i>Speaker dependent</i>	small large	1 4	small large	4 5	small large	5 6
<i>Multi speaker</i>	small large	2 4	small large	4 5	small large	6 7
<i>Speaker independent</i>	small large	3 5	small large	4 8	small large	5 10

ASR flow-chart (60 's)



The « analytical » way of doing ASR
Very poor efficiency

ASR flow-chart (70 's)

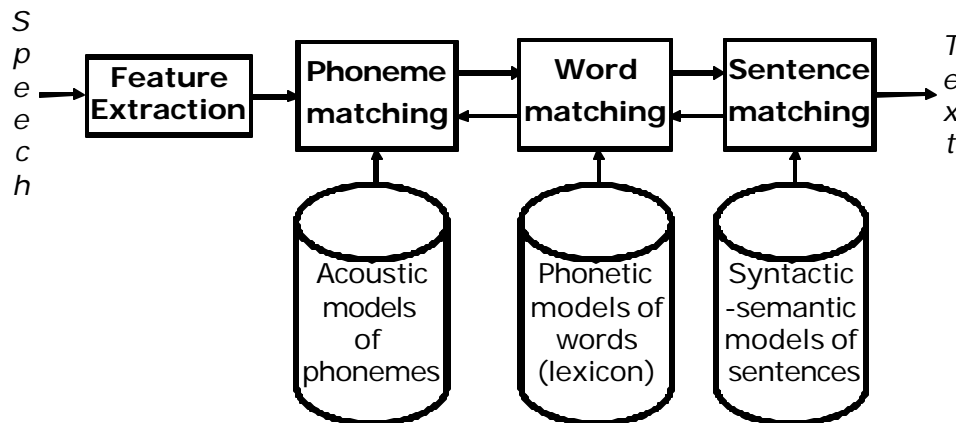


The **instance-based approach** (DTW)
OK for small vocabulary, speaker depdt

Contents

- Introduction
- **Feature extraction**
- Instance-based approach (DTW)
- Model-based approach (HMM, HMM/ANN)
 - Acoustic model
 - Phonetic model
 - Language model

Today 's ASR flow-chart

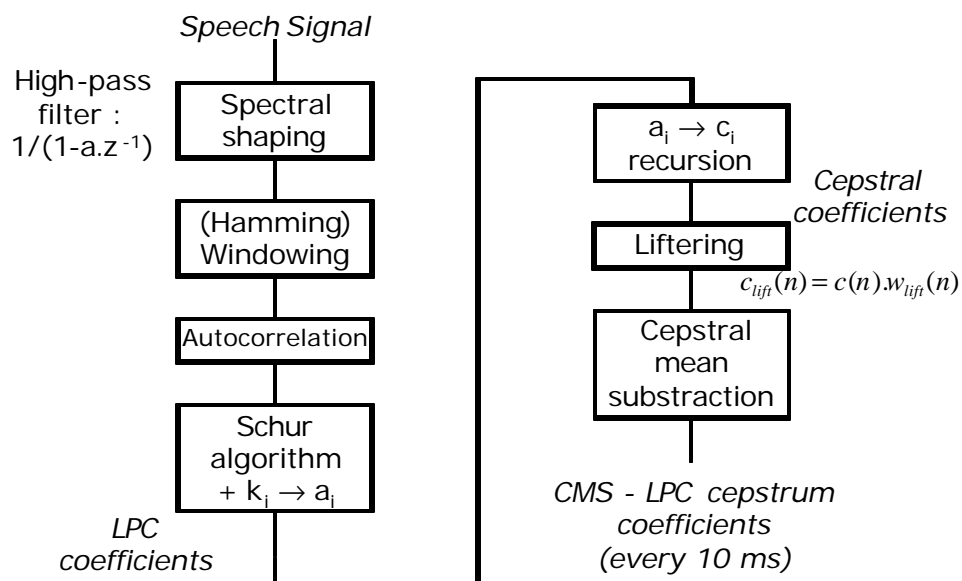


Phoneme-based approach using statistical models (HMM or HMM/ANN) for acoustics and linguistics: Large vocabulary, speaker indepd

Speech models for ASR

- Ideal properties of parameters:
 - Invariant across the speakers for the same sounds
 - Good discrimination between speech sounds
 - Robust to noise
- Types of features used in practice:
 - LPC-based features : **cepstrum coefficients**
 - Frequency warped spectral features: idem + use a non-linear frequency axis to mimic the human auditory system. (e.g. **PLP analysis, MEL - cepstrum**).
 - Auditory features : outputs of auditory models of the cochlea and auditory nerves

CMS-LPC coefficients



Enriching the feature set

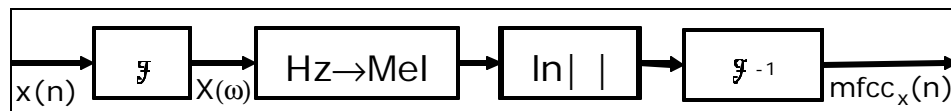
- $c_x(n)$ can be anything for silence frames
⇒ add **energy** σ_x of $x(n)$
- Add F_0 ? tried, without success
- Add **delta features** for better accounting for (frame-to-frame) spectral dynamics:

$$\Delta c_x^i(n) = c_x^i(n) - c_x^{i-1}(n)$$

$$\Delta \sigma_x^i(n) = \sigma_x^i(n) - \sigma_x^{i-1}(n)$$

MFCC - PLP

- Mel-Frequency Cepstrum Coefficients



- Perceptual Linear Prediction coefficients
 - apply frequency warping on the spectrum of $x(n)$
 - Bark scale : based on auditory critical bands (non linearly spaced in Hz; linearly spaced in Barks)
 - Apply LP model on the result

Contents

- Introduction
- Feature extraction
- **Instance-based approach (DTW)**
- Model-based approach (HMM, HMM/ANN)
 - Acoustic model
 - Phonetic model
 - Language model

Instance-based ASR

- Unknown utterance $\mathbf{X} = \{x_1, x_2, \dots, x_N\}$
(with $x_l = [x_{l1}, x_{l2}, \dots, x_{ld}]^T$)
- Known utterances $\mathbf{Y}^1 = \{y^1_1, y^1_2, \dots, y^1_{J(1)}\}$
 $\mathbf{Y}^2 = \{y^2_1, y^2_2, \dots, y^2_{J(2)}\}$
...
 $\mathbf{Y}^K = \{y^K_1, y^K_2, \dots, y^K_{J(K)}\}$
- Compute $D(\mathbf{X}, \mathbf{Y}^k)$ for $k=1 \dots M$
- Recognize:

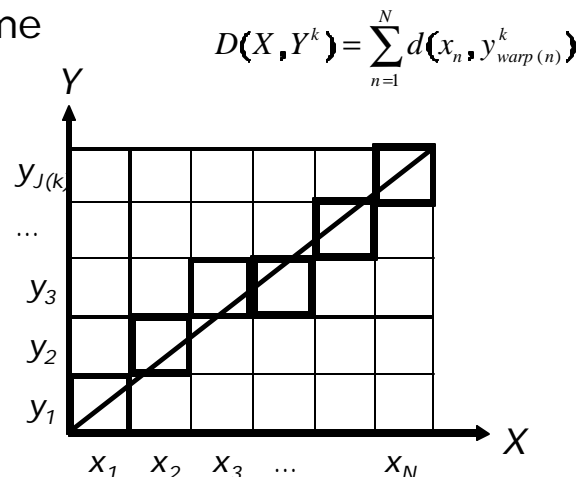
$$\mathbf{X} = \mathbf{Y}^{best}$$

with $D(\mathbf{X}, \mathbf{Y}^{best}) \leq D(\mathbf{X}, \mathbf{Y}^k)$ for $k=1 \dots M$

- OK for **spkr-dpndt isolated word reco.**

Global distance $D(\mathbf{X}, \mathbf{Y}^k)$?

- Linear time warping



- Not realistic for long words or phrases

Global distance $D(\mathbf{X}, \mathbf{Y}^k)$?

- Local distance: $d(x_n, y_j^k)$?

– Euclidian distance: $d(x_n, y_j^k) = (x_n - y_j^k)^T (x_n - y_j^k) = \sqrt{\sum_{i=1}^d (x_{ni} - y_{ji}^k)^2}$

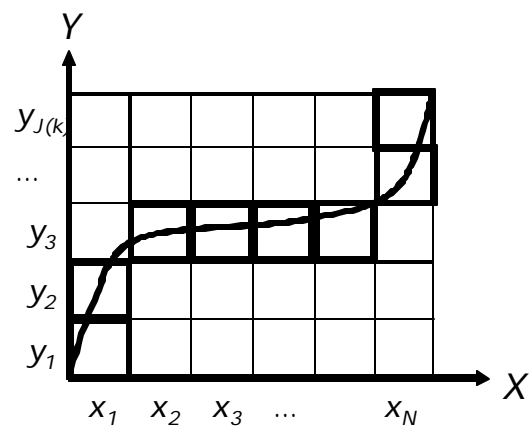
– Mahalanobis distance: $d(x_n, y_j^k) = (x_n - y_j^k)^T \Sigma^{-1} (x_n - y_j^k)$

– Itakura distance (LPC-based)

– ...

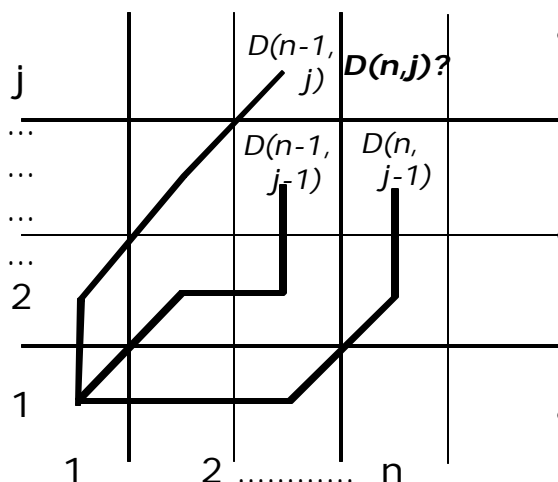
Global distance $D(\mathbf{X}, \mathbf{Y}^k)$?

- Non linear time warping



- Best path?

Dynamic Time Warping (DTW)



- $D(n, j)$ = accumulated distance from $(1, 1)$ to (n, j)
- $D(n, j) = \min$ of
 - $d(n, j) + D(n-1, j-1)$
 - $d(n, j) + D(n, j-1)$
 - $d(n, j) + D(n-1, j)$
- $D(\mathbf{X}, \mathbf{Y}^k) = D(N, J(k))$

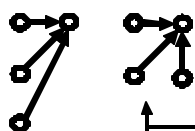
Contents

- Introduction
- Feature extraction
- Instance-based approach (DTW)
- **Model-based approach (HMM, HMM/ANN)**
 - Acoustic model
 - Phonetic model
 - Language model

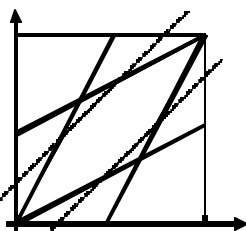
Dynamic Time Warping (DTW)

- Possible paths are constrained

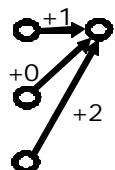
– by allowable steps; ex:



– by additional global constraints



- Penalties can be added to steps, so as to avoid too much deviation from diagonal



Model-based ASR

- Unknown utterance

$$\mathbf{X} = \{x_1, x_2, \dots, x_N\}$$

(with $x_l = [x_{l1}, x_{l2}, \dots, x_{ld}]^T$)

Models of known utterances:

$$\mathbf{M}_1, \mathbf{M}_2, \dots, \mathbf{M}_J$$

- **BAYESIAN (MAP) CLASSIFICATION:**

– Compute $P(\mathbf{M}_j | \mathbf{X})$ for $j=1 \dots J$

– Recognize:

$$\mathbf{X} = \mathbf{M}_{best}$$

with $P(\mathbf{M}_{best} | \mathbf{X}) \geq P(\mathbf{M}_j | \mathbf{X})$ for $j=1 \dots J$

Model-based ASR

- $P(\mathbf{M}_j|\mathbf{X})$ = « **posterior** probability of \mathbf{M}_j »
not easy to compute
- Bayes rule:
$$P(\mathbf{M}_j|\mathbf{X}) = \frac{P(\mathbf{X}|\mathbf{M}_j) \cdot P(\mathbf{M}_j)}{P(\mathbf{X})}$$
- $P(\mathbf{X}|\mathbf{M}_j)$ = « **likelihood** of \mathbf{X} »
- $P(\mathbf{M}_j)$ = « **prior** probability of \mathbf{M}_j »
- $P(\mathbf{X})$ = constant

$$\max P(\mathbf{M}_j|\mathbf{X}) = \max [P(\mathbf{X}|\mathbf{M}_j) \cdot P(\mathbf{M}_j)]$$

Model-based ASR

$$\max P(\mathbf{M}|\mathbf{X}) = \max [P(\mathbf{X}|\mathbf{M}) \cdot P(\mathbf{M})]$$

- \mathbf{M} is a sequence of words ($\mathbf{W}_1, \mathbf{W}_2, \dots, \mathbf{W}_K$)
- $(\mathbf{W}_1, \mathbf{W}_2, \dots, \mathbf{W}_K)$ may have several phonetic transcriptions \mathbf{P}_l ($l=1 \dots L$)
- $$P(\mathbf{X}|\mathbf{M}) = P(\mathbf{X}|\mathbf{P}_1)P(\mathbf{P}_1|\mathbf{M}) + P(\mathbf{X}|\mathbf{P}_2)P(\mathbf{P}_2|\mathbf{M}) \\ \dots + P(\mathbf{X}|\mathbf{P}_L)P(\mathbf{P}_L|\mathbf{M})$$

$P(\mathbf{X}|\mathbf{P}_l)$ ← Acoustic model

$P(\mathbf{P}_l|\mathbf{M})$ ← Phonetic model

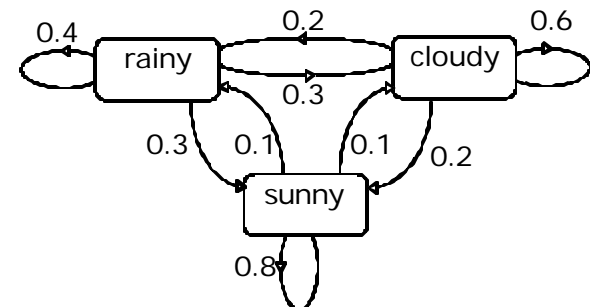
$P(\mathbf{M})$ ← Language model

Contents

- Introduction
- Feature extraction
- Instance-based approach (DTW)
- Model-based approach (HMM, HMM/ANN)
 - **Acoustic model**
 - Phonetic model
 - Language model

Markov chain

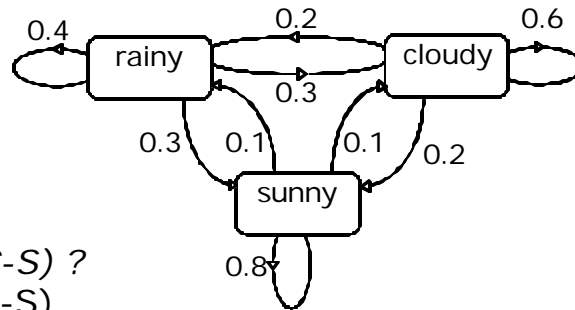
- = stochastic finite state automaton
 - ex:



- Defined by its topology (states, interconnections), and by transition probs.

Markov chain

- Use?

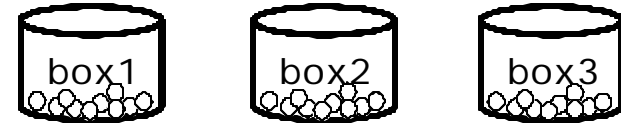


$$\begin{aligned}
 & - P(S-S-S-R-R-S-C-S) ? \\
 & P(S-S-S-R-R-S-C-S) \\
 & = P(S) \cdot P(S|S) \cdot P(S|S) \cdot P(R|S) \cdot P(R|R) \dots \\
 & = (1.0)(0.8)(0.8)(0.1)(0.4) \dots
 \end{aligned}$$

- **1 state** « **1 observation**

Try state=phoneme and observation= x_n ?
 x_n is observable with several phonemes

Hidden Markov Model (HMM)



$$\begin{aligned}
 P(r,b,r|Model) = & P(r,b,r|B1,B1,B1) P(B1,B1,B1) \\
 & + P(r,b,r|B1,B1,B2) P(B1,B1,B2) \\
 & + P(r,b,r|B1,B1,B3) P(B1,B1,B3) \\
 & + P(r,b,r|B1,B2,B1) P(B1,B2,B1) \\
 & + P(r,b,r|B1,B2,B2) P(B1,B2,B2) \\
 & \dots \\
 & + P(r,b,r|B3,B3,B3) P(B3,B3,B3)
 \end{aligned}$$

$$\begin{aligned}
 \text{with } P(r,b,r|Bi,Bj,Bk) &= P(r|Bi) P(b|Bj) P(r|Bk) \\
 P(Bi,Bj,Bk) &= P(Bi) P(Bj|Bi) P(Bk|Bj)
 \end{aligned}$$

Hidden Markov Model (HMM)



$$\begin{aligned}
 & P(r|B1) \\
 & P(b|B1) \\
 & P(g|B1)
 \end{aligned}$$



$$\begin{aligned}
 & P(r|B3) \\
 & P(b|B3) \\
 & P(g|B3)
 \end{aligned}$$



$$\begin{aligned}
 & P(r|B3) \\
 & P(b|B3) \\
 & P(g|B3)
 \end{aligned}$$

Emission probabilities

$$\begin{aligned}
 & P(B1|B1), P(B2|B1), P(B3|B1) \\
 & P(B1|B2), P(B2|B2), P(B3|B2) \\
 & P(B1|B3), P(B2|B3), P(B3|B3)
 \end{aligned}$$

Transition probabilities

Double, embedded stochastic process:

- choose box using transition probs.
- choose ball using emission probs.

Hidden Markov Model (HMM)

- **Estimation** problem : $P(\mathbf{X}|\mathbf{M})$?

$O(N_{states}^{N_{obs}})$ computations???

Dynamic programming

Baum-Welch

(all paths)

Viterbi

(best path only)

- **Training** problem: best emission and transition probs, given model topology and data
- **Decoding** problem: best sequence of states
 → **Viterbi** again

Training HMMs

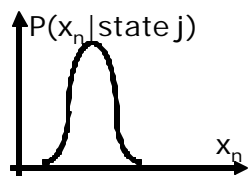


- Data: « brbrrbgbrbbgbggbrbgbgrggbrbogg... »
- Emission and transition probabilities?
If states were known: counting
- **EM** (expectation-maximization) **algorithm**
 - Initialize Probs. (first guess if possible): M^0
 - Decode the data with $M^0 \rightarrow$ states
 - Re-estimate Probs. by counting: M^1 until $M^j \approx M^{j+1}$

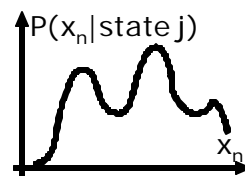
HMMs for ASR

- Observation = $x_n \Rightarrow P(x_n/\text{state}) = \textbf{continuous}$
Cannot be estimated by counting
Estimated via the p.d.f. of a distribution

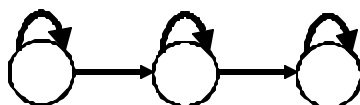
– ex: Gaussian



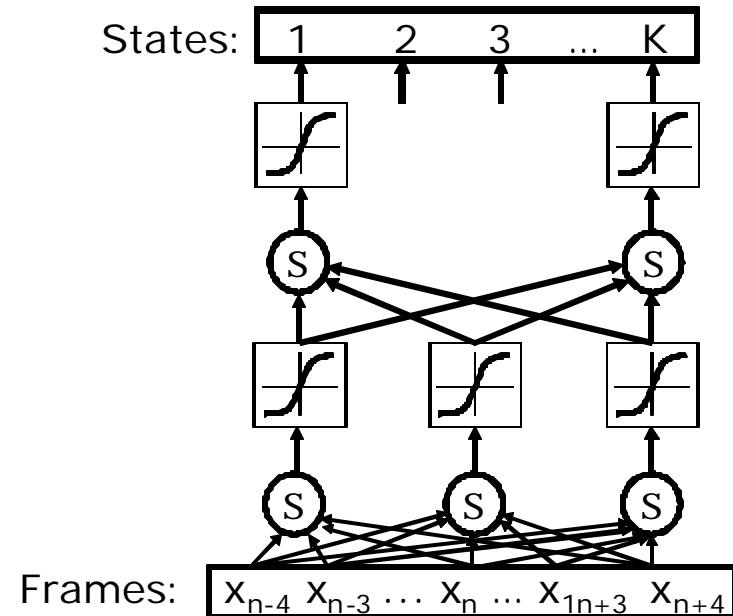
Multi-Gaussian



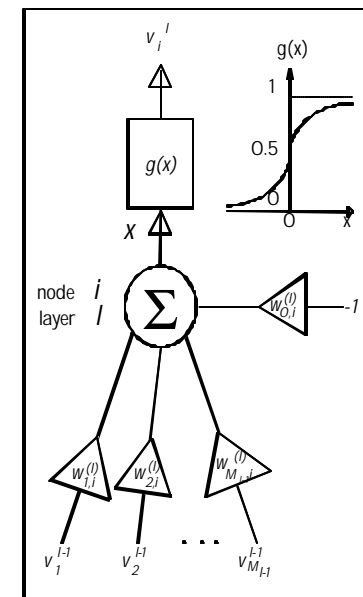
- In practice, each phoneme is modeled as 3 states (*Bakis model*)



HMM/ANN hybrids for ASR



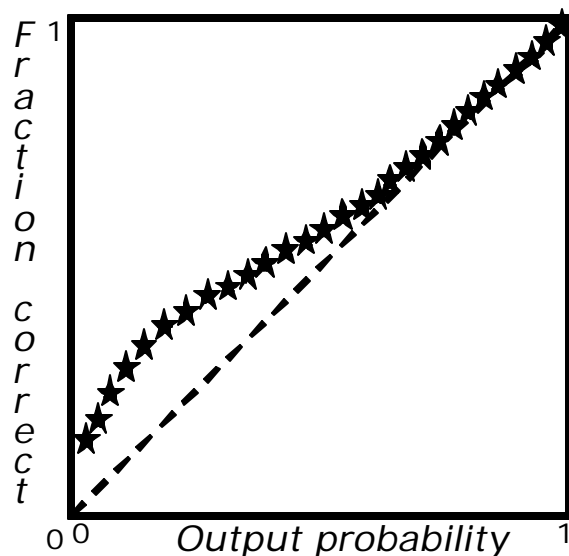
(The perceptron)



HMM/ANN hybrids for ASR

- Multi-Layer Perceptrons (MLP) can be used for estimating $P(x_n|\text{state})$

[Bourlard & Wellekens 90]



Contents

- Introduction
- Feature extraction
- Instance-based approach (DTW)
- Model-based approach (HMM, HMM/ANN)
 - Acoustic model
 - Phonetic model**
 - Language model

HMM/ANN hybrids for ASR

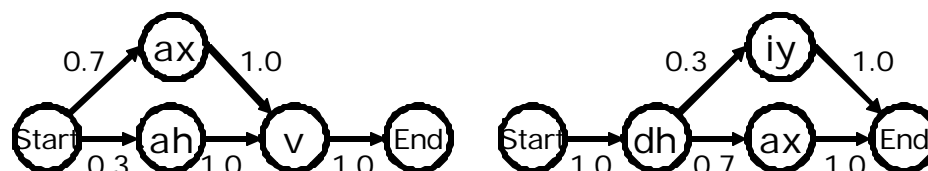
- $P(x_n|\text{state})$ is estimated **without any hypothesis on its p.d.f.** (as opposed to even multi-gaussians)
- MLPs are **discriminant**: at training time, not only the correct output (state) is maximally selected, but also the other outputs are maximally rejected
- Less parameters to train (the same MLP is used for prediction all states)

Phonetic model

- $$P(\mathbf{X}|\mathbf{M}) = P(\mathbf{X}|\mathbf{P}_1)P(\mathbf{P}_1|\mathbf{M}) + P(\mathbf{X}|\mathbf{P}_2)P(\mathbf{P}_2|\mathbf{M}) + \dots + P(\mathbf{X}|\mathbf{P}_L)P(\mathbf{P}_L|\mathbf{M})$$

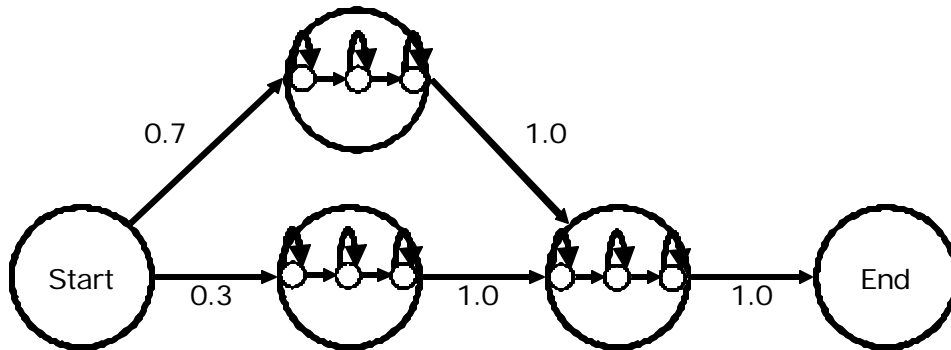
$P(\mathbf{P}_i|\mathbf{M})?$

- Using Markov chains: « of », « the »



Phonetic model

- Each phoneme is itself an HMM:



- Phonetic transition probabilities can be trained simultaneously with acoustic probs
- $P(\mathbf{X}|\mathbf{M})$ can be estimated in one shot

Contents

- Introduction
- Feature extraction
- Instance-based approach (DTW)
- Model-based approach (HMM, HMM/ANN)
 - Acoustic model
 - Phonetic model
 - Language model**

Language model

- $P(\mathbf{M})$ is actually $P(\mathbf{M}|\text{some language model})$
- 3 problems (cf. acoustic models)
 - Training** of the language model
 - Estimation** of $P(\mathbf{M}|\text{language model})$
 - Decoding**: how to integrate $P(\mathbf{M})$ in the recognition process?
- Notice that:

$$P(\mathbf{M}) = P(\mathbf{W}_1, \mathbf{W}_2, \dots, \mathbf{W}_K)$$

$$= \prod_{k=1}^K P(\mathbf{W}_k | \mathbf{W}_{k-1}, \mathbf{W}_{k-2}, \dots, \mathbf{W}_1)$$

word-pair model

- A sentence is admitted only if all word pairs have been encountered at least one in a very large text corpus:

$$P(\mathbf{M}) = 1 \quad \text{if} \quad \exists(\mathbf{W}_k, \mathbf{W}_{k-1}) \quad \text{for all } k$$

$$= 0 \quad \text{otherwise}$$

- Estimation and training are trivial
- Decoding: see n -grams
- Much too simple; banishes lots of well-formed sentences

***n*-gram models**

- **Hyp:** the probability of having a word in a sentence does not depend more on all the words of the sentence than on the n previous words:

$$P(\mathbf{M}) = \prod_{k=1}^K P(\mathbf{W}_k | \mathbf{W}_{k-1}, \mathbf{W}_{k-2}, \dots, \mathbf{W}_1)$$

$$= \prod_{k=1}^K P(\mathbf{W}_k | \mathbf{W}_{k-1}, \mathbf{W}_{k-2}, \dots, \mathbf{W}_{k-n})$$

- ex: **estimation with a trigram** ($n=2$):

$P(\text{the weather is nice}) = P(\text{the} | _, _)$

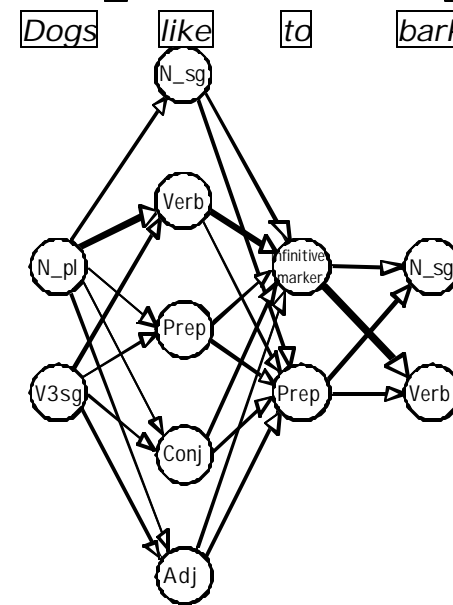
$P(\text{weather} | \text{the}, _)$ $P(\text{is} | \text{weather}, \text{the})$

$P(\text{nice} | \text{is}, \text{weather})$

***n*-gram training**

- $P(\mathbf{W}_k | \mathbf{W}_{k-1}, \dots, \mathbf{W}_{k-n})$?
- Count frequency of occurrence in a very, very, very large corpus
- If vocabulary \gg or $n > 1$, use *smoothing techniques*:
 - **Back-off**: if $P(\mathbf{W}_k | \mathbf{W}_{k-1}, \dots, \mathbf{W}_{k-n})$ cannot be estimated, try to use a combination of $n-i$ -grams ($i=1..n$)
 - Use **part-of-speech** instead of words for context :
 $P(\mathbf{W}_k | \mathbf{W}_{k-1}, \dots, \mathbf{W}_{k-n}) \gg P(\mathbf{W}_k | \text{pos}(\mathbf{W}_{k-1}), \dots, \text{pos}(\mathbf{W}_{k-n}))$
- In practice : $n \leq 2$: **trigrams**

Seeing words as pos



***n*-gram decoding**

- $\max P(\mathbf{M}_j | \mathbf{X}) = \max [P(\mathbf{X} | \mathbf{M}_j) \cdot P(\mathbf{M}_j)]$
- It has been assumed that \mathbf{M}_j was known, while \mathbf{M}_j should obviously depend on \mathbf{X} itself! (try to recognize the most probable sentence first)
- **Depth-first search**:
 - at each frame i , the decoder stores on a *stack* the list of most likely word sequences up to frame i
 - these word sequences are tested first for the estimation of the acoustic score for frame $i+1$

***n*-gram results**

"That this simple approach is so successful is a source of considerable irritation to me and to some of my colleagues. We have evidence that better language models are obtainable, we think we know many weaknesses of the trigram model, and yet, when we devise more or less subtle methods of improvement, we come up short."

F. Jelinek, « Up from trigrams », 1993

Today 's error rates

- Importance of the language model:
 - Ressource Management:
 - without LM: 85% words; with LM: 97%
 - ⇒ The last 3% might still be a *language* problem
- These are laboratory systems, working on read speech!
Real life systems, spontaneous speech: -30% :-)
- Current issues :

Robustness	Spkr adaptation	Language models
-------------------	------------------------	------------------------

Conclusion

Type	Task	Mode	Vocabulary	errorrate
Isolated words	Equiprobable words	Sp. Depdt	10 digits	0%
		Sp. Indepdt	39 ascii	4.5%
		Sp. Indepdt	1109 basic English	4.3%
		Sp. Indepdt	10 digits	0.1%
		Sp. Indepdt	39 ascii	7.0%
		Sp. Indepdt	1218 names	4.7%
Connected words	Sequence of digits id.	Sp. Depdt	10 digits	0.1%
	Flight reservation	Sp. Indepdt	11 digits	0.2%
		Sp. Depdt	129 words	0.1%
Continuous speech	Ressource management (perplexity 60)	Sp. Indepdt	991 words	3.0%
	Airline travel information system (perplexity 25)	Sp. Indepdt	1800 words	3.0%
	Wall street journal (perplexity 145)	Sp. Indepdt	20000 words	12.0%

Recognition is not understanding

