Speech Synthesis

Simon King University of Edinburgh



Speech signal analysis and modelling

- <u>analysis: generalising source + filter</u> to excitation + spectral envelope •
- •

modelling: representing speech parameters in a form suitable for statistical modelling

Orientation

- First part: speech signal **analysis**
- epochs
- F0
- spectral envelope
- <u>Second part</u>: speech signal **modelling**
 - speech parameters
- representations suitable for modelling
- converting back to a waveform



What you should already know

- <u>speech signals</u>
 - F0 & harmonics
 - vocal tract frequency response (formants)
- <u>source-filter model</u>
 - source: pulse train or noise
 - filter: a set of resonances
 - e.g., linear predictive (LP)
- Fourier analysis
 - magnitude & phase spectra



Speech signal analysis

- what we need to analyse
- epoch detection ('pitch marking')
- F0 estimation ('pitch tracking')
- spectral envelope estimation

Generalising the source-filter model concept

- Often, we don't really need the 'true' source and filter
- We just need to work with the speech **signal**, so that we can
 - measure
 - individual properties: e.g., F0 for use in the join cost
 - modify
 - phonetic identity
 - prosody
 - <u>manipulate</u>
 - waveforms: e.g., to smoothly concatenate candidates from the database

k independently



frequency

8kHz

Epoch detection vs F0 estimation

- Two different things, for different purposes
- - pitch-synchronous signal processing
 - TD-PSOLA
 - or simply just overlap-add joining of units
 - a few vocoders operate pitch synchronously
- <u>F0 estimation</u> (also known as pitch determination, **pitch tracking**, F0 tracking, ...)
 - a component of the join cost in all unit selection systems
 - used in the target cost, for systems that predict F0 targets (ASF)
 - a parameter for most (probably all) vocoders

• Epoch detection (also known as **pitch marking**, Glottal Closure Instant (GCI) detection)

Epoch detection vs F0 estimation



time

Speech signal analysis

- what we need to analyse
- <u>epoch detection ('pitch marking')</u>
- F0 estimation ('pitch tracking')
- spectral envelope estimation

What we need epoch detection for: PSOLA





Taylor - figure 12.22



A simple algorithm for epoch detection

- <u>Goal</u>
 - find a single, consistent location within each pitch period of the speech waveform
- <u>Plan</u>
 - make the problem simpler by removing all frequencies other than FO
 - find the main peak in each period
 - but, peak picking turns out to be hard
 - so, convert problem to detecting zero crossings, which is easy

A simple algorithm for epoch detection





frequency

8kHz



A simple algorithm for epoch detection

$\mathsf{MMMMMMMMMM}$



A simple algorithm for epoch detection





A simple algorithm for epoch detection

- Summary of method:
 - preprocess
 - remove unwanted frequencies with a low-pass filter
 - peak picking
 - differentiate
 - smooth, to remove spurious low-amplitude variations
 - find zero crossings (from positive to negative)
 - postprocess
 - correct for time offset e.g., to align pitchmark with largest peak in each period

Speech signal analysis

- what we need to analyse
- epoch detection ('pitch marking')
- F0 estimation ('pitch tracking')
- spectral envelope estimation

Orientation

- <u>epoch detection</u>
- pitch marking, Glottal Closure Instant (GCI) detection

- <u>F0 estimation</u>
- pitch determination, pitch tracking, F0 tracking, ...

So, can we obtain an estimate of F0 from the epochs?



for **parameterising** speech signals

าร

_____Obtaining an estimate of F0 from the detected epochs



Epoch detection vs F0 estimation



time



		IW	1114
	1	h	And



Cross-correlation (also known as "modified autocorrelation")

 $r_t(\tau)$ autocorrelation function of lag τ time index integration window size W





		hhhhh.
V		MAMMAN
	τ	





Cross-correlation, as a function of ${\cal T}$

The autocorrelation method

- We search for a peak in the (modified) autocorrelation function.
- There will be a large peak at a lag of 0, another at the pitch period and then every exact multiple of the pitch period
- Pick the highest non-zero-lag peak over **some search range**
 - the corresponding lag = the pitch period (measured in samples)
- Not always as easy at that sounds:
 - real signals are **not perfectly periodic**
 - formants will lead to some waveform self-similarity at lags other than exact multiples of the pitch period

Problems with autocorrelation: peak picking is **hard**

- We need to choose **the search range** carefully.
- If the upper limit is too high, we may choose a peak at too great a lag
- overestimate the pitch period = underestimate F0 by a factor (e.g., pitch halving) • If lower limit is too low we may choose the zero-lag peak
- Typically, pitch estimation algorithms are based on autocorrelation or cross-correlation, then add various pre- and post-processing mechanisms to deal with these problems

Autocorrelation is not enough: **pre**-processing

- low-pass filtering the speech waveform
 - often combined with downsampling the waveform (reduces computational cost) • removes vocal tract information (e.g., formants) and unvoiced sounds









frequency

8kHz



Auto-correlation is not enough: **pre**-processing

spectral flattening

- spectral envelope shape in the first place)
- inverse filtering is one way to do this
- harder

 - spectrum

• this is a general term for removing the vocal tract information (which is the cause of the

• Taylor suggests that inverse filtering introduces artefacts that actually make FO estimation

• I presume he means phase distortions - these make the residual less like an impulse train • Nevertheless, some pitch tracking algorithms do use inverse filtering to flatten the







frequency

8kHz



Figure 2 from David Talkin "A Robust Algorithm for Pitch Tracking (RAPT)" in Speech Coding and Synthesis, W. B. Kleijn and K. K. Palatal (eds), pages 497-518 Elsevier Science B.V., 1995

sampling rate = 8 kHz



Auto-correlation is not enough: **post**-processing

dynamic programming



Figure from YAAPT webpage http://ws2.binghamton.edu/zahorian/yaapt.htm

Pre-processing + autocorrelation + post-processing

- This is the typical architecture of many F0 estimation algorithms
- There tend to be a lot of parameters to tune

Constant	Meaning
F0 _{min}	minimum F0
Flores	maximum F0
t :	analysis frame
10	correlation wi
CAND.TR	minimum aco
LAG.WT	linear lag tap-
FREQ.WT	cost factor for
VTRAN.C	fixed voicing-s
VTR.A.C	delta amplitu
VTR.S.C	delta spectrus
VO_BIAS	bias to encour
DOUBL.C	cost of exact I
A.FACT	term to decrea
N.CANDS	max number

The tuneable parameters in RAPT - in David Talkin "A Robust Algorithm for Pitch Tracking (RAPT)" in Speech Coding and Synthesis, W. B. Kleijn and K. K. Palatal (eds), pages 497-518, Elsevier Science B.V., 1995

	Value
to search for (Hz)	50
to search for (Hz)	500
step size (sec)	.01
ndow size (sec)	0075
ptable peak value in NCCF	
r factor for NCCF	.3
F0 change	.02
tate transition cost	.005
de modulated transition cost	.5
n modulated transition cost	.5
age voiced hypotheses	0.0
F0 doubling or halving	.35
ase \$ of weak signals	10000
of hypotheses at each frame	20
use \$\$ of weak signals of hypotheses at each frame	1000 20



Alternatives to autocorrelation

cepstral domain methods



Michael Noll, "Cepstrum Pitch Determination" J. Acoust. Soc. Am. 41, 293, 1967

Alternatives to autocorrelation

comb filtering

- an adaptive filter that eliminates the harmonics (at multiples of FO)
- of energy from the signal (i.e., the output is as small as possible)



• the response of the filter is adaptively varied so that it removes the maximum amount

Y.K. Jang, J.F. Chicharo and B. Ribbum "Pitch Detection And Estimation Using Adaptive IIR Comb Filtering", in Proc SST 1992, pages 54-59, Brisbane, Australia, 1992



Alternatives to autocorrelation

probabilistic methods

- still uses pre-processing followed by autocorrelation
- train a classifier using supervised learning needs ground truth training data
- still needs some post-processing (dynamic programming)



Byung Suk Lee, Daniel P. W. Ellis "Noise Robust Pitch Tracking by Subband Autocorrelation Classification" in Proc. Interspeech 2012, September, Portland, OR, USA

Evaluation

- What is **ground truth**?
 - hand-labelled (or hand-corrected) F0 contours
 - Laryngograph (also known as an "Electroglottograph" or EGG) recordings
 - various **public databases** available
 - e.g., http://www.cstr.ed.ac.uk/research/projects/fda/
- What **type of errors**?
 - voicing status errors (in all speech)
 - F0 error (in voiced speech)





www.fon.hum.uva.nl



What about different voice qualities, such as creaky?

- <u>F0 estimation</u> algorithms
 - usually assume perfect periodicity
 - therefore will perform poorly on creaky voice
- epoch detection algorithms
 - vary in their ability to handle different voice qualities
- overall, we expect it will be harder to vocode some voice qualities

John Kane, Christer Gobl "Evaluation of glottal closure instant detection in a range of voice qualities" in Speech Communication 55(2), 2013, pages 295–314



Speech signal analysis

- what we need to analyse
- epoch detection ('pitch marking')
- F0 estimation ('pitch tracking')
- <u>spectral envelope estimation</u>



frequency

8kHz

Interference from the source

- Kawahara et al state:
 - time domain."
- and
 - spectrum shows periodic variation in the frequency domain."

Hideki Kawahara, Ikuyo Masuda-Katsuse and Alain de Cheveigné "Restructuring speech representations using" a pitch-adaptive time-frequency smoothing and an instantaneous-frequency-based FO extraction: Possible role of a repetitive structure in sounds" in Speech Communication, Volume 27, Issues 3–4, April 1999, Pages 187–207. DOI: 10.1016/S0167-6393(98)00085-5

• "When the length of a time window for spectral analysis is **comparable** to the fundamental period of the signal repetition, the resultant power spectrum shows periodic variation in the

• "When the length of a time window **spans several repetitions**, the resultant power





"When the length of a time window for spectral analysis is comparable to the fundamental period of the signal repetition, the resultant power spectrum shows periodic variation in





32768

"When the length of a time window spans several repetitions, the resultant power spectrum shows periodic variation in the frequency domain."



The STRAIGHT vocoder

- <u>Analysis phase</u> (the most important part of STRAIGHT)
 - F0 adaptive window to minimise the interference with spectral envelope
 - standard FFT analysis
 - amplitudes
 - we'll cover in a little later)
- <u>Synthesis phase</u> (generating a waveform)
 - coming later, when we consider **modelling** speech signals

• interpolation in frequency to extract **smooth** spectral envelope from harmonic

• estimation of ratio between **periodic and aperiodic energy** at each frequency (which



window size: 256 samples (16 ms, or approx 2 pitch periods)



STRAIGHT spectral envelope + Mel cepstral analysis



figure from Heiga Zen

Speech signal modelling

- speech parameters
- representations suitable for modelling
- converting back to a waveform

Speech signal modelling

- <u>speech parameters</u>
- representations suitable for modelling
- converting back to a waveform

Orientation

- <u>So far</u>: speech signal **analysis**
- epochs
- F0
- spectral envelope
- <u>Coming up</u>: speech signal **modelling**
 - speech parameters
 - representations suitable for modelling
 - converting back to a waveform



Orientation

- <u>So far</u>: speech signal **analysis**
 - epochs
 - F0
 - spectral envelope
- <u>Coming up</u>: speech signal **modelling**
 - speech parameters
 - representations suitable for modelling
 - converting back to a waveform



- fundamental frequency (F0)
- aperiodic energy

Speech signal modelling

- speech parameters
- <u>representations suitable for modelling</u>
- converting back to a waveform

Representations of the speech parameters that are suitable for modelling

- Vocoder is essentially a source-filter model
 - except we use an excitation signal + spectral envelope, not the "true" source+filter
- <u>excitation signal</u>
 - a periodic signal (e.g., a pulse train) at a frequency of FO
 - switched on and off by a voiced/unvoiced (V/UV) decision
- <u>spectral envelope</u>
 - we need a **representation** that is amenable to statistical modelling
- <u>aperiodic energy</u>
 - spectrally-shaped noise

Representations of the speech parameters that are suitable for modelling

- We want parameters that are
 - fixed in number (per frame) and as low dimensional as possible
 - at a **fixed** frame rate
 - a good **separation** of prosodic and segmental identity aspects of speech
 - so that we can model (and/or modify) either of them independently
 - well behaved and stable, when we perturb them (e.g., by averaging, or modelling error)
 - consecutive frames within a single speech sound
 - frames pooled from several similar sounds
- and for statistical modelling, we may additionally like to have
 - statistically **uncorrelated** parameters (to avoid having to model covariance)



What does STRAIGHT actually produce?

• ... and is it suitable for modelling?

- <u>smooth spectral envelope</u>
- <u>FO</u>
- <u>non-periodicity</u>
 - in other words, aperiodic energy



Figure: Hideki Kawahara

What does STRAIGHT actually produce?

- <u>smooth spectral envelope</u>
- high resolution (same as FFT)
- highly-correlated parameters
- probably **not** suitable for statistical modelling
 - at least, not with diagonalcovariance Gaussians

Figure: Hideki Kawahara

0



Improving the representation of the spectral envelope

- warp frequency scale
- decorrelate
- reduce dimensionality



Figure: Hideki Kawahara

Representing the spectral envelope as the Mel-cepstrum

- <u>warp the frequency scale</u>
 - instead of lossy discrete filterbank, use a **continuous** function (all-pass filter)
- <u>decorrelate</u>
 - convert from spectrum to cepstrum
- <u>reduce dimensionality</u>
 - **truncate** the cepstrum
 - in ASR, we kept the first 12 coefficients
 - in synthesis, we'll use a lot more, perhaps the first 40-60 coefficients

• Not quite the same as the MFCCs we use in ASR, but basically the **same motivation**

What does STRAIGHT actually produce?

- <u>aperiodic energy</u>
 - effectively the ratio between periodic and aperiodic energy, at each frequency
- high resolution (same as FFT)
- highly-correlated parameters

Figure: Hideki Kawahara

frequency

8kHz

Improving the representation of the aperiodic energy

- <u>aperiodic energy</u>
- reduce dimensionality
 - simply reduce resolution by averaging across broad frequency **bands**
 - e.g., between 5 and 25 bands (on a Mel scale, of course)

Figure: Hideki Kawahara

Speech signal modelling

- speech parameters
- representations suitable for modelling
- <u>converting back to a waveform</u>

STRAIGHT analysis and synthesis

Figure: Hideki Kawahara

Excitation signal

figure from Heiga Zen

What next?

- We have decomposed speech into
- F0, plus a V/UV decision
- smooth spectral envelope, parameterised as the Mel-cepstrum
- band aperiodicity parameters

- We've seen how to reconstruct the waveform
- Now we can insert a statistical model between the analysis and synthesis parts

Figures: Hideki Kawahara

What next?

Figures: Hideki Kawahara

)