

Autotune

Giacomo Radicchi

Maggio 2025

Indice

1	Introduzione	2
2	Cenni teorici	2
2.1	Concetto di pitch	2
2.2	Frequenza di una nota	2
2.3	Discrete Fourier Transform	3
2.4	Short-Time Fourier Transform	3
2.5	Autocorrelazione	4
3	Modifica del pitch	4
3.1	Resampling dell'audio	4
3.2	Time stretch con Phase Vocoder	7
3.3	Pitch shift	10
4	Pitch detection	12
4.1	Utilizzo dell'autocorrelazione	12
4.2	Squared Difference Function	15
5	Algoritmo di Tuning Completo	17
5.1	Segmentazione e Finestratura	17
5.2	Stima e Correzione del Pitch	17
5.3	Pitch Shifting	17
5.4	Overlap-Add con Allineamento	17
5.5	Normalizzazione Finale	19
5.6	Output	19
5.7	Implementazione	19
6	Conclusioni e sviluppi futuri	22

1 Introduzione

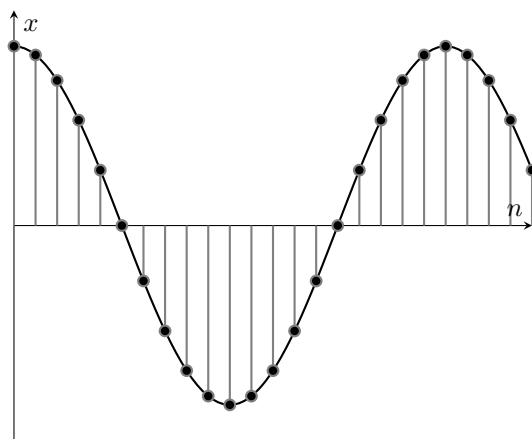
L'*autotune* è uno degli strumenti più utilizzati nella musica contemporanea. Nato alla fine degli anni '90, il suo impatto sulla produzione musicale è stato talmente profondo da trasformarsi in una vera e propria firma sonora, soprattutto nei generi pop, trap e hip-hop. La sua funzione principale è quella di correggere l'intonazione delle linee vocali, adattandole a una scala musicale di riferimento. In ogni istante, l'autotune rileva l'intonazione effettiva della voce e la modifica quel tanto che basta affinché corrisponda alla nota più vicina all'interno della scala selezionata.

Nel presente lavoro, si analizzeranno i fondamenti teorici alla base dell'elaborazione dell'intonazione, soffermandosi in particolare su tecniche come il time-stretching, il pitch-shifting, l'utilizzo della STFT e dell'autocorrelazione. Si affronterà anche un'implementazione in ambiente *MATLAB* e *Python* di alcuni algoritmi base utili alla realizzazione di un semplice sistema di autotune.

2 Cenni teorici

2.1 Concetto di pitch

Immaginiamo di avere come sequenza una sinusoide del tipo $x[n] = \cos(2\pi f_o n)$



Nello specifico, f_o rappresenta la frequenza in Hertz (**Hz**) che determina l'intonazione, o **pitch**, del segnale. Può anche essere visto come il reciproco del periodo del segnale $T = \frac{1}{f_o}$, relazione che ci tornerà più utile in avanti.

2.2 Frequenza di una nota

Ad ogni nota musicale è associata la propria frequenza. In realtà, è necessario definirla solamente per una nota, tipicamente il La - o A in inglese -, per determinare tutte le altre 12 note della scala cromatica (Do, Do#, Re, Re#, Mi, Fa, Sol, Sol#, La, La#, Si). Questo perché per poter ricavare la frequenza della

nota successiva, vale a dire aumentare il pitch di un **semitono**, è necessario moltiplicarlo per $2^{\frac{1}{12}}$. Da notare che la scala delle note è periodica di periodo 12, perciò avanzare di 12 semitoni - dunque moltiplicare per 2 - equivale a prendere la stessa nota ma di un'**ottava** superiore. La frequenza standard attribuita al La è quella di **440** Hz e da questa si possono ricavare le rimanenti. Per esempio, se si volesse calcolare il pitch del Si (B), basta moltiplicare quello del La (A) per $2^{\frac{2}{12}}$, dunque $440 \text{ Hz} \times 1.122 = 493.883 \text{ Hz}$.

2.3 Discrete Fourier Transform

Come vedremo successivamente, è necessario riuscire a passare dal dominio del tempo a dominio della frequenza. Per fare ciò, si applica la **DFT** (Discrete Fourier Transform) la cui formula è la seguente:

$$X[k] = \sum_{n=0}^{N-1} x[n] e^{-j2\pi \frac{k}{N} n} \quad (1)$$

dove N è la dimensione della sequenza X[n]. Da notare che X[k] è ottenuta per campionamento di uno spettro periodico di periodo 1, nonché della **DTFT** (Discrete-Time Fourier Transform), dunque la DFT sarà anch'essa periodica ma di periodo N.

2.4 Short-Time Fourier Transform

La **STFT** (Short-Time Fourier Transform) ci da come informazione come varia la trasformata di Fourier al variare del tempo. In particolare, la sequenza originale viene suddivisa in una serie di sotto-sequenze, o frame, ottenute mediante il prodotto con una specifica window function (ad esempio: Hann, Hamming, Blackman, ecc.):

$$x_i[n] = x[n + i(N - L)] w_N[n] \quad (2)$$

dove:

- N è la lunghezza della finestra, o **frame length**
- L è il numero di campioni sovrapposti tra un frame e il successivo, ovvero l'**hop length**
- $x_i[n]$ è la **window function** di lunghezza N

Da notare che le sotto sequenze sono sovrapposte per evitare discontinuità nel calcolo della STFT.

Per ogni frame $x_i[n]$, viene calcolata la Discrete Fourier Transform:

$$X_i[k] = \sum_{n=0}^{N-1} x_i[n] e^{-j2\pi \frac{k}{N} n} \quad (3)$$

2.5 Autocorrelazione

Dato un segnale $x(t)$, l'autocorrelazione misura la somiglianza del segnale con se stesso al variare di un ritardo temporale τ :

$$R_{xx}^{(E)}(\tau) = \int_{-\infty}^{\infty} x^*(t) x(t + \tau) dt \quad (4)$$

Dove $x^*(t)$ rappresenta il complesso coniugato di $x(t)$. Nel caso di segnali reali, in cui il complesso coniugato coincide con il segnale originale, la formula si semplifica in:

$$R_{xx}^{(E)}(\tau) = \int_{-\infty}^{\infty} x(t) x(t + \tau) dt \quad (5)$$

Nel caso dei segnali discreti, l'integrale si riduce ad una sommatoria:

$$R_{xx}^{(E)}[l] = \sum_{n=-\infty}^{\infty} x[n] x[n + l] \quad (6)$$

In pratica, l'autocorrelazione consente di evidenziare la periodicità presente nel segnale: per un segnale periodico, i valori di $R_{xx}^{(E)}(\tau)$ mostrano picchi nei ritardi corrispondenti al periodo fondamentale e ai suoi multipli. Questo rende l'autocorrelazione uno strumento efficace per stimare la frequenza fondamentale, o pitch, di una linea vocale.

3 Modifica del pitch

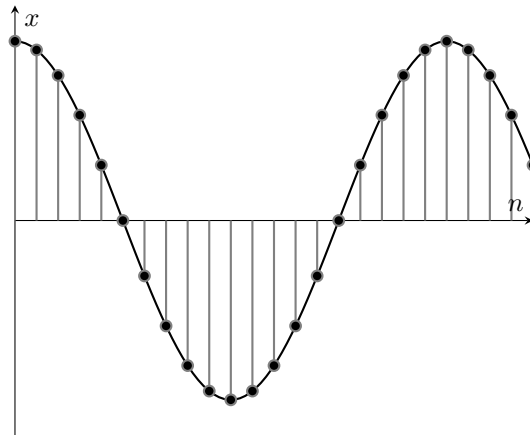
La modifica dell'intonazione di una traccia audio rappresenta una delle problematiche più discusse nella modellazione audio. Come è possibile riuscire ad aumentare la frequenza di un qualsiasi audio a piacimento?

3.1 Resampling dell'audio

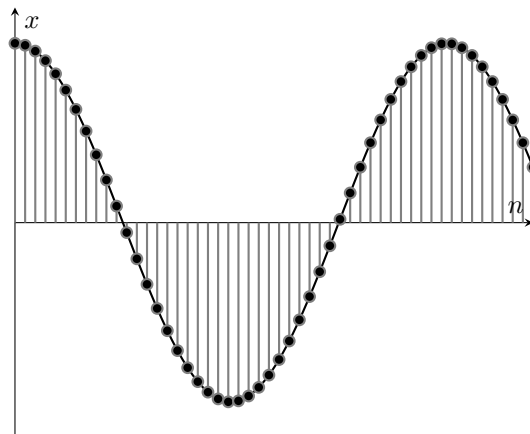
Un primo metodo che permette di fare ciò è quello di sovracampionare o sottocampionare il segnale, ovvero fare del **resampling**. Queste pratiche danno luogo, rispettivamente, ad una diminuzione o aumento del pitch.

Immaginiamo di voler dimezzare la frequenza della traccia audio, nonché abbassare il pitch di un'ottava. Per fare ciò, è necessario sovracampionare il nostro segnale dimezzando la distanza tra ogni singolo campione e, dunque, raddoppiare il numero di campioni presenti all'interno dell'audio.

Audio originale:



Audio sovracampionato:



Potremmo implementare ciò attraverso un semplice programma:

- MATLAB:

```

1 function x_out = resample(file_name, pitch_factor)
2     [data, sr] = audioread(file_name);
3
4     % lavoriamo con audio mono per semplicità
5     if size(data, 2) > 1
6         data = data(:, 1);
7     end
8
9     og_length = size(data, 1);
10
11     % genero l'asse dei tempi originale:
12     old_t = 1 : og_length;
13     % genero l'asse dei tempi del nuovo audio:

```

```

14     new_t = 1 : pitch_factor : og_length;
15
16     % nuova sequenza ottenuta per interpolazione lineare
17     :
18     x_out = interp1(old_t, data, new_t);
19
20     % salvataggio audio:
21     audiowrite("resampled_audio.wav", x_out, sr);
22
23     x_out = data;
24 end

```

- Python:

```

1 import numpy as np
2 import soundfile as sf
3
4 audio, sr = sf.read("audio.wav")
5
6 # lavoriamo con audio mono per semplicità
7 if len(np.shape(audio)) > 1:
8     audio = audio[:, 0]
9
10 og_length = np.shape(audio)[0]
11 pitch_factor = 0.5
12
13 # genero l'asse dei tempi originale:
14 old_t = np.arange(0, og_length)
15 # genero l'asse dei tempi del nuovo audio:
16 new_t = np.arange(0, og_length, pitch_factor)
17
18 # nuova sequenza ottenuta per interpolazione lineare:
19 new_audio = np.interp(new_t, old_t, audio)
20
21 # salvataggio audio
22 sf.write("pitched_audio1.wav", new_audio, sr)

```

Nell'esempio fornito, i campioni sono stati interpolati linearmente. Questo tipo di ricostruzione non è del tutto preciso e può portare ad artefatti o suoni indesiderati. L'ideale sarebbe attuare un'interpolazione cardinale sfruttando la formula tratta dal teorema del campionamento di Shannon-Nyquist:

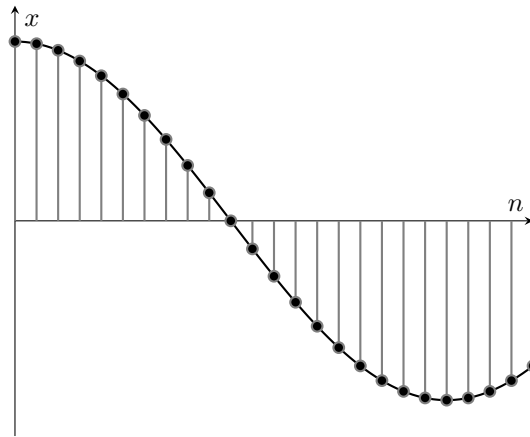
$$x(t) = \sum_{n=-\infty}^{\infty} x(nT_c) \operatorname{sinc} \left[\frac{\pi}{T_c} (t - nT_c) \right] \quad (7)$$

Per questioni di efficienza, verranno utilizzate le funzioni "interp1" per MATLAB e "scipy.signal.resample" della libreria scipy per Python.

Ovviamente, la frequenza di campionamento, o **sample rate**, con la quale si andrà a salvare la nuova traccia sovracampionata dovrà rimanere la stessa

precedente, altrimenti l'audio riprodotto rimarrà totalmente identico o, nel caso di sottocampionamento, potrebbe risultare distorto a causa del fenomeno di *aliasing*.

Audio sovracampionato con sample rate originale:



Mantenendo lo stesso sample rate, tuttavia, l'audio non solo cambierà frequenza, ma anche durata, cosa che si vorrebbe evitare.

3.2 Time stretch con Phase Vocoder

Arrivati a questo punto, vorremmo poter svincolare la durata dell'audio dalla sua intonazione. È possibile farlo tramite un time-stretcher, che consente di modificare la durata dell'audio senza alterarne l'intonazione, sfruttando un approccio basato sul **phase vocoder**: esso lavora nel dominio della frequenza, andando a modificare la fase di ogni singolo spettro ottenuto tramite **Short-Time Fourier Transform (STFT)**.

L'idea alla base del time stretching con un phase vocoder è quella di costruire una nuova STFT interpolando opportunamente quella originale (sia in modulo che in fase), e successivamente ricostruire la sequenza nel dominio del tempo mediante la **Inverse Short-Time Fourier Transform (I-STFT)**. Tale procedura permette di ottenere una nuova sequenza:

- di lunghezza diversa rispetto all'originale
- ma con stesso pitch

Tutto ciò può essere benissimo gestito sia su MATLAB che su Python utilizzando, rispettivamente, le funzioni built-in o importando la libreria SciPy.

- MATLAB:

```
1 function x_stretched = time_stretch(x, sr,
2   stretch_factor)
3     % Short-Time Fourier Transform
```

```

4
5     hop_length = 512;
6     nfft = 2048;
7     window = hann(nfft, "periodic");
8
9     [S, ~, T] = stft(x, sr, 'Window', window, ...
10         'OverlapLength', nfft - hop_length, ...
11         'FFTLenght', nfft);
12
13     % Calcolo del modulo e della fase
14     magnitude = abs(S);
15     phase = angle(S);
16
17     % Generazione time steps
18     % se stretch_factor > 1, l'audio sara' piu' breve
19     % se stretch_factor < 1, l'audio sara' piu' lungo
20     time_steps = 1 : stretch_factor : size(T, 1);
21
22     % Inizializzazione nuova STFT
23     new_S = zeros(size(S,1), length(time_steps));
24
25     % l'accumulatore svolge il ruolo dell'integrale
26     phase_accumulator = phase(:, 1);
27     og_stft_length = size(S, 2);
28     i = 1;
29     for step = time_steps
30         index = floor(step);
31         fraction = step - index;
32
33         % interpolazione tra i moduli
34         mag1 = magnitude(:, uint64(index));
35         mag2 = magnitude(:, min(uint64(index) + 1,
36             og_stft_length));
37         interpolated_magnitude = (1 - fraction) *
38             mag1 + fraction * mag2;
39
40         % calcolo della frequenza istantanea (derivata)
41         phase1 = phase(:, index);
42         phase2 = phase(:, min(index + 1, og_stft_length)
43             );
44         instant_freq = phase2 - phase1;
45
46         % wrapping della frequenza istantanea
47         instant_freq = mod(instant_freq + pi, 2 * pi) -
48             pi;
49
50         % aggiornamento dell'accumulatore (integrale)
51         phase_accumulator = phase_accumulator +
52             instant_freq;

```



```

51         new_S(:, i) = interpolated_magnitude .*
52         exp(1j * phase_accumulator);
53
54         i = i + 1;
55     end
56     x_stretched = real(istft(new_S, sr, 'Window', window
57         , ...
58         'OverlapLength', nfft - hop_length, ...
59         'FFTLength', nfft));
60
61     % dissolvenza di volume all'inizio e alla fine:
62     L = length(x_stretched(1:end));
63     smooth_window = zeros(L, 1);
64     dt = 0.1; %sec
65     t1 = dt * sr;
66     t2 = L - t1;
67     smooth_window(1:uint32(t1)) = (0:t1-1)/t1;
68     smooth_window(uint32(t1):L-uint32(t1)-1) = 1;
69     smooth_window(L-uint32(t1):L) = (L:-1:t2)/t1 - L/t1
70         +1;
71
72     x_stretched = smooth_window.*x_stretched;
73
74     % normalizzazione per evitare clipping
75     max_value = max(abs(x_stretched));
76     if max_value > 0.99
77         x_stretched = x_stretched/max_value;
78     end
79 end

```

- Python:

```

1  import numpy as np
2  import scipy as sp
3
4  def time_stretch(x, fs, stretch_factor, frame_length
5      =2048,
6      hop_length=128):
7      # Short-Time Fourier Transform
8      f, t, S = sp.signal.stft(x, fs=fs, nperseg=frame_length,
9      noverlap=frame_length-hop_length, window=np.hanning(
10         frame_length))
11
12     # Calcolo del modulo e della fase:
13     magnitude = np.abs(S)
14     phase = np.angle(S)
15
16     % Generazione time steps

```

```

15 % se stretch_factor > 1 , l'audio sara' piu' breve
16 % se stretch_factor < 1 , l'audio sara' piu' lungo
17 time_steps = np.arange(0, len(t), stretch_factor)
18
19 # Inizializzazione nuova STFT
20 new_S = np.zeros(shape=(np.shape(S)[0], len(time_steps)))
21
22 ,
23 dtype=np.complex64)
24
25 # l'accumulatore svolge il ruolo dell'integrale
26 phase_accumulator = phase[:, 0]
27
28 for i, step in enumerate(time_steps):
29     index = int(np.floor(step))
30     fraction = step - index
31
32     # interpolazione tra i moduli:
33     mag1 = magnitude[:, index]
34     mag2 = magnitude[:, min(index + 1,
35                             np.shape(magnitude)[1] - 1)]
36     interpolated_magnitude = (1 - fraction) *
37     mag1 + fraction * mag2
38
39     # calcolo della frequenza istantanea (derivata):
40     phase1 = phase[:, index]
41     phase2 = phase[:, min(index + 1,
42                             np.shape(magnitude)[1] - 1)]
43     instant_freq = phase2 - phase1
44
45     # wrapping della frequenza istantanea
46     instant_freq = np.mod(instant_freq + np.pi,
47                             2 * np.pi) - np.pi
48
49     # aggiornamento dell'accumulatore (integrale)
50     phase_accumulator += instant_freq #stretch_factor
51
52     new_S[:, i] = interpolated_magnitude
53     * np.exp(1j * phase_accumulator)
54
55 _, x_stretched = sp.signal.istft(new_S, fs=fs, nperseg=
56     frame_length,
57     noverlap=frame_length-hop_length)
58 return x_stretched

```

3.3 Pitch shift

Dato che l'obiettivo principale è quello di modificare solo e unicamente la frequenza, è possibile un *pitch shift* attraverso una combinazione di **time stretch** e **resampling** per raggiungere tale obiettivo. Qualora volessimo moltiplicare

la frequenza f_o di un fattore pari a k , posto L pari alla durata della traccia, la procedura da seguire è la seguente:

- **time stretch** di un fattore $\frac{1}{k} \Rightarrow \text{durata} = \frac{L}{k}$; frequenza = f_o
- **resample** con step pari a $k \Rightarrow \text{durata} = L\frac{k}{k} = L$; frequenza = kf_o

Al termine di ciò, la durata del nuovo audio coinciderà con quella dell'audio originale - o quasi, ci potrebbero essere degli errori di approssimazione dovuti al calcolo di $\frac{1}{k}$, dunque sarà necessario ricampionare un'ulteriore volta per mantenere lo stesso numero di campioni presenti inizialmente - ma l'intonazione sarà del tutto nuova. A livello implementativo, tale tecnica può essere realizzata in poche righe di codice:

- MATLAB:

```

1 function x_pitched = pitch_shifting(x, sr, pitch_factor)
2     x_stretched = time_stretch(x, sr, 1/pitch_factor);
3
4     % viene fatto uso della funzione resample
5     % precedentemente implementata:
6     x_pitched = resample(x_stretched, sr, pitch_factor);
7
8     % la lunghezza del segnale in ingresso x deve essere
9     % uguale a quella del segnale in uscita x_pitched:
10    old_t = 1:length(x_pitched);
11    new_t = 1:length(x);
12    x_pitched = interp1(old_t, x_pitched, new_t);
13 end

```

- Python:

```

1 def pitch_shifting(x, fs, pitch_factor,
2 frame_length=2048, hop_length=512):
3     x_stretched = time_stretch(x, fs, 1/pitch_factor,
4 frame_length, hop_length)
5     stretched_length = len(x_stretched)
6
7     t_old = np.arange(0, stretched_length)
8     t_new = np.arange(0, stretched_length, pitch_factor)
9     x_pitched = np.interp(t_new, t_old, x_stretched)
10
11    t_old = np.arange(0, len(x_pitched))
12    t_new = np.arange(0, len(x))
13    x_pitched = np.interp(t_new, t_old, x_pitched)
14
15    return x_pitched

```

4 Pitch detection

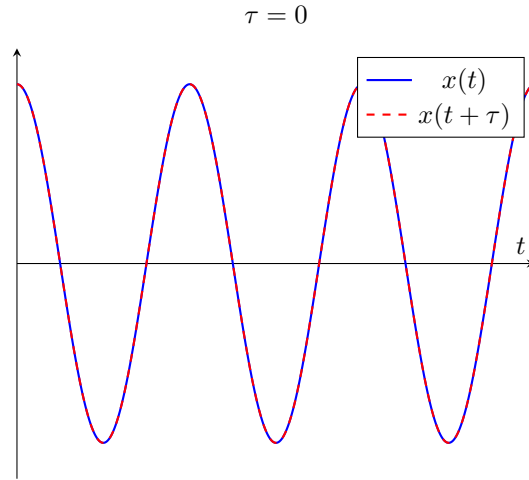
4.1 Utilizzo dell'autocorrelazione

A questo punto, è necessario sviluppare un algoritmo in grado di identificare con precisione l'intonazione di una linea vocale, al fine di correggere e adattare la nota desiderata. Ciò potrebbe essere implementato nel dominio del tempo attraverso l'autocorrelazione. Noto che la frequenza fondamentale f_0 è il reciproco del periodo T , ci basterà calcolare il periodo per risalire alla frequenza fondamentale.

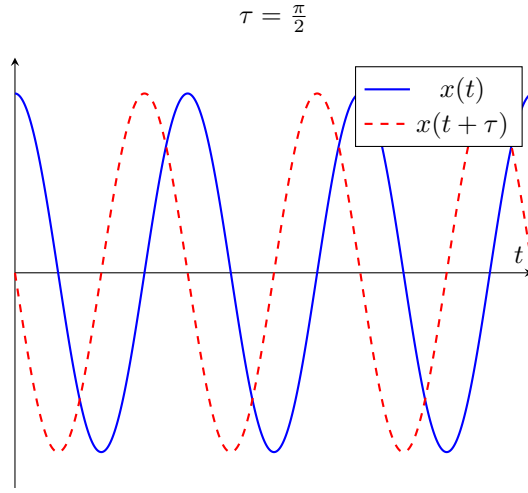
Per introdurre il concetto, consideriamo un segnale di test costituito da un coseno puro:

$$x(t) = \cos(t)$$

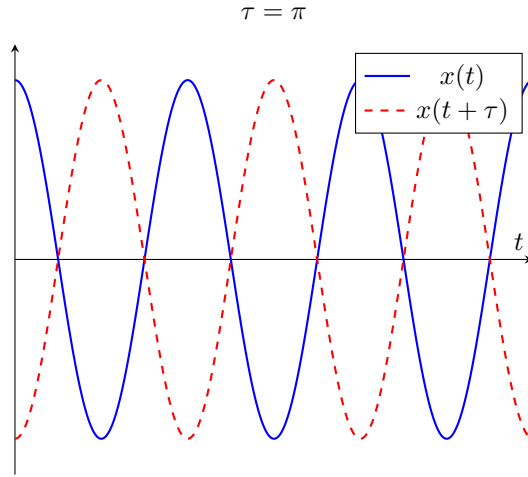
e calcoliamone l'autocorrelazione per $\tau \in \{0, \frac{\pi}{2}, \pi, \frac{3}{2}\pi, 2\pi\}$. Per semplicità, eseguiremo l'integrale su un intervallo finito.



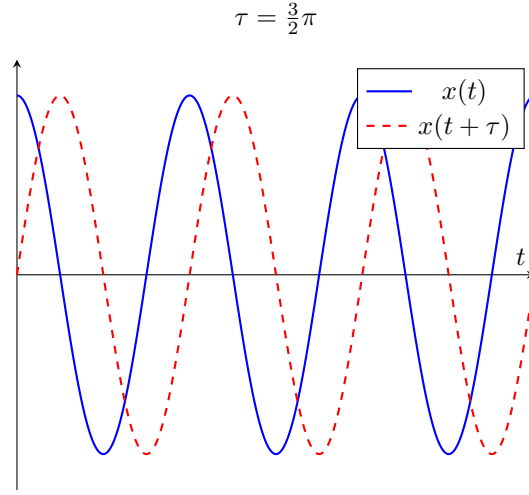
$$R_{xx}^{(E)}(0) = \int_0^{2\pi} x(t) x(t+0) dt = \int_0^{2\pi} x^2(t) dt = \pi \simeq 3.1416 \quad (8)$$



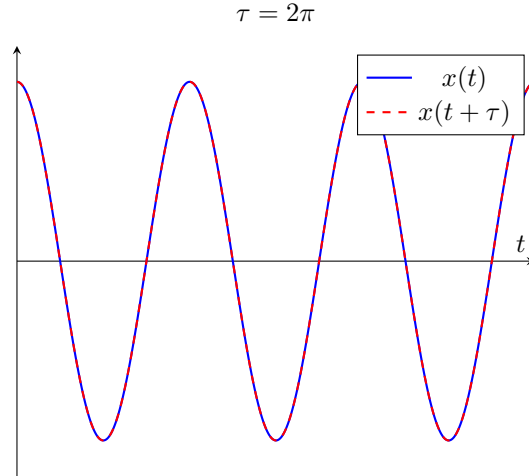
$$R_{xx}^{(E)}\left(\frac{\pi}{2}\right) = \int_0^{2\pi} x(t) x\left(t + \frac{\pi}{2}\right) dt = 0 \quad (9)$$



$$R_{xx}^{(E)}(\pi) = \int_0^{2\pi} x(t) x(t + \pi) dt = -\pi \simeq -3.1416 \quad (10)$$



$$R_{xx}^{(E)}\left(\frac{3}{2}\pi\right) = \int_0^{2\pi} x(t) x\left(t + \frac{3}{2}\pi\right) dt = 0 \quad (11)$$



$$R_{xx}^{(E)}(2\pi) = \int_0^{2\pi} x(t) x(t + 2\pi) dt = \pi \simeq 3.1416 \quad (12)$$

Alla luce di ciò, quando τ coincide con il periodo, l'autocorrelazione raggiunge lo stesso valore dell'origine, generando un massimo. Il primo massimo relativo della linea melodica permette quindi di determinare la frequenza fondamentale.

Nella pratica, però, le cose sono diverse. Innanzitutto, lavoriamo con segnali numerici, di lunghezza finita e discretizzati, per cui il calcolo si riduce a una

sommatoria. Inoltre, i segnali reali non sono perfette sinusoidi: per questo motivo, il metodo teorico visto in precedenza non risulta applicabile in modo diretto.

4.2 Squared Difference Function

Per ottenere una maggiore precisione nel calcolo, invece di valutare l'area del prodotto di due funzioni traslate, risulta più efficace calcolare l'area della loro differenza al quadrato. Questa funzione prende il nome di *Funzione di Differenza al Quadrato* o *Squared Difference Function* (**SDF**). Quando i due segnali si sovrappongono, la loro differenza è minima. La SDF si definisce come:

$$d(\tau) = \frac{1}{N} \sum_{n=0}^{N-1} (x[n] - x[n + \tau])^2 \quad (13)$$

dove:

- N è la lunghezza della finestra considerata;
- il fattore di normalizzazione $\frac{1}{N}$ permette di compensare la riduzione della lunghezza utile della finestra all'aumentare di τ , evitando che la funzione sia influenzata da questo effetto.

Una volta calcolata la SDF, sarà sufficiente individuare il primo minimo locale per determinare il valore di τ corrispondente alla migliore sovrapposizione tra i due segnali. Infine, dopo aver trovato τ , è necessario ricavare il periodo T noto che:

- $T = \frac{1}{f_0}$
- $f_0 = \frac{\tau}{sr}$, con sr pari alla frequenza di campionamento o **sample rate**

dunque:

$$T = \frac{sr}{\tau} \quad (14)$$

Ciò può essere implementato in questo modo:

- MATLAB:

```

1 function f0 = find_pitch(sample, fs)
2     % Intervallo di frequenze vocali (Hz)
3     voice_range = [80, 200];
4     min_lag = floor(fs / voice_range(2));
5     max_lag = floor(fs / voice_range(1));
6
7     % Calcolo delle differenze quadrate normalizzate
8     diff_areas = zeros(1, max_lag - min_lag + 1);
9     for tau = min_lag:max_lag
10         diff_areas(tau - min_lag + 1) = ...

```

```

11         get_difference_area(sample, sample(tau+1:end));
12     end
13
14     % Ricerca del primo minimo relativo
15     delta_lag = 0;
16     for tau = 2:length(diff_areas)-1
17         if diff_areas(tau) < diff_areas(tau-1) ...
18             && diff_areas(tau) < diff_areas(tau+1)
19             delta_lag = tau;
20             break;
21         end
22     end
23
24     % Se non trovato, prendo il minimo globale
25     if delta_lag == 0
26         [~, delta_lag] = min(diff_areas);
27     end
28
29     % Conversione in frequenza fondamentale
30     f0 = fs / (delta_lag + min_lag - 1);
31 end
32
33
34 function diff = get_difference_area(x1, x2)
35     L = min(length(x1), length(x2));
36     diff = sum((x1(1:L) - x2(1:L))^2) / L;
37 end

```

- Python:

```

1 def find_pitch(sample, fs):
2     voice_range = [80, 200] # Hz
3     min_lag = fs // voice_range[1]
4     max_lag = fs // voice_range[0]
5
6     diff_areas = [get_difference_area(sample,
7     sample[tau:]) for tau in range(min_lag, max_lag)]
8     mean = np.sum
9
10    # il primo minimo relativo trovato sara'
11    # il tau ricercato
12    delta_lag = 0
13    for tau in range(1, len(diff_areas) - 1):
14        if diff_areas[tau] < diff_areas[tau-1]
15            and diff_areas[tau] < diff_areas[tau+1]:
16            delta_lag = tau
17
18    if delta_lag == 0:
19        delta_lag = np.argmin(diff_areas)
20

```



```

21     return fs / (delta_lag + min_lag)
22
23
24 def get_difference_area(x1, x2):
25     length = min(len(x1), len(x2))
26     return np.sum((x1[:length] - x2[:length])
27                   ** 2) / length

```

5 Algoritmo di Tuning Completo

L'algoritmo finale di **tuning** implementa in maniera completa il processo di correzione dell'intonazione di un segnale audio discreto $x[n]$ campionato a frequenza f_s . Esso combina i metodi introdotti nelle sezioni precedenti: segmentazione in finestre, stima della frequenza fondamentale, ricerca della nota da raggiungere, pitch shifting e ricostruzione mediante overlap-add.

5.1 Segmentazione e Finestratura

Il segnale viene suddiviso in frame di lunghezza $L = 2048$ campioni, con hop size pari a $H = L/4$. Ogni frame è moltiplicato per una finestra di Hann $\omega[n]$ al fine di ridurre effetti di discontinuità.

5.2 Stima e Correzione del Pitch

Per ciascun frame viene stimata la frequenza fondamentale \hat{f}_0 mediante la funzione già precedentemente implementata **find_pitch**. La nota più vicina appartenente alla scala scelta S viene individuata come

$$f^* = \text{find_closest_note}(\hat{f}_0, S).$$

Il fattore di correzione è quindi

$$\alpha = \frac{f^*}{\hat{f}_0}.$$

5.3 Pitch Shifting

Il frame viene trasformato tramite la funzione di pitch shifting che è stata realizzata, secondo la relazione:

$$y_i[n] = \text{pitch_shift}(x_i[n], f_s, \alpha),$$

5.4 Overlap-Add con Allineamento

I frame processati vengono ricombinati nel segnale di output $x_{out}[n]$ con tecnica di overlap-add. Per garantire continuità temporale, viene stimato un ritardo τ

che massimizza la correlazione tra frame adiacenti:

$$x_{out}[n] = x_{out}[n] + y_i[n - \tau].$$

In questo modo:

- MATLAB:

```
1 function best_lag = find_best_lag(current_sample, frame,
   start_index, n_overlap)
2     % Calcola il lag massimo consentito
3     max_lag = floor(n_overlap / 8);
4
5     % Definisce la regione di ricerca
6     search_start = max(1, start_index - max_lag);
7     search_end = min(length(current_sample), start_index
   + length(frame) + max_lag);
8
9     search_region = current_sample(search_start:
   search_end);
10
11     % Se la regione troppo corta, ritorna 0
12     if length(search_region) < length(frame)
13         best_lag = 0;
14         return;
15     end
16
17     % Calcola la correlazione incrociata
18     correlation = xcorr(search_region, frame, 'none');
19
20     % Trova l'offset migliore
21     [~, best_offset] = max(correlation);
22
23     % Calcola il lag relativo alla posizione originale
24     best_lag = (search_start + best_offset - 1) -
   start_index;
25 end
```

- Python:

- MATLAB:

```
1 def find_best_lag_2(current_sample, frame, start_index,
   n_overlap):
2     max_lag = n_overlap // 8
3
4     # Estrai la regione di interesse dal segnale
   corrente
5     search_start = max(0, start_index - max_lag)
6     search_end = min(len(current_sample), start_index +
   len(frame) + max_lag)
```

```

7     search_region = current_sample[search_start:
      search_end]
8
9     if len(search_region) < len(frame):
10         return 0
11
12     # Usa correlazione crociata
13     correlation = np.correlate(search_region, frame,
      mode='valid')
14     best_offset = np.argmax(correlation)
15
16     # Converti in lag relativo alla posizione originale
17     best_lag = search_start + best_offset - start_index
18
19     return best_lag

```

5.5 Normalizzazione Finale

Infine, il segnale viene normalizzato per evitare fenomeni di clipping:

$$x_{out} \leftarrow \frac{x_{out}}{\max(x_{out})}.$$

5.6 Output

La funzione restituisce x_{out} , un segnale in cui ogni frame è stato intonato alla nota più vicina della scala musicale selezionata, con transizioni continue e senza artefatti percettibili.

5.7 Implementazione

- MATLAB:

```

1 function x_tuned = tuning(data, sr, scale, note_dict, A)
2     sample_length = numel(data);
3     frame_length = 2048;
4     hop_length = frame_length / 4;
5
6     x_tuned = zeros(size(data));
7
8     % overlap add
9     window = hamming(frame_length);
10
11     for ii = (1: hop_length: sample_length)
12         if ii + frame_length >= sample_length
13             frame = zeros(frame_length, 1);
14             frame(1:sample_length - (ii-1)) = data(ii:
      sample_length);
15         else

```

```

16         frame = data(ii : ii + frame_length-1);
17     end
18
19
20     % multiplying frame by window function
21     frame = frame.* window;
22
23
24     % finding frame pitch
25     pitch = find_pitch(frame, sr);
26
27     goal_pitch = find_closest_note(pitch, note_dict(
        scale), note_dict, A);
28
29     %shift_factor = goal_pitch/pitch;
30     shift_factor = 0.9;
31     frame = pitch_shifting(frame, sr, shift_factor);
32     tau = 0;
33     if ii ~= 1
34         tau = find_best_lag(x_tuned, frame, ii,
            frame_length - hop_length);
35     end
36
37     % add
38
39     if ii + tau + frame_length >= sample_length
40         x_tuned(ii + tau : end) = x_tuned(ii + tau :
            end) + frame(1: sample_length - ii -
            tau + 1);
41     else
42         temp = x_tuned(ii + tau : ii + tau +
            frame_length-1) + frame;
43         x_tuned(ii + tau : ii + tau + frame_length
            -1) = temp;
44     end
45 end
46
47 % Normalize output
48 max_value = max(x_tuned);
49 clip = 1;
50 if max_value > clip
51     x_tuned = x_tuned / (max_value/clip);
52 end
53 end

```

- Python:

```

1 def tuning(x, fs, scale="chromatic"):
2     sample_length = len(x)
3

```

```

4     frame_length = int(2048)
5     hop_length = frame_length // 4
6
7     x_out = np.zeros_like(x)
8
9     # Overlap-add
10    window = np.hanning(frame_length)
11    for i in range(0, sample_length, hop_length):
12        if i + frame_length >= sample_length:
13            frame = np.zeros(shape=(frame_length,))
14            frame[: sample_length - i] = x[i :
15                sample_length]
16        else:
17            frame = x[i : i + frame_length]
18
19        # multiplying frame by window function
20        frame = frame * window
21
22        # finding frame pitch
23        pitch = find_pitch(frame, fs)
24
25        goal_pitch = find_closest_note(pitch, dict_scale
26            [scale])
27
28        # shifting to the aim pitch
29        frame = pitch_shifting(frame, fs, goal_pitch/
30            pitch)
31
32        tau = find_best_lag(x_out, frame, i,
33            frame_length - hop_length)
34
35        # Add
36        if i + tau + frame_length >= sample_length:
37            x_out[i + tau : sample_length] = x_out[i +
38                tau : sample_length] + frame[:
39                    sample_length - i - tau]
40        else:
41            x_out[i + tau : i + tau + frame_length] =
42                x_out[i + tau : i + tau + frame_length] +
43                    frame
44
45    # Normalize output
46    max_value = np.max(x_out)
47    clip = 1
48    if max_value > clip:
49        x_out /= (max_value/clip)
50
51    return x_out

```

6 Conclusioni e sviluppi futuri

L'implementazione descritta costituisce una versione prototipale di un algoritmo di *autotuning*. Pur dimostrando l'efficacia concettuale della tecnica di segmentazione, rilevamento della frequenza fondamentale e ricampionamento con overlap-add, essa rimane inevitabilmente una versione semplificata rispetto a soluzioni di livello professionale.

In particolare, si evidenziano alcune limitazioni che potranno costituire linee guida per sviluppi futuri:

- **Resampling sinc-based:** l'attuale approccio al cambio di pitch si affida a trasformazioni di libreria che non garantiscono un'ottimale qualità del segnale. Un ricampionamento sinc-based permetterebbe di preservare maggiormente la fedeltà spettrale.
- **Mantenimento delle formanti:** l'algoritmo non si occupa della conservazione delle formanti vocali, con il rischio di introdurre artefatti timbrici percepibili, soprattutto su segnali vocali complessi.
- **Efficienza computazionale:** l'implementazione è stata sviluppata in linguaggi interpretati (MATLAB e Python), con l'uso estensivo di cicli `for`. Ciò comporta tempi di esecuzione non ottimali, che potrebbero essere sensibilmente ridotti mediante l'uso di tecniche di vectorization o con una successiva traduzione in linguaggi compilati.

In sintesi, il lavoro fornisce una base solida e didatticamente chiara per comprendere i meccanismi fondamentali del pitch shifting automatico, pur lasciando ampi margini di miglioramento sia dal punto di vista della qualità percettiva del segnale che dell'efficienza implementativa.