



UNIVERSITÀ
DEGLI STUDI
FIRENZE

Scuola di Scienze Matematiche, Fisiche e Naturali
Corso di Laurea in Informatica

Tesi di Laurea

LIBRERIA PYTHON PER L'OCR DI
DOCUMENTI D'IDENTITÀ

PYTHON LIBRARY TO PERFORM OCR ON
IDENTITY DOCUMENTS

ALESSIO FALAI

Relatore: *Maria Cecilia Verri*

Anno Accademico 2018-2019

INDICE

Introduzione	7
1 Prerequisiti matematici nell'ambito della computer vision	9
1.1 Immagini digitali	9
1.2 Convoluzione	9
1.3 Operazioni morfologiche	10
2 Binarizzazione di immagini	11
2.1 Introduzione	11
2.2 Sogliatura globale	11
2.3 Sogliatura adattativa (o locale)	11
2.4 Sogliatura di Otsu	11
2.5 Algoritmo di sogliatura proposto	11
3 Template matching	13
3.1 Introduzione	13
3.2 SIFT	13
3.3 SURF	13
3.4 AKAZE	13
3.5 BRISK	13
3.6 HED	13
4 Text detection	15
4.1 Introduzione	15
4.2 EAST	15
4.3 Algoritmo proposto	15
5 Ulteriori migliorie	17
5.1 Motori OCR	17
5.2 Benchmarks	17
5.3 Containerizzazione con Docker	17
5.4 Chiamate parametrizzabili sui campi del documento	17
5.5 Postprocessing	17
Conclusioni e sviluppi futuri	19

2	Indice
---	--------

Ringraziamenti	21
--------------------------	----

ELENCO DELLE FIGURE

"If a machine is expected to be infallible, it cannot also be intelligent"
— Alan Turing

INTRODUZIONE

Questa tesi vuole essere un riassunto delle attività svolte in un tirocinio curriculare presso l'azienda QI-LAB di Firenze nel periodo *marzo - giugno 2019*. L'offerta di tirocinio prevedeva l'apporto di migliorie, in termini di accuratezza ed efficienza, riguardanti una libreria per il riconoscimento ottico dei caratteri di alcuni campi di documenti d'identità, attualmente solo italiani. La libreria, denominata QI-OCR, era inizialmente il risultato del lavoro compiuto da un precedente tirocinante, *Emilio Cecchini*, che si era occupato dell'implementazione del framework di base, sul quale ho avuto il piacere di lavorare io stesso. QI-OCR viene distribuita come pacchetto *pip* ed è stata scritta con linguaggio Python, in versione 3.6.8, anche se compatibile con la versione minor successiva, ovvero la 3.7.x, nonché l'ultima versione stabile rilasciata, nella data in cui mi trovo a scrivere questo testo. QI-OCR prevedeva la suddivisione in 4 moduli principali:

1. **Preprocessing:** Si occupa della localizzazione del documento all'interno dell'immagine, del ritaglio dei contorni e del raddrizzamento, tramite l'algoritmo di template matching *SIFT*.
2. **Docfields:** Si occupa di effettuare un ritaglio statico, mediante coordinate predeterminate, dei campi d'interesse del documento già centrato.
3. **OCR:** Si occupa di effettuare una binarizzazione dell'immagine mediante un algoritmo di sogliatura globale *ad-hoc* e di fornire l'immagine prodotta in input a un motore di OCR.
4. **Postprocessing:** Si occupa di applicare tecniche di analisi sintattica e semantica sui risultati ritornati dal software OCR.

Il lavoro da me svolto riguarda il miglioramento della libreria QI-OCR, mediante tecniche che verranno descritte successivamente. A tal proposito, la tesi è suddivisa in cinque capitoli. Il capitolo 1 presenta una descrizione dei prerequisiti matematici necessari a comprendere correttamente le spiegazioni dei metodi presenti nei capitoli successivi. Il capitolo 2 riguarda lo studio e l'implementazione di tecniche di binarizzazione e

segmentazione di immagini, ovvero soluzioni necessarie per la "pulizia" di una generica immagine contenente testo, che permettano di ottenere idealmente il solo testo di colore nero e tutto il resto di colore bianco. Il capitolo 3 riguarda invece lo studio di algoritmi di segmentazione basata dal riscontro sul modello, in cui vengono introdotti metodi alternativi al ben noto SIFT per l'individuazione della posizione di un documento all'interno di un'immagine. Il capitolo 4 introduce l'applicazione di tecniche avanzate di *text-detection* con lo scopo di colmare problemi derivanti da documenti che non rispettano alcun tipo di template prefissato, come ad esempio la vecchia carta d'identità cartacea italiana, che risulta comunque rilasciabile, in alternativa alla nuova carta d'identità elettronica italiana, in casi di estrema esigenza o per tutti i cittadini che non possono recarsi per motivi di handicap presso il municipio. Infine, il capitolo 5 descrive ulteriori migliorie e innovazioni introdotte, sia per quanto riguarda il lato sviluppo/distribuzione (containerizzazione, benchmarks, ...), sia per quanto riguarda l'esperienza utente (OCR parametrizzabile sui campi del documento, implementazione di diversi motori di OCR, ...).

PREREQUISITI MATEMATICI NELL'AMBITO DELLA COMPUTER VISION

1.1 IMMAGINI DIGITALI

Un'immagine digitale I con risoluzione $m \times n$ è una matrice di valori interi di dimensione $n \times m$, che può essere matematicamente interpretata come una funzione semplice (nè iniettiva nè suriettiva) $I: \mathbb{N} \times \mathbb{N} \supseteq X \rightarrow \mathbb{P}^c$, che si occupa di mappare una coppia ordinata $(u, v) \in \mathbb{N} \times \mathbb{N}$, con $u, v \in \mathbb{N}$, in una c -upla $p = (p_1, p_2, \dots, p_c)$, in cui ciascun valore p_i di p rappresenta il valore del pixel associato nel canale i , con $0 \leq p_i \leq 2^b - 1$, dove b è il numero fissato di bit utilizzati per rappresentare ciascun pixel e 2^b indica il massimo numero di colori o di livelli di grigio utilizzabili. Per esempio, ipotizzando di utilizzare canali colore standard, nel caso in cui si lavori con immagini in scala di grigi si avrà $c = 1$, mentre nel caso in cui si lavori con immagini *RGB* si avrà $c = 3$ ($c = 4$ per immagini *RGBA* o *CMYK*). Inoltre, supponendo di utilizzare una profondità pari a $b = 8$ sarà possibile specificare interi compresi tra 0 e 255 per i valori di ciascun pixel. Infine, indichiamo la risoluzione dell'immagine come la coppia $(m, n) \in \mathbb{N} \times \mathbb{N}$: $(n, m) = \max_{(x,y) \in X} (x, y)$.

1.2 CONVOLUZIONE

La *convoluzione* è un'operazione alla base dell'*image processing*, grazie alla quale è possibile andare ad analizzare ed eventualmente accentuare o alleviare diversi aspetti di una data immagine, come la sfocatura, la nitidezza, i contorni e molto altro. L'operazione, che viene descritta nella forma discreta, prende dunque in input due immagini digitali, ovvero l'immagine originale I (normalmente un singolo canale) e l'immagine K , denominata *kernel* o *elemento strutturante*, e intuitivamente si occupa di far scorrere la matrice K sull'immagine I , generalmente a partire dal bordo

in alto a sinistra, effettuando una somma pesata dei valori di I dati dalle proiezioni delle posizioni correntemente analizzate dalla matrice K , in cui i pesi sono dati proprio dagli elementi di K . Solitamente la dimensione della matrice K è molto minore della dimensione dell'immagine originale e spesso K è una matrice quadrata $n \times n$, con n dispari. Più formalmente, l'operazione di convoluzione $I * K$ in un punto (i, j) è data da:

$$I'(i, j) = \sum_{x=1}^n \sum_{y=1}^n (I(i + \left\lceil \frac{n}{2} \right\rceil - x, j + \left\lceil \frac{n}{2} \right\rceil - y) \cdot K(x, y)) \quad (1.1)$$

Nel caso in cui nella formula 1.1 i due $-$ venissero sostituiti da due $+$ si otterrebbe la formulazione della *correlazione* $I \otimes K$ in un punto (i, j) :

$$I'(i, j) = \sum_{x=1}^n \sum_{y=1}^n (I(i + \left\lceil \frac{n}{2} \right\rceil + x, j + \left\lceil \frac{n}{2} \right\rceil + y) \cdot K(x, y)) \quad (1.2)$$

Convoluzione e *correlazione* possono risultare equivalenti se per passare da un'operazione all'altra si effettua un ribaltamento sui due assi, orizzontale e verticale, del filtro (graficamente K viene ruotata di 180°). Dunque le due operazioni risultano identiche nel caso in cui la matrice K sia simmetrica. In particolare, si preferisce utilizzare la convoluzione quando sono necessarie le proprietà commutativa, associativa e distributiva. Inoltre, dato che spesso il risultato delle operazioni descritte viene utilizzato come valore di intensità di pixel, tale valore viene normalizzato, dividendolo per la somma dei pesi del filtro. La complessità computazionale delle operazioni, data un'immagine di dimensione $n \times n$ e un kernel di dimensione $m \times m$, è piuttosto elevata, dato che richiede un numero di moltiplicazioni pari a $m^2 \cdot n^2$ e altrettante somme. Di seguito un esempio di convoluzione:

1.3 OPERAZIONI MORFOLOGICHE

BINARIZZAZIONE DI IMMAGINI

2.1 INTRODUZIONE

2.2 SOGLIATURA GLOBALE

2.3 SOGLIATURA ADATTATIVA (O LOCALE)

2.4 SOGLIATURA DI OTSU

2.5 ALGORITMO DI SOGLIATURA PROPOSTO

3

TEMPLATE MATCHING

3.1 INTRODUZIONE

3.2 SIFT

3.3 SURF

3.4 AKAZE

3.5 BRISK

3.6 HED

TEXT DETECTION

4.1 INTRODUZIONE

4.2 EAST

4.3 ALGORITMO PROPOSTO

5

ULTERIORI MIGLIORIE

5.1 MOTORI OCR

5.2 BENCHMARKS

5.3 CONTAINERIZZAZIONE CON DOCKER

5.4 CHIAMATE PARAMETRIZZABILI SUI CAMPI DEL DOCUMENTO

5.5 POSTPROCESSING

CONCLUSIONI E SVILUPPI FUTURI

RINGRAZIAMENTI

BIBLIOGRAFIA

- [1] Wilhelm Burger, Mark J. Burge - *Digital Image Processing - An Algorithmic Introduction Using Java, Second Edition*