

Západočeská univerzita v Plzni
Fakulta aplikovaných věd
Katedra informatiky a výpočetní techniky

Bakalářská práce

Automatická anotace obrázků

Místo této strany bude
zadání práce.

Prohlášení

Prohlašuji, že jsem bakalářskou práci vypracovala samostatně a výhradně s použitím citovaných pramenů.

V Plzni dne 24. května 2017

Kateřina Kratochvílová

Abstract

The text of the abstract (in English). It contains the English translation of the thesis title and a short description of the thesis.

Abstrakt

Text abstraktu (česky). Obsahuje krátkou anotaci (cca 10 řádek) v češtině. Budete ji potřebovat i při vyplňování údajů o bakalářské práci ve STAGu. Český i anglický abstrakt by měly být na stejné stránce a měly by si obsahem co možná nejvíce odpovídat (samozřejmě není možný doslovný překlad!).

Obsah

1	Úvod	6
2	JEC Joint Equal Contribution	7
2.1	Příznaky	7
2.1.1	Barva	7
2.1.2	textura	9
2.2	Vzdálenosti	10
2.3	Kombinace vzdáleností	10
2.4	Přenesení klíčových slov	11
3	Testovací databáze	13
3.1	iaprtc12	13
3.2	ČTK	13
3.3	ESP	14
4	Návrh systému	15
5	Použité programové prostředky	16
5.1	OpenCV	16
6	Vyhodnocení výsledků	17
6.1	Srovnání výsledků	17
6.1.1	iaprtc12	17
7	Závěr	18
8	Uživatelská dokumentace	19
	Literatura	20

1 Úvod

V dnešní době, kdy je svět přesycen obrázky v digitální podobě, není vůbec snadné nalézt obrázek zobrazující požadovaný obsah. Naneštěstí počítače nedokáží vnímat obraz jako lidé, vnímají totiž obrazy jako sérii binárních informací. Přitom počítače a jejich práce s obrazy by se dala využít v mnoha oborech jako je lékařství nebo doprava. Na základě toho vyplouvá na povrch problém jak spravovat digitální obrázky a efektivně mezi nimi vyhledávat. Prostřednictvím klíčových slov přiřazených k obrázkům se dá problém vyhledávání zjednodušit. Přiřazení klíčových slov probíhá pomocí procesu automatické anotace obrázků. Klíčová slova přiřazená k obrázku by měla vyjadřovat jeho obsah (například les, strom). Při reálném použití můžeme ovšem narazit na problém při zadávání abstraktních slov, například šťastná rodina.

Pro automatickou anotaci obrázků se používá strojové učení. Můžeme ji rozdělit na dvě části. V první části získáme klíčové příznaky ve druhé už je samotná anotace, tedy přidělení klíčových slov. Abychom tento postup mohli provést v praxi, musíme nejdřív klasifikátor natrénovat pomocí trénovací množiny. Trénovací množina je množina obrázků, která již má ke každému obrázku přidána metadata s klíčovými slovy připravenými od lidí. Vybrané obrázky v trénovací množině musí být různorodé, aby anotace probíhala správně. Pojem automatická anotace obrázků je jednoduše řečeno proces, při kterém jsou k obrázku automaticky přiřazena metadata, která obsahují klíčová slova.

Cílem práce je navrhnout a implementovat software umožňující automatickou anotaci obrázků. Konkrétně se bude zabývat metodou JEC. Práce bude využívat nízkourovňové příznaky, konkrétně barvu a texturu. Metodu budeme zkoušet na standardních datech, následně se budeme snažit výsledky vylepšit, a porovnat s další metodou a literaturou.

2 JEC Joint Equal Contribution

Tato metoda je založena na hypotéze, že podobné obrázky mají podobná klíčová slova. Pomocí metody hledání nejbližších sousedů (dále jen KNN) najdeme K nejpodobnějších obrázků. Přičemž klíčová slova od jednotlivých sousedů jsou posuzována odlišně a to právě na základě toho o kolik se s testovaným obrázkem liší. Metoda je postavena na dvou typech příznaků - barevných a texturových. [1]

2.1 Příznaky

Barva a textura jsou považovány za dva nejdůležitější nizkoúrovňové příznaky pro obrázkovou reprezentaci. Nejběžnější barevné deskriptory jsou barevné histogramy, které jsou často využívány pro porovnávání a indexování obrázků, zejména z důvodu jejich efektivnosti a snadného výpočtu. K vytvoření texturových příznaků se používají Haarovy a Gaborovy wavelety a to především z důvodu že jsou efektivní při vytváření rozptýlených diskriminativních obrázkových rysů.

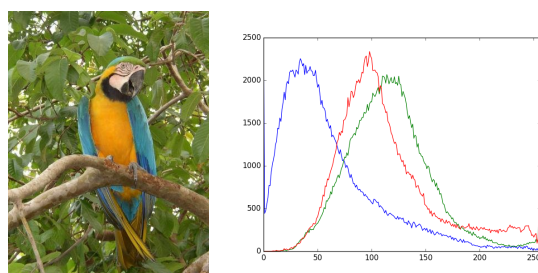
2.1.1 Barva

U digitálního obrazu je barva reprezentovaná n -rozměrným vektorem. Jeho velikost a význam jednotlivých složek (tzv. barevných kanálů) závisí na příslušném barevném prostoru. Počet bitů použitých k uložení buď celého vektoru nebo jeho jednotlivých složek se nazývá barevná hloubka (totožně bitová hloubka). Obvykle se můžeme setkat s hodnotami 8, 12, 14 a 16 bitů na kanál.

V použité metodě získáme vlastnosti z obrázků ve třech rozdílných barevných prostorech: RGB, HSV a LAB. RGB (Red, Green, Blue) je nejpoužívanější barevný prostor pro zachycení obrázu nebo jeho zobrazení. Oproti tomu HSV (Hue, Saturation and Value) se snaží zachytit barevný model tak jak ho vnímá lidské oko, ale zároveň se snaží zůstat jednoduchý na výpočet. Hue znamená odstín barvy (měří se jako poloha na standardním barevném kole $0^\circ - 360^\circ$), saturation je sytost barvy (množství šedi v poměru k odstínu 0% šedá barva - 100% plně sytá barva) a value je hodnota jasu nebo také množství bílého světla (relativní světlost nebo tmavost barvy). Některé

kombinace hodnot H, S a V mohou dávat nesmyslné výsledky. RGB je závislý na konkrétním zařízení, nemůže dosáhnout celého rozsahu barev, které vidí lidské oko, zatímco barevný model LAB je schopný obsáhnout celé viditelné spektrum a navíc je nezávislý na zařízení. L (ve zkratce LAB) značí Luminanci (jas dosahuje hodnot 0 - 100, kde 0 je černá a 100 je bílá). Zbylé A a B jsou dvě barvonosné složky, kdy A je ve směru červeno/zeleném a B se pohybuje ve směru modro/žlutém.

Pro RGB, HSV i LAB použijeme barevnou hloubku 16 bitů na kanál histogramu v jejich příslušném barevném prostoru.



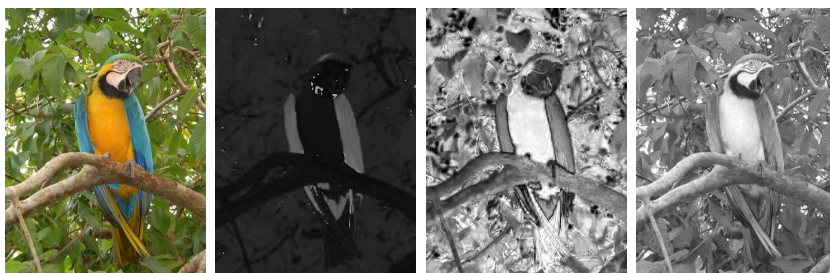
Obrázek 2.1: RGB histogram - zastoupení jednotlivých složek v obrázku



Obrázek 2.2: Barevný prostor RGB a jeho jednotlivé složky v pořadí R, G, B



Obrázek 2.3: Barevný prostor LAB a jeho jednotlivé složky v pořadí L, A, B



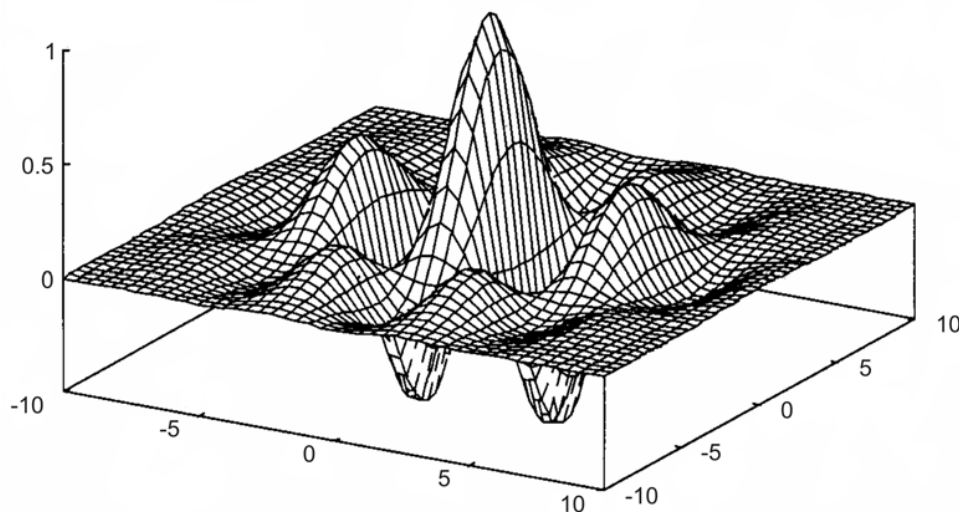
Obrázek 2.4: Barevný prostor HSV a jeho jednotlivé složky v pořadí H, S, V

2.1.2 textura

Jako reprezentace textur budou použity Gabor a Haar wavelety.

Gabor

Gabor filter je lineární filter používaný pro detekci hran. Frekvence a orientace reprezentující Gabor filter je podobná lidskému vnímání a jsou zvláště vhodné pro reprezentaci textury a rozlišování.



Obrázek 2.5: Gabor vlnka je tvořena kombinací dvou cosinových funkcí, pro každou osu, které mají v základu odlišné frekvence.[2]

U obrazku bylo: Gáborova vlnka je tvořena vynásobením dvou cosinovek (které mají obecně různé frekvence) pro každou osu a výsledek se pak násobí dvourozměrnou Gaussovou funkcí

určuje se velikost filtru (značeno jako $ksize$). Sigma standartní odchylka Gausovy funkce použité v gaborově filtru. theta je orientace normálu na

paralelní pruhy gaborovy funkce. lambda je vlnová délka sinusového faktoru (vlnova delka sinusovky), gamma - prostorový poměr stran (elipsicityda), když je gamma 1 je gaussian do kruhu, psi je fazovy posun (určije jestli nám vrátí realnou nebo imaginární část),

$$g(x, y; \lambda, \theta, \psi, \sigma, \gamma) = \exp\left(-\frac{x^2 + \gamma^2 y^2}{2\sigma^2}\right) \exp\left(i\left(2\pi\frac{x}{\lambda} + \psi\right)\right) \quad (2.1)$$

Každý obrázek bude filtrován na třech škálách (lambda = 2, 2*sqrt(2), 4) a čtyřech orientacích (theta = 0, 45, 90 a 135 stupňů)

Gabor filtr reaguje na hrany a texturové změny.

Obrázky jsou filtrovány za použité reálné části z různých odlišných Gabor filtrů jader. Průměrný a rozptyl filtrovaných snímků jsou pak použity jako funkce pro klasifikaci, která je založena na nejméně kvadratické chyby pro jednoduchost.

TODO Dohledat Haar a Gabor wavelety, přidat vzorečky a zase klidně i obrázky

2.2 Vzdálenosti

vzdálenost vs podobnost

K určení příslušné vzdálenosti se můžeme setkat se čtyřmi měřítky vzdálenosti pro histogramy a rozdělení Kullback-Leibler divergence KL - divergence, χ^2 statistika, L1 - vzdálenost a L2 - vzdálenost. Na RGB a HSV je nejlépejší použít L1 zatímco pro LAB je nejvhodnější KL - divergence.

Problém s KL - divergencí nastává pouze tehdy, když se histogramy nebudou shodovat v nulách. Jeden předpoklad pro fungování tohoto vzorce je totiž že když je $Q(i) = 0$ tak zároveň musí být i $P(i) = 0$.

$$D_{KL}(P||Q) = \sum_i P(i) \log_e\left(\frac{P(i)}{Q(i)}\right) \quad (2.2)$$

2.3 Kombinace vzdáleností

Nejrozumnějším přístupem ke zkombinování vzdáleností od různých deskriptorů je aby jednotlivé vzdálenosti přispívali rovnocenně. Z tohoto důvodu je potřeba vzdálenosti přeškálovat na jednotné měřítko.

Označme si I_i jako i -tý obrázek a řekněme, že máme N jeho příznaků f_i^1, \dots, f_i^N . Nadefinujme si $d_{(i,j)}^k$ jako vzdálenost mezi příznaky f_i^k a f_j^k . Chtěli bychom zkombinovat všechny vzdálenosti příznaků mezi obrázky I_i a I_j tedy

$d_{(i,j)}^k$, $k = 1, \dots, N$. Vzdálenosti nám ale v praxi nevyjdou tak aby měli stejný poměr na výsledku, proto předtím než vzdálenosti zkombinujeme musíme je normalizovat do jednotné formy. Získáme maximální a minimální hodnotu pro každý příznak a na základě toho hodnotu přeškálujeme na interval od 0 do 1. Jestliže označíme přeškálovanou vzdálenost jako $\tilde{d}_{(i,j)}^k$ následně můžeme označit kompletní vzdálenost mezi obrázky I_i a I_j jako (2.3) Joint Equal Contribution (JEC).

$$JEC = \sum_{k=1}^N \frac{\tilde{d}_{(i,j)}^k}{N} \quad (2.3)$$

2.4 Přenesení klíčových slov

Pro přenesení klíčových slov používáme metodu, kdy přeneseme n klíčových slov k dotazovanému obrázku \tilde{I} od K nejbližších sousedů z trénovací sady. Mějme $I_i, i = 1, \dots, K$, těchto K nejbližších sousedů seřadíme podle vzrůstající vzdálenosti (tzn. že I_1 je nejvíce podobný obrázek). Počet klíčových slov k danému I_i je označen jako $|I_i|$. Dále jsou popsány jednotlivé kroky algoritmu na přenesení klíčových slov.

1. Seřadíme klíčová slova z I_1 podle jejich frekvence výskytu v trénovací sadě.
2. Ze všech $|I_1|$ klíčových slov z I_1 přeneseme n nejvýše umístěná klíčová slova do dotazovaného \tilde{I} . Když $|I_1| < n$ pokračujte na krok 3.
3. Seřadíme klíčová slova sousedů od I_2 do I_K podle dvou faktorů
 - (a) výskytu v trénovací sadě s klíčovými slovy přenesených v kroku 2
 - (b) místní frekvence (tj. jak často se vyskytují jako klíčová slova u obrázků I_2 až I_K). Vybereme nejvíce vyskytující $n - |I_1|$ klíčových slov převedených do \tilde{I} .

Tento algoritmus pro přenos klíčových slov je poněkud odlišný od algoritmů, které se běžně používají. Jeden z běžně užívaných funguje na principu, že klíčová slova jsou vybrána od všech sousedů (se všemi sousedy je zacházeno stejně bez ohledu na to jak jsou danému obrázku podobní), jiný užívaný algoritmus k sousedům přistupuje váženě (každý soused má jinou váhu) a to na základě jejich vzdálenosti od testovaného obrázku. Při testování se ovšem

ukázalo, že tyto přímé přístupy přináší horší výsledky v porovnání s použitým dvoufaktorovým algoritmem pro přenos klíčových slov.

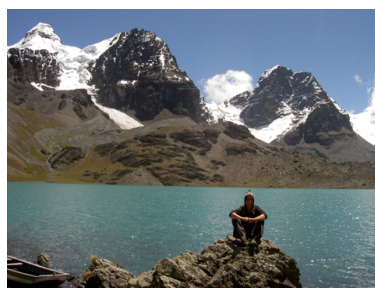
V souhrnu použitá metoda je složenina ze dvou složenin a to obrázkové vzdálenosti (JEC) a výše popsaným algoritmem na přenášení klíčových slov.

3 Testovací databáze

Pro natrénování a následné testování byla použita data z databází ČTK, ESP a IAPRC. Kolekce obrázků na natrénování musí být pečlivě vybrána aby zahrnovala co možná největší okruh z různých témat.

3.1 iaprtc12

Data iaprtc12 obsahují 20 000 obrázků ve formátu jpg s celkovým počtem 291 klíčových slov. Ke každému obrázku jsou přiložena metadata ve formátu XML, která obsahují informace o obrázku v různých jazycích. Kromě angličtiny je tam i například španělština nebo němčina. V metadatatech ovšem nenajdeme klíčová slova tak jak bychom si je představovali, ale v různých tagách nalezneme například titulek obrázku, který může vypadat například The Plaza de Armas, a v tagu description je například a woman and a child are walking over the square. Klíčová slova musela být z těchto dat získána. Spolu s databází jsme získali klíčová slova byla extrahována. K jednomu obrázku je v průměru přiřazeno 5.7 klíčových slov. Pro trénování bylo použito 17664 obrázků, na následné testování jich bylo použito 1960.



Obrázek 3.1: Ukázka obrázku s klíčovými slovy: front lake man mountain rock sky summit

3.2 ČTK

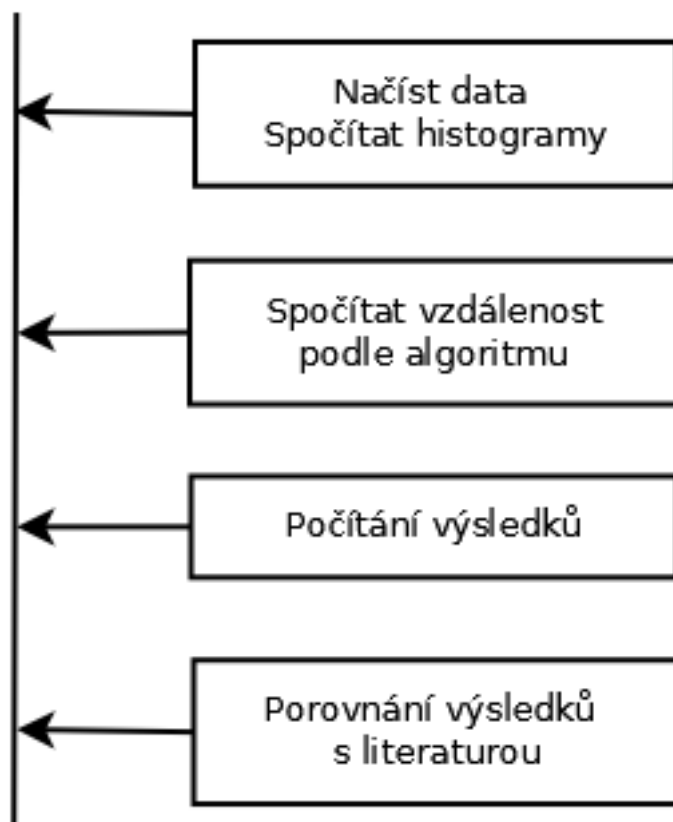
Data od ČTK obsahují 3 383 obrázků. Ke každému obrázku jsou přidána metadata ve formátu XML, která obsahují klíčová slova a další informace o obrázku.

3.3 ESP

Data od ESP obsahují 67 796 obrázků ve formátu jpg. Ke každému obrázku je přiřazen soubor ve formátu desc, který obsahuje anglické anotace.

4 Návrh systému

Systém byl navržen jako modulový a to z důvodu snadné obměny některé z částí, což je výhodné zejména pokud bychom potřebovali například spočítat vzdálenosti vektorů podle jiného algoritmu.



Obrázek 4.1: Návrh systému

5 Použité programové prostředky

Program byl navržen na operační systému Linux. Jako programovací jazyk byl zvolen Python a to z důvodu jeho jednoduchého použití, což je na prototyp, jako je tento velice výhodné na časovou náročnost. Program využívá knihovnu OpenCV 3.1.

5.1 OpenCV

OpenCV (Open source computer vision) je knihovna vydávána pod licencí BSD a je volně k dispozici jak pro akademické účely, tak pro komerční použití. Je vhodná pro použití v C++, C, Python a Javě. Podporuje operační systémy Windows, Linux, Mac OS, iOS a Android.

Knihovna byla navržena pro výpočetní efektivitu v oblasti počítačového vidění a zpracování obrazu se zaměřením na zpracování obrazu v reálném čase. Z důvodu optimalizace byla napsána v C/C++.

Knihovnu OpenCV je možné stáhnout na adrese: <http://opencv.org/>

6 Vyhodnocení výsledků

Zpracování výsledků probíhá jako porovnání anotací přidělených člověkem s anotacemi přidělenými klasifikátorem. Označme si w_{auto} jako počet obrázků, kterým bylo dané slovo přiřazeno klasifikátorem, w_{human} počet obrázků, kterým bylo dané slovo přiřazeno člověkem a $w_{correctly}$ jako počet obrázků, kterým bylo slovo přiřazeno správně. U klasifikátorů se počítá precision (přesnost) a recall (úplnost) pro každé slovo v testovací sadě.

Recall (6.1) je počet obrázků správně anotovaných s daným slovem děleno počtem obrázků, kterým bylo toto slovo přiděleno v anotaci člověkem. Precision (6.2) je počet správně anotovaných obrázků s tímto slovem děleno celkovým počtem anotovaných obrázků s tímto slovem (správně nebo ne). [6]

$$Rec = \frac{w_c}{w_h} \quad (6.1)$$

$$Prec = \frac{w_c}{w_a} \quad (6.2)$$

6.1 Srovnání výsledků

6.1.1 iaprtc12

	Přesnost (%)	Úplnost (%)	Nenulových slov
RGB	14	8	167
LAB	12	7	148
HSV	16	10	181
RGB, LAB, HSV	17	11	178
POEM	0	0	0
RGB, LAB, HSV, POEM	0	0	0
Barevný POEM	0	0	0

7 Závěr

V teoretické části byly popsány nízkourovňové příznaky barva a textura. Byla rozebrána metoda JEC, která bude v bakalářské práci implementována. Seznámili jsme se s knihovnou OpenCV, prostudovali obrázky a přiložená metadata od ČTK, ESP a iaprtc12.

8 Uživatelská dokumentace

popsání jak vypadá zdrojový soubor který to zere, nejdriv cesta k souboru
a pak jeho klicovy slova

Literatura

- [1] AMEESH MAKADIA, S. K. V. P. A new baseline for image annotation. *Commun. ACM*. July 1961, 4, 7, s. 321. ISSN 0001-0782. doi: 10.1145/366622.366644. Dostupné z: <http://doi.acm.org/10.1145/366622.366644>.
- [2] CRUSE, H. Neural Networks as Cybernetic Systems. *Commun. ACM*. July 1961, 4, 7, s. 321. ISSN 0001-0782. doi: 10.1145/366622.366644. Dostupné z: <http://doi.acm.org/10.1145/366622.366644>.
- [3] HOARE, C. A. R. Algorithm 64: Quicksort. *Commun. ACM*. July 1961, 4, 7, s. 321. ISSN 0001-0782. doi: 10.1145/366622.366644. Dostupné z: <http://doi.acm.org/10.1145/366622.366644>.
- [4] *Class Graphics2D* [online]. Oracle, 2016. [cit. 2016/03/09]. Java SE Documentation. Dostupné z: <https://docs.oracle.com/javase/7/docs/api/java/awt/Graphics2D.html>.
- [5] KNUTH, D. E. *The Art of Computer Programming, Volume 2 (3rd Ed.): Seminumerical Algorithms*. Addison-Wesley Longman Publishing Co., Inc., 1997. ISBN 0-201-89684-2.
- [6] V. LAVRENKO, J. J. R. A model for learning the semantics of pictures. *Commun. ACM*. July 1961, 4, 7, s. 321. ISSN 0001-0782. doi: 10.1145/366622.366644. Dostupné z: <http://doi.acm.org/10.1145/366622.366644>.