

Àâòïìàòè÷åñêàÿ êëàññèôèêàöèÿ øðèôòà ñòàðïïå÷àòíûõ òåêñòïâ

Àâòîð Àëåäèìø Ñåðåâåâà÷ þæèéâ
25.09.2009 å.
Íñéäåíå Íáíâéåíèå 05.10.2009 å.

The article describes an approach of font classification for old-printed texts and manuscripts. This is necessary if we have an unknown antique book and want to determine the printing press, place and time of issue, etc. The proposed approach is based on partitioning the page into rows and letters, and then makes a consistent comparison with the available font samples. Also the technique for the automatic font samples generating is described.

Âååääåíèå

Êiāää à ðóöèè -äéïäåéö, çàìèìäþùåíöñý
ñòàðïíå -àòiúù è òåëñòàìè, iïïäääååò íåèçååñòiáÿ ñòàðèíàÿ êíèää, eëéáí åå
ôðåäåíäiòû, åíçíèéååò iöfåéäià - èåé ïðåäååéèòù åðåíÿ è iåñòi åå iå -àòè, åñëè
ýòà èíöiðåéöÿ ií êåéèi -òi iöde -éíàí iåäñòòiíà. iñüöiúé èññéåäiåàòåëü, íåñiiíäií,
iñæåò iödeiåðíi ïðåäååéèòù ýòi, iññíåüåäyñü íå òåö êíèääò è ðóëítèñýö, ðòi ií
åèääéé ðäiåå, ií åäðåèëüíüé áiàéec øðèòòà, ñòèëÿ íaièñàíèÿ òðååáóåò ååñüìà
êòiñiøèéàíé ðäiåòù è íaiàéüö çàòðåò åðåíäié.

lîiñ÷ü â ðåøâíèè òàèéíè çääà÷è lîiñëà áû
ñiäöèäüäiy êiñüþòåðíay ñeñòåìà, êiòîðäy óíååò ñðàâíèåàòü áóêåû íåèçååñòíí
ñòðåíèöü ñ lîiñäî÷èñëåííùîè íåðäçöäiè ðøèòòòíâ èç ðàçíûò èñòî÷íèñâ, êiòîðûå
ôðåíÿòny â åå áàçå. Èñòîäy èç ýoññ, lîiññ ñòîðèñíåàòü lîiñññáû òðåáîâàíéy ê
òàèéíè ñeñòåìà:

Êæäüé íáðåçåö ðøðéòðà äíéæåí èíåðöü íáðæéíðòíåðéþ, òåéóþ êåé; ðéíïåðåðéý, áíðíä/ñòðåðåíä, áíä éçcääíéý è ó.ä.

Ãîëæíà áûòü âïçÿæíñòü ëåâëåïâ çàíåñåíèÿ â áàçó äàííûõ
íïðåïâ íàðàçöà øðèòöà.

Âûääà÷à ðåçóëüòàòîâ äïëæíà áûòü ðàíæèðîâàííé ïî êàæäîíó
jóíéòò jåòàèíòòjåòèè, låíðøjåð:

Ø Òøïïâðàôèëÿ - «Êðàêîâ» (87%), «Êèðèëë è îåôïäëé» (75%);

^o
Ãïä èçääíèÿ - 1764 (95%), 1795 (84%); ...

Êëàññèôèêàöèÿ øðèôòà íå äîëæíà çàíèìàòü íííâî âðåìåíè.

Íáçîð
ëèòåðàòóðû

Ííèññúâåâåìàÿ íðíáéâà à ðííññéòñý ê ééàññó
çàääà÷ ðâññíçíåâåíèÿ láðäçïà, à èíàííí - OFR (Optical Font Recognition). Ýòà íáéàñòü
ññáé÷àñ ãéòðéâåí ðâçåéâååòñý è íáðíäèò ññåíà íðíéíáíéå áî ííññéò ëíííåð÷åññéò
íðíáðéòåð, ðàééò ëàé FineReader, CuneiForm è ííññéò åððåéò. Ê ñíçæàéåíþ, êéàñññéòééàòíð
øðéòòíâ åôíäÿùéé â ñññòåå ýòëò ññòòåì, íå ííðíäèò åéÿ ðåðåíéÿ íåðåé çàääà÷è,
ò.é. «çàòí÷åí» ííä ååññüíà óçéóþ íáéàñòü - ííðíäðåòù íáéáíéåå ííðíäÿùóþ
ååðíéòòðò è ñòëëü íå÷åðåíéÿ èç í÷åíü íåðåíè÷åíííå íááíðà. íáú÷íí ýòí:
ííññéðéííüé øðéòò (Courier), øðéòò ñ
çàñå÷éàìè (Times New Roman), øðéòò ååç
çàñå÷åé (Arial), íëþñ èò ííññéðéííüé: bold è italic. Äéÿ çàääà÷è
ðâññíçíåâåíèÿ ðâññòåí (OCR) áíéüòååí è íå ðåðåååðåñý, ííòííó ÷òí ííàååéÿþþåå
÷èññéí ññåðåíåííüó ãíéóíåíóíà èñííëüçþþò åúøåíåðå÷èññéåííüå øðéòòù.

Íáðaðòèíý è éèðaððaðóða - á Íoáéëéàöðéýð
âññòðåð-á-ðòðvñ äåà ïññíññúð ïäðññäà è ðàññíñçàññàíèþ øðèðòða:

Âüääëäíéå
ääëäüüüöö ïðeçìäéïå èç ðåéñòåüüöö áéëéïå, òàéèöö éëæé ñëïåà, ñòðíéè è èö
ääëüüåéøèé äíåëëç. Ýöi ïïçåéëýåöö ïðåååëëýöü ñòëëü íå-÷åðòåíéÿ ððéòòå (bold, italic), ååí ðàçìåð,
íåéé-÷éå çåñå-÷åé. Éåéé ïðååéëç, ðåééïå ïïðöñïå íå ððååáóåò áåçü íåðåçöñå, íïæåò
ïðåååëëýöü íåñééüüéí ïñííåüüöö éëæññíå ððéòòåíå è íå çåäéñéò ìò ýçûéå (â
áéëüøéíñòåå ñéó-÷ååå), ÷òí åääëäåò ååí ïòëíåëëüüíùí äëëÿ èñííëüçíåàíéÿ å ñèñòåíàò
ðåñíçíååäíéÿ òåññòå OCR.

Ê iðèìâðó, â ðàâìòå Shi, Pavlidis [1] äëÿ
 ðàñïïçíåâàíèÿ ðøðøòà èñïïëüçóåòñÿ ðÿä åéïåàëüûò äëÿ âñåé ñòðàíèòû iðèçíàéîâ,
 òàéèò ëåê: åèñòîðàìà ðàñïðåâåëáíèÿ äééí ñéïå (â ièéñåëÿô), ååéè÷-èà
 iàæñòðî-ûò, iàæåóéâåíûò èíòåðâåéîâ. Yong Zhu [2] iðåâæàåâåò ðàñïïçíåâàíèå
 ðøðøòà, ïñïâåííà íà áíàééçå òåéñòóðû
 èçíåðàæåíèÿ âñåé ñòðàíèòû. Äéÿ èçåéå÷-åíèÿ ýòèò iðèçíàéîâ èñïïëüçóåòñÿ
 iííåéàíæüûò ôéëüòð Áåâîðà.

Âúääåéåíèå
eëíèåéüíûö ïðèçíàéîâ. Çääñü óæå àíàéèçèðóþöñý ìöääéüíûå áóêåû è eð ýëåíåíòû. Òàééé
ïíäöñä ÷ðâñöðåéòåéåí è íååíéüøèì èçíàíåíèÿì á ðøðòå, ÷òî ïçáíéÿåò õïðíøí ðàçääéÿöü
áéèçéèåí ïí áðòåíèþ ððéòööù è áíéüøå âññåñí ïíäöñäèò äéÿ íáøåé çäåà-è.

Â ðàáìòå Cooperman'à [3] èññéåäóþöñý ïöåíèé ñâíéñòå øðëèôòå àëÿ ñèñòåì OCR. Â êà÷ñòå âåðåêòîðíà ýòëõ ñâíéñòå èññéüçóþöñý èëåäëüíûå ïðèçíàèè, òåêèå êàë: íàëè÷ëå çàñå÷åê, iëîòíññòü è ò.ä. Â ðàáìòå Zramdini, Ingold [4] ïðääëåäåàòñý ñòaoðèñòå÷åññéé ïïðöïà àëÿ ðåøäíèéý çàäà÷è êëàññèòèàòëè øðëèôòå, íñííàáííûé íà ãûäåëåíèé èëåäëüíûö ïðèçíàèîâ. Ñòîæèé íàòïà èññéüçóåðñý è å ðàáìòå [5]

Íðåääàðèòåëüíàÿ
Íáðàáîòëà

Ñòàðññà÷àòíûé è, à ïññááííññòè, ñòàðéííûé
ðóéííà÷àòíûé (hand-writing) òâéñò ìèååò ðÿä ïññááííññòåé, êîòìðûå íà ïïçâíëýþò íaiðýìóþ íðèìåíèòü
ìåòíàú, íðåäëääååíùå á ïøáéëéàðëéý ïi OFR. Â ÷àñòññòè, íåæáóéååíùå è íåæñòðí÷íùå ðàññòëýéý áàðüèðóþòñý á
äññòàðí÷í ðøðíëéò ððåååéëåð ååæå íà ïññé ñòðåíèòå. Êññà òíñ, òåññò è íååò
áíëüøíå êíèé÷åñòå áåæíûò äëÿ éeññøêåöéè ððèòå ãæàéðòè÷åññèò ÿéåíåíòíå, íà
ñòðåíèòå ÷àñòí áñòü ðèññóíéè, íñíåòéè. Òàéæå èçíåðåæåíèå ñíåðæòåò áåôåéòü,
åñçååííùå äëéòåëüíùå ððåíåíéåí, êîòìðûå íñæíí ññéíàíí ðàçääéëòü íà 2 òéíà. È
íåðåííò íòíåñàí ååðååéòù ñòðåíèò íééæè, íiyâéåøèåñý á ðåçóëüòåòå äíëäíàí
ñòðåíàíéÿ, äåéñòåèÿ áéëæíññòè,
òåííåðåòðû, íñðåæåíèÿ ððéåééíò íòååëüíûò ñòðåíèò, áûöååòåíèå áóéå,
íåðååííåðíûé öååò áóíååè, êðóííûå è íåééèå íÿòíà è ó.ä. Èí áòíðíò òéíò íñæíí
íòíåñòè áåôåéòü, áïçíééøèå íðè íòéòðíåéå, ýòí: íåðååííåðíàÿ ýðéññòù è
éííòðåñòíññòù èçíåðåæåíèÿ (÷àñòí íðÿäëýåòñý íðè ñúåìéå òéòðíåùí
òíòíàíðåðòí), íññåå÷éååíèå íaäíèñåé ñ íåðåòííùå ñòðåííùå èéñòà, òéòðíåùíé øóí. Íðèìåð
óåéíàí èçíåðåæåíèÿ ííæçáí íà Ðéññóíéè.

Đèññóíîê 1. Ôðâäåíàíò ñòàðéííäí òåêñòà (Ñèííäè Çèëàíòñà ííàñòûðý, íàöèííàëüíûé àðôèå ðâññóáëèèè òàðåðñòàí)

Đèñóíîê 2. Íđåäâàðèòåëüíàÿ íáðàáðòéà èc íáðàæåíèÿ

Níñòàâëåíèå
øðèôòà

Íáíçíà-èì íeéñåëè áéíàðííáí èçíáðàæåíèý êàé
 Aij, áäå i=1..x, j=1..y, áäå x è y - øéðèíà è áûñîòà
 ñíòåðåñòðóåáíí. Íeéñåëè áíá ðåéñòåúö áéíèíà íðèíàì çà áåéëüå, ò.ê. ðèñóíèè è
 êåðòèíèé äëÿ íàøáé cäàä-è íá áàæíú. Íáíçíà-èì dres êàé ðàcðåøåíèå ñêàíèðíàáíèý á dpi.

Ióñòü íáðàçåö øðèôòà - ýòí íááíð øàáëíííâ
 áóêâ, áóíáyueò â ááíñ ñíñòáâ, â éíëë÷ðòòâ ïðèìáðíí ðàáíúí ÷èñëó áóêâ áàíííâ
 àëòâàèòà (íí÷åíò «íðèìáðíí», áóäâò íáúýñíáíí íèæâ). Iéñëáëè øàáëííà áóêâû èç
 íáðàçöà øðèôòà íáíçíá÷èí êâé , z=1..k, ãää k - éíëë÷ðòòâ
 øàáëíííâ á íáðàçöå øðèôòà, i=1..mz, j=1..nz, ãää mz è nz - øèðèíà è áûñòáà øàáëííà
 ñíñòáâòñòáâíí (Ðèñóííê 3).

Ðèñóííê 3. Âèä øàáëííà áóêâû

Íáðàçöîâ øðèôòâ òàéæå áûòü íííâ.

Äëý òíáí ÷òíáû êëàññèòëëòíâòò ðòëòò ìà íáèçâåñòííé ñòðàíèëå, íááóíäèíí íáéòè
 íàëëó÷ðåâ ñíñòáâòñòáâ ñ íáíèí èç íáðàçöîâ øðèôòà. Äëý ýòíáí ëàæäóþ áóêâó
 ñòðàíèëü áóâái ñòðàíèâàòü ñ êàæäúí øàáëíííâ áóêâû áñâõ íáðàçöîâ øðèôòà ñ
 íííùþ õòíëëëè èíððåéëòè, ííèñàíííé á [8].

,

ãää i=1, 2,..., mz, j=1, 2,..., nz, - ñòðâäíâå çíà÷áíèå

íèñëáëé á øàáëííâ (âû÷èñëýâííâ òíëüêí íäéí ðàç), - ñòðâäíâå çíà÷áíèå ýéâíâíòíâ
 èçíáðàæåíéý A á íáéàñòè, ñíâíâàþùâé ñ øàéòùèí ííëíæåíèåí B, à ñóííèðíâàíéå áâåäåòñý íí áñâí ìáðàí
 èíððåéíàò, íáúèí äëý A è B. Èíýôôèëåíó
 èíððåéëòè íáíýâòñý íò -1 áí 1 è íá çàâèñòè íò èçíáíâíéý íáñòòàáà àííëèòóä A è B.

Èòíáíâäý èíððåéëòè äëý íáðàçöà øðèôòà
 áóâåò ñíðåäåëëòñý êâé:

Íáðàçåö øðèôòà (éëë íáñéíëüéí íáðàçöîâ) ñ
 íàéñèíàëüíûí èíýôôèëéáíòíí èíððåéëòè áóâåò ÿâëëòñý íáéëó÷ðèí ñíñòáâòñòáâéâí ñí
 øðèôòíí íáèçâåñòííé ñòðàíèëü.

Nâáíâíòàöëý
 èçíáðàæåíéý íá ñòðîéè è áóêâû

Èòàé, íù ííðåäåëëè íáðàçöà øðèôòà. Íí äëý òíáí ÷òíáû ááí ïðèìáíèòü, íáíáóíäèíí íá
 èçíáðàæåíéè ñòðàíèëü (íá èíòíðíí èíáþòñý òíëüêí òâéñòíâû áéíéè áâç ðèñóíéíâ è
 éâðòèíí) áûâåäëëòü ñòðîéè, à çàòâí è íòáâëüíû áóêâû.

Äëý ñâáíâíòàöëè ñòðîé âíñííëüçóâíñý
 ñëåäóþùèí ííñòðíí: ííñòðíè íâðòèëàëüíóþ íðíâëëòþ áèíàðííâ èçíáðàæåíéý íí
 ñëåäóþùâé ôíðíóéå:

, äääå j=1..y

Íñééá ýóíáí nääéäæè ííéó÷áííóþ íðíáéöéþ
íóàí óñðåäíáíéý éàæäíáí ýéåíáíòå áåééòíðå ñ åáí níñääíéíè çíà-åíéýíè. Đåçóëüòàò
ííéàçáí íá Đèñóíéå 4.

Đèñóíîê 4. Âûääåëåíèå ñòðîê

Ääääå äüäääëýþöñý ñòðîèè ïì òåí íåïðåðüâåíû
ïïëñàì, äää çíà÷åíèå ωj > λline, äää λline - íåéîòïðûé êîýôôèöèåíò. Çàòåì èç ïïëó÷åííûõ ñòðîè
áíàëíäè÷íù íåðàçîì äüäääëýþöñý ìøäåëüíû áóêåû. Äëÿ ÷àñòè ñèìäåíâ, êîïðûå
ñïåääðæàò íåïíñâýçíûå íåéäñòè (íàïðèàð, áóêåà «Ù») ïïåóò áûòü íåïðåâèëüí
ïïðåäåäåéäíû áðàíèöû. Í yòí îàéí åéèýåò íà êííå÷íû ðåçóëüòàò, ó.ê. äëÿ
êëàññèôèêåöèè ñðèòòà äïñòåòí÷í ëñïëüçíàéí òïëüéí ÷àñòè áóêå.

Íáó÷åíèå

Òâïáðü ïñòåâðñý áùà ïäíà ïðíáéäìà - èäè
ïñäíòïâèòü íáðàçöû øðèòòà. lïæíí ñääëæòü ýòî áðó÷íóþ, áûðåçàÿ øàáéííû áóêâ èç
èçíáðàæåíéÿ, íí ýòî áóäåò áâññùà òðóáíàíè è çàéíàò iííñí áðåíåíè, ííyòïíó
íåññòäèìà áàòïíàøçäöéÿ áäííñí ïðñòåññà.

Íðåäéàåòñý ñëåäóþùèé íïäöïä - íà áôïä
íïàåòñý íåñéïëüéí ñòðàíéö ñ óäì øðéòòíí, éîòîðûé íðåäíëàåòñý çàíåñòè á
êà÷åñòå ìåðàçöå. Áñä ñòðàíéöû, èðííà íàííé, èñííëüçóþòñý äëÿ ôîðíèòíàíéý íåðàçöïå
áóéå, íàçíåíàí èø íàó÷åþùåé áúåíðééí. Á íñòàåòàÿñý íaià ñòðàíéöà - éííöðíëüíàÿ, äëÿ
óóí÷íåíéý ñòðíèòíåííûó íåðàçöïå.

Äëÿ ñòðàíèö îáó÷àþùåé âúáîðèè ääëæþþñý
 ýóàïü íðåäâàðèòåëüíé îáðàáîòè è âúáâåëåíèÿ ñòðîé/ñèïáîéïå ëìèñàííûå âúøå. Âúäåëåííûå
 áóêåû çàïíèíàþþñý êåé øàáéííû. Íñîéíèüéò ñðåäíýý ñòðàíèöà ñòàðíà÷-àòííàí
 óâéñòå ñíàáðæò, êåé íðåâéèí, 200-1000 ñèïáîéïå, ñéâåíâðåðåéüí, ñòîéüéò æå íü
 ííé-èí è øàáéííà áóêå. Ñíòóâåòñòååíí, äëÿ èàæäíàí ñèïáîéà, à ñòðåáíàí, íü
 èíàái 6-30 øàáéííà á èàæäéíè ñòðàíèöå. Ôàéíà áíèüòíà éíèé-åñòåí íáæåëàðåéüí,
 ííþòíò íåáðíàéí ñíèðåðèòöù èò -èñéèí. Ýòî áóäàí ääëæàöù à áåâ ýóàïà, íà íåðåíí

Iú m̄iði ñðåâáíéâáí áñá áûâáæéâíú ðøáæéíí ñ m̄iñuþ óóíéøè èíððæýøè, èñtílœüçíâáíí áûøð. Áñéè æëý t-áððæáíí íæðu çíà-âáíéâ ëýøðøééâíòå íáíüðå &equal, òí íæíó èç êííèé óääæýäí.

Íñíñéå òàééîé íðiòååóòðû óäåæáíéý áóáéèéàòíâ ó
íàñ èíàåòðñý óæå åíñòàòí-í íðeåíáíäý éíñééåéöéý òàáæííâ ñèíâíéíâ ððéòòà. Í òàí
íñåóò ñíäåðæàòùñý íñòíðííéå yéäåíåòû, íàíðèíåð ñíëéíòèåñý èéé íñåðåæäåííúå
áóééåû, yéäåíåòû ðèñíóéíâ, ééyéñ í èíðí-åå. Áéý òíâí ÷òíâû åúå áíéüøð íñåùñèòü
éà-÷åñòåí òaaééííâ, iù íðíèçâíäé íðiòååóòð ðåñíçíååàíéý ððéòòà ñ íñíñùþ ýóèò
òaaééííâ íà íñíñé íñòàåðåéñý éíñòðíéüñíé ñòðåíéòå. Íñíñéíüéó ððéòò íà íååéò ñòðåíéòå
ééååíòé-íüé, ðí òå òååééííû, éíòíðû íà ñíäåðæàòñý íà éíòíðíéüñíé, ñòðåíéòå íà
íòíñyöñy é ñèíâíéäí, ñéååíåðåéüñí, éó ííæíí óååééýóü. Íñíñéå òåééíé íðiòååóòðû iù
íñéó-èéé ååñüìà ðíðíþþ áéåæéíòååéó íåðåçöíâ åäíííâ ððéòòà.

Çàêëþ÷åíèå

Áûëà ïìèñàíà ìáóíäéèà äëÿ àâòîìàòè÷åñêíé
ééàññéôéèàöèè ñðèôòà ñòàðïíà÷àóííàí òâéñòà, ïðåäéíæåíû àéäîðèòìû è ìáóíäû äëÿ
ññâíàíòàöèè éçjíàðàæåíéÿ íà ñòðîè è ñèíàíèû. ïìèñàíà ìáóíäéèà äëÿ
àâòîìàòè÷åñêííàí ôîðîéðíâáíéÿ íàðacöà ñðèôòà.

Àëäîðèòù è ïäöñâû áûëè ðåàëëçîâáíû â áëää
ïðíâðàìííé ñèñòåìû. Íà êíîôðâíöèè òàêæå áóäåò ïðäàìíñòðèðîâàíà ðàáîòà ñàïíé
ñèñòåìû è ïíèò÷áíû ïðàëòè÷áññèå ðåçóëüòàòù.

Eèòåðàòóðà

- [1] H. Shi
and T. Pavlidis, "Font Recognition and Contextual Processing for More Accurate Text Recognition," ICDAR'97, pp. 39-44, Ulm, Germany, Aug. 1997.
 - [2] Yong Zhu,
Font Recognition Based on Global Texture Analysis. IEEE Transactions PAMI. October 2001 (vol. 23 no. 10) pp. 1192-1200.
 - [3] R.
Cooperman, "Producing Good Font Attribute Determination Using Error-Prone Information," SPIE, vol. 3,027, pp. 50-57, 1997.
 - [4] A.
Zramdini and R. Ingold, "Optical Font Recognition Using Typographical Features," IEEE Trans. Pattern Anal. Machine Intell., vol. 20, no. 8, pp.877-882, 1998.
 - [5] A.
Schreyer, P. Suda and G. Maderlechner, "Font Style Detection in Documents Using Textons", Proc. of 3rd IAPR Document Analysis Systems Workshop, Nagano, Japan, 1998.

[6] Nïëîâüââ Å.Ä., Þæèêîâ Å.Ñ.

Åâòîàòèçèðîâàíàÿ ñèñòâà ìáðàáòèè è ðåñòàâðàöèè èçîáðàæåíèé ñòàðñåðàòíûõ
òâèñòâà è ðóëññâé. Åâñòíèê ÈÃÒÓ èì. Òóññâé. Åü.3. - Èàçàíü: Èç-âî ÈÃÒÓ,
2006. - ñ.28-30

[7] Å.Ñ. Þæèêîâ. "Nââìåòàöèÿ
èçîáðàæåíèé ñòðàíèö äðââíèö ðóëññâé". Náñíèê òðóäâà êñññâðåíòèè
"RCDL-2007". Iáðâññâé-Çàëåññâé. Òñ 1, c. 236-240, 2007

[8] Đ. Åññâé. Öèôðîâàÿ ìáðàáòèà èçîáðàæåíèé. Iññâà: Ôåñññôâðà,
2005. - ñ. 996-997.