

Âáòíàòè÷ãñêÿ êëãññèòèèàòèÿ øðèòòà ñòàðííá÷àòíúõ òáêñòíá

Ááòíð Áéááèìø Ñáðáááè÷ Pæèéíá
25.09.2009 á.
Íñéááíáá íáííéáíéá 05.10.2009 á.

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.

Âááááíéá

Éíááá á ðóèè ÷ãñíááéó, çáíèìàðùáíóñÿ ñòàðííá÷àòíúíè òáêñòíáè, ïííááááò íáèçááñòíáÿ ñòàðèííáÿ éíéáá, èèáí áá òðááíáíóú, áíçíéèááò íðíáéáíá - éáé ïðáááéèòú áðáíÿ è íáñòí áá íá÷àòè, áñèè ÿòà éíðíðíàòèÿ ïí èáèè-òí íðè÷éíáí íááíñòóííá. Ííúòíúé èññéááíáàòáèú, íáñííáííí, ííæáò íðèíáðíí ïðáááéèòú ÿòí, íñííáúááÿñú íá òáò éíéááò è ðóéííèñÿó, ÷òí íí áèááé ðáíáá, íí ááòáèúíúé áíáèèç øðèòòà, ñòèèÿ íáíèñáíéÿ òðááóáò ááñúíá èðííòèèáíé ðááíóú è íáíáèúò çàòðàò áðáíáíé.

Ííí÷ú á ðáðáíèè òáéíé çááá÷è ïíéèá áú ñíáòèáèúíáÿ éííúðòáðíáÿ ñèñòáíá, éíòíðáÿ óíááò ñðááíéááòú áóéáú íáèçááñòíé ñòðáíéòú ñ íííáí÷èñéáííúé íáðáçòáíè øðèòòá èç ðáçíúò èñòí-íééíá, éíòíðúá òðáíÿòñÿ á áá ááçá. Èñòíáÿ èç ÿòíáí, ííæáí ñòíðíèðíáàòú íñíáíúá òðááíáíéÿ è òáéíé ñèñòáíá:

Éàæáúé íáðáçáò øðèòòà áíéæáí èíáòú íáòáéíðíðíàòèð, òáèòð éáè: òèííáðáòèÿ, áíðíá/ñòðáíá, áíá èçááíéÿ è ò.á.

Áíéæíá áúòú áíçííæíñòú éááéíáí çáíáñáíéÿ á áàçó ááííúò ííáíá íáðáçòà øðèòòà.

Áúáá÷á ðáçòèúòáòíá áíéæíá áúòú ðáíæèðíááííé ïí èáæáííó íóíéòò íáòáéíðíðíàòèè. Íáíðèíáð:

o
Oeiadaoeuy - «Edeia» (87%), «Edeee e laodiae» (75%);

o
Aia ecaaeuy - 1764 (95%), 1795 (84%); ...

.
Eeanhedeaeoey edeoda ia aieeia caieiaou iiaa adaiie.

Iaci
eeadaoou

Ienuaaiay idiaeia ioiineony e eeanho
caaa- daniaciaaeuy iadacia, a eiaii - OFR (Optical Font Recognition). Yoa iaeanou
nae-an aeoeaii dacaeaaaony e iadiae naita ideiaiea ai iiaeo eiiad-aneeo
idiaoedao, oaeeo eae FineReader, CunieForm e iiaeo adoeo. E niaaeaiet, eeanhedeaeoio
edeodi adiauee a nioaa yoeo nehoai, ia iadiae ae dcaeuy iaee caaa-e,
o.e. «caoi-ai» ita anuia ocoep iaeanou - itadiaou iaeeiaa itoiauyop
adioedoo e noeu ia-adoaiey ec i-aiu iadie-aiia iaadia. Iau-i yoi:
iiiiedeiue edeod (Courier), edeod n
canea-eae (Times New Roman), edeod aac
canea-ae (Arial), iepn eo itaiiaana: bold e italic. Aey caaa-e
daniaciaaeuy oaeioia (OCR) aieuaae e ia odaaaoony, itoio -oi itaapeyua
-eneti niadiaiue aeioiaie enieuyopo auaiada-eneaiua edeodu.

Iadaoeiny e eeadaooua - a ioaeaeoeyo
anoa-apony aaa niaiuo itoiaa e daniaciaaeuy edeoda:

Auaaeia
aeiaeuio ideiaiey ec oaeioiauo aeieia, oaeeo eae neia, noiee e eo
aeuiaee aeiaey. Yoi itcaieyo itaapeyuo noeu ia-adoaiey edeoda (bold, italic), aai daciao,
iaee-ea cana-ae. Eae idaeie, oaeie itoia ia odaaao aacy iadacoia, itaao
itdaapeyuo ianeieui niaiuo eeanha edeodia e ia caaeneo io yuea (a
aieueiehoaa neo-aaa), -oi aaeaa aai itoiaeuio ae enieuciaiey a nehoiaio
daniaciaaeuy oaeioia OCR.

E ideiaod, a daiaoa Shi, Pavlidis [1] aey
daniaciaaeuy edeoda enieuyoaony oya aeiaeuio aey ana noiaieou ideiaiea,
oaeeo eae: aenoiadai daniaciaaeuy aei neia (a ieneayeo), ae-e-ia
iaehodi-iou, iaaoeaeaiuo eioadaeia. Yong Zhu [2] idaeaaaao daniaciaaeia
edeoda, itiaaia ia aiaeyca oaeiooou
eciadaaeiey ana noiaieou. Aey ecae-aeuy yoeo ideiaiea enieuyoaony
itiaeiaeuie eeuo adaiia.

Auaaeia
eiaeuio ideiaiea. Caa nu oaa aeeycedopony itaapeyua aoeau e eo yeaiaou. Oaeie
itadia -oanadaeai e iaieuei eciaiey i adedeoa, -oi itcaieyo oiooi dacapeyuo
aeeyea it ia-adoaiey edeodu e aieua ana itadiae aey iaee caaa-e.

A daiaoa Cooperman'a [3] enieaapony ioiee naitenoa edeoda aey nehoai OCR. A
e-a-anoaa adaeodia yoeo naitenoa enieuyopony eiaeuia ideiaiee, oaeeo eae:
iaee-ea cana-ae, itoioho e o.a. A daiaoa Zramdini, Ingold [4] idaeaaaony
noaene-ane e itoia ae dcaeuy caaa-e eeanhedeaeoe edeoda, itiaaiue ia
auaeaeie eiaeuio ideiaiea. Noiiae ioia enieuyoaony e a daiaoa [5]

Íðáááàðèòáèúíáÿ
íáðááíòèà

Nòàðííá-àòíúé è, á íñíááííñòè, ñòàðèííúé ðòèííá-àòíúé (hand-writing) òáèñò èíááò ðÿá íñíááíííñòáé, èíòíðúá íá ïçáíéÿðò íáíðÿíòð ïðèíáíèòú íáòíáú, íðááèèàáááíúá á íóáèèèàòèÿò ïí OFR. Á -áñòííñòè, íáæáóéááííúá è íáæñòðí-íúá ðàññòíÿíéÿ áàðüèðòðòñÿ á áíñòáòí-íí ðèðíéèò ïðááèèàò áàæá íá íáííé ñòðáíéòá. Èðííá òíáí, òáèñò èíááò áíèüòíá èíèè-áñòáí áàæíúò äèÿ èèàññèòèèàòèè ððèòòà àèàèðòè-áñèèò ÿéáíáíòíá, íá ñòðáíéòáò -áñòí áñòò ðèñóíéè, ïííáòèè. Òáèæá èçíáðáæáíéá ñíááðæòè ááòáèòú, áúçááííúá áèèòáèúíúí òðáíáíéáí, èíòíðúá ïíæíí òñèíáíí ðàçááèèòú íá 2 òèíá. È íáðáííó íóíáñáí ááòáèòú ñáíèò ñòðáíéò èíèáè, ïíÿáèáøèèñÿ á ðàçóèüòáòá áíèáíáí òðáíáíéÿ, ááèñòáèÿ áèàæííñòè, òáííáðáòòðú, ïíðáæáíéÿ áðèáèíí ïðááèúíúò ñòðáíéò, áúòááòáíéá áóéá, íáðááííáðíúé òááò áóíááè, èðòííúá è íáèèèá ïÿòíá è ò.á. Èí áòíðíò òèíó ïíæíí íóíáñòè ááòáèòú, áíçíèèøèá ïðè ïðèòðíáèá, ÿòí: íáðááííáðíáÿ ÿðèíñòú è èííòðáñòííñòú èçíáðáæáíéÿ (-áñòí ïðíÿáèÿáòñÿ ïðè ñúáíéá òèòðíáíí òíòíáíáðáòíí), ïðíñáá-èááíéá íááíèñáé ñ íáðáòíé ñòðííú èèñòà, òèòðíáíé øóí. Íðèíáð òáèíáí èçíáðáæáíéÿ ïíèàçáí íá ðèñóíéá 1.

ðèñóííè 1. Óðááíáíò ñòàðèííáí òáèñòà (Ñèííáèè Çèèáíòíáá ïíáñòóðÿ, Íáøèíáèúíúé áððèá ðáñíóáèèèè Òàòáðñòáí)

Áñá ÿòè íñíááííñòè ïíáóò ïíáøáòú ïðíáññó èèàññèòèèàòèè ððèòòà. Á ñáÿçè ñ ÿòèí, èçíáðáæáíéá íáíáòíáèíí ïðáááàðèòáèúíí íáðááíòáòú è ïíáííóíáèòú. Áèÿ ÿòíáí áíñíèüçóáíñÿ áèáíðèòáíè, ïíèñáííúè á [6] - á ðàçóèüòáòá èò ðááíòú èíááí áéíáðííá èçíáðáæáíéá (ííí ñíñòíèò òíèüèí èç -áðíúò è ááèüò ïèññáèé), í-èúáííá íò øóíá è ïííáò, à òáèæá ïíááðíóòíá íá íáíáòíáèíúé óáíè, áñèè ñòðáíéòá ïðè ïðèòðíáèá áúèá íá íáðáèèáèúíá èðàð ñèáíéðòðúááí òñòðíèñòá. Ááèáá ñ ïíííúð íáòíáá, íðááèíáííáí á [7], áúááèÿáí òáèñòáíúè áéíè. Íááíèüøèá íáòí-ííñòè á ïðáááèéáèè áðáíèò òáèñòà, èíòíðúá áíçíèèàðò íá ñèíáíúò èçíáðáæáíéÿò, ïíæíí íá ïðèíèáòú áí áíèíáíéá, ïíòííò -òí íáøá òáèü - ÿòí íá ðàñííçíáááíéá òáèñòà, ïíÿòííò ïðááðÿ íáèíòíðúò áóéá íá áðáíèòáò ïðáèòè-áñèè íá ïíáèèÿáò íá èííá-íúé ðàçóèüòáò. Èòíá ïðáááàðèòáèúííé íáðááíòèè ïíèàçáí íá ðèñóíéá 2.

ðèñóííè 2. Íðáááàðèòáèúíáÿ
íáðááíòèà èçíáðáæáíéÿ

Ñííñòááèáíéá
ððèòòà

Íáíçíá-èí ïèèñáèè áéíáðííáí èçíáðáæáíéÿ èáè Aij, ááá í=1..x, j=1..y, ááá x è y - øèðéíá è áúñíòà ñííòááòñòááíí. Íèèñáèè áíá òáèñòáíúò áéíéíá íðèíáí çà ááèüá, ò.é. ðèñóíéè è èáððèéíèè äèÿ íáøáè çááá-è íá áàæíú. Íáíçíá-èí dres èáè ðàçðáøáíéá ñèáíèðíááíéÿ á dpi.

Ἰόνου ἰάδαζαὸ ὀδεὸδὰ - γοῖ ἰαίῖδ ὀααῖῖῖῖ
ἀόεῖ, ἀοῖῖῖῖῖ ἂ ἀῖ ἡῖῖῖῖ, ἂ εῖῖῖῖ-ἂῖῖῖῖ ἰδεῖῖῖῖ ὀαῖῖῖῖ ὀ-εῖῖῖῖ ἀόεῖ ἂῖῖῖῖ
ἀεὸῖῖῖῖ (ἰῖ-ἂῖῖ «ἰδεῖῖῖῖ», ἀόῖῖῖ ἰαῖῖῖῖῖ ἰεῖῖ). ἰεῖῖῖῖ ὀααῖῖῖῖ ἀόεῖῖ ἔϥ
ἰάδαζαὸ ὀδεὸδὰ ἰαῖῖῖῖ-εῖ ἔῖῖ, z=1..k, ἂῖῖ k - εῖῖῖῖ-ἂῖῖῖῖ
ὀααῖῖῖῖ ἂ ἰάδαζαὸ ὀδεὸδὰ, i=1..mz, j=1..nz, ἂῖῖ mz ἔ nz - ὀεὸεῖῖ ἔ ἂῖῖῖῖ ὀααῖῖῖῖ
ῖῖῖῖῖῖῖῖῖ (Ἐῖῖῖῖῖ 3).

Ἐῖῖῖῖῖ 3. Ἀεὸ ὀααῖῖῖῖ ἀόεῖῖ

Ἰάδαζοῖῖ ὀδεὸδῖῖ ὀῖῖῖ ἰῖῖῖ ἂῖῖῖ.
Ἀεῖ ὀῖῖῖ ὀ-ὀῖῖῖ ἔῖῖῖῖῖῖῖῖῖῖῖ ὀδεὸδ ἰῖ ἰῖῖῖῖῖῖῖῖῖ ἡῖῖῖῖῖῖῖῖ, ἰῖῖῖῖῖῖῖ ἰῖῖῖῖ
ἰῖῖῖῖῖῖῖ ῖῖῖῖῖῖῖῖῖῖ ἡ ἰῖῖῖῖ ἔϥ ἰάδαζοῖῖ ὀδεὸδὰ. Ἀεῖ γοῖῖῖ ἔῖῖῖῖ ἂῖῖῖῖ
ἡῖῖῖῖῖῖῖ ἂῖῖῖῖ ἡῖῖῖῖῖῖῖῖ ἡ ἔῖῖῖῖῖ ὀααῖῖῖῖῖ ἂῖῖῖῖ ἂῖῖῖῖ ἰάδαζοῖῖ ὀδεὸδὰ ἡ
ῖῖῖῖῖῖῖ ὀῖῖῖῖῖῖ ἔῖῖῖῖῖῖῖῖ, ῖῖῖῖῖῖῖ ἂ [8].

ἂῖῖ i=1, 2,..., mz, j=1, 2,..., nz, - ἡῖῖῖῖῖ ῖῖῖῖ-ἂῖῖῖῖ
ἰεῖῖῖῖῖ ἂ ὀααῖῖῖῖ (ἂῖῖ-εῖῖῖῖῖῖ ὀῖῖῖῖ ἰῖῖῖ ὀῖῖ), - ἡῖῖῖῖῖ ῖῖῖῖ-ἂῖῖῖῖ ῖῖῖῖῖῖῖ
ἔϥἰῖῖῖῖῖῖ A ἂ ἰῖῖῖῖῖ, ἡῖῖῖῖῖῖῖῖ ἡ ὀῖῖῖῖῖῖ ῖῖῖῖῖῖῖ B, ἂ ἡῖῖῖῖῖῖῖῖ ἂῖῖῖῖῖ ῖῖ ἂῖῖῖ ἰῖῖῖ
ἔῖῖῖῖῖῖ, ἰῖῖῖῖ ἂῖῖ A ἔ B. Ἐῖῖῖῖῖῖῖῖ
ἔῖῖῖῖῖῖῖ ἰῖῖῖῖῖῖ ῖῖ -1 ἂῖ 1 ἔ ἰῖ ῖῖῖῖῖῖ ῖῖ ἔϥἰῖῖῖῖῖ ἰῖῖῖῖῖῖ ἂῖῖῖῖῖῖ A ἔ B.

Ἐῖῖῖῖῖῖ ἔῖῖῖῖῖῖῖ ἂῖῖ ἰάδαζαὸ ὀδεὸδὰ
ἀόῖῖῖ ῖῖῖῖῖῖῖῖῖ ἔῖῖ:

Ἰάδαζαὸ ὀδεὸδὰ (ἔῖῖ ἰῖῖῖῖῖῖ ἰάδαζοῖῖ) ἡ
ἰῖῖῖῖῖῖῖῖῖ ἔῖῖῖῖῖῖῖῖῖῖ ἔῖῖῖῖῖῖῖῖ ἂῖῖῖῖ ῖῖῖῖῖῖῖ ἰῖῖῖῖῖῖ-ὀῖῖ ῖῖῖῖῖῖῖῖῖ ἡ
ὀδεὸδῖῖ ἰῖῖῖῖῖῖῖ ἡῖῖῖῖῖῖ.

Ἰῖῖῖῖῖῖῖ
ἔϥἰῖῖῖῖῖῖ ἰῖ ἡῖῖῖῖ ἔ ἂῖῖῖῖ

Ἐῖῖῖ, ἰῖ ῖῖῖῖῖῖῖῖ ἰῖῖῖ ἂῖῖ ἰῖῖῖῖῖῖῖ
ἰῖῖῖῖῖῖῖῖ ἰάδαζαὸ ὀδεὸδὰ. ἰῖ ἂῖῖ ὀῖῖῖ, ὀ-ὀῖῖῖ ἂῖ ἰῖῖῖῖῖῖῖ, ἰῖῖῖῖῖῖῖ ἰῖ
ἔϥἰῖῖῖῖῖῖ ἡῖῖῖῖῖῖ (ἰῖ ἔῖῖῖῖῖ ἔῖῖῖῖῖ ὀῖῖῖῖ ὀῖῖῖῖῖῖῖ ἂῖῖῖῖ ἂῖῖ ὀῖῖῖῖῖῖ ἔ
ἔῖῖῖῖῖῖ) ἂῖῖῖῖῖῖ ἡῖῖῖῖῖ, ἂ ῖῖῖῖ ἔ ἰῖῖῖῖῖῖῖ ἂῖῖῖῖ.

Ἀεῖ ἡῖῖῖῖῖῖῖῖ ἡῖῖῖῖ ἂῖῖῖῖῖῖῖῖῖ
ἡῖῖῖῖῖῖῖ ῖῖῖῖῖῖ: ῖῖῖῖῖῖῖ ἂῖῖῖῖῖῖῖῖῖ ῖῖῖῖῖῖῖ ἂῖῖῖῖῖῖῖ ἔϥἰῖῖῖῖῖῖ ῖῖ
ἡῖῖῖῖῖῖῖ ὀῖῖῖῖῖῖ:

íú ììãðíí ñðãáíéããáí áñá áúãáéáííúá ðãáéííú ñ ììíúúþ óóíéòéè éíððãéýðèè,
èñííéúçíãáííé áúðá. Áñèè äéý í-áðãáííé ìãðú çíà-áíéã éíýóðèòéáíòà íáíúðá λequal, òí íáíó èç éííéé óããéýáí.

Ìíñéã òãéíé ìðíòããóðú óããéáíéý áóáéèèàòíã ó
íãñ èíããòñý óãã áíñòàòí-íí ìðéáííáý éíééãéòéý ðãáéííã ñèíáíéíã ððèòòã. Íí òàí
ííãóò ñíããðãèòñý ìñòíðííéã ýéáíáíòú, íáíðèíãð ñèèèðèãñý èèè ìãðãããáííúá
áóéãú, ýéáíáíòú ðèñíóíéíã, èéýéñ è ìðí-ãã. Äéý òíãí -òíãú áúã áíéúðã ìãúñèòú
èã-ãñòáí ðãáéííã, íú ìðíéçãíãè ìðíòããóðó ðãñíçíããáíéý ððèòòã ñ ììíúúþ ýòèð
ðããéííã íã íáíé ñòããðãéñý éííòðíéúííé ñòðáíéòã. Ìíñéíéúéó ððèòò íã íããéð ñòðáíéòã
èããíòè-íúé, òí òã ðããéííú, éíòíðúá íã ñíããðãèòñý íã éííòðíéúííé, ñòðáíéòã íã
íóííñýòñý è ñèíáíéíã, ñéãáíããòãéúíí, èð ìãíí óããéýòú. Ìíñéã òãéíé ìðíòããóðú íú
ííéó-èèè áãñúá òíðíòðð áéãéèíòãéó íãðãçòíã ááííãí ððèòòã.

Çãèèþ-áíéã

Áúéã ìíèñáíã íãòíãéèã äéý áãòíãòè-ãñéíé
èèãññéòéèãòéè ððèòòã ñòãðííã-òòííã òãéñòã, ìðããéíããú áéãíðèòú è íãòíãú äéý
ñããíãíòãòéè èçíãðããáíéý íã ñòðíéè è ñèíáíéú. Ìíèñáíã íãòíãéèã äéý
ããòíãòè-ãñéíãí òíðíéðíããáíéý íãðãçòã ððèòòã.

Áéãíðèòú è ìãòíãú áúéè ðããéèçíããú á àéãã
ìðíãðáííé ñèñòáíú. Íã éííòãðáííé èòããá òðíããíñòðèðíããá ðããíòã ñáííé
ñèñòáíú è ìíéó-áííúã ìðãèòè-ãñéèã ðãçóéúòãòú.

Èèòãðãòòðã

[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] Ñííéíáúáá Á.Á., Þæèéíá Á.Ñ.

Ááòííàòèçèðíááíáy ñèñòáíà íáðááíòéè è ðáñòááðàðèè èçíáðáæáíéé ñòáðííá-àóíúð òáèñòíá è ðóèííèñáé. Ááñòíéé ÉÁÓÓ èì. Óóííéááá. Áúí.3. - Éçáíú: Éç-áí ÉÁÓÓ, 2006. - ñ.28-30

[7] Á.Ñ. Þæèéíá. "Ñááíáíòàòèý

èçíáðáæáíéé ñòðáíèð äðááíèð ðóèííèñáé". Ñáíðíéé òðóáíá éííòáðáíòéè "RCDL-2007". Íáðáñéàáèü-Çàèáññéé. Òíí 1, c. 236-240, 2007

[8] Ð. Áííñáèñ. Õèððíááy íáðááíòèà èçíáðáæáíéé. Ííñéá: Õáðííñòáðà, 2005. - ñ. 996-997.