

Íá íäíí iåòíäå ñòàòèñòè÷åñéíé ôèëüòðàöèè òåêñòíâíé èíóîðìàöèè

Àâòìð Ñòàïèñéàå Âàñèëüâåè÷ lì÷âñíà
18.07.2008 å.
lìñéäåíå ïáñíæäíèå 20.07.2008 å.

Òåêñò ïå÷àòííäî èçääàíèÿ â ôîðìàòå PDF

Âååääåíèå

Â ääiííé ñoåòöüä iðiâiæöñý íáçïð ñoåòëññòð-âññéèð íåðiäiâi áiâæçä ðåâññòðiâiñíûð ýçüñéâð, à òåðéæð ïiñéâçäiûð âiçíññøññòð

1. Эðàòèèé íáçïð ñòàòèñòè÷åñèèõ íåðîäíà èññëåäíàíèÿ òåêñòîâíé èíôîðìàöè

Øðøðiðiðið ðáññiðiðiðiðááíéå ñòðåðeññòðe-ñáññééå íåðiñü íáøðeé á áóíàíéðøðáðiñú íáéññòðyç cíáíéé, á ÷-ñáññiðiñòðe, á èññiðiðeé è èéðøðáðøðoðáð ááðiñü ñòðaâ èéðøðáðøðiñíâ iðfèçâááíéy “Øðøðeé Áíí” í. Á. Øðiðiðiðiâ [Óññáíé 1985]. Áíññüðéññóâí ýòðe íåðiñü íáññéé åññéé ñòðaâ íåññéé.

Í·áâéáíí, ÷óí éíéè·áñòááííáý íoáíéá èíáéáéáóáéúíúó íoéè·éòáéúíúó íñíááíííñòáé íðíéçááááéý óíáí ééé èííáí áâòíðá ááñüíá íáñ

Íðé áúáíðá óáó ééé éíúó ñòáòéñòé·áñééó óàðáéòáðéñòéé, èñííéúçóáíúó á éá·áñòáá áíááðéáíóúó, íáíáóíáéíí, ÷óíáú ííé óáíáéá

– éíóááðáéúíñòé (íáíáùáíéá ííñòé áðóíííü ííéáçáóáéé);

– íñòíýíñòáó (äéý íðáááéáííáíí áâòíðá íá áðóíííü íðíéçáááéé);

– äéáíáçííó éçíáíáéé (äéý áðóíí áâòíðíá).

Óàéíá ñí·áòáíéá áñáó óáó óáó·éñéáííúó óñéíáéé ííçáíéýáó áíáíðéóü íáéé·éé íáéíðíá áâòíðñéíá áíááðéáíóá.

Á éá·áñòáá áíéé·áñòááííúó óàðáéòáðéñòéé óáéñòíá íðááéáááþòñý íéááóþùéá:

– äééíá íðááéíáéé (ñðááíá áéñéíí áñéíá áíðááéíáéé, ííñ·éòáííá äéý éáæáíé áúáíðéé);

– äééíá ñéíá (ñðááíá áéñéíí áñéíá áñéíá áñéíá áíðááéíáéé);

– ÷áñòíóá óííóðááéáéý íéóæááíúó ñéíá íðááéíá áíþçíá, ÷áñòéö;

– ÷áñòíóá óííóðááéáéý íóùáñòáéóáéúíúó;

– ÷àñòòà óïòðåáëåíèÿ ãëàãîëîâ;

– ÷àñòîòà óíïòðåáäéåíèÿ íðèëàäàòåëüíûð;

– ðàñòîòà óïîòðåáëåíèÿ ïðåäëîää «â»;

– ÷àñòòà óííòðåáëåíèÿ ÷àñòèöû «íå»;

– ðíèéè-ðñòðâî ñëóäåáíûð ñëîïà â ïðåäåéíæåíèè (ñðåäåíàð ÷ èñëëè ñïþçâ, ïðåäåéíàð è ÷-àñòðö â ïðåäåéíæåíèè äéÿ èåæåíè åúâ

Ñéääöåò ìòiåòèöü, ÷òi à óéäçäííüö ðäáòäö [Óäðëí] iðéäåäåííüå ñòäòèñòè·åññéäà õäðäéòåðéñòèé èñíiøüçíåäèñü èëöü äëÿ åðóíí àåòïðîå. Å òi æå åðäìy iðè ðåðäíèè çääà· àíàëèçà åàæíüì ýäëÿåöñý ñíèðåùäíèå íáúåìà èññöïåíé òåéñòåíé èíóïðìàöèé ååq

2. Àëäîðèòì ôèëüòðàöèè

Â ðàáîòà [lî-âíî è àð. 2004] iðè ñòðöðåíè è ñòàðèñòð-ðåññé è íóîðìàðëíííé iñååëë èññüçåàëèñü òåññòû ï ýéíííè-ðåññé, ýéíë 10 ñòðåíèò ìàøëííèñíäî òåññòà ôîðìàòà À4, ñíååðæàùèò iñðÿäà 300 iñðääëíæåíè èäæäúé).

Íðäääíàðòíí ááííñûð êññéäåíàíéé ýäéëýåðñý áúÿâëåíéå íáíáùåíñûð ñòàòèñòè÷åññéèõ õàðàéòåðèñòèê ðàññïðåäåëåíéý êîëè÷åñòàà ñòðåøåíéå ñéäåäóþüèõ çäåäà:

– ó láíslúðráleà láuðáâí láuðáâí àálaëëçéððóâíâí òâéñòâ ïðé ñíððáâíâíè ââí ñâílæðóâ-âññéïé ñíñðoaâëýþþuâé;

– àâòïìàòè÷åñêîå ðåôåðèðîâàíèå;

– äüyäéäiéå ñäööèöðe ãñéëö ðöðäiaíöïä òäéñöà, öäïäéäöâïðýþùëö çääääííü èöðöàäðéýi;

– iñääñòñâêà òåêñòà äëÿ äàëüíåéøèõ ýòàïñâ àíàëèçà.

Âúyäéäííñá ñíòàðòèñòé÷-âññééå ðäàðåéòåðéèñòéééò lïäóò áùñöü èññïñéüçñåíàííú è íðé ðäåøäíéè çäåäà÷ ñéíòåçå òåéñòå.

À róðóðánñá èññéðåáðáié è úúëè áúðíðéáíù yéññáððéíáðòù, láðáðaaéáíùá lá éçóð-áíéà ðíñéè íððåáðéíáðé ðàçñííé äéëéíù. lá ððéè. 1 íñéacá

Đèñ. 1. Đàñïðåäåëåíèå ïðåäëîæåíèé iî êîëè÷âñòâó ñëîâ à ïðåäëîæåíèé

à ãòðåòòéå à ñ iñè ðäëéíàò ðòëíàðåíí àééò-áñòóà íòðåäééíðééé, à ñ iñè ááñòöéññ — èíòåðåàéüñ áðóíí à ñ èññéò ñíéíà á ïðåäééíæ

Է՞նդապէստի առաջարկը հայտնաբերություն է կազմում և առաջարկ է առաջնային առողջապահության համար:

Äëy öäëåé iàïðàâåëåííé ôëëüòðàöèè iïäóò áûòü èñïëüçîàíú è äðóäèå iáîáùåííûå õàðàêòåðèñòèéè, í êîòïðûô øëà ðå÷ü áûøå.

Çàêëþ÷åíèå

Ñòàðèñòè÷åñêèå ìåðîäû àíàëëçà ðóññèíýçû÷íûô òåêñòà ëíäóò ñ óñïäõî îðèìåíýöüñÿ äéÿ ðåøáíèÿ ðàçííâðàçíûô çàäà÷ íàðåáàòè õòëëüöðàòëÿ íïçâíëÿåò ëíâðèðîâàòü ýéàíàíòàè òåêñòà, íàïðèìåð, ãèàâîé, ñòðàíèöåé, àáçàòåí.

Summary

The paper presents a review of the statistic methods of analysis of texts in natural languages and shows the possibilities of filtration of the text information on the basis of one of the most informative statistic text characteristics. There is developed a method of filtration of the texts in the Russian language that can be applied to the solution of various tasks of text processing like the decrease of the volume of textual information, writing essays, determination of the semantic component of the text units etc.

Nïèñîê ëèòåðàòóðû

ìàòåìàòè÷åñêèå 1977 ìàòåìàòè÷åñêèå ìàòåäû à èñòîðèéí-ýéíí-ìè÷åñêèõ è èñòîðèéí-êóëüööðûõ èññëåäíâàíéþ. Ì., 1977.

ìàòåìàòè÷åñêèå 1985 ìàòåìàòè÷åñêèå ìàòîäû è ÝÂÌ à èñòîðë÷åñêèõ èññëåâîàìèþ. Ì., 1985.

lì÷åíîâ è äð. 2005 lì÷åíîâ, Ñ. Å. lõðèlåíáíèå ñòàòèñòè÷åñêèõ låòíäîâ iðè àíäëèçå òåêñòåíé èíôîðìàöèè / Ñ. Å. lì÷åíîâ, Å. I. Äéää

Íññâññéè è äð. 1989 Íññâññéè, Ä. Ä. Ñòàòèñòè÷åññéèå äóáëëèåòù â óïñðýäî÷åíñûõ ñïèñèåò ñ ðàçáèåíèå / Ä. Ä. Íññâññéè, Ä. Ä. Ñòàòèñòè÷åññéèå ïíññéåäíàìèÿ. Ì., 1989. № 138–148. (Íàö÷. Ññâåò î ëíññéåññéè ïðíæäíà «Èéååðíåòèè», Àí ÑÑÑÐ).

Ôñâåíêî 1980 Ôñâåíêî, À. Ò. Íåêîòïðûå ñòàòèñòè÷åñêèå çàêîíñåðíñòè ðàñïðåäåëåíÿ ÿ èíòíñòè è íóïðàöèè à òåêñòàõ ñî øêàëíé

Ôííáíê 1983 Ôííáíê, Á. Ò. Àâòïðñêèé èíâàðèàíò ðóññêèõ ëèòåðàòóðíûõ òåêñòíâ // Íåòíäû êíëè-åñòâåííñá àíàëèçà òåêñòíâ íáñ. 86–109.

Ôííáíê 1985 Ôííáíê, Á. Ò. Èíôïðìàòèâíûå ôóíêöè è ñâýçàííûå ñ íèìè ñòàòèñòè-åñêèå çàéííñíàðñòè / Á. Ò. Ôííáíê // Ñòàòèñòè

Õàðèí Õàðèí, I. Í. Èññëåäíâàíèå ïðèíðèíâ ñâàíòè-åñêíâí ïíèñêà òåêñòíâíé èíôïðìàòè è íà íñííâå èñííëüçíàíèÿ èíòåëéåéòóàëüíûõ