

Ååêòîðíàÿ ïääåëü îðåäñòàâëåíèÿ òåêñòîâîé èíóîðìàöèè

Àâòîð Ñòàíèñèåâ Äàñèüåâè÷ ëí÷åíïâ
18.07.2008 å.
ëíñëåäíåâ íáíïåäéåíèå 20.07.2008 å.

Òåêñò ïå÷àòííäî èçääàíèÿ â ôîðìàòå PDF

Âååäåíèå

Íñîáóþ àèòóàëüíñòö ïðèïáðåòàåò ðàçöàáìøèà ìàðïäíâ èçâëå÷åìèÿ è ôîðìèðíâàíèÿ ííâûö çíàíèé, íåíàöïäèìûö äëÿ ðåðåíèÿ êíñêðåòí êíñïëåéñííà ñèñòàííà èñïñëüçâàíèå ðàçèè÷íûö ëèíâàèñòè÷åññèò ïäðïäíâ è ìàðïäíâ èñêòññòåâííà ëìòåðåéëåêòà, íàïðàâëåíûö íà èäåþ, çàëíæåíóþ ààòòíî.

Íírueüçláàòåðáëü íá ãññääää ííæåò ðí-íí ñíòðíóëëðíláàðü ííñéñéñáûé çàïðíñ íá ííéó-ðáíéà èí ëíòðíàëöèë, èíòðíðäý álöð íáíáññíàëìà. Áíéëåå ðíññí, Òðóðáíññòë, ñâýçáíñúá ñ ðåðáíéàë ýóëé çàëäà-è, çåééþ-áþòñý á ííññíáðàçèë áíçíññíñúð ðíòí ãúðàæáíéý íäííé e ðíë æå èäéëè, üññéé

Â ääiííé ñòåòòüä ðàññíàòðéäþöñü íåéòïòðñä ïïäöñäü ê ðäøäíèþ öéâçäííûò ïðíáëä. ïñííäííà áíéíàíéä óääëäíí íåðííäí áåéòïðííàí íðå

1. Îáçíð ìåòîäíâ âåêèòíðííâí îðåäñòàâëåíèý òåêñòíâ

Â êííøå 80 õ áâäââ â ðàâîòåðå ñàëöííà [Salton et al. 1994] áúëæ iðåäëíæåíà âåéòíðíàÿ iíäåéü êàé àëüøåðíàòèåà åéñè-åñéííó ááñéííó ñíløåðà ñéíà è ñíløåðòñòååííí âåéòíðà â åéñè-åñéíí iðíñòðåíñòåå. Â iðíøåññå iíèñéà ÷àñòòíóúé iíðòðåò çàïðíñà ðàññìàòðèåà ðåéåâåíóúå áééòíåíóú.

Â Áíëåå ïðíäâèíóòûõ áâéòïðíûõ ííäåéýõ ðàcìàðíñòü ïðíñòðàíñòâà ñíéðàùàåòñÿ ìòáðàñûâàíèåì íàéáíëåå ðàñïðíñòðàíåíûõ èéè ðâ

Ãëàâíúi áiñòïèíñòâi ââéòïðíé iñäåéè ýâéýâòñý âíçiiæíñòü iñèñêà è ðàíæèðiâàíèý áîéóiaíòâi iñäåáèp, òi âñòü iñ èô áéèçíñòè óâíâæòâi ðèòâåëüíûìè, ÷òi iññáâi iñiyâéýâòñý, êtääà çàïðiñ iñäååðæèò iñèè-âñòâi ñëiñâ. Äey iñéó-âíèý èô-ðæé ðâéâåâiñòü Latent Semantic Indexing (LSI). Iñäåéü èñiñëüçâàë Singular Value Decomposition (SVD) äey iñðâññâàòâi iñ ðàçðâæâiñíé iñòðèòöü

LSI iñèâçàëà çíà-èòâåëüíâ iñäåññòâi â ðâçóëüòâòâo ñòñêà è ñòâåâiâièp ñ ëâéñè-âññèè iñòâiñâ, iñíàéñ ñëâæíñòü iñäåéè ÷âñòâoâi ñèéé [Salton 1989]. Íñáà èç iñèâéâå ðââiññiñâiñòü ñèñòâi iññâ LSI áuëà ñíçâàiâ â Áâðéèè à 1995 ãñäo iñééñâ Áâððé

Íñèññâàâiâi ìèæâ ñèñòâiâ èñiñëüçóâo ñââåððâiñíí äðóâóþ èíòâðiðâòâoëp iñiyòèý ââéòïðíé iñäåéè òâéñòâ, â êiòïðié iñ iñðèiâiýþò

2. Áâéòïðiâi ñòññâ ñòââéâièý òâéñòâiâi èíòâðiâoëè

Â äâiñíé ðââiòâ òâéñò ðâññiâòðèâââòñý êââ ñòðóêòâðâ, òi âñòü êââ ñââiêoíññòü iñäåéüíûò âçâèiññâýçâiñûò iñðâæiæâièé, iñââiòðiñ, ÷âðâç iññâññòâi iñäåéâéâ ðâçññâ ñââiññâ, iñðâååëýâiñûò iñäåéüíûìè iñðâæiæâièé, àáçâoâiâ, iñðâññâòâiâ, ãëâââiâ è ò. iññâññâ

Íèæâ iñðèâiâyòñý iñðèiâðû iñòâiâòâ-âññèè èíòâðiðâòâoëè ââéòïðíé iñäåéè iñâiññâiâi çâeññâiâi yëââiâiòâ òâéñòâ, ñññòññââi, iññâññâ

$$G = \{G_1, G_2, \dots, G_i, \dots, G_n\}$$

$$Vg = \{Vg_1, Vg_2, \dots, Vg_i, \dots, Vg_n\},$$

âââ G —; iññâññòâi âëââ; Gi —; i àý âëâââ, i = 1, ..., n; Vg —; iññâññòâi ââéòïðiâ õâéâé âëââ; Vgi —;

Â ñâiþ iñðâæü

$$Ai = \{Ai_1, Ai_2, \dots, Ai_j, \dots, Ai_m\}$$

$$Vai = \{Vai_1, Vai_2, \dots, Vai_j, \dots, Vai_m\},$$

ãääå Ái — iííæåñòáí àáçàöåâ i íé ãëàâû; Áij — j ûé àáçàö i íé ãëàâû, j = 1 … m; Vái — iííæåñòáí áåéèòå ãëàâû.

Íàòåìàòè÷åñêàÿ èíòåðïðåòàöèÿ áåéèòïðíé iíäåéè ãëÿ íðåäéíæåíèé áúðàæàåòñÿ á åèäå:

Pij = {Pij1, Pij2 ,…, Pijh ,…, Pk}

Vðij = {Vðij1, Vðij2, …, Vðijh,..., Vðijk},

ãääå Ðij — iííæåñòáí íðåäéíæåíèé i íé ãëàâû j ãí àáçàö; Ðijh — h íå íðåäéíæåíèå i íé ãëàâû j ãí àáçàö, h = 1 … Vpjh — áåéòïð öäéè h íå íðåäéíæåíèÿ i íé ãëàâû j ãí àáçàö.

Èç íðåäñòåâæåíííá ñíèñàíèÿ áèäíí, ÷òî êàæäííó ýëåìåíòó (ôðàäìåíòó) òåéñòà ñòàâèòñÿ á ñííòåâòñòåèå íåéíòïðûé áåéèòïð öåéè.

Èåé èçåâñòíí, ñíùñéíâû è áðàììàòè÷åñêè õåíòïí íðåäéíæåíèÿ íáû÷íí ýäéÿåòñÿ ñêàçóåííå, áúðàæåíííå áëàäåíèíí (ííèíçíà÷íûí èéè çàâèñèíûè ñéíâàìè èéè áåç íèð) ííäóó áûñòóìàòü á êà÷åñòåå çàéíí÷åíííå íðåäéíæåíèÿ. Íðèìåðàìè òàéèø íðåäéíæåíèé, íàçûâàåíííá.

Áåäöåòòü íåðâåíâ. Íí÷ü. Ííåäåéüíèé. Í÷åðòàíüý ñòîéèöü áí íæéå. (Àðìàòòå).

Èðííå èíåííûõ íðåäéíæåíèé íæíí ðàññìàòðèåàòü è íaiííéüå, êîòïðûå íáðàçóþòñÿ èç ííèíûõ íóòåíííüõ ñíèðàùåíèé. Íàïðèìåð

Íòåâðòéó! (âìåñòí Äàé íòåâðòéó!).

Íðè÷éíû òàéíâí ñíèðàùåíèÿ, íàçûâàåíííá ýëéèíñèñíí, ííäóó áûòü ðàçíííáðàçíû, íí íáû÷íí ñíèðàùåàåòñÿ òà ÷àñòü íðåäéíæåíèÿ, êîòïðà

Èáé óæå ìòìå÷æéñü áûøå, èàæäíå ïðåäëíæåíèå íáñåò â ñååå ïðåäääëåííùé ñìùñé, çàéëäüåàåìùé àåòòðíí, è íáññá÷èååò ïðåä

Â ñéö÷å äæäíå ïðåäëíæåíèå èååò ññòåòñòåóþùéé ååéòòð öåëè. Òàéèè íáðàçíí, òåéñò ïðåäääëèòü êåé ñòðóéòó:

Vbegin — íà÷æéñü õåëü, áûðàæäíà ðåðåç íà÷æéñü ñçàäàííùé èíñðäèíàòàíè;

Vend èíñá÷íà õåëü, áûðàæäíà ðåðåç èíñá÷íùé ååéòòð ñçàäàííùé èíñðäèíàòàíè;

Z åèä ñâýçè íæäó íà÷æéñü ãåéòòðíí è èíñá÷íùé ååéòòðíí.

Â èà÷åñòå ååéòòðíà ññòåòñòåóþùé, èíñðäèíà ïðåäääëåííùé ñéïåà, ííýòéé, èíñííùå ãðóííù, íòääëüíùå ïðåäëíæåíèé, ñìùñëíåñå åðóííù

Íññéüé òðè ññòåòñòåóþùé, èíñðäèíà ïðåäääëåííùé, èíñðäèíà ïðåäääëåííùé ñéïåà, ííýòéé, èíñííùå ãðóííù, íòääëüíùå ïðåäëíæåíèé, ñìùñëíåñå åðóííù

1) ïðññòíé ååéòòð: = () èéè V = ();

2) íóëåâíé ååéòòð: = (\emptyset);

3) íéíùé ååéòòð: = (,) ñí ñâýçüþ Z;

4) íóññòíé ååéòòð: = (,) áåç ñâýçè Z;

5) ëåâûé âåêòíð: = ();

6) iðàâûé âåêòíð: = ().

Â ñâíþ í-åðåäü âåêòíðû è ííãóò ñíñòíýöü èç ííäâåêòíðâ, êàê íòäåëüíûõ ñàíñòíýöåëüíûõ ÷àñòåé êííðæéíàò, iðèíàäéåæàùèô ò

Êàæääý êííðæéíàò àìååò ñâíè àòðèáóòû atr. Åòðèáóòàìè ííãóò ýâëýöüñý âðåìåííûå èëè iðíñòðàíñòååííûå õàðàéòåðèñòëëë êííðæéíàò

Ñíñòàâ êííðæéíàò âåêòíðà ííðåäåëýåòñý ñëíæíñòüþ ííñòðåíèý íðåäëíæåíèý. Â íáúåì ñëó÷àå ñíñä÷éíåíñòü íòäåëüíûõ ÷àñòåé ííñòðåíèý

Ðàññííòðèì iðèíàíåíèå ííñàííé âûøå âåêòíðîé ííäåëè íà êííðåòíí iðèíåðå.

Íðíàíåëèçèðóàì ñëåäóþùåå íðåäëíæåíèå.

Â íâñå âðåìåíà ëþäè ñòàëëèâàþòñý ñ íáíèìè è òåìè æå íðíáëåìàìè ýéíííèéè.

Äàíííá íðåäëíæåíèå ííæåò áûòü íðåäñòàåéåíí â âåêòíðîé õíðìå: Vp(x1; y1) ñí ñâýçüþ âèäà z1, èëè å óíðíùåííé ôíðìå Vp(x1; y1) (

ãäå êííðæéíàòà x1 = {ëþäè};

êííðæéíàòà y1 = {íðíáëåìû ýéíííèéè};

ñâýçü âèäà z1 = (ñòàëëëâàþòñy).

Iðè ýòî àòðëáóòàìè êñðäëíàòû ñ1 ýâëýþòñy atrx = (âî âñå âðåìáà), à àòðëáóòàìè êñðäëíàòû y1 ýâëýþòñy atry = (iáie è òå æå)

Ià iññâå ååëòññá iñðäëñòàâëåíèÿ iññóò áûòü ðåøåíû íåëòññá iñðäëíàòû ñðåáñòàâëåíèÿ èíññòàâëåíèÿ

– ñîëðàùåíèå íáúåìà èññññé èíññòàâëåíèÿ èíññòàâëåíèÿ äëÿ åññññé iñðäëíàòû ñðåáñòàâëåíèÿ èíññòàâëåíèÿ

– ñèíòåç òåëñòà ñ èñññüçñâàíèå èíññòàâëåíèÿ èçåëåâàåíèé èç áàç çíàíèé.

Â ñëåäöþùåi iñðäàñðàôå ñðñññàòðèâàåòñy ååññòðè÷åññàÿ èíòåðñðåðàöèÿ ñðåáñòàâëåíèÿ èíññòàâëåíèÿ

3. Iðèíåíàíèå ñðåáññòàâëåíèÿ èíññòàâëåíèÿ iñðäëñòàâëåíèÿ èíññòàâëåíèÿ èíññòàâëåíèÿ

Ðàññññòðåíàÿ áûøå ååëòññá iññóò áûòü ñðåáñòàâëåíèÿ ñðåáñòàâëåíèÿ èíññòàâëåíèÿ èíññòàâëåíèÿ èíññòàâëåíèÿ èíññòàâëåíèÿ

Ià ðèñ. 1 iñðäëñòàâëåíà óïðñùåíàÿ èíòåðñðåðàöèÿ ååëòññá iñðäëñòàâëåíèÿ ñðåáñòàâëåíèÿ èíññòàâëåíèÿ èíññòàâëåíèÿ èíññòàâëåíèÿ

Ðèñ. 1. Ååëòññá iñðäëñòàâëåíèå ñðåáññòàâëåíèå

Ià iñðäëñòàâëåíí ðèñññé åëàçàíû ñðè ååëòññá Vp1(x1, y1) (z1); Vp2(x2, y2) (z2); Vp3(x3, y3) (z3) è èõ iñðåññéèÿ íà iññññòû x, y, atr.

Àòðèáóòú ñíäóò èìåòü âðåìáííúå, iðñòðàíñòâåííúå è äðóæå èçìåðýåíúå õàðàéòåðèñòèéè.

Êñðäèíàòú õi ñíðåäåéýþò íà÷àëüíúå êñðäèíàòú âåêòîðà. Êñðäèíàòú yi ñíðåäåéýþò êííå÷íúå êñðäèíàòú âåêòîðà.

Èññïäý èç iðåäûäóùåâí ñíèñàíéý âåêòîð V ñíðåäåéýåò êííå÷íóþ öåëü ðàññìàòðèåàåíé åäèíèöü òåêñòà è èìååò ñòðóêòóðó âåêò

Ðàññìòðèì èññïäüçîàíéå ñíèñàíííé ñíäåéè íà iðèìåðå.

Íóñòü çàäàíú òðè âåêòîðà:

Vp1(x1, y1) (z1); Vp2(x1, y2) (z2); Vp3(x1, y3) (z3).

Íà ðèñ. 2 ñíèàçàíú íåêòîðûå âíçííæíúå âàðèàíòû ãåñìàòðè÷åñéíé èíòåðïðåòàöèè áçàèíäåéñòåèý òðåð ãåêòîðâ.

à

á

â

Đèñ. 2. Åâññâòðè÷åñêàÿ èíòåðïðåòàöèÿ âçàèñäåéñòâèÿ òðåô âåêòîðîâ

Iøèìåð, iøåäñòàâëåííûé íà ðèñ. 2à, iiêàçûâàåò, ÷òî êññðäèíàòà õ1 âåêòîðà iøåäñòàâëÿò ñâáîé èåðàðõèþ iiýòèé. Íaïøèìåð, õàðà

Íà ðèñóíéå 2á iøåäñòàâëåíà åâññâòðè÷åñêàÿ èíòåðïðåòàöèÿ âçàèñäåéñòâèÿ äðóñé ãðóññû èç òðåô âåêòîðîâ:

Vp1(x1, y1) (z1); Vp2(x2, y2) (z2); Vp3(x3, y3) (z3),

ãäå y1 = x2, y2 = x3.

Ôàêòè÷åñêè ñâáîñóíñòü ýòèõ òðåô âåêòîðîâ iøåäñòàâëÿò íâéòîðûé ðåçóëüòèðóþùèé âåêòîð V(x1, y3) (z'); ñññòâåòñòâó

Nëåäóþùèé iøèìåð (ðèñ. 2â) èëëþñòðèðóåò íåçàâèñèñòü öåëåé â iøèåäåäåíííí íàáîðå âåêòîðîâ:

Vp1(x1, y1) (z1); Vp2(x1, y2) (z2); Vp3(x3, y3) (z3),

ãääå y2 = x3.

Äðóäèì ïðèìáíáíéáí áåêòíðííé ïäääéè ýâëýåòñý áíçííæíñòü ðåàëèçàöèè ñèíòåçà òåêñòíâíé èíôíðíàöèè.

Íðåäííæèì, ÷òí ðåøàåòñý çàääà÷à, ñâýçáííàÿ ñ ðàñêðûòèáí ïíýòèÿ õ1. Â ýòíí ñëó÷àå, áåêòíð öåëè äëÿ ïèñàíèÿ ïðåäääéåíûõ ïðéíñòðóèðíåíèÿ ïäöåëåé.

Äñíñòèì, á áàçå çíàíèé ïíýòèå õ1 ïðåäääéåííá íà ííæåñòåå åíòíëíäéé ÷åðåç áåêòíð Vp0(x1, y1) (z1). Â ñâíþ í÷åðåäü ïäääéòíð y1 ðàñêðûòèÿ ïíýòèÿ õ1. Íåðàíèçì ðàçååðòûåàíèÿ áåêòíðà äëÿ ïèñàíèÿ ïðòåññíà è ýâëåíèé ïæååò áúòü áâîÿèè: éèáí íà íñííåå ñíóúéå.

Íà íñííåå ïðåäëíæåííé ïäääéè ðàçðàáíòàíà òåôííëíäé ìáðàáíòéè òåêñòíâíé èíôíðíàöèè íà íñííåå áåêòíðííé ïäääéè òåêñòà.

Çàêëþ÷åíèå

Ðàññííòðåííûå á äàíííé ñòàòüå íñííåíúå ïíëíæåíèÿ òåôííëíäéè áåêòíðííáí ïðåäñòååëåíèÿ òåêñòíâíé èíôíðíàöèè è àåòííàöèçàöèÿ ýòíí

– ñíçääíèÿ ïðòåññèíàëüíûõ ñèñòåí è áàç çíàíèé;

– ïäääðæéè ïðòåññèíàëüííé äåëòåëüíñòè ðàáíòíèéåíà ðàçëè÷íûõ ìòðàñëåé;

– iñáùáñòáíéý óðíáíý êííáòáíóèè ñiáöèáéèñòíâ çà ñ÷áð ñéó÷áíéý áíçííæíñòé áúñòðíáí áíáéèçà è iðááñòááéáíéý á óáíáíé

– iðíááááíéý ñéíóáçà òáéñòáñòáí ãíéóíåíóíâ ñ ðàçéè-ííé ñòáíáíüþ íáíáùáíéý èíóíðíàöèè;

– áâòííàòéçàöèè iðíóáññíâ óíðíèðíááíéý ñéñòáíû ííóíéíæé á òíé èéè èíé iðíóáññéíáéüíé íáéáñòé;

– iðíááááíéý íáíðááéáíííííâí ñéñêá è óééüòðàöèè òáéñòáñòáí ãíéóíåíóíâ;

– áâòííàòé-áñéíâí ðåôåðèðíááíéý òáéñòíâ ãíéóíåíóíâ.

Summary

The paper considers the approaches associated with the vector representation of the textual information. The particularity of the approach under consideration is in the determination of the goal functions of separate sentences and representation of them in the form of some local vectors on which basis a global vector is built that determines the semantic component of the text on the whole. Various aspects of application of the proposed approach are considered.

Ñiñéñíé ëèòáðàòóðû

Àðóòþííâ 2005 Àðóòþííâ, I. Ä. Íðåäéíæáíéá è ááí ñiñéñíé / I. Ä. Àðóòþííâ. I. : ÓÐÑÑ, 2005.

Íí-áííâ è äð. 2005 Íí-áííâ, N. Ä. Íðèíáíáíéá ñòáðèñòé-áñéèé ñåðíáíâ äéý ñåìáíðé-áñéíâí áíáéèçà òáéñòá / N. Ä. Íí-áííâ, Á. I. Äé-

Éàðàóöíâ è äð. 1982 Éàðàóöíâ, P. I. Ðóññééé ñåìáíðé-áñéèé ñéíáàðü. Ííûò áâòííàòé-áñéíâí ñíñòðíáíéý òáçàóðóñà: ìò ìííýøéý

Ðóáàøééí è äð. 1998 Ðóáàøééí, Á. Ø. Ñåìáíðé-áñéèé (êííóáíðóáéüíûé) ñéíáàðü äéý èíóíðíàöèííûó òåðííéíæé. x. 1 / Á. Ø. Ðó-

Ñíñéðéí è äð. 2005 Ñíñéðéí, Á. Á. Íðíàêò ÄÈÀËÈÍÃ, COM-íáúâêò Goldrml / Á. Á. Ñíñéðéí, Á. Á. Íáíêðàòíâ. Í. : Äèàëíã, 2005.

Ôéíí 1999 Ôéíí Á. É. Í ðíëè ìàøèííáí íáó÷åíèý á èíðåëëåêòóàëüíûõ ñèñòåìàõ // ÍÒÈ. Ñåð. 2. 1999. ¹ 12. Ñ. 1–3.

Salton 1989 G. Salton. Automatic Text Processing. Addison-Wesley Publishing Company, Inc., Reading, MA, 1989.

Salton et al. 1994 G. Salton, J. Allan, and C. Buckley. Automatic structuring and retrieval of large text files. Communications of the ACM, 37(2), February 1994.

Todd et al. Todd A. Letsche and Michael W. Berry. Large-Scale Information Retrieval with Latent Semantic Indexing. URL: <http://www.cs.utk.edu/~berry/sc95/sc95.html>.