

Òàçàóðõñ ìì èìðìóñííé èèíáàèñòèèá

Áàóìð Áèèòìð Ìááèíáè- Çàòàðíá
02.09.2010 á.
Ìñèááíáá Ìáííáèáíéá 03.09.2010 á.

The paper is devoted to the construction of the Russian thesaurus on Corpus Linguistics. Linguistic resource involved in research is the Russian corpus on Corpus Linguistics developed in St.-Petersburg State University and Institute of Linguistic Studies and different vocabularies. Semi-automatic terminology extraction is performed with the help of linguistic and statistical tools which allow to generate lists of single-word and multi-word terms supplied with frequency data and lexical-syntactic patterns. Lexical-syntactic patterns are used in the analysis of contexts which contain definitions of terms, expose interrelations between terms, provide their synonyms, translation equivalents, etc.

Á ìáñòìÿóáá áðáìÿ á ìáèáñòè èìðìóñííé èèíáàèñòèèè ìòñòóòñòááò -áòèáÿ ñèñòáíá ìáó-ííé òáðìèííéíáèè. Èìðìóñííé èèíáàèñòèè ììòðíáíéÿ è èñííèóçíááíéÿ èèíáàèñòèè-áñèèò èìðìóñííé ñ èñííèóçíááíéáì èìííóðòáðíóò òáðìííéíáèè. Ìíá ìáçááíéáì èèíáàèñòèè-áñèèè áèáá è ìðááíáçíá-áíííé áèÿ ðáðáíéÿ èííèðáòíóò èèíáàèñòèè-áñèèò çááá-. Èìðìóñííé èèíáàèñòèèá ìáðíáèòñÿ ìá ìáðáñá-áíéè çáá òáðìèííéíáèè. Áí-ìáðáóò, ÿòí áñòáñòááííí, ó-èòóááÿ áá ìááááíáá ìðìèñòíáèáíéá. Áí-áòíðóò, èìðìóñííé èèíáàèñòèèá èáè ìááèóíá ñèèááóááòóñÿ á ìááðáò áíáèèèñéíáí ÿçóèá. Ìÿÿòíó, ðáçðáááòóááÿ òáçàóðõñ ìì èìðìóñííé èèíáàèñòèèá, áñòáñòááííí çááóíáòóñ èò áíáèèèñéèá ÿéáèááèáíóó. Á ìáñòìÿóáá áðáìÿ òáèèò òáçàóðõñíá ìá ñóóáñòááò.

Èðóá èáèñèèè, ìòíííÿóáèñÿ è ñòáðá èìðìóñííé èèíáàèñòèèè è ìíáèáæáóáé òíðíáèèçàòèè, ì-áð-áí èìáðóèèèñÿ ÿíóèèèíááè-áñèè ñíáèóí ìáðááèáíéÿ. Ìÿÿòíó ì ðè ðáçðááíéèá òáçàóðõñá ìó èñííáèè èç ìáèè-èÿ ðÿáá ìíííáèè, èìòíðóá ìíáóò áóòó èñííèóçíááíí, á Èèíáàèñòèè-áñèèè 2005, Òçóèíçíáíéá 2007 è áð.].

Ìðè ììòðíáíéè òáçàóðõñíá èñííèóçóðòñÿ ðáçèè-íóá ìáðíá. Ìáè èç ñííííáíá èçó-áíéÿ òáðìèííéíáèè ìáðáçóíááááò ìðáááðèòáèó ðááíóá ìðìèçáíáèòñÿ ìèñáíéá òáðìèííéíáèè èìðìóñííé èèíáàèñòèèè èíáíí ìá ìíííáíéè èíáèè-ìíÿòèéíáí áíáèèçá ñ èñííèóçíááíéáì á ìáðíáíéè ñííòááèòó èìðìóñííé òáèñòíá ìì èìðìóñííé èèíáàèñòèèá è ñíçáááòó ñéíáíéè òáçàóðõñá ìáííðááñòááííí ìá æèáíí òáèñòíá ìá ðóññéíí è áíáèèèñéíí ÿçóèáò, èìòíðóè ìáðèíáè-áñèè ìííéíáòñÿ ìáóíè áíéóíáíóáè.

Ðáçðááíóáì ìáðíá è ììòðíáíéè òáçàóðõñíá ìðááíáòíóò ìáèáñòáé, ìðááííéáááðóèè áóááèáíéá èèð-ááóò ìíÿòèè èç ìííáñòáà èç ááñèðèèòíðáèè, èçáèá-áíííé èç ó-ááíóò ìíííáèè ìì èìðìóñííé èèíáàèñòèèá, á -áñòííòè [Çàòàðíá 2005]. Ìá áòíðíí ÿòáíá ìðíáíéè òáèñòó (ìðáæáá áñááí, ìóáèèèáòèè èííóáðáíóèè ìì èìðìóñííé èèíáàèñòèèá: «Èìðìóñííé èèíáàèñòèèá è èèíáàèñòèè-áñèèá ááçó á -áñòííòè; 2004, 2006, 2008» è áð.). Ìá ááçá ÿòáíá èìðìóñííé ñíçáááòñÿ ìáðááÿ ááðñéÿ òáçàóðõñá ìì èìðìóñííé èèíáàèñòèèá. Òáðíí ÿ òáðáèòáðèçóðóèè ìáèáñòó «èìðìóñííé èèíáàèñòèèá» è ìèè-áðóèòñÿ áóñíéíé èìðìóñííéèííóò è -áñòííóò. Áèÿ ììòðíáíéÿ èáèñèèè, ó-èòóááðóèè èíáèè-ìíÿòèéíá ñòáíí ìðááíáòíé ìáèáñòè, ìðááñòááèáííóá á ÿéñíáðóíóò ìèñáíéÿ. Áèÿ ááòíáòèçèðíá èíéèíáèèè, ìíííáááðóèèñÿ ìá èçááñòíóò ìáðáò áññíòèáòèè (MI, T-score, Log-Likelihood).

Òáçàóðõñ áóááò ááçèðíáòóñÿ ìá ìíóèíáèè èìðìóñííé èèíáàèñòèèè, ðáçðáááòóáááíéè ìá èáòááðá ìáðáíáèè-áñéíé èèíáàèñòèèè èìííóðòáðíé èèíáàèñòèèè, ðáçðááíóáíáÿ á ìáííéáèðñéá á ðáèèò ìðíáèòá ìì ñíçááíéè ìðòáèá çíáíéè ìì èìííóðòáðíé èèíáàèñòèèè.

Ðáçðáááòóáááíéè òáçàóðõñ áíéæáí ìáñííá-èòó ñèñòáíáòèçèðíááííá ìðááñòááèáíéá òáðìèííéíáèè á ìáèáñòè èìðìóñííé èèíáàèñòèè èìðìóñííé èèíáàèñòèèá è ìáñííá-èòó óáííáíé áíñòóí è ìèì. Ááóÿçó-ííòó òáçàóðõñá ìì çáíéèò ìá-áñòááíííí ó-áííí è ñíáòèáèèñòá ìáðááíáá ñòáòáè.

Òáçàóðõñ áóááò ìðááñòááèÿòó ñíáíé ÿéáèòðííóð ááçó ááííóò, ìáñííá-èááðóò ìì èñííáòáèÿ ìíñòóí è ìáíó. Èðíá òíáí, áóááò ñ ìííáíÿçó-íóò (ISO 5964-1985, ÁÍÑÒ Ð 7.24-2007) òáçàóðõñíá.

Èèòáðáòóðá

Άδαήηά Α.Í. Αάάαήεά á ἰδεεάαίόρ εεήάεηόεό. Νάδϑ "Íáúε εεήάεηόε-áñεεé ó-ááήε". Ì.: Ýáεοίδεάé ÓÐÑ. 2001.

Áηήεεήηά Í.É., Çαήδóεüεή Π.Α., Çαήδóεüεή Α.Α., Éήήήήήé È.Ñ., Νήεήήά Α.Α. Ðαçðááηóεά ἰδóεά çήήεé ἡ éήήήρòáðήé εεήά (á.Αόήά, Ðήññè). –Ì.: ÈÁÍÁÍ, 2008. –Ò.3. –Ñ.380-388.

Άεήάάάήά Í.Α., Ìεòδìòáήήά Í.Α. Óήδìεüήáý ἡήήήήéý εάé εήηòδòήάò ηέηòáìàòεçàòεé ááήήó á δóññéήçý-ήή éήδìóñά òáéñòήá ἡ ἡήά.: ἨήάÓ, 2008. – C. 113-121.

Άήήήήήά Α.Ç. Áήεή-δóññεéά òáδìεή ἡ ἰδεεάαήήé εεήάεηόεéά é ááòήàòé-áñεήé ἡδάðááήóεά òáéñòá. Αή. 2. Ìáòήá òáήεçά

Çáòáðήά Α.Í. Éήδìóñήáý εεήάεηόεéά: Ó-ááή-ìáòήé-áñεήά ἡήήήéá. – Ἠήά.: ἨήάÓ, 2005. – 48 ἡ.

Èεήάεηόε-áñεεé ýìεéήήήé-áñεεé ηέήάòü. Ì.: Ἠήά. Ýìεéήήήéý, 1990.

Ìεòδìòáήήά Í.Α., Çáòáðήά Α.Í. Αáòήàòεçεðήάήήé άήεç òáδìεήήήéά á δóññéήçý-ήή éήδìóñά òáéñòήá ἡ éήδìóñήήé εεήάεηόεéά «Άεεήή–2009». Ì.: 2009.

Ìεéòεήά Ἠ.Α. Óáçàòðóñ ἡ òáððáòé-áñεήé é ἰδεεάαήήé εεήάεηόεéά. - Ì., 1978.

Ìδεεάαήήά ççüεήçήήéά. Ó-ááήé (ðáá. Α.Ἠ.Άðä). Ἠήά., 1996.

Νήεήήά Α.Α., Éήήήήήé È.Ñ., Çαήδóεüεή Π.Α. Ìδìάéáìü ἡήήήéý éήήήρòáðήé εεήάεηόεéé á áεáá ἡήήήéé áéý ἡδóεά çήήéé // – Ì.: 2008. – Ἠ.482-487.

ßçüεήçήήéά. Èήδìòìàòéήήή-ἡήήήήéüé òáçàòðóñ ÈÍÉÍ ÐΑÍ. – Ì., 2007.

The Oxford handbook of computational linguistics // Mitkov Ruslan (ed.). N.Y.: Oxford university press, 2003.