

КЛАССИФИКАЦИЯ СКРЕП

N00. Разрывные скрепы

N01. ИДЕНТИФИКАЦИЯ

N02. СРАВНЕНИЕ

N03. ПРИЧИНА

N04. ВОЗМОЖНОСТЬ

N05. ПРОТИВОПОСТАВЛЕНИЕ, ПРОТИВИТЕЛЬНОСТЬ, УСТУПИТЕЛЬНОСТЬ

N06. ОБОБЩЕНИЕ-КОНКРЕТИЗАЦИЯ

N07. ОБОСОБЛЕНИЕ-ОГРАНИЧЕНИЕ

N08. ЦЕЛЬ

N09. ВРЕМЯ

N10а. ПРЕДСКАЗУЕМОСТЬ, СЧЁТНОСТЬ, ПОСТОЯНСТВО, УСТОЙЧИВОСТЬ

N10б. НЕПРЕДСКАЗУЕМОСТЬ

N11. УСЛОВИЕ

N12. ВВОД ТЕМЫ (ИСТОЧНИК ИНФОРМАЦИИ)

N13. МЕСТО

N14. ОБЪЕДИНЕНИЕ

ОСНОВНЫЕ ОПРЕДЕЛЕНИЯ

C_j – коннектор _ номер _ j

T_k – текст _ номер _ k

$L(T_k)$ – длина _ текста _ T_k _ (количество _ словоформ)

$n_j(T_k)$ – количество _ появлений _ коннектора _ C_j _ в _ тексте _ T_k

частота коннектора C_j в тексте T_k :

$$f_j(T_k) = \frac{n_j(T_k) * 1000000}{L_k} =$$

$$\frac{\text{количество появлений коннектора } C_j \text{ в тексте } T_k, \text{ умноженное на } 1000000}{\text{длина текста } T_k}$$

средняя частота коннектора C_j в текстах некоторого автора =

$$f_j^* = \frac{\text{суммарное количество появлений коннектора } C_j \text{ в текстах автора}}{\text{суммарная длина текстов автора}} * 1000000 =$$

$$\frac{n_j(T_1) + n_j(T_2) + \dots + n_j(T_n)}{L(T_1) + L(T_2) + \dots + L(T_n)} * 1000000 = \frac{\sum_k n_j(T_k)}{\sum_k L(T_k)} * 1000000 = \sum_k f_j(T_k) * \frac{L_k}{\sum_m L(T_m)}$$

Для каждого текста вычисляется последовательность частот:

$$f_1(T_k), f_2(T_k), f_3(T_k), \dots$$

Мы можем рассматривать эти числа как координаты точки T_k в многомерном пространстве.

Текстам каждого автора соответствует некоторое множество – облако точек. Удобно считать, что в каждой точке сосредоточена масса, равная объёму текста. Тогда числа f_j^* – это координаты центра тяжести облака точек. Мы можем для наглядности считать, что точки каждого автора закрашены своим цветом.

Суммы $\sum_j f_j(T)$ для разных текстов одного и того же автора могут изменяться в довольно широком интервале. При этом интервалы, соответствующие текстам разных авторов, могут сильно перекрываться.

РАССТОЯНИЕ МЕЖДУ ТЕКСТАМИ

Расстояние между текстами T_a и T_b =

$$\sum_j |f_j(T_a) - f_j(T_b)|$$

Расстояние между текстом и множеством отличных от него текстов некоторого автора

$$\sum_j |f_j(T) - f_j^*|$$

Расстояние от текста до облака точек, соответствующих текстам некоторого автора, определяется как расстояние от точки, соответствующей тексту, до центра тяжести облака точек. При этом точка, соответствующая исследуемому тексту, исключается из облака точек, соответствующих текстам автора, если она ему принадлежала.

Таблица авторов

ТАБЛИЦА ПОДГОТОВЛЕННЫХ ДЛЯ ОБРАБОТКИ ТЕКСТОВ

Номер текста	Номер предложения	Номер слова	Слово	Знаки препинания после слова
92	1	1	Повести	,
92	1	2	Изданные	
92	1	3	пасичником	
92	1	4	Рудым	
92	1	5	Паньком	
92	2	1	Часть	
92	2	2	ПЕРВАЯ	
92	3	1	Предисловие	«
92	4	1	Это	
92	4	2	что	
92	4	3	за	
92	4	4	невидадь	: «
92	4	5	Вечера	
92	4	6	на	
92	4	7	хуторе	
92	4	8	близ	
92	4	9	Диканьки	»?

ПРИМЕР ПЕРЕСЕЧЕНИЯ КОННЕКТОРОВ

933	Нет ₁ , матушка ₂ не ₃ обижу ₄ , – говорил ₅ он ₆ , а₇ между₈ тем₉ отирал ₁₀ рукою ₁₁ пот ₁₂ , который ₁₃ в ₁₄ три ₁₅ ручья ₁₆ катился ₁₇ по ₁₈ лицу ₁₉ его ₂₀ .	7; 8; 9;	1	ТОТ (2-15), КОТОРЫЙ	9; 13;	0
				МЕЖДУ ТЕМ	8; 9;	2
				НЕ (2-15), А	3; 7;	1