

Отчет по заданию 1

“Сглаживание”

Основные задачи

- Взять отрывок из стихотворения “Дом, который построил Джек”
- Разбить стихотворение на обучающую и тестовую части
- Убрать из текста знаки препинания, заменив при этом дефис пробелом
- Разбить обучающую выборку на униграммы и биграммы, подсчитать их вероятности
- Провести для тестовой выборки аддитивное сглаживание и сглаживание с подбором параметра λ согласно закону Лидстоуна
- Подсчитать перплексию для униграмм и биграмм

Униграммы

Обучение на обучающем множестве

Исходная вероятность униграммы:

$$P(w_i) = \frac{c_i}{N}$$

Исходные вероятности

#	Aa Words	# Probabilities
0	<u>во</u> т	0.02702702702702703
1	до <u>м</u>	0.02702702702702703
2	<u>ко</u> то <u>р</u> ый	0.08108108108108109
3	<u>по</u> стро <u>и</u> л	0.08108108108108109
4	дже <u>к</u>	0.08108108108108109
5	<u>а</u>	0.05405405405405406
6	<u>э</u> то	0.05405405405405406
7	<u>п</u> шени <u>ц</u> а	0.02702702702702703

#	Aa Words	# Probabilities
8	<u>которая</u>	0.08108108108108109
9	<u>в</u>	0.10810810810810811
10	<u>темном</u>	0.05405405405405406
11	<u>чулане</u>	0.05405405405405406
12	<u>хранится</u>	0.05405405405405406
13	<u>доме</u>	0.05405405405405406
14	<u>веселая</u>	0.02702702702702703
15	<u>птица</u>	0.02702702702702703
16	<u>синица</u>	0.02702702702702703
17	<u>часто</u>	0.02702702702702703
18	<u>ворует</u>	0.02702702702702703
19	<u>пшеницу.</u>	0.02702702702702703

Вероятности при аддитивном сглаживании (сглаживание Лапласа)

$$P_{Lap}(w_i) = \frac{c_i + 1}{N + V}$$

Сглаживание Лапласа для униграмм

#	Aa Words	# Probabilities
0	<u>вот</u>	0.03508771929824561
1	<u>кот</u>	0.017543859649122806
2	<u>который</u>	0.07017543859649122
3	<u>пугает</u>	0.017543859649122806
4	<u>и</u>	0.017543859649122806
5	<u>ловит</u>	0.017543859649122806
6	<u>синицу.</u>	0.017543859649122806
7	<u>которая</u>	0.07017543859649122
8	<u>часто</u>	0.03508771929824561
9	<u>ворует</u>	0.03508771929824561
10	<u>пшеницу.</u>	0.03508771929824561

#	Aa Words	# Probabilities
11	<u>в</u>	0.08771929824561403
12	<u>темном</u>	0.05263157894736842
13	<u>чулане</u>	0.05263157894736842
14	<u>хранится</u>	0.05263157894736842
15	<u>доме</u>	0.05263157894736842
16	<u>построил</u>	0.07017543859649122
17	<u>джек</u>	0.07017543859649122

Сглаживание с подбором λ (закон Лидстоуна)

$$P_{Lid}(w_i) = \frac{c_i + \lambda}{N + B\lambda}$$

$$\lambda = 0.25$$

Закон Линдстоуна с параметром 0.25

#	Aa Words	# Probabilities
0	<u>вот</u>	0.02976190476190476
1	<u>кот</u>	0.005952380952380952
2	<u>который</u>	0.07738095238095238
3	<u>пугает</u>	0.005952380952380952
4	<u>и</u>	0.005952380952380952
5	<u>ловит</u>	0.005952380952380952
6	<u>синицу.</u>	0.005952380952380952
7	<u>которая</u>	0.07738095238095238
8	<u>часто</u>	0.02976190476190476
9	<u>ворует</u>	0.02976190476190476
10	<u>пшеницу.</u>	0.02976190476190476
11	<u>в</u>	0.10119047619047619
12	<u>темном</u>	0.05357142857142857
13	<u>чулане</u>	0.05357142857142857
14	<u>хранится</u>	0.05357142857142857
15	<u>доме</u>	0.05357142857142857

#	Aa Words	# Probabilities
16	<u>построил</u>	0.07738095238095238
17	<u>джек</u>	0.07738095238095238

$$\lambda = 0.5$$

Закон Линдстоуна с параметром 0.5

#	Aa Words	# Probabilities
0	<u>вог</u>	0.031914893617021274
1	<u>кот</u>	0.010638297872340425
2	<u>который</u>	0.07446808510638298
3	<u>пугает</u>	0.010638297872340425
4	<u>и</u>	0.010638297872340425
5	<u>ловит</u>	0.010638297872340425
6	<u>синицу.</u>	0.010638297872340425
7	<u>которая</u>	0.07446808510638298
8	<u>часто</u>	0.031914893617021274
9	<u>ворует</u>	0.031914893617021274
10	<u>пшеницу.</u>	0.031914893617021274
11	<u>в</u>	0.09574468085106383
12	<u>темном</u>	0.05319148936170213
13	<u>чулане</u>	0.05319148936170213
14	<u>хранится</u>	0.05319148936170213
15	<u>доме</u>	0.05319148936170213
16	<u>построил</u>	0.07446808510638298
17	<u>джек</u>	0.07446808510638298

$$\lambda = 0.75$$

Закон Линдстоуна с параметром 0.75

#	Aa Words	# Probabilities
0	<u>вог</u>	0.03365384615384615
1	<u>кот</u>	0.014423076923076924
2	<u>который</u>	0.07211538461538461

#	Aa Words	# Probabilities
3	<u>пугает</u>	0.014423076923076924
4	<u>и</u>	0.014423076923076924
5	<u>ловит</u>	0.014423076923076924
6	<u>синицу.</u>	0.014423076923076924
7	<u>которая</u>	0.07211538461538461
8	<u>часто</u>	0.03365384615384615
9	<u>ворует</u>	0.03365384615384615
10	<u>пшеницу.</u>	0.03365384615384615
11	<u>в</u>	0.09134615384615384
12	<u>темном</u>	0.052884615384615384
13	<u>чулане</u>	0.052884615384615384
14	<u>хранится</u>	0.052884615384615384
15	<u>доме</u>	0.052884615384615384
16	<u>построил</u>	0.07211538461538461
17	<u>джек</u>	0.07211538461538461

Перплексия

+1 : 16.175358888003007

$\lambda = 0.25$: 20.976732603415496

$\lambda = 0.5$: 18.23156241301383

$\lambda = 0.75$: 16.946076133700398

Биграммы

Обучение на обучающем множестве

$$P(w_n|w_{n-1}) = \frac{c(w_n w_{n-1})}{c(w_{n-1})}$$

Исходные вероятности

#	Aa Words	# Probabilities
0	('Вот', '_дом')	1

#	Aa Words	# Probabilities
1	(<u>дом</u> ', ' <u>который</u> ').	1
2	(<u>который</u> ', ' <u>построил</u> ').	1
3	(<u>построил</u> ', ' <u>джек</u> ').	1
4	(<u>джек</u> ', ' <u>а</u> ').	0.6666666666666666
5	(<u>а</u> ', ' <u>это</u> ').	1
6	(<u>это</u> ', ' <u>пшеница</u> ').	0.5
7	(<u>пшеница</u> ', ' <u>которая</u> ').	1
8	(<u>которая</u> ', ' <u>в</u> ').	0.6666666666666666
9	(<u>в</u> ', ' <u>темном</u> ').	0.5
10	(<u>темном</u> ', ' <u>чулане</u> ').	1
11	(<u>чулане</u> ', ' <u>хранится</u> ').	1
12	(<u>хранится</u> ', ' <u>в</u> ').	1
13	(<u>в</u> ', ' <u>доме</u> ').	0.5
14	(<u>доме</u> ', ' <u>который</u> ').	1
15	(<u>это</u> ', ' <u>веселая</u> ').	0.5
16	(<u>веселая</u> ', ' <u>птица</u> ').	1
17	(<u>птица</u> ', ' <u>синица</u> ').	1
18	(<u>синица</u> ', ' <u>которая</u> ').	1
19	(<u>которая</u> ', ' <u>часто</u> ').	0.3333333333333333
20	(<u>часто</u> ', ' <u>ворует</u> ').	1
21	(<u>ворует</u> ', ' <u>пшеницу</u> ').	1
22	(<u>пшеницу</u> ', ' <u>которая</u> ').	1

Аддитивное сглаживание

$$P_{Lap}(w_n|w_{n-1}) = \frac{c(w_n w_{n-1}) + 1}{c(w_{n-1}) + V^2}$$

Сглаживание Лапласа для биграмм

#	Aa Words	# Probabilities
0	(<u>вот</u> ', ' <u>кот</u> ').	0.001890359168241966
1	(<u>кот</u> ', ' <u>который</u> ').	0.001890359168241966

#	Aa Words	# Probabilities
2	('который', 'пугает').	0.001890359168241966
3	('пугает', 'и').	0.001890359168241966
4	('и', 'ловит').	0.001890359168241966
5	('ловит', 'синицу').	0.001890359168241966
6	('синицу', 'которая').	0.001890359168241966
7	('которая', 'часто').	0.0037593984962406013
8	('часто', 'ворует').	0.0037735849056603774
9	('ворует', 'пшеницу').	0.0037735849056603774
10	('пшеницу', 'которая').	0.0037735849056603774
11	('которая', 'в').	0.005639097744360902
12	('в', 'темном').	0.005628517823639775
13	('темном', 'чулане').	0.005649717514124294
14	('чулане', 'хранится').	0.005649717514124294
15	('хранится', 'в').	0.005649717514124294
16	('в', 'доме').	0.005628517823639775
17	('доме', 'который').	0.005649717514124294
18	('который', 'построил').	0.007518796992481203
19	('построил', 'джек').	0.007518796992481203

Сглаживание с подбором λ

Общая формула подсчета вероятности:

$$P_{Lid}(w_1 \dots w_n) = \frac{C(w_1 \dots w_n) + \lambda}{N + B\lambda}, \quad B = V$$

где V — словарь n-грамм.

$$\lambda = 0.25$$

Закон Линдстоуна при безусловном подсчете частотности биграмм с параметром 0.25

#	Aa Words	# Probabilities
0	('вот', 'кот').	0.005988023952095809
1	('кот', 'который').	0.005988023952095809

#	Aa Words	# Probabilities
2	(<u>'который'</u> , ' <u>пугает'</u>).	0.005988023952095809
3	(<u>'пугает'</u> , ' <u>и'</u>).	0.005988023952095809
4	(<u>'и'</u> , ' <u>ловит'</u>).	0.005988023952095809
5	(<u>'ловит'</u> , ' <u>синицу'</u>).	0.005988023952095809
6	(<u>'синицу'</u> , ' <u>которая'</u>).	0.005988023952095809
7	(<u>'которая'</u> , ' <u>часто'</u>).	0.029940119760479042
8	(<u>'часто'</u> , ' <u>ворует'</u>).	0.029940119760479042
9	(<u>'ворует'</u> , ' <u>пшеницу'</u>).	0.029940119760479042
10	(<u>'пшеницу'</u> , ' <u>которая'</u>).	0.029940119760479042
11	(<u>'которая'</u> , ' <u>в'</u>).	0.05389221556886228
12	(<u>'в'</u> , ' <u>темном'</u>).	0.05389221556886228
13	(<u>'темном'</u> , ' <u>чулане'</u>).	0.05389221556886228
14	(<u>'чулане'</u> , ' <u>хранится'</u>).	0.05389221556886228
15	(<u>'хранится'</u> , ' <u>в'</u>).	0.05389221556886228
16	(<u>'в'</u> , ' <u>доме'</u>).	0.05389221556886228
17	(<u>'доме'</u> , ' <u>который'</u>).	0.05389221556886228
18	(<u>'который'</u> , ' <u>построил'</u>).	0.07784431137724551
19	(<u>'построил'</u> , ' <u>джек'</u>).	0.07784431137724551

$$\lambda = 0.5$$

Закон Линдстоуна при безусловном подсчете частотности биграмм с параметром 0.5

#	Aa Words	# Probabilities
0	(<u>'вот'</u> , ' <u>кот'</u>).	0.010526315789473684
1	(<u>'кот'</u> , ' <u>который'</u>).	0.010526315789473684
2	(<u>'который'</u> , ' <u>пугает'</u>).	0.010526315789473684
3	(<u>'пугает'</u> , ' <u>и'</u>).	0.010526315789473684
4	(<u>'и'</u> , ' <u>ловит'</u>).	0.010526315789473684
5	(<u>'ловит'</u> , ' <u>синицу'</u>).	0.010526315789473684
6	(<u>'синицу'</u> , ' <u>которая'</u>).	0.010526315789473684
7	(<u>'которая'</u> , ' <u>часто'</u>).	0.031578947368421054

#	Aa Words	# Probabilities
8	('часто', 'ворует').	0.031578947368421054
9	('ворует', 'пшеницу').	0.031578947368421054
10	('пшеницу', 'которая').	0.031578947368421054
11	('которая', 'в').	0.05263157894736842
12	('в', 'темном').	0.05263157894736842
13	('темном', 'чулане').	0.05263157894736842
14	('чулане', 'хранится').	0.05263157894736842
15	('хранится', 'в').	0.05263157894736842
16	('в', 'доме').	0.05263157894736842
17	('доме', 'который').	0.05263157894736842
18	('который', 'построил').	0.07368421052631578
19	('построил', 'джек').	0.07368421052631578

$$\lambda = 0.75$$

Закон Линдстоуна при безусловном подсчете частотности биграмм с параметром 0.75

#	Aa Words	# Probabilities
0	('вот', 'кот').	0.014084507042253521
1	('кот', 'который').	0.014084507042253521
2	('который', 'пугает').	0.014084507042253521
3	('пугает', 'и').	0.014084507042253521
4	('и', 'ловит').	0.014084507042253521
5	('ловит', 'синицу').	0.014084507042253521
6	('синицу', 'которая').	0.014084507042253521
7	('которая', 'часто').	0.03286384976525822
8	('часто', 'ворует').	0.03286384976525822
9	('ворует', 'пшеницу').	0.03286384976525822
10	('пшеницу', 'которая').	0.03286384976525822
11	('которая', 'в').	0.051643192488262914
12	('в', 'темном').	0.051643192488262914
13	('темном', 'чулане').	0.051643192488262914

#	Aa Words	# Probabilities
14	('чулане', 'хранится').	0.051643192488262914
15	('хранится', 'в').	0.051643192488262914
16	('в', 'доме').	0.051643192488262914
17	('доме', 'который').	0.051643192488262914
18	('который', 'построил').	0.07042253521126761
19	('построил', 'джек').	0.07042253521126761

Перплексия

+1 : 273.7287283373542

$\lambda = 0.25$: 139.09749617838835

$\lambda = 0.5$: 200.17679598819893

$\lambda = 0.75$: 242.19156910489474

Сглаживание с подбором λ . Подсчет частотности биграмм в зависимости от первого слова в биграмме

В ходе выполнения задания было выяснено, что существует еще один подход к сглаживанию с параметром λ , когда мы не просто подсчитываем вероятность появления биграммы как независимого объекта, а определяем частотность появления в биграмме слова после конкретного слова (формула почти аналогична формуле для сглаживания +1)

Общая формула подсчета вероятности:

$$P_{Lid}(w_n|w_1...w_{n-1}) = \frac{C(w_1...w_n) + \lambda}{C(w_1...w_{n-1}) + B\lambda}, \quad B = V^2$$

где V — словарь n-грамм.

$\lambda = 0.25$

Закон Линдстоуна с параметром 0.25

#	Aa Words	# Probabilities
0	('вот', 'кот').	0.001890359168241966
1	('кот', 'который').	0.001890359168241966

#	Aa Words	# Probabilities
2	(<u>который</u> , <u>пугает</u>).	0.001890359168241966
3	(<u>пугает</u> , <u>и</u>).	0.001890359168241966
4	(<u>и</u> , <u>ловит</u>).	0.001890359168241966
5	(<u>ловит</u> , <u>синицу</u>).	0.001890359168241966
6	(<u>синицу</u> , <u>которая</u>).	0.001890359168241966
7	(<u>которая</u> , <u>часто</u>).	0.009242144177449169
8	(<u>часто</u> , <u>ворует</u>).	0.009380863039399626
9	(<u>ворует</u> , <u>пшеницу</u>).	0.009380863039399626
10	(<u>пшеницу</u> , <u>которая</u>).	0.009380863039399626
11	(<u>которая</u> , <u>в</u>).	0.0166358595194085
12	(<u>в</u> , <u>темном</u>).	0.01651376146788991
13	(<u>темном</u> , <u>чулане</u>).	0.01675977653631285
14	(<u>чулане</u> , <u>хранится</u>).	0.01675977653631285
15	(<u>хранится</u> , <u>в</u>).	0.01675977653631285
16	(<u>в</u> , <u>доме</u>).	0.01651376146788991
17	(<u>доме</u> , <u>который</u>).	0.01675977653631285
18	(<u>который</u> , <u>построил</u>).	0.024029574861367836
19	(<u>построил</u> , <u>джек</u>).	0.024029574861367836

$$\lambda = 0.5$$

Закон Линдстоуна с параметром 0.5

#	Aa Words	# Probabilities
0	(<u>вот</u> , <u>кот</u>).	0.001890359168241966
1	(<u>кот</u> , <u>который</u>).	0.001890359168241966
2	(<u>который</u> , <u>пугает</u>).	0.001890359168241966
3	(<u>пугает</u> , <u>и</u>).	0.001890359168241966
4	(<u>и</u> , <u>ловит</u>).	0.001890359168241966
5	(<u>ловит</u> , <u>синицу</u>).	0.001890359168241966
6	(<u>синицу</u> , <u>которая</u>).	0.001890359168241966
7	(<u>которая</u> , <u>часто</u>).	0.005607476635514018
8	(<u>часто</u> , <u>ворует</u>).	0.005649717514124294

#	Aa Words	# Probabilities
9	(<u>ворует</u> , <u>пшеницу</u>).	0.005649717514124294
10	(<u>пшеницу</u> , <u>которая</u>).	0.005649717514124294
11	(<u>которая</u> , <u>в</u>).	0.009345794392523364
12	(<u>в</u> , <u>темном</u>).	0.00931098696461825
13	(<u>темном</u> , <u>чулане</u>).	0.009380863039399626
14	(<u>чулане</u> , <u>хранится</u>).	0.009380863039399626
15	(<u>хранится</u> , <u>в</u>).	0.009380863039399626
16	(<u>в</u> , <u>доме</u>).	0.00931098696461825
17	(<u>доме</u> , <u>который</u>).	0.009380863039399626
18	(<u>который</u> , <u>построил</u>).	0.013084112149532711
19	(<u>построил</u> , <u>джек</u>).	0.013084112149532711

$$\lambda = 0.75$$

Закон Линдстоуна с параметром 0.75

#	Aa Words	# Probabilities
0	(<u>вот</u> , <u>кот</u>).	0.001890359168241966
1	(<u>кот</u> , <u>который</u>).	0.001890359168241966
2	(<u>который</u> , <u>пугает</u>).	0.001890359168241966
3	(<u>пугает</u> , <u>и</u>).	0.001890359168241966
4	(<u>и</u> , <u>ловит</u>).	0.001890359168241966
5	(<u>ловит</u> , <u>синицу</u>).	0.001890359168241966
6	(<u>синицу</u> , <u>которая</u>).	0.001890359168241966
7	(<u>которая</u> , <u>часто</u>).	0.004377736085053158
8	(<u>часто</u> , <u>ворует</u>).	0.0043997485857950975
9	(<u>ворует</u> , <u>пшеницу</u>).	0.0043997485857950975
10	(<u>пшеницу</u> , <u>которая</u>).	0.0043997485857950975
11	(<u>которая</u> , <u>в</u>).	0.0068792995622263915
12	(<u>в</u> , <u>темном</u>).	0.006862133499688085
13	(<u>темном</u> , <u>чулане</u>).	0.006896551724137931
14	(<u>чулане</u> , <u>хранится</u>).	0.006896551724137931
15	(<u>хранится</u> , <u>в</u>).	0.006896551724137931

#	Aa Words	# Probabilities
16	('в', 'доме').	0.006862133499688085
17	('доме', 'который').	0.006896551724137931
18	('который', 'построил').	0.009380863039399626
19	('построил', 'джек').	0.009380863039399626

Перплексия

+1 : 273.7287283373542

$\lambda = 0.25$: 139.09749617838835

$\lambda = 0.5$: 200.17679598819893

$\lambda = 0.75$: 242.19156910489474

Заметим, что перплексии двух различных вариаций закона Лидстоуна совпадают при одинаковых λ .