

# Языковые модели

Павел Браславский

*использованы слайды Dan Jurafsky*

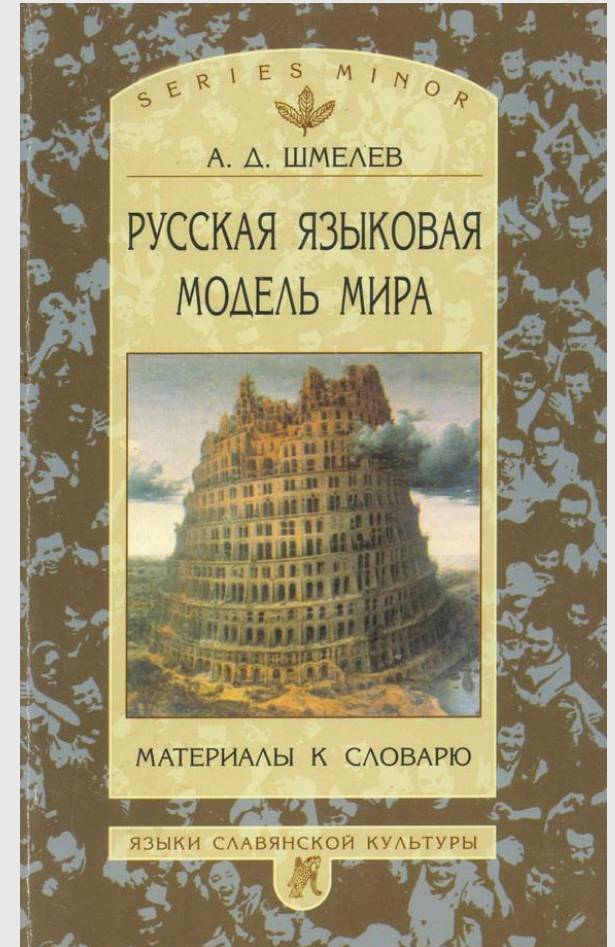
# Терминологическое замечание

(statistical) language model (LM)



языковая модель

модель языка



# **ВВЕДЕНИЕ**

# Языковые модели

Задача: определить вероятность последовательности слов

- распознавание речи  
*скрипка лиса vs. скрип колеса*
- исправление опечаток  
*курсовая работа vs. курсовая робота*
- генерация текста (в том числе спама 😊)
- машинный перевод  
*крепкий чай vs. сильный чай*
- ...



почему люди|занимаются

почему люди занимаются бизнесом

почему люди не летают как птицы монолог катерины

почему люди путешествуют

почему люди завидуют

почему люди икают

почему люди осваивали новые земли окружающий мир 4 класс

почему люди воюют

почему люди заботятся о скворцах

почему люди становятся вегетарианцами

почему люди приобретают привычку к курению

городская и культурная информация

113

# ВРЕМЯ И СТЕКЛО

ДВОРЕЦ  
МОЛОДЕЖИ  
5 ОКТЯБРЯ

ЗВЕЗДЫ  
РУССКОГО  
РАДИО

16+

РУССКОЕ  
РАДИО

"время и стекло"

All

Videos

Images

News

About 2,100,000 results (0.30 seconds)

"время истекло"

All

Videos

Images

News

About 291,000 results (0.29 seconds)

# Предсказание следующего слова

- яблоко от яблони...
- круглый *год/сирота/...*
- выпил *чаю/воды/стакан/...*

309 и  
192 стакан  
138 еще  
120 рюмку  
86 залпом  
82 с  
69 водки  
66 ее  
66 за  
63 две  
55 его  
48 и  
45 а  
43 чаю  
42 бы  
42 в  
39 два  
39 свой

# Модели на основе n-грамм

- Как вычислить вероятность новой последовательности?

$$P(S) = P(w_1, w_2, \dots, w_n)$$

- Произведение вероятностей:

$$P(S) = P(w_1)P(w_2 | w_1) P(w_3 | w_1, w_2) \dots P(w_n | w_1, w_2, \dots, w_{n-1})$$

- Предположим, что это Марковский процесс (т.е. ограничим влияние «истории»)
  - униграммная  $P(w_n)$
  - биграммная  $P(w_n | w_{n-1})$
  - триграммная  $P(w_n | w_{n-2}, w_{n-1})$

# Модели на основе n-грамм –2

- моделируют *локальные* зависимости
- n-граммы более высоких порядков: мало данных, нет «обобщения»



# ОЦЕНКА ВЕРОЯТНОСТЕЙ

# Оценка вероятностей биграмм

- Метод максимального правдоподобия (Maximum Likelihood Estimate)

$$P(w_i | w_{i-1}) = \frac{c(w_{i-1}, w_i)}{c(w_{i-1})}$$

# Пример

<s>иду шагаю по москве</s>

<s>иду на вы</s>

<s>шагаю босиком по улице</s>

$$P(\text{иду} \mid \langle s \rangle) = 0.67$$

$$P(\text{шагаю} \mid \langle s \rangle) = 0.33$$

$$P(\text{шагаю} \mid \text{иду}) = 0.5$$

$$P(\text{москве} \mid \text{по}) = 0.5$$

$$P(\text{по} \mid \text{босиком}) = 1.0$$

...

# Berkeley Restaurant Project corpus

9,222 предложений

can you tell me about any good cantonese restaurants close by

mid priced thai food is what i'm looking for

tell me about chez panisse

can you give me a listing of the kinds of food that are available

i'm looking for a good place to eat breakfast

when is caffe venezia open during the day

# Частоты биграмм

	i	want	to	eat	chinese	food	lunch	spend
i	5	827	0	9	0	0	0	2
want	2	0	608	1	6	6	5	1
to	2	0	4	686	2	0	6	211
eat	0	0	2	0	16	2	42	0
chinese	1	0	0	0	0	82	1	0
food	15	0	15	0	1	4	0	0
lunch	2	0	0	0	0	1	0	0
spend	1	0	1	0	0	0	0	0

# Вероятности биграмм

i	want	to	eat	chinese	food	lunch	spend
2533	927	2417	746	158	1093	341	278

	i	want	to	eat	chinese	food	lunch	spend
i	0.002	0.33	0	0.0036	0	0	0	0.00079
want	0.0022	0	0.66	0.0011	0.0065	0.0065	0.0054	0.0011
to	0.00083	0	0.0017	0.28	0.00083	0	0.0025	0.087
eat	0	0	0.0027	0	0.021	0.0027	0.056	0
chinese	0.0063	0	0	0	0	0.52	0.0063	0
food	0.014	0	0.014	0	0.00092	0.0037	0	0
lunch	0.0059	0	0	0	0	0.0029	0	0
spend	0.0036	0	0.0036	0	0	0	0	0

# Вероятность предложения

$$\begin{aligned} P(\langle s \rangle \text{ I want english food } \langle /s \rangle) &= \\ &P(\text{I} | \langle s \rangle) \\ &\times P(\text{want} | \text{I}) \\ &\times P(\text{english} | \text{want}) \\ &\times P(\text{food} | \text{english}) \\ &\times P(\langle /s \rangle | \text{food}) \\ &= .000031 \end{aligned}$$

# О чем говорят вероятности?

- $P(\text{english} | \text{want}) = .0011$
- $P(\text{chinese} | \text{want}) = .0065$
- $P(\text{to} | \text{want}) = .66$
- $P(\text{eat} | \text{to}) = .28$
- $P(\text{food} | \text{to}) = 0$
- $P(\text{want} | \text{spend}) = 0$
- $P(i | \langle s \rangle) = .25$



# Практические соображения

- лучше оперировать логарифмами вероятностей
  - избежать переполнения
  - сложение быстрее, чем умножение

$$\begin{aligned}\log(p_1 \times p_2 \times p_3 \times p_4) &= \\ &= \log p_1 + \log p_2 + \log p_3 + \log p_4\end{aligned}$$

# Оценка

- «внешняя» оценка: насколько хорошо модель помогает решить задачу
- «внутренняя» оценка модели : коэффициент неопределенности (перплексия, perplexity):
  - хорошая модель приписывает более высокую вероятность последовательности (предложению), которая действительно встречается в тексте

# Перплексия

- Обратная вероятность тестового набора, нормализованная на длину

$$\begin{aligned} PP(W) &= P(w_1 w_2 \dots w_N)^{-\frac{1}{N}} \\ &= \sqrt[N]{\frac{1}{P(w_1 w_2 \dots w_N)}} \end{aligned}$$

- связана с энтропией

$$PP(W) = \sqrt[N]{\prod_{i=1}^N \frac{1}{P(w_i | w_1 \dots w_{i-1})}}$$

$$PP(W) = \sqrt[N]{\prod_{i=1}^N \frac{1}{P(w_i | w_{i-1})}}$$

# Ниже перплексия = лучше модель

- Обучение на 38 млн. слов, тестирование на 1.5 млн., WSJ

	Unigram	Bigram	Trigram
Perplexity	962	170	109

**ДААННЫЕ**



[главная](#)

[архив новостей](#)

[поиск в корпусе](#)

[что такое корпус?](#)

[состав и структура](#)

[статистика](#)

[графики](#)

[частоты](#)

[морфология](#)

[обороты](#)

[синтаксис](#)

[семантика](#)

## Частоты словоформ и словосочетаний

Вы можете скачать архивы с текстовыми файлами, содержащими частоты словоформ и словосочетаний в основном корпусе.

При подсчёте учитывался регистр букв, а также знаки препинания.  
Общий объём корпуса – 192689044 словоформы.

Словоформы	<a href="#">zip-архив</a> (5,5 Мб, обрезаны по частоте 3)	<a href="#">топ-100</a>
2-граммы	<a href="#">zip-архив</a> (39 Мб, обрезаны по частоте 3)	<a href="#">топ-100</a>
3-граммы	<a href="#">zip-архив</a> (31 Мб, обрезаны по частоте 3)	<a href="#">топ-100</a>
4-граммы	<a href="#">zip-архив</a> (44 Мб, обрезаны по частоте 2)	<a href="#">топ-100</a>
5-граммы	<a href="#">zip-архив</a> (28 Мб, обрезаны по частоте 2)	<a href="#">топ-100</a>
6-граммы		<a href="#">топ</a>

# Google Books n-gram viewer



## Google Books Ngram Viewer

Graph these comma-separated phrases:   case-insensitive

between  and  from the corpus  with smoothing of

[Search lots of books](#)

[G+ Share](#)

[Tweet](#)

[Embed Chart](#)



(click on line/label for focus)

# Google Books n-grams

Russian

Version 20120701

[total counts](#)

**1-grams** [0](#) [1](#) [2](#) [3](#) [4](#) [5](#) [6](#) [7](#) [8](#) [9](#) [a](#) [b](#) [c](#) [d](#) [e](#) [f](#) [g](#) [h](#) [i](#) [j](#) [k](#) [l](#) [m](#) [n](#) [o](#) [other](#) [p](#) [pos](#) [punctuation](#) [q](#) [r](#) [s](#) [t](#) [u](#) [v](#) [w](#) [x](#) [y](#) [z](#)

**2-grams** [0](#) [1](#) [2](#) [3](#) [4](#) [5](#) [6](#) [7](#) [8](#) [9](#) [ADJ](#) [ADP](#) [ADV](#) [CONJ](#) [NOUN](#) [NUM](#) [PRT](#) [VERB](#) [a](#) [aa](#) [ab](#) [ac](#) [ad](#) [ae](#) [af](#) [ag](#) [ah](#) [ai](#) [aj](#) [ak](#) [al](#) [am](#) [an](#) [ao](#) [ap](#) [aq](#) [ar](#) [as](#) [at](#) [au](#) [av](#) [aw](#) [ax](#) [ay](#) [az](#) [b](#) [ba](#) [bb](#) [bc](#) [bd](#) [be](#) [bf](#) [bg](#) [bh](#) [bi](#) [bj](#) [bk](#) [bl](#) [bm](#) [bn](#) [bo](#) [bp](#) [br](#) [bs](#) [bt](#) [bu](#) [bv](#) [bx](#) [by](#) [bz](#) [c](#) [ca](#) [cb](#) [cc](#) [cd](#) [ce](#) [cf](#) [cg](#) [ch](#) [ci](#) [cj](#) [ck](#) [cl](#) [cm](#) [cn](#) [co](#) [cp](#) [cq](#) [cr](#) [cs](#) [ct](#) [cu](#) [cv](#) [cx](#) [cy](#) [cz](#) [d](#) [da](#) [db](#) [dc](#) [dd](#) [de](#) [df](#) [dg](#) [dh](#) [di](#) [dj](#) [dk](#) [dl](#) [dm](#) [dn](#) [do](#) [dp](#) [dr](#) [ds](#) [dt](#) [du](#) [dv](#) [dw](#) [dx](#) [dy](#) [dz](#) [e](#) [ea](#) [eb](#) [ec](#) [ed](#) [ee](#) [ef](#) [eg](#) [eh](#) [ei](#) [ej](#) [ek](#) [el](#) [em](#) [en](#) [eo](#) [ep](#) [eq](#) [er](#) [es](#) [et](#) [eu](#) [ev](#) [ew](#) [ex](#) [ey](#) [ez](#) [f](#) [fa](#) [fb](#) [fc](#) [fd](#) [fe](#) [ff](#) [fg](#) [fh](#) [fi](#) [fj](#) [fk](#) [fl](#) [fm](#) [fn](#) [fo](#) [fp](#) [fr](#) [fs](#) [ft](#) [fu](#) [fv](#) [fx](#) [fy](#) [fz](#) [g](#) [ga](#) [gb](#) [gc](#) [gd](#) [ge](#) [gf](#) [gg](#) [gh](#) [gi](#) [gj](#) [gk](#) [gl](#) [gm](#) [gn](#) [go](#) [gp](#) [gr](#) [gs](#) [gt](#) [gu](#) [gv](#) [gx](#) [gy](#) [gz](#) [h](#) [ha](#) [hb](#) [hc](#) [hd](#) [he](#) [hf](#) [hg](#) [hh](#) [hi](#) [hj](#) [hk](#) [hl](#) [hm](#) [hn](#) [ho](#) [hp](#) [hr](#) [hs](#) [ht](#) [hu](#) [hv](#) [hx](#) [hy](#) [hz](#) [i](#) [ia](#) [ib](#) [ic](#) [id](#) [ie](#) [if](#) [ig](#) [ih](#) [ii](#) [ij](#) [ik](#) [il](#) [im](#) [in](#) [io](#) [ip](#) [iq](#) [ir](#) [is](#) [it](#) [iu](#) [iw](#) [ix](#) [iy](#) [iz](#) [j](#) [ja](#) [jb](#) [jc](#) [jd](#) [je](#) [if](#) [ig](#) [ih](#) [ii](#) [ij](#) [ik](#) [il](#) [im](#) [in](#) [io](#) [ip](#) [ir](#) [is](#) [it](#) [iu](#) [iw](#) [ix](#) [iy](#) [iz](#) [k](#) [ka](#) [kb](#) [kc](#) [kd](#) [ke](#) [kf](#) [kg](#) [kh](#) [ki](#) [kj](#) [kk](#) [kl](#) [km](#) [kn](#) [ko](#) [kp](#) [kr](#) [ks](#) [kt](#) [ku](#) [kv](#) [kx](#) [ky](#) [kz](#) [l](#) [la](#) [lb](#) [lc](#) [ld](#) [le](#) [lf](#) [lg](#) [lh](#) [li](#) [lj](#) [lk](#) [ll](#) [lm](#) [ln](#) [lo](#) [lp](#) [lr](#) [ls](#) [lt](#) [lu](#) [lv](#) [lw](#) [lx](#) [ly](#) [lz](#) [m](#) [ma](#) [mb](#) [mc](#) [md](#) [me](#) [mf](#) [mg](#) [mh](#) [mi](#) [mj](#) [mk](#) [ml](#) [mm](#) [mn](#) [mo](#) [mp](#) [mr](#) [ms](#) [mt](#) [mu](#) [mv](#) [mw](#) [mx](#) [my](#) [mz](#) [n](#) [na](#) [nb](#) [nc](#) [nd](#) [ne](#) [nf](#) [ng](#) [nh](#) [ni](#) [nj](#) [nk](#) [nl](#) [nm](#) [nn](#) [no](#) [np](#) [nr](#) [ns](#) [nt](#) [nu](#) [nv](#) [nw](#) [nx](#) [ny](#) [nz](#) [o](#) [oa](#) [ob](#) [oc](#) [od](#) [oe](#) [of](#) [og](#) [oh](#) [oi](#) [oj](#) [ok](#) [ol](#) [om](#) [on](#) [oo](#) [op](#) [or](#) [os](#) [ot](#) [other](#) [ou](#) [ov](#) [ow](#) [ox](#) [oy](#) [oz](#) [p](#) [pa](#) [pb](#) [pc](#) [pd](#) [pe](#) [pf](#) [pg](#) [ph](#) [pi](#) [pj](#) [pk](#) [pl](#) [pm](#) [pn](#) [po](#) [pp](#) [pr](#) [ps](#) [pt](#) [pu](#) [punctuation](#) [pv](#) [px](#) [py](#) [pz](#) [q](#) [qa](#) [qi](#) [qn](#) [qu](#) [qx](#) [r](#) [ra](#) [rb](#) [rc](#) [rd](#) [re](#) [rf](#) [rg](#) [rh](#) [ri](#) [rj](#) [rk](#) [rl](#) [rm](#) [rn](#) [ro](#) [rp](#) [rr](#) [rs](#) [rt](#) [ru](#) [rv](#) [rw](#) [rx](#) [ry](#) [rz](#) [s](#) [sa](#) [sb](#) [sc](#) [sd](#) [se](#) [sf](#) [sg](#) [sh](#) [si](#) [sj](#) [sk](#) [sl](#) [sm](#) [sn](#) [so](#) [sp](#) [sq](#) [sr](#) [ss](#) [st](#) [su](#) [sv](#) [sw](#) [sx](#) [sy](#) [sz](#) [t](#) [ta](#) [tb](#) [tc](#) [td](#) [te](#) [tf](#) [tg](#) [th](#) [ti](#) [tj](#) [tk](#) [tl](#) [tm](#) [tn](#) [to](#) [tp](#) [tr](#) [ts](#) [tt](#) [tu](#) [tv](#) [tw](#) [tx](#) [ty](#) [tz](#) [u](#) [ua](#) [ub](#) [uc](#) [ud](#) [ue](#) [uf](#) [ug](#) [uh](#) [ui](#) [uj](#) [uk](#) [ul](#) [um](#) [un](#) [uo](#) [up](#) [ur](#) [us](#) [ut](#) [uu](#) [uv](#) [ux](#) [uy](#) [uz](#) [v](#) [va](#) [vb](#) [vc](#) [vd](#) [ve](#) [vf](#) [vg](#) [vh](#) [vi](#) [vj](#) [vk](#) [vl](#) [vm](#) [vn](#) [vo](#) [vp](#) [vr](#) [vs](#) [vt](#) [vu](#) [vv](#) [vx](#) [vy](#) [vz](#) [w](#) [wa](#) [wb](#) [wc](#) [we](#) [wh](#) [wi](#) [wj](#) [wl](#) [wm](#) [wn](#) [wo](#) [wr](#) [ws](#) [wt](#) [wu](#) [ww](#) [wx](#) [wy](#) [x](#) [xa](#) [xb](#) [xc](#) [xd](#) [xe](#) [xf](#) [xg](#) [xh](#) [xi](#) [xj](#) [xk](#) [xl](#) [xm](#) [xn](#) [xo](#) [xp](#) [xr](#) [xs](#) [xt](#) [xu](#) [xv](#) [xx](#) [xy](#) [y](#) [ya](#) [yb](#) [yc](#) [yd](#) [ye](#) [yf](#) [yg](#) [yh](#) [yi](#) [yj](#) [yk](#) [yl](#) [ym](#) [yn](#) [yo](#) [yp](#) [yr](#) [ys](#) [yt](#) [yu](#) [yv](#) [yx](#) [yz](#) [z](#) [za](#) [zb](#) [zc](#) [zd](#) [ze](#) [zf](#) [zg](#) [zh](#) [zi](#) [zj](#) [zk](#) [zl](#) [zm](#) [zn](#) [zo](#) [zp](#) [zr](#) [zs](#) [zt](#) [zu](#) [zv](#) [zw](#) [zx](#) [zy](#) [zz](#)

**3-grams** [0](#) [1](#) [2](#) [3](#) [4](#) [5](#) [6](#) [7](#) [8](#) [9](#) [ADJ](#) [ADP](#) [ADV](#) [CONJ](#) [NOUN](#) [NUM](#) [PRT](#) [VERB](#) [a](#) [aa](#) [ab](#) [ac](#) [ad](#) [ae](#) [af](#) [ag](#) [ah](#) [ai](#) [aj](#) [ak](#) [al](#) [am](#) [an](#) [ao](#) [ap](#) [aq](#) [ar](#) [as](#) [at](#) [au](#) [av](#) [aw](#) [ax](#) [ay](#) [az](#) [b](#) [ba](#) [bb](#) [bc](#) [bd](#) [be](#) [bf](#) [bg](#) [bh](#) [bi](#) [bj](#) [bk](#) [bl](#) [bm](#) [bn](#) [bo](#) [bp](#) [br](#) [bs](#) [bt](#) [bu](#) [bv](#) [bx](#) [by](#) [bz](#) [c](#) [ca](#) [cb](#) [cc](#) [cd](#) [ce](#) [cf](#) [cg](#) [ch](#) [ci](#) [cj](#) [ck](#) [cl](#) [cm](#) [cn](#) [co](#) [cp](#) [cq](#) [cr](#) [cs](#) [ct](#) [cu](#) [cv](#) [cx](#) [cy](#) [cz](#) [d](#) [da](#) [db](#) [dc](#) [dd](#) [de](#) [df](#) [dg](#) [dh](#) [di](#) [dj](#) [dk](#) [dl](#) [dm](#) [dn](#) [do](#) [dp](#) [dr](#) [ds](#) [dt](#) [du](#) [dv](#) [dw](#) [dx](#) [dy](#) [dz](#) [e](#) [ea](#) [eb](#) [ec](#) [ed](#) [ee](#) [ef](#) [eg](#) [eh](#) [ei](#) [ej](#) [ek](#) [el](#) [em](#) [en](#) [eo](#) [ep](#) [eq](#) [er](#) [es](#) [et](#) [eu](#) [ev](#) [ex](#) [ey](#) [ez](#) [f](#) [fa](#) [fb](#) [fc](#) [fd](#) [fe](#) [ff](#) [fg](#) [fh](#) [fi](#) [fj](#) [fk](#) [fl](#) [fm](#) [fn](#) [fo](#) [fp](#) [fr](#) [fs](#) [ft](#) [fu](#) [fv](#) [fx](#) [fy](#) [fz](#) [g](#) [ga](#) [gb](#) [gc](#) [gd](#) [ge](#) [gf](#) [gg](#) [gh](#) [gi](#) [gj](#) [gk](#) [gl](#) [gm](#) [gn](#) [go](#) [gp](#) [gr](#) [gs](#) [gt](#) [gu](#) [gv](#) [gx](#) [gy](#) [gz](#) [h](#) [ha](#) [hb](#) [hc](#) [hd](#) [he](#) [hf](#) [hg](#) [hh](#) [hi](#) [hj](#) [hk](#) [hl](#) [hm](#) [hn](#) [ho](#) [hp](#) [hr](#) [hs](#) [ht](#) [hu](#) [hv](#) [hx](#) [hy](#) [hz](#) [i](#) [ia](#) [ib](#) [ic](#) [id](#) [ie](#) [if](#) [ig](#) [ih](#) [ii](#) [ij](#) [ik](#) [il](#) [im](#) [in](#) [io](#) [ip](#) [iq](#) [ir](#) [is](#) [it](#) [iu](#) [iw](#) [ix](#) [iy](#) [iz](#) [j](#) [ja](#) [jb](#) [jc](#) [jd](#) [je](#) [if](#) [ig](#) [ih](#) [ii](#) [ij](#) [ik](#) [il](#) [im](#) [in](#) [io](#) [ip](#) [ir](#) [is](#) [it](#) [iu](#) [iw](#) [ix](#) [iy](#) [iz](#) [k](#) [ka](#) [kb](#) [kc](#) [kd](#) [ke](#) [kf](#) [kg](#) [kh](#) [ki](#) [kj](#) [kk](#) [kl](#) [km](#) [kn](#) [ko](#) [kp](#) [kr](#) [ks](#) [kt](#) [ku](#) [kv](#) [kx](#) [ky](#) [kz](#) [l](#) [la](#) [lb](#) [lc](#) [ld](#) [le](#) [lf](#) [lg](#) [lh](#) [li](#) [lj](#) [lk](#) [ll](#) [lm](#) [ln](#) [lo](#) [lp](#) [lr](#) [ls](#) [lt](#) [lu](#) [lv](#) [lx](#) [ly](#) [lz](#) [m](#) [ma](#) [mb](#) [mc](#) [md](#) [me](#) [mf](#) [mg](#) [mh](#) [mi](#) [mj](#) [mk](#) [ml](#) [mm](#) [mn](#) [mo](#) [mp](#) [mr](#) [ms](#) [mt](#) [mu](#) [mv](#) [mw](#) [mx](#) [my](#) [mz](#) [n](#) [na](#) [nb](#) [nc](#) [nd](#) [ne](#) [nf](#) [ng](#) [nh](#) [ni](#) [nj](#) [nk](#) [nl](#) [nm](#) [nn](#) [no](#) [np](#) [nr](#) [ns](#) [nt](#) [nu](#) [nv](#) [nw](#) [nx](#) [ny](#) [nz](#) [o](#) [oa](#) [ob](#) [oc](#) [od](#) [oe](#) [of](#) [og](#) [oh](#) [oi](#) [oj](#) [ok](#) [ol](#) [om](#) [on](#) [oo](#) [op](#) [or](#) [os](#) [ot](#) [other](#) [ou](#) [ov](#) [ow](#) [ox](#) [oy](#) [oz](#) [p](#) [pa](#) [pb](#) [pc](#) [pd](#) [pe](#) [pf](#) [pg](#) [ph](#) [pi](#) [pj](#) [pk](#) [pl](#) [pm](#) [pn](#) [po](#) [pp](#) [pr](#) [ps](#) [pt](#) [pu](#) [punctuation](#) [pv](#) [px](#) [py](#) [pz](#) [q](#) [qa](#) [qn](#) [qu](#) [qx](#) [r](#) [ra](#) [rb](#) [rc](#) [rd](#) [re](#) [rf](#) [rg](#) [rh](#) [ri](#) [rj](#) [rk](#) [rl](#) [rm](#) [rn](#) [ro](#) [rp](#) [rr](#) [rs](#) [rt](#) [ru](#) [rv](#) [rw](#) [rx](#) [ry](#) [rz](#) [s](#) [sa](#) [sb](#) [sc](#) [sd](#) [se](#) [sf](#) [sg](#) [sh](#) [si](#) [sj](#) [sk](#) [sl](#) [sm](#) [sn](#) [so](#) [sp](#) [sq](#) [sr](#) [ss](#) [st](#) [su](#) [sv](#) [sw](#) [sx](#) [sy](#) [sz](#) [t](#) [ta](#) [tb](#) [tc](#) [td](#) [te](#) [tf](#) [tg](#) [th](#) [ti](#) [tj](#) [tk](#) [tl](#) [tm](#) [tn](#) [to](#) [tp](#) [tr](#) [ts](#) [tt](#) [tu](#) [tv](#) [tw](#) [tx](#) [ty](#) [tz](#) [u](#) [ua](#) [ub](#) [uc](#) [ud](#) [ue](#) [uf](#) [ug](#) [uh](#) [ui](#) [uj](#) [uk](#) [ul](#) [um](#) [un](#) [uo](#) [up](#) [ur](#) [us](#) [ut](#) [uu](#) [uv](#) [ux](#) [uy](#) [uz](#) [v](#) [va](#) [vb](#) [vc](#) [vd](#) [ve](#) [vf](#) [vg](#) [vh](#) [vi](#) [vj](#) [vk](#) [vl](#) [vm](#) [vn](#) [vo](#) [vp](#) [vr](#) [vs](#) [vt](#) [vu](#) [vv](#) [vx](#) [vy](#) [vz](#) [w](#) [wa](#) [wb](#) [wc](#) [we](#) [wh](#) [wi](#) [wj](#) [wl](#) [wm](#) [wn](#) [wo](#) [wr](#) [ws](#) [wt](#) [wu](#) [ww](#) [wx](#) [x](#) [xa](#) [xc](#) [xd](#) [xe](#) [xf](#) [xg](#) [xh](#) [xi](#) [xj](#) [xk](#) [xl](#) [xm](#) [xn](#) [xo](#) [xp](#) [xr](#) [xs](#) [xt](#) [xu](#) [xv](#) [xx](#) [xy](#) [y](#) [ya](#) [yb](#) [yc](#) [yd](#) [ye](#) [yf](#) [yg](#) [yh](#) [yi](#) [yj](#) [yk](#) [yl](#) [ym](#) [yn](#) [yo](#) [yp](#) [yr](#) [ys](#) [yt](#) [yu](#) [yv](#) [yx](#) [yz](#) [z](#) [za](#) [zb](#) [zc](#) [zd](#) [ze](#) [zf](#) [zg](#) [zh](#) [zi](#) [zj](#) [zk](#) [zl](#) [zm](#) [zn](#) [zo](#) [zp](#) [zr](#) [zs](#) [zt](#) [zu](#) [zv](#) [zw](#) [zx](#) [zy](#) [zz](#)

**4-grams** [0](#) [1](#) [2](#) [3](#) [4](#) [5](#) [6](#) [7](#) [8](#) [9](#) [ADJ](#) [ADP](#) [ADV](#) [CONJ](#) [NOUN](#) [NUM](#) [PRT](#) [VERB](#) [a](#) [aa](#) [ab](#) [ac](#) [ad](#) [ae](#) [af](#) [ag](#) [ah](#) [ai](#) [aj](#) [ak](#) [al](#) [am](#) [an](#) [ao](#) [ap](#) [aq](#) [ar](#) [as](#) [at](#) [au](#) [av](#) [aw](#) [ax](#) [ay](#) [az](#) [b](#) [ba](#) [bb](#) [bc](#) [bd](#) [be](#) [bf](#) [bg](#) [bh](#) [bi](#) [bj](#) [bk](#) [bl](#) [bm](#) [bn](#) [bo](#) [bp](#) [br](#) [bs](#) [bt](#) [bu](#) [bv](#) [bx](#) [by](#) [bz](#) [c](#) [ca](#) [cb](#) [cc](#) [cd](#) [ce](#) [cf](#) [cg](#) [ch](#) [ci](#) [cj](#) [ck](#) [cl](#) [cm](#) [cn](#) [co](#) [cp](#) [cq](#) [cr](#) [cs](#) [ct](#) [cu](#) [cv](#) [cx](#) [cy](#) [cz](#) [d](#) [da](#) [db](#) [dc](#) [dd](#) [de](#) [df](#) [dg](#) [dh](#) [di](#) [dj](#) [dk](#) [dl](#) [dm](#) [dn](#) [do](#) [dp](#) [dr](#) [ds](#) [dt](#) [du](#) [dv](#) [dw](#) [dx](#) [dy](#) [dz](#) [e](#) [ea](#) [eb](#) [ec](#) [ed](#) [ee](#) [ef](#) [eg](#) [eh](#) [ei](#) [ej](#) [ek](#) [el](#) [em](#) [en](#) [eo](#) [ep](#) [eq](#) [er](#) [es](#) [et](#) [eu](#) [ev](#) [ex](#) [ey](#) [ez](#) [f](#) [fa](#)





# Cognitive Services

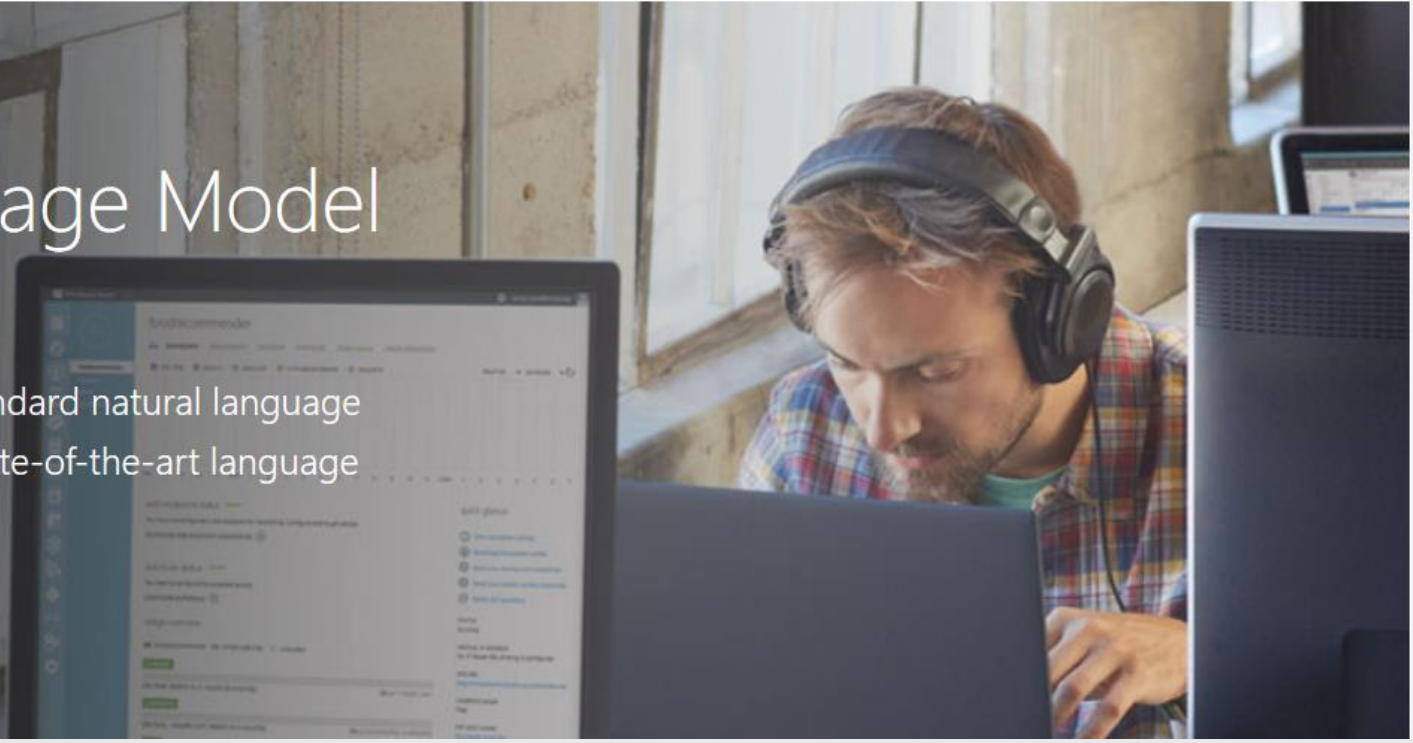
My account

- Home
- APIs ▾
- Applications
- Developers ▾
- Pricing

## Web Language Model API

Automate a variety of standard natural language processing tasks using state-of-the-art language modeling APIs.

Get started for free



## Joint probabilities

Calculate how often a particular sequence of words appear together.



## Conditional probabilities

Given a sequence of words, calculate how often a particular word tends to follow.



## Next word completions

Given a sequence of words, get the list of words most likely to follow.



# **СГЛАЖИВАНИЕ ЛАПЛАСА**

# Новые последовательности

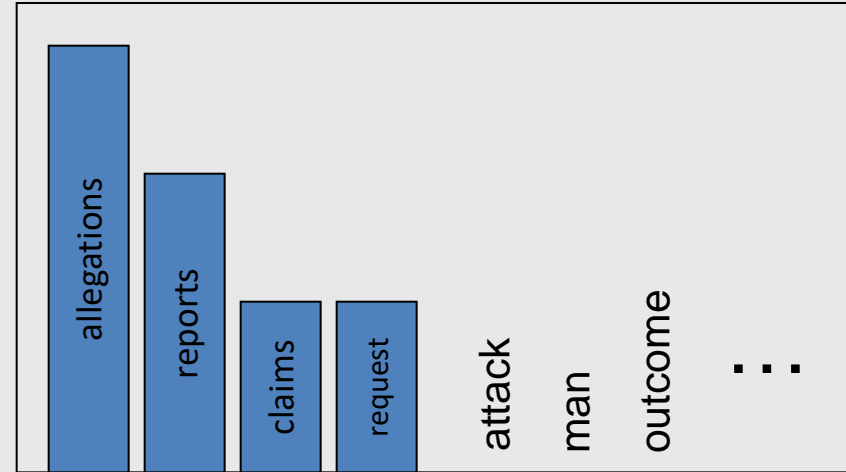
- языковая модель должна *обобщать* (а не повторять) данные, на которых она обучалась
- всегда будут новые последовательности слов, которые не встречались в корпусе для обучения
- «новые» n-граммы (в тестовых данных, но не в данных для обучения) «ломают» модель

# Сглаживание

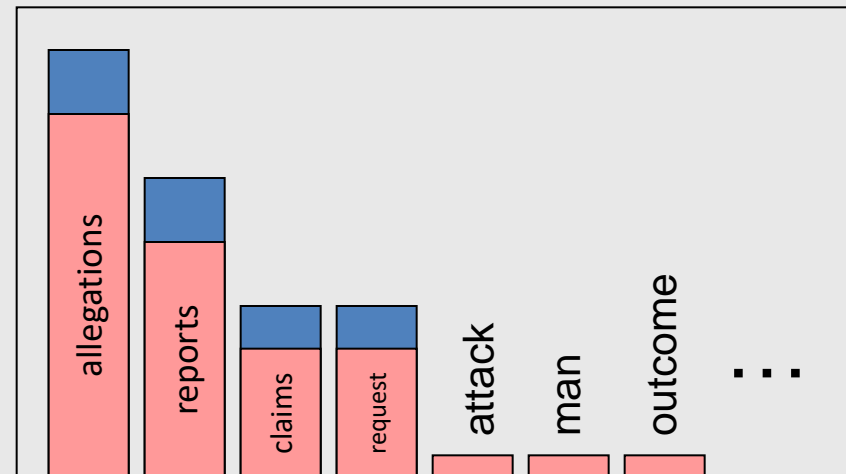
- «зарезервировать» часть вероятностной массы для событий, которые еще не встречались

# Пример

$P(w \mid \text{denied the})$   
3 allegations  
2 reports  
1 claims  
1 request  
7 total



$P(w \mid \text{denied the})$   
2.5 allegations  
1.5 reports  
0.5 claims  
0.5 request  
**2 other**  
7 total



# Сглаживание Лапласа

- или сглаживание «+1»
- Предположим, что мы встречали каждую пару слов на один раз больше

- оценка максимального правдоподобия: 
$$P_{MLE}(w_i | w_{i-1}) = \frac{c(w_{i-1}, w_i)}{c(w_{i-1})}$$

- оценка «+1»: 
$$P_{Add-1}(w_i | w_{i-1}) = \frac{c(w_{i-1}, w_i) + 1}{c(w_{i-1}) + V}$$

	i	want	to	eat	chinese	food	lunch	spend
i	5	827	0	9	0	0	0	2
want	2	0	608	1	6	6	5	1
to	2	0	4	686	2	0	6	211
eat	0	0	2	0	16	2	42	0
chinese	1	0	0	0	0	82	1	0
food	15	0	15	0	1	4	0	0
lunch	2	0	0	0	0	1	0	0
spend	1	0	1	0	0	0	0	0

	i	want	to	eat	chinese	food	lunch	spend
i	6	828	1	10	1	1	1	3
want	3	1	609	2	7	7	6	2
to	3	1	5	687	3	1	7	212
eat	1	1	3	1	17	3	43	1
chinese	2	1	1	1	1	83	2	1
food	16	1	16	1	2	5	1	1
lunch	3	1	1	1	1	2	1	1
spend	2	1	2	1	1	1	1	1

[Dan Jurafsky]



# Сглаженные вероятности

$$P^*(w_n | w_{n-1}) = \frac{C(w_{n-1}w_n) + 1}{C(w_{n-1}) + V}$$

	i	want	to	eat	chinese	food	lunch	spend
i	0.0015	0.21	0.00025	0.0025	0.00025	0.00025	0.00025	0.00075
want	0.0013	0.00042	0.26	0.00084	0.0029	0.0029	0.0025	0.00084
to	0.00078	0.00026	0.0013	0.18	0.00078	0.00026	0.0018	0.055
eat	0.00046	0.00046	0.0014	0.00046	0.0078	0.0014	0.02	0.00046
chinese	0.0012	0.00062	0.00062	0.00062	0.00062	0.052	0.0012	0.00062
food	0.0063	0.00039	0.0063	0.00039	0.00079	0.002	0.00039	0.00039
lunch	0.0017	0.00056	0.00056	0.00056	0.00056	0.0011	0.00056	0.00056
spend	0.0012	0.00058	0.0012	0.00058	0.00058	0.00058	0.00058	0.00058

# Восстановленные частоты

$$c^*(w_{n-1}w_n) = \frac{[C(w_{n-1}w_n) + 1] \times C(w_{n-1})}{C(w_{n-1}) + V}$$

	i	want	to	eat	chinese	food	lunch	spend
i	3.8	527	0.64	6.4	0.64	0.64	0.64	1.9
want	1.2	0.39	238	0.78	2.7	2.7	2.3	0.78
to	1.9	0.63	3.1	430	1.9	0.63	4.4	133
eat	0.34	0.34	1	0.34	5.8	1	15	0.34
chinese	0.2	0.098	0.098	0.098	0.098	8.2	0.2	0.098
food	6.9	0.43	6.9	0.43	0.86	2.2	0.43	0.43
lunch	0.57	0.19	0.19	0.19	0.19	0.38	0.19	0.19
spend	0.32	0.16	0.32	0.16	0.16	0.16	0.16	0.16

# Сравнение

	i	want	to	eat	chinese	food	lunch	spend
i	5	827	0	9	0	0	0	2
want	2	0	608	1	6	6	5	1
to	2	0	4	686	2	0	6	211
eat	0	0	2	0	16	2	42	0
chinese	1	0	0	0	0	82	1	0
food	15	0	15	0	1	4	0	0
lunch	2	0	0	0	0	1	0	0
spend	1	0	1	0	0	0	0	0

	i	want	to	eat	chinese	food	lunch	spend
i	3.8	527	0.64	6.4	0.64	0.64	0.64	1.9
want	1.2	0.39	238	0.78	2.7	2.7	2.3	0.78
to	1.9	0.63	3.1	430	1.9	0.63	4.4	133
eat	0.34	0.34	1	0.34	5.8	1	15	0.34
chinese	0.2	0.098	0.098	0.098	0.098	8.2	0.2	0.098
food	6.9	0.43	6.9	0.43	0.86	2.2	0.43	0.43
lunch	0.57	0.19	0.19	0.19	0.19	0.38	0.19	0.19
spend	0.32	0.16	0.32	0.16	0.16	0.16	0.16	0.16

[Dan Jurafsky]

# Сглаживание «+1»

- плохо работает для языковых моделей
  - более мягкий вариант «+ $\alpha$ »:

$$P_{Add-1}(w_i | w_{i-1}) = \frac{c(w_{i-1}, w_i) + \alpha}{c(w_{i-1}) + \alpha V}$$

- используется для классификации текстов
  - и там, где меньше нулевых оценок

# ОТКАТ И ИНТЕРПОЛЯЦИЯ

# Откат и интерполяция

- Откат (backoff)
  - переход к  $n$  с более низким  $n$ , для которых достаточно данных ( $3 \rightarrow 2 \rightarrow 1$ )
- Интерполяция
  - смесь униграмм, биграмм, триграмм

# Линейная интерполяция

$$\hat{P}(w_n|w_{n-1}w_{n-2}) = \lambda_1 P(w_n|w_{n-1}w_{n-2}) \\ + \lambda_2 P(w_n|w_{n-1}) \\ + \lambda_3 P(w_n)$$

$$\sum_i \lambda_i = 1$$

- $\lambda$  могут зависеть от контекста

$$\hat{P}(w_n|w_{n-2}w_{n-1}) = \lambda_1 (w_{n-2}^{n-1}) P(w_n|w_{n-2}w_{n-1}) \\ + \lambda_2 (w_{n-2}^{n-1}) P(w_n|w_{n-1}) \\ + \lambda_3 (w_{n-2}^{n-1}) P(w_n)$$

# Для очень больших корпусов (веб)

- Наивный откат (*stupid backoff*) [Brants *et al.* 2007]
- Использование относительных частот

$$S(w_i | w_{i-k+1}^{i-1}) = \begin{cases} \frac{\text{count}(w_{i-k+1}^i)}{\text{count}(w_{i-k+1}^{i-1})} & \text{если } \text{count}(w_{i-k+1}^i) > 0 \\ 0.4S(w_i | w_{i-k+2}^{i-1}) & \text{иначе} \end{cases}$$

$$S(w_i) = \frac{\text{count}(w_i)}{N}$$



# **СГЛАЖИВАНИЕ ГУДА-ТЬЮРИНГА**

# Идея метода: рыбалка

- пример Josh Goodman
- Улов: 10 карпов, 3 окуня, 2 сига, 1 форель, 1 лосось, 1 угорь = 18 рыб
- Какова вероятность, что следующая рыба – форель?  $1/18$
- Какова вероятность поймать новую рыбу (сома или щуку)?
  - Предположим, что вероятность поймать новую равна вероятности поймать рыбу, которая до сих пор попадалась только один раз:  $3/18$
- Тогда какова вероятность поймать форель?  $< 1/18$

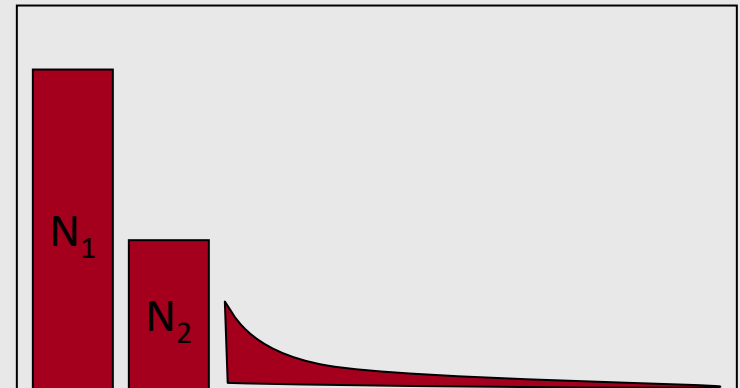
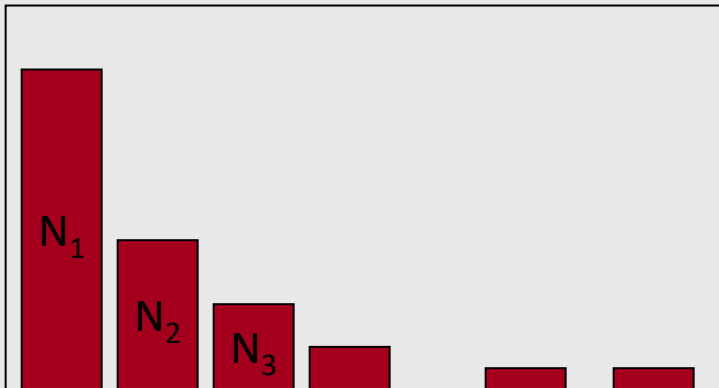
# Пример

$$P_{GT}^* \text{ (нулевая вероятность)} = \frac{N_1}{N} \quad c^* = \frac{(c+1)N_{c+1}}{N_c}$$

- не видели (сом или щука)
  - $c = 0$ :
  - MLE  $p = 0/18 = 0$
  - $P_{GT}^* \text{ (не видели)} = N_1/N = 3/18$
- видели однажды (форель)
  - $c = 1$
  - MLE  $p = 1/18$
  - $C^* \text{ (форель)} = 2 * N_2/N_1 = 2 * 1/3 = 2/3$
  - $P_{GT}^* \text{ (форель)} = 2/3 / 18 = 1/27$

# Аппроксимация низких частот

- В области больших  $k$   $N_k$  идут не подряд (высокие частоты – редкие) и не позволяют использовать формулу напрямую
- Аппроксимация для больших  $k$  с помощью степенной функции



# Пример

- [Church and Gale, 1991]
- новостной корпус 22 млн. слов

$$c^* = \frac{(c + 1)N_{c+1}}{N_c}$$

Count c	Good Turing c*
0	.0000270
1	0.446
2	1.26
3	2.24
4	3.24
5	4.22
6	5.19
7	6.21
8	7.24
9	8.25

# Другие подходы

- Witten-Bell:  
учет разнообразия продолжения  
(низкое разнообразие слов, которые следуют  
за «високосный», по сравнению с «жадному»  
→ новая биграмма с «идиотический» более  
вероятна)
- Kneser-Neu:  
«в другую сторону» – учет вероятности слова  
быть продолжением (биграммы),  
разнообразия «истории»

# Оценка

Perplexity for language models trained on the Europarl corpus

<b>Smoothing method</b>	<b>bigram</b>	<b>trigram</b>	<b>4-gram</b>
Good-Turing	96.2	62.9	59.9
Witten-Bell	97.1	63.8	60.4
Modified Kneser-Ney	95.4	61.6	58.6
Interpolated Modified Kneser-Ney	94.5	59.3	54.0

# НЕЙРОННЫЕ ЯЗЫКОВЫЕ МОДЕЛИ



# Преимущества

- не надо сглаживать
- учет более длинных «историй»
- обобщение

