

Introduction to Information Retrieval

لغت نامه‌ها و بازیابی مقاوم
علی قنبری سرخی

تفاوت نوع (type) / نشانه (token)

- **نشانه:** نمونه‌ای از دنباله کاراکترها (Word یا Term) در اسناد است که با هم به عنوان یک واحد معنایی برای پردازش
- **نوع:** دسته تمامی نشانه‌هاست که دارای دنباله کاراکتر یکسان هستند.

■ مثال:

- In June, the dog likes to chase the cat in the barn.
- 12 word tokens, 9 word types

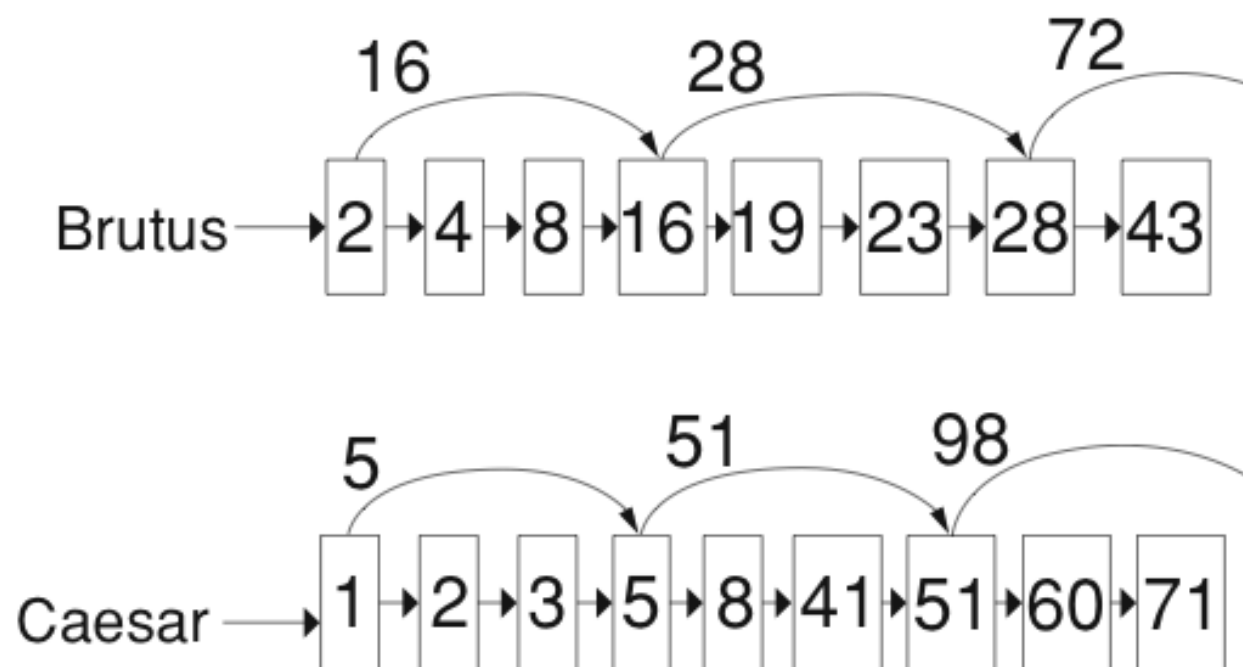
مشکلات نشانه گذاری

- جداکننده‌ها (delimiters) چه هستند؟ فضا خالی؟ آپوستروف؟ خط ربط؟
- در این موارد، بعضی از مواقع آنها را محدود می‌کنند. ولی در بعضی از موارد نمی‌توان محدود کرد.
- در بعضی از زبان‌ها فاصله ندارند! (برای نمونه زبان چینی)
- کلمات مرکب که فاصله ندارند (برای نمونه زبان‌های آلمانی، هلندی و سوئدی)
- مثال کلمات مرکب:
- (*Lebensversicherungsgesellschaftsangestellter*)

مشکلات دسته‌بندی هم‌ارزی

- یک عبارت، یک دسته هم‌ارز از نشانه‌ها است.
- چگونه ما می‌توانیم یک کلاس هم‌ارز را تعریف کنیم؟
- شماره‌ها (3/20/91 vs. 20/3/91)
- نوشتن با حروف بزرگ/غیر حساس کردن به حروف کوچک و بزرگ
- ریشه‌گیری، ریشه‌گیر Porter
- تحلیل ریخت‌شناسی: inflectional vs. Derivational
- مسائل دسته‌بندی هم‌ارزی در زبانه‌های دیگر
- ریخت‌شناسی بسیار پیچیده‌تر نسبت به زبان انگلیسی
- فلاندی: یک فعل تنها می‌تواند ۱۲۰۰۰ فرم متفاوت داشته باشد.
- لهجه، اعراب

اشاره گر پرس



شاخص‌های موقعیتی

- لیست‌های پست‌ها در شاخص بدون موقعیت: هر پست فقط یک شماره سند (docID) است.
- لیست‌های پست‌ها در شاخص موقعیتی: هر پست یک شماره سند (docID) و یک لیست از موقعیت‌ها است.

■ پرس و جوی نمونه “ $to_1 be_2 or_3 not_4 to_5 be_6$ ”

TO, 993427:

< 1: <7, 18, 33, 72, 86, 231>;
 2: <1, 17, 74, 222, 255>;
 4: <8, 16, 190, 429, 433>;
 5: <363, 367>;
 7: <13, 23, 191>; . . . >

BE, 178239:

< 1: <17, 25>;
 4: <17, 191, 291, 430, 434>;
 5: <14, 19, 101>; . . . > Document 4 is a match!

-
- با شاخص موقعیت، ما می‌توانیم به پرس و جویهای اصطلاح پاسخ دهد.
 - با شاخص موقعیت، ما می‌توانیم به پرس و جویهای مجاورت پاسخ دهد.

Take-away

- **بازیابی مقاوم:** در صورتی که انطباق دقیق بین عبارت پرس و جو و عبارت استادوجود نداشته باشد چه کاری باید کرد:
- پرس و جویهای جایگزین
- تصحیح املایی

شاخص وارونه

For each term t , we store a list of all documents that contain t .

BRUTUS →

1	2	4	11	31	45	173	174
---	---	---	----	----	----	-----	-----

CAESAR →

1	2	4	5	6	16	57	132	...
---	---	---	---	---	----	----	-----	-----

CALPURNIA →

2	31	54	101
---	----	----	-----

⋮


dictionary


postings

شاخص وارونه

For each term t , we store a list of all documents that contain t .

BRUTUS → 1 2 4 11 31 45 173 174

CAESAR → 1 2 4 5 6 16 57 132 ...

CALPURNIA → 2 31 54 101

⋮

dictionary

postings

دیکشنری‌ها (لغت نامه‌ها)

- لغت نامه، ساختار داده برای ذخیره واژگان عبارت است.
- واژگان عبارت: داده
- لغت نامه: ساختار داده برای ذخیره واژگان داده است.

لغت نامه به عنوان آرایه‌ای از ورودی‌های با طول ثابت

- برای هر عبارت، ما نیاز به ذخیره یک جفت عبارت داریم:

- تعداد تکرار اسناد

- اشاره‌گر به لیست پست‌ها

- فرض کنید که در حال حاضر ما می‌توانیم این اطلاعات را در ورودی با طول ثابت ذخیره کنیم.

- فرض کنید ما این ورودی‌ها را در آرایه ذخیره می‌کنیم.

لغت نامه به عنوان آرایه‌ای از ورودی‌های با طول ثابت

term	document frequency	pointer to postings list
a	656,265	→
aachen	65	→
...
zulu	221	→

space needed: 20 bytes 4 bytes 4 bytes

- چطور ما یک عبارت پرس و جو q_i را در این آرایه در مرحله پرس و جو جستجو می‌کنیم؟
- به این معنا که: ما چه ساختار داده‌ای برای قرار دادن ورودی (سطر) در آرایه که q_i ذخیره شده، استفاده می‌کنیم

ساختار داده‌ها برای جستجو عبارت

- دو کلاس اصلی ساختار داده‌ها: درخت‌ها و درهم سازی‌ها
- بعضی از سیستم‌های بازیابی اطلاعات از درهم سازی و بعضی‌ها از درخت‌ها استفاده می‌کنند
- معیارهای استفاده از درهم سازی‌ها در مقابل درخت‌ها
 - چه تعداد کلید خواهیم داشت؟
 - آیا این تعداد ثابت باقی می‌ماند یا به شدت تغییر می‌کند- و در مورد تغییرات، تنها کلیدهای جدید درج شده را داریم، یا برخی از کلیدها در لغت‌نامه باید حذف شوند؟
 - فراوانی نسبی دسترسی به کلیدهای مختلف چیست؟

درهم سازی ها

- هر عبارت مجموعه واژگان (کلید) به یک عدد صحیح روی فضایی که به اندازه کافی بزرگ است درهم سازی می شود.
- سعی می شود از تصادم اجتناب شود.
- در زمان پرس و جو مراحل زیر انجام می شود
 - درهم سازی عبارت پرس و جو
 - حل تصادم ها
 - قرار دادن ورودی در آرایه با طول ثابت

درهم سازی ها

- **نقاط قوت:** جستجو در درهم سازی سریع تر از جستجو در درخت می باشد

- زمان جستجو ثابت است

- **نقاط ضعف:**

- هیچ راه آسانی برای پیدا کردن تغییرات کوچک وجود ندارد (برای نمونه عبارات لهجه دار و بی لهجه *resume* vs. *résumé*)

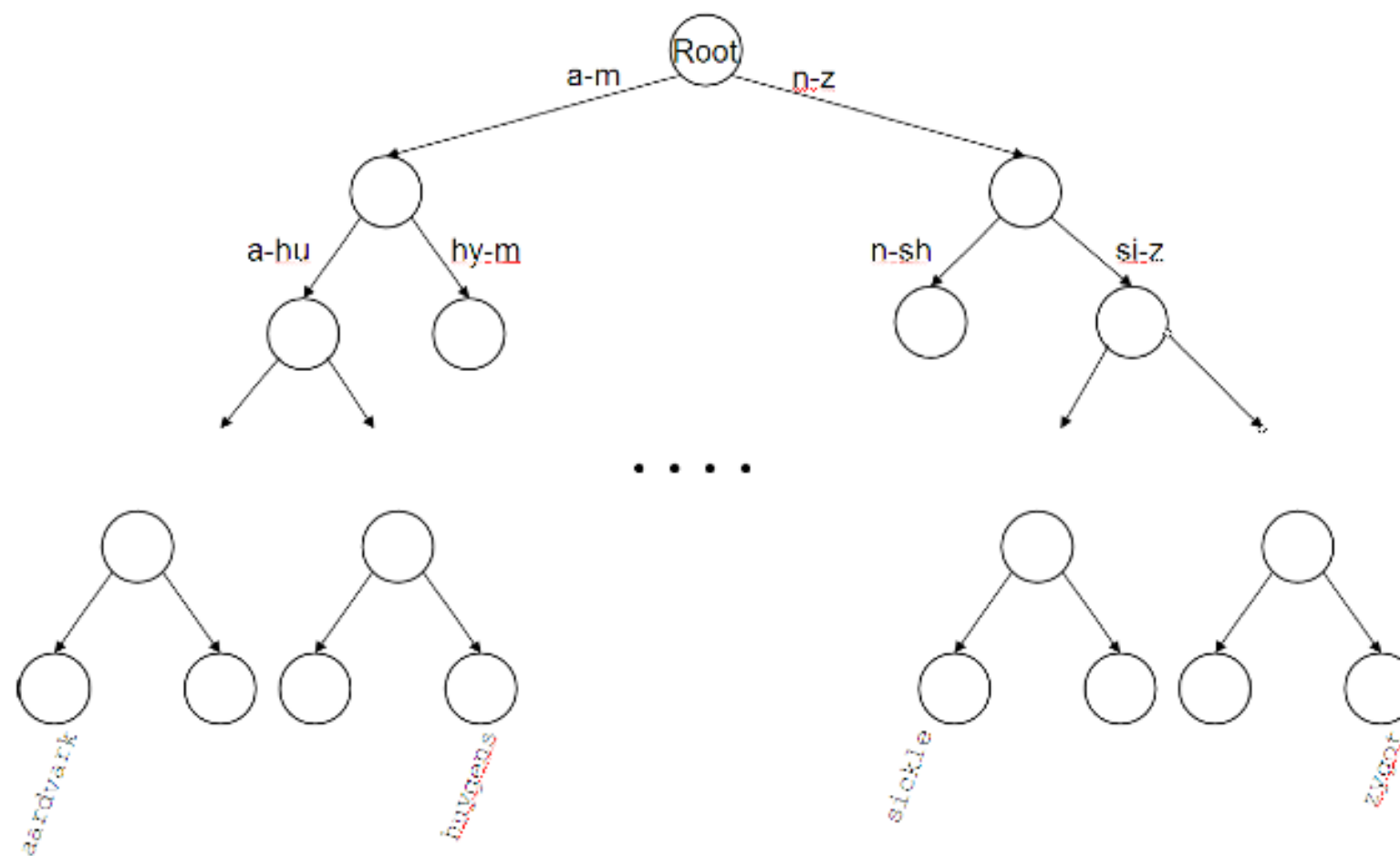
- جستجو پیشوندی انجام نمی شود (برای نمونه تمام عباراتی که با *automat* شروع می شود را نمی توان جستجو کرد)

- در صورتی که عبارت واژگان در حال رشد کردن باشد نیاز به دوباره درهم سازی برای هر چیز به صورت دوره ای می باشد. (مثلا وب)

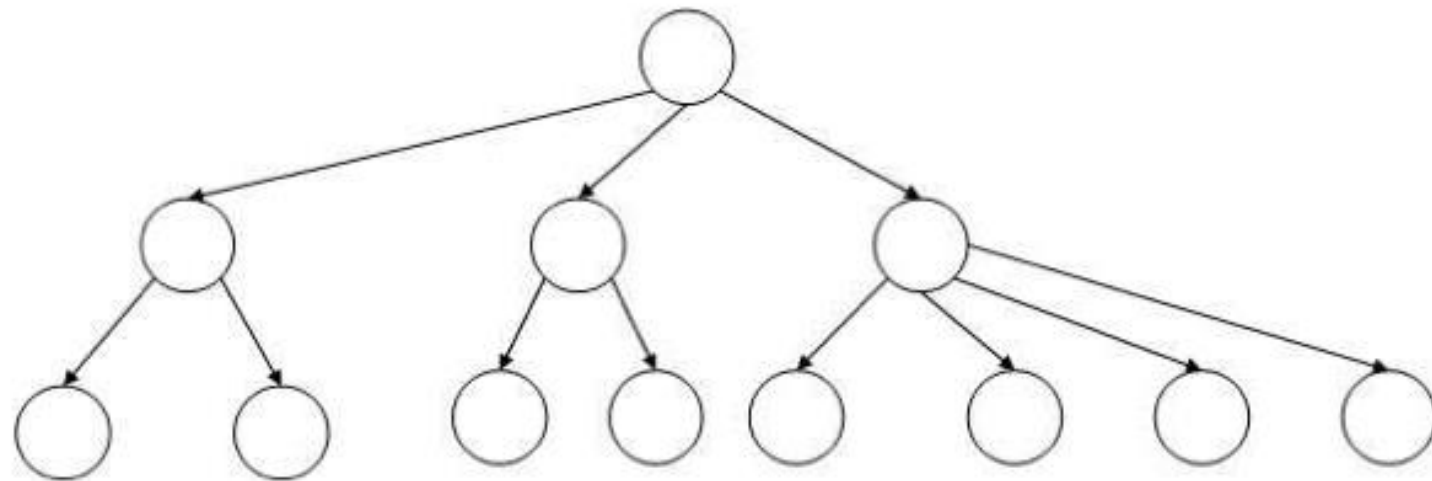
درخت‌ها

- درخت‌ها مشکلات پیشوندی را حل کردند (پیدا کردن همه ی عبارات که با automat شروع می شوند)
- ساده‌ترین درخت: درخت باینری
- جستجو کمی آهسته‌تر از درهم‌سازی ها هست: هزینه جستجو $O(\log M)$ ، در اینجا M اندازه واژگان می‌باشد.
- هزینه $O(\log M)$ فقط برای درخت‌های متوازن می‌باشد
- توازن مجدد درخت‌های باینری بسیار پرهزینه می‌باشد.
- B-trees، مسئله توازن مجدد را کاهش می‌دهند.
- تعریف B-tree: یک درخت جستجو است که در آن هر گره داخلی دارای تعدادی فرزند در بازه $[a, b]$ است. که در آن a و b اعداد صحیح مثبت هستند. برای نمونه درخت باینری $[2, 4]$.

درخت باینری



B-tree



پرس و جوهای جایگزین

- mon^* : پیدا کردن همه‌ی اسنادی که شامل، عبارتی هستن که با mon شروع می‌شوند.
- به راحتی با لغت‌نامه B-tree انجام می‌شود: همه‌ی عباراتی به مانند t ، که در بازه $mon \leq t < moo$ هستند بازیابی می‌شوند.
- mon^* : پیدا کردن همه‌ی اسنادی که شامل، عبارتی هستن که به mon ختم می‌شوند.
- حفظ یک درخت اضافی برای عبارات عقب‌گرد
- سپس همه‌ی عبارات t که در بازه $nom \leq t < non$ هستن بازیابی می‌شوند
- نتیجه: یک مجموعه از عبارات که با پرس و جوی جایگزین انطباق دارند
- سپس اسنادی که شامل هر یک از این عبارات هستند بازیابی می‌شوند.

کنترل * در وسط عبارات

- نمونه: $m*nchen$
- ما می‌توانیم $m*$ و $*nchen$ را در B-tree جستجو نماییم و مجموعه‌های دو عبارت را اشتراک بگیریم
- بسیار پر هزینه می‌باشد.
- جایگزین: شاخص جایگردانی
- ایده پایه: چرخش هر پرس و جو جایگزین، بنابراین $*$ در انتها اتفاق می‌افتد.
- هریم از این چرخش‌ها در لغت‌نامه (B-tree) ذخیره می‌شود.

شخص جایگردانی

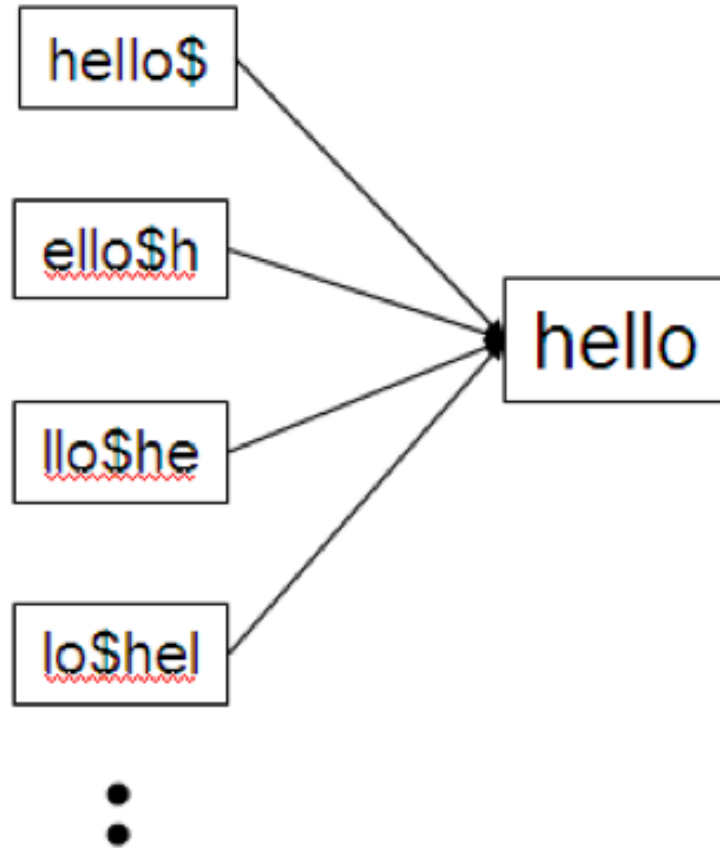
- برای عبارت HELLO: عبارات

hello\$, ello\$h, llo\$he, lo\$hel, and o\$hell

به B-tree اضافه می‌شوند.

- در اینجا \$ یک نماد ویژه است

Permuterm \rightarrow term mapping



شاخص جایگردانی

■ برای عبارت HELLO: ما عبارات hello\$, ello\$h, llo\$he, lo\$hel, and o\$hell را ذخیره کرده ایم.

■ پرس و جوها:

- For X, look up X\$
- For X*, look up X*\$
- For *X, look up X\$*
- For *X*, look up X*
- For X*Y, look up Y\$X*
- Example: For hel*o, look up o\$hel*

■ بهتر است بجای استفاده از شاخص جایگردانی از درخت جایگردانی استفاده شود. ولی شاخص جایگزینی رایج تر است.

■ شاخص جایگردانی به ما امکان می دهد دوران هایی با پیشوند مشخص را انتخاب کنیم.

جستجو کردن در شاخص جایگردانی

- چرخش پرس و جو جایگزین به راست
- استفاده از جستجو B-tree به مانند قبل
- مشکل: لغت نامه جایگزینی بسیار گسترده شده است و شامل تمامی دوران‌های هر عبارات است.
- بیش از ۴ برابر اندازه لغت نامه در مقایسه B-tree (عدد تجربی)

شاخص‌های k-گرمی

- این شاخص بسیار کارآمدتر از شاخص جایگردانی است.
- لغت نامه شامل تمام کارکترهای k-گرم‌هایی است که در هر عبارت مجموعه واژگان رخ می‌دهد (یک k-گرمی دنباله‌ای از k کاراکتر است).
- ۲-گرمی به اصطلاح bigrams گفته می‌شود.
- مثال:
- "April is the cruelest month"
 - bigrams: \$a ap pr ri il l\$ \$i is s\$ \$t th he e\$ \$c cr ru ue el le es st t\$ \$m mo on nt h\$
 - \$ یک نماد مرزی ویژه به مانند قبل می‌باشد.
 - یک شاخص وارونه از ۲-گرمی‌ها برای عباراتی که شامل ۲-گرم هستند را حفظ می‌کند.

لیست پست‌ها در یک شاخص وارونه ۳- گرمی



فرآیند عبارتهای جایگزین در شاخص bigram

- حالا پرس وجوی mon^* می تواند به عنوان $mon AND mo AND on$ اجرا شود.
- تمام عباراتی که با پیشوند mon می باشند را در اختیار ما قرار می دهد.
- اما دارای "مثبت کاذب" به مانند $Moon$ می باشد.
- ما باید این عبارات را در برابر پرس و جو پس فیلتر ($postfilter$) کنیم
- سپس عبارات باقی مانده (زنده مانده) در شاخص وارونه شده عبارت-سند جستجو می شود.
- شاخص k -گرمی در مقابل شاخص جایگزینی
- شاخص k -گرمی بسیار از لحاظ فضا کارآمدتر است.
- شاخص جایگزینی به پسافیلتر نیازی ندارد.

تمرین

- Google یک حمایت محدود برای پرس و جوهای جایگزینی دارد.
- برای مثال، همچنین پرس و جویی به مانند [gen* universit*] عملکرد خیلی خوبی در Google ندارند.
- هدف: شما می‌خواهید عباراتی به مانند the University of Geneva را جستجو نمائید اما نمی‌دانیم کدام لهجه برای استفاده از کلمات فرانسوی برای عبارت university and Geneva استفاده می‌شود.
- بر اساس جستجو Google در سال 2010-04-29: "توجه داشته باشید که عملگر * برای کل کلمات نه قسمتی از کلمات عمل می‌کند"
- اما این مورد به صورت کاملاً درست نمی‌باشد. تلاش کنید عبارات [pythag*] and [m*nchen]
- تمرین: چرا Google به صورت کامل از پرس و جوهای جایگزینی حمایت نمی‌کند؟

فرآیند پرس و جوهای جایگزینی در شاخص عبارت-سند

- مسئله ۱: ما باید به صورت بالقوه یک تعداد زیادی از پرس و جوهای بولی را اجرا نمائیم.
- Most straightforward semantics: Conjunction of disjunctions
- For [gen* universit*]: geneva university OR geneva université OR genève university OR genève université OR general universities OR . . .
- بسیار پرهزینه
- مسئله ۲: کاربران نسبت به نوع نفرت دارند.
- اگر پرس و جوی به صورت اختصار وجود داشته باشد برای مثال اجازه داشته باشید از [pyth* theo*] برای [pythagoras' theorem] استفاده نمائید بسیار از کاربران از آنها استفاده خواهند کرد.
- این عمل هزینه پاسخدهی به پرس و جوها را به صورت قابل توجه افزایش می دهد.
- تا حدودی توسط پیشنهادات گوگل رفع می شود.

تصحیح املائی

- دو اصل استفاده می‌شود:
- اسنادی که شاخص‌گذاری شده تصحیح شوند.
- پرس و جویهای کاربران تصحیح شود.
- دو روش متفاوت برای تصحیح املائی وجود دارد:
- تصحیح املائی عبارت مجزا
- هر عبارت را به صورت غلط املائی چک کنید
- تلاش می‌کنیم تا یک عبارت واحد پرس و جو را تصحیح کنیم - حتی زمانی که پرس و جو چند عبارتی داشته باشیم.
برای مثال: an asteroid that fell **form** the sky
- تصحیح املائی حساس به متن
- به عبارات پیرامون نگاه کن
- در این صورت می‌توان خطای بالا form/from را تصحیح کرد.

تصحیح اسناد

- در این کلاس، ما علاقه‌ای به تعامل تصحیح خطا اسناد (برای مثال MS Word) را نداریم.
- در بازیابی اطلاعات، ما تصحیح اسناد اصلی برای اسناد OCR شده را استفاده می‌کنیم.
- منظور از OCR = optical character recognition
- فلسفه اصلی در IR :
- تغییر ندادن اسناد

تصحیح پرس و جو

- ابتدا: تصحیح املاي عبارت مجزا
- فرض ۱: یک لیست از "عبارات صحیح" که در آن املاي درست آمده است.
- فرض ۲: یک راهی برای محاسبه فاصله بین عبارت اشتباه (با غلط املاي) و عبارت صحیح وجود دارد.
- ساده‌ترین الگوریتم تصحیح املاي، عبارت "صحیح" که کوچکترین فاصله به عبارت اشتباه دارد را بر می‌گرداند.
- Example: informaton → information
- ما می‌توانیم برای لیستی از عبارات صحیح، از واژگانی از تمام عباراتی که در مجموعه ما رخ می‌دهند استفاده نمائیم
- چرا این مشکل‌ساز است؟

استفاده کردن واژگان عبارت

- یک لغت نامه استاندارد (Webster's, OED etc.)
- یک لغت نامه خاص منظور صنعتی (مثلا برای سیستم های بازیابی اطلاعات خاص منظوره شده اند)
- واژگان عبارت از مجموعه که به صورت صحیح وزندار شده اند.

فاصله بین کلمات اشتباه و واژه "صحیح"

- ما چندین گزینه را برای محاسبه فاصله مطالعه می کنیم.

- فاصله ویرایشی و فاصله Levenshtein

- فاصله ویرایش وزندار

- k-gram overlap

فاصله ویرایشی

- فاصله ویرایشی بین دو رشته $S1$ و $S2$
- حداقل تعداد عملیات ویرایشی مورد نیاز برای تبدیل $s1$ به $s2$ است.
- فاصله Levenshtein
 - عملیات ویرایش قابل قبول، درج (insert)، حذف (delete) و جایگزینی (replace) است.
- Levenshtein distance *dog-do*: 1
- Levenshtein distance *cat-cart*: 1
- Levenshtein distance *cat-cut*: 1
- Levenshtein distance *cat-act*: 2
- Damerau-Levenshtein شامل انتقال (transposition) به عنوان چهارمین عملیات ویرایشی قابل قبول می باشد.
 - Damerau-Levenshtein distance *cat-act*: 1

محاسبه فاصله Levenshtein

		f	a	s	t
	0	1	2	3	4
c	1	1	2	3	4
a	2	2	1	2	3
t	3	3	2	2	2
s	4	4	3	2	3

الگوریتم فاصله Levenshtein

LEVENSHTEINDISTANCE(s_1, s_2)

```
1  for  $i \leftarrow 0$  to  $|s_1|$ 
2  do  $m[i, 0] = i$ 
3  for  $j \leftarrow 0$  to  $|s_2|$ 
4  do  $m[0, j] = j$ 
5  for  $i \leftarrow 1$  to  $|s_1|$ 
6  do for  $j \leftarrow 1$  to  $|s_2|$ 
7      do if  $s_1[i] = s_2[j]$ 
8          then  $m[i, j] = \min\{m[i-1, j]+1, m[i, j-1]+1, m[i-1, j-1]\}$ 
9          else  $m[i, j] = \min\{m[i-1, j]+1, m[i, j-1]+1, m[i-1, j-1]+1\}$ 
10 return  $m[|s_1|, |s_2|]$ 
```

Operations: insert (cost 1), delete (cost 1), replace (cost 1), copy (cost 0)

الگوریتم فاصله Levenshtein

LEVENSHTEINDISTANCE(s_1, s_2)

```
1  for  $i \leftarrow 0$  to  $|s_1|$ 
2  do  $m[i, 0] = i$ 
3  for  $j \leftarrow 0$  to  $|s_2|$ 
4  do  $m[0, j] = j$ 
5  for  $i \leftarrow 1$  to  $|s_1|$ 
6  do for  $j \leftarrow 1$  to  $|s_2|$ 
7      do if  $s_1[i] = s_2[j]$ 
8          then  $m[i, j] = \min\{m[i-1, j]+1, m[i, j-1]+1, m[i-1, j-1]\}$ 
9          else  $m[i, j] = \min\{m[i-1, j]+1, m[i, j-1]+1, m[i-1, j-1]+1\}$ 
10 return  $m[|s_1|, |s_2|]$ 
```

Operations: insert (cost 1), delete (cost 1), replace (cost 1), copy (cost 0)

الگوریتم فاصله Levenshtein

LEVENSHTEINDISTANCE(s_1, s_2)

```
1  for  $i \leftarrow 0$  to  $|s_1|$ 
2  do  $m[i, 0] = i$ 
3  for  $j \leftarrow 0$  to  $|s_2|$ 
4  do  $m[0, j] = j$ 
5  for  $i \leftarrow 1$  to  $|s_1|$ 
6  do for  $j \leftarrow 1$  to  $|s_2|$ 
7      do if  $s_1[i] = s_2[j]$ 
8          then  $m[i, j] = \min\{m[i-1, j]+1, m[i, j-1]+1, m[i-1, j-1]\}$ 
9          else  $m[i, j] = \min\{m[i-1, j]+1, m[i, j-1]+1, m[i-1, j-1]+1\}$ 
10 return  $m[|s_1|, |s_2|]$ 
```

Operations: insert (cost 1), delete (cost 1), replace (cost 1), copy (cost 0)

الگوریتم فاصله Levenshtein

LEVENSHTEINDISTANCE(s_1, s_2)

```
1  for  $i \leftarrow 0$  to  $|s_1|$ 
2  do  $m[i, 0] = i$ 
3  for  $j \leftarrow 0$  to  $|s_2|$ 
4  do  $m[0, j] = j$ 
5  for  $i \leftarrow 1$  to  $|s_1|$ 
6  do for  $j \leftarrow 1$  to  $|s_2|$ 
7      do if  $s_1[i] = s_2[j]$ 
8          then  $m[i, j] = \min\{m[i-1, j]+1, m[i, j-1]+1, m[i-1, j-1]\}$ 
9          else  $m[i, j] = \min\{m[i-1, j]+1, m[i, j-1]+1, m[i-1, j-1]+1\}$ 
10 return  $m[|s_1|, |s_2|]$ 
```

Operations: insert (cost 1), delete (cost 1), **replace (cost 1)**, copy (cost 0)

الگوریتم فاصله Levenshtein

LEVENSHTEINDISTANCE(s_1, s_2)

```
1  for  $i \leftarrow 0$  to  $|s_1|$ 
2  do  $m[i, 0] = i$ 
3  for  $j \leftarrow 0$  to  $|s_2|$ 
4  do  $m[0, j] = j$ 
5  for  $i \leftarrow 1$  to  $|s_1|$ 
6  do for  $j \leftarrow 1$  to  $|s_2|$ 
7      do if  $s_1[i] = s_2[j]$ 
8          then  $m[i, j] = \min\{m[i-1, j]+1, m[i, j-1]+1, m[i-1, j-1]\}$ 
9          else  $m[i, j] = \min\{m[i-1, j]+1, m[i, j-1]+1, m[i-1, j-1]+1\}$ 
10 return  $m[|s_1|, |s_2|]$ 
```

Operations: insert (cost 1), delete (cost 1), replace (cost 1), **copy**
(cost 0)

مثال فاصله Levenshtein

		f	a	s	t
	<div><div></div><div>0</div></div>	<div><div>1</div><div>1</div></div>	<div><div>2</div><div>2</div></div>	<div><div>3</div><div>3</div></div>	<div><div>4</div><div>4</div></div>
c	<div><div>1</div><div>1</div></div>	<div><div>1</div><div>2</div><div>2</div><div>1</div></div>	<div><div>2</div><div>3</div><div>2</div><div>2</div></div>	<div><div>3</div><div>4</div><div>3</div><div>3</div></div>	<div><div>4</div><div>5</div><div>4</div><div>4</div></div>
a	<div><div>2</div><div>2</div></div>	<div><div>2</div><div>2</div><div>3</div><div>2</div></div>	<div><div>1</div><div>3</div><div>3</div><div>1</div></div>	<div><div>3</div><div>4</div><div>2</div><div>2</div></div>	<div><div>4</div><div>5</div><div>3</div><div>3</div></div>
t	<div><div>3</div><div>3</div></div>	<div><div>3</div><div>3</div><div>4</div><div>3</div></div>	<div><div>3</div><div>2</div><div>4</div><div>2</div></div>	<div><div>2</div><div>3</div><div>3</div><div>2</div></div>	<div><div>2</div><div>4</div><div>3</div><div>2</div></div>
s	<div><div>4</div><div>4</div></div>	<div><div>4</div><div>4</div><div>5</div><div>4</div></div>	<div><div>4</div><div>3</div><div>5</div><div>3</div></div>	<div><div>2</div><div>3</div><div>4</div><div>2</div></div>	<div><div>3</div><div>3</div><div>3</div><div>3</div></div>

هر سلول ماتریس Levenshtein

هزینه دریافتی از همسایه بالا سمت چپ (copy or replace)	هزینه دریافت از همسایه بالایی (delete)
هزینه دریافت از همسایه چپ (insert)	کمینه سه عنصر این ماتریس؛ ارزانترین راه برای رسیدن به اینجا

مثال فاصله Levenshtein

		f	a	s	t
	<div><div></div><div>0</div></div>	<div><div>1</div><div>1</div></div>	<div><div>2</div><div>2</div></div>	<div><div>3</div><div>3</div></div>	<div><div>4</div><div>4</div></div>
c	<div><div>1</div><div>1</div></div>	<div><div>1</div><div>2</div><div>2</div><div>1</div></div>	<div><div>2</div><div>3</div><div>2</div><div>2</div></div>	<div><div>3</div><div>4</div><div>3</div><div>3</div></div>	<div><div>4</div><div>5</div><div>4</div><div>4</div></div>
a	<div><div>2</div><div>2</div></div>	<div><div>2</div><div>2</div><div>3</div><div>2</div></div>	<div><div>1</div><div>3</div><div>3</div><div>1</div></div>	<div><div>3</div><div>4</div><div>2</div><div>2</div></div>	<div><div>4</div><div>5</div><div>3</div><div>3</div></div>
t	<div><div>3</div><div>3</div></div>	<div><div>3</div><div>3</div><div>4</div><div>3</div></div>	<div><div>3</div><div>2</div><div>4</div><div>2</div></div>	<div><div>2</div><div>3</div><div>3</div><div>2</div></div>	<div><div>2</div><div>4</div><div>3</div><div>2</div></div>
s	<div><div>4</div><div>4</div></div>	<div><div>4</div><div>4</div><div>5</div><div>4</div></div>	<div><div>4</div><div>3</div><div>5</div><div>3</div></div>	<div><div>2</div><div>3</div><div>4</div><div>2</div></div>	<div><div>3</div><div>3</div><div>3</div><div>3</div></div>

فاصله ویرایشی وزندار

- وزن یک عملیات بستگی به کارکترهای درگیر دارد.
- برای مدیریت کردن خطاهای کیبورد. برای نمونه: احتمال اینکه m به صورت اشتباه n تایپ شود بیشتر از q می باشد.
- بنابراین، جایگزینی (replacing)، m به وسیله n یک فاصله ویرایشی کوچکتر نسبت به m به وسیله q دارد.
- در اینجا ما به یک ماتریس وزندار به عنوان ورودی نیاز داریم.
- برنامه نویسی پویا برای کنترل وزن ها اصلاح می شود.

استفاده از فاصله ویرایشی تصحیح املایی

- با توجه به پرس و جو:
- ابتدا تمام توالی‌های کارکترها را در یک فاصله ویرایشی از پیش تعیین شده (احتمالا وزندار) ثبت کنید
- اشتراک این مجموعه را با لیست کلمات "صحیح" خود بدست آورید.
- سپس عبارتی را از این اشتراک‌گیری به کاربر پیشنهاد دهید.

تمرین

- ماتریش فاصله Levenshtein را برای OSLO – SNOW محاسبه کنید؟
- چرا عملیات ویرایشی Levenshtein عبارت cat را به catcat تبدیل می کند؟

		s	n	o	w
	<div><div></div><div>0</div></div>	<div><div>1</div><div>1</div></div>	<div><div>2</div><div>2</div></div>	<div><div>3</div><div>3</div></div>	<div><div>4</div><div>4</div></div>
o	<div><div>1</div><div>1</div></div>				
s	<div><div>2</div><div>2</div></div>				
l	<div><div>3</div><div>3</div></div>				
o	<div><div>4</div><div>4</div></div>				

		s	n	o	w
	<div><div></div><div>0</div></div>	<div><div>1</div><div>1</div></div>	<div><div>2</div><div>2</div></div>	<div><div>3</div><div>3</div></div>	<div><div>4</div><div>4</div></div>
o	<div><div>1</div><div>1</div></div>	<div><div>1</div><div>2</div><div>2</div><div>?</div></div>			
s	<div><div>2</div><div>2</div></div>				
l	<div><div>3</div><div>3</div></div>				
o	<div><div>4</div><div>4</div></div>				

			s	n	o	w
		0	11	22	33	44
o		11	1221			
s		22				
l		33				
o		44				

		s	n	o	w
	<div><div></div><div>0</div></div>	<div><div>11</div></div>	<div><div>22</div></div>	<div><div>33</div></div>	<div><div>44</div></div>
o	<div><div>11</div></div>	<div><div>1221</div></div>	<div><div>232?</div></div>		
s	<div><div>22</div></div>				
l	<div><div>33</div></div>				
o	<div><div>44</div></div>				

		s	n	o	w
	<div><div></div><div>0</div></div>	<div><div>1</div><div>1</div></div>	<div><div>2</div><div>2</div></div>	<div><div>3</div><div>3</div></div>	<div><div>4</div><div>4</div></div>
o	<div><div>1</div><div>1</div></div>	<div><div>1</div><div>2</div><div>2</div><div>1</div></div>	<div><div>2</div><div>3</div><div>2</div><div>2</div></div>		
s	<div><div>2</div><div>2</div></div>				
l	<div><div>3</div><div>3</div></div>				
o	<div><div>4</div><div>4</div></div>				

		s	n	o	w
	<div><div></div><div>0</div></div>	<div><div>1</div><div>1</div></div>	<div><div>2</div><div>2</div></div>	<div><div>3</div><div>3</div></div>	<div><div>4</div><div>4</div></div>
o	<div><div>1</div><div>1</div></div>	<div><div>1</div><div>2</div><div>2</div><div>1</div></div>	<div><div>2</div><div>3</div><div>2</div><div>2</div></div>	<div><div>2</div><div>4</div><div>3</div><div>?</div></div>	
s	<div><div>2</div><div>2</div></div>				
l	<div><div>3</div><div>3</div></div>				
o	<div><div>4</div><div>4</div></div>				

			s	n	o	w
	<div><div></div><div>0</div></div>	<div><div>1</div><div>1</div></div>	<div><div>2</div><div>2</div></div>	<div><div>3</div><div>3</div></div>	<div><div>4</div><div>4</div></div>	
o	<div><div>1</div><div>1</div></div>	<div><div>1</div><div>2</div><div>2</div><div>1</div></div>	<div><div>2</div><div>3</div><div>2</div><div>2</div></div>	<div><div>2</div><div>4</div><div>3</div><div>2</div></div>		
s	<div><div>2</div><div>2</div></div>					
l	<div><div>3</div><div>3</div></div>					
o	<div><div>4</div><div>4</div></div>					

			s	n	o	w
	<div><div></div><div>0</div></div>	<div><div>1</div><div>1</div></div>	<div><div>2</div><div>2</div></div>	<div><div>3</div><div>3</div></div>	<div><div>4</div><div>4</div></div>	
o	<div><div>1</div><div>1</div></div>	<div><div>1</div><div>2</div><div>2</div><div>1</div></div>	<div><div>2</div><div>3</div><div>2</div><div>2</div></div>	<div><div>2</div><div>4</div><div>3</div><div>2</div></div>	<div><div>4</div><div>5</div><div>3</div><div>?</div></div>	
s	<div><div>2</div><div>2</div></div>					
l	<div><div>3</div><div>3</div></div>					
o	<div><div>4</div><div>4</div></div>					

		s	n	o	w
	<div><div></div><div>0</div></div>	<div><div>1</div><div>1</div></div>	<div><div>2</div><div>2</div></div>	<div><div>3</div><div>3</div></div>	<div><div>4</div><div>4</div></div>
o	<div><div>1</div><div>1</div></div>	<div><div>1</div><div>2</div><div>2</div><div>1</div></div>	<div><div>2</div><div>3</div><div>2</div><div>2</div></div>	<div><div>2</div><div>4</div><div>3</div><div>2</div></div>	<div><div>4</div><div>5</div><div>3</div><div>3</div></div>
s	<div><div>2</div><div>2</div></div>				
l	<div><div>3</div><div>3</div></div>				
o	<div><div>4</div><div>4</div></div>				

		s	n	o	w
	<div><div></div><div>0</div></div>	<div><div>11</div></div>	<div><div>22</div></div>	<div><div>33</div></div>	<div><div>44</div></div>
o	<div><div>11</div></div>	<div><div>12</div><div>21</div></div>	<div><div>23</div><div>22</div></div>	<div><div>24</div><div>32</div></div>	<div><div>45</div><div>33</div></div>
s	<div><div>22</div></div>	<div><div>12</div><div>3?</div></div>			
l	<div><div>33</div></div>				
o	<div><div>44</div></div>				

			s		n		o		w	
		<div><div></div><div>0</div></div>	<div><div>1</div><div>1</div></div>		<div><div>2</div><div>2</div></div>		<div><div>3</div><div>3</div></div>		<div><div>4</div><div>4</div></div>	
o		<div><div>1</div><div>1</div></div>	<div><div>1</div><div>2</div><div>2</div><div>1</div></div>		<div><div>2</div><div>3</div><div>2</div><div>2</div></div>		<div><div>2</div><div>4</div><div>3</div><div>2</div></div>		<div><div>4</div><div>5</div><div>3</div><div>3</div></div>	
s		<div><div>2</div><div>2</div></div>	<div><div>1</div><div>2</div><div>3</div><div>1</div></div>							
l		<div><div>3</div><div>3</div></div>								
o		<div><div>4</div><div>4</div></div>								

		s	n	o	w
	<div><div></div><div>0</div></div>	<div><div>11</div></div>	<div><div>22</div></div>	<div><div>33</div></div>	<div><div>44</div></div>
o	<div><div>11</div></div>	<div><div>1221</div></div>	<div><div>2322</div></div>	<div><div>2432</div></div>	<div><div>4533</div></div>
s	<div><div>22</div></div>	<div><div>1231</div></div>	<div><div>232?</div></div>		
l	<div><div>33</div></div>				
o	<div><div>44</div></div>				

		s	n	o	w
	<div><div></div><div>0</div></div>	<div><div>11</div></div>	<div><div>22</div></div>	<div><div>33</div></div>	<div><div>44</div></div>
o	<div><div>11</div></div>	<div><div>1221</div></div>	<div><div>2322</div></div>	<div><div>2432</div></div>	<div><div>4533</div></div>
s	<div><div>22</div></div>	<div><div>1231</div></div>	<div><div>2322</div></div>		
l	<div><div>33</div></div>				
o	<div><div>44</div></div>				

		s	n	o	w
	<div><div></div><div>0</div></div>	<div><div>11</div></div>	<div><div>22</div></div>	<div><div>33</div></div>	<div><div>44</div></div>
o	<div><div>11</div></div>	<div><div>12</div><div>21</div></div>	<div><div>23</div><div>22</div></div>	<div><div>24</div><div>32</div></div>	<div><div>45</div><div>33</div></div>
s	<div><div>22</div></div>	<div><div>12</div><div>31</div></div>	<div><div>23</div><div>22</div></div>	<div><div>33</div><div>3?</div></div>	
l	<div><div>33</div></div>				
o	<div><div>44</div></div>				

		s	n	o	w
	<div><div></div><div>0</div></div>	<div><div>11</div></div>	<div><div>22</div></div>	<div><div>33</div></div>	<div><div>44</div></div>
o	<div><div>11</div></div>	<div><div>12</div><div>21</div></div>	<div><div>23</div><div>22</div></div>	<div><div>24</div><div>32</div></div>	<div><div>45</div><div>33</div></div>
s	<div><div>22</div></div>	<div><div>12</div><div>31</div></div>	<div><div>23</div><div>22</div></div>	<div><div>33</div><div>33</div></div>	
l	<div><div>33</div></div>				
o	<div><div>44</div></div>				

		s	n	o	w
	<div><div></div><div>0</div></div>	<div><div>1</div><div>1</div></div>	<div><div>2</div><div>2</div></div>	<div><div>3</div><div>3</div></div>	<div><div>4</div><div>4</div></div>
o	<div><div>1</div><div>1</div></div>	<div><div>1</div><div>2</div><div>2</div><div>1</div></div>	<div><div>2</div><div>3</div><div>2</div><div>2</div></div>	<div><div>2</div><div>4</div><div>3</div><div>2</div></div>	<div><div>4</div><div>5</div><div>3</div><div>3</div></div>
s	<div><div>2</div><div>2</div></div>	<div><div>1</div><div>2</div><div>3</div><div>1</div></div>	<div><div>2</div><div>3</div><div>2</div><div>2</div></div>	<div><div>3</div><div>3</div><div>3</div><div>3</div></div>	<div><div>3</div><div>4</div><div>4</div><div>?</div></div>
l	<div><div>3</div><div>3</div></div>				
o	<div><div>4</div><div>4</div></div>				

		s	n	o	w
	<div><div></div><div>0</div></div>	<div><div>11</div></div>	<div><div>22</div></div>	<div><div>33</div></div>	<div><div>44</div></div>
o	<div><div>11</div></div>	<div><div>1221</div></div>	<div><div>2322</div></div>	<div><div>2432</div></div>	<div><div>4533</div></div>
s	<div><div>22</div></div>	<div><div>1231</div></div>	<div><div>2322</div></div>	<div><div>3333</div></div>	<div><div>3443</div></div>
l	<div><div>33</div></div>				
o	<div><div>44</div></div>				

		s	n	o	w
	<div><div></div><div>0</div></div>	<div><div>11</div></div>	<div><div>22</div></div>	<div><div>33</div></div>	<div><div>44</div></div>
o	<div><div>11</div></div>	<div><div>1221</div></div>	<div><div>2322</div></div>	<div><div>2432</div></div>	<div><div>4533</div></div>
s	<div><div>22</div></div>	<div><div>1231</div></div>	<div><div>2322</div></div>	<div><div>3333</div></div>	<div><div>3443</div></div>
l	<div><div>33</div></div>	<div><div>324?</div></div>			
o	<div><div>44</div></div>				

		s	n	o	w
	<div><div></div><div>0</div></div>	<div><div>11</div></div>	<div><div>22</div></div>	<div><div>33</div></div>	<div><div>44</div></div>
o	<div><div>11</div></div>	<div><div>1221</div></div>	<div><div>2322</div></div>	<div><div>2432</div></div>	<div><div>4533</div></div>
s	<div><div>22</div></div>	<div><div>1231</div></div>	<div><div>2322</div></div>	<div><div>3333</div></div>	<div><div>3443</div></div>
l	<div><div>33</div></div>	<div><div>3242</div></div>			
o	<div><div>44</div></div>				

		s	n	o	w
	<div><div></div><div>0</div></div>	<div><div>11</div></div>	<div><div>22</div></div>	<div><div>33</div></div>	<div><div>44</div></div>
o	<div><div>11</div></div>	<div><div>1221</div></div>	<div><div>2322</div></div>	<div><div>2432</div></div>	<div><div>4533</div></div>
s	<div><div>22</div></div>	<div><div>1231</div></div>	<div><div>2322</div></div>	<div><div>3333</div></div>	<div><div>3443</div></div>
l	<div><div>33</div></div>	<div><div>3242</div></div>	<div><div>233?</div></div>		
o	<div><div>44</div></div>				

			s		n		o		w	
		<div><div></div><div>0</div></div>	<div><div>1</div><div>1</div></div>		<div><div>2</div><div>2</div></div>		<div><div>3</div><div>3</div></div>		<div><div>4</div><div>4</div></div>	
o		<div><div>1</div><div>1</div></div>	<div><div>1</div><div>2</div><div>2</div><div>1</div></div>		<div><div>2</div><div>3</div><div>2</div><div>2</div></div>		<div><div>2</div><div>4</div><div>3</div><div>2</div></div>		<div><div>4</div><div>5</div><div>3</div><div>3</div></div>	
s		<div><div>2</div><div>2</div></div>	<div><div>1</div><div>2</div><div>3</div><div>1</div></div>		<div><div>2</div><div>3</div><div>2</div><div>2</div></div>		<div><div>3</div><div>3</div><div>3</div><div>3</div></div>		<div><div>3</div><div>4</div><div>4</div><div>3</div></div>	
l		<div><div>3</div><div>3</div></div>	<div><div>3</div><div>2</div><div>4</div><div>2</div></div>		<div><div>2</div><div>3</div><div>3</div><div>2</div></div>					
o		<div><div>4</div><div>4</div></div>								

		s	n	o	w
	<div><div></div><div>0</div></div>	<div><div>11</div></div>	<div><div>22</div></div>	<div><div>33</div></div>	<div><div>44</div></div>
o	<div><div>11</div></div>	<div><div>12</div><div>21</div></div>	<div><div>23</div><div>22</div></div>	<div><div>24</div><div>32</div></div>	<div><div>45</div><div>33</div></div>
s	<div><div>22</div></div>	<div><div>12</div><div>31</div></div>	<div><div>23</div><div>22</div></div>	<div><div>33</div><div>33</div></div>	<div><div>34</div><div>43</div></div>
l	<div><div>33</div></div>	<div><div>32</div><div>42</div></div>	<div><div>23</div><div>32</div></div>	<div><div>34</div><div>3?</div></div>	
o	<div><div>44</div></div>				

		s	n	o	w
	<div><div></div><div>0</div></div>	<div><div>11</div></div>	<div><div>22</div></div>	<div><div>33</div></div>	<div><div>44</div></div>
o	<div><div>11</div></div>	<div><div>12</div><div>21</div></div>	<div><div>23</div><div>22</div></div>	<div><div>24</div><div>32</div></div>	<div><div>45</div><div>33</div></div>
s	<div><div>22</div></div>	<div><div>12</div><div>31</div></div>	<div><div>23</div><div>22</div></div>	<div><div>33</div><div>33</div></div>	<div><div>34</div><div>43</div></div>
l	<div><div>33</div></div>	<div><div>32</div><div>42</div></div>	<div><div>23</div><div>32</div></div>	<div><div>34</div><div>33</div></div>	
o	<div><div>44</div></div>				

		s	n	o	w
	<div><div></div><div>0</div></div>	<div><div>11</div></div>	<div><div>22</div></div>	<div><div>33</div></div>	<div><div>44</div></div>
o	<div><div>11</div></div>	<div><div>12</div><div>21</div></div>	<div><div>23</div><div>22</div></div>	<div><div>24</div><div>32</div></div>	<div><div>45</div><div>33</div></div>
s	<div><div>22</div></div>	<div><div>12</div><div>31</div></div>	<div><div>23</div><div>22</div></div>	<div><div>33</div><div>33</div></div>	<div><div>34</div><div>43</div></div>
l	<div><div>33</div></div>	<div><div>32</div><div>42</div></div>	<div><div>23</div><div>32</div></div>	<div><div>34</div><div>33</div></div>	<div><div>44</div><div>4?</div></div>
o	<div><div>44</div></div>				

			s	n	o	w
	<div><div></div><div>0</div></div>	<div><div>1</div><div>1</div></div>	<div><div>2</div><div>2</div></div>	<div><div>3</div><div>3</div></div>	<div><div>4</div><div>4</div></div>	
o	<div><div>1</div><div>1</div></div>	<div><div>1</div><div>2</div><div>2</div><div>1</div></div>	<div><div>2</div><div>3</div><div>2</div><div>2</div></div>	<div><div>2</div><div>4</div><div>3</div><div>2</div></div>	<div><div>4</div><div>5</div><div>3</div><div>3</div></div>	
s	<div><div>2</div><div>2</div></div>	<div><div>1</div><div>2</div><div>3</div><div>1</div></div>	<div><div>2</div><div>3</div><div>2</div><div>2</div></div>	<div><div>3</div><div>3</div><div>3</div><div>3</div></div>	<div><div>3</div><div>4</div><div>4</div><div>3</div></div>	
l	<div><div>3</div><div>3</div></div>	<div><div>3</div><div>2</div><div>4</div><div>2</div></div>	<div><div>2</div><div>3</div><div>3</div><div>2</div></div>	<div><div>3</div><div>4</div><div>3</div><div>3</div></div>	<div><div>4</div><div>4</div><div>4</div><div>4</div></div>	
o	<div><div>4</div><div>4</div></div>					

		s	n	o	w
	<div><div></div><div>0</div></div>	<div><div>11</div></div>	<div><div>22</div></div>	<div><div>33</div></div>	<div><div>44</div></div>
o	<div><div>11</div></div>	<div><div>12</div><div>21</div></div>	<div><div>23</div><div>22</div></div>	<div><div>24</div><div>32</div></div>	<div><div>45</div><div>33</div></div>
s	<div><div>22</div></div>	<div><div>12</div><div>31</div></div>	<div><div>23</div><div>22</div></div>	<div><div>33</div><div>33</div></div>	<div><div>34</div><div>43</div></div>
l	<div><div>33</div></div>	<div><div>32</div><div>42</div></div>	<div><div>23</div><div>32</div></div>	<div><div>34</div><div>33</div></div>	<div><div>44</div><div>44</div></div>
o	<div><div>44</div></div>	<div><div>43</div><div>5?</div></div>			

		s	n	o	w
	<div><div></div><div>0</div></div>	<div><div>11</div></div>	<div><div>22</div></div>	<div><div>33</div></div>	<div><div>44</div></div>
o	<div><div>11</div></div>	<div><div>12</div><div>21</div></div>	<div><div>23</div><div>22</div></div>	<div><div>24</div><div>32</div></div>	<div><div>45</div><div>33</div></div>
s	<div><div>2</div><div>2</div></div>	<div><div>12</div><div>31</div></div>	<div><div>23</div><div>22</div></div>	<div><div>33</div><div>33</div></div>	<div><div>34</div><div>43</div></div>
l	<div><div>3</div><div>3</div></div>	<div><div>32</div><div>42</div></div>	<div><div>23</div><div>32</div></div>	<div><div>34</div><div>33</div></div>	<div><div>44</div><div>44</div></div>
o	<div><div>4</div><div>4</div></div>	<div><div>43</div><div>53</div></div>			

		s	n	o	w
	<div><div></div><div>0</div></div>	<div><div>11</div></div>	<div><div>22</div></div>	<div><div>33</div></div>	<div><div>44</div></div>
o	<div><div>11</div></div>	<div><div>12</div><div>21</div></div>	<div><div>23</div><div>22</div></div>	<div><div>24</div><div>32</div></div>	<div><div>45</div><div>33</div></div>
s	<div><div>22</div></div>	<div><div>12</div><div>31</div></div>	<div><div>23</div><div>22</div></div>	<div><div>33</div><div>33</div></div>	<div><div>34</div><div>43</div></div>
l	<div><div>33</div></div>	<div><div>32</div><div>42</div></div>	<div><div>23</div><div>32</div></div>	<div><div>34</div><div>33</div></div>	<div><div>44</div><div>44</div></div>
o	<div><div>44</div></div>	<div><div>43</div><div>53</div></div>	<div><div>33</div><div>4?</div></div>		

		s	n	o	w
	<div><div></div><div>0</div></div>	<div><div>11</div></div>	<div><div>22</div></div>	<div><div>33</div></div>	<div><div>44</div></div>
o	<div><div>11</div></div>	<div><div>12</div><div>21</div></div>	<div><div>23</div><div>22</div></div>	<div><div>24</div><div>32</div></div>	<div><div>45</div><div>33</div></div>
s	<div><div>22</div></div>	<div><div>12</div><div>31</div></div>	<div><div>23</div><div>22</div></div>	<div><div>33</div><div>33</div></div>	<div><div>34</div><div>43</div></div>
l	<div><div>33</div></div>	<div><div>32</div><div>42</div></div>	<div><div>23</div><div>32</div></div>	<div><div>34</div><div>33</div></div>	<div><div>44</div><div>44</div></div>
o	<div><div>44</div></div>	<div><div>43</div><div>53</div></div>	<div><div>33</div><div>43</div></div>		

		s	n	o	w
	<div><div></div><div>0</div></div>	<div><div>11</div></div>	<div><div>22</div></div>	<div><div>33</div></div>	<div><div>44</div></div>
o	<div><div>11</div></div>	<div><div>1221</div></div>	<div><div>2322</div></div>	<div><div>2432</div></div>	<div><div>4533</div></div>
s	<div><div>22</div></div>	<div><div>1231</div></div>	<div><div>2322</div></div>	<div><div>3333</div></div>	<div><div>3443</div></div>
l	<div><div>33</div></div>	<div><div>3242</div></div>	<div><div>2332</div></div>	<div><div>3433</div></div>	<div><div>4444</div></div>
o	<div><div>44</div></div>	<div><div>4353</div></div>	<div><div>3343</div></div>	<div><div>244?</div></div>	

		s	n	o	w
	<div><div></div><div>0</div></div>	<div><div>11</div></div>	<div><div>22</div></div>	<div><div>33</div></div>	<div><div>44</div></div>
o	<div><div>11</div></div>	<div><div>12</div><div>21</div></div>	<div><div>23</div><div>22</div></div>	<div><div>24</div><div>32</div></div>	<div><div>45</div><div>33</div></div>
s	<div><div>22</div></div>	<div><div>12</div><div>31</div></div>	<div><div>23</div><div>22</div></div>	<div><div>33</div><div>33</div></div>	<div><div>34</div><div>43</div></div>
l	<div><div>33</div></div>	<div><div>32</div><div>42</div></div>	<div><div>23</div><div>32</div></div>	<div><div>34</div><div>33</div></div>	<div><div>44</div><div>44</div></div>
o	<div><div>44</div></div>	<div><div>43</div><div>53</div></div>	<div><div>33</div><div>43</div></div>	<div><div>24</div><div>42</div></div>	

		s	n	o	w
	<div><div></div><div>0</div></div>	<div><div>11</div></div>	<div><div>22</div></div>	<div><div>33</div></div>	<div><div>44</div></div>
o	<div><div>11</div></div>	<div><div>12</div><div>21</div></div>	<div><div>23</div><div>22</div></div>	<div><div>24</div><div>32</div></div>	<div><div>45</div><div>33</div></div>
s	<div><div>22</div></div>	<div><div>12</div><div>31</div></div>	<div><div>23</div><div>22</div></div>	<div><div>33</div><div>33</div></div>	<div><div>34</div><div>43</div></div>
l	<div><div>33</div></div>	<div><div>32</div><div>42</div></div>	<div><div>23</div><div>32</div></div>	<div><div>34</div><div>33</div></div>	<div><div>44</div><div>44</div></div>
o	<div><div>44</div></div>	<div><div>43</div><div>53</div></div>	<div><div>33</div><div>43</div></div>	<div><div>24</div><div>42</div></div>	<div><div>45</div><div>3?</div></div>

		s	n	o	w
	$\frac{\quad}{0}$	$\frac{1}{1}$	$\frac{2}{2}$	$\frac{3}{3}$	$\frac{4}{4}$
o	$\frac{1}{1}$	$\frac{1}{2}$ $\frac{2}{1}$	$\frac{2}{2}$ $\frac{3}{2}$	$\frac{2}{3}$ $\frac{4}{2}$	$\frac{4}{3}$ $\frac{5}{3}$
s	$\frac{2}{2}$	$\frac{1}{3}$ $\frac{2}{1}$	$\frac{2}{2}$ $\frac{3}{2}$	$\frac{3}{3}$ $\frac{3}{3}$	$\frac{3}{4}$ $\frac{4}{3}$
l	$\frac{3}{3}$	$\frac{3}{4}$ $\frac{2}{2}$	$\frac{2}{3}$ $\frac{3}{2}$	$\frac{3}{3}$ $\frac{4}{3}$	$\frac{4}{4}$ $\frac{4}{4}$
o	$\frac{4}{4}$	$\frac{4}{5}$ $\frac{3}{3}$	$\frac{3}{4}$ $\frac{3}{3}$	$\frac{2}{4}$ $\frac{4}{2}$	$\frac{4}{3}$ $\frac{5}{3}$

		s	n	o	w
	<div><div></div><div>0</div></div>	<div><div>11</div></div>	<div><div>22</div></div>	<div><div>33</div></div>	<div><div>44</div></div>
o	<div><div>11</div></div>	<div><div>12</div><div>21</div></div>	<div><div>23</div><div>22</div></div>	<div><div>24</div><div>32</div></div>	<div><div>45</div><div>33</div></div>
s	<div><div>22</div></div>	<div><div>12</div><div>31</div></div>	<div><div>23</div><div>22</div></div>	<div><div>33</div><div>33</div></div>	<div><div>34</div><div>43</div></div>
l	<div><div>33</div></div>	<div><div>32</div><div>42</div></div>	<div><div>23</div><div>32</div></div>	<div><div>34</div><div>33</div></div>	<div><div>44</div><div>44</div></div>
o	<div><div>44</div></div>	<div><div>43</div><div>53</div></div>	<div><div>33</div><div>43</div></div>	<div><div>24</div><div>42</div></div>	<div><div>45</div><div>33</div></div>

			s		n		o		w	
		<hr/> 0	<hr/> 1	1	<hr/> 2	2	<hr/> 3	3	<hr/> 4	4
o		1	1	2	2	3	2	4	4	5
	<hr/>	1	2	1	2	2	3	2	3	3
s		2	1	2	2	3	3	3	3	4
	<hr/>	2	3	1	2	2	3	3	4	3
l		3	3	2	2	3	3	4	4	4
	<hr/>	3	4	2	3	2	3	3	4	4
o		4	4	3	3	3	2	4	4	5
	<hr/>	4	5	3	4	3	4	2	3	3

How do

I read out the editing operations that transform OSLO into SNOW?

		s	n	o	w
	<div><div></div><div>0</div></div>	<div><div>11</div></div>	<div><div>22</div></div>	<div><div>33</div></div>	<div><div>44</div></div>
o	<div><div>11</div></div>	<div><div>12</div><div>21</div></div>	<div><div>23</div><div>22</div></div>	<div><div>24</div><div>32</div></div>	<div><div>45</div><div>33</div></div>
s	<div><div>22</div></div>	<div><div>12</div><div>31</div></div>	<div><div>23</div><div>22</div></div>	<div><div>33</div><div>33</div></div>	<div><div>34</div><div>43</div></div>
l	<div><div>33</div></div>	<div><div>32</div><div>42</div></div>	<div><div>23</div><div>32</div></div>	<div><div>34</div><div>33</div></div>	<div><div>44</div><div>44</div></div>
o	<div><div>44</div></div>	<div><div>43</div><div>53</div></div>	<div><div>33</div><div>43</div></div>	<div><div>24</div><div>42</div></div>	<div><div>45</div><div>33</div></div>

cost	operation	input	output
1	insert	*	w

		s	n	o	w
	<div><div></div><div>0</div></div>	<div><div>11</div></div>	<div><div>22</div></div>	<div><div>33</div></div>	<div><div>44</div></div>
o	<div><div>11</div></div>	<div><div>12</div><div>21</div></div>	<div><div>23</div><div>22</div></div>	<div><div>24</div><div>32</div></div>	<div><div>45</div><div>33</div></div>
s	<div><div>22</div></div>	<div><div>12</div><div>31</div></div>	<div><div>23</div><div>22</div></div>	<div><div>33</div><div>33</div></div>	<div><div>34</div><div>43</div></div>
l	<div><div>33</div></div>	<div><div>32</div><div>42</div></div>	<div><div>23</div><div>32</div></div>	<div><div>34</div><div>33</div></div>	<div><div>44</div><div>44</div></div>
o	<div><div>44</div></div>	<div><div>43</div><div>53</div></div>	<div><div>33</div><div>43</div></div>	<div><div>24</div><div>42</div></div>	<div><div>45</div><div>33</div></div>

cost	operation	input	output
0	(copy)	o	o
1	insert	*	w

		s		n		o		w	
		<div><div></div><div>0</div></div>	<div><div>1</div><div>1</div></div>	<div><div>2</div><div>2</div></div>	<div><div>3</div><div>3</div></div>	<div><div>4</div><div>4</div></div>			
o		<div><div>1</div><div>1</div></div>	<div><div>1</div><div>2</div><div>2</div><div>1</div></div>	<div><div>2</div><div>3</div><div>2</div><div>2</div></div>	<div><div>2</div><div>4</div><div>3</div><div>2</div></div>	<div><div>4</div><div>5</div><div>3</div><div>3</div></div>			
s		<div><div>2</div><div>2</div></div>	<div><div>1</div><div>2</div><div>3</div><div>1</div></div>	<div><div>2</div><div>3</div><div>2</div><div>2</div></div>	<div><div>3</div><div>3</div><div>3</div><div>3</div></div>	<div><div>3</div><div>4</div><div>4</div><div>3</div></div>			
l		<div><div>3</div><div>3</div></div>	<div><div>3</div><div>2</div><div>4</div><div>2</div></div>	<div><div>2</div><div>3</div><div>3</div><div>2</div></div>	<div><div>3</div><div>4</div><div>3</div><div>3</div></div>	<div><div>4</div><div>4</div><div>4</div><div>4</div></div>			
o		<div><div>4</div><div>4</div></div>	<div><div>4</div><div>3</div><div>5</div><div>3</div></div>	<div><div>3</div><div>3</div><div>4</div><div>3</div></div>	<div><div>2</div><div>4</div><div>4</div><div>2</div></div>	<div><div>4</div><div>5</div><div>3</div><div>3</div></div>			

cost	operation	input	output
1	replace	l	n
0	(copy)	o	o
1	insert	•	w

		s	n	o	w
	<u> </u> 0	<u> 1 </u> 1	<u> 2 </u> 2	<u> 3 </u> 3	<u> 4 </u> 4
o	<u> 1 </u> 1	<u> 1 2 </u> 2 1	<u> 2 3 </u> 2 2	<u> 2 4 </u> 3 2	<u> 4 5 </u> 3 3
s	<u> 2 </u> 2	<u> 1 2 </u> 3 1	<u> 2 3 </u> 2 2	<u> 3 3 </u> 3 3	<u> 3 4 </u> 4 3
l	<u> 3 </u> 3	<u> 3 2 </u> 4 2	<u> 2 3 </u> 3 2	<u> 3 4 </u> 3 3	<u> 4 4 </u> 4 4
o	<u> 4 </u> 4	<u> 4 3 </u> 5 3	<u> 3 3 </u> 4 3	<u> 2 4 </u> 4 2	<u> 4 5 </u> 3 3

cost	operation	input	output
0	(copy)	s	s
1	replace	l	n
0	(copy)	o	o
1	insert	*	w

		s	n	o	w
	<div><div></div><div>0</div></div>	<div><div>1</div><div>1</div></div>	<div><div>2</div><div>2</div></div>	<div><div>3</div><div>3</div></div>	<div><div>4</div><div>4</div></div>
o	<div><div>1</div><div>1</div></div>	<div><div>1</div><div>2</div><div>2</div><div>1</div></div>	<div><div>2</div><div>3</div><div>2</div><div>2</div></div>	<div><div>2</div><div>4</div><div>3</div><div>2</div></div>	<div><div>4</div><div>5</div><div>3</div><div>3</div></div>
s	<div><div>2</div><div>2</div></div>	<div><div>1</div><div>2</div><div>3</div><div>1</div></div>	<div><div>2</div><div>3</div><div>2</div><div>2</div></div>	<div><div>3</div><div>3</div><div>3</div><div>3</div></div>	<div><div>3</div><div>4</div><div>4</div><div>3</div></div>
l	<div><div>3</div><div>3</div></div>	<div><div>3</div><div>2</div><div>4</div><div>2</div></div>	<div><div>2</div><div>3</div><div>3</div><div>2</div></div>	<div><div>3</div><div>4</div><div>3</div><div>3</div></div>	<div><div>4</div><div>4</div><div>4</div><div>4</div></div>
o	<div><div>4</div><div>4</div></div>	<div><div>4</div><div>3</div><div>5</div><div>3</div></div>	<div><div>3</div><div>3</div><div>4</div><div>3</div></div>	<div><div>2</div><div>4</div><div>4</div><div>2</div></div>	<div><div>4</div><div>5</div><div>3</div><div>3</div></div>

cost	operation	input	output
1	delete	o	*
0	(copy)	s	s
1	replace	l	n
0	(copy)	o	o
1	insert	*	w

		c	a	t	c	a	t
	<u> </u> 0	<u>1</u> 1	<u>2</u> 2	<u>3</u> 3	<u>4</u> 4	<u>5</u> 5	<u>6</u> 6
c	<u> </u> 1	<u>0</u> 2 2 0	<u>2</u> 3 1 1	<u>3</u> 4 2 2	<u>3</u> 5 3 3	<u>5</u> 6 4 4	<u>6</u> 7 5 5
a	<u> </u> 2	<u>2</u> 1 3 1	<u>0</u> 2 2 0	<u>2</u> 3 1 1	<u>3</u> 4 2 2	<u>3</u> 5 3 3	<u>5</u> 6 4 4
t	<u> </u> 3	<u>3</u> 2 4 2	<u>2</u> 1 3 1	<u>0</u> 2 2 0	<u>2</u> 3 1 1	<u>3</u> 4 2 2	<u>3</u> 5 3 3

		c	a	t	c	a	t
	$\begin{array}{ c } \hline 0 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 1 & 1 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 2 & 2 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 3 & 3 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 4 & 4 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 5 & 5 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 6 & 6 \\ \hline \end{array}$
c	$\begin{array}{ c } \hline 1 \\ \hline 1 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 0 & 2 \\ \hline 2 & 0 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 2 & 3 \\ \hline 1 & 1 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 3 & 4 \\ \hline 2 & 2 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 3 & 5 \\ \hline 3 & 3 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 5 & 6 \\ \hline 4 & 4 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 6 & 7 \\ \hline 5 & 5 \\ \hline \end{array}$
a	$\begin{array}{ c } \hline 2 \\ \hline 2 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 2 & 1 \\ \hline 3 & 1 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 0 & 2 \\ \hline 2 & 0 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 2 & 3 \\ \hline 1 & 1 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 3 & 4 \\ \hline 2 & 2 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 3 & 5 \\ \hline 3 & 3 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 5 & 6 \\ \hline 4 & 4 \\ \hline \end{array}$
t	$\begin{array}{ c } \hline 3 \\ \hline 3 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 3 & 2 \\ \hline 4 & 2 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 2 & 1 \\ \hline 3 & 1 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 0 & 2 \\ \hline 2 & 0 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 2 & 3 \\ \hline 1 & 1 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 3 & 4 \\ \hline 2 & 2 \\ \hline \end{array}$	$\begin{array}{ c c } \hline 3 & 5 \\ \hline 3 & 3 \\ \hline \end{array}$

cost	operation	input	output
1	insert	*	c
1	insert	*	a
1	insert	*	t
0	(copy)	c	c
0	(copy)	a	a
0	(copy)	t	t

		c		a		t		c		a		t		
		0	1	1	2	2	3	3	4	4	5	5	6	6
c		1 1	0 2	2 0	2 1	3 1	3 2	4 2	3 3	5 3	5 4	6 4	6 5	7 5
a		2 2	2 3	1 1	0 2	2 0	2 1	3 1	3 2	4 2	3 3	5 3	5 4	6 4
t		3 3	3 4	2 2	2 3	1 1	0 2	2 0	2 1	3 1	3 2	4 2	3 3	5 3

cost	operation	input	output
0	(copy)	c	c
1	insert	*	a
1	insert	*	t
1	insert	*	c
0	(copy)	a	a
0	(copy)	t	t

		c	a	t	c	a	t
	<div><div></div><div>0</div></div>	<div><div>11</div></div>	<div><div>22</div></div>	<div><div>33</div></div>	<div><div>44</div></div>	<div><div>55</div></div>	<div><div>66</div></div>
c	<div><div>11</div></div>	<div><div>02</div><div>20</div></div>	<div><div>23</div><div>11</div></div>	<div><div>34</div><div>22</div></div>	<div><div>35</div><div>33</div></div>	<div><div>56</div><div>44</div></div>	<div><div>67</div><div>55</div></div>
a	<div><div>22</div></div>	<div><div>21</div><div>31</div></div>	<div><div>02</div><div>20</div></div>	<div><div>23</div><div>11</div></div>	<div><div>34</div><div>22</div></div>	<div><div>35</div><div>33</div></div>	<div><div>56</div><div>44</div></div>
t	<div><div>33</div></div>	<div><div>32</div><div>42</div></div>	<div><div>21</div><div>31</div></div>	<div><div>02</div><div>20</div></div>	<div><div>23</div><div>11</div></div>	<div><div>34</div><div>22</div></div>	<div><div>35</div><div>33</div></div>

cost	operation	input	output
0	(copy)	c	c
0	(copy)	a	a
1	insert	*	t
1	insert	*	c
1	insert	*	a
0	(copy)	t	t

		c		a		t		c		a		t		
		0	1	1	2	2	3	3	4	4	5	5	6	6
c		1 1	0 2	2 1	3 1	4 2	5 3	6 4	7 5	8 6	9 7	10 8	11 9	
a		2 2	1 3	0 2	2 1	3 1	4 2	5 3	6 3	7 4	8 5	9 6	10 7	
t		3 3	2 4	1 3	0 2	2 1	3 1	4 2	5 3	6 3	7 4	8 5	9 6	

cost	operation	input	output
0	(copy)	c	c
0	(copy)	a	a
0	(copy)	t	t
1	insert	*	c
1	insert	*	a
1	insert	*	t

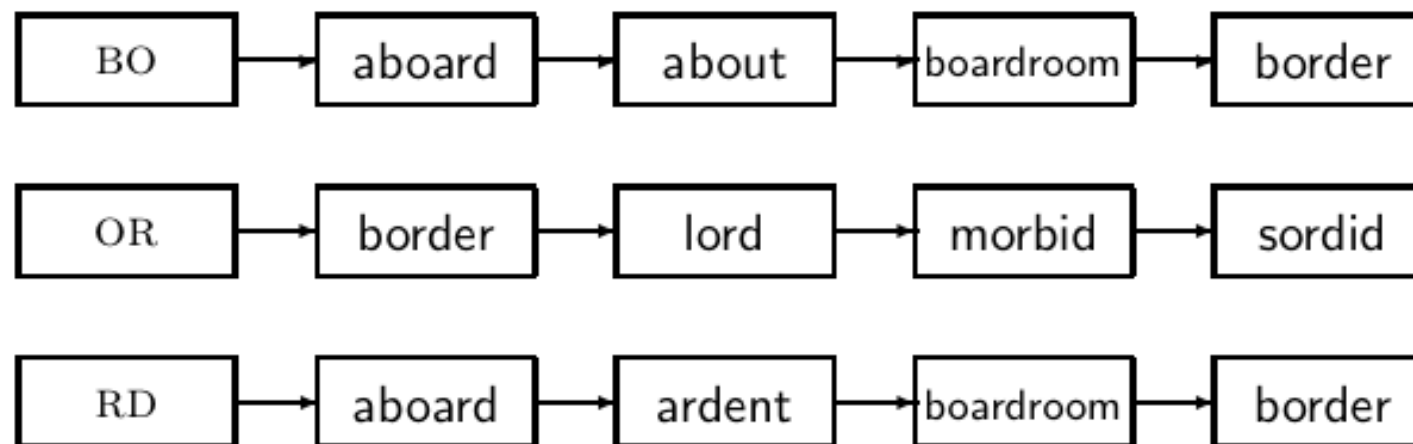
تصحیح املائی

- حالا که ما می‌توانیم فاصله ویرایشی را محاسبه کنیم:
- چطور این فاصله را برای تصحیح املائی کلمه مجزا استفاده کنیم؟
- شاخص‌های k -گرمی برای تصحیح املائی کلمه مجزا
- تصحیح املائی حساس به متن
- مسائل عمومی

شاخص‌های k -گرمی برای تصحیح املائی

- همه‌ی k -گرمی‌ها را در عبارت پرس‌وجو محاسبه و ثبت می‌کنیم.
- نمونه: شاخص دو-گرمی (bigram) کلمه با اشتباه املائی `bordroom` به صورت زیر می‌باشد.
 - Bigrams: `bo, or, rd, dr, ro, oo, om`
- از شاخص k -گرمی برای بازیابی کلمات "صحیح" که با k -گرمی‌های عبارت پرس‌وجو انطباق پیدا می‌کنند استفاده می‌شود.
- حد‌آستانه با تعداد تطبیفات k -گرمی‌ها
- به عنوان مثال: فقط عبارات واژگانی که بیشتر از ۳، k -گرم‌ها متفاوت هستند

شاخص‌های k-گرمی برای تصحیح املایی: *bordroom*



- شاخص ۲-گرمی - پست ها (بخشی از آنها) را برای سه ۲-گرمی در پرس و جوی board را نشان می دهد.
- فرض کنید می خواستیم عبارات مجموعه واژگانی را که شامل حداقل دو تا از این سه ۲-گرمی است، بازیابی کنیم.
- یک اسکن واحد از پست ها به ما اجازه می دهد تا تمامی چنین عباراتی را محاسبه کنیم.
- در این مثال، عبارات boardroom، aboard و border در نظر گرفته می شوند.

تصحیح املائی حساس به متن

- نمونه: *an asteroid that fell **form** the sky*
- چگونه می‌توان form را اینجا تصحیح کرد؟
- ایده اول: تصحیح املائی مبتنی بر ضربه (hit-based)
- در این روش بازیابی عبارات "صحیح" به مانند بازیابی هر عبارت پرس‌وجو (بخش قبل)
- در این روش اسناد اندکی بازیابی می‌شوند.
- برای flew form munich:
- flea for flew
- from for form
- munch for munich

تصحیح املایی حساس به متن

- حالا همه‌ی نتایج عبارات ممکن را به عنوان پرس‌وجوهایی با یک کلمه "ثابت" در یک زمان امتحان کنید
 - Try query “flea form munich”
 - Try query “flew from munich”
 - Try query “flew form munch”
 - The correct query “flew from munich” has the most hits.
- فرض کنید ما ۷ گزینه برای flew، ۲۰ برای form و ۳ تا برای munich
- چند عبارت تصحیح شده باید شمارش شود؟
- اگر برای یک عبارت واحد تصحیح‌های بسیاری پیدا شود، این محاسبات می‌تواند بسیار پرهزینه باشد چون ممکن است با تعداد زیادی از ترکیب‌های عبارات جایگزین مواجه شویم.

تصحیح املایی حساس به متن

- الگوریتم مبتنی بر ضربه (hit-based) در اینجا فقط معرفی شده است و کارآمدی خوبی ندارد.
- جایگزین کارآمدتر: نگاهی به "مجموعه" پرس وجوها ، نه اسناد

Soundex

- Soundex اساسا برای پیدا کردن آوایی (**phonetic**) می باشد
- نمونه: *chebyshev / tchebyscheff*
- الگوریتم:
- هر عبارت را به صورت کاهش یافته ۴ کارکتری شاخص کنید.
- با عبارات پرس و جو نیز همین کار را انجام دهید.
- هنگامی که پرس و جو نیاز به تطبیق Soundex دارد، این شاخص را جستجو کنید.

الگوریتم Soundex

۱. اولین عبارت را حذف کنید.

۲. تمامی رخدادهای حروف A, E, I, O, U, H, W, Y را به '0' تغییر دهید.

۳. حروف را به ترتیب زیر به ارقام تغییر دهید:

- B, F, P, V to 1
- C, G, J, K, Q, S, X, Z to 2
- D, T to 3
- L to 4
- M, N to 5
- R to 6

۴. به طور تکرار شونده، یکی از جفت ارقام یکسان متوالی را حذف کنید.

۵. تمامی صفرها را از رشته‌ی حاصل شده حذف کنید. رشته‌ی حاصل شده را با یک دنباله از صفر پر کنید و چهار موقعیت اول را برگردانید که شامل حرفی است که به دنبال آن سه رقم وجود دارد.

مثال Soundex برای HERMAN

- Retain H
- *ERMAN* → *ORMON*
- *ORMON* → *06505*
- *06505* → *06505*
- *06505* → *655*
- Return *H655*
- Note: *HERMANN* will generate the same code