

Часть III

Структурный подход в распознавании образов (2)

Разделы

Введение в синтаксический подход к распознаванию образов

Левенштейновская аппроксимация слова словом из регулярного языка

Введение в проблему

Основной результат

Алгоритм решения задачи РЛ

Что ещё почитать

Разделы

Введение в синтаксический подход к распознаванию образов

Левенштейновская аппроксимация слова словом из регулярного языка

Введение в проблему

Основной результат

Алгоритм решения задачи РЛ

Что ещё почитать

- └ Левенштейновская аппроксимация слова словом из регулярного языка

- └ Введение в проблему

Разделы

Введение в синтаксический подход к распознаванию образов

Левенштейновская аппроксимация слова словом из регулярного языка

- Введение в проблему

- Основной результат

- Алгоритм решения задачи РЛ

Что ещё почитать

└ Леვენштейновская аппроксимация слова словом из регулярного языка

└ Введение в проблему

Отличие последовательности \bar{x} от языка L

Пусть заданы

- ▶ L — регулярный язык, задаваемый генерирующим автоматом $\tilde{A} = \langle X, K, \varphi, P, \psi \rangle$;
- ▶ $d: X^* \times X^* \rightarrow \mathbb{R}$ — функция *отличия слова* $\bar{x}_2 \in X^*$ от слова $\bar{x}_1 \in X^*$.

Отличие не обязательно (1) симметричная функция и (2) образует метрику на множестве последовательностей.

Задача: построить алгоритм, который для каждого слова $\bar{x} \in X^*$ и каждого регулярного языка $L \subset X^*$ вычисляет *отличие слова \bar{x} от языка L* — число

$$D(\bar{x}) = \min_{\bar{y} \in L} \{ d(\bar{y}, \bar{x}) \}.$$

Далее рассматривается случай, когда функция d принадлежит классу т.н. *леვენштейновских функций*.

- └ Леვენштейновская аппроксимация слова словом из регулярного языка

- └ Введение в проблему

В.И. Леვენштейн



Владимир Иосифович Леვენштейн

(20.05.1935) — российский учёный-математик,
д.ф.-м.н, вед.н.с. ИПМ им. М. В. Келдыша.

В 1965 г. ввёл понятие расстояния
редактирования для
0-1 последовательностей.

В 2006 году получил престижную
награду США — Медаль Ричарда Хэмминга.

Пример: чтобы перевести слово КОНЬ в слово КОТ нужно совершить одно удаление и одну замену, соответственно расстояние Леვენштейна составляет 2:

КОНЬ $\xrightarrow{\text{заменяем Н на Т}}$ КОТЬ $\xrightarrow{\text{удаляем Ъ}}$ КОТ

└ Леვენштейновская аппроксимация слова словом из регулярного языка

└ Введение в проблему

Операции редактирования слов

Определим 3 операции посимвольного редактирования слов и действительные неотрицательные функции их стоимости (все символы $x, x' \in X$, все слова $\bar{x} \in X^*$):

INsert (вставка) преобразует слово $\bar{x}_1\bar{x}_2$ в слово $\bar{x}_1x\bar{x}_2$,
стоимость — $in(x)$;

CHange (замена) преобразует слово $\bar{x}_1x\bar{x}_2$ в слово $\bar{x}_1x'\bar{x}_2$,
стоимость — $ch(x, x')$;

DElete (исключение) преобразует слово $\bar{x}_1x\bar{x}_2$ в слово $\bar{x}_1\bar{x}_2$,
стоимость — $de(x)$.

Стоимость последовательности операций = сумма стоимостей операций, входящих в эту последовательность; стоимость пустой последовательности операций = 0.

└ Левенштейновская аппроксимация слова словом из регулярного языка

└ Введение в проблему

Операции редактирования слов

Определение

Стоимость $d(\bar{x}_1, \bar{x}_2)$ самой дешёвой последовательности редакторских операций, преобразующей $\bar{x}_1 \rightarrow \bar{x}_2$ называют левенштейновым отличием слова \bar{x}_2 от слова \bar{x}_1 .

Нахождение стоимости $d(\bar{x}_1, \bar{x}_2)$ — непростая задача: существует бесконечное количество последовательностей редакторских операций преобразования $\bar{x}_1 \rightarrow \bar{x}_2$.

Левенштейновы функции стоимости не обязательно задают метрику на множестве слов: они могут не удовлетворять условию треугольника и быть несимметричными.

- └ Левенштейновская аппроксимация слова словом из регулярного языка

- └ Введение в проблему

Операции редактирования слов...

Определение левенштейновского отличия провоцирует предположения о некоторых свойствах, кажущихся очевидными, но не имеющих места в действительности.

Основанные на таких предположениях алгоритмы решают только некоторые задачи из указанного класса, а в общем случае их оптимальность не гарантируется и они выдают *псевдорешение*.

Поэтому нужно точно описать подкласс задач, решаемый таким алгоритмом правильно.

Известный и часто применяемый алгоритм вычисления функций Левенштейна является ярким примером такого псевдорешения.

└ Левенштейновская аппроксимация слова словом из регулярного языка

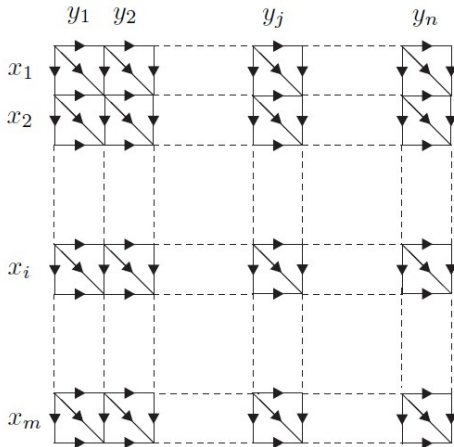
└ Введение в проблему

Наивный алгоритм вычисления $d(\bar{y}, \bar{x})$

$\bar{x} = (x_1, \dots, x_i, \dots, x_m)$, $\bar{y} = (y_1, \dots, y_j, \dots, y_n)$ — слова из X^* .

Алгоритм вычисления стоимости $d(\bar{y}, \bar{x})$ преобразования

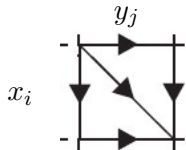
$\bar{y} \rightarrow \bar{x}$ зададим на графе $\Gamma_0(\bar{y}, \bar{x})$:



└ Левенштейновская аппроксимация слова словом из регулярного языка

└ Введение в проблему

Наивный алгоритм вычисления $d(\bar{y}, \bar{x}) \dots$



i -ой строке приписан i -й символ слова \bar{x} .

j -му столбцу приписан j -й символ слова \bar{y} .

Рёбра (стрелки) задают множество допустимых путей из левого верхнего угла

графа в правый нижний и каждому пути соответствует последовательность *редакторских операций* преобразования $\bar{y} \rightarrow \bar{x}$ слов: прохождению

- ▶ вертикальной стрелке в i -й строке соответствует вставка символа x_i ;
- ▶ горизонтальной стрелке j -м столбце соответствует исключение символа y_j ;
- ▶ диагональной стрелке в i -й и j -м столбце соответствует замена $y_j \mapsto x_i$ символов.

- └ Леვენштейновская аппроксимация слова словом из регулярного языка

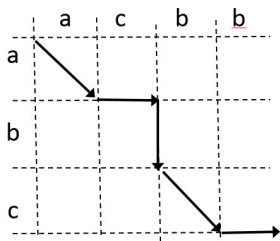
- └ Введение в проблему

Наивный алгоритм вычисления $d(\bar{y}, \bar{x})$...

Слова обрабатывается последовательно «слева направо»: на каждом шаге преобразуемое слово имеет вид $\bar{x}'\bar{y}'$, где \bar{x}' — префикс слова \bar{x} , а \bar{y}' — окончание слова \bar{y} .

Пример. Пусть $X = \{a, b, c\}$.

Преобразуем по данному алгоритму слово $\bar{y} = acbb$ в слово $\bar{x} = abc$ по указанному пути:



a	c	b	b
a	c	b	b
a	b	b	
a	b	b	b
a	b	c	b
a	b	c	

└ Лебенштейновская аппроксимация слова словом из регулярного языка

└ Введение в проблему

Наивный алгоритм вычисления $d(\bar{y}, \bar{x})\dots$

Присвоим стрелкам неотрицательные длины: длина каждой

- ▶ вертикальной стрелки в i -й строке равна стоимости $in(x_i)$;
- ▶ горизонтальной стрелки в j -м столбце — стоимости $de(y_j)$;
- ▶ наклонной стрелки в i -ой строке и j -ом столбце — стоимости $ch(y_j, x_i)$.

Длина каждого пути = стоимость последовательности операций, которую этот путь представляет.

Вроде бы: поиск оптимальной слова операций преобразования $\bar{y} \rightarrow \bar{x}$ сводится к поиску кратчайшего пути на графе от верхнего левого угла к правому нижнему?

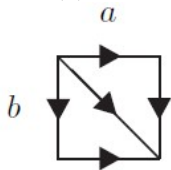
Покажем **ошибочность** данного предположения.

└ Левенштейновская аппроксимация слова словом из регулярного языка

└ Введение в проблему

Наивный алгоритм вычисления $d(\bar{y}, \bar{x}) \dots$

Пример. Пусть в алфавите $X = \{a, b, c, d\}$ даны слова $\bar{x} = (b)$, $\bar{y} = (b)$; вычислить $d(\bar{y}, \bar{x})$.



Ответ. По графу $\Gamma_0(\bar{y}, \bar{x})$: $d(\bar{y}, \bar{x}) = \min \{ ch(a, b), de(a) + in(b), in(b) + de(a) \}$.

Но имеется и много других вариантов, например:

$$a \xrightarrow{ch(a,c)} c \xrightarrow{de(c)} d \xrightarrow{in(d)} b$$

Граф Γ_0 представляет лишь очень малую часть всех возможных редакторских операций преобразования $a \rightarrow b$.

\therefore наивный алгоритм решает задачу правильно только в случае, когда он содержит самую дешевую последовательность.

Вывод: задача определения левенштейнового отличия слова от заданного языка должна быть строго формализована.

└ Левенштейновская аппроксимация слова словом из регулярного языка

└ Введение в проблему

Граф преобразований слов

Определим бесконечный направленный мультиграф Γ :

- ▶ $V(\Gamma) \leftrightarrow X^*$ — множество всех слов конечной длины над конечным алфавитом X .
- ▶ $E(\Gamma)$ составляют дуги 3 типов *in*, *de* и *ch*: две вершины, соответствующие последовательностям вида

$\bar{x}_1\bar{x}_2$ и $\bar{x}_1x\bar{x}_2$ соединяют дуги типа *in* длины $in(x)$ от первой вершины ко второй и типа *de* длины $de(x)$ от второй вершины к первой;

$\bar{x}_1y\bar{x}_2$ и $\bar{x}_1x\bar{x}_2$ соединяют дуги типа *ch* длины $ch(y,x)$ от первой вершины ко второй и длины $ch(x,y)$ от второй вершины к первой.

Левенштейново отличие — функция $d : X^* \times X^* \rightarrow \mathbb{R}_{\geq 0}$, значение которой $d(\bar{y}, \bar{x})$ для пары \bar{y}, \bar{x} есть длина кратчайшего пути в Γ от вершины \bar{y} до вершины \bar{x} .

Свойства функций Леуенштейна

Лемма (о порядке редакторских операций)

Для любых двух слов \bar{y} и \bar{x} из X^* существует кратчайший путь

- ▶ начинающийся последовательностью стрелок типа in ,
- ▶ за которой следует последовательность стрелок типа ch ,
- ▶ и завершающийся последовательностью стрелок типа de ,

причём любая из этих последовательностей может быть пустой.

Доказательство.

Пусть $\bar{y}, \bar{x}_1, \bar{x}_2, \dots, \bar{x}_n, \bar{x}$ — кратчайший путь от вершины \bar{y} к вершине \bar{x} , который не обладает указанным свойством, что обнаруживается при последовательных переходах от \bar{x}_{i-1} к \bar{x}_i и от \bar{x}_i к \bar{x}_{i+1} :

$$\bar{x}_{i-1} \longrightarrow \bar{x}_i \longrightarrow \bar{x}_{i+1}$$

└ Левенштейновская аппроксимация слова словом из регулярного языка

└ Введение в проблему

Свойства функций Левенштейна: $in - ch - de$

Это может произойти только в трёх следующих ситуациях.

① Дуга $(\bar{x}_{i-1}, \bar{x}_i)$ имеет тип ch , а дуга $(\bar{x}_i, \bar{x}_{i+1})$ — тип in .

Это значит, что слова $\bar{x}_{i-1}, \bar{x}_i, \bar{x}_{i+1}$ имеют один из двух следующих видов:

$$\left\{ \begin{array}{l} \bar{x}_{i-1} = \bar{x}'y\bar{x}''\bar{x}''', \\ \bar{x}_i = \bar{x}'x\bar{x}''\bar{x}''', \\ \bar{x}_{i+1} = \bar{x}'x\bar{x}''z\bar{x}''', \end{array} \right. \quad \text{или} \quad \left\{ \begin{array}{l} \bar{x}_{i-1} = \bar{x}'\bar{x}''y\bar{x}''', \\ \bar{x}_i = \bar{x}'\bar{x}''x\bar{x}''', \\ \bar{x}_{i+1} = \bar{x}'z\bar{x}''x\bar{x}'''. \end{array} \right.$$

В первом случае заменяем вершину \bar{x}_i на $\bar{x}'y\bar{x}''z\bar{x}'''$, а во втором — на $\bar{x}'z\bar{x}''y\bar{x}'''$.

В обоих случаях получим новый путь из вершины \bar{y} в вершину \bar{x} с той же длиной, но в этом пути первая рассматриваемая стрелка будет иметь тип in , а вторая — тип ch .

└ Левенштейновская аппроксимация слова словом из регулярного языка

└ Введение в проблему

Свойства функций Левенштейна: $in - ch - de...$

② Дуга $(\bar{x}_{i-1}, \bar{x}_i)$ имеет тип de , а дуга $(\bar{x}_i, \bar{x}_{i+1})$ — тип in .

Это значит, что слова $\bar{x}_{i-1}, \bar{x}_i, \bar{x}_{i+1}$ имеют один из двух следующих видов $(\bar{x}', \bar{x}'', \bar{x}''' \in X^*, x, y \in X)$:

$$\left\{ \begin{array}{l} \bar{x}_{i-1} = \bar{x}'y\bar{x}''\bar{x}''', \\ \bar{x}_i = \bar{x}'\bar{x}''\bar{x}''', \\ \bar{x}_{i+1} = \bar{x}'\bar{x}''x\bar{x}''', \end{array} \right. \quad \text{или} \quad \left\{ \begin{array}{l} \bar{x}_{i-1} = \bar{x}'\bar{x}''y\bar{x}''', \\ \bar{x}_i = \bar{x}'\bar{x}''\bar{x}''', \\ \bar{x}_{i+1} = \bar{x}'x\bar{x}''\bar{x}'''. \end{array} \right.$$

В первом случае вершину \bar{x}_i заменяем на $\bar{x}'y\bar{x}''x\bar{x}'''$, а во втором — на $\bar{x}'x\bar{x}''y\bar{x}'''$.

В обоих случаях получим новый путь с той же длиной, но в этом пути первая рассматриваемая стрелка будет иметь тип in , а вторая — тип de .

└ Левенштейновская аппроксимация слова словом из регулярного языка

└ Введение в проблему

Свойства функций Левенштейна: $in - ch - de...$

③ Дуга $(\bar{x}_{i-1}, \bar{x}_i)$ имеет тип de , а дуга $(\bar{x}_i, \bar{x}_{i+1})$ — тип ch .
Это возможно в двух случаях:

$$\left\{ \begin{array}{l} \bar{x}_{i-1} = \bar{x}'x\bar{x}''y\bar{x}''', \\ \bar{x}_i = \bar{x}'\bar{x}''y\bar{x}''', \\ \bar{x}_{i+1} = \bar{x}'x\bar{x}''z\bar{x}''', \end{array} \right. \quad \text{или} \quad \left\{ \begin{array}{l} \bar{x}_{i-1} = \bar{x}'y\bar{x}''x\bar{x}''', \\ \bar{x}_i = \bar{x}'y\bar{x}''\bar{x}''', \\ \bar{x}_{i+1} = \bar{x}'z\bar{x}''\bar{x}'''. \end{array} \right.$$

Заменяем в первом случае вершину \bar{x}_i на $\bar{x}'x\bar{x}''z\bar{x}'''$, а во втором — на $\bar{x}'z\bar{x}''x\bar{x}'''$.

В обоих случаях получим новый путь с той же длиной, но в этом пути первая рассматриваемая стрелка будет иметь тип ch , а вторая — тип de .

Последовательно изменяя путь от \bar{y} до \bar{x} по указанным правилам, найдем путь, в котором ни одна из указанных трех ситуаций не встретится.

└ Леვენштейновская аппроксимация слова словом из регулярного языка

└ Введение в проблему

Эквивалентное определение левенштейнова отличия

— на основе доказанной леммы.

① Определим три частных отличия d_{in} , d_{ch} и d_{de} по построенному кратчайшему пути: они равны длинам пути по стрелкам соответствующего типа, при этом если $\bar{y} = \bar{x}$, то полагаем все отличия равными 0, а если пути по стрелкам данного типа нет, то полагаем соответствующее отличие равным ∞ .

② Введём новые «длинные» стрелки в графе Γ (ранее введённые — «короткие»):

- ▶ если \bar{x} получена из \bar{y} вставкой некоторых символов, то введём в граф Γ две длинные стрелки: типа in от \bar{y} до \bar{x} длины $d_{in}(\bar{y}, \bar{x})$ и типа de от \bar{x} до \bar{y} длины $d_{de}(\bar{x}, \bar{y})$;
- ▶ если \bar{x} и \bar{y} имеют одинаковую длину, введём в граф Γ длинную стрелку типа ch от \bar{y} до \bar{x} длины $d_{ch}(\bar{y}, \bar{x})$.

└ Левенштейновская аппроксимация слова словом из регулярного языка

└ Введение в проблему

Эквивалентное определение левенштейнова отличия...

Ясно, что

- ▶ длина кратчайшего пути от \bar{y} до \bar{x} по коротким стрелкам = длине кратчайшего пути по длинным стрелкам,
- ▶ и один из этих кратчайших путей содержит не более, чем три длинные стрелки типа *in*, *ch* и *de*, идущие друг за другом в указанном порядке.

Математическое представление этого высказывания —

$$d(\bar{y}, \bar{x}) = \min_{\bar{z}_1 \in X^*} \min_{\bar{z}_2 \in X^*} \{ d_{in}(\bar{y}, \bar{z}_1) + d_{ch}(\bar{z}_1, \bar{z}_2) + d_{de}(\bar{z}_2, \bar{x}) \}$$

не годится для конструктивного вычисления $d(\bar{y}, \bar{x})$, но полезно, т.к. раскладывает понятие левенштейнова отличия на три частных понятия: его вычисление сводится к отысканию двух вспомогательных слов \bar{z}_1 и \bar{z}_2 и $|\bar{y}| \leq |\bar{z}_1| = |\bar{z}_2| \geq |\bar{x}|$.

- └ Левенштейновская аппроксимация слова словом из регулярного языка

- └ Основной результат

Разделы

Введение в синтаксический подход к распознаванию образов

Левенштейновская аппроксимация слова словом из регулярного языка

- Введение в проблему

- Основной результат

- Алгоритм решения задачи РЛ

Что ещё почитать

└ Левенштейновская аппроксимация слова словом из регулярного языка

└ Основной результат

Задача РЛ вычисления расстояния Левинштейна

Пусть для конечного алфавита X заданы:

- ▶ регулярный язык L как подмножество $\{(x_1, \dots, x_n)\}$ слов X^* , генерируемых автоматом $\tilde{A} = \langle X, K, \varphi, P, \psi \rangle$, т.е. для которых справедливо утверждение

$$\bigvee_{k_0 \in K} \dots \bigvee_{k_n \in K} \varphi(k_0) \& \bigwedge_{i=1}^n P(k_{i-1}, x_i, k_i) \& \psi(k_n);$$

- ▶ три функции in , ch и de на X^* , $(X^*)^2$ и X^* соответственно, определяющие левенштейново отличие слов $d : (X^*)^2 \rightarrow \mathbb{R}_{\geq 0}$.

Задача РЛ: создать алгоритм, который для каждого слова $\bar{x} \in X^*$ и каждой шестерки функций $(\varphi, P, \psi, in, ch, de)$ вычисляет левенштейновское отличие

$$D(\bar{x}) = \min_{\bar{y} \in L} d(\bar{y}, \bar{x}).$$

└ Леვენштейновская аппроксимация слова словом из регулярного языка

└ Основной результат

Особенности задачи РЛ

— нахождение $D(\bar{x}) = \min_{\bar{y} \in L} d(\bar{y}, \bar{x})$.

Обычно рассматривают многомерные задачи нахождения $\sup_x f(x)$ — оптимизации функций целевых функций $f(x)$, в которых

- ▶ количество переменных, может быть, большое заранее задано, в то время, как задача РЛ требуется найти слово с неизвестной заранее длиной из бесконечного множества;
- ▶ функция $f(x)$, заданными в виде явной формулы, а задаче РЛ функция $d(\bar{y}, \bar{x})$, не представлена в виде вспомогательной задачи её поиска.

Основная задача — оптимизация найденной в результате решения вспомогательной задачи функции $d(\bar{y}, \bar{x})$.

└ Левенштейновская аппроксимация слова словом из регулярного языка

└ Основной результат

Эквивалентное представление левенштейнова отличия

Теорема

Пусть X и K — два конечных множества,

$$\varphi : K \rightarrow \{0, 1\}, \quad P : K \times X \times K \rightarrow \{0, 1\}, \quad \psi : K \rightarrow \{0, 1\}$$

— три функции, определяющие регулярный язык $L \subset X^*$ как множество слов (x_1, x_2, \dots, x_n) , для которых истинен предикат

$$\bigvee_{k_0 \in K} \dots \bigvee_{k_n \in K} \varphi(k_0) \& \bigg\&_{i=1}^n P(k_{i-1}, x_i, k_i) \& \psi(k_n).$$

Пусть также $in : X \rightarrow \mathbb{R}$, $ch : X \times X \rightarrow \mathbb{R}$, $de : X \rightarrow \mathbb{R}$ — три неотрицательные функции, определяющие левенштейновы отличия $d : X^* \times X^* \rightarrow \mathbb{R}$ и $D : X^* \rightarrow \mathbb{R}$,

$$D(\bar{x}) = \min_{\bar{y} \in L} d(\bar{y}, \bar{x}).$$

└ Левенштейновская аппроксимация слова словом из регулярного языка

└ Основной результат

Эквивалентное представление левенштейнова отличия...

Теорема (продолжение)

Тогда для каждой шестерки функций $(\varphi, P, \psi, in, ch, de)$ существует такая пара функций P', ψ' , что равенство

$$D(\bar{x}) = \min_{k_0 \in K} \dots \min_{k_n \in K} \left\{ \varphi(k_0) + \sum_i^n P'(k_{i-1}, x_i, k_i) + \psi'(k_n) \right\}.$$

выполняется для любого слова $(x_1, x_2, \dots, x_n) \in X^*$.

└ Левенштейновская аппроксимация слова словом из регулярного языка

└ Основной результат

Следствия и выводы

- ▶ Вычисление числа $D(\bar{x}) = \min_{\bar{y} \in L} d(\bar{y}, \bar{x})$, вопреки всей сложности его определения, имеет сложность $O(|K|^2 n)$ что и при распознавании принадлежности слова \bar{x} обыкновенному, не штрафному, регулярному языку.
- ▶ Суммарная сложность определения $\bar{y} \in L$ и вычисление $d(\bar{y}, \bar{x})$ имеет порядок $O(|K|^2 |\bar{y}| + |\bar{y}| |\bar{x}|)$.
Т.е. вычисление $D(\bar{x})$ имеет сложность $O(|K|^2 |\bar{x}|)$, которая не зависит от $|\bar{y}|$, на которой достигается минимум, и более того, меньше, чем сложность вычисления числа $d(\bar{y}, \bar{x})$ для некоторых слов $\bar{y} \in L$.

Эти вытекающие из теоремы замечательные свойства априори неправдоподобны.

└ Левенштейновская аппроксимация слова словом из регулярного языка

└ Основной результат

Следствия и выводы...

- ▶ *Редакционным предписанием* называется последовательность действий, необходимых для получения из первой строки второй кратчайшим образом. К рассмотренным действиям (вставить, заменить, удалить) добавляют MATch) — совпадение. Найти только расстояние Левенштейна — более простая задача, чем найти ещё и редакционное предписание.
- ▶ Если к списку разрешённых операций добавить транспозицию (два соседних символа меняются местами), приходим к понятию *расстояния Дамерау — Левенштейна*. Расстояние Дамерау — Левенштейна используется —
 - ▶ при анализе текстов: Дамерау показал, что 80% ошибок при наборе текста человеком являются транспозициями;
 - ▶ в биоинформатике.

- └ Левенштейновская аппроксимация слова словом из регулярного языка

- └ Алгоритм решения задачи РЛ

Разделы

Введение в синтаксический подход к распознаванию образов

Левенштейновская аппроксимация слова словом из регулярного языка

- Введение в проблему

- Основной результат

- Алгоритм решения задачи РЛ

Что ещё почитать

└ Левенштейновская аппроксимация слова словом из регулярного языка

└ Алгоритм решения задачи РЛ

Исходная постановка задачи (повторение)

Пусть даны конечные множества X и K и функции $\varphi : K \rightarrow \{0, \infty\}$, $P : K \times X \times K \rightarrow \{0, \infty\}$, $\psi : K \rightarrow \{0, \infty\}$, задающие множество $L \in X^*$ слов, для которых существует такая слово (k_0, k_1, \dots, k_n) , что

- ① $\varphi(k_0) = 0$,
- ② $P(k_{i-1}, x_i, k_i) = 0$, $i = \overline{1, n}$,
- ③ $\psi(k_n) = 0$.

Пусть также даны также три функции $in : X \rightarrow \mathbb{R}$, $ch : X \times X \rightarrow \mathbb{R}$, $de : X \rightarrow \mathbb{R}$, которые определяют функцию $d : X \times X \rightarrow \mathbb{R}$, значения которой $d(\bar{y}, \bar{x})$ есть левенштейново отличие слова \bar{x} от слова \bar{y} .

Задача: построить алгоритм, который для любой поданной на его вход слова $\bar{x} \in X^*$ вычисляет число

$$D(\bar{x}) = \min_{\bar{y} \in L} d(\bar{y}, \bar{x}).$$

└ Левенштейнская аппроксимация слова словом из регулярного языка

└ Алгоритм решения задачи РЛ

Алгоритм решения

Алгоритм вычисления $D(\bar{x})$ состоит из двух частей.

- I. Первая часть — подготовительная и заключается в построении функций P' и ψ' , которые получают из функций P, ψ, in, ch, de .
Вычисления в этой части не зависят от входной слова \bar{x} и для данного множества L и данной левенштейновой функции выполняются только один раз.
- II. Вторая часть вычислений зависит от входной слова и состоит собственно в вычислении $D(\bar{x})$.

└ Леვენштейновская аппроксимация слова словом из регулярного языка

└ Алгоритм решения задачи РЛ

Алгоритм решения: часть I.

① Построить функцию $P_1 : K \times X \times K \rightarrow \mathbb{R}$:

$$P_1(k', y, k'') = \begin{cases} \min \{ P(k', y', k''), in(y) \}, & \text{если } k' = k'', \\ P(k', y', k''), & \text{иначе.} \end{cases}$$

Числа $P_1(k', y, k'')$ обозначают стоимость самого дешевого добавления символа y в **конец строки** при условии, что до этого автомат находился в состоянии k' , а после будет находится в состоянии k'' . При

$k' \neq k''$ сгенерируемый автоматом символ y дописывается к строке, стоимость дописывания — $P(k', y', k'')$;

$k' = k'' = k$ из этих двух возможностей: (1) символ y генерируется автоматом, стоимость $P(k, y, k)$ и (2) символ y появляется в результате редакторской операции вставки со стоимостью $in(y)$ — выбирается более дешевая.

└ Левенштейновская аппроксимация слова словом из регулярного языка

└ Алгоритм решения задачи РЛ

Алгоритм решения: часть I...

② Построить функцию $ch^* : X \times X \rightarrow \mathbb{R}$, например, следующим образом: сначала положим $ch^*(x, y) = 0$ при $x = y$ и $ch^*(x, y) = ch(x, y) = 0$, иначе; затем эти числа многократно преобразуются оператором

$$ch^*(x, y) = \min_{z \in X} \{ ch^*(x, z) + ch^*(z, y) \}.$$

Число ch^* есть стоимость самой дешевой слова замен (возможно, пустой), превращающей символ x в символ y .

└ Леვენштейновская аппроксимация слова словом из регулярного языка

└ Алгоритм решения задачи РЛ

Алгоритм решения: часть I...

③ Построить функцию $P_2 : K \times X \times K \rightarrow \mathbb{R}$:

$$P_2(k', x, k'') = \min_{y \in X} \{ P_1(k', y, k'') + ch^*(y, x) \}.$$

Число $P_2(k', x, k'')$ — это стоимость самого дешевого способа дописывания символа x в конец слова при условии, что до этого дописывания автомат находился в состоянии k' , а после дописывания — в состоянии k'' .

Добавленный символ мог быть сгенерирован автоматом или вставлен редактором.

После этого с этим символом выполняется (или не выполняется) сколь угодно длинная слово замен, результатом которой является символ x , т.е. число $P_2(k', x, k'')$ — результат оптимального выбора из довольно большого множества.

└ Леვენштейновская аппроксимация слова словом из регулярного языка

└ Алгоритм решения задачи РЛ

Алгоритм решения: часть I...

④ Построить вспомогательную функцию $q : K \times K \rightarrow \mathbb{R}$:

$$q(k', k'') = \begin{cases} 0, & \text{если } k' = k'', \\ \min_{x \in K} \{ P_2(k', x, k'') + de(x) \}, & \text{иначе.} \end{cases}$$

Число $q(k', k'')$ есть стоимость самого дешевого процесса, в результате которого автомат переходит из состояния k' в состояние k'' , причем такого, что после окончания процесса слово символов оказывается такой же, как до его начала, хотя в течение самого процесса к этой слова и дописывался какой-то символ.

Этот процесс состоит из следующих частей:

- ▶ добавление символа в конец слова, сгенерировав его автоматом или выполнив редакторскую вставку;
- ▶ слово (возможно, пустая) замен добавленного символа;
- ▶ исключение символа из слова.

└ Леуенштейновская аппроксимация слова словом из регулярного языка

└ Алгоритм решения задачи РЛ

Алгоритм решения: часть I...

⑤ Построить вспомогательную функцию $q^* : K \times K \rightarrow \mathbb{R}$:
сначала выполнить присвоение

$$q^*(k', k'') = \begin{cases} 0, & \text{если } k' = k'', \\ q(k', k''), & \text{иначе.} \end{cases}$$

Затем многократно выполнить оператор

$$q^*(k', k'') = \min_{k \in K} \{ q^*(k', k) + q^*(k, k'') \}.$$

Число $q^*(k', k'')$ подобно числу $q(k', k'')$, а различие в том, что последнее число есть стоимость самого дешевого способа перехода из состояния, при котором разрешено один и только один раз дописать символ в слово, выполнить с ним ряд замен и исключить, а первое — стоимость самого дешевого способа перехода из состояния k' в состояние k'' , при котором в слово дописывается какое угодно количество символов, включая нулевое, выполняются с ними какие угодно замены (или не выполняются их) и в конечном итоге все эти символы исключаются.

└ Левенштейновская аппроксимация слова словом из регулярного языка

└ Алгоритм решения задачи РЛ

Алгоритм решения: часть I...

⑥ Построить функцию $P : K \times X \times K \rightarrow \mathbb{R}$:

$$P'(k', x, k'') = \min_{k \in K} \{ q^*(k', k) + P_2(k, x, k'') \}.$$

Число $P'(k', x, k'')$ есть стоимость самого дешевого процесса, в результате которого автомат переходит из состояния k' в состояние k'' , а слово символов наращивается символом x . Во время этого процесса автомат генерирует какую-то слово символов, с каждым из них выполняются какие-то замены, и в конечном итоге все они исключаются.

Затем в слово дописывается один символ, либо позволяя автомату сгенерировать его, либо вставляя его редактором, с этим символом выполняется ряд замен, и в конечном итоге он остается в слова.

└ Левенштейновская аппроксимация слова словом из регулярного языка

└ Алгоритм решения задачи РЛ

Алгоритм решения: часть I...

⑦ Вычислить числа $\psi'(k)$:

$$\psi'(k) = \min_{k' \in K} \{ q^*(k, k') + \psi(k') \}.$$

Число $\psi'(k)$ есть стоимость самого дешевого процесса следующего класса: автомат, находящийся в состоянии k , генерирует какую-то слово символов и оказывается в состоянии k' , в котором он прекращает работу, а затем каждый сгенерированный символ произвольное количество раз заменяется другим символом и в конечном итоге исключается.

└ Леვენштейновская аппроксимация слова словом из регулярного языка

└ Алгоритм решения задачи РЛ

Алгоритм решения: часть II

Вычислительная процедура вычисления $D(\bar{x})$ имеет следующий вид: для поданной на вход алгоритма слова $\bar{x} = (x_0, x_1, \dots, x_n)$ вычисляются числа $f_i(k)$, $k \in K$, $i = 0, 1, \dots, n$, и число $D(\bar{x})$ по формулам

$$f_0(k_0) = \varphi(k_0), \quad k_0 \in K,$$

$$f_i(k_i) = \min_{k_{i-1} \in K} \{ f_{i-1}(k_{i-1}) + P'(k_{i-1}, x_i, k_i) \}, \quad k_i \in K, \quad i = \overline{1, n},$$

$$D(\bar{x}) = \min_{k_n \in K} \{ f_n(k_n) + \psi'(k_n) \}.$$

└ Левенштейновская аппроксимация слова словом из регулярного языка

└ Алгоритм решения задачи РЛ

Алгоритм решения: замечания

«Что касается убедительности обоснованности [описанного алгоритма] ... то нам представляется почти безнадежным сделать это с помощью одних лишь так называемых разумных соображений, выраженных словами естественного языка».

Рассмотренный алгоритм реализует метод ближайшего соседа в случае, когда L есть регулярный язык, а d — есть функция Левенштейна.

Допустим, что на основании слова $\bar{x} = (x_1, x_2, \dots, x_i, \dots, x_n)$ наблюдений нужно определить слово $\bar{k} = (k_0, k_1, \dots, k_n)$ состояний. Все описанные в лекциях алгоритмы сходны в том, что решение о всей слова \bar{k} целиком принимается на основании всей слова \bar{x} .

Известны результаты по левенштейновой аппроксимации в более общем случае контекстно-свободных языков.

Разделы

Введение в синтаксический подход к распознаванию образов

Левенштейновская аппроксимация слова словом из регулярного языка






Введение в проблему

Основной результат

Алгоритм решения задачи РЛ

Что ещё почитать

Литература I

-  *Сойфер В.А. (ред.)* Методы компьютерной обработки изображений. — М.: Физматлит, 2003.
-  *Структурное* распознавание образов. Учебно-методическое пособие для вузов. / Составитель Н. М. Новикова. — Воронеж: Издат.-полиграфич. центр Воронежского государств. унив-та, 2008.
-  *Шлезингер М., Главач В.* Десять лекций по статистическому и структурному распознаванию. — К.: Наукова думка, 2004.
-  *Фу К.* Структурные методы в распознавании образов. — М.: Мир, 1997.
-  *Хэмминг Р. В.* Цифровые фильтры. — М.: Сов. радио, 1980.