Processamento de Texto: o básico

Prof. Dr. Thiago de Paulo Faleiros

Roteiro da apresentação

- Expressão Regulares
- Palavras e Corpora
- Minimum Edit Distance

Expressão Regular são úteis em todos os lugares

- Toda tarefa de processamento de texto
 - Não é uma solução geral em PLN (por esse motivo utilizamos grande sistemas de PLN que serão vistos ao longo desta disciplina).
 - Muito útil como parte dos sistemas de PLN (ex., pré-processamento ou formatação de texto).
- Necessária para análise de dados do tipo texto
- Ferramenta amplamente utilizada na indústria e academia

Expressão Regular

- Uma linguagem formal para especificar cadeias de caracteres
- Como podemos pesquisar menções deste animal fofo no texto?
- woodchuck
- woodchucks
- Woodchuck
- Woodchucks
- Groundhog
- groundhogs



Expressão Regular: Disjunções

Letras dentro de cochetes []

Pattern	Matches
[wW]oodchuck	Woodchuck, woodchuck
[123456789]	Qualquer dígito

Intervalos usando traço [A – Z]

Pattern	Matches	
[A – Z]	Letras maiúsculasa	Thiago de Paulo
[a – z]	Letras minúsculas	estou impaciente
[0 - 9]	Um simples dígito	Capítulo 1: Expressão Regulares

Expressão Regulares: Negação em Disjução

```
Circunflexo ('^{^{\prime}} – carat) como primeiro caracter em [ ] é a negação da lista
```

- Nota: ^ significa negação só quando ele é o primeiro em []
- Caracteres especiais (., * , +, ?) perde seu significado dentro da lista []

Pattern	Matches	Exemplos
[^A - Z]	Não Letras maiúsculasa	Thiago de Paulo
[^Ss]	Nenhum 'S' ou 's'	Estou impaciente
[^.]	Não é um ponto	Capítulo 1: Expressão Regulares
[e^]	Nenhum 'e' ou '^'	tudo ^ bem

Expressão Regulares: Atalhos convenientes

Pattern	Expansion	Matches	Exemplos
$\setminus d$	[0-9]	Qualquer dígito	Fahreneit 451
$\setminus D$	[^0-9]	Qualquer não dígito	100 <mark>k</mark> m
\w	[a-zA-Z0-9_]	Qualquer alfanumé-	Thiago
		rico ou _	
$\setminus W$	[^\w]	Nenhum alfanumérico	Olhe!
		ou _	
\s	$[\r \r \]$	Espaço em branco	Olhe⊔para cima!
		(epaço,tab)	
\S	[^\s]	Não espaço em	Thiago Faleiros
		branco	

Expressão Regulares: Mais Disjunções

Marmota, em inglês pode ser *Groundhog* ou *Woodchuck* Utilize o símbolo '|' para disjunção

Pattern	Matches
groundhog woodchuck	woodchuck
yours mine	yours
a b c	= [abc]
[gG]groundhog [Ww]oodchuck	Woodchuck





Caracteres Curinga: ., ?, *, +

Pattern	Matches	Exemplo
beg.n	Qualquer caracter	begin begun beg3n beg⊔n
woodchucks?	s é opcional	woodchuck woodchucks
to*	0 ou mais do ca-	t to too tooo
	racter anterior	
to+	1 ou mais do ca-	to too tooo
	racter anterior	

Expressão Regular: Âncoras ^ \$

Pattern	Matches	
^[A–Z]	Palo Alto	
^[^A–Za–z]	'Hellow'	
\.\$	The end.	
.\$	The end?	The end!

Nota sobre Expressão Regulares em Python

- Regex e Python usam "\" para caracteres especiais. Você precisa de um "\" extra!
 - "\\d+" para procurar 1 ou mais dígitos
 - "\n" em Python significa nova linha. Precisaria de "\\n" para dois caracteres.
- Utilize raw string notation para regex:
 - r"[tT]he"
 - r"\d+" para um ou mais dígitos
 - em vez de "\\d+"

Processo iterativo para escrita de Regex

Encontre todas as instâncias da palavra "the" no texto.

- the
 - Faltou a letra maiúscula
- [tT]he
 - Pode encontrar palavras como Theology

Falso positivo e falso negativo

- O processo foi baseado na correção de dois erros:
 - Não encontrar coisas que deveriam encontrar (The) Falso Negativos
 - Encontrar strings que n\u00e3o deveriam ser encontradas (there, then, other) – Falso Positivos

Caracterizando o trabalho em PLN

Em PLN nós sempre tratamos com esses erros Reduzir a taxa de erro para uma aplicação frequentemente envolve dois esforços antagônicos:

- Aumentar a cobertura (ou recall) Minimizando os falsos negativos
- Aumentando a acurácia (ou precisão) Minimizando os falsos positivos

Expressões Regulares são importantes!

Amplamente usadas na academia e Indústria

- Parte da maioria das tarefas de processamento de texto, mesmo para grandes modelos de redes neurais.
 - Incluindo formatação e pré-processamento de texto
- Muito útil para análise de dados do tipo texto

Substituição

Substituição em Python e UNIX:

- s/regexp1/pattern/
 - s/colour/color/

Capturando grupos

- Queremos colocar ângulos (<>) em volta dos números:
 - the 35 boxes \rightarrow the <35> boxes
- Use parênteses () para capturar um padrão em indexado sequencialmente pela ocorrência (1, 2, 3, ...)
- Use \1 para referir ao conteúdo do primeiro padrão
 - $s/([0-9]+)/<\backslash 1>/$

Capturando grupos: registros múltiplos

```
/the (.*)er they (.*), the \1er we \2/
```

- Matches
 - the faster they ran, the faster we ran
- Mas não
 - the faster they ran, the faster we ate

Suponha que não queremos capturar

Parênteses tem função dupla: agrupa e captura termos

- Não capturando grupos: adicione ?: depois do parênteses: /(?:some|a few) (people|cats) like some \1/
- matches
 - some cats like some cats
- mas não
 - some cats like some some

Assertivas de Antevisão

- (?= pattern) é verdade se o padrão casa, mas sem incluir na correspondência final
- (?! pattern) verdade se o padrão não casa.
 - Como casar, no começo de uma linha, qualquer palavra simples que não comece com "Volcano":

$$/^(?!Volcano)[A-Za-z]+/$$

Quantas palavras na sentença?

"I do uh main- mainly business data processing"

Fragmentos e pausas

"Seuss's cat in the hat is different from other cats!"

- Lemma: mesmo stem, part-of-speech e radical
 - cat e cats = mesmo lema
- Wordform: a palavra flexionada completa
 - cat e cats = diferentes wordforms

Quantas palavras na sentença?

they lay back on the San Francisco grass and looked at the stars and their

- Tipo: um elemento do vocabulário
- Token: uma instância do tipo no texto

Quantas palavras?

- 15 tokens (ou 14)
- 13 tipos (ou 12 ou 11?)

Quantas palavras na sentença?

- N = número de tokens
- V = vocabulário = conjunto de tipos, |V| é o tamanho do vocabulário
- Heaps Law = Herdan's Law = $|V| = kN^{\beta}$ onde $0.67 < \beta < 0.75$.
 - O tamanho do vocabulário cresce maior do que a raíz quadrada do número de tokens

	Tokens=N	Tipos= V
Conversas telefônicas	2.4 milhôes	20 mil
Shakespeare	884 mil	31 mil
COCA	440 milhôes	2 milhôes
Google N-grams	1 trilhão	13+ milhões



Corpora

As palavras vem de algum lugar! Texto é produzido por:

- um escrito específico
- em um determinado momento
- de um determinado gênero
- em uma linguagem específica
- para uma determinada função

Corpora varia conforme dimensões

- linguagens 7097 linguagens no mundo
- Variedade Variações de linguagens
 - iont (I don't)
- Múltiplas linguagens:
 Por primeira vez veo a @username actually being hateful! it was beautiful:)
- Gênero: notícias, ficção, artigos científicos, Wikipedia
- Autor demográfico: idade, gênero, etinia,...

Tabela do Corpus

- Motivação:
 - Qual o motivo da coleta do corpus?
 - Por quem?
 - Quem financiou?
- Situação: Em que situação o texto foi escrito?
- Processo de coleta: Se é uma subamostra? Como foi amostrada? Existe pré-processamento?
 - processo de anotação, variedade linguística, demográfica, etc...

Normalização do texto

Toda tarefa de PLN requere normalização de texto:

- Tokenização (segmentação) de palavras
- Normalização do formato das palavras
- Segmentação das sentenças

Segmentação baseada em espaço

A forma mais simples de tokenização

- Para linguagens que utilização espaço em branco entre palavras
 - Arábico, Grego, latin, etc.
- Segmenta tokens entre instâncias de espaços em branco
- o comando "tr" em UNIX
 \$ tr -sc 'A-Za-z' '\n' < shakes.txt | sort | uniq -c</pre>

Problemas com a tokenização

Não pode remover cegamente as pontuações:

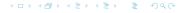
- m.p.h, Ph.D, AT&T, cap'n
- preços (\$45.55)
- datas (01/02/06)
- URLs (http://unb.br)
- hashtags (#nlproc)
- e-mails (thiagodepaulo@unb.br)

Clítico: uma palavra que não depende de outra foneticamente ou gramaticamente

• "are" (we're), "je" (j'ai), "le" (l'honneur), me (viu-me)

Expressões com múltiplas palavras são palavras únicas?

New York, Rock'n roll



Tokenização com o NLTK

```
>>> text = 'That U.S.A. poster-print costs $12.40...'
>>> pattern = r'''(?x)  # set flag to allow verbose regexps
...  ([A-Z]\.)+  # abbreviations, e.g. U.S.A.
... | \w+(-\w+)*  # words with optional internal hyphens
... | \$?\d+(\.\d+)?%?  # currency and percentages, e.g. $12.40, 82%
... | \.\.\.  # ellipsis
... | [][.,;"'?():-_']  # these are separate tokens; includes ], [
... '''
>>> nltk.regexp_tokenize(text, pattern)
['That', 'U.S.A.', 'poster-print', 'costs', '$12.40', '...']
```

Tokenização de linguagens sem espaço

Várias linguagens (como o Chinês, Japonês, Tailandês) não usam espaço para separar as palavras.

Como decidir onde ficam os limites dos tokens?

Tokenização de palavras em Chinês

- Palavras em chinês são compostas por caracteres chamados "hanzi".
- Cada um representa uma unidade semântica chamada morfema.
- Cada palavra tem em média 2.4 morfemas.
- Mas decidir o que vale como uma palavra é complexo.

Como fazer tokenização de palavras em Chinês?

```
3 words?
姚明 进入 总决赛
YaoMing reaches finals
5 words?
姚 明 进入 总  决赛
Yao Ming reaches overall finals
```

姚 明 进 入 总 决 赛 Yao Ming enter enter overall decision game

7 characters? (don't use words at all):

姚明进入总决赛"Yao Ming reaches the finals"

Tokenização

Em chinês é comum apenas tratar cada caracter como um token.

• O processo de segmentação é simples

Em outras linguagens (Como o japonês), é necessário um processo complexo de segmentação.

 O algoritmo padrão utiliza modelo de redes neurais sequenciais treinado por um técnicas de aprendizagem de máquina supervisionada.

Outras opções de Tokenização

Em vez de

- Segmentação por espaços em branco
- Segmentação por caracteres simples

Usar os próprios dados para dizer como tokenizar.

Tokenização em subpalavras

Tokens podem ser partes da palavra ou toda a palavra.

Tokenização em subpalavras

Três algoritmos comuns:

- Byte-Pair Encoding (BPE) (Sennrich et al., 2016)
- Unigram language modeling tokenization (Kudo, 2018)
- WordPiece (Schuster and Nakajima, 2012)

Todos tem duas partes:

- Apender o que s\(\tilde{a}\)o tokens a partir do corpus cru e induzir um vocabul\(\tilde{a}\)rio (um conjunto de tokens)
- Um segmentador de tokens que pega uma sentença de teste e tokeniza ela de acordo com o vocabulário criado.

Byte Pair Encoding (BPE)

Seja o vocabulário o conjunto de todos os caracteres = $\{A, B, C, D, \dots, a, b, c, d, \dots\}$. Repita:

- Escolha dois símbolos que mais ocorrem frequentemente adjacentes (digamos 'A' e 'B').
- Adicione um novo símbolo da união, formando 'AB', no vocabulário.
- Substitua cada adjacência dos símbolos 'A' e 'B' no corpus pelo novo símbolo 'AB'.

Faça esse processo até k uniões de símbolos serem realizadas.



```
function BYTE-PAIR ENCODING(strings C, number of merges k) returns vocab V
```

```
V \leftarrow all unique characters in C # initial set of tokens is characters for i=1 to k do # merge tokens til k times t_L, t_R \leftarrow Most frequent pair of adjacent tokens in C # make new token by concatenating V \leftarrow V + t_{NEW} # update the vocabulary Replace each occurrence of t_L, t_R in C with t_{NEW} # and update the corpus return V
```

- A maioria das subpalavras estão separadas por espaço
- Então é comumente adicionado um símbolo '_' de fim-de-palavra antes do espaço no corpus de treinamento.

Corpus original (e muito fascinante :P)

low low low low low lowest lowest newer newer newer newer newer newer newer newer wider wider new new

Adiciona o token de fim-de-palavra, resultando no seguinte vocabulário:

Vocabulário:

corpus

- 5 low_
- 2 lowest_
- 6 newer_
- 3 wider_
- 2 new_

Unindo e r para er corpus

- 5 low_
- 2 lowest_
- 6 newer_
- 3 wider
- 2 new_

vocabulário

_, d, e, i, l, n, o, r, s, t, w

vocabulário

_, d, e, i, l, n, o, r, s, t, w, er

```
corpus
5 low
2 lowest_
6 newer_
3 wider_
2 new_
 Unindo er _ para er_
corpus
   low_
2 lowest_
6 newer_
3 wider_
```

n e w

vocabulário

```
_, d, e, i, l, n, o, r, s, t, w, er
```

vocabulário

```
_, d, e, i, l, n, o, r, s, t, w, er, er_
```

corpus

- 5 low_
- 2 lowest_
- 6 newer_
- 3 wider_
- 2 new_

Unindo **n** e para **ne** corpus

- 5 low
- 2 lowest_
- 6 ne w er_
- 3 wider
- 2 ne w _

vocabulário

__, d, e, i, l, n, o, r, s, t, w, er, er__

vocabulário

, d, e, i, l, n, o, r, s, t, w, er, er, ne

- Unindo ne w para new
- Unindo i o para lo
- Unindo lo w para low
- Unindo new er_ para newer_
- Unindo low __ para low__

corpus

Propriedades dos tokens do BPE

- Normalmente incluem palavras frequentes
- Incluem também subpalavras frequentes
 - Que são morfemas, como -est ou -er
- Um morfema é a menor unidade portadora de significado de uma língua
 - 'unlikeliest' tem 3 morfemas 'un-', 'likely' e '-est'.

Normalização de palavras

Colocando palavras/tokens em um formato padrão

- U.S.A. ou USA
- uhhuh ou hu-huh
- Fed ou fed
- am, is, be, are

Capitalizar palavras

- Aplicações de Recuperação de Informação: reduz tudo para minúsculo
 - Os usuários tendem a usar letras minúsculas
 - Casos possíveis de palavras maiúsculas no meio de sentenças
 - General Motors
 - Fed vs fed
 - SAIL vs sail
- Para análise de sentimentos, Mineração de texto, extração de informação
 - Palavras maiúsculas podem ser importantes



Lematização

- Representa todas as palavras com os seus lemas e radicais
 - am, are, is \rightarrow be
 - ullet car, cars, cars' o car
 - ullet quiero, quieres o querer
 - correram → correr
 - \bullet maiores \rightarrow grande
- "Ele está lendo histórias de detetive."
 - \rightarrow "Ele estar ler história de detetive."

Lematização é feita por análise morfológica

Morfemas:

- Menores unidades de significação que compõem as palavras
- Stems: As unidades centrais portadoras de significado
- Afixos: Partes que aderem ao stems, geralmente com funções gramaticais

Analisador Morfológico (Morphological Parsers):

- 'cats' tem dois morfemas 'cat' e 's'
- 'amarão' tem o morfema 'amar' e características morfológicas terceira pessoa no plural e futuro do presente.

Stemming

Reduzindo os termos para os stems, removendo sufixo e flexões, deixando apenas as raízes que transmitem o sentido básico das palavras

- "Este não era o mapa que encontramos no baú de Billy Bones, mas uma cópia precisa, completa em todas as coisas nomes e alturas e sondagens - com a única exceção das cruzes vermelhas e das notas escritas."
- "Este não era o mapa que encontr no baú de Billy Bones, mas uma cóp precis, complet em tod as cois - nom e altur e sond - com a únic exceç das cruz vermelh e das not escrit."

Porter Stemmer

- O algoritmo de Porter divide o processo de stemming em várias etapas.
- Cada etapa contém regras que identificam padrões específicos nos sufixos das palavras.
- A aplicação dessas regras é sequencial, e, em cada etapa, certas condições precisam ser satisfeitas para que uma substituição ocorra.

Algumas regras:

- ATIONAL \rightarrow ATE (ex.: relational \rightarrow relate)
- ING $ightarrow \epsilon$ (ex.: motoring ightarrow motor, se o stem contém vogal)
- SSES \rightarrow SS (ex.: grasses \rightarrow grass)



Stemming em Português

- Porter Stemmer adaptado para o português
 - Snowball Stemmer para português é uma extensão do algoritmo de Porter.
 - Mais simples e rápido!
- RSLP (Removedor de Sufixos da Língua Portuguesa)
 - O RSLP segue um processo baseado em várias fases, onde diferentes tipos de sufixos são removidos
 - Considerado o algoritmo mais adequado para stemming em português devido à sua complexidade e tratamento detalhado dos sufixos da língua.

Segmentação de Sentenças

- símbolos como '!' e '?' são muito ambíguas, mas o ponto final '.' é muito mais ambiguo.
 - Pode ser os limites de uma sentença
 - Abreviação (Dr., Prof.)
 - Números (.02%, 3.14)
- Algoritmos comuns: Tokenizar primeiro depois utilize regras ou algoritmos de ML para classificar pontos como (a) parte da palavra ou (b) limites da sentença.
 - Utilizar dicionários com abreviações pode ajudar
- Segmentação de sentenças podem frequentemente ser feito por regras baseadas nas tokenizações.



O quanto similar são duas strings?

- Corretor automático
- Biologia Computacional
- Tradução de Máquina
- Extração de Informação
- Reconhecimento de fala

Edit Distance

- O número mínimo de edições entre duas strings
- O número de operação de edições
 - Inserção
 - Deleção
 - Substituição
- Necessário para transformar uma string em outra

Minimum Edit Distance

Alinhando duas Strings

Minimum Edit Distance

Alinhando duas Strings



- Se cada operação tem custo 1
 - Distância será 5
- Se substituição tem custo 2 (Distância de Levenshtein)
 - Distância será 8

Outros usos da Edit Distance em PLN

- Avaliação de Tradução de Máquina e reconhecimento de fala
 - Spokesman confirms senior government adviser was shot
 - Spokesman said the senior adviser was shot dead
- Extração de Entidades Nomeadas e Conferência de Entidades
 - IBM Inc. announced today
 - IBM profits
 - Stanford President John Hennessy announced yesterday
 - for Stamford University President John Hennesy

Como encontrar a Mínima Distância de Edição?

- Procurando um caminho (sequência de edições) da string inicial até a final
 - Estado inicial: A palavra que estamos transformando
 - Operador: inserção, deleção, substituição
 - Estabo objetivo: A palavra que estamos querendo obter
 - Custo do caminho: O que queremos minimizar: o número de edições

Edições Mínimas como busca

- O espaço de todas as sequências de edições é enorme!
 - Não podemos nos dar ao luxo de navegar ingenuamente
 - Vários caminhos distintos acabam no mesmo estado
 - Não precisamos registrar todos
 - Apenas o caminho mais curto para cada estado

Definindo a Distância mínima de Edição

- Para duas Strings
 - X de tamanho n
 - Y de tamanho m
- Nos definimos D(i, j)
 - A distância entre X[1 ... i] e Y[1 ... j]
 - *i.e.*, os primeiros i caracteres de X e os primeiros j de Y.
 - A distância de edição entre X e Y é então D(n, m).

Programação Dinâmica para Minimum Edit Distance

- **Programação Dinâmica:** Computação de uma tabela que armazena os valores de D(n, m)
- Resolvendo problemas pela combinação de soluções de subproblemas
- Bottom-up
 - Computamos D(i,j) para pequenos valores de i e j
 - Some grandes valores de D(i,j) baseado nos valores menores computados anteriormente
 - *i.e.*, compute D(i, j) para todo i(0 < i < n) e j(0 < j < m).

Definindo Distância de Edições (Levenshtein)

Inicialização

$$D(i,0) = i$$
$$D(0,j) = j$$

• Relação de Recorrência

Para cada
$$i = 1 \dots m$$

Para cada $i = 1 \dots n$

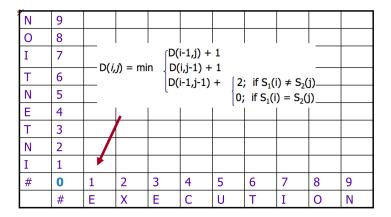
$$D(i,j) = min \left\{ \begin{array}{l} D(i-1,j)+1 \\ D(i,j-1)+1 \\ D(i-1,j-1)+2 \text{ se } X(i) \neq X(j) \\ D(i-1,j-1)+0 \text{ se } X(i) = Y(j) \end{array} \right.$$

2024/2

• Finalização D(n,m) é a distância



N	9									
0	8									
I	7									
Т	6									
N	5									
Е	4									
Т	3									
N	2									
I	1									
#	0	1	2	3	4	5	6	7	8	9
	#	Е	Χ	Е	С	U	Т	I	0	N



N	9									
0	8									
I	7									
Т	6									
N	5									
Е	4									
Т	3									
N	2									
I	1									
#	0	1	2	3	4	5	6	7	8	9
	#	Е	Χ	Е	С	U	Т	I	0	N

N	9	8	9	10	11	12	11	10	9	8
0	8	7	8	9	10	11	10	9	8	9
I	7	6	7	8	9	10	9	8	9	10
Т	6	5	6	7	8	9	8	9	10	11
N	5	4	5	6	7	8	9	10	11	10
Е	4	3	4	5	6	7	8	9	10	9
Т	3	4	5	6	7	8	7	8	9	8
N	2	3	4	5	6	7	8	7	8	7
I	1	2	3	4	5	6	7	6	7	8
#	0	1	2	3	4	5	6	7	8	9
	#	Е	Χ	Е	С	U	Т	I	0	N

Computando o alinhamento

- Distância de Edição não é suficiente
 - Frequentemente é necessário alinhar cada caracter de cada string
- Isso é feito fazendo o "backtrace"
- Toda vez que entramos em uma célula, lembre-se de onde veio
- Quando alcançado o fim,
 - Volte o caminho do canto da parte superior direita para definir o alinhamento

n	9	↓8	∠ ←↓9	∠←↓ 10	∠ ← ↓ 11	∠ ← ↓ 12	↓ 11	↓ 10	↓9	∠8	
0	8	↓ 7	∠←↓ 8	∠←↓ 9	∠←↓ 10	∠←↓ 11	↓ 10	↓9	∠ 8	← 9	
i	7	↓ 6	∠←↓ 7	∠←↓ 8	∠←↓ 9	∠←↓ 10	↓9	∠ 8	← 9	← 10	
t	6	↓ 5	∠ - ↓6	∠←↓ 7	∠←↓ 8	∠←↓ 9	∠ 8	← 9	← 10	← ↓ 11	
n	5	↓ 4	∠ 	∠←↓ 6	∠←↓ 7	∠←↓ 8	<u>/</u> ←↓9	∠ - ↓ 10	∠←↓ 11	∕↓ 10	
e	4	∠3	← 4	∠ ← 5	← 6	← 7	<i>←</i> ↓ 8	∠←↓ 9	∠←↓ 10	↓9	
t	3	∠ 4	∠ ← ↓ 5	∠←↓ 6	∠←↓ 7	∠←↓ 8	∠7	←↓ 8	∠←↓ 9	↓ 8	
n	2	∠ ← ↓ 3	∠ - ↓4	∠←↓ 5	∠←↓ 6	∠←↓ 7	∠←↓ 8	↓ 7	∠←↓ 8	∠7	
i	1	∠←↓ 2	∠ - ↓3	∠←↓ 4	∠ ← ↓ 5	∠←↓ 6	∠ ← ↓ 7	∠ 6	← 7	← 8	
#	0	1	2	3	4	5	6	7	8	9	
	#	e	X	e	С	u	t	i	0	n	

Adicionando o Backtrace

Condição Base:

$$D(i,0)=i D(0,j)=j$$

• Relação de Recorrência

Para cada
$$i = 1 \dots m$$

Para cada $j = 1 \dots n$

$$D(i,j) = min \begin{cases} D(i-1,j) + 1 \\ D(i,j-1) + 1 \\ D(i-1,j-1) + 2 \text{ se } X(i) \neq X(j) \\ D(i-1,j-1) + 0 \text{ se } X(i) = Y(j) \end{cases}$$

$$ptr(i,j) = min \begin{cases} ESQUERDA & inserção \\ DIREITA & deleção \\ DIAGONAL & substituição \end{cases}$$

Resultado do Backtrace

• As duas Strings e seu alinhamento:



Desempenho

• Tempo: O(nm)

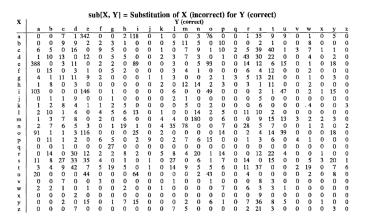
• Espaço: O(nm)

• Backtrace O(n+m)

Distância de Edição ponderada

- Qual o motivo para adicionar pesos na computação da Edit Distance?
 - Corretor gramatial: Algumas letras são mais prováveis de serem erradas do que outras
 - Biologia: Certos tipos de deleção ou inserção são mais prováveis do que outras

Matriz de Confusão para corretor gramatical



Adicionando o Backtrace

Inicialização:

$$D(0,0) = 0$$

 $D(i,0) = D(i-1,0) + del[X(i)]$ $1 < i \le n$
 $D(0,j) = D(0,j-1) + ins[Y(j)]$ $1 < j \le m$

Relação de Recorrência

$$D(i,j) = \min \left\{ \begin{array}{ll} D(i-1,j) + 1 & +del[X(i)] \\ D(i,j-1) + 1 & +ins[X(i)] \\ D(i-1,j-1) + 2 & +sub[X(i),Y(j)] \end{array} \right.$$

• Finalização D(n, m) é a distância

