



Processamento de Linguagem Natural

Experimentos em PLN e Morfologia Computacional

Marlo Souza¹

¹Universidade Federal da Bahia - Brasil

23 de julho de 2024



Sumário

*Datasets e
Corpora*

Morfologia
Computacional e
Morfoossintaxe

Datasets e Corpora

Morfologia Computacional e Morfoossintaxe



3

*Datasets e
Corpora*

Morfologia
Computacional e
Morfoossintaxe

26

Datasets e Corpora



Dataset

*Datasets e
Corpora*

Morfologia
Computacional e
Morfossintaxe

4

Um dataset é um conjunto de dados, i.e. elementos que, organizados e tratados, produzem informação. Em PLN, *datasets* são compostos de dados linguísticos, que podem ser:



Dataset

Datasets e
Corpora

Morfologia
Computacional e
Morfossintaxe

4

Um dataset é um conjunto de dados, i.e. elementos que, organizados e tratados, produzem informação. Em PLN, *datasets* são compostos de dados linguísticos, que podem ser:

- ▶ Palavras com suas classificações gramaticais



Dataset

Datasets e
Corpora

Morfologia
Computacional e
Morfo sintaxe

4

Um dataset é um conjunto de dados, i.e. elementos que, organizados e tratados, produzem informação. Em PLN, *datasets* são compostos de dados linguísticos, que podem ser:

- ▶ Palavras com suas classificações gramaticais
- ▶ Palavras ou unidades maiores, que podem ser classificadas como pessoa



Dataset

Datasets e
Corpora

Morfologia
Computacional e
Morfoossintaxe

4

Um dataset é um conjunto de dados, i.e. elementos que, organizados e tratados, produzem informação. Em PLN, *datasets* são compostos de dados linguísticos, que podem ser:

- ▶ Palavras com suas classificações gramaticais
- ▶ Palavras ou unidades maiores, que podem ser classificadas como pessoa
- ▶ Pares de frases classificadas quanto ao seu grau de similaridade



Corpus

Datasets e
Corpora

Morfologia
Computacional e
Morfoossintaxe

5

Um corpus é um *conjunto de dados linguísticos organizado*. O material que compõe um corpus (os textos) é coletado *com algum propósito* e *contém exemplos reais do uso da língua*.



Corpus

Datasets e
Corpora

Morfologia
Computacional e
Morfossintaxe

5

Um corpus é um *conjunto de dados linguísticos organizado*. O material que compõe um corpus (os textos) é coletado *com algum propósito* e *contém exemplos reais do uso da língua*.

Um corpus, para ser considerado um dataset linguístico no PLN, precisa ter tamanho o suficiente para permitir avaliar de maneira confiável o resultado de uma análise automática.



Corpus

Datasets e
Corpora

Morfologia
Computacional e
Morfofossintaxe

5

Um corpus é um *conjunto de dados linguísticos organizado*. O material que compõe um corpus (os textos) é coletado *com algum propósito e contém exemplos reais do uso da língua*.

Um corpus, para ser considerado um dataset linguístico no PLN, precisa ter tamanho o suficiente para permitir avaliar de maneira confiável o resultado de uma análise automática.

A diferença pode ser também a nível de utilidade quando um corpus anotado passa a ser usado como material para treinar sistemas, ele também pode ser visto como um dataset.



Anotação Linguística e dados rotulados

Datasets e
Corpora

Morfologia
Computacional e
Morfoossintaxe

6

A anotação é uma atividade de classificação: temos um conjunto de classes (as etiquetas) – também chamado de *tagset* – previamente definidas e critérios que guiam a classificação.

26



Anotação Linguística e dados rotulados

Datasets e
Corpora

Morfologia
Computacional e
Morfossintaxe

6

A anotação é uma atividade de classificação: temos um conjunto de classes (as etiquetas) – também chamado de *tagset* – previamente definidas e critérios que guiam a classificação.

1. TRISTE_SUBST_[0] é_V_[0] uma_ART_[0] palavra_SUBST_[0] de_PREP_[0] 6_NUM_[0] letras_SUBST_[0]
[frase neutra]
2. Um_NUM_[0] dos_PREP+ART_[0] livros_SUBST_[0] mais_ADV_[0] tristes_ADJ_[0] que_PRON-Rel_[0] já
ADV[0] !l_V_[0] [frase neutra?]
3. Sofri_V_[0] com_PREP_[0] a_ART_[0] protagonista_SUBST_[0] a_PREP_[0] cada_PRON_[0] nova_ADJ_[0]
página_SUBST_[0]; sofri_V_[0] quando_ADV_[0] o_ART_[0] livro_SUBST_[0] acabou_V_[0] [frase
positiva]
4. Nunca_ADV_[0] sofri_V_[0] tanto_ADV_[0] para_PREP_[0] ler_V_[0] um_ART_[0] livro_SUBST_[0] [frase
negativa]



Anotação Linguística e dados rotulados

Datasets e
Corpora

Morfologia
Computacional e
Morfossintaxe

6

A anotação é uma atividade de classificação: temos um conjunto de classes (as etiquetas) – também chamado de *tagset* – previamente definidas e critérios que guiam a classificação.

1. TRISTE_SUBST_[0] é_V_[0] uma_ART_[0] palavra_SUBST_[0] de_PREP_[0] 6_NUM_[0] letras_SUBST_[0]
[frase neutra]
2. Um_NUM_[0] dos_PREP+ART_[0] livros_SUBST_[0] mais_ADV_[0] tristes_ADJ_[0] que_PRON-Rel_[0] já
ADV[0] !l_V_[0] [frase neutra?]
3. Sofri_V_[0] com_PREP_[0] a_ART_[0] protagonista_SUBST_[0] a_PREP_[0] cada_PRON_[0] nova_ADJ_[0]
página_SUBST_[0]; sofri_V_[0] quando_ADV_[0] o_ART_[0] livro_SUBST_[0] acabou_V_[0] [frase
positiva]
4. Nunca_ADV_[0] sofri_V_[0] tanto_ADV_[0] para_PREP_[0] ler_V_[0] um_ART_[0] livro_SUBST_[0] [frase
negativa]

É fundamental ter em mente que a anotação é uma atividade interpretativa por parte do (grupo de) anotador(es).

26



Datasets: pra quê?

*Datasets e
Corpora*

Morfologia
Computacional e
Morfossintaxe

7

A existência de datasets linguísticos, ou corpora padrão ouro, é fundamental para uma série de tarefas e aplicações de PLN por três motivos principais:

26



Datasets: pra quê?

*Datasets e
Corpora*

Morfologia
Computacional e
Morfofossintaxe

7

A existência de datasets linguísticos, ou corpora padrão ouro, é fundamental para uma série de tarefas e aplicações de PLN por três motivos principais:

- São necessários para treinar modelos

26



Datasets: pra quê?

*Datasets e
Corpora*

Morfologia
Computacional e
Morfossintaxe

7

A existência de datasets linguísticos, ou corpora padrão ouro, é fundamental para uma série de tarefas e aplicações de PLN por três motivos principais:

- ▶ São necessários para treinar modelos
- ▶ São necessários para avaliar e comparar modelos

26



Datasets: pra quê?

*Datasets e
Corpora*

Morfologia
Computacional e
Morfofossintaxe

7

A existência de datasets linguísticos, ou corpora padrão ouro, é fundamental para uma série de tarefas e aplicações de PLN por três motivos principais:

- ▶ São necessários para treinar modelos
- ▶ São necessários para avaliar e comparar modelos
- ▶ São necessários para estabelecer uma interpretação comum do fenômeno e determinar caminhos futuros de desenvolvimento

26



Características de um bom dataset linguístico

*Datasets e
Corpora*

Morfologia
Computacional e
Morfoossintaxe

8

Uma avaliação adequada, justa, abrangente, detalhada, sistemática e transparente é um passo essencial ao se desenvolver, construir, analisar, comparar e usar tecnologias de linguagem.



Características de um bom dataset linguístico

Datasets e
Corpora

Morfologia
Computacional e
Morfo sintaxe

8

Uma avaliação adequada, justa, abrangente, detalhada, sistemática e transparente é um passo essencial ao se desenvolver, construir, analisar, comparar e usar tecnologias de linguagem.

Toda avaliação está situada em um contexto e deve ser feita de acordo com ele. Assim, apesar de determinadas práticas avaliativas serem difundidas dentro da área, não existe uma forma correta e definitiva de se avaliar um sistema, apenas questões que devem ser consideradas nas práticas de avaliação.



Avaliação de sistemas de PLN

Datasets e
Corpora

Morfologia
Computacional e
Morfoossintaxe

9

As avaliações de sistemas podem ser classificadas em relação às seguintes dimensões:

- ▶ manual ou automática
- ▶ intrínseca e extrínseca, ou direta e indireta
- ▶ formativa ou sumativa
- ▶ por usuários ou por especialistas
- ▶ componente ou ponta-a-ponta
- ▶ transparente ou às escuras
- ▶ investigativa ou experimental
- ▶ qualitativa ou quantitativa
- ▶ objetiva ou subjetiva
- ▶ supervisionada ou não-supervisionada

26



Leaderboards

Datasets e
Corpora

Morfologia
Computacional e
Morfoossintaxe

10

Uma prática popular se resume a tentar melhorar uma métrica em um leaderboard.

Model	▲ Clemscore ▲	% Played ▲	Quality Score ▲
claude-2.1	36.38	83.08	43.79
claude-2	33.71	82.12	41.05
gpt-3.5-turbo-0613	32.53	91.96	35.37
gpt-3.5-turbo-1106	30.45	77.12	39.49
openchat_3.5	19.72	57.57	34.26
sheep-duck-llama-2-70b-v1.1	17.12	40.82	41.93
Yi-34B-Chat	16.77	63.76	26.3
WizardLM-70b-v1.0	16.7	51.65	32.34
Mixtral-8x7B-Instruct-v0.1	16.53	57.68	28.66
tulu-2-dpo-70b	15.9	54.49	29.18
claude-instant-1.2	15.44	59.61	25.91
Codellama-34b-Instruct-hf	10.34	23.96	43.15

Figura: Leaderboard extraído do benchmark clembench (Madureira,2024)

26



Leaderboards

*Datasets e
Corpora*

Morfologia
Computacional e
Morfoossintaxe

11

Embora esse paradigma ajude a fomentar um progresso mensurável, ele vem acompanhado de muitas críticas: foco na métrica não no uso real da tecnologia; corrupção da métrica, etc. Ainda assim, avaliações conjuntas de tarefas podem trazer consideráveis evoluções para sua solução, especialmente para tarefas novas.

26



Elementos de uma avaliação experimental em PLN

- ▶ Hipóteses e Experimentos: avaliações são projetadas com base numa metodologia experimental na qual tentamos validar uma hipótese de forma controlada
- ▶ Material de referência: dados rotulados ou históricos são comumente necessários para permitir a avaliação da performance do sistema em relação aos objetivos e hipóteses pré-estabelecidas
- ▶ Metodologia de avaliação: é preciso garantir que o esforço de avaliação seja correto (i.e. confiável) e replicável.
- ▶ Métricas experimentais: conjunto de medidas que são utilizadas na avaliação quantitativas de sistema para representar sua performance



Partição dos dados

Datasets e
Corpora

Morfologia
Computacional e
Morfossintaxe

13

Dados que seja utilizados no processo de desenvolvimento da solução não devem ser usados para sua avaliação, uma vez que o sistema foi desenvolvido especificamente para tratar desses caos. Assim, é necessário realizar uma partição dos dados, de forma que os dados de avaliação não tenham sido observados pelo sistema ou pelos projetistas durante o processo de desenvolvimento.

26



Partição dos dados

Datasets e
Corpora

Morfologia
Computacional e
Morfossintaxe

13

Dados que seja utilizados no processo de desenvolvimento da solução não devem ser usados para sua avaliação, uma vez que o sistema foi desenvolvido especificamente para tratar desses caos. Assim, é necessário realizar uma partição dos dados, de forma que os dados de avaliação não tenham sido observados pelo sistema ou pelos projetistas durante o processo de desenvolvimento. Técnicas de partição de dados

- ▶ Partição dos dados: treino, teste, desenvolvimento
- ▶ Validação cruzada

26



Baselines e Validação estatística

Datasets e
Corpora

Morfologia
Computacional e
Morfofossintaxe

14

Comumente devemos comparar a performance do sistema proposto com uma solução inicial básica ou pré-existente para garantir que a performance do sistema justifica sua aplicação. A tais sistemas chamamos de *baseline*.

26



Baselines e Validação estatística

Datasets e
Corpora

Morfologia
Computacional e
Morfo sintaxe

14

Comumente devemos comparar a performance do sistema proposto com uma solução inicial básica ou pré-existente para garantir que a performance do sistema justifica sua aplicação. A tais sistemas chamamos de *baseline*. A elaboração de hipóteses experimentais usualmente envolvem a comparação da performance do sistema em relação aos *baselines* escolhidos e a validação das mesmas se dá através de testes estatísticos.

26



Análise de Erros

Datasets e
Corpora

Morfologia
Computacional e
Morfoossintaxe

15

A simples validação estatística da hipótese não nos fornece base para compreender o funcionamento do sistema proposto - do contrário recaímos num reducionismo (como discutido anteriormente sobre os leaderboards).

26



Análise de Erros

Datasets e
Corpora

Morfologia
Computacional e
Morfoossintaxe

15

A simples validação estatística da hipótese não nos fornece base para compreender o funcionamento do sistema proposto - do contrário recaímos num reducionismo (como discutido anteriormente sobre os leaderboards). É preciso então investigar de formas diferentes, de preferência de forma qualitativa, a performance do sistema, levando em consideração aspectos não quantificados.

26



Análise de Erros

Datasets e
Corpora

Morfologia
Computacional e
Morfoossintaxe

15

A simples validação estatística da hipótese não nos fornece base para compreender o funcionamento do sistema proposto - do contrário recaímos num reducionismo (como discutido anteriormente sobre os leaderboards).

É preciso então investigar de formas diferentes, de preferência de forma qualitativa, a performance do sistema, levando em consideração aspectos não quantificados.

Uma tal forma é a análise qualitativa de erros, que nos permite entender de forma minuciosa o comportamento do sistema e quais suas limitações de aplicação. Ela consiste em analisar manualmente as saídas do sistema e os erros cometidos pelo mesmo a fim de levantar hipóteses para melhorias futuras.

26



*Datasets e
Corpora*

Morfologia
Computacional e
Morfossintaxe

16

Morfologia Computacional e Morfossintaxe

26

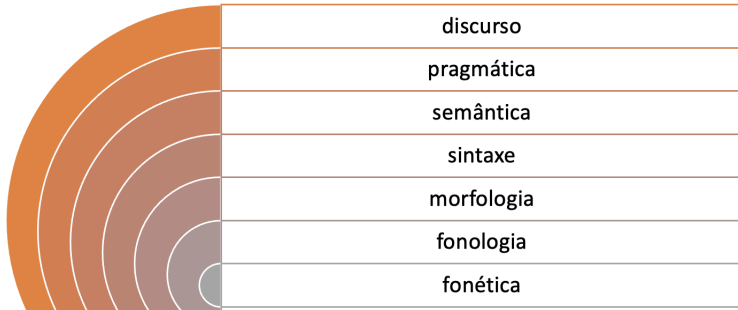


Níveis de Análise Linguística

Datasets e
Corpora

Morfologia
Computacional e
Morfofossintaxe

17



26



Morfologia

*Datasets e
Corpora*

Morfologia
Computacional e
Morfossintaxe

18

A morfologia é a área de linguística que estuda a estrutura e o processo de formação das palavras de uma língua.

26



Morfologia

Datasets e
Corpora

Morfologia
Computacional e
Morfossintaxe

18

A morfologia é a área de linguística que estuda a estrutura e o processo de formação das palavras de uma língua.

casa $\xrightarrow{\text{diminutivo}}$ casinha

26



Morfologia

Datasets e
Corpora

Morfologia
Computacional e
Morfossintaxe

18

A morfologia é a área de linguística que estuda a estrutura e o processo de formação das palavras de uma língua.

casa $\xrightarrow{\text{diminutivo}}$ casinha

A menor unidade dotada de significado é chamada de **morfema** (e.g., inha). Em uma explicação simples, podemos dizer que os morfemas são os pedacinhos que se juntam para formar as palavras

26



Palavras e tokens

*Datasets e
Corpora*

Morfologia
Computacional e
Morfossintaxe

19

O conceito de palavra em linguística é um conceito complexo.

26



Palavras e tokens

Datasets e
Corpora

Morfologia
Computacional e
Morfoossintaxe

19

O conceito de palavra em linguística é um conceito complexo. Segundo Cabré (1999, p. 20), as palavras são as unidades de referência da realidade empregadas pelos falantes.

26



Palavras e tokens

Datasets e
Corpora

Morfologia
Computacional e
Morfoossintaxe

19

O conceito de palavra em linguística é um conceito complexo. Segundo Cabré (1999, p. 20), as palavras são as unidades de referência da realidade empregadas pelos falantes.

O léxico constitui o léxico consiste no conjunto das palavras de uma língua e dos padrões que possibilitam a criatividade do falante.

26



Palavras e tokens

Datasets e
Corpora

Morfologia
Computacional e
Morfofossintaxe

19

O conceito de palavra em linguística é um conceito complexo. Segundo Cabré (1999, p. 20), as palavras são as unidades de referência da realidade empregadas pelos falantes.

O léxico constitui o léxico consiste no conjunto das palavras de uma língua e dos padrões que possibilitam a criatividade do falante.

Em PLN, empregamos o conceito de *token*, ou palavra computacional, que significa qualquer sequência de caracteres à qual se atribui um valor. Exemplo: a sentença

“Eu sempre viajo para Campinas, para Salvador e para Belém”

possui os tokens (“eu”, “sempre”, “viajo”, “para”, “Campinas”, “,”, “Salvador”, “e”, “Belém” e “.”)

26



A Morfologia Computacional

*Datasets e
Corpora*

Morfologia
Computacional e
Morfoossintaxe

20

A morfologia computacional é a área do PLN que estuda técnicas para a análise de propriedades morfológicas do texto ou da língua.

26



A Morfologia Computacional

Datasets e
Corpora

Morfologia
Computacional e
Morfofossintaxe

20

A morfologia computacional é a área do PLN que estuda técnicas para a análise de propriedades morfológicas do texto ou da língua.

São tarefas comumente estudadas no processamento morfológico:

- ▶ a sentenciação
- ▶ a tokenização
- ▶ a tokenização em subpalavras (descoberta de morfemas ou de morfotokens)
- ▶ a normalização
- ▶ a análise morfofossintática
- ▶ a análise de atributos morfológicos

26



A Sentenciação

*Datasets e
Corpora*

*Morfologia
Computacional e
Morfofossintaxe*

21

A sentenciação (ou sentenciamento) é o processo de segmentação do texto em sentenças, ou seja, é o processo de identificação de unidades textuais de processamento onde se definem os limites de cada sentença.

26



A Sentenciação

*Datasets e
Corpora*

*Morfologia
Computacional e
Morfofossintaxe*

21

A sentenciação (ou sentenciamento) é o processo de segmentação do texto em sentenças, ou seja, é o processo de identificação de unidades textuais de processamento onde se definem os limites de cada sentença.

Estratégias comuns empregadas para a sentenciação:

- ▶ Abordagens baseadas em regras
- ▶ Abordagens baseadas em aprendizado de máquina supervisionado
- ▶ Abordagens baseadas em aprendizado de máquina não supervisionado,

26



A Sentenciação

Datasets e
Corpora

Morfologia
Computacional e
Morfoossintaxe

21

A sentenciação (ou sentenciamento) é o processo de segmentação do texto em sentenças, ou seja, é o processo de identificação de unidades textuais de processamento onde se definem os limites de cada sentença.

Estratégias comuns empregadas para a sentenciação:

- ▶ Abordagens baseadas em regras
- ▶ Abordagens baseadas em aprendizado de máquina supervisionado
- ▶ Abordagens baseadas em aprendizado de máquina não supervisionado,

Fui à clínica do Dr. Nilo.

26



A Tokenização

*Datasets e
Corpora*

Morfologia
Computacional e
Morfoossintaxe

22

A separação em unidades linguísticas mínimas é denominada tokenização. No caso do português é feita partindo da separação das palavras através de delimitadores.

26



A Tokenização

Datasets e
Corpora

Morfologia
Computacional e
Morfossintaxe

22

A separação em unidades linguísticas mínimas é denominada tokenização. No caso do português é feita partindo da separação das palavras através de delimitadores. Outra tarefa frequente da tokenização é a separação de palavras contraídas, por exemplo, a palavra “da” é separada em dois tokens: “de”+“a”

26



A Tokenização de subpalavras

*Datasets e
Corpora*

Morfologia
Computacional e
Morfofossintaxe

23

Uma tarefa relacionada à tokenização, denominada tokenização de subpalavras, busca identificar fragmentos de palavras que possuem conteúdo relevante no processamento morfológico de uma palavra. Apesar de ter se popularizado recentemente com sua utilização em modelos de linguagem, é estudado no contexto de morfologia computacional não supervisionada para identificação de morfemas.

26



A Normalização

*Datasets e
Corpora*

Morfologia
Computacional e
Morfoossintaxe

24

A normalização é a tarefa que converte as palavras para alguma forma padrão. São exemplos de normalização: conversão de versões abreviadas de palavras, conversão para caracteres minúsculos , lematização (e.g., estabelecendo que “somos” é uma conjugação do verbo “ser”) e radicalização (e.g., estabelecendo que “retrabalho” tem o radical “trabalho”).

26



A Análise Morfossintática

Datasets e
Corpora

Morfologia
Computacional e
Morfossintaxe

25

A análise morfossintática, ou PoS tagging, é uma técnica fundamental na área de PLN que envolve a atribuição de etiquetas gramaticais a cada palavra em um texto, com base na sua classe gramatical e em suas características morfológicas.

		DET	NOUN	ADJ	VERB	ADV	PUNCT		
<p>	<s>	O	gato	preto	correu	rapidamente	.	</s>	</p>

Figura: Exemplo de anotação de PoS (Finatto et al, 2024)

26



Anotação de atributos morfológicos

Datasets e
Corpora

Morfologia
Computacional e
Morfossintaxe

26

A anotação de atributos morfológicos, também conhecida como atribuição de features morfológicas (ou somente feats), é uma importante tarefa que envolve a marcação ou identificação de informações específicas sobre as características gramaticais e morfológicas de palavras em um texto. Esses atributos morfológicos incluem características como número, gênero, modo, tempo, pessoa e outras informações semelhantes.

		Gender=Masc Number=Sing	Gender=Masc Number=Sing	Gender=Masc Number=Sing	Mood=Ind Number=Sing Person=3 Tense=Past		
<p>	<s>	O	GATO	PRETO	CORRER	RAPIDAMENTE	.
		DET	NOUN	ADJ	VERB	ADV	PUNCT
		O	gato	preto	correu	rapidamente	.
							</s> </p>

Figura: Exemplo de anotação de features morfológicas (Finatto et al, 2024)

26