

# Formação Cientista de Dados

---

NLP

# Corpus

- Conjunto de dados (texto não estruturado) em linguagem natural

1;Na era da informação e conhecimento, analisar dados não é uma atividade qualquer, para empresas e governos, é uma questão de sobrevivência. Num mundo globalizado e cada vez mais dependente de dados, como se extrai informação e conhecimento de dados? Implementando projetos de Big Data. Um projeto, segundo definição do PMBOK (2013) é um esforço temporário para produzir um resultado único, com começo, meio e fim, com recursos, orçamento, cronograma e equipe dedicados para a produção de um resultado específico. Um projeto de Big Data é um projeto de informação e conhecimento relevantes para a tomada de decisão.

3;Tudo isso para dizer que, de uma maneira geral, um projeto de Big Data nada mais é do que um projeto qualquer, com início e fim, escopo, orçamento, cronograma entre outros aspectos.

4;Nesta obra vamos usar alguns termos que precisamos definir antes, de forma a não haver um entendimento ambíguo do seu significado:

5;Projeto de análise de dados ou de Big Data: Nesta obra vamos usar projetos de análise de dados ou de Big Data como uma definição genérica para qualquer tipo de projeto que envolva a análise de dados.

6;Dados de origem: analisar dados requer coletar dados de algum lugar. Vamos usar dados de origem para nos referirmos a dados ou conjuntos de dados que são utilizados para produzir um resultado específico.

7;Dados de Staging: Muitos processos de análises de dados possuem uma etapa intermediária, entre os dados de origem e o resultado dos dados já consolidados. Esta etapa intermediária é usada para armazenar dados antes de serem processados.

8;Dados de destino: aqui estaremos sempre nos referindo ao resultado: o cubo, a previsão de uma análise preditiva, o relatório, o arquivo gerado, entre outros.

9;O PMBOK nos ensina que nem todos os seus processos são obrigatórios. Tampouco temos que aplicar os processos sempre com a mesma intensidade. Quais processos usar e em qual ordem depende do projeto.

10;Se você está lendo esta obra provavelmente já ouviu e leu muito sobre Big Data, e não será nosso objetivo aqui apresentar este conceito. Porém, muitos entendimentos são apresentados e precisamos esclarecer alguns pontos.

11;Falamos em seção anterior mas vamos repetir: cabe ao gerente de projeto decidir quais processos e com qual intensidade eles serão aplicados, conforme a complexidade do projeto.

12;Desde a pré-história o homem analisa dados. A análise de dados eletrônicos é um evento mais recente, com pouco mais de 70 anos. A capacidade de usar um computador para analisar dados é uma característica da era da informação.

13;Mas a análise de dados só começou a tomar força na década de 90, foi quando os grandes armazéns de dados, ou data warehouse começaram a se tornar mais comuns para apoiar a análise de dados.

14;Mas o que diferencia um projeto de análise de dados tradicional, como os popularizados nos anos 90, de um projeto de Big Data? Primeiro, os “Vs” que falamos na seção anterior.

15;Velocidade: a velocidade diz respeito não somente a da produção do dado em si, mas a velocidade do processamento e produção de informação e conhecimento, visto que o valor da informação é dinâmico e muda rapidamente.

16;Volume: projetos tradicionais eram construídos em armazéns de dados contendo por volta de terabytes de dados. Projetos de Big Data são de petabytes ou mais.

17;Variedade: projetos tradicionais carregavam dados estruturados de sistemas de operação tradicionais, em modelos relacionais, hierárquicos ou de rede. Projetos de Big Data podem carregar dados não estruturados.

18;Mas além dos “Vs” existem outras diferenças significativas que devem ser consideradas em projetos de B&A. Vamos ver alguns:

19;Primeiro, do ponto de vista de arquitetura: projetos tradicionais tem uma arquitetura centralizada, enquanto Big Data é distribuída. Se o projeto tradicional precisar crescer, a arquitetura precisa ser redesenhada.

20;Em projetos tradicionais, existe uma grande preocupação em só carregar dados em que, a partir de uma análise prévia, se vê valor. Estes dados são tratados e carregados em um formato específico.

21;Outra forma que podemos olhar uma solução de Big Data é sob sua arquitetura básica. Neste contexto, temos quatro elementos: fontes de dados, carga, armazenamento, análise e visualização.

22;Quanto as fontes de dados, podemos ter nos dois casos os mesmos elementos: dados estruturados ou não estruturados. Porém, projetos de Big Data tem mais presentes fontes de dados não estruturados.

# Annotations

- Processo mais importante em NLP
- Coloca anotações no texto, como flexões, dependências etc.

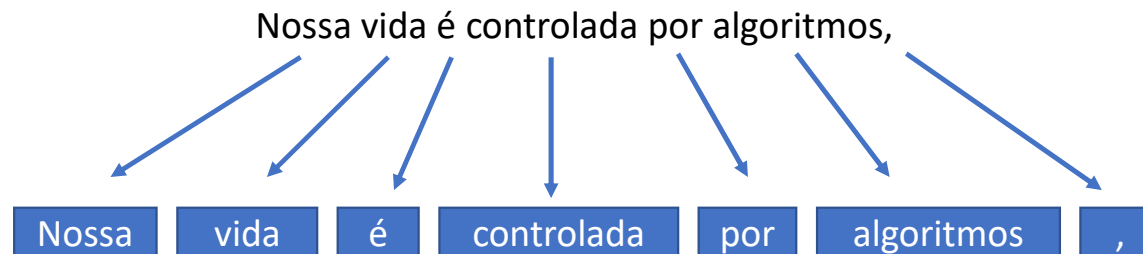
```
1\n# newpar\n# sent_id = 1\n# text = Nossa vida é controlada por algoritmos, disse artista e professor de artes digitais de uma universidade
```

americana\|n1\|tNossa\|t\_ \|tDET\|tDET\|t\_ \|t2\|tdet:poss\|t\_ \|t\_ \|n2\|tvida\|t\_ \|tNOUN\|tNOUN\|t\_ \|t4\|tnsubj:pass\|t\_ \|t\_ \|n3\|té\|t\_ \|tAUX\|tAUX\|t\_ \|t4\|taux:pass\|t\_ \|t\_ \|n4\|tcontrolada\|t\_ \|tVERB\|tVERB\|t\_ \|t8\|tccomp\|t\_ \|t\_ \|n5\|tpor\|t\_ \|tADP\|tADP\|t\_ \|t6\|tcase\|t\_ \|t\_ \|n6\|talgoritmos\|t\_ \|tNOUN\|tNOUN\|t\_ \|t4\|tnmod\|t\_ \|tSpaceAfter=No\|n7\|t\_ \|t\_ \|tPUNCT\|t\_ \|t\_ \|t8\|tpunct\|t\_ \|t\_ \|n8\|tdisse\|t\_ \|tVERB\|tVERB\|t\_ \|t0\|troot\|t\_ \|t\_ \|n9\|tartista\|t\_ \|tNOUN\|tNOUN\|t\_ \|t8\|tnsubj\|t\_ \|t\_ \|n10\|te\|t\_ \|tCCONJ\|tCONJ\|t\_ \|t11\|tc\|t\_ \|t\_ \|n11\|tprofessor\|t\_ \|tNOUN\|tNOUN\|t\_ \|t9\|tconj\|t\_ \|t\_ \|n12\|tde\|t\_ \|tADP\|tADP\|t\_ \|t13\|tcase\|t\_ \|t\_ \|n13\|tartes\|t\_ \|tNOUN\|tNOUN\|t\_ \|t17\|tnmod\|t\_ \|t\_ \|n14\|tdigitais\|t\_ \|tADJ\|tADJ\|t\_ \|t13\|tamod\|t\_ \|t\_ \|n15\|tde\|t\_ \|tADP\|tADP\|t\_ \|t17\|tcase\|t\_ \|t\_ \|n16\|tuma\|t\_ \|tDET\|tDET\|t\_ \|t17\|tdet\|t\_ \|t\_ \|n17\|tuniversidade\|t\_ \|tNOUN\|tNOUN\|t\_ \|t9\|tnmod\|t\_ \|t\_ \|n18\|tamericana\|t\_ \|tADJ\|tADJ\|t\_ \|t17\|tamod\|t\_ \|tSpacesAfter=\\|n\|n\|n"



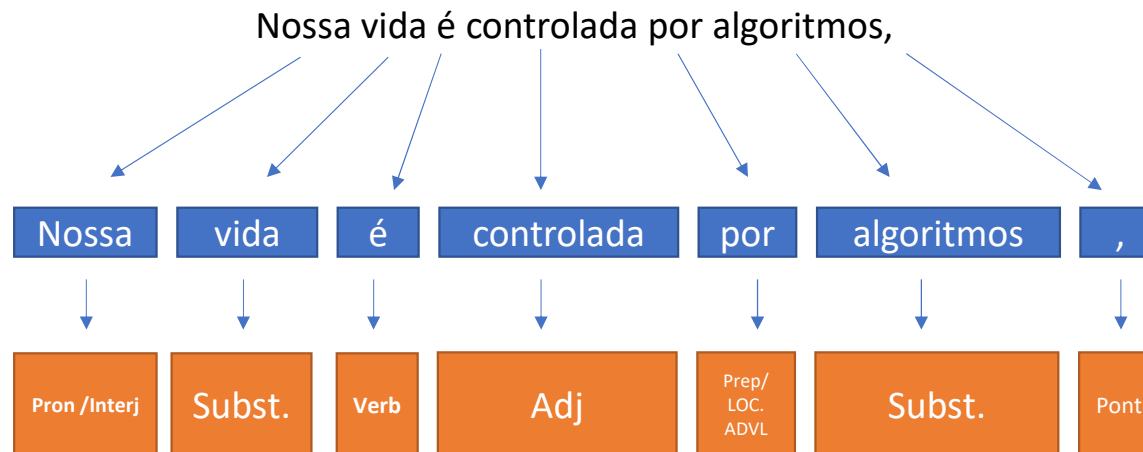
# Tokenization

- Processo de separar a sentença em suas partes: palavras, pontos, símbolos etc.



# Parts-of-Speech Tagging (POS)

Adiciona tags a cada token, como por exemplo, se é verbo, substantivo, adjetivo etc.

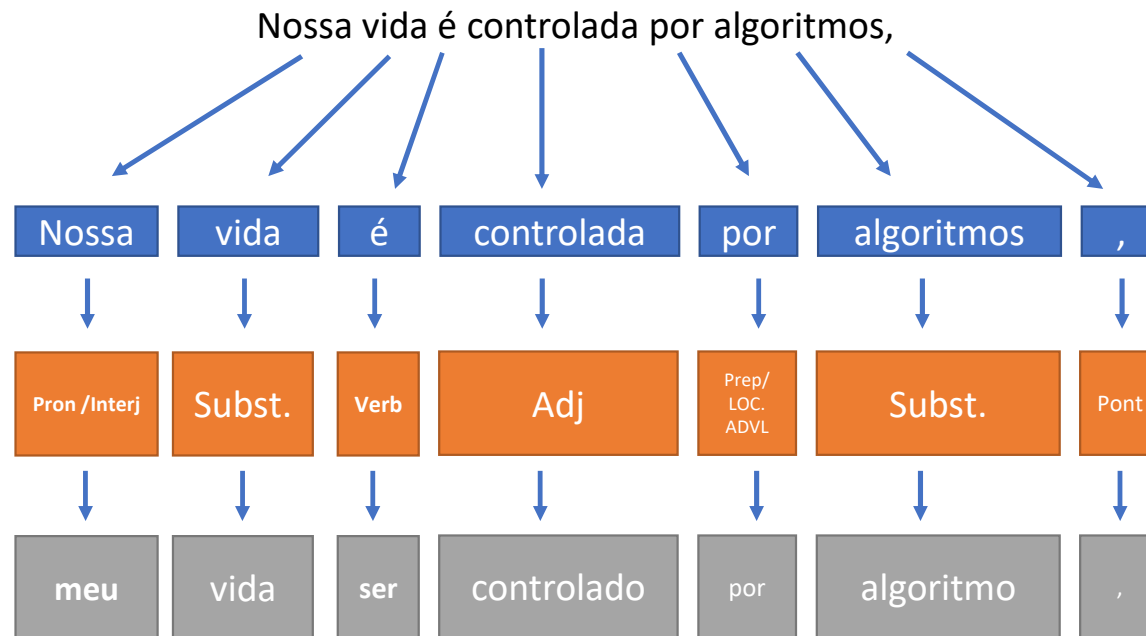


# POS Tagging

ABREVIACÃO	SIGNIFICADO	EXEMPLO
PROPN	Nome Próprio	José, Maria
VERB	verbo	Andar, Dirigir
ADP	Adposição	De, em, durante
DET	Determinante	A, Aquela, muitas
NOUN	Substantivo	Casa, carro
PUNCT	Pontuação	,,;
ADJ	Adjetivo	Infeliz, apavorado, brasileiro
CCONJ	Conjunções Coordenativas	E, nem, mas, entretanto
SCONJ	Conjunções Subordinativas	Embora, mesmo que, uma vez que
AUX	Verbos Auxiliares	Ser, estar, ter
PART	Funções de Partícula	Se, que
PRON	Pronomes	Meu, minha, meus, os quais
NUM	Números	10, vinte
ADV	Advérbios	Tarde, aqui, mal
SYM	Sinais Gráficos	~, ", '
INTJ	Interjeição	Ah, droga, psiu, hum
X	Outros	

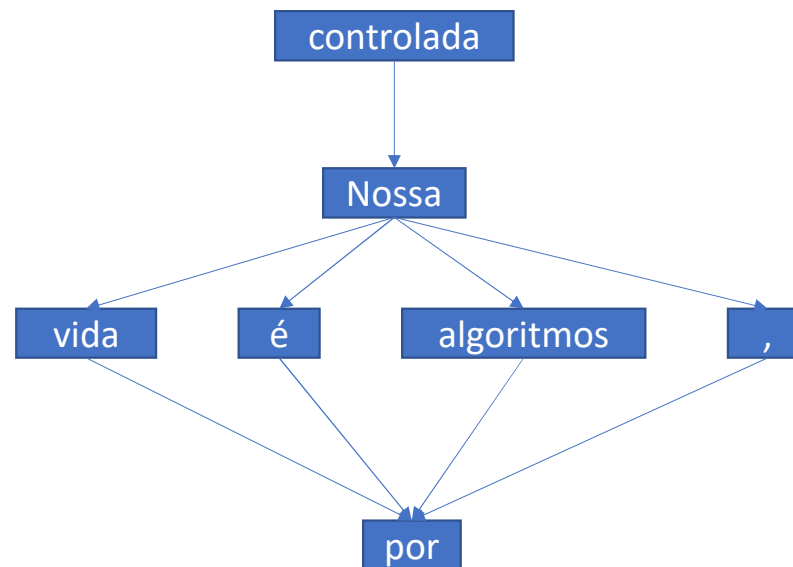
# Lemmatizing (Lemma)

- Traz a palavra na sua flexão, de modo que possam ser analisadas juntas

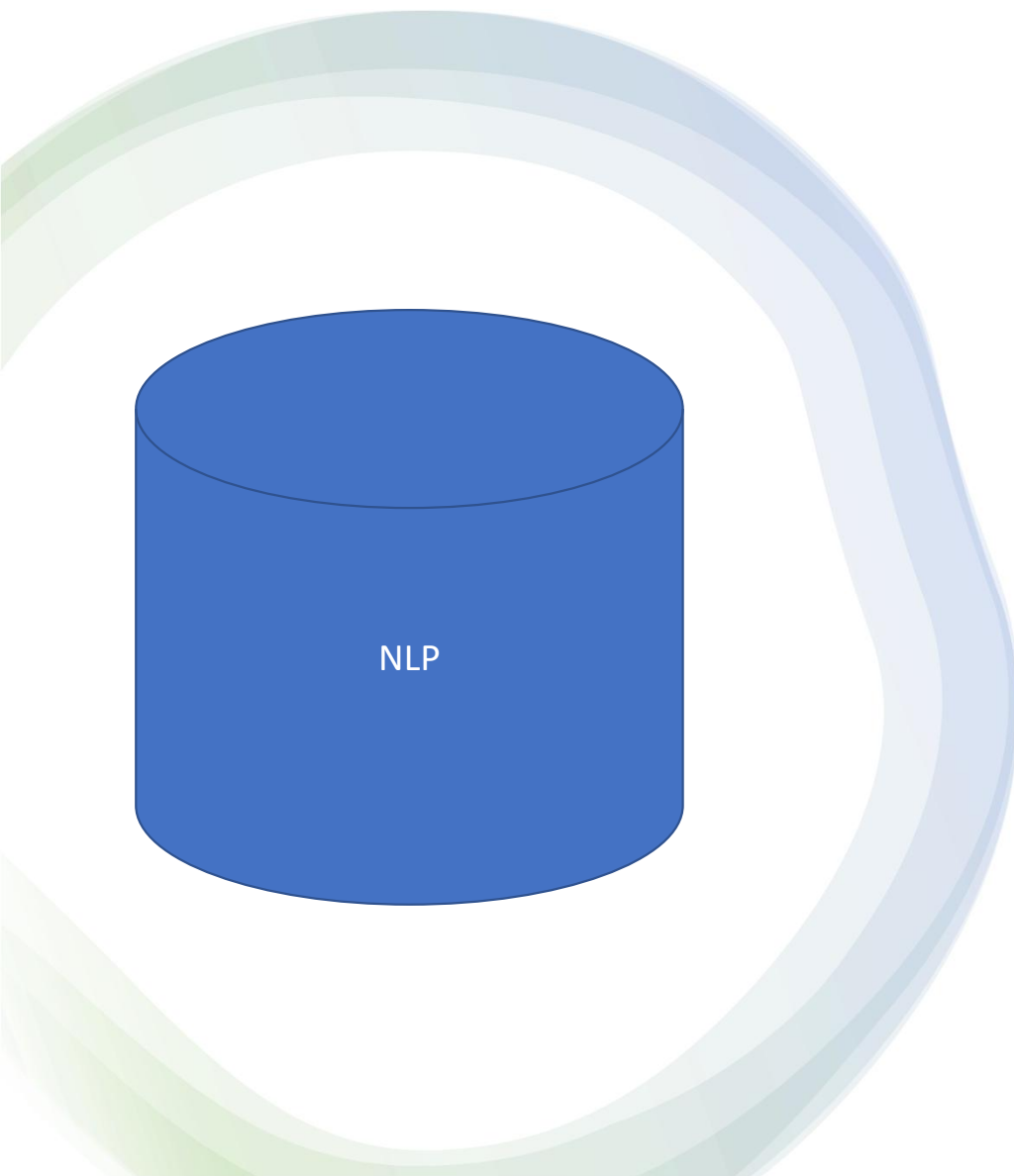


# Dependency Parsing

Encontra relação entre palavras “pais” e “filhos”



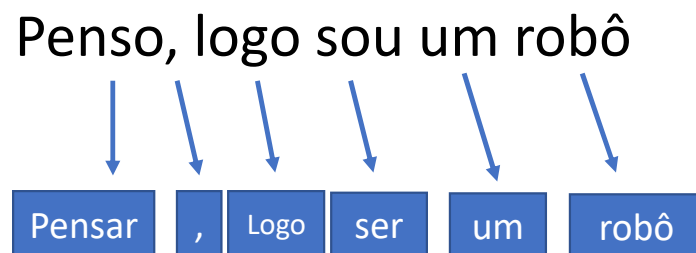




## Modelo

- Análise
  - Verbo? Substantivo? Quais são as flexões?  
Quais as dependências?
- Um modelo é um banco de dados linguístico
- Específico de cada idioma
- Maioria das plataformas de NLP tem seus próprios modelos (ou usam de terceiros)
- Você pode criar o seu!

Pergunta: defina o processo abaixo:



- ☐ Tokenization
- ☐ Parts-of-Speech Tagging
- ☒ Lemmatizing
- ☐ Dependency Parsing

