



Automated Text Generation & Data-Augmentation for Medicine, Finance, Law, and E-Commerce

September 7. 2022

Christian Kasim Loan
Lead Data Scientist, ML and Big Data Engineer
christian@johnsnowlabs.com



Introducing Spark NLP



Total downloads 33,767,472

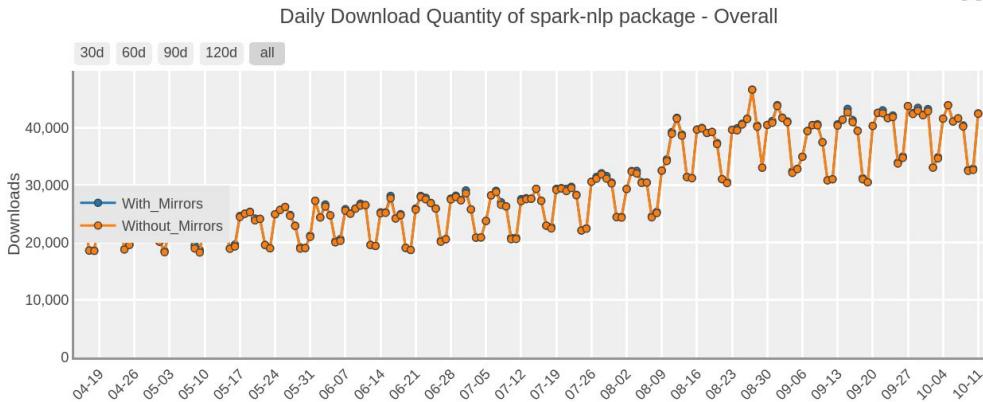
Total downloads - 30 days 2,356,291

Total downloads - 7 days 562,523

downloads 33M

downloads/month 2M

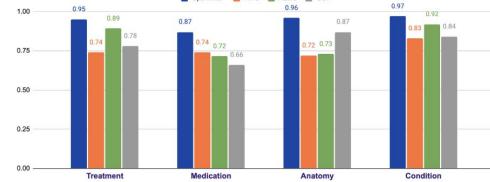
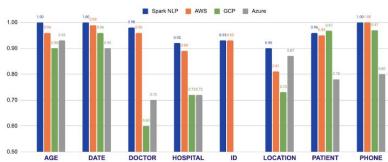
downloads/week 562k



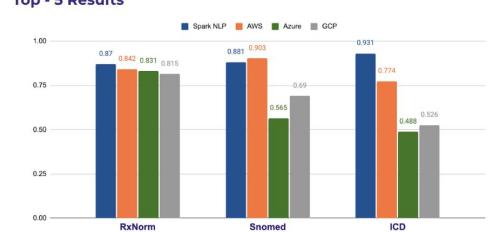
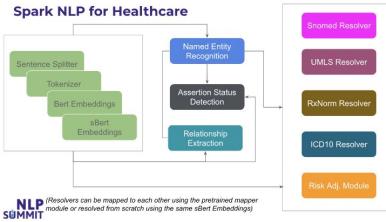
- Spark NLP is an open-source natural language processing library, built on top of Apache Spark and Spark ML. (initial release: Oct 2017)
 - A single unified solution for all your NLP needs
 - Take advantage of transfer learning and implementing the latest and greatest SOTA algorithms and models in NLP research
 - The most widely used NLP library in industry (5 yrs in a row)
 - Delivering a mission-critical, enterprise grade NLP library (used by multiple Fortune 500)
 - The most scalable, accurate and fastest library in NLP history

Most Accurate in the Industry

De-Identification Benchmarks (en)

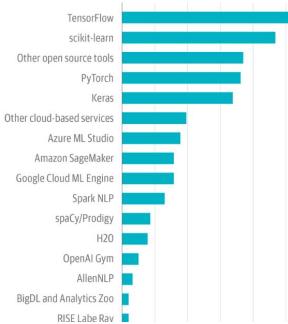


Spark NLP for Healthcare

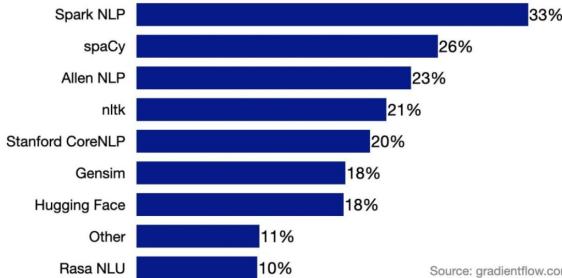


Spark NLP in Industry

Which of the following AI tools do you use?



Which NLP libraries does your organization use?



Source: gradientflow.co

TRUSTED BY



Recognized by the Technology Experts



Clinical Entity Recognition

Suspect diabetes	SNOMED-CT	Present
Unstable angina	ICD-10	Present
Hypertension	ICD-10	Present
Father with Alzheimer's	FAMILY	Present

Algorithms

Extract Knowledge	Entity Linker Entity Disambiguator Document Classifier Contextual Parser
De-Identity Text	Structured Data Unstructured Text Obfuscator Generalizer
Medical Transformers	350-BERT-Clinical ClinicalBERT Glove-Med Glove-ICD-O
Linked Medical Terminologies	SNOMED-CT CPT UMLS ICD-10-CM ICD-10-PCS ICD-O LOINC

Content

700+ Pretrained Models

Clinical:	Drug, Symptoms, Treatments, Procedures, Tests, Labs, Sessions
Anatomy:	Cell, Organ, Tissue, Cell Structure, Organ, Tissue, Gene, Chemical
Drugs:	Name, Dose, Strength, Route, Duration, Frequency, Potency, Active Ingredients, Side Effects
Demographics:	Age, Gender, Height, Weight, Address, University Status, Vital Signs
Risk Factors:	Smoking, Obesity, Diabetes, Hypertension, Resistance, Risk
Sensitive Data:	Patient Name, Address, Phone, Electronic Prescribers, Identifiers

Find in Text

- Text Matcher
- Regex Matcher
- Date Matcher
- Number Matcher

Trainable & Tunable

Scalable to a Cluster

Fast Inference

Hardware Optimized

Community

NLP SUMMIT

Spark NLP Modules Healthcare and Public

Biomedical Named Entity Recognition at Scale

Veysel Kocaman
John Snow Labs Inc.
16192 Coastal Highway
Lewes, DE, USA 19958
vey sel@johnsnowlabs.com

Abstract

Named entity recognition (NER) is a widely applicable natural language processing task and building block of question answering, topic modeling, information retrieval, etc. In this paper, we introduce NER models for extracting meaningful chunks from clinical notes and reports, which are then fed to downstream tasks like assertion status detection, entity resolution, and de-identification. We present a model implementing a Bi-LSTM-CNN-Char deep learning architecture on top of Apache Spark, to present a highly trainable NER model that obtains state-of-the-art performance on several benchmarks while being heavy contextual embeddings like BERT. This includes improving 4PC4MED to 93.72% (41.5% gain), 4PC4MED++ to 93.71% (41.5% gain), and 4PC4MED++-min to 93.2% (41.5% gain). In addition, this model is freely available within a production-grade code base as part of the open-source Spark NLP library. We also present our findings on how any Spark cluster (CPU, GPU, and Flink) for popular pre-trained languages such as Python, R, Scala, Java, and C can be extended to support our human languages with no code changes.

Anonymous NAACL-HLT 2021 submission

Improving Clinical Document Understanding on COVID-19 Research with Spark NLP

Veysel Kocaman, David Talby
John Snow Labs Inc.
16192 Coastal Highway
Lewes, DE, USA 19958
vey sel, david@johnsnowlabs.com

Abstract

Following the global COVID-19 pandemic, the number of scientific papers studying the virus has grown massively, leading to increased interest in automated literature review. We present a clinical document understanding pipeline designed in three ways. First, it can recognize over 100 different entity types including standard determinants of health, anatomy, risk factors, and clinical outcomes. Second, the pipeline includes named entity recognition based on a modified Bi-LSTM-CNN-Char DL architecture built on top of Apache Spark. This algorithm establishes new state-of-the-art accuracy on 7 and 8 entity types in biomedical NER benchmarks and 3 clinical document understanding challenges: 2010 i2b2/VA clinical concept extraction, 2014 n2c2 de-identification, and 2018 n2c2 medication extraction. Moreover, clinical NER models trained using this implementation

Peer-reviewed conference papers on Spark NLP NER

Clinical Entity Linking

Assertion Status

Relation Extraction

Entity Recognition

Information Extraction

Spelling & Grammar

Text Classification

Translation

Summarization

Question Answering

Emotion Detection

7,000+
Pre-trained Pipelines, Models & Transformers

250+
Languages

Understand Grammar

Find in Text

- Text Matcher
- Regex Matcher
- Date Matcher
- Character Matcher
- Question Answering

Trainable & Tunable

Scalable to a Cluster

Fast Inference

Hardware Optimized

Community

NLP SUMMIT

The most scalable, accurate and fastest arsenal of NLP transformers in history

Core NLP Transformer Models

GloveBERT Embeddings	BertEmbeddings	RoBERTaEmbeddings	XlmRobertaEmbeddings
DeBERTaEmbeddings	LongformerEmbeddings	BertEmbeddings	XlmRobertaEmbeddings
DeBERTaForQuestionAnswering	LongformerForTokenClassification	BertForTokenClassification	XlmRobertaForTokenClassification
GPT2Transformer	LongformerForSequenceClassification	BertForSequenceClassification	XlmRobertaForSequenceClassification
MarianTransformer	LongformerForQuestionAnswering	BertForQuestionAnswering	XlmRobertaForQuestionAnswering
T5Transformer	LongformerForEncoderDecoder	BertForEncoderDecoder	XlmRobertaForEncoderDecoder
RobertaEmbeddings	XlnetEmbeddings	BertEmbeddings	AlbertEmbeddings
RobertaForTokenClassification	XlnetForTokenClassification	BertForTokenClassification	AlbertForTokenClassification
RobertaForSequenceClassification	XlnetForSequenceClassification	BertForSequenceClassification	AlbertForSequenceClassification
RobertaForQuestionAnswering	XlnetForQuestionAnswering	BertForQuestionAnswering	AlbertForQuestionAnswering

Subflavored NLP Transformer Extensions

SciBertEmbeddings	ElectraForQuestionAnswering	MuRILBERTForQuestionAnswering	MuRILBERTEmbeddings	LanuseCenteveciembbedings
IndoBERTForQuestionAnswering	DeBERTaV3Embeddings	MuRILBERTSequenceEmbeddings	SciBertForQuestionAnswering	SapBERTForQuestionAnswering
ElectraEmbeddings	BiBERTForQuestionAnswering	CovidBERTForQuestionAnswering	MinILMForQuestionAnswering	MacBERTForQuestionAnswering
BioFormerForQuestionAnswering	RedditGPTEmbeddings	CovidDistilBERTForQuestionAnswering	LinkBERTForQuestionAnswering	BroderEmbeddings

Optimized, Tested, Supported Integrations

jupyter

ANACONDA

kubernetes

Scala

R

CLOUDERA

oneAPI

comet

Amazon SageMaker

mlflow

Cloud Dataproc

intel

databricks

Google Cloud

Microsoft Azure

NVIDIA

Apache Spark

Java

python

Apache

Apache Synapse ML

Simple and Distributed Machine Learning

Legend:

- Transformer for Sequence to Sequence
- Transformer for Token Classification
- Transformer for Sequence Classification
- Transformer for Question Answering
- Transformer Sentence Embeddings
- Transformer Word Embeddings
- Non Transformer Embeddings

NLU: All 9000+ John Snow Labs models & features in 1 line of code

How does it work?

```
model= nlu.load(model)
```

- Returns a nlu pipeline object

```
model.predict(data)
```

- Returns a pandas DF

How does it work?

```
model = nlu.load('emotion')
```

- Returns a nlu pipeline object

```
model.predict('I love NLU!')
```

- Returns a pandas DF

EMOTION DETECTION

```
nlu.load('emotion').predict('I love NLU!')
```

sentence_embeddings	category_sentence	category_surprise	category_sadness	category_joy	category_fear	sentence	category	id
[0.027570432052016258, -0.05264767631883896, ...]	0	0.012899903	0.0015578865	0.9760173	0.0095249	I love NLU!	joy	1

NLU WORKS DIRECTLY ON TYPICAL PYTHON DATASETS

Strings

```
import nlu  
nlu.load('sentiment').predict('This is just one string')
```

Lists

```
import nlu  
nlu.load('sentiment').predict(['This is an array', ' of strings!'])
```

Pandas data frame

```
import nlu  
import pandas as pd  
data = {'text': ['This day sucks', 'I love this day', 'I dont like Sami']}  
text_df = pd.DataFrame(data)  
nlu.load('sentiment').predict(text_df)
```

Pandas series

```
import nlu  
import pandas as pd  
data = {'text': ['This day sucks', 'I love this day', 'I dont like Sami']}  
text_df = pd.DataFrame(data)  
nlu.load('sentiment').predict(text_df['text'])
```

Spark
Data Frame

Ray
Data Frame

Dask
Data Frame

Where is Text Generation used in the Industry?

- Natural Language Generation (NLG) is a sub-field of natural language processing (NLP)
 - Automatically generate **coherent** and **useful** text for human consumption
- NLG Systems in use
 - Google Translate
 - Dialogue Systems (Siri, Alexa, Cortana, Google Home, Bigsby)
 - Summarizers (Google Search Results)
 - Emails, News, Papers, Reports, Conversations
 - Creative writing
 - Stories, Poetry, Narratives
 - Fake News, Botnets



Automated Text Generation & Data-Augmentation for Medicine, Finance, Law, and E-Commerce

Explore GPT Parameters Play with Parameters

Sampling Methods

Play with temperature

GPT2 Text Generation Industry Use-Cases

Zero-Shot-Learning Applications & Examples

Data2Text for finance assets

Treatment Recommendations for Medicine

SEO Text / Content Writing / Copy Writing

Marketing Texts for Hoodies

Marketing Texts for Soap

Marketing Texts for beard products

Generate Product Recommendations

Top 10 Movie List

Top 10 Game List

Top 10 Book List

Song Lyrics

Last Christmas

Fresh prince of GPT

GPT Jackson

Real Slim GPT

GPT Rap God Bot

Fiction Writing

Fantasy Stories

Lovercraftian Horror Stories

Re-Write movie script

GPT for Programming Code

Scala

Python

Data Augmentation

Data Augmentation example for news Classifier Dataset

Vanilla Training

Base Dataset Preparation
Train a new News Classifier from the base dataset

Training Dataset Metrics

Test Dataset Metrics

Augmented Training

Augmented Dataset Creation
Train a new Model with Augmented Data

Evaluate Train Metrics for Augmented Model

Evaluate Test Metrics for Augmented Model

Let's Create some generic functions and test data augmentation on other datasets

Dataset Creation Function

Fit and Evaluate Function

Create Augmented Dataset Function

Compare Vanilla with Augmented Training Function

Let's test out the functions on a dataset

Data Augmentation on a Finance Dataset

Data Augmentation on a Medical Dataset

Data Augmentation on a Legal Dataset

There are many more models you can put to use in a line of code!

Text Generation in Action : Code Time

Dipping our toes in NLG with NLU using GPT2

GPT2 NLU Notebook

Why are generation repeating and what did we just do?

A quick theory crash course follows

```
[ ] gpt2_model['gpt2'].setDoSample(False)
generation_df = gpt2_model.predict("My Favorite Food is!")
print(generation_df.generated.values[0])
generation_df = gpt2_model.predict("My Favorite Food is!")
print(generation_df.generated.values[0])
```

My Favorite Food is!

I love the way the food is cooked. I love the texture. I like the way it
My Favorite Food is!

I love the way the food is cooked. I love the texture. I like the way it

Why are generation repeating?

```
[ ] gpt2_model['gpt2'].setDoSample(True)
gpt2_model['gpt2'].setMaxOutputLength(25)
generation_df = gpt2_model.predict("My Favorite Food is!")
print(generation_df.generated.values[0])
generation_df = gpt2_model.predict("My Favorite Food is!")
print(generation_df.generated.values[0])
```

My Favorite Food is!

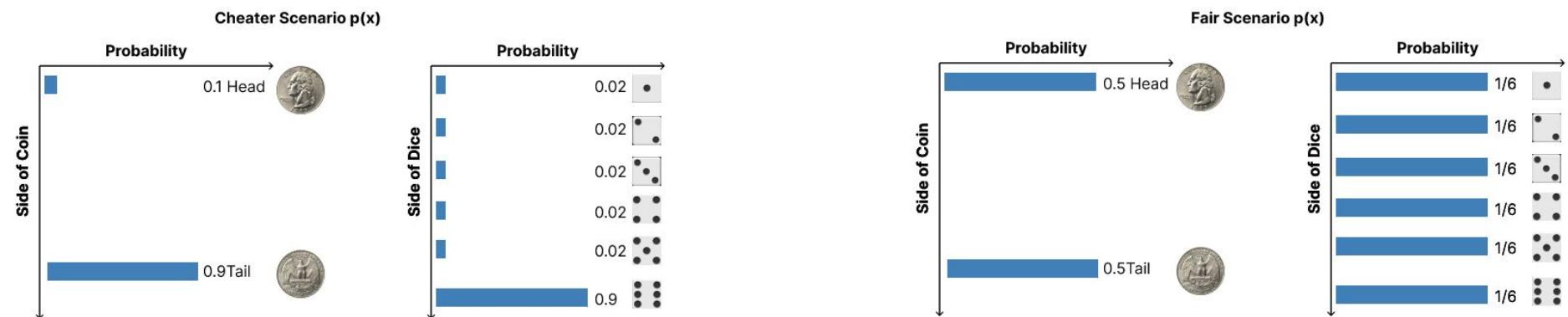
We have a great family who loves having family gatherings with our kids. From our family gatherings to
My Favorite Food is!

A recent "food for thought" article, entitled "How to Make a Food that is Easy

Why are generations not repeating?

What is a Probability Distribution?

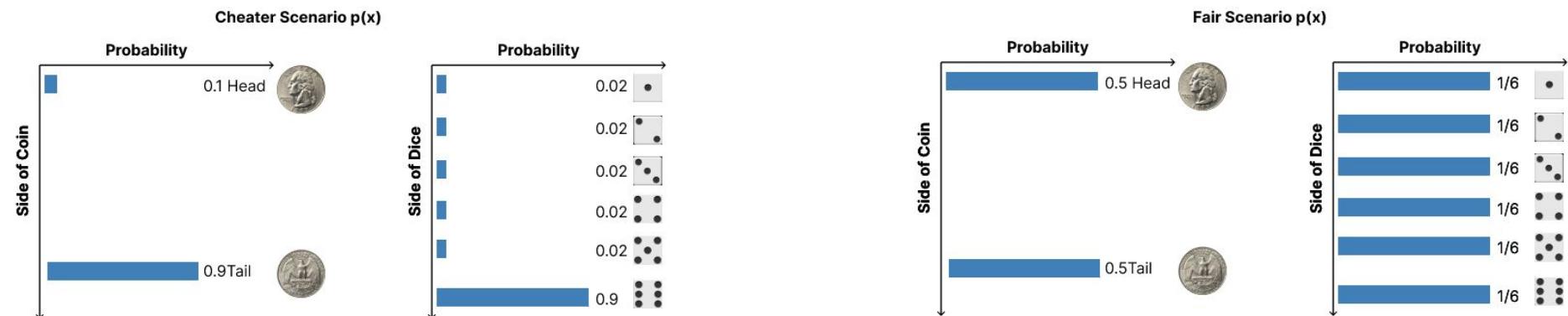
- A probability distribution is simply a function $p : \Omega \rightarrow [0, 1]$ which returns a probability for events x from a sample space Ω
- In a coin flip scenario, sample space is $\Omega = \{Heads, Tails\}$ and for every event $p(x) = 0.5$ if the coin is fair.
- In a 6-sided-Dice throw scenario, sample space is $\Omega = \{1, 2, 3, 4, 5, 6\}$ and for every $p(x) = 1/6$ if the dice is fair.



What is drawing from a probability distribution?

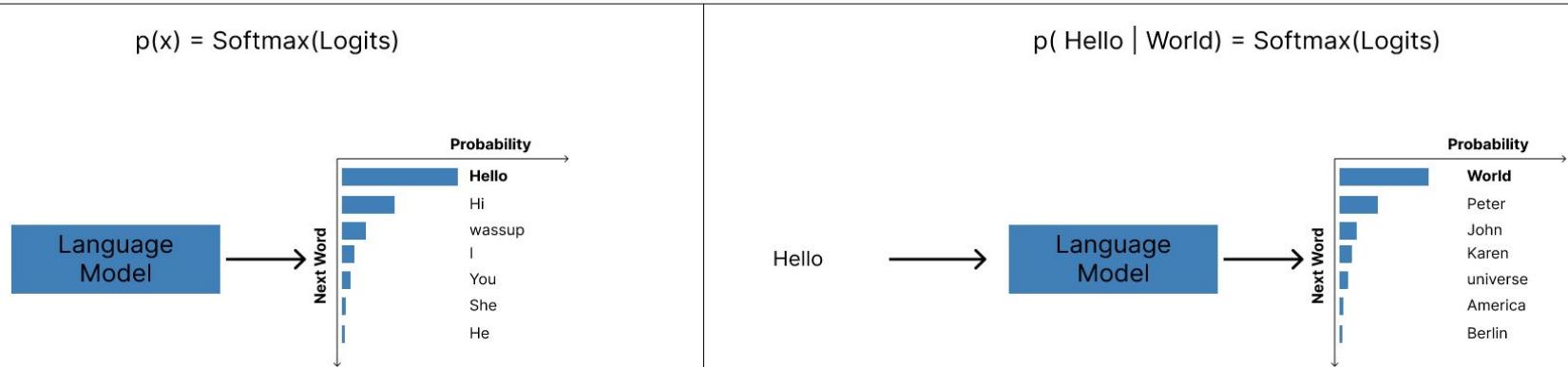
Also known as sampling from a distribution

- If we know the probability distribution of a sample space, we can **simulate** it
- We simulate by **drawing** elements in Ω using some **sampling algorithm**



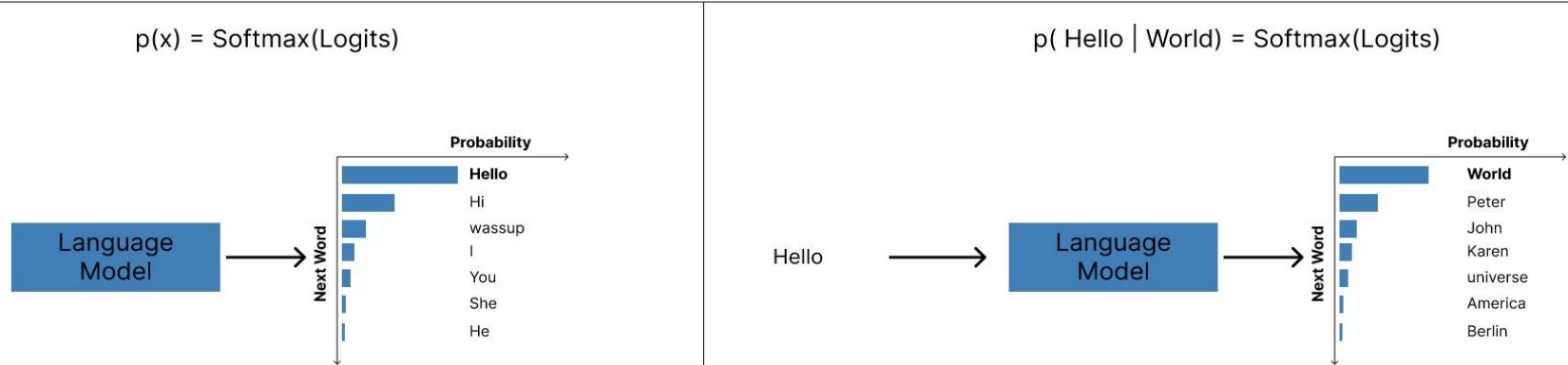
What is a Language Model?

- A Language Model is simply a probability distribution $p(x)$ where the sample space Ω is the **set of words in the human language**
- $p(\text{hello})$ should tell us the probability of the word **hello** occurring as **beginning of human written text**
- $p(\text{world} \mid \text{hello})$ should tell us the probability of the word **world** given that human written text started with **hello**

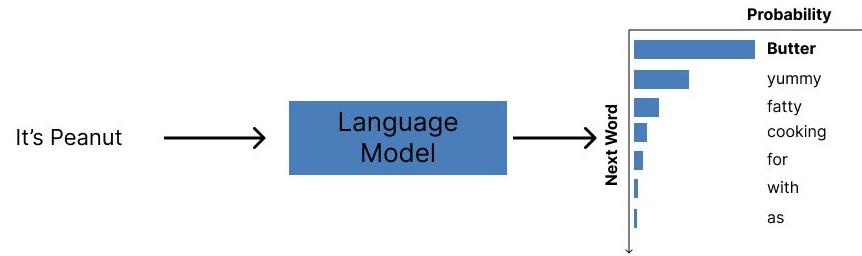


How does a Language Model generate Text?

- The function $p(x)$ is used to draw a word from the human language
- The function $p(x|y)$ is used to draw a word from the human language, given a start word
 - We can feed the outputs of $p(x|y)$ back to itself to generate one word at a time, conditioned on preceding words
 - **y** is referred to as **the prompt** we are giving to the model
- Finding an decent estimation of p has been solved by Transformer based deep learning architectures like Bert, GPT, T5, etc..
- When we use these Pretrained Language Models to generate text our main concern **shaping** and **sampling** from these distributions
- Final layer of generative NLP Transformer exposes one logit for every word in the training vocabulary.
 - Applying Softmax function to outputs of logits, yields probability distribution $p(x | y)$ which we sample from

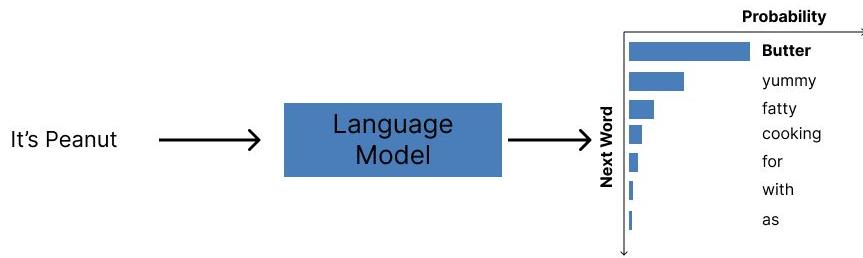


$$P(\text{Next-Word} \mid \text{It's, Peanut}) = \text{Softmax}(\text{Logits})$$

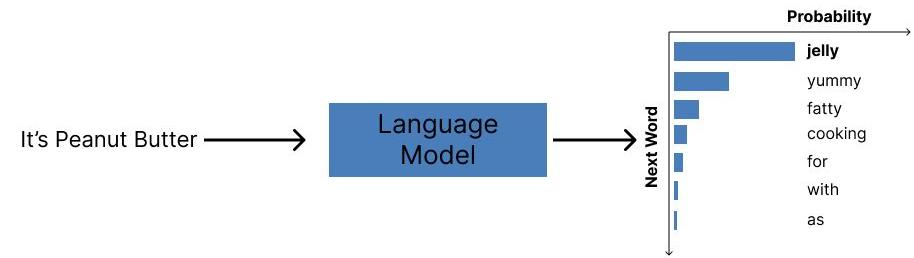


How does a Language Model generate Text?

$$P(\text{ Next-Word} \mid \text{It's, Peanut}) = \text{Softmax}(\text{Logits})$$

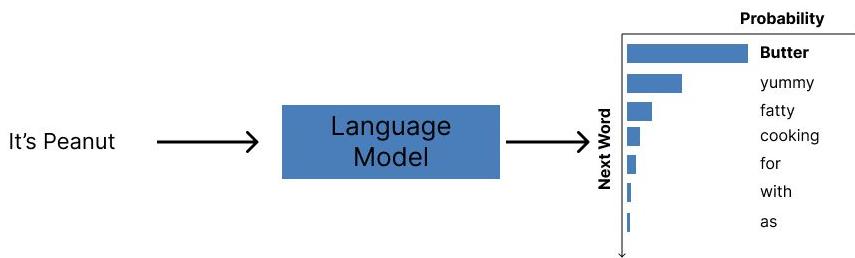


$$P(\text{ Next-Word} \mid \text{It's, Peanut, Butter}) = \text{Softmax}(\text{Logits})$$

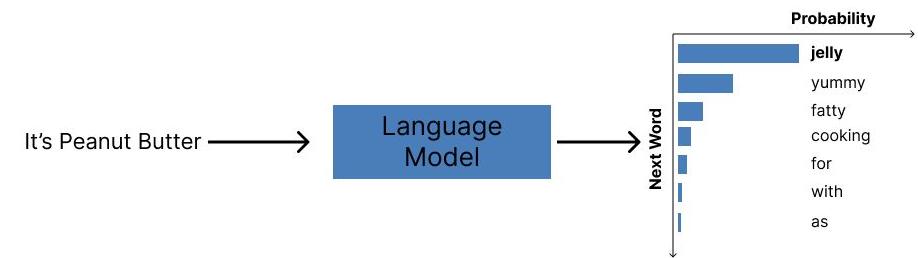


How does a Language Model generate Text?

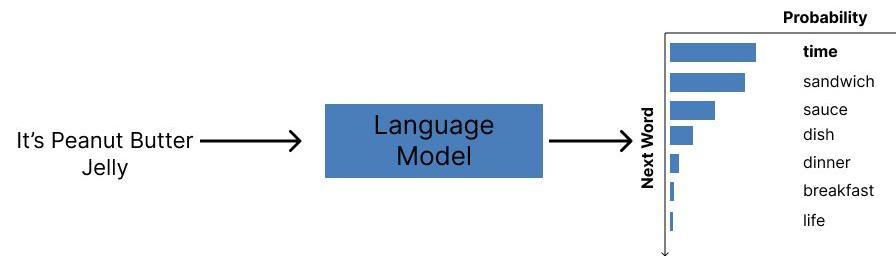
$$P(\text{ Next-Word} \mid \text{It's, Peanut}) = \text{Softmax}(\text{Logits})$$



$$P(\text{ Next-Word} \mid \text{It's, Peanut, Butter}) = \text{Softmax}(\text{Logits})$$

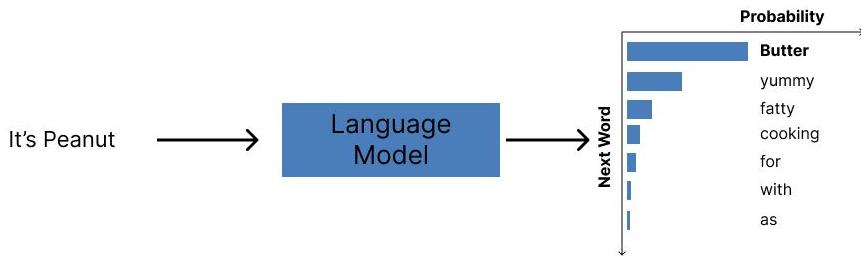


$$P(\text{ Next-Word} \mid \text{It's, Peanut, Butter,Jelly}) = \text{Softmax}(\text{Logits})$$

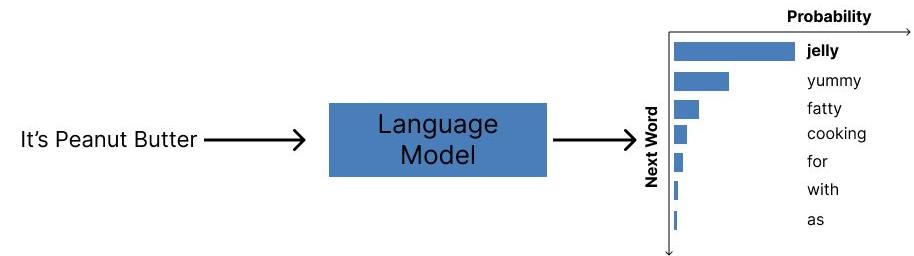


How does a Language Model generate Text?

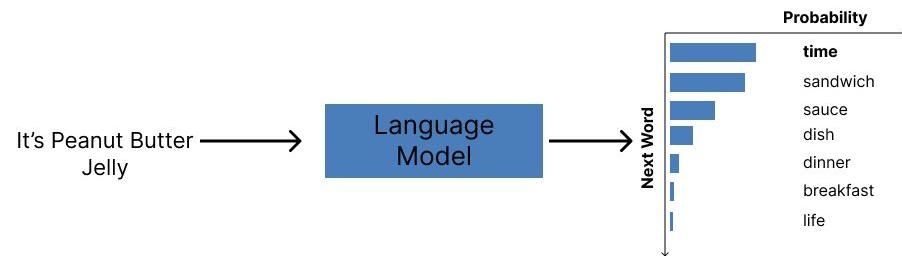
$$P(\text{ Next-Word} \mid \text{It's, Peanut}) = \text{Softmax}(\text{Logits})$$



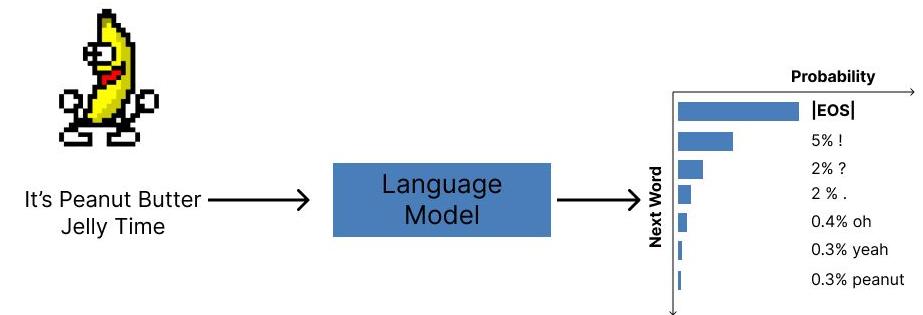
$$P(\text{ Next-Word} \mid \text{It's, Peanut, Butter}) = \text{Softmax}(\text{Logits})$$



$$P(\text{ Next-Word} \mid \text{It's, Peanut, Butter,Jelly}) = \text{Softmax}(\text{Logits})$$



$$P(\text{ Next-Word} \mid \text{It's, Peanut, Butter}) = \text{Softmax}(\text{Logits})$$



How does a Language Model generate Text?

What are some Sampling Algorithms and Parameters ?

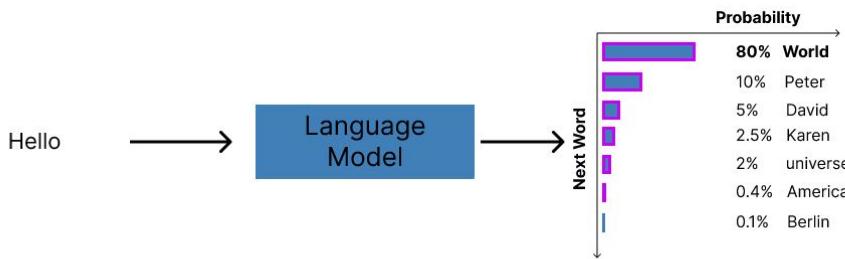
- **Parameters:**

- Repetition Penalty: Penalize generation candidates based on how often they already exist based See [CTRL](#)
- Max Repeated NGrams: Block some NGrams to occur more than K times
- Temperature: Reshape distributed given by Softmax
- Min/Max Output Length: Limit the length of generated

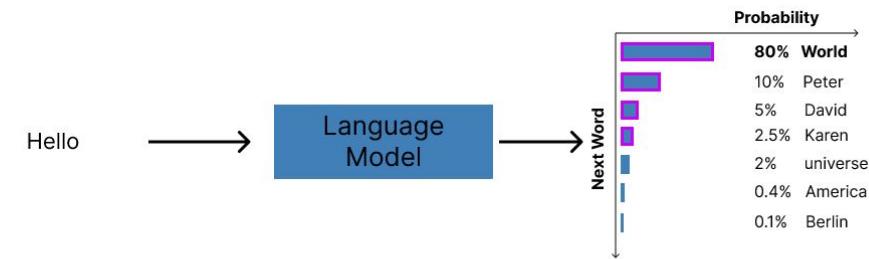
- **Sampling Algorithms:**

- ArgMax based : Pick word with highest probability of occurring next, Greedy Algorithm.
 - Makes generation deterministic, i.e. always generates the same text for the same prompt
 - When `.doSample(False)` our will use this algorithm and all previous examples used this algorithm
- Top-K Sampling: choose from **Top-K highest probability** words in distribution
- Top-P Sampling: choose from words with **summed probability** equal or greater than P

Top-K Sampling
Choose from first K word with highest probability
For K=6



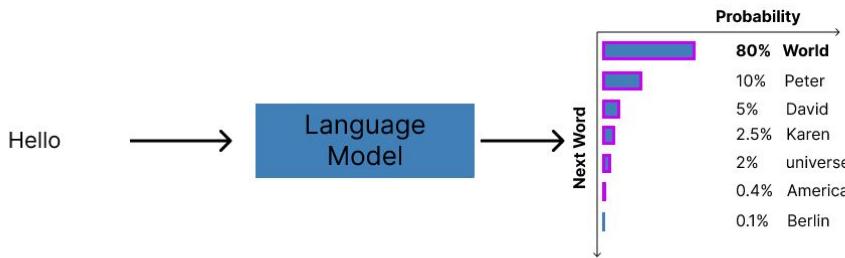
Top-P Sampling
Choose from first tokens whose collective probability is at least P
For P = 0.95



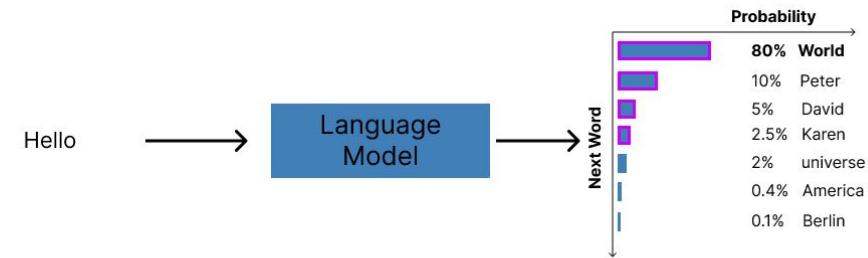
What's the difference between Top-P and Top-K Sampling?

- **Higher K and P:** makes choices more diverse and risky. Less logical and coherent
- **Lower K and P:** makes choices more generic and safe. Also more logical and coherent
- Softmax normalizes logits between 0 and 1. We can tweak the function with the T parameter
- **High T:** Make Distribution more uniform and diverse.
Associated with creativity and randomness of generations.
- **Low T:** Make Distribution concentrate around top words, more spiky and less diverse.
Associated with more logical and coherent generations.

Top-K Sampling
Choose from first K word with highest probability
For K=6



Top-P Sampling
Choose from first tokens whose collective probability is at least P
For P = 0.95



What is Prompt Engineering?

- Prompt Engineering is the art of **finding prompts** which shape the conditioned probability $p(x | y)$ such that only tokens relevant to a specific use case have a high chance of being sampled

^{37]} # Bad prompting, the input text we condition GPT2 yields bad output, it does not understand the pattern we want from the original input
gpt2_pipe.predict("Suggest me a good Sci-Fi movie")

index	document	generated
0	Suggest me a good Sci-Fi movie	generate Suggest me a good Sci-Fi movie. I'm not sure if I'm going to be able to do this, but I'm sure I'll be able. (I'm sure you're going to want to do it.) So, I'm not going to do that . . .

Show 25 ▾ per page
Like what you see? Visit the [data table notebook](#) to learn more about interactive tables.

Good prompting. help GPT2 out and by giving it a few samples in the prompt we condition it on
gpt2_pipe.predict("""Generate a top 10 movie list: \n1.The Matrix \n2.Terminator \n3. """)

index	document	generated
0	Generate a top 10 movie list: 1.The Matrix 2.Terminator 3.The Hunger Games 4.The Dark Knight Rises 5.The Lord of the Rings 6.The Hobbit 7.The Last Crusade 8.The Lion King 9.The Lego Movie 10.The LEGO Movie 11.The King of the Hill 12.The Jungle Book 13.The Legend of Zelda 14.The Little Mermaid 15.The Princess Bride 16.The Pirates of the Caribbean 17.The Twilight Saga 18	generate Generate a top 10 movie list: 1.The Matrix 2.Terminator 3.The Hunger Games 4.The Dark Knight Rises 5.The Lord of the Rings 6.The Hobbit 7.The Last Crusade 8.The Lion King 9.The Lego Movie 10.The LEGO Movie 11.The King of the Hill 12.The Jungle Book 13.The Legend of Zelda 14.The Little Mermaid 15.The Princess Bride 16.The Pirates of the Caribbean 17.The Twilight Saga 18

Show 25 ▾ per page

Like what you see? Visit the [data table notebook](#) to learn more about interactive tables.

Let's go back to the notebook and apply what we just learned

GPT2 NLU Notebook

Table of contents 

Automated Text Generation & Data-Augmentation for Medicine, Finance, Law, and E-Commerce

Explore GPT Parameters Play with Params

Sampling Methods

Play with temperature

GPT2 Text Generation Industry Use-Cases

Zero-Shot-Learning Applications & Examples

- Data2Text for finance assets
- Treatment Recommendations for Medicine

SEO Text / Content Writing / Copy Writing

- Marketing Texts for Hoodies
- Marketing Texts for Soap
- Marketing Texts for beard products

Generate Product Recommendations

- Top 10 Movie List
- Top 10 Game List
- Top 10 Book List

Song Lyrics

- Last Christmas
- Fresh prince of GPT
- GPT Jackson
- Real Slim GPT
- GPT Rap God Bot

Fiction Writing

- Fantasy Stories
- Lovercraftian Horror Stories
- Re-Write movie scripts

GPT for Programming Code

- Scala
- Python

Data Augmentation

Data Augmentation example for news Classifier Dataset

Vanilla Training

- Base Dataset Preparation
- Train a new News Classifier from the base dataset
- Training Dataset Metrics
- Test Dataset Metrics

Augmented Training

- Augmented Dataset Creation
- Train a new Model with Augmented Data
- Evaluate Train Metrics for Augmented Model
- Evaluate Test Metrics for Augmented Model

Let's create some generic functions and test data augmentation on other datasets

- Dataset Creation Function
- Fit and Evaluate Function
- Create Augmented Dataset Function
- Compare Vanilla with Augmented Training Function

Let's test out the functions on a dataset

Data Augmentation on a Finance Dataset

Data Augmentation on a Medical Dataset

Data Augmentation on a Legal Dataset

There are many more models you can put to use in 1 line of code!

What is Data Augmentation with Language Models?

- Since the largest Language Modes are trained on almost the entire Internet as text corpus, they can generate realistic text for almost any topic.
- We can **improve accuracy** of most NLP Classifiers by **augmenting our dataset** generated examples
 - Especially useful when only very little training data is available and creating new by hand is tedious
- The newly NLP classifier is likely much faster and smaller than the large Language Model
 - Essentially a distilled version or the Language Models knowledge
- How to **create more training** data and Augmented Dataset:
 - Use the label itself as prompt and suffix it with a chunk of the original text.
Try to use chunks which have enough semantic information to justify a prediction of the label.
One of the most simple methods is to use the first N token or characters as suffix to the label
- The Language Model could be used for training, but the trained model is much faster, efficient and possibly more accurate

How to create more training data and Augmented Dataset?

- There are dozens of techniques and approaches, see [A Survey of Data Augmentation Approaches for NLP](#)
- For our simple demo we will :
 - Use the label itself as prompt and suffix it with a chunk of the original text.
Try to use chunks which have enough semantic information to justify a prediction of the label.
One of the most simple methods is to use the first N token or characters as suffix to the label

DA Method	Ext.Know	Pretrained	Preprocess	Level	Task-Agnostic
SYNONYM REPLACEMENT (Zhang et al., 2015)	✓	✗	tok	Input	✓
RANDOM DELETION (Wei and Zou, 2019)	✗	✗	tok	Input	✓
RANDOM SWAP (Wei and Zou, 2019)	✗	✗	tok	Input	✓
BACKTRANSITION (Sennrich et al., 2016)	✗	✓	Depends	Input	✓
SCPN (Wieting and Choueiri, 2017)	✗	✓	const	Input	✓
SIMILAR TEXT EXCHANGE (Feng et al., 2019)	✗	✓	const	Input	✓
CONTEXTUAL AUG (Kobayashi, 2018)	✗	✓	-	Input	✓
LAMMADA (Amby-Tavor et al., 2020)	✗	✓	-	Input	✗
Geca (Andreas, 2019)	✗	✗	tok	Input	✗
SeqMixUp (Guo et al., 2020)	✗	✗	tok	Input	✗
SWITCHOUT (Wang et al., 2018b)	✗	✗	tok	Input	✗
EMIX (Jindal et al., 2020a)	✗	✗	-	Emb/Hidden	✓
SPEECHMIX (Jindal et al., 2020b)	✗	✗	-	Emb/Hidden	Speech/Audio
MIXTEXT (Chen et al., 2020c)	✗	✗	-	Emb/Hidden	✓
SUPEREDGRAPH (Chen et al., 2020b)	✗	✗	-	Input	✗
DTremMorph (Suhir and Steedman, 2018)	✗	✗	dep	Input	✓
Sub ^b (Shi et al., 2021)	✗	✗	dep	Input	Substructural
DAGA (Ding et al., 2020)	✗	✗	tok	Input+Label	✗
WN-HYPERS (Feng et al., 2020)	✓	✗	const+KWE	Input	✓
SYNTHETIC NOISE (Feng et al., 2020)	✗	✗	tok	Input	✓
UEDIN-MS (DA pair) (Grundkiewicz et al., 2019)	✓	✗	tok	Input	✓
NONCE (Gulordava et al., 2018)	✓	✗	const	Input	✓
XLD (Singh et al., 2019)	✗	✓	Depends	Input	✓
SeqMix (Zhang et al., 2020)	✗	✓	tok	Input+Label	✗
SLOT-SUB-LM (Lougee and Magnini, 2020)	✗	✓	tok	Input	✓
UBT & TBT (Vaihauer et al., 2019)	✗	✓	Depends	Input	✓
SOFT CONTEXTUAL DA (Gao et al., 2019)	✗	✓	tok	Emb/Hidden	✓
DATA DIVERSIFICATION (Nguyen et al., 2020)	✗	✓	Depends	Input	✓
DIPS (Kumar et al., 2019a)	✗	✓	tok	Input	✓
AUGMENTED SBERT (Thakur et al., 2021)	✗	✓	-	Input+Label	Sentence Pairs

Table 1: Comparing a selection of DA methods by various aspects relating to their applicability, dependencies, and requirements. *Ext.Know*, *KWE*, *tok*, *const*, and *dep* stand for External Knowledge, keyword extraction, tokenization, constituency parsing, and dependency parsing, respectively. *Ext.Know* refers to whether the DA method requires external knowledge (e.g. WordNet) and *Pretrained* if it requires a pretrained model (e.g. BERT). *Preprocess* denotes preprocessing required, *Level* denotes the depth at which data is modified by the DA, and *Task-Agnostic* refers to whether the DA method can be applied to different tasks. See Appendix B for further explanation.

Why does this work?

- **Universal Approximation Theorem** states: More data = more accuracy
Holds for Neural Networks and many ML models
- Leverage this when we augment dataset size with generative models
- Bad generation can make our model perform worse though, we need to be careful how and what is generated

Data Augmentation in Action

apply what we just learned

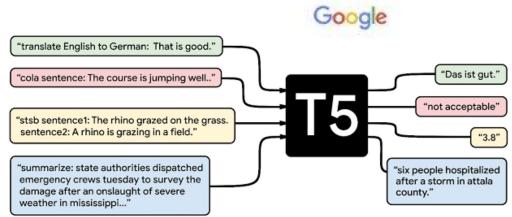
GPT2 NLU Notebook

Table of contents

- Automated Text Generation & Data-Augmentation for Medicine, Finance, Law, and E-Commerce
- Explore GPT Parameters Play with Params
 - Sampling Method
 - Play with temperature
- GPT2 Text Generation Industry Use-Cases
 - Zero-Shot-Learning Applications & Examples
 - DataText for finance assets
 - Treatment Recommendations for Medicine
 - SEO Text / Content Writing / Copy Writing
 - Marketing Texts for Hoodies
 - Marketing Texts for Soap
 - Marketing Texts for beard products
 - Generate Product Recommendations
 - Top 10 Movie List
 - Top 10 Game List
 - Top 10 Book List
 - Song Lyrics
 - Last Christmas
 - Fresh prince of GPT
 - GPT.Jackson
 - Real Slim GPT
 - GPT Rap God Bot
 - Fiction Writing
 - Fantasy Stories
 - Lovercraftian Horror Stories
 - Re-Write movie scripts
 - GPT for Programming Code
 - Scala
 - Python
- Data Augmentation
 - Data Augmentation example for news Classifier Dataset
 - Vanilla Training
 - Base Dataset Preparation
 - Train a new News Classifier from the base dataset
 - Training Dataset Metrics
 - Test Dataset Metrics
 - Augmented Training
 - Augmented Dataset Creation
 - Train a new Model with Augmented Data
 - Evaluate Train Metrics for Augmented Model
 - Evaluate Test Metrics for Augmented Model
 - Let's create some generic functions and test data augmentation on other datasets
 - Dataset Creation Function
 - Fit and Evaluate Function
 - Create Augmented Dataset Function
 - Compare Vanilla with Augmented Training Function
 - Let's test out the functions on a dataset
 - Data Augmentation on a Finance Dataset
 - Data Augmentation on a Medical Dataset
 - Data Augmentation on a Legal Dataset

There are many more models you can put to use in 1 line of code!

Text Generating Sequence to Sequence Transformers in NLU & Spark NLP



Task Name	Explanation
1.CoLA	Classify if a sentence is grammatically correct
2.RTE	Classify whether a statement can be deduced from a sentence
3.MNLI	Classify for a hypothesis and premise whether they contradict or contradict each other or neither of both (3 classes).
4.MRPC	Classify whether a pair of sentences is a re-phrasing of each other (semantically equivalent)
5.QNLI	Classify whether the answer to a question can be deduced from an answer candidate.
6.QQP	Classify whether a pair of questions is a re-phrasing of each other (semantically equivalent)
7.SST2	Classify the sentiment of a sentence as positive or negative
8.STSB	Classify the sentiment of a sentence on a scale from 1 to 5 (21 Sentiment classes)
9.CB	Classify for a premise and a hypothesis whether they contradict each other or not (binary).
10.COPA	Classify for a question, premise, and 2 choices which choice the correct choice is (binary).
11.MultiRc	Classify for a question, a paragraph of text, and an answer candidate, if the answer is correct (binary).
12.WIC	Classify for a pair of sentences and a disambiguous word if the word has the same meaning in both sentences.
13.WSCIDPR	Predict for an ambiguous pronoun in a sentence what it is referring to.
14.Summarization	Summarize text into a shorter representation.
15.SQuAD	Answer a question for a given context.
16.WMT1	Translate English to German
17.WMT2	Translate English to French
18.WMT3	Translate English to Romanian

Summarize, answer questions, generate SQL and more with [Google's T5](#)
Showcased in [this tutorial Webinar](#)

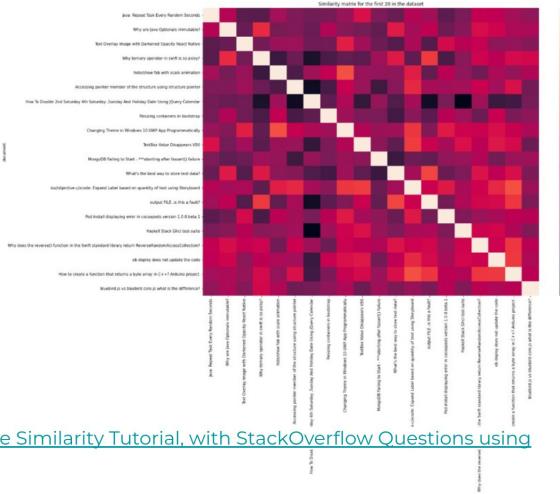


Generate long sequences of text, Augment Text and more with [Open AI's GPT2](#)
Showcased today



Translate between over 200 languages with [Microsoft's MarianNMT](#)
showcased on chinese in [this tutorial webinar](#)

Sentence Similarity With BERTology Embeds or T5



[Indepth and Easy Sentence Similarity Tutorial, with StackOverflow Questions using BERTology embeddings](#)

Explore 100+ Embeddings via 6 Similarity Metrics with 0 lines of code with NLU & Streamlit

- Compare Multiple similarity Metrics And Embeddings at the same time
 - Supported similarities :
 - Cosine
 - Cityblock
 - Euclidean
 - L2
 - L1
 - Manhattan

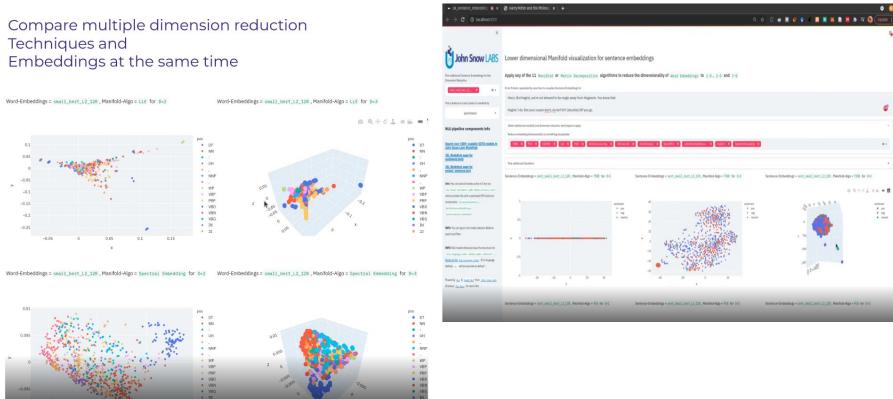


t-SNE Visualizations with NLU



1line of Python code for BERT, ALBERT, ELMO, ELECTRA, XLNET, GLOVE, Part of Speech with NLU and t-SNE

Explore 100+ Embeddings via 10+ Manifold and Matrix Decomposition with 0 lines of code with NLU & Streamlit



OCR Tables from PDFs, DOC, DOCX and PPT files in 1 Line of code with NLU 4.0



```
[1] nlu.load('pdf2table').predict('/content/tables.pdf')
```

	mpg	cyl	disp	hp	drat	wt	qsec	vs	am	gear
1	21.0	6	160.0	110	3.90	2.620	16.46	0	1	4
2	21.0	6	160.0	110	3.90	2.875	17.02	0	1	4
3	22.8	4	108.0	93	3.85	2.320	18.61	1	1	4
4	21.4	6	258.0	110	3.08	3.215	19.44	1	0	3
5	18.7	8	360.0	175	3.15	3.440	17.02	0	0	3
6	18.1	6	225.0	105	2.76	3.460	20.22	1	0	3
7	14.3	8	360.0	245	3.21	3.570	15.84	0	0	3
8	24.4	4	146.7	62	3.69	3.190	20.00	1	0	4
9	22.8	4	140.8	95	3.92	3.150	22.90	1	0	4
10	19.2	6	167.6	123	3.92	3.440	18.30	1	0	4
11	17.8	6	167.6	123	3.92	3.440	18.90	1	0	4
12	16.4	8	275.8	180	3.07	4.070	17.40	0	0	3
13	17.3	8	275.8	180	3.07	3.730	17.60	0	0	3
14	15.2	8	275.8	180	3.07	3.780	18.00	0	0	3
15	10.4	8	472.0	205	2.93	5.250	17.98	0	0	3
16	10.4	8	460.0	215	3.00	5.424	17.82	0	0	3
17	14.7	8	440.0	230	3.23	5.345	17.42	0	0	3
18	32.4	4	78.7	66	4.08	2.200	18.47	1	1	4
19	30.4	4	75.7	52	4.93	1.615	18.52	1	1	4

Overview of all OCR features in NLU

Overview of OCR Text Extractors

These models grab the text directly from your input file and returns it as a Pandas DataFrame

NLU Spell	Transformer Class
nlu.load('img2text')	ImageToText
nlu.load('pdf2text')	PdfToText
nlu.load('doc2text')	DocToText

Overview of OCR Table Extractors

These models grab all Table data from the files detected and return a [list of Pandas DataFrames](#), containing Pandas DataFrame for every table detected

NLU Spell	Transformer Class
nlu.load('pdf2table')	PdfToTextTable
nlu.load('ppt2table')	PptToTextTable
nlu.load('dec2table')	DocToTextTable

More info on

https://nlu.johnsnowlabs.com/docs/en/nlu_for_ocr

Image to Text

"The Old Pond" by Matsuo Bashō

An old silent pond

A frog jumps into the pond—

Splash! Silence again.

Sample image:

```
nlu.load('img2text').predict('path/to/haiku.png')
```

Output of IMG OCR:

text
"The Old Pond" by Matsuo Bashō
An old silent pond
A frog jumps into the pond—
Splash! Silence again.

PDF to Text

haiku.pdf

Sample PDF: ["Lighting One Candle" by Yosa Buson](#)

The light of a candle

Is transferred to another candle—

Spring twilight

Output of PDF OCR:

text
"Lighting One Candle" by Yosa Buson
The light of a candle
Is transferred to another candle—
Spring twilight

DOCX to text

haiku.docx - LibreOffice Writer

Sample DOCX:

"In a Station of the Metro" by Ezra Pound

The apparition of these faces in the crowd;

Petals on a wet, black bough.

Output of DOCX OCR:

text
"In a Station of the Metro" by Ezra Pound
The apparition of these faces in the crowd;
Petals on a wet, black bough.

More NLU Ressources

- [Join our Slack](#)
- [NLU Website](#)
- [NLU Github](#)
- [Many more NLU example tutorials](#)
- [Overview of every powerful nlu 1-liner](#)
- [Checkout the Modelshub for an overview of all models](#)
- [Checkout the NLU Namespace where you can find every model as a tabel](#)
- [Intro to NLU article](#)
- [In Depth and Easy Sentence Similarity Tutorial, with StackOverflow Questions using BERTology embeddings](#)
- [1 line of Python code for BERT, ALBERT, ELMO, ELECTRA, XLNET, GLOVE, Part of Speech with NLU and t-SNE](#)

Webinars and Videos

- [NLU & Streamlit Tutorial](#)
- [Crash course of the 50 + Medical Domains and the 200+ Healthcare models in NLU](#)
- [Multi Lingual NLU Webinar - Working with Chinese Text Data](#)
- [John Snow Labs NLU: Become a Data Science Superhero with One Line of Python code](#)
- [Python Web Def Conf - Python's NLU library: 4,000+ Models, 200+ Languages, State of the Art Accuracy, 1 Line of Code](#)
- [NYC/DC NLP Meetup with NLU](#)

Thank you.



christian@JohnSnowLabs.com



@ckl_it



In/Christian-Kasim-Loan



Medium.com/@Christian.Kasim.Loan