



UNITED STATES MILITARY ACADEMY
WEST POINT

Quickly Build Analyst NLP Capability in Your Organization

MAJ John Scudder
AORS 2023





- NLP permeates our daily lives
- How can you quickly learn and leverage NLP for your own use cases?
- **My Goal: Equip you with a framework and resources to bring NLP capability to analysts in your organization.**





- Background
- Methodology: Environment / Audience Considerations
- Lesson Plan
- Best Practices
- Resources for YOU!
 - Curated opensource resources
 - Code repository with Jupyter Notebooks and slides

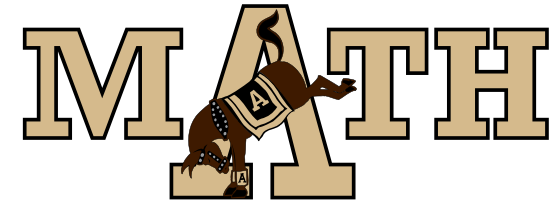




UNITED STATES MILITARY ACADEMY
WEST POINT

Background





- Problem set within USMA D/Math:
 - Existing research projects have a requirement to analyze unstructured text data → NLP knowledge
 - All D/Math Faculty have experience with w/ structured data
 - ~5 members w/ NLP experience out of > 90 faculty
- Solution: Combine local experience & offer 5x interactive lectures w/ runnable code
- I'm providing you the product we used → tailor it for your purposes!





UNITED STATES MILITARY ACADEMY
WEST POINT

Methodology



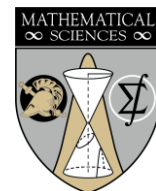


- Requirements
 - Collaborative coding environment
 - Common Python packages installed
 - Classroom for instruction
- Possible Solutions
 - [CoCalc](#)
 - [Google Colab](#)
 - Institutional Compute Cluster
 - Cloud-based (e.g., AWS, Azure)





- D/Math Population
 - Experienced Analysts / Researchers
 - Range of Python skills
- Method of Instruction
 - Demos and slides
 - Executable code to discuss
 - Jupyter Notebooks
 - Common dataset to our organization:
course-end feedback





UNITED STATES MILITARY ACADEMY
WEST POINT

Lesson Plan





Lesson Plan

1. NLP Tasks & Terminology; Preprocessing
2. Basic Lexical Analysis & Word Representations
3. Word Embeddings & Model Training
4. SOTA models & Python Packages ([Hugging Face](#))
5. Additional packages & Ongoing Research in Department





UNITED STATES MILITARY ACADEMY
WEST POINT

Lesson 1: NLP Tasks & Terminology; Preprocessing





Common NLP Tasks

Syntax

Morphology

Word Segmentation

Part-of-Speech Tagging

Parsing

Constituency

Dependency

Discourse

Summarization

Coreference Resolution

Semantics

Sentiment Analysis

Topic Modelling

Named Entity Recognition (NER)

Relation Extraction

Word Sense Disambiguation

Natural Language Understanding (NLU)

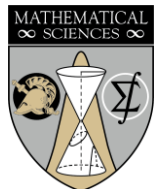
Natural Language Generation (NLG)

Machine Translation

Entailment

Question Answering

Language Modelling





Common NLP Tasks

Syntax

Morphology

Word Segmentation

Part-of-Speech Tagging

Parsing

Constituency

Dependency

Discourse

Summarization

Coreference Resolution

Semantics

Sentiment Analysis

Topic Modeling

Question Answering

Language Modelling

- Really slow wait. Took forever to get food.
- + Freshest ingredients ever. New favorite restaurant. Will be back!
- Found a hair in the food. Horrible.
- + Waited 6 months to get a reservation at this place. Totally worth it.





Tokenization

“I particularly thought getting up to the boards was helpful.”

→ [I, particularly, thought, getting, up, to, ...]

```
1 # tokenize our corpus
2 tokens = helpful.apply(lambda x: nltk.word_tokenize(str(x)))
3 print("NLTK Tokenization: \n", tokens.head())
4 tokens = helpful.apply(lambda x: TextBlob(str(x)).words)
5 print("TextBlob Tokenization:\n",tokens.head())
```

Python

NLTK Tokenization:

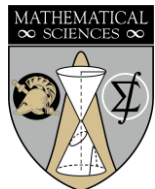
```
0  [!, !, !, 12354, In, class, examples, were, gr...
1  [I, particularly, thought, getting, up, to, th...
2  [every, time, we, would, go, to, boards, we, w...
3  [humor, that, motivated, learning]
4  [During, almost, every, class, period, ,, we, ...
```

Name: text, dtype: object

TextBlob Tokenization:

```
0  [12354, In, class, examples, were, grate, for,...
1  [I, particularly, thought, getting, up, to, th...
2  [every, time, we, would, go, to, boards, we, w...
3  [humor, that, motivated, learning]
4  [During, almost, every, class, period, we, got...
```

Name: text, dtype: object





UNITED STATES MILITARY ACADEMY
WEST POINT

Lesson 2:

Basic Lexical Analysis &

Word Representations





Basic Sentiment Analysis

“Mr. XXX is the best teacher in this course.”

→ {'neg': 0.0, 'neu': 0.656, 'pos': 0.344, 'compound': 0.6369}

“Mr. XXX is not the best teacher in this course.”

→ {'neg': 0.272, 'neu': 0.728, 'pos': 0.0, 'compound': -0.5216}

```
1 from vaderSentiment.vaderSentiment import SentimentIntensityAnalyzer #Valence Aware Dictionary for sEnt
2 analyzer = SentimentIntensityAnalyzer()
3 for sentence in blob.sentences:
4     print(sentence, analyzer.polarity_scores(sentence))
```

Python

Mr. XXX is the best teacher in this course. {'neg': 0.0, 'neu': 0.656, 'pos': 0.344, 'compound': 0.6369}

Mr. XXX is not the best teacher in this course. {'neg': 0.272, 'neu': 0.728, 'pos': 0.0, 'compound': -0.5216}



Simplistic Summarization: Keywords

What do students think the most helpful classroom experiences include?

```
1 import random
2 nouns = list()
3 blob = TextBlob(allResponses)
4 for word, tag in blob.tags:
5     if tag == 'NN': # Try noun? NN
6         nouns.append(word.lemmatize())
7
8 print("Students think that helpful classroom experiences include:")
9 for item in random.sample(nouns,10):
10     print(item)
```

```
Students think that helpful classroom experiences include:
room
use
instructor
assignment
test
board
```





UNITED STATES MILITARY ACADEMY
WEST POINT

Lesson 3: Word Embeddings & Model Training





Task: find similar comments based on cosine similarity of word2vec embeddings

Input: doing practice problems and having class examples were very helpful





Lesson 3: Word Embeddings & Model Training

Input: doing practice problems and having class examples were very helpful

**** word2vec Recommendations ****

- the various class practice problems were especially helpful
- doing problems in class was very helpful for my learning
- this path course offered a lot of opportunities to go to the boards during class and practice class examples this was extremely helpful to not only practice problems but to also have the instructor critique our work
- doing board problems were very helpful
- working through problems as a class was very helpful to me



UNITED STATES MILITARY ACADEMY
WEST POINT

Lesson 4: SOTA Models & Python Packages





Sentiment Analysis with a Pre-Trained LLM (SOTA)

```
1 MODEL_NAME = f"cardiffnlp/twitter-roberta-base-sentiment-latest"
2 device = torch.device(0 if torch.cuda.is_available() else -1)
3
4 classifier = pipeline("sentiment-analysis",model = MODEL_NAME, device = device) # device should make
5
6 posCom = "MA104 is my favorite class -- I especially enjoy the self-study days."
7 negCom = "This was the worst class I've ever had. The work was tedious and not related to my major."
8 print()
9 print("Comment: ", posCom,"\n",classifier(posCom))
10 print("Comment: ", negCom,"\n",classifier(negCom))
```



```
Comment: MA104 is my favorite class -- I especially enjoy the self-study days.
[{'label': 'positive', 'score': 0.9886686205863953}]
Comment: This was the worst class I've ever had. The work was tedious and not related to my major.
[{'label': 'negative', 'score': 0.9576689004898071}]
```



Where do I get SOTA models? Benchmark datasets?

The screenshot displays the Hugging Face Model Hub interface. The top navigation bar includes the Hugging Face logo, a search bar, and links for Models, Datasets, Spaces, Docs, Solutions, Pricing, Log In, and Sign Up. The left sidebar shows a 'Tasks' menu with categories like Multimodal, Computer Vision, and Natural Language Processing. The main content area is titled 'Models 82' and shows a search for 'rotten tomatoes'. It lists several models for sentiment classification, including:

- textattack/bert-base-uncased-rotten-tomatoes**: Text Classification, Updated May 20, 2021, 1.33k downloads, 1 heart.
- Movasaghi/finetuning-sentiment-rottentomatoes**: Text Classification, Updated May 29, 195 downloads.
- textattack/xlnet-base-cased-rotten-tomatoes**: Text Generation, Updated Jul 6, 2020, 136 downloads.
- textattack/distilbert-base-uncased-rotten-tomatoes**: Text Classification, Updated Jul 6, 2020, 67 downloads.
- textattack/roberta-base-rotten-tomatoes**: Text Classification, Updated May 20, 2021, 39 downloads.
- muhtasham/tiny-mlm-rotten_tomatoes-custom-tokenizer**: Fill-Mask, Updated Jan 7, 35 downloads, 1 heart.
- textattack/albert-base-v2-rotten-tomatoes**: Text Classification, Updated Jul 6, 2020, 29 downloads.
- RJZauner/distilbert_rotten_tomatoes_sentiment_class...**: Text Classification, Updated Aug 24, 2022, 27 downloads, 1 heart.
- klin1/rotten-tomatoes-model**: Text Classification, Updated Apr 18, 2022, 13 downloads.
- textattack/bert-base-uncased-rotten_tomatoes**: Fill-Mask, Updated May 20, 2021, 10 downloads.
- textattack/albert-base-v2-rotten_tomatoes**: Fill-Mask, Updated Jun 25, 2020, 7 downloads.
- mrm8488/distilroberta-finetuned-rotten_tomatoes-sen...**: Text Classification, Updated Aug 26, 2021, 5 downloads.
- AdapterHub/roberta-base-pf-rotten_tomatoes**: Text Classification, Updated Nov 24, 2021, 4 downloads.
- AdapterHub/bert-base-uncased-pf-rotten_tomatoes**: Text Classification, Updated Nov 24, 2021, 3 downloads.
- liliaciolite/rotten_tomatoes_sentiment_analysis**: Text Classification, Updated Jun 23, 3 downloads.
- xianzhew/distilbert-base-uncased_rotten_tomatoes**: Text Classification, Updated Nov 21, 2022, 2 downloads.

Model Hub: <https://huggingface.co/models>



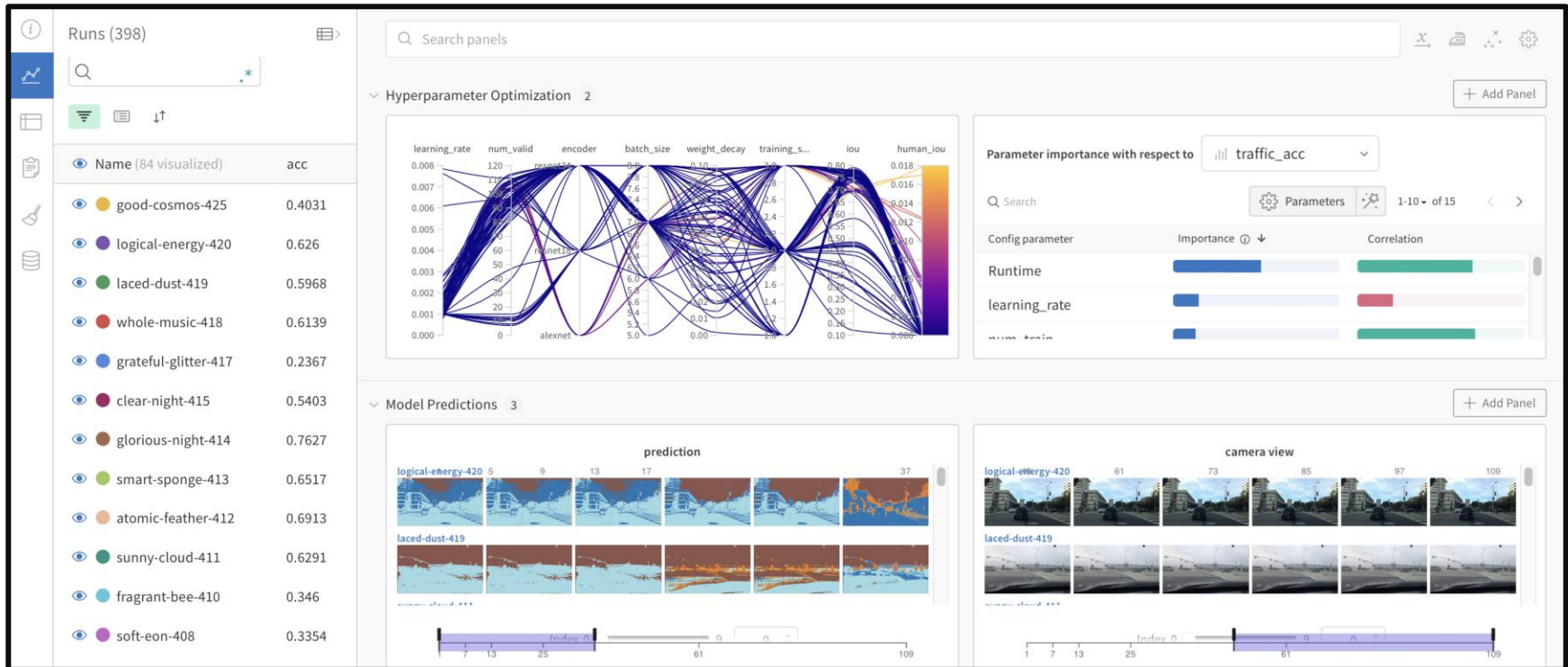
UNITED STATES MILITARY ACADEMY
WEST POINT

Lesson 5: Additional Packages & Ongoing Research in Department



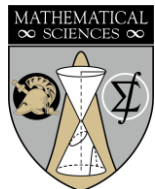


Why use Weights & Biases? **It automates your record keeping!**



[W&B Docs | Weights & Biases Documentation \(wandb.ai\)](https://docs.wandb.ai/)

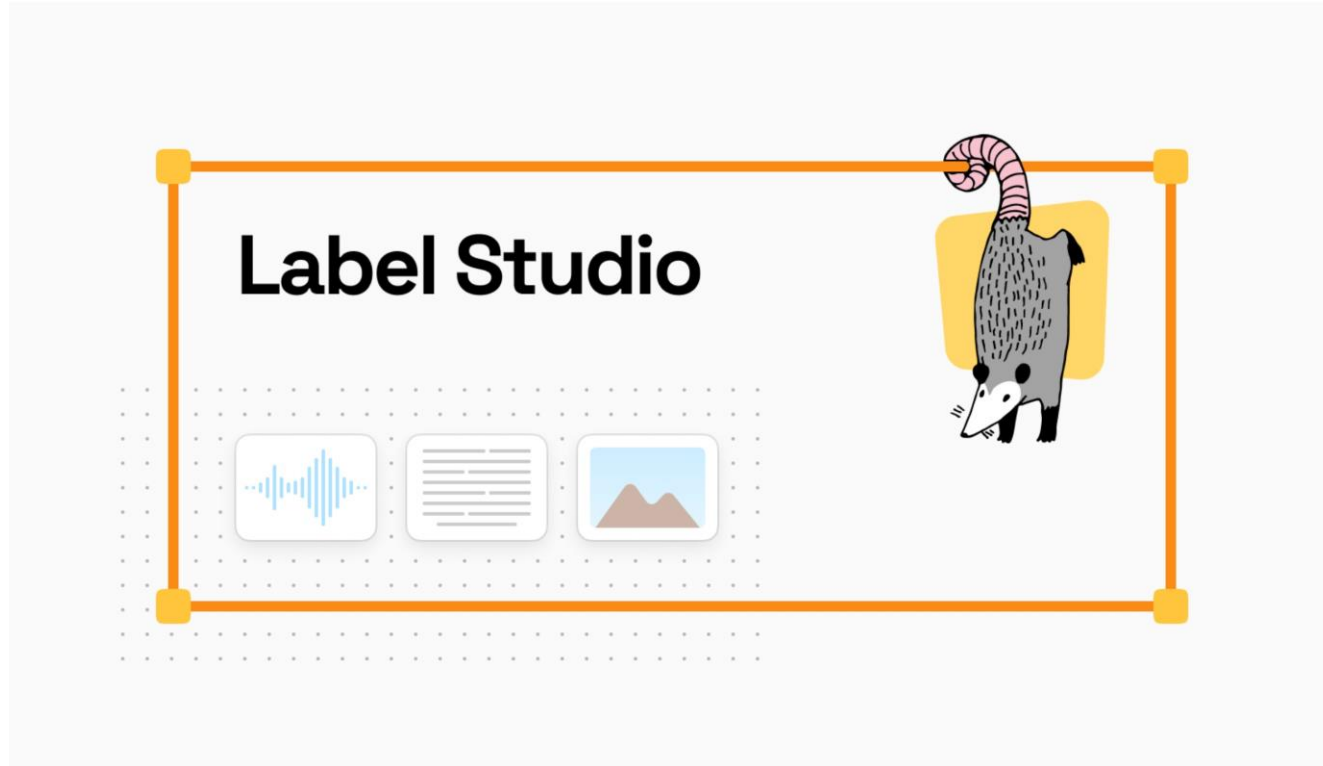
[Intro to Weights & Biases.ipynb - Colaboratory \(google.com\)](https://colab.research.google.com/github/WeightsandBiases/Intro-to-Weights-&-Biases/blob/master/Intro%20to%20Weights%20%26%20Biases.ipynb)





Lesson 5: Packages & Ongoing Research

Why use Label Studio? **It's an intuitive environment to label data!**



[Open Source Data Labeling | Label Studio](https://labelstudio.com/)



UNITED STATES MILITARY ACADEMY
WEST POINT

Best Practices





- Use a dataset familiar for your faculty
- Consider audience coding skills and desired coding environment
- Include orientation to using Graphics Processing Units (GPUs) for large language models (LLM)
- Provide executable code – live coding is slow





UNITED STATES MILITARY ACADEMY
WEST POINT

What are your next steps?





UNITED STATES MILITARY ACADEMY
WEST POINT

Resources for YOU

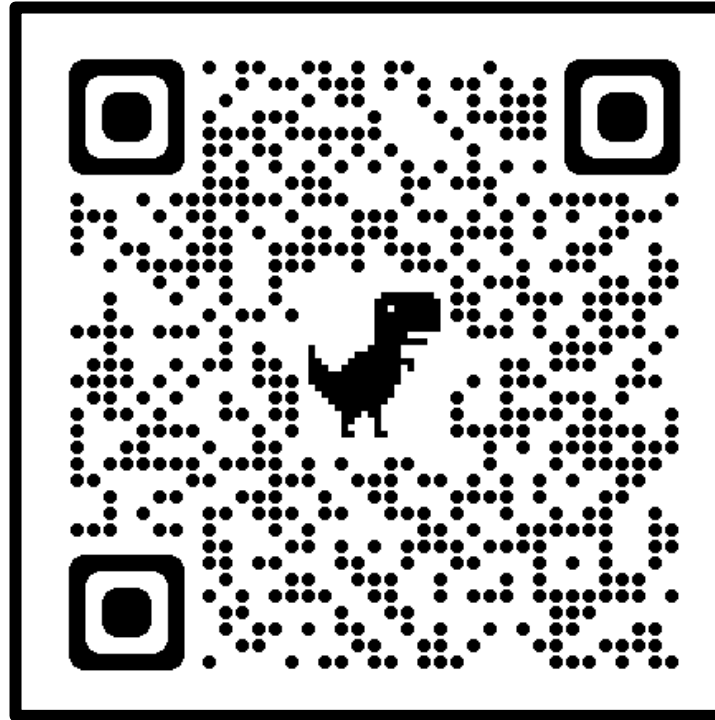




- **The [Hugging Face](#) community** (particularly their [Transformers](#) library) is one of the most popular communities/platforms for NLP/AI/ML. They have great tutorials and readable documentation. Hugging Face models and datasets are designed to work well with and/or build on other popular NLP packages [[nltk](#), [AllenNLP](#), [spacy](#)].
- **CS287: Deep Learning for NLP** is a course I took in grad school. It was only taught for one semester, but the lecture slides are all available online (for now) [[link - scroll to "Lectures"](#)]. There's also a [Supplemental Resources](#) page for NLP.
- If you want straightforward explanations of how the most used NLP models work, [Jay Alammar has a blog](#) with great visuals and explanations for word2vec (2013) and BERT (2018+) toward the bottom of his blog posts.
- Sebastian Ruder [[his site](#)] [[his Google Scholar](#)] has a great blog that covers new ideas in NLP, NLP Progress, his highlights from NLP conferences...etc. In 2018 he gave a brief timeline of NLP. <https://ruder.io/a-review-of-the-recent-history-of-nlp/>.



Questions?



Scan QR code for [GitHub repository](#) access!

Email: john.scudder@westpoint.edu





UNITED STATES MILITARY ACADEMY
WEST POINT





UNITED STATES MILITARY ACADEMY
WEST POINT

Backup Slides





- 15 March @ 1500 in TH120
 - Preprocessing, Tokenization, Part of Speech Tagging
- 22 March @ 1000 in TH120
 - Basic lexical analysis – sentiment, cosine similarity
 - Representations – N-Grams, Bag of Words
- 5 April @ 1100 in TH120
 - TF-IDF & word2vec
 - Visualizations
 - Training models
- 19 April @ 1000 in TH120
 - State-of-the-art models: transformers, BERT
 - Hugging Face packages
- 9 May @ 1000 in TH120
 - 2x packages: Weights & Biases + Label Studio
 - Hone your craft! Ongoing D/Math research with NLP





UNITED STATES MILITARY ACADEMY
WEST POINT

LESSON 1





UNITED STATES MILITARY ACADEMY
WEST POINT

Common Tasks





Syntax

Morphology

Word Segmentation

Part-of-Speech Tagging

Parsing

Constituency

Dependency

Discourse

Summarization

Coreference Resolution

Semantics

Sentiment Analysis

Topic Modelling

Named Entity Recognition (NER)

Relation Extraction

Word Sense Disambiguation

Natural Language Understanding (NLU)

Natural Language Generation (NLG)

Machine Translation

Entailment

Question Answering

Language Modelling





Syntax

Morphology

Word Segmentation

Part-of-Speech Tagging

Parsing

Constituency

Dependency

Discourse

Summarization

Coreference Resolution

Semantics

Sentiment Analysis

Topic Modelling



“The U.S. secretly modified the Himars rocket launchers it gave Ukraine to keep Kyiv from firing missiles into Russia.”

Question Answering

Language Modelling

NLU)
)



Syntax

Morphology

Word Segmentation

Part-of-Speech Tagging

Parsing

Constituency

Dependency

Discourse

Summarization

Coreference Resolution

Semantics

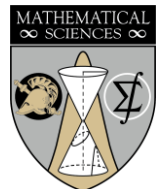
Sentiment Analysis

Topic Modeling

- − Really slow wait. Took forever to get food.
- + Freshest ingredients ever. New favorite restaurant. Will be back!
- − Found a hair in the food. Horrible.
- + Waited 6 months to get a reservation at this place. Totally worth it.

Question Answering

Language Modelling





Syntax

Morphology

Word Segmentation

Part-of-Speech Tagging

Parsing

Constituency

Dependency

Discourse

Summarization

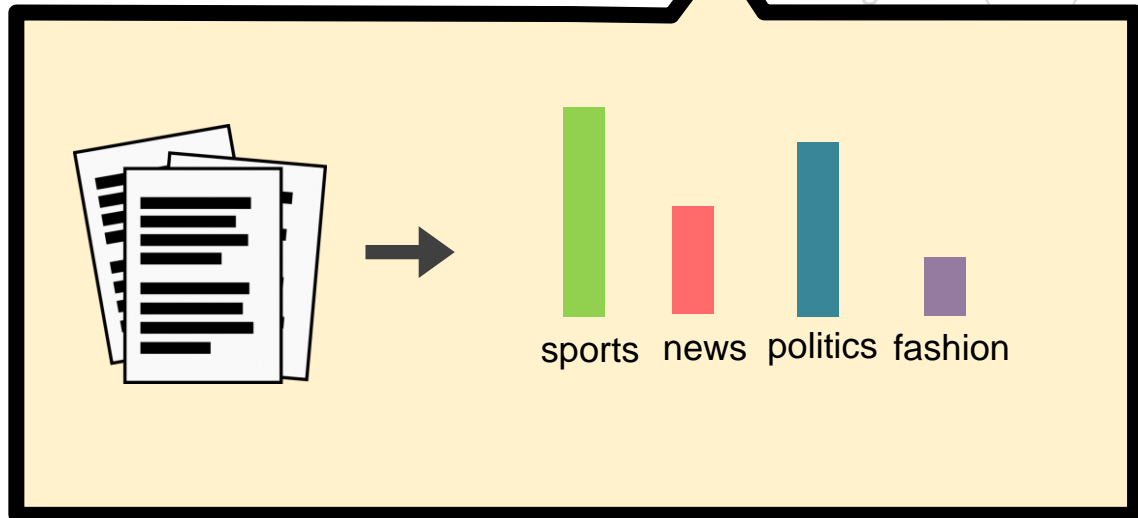
Coreference Resolution

Semantics

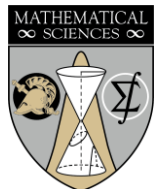
Sentiment Analysis

Topic Modelling

Named Entity Recognition (NER)



Language Modelling





Syntax

Morphology

Word Segmentation

Part-of-Speech

Parsing

Constituency

Dependency

“Alexa, play White Christmas by Bing Crosby”



“Alexa, **play** **White Christmas** by **Bing Crosby**”
INTENT **SONG** **ARTIST**

Named Entity Recognition (NER)

Semantic

Discourse

Summarization

Coreference Resolution

Natural Language Understanding (NLU)

Natural Language Generation (NLG)

Machine Translation

Entailment

Question Answering

Language Modelling





Syntax

Morphology

Word Segmentation

Part-of-Speech Tagging

Parsing

Constituency

Dependency

Semantics

Sentiment Analysis

Topic Modelling

Named Entity Recognition (NER)

Dzień dobry

POLISH



Good morning

ENGLISH

Discourse

Summarization

Coreference Resolution

Machine Translation

Entailment

Question Answering

Language Modelling





UNITED STATES MILITARY ACADEMY
WEST POINT

Hands-on Coding





UNITED STATES MILITARY ACADEMY
WEST POINT

LESSON 2





UNITED STATES MILITARY ACADEMY
WEST POINT

Terminology





How do we get *any* system to process, “understand”, leverage language?

- **Representation**: how we transform symbolic meaning (e.g., words, signs, braille, speech audio) into something the computer can use
- **Modelling**: given these represented symbols, how we use them to model the task at hand





Representing images is relatively easy. There is a meaningful relationship between the byte values and color.



170 **33** **71**
r **g** **b**



255 **33** **71**
r **g** **b**





In contrast – words are represented by strings and there are no meaningful relationships between the byte values and language.

h a t e

68	61	74	65
----	----	----	----

a t e

61	74	65
----	----	----

hate and **ate**. No relation but similar byte values.

H a t

48	61	74
----	----	----

h a t

68	61	74
----	----	----

Hat and **hat**. Identical concept but different byte values.





A **Language Model** estimates the probability of any sequence of words

Let X = “Jack was late for class”

w_1 w_2 w_3 w_4 w_5

$$P(X) = P(\text{“Jack was late for class”}) = 2.3 \times 10^{-14}$$





Scenario: assume we have a finite vocabulary V

V^* represents the **infinite set** of strings/sentences that we could construct

e.g., $V^* = \{a, a \text{ dog}, a \text{ frog}, \text{dog } a, \text{dog dog}, \text{frog dog}, \text{frog } a \text{ dog}, \dots\}$

Data: we have a training set of sentences $x \in V^*$

Problem: estimate a probability distribution:

$$\sum_{x \in V^*} p(x) = 1$$

$$p(\textit{the}) = 10^{-2}$$

$$p(\textit{the}, \textit{sun}, \textit{okay}) = 2.5 \times 10^{-13}$$

$$p(\textit{waterfall}, \textit{the}, \textit{icecream}) = 3.2 \times 10^{-18}$$





A word **token** is a specific occurrence of a word in a text

A word **type** refers to the general form of the word, defined by its lexical representation

If our corpus were just “I ran and ran and ran”, you’d say we have:

- 6x word **tokens** [I, ran , and , ran , and , ran]
- 3x word **types**: {I, ran, and}





Bigram Model Example – condition each word on its immediate predecessor

Let $X = \text{"Jack was late for class"}$

w_1 w_2 w_3 w_4 w_5

$$P(w' | w) = P("w, w'") = \frac{n_{w, w'}(d)}{n_{w, w^*}(d)}$$

$$P(\text{class} | \text{for}) = P(\text{for, class}) = \frac{12}{3,000}$$

Let's say our corpus d has 100,000 words

word	# occurrences
Jack	15
was	1,000
late	400
for	3,000
class	350

$$n_{w^*}(d) = 100,000$$

$n_{w, w'}(d) = \#$ of times words w and w' appear together as a bigram in d

$n_{w, w^*}(d) = \#$ of times word w is the first token of a bigram in d





Simple idea: let's represent each document as a feature vector, which can serve as the input to any of your favorite supervised ML models

Let's say our dataset's entire *vocabulary* is just 10 words. Each unique word can have its own dimension (feature index). Each document's vector has a 1 if the word is present. Otherwise, 0.

e.g., “the dog jumped” is represented as

[1	1	0	0	0	0	1	0	0	0]
	dog	the	quick	went	brown	a	jumped	fast	over	store	





Imagine a document is a sports broadcast transcript, which concerns a few teams but mostly discusses the local home team, the Cubs. We have no indication of *how much* the document is about the Cubs.

[1	1	1	1	0	1	0	0	0	1]
	baseball	chicago	cubs	the	wrigley	padres	shohei	mvp	homerun	crowd	

Using a count-based approach, now we can see that it's much more about the **Chicago Cubs** than the **Padres**.

[2	9	17	8	0	2	0	0	0	2]
	baseball	chicago	cubs	the	wrigley	padres	shohei	mvp	homerun	crowd	





Weaknesses:

- Flattened view of the document
- Context-insensitive (“the horse ate” = “ate the horse”)
- Curse of Dimensionality (vocab could be over 100k)
- Orthogonality: no concept of semantic similarity at the word-level
 - e.g., $d(\text{dog}, \text{cat}) = d(\text{dog}, \text{chair})$





UNITED STATES MILITARY ACADEMY
WEST POINT

Hands-on Coding





UNITED STATES MILITARY ACADEMY
WEST POINT

LESSON 3





TF (term frequency) = f_{w_i} = # times word w_i appeared in the document

IDF (inverse document frequency) = $\log \left(\frac{\# \text{ docs in corpus}}{\# \text{ docs containing } w_i} \right)$

$$\text{TFIDF} = f_{w_i} * \log \left(\frac{\# \text{ docs in corpus}}{\# \text{ docs containing } w_i} \right)$$

- TF-IDF rewards “rare” words that frequently appear in a few documents. Words that appear in every document in the corpus have a TF-IDF score = 0.
- Hans Peter Luhn (1957) credited with TF; Karen Spärck Jones (1972) credited with IDF.



Thou shalt not make a machine in the likeness of a human mind

thou	shalt	not	make	a	machine	in	the	...
------	-------	-----	------	---	---------	----	-----	-----

thou	shalt	not	make	a	machine	in	the	...
------	-------	-----	------	---	---------	----	-----	-----

thou	shalt	not	make	a	machine	in	the	...
------	-------	-----	------	---	---------	----	-----	-----

thou	shalt	not	make	a	machine	in	the	...
------	-------	-----	------	---	---------	----	-----	-----

thou	shalt	not	make	a	machine	in	the	...
------	-------	-----	------	---	---------	----	-----	-----

input word	target word
not	thou
not	shalt
not	make
not	a
make	shalt
make	not
make	a
make	machine
a	not
a	make
a	machine
a	in
machine	make
machine	a
machine	in
machine	the
in	a
in	machine
in	the
in	likeness

Graphic:

<https://jalammar.github.io/illustrated-word2vec/>

Original Papers:

<https://arxiv.org/abs/1301.3781>

<https://arxiv.org/abs/1309.4168>



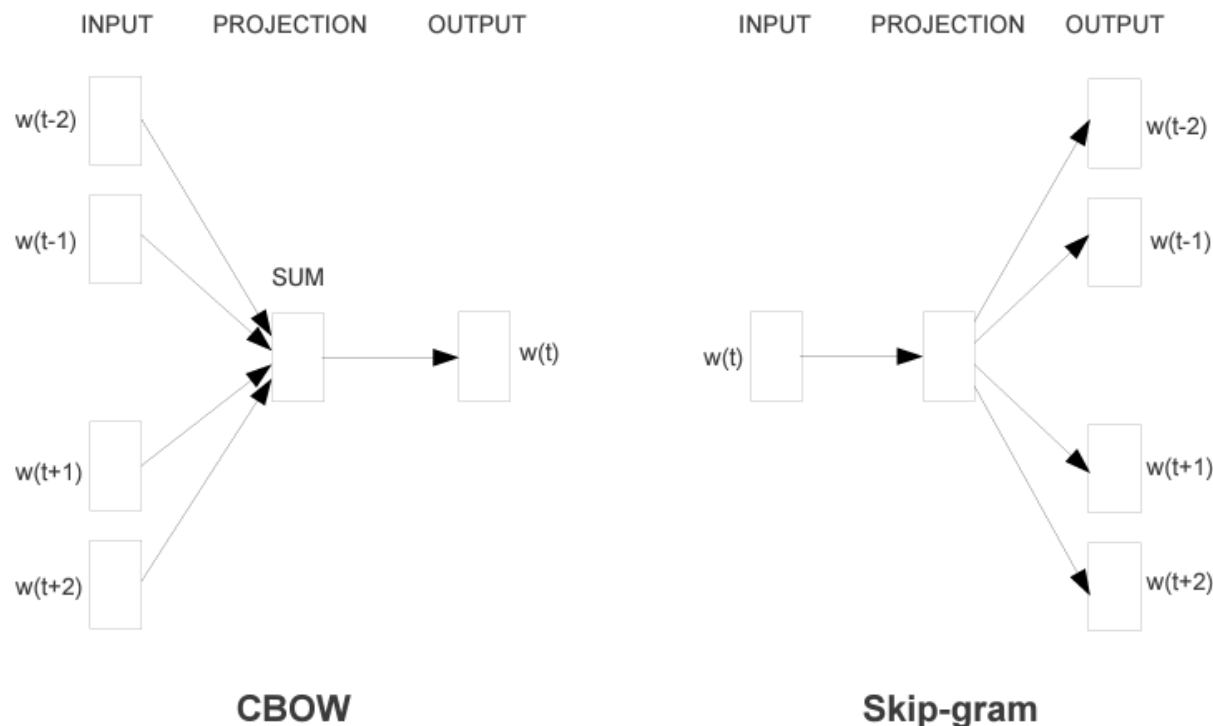


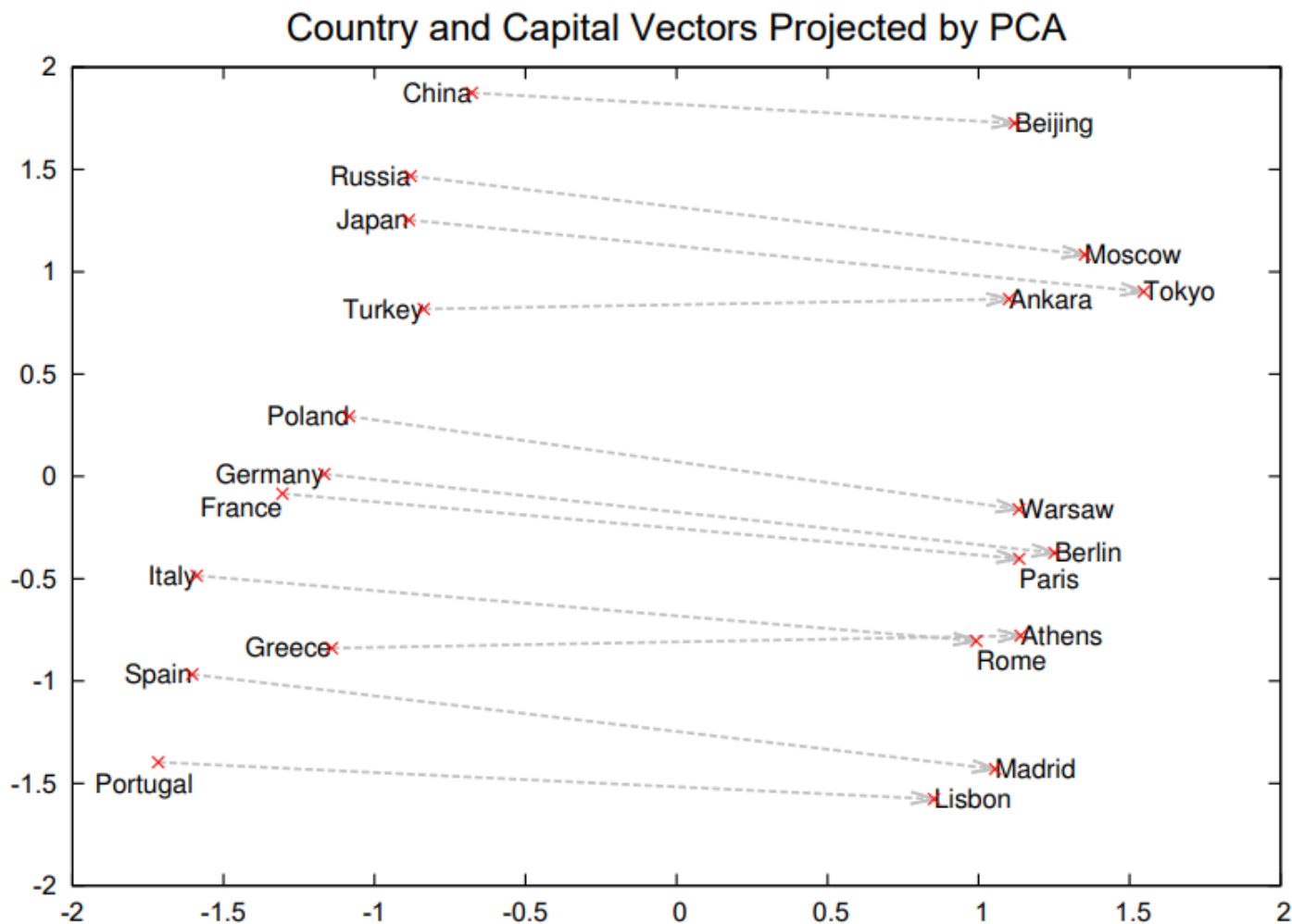
Figure 1: New model architectures. The CBOW architecture predicts the current word based on the context, and the Skip-gram predicts surrounding words given the current word.

Original Papers:

<https://arxiv.org/abs/1301.3781>

<https://arxiv.org/abs/1309.4168>





Original Papers:

<https://arxiv.org/abs/1301.3781>

<https://arxiv.org/abs/1309.4168>

<https://arxiv.org/abs/1310.4546>





Document 1	Document Vector Embedding --->	0.4, -0.0375, 0.725, ...
Word Embedding	the	0.4, 0.25, 0.1, ...
Word Embedding	board	0.9, -0.42, 0.01, ...
Word Embedding	was	0.6, 0.34, -0.21, ...
Word Embedding	great	-.3, -0.32, 3, ...

Document 2	Document Vector Embedding --->	0.285, 0.23, -0.34, ...
Word Embedding	the	0.4, 0.25, 0.1, ...
Word Embedding	elephant	0.34, 0.23, -0.3, ...
Word Embedding	was	0.6, 0.34, -0.21, ...
Word Embedding	big	-0.2, 0.1, 0.07, ...



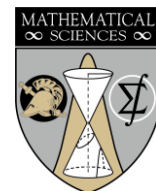


```
: print("Length of the vector: ",len(wv['board']), "\n Embedding for BOARD: \n",wv['board'])
```

Length of the vector: 300

Embedding for BOARD:

```
[-0.14453125 -0.25976562 -0.01611328 -0.01074219 -0.01281738 -0.34765625  
 0.10839844 0.00340271 0.07080078 0.04199219 0.0456543 -0.14160156  
-0.03808594 -0.19335938 -0.30273438 0.09619141 0.0703125 -0.11425781  
-0.02709961 0.01306152 -0.09863281 0.22070312 0.00118256 0.1328125  
0.02783203 0.14453125 -0.21386719 0.30664062 -0.20117188 -0.29101562  
0.07080078 -0.07861328 -0.07958984 -0.06738281 0.17675781 -0.23730469  
0.171875 0.31445312 0.13378906 -0.12109375 -0.09423828 0.13671875  
0.0390625 -0.09619141 0.07666016 -0.12695312 0.19140625 -0.04907227  
0.04589844 0.21679688 -0.00778198 0.08886719 0.05566406 -0.0859375  
-0.08349609 0.0559082 -0.17382812 0.04443359 0.10644531 0.00653076  
0.09863281 0.10205078 -0.15429688 0.09130859 0.06933594 0.06030273  
0.00701904 0.19433594 0.140625 0.19238281 0.01043701 -0.140625  
0.33984375 0.09326172 0.00125122 0.19628906 0.03613281 -0.09277344  
0.25 0.08496094 0.11083984 -0.05419922 0.09912109 0.09716797  
0.28710938 -0.11914062 -0.1640625 -0.05810547 0.15820312 0.11816406  
0.30078125 -0.00357056 0.02087402 0.20703125 0.0378418 -0.07324219  
-0.24707031 0.22558594 -0.06201172 0.12011719 0.0612793 0.00628662  
0.0267334 0.22949219 -0.06494141 -0.09521484 -0.02026367 -0.171875  
0.00494385 -0.2734375 0.0300293 -0.18554688 0.04321289 -0.22070312  
0.3203125 -0.10986328 0.12011719 0.07666016 0.22070312 -0.16503906  
-0.06176758 -0.11962891 -0.15039062 0.0859375 -0.13769531 -0.11962891  
-0.25976562 -0.08496094 0.05151367 0.04101562 0.06591797 -0.04907227  
-0.18457031 -0.11865234 0.02941895 0.26171875 0.04516602 0.11425781  
0.06640625 0.14160156 -0.12597656 -0.04467773 0.06005859 -0.08789062  
0.04150391 -0.02868652 -0.14160156 -0.0378418 0.34570312 0.17382812  
-0.08154297 0.05883789 -0.23632812 0.24511719 -0.06542969 0.0144043  
0.12988281 0.08105469 0.20019531 0.14941406 -0.0559082 0.00105286  
0.16308594 0.03198242 0.03173828 -0.18847656 -0.02893066 0.02001953  
0.26171875 0.12792969 0.10400391 -0.00041008 0.00915527 -0.13085938  
0.21073656 0.21679688 0.10000000 0.01055000 0.00531000 0.13193500
```





UNITED STATES MILITARY ACADEMY
WEST POINT

Hands-on Coding





LESSON 4

[Link to Behemoth Google Colab Notebook!](#)

Separate Notebooks for Exploring vs. Training a Model

- [Part 1: Explore Hugging Face Datasets and Models](#)
- [Part 2: Train a Model](#)





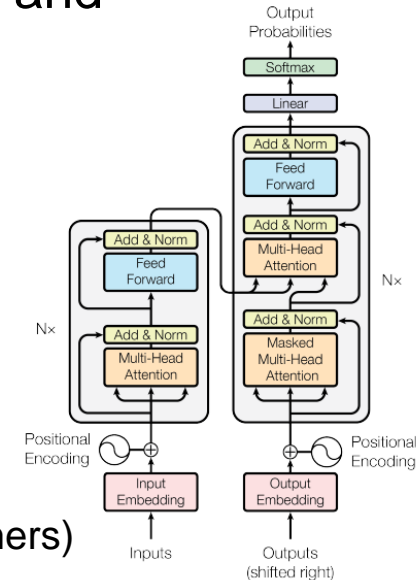
UNITED STATES MILITARY ACADEMY
WEST POINT®

Transformers ...and BERT





- Transformers are an **architecture** that enables accurate and efficient transfer learning
 - Keeps track ONLY of model weights, biases, configuration
 - Customizable
- Encoder / Decoder
 - Encoder (extract meaning of input – all at once):
 - Document classification
 - Named entity recognition
 - Extractive Q&A
 - Ex: BERT (Bidirectional Encoder Representations from Transformers)
 - Decoder (access meaning of input – sequential):
 - Next word prediction
 - Text Generation
 - Encoder / Decoder (extract features, then produce sequential output):
 - Text Summarization
 - Translation





What do Transformers keep track of?

```
fineTunedModel /
├── config.json
└── pytorch_model.bin

▼ root:
  _name_or_path: "cardiffnlp/twitter-roberta-base-sentiment-latest"
  ▼ architectures: [ ] 1 item
    0: "RobertaForSequenceClassification"
    attention_probs_dropout_prob: 0.1
    bos_token_id: 0
    classifier_dropout: null
    eos_token_id: 2
    gradient_checkpointing: false
    hidden_act: "gelu"
    hidden_dropout_prob: 0.1
    hidden_size: 768
  ▼ id2label:
    0: "NEGATIVE"
    1: "NEUTRAL"
    2: "POSITIVE"
  initializer_range: 0.02
  intermediate_size: 3072
  ▼ label2id:
    NEGATIVE: 0
    NEUTRAL: 1
    POSITIVE: 2
  layer_norm_eps: 0.0001
  max_position_embeddings: 514
  model_type: "roberta"
  num_attention_heads: 12
  num_hidden_layers: 12
  pad_token_id: 1
  position_embedding_type: "absolute"
  problem_type: "single_label_classification"
  torch_dtype: "float32"
  transformers_version: "4.28.1"
  type_vocab_size: 1
  use_cache: true
  vocab_size: 50265
```

```
RobertaConfig {
  "attention_probs_dropout_prob": 0.1,
  "bos_token_id": 0,
  "classifier_dropout": null,
  "eos_token_id": 2,
  "hidden_act": "gelu",
  "hidden_dropout_prob": 0.1,
  "hidden_size": 768,
  "initializer_range": 0.02,
  "intermediate_size": 3072,
  "layer_norm_eps": 1e-12,
  "max_position_embeddings": 512,
  "model_type": "roberta",
  "num_attention_heads": 12,
  "num_hidden_layers": 12,
  "pad_token_id": 1,
  "position_embedding_type": "absolute",
  "transformers_version": "4.28.1",
  "type_vocab_size": 2,
  "use_cache": true,
  "vocab_size": 30522
}
```



- Common Architecture / Ecosystem enables **transfer learning**
 - “Reuse a pre-trained model as the starting point for a model on an adjacent task” – Layman’s Terms
 - Repurposing on a related task is MUCH quicker
- Most language models begin with an MLM, trained at GREAT expense... then add a classifier.
 - MLMs vary based on tokenizers (BPE vs. WordPiece) and architecture (number of hidden layers, size of hidden layers, etc.)
 - You can even repurpose repurposed models (what we’ll do today)
- Model Hub: <https://huggingface.co/models>



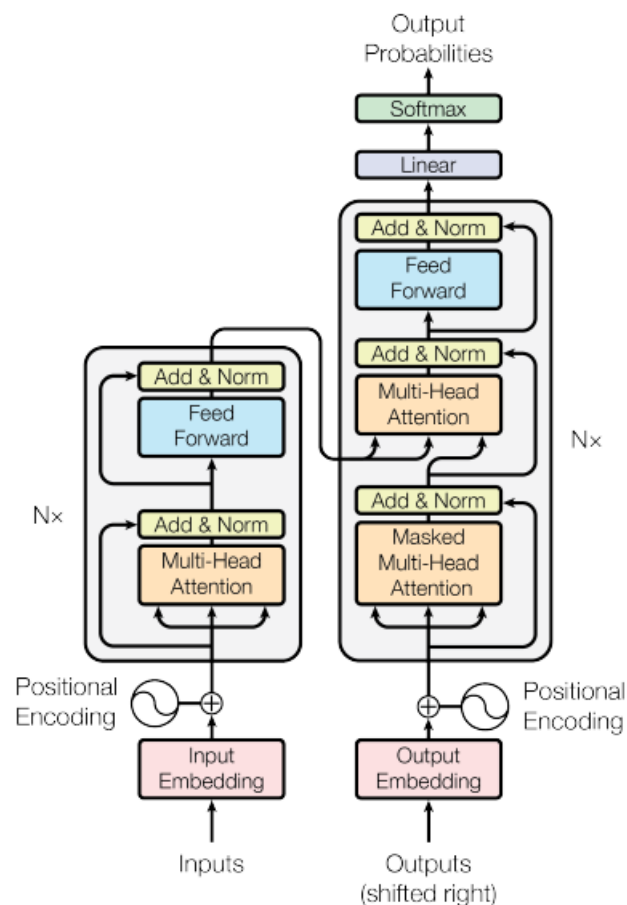


- Up to 512 tokens per sample
- Input Embedding (vocab size x embedding dim)
 $30522 \times 768 + \text{others} = 24\text{M}$
- 12 Layers: $7.1\text{M} \times 12 = 85\text{M}$
 - “Attention Head”: 7.1M
- Output Layer: 0.5M

 config.json

 pytorch_model.bin

[Computing Parameters - Stack Overflow](#)





Batches, Samples, and Epochs, Oh, my!

	workend	randomi	consent immun
1	A0017268R9SKD8U2Y3F	18097	1 sounds abso
2	A105GUJGUWPAOK	95789	1 od was high c
3	A106MX68VHW6T	50999	1 ne else has a
4	A113XP6TXPMO2	53169	1 ho i eat ex:re
5	A115PR5CQ03UJ	68906	1 less than i di
6	A11HD54H2OLIEP	55634	1 food, it isn't
7	A11LS06D7BMY99	45489	1 ds, my own la
8	A11NZXYG9AGJ30	57064	1 larent about ti
9	A120PN2KQW0VPY	31075	1 en creeping u
10	A121T4P5NKCHCS	75709	1 ces in the thr
11	A128V4VDTJJ9MS	56050	1 heat and milk
12	A12PNTH0DSUF8	95414	1 eating snack
13	A12TZ1W7		
14	A12		
15	A1		
16	A13		
17	A1		
18	A13		
19	A1		
20	A13		
21	A1		

Samples

- “Rows of data”
- documents or records – not tokens

1		1 A0017268R9SKD8U2Y3
2		
3	1	1 A0017268R9SK
4	2	
5	3	1 A0017
6	4	2 A
7	5	3
8		
9		

Batches

- Enables efficient parallelism
- Batch Size = Samples / Batches
- Updates model weights each batch

Epochs

- How many times should the model see all data?
- Forward / Backward pass



UNITED STATES MILITARY ACADEMY
WEST POINT

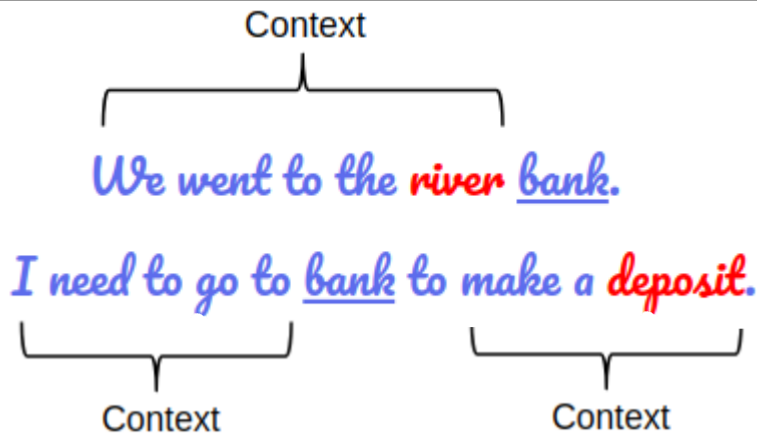
Embedding Differences





Word Embeddings vs. Contextualized Embeddings

- word embeddings: words have the **same** vector representation **regardless** of context.
- contextualized embeddings: words have a **different** vector representation **dependent** on their context



[Analytics Vidhya: Demistifying BERT.](#)

- Word2Vec: 'bank' vector is the same



- BERT: 'bank' vector is different.





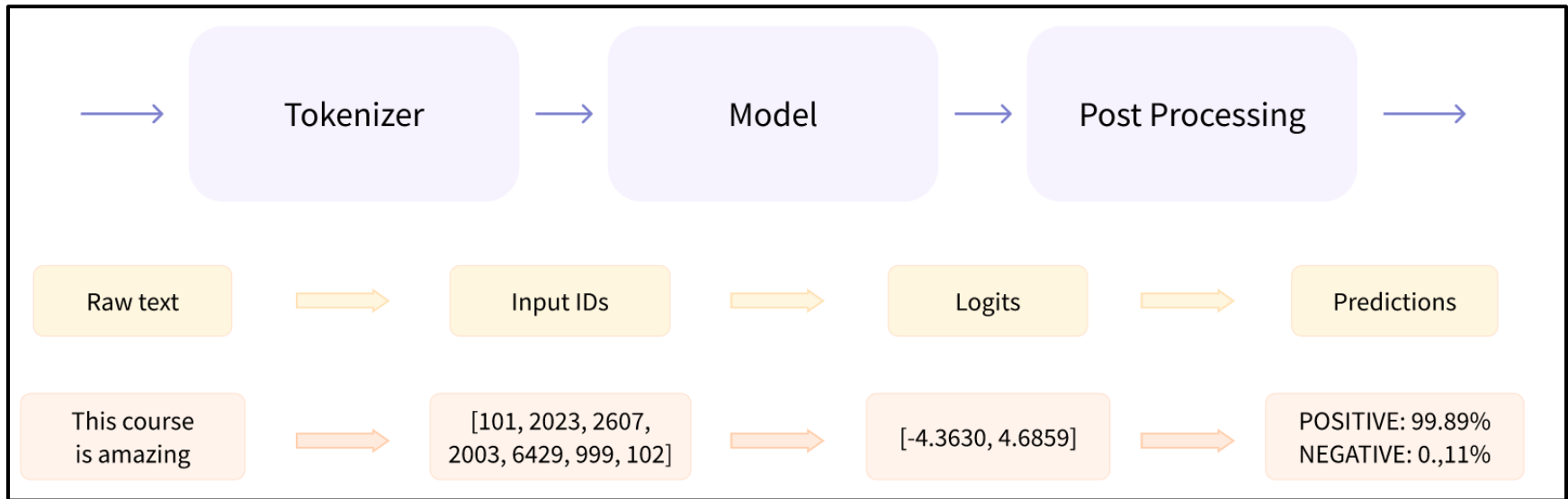
UNITED STATES MILITARY ACADEMY
WEST POINT

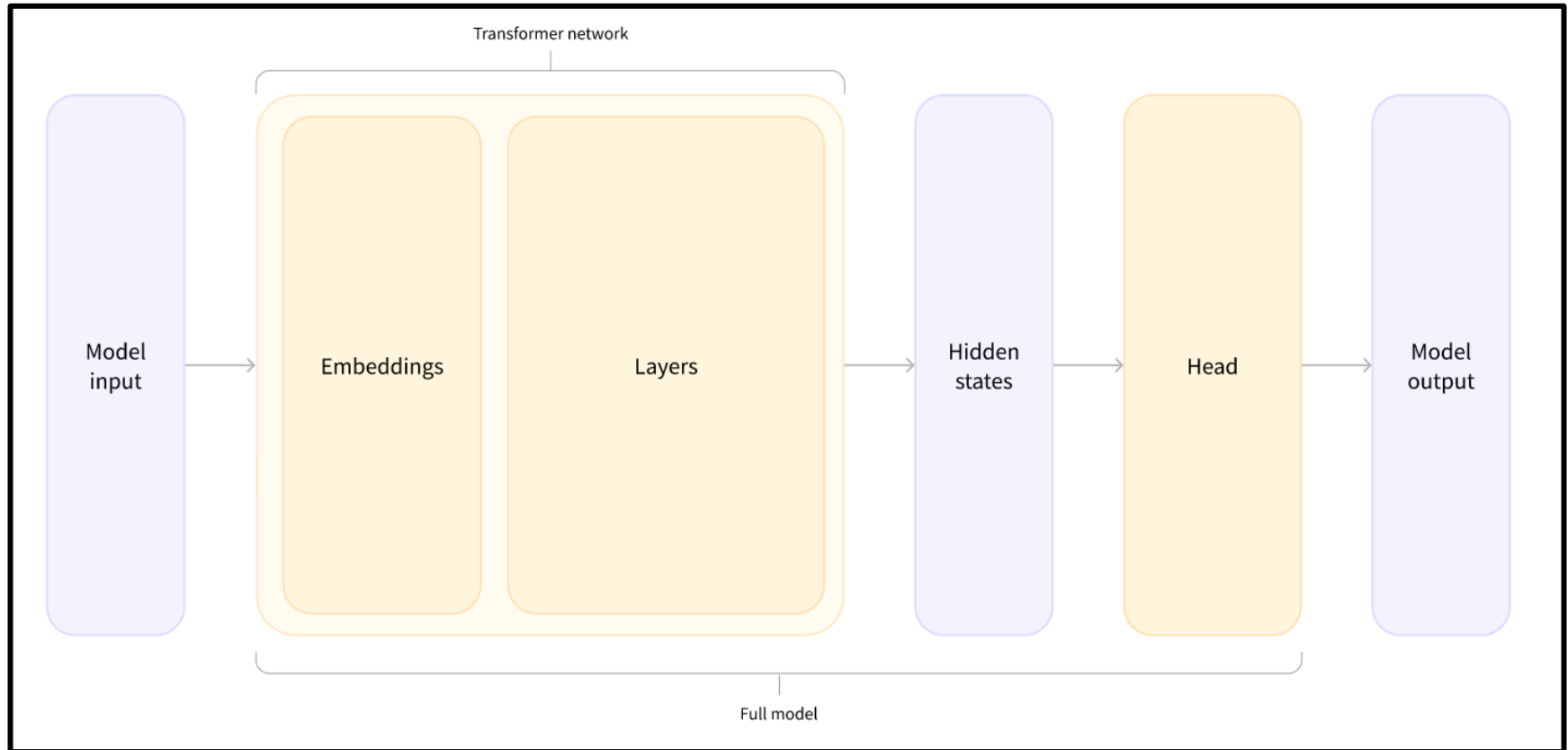
Text Processing Pipeline





Behind the pipeline - Hugging Face Course







UNITED STATES MILITARY ACADEMY
WEST POINT

Hands-on Coding





UNITED STATES MILITARY ACADEMY
WEST POINT

LESSON 5





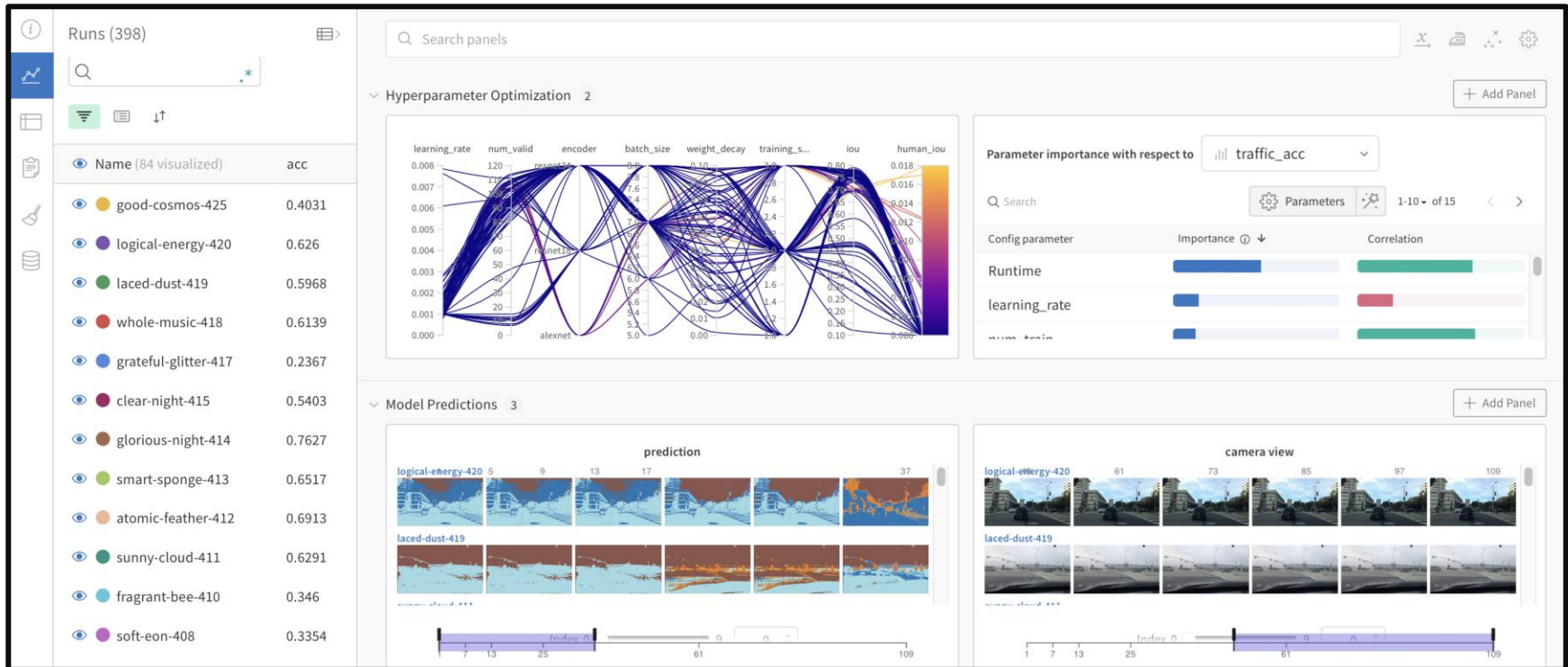
UNITED STATES MILITARY ACADEMY
WEST POINT

Weights and Biases



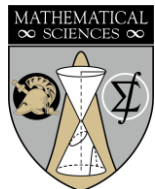


Why? **It automates your record keeping!**



[W&B Docs | Weights & Biases Documentation \(wandb.ai\)](https://docs.wandb.ai/)

[Intro to Weights & Biases.ipynb - Colaboratory \(google.com\)](https://colab.research.google.com/github/WeightsandBiases/Intro-to-Weights-and-Biases/blob/master/Intro%20to%20Weights%20&%20Biases.ipynb)





Three calls to start logging!

1. [wandb.init - Documentation](#)
 1. `name` – a label for your run
 2. `project` – where you want to group your runs
 3. `entity` – username (important for collaboration)
 4. `tags`
 5. `notes`
 6. `config` – dictionary to store information about your run
2. [wandb.config - Documentation](#)
 - a. dictionary of tracked items
 - b. `argparse` compatible
3. [wandb.log - Documentation](#) – logging. This is the best part!





- Create a Weights & Biases account
- Create a project “basic-demo”
- In terminal:
 - `pip install wandb`
 - `wandb login` → use API key from “User Settings” to login

1. Quickstart Demo from W&B

- [Quickstart Colab Demo from W&B](#)

2. Integrate with `Trainer` object from Hugging Face

- [Modified Colab demo from last lecture](#)
- [Guide from W&B](#)





UNITED STATES MILITARY ACADEMY
WEST POINT

Hands-on Coding






Fill in the blank...

Virtually ALL machine learning techniques require lots of _____ to work well.

So... for a custom project, using custom data, how do we do that, exactly?





- Many techniques – but I recommend:
<https://labelstud.io/>  **Label Studio**
- Why?
 - Distributed labelling (more than 1 at a time)
 - Now hosted on ACI's WIRE
<https://icsarl.westpoint.edu/wire>
 - Versatile
 - Text? Sure.
 - Image? Sure.
 - Audio? Sure.
 - Let's take it for a test drive!



- Topic modeling is an unsupervised machine learning technique to reduce dimensionality of your corpus.*
 - Start with 1000 unique documents
 - Group / Cluster ones that are similar
 - Think PCA for text data
- Two main techniques:
 - Lexical (LDA)
 - Transformer (BERTopic)
- Let's take this for a test drive on our course end feedback.

* An alternative to topic modeling is topic classification model (supervised)





UNITED STATES MILITARY ACADEMY
WEST POINT

Research Opportunities in D/Math!





UNITED STATES MILITARY ACADEMY
WEST POINT

Questions?





UNITED STATES MILITARY ACADEMY
WEST POINT

Additional Slides





1. Old Demos

- a. Colab Notebooks:
 - i. [Fine-tuning a model in Google Colab](#) [training]
 - ii. [Using a fine-tuned model in Google Colab](#)
- b. Hugging Face: [Hugging Face – The AI community building the future.](#)

2. NLP

- a. **CS287: Deep Learning for NLP** is a course Jack took in grad school. It was only taught for one semester, but the lecture slides are all available online (for now) [[link - scroll to "Lectures"](#)]. There's also a helpful [Supplemental Resources](#) page for NLP.
- b. [Hugging Face](#) (particularly their [Transformers](#) library) is one of the most popular communities/platforms for NLP/AI/ML. They have great tutorials and readable documentation. Hugging Face models and datasets are designed to work well with and/or build on other popular NLP packages [[nlk](#), [AllenNLP](#), [spacy](#)].

3. Python

- a. If you want to use Python without downloading it – [Google Colab](#) is a free online version of a Jupyter notebook with incredible functionality.
- b. Jack uses [VSCode](#) as my IDE for Python.
- c. Most individuals use the Anaconda Distribution for Windows. You can find instructions here: [[link](#)]
- d. Jack paid ~\$15 back in 2020 for a Udemy course called "Complete Python Bootcamp". [[link](#) – the video to install Python is available as a free preview Python Setup>Installing Python (Step by Step)]. The course lecture slides and materials are freely available. [[slides](#)] [[github repo w/ Jupyter Notebooks](#)].
- e. Google also has a free Python Course. [[link](#)]



1. Sebastian Ruder [[his site](#)] [[his Google Scholar](#)] has a great blog that covers new ideas in NLP, NLP Progress, his highlights from NLP conferences...etc. In 2018 he gave a brief timeline of NLP. <https://ruder.io/a-review-of-the-recent-history-of-nlp/>.
2. If you want straightforward explanations of how the most used NLP models work, [Jay Alammam has a blog](#) with great visuals and explanations for word2vec (2013) and BERT (2018+) toward the bottom of his blog posts.

