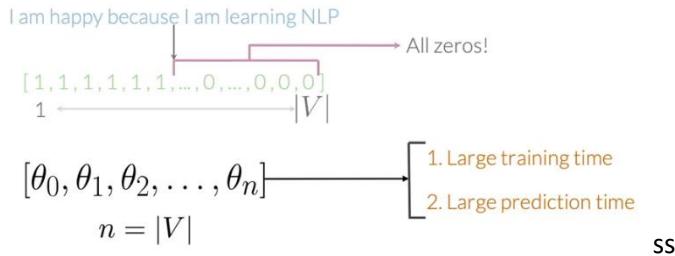


Sentiment Analysis

- Classify positive and negative sentiment
- Process:
 - o Logistic Regression classifier to classify the sentiment of a sentence from its extracted features.

Vocabulary & Feature Extraction



- Vector representation:
 - o Sparse representation
- Process:
 - o Build a vocabulary
 - Go through all words in tweets
 - Save each unique word encountered
 - o Feature extraction:
 - Check the present of vocabulary words in a tweet
 - Assign a value of 1 for presence, 0 for absence.
- Problem:
 - o Total features in the vector equal the size of the vocabulary
 - o $n+1$ total parameters to learn, $n = |V|$
 - o Majority of features will likely be 0 for each tweet

Negative and Positive Frequencies

Word frequency in classes		
Vocabulary	PosFreq (1)	NegFreq (0)
I	3	3
am	3	3
happy	2	0
because	1	0
learning	1	1
NLP	1	1
sad	0	2
not	0	1

- Simply number of times the words show up in the corresponding class

Feature Extraction with Frequencies

- freqs: dictionary mapping from (word, class) to frequency.

$$X_m = [1, \sum_w \text{freqs}(w, 1), \sum_w \text{freqs}(w, 0)]$$

↓ ↓ ↓ ↓
 Features of Bias Sum Pos. Sum Neg.
 ○ tweet m Frequencies Frequencies

- Word refers to unique, no duplicates

- Example:

$$\sum_w \text{freqs}(w, 0)]$$

I am sad, I am not learning NLP

-

- The value of the sum of positive frequencies = 8
- The value of the sum of negative frequencies = 11
- $X_M = [1, 8, 11]$

Vocabulary	PosFreq (1)	NegFreq (0)
I	3	3
am	3	3
happy	2	0
because	1	0
learning	1	1
NLP	1	1
sad	0	2
not	0	1

Preprocessing

- Stop words and punctuation
 - In a sentences, remove all stop words like "and", "has" ... and punctuations
 - Example:

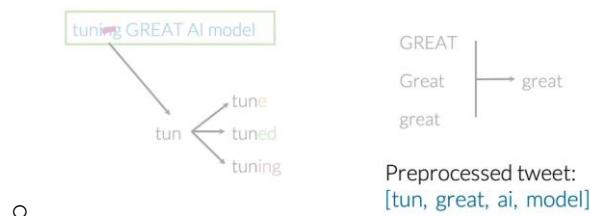
	Stop words	Punctuation
@YMourri @AndrewYNg tuning GREAT AI model https://deeplearning.ai!!!	and is a at has for of	,
@YMourri @AndrewYNg tuning GREAT AI model https://deeplearning.ai		.
		:
		!
		"
		'

- Remove all handles and URLs:

~~@YMourri @AndrewYNg tuning GREAT AI model~~
~~<https://deeplearning.ai>~~

- Stemming and lowercasing

- Stemming:
 - Convert every word to its stem.
 - Tuning → Tun
- Lowercaseing:
 - GREAT → great
- Tuning GREAT AI model → tun great ai model



Preprocessing lab:

```
1. import nltk                                     # Python library for NLP
2. from nltk.corpus import twitter_samples        # sample Twitter dataset from NLTK
3. from nltk.corpus import stopwords             # module for stop words that come with NLTK
4. from nltk.stem import PorterStemmer          # module for stemming
5. from nltk.tokenize import TweetTokenizer       # module for tokenizing strings
6. import re                                       # library for regular expression operations
7. import string                                    # for string operations
8.
9. nltk.download('twitter_samples')
10.
11. all_positive_tweets = twitter_samples.strings('positive_tweets.json') # a list of string(tweet)
12. all_negative_tweets = twitter_samples.strings('negative_tweets.json')
13.
14. all_positive_tweets[random.randint(0,5000)]
15. --> "#KFCkitchentours Kitchen is so clean. I'm totally amazed :) @KFC_India"
16. all_negative_tweets[random.randint(0,5000)]
17. --> "Athabasca glacier was there in #1948 :-( #athabasca #glacier #jasper #jaspernationalpark #alberta #explorealberta
#... http://t.co/dZZdqmf7Cz"
18. tweet = all_positive_tweets[2277]
19.
20. #Preprocess raw text for sentiment analysis
21.     nltk.download('stopwords')
22.     tweet2 = re.sub(r'^RT[\s]+', '', tweet)
23.     tweet2 = re.sub(r'https?://[^\s\n\r]+', '', tweet2)
24.     tweet2 = re.sub(r'#', '', tweet2)
25.     print(tweet2)
26. #- Tokenizing the string && lowercasing
27.     tokenizer = TweetTokenizer(preserve_case=False, strip_handles=True, reduce_len=True)
28.     tweet_tokens = tokenizer.tokenize(tweet2)
29.     #tweet2:
30.         #My beautiful sunflowers on a sunny Friday morning off :) sunflowers favourites happy Friday off...
31.     #tweet_tokens:
32.         #['my', 'beautiful', 'sunflowers', 'on', 'a', 'sunny', 'friday', 'morning', 'off', ':)', 'sunflowers',
33.         'favourites', 'happy', 'friday', 'off', ...]
34. #- Removing stop words and punctuation
35.     stopwords_english = stopwords.words('english')
36.     tweets_clean = []
37.     for word in tweet_tokens: # Go through every word in your tokens list
38.         if (word not in stopwords_english and # remove stopwords
39.             word not in string.punctuation): # remove punctuation
40.             tweets_clean.append(word)
41.     #tweets_clean (removed stop words and punctuation)
42.     #['beautiful', 'sunflowers', 'sunny', 'friday', 'morning', ':)', 'sunflowers', 'favourites', 'happy', 'friday',
43.     ...]
43. #- Stemming
44.     stemmer = PorterStemmer()
45.     tweets_stem = []
46.     for word in tweets_clean:
47.         stem_word = stemmer.stem(word) # stemming word
48.         tweets_stem.append(stem_word) # append to the list
49.     #tweets_stem:
50.         #['beauti', 'sunflow', 'sunni', 'friday', 'morn', ':)', 'sunflow', 'favourit', 'happi', 'friday', ...]
51.
52. --> make a single function process_tweet that does all above
```

Putting it all together

- Example: I am Happy Because i am learning NLP @deeplearning
 - o → Preprocessing: [happy, learn, nlp]
 - o → feature extraction:
 - [bias, sum positive frequencies, sum negative frequencies]
 - → [1,4,2]
- In practice,
 - o do preprocessing for a set of m raw tweets one by one
 - o extract feature for them → put them together in a X_matrix

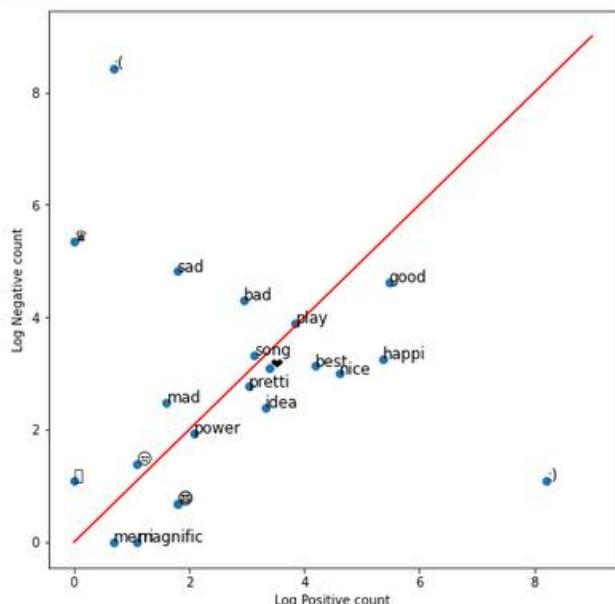
$$\begin{bmatrix} 1 & X_1^{(1)} & X_2^{(1)} \\ 1 & X_1^{(2)} & X_2^{(2)} \\ \vdots & \vdots & \vdots \\ 1 & X_1^{(m)} & X_2^{(m)} \end{bmatrix} \longrightarrow \begin{array}{l} [1, 40, 20] \\ [1, 20, 50], \\ \dots \\ [1, 5, 35] \end{array}$$

Word frequencies Lab:

```

1. def build_freqs(tweets, ys):
2.     freqs = {}
3.     for y, tweet in zip(yslist, tweets):
4.         for word in process_tweet(tweet):
5.             pair = (word, y)
6.             freqs[pair] = freqs.get(pair, 0) + 1 #example: ('followfriday', 1.0): 25
7.     return freqs
8. keys = ['happi', 'merri', 'nice', 'good', 'bad', 'sad', 'mad', '...']
9. data = []
10. for word in keys:
11.     pos = 0
12.     neg = 0
13.     if (word, 1) in freqs:
14.         pos = freqs[(word, 1)]
15.     if (word, 0) in freqs:
16.         neg = freqs[(word, 0)]
17.     data.append([word, pos, neg])
18. #data: [['happi', 212, 25],[...],['nice', 99, 19]]

```



Logistic Regression

- Training:
 - o Decision boundary:
 - $z = \theta^T x = 0$
 - o $x = [1, \text{pos}, \text{neg}]$
 - o $z(\theta, x) = \theta_0 + \theta_1 * \text{pos} + \theta_2 * \text{neg} = 0$
 - o $\text{neg} = (-\theta_0 - \theta_1 * \text{pos})/\theta_2$
- Testing
 - o $\text{pred} = h(X_{\text{val}}, \theta) \geq 0.5$
 - h : sigmoid
- Accuracy:
 - o $\frac{\sum_{i=1}^m (pred^{(i)} == y_{\text{val}}^{(i)})}{m}$

Review:

$$\begin{aligned}
 h(x)' &= \frac{1}{1+e^{-x}}' = \frac{-(1+e^{-x})'}{(1+e^{-x})^2} = \frac{-1' - (e^{-x})'}{(1+e^{-x})^2} = \frac{0 - (-x)'(e^{-x})}{(1+e^{-x})^2} = \frac{-(-1)(e^{-x})}{(1+e^{-x})^2} = \frac{e^{-x}}{(1+e^{-x})^2} \\
 &= \frac{1}{1+e^{-x}} \frac{e^{-x}}{1+e^{-x}} = h(x) \frac{1+1+e^{-x}}{1+e^{-x}} = h(x) \frac{1+e^{-x}}{1+e^{-x}} - \frac{1}{1+e^{-x}} = h(x)(1-h(x)) \\
 \frac{\partial}{\partial \theta_j} J(\theta) &= \frac{\partial}{\partial \theta_j} \frac{-1}{m} \sum_{i=1}^m y^{(i)} \log(h(x^{(i)}, \theta)) + (1-y^{(i)}) \log(1-h(x^{(i)}, \theta)) \\
 &= -\frac{1}{m} \sum_{i=1}^m y^{(i)} \frac{\partial}{\partial \theta_j} \log(h(x^{(i)}, \theta)) + (1-y^{(i)}) \frac{\partial}{\partial \theta_j} \log(1-h(x^{(i)}, \theta)) \\
 &= -\frac{1}{m} \sum_{i=1}^m \frac{y^{(i)} \frac{\partial}{\partial \theta_j} h(x^{(i)}, \theta)}{h(x^{(i)}, \theta)} + \frac{(1-y^{(i)}) \frac{\partial}{\partial \theta_j} (1-h(x^{(i)}, \theta))}{1-h(x^{(i)}, \theta)} \\
 &= -\frac{1}{m} \sum_{i=1}^m \frac{y^{(i)} \frac{\partial}{\partial \theta_j} h(x^{(i)}, \theta)}{h(x^{(i)}, \theta)} + \frac{(1-y^{(i)}) \frac{\partial}{\partial \theta_j} (1-h(x^{(i)}, \theta))}{1-h(x^{(i)}, \theta)} \\
 &= -\frac{1}{m} \sum_{i=1}^m \frac{y^{(i)} h(x^{(i)}, \theta) (1-h(x^{(i)}, \theta)) \frac{\partial}{\partial \theta_j} \theta^T x^{(i)}}{h(x^{(i)}, \theta)} + \frac{-(1-y^{(i)}) h(x^{(i)}, \theta) (1-h(x^{(i)}, \theta)) \frac{\partial}{\partial \theta_j} \theta^T x^{(i)}}{1-h(x^{(i)}, \theta)} \\
 &= -\frac{1}{m} \sum_{i=1}^m \frac{y^{(i)} h(x^{(i)}, \theta) (1-h(x^{(i)}, \theta)) \frac{\partial}{\partial \theta_j} \theta^T x^{(i)}}{h(x^{(i)}, \theta)} - \frac{(1-y^{(i)}) h(x^{(i)}, \theta) (1-h(x^{(i)}, \theta)) \frac{\partial}{\partial \theta_j} \theta^T x^{(i)}}{1-h(x^{(i)}, \theta)} \\
 &= -\frac{1}{m} \sum_{i=1}^m y^{(i)} (1-h(x^{(i)}, \theta)) x_j^{(i)} - (1-y^{(i)}) h(x^{(i)}, \theta) x_j^{(i)} \\
 &= -\frac{1}{m} \sum_{i=1}^m y^{(i)} (1-h(x^{(i)}, \theta)) - (1-y^{(i)}) h(x^{(i)}, \theta) x_j^{(i)} \\
 &= -\frac{1}{m} \sum_{i=1}^m y^{(i)} - y^{(i)} h(x^{(i)}, \theta) - h(x^{(i)}, \theta) + y^{(i)} h(x^{(i)}, \theta) x_j^{(i)} \\
 &= -\frac{1}{m} \sum_{i=1}^m y^{(i)} - h(x^{(i)}, \theta) x_j^{(i)} \\
 &= \frac{1}{m} \sum_{i=1}^m h(x^{(i)}, \theta) - y^{(i)} x_j^{(i)}
 \end{aligned}$$

The vectorized version:

$$\nabla J(\theta) = \frac{1}{m} \cdot X^T \cdot (H(X, \theta) - Y)$$

Frequencies to probability:

$P(w_i | \text{class})$

word	Pos	Neg
I	3	3
am	3	3
happy	2	1
because	1	0
learning	1	1
NLP	1	1
sad	1	2
not	1	2
Nclass	13	12

word	Pos	Neg
I	0.24	0.25
am	0.24	0.25
happy	0.15	0.08
because	0.08	0.01
learning	0.08	0.08
NLP	0.08	0.08
sad	0.08	0.17
not	0.08	0.17

- The columns summed to one

Naïve Bayes:

Tweet: I am happy today; I am learning.

$$\prod_{i=1}^m \frac{P(w_i | pos)}{P(w_i | neg)}$$

$$\frac{0.20}{0.20} * \frac{0.20}{0.20} * \frac{0.14}{0.10} * \frac{0.20}{0.20} * \frac{0.20}{0.20} * \frac{0.10}{0.10}$$

$$\prod_{i=1}^m \frac{P(w_i | pos)}{P(w_i | neg)} = \frac{0.14}{0.10} = 1.4 > 1 \rightarrow \text{Positive}$$

- Problem:

- o Number of negative and positive class may be different
 \rightarrow add a prior ratio:

$$\text{ratio}(w_i) = \frac{P(w_i | \text{Pos})}{P(w_i | \text{Neg})}$$

$$\approx \frac{\text{freq}(w_i, 1) + 1}{\text{freq}(w_i, 0) + 1}$$

$$\frac{P(pos)}{P(neg)} \prod_{i=1}^m \frac{P(w_i | pos)}{P(w_i | neg)}$$

- o Products bring risk of underflow
 \rightarrow use log-likelihood

$$\log \frac{P(tweet | pos)}{P(tweet | neg)} = \log(P(tweet | pos)) - \log(P(tweet | neg))$$

$$\text{positive} = \log(P(tweet | pos)) = \sum_{i=0}^n \log P(W_i | pos)$$

$$\text{negative} = \log(P(tweet | neg)) = \sum_{i=0}^n \log P(W_i | neg)$$

Summing the Lambdas

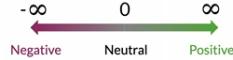
doc: I am happy because I am learning.

$$\lambda(w) = \log \frac{P(w | pos)}{P(w | neg)}$$

$$\lambda(\text{happy}) = \log \frac{0.09}{0.01} \approx 2.2$$

word	Pos	Neg	λ
I	0.05	0.05	0
am	0.04	0.04	0
happy	0.09	0.01	2.2
because	0.01	0.01	0
learning	0.03	0.01	1.1
NLP	0.02	0.02	0
sad	0.01	0.09	-2.2
not	0.02	0.03	-0.4

- o Inference time:



Laplacian Smoothing:

		Pos	Neg			Pos	Neg
word	1	3	3			1	0.19
am	3	3					
happy	2	1					
because	1	0					
learning	1	1					
NLP	1	1					
sad	1	2					
not	1	2					
Nclass	13	12					

$P(I|Pos) = \frac{3+1}{13+8}$

$V = 8$

Training Naïve Bayes:

- Get or annotate a dataset with positive and negative tweets
- Preprocess the tweets: $\text{process_tweet(tweet)} \rightarrow [w_1, w_2, \dots]$
- Compute $\text{freq}(w, \text{class})$
- Get $P(w | pos)$, $P(w | neg)$
- Get $\lambda(w)$

$$\lambda(w) = \log \frac{P(w|pos)}{P(w|neg)}$$

- Compute $\text{logprior} = \log(P(\text{pos})/P(\text{neg}))$

Testing Naïve Bayes on unseen data:

- X_{val} Y_{val} λ log_prior
- Score = $\text{predict}(X_{\text{val}}, \lambda, \text{logprior})$
- pred = score > 0
- Accuracy:

$$\circ \quad \frac{1}{m} \sum_{i=1}^m (\text{pred}_i == Y_{\text{val},i})$$

Applications of Naïve Bayes

- Author identification

$$\frac{P(\text{book}|\text{book})}{P(\text{not book}|\text{book})}$$

- Spam filtering

$$\frac{P(\text{spam}|\text{email})}{P(\text{nonspam}|\text{email})}$$

- Information retrieval

$$P(\text{document}_k|\text{query}) \propto \prod_{i=0}^{|\text{query}|} P(\text{query}_i|\text{document}_k)$$

○ Retrieve document if $P(\text{document}_k|\text{query}) > \text{threshold}$

- Word disambiguation

$$\frac{P(\text{river}|\text{text})}{P(\text{money}|\text{text})}$$

Bank:



○

- This method is usually used as a simple baseline. It is also really fast

Naïve Bayes Assumptions

- Words in a sentence (predictors or features) are independent of each other.
- Example Sentence: "It is sunny and hot in the Sahara desert."
 - Naïve Bayes would treat "sunny" and "hot" as independent when they are often correlated.
 - Could lead to incorrect estimations of conditional probabilities.
 - Naïve Bayes make the independence assumption and is affected by the word frequencies in the corpus. For example:



"It is sunny and hot in the Sahara desert."



"It's always cold and snowy in ___."

○

- The naïve model will assign equal weight to the words "spring, summer, fall, winter".

Reliance on Training Data Distribution

- Naïve Bayes is heavily dependent on the distribution of the training dataset.
- Ideal training datasets should reflect the actual proportions of classes (e.g., positive and negative tweets) found in a random sample.
- Many annotated datasets are artificially balanced, which does not match the real distribution (e.g., the prevalence of positive over negative tweets).

Issues due to Dataset Imbalance

- Artificially balanced datasets → overly optimistic or pessimistic models.
- Reasons for imbalance include platform content bans or user content muting, affecting the negative tweet count.

Analyzing NLP Errors

- Key to error analysis is examining the processed text version.

Common Sources of Errors

- Semantic Loss During Preprocessing:
 - o Example: "my beloved grandmother :(" becomes "beloved grandmother" after removing punctuation, altering the sentiment.
- Word Order Significance:
 - o The order of words can affect the meaning of a sentence, which may not be captured by some NLP models.
 - o Example: "I am not happy because I did not go" can be misinterpreted if the word "not" is overlooked.
- Language Quirks:
 - o Nuances like sarcasm, irony, and euphemism, which are naturally understood by humans, can confuse models like Naïve Bayes.
 - o Example: "This is a ridiculously powerful movie." might be misinterpreted due to the word "ridiculously."

Preprocessing Considerations

- Importance of preserving sentiment-indicating punctuation and negations.
- Check the results of preprocessing to ensure the model receives accurate data.

Model Limitations & Solutions

- Naïve Bayes classifiers can miss nuances due to their reliance on word frequency and the independence assumption.
- The importance of word order and handling of negations ("not") will be addressed later.

Adversarial Attacks

- Adversarial attacks describe situations where the model is tricked by language phenomena.
- Example: Positive statements with negative words due to sarcasm or emphasis can lead to misclassification.

Naïve Bayes Lab

- Identify the number of classes:

$$P(D_{pos}) = \frac{D_{pos}}{D}$$

$$P(D_{neg}) = \frac{D_{neg}}{D}$$

- Prior and Logprior

- o Prior: the ratio of the probability:

$$\frac{P(D_{pos})}{P(D_{neg})}$$

- o Logprior:

$$\text{logprior} = \log\left(\frac{P(D_{pos})}{P(D_{neg})}\right) = \log\left(\frac{D_{pos}}{D_{neg}}\right)$$

$$\text{logprior} = \log(P(D_{pos})) - \log(P(D_{neg})) = \log(D_{pos}) - \log(D_{neg})$$

- Positive and Negative Probability of a word:

$$P(W_{pos}) = \frac{freq_{pos} + 1}{N_{pos} + V}$$

$$P(W_{neg}) = \frac{freq_{neg} + 1}{N_{neg} + V}$$

- o

```

1. for pair in freqs.keys():
2.     if pair[1] > 0:
3.         N_pos += freqs[pair]
4.     else:
5.         N_neg += freqs[pair]

1. freq_pos = freqs.get((word,1),0)
2. freq_neg = freqs.get((word,0),0)
3. p_w_pos = (freq_pos+1)/(N_pos+V)
4. p_w_neg = (freq_neg+1)/(N_neg+V)

```

- Log-likelihood:

$$\log\left(\frac{P(W_{pos})}{P(W_{neg})}\right)$$

$$1. \text{loglikelihood[word]} = \text{np.log}(p_w_pos/p_w_neg)$$

- Make prediction:

$$p = \text{logprior} + \sum_i^N (\text{loglikelihood}_i)$$

```

1. for word in word_l:
2.     if word in loglikelihood:
3.         p += loglikelihood[word]

```

Vector space models

- Usage:
 - o identify similar meanings between sentences with different words
 - Use cases: question answering, paraphrasing, and summarization
 - o capture dependencies between words (e.g., "buyand "sell")
- Application:
 - o Information extraction ("answering who, what, where, how, etc.")
 - o Machine Translation
 - o Chatbot

Word-By-Word design

- Number of times two words occur together within a certain distance k
- Example:
 - o $k = 2$
 - 2 words preceding/following it in a sentence will be used to determine co-occurrences.
 - o "I like simple data", "I prefer simple raw data".

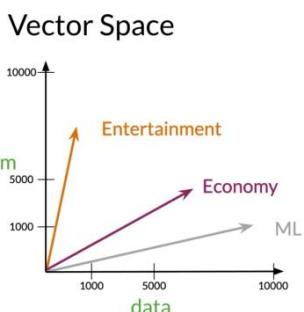
	simple	Raw	like	I
data	2	1	1	0

Word-By-Document Design

- Number of times a word occurs within a certain category

Word/Type	Entertainment	Economy	Machine Learning
data	500	6620	9320
film	7000	4000	1000

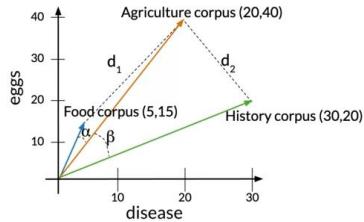
Vector Space



- Measures of "similarity" by Angle, Distance
 - o Calculate the euclidean distance between the corpus
- Economy and ML are more alike than with entertainment

Problem with using Euclidean distance

- Euclidean distance is affected by the size of the document corpora.



Cosine Similarity:

- Intuition:
 - o Use the angle between the documents.
- Equation:

$$\hat{v} \cdot \hat{w} = \|\hat{v}\| \|\hat{w}\| \cos(\beta)$$

$$\cos(\beta) = \frac{\hat{v} \cdot \hat{w}}{\|\hat{v}\| \|\hat{w}\|}$$

$$= \frac{(20 \times 30) + (40 \times 20)}{\sqrt{20^2 + 40^2} \times \sqrt{30^2 + 20^2}}$$

- o $\beta = 90^\circ \rightarrow \text{cosine similarity} = 0$
- o $\beta = 0^\circ \rightarrow \text{cosine similarity} = 1$

Manipulating words in vector space

- to infer unknown relations among words



- Predict Russia's capital with the closest one

```

1. capital = vec('France') - vec('Paris')
2. country = vec('Madrid') + capital
3. def find_closest_word(v, k = 1):
4.     diff = embedding.values - v
5.     delta = np.sum(diff * diff, axis=1)
6.     i = np.argmin(delta)
7.     return embedding.iloc[i].name
8. find_closest_word(country)

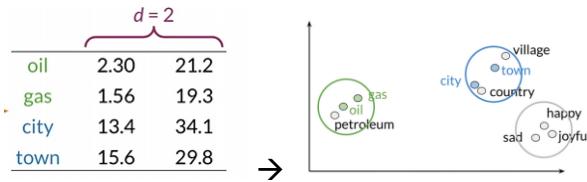
```

- Problem:

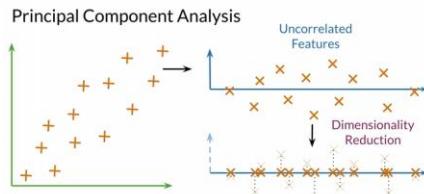
- o Often at very high dimension, unable to visualize clearly

	$d > 2$		
oil	0.20	...	0.10
gas	2.10	...	3.40
city	9.30	...	52.1
town	6.20	...	34.3

- o \rightarrow reduce the dimension by PCA



PCA



- How to get uncorrelated features
 - o Eigenvector: Uncorrelated features for your data
 - o Eigenvalue: the amount of information retained by each feature
- Algorithm:
 - o Mean normalize data
 - o Get Covariance Matrix
 - o Perform SVD on the Covariance matrix → eigenvectors(U) and eigenvalues(S)
 - Eigenvectors give the direction of uncorrelated features
 - Eigenvalues are the variance of the new features
 - o Dot product to project data: $X' = XU[:,0:2]$
 - Gives the projection on uncorrelated features
 - o Percentage of retained variance

$$\frac{\sum_{i=0}^1 S_{ii}}{\sum_{j=0}^d S_{jj}}$$
 -

PCA algorithm:

```

1. def compute_pca(X, n_components=2):
2.     X_demeaned = X - np.mean(X, axis=0)
3.     covariance_matrix = np.cov(X_demeaned, rowvar = False)
4.     eigen_vals, eigen_vecs = np.linalg.eigh(covariance_matrix)
5.     idx_sorted = np.argsort(eigen_vals)
6.     idx_sorted_decreasing = idx_sorted[::-1]
7.     eigen_vals_sorted = eigen_vals[idx_sorted_decreasing]
8.     eigen_vecs_sorted = eigen_vecs[:,idx_sorted_decreasing]
9.     eigen_vecs_subset = eigen_vecs_sorted[:, :n_components]
10.    X_reduced = np.dot(eigen_vecs_subset.T, X_demeaned.T).T
11.    return X_reduced

```

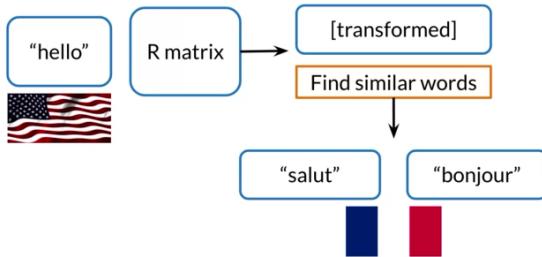
Transforming word vectors

- Machine computes word embeddings for each language and learns to align them.
- Steps:
 - o Retrieve the English word embedding (e.g., "cat") [1,0,1].
 - o Transform this embedding into the French word vector space.
 - o Search for the closest French word vector to find translation candidates.
 - o Ideally, the machine will identify "chat" as the French equivalent for "cat".

Align word vectors

- $$\begin{matrix} \left[\begin{matrix} ["\text{cat}"] \\ [...] \\ ["\text{zebra}"] \end{matrix} \right] & \mathbf{X} \end{matrix} \quad \mathbf{X}\mathbf{R} \approx \mathbf{Y} \quad \begin{matrix} \left[\begin{matrix} ["\text{chat}"] \\ [...] \\ ["\text{z\u00e9bresse}"] \end{matrix} \right] & \mathbf{Y} \end{matrix}$$
- Finding the transformation matrix R:
 - o Initialize R randomly
 - o Iteratively:
 - Use Frobenius norm
 - $\|\mathbf{A}\|_F \equiv \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2}$
 - ```
1. A = np.array([[2,2],[2,2]])
2. A_squared = np.square(A)
3. A_Frobenious = np.sqrt(np.sum(A_squared))
4. A_Frobenious #4.0
```
      - Loss =  $\|\mathbf{XR} - \mathbf{Y}\|_F$
      - Easier to work with frobenius norm squared:
        - o  $\|\mathbf{XR} - \mathbf{Y}\|_F^2$
        - o  $\mathbf{A} = \begin{pmatrix} 2 & 2 \\ 2 & 2 \end{pmatrix}$
        - o  $\|\mathbf{A}\|_F^2 = (\sqrt{2^2 + 2^2 + 2^2 + 2^2})^2$
        - o  $\|\mathbf{A}\|_F^2 = 16$
      - Update R by subtracting the gradient
 
$$g = \frac{d}{dR} \text{Loss}$$
        - $R = R - \alpha g$
        - With frobenious norm squared
 
$$\text{Loss} = \|\mathbf{XR} - \mathbf{Y}\|_F^2$$
          - $g = \frac{d}{dR} \text{Loss} = \frac{2}{m} (\mathbf{X}^T(\mathbf{XR} - \mathbf{Y}))$

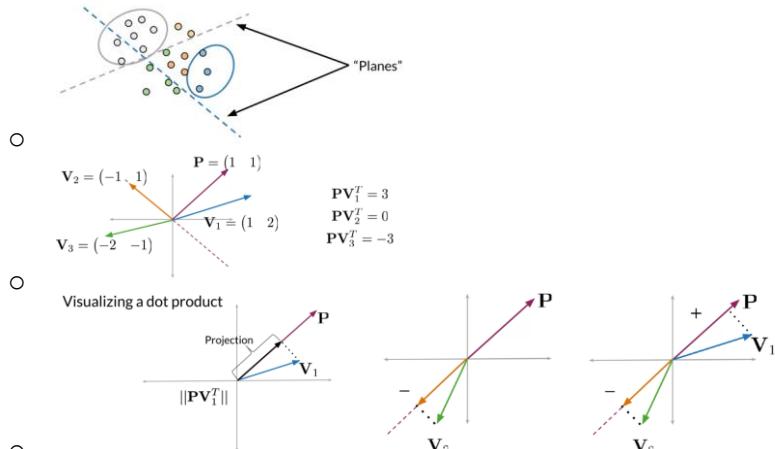
How to find similar word vectors?



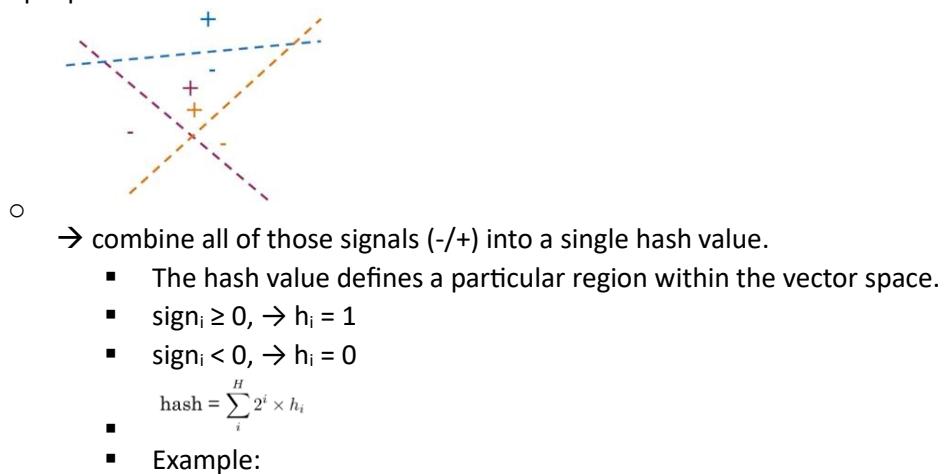
- K-nearest neighbors
  - o Key operation to find a matching word.
  - o Use locality sensitive hashing

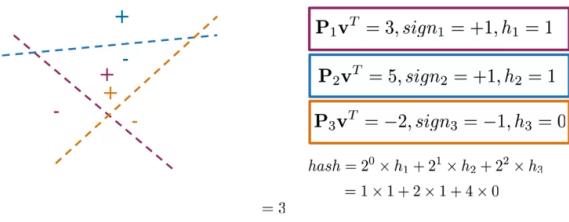
Locality sensitive hashing

- a hashing method that cares very deeply about assigning items based on where they're located in vector space
- Planes:



- Multiple planes:



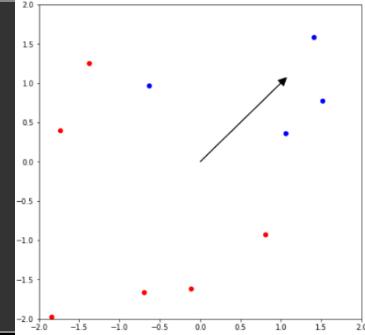


## Hash function and multiplanes lab

```

1. P = np.array([[1, 1]]) # Define a single plane.
2. fig, ax1 = plt.subplots(figsize=(8, 8))
3.
4. plot_vectors([P], axes=[2, 2], ax=ax1)
5. for i in range(0, 10):
6. v1 = np.array(np.random.uniform(-2, 2, 2))
7. side_of_plane = np.sign(np.dot(P, v1.T))
8. if side_of_plane == 1:
9. ax1.plot([v1[0]], [v1[1]], 'bo')
10. else:
11. ax1.plot([v1[0]], [v1[1]], 'ro')
12. plt.show()

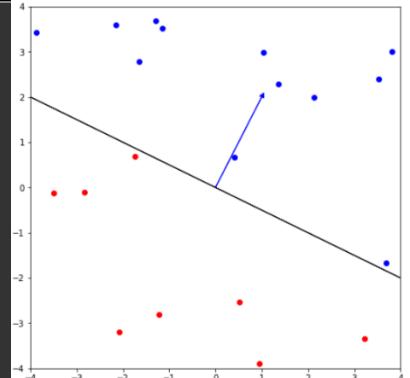
```



```

1. P = np.array([[1, 2]])
2. PT = np.dot([[0, 1], [-1, 0]], P.T).T
3.
4. fig, ax1 = plt.subplots(figsize=(8, 8))
5. plot_vectors([P], colors=['b'], axes=[2, 2], ax=ax1)
6. plot_vectors([PT*4,PT*-4],colors=['k','k'],axes=[4,4],ax=ax1)
7. for i in range(0, 20):
8. v1 = np.array(np.random.uniform(-4, 4, 2))
9. side_of_plane = np.sign(np.dot(P, v1.T))
10. if side_of_plane == 1:
11. ax1.plot([v1[0]], [v1[1]], 'bo')
12. else:
13. ax1.plot([v1[0]], [v1[1]], 'ro')
14. plt.show()

```

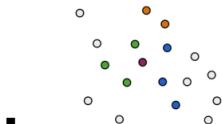


```

1. def side_of_plane(P, v):
2. dotproduct = np.dot(P, v.T)
3. sign_of_dot_product = np.sign(dotproduct)
4. sign_of_dot_product_scalar = sign_of_dot_product.item()
5. return sign_of_dot_product_scalar
1. def hash_multi_plane(P_l, v):
2. hash_value = 0
3. for i, P in enumerate(P_l):
4. sign = side_of_plane(P,v)
5. hash_i = 1 if sign >=0 else 0
6. hash_value += 2**i * hash_i
7. return hash_value
1. def side_of_plane_matrix(P, v):
2. dotproduct = np.dot(P, v.T)
3. sign_of_dot_product = np.sign(dotproduct)
4. return sign_of_dot_product
1. def hash_multi_plane_matrix(P, v, num_planes):
2. sides_matrix = side_of_plane_matrix(P, v) # Get the side of planes for P and v
3. hash_value = 0
4. for i in range(num_planes):
5. sign = sides_matrix[i].item() # Get the value inside the matrix cell
6. hash_i = 1 if sign >=0 else 0
7. hash_value += 2**i * hash_i # sum 2^i * hash_i
8. return hash_value

```

## Approximate nearest neighbors

- Vector Space Partitioning
  - o Planes as Dividers: A few planes can divide the vector space into distinct regions.
  - o Optimal Division Uncertainty: It is unknown which set of planes best divides the vector space.
- Multiple Sets of Random Planes
  - o Creating Multiple Sets: Instead of one set of planes, multiple sets of random planes are generated.
- Implementation for ANN
  - o Vector Space Example: A vector space is considered with a magenta dot representing a transformed English word into a French word vector.
  - o Hash Pockets Assignment: Different sets of random planes are used to assign similar vectors (e.g., French word vectors) to the same hash pockets.
    - One set pairs magenta and green vectors.
    - Another set pairs red and blue vectors.
    - A third set pairs magenta and orange vectors.
  - o Robust Searching: Multiple sets of random planes make the search for nearest neighbors more robust by considering various hash pockets.
- Approximation vs. Precision
  - o ANN Definition: Searching only a subset of the entire vector space for potential nearest neighbors, leading to an approximate, rather than exact, k-NN.
  - o Trade-off: There is a trade-off of some precision for significantly increased efficiency in the search process.

## Document representation

| Document representation |              |
|-------------------------|--------------|
| I love learning!        | [?, ?, ?]    |
| I                       | [1, 0, 1]    |
| love                    | + [-1, 0, 1] |
| learning                | + [1, 0, 1]  |
| I love learning!        | = [1, 0, 3]  |

## Naïve Machine Translation and LSH

|                                                                                                                                                                                                                                                                           |                                                                                                                                                                                                                                                                          |
|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| English word embedding $\mathbf{X}$<br>$\begin{pmatrix} 2 & 9 & \dots & \dots & 1 & 3 \\ 1 & 2 & \dots & \dots & 0 & 3 \\ 9 & 8 & \dots & \dots & 7 & 5 \\ 1 & 0 & \dots & \dots & 9 & 3 \\ 4 & 9 & \dots & \dots & 4 & 0 \\ 2 & 1 & \dots & \dots & 0 & 5 \end{pmatrix}$ | French word embedding $\mathbf{Y}$<br>$\begin{pmatrix} 1 & 9 & \dots & \dots & 9 & 3 \\ 4 & 2 & \dots & \dots & 0 & 8 \\ 9 & 6 & \dots & \dots & 6 & 5 \\ 1 & 9 & \dots & \dots & 9 & 3 \\ 4 & 2 & \dots & \dots & 4 & 0 \\ 2 & 8 & \dots & \dots & 6 & 5 \end{pmatrix}$ |
|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|

```

1. def get_matrices(en_fr, french_vecs, english_vecs):
2. X_l = list()
3. Y_l = list()
4. english_set = set(english_vecs)
5. french_set = set(french_vecs)
6. french_words = set(en_fr.values())
7. for en_word, fr_word in en_fr.items():
8. if fr_word in french_set and en_word in english_set:
9. en_vec = english_vecs[en_word]
10. fr_vec = french_vecs[fr_word]
11. X_l.append(en_vec)
12. Y_l.append(fr_vec)
13. X = np.vstack(X_l)
14. Y = np.vstack(Y_l)
15. return X, Y

```

Translation

Problem:

$$\arg \min_{\mathbf{R}} \|\mathbf{X}\mathbf{R} - \mathbf{Y}\|_F$$

Frobenius norm:

$$\|\mathbf{A}\|_F \equiv \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2}$$

Loss function:

$$\frac{1}{m} \|\mathbf{X}\mathbf{R} - \mathbf{Y}\|_F^2$$

Compute loss:

$$L(\mathbf{X}, \mathbf{Y}, \mathbf{R}) = \frac{1}{m} \sum_{i=1}^m \sum_{j=1}^n (a_{ij})^2$$

```

1. def compute_loss(X, Y, R):
2. m = X.shape[0]
3. diff = np.dot(X, R) - Y
4. diff_squared = np.square(diff)
5. sum_diff_squared = np.sum(diff_squared)
6. loss = sum_diff_squared / m
7. return loss

```

Compute gradient:

```

1. def compute_gradient(X, Y, R):
2. m = X.shape[0]
3. gradient = np.dot(X.T, (np.dot(X,R) - Y)) * 2/m
4. return gradient

```

Update R:

$$R_{\text{new}} = R_{\text{old}} - \alpha g$$

```

1. def align_embeddings(X, Y, train_steps=100, learning_rate=0.0003, verbose=True,
2. compute_loss=compute_loss, compute_gradient=compute_gradient):
3. np.random.seed(129)
4. R = np.random.rand(X.shape[1], X.shape[1])
5. for i in range(train_steps):
6. if verbose and i % 25 == 0:
7. print(f"loss at iteration {i} is: {compute_loss(X, Y, R):.4f}")
8. gradient = compute_gradient(X, Y, R)
9. R -= learning_rate*gradient

```

Testing the translation by K-Nearest Neighbor Algorithm

- Use 1-NN with eR as input → search for an embedding f (as a row) in the matrix Y which is the closest to the transformed vector eR.

Cosine Similarity:

- Cosine of the angle between vector u and v.

$$\cos(u, v) = \frac{u \cdot v}{\|u\| \|v\|}$$

Nearest\_neighbor

```

1. def nearest_neighbor(v, candidates, k=1, cosine_similarity=cosine_similarity):
2. similarity_l = []
3. for row in candidates:
4. cos_similarity = cosine_similarity(v, row)
5. similarity_l.append(cos_similarity)
6. sorted_ids = np.argsort(similarity_l)
7. sorted_ids = sorted_ids[::-1]
8. k_idx = sorted_ids[:k]
9. return k_idx

```

Test vocabulary

```

1. def test_vocabulary(X, Y, R, nearest_neighbor=nearest_neighbor):
2. pred = np.dot(X, R)
3. num_correct = 0
4. for i in range(len(pred)):
5. pred_idx = nearest_neighbor(pred[i], Y)
6. if pred_idx == i:
7. num_correct += 1
8. accuracy = num_correct / pred.shape[0]
9. return accuracy

```

Get document embedding

```

1. def get_document_embedding(tweet, en_embeddings, process_tweet=process_tweet):
2. doc_embedding = np.zeros(300)
3. processed_doc = process_tweet(tweet)
4. for word in processed_doc:
5. doc_embedding = doc_embedding + en_embeddings.get(word, 0)
6. return doc_embedding

```

## Get document vecs

- Store all tweet embeddings into a dictionary

```

1. def get_document_vecs(all_docs, en_embeddings,
get_document_embedding=get_document_embedding):
2. ind2Doc_dict = {}
3. document_vec_l = []
4. for i, doc in enumerate(all_docs):
5. doc_embedding = get_document_embedding(doc, en_embeddings)
6. ind2Doc_dict[i] = doc_embedding
7. document_vec_l.append(doc_embedding)
8. document_vec_matrix = np.vstack(document_vec_l)
9. return document_vec_matrix, ind2Doc_dict

```

## Locality sensitive hashing (LSH)

Let's say your data points are plotted like this:

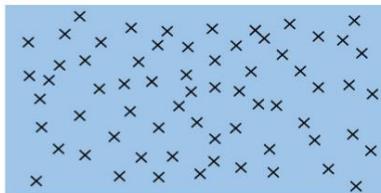


Figure 3

You can divide the vector space into regions and search within one region for nearest neighbors of a given vector.

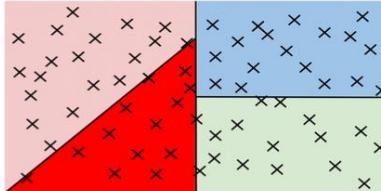


Figure 4

## Hash value of vector

$$\text{hash} = \sum_{i=0}^{N-1} (2^i \times h_i)$$

```

1. def hash_value_of_vector(v, planes):
2. dot_product = np.dot(v, planes)
3. sign_of_dot_product = np.sign(dot_product)
4. h = sign_of_dot_product >= 0
5. h = np.squeeze(h)
6. hash_value = 0
7. n_planes = planes.shape[1]
8. for i in range(n_planes):
9. hash_value += np.power(2, i) * h[i]
10. hash_value = int(hash_value)
11. return hash_value

```

## Make hash table

```
1. def make_hash_table(vecs, planes, hash_value_of_vector=hash_value_of_vector):
2. num_of_planes = planes.shape[1]
3. num_buckets = 2**num_of_planes
4. hash_table = {i: [] for i in range(num_buckets)}
5. id_table = {i: [] for i in range(num_buckets)}
6. for i, v in enumerate(vecs):
7. h = hash_value_of_vector(v, planes)
8. hash_table[h].append(v)
9. id_table[h].append(i)
10. return hash_table, id_table
```

## Creating all hash tables

```
1. def create_hash_id_tables(n_universes):
2. hash_tables = []
3. id_tables = []
4. for universe_id in range(n_universes): # there are 25 hashes
5. print('working on hash universe #:', universe_id)
6. planes = planes_l[universe_id]
7. hash_table, id_table = make_hash_table(document_vecs, planes)
8. hash_tables.append(hash_table)
9. id_tables.append(id_table)
10. return hash_tables, id_tables
11. hash_tables, id_tables = create_hash_id_tables(N_UNIVERSES)
```

## Approximate KNN

```
1. def approximate_knn(doc_id, v, planes_l, hash_tables, id_tables, k=1,
num_universes_to_use=25, hash_value_of_vector=hash_value_of_vector):
2. vecs_to_consider_l = list()
3. ids_to_consider_l = list()
4. ids_to_consider_set = set()
5. for universe_id in range(num_universes_to_use):
6. planes = planes_l[universe_id]
7. hash_value = hash_value_of_vector(v, planes)
8. hash_table = hash_tables[universe_id]
9. document_vectors_l = hash_table[hash_value]
10. id_table = id_tables[universe_id]
11. new_ids_to_consider = id_table[hash_value]
12. for i, new_id in enumerate(new_ids_to_consider):
13. if doc_id == new_id:
14. continue
15. if new_id not in ids_to_consider_set:
16. document_vector_at_i = document_vectors_l[i]
17. vecs_to_consider_l.append(document_vector_at_i)
18. ids_to_consider_l.append(new_id)
19. ids_to_consider_set.add(new_id)
20. print("Fast considering %d vecs" % len(vecs_to_consider_l))
21. vecs_to_consider_arr = np.array(vecs_to_consider_l)
22. nearest_neighbor_idx_l = nearest_neighbor(v, vecs_to_consider_arr, k=k)
23. nearest_neighbor_ids = [ids_to_consider_l[idx]
24. for idx in nearest_neighbor_idx_l]
25. return nearest_neighbor_ids
```

## Overview of Autocorrect

- How it works:
  - o Identify a misspelled word
 

1. If word not in vocab:
    2. Misspelled = True
  - o Find strings n edit distance away
    - Number of edit operation away one string from another
    - Insert, delete, replace(a letter), switch (swap 2 adjacent letters)
  - o Filter candidates
  - o Calculate word probabilities
    - Make a frequency table storing probability for each word
      - Number of times the word appears/total size of the corpus

## Minimum edit distance

- Edit cost:
  - o Insert(1), Delete(1), Replace(2)
  - o Examples:
    - play → stay:
      - p → s: replace(2)
      - l → t: replace(2)
    - total edit distance = 2+2 = 4
- dynamic programming

○

|   | 0 | 1 | 2 | 3 | 4 |   |
|---|---|---|---|---|---|---|
| 0 | # | 0 | 1 | 2 | 3 | 4 |
| 1 | p | 1 | 2 | 3 | 4 | 5 |
| 2 | I | 2 | 3 | 4 | 5 | 6 |
| 3 | a | 3 | 4 | 5 | 4 | 5 |
| 4 | y | 4 | 5 | 6 | 5 | 4 |

○  $D[i, j] =$

$$\min \begin{cases} D[i - 1, j] + \text{del\_cost} \\ D[i, j - 1] + \text{ins\_cost} \\ D[i - 1, j - 1] + \begin{cases} \text{rep\_cost}; & \text{if } \text{src}[i] \neq \text{tar}[j] \\ 0; & \text{if } \text{src}[i] = \text{tar}[j] \end{cases} \end{cases}$$

## Auto Correct Lab

- Edit Distance
  - $$P(c|w) = \frac{P(w|c) \times P(c)}{P(w)}$$

- Data Preprocessing

```
1. def process_data(file_name):
2. words = [] # return this variable correctly
3. with open(file_name) as f:
4. file_name_data = f.read()
5. file_name_data = file_name_data.lower()
6. words = re.findall('\w+', file_name_data)
7. return words
```

- Get count

```
1. def get_count(word_l):
2. word_count_dict = {} # fill this with word counts
3. word_count_dict = Counter(word_l)
4. return word_count_dict
```

- Get probability of word count

$$\circ \quad P(w_i) = \frac{C(w_i)}{M}$$

```
1. def get_probs(word_count_dict):
2. probs = {} # return this variable correctly
3. m = sum(word_count_dict.values())
4. for key in word_count_dict.keys():
5. probs[key] = word_count_dict[key] / m
6. return probs
```

- Delete letter

```
1. def delete_letter(word, verbose=False):
2. delete_l = []
3. split_l = []
4. for c in range(len(word)):
5. split_l.append((word[:c], word[c:]))
6. for a,b in split_l:
7. delete_l.append(a+b[1:])
8. if verbose: print(f"input word {word}, \nsplit_l = {split_l}, \ndelete_l = {delete_l}")
9. return delete_l
```

- Switch letter

```
1. def switch_letter(word, verbose=False):
2. switch_l = []
3. split_l = []
4. word_len = len(word)
5. for c in range(word_len):
6. split_l.append((word[:c], word[c:]))
7. switch_l = [a + b[1] + b[0] + b[2:] for a,b in split_l if len(b) >= 2]
8. if verbose: print(f"Input word = {word} \nsplit_l = {split_l} \nswitch_l = {switch_l}")
9. return switch_l
```

- Replace

```
1. def replace_letter(word, verbose=False):
2. letters = 'abcdefghijklmnopqrstuvwxyz'
3. replace_l = []
4. split_l = []
5. for c in range(len(word)):
6. split_l.append((word[0:c], word[c:]))
7. replace_l = [a + l + (b[1:] if len(b) > 1 else '') for a,b in split_l if b for l in letters]
8. replace_set = set(replace_l)
9. replace_set.remove(word)
10. replace_l = sorted(list(replace_set))
11. if verbose: print(f"Input word = {word} \nsplit_l = {split_l} \nreplace_l = {replace_l}")
12. return replace_l
```

- Insert

```
1. def insert_letter(word, verbose=False):
2. letters = 'abcdefghijklmnopqrstuvwxyz'
3. insert_l = []
4. split_l = []
5. for c in range(len(word)+1):
6. split_l.append((word[0:c], word[c:]))
7. insert_l = [a + l + b for a,b in split_l for l in letters]
8. if verbose: print(f"Input word {word} \nsplit_l = {split_l} \ninsert_l = {insert_l}")
9. return insert_l
```

- Edit one letter → get all possible edits

```
1. def edit_one_letter(word, allow_switches = True):
2. edit_one_set = set()
3. edit_one_set = set(insert_letter(word)+delete_letter(word)+replace_letter(word))
4. if(allow_switches):
5. edit_one_set = edit_one_set.union(set(switch_letter(word,verbose=False)))
6. return set(edit_one_set)
```

- Edit two letters

```
1. def edit_two_letters(word, allow_switches = True):
2. edit_two_set = set()
3. edit_one_set = edit_one_letter(word,allow_switches=allow_switches)
4. edit_two_set = edit_one_set
5. for w in edit_one_set:
6. edit_two_set = edit_two_set.union(set(edit_one_letter(w,allow_switches=allow_switches)))
7. return set(edit_two_set)
```

- Get correction

```

1. def get_corrections(word, probs, vocab, n=2, verbose = False):
2. suggestions = []
3. n_best = []
4. temp=[]
5. if word in probs.keys():
6. temp+=[(word,probs[word])]
7. temp+=sorted([(w,probs[w]) for w in edit_one_letter(word) if w in probs],key = lambda x : -x[1])
8. temp+=sorted([(w,probs[w]) for w in edit_two_letters(word) if w in probs],key = lambda x: -x[1])
9. suggestions = set([sugg[0] for sugg in temp[:n]])
10. n_best = [sugg for sugg in temp[:n]]
11. if verbose: print("entered word = ", word, "\nsuggestions = ", suggestions)
12. return n_best

```

### Minimum edit distance

#### Initialization

$$\begin{aligned} D[0,0] &= 0 \\ D[i,0] &= D[i-1,0] + \text{del\_cost}(\text{source}[i]) \\ D[0,j] &= D[0,j-1] + \text{ins\_cost}(\text{target}[j]) \end{aligned}$$

#### Per Cell Operations

$$D[i,j] = \min \begin{cases} D[i-1,j] + \text{del\_cost} \\ D[i,j-1] + \text{ins\_cost} \\ D[i-1,j-1] + \begin{cases} \text{rep\_cost}; & \text{if } \text{src}[i] \neq \text{tar}[j] \\ 0; & \text{if } \text{src}[i] = \text{tar}[j] \end{cases} \end{cases}$$

```

1. def min_edit_distance(source, target, ins_cost = 1, del_cost = 1, rep_cost = 2):
2. m = len(source)
3. n = len(target)
4. D = np.zeros((m+1, n+1), dtype=int)
5. for row in range(0,m+1): # Replace None with the proper range
6. D[row,0] = row*ins_cost
7. for col in range(0,n+1): # Replace None with the proper range
8. D[0,col] = col*del_cost
9. for row in range(1,m+1):
10. for col in range(1,n+1):
11. r_cost = rep_cost
12. if source[row-1]==target[col-1]:
13. r_cost = 0
14. D[row,col] = min(D[row-1,col-1]+r_cost , D[row-1,col]+del_cost , D[row,col-1]+ins_cost)
15. med = D[m,n]
16. return D, med

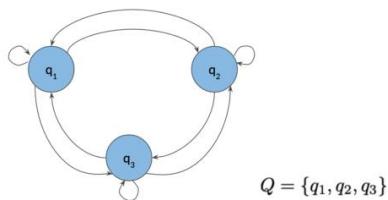
```

## Part of Speech Tagging

| lexical term | tag | example            |
|--------------|-----|--------------------|
| noun         | NN  | something, nothing |
| verb         | VB  | learn, study       |
| determiner   | DT  | the, a             |
| w-adverb     | WRB | why, where         |
| ...          | ... |                    |

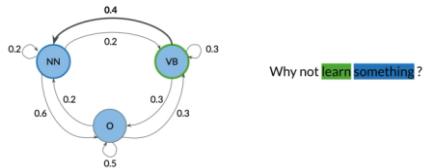
## Markov Chains

- Stochastic Model: Describes a sequence where each event's probability is dependent only on the state of the previous event.
  - o "Stochastic": Indicates randomness within the process.



- Transition Probabilities: Quantifies the likelihood of moving from one state to another.
- States: Represents the current condition or moment in the chain.

## Markov Chains and POS Tags



- Using transition matrix:

|                 | NN  | VB  | O   |
|-----------------|-----|-----|-----|
| $\pi$ (initial) | 0.4 | 0.1 | 0.5 |
| NN (noun)       | 0.2 | 0.2 | 0.6 |
| VB (verb)       | 0.4 | 0.3 | 0.3 |
| O (other)       | 0.2 | 0.3 | 0.5 |

- o Written as an actual matrix:

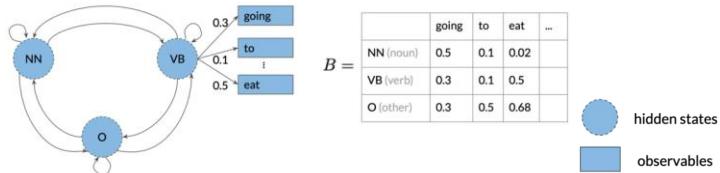
$$A = \begin{pmatrix} 0.4 & 0.1 & 0.5 \\ 0.2 & 0.2 & 0.6 \\ 0.4 & 0.3 & 0.3 \\ 0.2 & 0.3 & 0.5 \end{pmatrix}$$

- Transition matrix A given some states Q:

| States                    | Transition matrix                                                                                                          |
|---------------------------|----------------------------------------------------------------------------------------------------------------------------|
| $Q = \{q_1, \dots, q_N\}$ | $A = \begin{pmatrix} a_{1,1} & \dots & a_{1,N} \\ \vdots & \ddots & \vdots \\ a_{N+1,1} & \dots & a_{N+1,N} \end{pmatrix}$ |

## Hidden Markov Models

- Implies that states are hidden
- Emission probabilities:
  - o Describe the transition from the hidden states of the hidden markov



- o The last row (O) is incorrect, the sum of each row is equal to 1.

$$\sum_{j=1}^V b_{ij} = 1$$

|                                     |                                                                                                                                                 |                                                                                                                                       |
|-------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------|
| States<br>$Q = \{q_1, \dots, q_N\}$ | Transition matrix<br>$A = \begin{pmatrix} a_{1,1} & \dots & a_{1,N} \\ \vdots & \ddots & \vdots \\ a_{N+1,1} & \dots & a_{N+1,N} \end{pmatrix}$ | Emission matrix<br>$B = \begin{pmatrix} b_{11} & \dots & b_{1V} \\ \vdots & \ddots & \vdots \\ b_{N1} & \dots & b_{NV} \end{pmatrix}$ |
|-------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------|

## Calculate probability

- Transition probability



- o Let occurrences of tag pairs =  $C(t_{i-1}, t_i)$

Calculate probabilities using the counts

$$P(t_i | t_{i-1}) = \frac{C(t_{i-1}, t_i)}{\sum_{j=1}^N C(t_{i-1}, t_j)}$$

- o Training:

- Given a set of corpus, preprocess →

<s> in a station of the metro  
<s> the apparition of these faces in the crowd :  
● <s> petals on a wet, black bough .  
● Initial probability

## Populating the transition matrix

|           | NN | VB | O |
|-----------|----|----|---|
| x         | 1  | 0  | 2 |
| NN (noun) | 0  | 0  | 6 |
| VB (verb) | 0  | 0  | 0 |
| O (other) | 6  | 0  | 8 |

Era Pound - 1913

- Transition probability:

$$P(t_i | t_{i-1}) = \frac{C(t_{i-1}, t_i)}{\sum_{j=1}^N C(t_{i-1}, t_j)}$$

|    | NN  | VB  | O   |
|----|-----|-----|-----|
| x  | 1+ε | 0+ε | 2+ε |
| NN | 0+ε | 0+ε | 6+ε |
| VB | 0+ε | 0+ε | 0+ε |
| O  | 6+ε | 0+ε | 8+ε |

$P(t_i | t_{i-1}) = \frac{C(t_{i-1}, t_i) + \boxed{ε}}{\sum_{j=1}^N C(t_{i-1}, t_j) + \boxed{N} * \boxed{ε}}$

- o with smoothing:

|    | NN     | VB     | O      |
|----|--------|--------|--------|
| x  | 0.3333 | 0.0003 | 0.6663 |
| NN | 0.0001 | 0.0001 | 0.9996 |
| VB | 0.3333 | 0.3333 | 0.3333 |
| O  | 0.4285 | 0.0000 | 0.5713 |

- May not need smooth for initial probability in real world application.

## Populating the Emission matrix

$$P(w_i|t_i) = \frac{C(t_i, w_i) + \epsilon}{\sum_{j=1}^V C(t_i, w_j) + N * \epsilon}$$

$$= \frac{C(t_i, w_i) + \epsilon}{C(t_i) + N * \epsilon}$$

|           | in | a | ... |
|-----------|----|---|-----|
| NN (noun) | 0  |   |     |
| VB (verb) | 0  |   |     |
| O (other) | 2  |   |     |

<s> in a station of the metro  
<s> the apparition of these faces in the crowd:  
<s> petals on a wet, black bough.

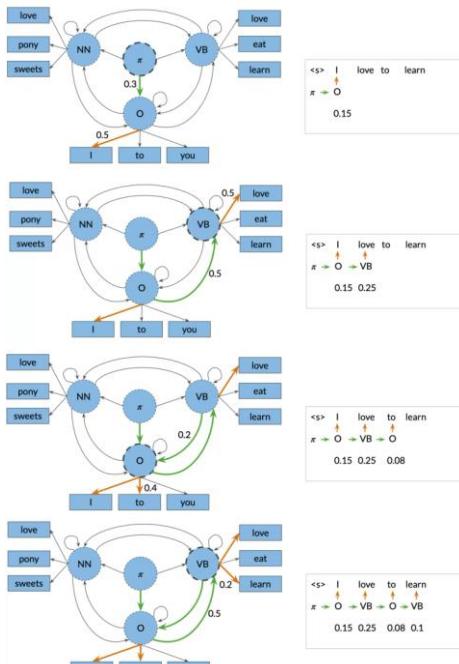
## Working with tags

```

1. tags = ['RB', 'NN', 'TO']
2. num_tags = len(tags)
3. transition_matrix = np.zeros((num_tags, num_tags))
4. transition_matrix
5. transition_counts = {
6. ('NN', 'NN'): 16241,
7. ('RB', 'RB'): 2263,
8. ('TO', 'TO'): 2,
9. ('NN', 'TO'): 5256,
10. ('RB', 'TO'): 855,
11. ('TO', 'NN'): 734,
12. ('NN', 'RB'): 2431,
13. ('RB', 'NN'): 358,
14. ('TO', 'RB'): 200
15. }
16. sorted_tags = sorted(tags)
17. for i in range(num_tags):
18. for j in range(num_tags):
19. tag_tuple = (sorted_tags[i], sorted_tags[j])
20. transition_matrix[i, j] = transition_counts.get(tag_tuple)
21. rows_sum = transition_matrix.sum(axis=1, keepdims=True)
22. transition_matrix = transition_matrix / rows_sum

```

## Viterbi Algorithm



- Probability for this sequence of hidden states:  $0.15 * 0.25 * 0.08 * 0.1 = 0.0003$

- Procedures:
  - o Initialization
  - o Forward pass
  - o Backward pass
- Two matrices populating:

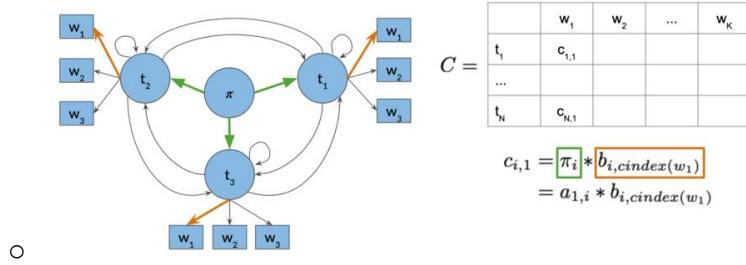
$$C = \begin{array}{|c|c|c|c|} \hline & w_1 & w_2 & \dots & w_K \\ \hline t_1 & & & & \\ \hline \dots & & & & \\ \hline t_N & & & & \\ \hline \end{array}$$
  

$$D = \begin{array}{|c|c|c|c|} \hline & w_1 & w_2 & \dots & w_K \\ \hline t_1 & & & & \\ \hline \dots & & & & \\ \hline t_N & & & & \\ \hline \end{array}$$

o

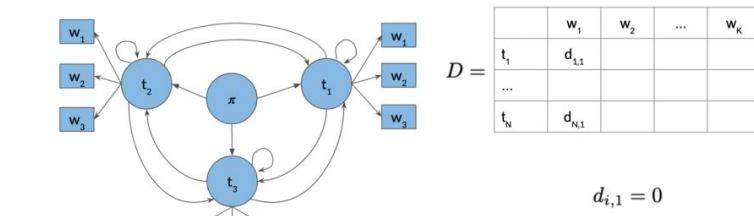
### Viterbi Algorithm: initialization

- C matrix: (Transition and Emission Probabilities)



- First column: the probability of transitions from startstates to the first tag  $t_i$  and  $w_1$ .

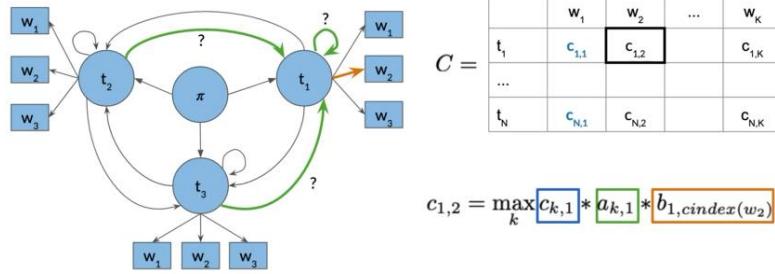
- D matrix:



- Store labels representing the different states traversed while finding the most likely sequence of POS tags for a given seqquence of words,  $w_1$  to  $w_k$ .
- First column: all 0 as there are no preceding POS tags.

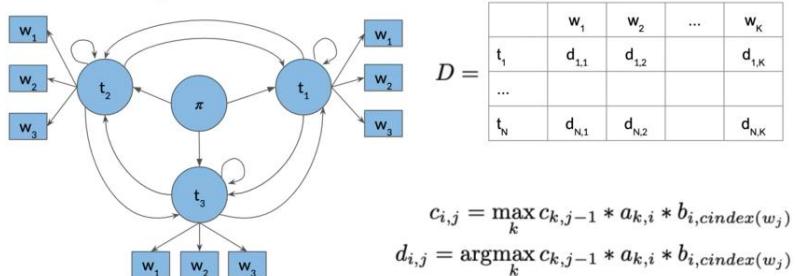
## Viterbi: Forward Pass

- C matrix:



- o  $b_{1,cindex(w_2)}$ : the emission probability from tag  $t_1$  to  $w_2$
- o  $a_{k,1}$ : the transition probability from the POS tag  $t_k$  to current tag  $t_1$
- o  $c_{k,1}$ : the probability for the preceding path you've traversed

- D matrix:



- o Argmax for k maximizing c

## Viterbi: Backward Pass

|                | w <sub>1</sub>   | w <sub>2</sub>   | ... | w <sub>K</sub>   |
|----------------|------------------|------------------|-----|------------------|
| t <sub>1</sub> | c <sub>1,1</sub> | c <sub>1,2</sub> |     | c <sub>1,K</sub> |
| ...            |                  |                  |     |                  |
| t <sub>N</sub> | c <sub>N,1</sub> | c <sub>N,2</sub> |     | c <sub>N,K</sub> |

$$s = \operatorname{argmax}_i c_{i,K}$$

- Example:

|                | w <sub>1</sub> | w <sub>2</sub> | w <sub>3</sub> | w <sub>4</sub> | w <sub>5</sub> |
|----------------|----------------|----------------|----------------|----------------|----------------|
| t <sub>1</sub> | 0.25           | 0.125          | 0.025          | 0.0125         | 0.01           |
| t <sub>2</sub> | 0.1            | 0.025          | 0.05           | 0.01           | 0.003          |
| t <sub>3</sub> | 0.3            | 0.05           | 0.025          | 0.02           | 0.0000         |
| t <sub>4</sub> | 0.2            | 0.1            | 0.000          | 0.0025         | 0.0003         |

$$s = \operatorname{argmax}_i c_{i,K} = 1$$

o

|                | w <sub>1</sub> | w <sub>2</sub> | w <sub>3</sub> | w <sub>4</sub> | w <sub>5</sub> |
|----------------|----------------|----------------|----------------|----------------|----------------|
| t <sub>1</sub> | 0              | 1              | 3              | 2              | 3              |
| t <sub>2</sub> | 0              | 2              | 4              | 1              | 3              |
| t <sub>3</sub> | 0              | 2              | 4              | 1              | 4              |
| t <sub>4</sub> | 0              | 4              | 4              | 3              | 1              |

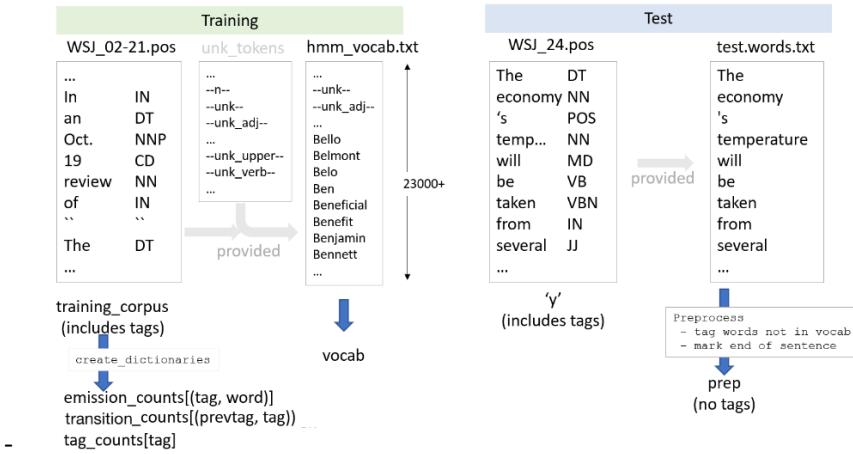
<s> w1 w2 w3 w4 w5

$\pi \leftarrow t_2 \leftarrow t_3 \leftarrow t_1 \leftarrow t_3 \leftarrow t_1$

- Implementation:

- o Use log probabilities instead.

## Parts-of-Speech Tagging (POS) lab



training\_corpus:

- ['In\tn', 'an\tDT', 'Oct.\tNNP', '19\tCD', 'review\tNN', ...]

voc\_l:

- ['!', '#', '\$', '%', '&', "", "", "40s", ..., 'zip', 'zombie', 'zone', 'zones', 'zoning', '{}', '}', '"]

vocab (dictionary):

```

:0
!:1
#:2
$:3
%:4
&:5
':6
...

```

y:

- ['The\tDT', 'economy\tNN', "'s\tPOS", 'temperature\tNN', ...]

Prep:

- ['The', 'economy', "'s", 'temperature', 'will', 'be', 'taken', 'from', 'several', '--unk--', ...]

Create dictionaries:

```

1. def create_dictionaries(training_corpus, vocab, verbose=True):
2. emission_counts = defaultdict(int)
3. transition_counts = defaultdict(int)
4. tag_counts = defaultdict(int)
5. prev_tag = '--s--'
6. i = 0
7. for word_tag in training_corpus:
8. i += 1
9. if i % 5000 == 0 and verbose:
10. print(f"word count = {i}")
11. word, tag = get_word_tag(word_tag, vocab)
12. transition_counts[(prev_tag, tag)] += 1
13. emission_counts[(tag, word)] += 1
14. tag_counts[tag] += 1
15. prev_tag = tag
16. return emission_counts, transition_counts, tag_counts

```

transition examples:  
 ('--s--', 'IN'), 5050  
 ('IN', 'DT'), 32364  
 ('DT', 'NNP'), 9044

emission examples:  
 ('DT', 'any'), 721  
 ('NN', 'decrease'), 7  
 ('NN', 'insider-trading'), 5

ambiguous word example:  
 ('RB', 'back') 304  
 ('VB', 'back') 20  
 ('RP', 'back') 84  
 ('JJ', 'back') 25  
 ('NN', 'back') 29  
 ('VBP', 'back') 4

```

1. print("transition examples: ")
2. for ex in list(transition_counts.items())[:3]:
3. print(ex)
4. print()
5. print("emission examples: ")
6. for ex in list(emission_counts.items())[200:203]:
7. print(ex)
8. print()
9. print("ambiguous word example: ")
10. for tup,cnt in emission_counts.items():
11. if tup[1] == 'back': print(tup, cnt)

transition examples:
('---', 'IN'), 5050
('IN', 'DT'), 32364
('DT', 'NNP'), 9044

emission examples:
('DT', 'any'), 721
('NN', 'decrease'), 7
('NN', 'insider-trading'), 5

ambiguous word example:
('RB', 'back') 304
('VB', 'back') 20
('RP', 'back') 84
('JJ', 'back') 25
('NN', 'back') 29
('VBP', 'back') 4

```

states:

- Parts-of-speech designations found in the training data
  - o ['#', '\$', '""', '(', ')', ',', '--s--', ':', '!', 'CC', 'CD', 'DT', 'EX', ...]
  - o <https://github.com/clips/MBSP/blob/master/tags.py>

## Predict POS

```

1. def predict_pos(prep, y, emission_counts, vocab, states):
2. num_correct = 0
3. all_words = set(emission_counts.keys())
4. total = 0
5. for word, y_tup in zip(prep, y):
6. y_tup_l = y_tup.split()
7. if len(y_tup_l) == 2:
8. true_label = y_tup_l[1]
9. else:
10. continue
11. count_final = 0
12. pos_final = ''
13. if word in vocab:
14. for pos in states:
15. key = (pos, word)
16. if key in emission_counts:
17. count = emission_counts[key]
18. if count > count_final:
19. count_final = count
20. pos_final = pos
21. if pos_final == true_label:
22. num_correct += 1
23. total += 1
24. accuracy = num_correct / total
25. return accuracy

```

## HMM for POS

- Constructing "A": transition probabilities matrix

$$\circ \quad P(t_i|t_{i-1}) = \frac{C(t_{i-1}, t_i) + \alpha}{C(t_{i-1}) + \alpha * N}$$

```

1. def create_transition_matrix(alpha, tag_counts, transition_counts):
2. all_tags = sorted(tag_counts.keys())
3. num_tags = len(all_tags)
4. A = np.zeros((num_tags, num_tags))
5. trans_keys = set(transition_counts.keys())
6. for i in range(num_tags):
7. for j in range(num_tags):
8. count = 0
9. key = (all_tags[i], all_tags[j])
10. if key in trans_keys:
11. count = transition_counts[key]
12. count_prev_tag = tag_counts[key[0]]
13. A[i, j] = (count+alpha)/(count_prev_tag+num_tags*(alpha))
14. return A

```

- Constructing "B": emission probabilities matrix

$$\circ \quad P(w_i|t_i) = \frac{C(t_i, word_i) + \alpha}{C(t_i) + \alpha * N}$$

```

1. def create_emission_matrix(alpha, tag_counts, emission_counts, vocab):
2. num_tags = len(tag_counts)
3. all_tags = sorted(tag_counts.keys())
4. num_words = len(vocab)
5. B = np.zeros((num_tags, num_words))
6. emis_keys = set(list(emission_counts.keys()))
7. for i in range(num_tags):
8. for j in range(num_words):
9. count = 0
10. key = (all_tags[i], vocab[j]) # tuple of form (tag,word)
11. if key in emission_counts:
12. count = emission_counts[key]
13. count_tag = tag_counts[key[0]]
14. B[i, j] = (count+alpha)/(count_tag+alpha*num_words)
15. return B

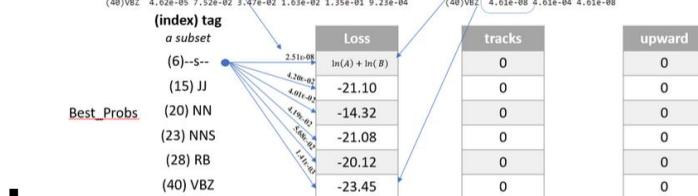
```

- Viterbi Initialization:

- o Two matrices: best\_probs, best\_paths(all zeros)

- $\text{best\_probs}[i, 0] = \log(A[s_{idx}, i]) + \log(B[i, vocab[corpus[0]]])$

| A – Transitions Matrix (subset) |          |          |          | B – Emissions Matrix (subset) |          |          |        |
|---------------------------------|----------|----------|----------|-------------------------------|----------|----------|--------|
| -5--                            | 23       | NN       | IHS      | RB                            | VBD      | LQSS     | tracks |
| (15)JJ                          | 6.54e-05 | 4.47e-02 | 4.50e-02 | 4.39e-02                      | 5.68e-02 | 1.41e-03 | upward |
| (28)NN                          | 4.96e-04 | 5.79e-03 | 1.20e-02 | 7.78e-02                      | 1.83e-02 | 4.37e-03 |        |
| (28)IHS                         | 4.96e-04 | 5.79e-03 | 1.20e-02 | 7.78e-02                      | 1.83e-02 | 8.76e-03 |        |
| (28)RB                          | 4.52e-04 | 1.83e-01 | 1.10e-02 | 4.13e-03                      | 7.31e-02 | 3.95e-02 |        |
| (48)VBD                         | 4.62e-05 | 7.52e-02 | 3.47e-02 | 1.63e-02                      | 1.35e-02 | 9.23e-02 |        |

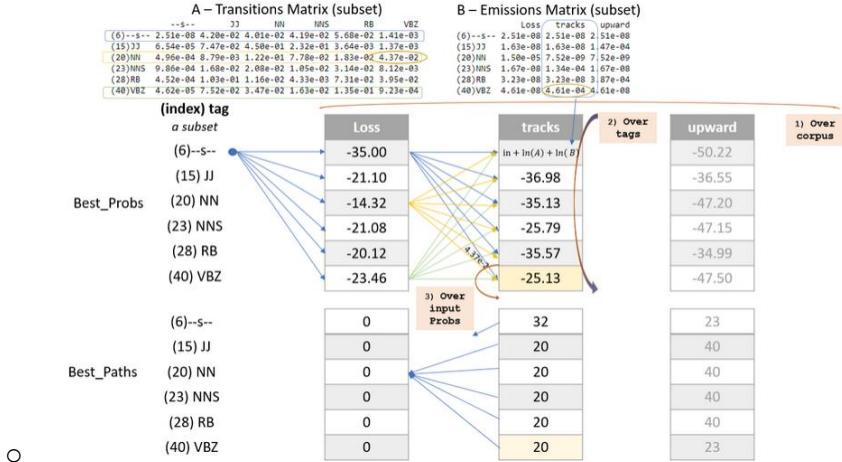


```

1. def initialize(states, tag_counts, A, B, corpus, vocab):
2. num_tags = len(tag_counts)
3. best_probs = np.zeros((num_tags, len(corpus)))
4. best_paths = np.zeros((num_tags, len(corpus)), dtype=int)
5. s_idx = states.index("--S--")
6. for i in range(num_tags):
7. best_probs[i, 0] = np.log(A[s_idx, i])+np.log(B[i,vocab[corpus[0]]]))
8. return best_probs, best_paths

```

- Viterbi Forward:

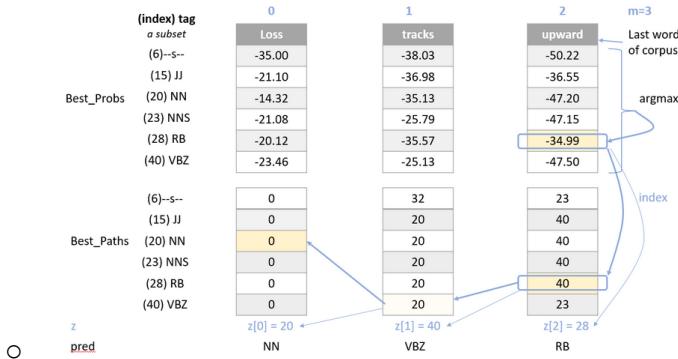


```

1. def viterbi_forward(A, B, test_corpus, best_probs, best_paths, vocab, verbose=True):
2. num_tags = best_probs.shape[0]
3. for i in range(1, len(test_corpus)):
4. if i % 5000 == 0 and verbose:
5. print("Words processed: {:>8}".format(i))
6. for j in range(num_tags):
7. best_prob_i = float("-inf")
8. best_path_i = None
9. for k in range(num_tags):
10. prob = best_probs[k, i-1] + math.log(A[k,j]) + math.log(B[j, vocab[test_corpus[i]]])
11. if prob > best_prob_i:
12. best_prob_i = prob
13. best_path_i = k
14. best_probs[j,i] = best_prob_i
15. best_paths[j,i] = best_path_i
16. return best_probs, best_paths

```

- Viterbi Backward:



```

1. def viterbi_backward(best_probs, best_paths, corpus, states):
2. m = best_paths.shape[1]
3. z = [None] * m
4. num_tags = best_probs.shape[0]
5. best_prob_for_last_word = float('-inf')
6. pred = [None] * m
7. for k in range(num_tags):
8. if best_probs[k, -1] > best_prob_for_last_word:
9. best_prob_for_last_word = best_probs[k, -1]
10. z[m-1] = k
11. pred[m-1] = states[z[m-1]]
12. for i in range(m-1, 0, -1):
13. pos_tag_for_word_i = z[i]
14. z[i-1] = best_paths[pos_tag_for_word_i, i]
15. pred[i-1] = states[z[i-1]]
16. return pred

```

Compute Accuracy:

```
1. def compute_accuracy(pred, y):
2. num_correct = 0
3. total = 0
4. for prediction, y in zip(pred, y):
5. word_tag_tuple = y.split()
6. if len(word_tag_tuple)!=2:
7. continue
8. word, tag = word_tag_tuple[0], word_tag_tuple[1]
9. if tag==prediction:
10. num_correct += 1
11. total += 1
12. return num_correct/total
```

## AutoComplete: N-Grams

### Sequence notation

Corpus: This is great ... teacher drinks tea.  
 $w_1 \ w_2 \ w_3 \dots w_{498} \ w_{499} \ w_{500}$   $m = 500$

$$w_1^m = w_1 \ w_2 \dots w_m$$

$$w_1^3 = w_1 \ w_2 \ w_3$$

$$\dots \quad w_{m-2}^m = w_{m-2} \ w_{m-1} \ w_m$$

### N-grams

- An N-gram is a sequence of N words (can be characters)

- Example:

- Corpus: I am happy because I am learning

- Unigrams: {I, am, happy, ...}

- Probability:

Corpus: I am happy because I am learning

Size of corpus m = 7

$$P(I) = \frac{2}{7}$$

$$P(happy) = \frac{1}{7}$$

$$\text{Probability of unigram: } P(w) = \frac{C(w)}{m}$$

○

- Bigrams: {I am, am happy, ...}

- Probability:

Corpus: I am happy because I am learning

$$P(am|I) = \frac{C(I am)}{C(I)} = \frac{2}{2} = 1$$

$$P(happy|I) = \frac{C(I happy)}{C(I)} = \frac{0}{2} = 0 \quad \text{X} \ I \ happy$$

$$P(learning|am) = \frac{C(am learning)}{C(am)} = \frac{1}{2}$$

$$\text{Probability of a bigram: } P(y|x) = \frac{C(x \ y)}{\sum_w C(x \ w)} = \frac{C(x \ y)}{C(x)}$$

○

- Trigrams: {I am happy, am happy because, ...}

- Probability:

Corpus: I am happy because I am learning

$$P(happy|I am) = \frac{C(I am happy)}{C(I am)} = \frac{1}{2}$$

$$\text{Probability of a trigram: } P(w_3|w_1^2) = \frac{C(w_1^2 w_3)}{C(w_1^2)}$$

$$C(w_1^2 w_3) = C(w_1 w_2 w_3) = C(w_1^3)$$

○

- N-gram probability:

$$P(w_N|w_1^{N-1}) = \frac{C(w_1^{N-1} w_N)}{C(w_1^{N-1})}$$

$$C(w_1^{N-1} w_N) = C(w_1^N)$$

●

## Sequence Probability:

### Version 1:

- $P(I \text{ ate an apple}) = P(I) * P(\text{ate}|I) * P(\text{an}|I \text{ ate}) * P(\text{apple}|I \text{ ate an})$
- Problem:
  - o Corpus almost never contains the exact sentence we're interested in or even its longer subsequences.

### Approximation of sequence probability

- The product of the probabilities of bigrams.
- $P(\text{apple}|I \text{ ate an}) \approx P(\text{apple} | \text{ an})$

### Markov assumption: only last N words matter

- Bigram:
  - o  $P(w_n|w_1^{n-1}) \approx P(w_n|w_{n-1})$
- N-gram:
  - o  $P(w_n|w_1^{n-1}) \approx P(w_n|w_{n-N+1}^{n-1})$
- Entire sentence with bigram:
 
$$P(w_1^n) \approx \prod_{i=1}^n P(w_i|w_{i-1})$$
  - o ...

### Starting and ending sentences

- Start of sentence token <s>
  - o For sentence: I ate an apple
  - o Bigram: <s>I ate an apple
    - $P(<s> \text{I ate an apple}) \approx P(I | <s>) * P(\text{ate} | I) * P(\text{an} | \text{ate}) * P(\text{apple} | \text{an})$
  - o Trigram: <s> <s> I ate an apple
    - $P(w_1 | <s> <s>) P(w_2 | <s> w_1) \dots P(w_n | w_{n-2} w_{n-1})$
  - o N-gram model: add N-1 start tokens <s>
- End of sentence token </s>
  - o Motivation:
    - $P(y|x) = \frac{C(x,y)}{\sum_w C(x,w)} = \frac{C(x,y)}{C(x)}$
    - Does not work when:
      - Corpus: <s> Lyn<sup>drinks</sup> chocolate  
<s> John<sup>drinks</sup>
      - Sentences of length 3:  $P(<s> \text{ yes yes yes}) = \dots$   
 $\begin{array}{ll} <s> \text{ yes no} & <s> \text{ yes yes yes} \\ <s> \text{ yes yes} & <s> \text{ yes yes no} \\ <s> \text{ no no} & \dots \end{array}$   
 $P(<s> \text{ yes yes no}) = \dots$   
 $\dots = \dots$   
 $P(<s> \text{ no no no}) = \dots$
      - $\sum_w C(drinks, w) = 1$   
 $C(drinks) = 2$
      - $\sum_{3 \text{ word}} P(\dots) = 1$

- Solution:
  - Bigram
    - $<s> \text{ the teacher} \text{ drinks tea} \Rightarrow <s> \text{ the teacher} \text{ drinks tea} </s>$
    - $P(\text{the}|<s>) P(\text{teacher}|\text{the}) P(\text{drinks}|\text{teacher}) P(\text{tea}|\text{drinks}) P(</s>|\text{tea})$
  - Corpus: <s> Lyn<sup>drinks</sup> chocolate </s>  
<s> John<sup>drinks</sup> </s>
    - $\sum_w C(drinks, w) = 2$
    - $C(drinks) = 2$

## Example - bigram

Corpus

```
<s> Lyn drinks chocolate </s>
<s> John drinks tea </s>
<s> Lyn eats chocolate </s>
```

$$P(John|<s>) = \frac{1}{3}$$

$$P(chocolate|eats) = \frac{1}{1}$$

$$P(sentence) = \frac{2}{3} * \frac{1}{2} * \frac{1}{2} * \frac{2}{2} = \frac{1}{6}$$

$$P(</s>|tea) = \frac{1}{1}$$

$$P(Lyn|<s>) = ? = \frac{2}{3}$$

•

## N-grams Corpus Preprocessing

```
1. corpus = "Learning% makes 'me' happy. I am happy be-cause I am learning! :)"
2. corpus = corpus.lower()
3. corpus = re.sub(r"[^a-zA-Z0-9.?]+", "", corpus)
4. # learning makes me happy. i am happy because i am learning!
5. input_date="Sat May 9 07:33:35 CEST 2020"
6. date_parts = input_date.split(" ")
7. # ['Sat', 'May', ' ', '9', '07:33:35', 'CEST', '2020']
8. time_parts = date_parts[4].split(":")
9. # ['07', '33', '35']
10. sentence = 'i am happy because i am learning.'
11. tokenized_sentence = nltk.word_tokenize(sentence)
12. # ['i', 'am', 'happy', 'because', 'i', 'am', 'learning', '.']
13. sentence = ['i', 'am', 'happy', 'because', 'i', 'am', 'learning', '.']
14. word_lengths = [(word, len(word)) for word in sentence]
15. # [('i', 1), ('am', 2), ('happy', 5), ('because', 7), ..., ('.', 1)]
```

Trigram:

```
1. def sentence_to_trigram(tokenized_sentence):
2. for i in range(len(tokenized_sentence) - 3 + 1):
3. trigram = tokenized_sentence[i : i + 3]
4. print(trigram)
```

## The N-gram Language Model

- Count matrix

- Rows: unique corpus (N-1)-grams
  - Columns: unique corpus words

Bigram count matrix

Corpus: <s> I study I learn </s>

|       | <s> | </s> | I | study | learn |
|-------|-----|------|---|-------|-------|
| <s>   | 0   | 0    | 1 | 0     | 0     |
| </s>  | 0   | 0    | 0 | 0     | 0     |
| I     | 0   | 0    | 0 | 1     | 1     |
| study | 0   | 0    | 1 | 0     | 0     |
| learn | 0   | 1    | 0 | 0     | 0     |



- Probability matrix

$$P(w_n | w_{n-N+1}^{n-1}) = \frac{C(w_{n-N+1}^{n-1}, w_n)}{C(w_{n-N+1}^{n-1})}$$



$$\text{sum}(row) = \sum_{w \in V} C(w_{n-N+1}^{n-1}, w) = C(w_{n-N+1}^{n-1})$$



Count matrix (bigram)

|       | <s> | </s> | I | study | learn | sum |
|-------|-----|------|---|-------|-------|-----|
| <s>   | 0   | 0    | 1 | 0     | 0     | 1   |
| </s>  | 0   | 0    | 0 | 0     | 0     | 0   |
| I     | 0   | 0    | 0 | 1     | 1     | 2   |
| study | 0   | 0    | 1 | 0     | 0     | 1   |
| learn | 0   | 1    | 0 | 0     | 0     | 1   |



Probability matrix

|       | <s> | </s> | I | study | learn |
|-------|-----|------|---|-------|-------|
| <s>   | 0   | 0    | 1 | 0     | 0     |
| </s>  | 0   | 0    | 0 | 0     | 0     |
| I     | 0   | 0    | 0 | 0.5   | 0.5   |
| study | 0   | 0    | 1 | 0     | 0     |
| learn | 0   | 1    | 0 | 0     | 0     |

- Language model

Sentence probability:

<s> I learn </s>

$$P(\text{sentence}) =$$

$$P(I|<s>)P(\text{learn}|I)P(</s>|\text{learn}) =$$

$$1 \times 0.5 \times 1 =$$

$$0.5$$



o Use log probability to prevent numerical underflow.

- Generative Language model:

Corpus:

1. <s>, Lyn or (<s>, John)?
2. (Lyn.eats) or (Lyn.drinks) ?
3. (drinks.tea) or (drinks.chocolate) ?
4. (tea,</s>) - always



o Algorithm:

- Choose sentence start.
- Choose next bigram starting with previous word.
- Continue until </s> is picked.

## Language Model Evaluation

- Train/val/test split

- o Ratio:

- For smaller corpora: 80:10:10
    - For larger corpora: 98:1:1

- o Two splitting methods:

- Split the corpus by choosing longer continuous segments like Wikipedia articles.
    - Randomly choose short sequences of words.

- Perplexity

- o A measure of the complexity in a sample of texts.

- o Used to tell us whether a set of sentences look like they were written by humans rather than by a simple program choosing words at random.

$$PP(W) = P(s_1, s_2, \dots, s_m)^{-\frac{1}{m}}$$

$W \rightarrow$  test set containing  $m$  sentences  $s$   
 $s_i \rightarrow$  i-th sentence in the test set, each ending with </s>  
 $m \rightarrow$  number of all words in entire test set  $W$  including  
</s> but not including <s>

human: lower, machine: higher

- o Higher language model estimates the probability of the test set  $\rightarrow$  lower

- Perplexity example:
  - E.g. m=100
 
$$P(W) = 0.9 \Rightarrow PP(W) = 0.9^{-\frac{1}{100}} = 1.00105416$$

$$P(W) = 10^{-250} \Rightarrow PP(W) = (10^{-250})^{-\frac{1}{100}} \approx 316$$
  - o Smaller perplexity = better model
  - o Character level models  $PP < \text{word-based models } PP$
- Perplexity for bigram models

$$PP(W) = \sqrt[m]{\prod_{i=1}^m \prod_{j=1}^{|s_i|} \frac{1}{P(w_j^{(i)} | w_{j-1}^{(i)})}}$$

$w_j^{(i)} \rightarrow j\text{-th word in } i\text{-th sentence}$

- concatenate all sentences in W

$$PP(W) = \sqrt[m]{\prod_{i=1}^m \frac{1}{P(w_i | w_{i-1})}}$$

$w_i \rightarrow i\text{-th word in test set}$

- Log perplexity

$$\log PP(W) = -\frac{1}{m} \sum_{i=1}^m \log_2(P(w_i | w_{i-1}))$$

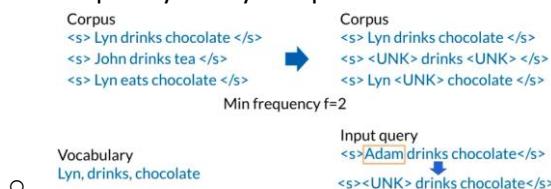
- Example:

|                                                                                                                                                                                                                                                                                                  |
|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <b>Unigram</b><br>Months the my and issue of year foreign new exchange's september were recession ex-change new endorsed a acquire to six executives                                                                                                                                             |
| <b>Bigram</b><br>Last December through the way to preserve the Hudson corporation N. B. E. C. Taylor would seem to complete the major central planners one point five percent of U. S. E. has already old M. X. corporation of living on information such as more frequently fishing to keep her |
| <b>Trigram</b><br>They also point to ninety nine point six billion dollars from two hundred four oh six three percent of the rates of interest stores as Mexico and Brazil on market conditions                                                                                                  |

- o Unigram: 962, Bigram: 170, Trigram: 109

## Out of vocabulary words

- Using <UNK> in corpus:
  - o Create vocabulary V, replace any word in corpus and not in V by <UNK>, count the probabilities with <UNK> as with any other word.
- Create vocabulary V:
  - o Criteria: min word frequency f / Max |V|, including words by frequency
  - o Use <UNK> sparingly
  - o Perplexity – only compare LMs with the same V



## Smoothing

- Missing N-grams in training corpus
  - o Problem: N-grams made of known words still might be missing in the training corpus
    - “John”, “eats” in corpus, but “John eats” missing
- Add-one smoothing (Laplacian smoothing)
  - o 
$$P(w_n|w_{n-N+1}^{n-1}) = \frac{C(w_{n-N+1}^{n-1}, w_n)}{C(w_{n-N+1}^{n-1})}$$
 ← Can be 0
- Add-k smoothing
  - o 
$$P(w_n|w_{n-1}) = \frac{C(w_{n-1}, w_n) + k}{\sum_{w \in V} (C(w_{n-1}, w) + k)} = \frac{C(w_{n-1}, w_n) + k}{C(w_{n-1}) + k * V}$$
- Kneser-Ney smoothing
- Good-Turing smoothing

## Backoff

- If N-gram missing => use (N-1)-gram
  - o Probability discounting e.g. Katz backoff
  - o “Stupid” backoff

Corpus

<s> Lyn drinks chocolate </s>       $P(chocolate|John \text{ } drinks) = ?$   
 ↓  
 <s> John drinks tea </s>  
 <s> Lyn eats chocolate </s>       $0.4 \times P(chocolate|drinks)$

## Interpolation

- Use all (assign different portions to different grams)

$$\hat{P}(w_n|w_{n-2} \dots w_{n-1}) = \lambda_1 \times P(w_n|w_{n-2} \dots w_{n-1}) + \lambda_2 \times P(w_n|w_{n-1}) + \lambda_3 \times P(w_n) \quad \sum_i \lambda_i = 1$$

- Example:

$$\hat{P}(chocolate|John \text{ } drinks) = 0.7 \times P(chocolate|John \text{ } drinks) + 0.2 \times P(chocolate|drinks) + 0.1 \times P(chocolate)$$

## Building the language model

### Count matrix

```
1. def single_pass_trigram_count_matrix(corpus):
2. bigrams = []
3. vocabulary = []
4. count_matrix_dict = defaultdict(dict)
5. for i in range(len(corpus) - 3 + 1):
6. trigram = tuple(corpus[i : i + 3])
7. bigram = trigram[0 : -1]
8. if not bigram in bigrams:
9. bigrams.append(bigram)
10. last_word = trigram[-1]
11. if not last_word in vocabulary:
12. vocabulary.append(last_word)
13. if (bigram, last_word) not in count_matrix_dict:
14. count_matrix_dict[bigram, last_word] = 0
15. count_matrix_dict[bigram, last_word] += 1
16. count_matrix = np.zeros((len(bigrams), len(vocabulary)))
17. for trigram_key, trigam_count in count_matrix_dict.items():
18. count_matrix[bigrams.index(trigram_key[0]), \
19. vocabulary.index(trigram_key[1])] \
20. = trigam_count
21. count_matrix = pd.DataFrame(count_matrix, index=bigrams, columns=vocabulary)
22. return bigrams, vocabulary, count_matrix
```

```
1. corpus = ['i', 'am', 'happy', 'because', 'i', 'am', 'learning', '.']
2. bigrams, vocabulary, count_matrix = single_pass_trigram_count_matrix(corpus)
3. print(count_matrix)
```

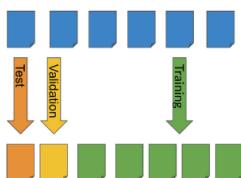
|                  | happy | because | i   | am  | learning | .   |
|------------------|-------|---------|-----|-----|----------|-----|
| (i, am)          | 1.0   | 0.0     | 0.0 | 0.0 | 1.0      | 0.0 |
| (am, happy)      | 0.0   | 1.0     | 0.0 | 0.0 | 0.0      | 0.0 |
| (happy, because) | 0.0   | 0.0     | 1.0 | 0.0 | 0.0      | 0.0 |
| (because, i)     | 0.0   | 0.0     | 0.0 | 1.0 | 0.0      | 0.0 |
| (am, learning)   | 0.0   | 0.0     | 0.0 | 0.0 | 0.0      | 1.0 |

### Probability matrix

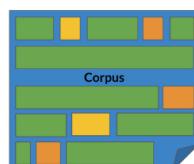
```
1. row_sums = count_matrix.sum(axis=1)
2. prob_matrix = count_matrix.div(row_sums, axis=0)
3. print(prob_matrix)
```

|                  | happy | because | i   | am  | learning | .   |
|------------------|-------|---------|-----|-----|----------|-----|
| (i, am)          | 0.5   | 0.0     | 0.0 | 0.0 | 0.5      | 0.0 |
| (am, happy)      | 0.0   | 1.0     | 0.0 | 0.0 | 0.0      | 0.0 |
| (happy, because) | 0.0   | 0.0     | 1.0 | 0.0 | 0.0      | 0.0 |
| (because, i)     | 0.0   | 0.0     | 0.0 | 1.0 | 0.0      | 0.0 |
| (am, learning)   | 0.0   | 0.0     | 0.0 | 0.0 | 0.0      | 1.0 |

- Continuous text



- Random short sequences



## LM – AutoComplete

```
1. import math
2. import random
3. import numpy as np
4. import pandas as pd
5. import nltk
6. nltk.download('punkt')
7. import w3_unittest
8. nltk.data.path.append('..')
```

- Load and preprocess data
  - o Load and tokenize data

```
1. #data - type: <class: 'str'>
2. #3335477 number of letters

1. def split_to_sentences(data):
2. sentences = data.split("\n")
3. sentences = [s.strip() for s in sentences]
4. sentences = [s for s in sentences if len(s) > 0]
5. return sentences

1. def tokenize_sentences(sentences):
2. tokenized_sentences = []
3. for sentence in sentences: # complete this line
4. sentence = sentence.lower()
5. tokenized = nltk.word_tokenize(sentence)
6. tokenized_sentences.append(tokenized)
7. return tokenized_sentences

1. def get_tokenized_data(data):
2. sentences = split_to_sentences(data)
3. tokenized_sentences = tokenize_sentences(sentences)
4. return tokenized_sentences
```

- o Split the sentences into train and test sets.

```
1. tokenized_data = get_tokenized_data(data)
2. random.seed(87)
3. random.shuffle(tokenized_data)
4. train_size = int(len(tokenized_data) * 0.8)
5. train_data = tokenized_data[0:train_size]
6. test_data = tokenized_data[train_size:]
```

- o Replace words with a low frequency by an unknown marker <UNK>

```
1. def count_words(tokenized_sentences):
2. word_counts = {}
3. for sentence in tokenized_sentences: # complete this line
4. for token in sentence: # complete this line
5. if token not in word_counts.keys():
6. word_counts[token] = 1
7. else:
8. word_counts[token] += 1
9. return word_counts

1. def get_words_with_nplus_frequency(tokenized_sentences,
count_threshold):
2. closed_vocab = []
3. word_counts = count_words(tokenized_sentences)
4. for word, cnt in word_counts.items(): # complete this line
5. if word_counts[word]>=count_threshold:
6. closed_vocab.append(word)
7. return closed_vocab
```

```

1. def replace_oov_words_by_unk(tokenized_sentences, vocabulary,
unknown_token=<unk>):
2. vocabulary = set(vocabulary)
3. replaced_tokenized_sentences = []
4. for sentence in tokenized_sentences:
5. replaced_sentence = []
6. for token in sentence:
7. if token in vocabulary:
8. replaced_sentence.append(token)
9. else:
10. replaced_sentence.append(unknown_token)
11. replaced_tokenized_sentences.append(replaced_sentence)
12. return replaced_tokenized_sentences

1. def preprocess_data(
train_data, test_data, count_threshold, unknown_token=<unk>,
get_words_with_nplus_frequency=get_words_with_nplus_frequency,
replace_oov_words_by_unk=replace_oov_words_by_unk
):
2. vocabulary = get_words_with_nplus_frequency(train_data,
count_threshold)
3. train_data_replaced = replace_oov_words_by_unk(train_data,
vocabulary,unknown_token)
4. test_data_replaced = replace_oov_words_by_unk(test_data,
vocabulary,unknown_token)
5. return train_data_replaced, test_data_replaced, vocabulary

```

- Develop N-gram based language models

- o Compute the count of n-grams from a given data set

```

1. def count_n_grams(data, n, start_token='<s>', end_token = '<e>'):
2. n_grams = {}
3. for sentence in data: # complete this line
4. sentence = [start_token] * n+ sentence + [end_token]
5. sentence = tuple(sentence)
6. m = len(sentence) if n==1 else len(sentence)-n+1
7. for i in range(m): # complete this line
8. n_gram = sentence[i:i+n]
9. if n_gram in n_grams.keys(): # complete this line
10. n_grams[n_gram] += 1
11. else:
12. n_grams[n_gram] = 1
13. return n_grams

```

- o Estimate the conditional probability of a next word with k-smoothing

$$\hat{P}(w_t | w_{t-n} \dots w_{t-1}) = \frac{C(w_{t-n} \dots w_{t-1}, w_t) + k}{C(w_{t-n} \dots w_{t-1}) + k|V|}$$

```

1. def estimate_probability(word, previous_n_gram,
n_gram_counts, n_plus1_gram_counts, vocabulary_size, k=1.0):
2. previous_n_gram = tuple(previous_n_gram)
3. previous_n_gram_count = n_gram_counts[previous_n_gram] if
previous_n_gram in n_gram_counts else 0
4. denominator = previous_n_gram_count + k*vocabulary_size
5. n_plus1_gram = previous_n_gram+(word,)
6. n_plus1_gram_count = n_plus1_gram_counts[n_plus1_gram] if
n_plus1_gram in n_plus1_gram_counts else 0
7. numerator = n_plus1_gram_count + k
8. probability = numerator/denominator
9. return probability

```

```

1. def estimate_probabilities(previous_n_gram, n_gram_counts,
n_plus1_gram_counts, vocabulary, end_token='<e>',
unknown_token='<unk>', k=1.0):
2. previous_n_gram = tuple(previous_n_gram)
3. vocabulary = vocabulary + [end_token, unknown_token]
4. vocabulary_size = len(vocabulary)
5. probabilities = {}
6. for word in vocabulary:
7. probability = estimate_probability(word, previous_n_gram,
n_gram_counts,
8. n_plus1_gram_counts,
9. vocabulary_size, k=k)
10. probabilities[word] = probability
11. return probabilities

1. def make_count_matrix(n_plus1_gram_counts, vocabulary):
2. vocabulary = vocabulary + [<e>, <unk>]
3. n_grams = []
4. for n_plus1_gram in n_plus1_gram_counts.keys():
5. n_gram = n_plus1_gram[0:-1]
6. n_grams.append(n_gram)
7. n_grams = list(set(n_grams))
8. row_index = {n_gram:i for i, n_gram in enumerate(n_grams)}
9. col_index = {word:j for j, word in enumerate(vocabulary)}
10. nrow = len(n_grams)
11. ncol = len(vocabulary)
12. count_matrix = np.zeros((nrow, ncol))
13. for n_plus1_gram, count in n_plus1_gram_counts.items():
14. n_gram = n_plus1_gram[0:-1]
15. word = n_plus1_gram[-1]
16. if word not in vocabulary:
17. continue
18. i = row_index[n_gram]
19. j = col_index[word]
20. count_matrix[i, j] = count
21. count_matrix = pd.DataFrame(count_matrix, index=n_grams, columns=vocabulary)
22. return count_matrix

1. def make_probability_matrix(n_plus1_gram_counts, vocabulary, k):
2. count_matrix = make_count_matrix(n_plus1_gram_counts, unique_words)
3. count_matrix += k
4. prob_matrix = count_matrix.div(count_matrix.sum(axis=1), axis=0)
5. return prob_matrix

```

```

sentences = [['i', 'like', 'a', 'cat'],
['this', 'dog', 'is', 'like', 'a', 'cat']]
unique_words = list(set(sentences[0] + sentences[1]))
bigram_counts = count_n_grams(sentences, 2)
print('bigram counts')
display(make_count_matrix(bigram_counts, unique_words))
bigram counts

```

|        | i   | this | is  | like | dog | a   | cat | <e> | <unk> |
|--------|-----|------|-----|------|-----|-----|-----|-----|-------|
| (dog)  | 0.0 | 0.0  | 1.0 | 0.0  | 0.0 | 0.0 | 0.0 | 0.0 | 0.0   |
| (<e>)  | 1.0 | 1.0  | 0.0 | 0.0  | 0.0 | 0.0 | 0.0 | 0.0 | 0.0   |
| (a)    | 0.0 | 0.0  | 0.0 | 0.0  | 0.0 | 2.0 | 0.0 | 0.0 | 0.0   |
| (like) | 0.0 | 0.0  | 0.0 | 0.0  | 2.0 | 0.0 | 0.0 | 0.0 | 0.0   |
| (cat)  | 0.0 | 0.0  | 0.0 | 0.0  | 0.0 | 0.0 | 2.0 | 0.0 | 0.0   |
| (i)    | 0.0 | 0.0  | 0.0 | 1.0  | 0.0 | 0.0 | 0.0 | 0.0 | 0.0   |
| (this) | 0.0 | 0.0  | 0.0 | 0.0  | 1.0 | 0.0 | 0.0 | 0.0 | 0.0   |
| (is)   | 0.0 | 0.0  | 0.0 | 1.0  | 0.0 | 0.0 | 0.0 | 0.0 | 0.0   |

```

sentences = [['i', 'like', 'a', 'cat'],
['this', 'dog', 'is', 'like', 'a', 'cat']]
unique_words = list(set(sentences[0] + sentences[1]))
bigram_counts = count_n_grams(sentences, 2)
print('bigram probabilities')
display(make_probability_matrix(bigram_counts, unique_words, k=1))
bigram probabilities

```

|        | i        | this     | is       | like     | dog      | a        | cat      | <e>      | <unk>    |
|--------|----------|----------|----------|----------|----------|----------|----------|----------|----------|
| (dog)  | 0.100000 | 0.100000 | 0.200000 | 0.100000 | 0.100000 | 0.100000 | 0.100000 | 0.100000 | 0.100000 |
| (<e>)  | 0.181818 | 0.181818 | 0.090909 | 0.090909 | 0.090909 | 0.090909 | 0.090909 | 0.090909 | 0.090909 |
| (a)    | 0.090909 | 0.090909 | 0.090909 | 0.090909 | 0.090909 | 0.090909 | 0.272727 | 0.090909 | 0.090909 |
| (like) | 0.090909 | 0.090909 | 0.090909 | 0.090909 | 0.090909 | 0.272727 | 0.090909 | 0.090909 | 0.090909 |
| (cat)  | 0.090909 | 0.090909 | 0.090909 | 0.090909 | 0.090909 | 0.090909 | 0.090909 | 0.272727 | 0.090909 |
| (i)    | 0.100000 | 0.100000 | 0.100000 | 0.200000 | 0.100000 | 0.100000 | 0.100000 | 0.100000 | 0.100000 |
| (this) | 0.100000 | 0.100000 | 0.100000 | 0.100000 | 0.200000 | 0.100000 | 0.100000 | 0.100000 | 0.100000 |
| (is)   | 0.100000 | 0.100000 | 0.100000 | 0.200000 | 0.100000 | 0.100000 | 0.100000 | 0.100000 | 0.100000 |

Evaluate the N-gram models by computing the perplexity score

$$PP(W) = \sqrt[N-1]{\prod_{t=n}^{N-1} \frac{1}{P(w_t|w_{t-n} \dots w_{t-1})}}$$

(array indexing starts from 0)

```

1. def calculate_perplexity(sentence, n_gram_counts, n_plus1_gram_counts,
vocabulary_size, start_token='<s>', end_token = '<e>', k=1.0):
2. n = len(list(n_gram_counts.keys())[0])
3. sentence = [start_token] * n + sentence + [end_token]
4. sentence = tuple(sentence)
5. N = len(sentence)
6. product_pi = 1.0
7. for t in range(n, N):
8. n_gram = sentence[t-n:t]
9. word = sentence[t]
10. probability =
estimate_probability(word,n_gram,n_gram_counts,n_plus1_gram_counts,vocabulary_size,k)
11. product_pi *= 1/probability
12. perplexity = (product_pi)**(1/N)
13. return perplexity

```

Auto-complete System

```

1. def suggest_a_word(previous_tokens, n_gram_counts, n_plus1_gram_counts,
vocabulary, end_token='<e>', unknown_token="<unk>", k=1.0, start_with=None):
2. n = len(list(n_gram_counts.keys())[0])
3. previous_tokens = ['<s>'] * n + previous_tokens
4. previous_n_gram = previous_tokens[-n:]
5. probabilities = estimate_probabilities(previous_n_gram,
6. n_gram_counts, n_plus1_gram_counts,
7. vocabulary, k=k)
8. suggestion = None
9. max_prob = 0
10. for word, prob in probabilities.items(): # complete this line
11. if start_with!=None: # complete this line with the proper condition
12. if not word.startswith(start_with):
13. continue
14. if prob > max_prob: # complete this line with the proper condition
15. suggestion = word
16. max_prob = prob
17. return suggestion, max_prob
18.
19. def getSuggestions(previous_tokens, n_gram_counts_list, vocabulary, k=1.0,
start_with=None):
20. model_counts = len(n_gram_counts_list)
21. suggestions = []
22. for i in range(model_counts-1):
23. n_gram_counts = n_gram_counts_list[i]
24. n_plus1_gram_counts = n_gram_counts_list[i+1]
25. suggestion = suggest_a_word(previous_tokens, n_gram_counts,
26. n_plus1_gram_counts, vocabulary,
27. k=k, start_with=start_with)
28. suggestions.append(suggestion)
29. return suggestions

```

## Basic word representation

| Word  | Number |
|-------|--------|
| a     | 1      |
| able  | 2      |
| about | 3      |
| ...   | ...    |
| hand  | 615    |
| ...   | ...    |
| happy | 621    |
| ...   | ...    |
| zebra | 1000   |

- the ordering doesn't make semantic sense

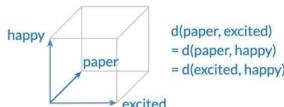
hand < happy < zebra  
615 ?! 621 ?! 1000

○

## One-hot-vectors

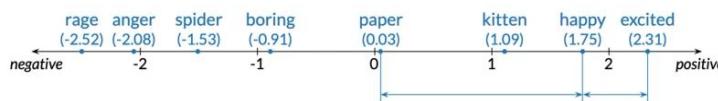
| Word  | Number | "happy"         |
|-------|--------|-----------------|
| a     | 1      | 1 [ 0 ]         |
| able  | 2      | 2 [ 0 ]         |
| about | 3      | 3 [ 0 ]         |
| ...   | ...    | ... [ 0 ]       |
| hand  | 615    | 615 [ 0 ]       |
| ...   | ...    | ... [ 0 ]       |
| happy | 621    | 621 [ 1 ] happy |
| ...   | ...    | ... [ 0 ]       |
| zebra | 1000   | 1000 [ 0 ]      |

- Advantages:
  - No implied ordering
  - Simple
- Disadvantages:
  - Huge (one million word → one million rows)
  - No embedded meaning



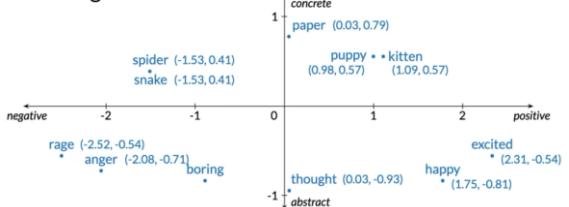
## Word embeddings

### Meaning as vectors



- Happy is more related to excited than paper

### Meaning as vectors



- Possible for two words to be on the same point (e.g. snake, spider)

### Advantages:

- Low dimension (~100)
- Embed meaning
  - Semantic distance
  - Analogies

## How to create word embedding

Prepare:

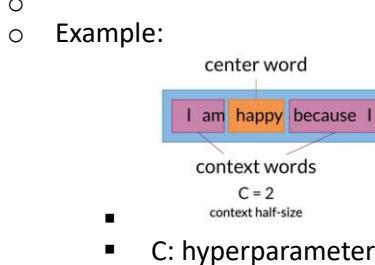
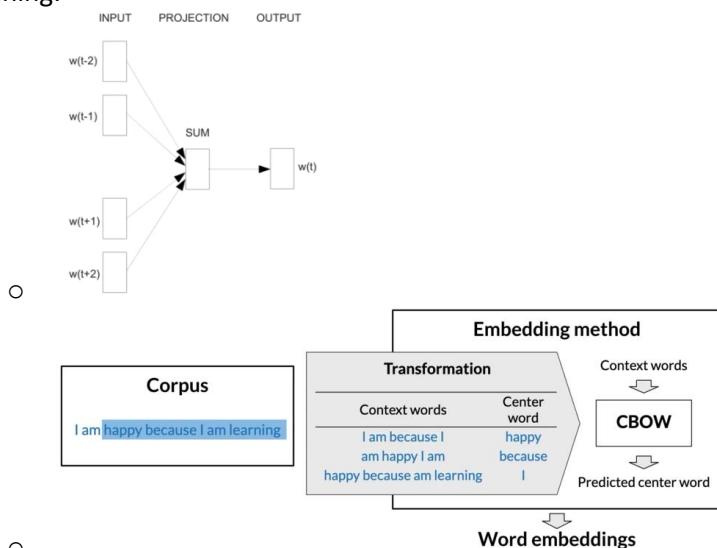
- Corpus:
  - o Words in context.
  - o Could be general-purpose or specialized.
  - o Need to be transformed into machine-readable like vectors before feeding into the model.
- Embedding method:
  - o Machine learning model
    - Performs a learning task → define meaning → word embeddings
    - Self-supervised (unsupervised + supervised)
      - The corpus is unlabeled, but it itself provides the necessary context which would ordinarily make up the labels.
      - Dimension of word embedding is one of the hyperparameters.

Methods:

- word2vec
  - o Continuous bag-of-words (CBOW)
  - o Continuous skip-gram / Skip-gram with negative sampling (SGNS)
- Global Vectors (GloVe)
- fastText
  - o Supports out-of-vocabulary (OOV) words
- Advanced:
  - o Deep learning, contextual embeddings
  - o BERT
  - o ELMo
  - o GPT

## Continuous Bag-of-Words (CBOW) Model

- Center word prediction
  - o Objective: predict the missing based on the surrounding words
- Rationale:
  - o The little \_\_\_\_\_ is barking
    - Dog, puppy, terrier, ...
- Training:



## Cleaning and Tokenization

- Corpus should be case-insensitive:
  - o "The" == "the" == "THE" → lowercase/upper case
- Punctuation:
  - o , ! . ? → .      " ‘ ’ ” → Ø      ... !! ??? → .
- Numbers:
  - o 1 2 3 5 8 → Ø      3.14159 90210 → as is/<NUMBER>
- Special characters:
  - o ∇ \$ € § ¶ \*\* → Ø
- Special words:
  - o Hashtag, emojis

## Sliding window of words in python

```

1. def get_windows(words, C):
2. i = C
3. while i < len(words) - C:
4. center_word = words[i]
5. context_words = words[(i-C):i] + words[(i+1):(i+C+1)]
6. yield context_words, center_word
7. i += 1

```

## Transforming words into vectors

- Corpus → vocabulary → one-hot vector
 

|                                                       |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              |                                                       |                                                       |                                                       |   |          |                                                       |                                                       |                                                       |                                                       |                                                       |
|-------------------------------------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------|-------------------------------------------------------|-------------------------------------------------------|---|----------|-------------------------------------------------------|-------------------------------------------------------|-------------------------------------------------------|-------------------------------------------------------|-------------------------------------------------------|
| Corpus                                                | I am happy because I am learning                                                                                                                                                                                                                                                                                                                                                                                                                                                                             |                                                       |                                                       |                                                       |   |          |                                                       |                                                       |                                                       |                                                       |                                                       |
| Vocabulary                                            | am, because, happy, I, learning                                                                                                                                                                                                                                                                                                                                                                                                                                                                              |                                                       |                                                       |                                                       |   |          |                                                       |                                                       |                                                       |                                                       |                                                       |
| One-hot vector                                        | <table border="0"> <tr> <td>am</td> <td>because</td> <td>happy</td> <td>I</td> <td>learning</td> </tr> <tr> <td><math>\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}</math></td> <td><math>\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}</math></td> <td><math>\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}</math></td> <td><math>\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}</math></td> <td><math>\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}</math></td> </tr> </table> | am                                                    | because                                               | happy                                                 | I | learning | $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ | $\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ | $\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$ | $\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$ | $\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$ |
| am                                                    | because                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      | happy                                                 | I                                                     | learning                                              |   |          |                                                       |                                                       |                                                       |                                                       |                                                       |
| $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ | $\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$                                                                                                                                                                                                                                                                                                                                                                                                                                                        | $\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$ | $\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$ | $\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$ |   |          |                                                       |                                                       |                                                       |                                                       |                                                       |

- Take the average of individual one-hot-vectors

$$\text{O } \left( \begin{array}{c} \text{I} \\ \text{am} \\ \text{because} \\ \text{happy} \\ \text{I} \\ \text{learning} \end{array} \right) + \left( \begin{array}{c} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{array} \right) + \left( \begin{array}{c} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{array} \right) + \left( \begin{array}{c} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{array} \right) / 4 = \left( \begin{array}{c} 0.25 \\ 0.25 \\ 0 \\ 0.5 \\ 0 \end{array} \right)$$

- Final prepared training set

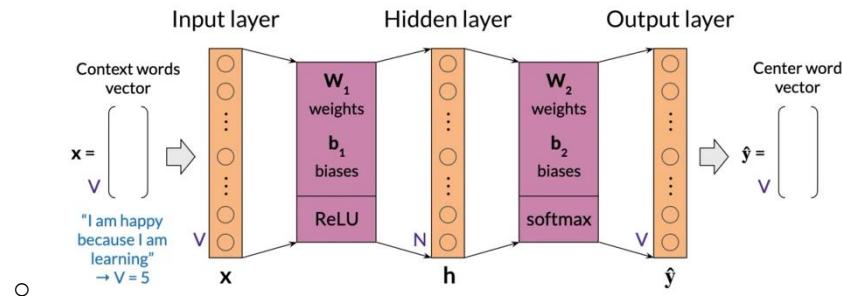
| Context words             | Context words vector        | Center word | Center word vector |
|---------------------------|-----------------------------|-------------|--------------------|
| I am because I            | [0.25; 0.25; 0; 0.5; 0]     | happy       | [0; 0; 1; 0; 0]    |
| am happy I am             | [0.5; 0; 0.25; 0.25; 0]     | because     | [0; 1; 0; 0; 0]    |
| happy because am learning | [0.25; 0.25; 0.25; 0; 0.25] | I           | [0; 0; 0; 1; 0]    |

## CBOW Model

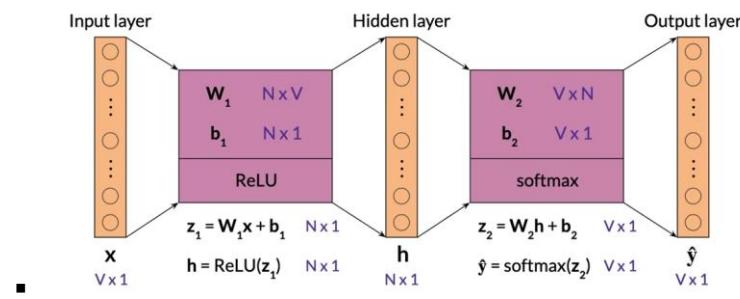
- Architecture

## Architecture of the CBOW model

Hyperparameters  
N: Word embedding size ...



- Dimensions:



- Column vectors

$$z_1 = W_1x + b_1 \quad z_1 = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix} \quad W_1 = \begin{bmatrix} N \times V \\ N \times 1 \end{bmatrix} \quad x = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix} \quad b_1 = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix} \quad N \times 1$$

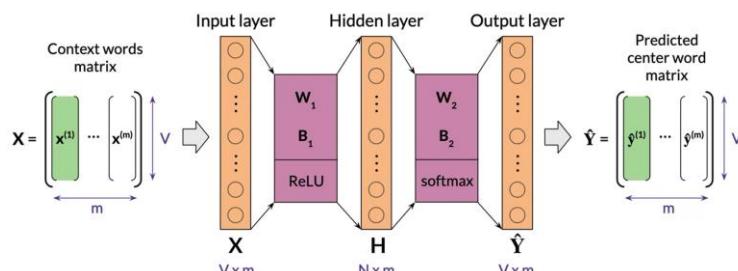
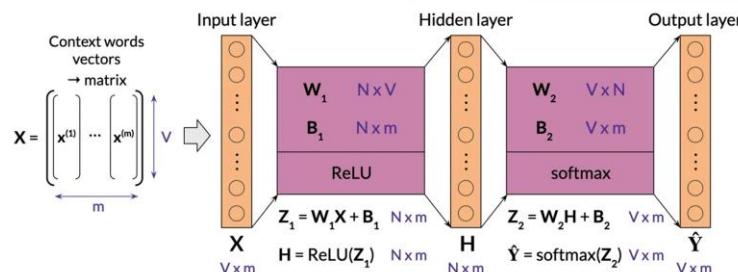
- Row vectors

$$z_1 = xW_1^T + b_1 \quad b_1 = [1 \times N] \quad W_1 = [N \times V] \quad b_1 = [1 \times N] \quad x = [1 \times V]$$

- Dimensions for batch input:

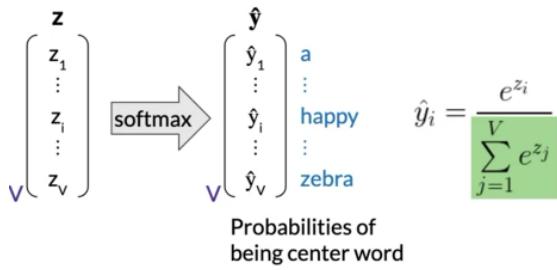
### Dimensions (batch input)

$[b_1] \rightarrow B_1 = \left[ \begin{array}{c} b_1 \\ \vdots \\ b_1 \end{array} \right] \underset{m}{\underbrace{\dots}} \underset{m}{\overbrace{\dots}} \quad N \text{ broadcasting}$



- Activation function:

- Input layer: ReLU
- Hidden layer: softmax



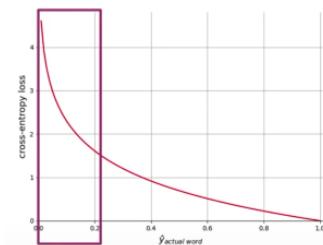
- Cost function:

  - o Cross-entropy loss

$$J = - \sum_{k=1}^V y_k \log \hat{y}_k$$

$$\begin{array}{c} \mathbf{y} \\ \left[ \begin{array}{c} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{array} \right] \end{array} \begin{array}{c} \mathbf{\hat{y}} \\ \left[ \begin{array}{c} 0.083 \\ 0.611 \\ 0.225 \\ 0.05 \\ 0 \end{array} \right] \end{array} \xrightarrow{\text{log}} \begin{array}{c} \log(\mathbf{\hat{y}}) \\ \left[ \begin{array}{c} -2.49 \\ -3.49 \\ -0.49 \\ -1.49 \\ -2.49 \end{array} \right] \end{array} \xrightarrow{\odot \mathbf{y}} \begin{array}{c} \mathbf{y} \odot \log(\mathbf{\hat{y}}) \\ \left[ \begin{array}{c} 0 \\ 0 \\ -0.49 \\ 0 \\ 0 \end{array} \right] \end{array} \xrightarrow{-\Sigma} J = 0.49$$

$$\begin{array}{c} \mathbf{y} \\ \left[ \begin{array}{c} 0 \\ 0 \\ 1 \\ 0 \end{array} \right] \end{array} \begin{array}{c} \mathbf{\hat{y}} \\ \left[ \begin{array}{c} 0.96 \\ 0.01 \\ 0.01 \\ 0.01 \end{array} \right] \end{array} \xrightarrow{\text{log}} \begin{array}{c} \log(\mathbf{\hat{y}}) \\ \left[ \begin{array}{c} -0.04 \\ -4.61 \\ -4.61 \\ -4.61 \end{array} \right] \end{array} \xrightarrow{\odot \mathbf{y}} \begin{array}{c} \mathbf{y} \odot \log(\mathbf{\hat{y}}) \\ \left[ \begin{array}{c} 0 \\ 0 \\ -4.61 \\ 0 \end{array} \right] \end{array} \xrightarrow{-\Sigma} J = 4.61$$



- Training process:

  - o Forward propagation

  - o Cost

$$J_{batch} = -\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^V y_j^{(i)} \log \hat{y}_j^{(i)}$$

$$J_{batch} = \frac{1}{m} \sum_{i=1}^m J^{(i)}$$

  - o Backpropagation and gradient descent

$$\frac{\partial J_{batch}}{\partial \mathbf{W}_1} = \frac{1}{m} (\mathbf{W}_2^\top (\hat{\mathbf{Y}} - \mathbf{Y}) \cdot \text{step}(\mathbf{Z}_1)) \mathbf{X}^\top$$

$$\frac{\partial J_{batch}}{\partial \mathbf{W}_2} = \frac{1}{m} (\hat{\mathbf{Y}} - \mathbf{Y}) \mathbf{H}^\top$$

$$\frac{\partial J_{batch}}{\partial \mathbf{b}_1} = \frac{1}{m} (\mathbf{W}_2^\top (\hat{\mathbf{Y}} - \mathbf{Y}) \cdot \text{step}(\mathbf{Z}_1)) \mathbf{1}_m^\top$$

$$\frac{\partial J_{batch}}{\partial \mathbf{b}_2} = \frac{1}{m} (\hat{\mathbf{Y}} - \mathbf{Y}) \mathbf{1}_m^T$$

$\mathbf{1}_m = [1, \dots, 1]$   
 $\mathbf{A} \cdot \mathbf{1}_m^\top = \left( \begin{array}{c} \mathbf{A} \\ \vdots \\ \mathbf{A} \end{array} \right) \left( \begin{array}{c} 1 \\ \vdots \\ 1 \end{array} \right) = \left[ \begin{array}{c} \mathbf{A} \\ \vdots \\ \mathbf{A} \end{array} \right] = \mathbf{\Sigma}$ 

```
import numpy as np
code to initialize matrix a omitted
np.sum(a, axis=1, keepdims=True)
```

Word Embeddings: Training the CBOW model:

```

1. # size of the word embedding vectors
2. N = 3
3. # size of the vocabulary
4. V = 5
5.
6. # predefined W1 (3 x 5), W2 (5 x 3), b1 (3 x 1), b2 (5 x 1)
7. # W1: (N x V), W2: (V x N), b1: (N x 1), b2: (V x 1)
8.
9.
10. words = ['i', 'am', 'happy', 'because', 'i', 'am', 'learning']
11. word2Ind, Ind2word = get_dict(words)
12. def get_windows(words, C):
13. i = C
14. while i < len(words) - C:
15. center_word = words[i]
16. context_words = words[(i - C):i] + words[(i+1):(i+C+1)]
17. yield context_words, center_word
18. i += 1
19. def word_to_one_hot_vector(word, word2Ind, V):
20. one_hot_vector = np.zeros(V)
21. one_hot_vector[word2Ind[word]] = 1
22. return one_hot_vector
23. def context_words_to_vector(context_words, word2Ind, V):
24. context_words_vectors = [word_to_one_hot_vector(w, word2Ind, V) for w in context_words]
25. context_words_vectors = np.mean(context_words_vectors, axis=0)
26. return context_words_vectors
27. def get_training_example(words, C, word2Ind, V):
28. for context_words, center_word in get_windows(words, C):
29. yield context_words_to_vector(context_words, word2Ind, V),
30. word_to_one_hot_vector(center_word, word2Ind, V)
31.
32. x = x_array.copy()
33. x.shape = (V, 1)
34. y = y_array.copy()
35. y.shape = (V, 1)
36.
37. def relu(z):
38. result = z.copy()
39. result[result < 0] = 0
40. return result
41. def softmax(z):
42. e_z = np.exp(z)
43. sum_e_z = np.sum(e_z)
44. return e_z / sum_e_z
45.
46. z1 = np.dot(W1, x) + b1
47. h = relu(z1)
48. z2 = np.dot(W2, h) + b2
49. y_hat = softmax(z2)
50.
51. def cross_entropy_loss(y_predicted, y_actual):
52. loss = -np.sum(np.dot(y_actual.T,np.log(y_predicted)))
53. return loss

```

Extracting word embedding vectors

Word Embedding:

```

1. words = ['i', 'am', 'happy', 'because', 'i', 'am', 'learning']
2. V = 5
3. word2Ind, Ind2word = get_dict(words)
4. W1 = np.array([[0.41687358, 0.08854191, -0.23495225, 0.28320538, 0.41800106],
5. [0.32735501, 0.22795148, -0.23951958, 0.4117634 , -0.23924344],
6. [0.26637602, -0.23846886, -0.37770863, -0.11399446, 0.34008124]])
7. W2 = np.array([[-0.22182064, -0.43008631, 0.13310965],
8. [0.08476603, 0.08123194, 0.1772054],
9. [0.1871551 , -0.06107263, -0.1790735],
10. [0.07055222, -0.02015138, 0.36107434],
11. [0.33480474, -0.39423389, -0.43959196]])
12. b1 = np.array([[0.09688219],
13. [0.29239497],
14. [-0.27364426]])

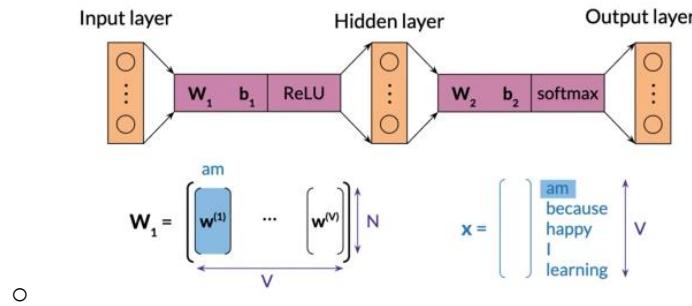
```

```

15. b2 = np.array([[0.0352008],
16. [-0.36393384],
17. [-0.12775555],
18. [-0.34802326],
19. [-0.07017815]])

```

- Option 1:

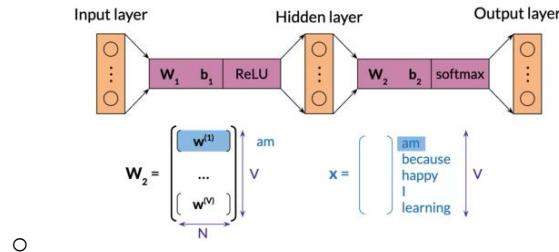


```

1. for word in word2Ind:
2. word_embedding_vector = W1[:, word2Ind[word]]
3. print(f'{word}: {word_embedding_vector}')

```

- Option 2:

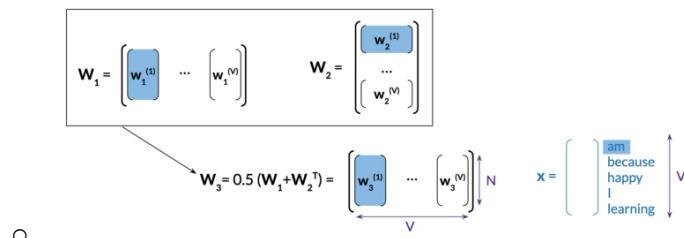


```

1. for word in word2Ind:
2. word_embedding_vector = W2.T[:, word2Ind[word]]
3. print(f'{word}: {word_embedding_vector}')

```

- Option 3:



```

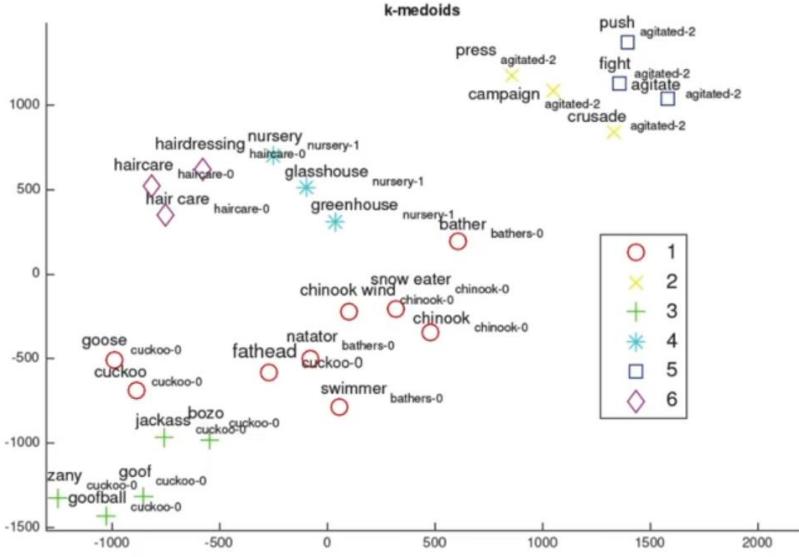
1. W3 = (W1+W2.T)/2
2. for word in word2Ind:
3. word_embedding_vector = W3[:, word2Ind[word]]
4. print(f'{word}: {word_embedding_vector}')

```

## Evaluating word embeddings: intrinsic evaluation

- Test relationships between words
  - o Analogies
    - Semantic analogies:
      - France is to Paris as Italy is to <?>
    - Syntactic analogies:
      - Seen is to saw as been is to <?>
    - There could be ambiguity ( can have many different answers)

- Clustering



### Evaluating Word Embedding: Extrinsic Evaluation

- Test word embedding on external task
  - E.g. named entity recognition, parts-of-speech tagging
- Advantages:
  - Evaluates actual usefulness of embeddings
- Disadvantages:
  - Time-consuming
  - More difficult to troubleshoot

### Word Embedding lab

Utils:

```

1. def get_batches(data, word2Ind, V, C, batch_size):
2. batch_x = []
3. batch_y = []
4. for x, y in get_vectors(data, word2Ind, V, C):
5. if len(batch_x) < batch_size:
6. batch_x.append(x)
7. batch_y.append(y)
8. else:
9. yield np.array(batch_x).T, np.array(batch_y).T
10. batch_x = []
11. batch_y = []

1. def compute_pca(data, n_components=2):
2. m, n = data.shape
3. data -= data.mean(axis=0)
4. R = np.cov(data, rowvar=False)
5. evals, evecs = linalg.eigh(R)
6. idx = np.argsort(evals)[::-1]
7. evecs = evecs[:, idx]
8. evals = evals[idx]
9. evecs = evecs[:, :n_components]
10. return np.dot(evecs.T, data.T).T

```

Load, tokenize and process the data:

```

1. with open('./data/shakespeare.txt') as f:
2. data = f.read()
3. data = re.sub(r'[,!?;-]', '.', data)
4. data = nltk.word_tokenize(data)

```

```
5. data = [ch.lower() for ch in data if ch.isalpha() or ch == '.']
```

Data:

- ['o', 'for', 'a', 'muse', 'of', 'fire', '.', 'that', 'would', 'ascend', 'the', ... ]

Compute the frequency distribution of the words in the dataset

```
1. fdist = nltk.FreqDist(word for word in data)
```

fdist.most\_common(20):

- [('.', 9630), ('the', 1521), ('and', 1394), ('i', 1257), ('to', 1159), ...]

Mapping words to indices and viceversa:

```
1. word2Ind, Ind2word = get_dict(data)
2. #word2Ind['king'] = 2745, Ind2word[2745] = 'king'
```

Initialize model

```
1. def initialize_model(N,V, random_seed=1):
2. np.random.seed(random_seed)
3. W1 = np.random.rand(N,V)
4. W2 = np.random.rand(V,N)
5. b1 = np.random.rand(N,1)
6. b2 = np.random.rand(V,1)
7. return W1, W2, b1, b2
```

Softmax:

```
1. def softmax(z):
2. yhat = np.exp(z)/np.sum(np.exp(z),0) #sum columns
3. return yhat
```

Forward prop:

Implement the forward propagation  $z$  according to equations (1) to (3).

$$h = W_1 X + b_1 \quad (1)$$

$$h = \text{ReLU}(h) \quad (2)$$

$$z = W_2 h + b_2 \quad (3)$$

For that, you will use as activation the Rectified Linear Unit (ReLU) given by:

$$f(h) = \max(0, h) \quad (6)$$

```
1. def forward_prop(x, W1, W2, b1, b2):
2. h = np.dot(W1, x) + b1
3. h = np.maximum(0, h)
4. z = np.dot(W2, h) + b2
5. return z, h
```

Cost function (cross-entropy)

```
1. def compute_cost(y, yhat, batch_size):
2. logprobs = np.multiply(np.log(yhat),y)
3. cost = - 1/batch_size * np.sum(logprobs)
4. cost = np.squeeze(cost)
5. return cost
```

Back prop:

$$\frac{\partial J_{batch}}{\partial \mathbf{W}_1} = \frac{1}{m} (\mathbf{W}_2^\top (\hat{\mathbf{Y}} - \mathbf{Y}) \cdot \text{step}(\mathbf{Z}_1)) \mathbf{X}^\top$$

$$\frac{\partial J_{batch}}{\partial \mathbf{W}_2} = \frac{1}{m} (\hat{\mathbf{Y}} - \mathbf{Y}) \mathbf{H}^\top$$

$$\frac{\partial J_{batch}}{\partial \mathbf{b}_1} = \frac{1}{m} (\mathbf{W}_2^\top (\hat{\mathbf{Y}} - \mathbf{Y}) \cdot \text{step}(\mathbf{Z}_1)) \mathbf{1}_m^\top$$

$$\frac{\partial J_{batch}}{\partial \mathbf{b}_2}$$

$$\mathbf{1}_m = [1, \dots, 1]$$

$$\mathbf{A} \cdot \mathbf{1}_m^\top = \begin{pmatrix} \text{[ ]} \\ \vdots \end{pmatrix} \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} = \begin{pmatrix} \Sigma \end{pmatrix}$$

```
import numpy as np
code to initialize matrix a omitted
np.sum(a, axis=1, keepdims=True)
```

```
1. def back_prop(x, yhat, y, h, W1, W2, b1, b2, batch_size):
2. z1 = np.dot(W1, x) + b1
3. l1 = np.dot(W2.T, (yhat-y))
4. l1[z1 < 0] = 0 # use "l1" to compute gradients below
5. grad_W1 = np.dot(l1, x.T) / batch_size
6. grad_W2 = np.dot((yhat-y), h.T) / batch_size
7. grad_b1 = np.sum(l1, axis=1, keepdims=True) / batch_size
8. grad_b2 = np.sum(yhat-y, axis=1, keepdims = True) / batch_size
9. return grad_W1, grad_W2, grad_b1, grad_b2
```

Gradient Descent:

```
1. def gradient_descent(data, word2Ind, N, V, num_iters, alpha=0.03,
2. random_seed=282, initialize_model=initialize_model,
3. get_batches=get_batches, forward_prop=forward_prop,
4. softmax=softmax, compute_cost=compute_cost,
5. back_prop=back_prop):
6. W1, W2, b1, b2 = initialize_model(N,V, random_seed=random_seed) #W1=(N,V) and W2=(V,N)
7. batch_size = 128
8. iters = 0
9. C = 2
10. for x, y in get_batches(data, word2Ind, V, C, batch_size):
11. z, h = forward_prop(x,W1,W2,b1,b2)
12. yhat = softmax(z)
13. cost = compute_cost(y,yhat,batch_size)
14. if ((iters+1) % 10 == 0):
15. print(f"iters: {iters + 1} cost: {cost:.6f}")
16. grad_W1, grad_W2, grad_b1, grad_b2 = back_prop(x,yhat,y,h,W1,W2,b1,b2,batch_size)
17. W1 = W1-alpha*grad_W1
18. W2 = W2-alpha*grad_W2
19. b1 = b1-alpha*grad_b1
20. b2 = b2-alpha*grad_b2
21. iters +=1
22. if iters == num_iters:
23. break
24. if iters % 100 == 0:
```

```

25. alpha *= 0.66
26. return W1, W2, b1, b2

```

### Visualization

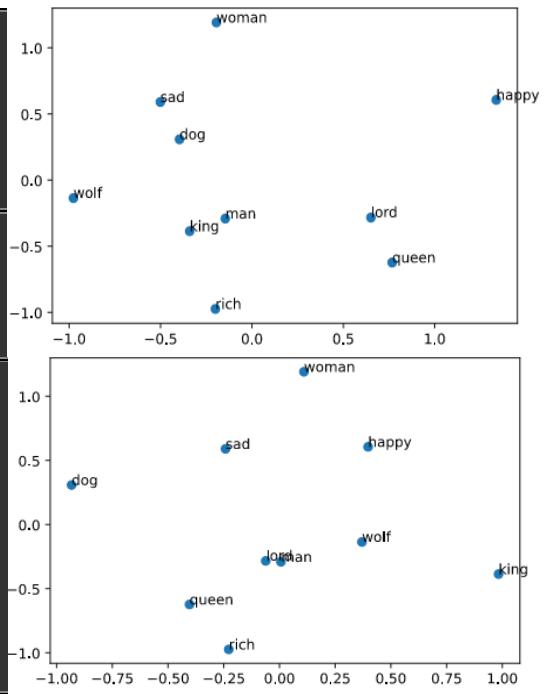
```

1. from matplotlib import pyplot
2. %config InlineBackend.figure_format = 'svg'
3. words = ['king', 'queen', 'lord', 'man', 'woman', 'dog',
4. 'wolf', 'rich', 'happy', 'sad']
5. embs = (W1.T + W2)/2.0
6. idx = [word2Ind[word] for word in words]
7. X = embs[idx, :]

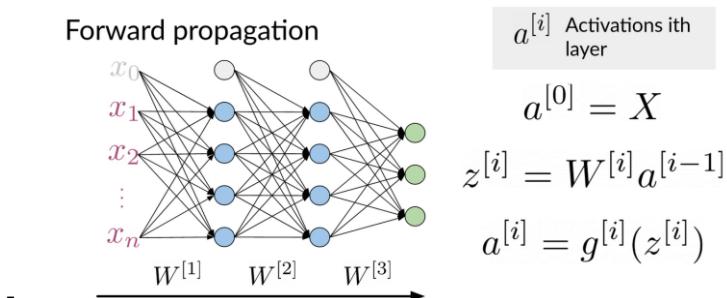
1. result= compute_pca(X, 2)
2. pyplot.scatter(result[:, 0], result[:, 1])
3. for i, word in enumerate(words):
4. pyplot.annotate(word, xy=(result[i, 0], result[i, 1]))
5. pyplot.show()

1. result= compute_pca(X, 4)
2. pyplot.scatter(result[:, 0], result[:, 1])
3. for i, word in enumerate(words):
4. pyplot.annotate(word, xy=(result[i, 0], result[i, 1]))
5. pyplot.show()

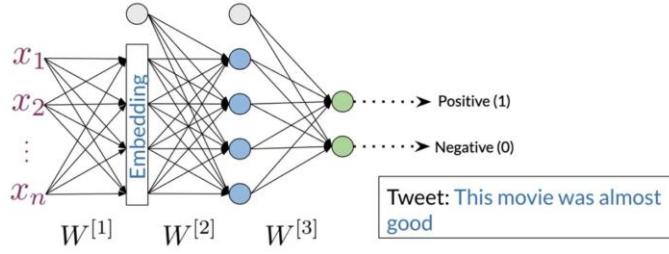
```



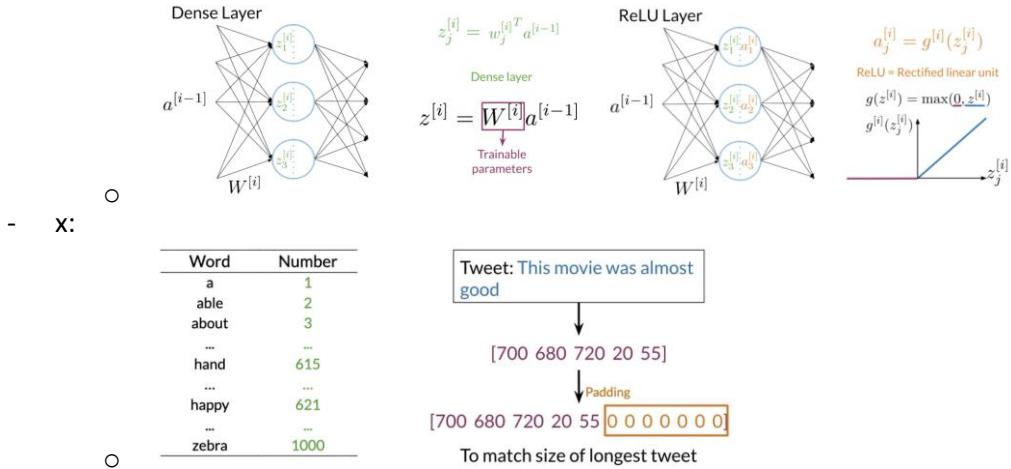
### Neural Network:



Neural Networks for sentiment analysis

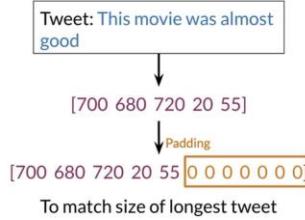


- Uses Dense layer and ReLU Layer



- X:

| Word  | Number |
|-------|--------|
| a     | 1      |
| able  | 2      |
| about | 3      |
| ...   | ...    |
| hand  | 615    |
| ...   | ...    |
| happy | 621    |
| ...   | ...    |
| zebra | 1000   |



## Embedding Layer

Embedding Layer

| Vocabulary | Index |        |        |
|------------|-------|--------|--------|
| I          | 1     | 0.020  | 0.006  |
| am         | 2     | -0.003 | 0.010  |
| happy      | 3     | 0.009  | 0.010  |
| because    | 4     | -0.011 | -0.018 |
| learning   | 5     | -0.040 | -0.047 |
| NLP        | 6     | -0.009 | 0.050  |
| sad        | 7     | -0.044 | 0.001  |
| not        | 8     | 0.011  | -0.022 |

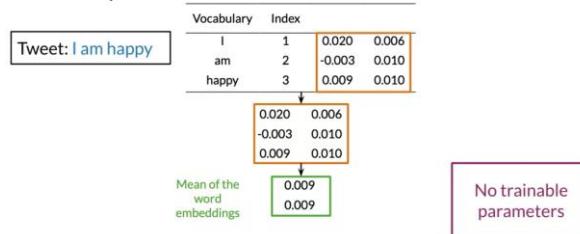
Trainable weights

Vocabulary x Embedding

- **Purpose:** Maps words to a high-dimensional space where words with similar meaning have similar representations.
- **Vocabulary:** Set of unique words from your text data.

- **Functionality:** Takes an index for each word and maps it to a fixed-size dense vector.
- **Dimensionality:** The size of the embedding (e.g., 2) is a hyperparameter.
- **Trainable Parameters:** The embedding weights are learned during training to optimize performance on the NLP task.
- **Benefits:** Allows the model to learn and improve word representations for the specific NLP task.

### Mean Layer



- **Purpose:** Reduces the dimensionality of the output from the embedding layer.
- **Functionality:** Computes the mean of the features across the word embeddings.
- **Parameters:** No trainable parameters; it simply averages the embeddings.
- **Output:** Returns a fixed-size vector regardless of the sequence length, preserving the number of features as the embedding size.
- **Use Case:** Useful for handling variable-length text inputs with a fixed-size layer downstream.

### Programming lab: Sentiment with Deep neural Networks

Import library & data:

```

1. import os
2. os.environ['TF_CPP_MIN_LOG_LEVEL'] = '3'
3. import numpy as np
4. import tensorflow as tf
5. import matplotlib.pyplot as plt
6. from sklearn.decomposition import PCA
7. from utils import load_tweets, process_tweet
8. %matplotlib inline
9.
10. all_positive_tweets, all_negative_tweets = load_tweets()

```

Process the tweets:

- tokenizing the sentence (splitting to words)
- removing stock market tickers like \$GE
- removing old style retweet text "RT"
- removing hyperlinks
- removing hashtags
- lowercasing
- removing stopwords and punctuation
- stemming

```

1. all_positive_tweets_processed = [process_tweet(tweet) for tweet in
all_positive_tweets]
2. all_negative_tweets_processed = [process_tweet(tweet) for tweet in
all_negative_tweets]
```

Split the dataset:

```

1. val_pos = all_positive_tweets_processed[4000:]
2. train_pos = all_positive_tweets_processed[:4000]
3. val_neg = all_negative_tweets_processed[4000:]
4. train_neg = all_negative_tweets_processed[:4000]
5. train_x = train_pos + train_neg
6. val_x = val_pos + val_neg
7. train_y = [[1] for _ in train_pos] + [[0] for _ in train_neg]
8. val_y = [[1] for _ in val_pos] + [[0] for _ in val_neg]
```

Build the vocabulary:

```

1. def build_vocabulary(corpus):
2. vocab = {'': 0, '[UNK]': 1}
3. index = 2
4. # For each tweet in the training set
5. for tweet in corpus:
6. # For each word in the tweet
7. for word in tweet:
8. # If the word is not in vocabulary yet, add it to vocabulary
9. if word not in vocab:
10. vocab[word] = index
11. index += 1
12. return vocab
13. vocab = build_vocabulary(train_x)
14. num_words = len(vocab)
```

Find the length of the longest tweet

```

1. def max_length(training_x, validation_x):
2. max_len = max(len(tweet_tokens) for tweet_tokens in training_x + validation_x)
3. return max_len
```

Pad the sequence → max\_len

```

1. def padded_sequence(tweet, vocab_dict, max_len, unk_token='[UNK]'):
2. unk_ID = vocab_dict[unk_token]
3. encoded_tweet = [vocab_dict.get(word, unk_ID) for word in tweet]
4. padded_tensor = encoded_tweet + [0] * (max_len - len(encoded_tweet))
5. return padded_tensor
6.
7. train_x_padded = [padded_sequence(x, vocab, max_len) for x in train_x]
8. val_x_padded = [padded_sequence(x, vocab, max_len) for x in val_x]
```

ReLU

```

1. def relu(x):
2. activation = np.maximum(0,x)
3. return activation
```

## Sigmoid

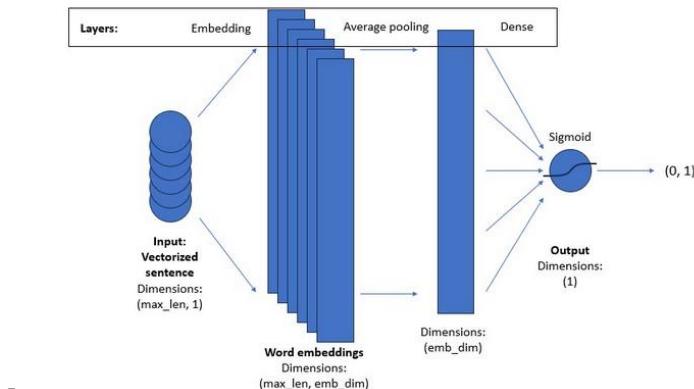
```
1. def sigmoid(x):
2. activation = 1/(1+np.exp(-x))
3. return activation
```

## Dense

```
1. class Dense():
2. def __init__(self, n_units, input_shape, activation, stdev=0.1, random_seed=42):
3. self.n_units = n_units
4. self.random_generator = np.random.default_rng(seed=random_seed)
5. self.activation = activation
6. w = self.random_generator.normal(scale=stdev, size = (input_shape[-1], self.n_units))
7. self.weights = w
8. def __call__(self, x):
9. return self.forward(x)
10. def forward(self, x):
11. dense = np.dot(x, self.weights)
12. dense = self.activation(dense)
13. return dense
```

## Model:

- `tf.keras.layers.Embedding`
  - o Turns positive integers (word indices) into vectors of fixed size. You can imagine it as creating one-hot vectors out of indices and then running them through a fully-connected (dense) layer.
- `tf.keras.layers.GlobalAveragePooling1D`
- `tf.keras.layers.Dense`
  - o Regular fully connected layer



## Model Implementation:

```
1. def create_model(num_words, embedding_dim, max_len):
2. tf.random.set_seed(123)
3. model = tf.keras.Sequential([
4. tf.keras.layers.Embedding(num_words, embedding_dim, input_length=max_len),
5. tf.keras.layers.GlobalAveragePooling1D(),
6. tf.keras.layers.Dense(1, activation='sigmoid')
7.])
8. model.compile(loss='binary_crossentropy',
9. optimizer='adam',
10. metrics=['accuracy'])
11. return model
12.
13. model = create_model(num_words=num_words, embedding_dim=16, max_len=max_len)
```

Training model:

```
1. train_x_prepared = np.array(train_x_padded)
2. val_x_prepared = np.array(val_x_padded)
3. train_y_prepared = np.array(train_y)
4. val_y_prepared = np.array(val_y)
5. history = model.fit(train_x_prepared, train_y_prepared, epochs=20,
validation_data=(val_x_prepared, val_y_prepared))
```

Evaluating model:

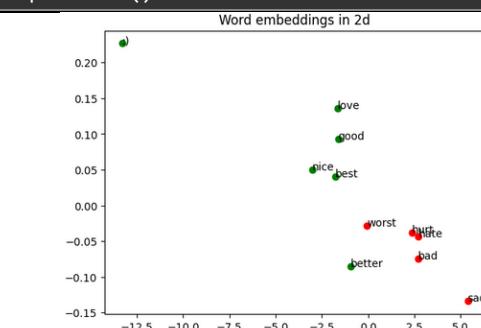
```
1. def plot_metrics(history, metric):
2. plt.plot(history.history[metric])
3. plt.plot(history.history[f'val_{metric}'])
4. plt.xlabel("Epochs")
5. plt.ylabel(metric.title())
6. plt.legend([metric, f'val_{metric}'])
7. plt.show()
8.
9. plot_metrics(history, "accuracy")
10. plot_metrics(history, "loss")
```

Predict the data:

```
1. example_for_prediction = np.append(val_x_prepared[0:10], val_x_prepared[-10:], axis=0)
2. model.predict(example_for_prediction)
3.
4. def get_prediction_from_tweet(tweet, model, vocab, max_len):
5. tweet = process_tweet(tweet)
6. tweet = padded_sequence(tweet, vocab, max_len)
7. tweet = np.array([tweet])
8. prediction = model.predict(tweet, verbose=False)
9. return prediction[0][0]
10.
11. unseen_tweet = '@DLAI @NLP_team_dlai OMG!!! what a daaay, wow, wow. This AsSiGnMeNT was gr8.'
12. prediction_unseen = get_prediction_from_tweet(unseen_tweet, model, vocab, max_len)
```

Word Embeddings

```
1. embeddings_layer = model.layers[0]
2. embeddings = embeddings_layer.get_weights()[0]
3. pca = PCA(n_components=2)
4. embeddings_2D = pca.fit_transform(embeddings)
5. neg_words = ['bad', 'hurt', 'sad', 'hate', 'worst']
6. pos_words = ['best', 'good', 'nice', 'love', 'better', ':)']
7. neg_n = [vocab[w] for w in neg_words]
8. pos_n = [vocab[w] for w in pos_words]
9. plt.figure()
10. plt.scatter(embeddings_2D[neg_n][:,0], embeddings_2D[neg_n][:,1], color = 'r')
11. for i, txt in enumerate(neg_words):
12. plt.annotate(txt, (embeddings_2D[neg_n][i,0], embeddings_2D[neg_n][i,1]))
13. plt.scatter(embeddings_2D[pos_n][:,0], embeddings_2D[pos_n][:,1], color = 'g')
14. for i, txt in enumerate(pos_words):
15. plt.annotate(txt, (embeddings_2D[pos_n][i,0], embeddings_2D[pos_n][i,1]))
16. plt.title('Word embeddings in 2d')
17. plt.show()
```

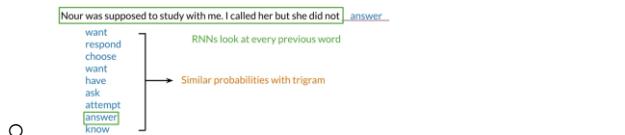


## Traditional Language model:

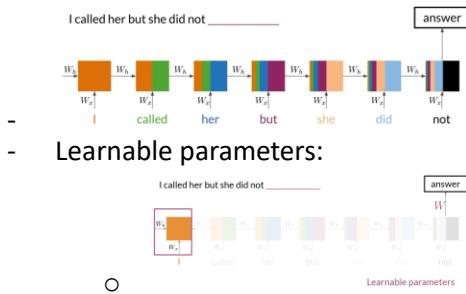
- N-Grams
  - o Large N-grams needed to capture dependencies between distant words
  - o Need a lot of space and RAM
  - o Example: Bigrams:
    - $P(w_1, w_2, w_3) = P(w_1) * P(w_2 | w_1) * P(w_3 | w_2)$

## Recurrent Neural Networks vs N-Grams:

- Example:



## RNNs Basic Structure:

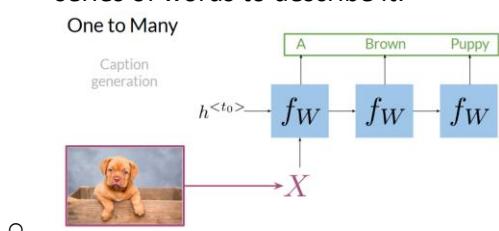


## Applications of RNNs

- One to One:
  - o A simple model that takes a single input to produce a single output.
  - o Example: Predicting the position of a soccer team on the leaderboard based on a list of scores from La Liga Santander.
  - o For such tasks, the advantage of RNNs over traditional neural networks is minimal.

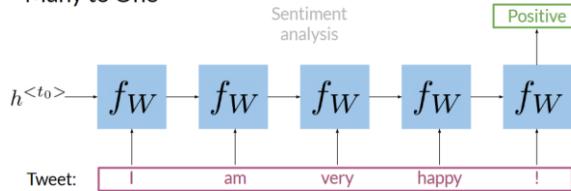


- One to Many:
  - o This architecture takes a single input to generate multiple outputs.
  - o Example Use: Image captioning where an RNN takes an image and outputs a series of words to describe it.



- Many to One:
  - o Multiple inputs leading to a single output.
  - o Example: Sentiment analysis where an RNN analyzes a sequence of words (like a tweet) to determine the overall sentiment.

Many to One

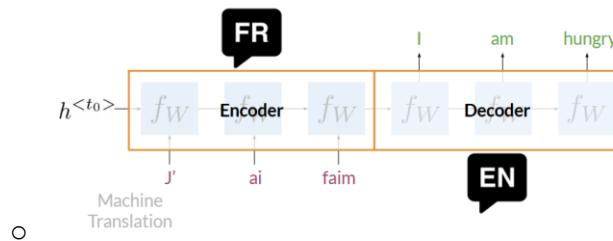


- o

- Many to Many:

- o Involves multiple inputs and multiple outputs.
- o Example: Machine translation where an RNN takes a series of words in one language and translates them into another.
- o Specific Architecture: The encoder-decoder model is commonly used for this task. The encoder processes the input sequence into a single representation capturing the sentence's meaning, which the decoder then translates into another language

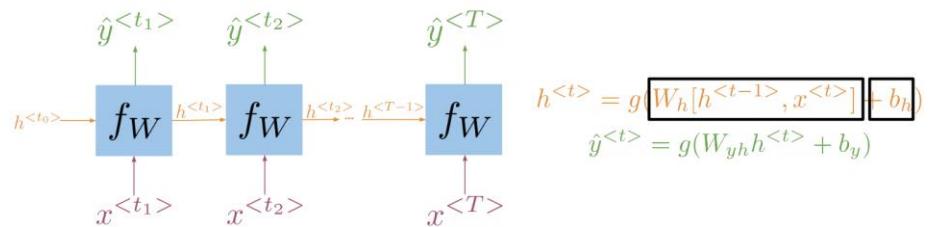
Many to Many



- o

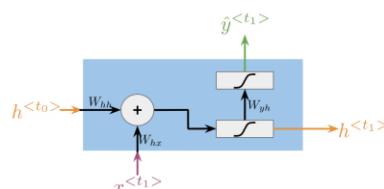
## Math in Simple RNNs

- A Vanilla RNN:



- o

- o Look into the first cell:



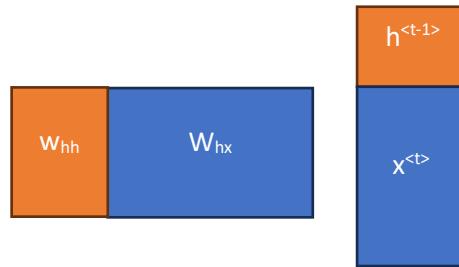
- o Shapes:

- If  $h^{<t>} \rightarrow 4 \times 1$ ,  $x^{<t>} = 10 \times 1$ 
  - $W_h$  will be  $4 \times 14$

$$W_h = [W_{hh} \mid W_{hx}]$$

Formula 1:  $h^{<t>} = g(W_h[h^{<t-1>}, x^{<t>}] + b_h)$

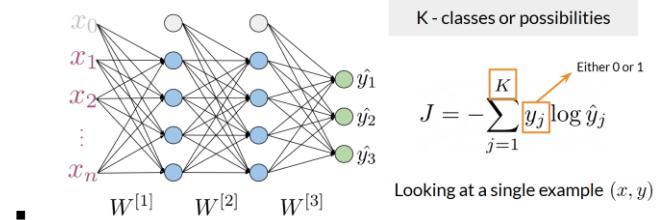
Formula 2:  $h^{<t>} = g(W_{hh}h^{<t-1>} + W_{hx}x^{<t>} + b_h)$



Formula 1  $\Leftrightarrow$  Formula 2:

Cross function for RNNs

- Cross Entropy Loss
  - o For K – classes or possibilities:

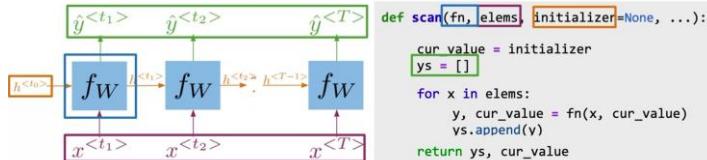


- o For Vanilla RNN:
  - Average with respect to time
    - there is just one value in every vector  $y^{<t>}$  different from 0.

$$J = -\frac{1}{T} \sum_{t=1}^T \sum_{j=1}^K y_j^{<t>} \log \hat{y}_j^{<t>}$$

Scan function in Tensorflow:

`tf.scan()` function



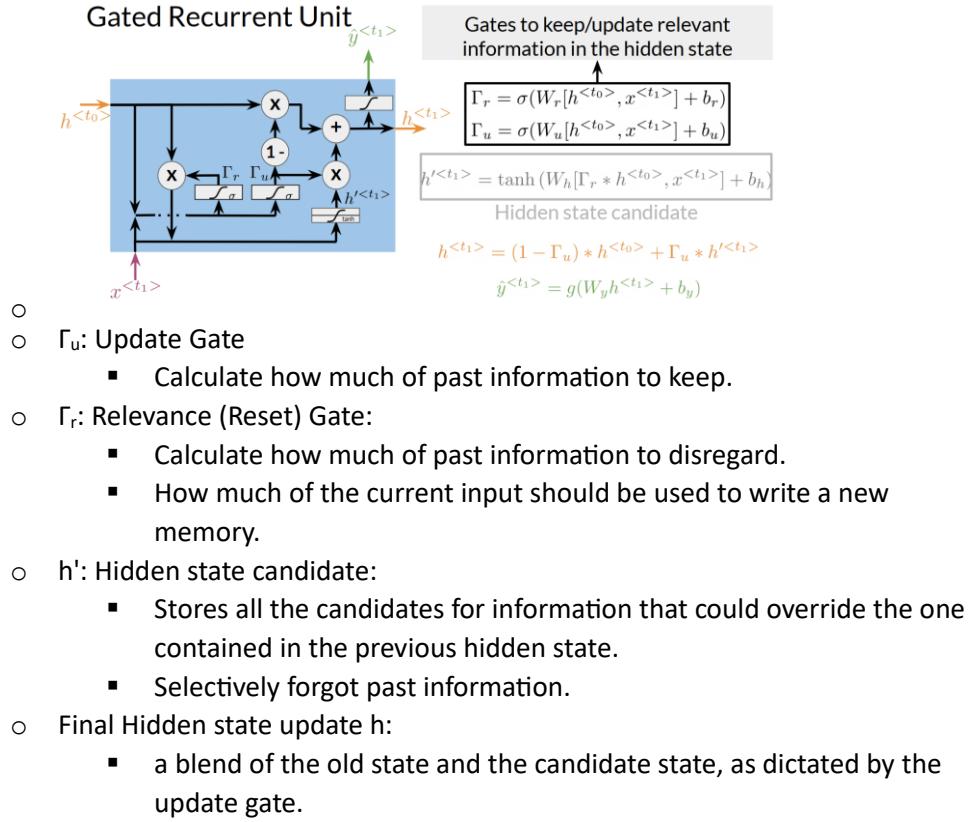
- `cur_value` corresponds to the hidden state in an RNN
- Framework like Tensorflow need this type of abstraction for parallel computations and GPU usage.

## Gated Recurrent Units

- Allows relevant information to be kept in the hidden state over long sequences.
- Example:

- "Ants are really interesting. They are everywhere."
- Plural
- Relevance and update gates to remember important prior information.

- Architecture:



## Vanilla RNN vs GRUs

| Vanilla RNN                                                                                                                                                                        | GRUs                                                                                                                                                                                                                                                                                                                                           |
|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <p><b>Vanilla RNN</b></p> $\begin{aligned}h^{<t_1>} &= g(W_h[h^{<t-1>}, x^{<t_1>}] + b_h) \\ \hat{y}^{<t_1>} &= g(W_y h^{<t_1>} + b_y)\end{aligned}$                               | <p><b>GRUs</b></p> $\begin{aligned}\Gamma_u &= \sigma(W_u[h^{<t_0>}, x^{<t_1>}] + b_u) \\ \Gamma_r &= \sigma(W_r[h^{<t_0>}, x^{<t_1>}] + b_r) \\ h'^{<t_1>} &= \tanh(W_h[\Gamma_r * h^{<t_0>}, x^{<t_1>}] + b_h) \\ h^{<t_1>} &= (1 - \Gamma_u) * h^{<t_0>} + \Gamma_u * h'^{<t_1>} \\ \hat{y}^{<t_1>} &= g(W_y h^{<t_1>} + b_y)\end{aligned}$ |
| <p><b>Disadvantages:</b></p> <ul style="list-style-type: none"> <li>- Updating the hidden state at every time step, for long sequence, the information tends to vanish.</li> </ul> | <p><b>Advantages:</b></p> <ul style="list-style-type: none"> <li>- Allow the network to learn what type of information to keep and when to override it.</li> </ul>                                                                                                                                                                             |

## Programming lab: Vanilla RNN and GRUs

Forward method for Vanillar RNN cell:

$$h^{<t>} = g(W_h[h^{<t-1>}, x^{<t>}] + b_h)$$

$$\hat{y}^{<t>} = g(W_{yh}h^{<t>} + b_y)$$

```
1. def forward_V_RNN(inputs, weights):
2. x, h_t = inputs
3. wh, bh = weights
4. h_t = np.dot(wh, np.concatenate([h_t, x])) + bh
5. h_t = sigmoid(h_t)
6. y = h_t
7. return y, h_t
```

Forward method for GRUs:

$$\Gamma_r = \sigma(W_r[h^{<t-1>}, x^{<t>}] + b_r)$$

$$\Gamma_u = \sigma(W_u[h^{<t-1>}, x^{<t>}] + b_u)$$

$$c^{<t>} = \tanh(W_h[\Gamma_r * h^{<t-1>}, x^{<t>}] + b_h)$$

$$h^{<t>} = \Gamma_u * c^{<t>} + (1 - \Gamma_u) * h^{<t-1>}$$

```
1. def forward_GRU(inputs, weights): # Forward propagation for a single GRU cell
2. x, h_t = inputs
3. # weights.
4. wu, wr, wc, bu, br, bc = weights
5. # Update gate
6. u = np.dot(wu, np.concatenate([h_t, x])) + bu
7. u = sigmoid(u)
8. # Relevance gate
9. r = np.dot(wr, np.concatenate([h_t, x])) + br
10. r = sigmoid(r)
11. # Candidate hidden state
12. c = np.dot(wc, np.concatenate([r * h_t, x])) + bc
13. c = np.tanh(c)
14. # New Hidden state h_t
15. h_t = u * c + (1 - u) * h_t
16. # We avoid implementation of y for clarity
17. y = h_t
18. return y, h_t
19.
20. forward_GRU([X[1], h_0], weights_GRU)[0]
```

Scan function:

```
1. def scan(fn, elems, weights, h_0): # Forward propagation for RNNs
2. h_t = h_0
3. ys = []
4. for x in elems:
5. y, h_t = fn([x, h_t], weights)
6. ys.append(y)
7. return ys, h_t
```

Go through all the elements from a sequence

```
1. # vanilla RNNs
2. ys, h_T = scan(forward_V_RNN, X, weights_vanilla, h_0)
3. # GRUs
4. ys, h_T = scan(forward_GRU, X, weights_GRU, h_0)
```

Create a GRU model in tensorflow:

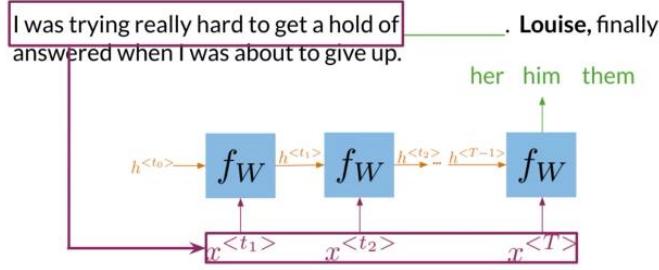
```
1. model_GRU = tf.keras.Sequential([
2. tf.keras.layers.GRU(256, return_sequences=True, name='GRU_1_returns_seq'),
3. tf.keras.layers.GRU(128, return_sequences=True, name='GRU_2_returns_seq'),
4. tf.keras.layers.GRU(64, name='GRU_3_returns_last_only'),
5. tf.keras.layers.Dense(10)
6.])
7. batch_size = 60
8. sequence_length = 50
9. word_vector_length = 40
10. input_data = tf.random.normal([batch_size, sequence_length, word_vector_length])
11. prediction = model_GRU(input_data)
12. model_GRU.summary()
```

Model: "sequential"

| Layer (type)                        | Output Shape  | Param # |
|-------------------------------------|---------------|---------|
| =====                               |               |         |
| GRU_1_returns_seq (GRU)             | (60, 50, 256) | 228864  |
| GRU_2_returns_seq (GRU)             | (60, 50, 128) | 148224  |
| GRU_3_returns_last_only (GRU)       | (60, 64)      | 37248   |
| dense (Dense)                       | (60, 10)      | 650     |
| =====                               |               |         |
| Total params: 414986 (1.58 MB)      |               |         |
| Trainable params: 414986 (1.58 MB)  |               |         |
| Non-trainable params: 0 (0.00 Byte) |               |         |

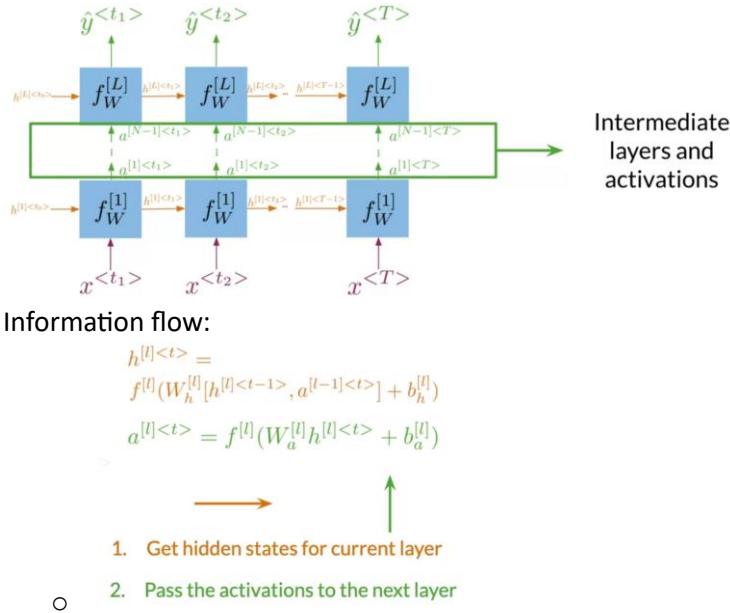
- 
- the model's first GRU layer takes the original 40-dimensional word vectors and transforms them into its own 256-dimensional representations.

Limitation of simple RNNs:



Deep RNNs:

- Just RNNs stack together



- Information flow:

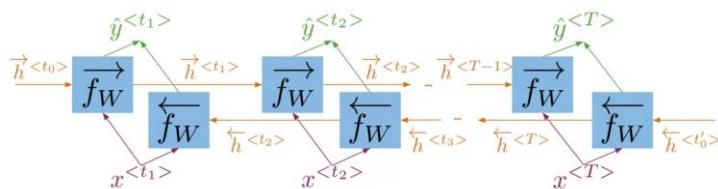
$$\begin{aligned} h^{[l]} < t > &= h^{[l]}(W_h^{[l]}[h^{[l-1]} < t-1 >, a^{[l-1]} < t >] + b_h^{[l]}) \\ a^{[l]} < t > &= a^{[l]}(W_a^{[l]}h^{[l]} < t > + b_a^{[l]}) \end{aligned}$$

→

1. Get hidden states for current layer

○ 2. Pass the activations to the next layer

Bi-directional RNNs



- information flows from the past and from the future independently
- Prediction:

$$\hat{y}^{<t>} = g(\vec{h}^{<t>}, \overleftarrow{h}^{<t>}) + b_y$$

Perplexity:

$$P(W) = \sqrt[N]{\prod_{i=1}^N \frac{1}{P(w_i | w_1, \dots, w_{i-1})}}$$

- Taking the logarithm:

$$= -\frac{1}{N} \sum_{i=1}^N \log P(w_i | w_1, \dots, w_{i-1})$$

○

## Assignment: Deep N-grams

Load library

```
1. import os
2. import traceback
3. os.environ['TF_CPP_MIN_LOG_LEVEL'] = '3'
4. import shutil
5. import numpy as np
6. import random as rnd
7. import tensorflow as tf
8. from tensorflow.keras import layers
9. from tensorflow.keras.layers import Input
10. from termcolor import colored
11. rnd.seed(32)
```

Load the Data:

```
1. dirname = 'data/'
2. filename = 'shakespeare_data.txt'
3. lines = []
4. counter = 0
5. with open(os.path.join(dirname, filename)) as files:
6. for line in files:
7. pure_line = line.strip().lower()
8. if pure_line:
9. lines.append(pure_line)
```

Create the vocabulary:

```
1. text = "\n".join(lines)
2. vocab = sorted(set(text))
3. vocab.insert(0, "[UNK]")
4. vocab.insert(1, "")
```

Convert a line to Tensor:

```
1. line = "Hello world!"
2. chars = tf.strings.unicode_split(line, input_encoding='UTF-8')
```

chars:

- `tf.Tensor([b'H' b'e' b'l' b'l' b'o' b' ' b'w'... b'!'], shape=(12,), dtype=string)`

StringLookup:

```
1. ids = tf.keras.layers.StringLookup(vocabulary=list(vocab), mask_token=None)(chars)
```

Ids:

- `tf.Tensor([34 59 66 66 69 4 77 69 72 66 58 5], shape=(12,), dtype=int64)`

Line to Tensor:

Takes in a single line and transforms each character into its unicode integer

```
1. def line_to_tensor(line, vocab):
2. chars = tf.strings.unicode_split(line, input_encoding='UTF-8')
3. ids = tf.keras.layers.StringLookup(vocabulary=list(vocab), mask_token=None)(chars)
4. return ids
```

Ids to Text:

Converting back from ids to text

```
1. def text_from_ids(ids, vocab):
2. chars_from_ids = tf.keras.layers.StringLookup(vocabulary=vocab, invert=True, mask_token=None)
3. return tf.strings.reduce_join(chars_from_ids(ids), axis=-1)
```

Predict the next character: Create the input and the output for the model:

```
1. def split_input_target(sequence):
2. input_text = sequence[:-1]
3. target_text = sequence[1:]
4. return input_text, target_text
```

split\_input\_target(list("Tensorflow")):

```
- ([T', 'e', 'n', 's', 'o', 'r', 'f', 'l', 'o'], [e', 'n', 's', 'o', 'r', 'f', 'l', 'o', 'w'])
```

Create batch dataset

```
1. def create_batch_dataset(lines, vocab, seq_length=100, batch_size=64):
2. BUFFER_SIZE = 10000
3. single_line_data = "\n".join(lines)
4. all_ids = line_to_tensor(single_line_data, vocab)
5. ids_dataset = tf.data.Dataset.from_tensor_slices(all_ids)
6. data_generator = ids_dataset.batch(seq_length + 1, drop_remainder=True)
7. dataset_xy = data_generator.map(split_input_target)
8. dataset = (
9. dataset_xy
10. .shuffle(BUFFER_SIZE)
11. .batch(batch_size, drop_remainder=True)
12. .prefetch(tf.data.experimental.AUTOTUNE)
13.)
14. return dataset
```

GRU Language Model (GRULM)

```
1. class GRULM(tf.keras.Model):
2. def __init__(self, vocab_size=256, embedding_dim=256, rnn_units=128):
3. super().__init__()
4. self.embedding = tf.keras.layers.Embedding(vocab_size, embedding_dim)
5. self.gru = tf.keras.layers.GRU(
6. units=rnn_units,
7. return_sequences=True,
8. return_state=True,
9.)
10. self.dense = tf.keras.layers.Dense(vocab_size, activation=tf.nn.log_softmax)
11. def call(self, inputs, states=None, return_state=False, training=False):
12. x = inputs
13. x = self.embedding(x, training=training)
14. if states is None:
15. states = self.gru.get_initial_state(x)
16. x, states = self.gru(x, initial_state=states, training=training)
17. x = self.dense(x, training=training)
18. if return_state:
19. return x, states
20. else:
21. return x
```

Using the model to predict the next character using the untrained model:

```
1. vocab_size = 82
2. embedding_dim = 256
3. rnn_units = 512
4. model = GRULM(
5. vocab_size=vocab_size,
6. embedding_dim=embedding_dim,
7. rnn_units = rnn_units)

1. for input_example_batch, target_example_batch in dataset.take(1):
2. example_batch_predictions = model(tf.constant([input_example_batch[0].numpy()]))
```

example\_batch\_prediction.shape:

- (1,100,82): the first sequence generated by the batch generator, predict 100 characters, 82 values for each predicted character.

Choose the next character by getting the index of the element with the highest likelihood:

```
1. last_character = tf.math.argmax(example_batch_predictions[0][99])
2. sampled_indices = tf.math.argmax(example_batch_predictions[0], axis=1)
3.
4. print("Input:\n", text_from_ids(input_example_batch[0], vocab))
5. print()
6. print("Next Char Predictions:\n", text_from_ids(sampled_indices, vocab))
```

Compile Model:

```
1. def compile_model(model):
2. loss = tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True)
3. opt = tf.keras.optimizers.Adam(learning_rate=0.00125)
4. model.compile(optimizer=opt, loss=loss)
5. return model
```

Training:

```
1. EPOCHS = 10
2. model = compile_model(model)
3. history = model.fit(dataset, epochs=EPOCHS)
```

Log\_perplexity:

```
1. def log_perplexity(preds, target):
2. PADDING_ID = 1
3. log_p = np.sum(preds * tf.one_hot(target, preds.shape[-1]), axis=-1)
4. non_pad = 1.0 - np.equal(target, PADDING_ID)
5. log_p = log_p * non_pad
6. log_ppx = np.sum(log_p, -1) / np.sum(non_pad, -1)
7. log_ppx = np.mean(log_ppx)
8. return -log_ppx
```

Loading the pretrained model:

```
1. vocab_size = len(vocab)
2. embedding_dim = 256
3. rnn_units = 512
4. model = GRULM(
5. vocab_size=vocab_size,
6. embedding_dim=embedding_dim,
7. rnn_units = rnn_units)
8. model.build(input_shape=(100, vocab_size))
9. model.load_weights('./model/')
```

Calculate the perplexity:

```
1. eval_text = "\n".join(eval_lines)
2. eval_ids = line_to_tensor([eval_text], vocab)
3. input_ids, target_ids = split_input_target(tf.squeeze(eval_ids, axis=0))
4. preds, status = model(tf.expand_dims(input_ids, 0), training=False, states=None,
 return_state=True)
5. log_ppx = log_perplexity(preds, tf.expand_dims(target_ids, 0))
6. print(f'The log perplexity and perplexity of your model are {log_ppx} and
 {np.exp(log_ppx)} respectively')
```

Output:

- The log perplexity and perplexity of your model are 1.2239635591264044 and 3.400639693684683 respectively

This model in its default form is deterministic and can result in repetitive and monotonous outputs. To make the Language model more dynamic and versatile:

- Introduce an element of randomness into its predictions.
- Employ random sampling, regulated by temperature.

```
1. def temperature_random_sampling(log_probs, temperature=1.0):
2. u = tf.random.uniform(minval=1e-6, maxval=1.0 - 1e-6, shape=log_probs.shape)
3. g = -tf.math.log(-tf.math.log(u))
4. return tf.math.argmax(log_probs + g * temperature, axis=-1)
```

GenerativeModel:

```
1. class GenerativeModel(tf.keras.Model):
2. def __init__(self, model, vocab, temperature=1.0):
3. super().__init__()
4. self.temperature = temperature
5. self.model = model
6. self.vocab = vocab
7. @tf.function
8. def generate_one_step(self, inputs, states=None):
9. input_ids = line_to_tensor(inputs, self.vocab)
10. predicted_logits, states = self.model(inputs=input_ids, states=states, return_state=True)
11. predicted_logits = predicted_logits[:, -1, :]
12. predicted_ids = temperature_random_sampling(predicted_logits, self.temperature)
13. predicted_chars = text_from_ids([predicted_ids], self.vocab)
14. return tf.expand_dims(predicted_chars, 0), states
15. def generate_n_chars(self, num_chars, prefix):
16. states = None
17. next_char = tf.constant([prefix])
18. result = [next_char]
19. for n in range(num_chars):
20. next_char, states = self.generate_one_step(next_char, states=states)
21. result.append(next_char)
22. return tf.strings.join(result)[0].numpy().decode('utf-8')
```

Generate a longer text:

```
1. tf.random.set_seed(np.random.randint(1, 1000))
2. gen = GenerativeModel(model, vocab, temperature=0.8)
3. import time
4. start = time.time()
5. print(gen.generate_n_chars(1000, "ROMEO "), '\n\n' + '_'*80)
6. print('\nRun time:', time.time() - start)
```

## RNNs and Vanishing Gradients

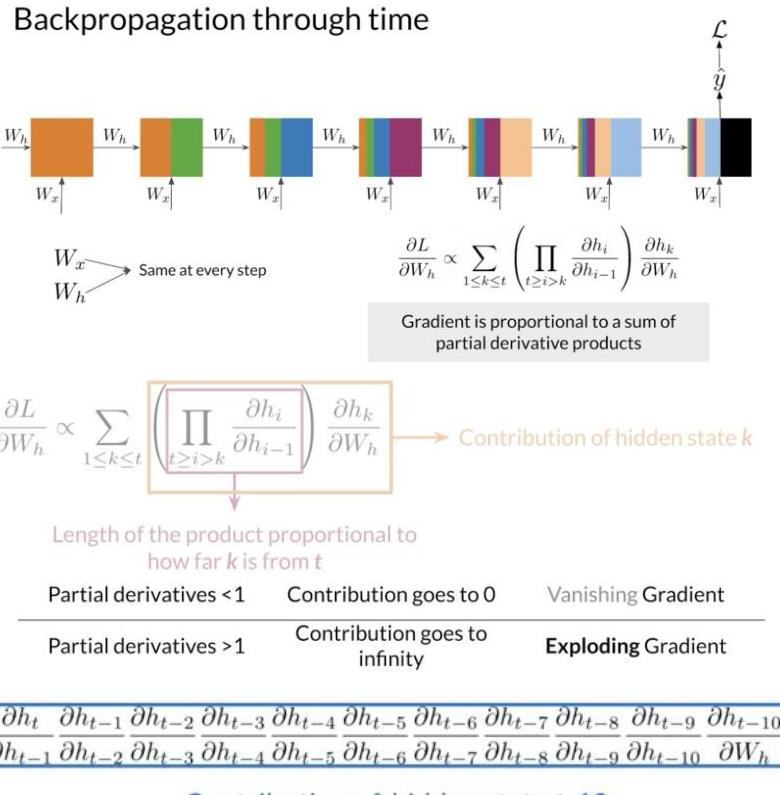
RNNs: Advantages:

- Captures dependencies within a short range
- Takes up less RAM than other n-gram models

RNNs: Disadvantages:

- Struggles to capture long term dependencies
- Prone to vanishing or exploding gradients

RNNs: Back propagation through time:



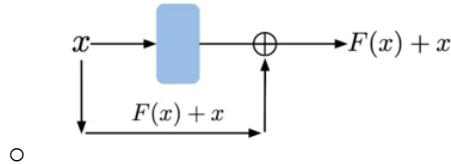
Solving for vanishing gradients problem:

- Identity RNN:
  - Identity RNN with ReLU activation
  - Initializing the weights to the identity matrix, using ReLU activation.
  - Interpretation:
    - Copy the previous hidden states and information from the current inputs and replace any negative values with zero.
  - Effect:
- Encourage the network to stay close to the values and the identity matrix.

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad -1 \longrightarrow 0$$

Solving for exploding gradients problem:

- Gradient clipping
- Skip connections:



Detailed Explanation for Identity RNN:

---

## Identity Matrix Initialization

In an RNN, the hidden state at time  $t$  is typically computed as:

$$h_t = \text{activation}(W * h_{t-1} + U * x_t + b)$$

where:

- $h_t$  is the hidden state at time  $t$ .
- $h_{t-1}$  is the hidden state at time  $t-1$ .
- $x_t$  is the input at time  $t$ .
- $W$  and  $U$  are the weight matrices for the hidden state and input respectively.
- $b$  is a bias vector.
- $\text{activation}$  is the activation function applied element-wise.

When initializing  $W$  to the identity matrix (and assuming no bias for simplicity), the recurrent part of the computation simplifies to:

$$h_t = \text{activation}(h_{t-1} + U * x_t)$$

If  $W$  is the identity matrix and  $b$  is zero,  $W * h_{t-1}$  results in  $h_{t-1}$  because multiplying by the identity matrix leaves the vector unchanged. This means that, in the absence of new inputs,  $h_t$  would be the same as  $h_{t-1}$ .

## ReLU Activation

ReLU, or Rectified Linear Unit, is an activation function defined as:

$$\text{ReLU}(x) = \max(0, x)$$

The derivative of ReLU is as follows:

- 0 for  $x < 0$
- 1 for  $x > 0$
- Undefined at  $x = 0$ , but typically handled as either 0 or 1 in practice.

Because the derivative of ReLU is 1 for positive inputs, when you backpropagate through a ReLU neuron that has a positive activation, the gradient does not diminish. This is particularly useful in combating the vanishing gradient problem because it means that as long as the neuron activates, it will not diminish the gradient during backpropagation.

## Mathematical Effect in Identity RNN

By combining the identity matrix with ReLU, the RNN has the following properties:

- Preservation of Gradient Magnitude:** During backpropagation, the gradient with respect to  $h_{(t-1)}$  is not scaled down by  $w$  because  $w$  is an identity matrix. This helps prevent the decay of gradients as they are propagated back through time.
  - Non-Negativity of Activations:** ReLU ensures that all negative values are set to zero and positive values pass through unchanged. This means that negative data does not contribute to the gradient, which may help in maintaining the stability of the gradient's magnitude over time.
  - Simplicity and Efficiency:** Since the identity matrix and ReLU are simple and involve minimal computation, this approach can be computationally efficient.
- 

Lab: Vanishing gradients and exploding gradients in RNNs

- For Vanilla RNNs, the gradient with respect to  $W_h$  is proportional to:

$$\circ \quad \frac{\delta L}{\delta W_h} \propto \sum_{1 \leq k \leq t} \left( \prod_{i \geq i > k} \frac{\delta h_i}{\delta h_{i-1}} \right) \frac{\delta h_k}{\delta W_h}$$

Activations & Partial Derivative:

- Hidden state:

$$\circ \quad h_i = \sigma(W_{hh}h_{i-1} + W_{hx}x_i + b_h)$$

Backpropagation through time (BPTT):

- We want to compute the partial derivative of  $h_i$  with respect to  $h_{i-1}$ .

By chain rule:

- Compute the derivative of the outer function:
  - o  $\sigma'(z) = \sigma(z)(1-\sigma(z))$
- Compute the derivative of the inner function:
  - o The inner function is the affine transformation:
    - $W_{hh}h_{i-1} + W_{hx}x_i + b_h$  with respect to  $h_{i-1}$
    - The derivative is simply  $W_{hh}$
- We get:
  - o  $\frac{\delta h_i}{\delta h_{i-1}} = \text{diag}(\sigma'(z))W_{hh}$ 
    - The diag function used to ensure that the element-wise derivative of the activation function is applied correctly during matrix multiplication.

Partial derivative:

$$\frac{\delta h_i}{\delta h_{i-1}} = W_{hh}^T \text{diag}(\sigma'(W_{hh}h_{i-1} + W_{hx}x_i + b_h))$$

Gradients:

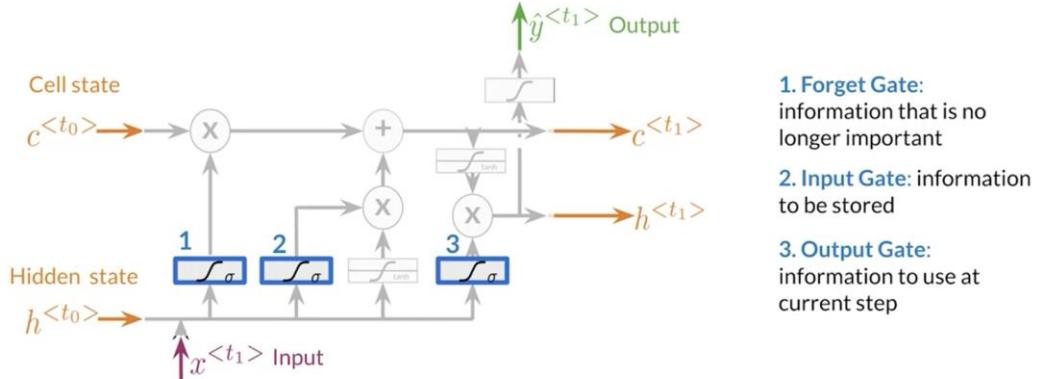
$$\prod_{i \geq i > k} \frac{\delta h_i}{\delta h_{i-1}} = \prod_{i \geq i > k} W_{hh}^T \text{diag}(\sigma'(W_{hh}h_{i-1} + W_{hx}x_i + b_h))$$

- Approaches 0 → vanishing gradient problems, vice versa.

## Introduction to LSTMs

- Learns when to remember and when to forget.

## LSTM Overview:



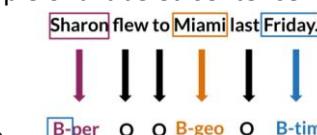
- Forget gate:
  - o Decides what to keep.
  - o  $f = \sigma(W_f [h_{t-1}; x_t] + b_f)$
- Input gate:
  - o Decides what to add.
  - o  $i = \sigma(W_i [h_{t-1}; x_t] + b_i)$
  - o Candidate values vector:  $\tilde{C}_t = \tanh(W_C \bullet [h_{t-1}, x_t] + b_C)$
  - o Cell state update:  $C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$ 
    - Forgetting the things decided by the forget gate & adding the new candidate values scaled by how much we decided to update each state value.
- Output gate:
  - o Decides what the next hidden state will be.
  - o  $o = \sigma(W_o [h_{t-1}; x_t] + b_o)$
- Tanh layer:
  - o Ensures the values in the network stay numerically stable by squeezing all values between -1 and 1.
- The output of LSTM unit:
  - o  $h_t = o_t \odot \tanh(c_t)$

## Application of LSTMs

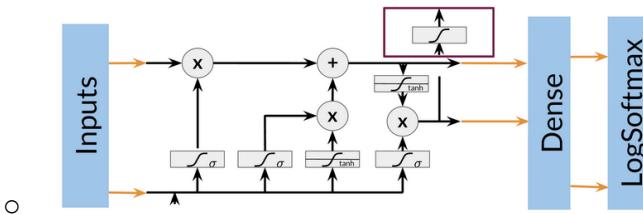
- Next-character prediction
- Image captioning
- Chatbots
- Music composition
- Speech recognition

## Name Entity Recognition

- Locates and extracts predefined entities from text.
- Can be places, organizations, names, time and dates.
- Types of Entities:
  - o Geographical (e.g. Thailand)
  - o Organization (e.g. Google)
  - o Geopolitical (e.g. Indian)
  - o Time indicator (e.g. December)
  - o Artifact (e.g. Egyptian statue)
  - o Person (e.g. Barack Obama)
- Example of a labeled sentence:
 

Sharon flew to Miami last Friday.  

  
 ○ B-per ○ O ○ B-geo ○ O B-tim
- Application:
  - o Search engine efficiency
  - o Recommendation engines
  - o Customer service
  - o Automatic trading

## Name Entity Recognition(NER) System:

- Training:
  - o Create a tensor for each input and its corresponding number
  - o Put them in a batch ==> 64, 128, 256, 512 ...
  - o Feed it into an LSTM unit
  - o Run the output through a dense layer
  - o Predict using a log softmax over K classes
- Architecture:
 

## Computing Accuracy:

- Pass test set through the model
- Get arg max across the prediction array
- Mask padded tokens making the array the same length
  - o Padded tokens are not part of the data and are just used to help us keep the same sequence length for more efficient batch processing. We should not include their loss.
- Compare outputs against test labels

## Name Entity Recognition Assignment:

### Library:

```
1. import os
2. os.environ['TF_CPP_MIN_LOG_LEVEL'] = '3'
3. import numpy as np
4. import pandas as pd
5. import tensorflow as tf
6. tf.keras.utils.set_random_seed(33)
```

### Data demo:

```
1. # display original kaggle data
2. data = pd.read_csv("data/ner_dataset.csv", encoding = "ISO-8859-1")
3. train_sents = open('data/small/train/sentences.txt', 'r').readline()
4. train_labels = open('data/small/train/labels.txt', 'r').readline()
5. del(data, train_sents, train_labels)
```

SENTENCE: Thousands of demonstrators have marched through London to protest the war in Iraq and demand the withdrawal of British troops from that country .

SENTENCE LABEL: 0 0 0 0 0 B-geo 0 0 0 0 0 B-geo 0 0 0 0 0 B-gpe 0 0 0 0 0

#### ORIGINAL DATA:

|   | Sentence #  | Word          | POS | Tag |
|---|-------------|---------------|-----|-----|
| 0 | Sentence: 1 | Thousands     | NNS | 0   |
| 1 | NaN         | of            | IN  | 0   |
| 2 | NaN         | demonstrators | NNS | 0   |
| 3 | NaN         | have          | VBP | 0   |
| 4 | NaN         | marched       | VBN | 0   |

### Importing the data:

```
1. def load_data(file_path):
2. with open(file_path, 'r') as file:
3. data = np.array([line.strip() for line in file.readlines()])
4. return data
5.
6. train_sentences = load_data('data/large/train/sentences.txt')
7. train_labels = load_data('data/large/train/labels.txt')
8. val_sentences = load_data('data/large/val/sentences.txt')
9. val_labels = load_data('data/large/val/labels.txt')
10. test_sentences = load_data('data/large/test/sentences.txt')
11. test_labels = load_data('data/large/test/labels.txt')
```

### Encoding the sentences:

```
1. def get_sentence_vectorizer(sentences):
2. tf.keras.utils.set_random_seed(33) ## Do not change this line.
3. sentence_vectorizer = tf.keras.layers.TextVectorization(stdize=None)
4. sentence_vectorizer.adapt(sentences)
5. vocab = sentence_vectorizer.get_vocabulary()
6. return sentence_vectorizer, vocab
7.
8. sentence_vectorizer, vocab = get_sentence_vectorizer(train_sentences)
```

### Encoding the labels

```
1. def get_tags(labels):
2. tag_set = set() # Define an empty set
3. for el in labels:
4. for tag in el.split(" "):
5. tag_set.add(tag)
6. tag_list = list(tag_set)
7. tag_list.sort()
8. return tag_list
9. tags = get_tags(train_labels)
```

tags:

- ['B-art', 'B-eve', 'B-geo', 'B-gpe', 'B-nat', 'B-org', 'B-per', 'B-tim', ..., 'I-tim', 'O']

Mapping the tags to positive integers:

```
1. def make_tag_map(tags):
2. tag_map = {}
3. for i,tag in enumerate(tags):
4. tag_map[tag] = i
5. return tag_map
6. tag_map = make_tag_map(tags)
```

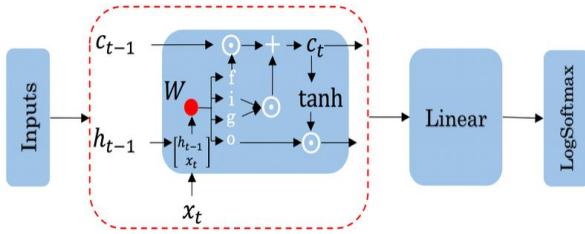
Label vectorizer: inputs a list of labels and a tag mapping and outputs their respective label ids via a tag map lookup.

```
1. def label_vectorizer(labels, tag_map):
2. label_ids = []
3. for element in labels:
4. tokens = element.split(' ')
5. element_ids = [tag_map[token] for token in tokens]
6. label_ids.append(element_ids)
7. label_ids = tf.keras.utils.pad_sequences(label_ids, padding='post', value=-1)
8. return label_ids
```

Building the dataset:

```
1. def generate_dataset(sentences, labels, sentence_vectorizer, tag_map):
2. sentences_ids = sentence_vectorizer(sentences)
3. labels_ids = label_vectorizer(labels, tag_map=tag_map)
4. dataset = tf.data.Dataset.from_tensor_slices((sentences_ids, labels_ids))
5. return dataset
6.
7. train_dataset = generate_dataset(train_sentences, train_labels, sentence_vectorizer, tag_map)
8. val_dataset = generate_dataset(val_sentences, val_labels, sentence_vectorizer, tag_map)
9. test_dataset = generate_dataset(test_sentences, test_labels, sentence_vectorizer, tag_map)
```

## Model architecture:



```

1. def NER(len_tags, vocab_size, embedding_dim = 50):
2. model = tf.keras.Sequential(name = 'sequential')
3. model.add(
4. tf.keras.layers.Embedding(
5. input_dim=vocab_size + 1,
6. output_dim=embedding_dim,
7. mask_zero=True
8.)
9.)
10. model.add(
11. tf.keras.layers.LSTM(
12. units=embedding_dim,
13. return_sequences=True
14.)
15.)
16. model.add(
17. tf.keras.layers.Dense(
18. units=len_tags,
19. activation=tf.nn.log_softmax
20.)
21.)
22. return model

```

## Masked loss and metrics:

```

1. def masked_loss(y_true, y_pred):
2. y_true = tf.reshape(y_true, [-1])
3. y_pred = tf.reshape(y_pred, [-1, tf.shape(y_pred)[-1]])
4. mask = tf.cast(tf.not_equal(y_true, -1), dtype=tf.float32)
5. loss_fn = tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True)
6. y_true_masked = tf.boolean_mask(y_true, mask)
7. y_pred_masked = tf.boolean_mask(y_pred, mask)
8. loss = loss_fn(y_true_masked, y_pred_masked)
9. return loss
10.
11. def masked_accuracy(y_true, y_pred):
12. y_true = tf.cast(y_true, tf.float32)
13. mask = tf.not_equal(y_true, -1.0)
14. mask = tf.cast(mask, tf.float32)
15. y_pred_class = tf.argmax(y_pred, axis=-1, output_type=tf.int32)
16. y_pred_class = tf.cast(y_pred_class, tf.float32)
17. matches_true_pred = tf.equal(y_true, y_pred_class)
18. matches_true_pred = tf.cast(matches_true_pred, tf.float32)
19. matches_true_pred *= mask
20. masked_acc = tf.reduce_sum(matches_true_pred) / tf.reduce_sum(mask)
21. return masked_acc

```

## Experiment:

```

1. true_labels = [0,1,2,0]
2. predicted_logits = [[0.1,0.6,0.3], [0.2,0.7,0.1], [0.1, 0.5,0.4], [0.4,0.4,0.2]]
3. print(masked_loss(true_labels, predicted_logits))

```

Model compile:

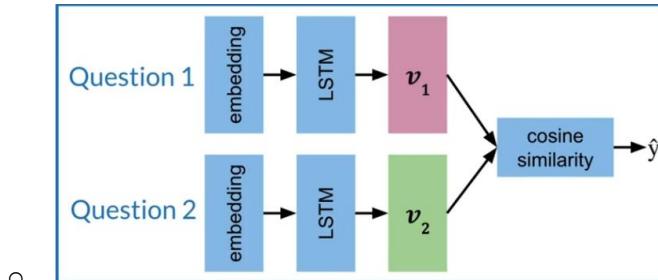
```
1. model.compile(optimizer=tf.keras.optimizers.Adam(0.01),
2. loss = masked_loss,
3. metrics = [masked_accuracy])
4. tf.keras.utils.set_random_seed(33)
5. BATCH_SIZE = 64
6. model.fit(train_dataset.batch(BATCH_SIZE),
7. validation_data = val_dataset.batch(BATCH_SIZE),
8. shuffle=True,
9. epochs = 2)
```

Make Prediction:

```
1. def predict(sentence, model, sentence_vectorizer, tag_map):
2. index_to_tag = {index: tag for tag, index in tag_map.items()}
3. sentence_vector = sentence_vectorizer([sentence])
4. predictions = model(sentence_vector)
5. predicted_indices = tf.argmax(predictions, axis=-1)
6. predicted_indices = tf.squeeze(predicted_indices)
7. pred = [index_to_tag[index] for index in predicted_indices.numpy() if index != -1]
8. return pred
1. sentence = "Peter Parker , the White House director of trade and manufacturing policy
of U.S , said in an interview on Sunday morning that the White House was working to
prepare for the possibility of a second wave of the coronavirus in the fall , though he
said it wouldn 't necessarily come"
2. predictions = predict(sentence, model, sentence_vectorizer, tag_map)
3. for x,y in zip(sentence.split(' '), predictions):
4. if y != '0':
5. print(x,y)
```

## Siamese Networks:

- Identify similarity between things
- Example: same:
  - o What is your age?
  - o How old are you?
- Applications:
  - o Handwritten checks
  - o Identify question duplicates.
  - o Check search engine queries.
- Architecture:



- o Contains two identical subnetworks (sister networks), merge to produce a similarity score between inputs.
- o The two subnetworks share identical parameters.
- o Process:
  - Transform each question into an embedding.
  - Feed embeddings into an LSTM layer to capture the meaning of each question.
  - Each LSTM outputs a vector, one for each question.
  - Compared the output vectors using cosine similarity → how similar the two questions are.

## Siamese Network lab:

```

1. import tensorflow as tf
2. from tensorflow.keras import layers
3. from tensorflow.keras.models import Model, Sequential
4. from tensorflow import math
5. import numpy
6. numpy.random.seed(10)

```

```

1. vocab_size = 500
2. model_dimension = 128
3. LSTM = Sequential()
4. LSTM.add(layers.Embedding(input_dim=vocab_size, output_dim=model_dimension))
5. LSTM.add(layers.LSTM(units=model_dimension, return_sequences = True))
6. LSTM.add(layers.AveragePooling1D())
7. LSTM.add(layers.Lambda(lambda x: math.l2_normalize(x)))
8. input1 = layers.Input((None,))
9. input2 = layers.Input((None,))
10. conc = layers.concatenate(axes=1)((LSTM(input1), LSTM(input2)))
11. Siamese = Model(inputs=(input1, input2), outputs=conc)

```

### Loss Function for Siamese Network:

- Anchor: the reference input we compare against
- Positive: A different input that is supposed to be similar to the anchor
- Negative: An input that is supposed to be dissimilar to the anchor
- Example:
  - o Anchor: How old are you?
  - o Positive: What is your age?
  - o Negative: Where are you from?

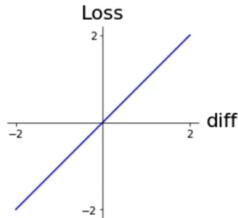
- Cosine Similarity:

$$\cos(v_1, v_2) = \frac{v_1 \cdot v_2}{\|v_1\| \|v_2\|}$$

$$s(v_1, v_2)$$

$$s(A, P) \approx 1$$

- o  $s(A, N) \approx -1$
- o To begin building a loss function:  $\text{diff} = s(A, N) - s(A, P)$



- Problem: can decrease towards negative values, can even go to negative infinity.

### Triplets

- 
- Goal: Minimize the difference between similarity of (A,N) and (A,P) with constraint With alpha margin
- $$\mathcal{L} = \begin{cases} 0; & \text{if } \text{diff} + \alpha \leq 0 \\ \text{diff} + \alpha; & \text{if } \text{diff} + \alpha > 0 \end{cases}$$
  - o Alpha margin is used to ensure difference between similarities is equal or close to the margin value.
  - o Allows you to have some safety.
  - o Controls how far  $\cos(A, P)$  is from  $\cos(A, N)$

### Triplet loss (simplified form):

- $$\mathcal{L}(A, P, N) = \max(\text{diff} + \alpha, 0)$$
- $$\max(-\cos(A, P) + \cos(A, N) + \alpha, 0)$$

### Triplet Selection:

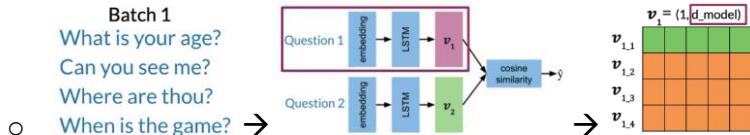
- Objective: Choose challenging triplets to efficiently train the model.
- Selection Process:
  - o Select a pair of duplicate questions (A and P).
  - o Select a non-duplicate question (N).
- Avoid Random Selection: Leads to trivial cases with loss close to zero.
- Hard Triplets:
  - o Select triplets that are difficult for the model to distinguish between P and N.
  - o The similarity between A and N is close to but less than between A and P.
  - o Focusing on hard triplets accelerates learning and improves model performance.

### Computing the cost

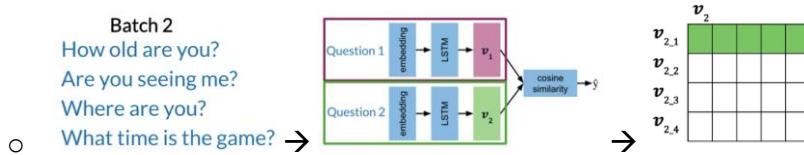
| Prepare the batches as follows: |                                                       |
|---------------------------------|-------------------------------------------------------|
| What is your age?               | How old are you? <span style="color: green;">✓</span> |
| Can you see me?                 | Are you seeing me?                                    |
| Where are thou?                 | Where are you?                                        |
| When is the game?               | What time is the game?                                |

b = 4

- Given batch 1, pass the batch into the model, get the output.



- Do the same for batch 2:



- Example of computing the cost:

$$s(v_1, v_2)$$

$$\begin{matrix}
 & & v_1 & \\
 & & -1 & -2 & -3 & -4 \\
 v_2 & -1 & 0.9 & -0.8 & 0.3 & -0.5 \\
 & -2 & -0.8 & 0.5 & 0.1 & -0.2 \\
 & -3 & 0.3 & 0.1 & 0.7 & -0.8 \\
 & -4 & -0.5 & -0.2 & -0.8 & 1.0
 \end{matrix}$$

- o The diagonal, the values are the similarities for all your positive example (Question duplicates)

$$\mathcal{L}(A, P, N) = \max(\text{diff} + \alpha, 0)$$

$$\text{diff} = s(A, N) - s(A, P)$$

$$\mathcal{J} = \sum_{i=1}^m \mathcal{L}(A^{(i)}, P^{(i)}, N^{(i)})$$

o

Modify the cost:

- Current lost:

|       |       | $s(v_1, v_2)$ |       |       |       |       |
|-------|-------|---------------|-------|-------|-------|-------|
|       |       | $v_1$         |       |       |       |       |
|       |       | $v_2$         | $v_1$ | $v_1$ | $v_1$ | $v_1$ |
| $v_1$ | $v_2$ | $v_1$         | $v_1$ | $v_1$ | $v_1$ | $v_1$ |
| $v_1$ | $v_2$ | -1            | 0.9   | -0.8  | 0.3   | -0.5  |
| $v_1$ | $v_2$ | -2            | -0.8  | 0.5   | 0.1   | -0.2  |
| $v_1$ | $v_2$ | -3            | 0.3   | 0.1   | 0.7   | -0.8  |
| $v_1$ | $v_2$ | -4            | -0.5  | -0.2  | -0.8  | 1.0   |

- Mean negative:
  - Mean of off-diagonal values in each row
- Cloest negative:
  - Off-diagonal value closest to (but less than) the value on diagonal in each row.

- Hard Negative Mining:

$$\mathcal{L}_{\text{Original}} = \max \underbrace{|s(A, N) - s(A, P)|}_{\text{diff}} + \alpha, 0$$

$$\mathcal{L}_1 = \max (\text{mean\_neg} - s(A, P) + \alpha, 0)$$

$$\mathcal{L}_2 = \max (\text{closest\_neg} - s(A, P) + \alpha, 0)$$

- - L1: helps the model converge faster during training by reducing noise by training on just the average of several observations
    - How? Noise is defined to be a small value that comes from a distribution that is centered around zero. The average of several noise values is usually zero.
    - The individual noise from the observations cancelled out.
  - L2: Creates a larger penalty by focusing on the most challenging examples.

Full lost:

$$\mathcal{L}_{\text{Full}}(A, P, N) = \mathcal{L}_1 + \mathcal{L}_2$$

$$\mathcal{J} = \sum_{i=1}^m \mathcal{L}_{\text{Full}}(A^{(i)}, P^{(i)}, N^{(i)})$$

Tensorflow implementation for Modified Triplet Loss:

```

1. v1 = np.array([1, 2, 3], dtype=float)
2. v2 = np.array([1, 2, 3.5], dtype=float) # notice the 3rd element is offset by 0.5
3. print("-- Inputs --")
4. print("v1 : ", v1)
5. print("v2 : ", v2, "\n")
6. def cosine_similarity(v1, v2):
7. numerator = tf.math.reduce_sum(v1*v2)
8. denominator = tf.math.sqrt(tf.math.reduce_sum(v1*v1) * tf.math.reduce_sum(v2*v2))
9. return numerator / denominator

```

```

1. v1_1 = np.array([1.0, 2.0, 3.0])
2. v1_2 = np.array([9.0, 8.0, 7.0])
3. v1_3 = np.array([-1.0, -4.0, -2.0])
4. v1_4 = np.array([1.0, -7.0, 2.0])
5. v1 = np.vstack([v1_1, v1_2, v1_3, v1_4])
6. v2_1 = v1_1 + np.random.normal(0, 2, 3)
7. v2_2 = v1_2 + np.random.normal(0, 2, 3)
8. v2_3 = v1_3 + np.random.normal(0, 2, 3)
9. v2_4 = v1_4 + np.random.normal(0, 2, 3)
10. v2 = np.vstack([v2_1, v2_2, v2_3, v2_4])
11. b = len(v1)
12. sim_1 = np.zeros([b, b])
13. for row in range(0, sim_1.shape[0]):
14. for col in range(0, sim_1.shape[1]):
15. sim_1[row, col] = cosine_similarity(v2[row], v1[col]).numpy()

1. def norm(x):
2. return tf.math.l2_normalize(x, axis=1)
3. sim_2 = tf.linalg.matmul(norm(v2), norm(v1), transpose_b=True)

```

### Hard Negative Mining:

```

1. sim_hardcoded = np.array(
2. [
3. [0.9, -0.8, 0.3, -0.5],
4. [-0.4, 0.5, 0.1, -0.1],
5. [0.3, 0.1, -0.4, -0.8],
6. [-0.5, -0.2, -0.7, 0.5],
7.]
8.)
9. sim = sim_hardcoded
10. b = sim.shape[0]
11. sim_ap = np.diag(sim)
12. sim_an = sim - np.diag(sim_ap)
13. mean_neg = np.sum(sim_an, axis=1, keepdims=True) / (b - 1)
14. mask_1 = np.identity(b) == 1
15. mask_2 = sim_an > sim_ap.reshape(b, 1)
16. mask = mask_1 | mask_2
17. sim_an_masked = np.copy(sim_an)
18. sim_an_masked[mask] = -2
19. closest_neg = np.max(sim_an_masked, axis=1, keepdims=True)

```

Implementation in tensorflow:

```
1. sim_hardcoded = np.array(
2. [
3. [0.9, -0.8, 0.3, -0.5],
4. [-0.4, 0.5, 0.1, -0.1],
5. [0.3, 0.1, -0.4, -0.8],
6. [-0.5, -0.2, -0.7, 0.5],
7.]
8.)
9. sim = sim_hardcoded
10. b = sim.shape[0]
11. sim_ap = tf.linalg.diag_part(sim)
12. sim_an = sim - tf.linalg.diag(sim_ap)
13. mean_neg = tf.math.reduce_sum(sim_an, axis=1) / (b - 1)
14. mask_1 = tf.eye(b) == 1
15. mask_2 = sim_an > tf.expand_dims(sim_ap, 1)
16. mask = tf.cast(mask_1 | mask_2, tf.float64)
17. sim_an_masked = sim_an - 2.0*mask
18. closest_neg = tf.math.reduce_max(sim_an_masked, axis=1)
```

### The Loss Functions

$$\mathcal{L}_1 = \max(\text{mean\_neg} - s(A, P) + \alpha, 0)$$

$$\mathcal{L}_2 = \max(\text{closest\_neg} - s(A, P) + \alpha, 0)$$

$$- \quad \mathcal{L}_{\text{Full}} = \mathcal{L}_1 + \mathcal{L}_2$$

```
1. alpha = 0.25
2. l_1 = tf.maximum(mean_neg - sim_ap + alpha, 0)
3. l_2 = tf.maximum(closest_neg - sim_ap + alpha, 0)
4. l_full = l_1 + l_2
5. cost = tf.math.reduce_sum(l_full)
```

## One shot learning

- $s(sig1, sig2) > \tau$  ✓ → Yes
- $s(sig1, sig2) \leq \tau$
  - Challenge:
    - o Adding new classes necessitates costly model retraining.
  - Solution:
    - o Recognize a new class from a single example without retraining
  - Method:
    - o Utilizes a learned similarity function to compare and evaluate signatures based on similarity scores.
  - Advantages:
    - o Significantly more efficient for dynamic datasets, such as in banking scenarios, where new signatures are regularly introduced.

### Training:

- The subnetworks have identical parameters so you only need to train one set of them, which are then shared by both.

### Testing:

- Convert each input into an array of numbers.
- Feed arrays into the model.
- Compare  $v_1, v_2$  using cosine similarity.
- Test against a threshold  $\tau$ .

Evaluating a Siamese model:

Libraries:

```
1. import numpy as np
2. import tensorflow as tf
3. import tensorflow.math as math
4. import tensorflow.linalg as linalg
```

Question vector:

```
1. q1 = np.load('../data/q1.npy') #shape: (512,64)
2. q2 = np.load('../data/q2.npy') #shape: (512,64)
3. y_test = np.load('../data/y_test.npy') #shape: (512,)
4. v1 = np.load('../data/v1.npy') #shape: (512,128)
5. v2 = np.load('../data/v2.npy') #shape: (512,128)
```

Calculating the accuracy:

```
1. accuracy = 0
2. batch_size = 512
3. threshold = 0.7
4. for j in range(batch_size):
5. d = math.reduce_sum(v1[j]*v2[j])
6. res = d > threshold
7. accuracy += tf.cast(y_test[j] == res, tf.int32)
8. accuracy = accuracy / batch_size
```

## Question duplicates

Libraries % data:

```
1. import os
2. os.environ['TF_CPP_MIN_LOG_LEVEL'] = '3'
3. import os
4. import numpy as np
5. import pandas as pd
6. import random as rnd
7. import tensorflow as tf
8. rnd.seed(34)
9. data = pd.read_csv("questions.csv")
10. N = len(data)
11. N_train = 300000
12. N_test = 10240
13. data_train = data[:N_train]
14. data_test = data[N_train:N_train + N_test]
15. del (data) # remove to free memory
```

Data split:

```
1. td_index = data_train['is_duplicate'] == 1
2. td_index = [i for i, x in enumerate(td_index) if x]
3. Q1_train = np.array(data_train['question1'][td_index])
4. Q2_train = np.array(data_train['question2'][td_index])
5. Q1_test = np.array(data_test['question1'])
6. Q2_test = np.array(data_test['question2'])
7. y_test = np.array(data_test['is_duplicate'])

- data_train['question1'][5]:
 o Astrology: I am a Capricorn Sun Cap moon and cap rising...what does that
 say about me?
- data_train['question2'][5]:
 o I'm a triple Capricorn (Sun, Moon and ascendant in Capricorn) What does
 this say about me?
- data_train['is_duplicate']
 o 1
```

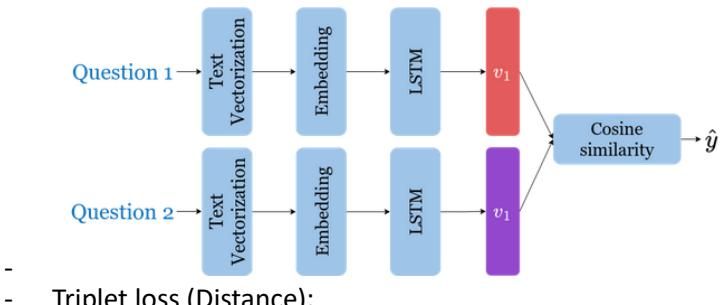
Splitting the data:

```
1. cut_off = int(len(Q1_train) * 0.8)
2. train_Q1, train_Q2 = Q1_train[:cut_off], Q2_train[:cut_off]
3. val_Q1, val_Q2 = Q1_train[cut_off:], Q2_train[cut_off:]
```

Learning question encoding:

```
1. tf.random.set_seed(0)
2. text_vectorization =
tf.keras.layers.TextVectorization(output_mode='int',split='whitespace',
standardize='strip_punctuation')
3. text_vectorization.adapt(np.concatenate((Q1_train,Q2_train)))
```

Defining the Siamese model:



- Triplet loss (Distance):

$$\mathcal{L}(A, P, N) = \max (\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \alpha, 0)$$

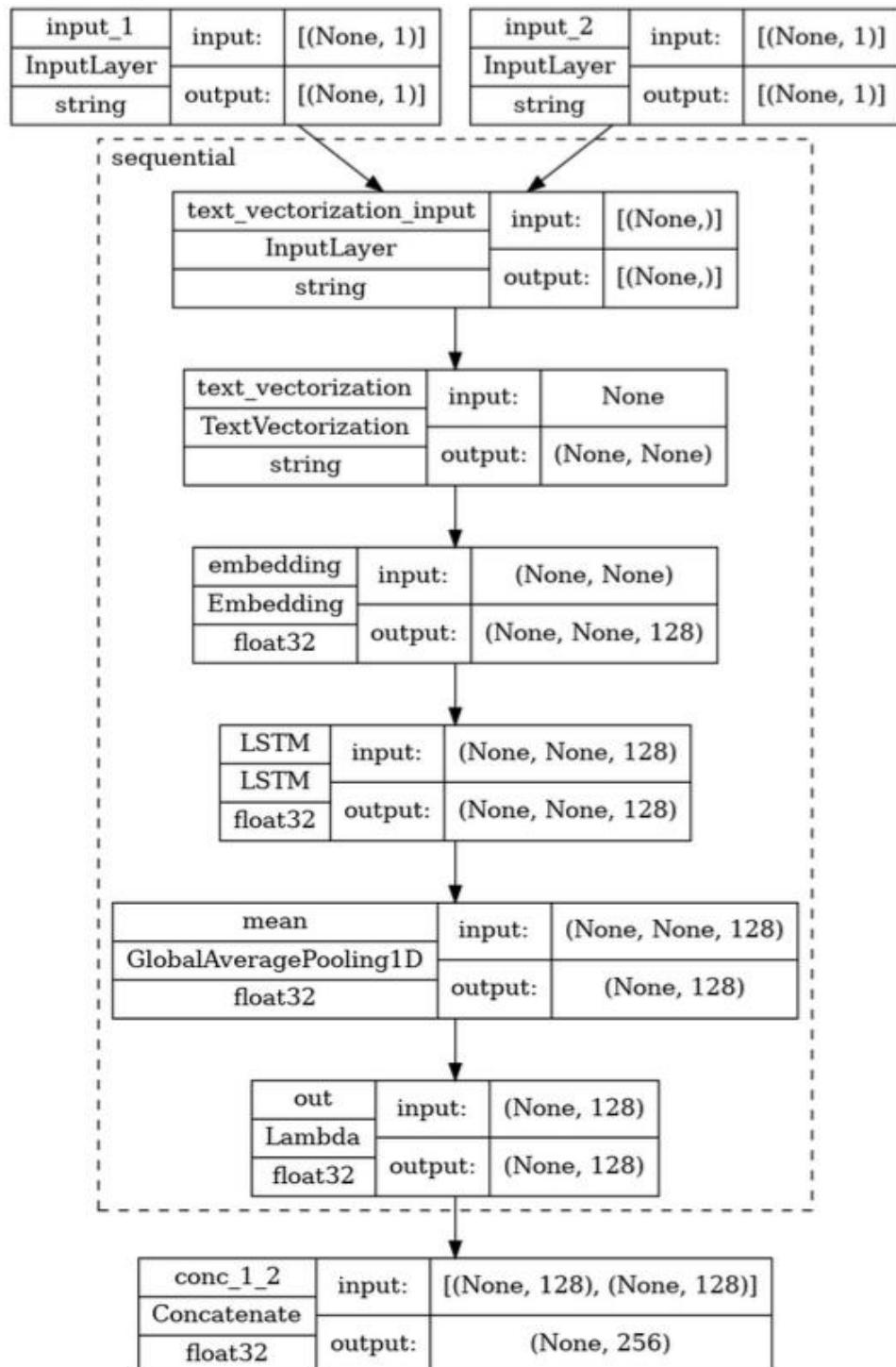
Siamese:

```
1. def Siamese(text_vectorizer, vocab_size=36224, d_feature=128):
2. branch = tf.keras.models.Sequential(name='sequential')
3. branch.add(text_vectorizer)
4. branch.add(tf.keras.layers.Embedding(input_dim=vocab_size, output_dim=d_feature, name='embedding'))
5. branch.add(tf.keras.layers.LSTM(d_feature, return_sequences=True, name="LSTM"))
6. branch.add(tf.keras.layers.GlobalAveragePooling1D(name="mean"))
7. branch.add(tf.keras.layers.Lambda(lambda x: tf.math.l2_normalize(x, axis=1), name='out'))
8. input1 = tf.keras.Input(shape=(1,), dtype=tf.string, name='input_1')
9. input2 = tf.keras.Input(shape=(1,), dtype=tf.string, name='input_2')
10. branch1 = branch(input1)
11. branch2 = branch(input2)
12. conc = tf.keras.layers.concatenate([branch1, branch2], axis=-1, name='conc_1_2')
13. return tf.keras.Model(inputs=[input1, input2], outputs=conc, name="SiameseModel")
```

```
1. model = Siamese(text_vectorization, vocab_size=text_vectorization.vocabulary_size())
2. model.build(input_shape=None)
```

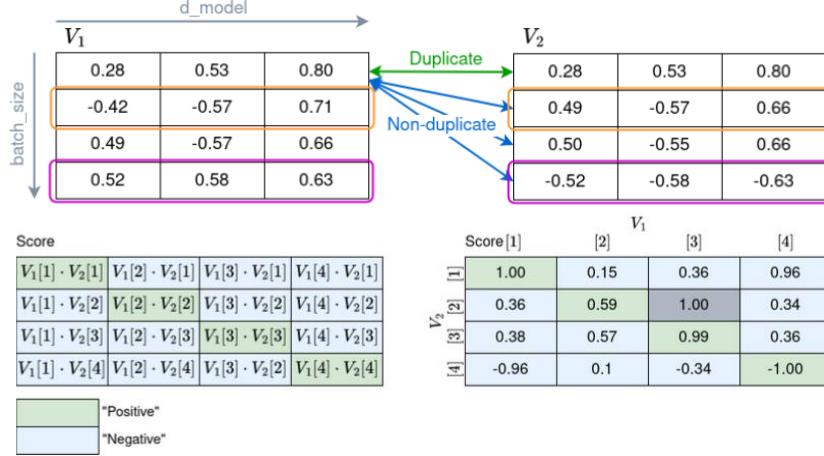
Plot model:

```
1. tf.keras.utils.plot_model(
2. model,
3. to_file="model.png",
4. show_shapes=True,
5. show_dtype=True,
6. show_layer_names=True,
7. rankdir="TB",
8. expand_nested=True)
```



Triplet Loss with Hard Negative Mining:

$$\begin{aligned} \text{Loss}_1(\mathcal{A}, \mathcal{P}, \mathcal{N}) &= \max(-\cos(A, P) + \text{mean}_{\text{neg}} + \alpha, 0) \\ \text{Loss}_2(\mathcal{A}, \mathcal{P}, \mathcal{N}) &= \max(-\cos(A, P) + \text{closest}_{\text{neg}} + \alpha, 0) \\ \text{Loss}(\mathcal{A}, \mathcal{P}, \mathcal{N}) &= \text{mean}(\text{Loss}_1 + \text{Loss}_2) \end{aligned}$$



Triplet Loss Function:

```

1. def TripletLossFn(v1, v2, margin=0.25):
2. scores = tf.linalg.matmul(v1, v2, transpose_b = True)
3. batch_size = tf.cast(tf.shape(v1)[0], scores.dtype)
4. positive = tf.linalg.diag_part(scores)
5. negative_zero_on_duplicate = scores - tf.linalg.diag(positive)
6. mean_negative = tf.math.reduce_sum(negative_zero_on_duplicate, axis=1) / (batch_size - 1)
7. mask_diagonal = tf.eye(batch_size) == 1
8. mask_greater_than_positive = negative_zero_on_duplicate > tf.expand_dims(positive, 1)
9. mask_exclude_positives = tf.cast(mask_diagonal | mask_greater_than_positive, scores.dtype)
10. negative_without_positive = negative_zero_on_duplicate - mask_exclude_positives * 2.0
11. closest_negative = tf.math.reduce_max(negative_without_positive, axis=1)
12. triplet_loss1 = tf.maximum(0.0, closest_negative - positive + margin)
13. triplet_loss2 = tf.maximum(0.0, mean_negative - positive + margin)
14. loss_full = triplet_loss1 + triplet_loss2
15. triplet_loss = tf.math.reduce_sum(loss_full)
16. return triplet_loss

1. def TripletLoss(labels, out, margin=0.25):
2. _, embedding_size = out.shape # get embedding size
3. v1 = out[:, :int(embedding_size/2)] # Extract v1 from out
4. v2 = out[:, int(embedding_size/2):] # Extract v2 from out
5. return TripletLossFn(v1, v2, margin=margin)

```

## Model:

```
1. def train_model(Siamese, TripletLoss, text_vectorizer, train_dataset, val_dataset,
d_feature=128, lr=0.01, train_steps=5):
2. vocab_size = len(text_vectorizer.get_vocabulary())
3. model = Siamese(text_vectorizer,
4. vocab_size = vocab_size,
5. d_feature = d_feature)
6. model.compile(loss=TripletLoss,
7. optimizer=tf.keras.optimizers.Adam(lr=lr)
8.)
9. history = model.fit(train_dataset,
10. epochs=train_steps,
11. validation_data=val_dataset
12.)
13. return model
```

## Classify:

```
1. def classify(test_Q1, test_Q2, y_test, threshold, model, batch_size=64, verbose=True):
2. test_gen = tf.data.Dataset.from_tensor_slices(((test_Q1, test_Q2), None)).batch(batch_size=batch_size)
3. pred = model.predict(test_gen)
4. v1, v2 = np.split(pred, 2, axis=1)
5. d = tf.reduce_sum(v1 * v2, axis=1) / (tf.linalg.norm(v1, axis=1) * tf.linalg.norm(v2, axis=1))
6. y_pred = tf.cast(d > threshold, tf.float64)
7. accuracy = tf.reduce_mean(tf.cast(tf.equal(y_pred, y_test), tf.float64))
8. cm = tf.math.confusion_matrix(y_test, y_pred)
9. return accuracy, cm
```

## Predict

```
1. def predict(question1, question2, threshold, model, verbose=False):
2. generator = tf.data.Dataset.from_tensor_slices(([question1], [question2]),
None)).batch(batch_size=1)
3. v1v2 = model.predict(generator)
4. v1, v2 = np.split(v1v2, 2, axis=1)
5. d = tf.reduce_sum(v1 * v2, axis=1) / (tf.linalg.norm(v1, axis=1) *
tf.linalg.norm(v2, axis=1))
6. res = d > threshold
7. if(verbose):
8. print("Q1 = ", question1, "\nQ2 = ", question2)
9. print("d = ", d.numpy())
10. print("res = ", res.numpy())
11. return res.numpy()[0]
12.
```