



Finding Rare Events in Text

Debanjana Banerjee

Senior Data Scientist, Walmart Global Tech

Sep 15th, 2021

Workshop,

Open Data Science Conference, India

Walmart  Global Tech



Hello, fellow explorers!



Debanjana Banerjee

Senior Data Scientist,
Walmart Global Tech





Your Instructor for the day!!

My first workshop
- please be kind!

- Born and bred in the world of Statistics (I worship the old gods!)
- 4+ years in the industry
- In love with the diverse ML use cases retail has to offer



Our Agenda today

1. Introduction & Examples
2. Rare Events in Text: Defining Characteristics
3. Q&A 
4. Text Mining
5. PU Learning
6. Q&A 
7. Break 
8. Semi Supervised Learners
9. iCASSTLe
10. Q&A 

View slides for this workshop here:

[github.com/debanjana-banerjee/Finding-Rare-Events-in-Text-ODSC-2021-](https://github.com/debanjana-banerjee/Finding-Rare-Events-in-Text-ODSC-2021)

Available after the session!



Wands at the ready!

We will communicate on the Event X Ai Platform

Channel: **wed-debanjana-banerjee-finding-rare-events-in-text**

What do we need to do before we get started?

1. Create a folder named **FRET** on your Google Drive with sufficient space (+300MB)
2. Download **GloVe** file shared and load to folder **FRET** on your Google drive (this may take a while!)
3. Two data files will be shared shortly. We will load those onto **the same folder** on Google Drive



Finding Rare Events in Text – what is it all about?



Fraud Detection

01

Ethics Compliance

02

Consumer Safety

03

Social Media Compliance

04

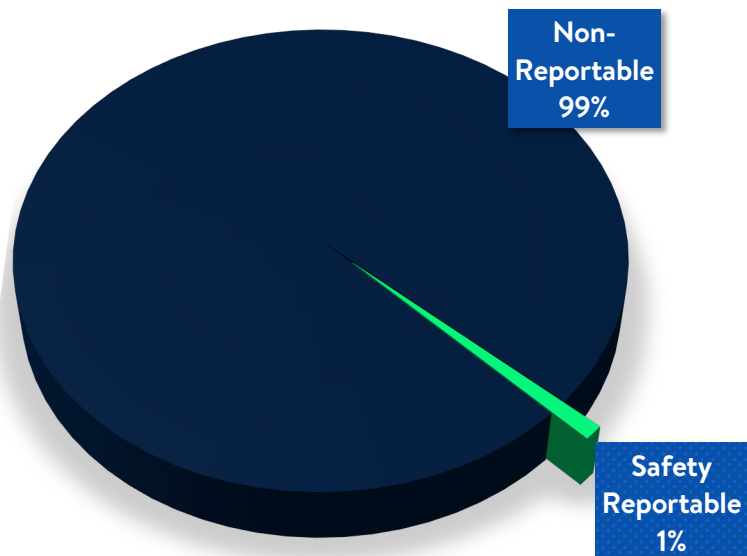
Rare Events in Text

Rare Events!!



Rare Events in Consumer Safety

Customer Reviews



Non-Reportable Reviews

- Product Enquiry
- Compliments
- User guidance
- Discounts
- Sharing product info, etc.



Safety-Reportable Reviews

- Freshness of Produce
- Allergen Ingredients
- Safety guidelines on electronics
- Slippery tubs
- Brake safety on bikes, etc.

Let's look at a few examples!



Rare Events are not the same as Anomalies

Anomalies

- Generally, do not adhere to specification (i.e., cannot be described by a common theme)
- Variability in rate of imbalance is limited (rate of imbalance is always very high)
- Calls for unsupervised learning
- Sample dependent

Rare Events

- Adhere to certain specification (i.e., usually described by a common theme)
- Allows variability in rate of imbalance (depending on degree of rarity)
- May pertain to supervised or, unsupervised learning
- Sample independent

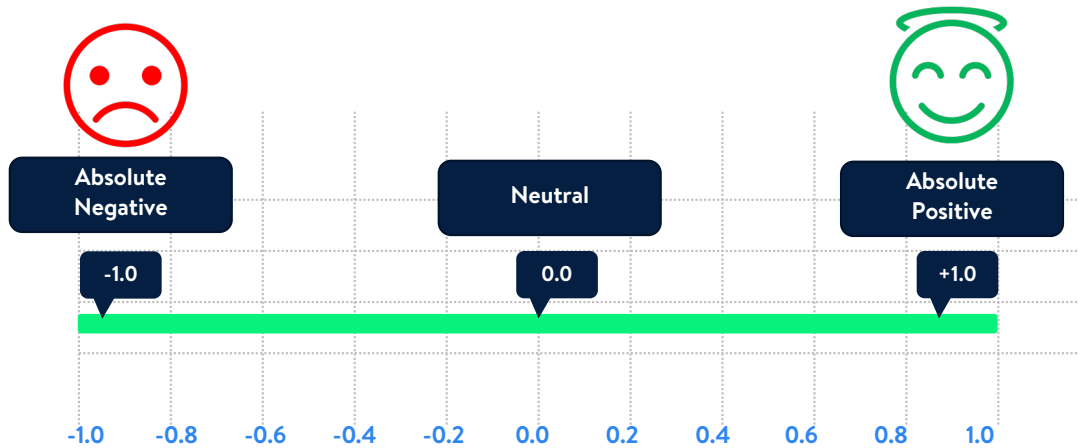


Rare Events in Text

Defining Characteristics



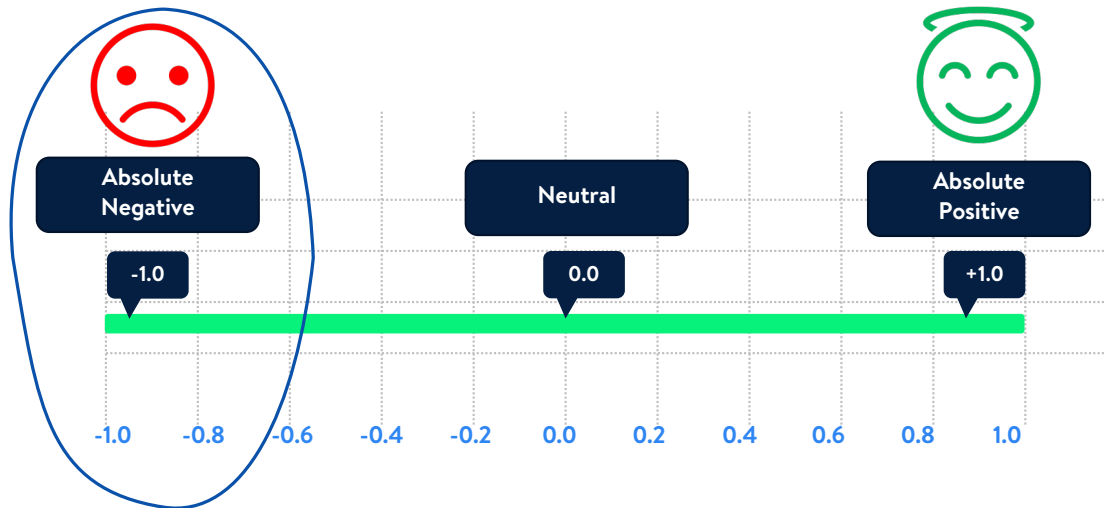
Sentiment Inclination



Depending on the specification, a rare event is likely to have inclination toward a polar sentiment



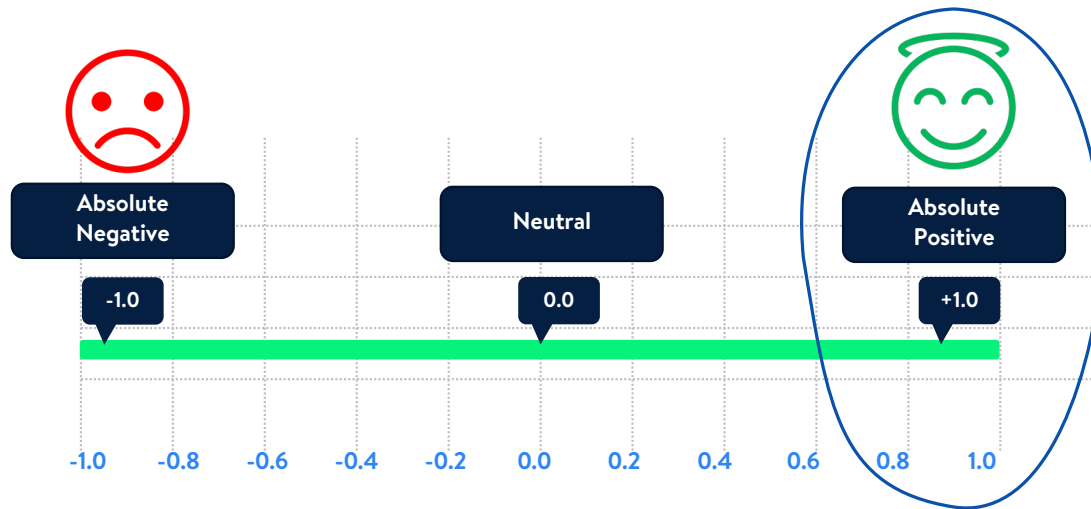
Sentiment Inclination



Depending on the specification, a rare event is likely to have inclination toward a polar sentiment



Sentiment Inclination



Depending on the specification, a rare event is likely to have inclination toward a polar sentiment



Sentiment Inclination



Depending on the specification, a rare event is likely to have inclination toward a polar sentiment



Safety-Reportable Reviews

“My salad was stale. The lettuce had blackened”

Sentiment Polarity: **-0.12**



Token Sensitivity



Review A

“Customer stated ABC microwave started smoking when she used it”



Review B

“Customer stated ABC microwave worked smoothly when she used it”



Token Sensitivity



Review A

“Customer stated ABC microwave started smoking when she used it”



Review B

“Customer stated ABC microwave worked smoothly when she used it”



Review C

“Customer stated he purchased the bike a year ago and the chains came off”



Review D

“Customer stated he purchased the bike a week ago and the chains came off”



Token Sensitivity



Review A

“Customer stated ABC microwave **started smoking** when she used it”



Review B

“Customer stated ABC microwave **worked smoothly** when she used it”



Review C

“Customer stated he purchased the bike **a year ago** and the chains came off”



Review D

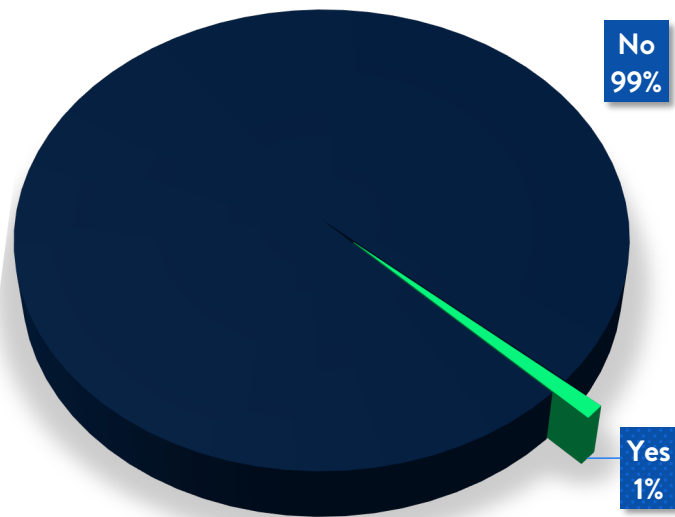
“Customer stated he purchased the **bike a week ago** and the chains came off”

For rare events in text, only certain tokens in the text dictate reportability or, non-reportability of the case



Data Availability

Is it a rare event?



- Rare Events are scarcely observed
 - Poor data availability

Quantity:

Number of available instances is low

Quality:

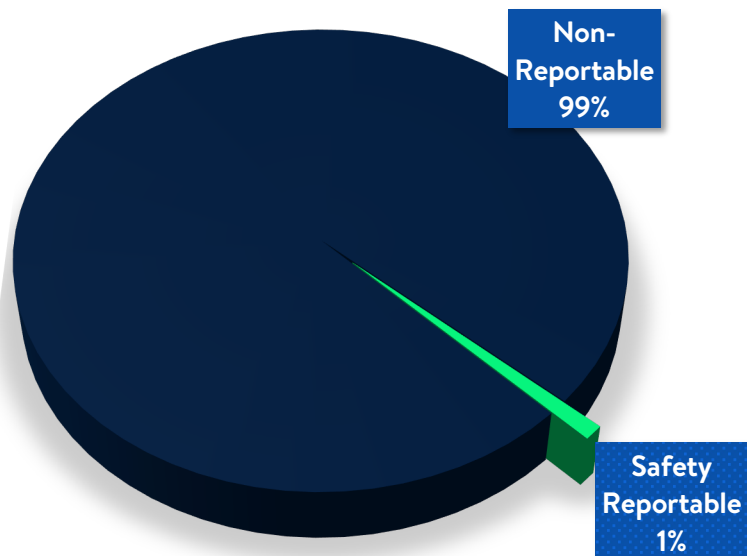
Available instances are not exhaustive for event specification

- Similar instances are repeated
- Not all sub-classes of the rare events are observed in the data



Data Availability

Consumer Safety



Non-Reportable Reviews

- Product Enquiry
- Compliments
- User guidance
- Discounts
- Sharing product info, etc.



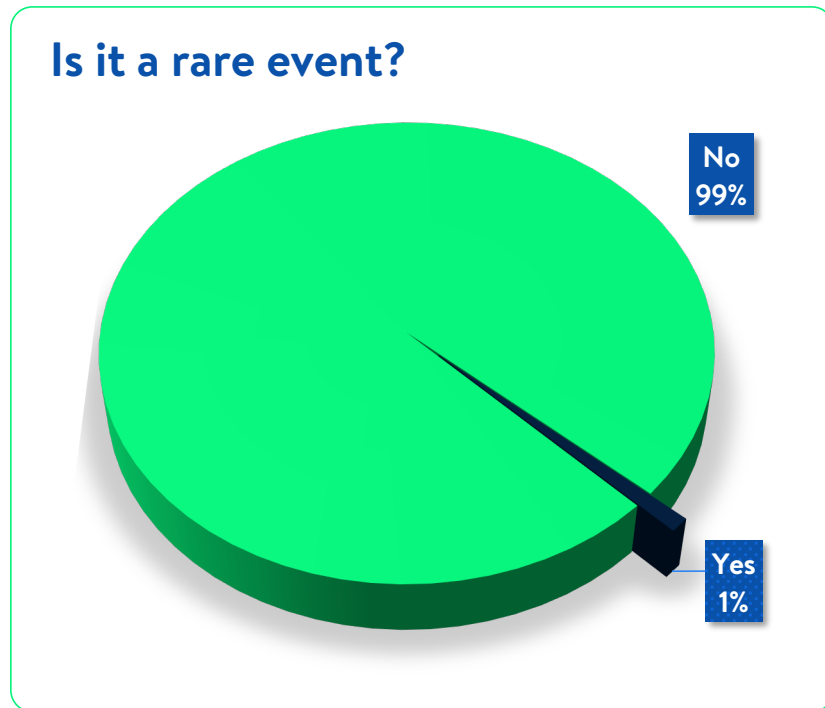
Safety-Reportable Reviews

- Freshness of Produce
- Allergen Ingredients
- Safety guidelines on electronics
- Slippery tubs
- Brake safety on bikes, etc

Of these, your data may show only a few kind!



No or, poor record of quality non-reportables



- Quality of Non-Reportables is important

Non-Reportable Case A

Customer stated ABC microwave **worked smoothly** when she used it

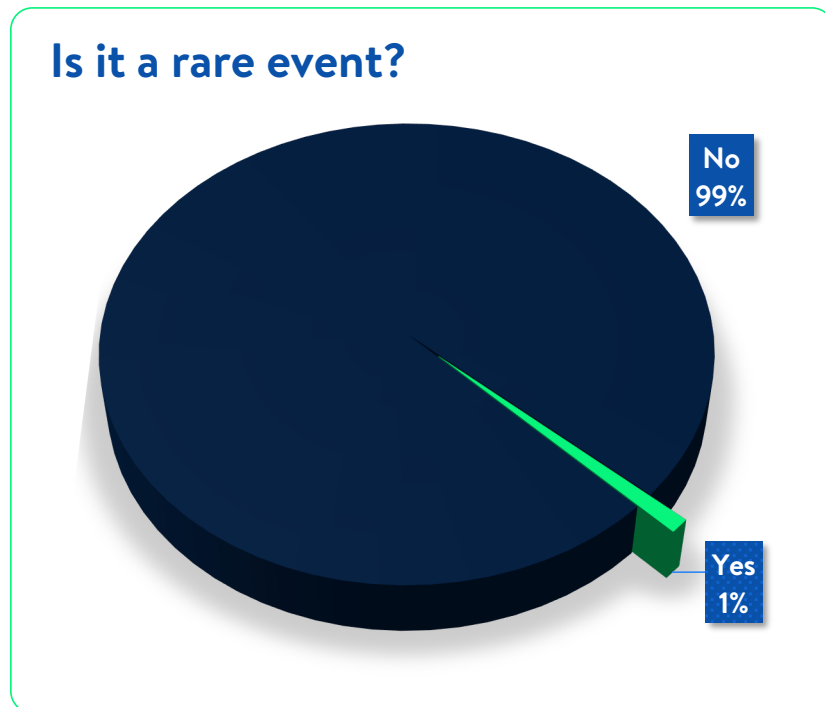


Non-Reportable Case B

Customer stated he purchased the bike **a year ago** and the chains came off



Data Quality Issues in Rare Event Extraction



Poor Quality and low quantity of Positives (Reportables)

Poor Quality or, no record of Negatives (Non-Reportables)

Positive Unlabeled Learning



Questions so far?

Please post on Event X Ai Platform

Channel: **wed-debanjana-banerjee-finding-rare-events-in-text**



Text Mining



Text Pre-Processing

- Text Cleaning
- Lemmatization
- Numeric Representation of Text

Please upload the following csv's to FRET folder in google drive

- FRET_Test.csv
- FRET_Positive.csv

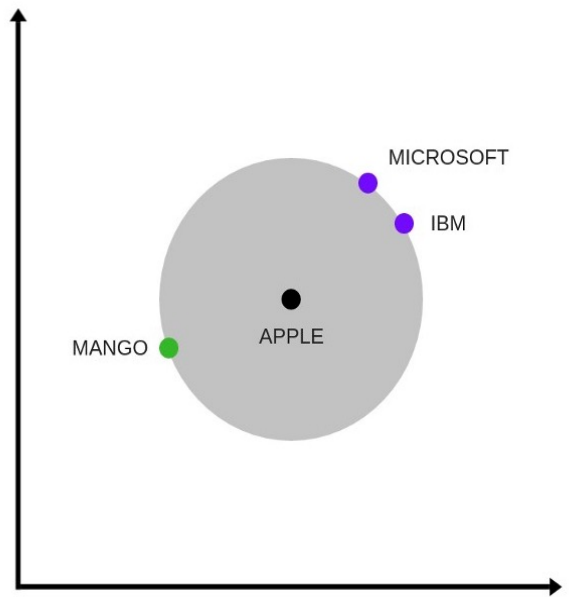
Let's try this hands-on!

Notebook link is shared on [wed-debanjana-banerjee-finding-rare-events-in-text](#)



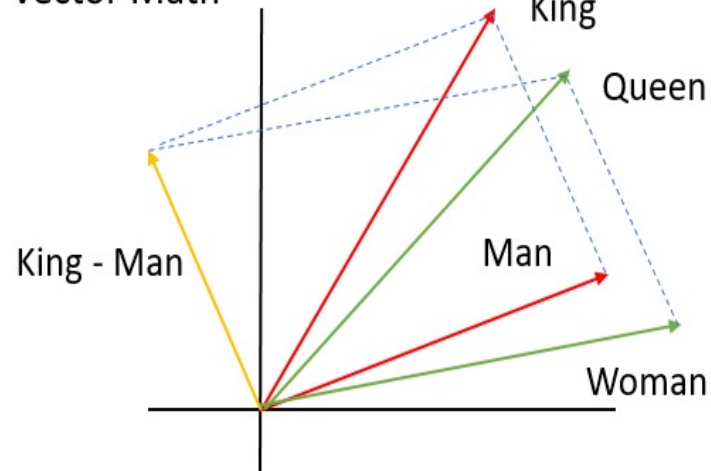
What are Text Embeddings?

Embeddings are numeric representation of meaningful words (or, phrases) such that their inter-relationships are preserved in the vector forms.



In the vector space above, the word 'Apple' is equidistant from the words 'Microsoft', 'IBM' and 'Mango', justifying two meanings of the word 'Apple'

Vector Math



In the vector space above, the distance between 'King' and 'Queen' is the same as that of 'Man' and 'Woman' capturing pairwise symmetry of words



Training your own text embeddings using a corpus

- Context Window
- Term Co-occurrence Matrix

Original Text 1

- The child got a rash from the diapers

Original Text 2

- The diapers did not fit my child

Cleaned Text 1

- child got rash diaper

Cleaned Text 2

- diaper not fit child

Vocab: child, got, rash, diaper, not, fit

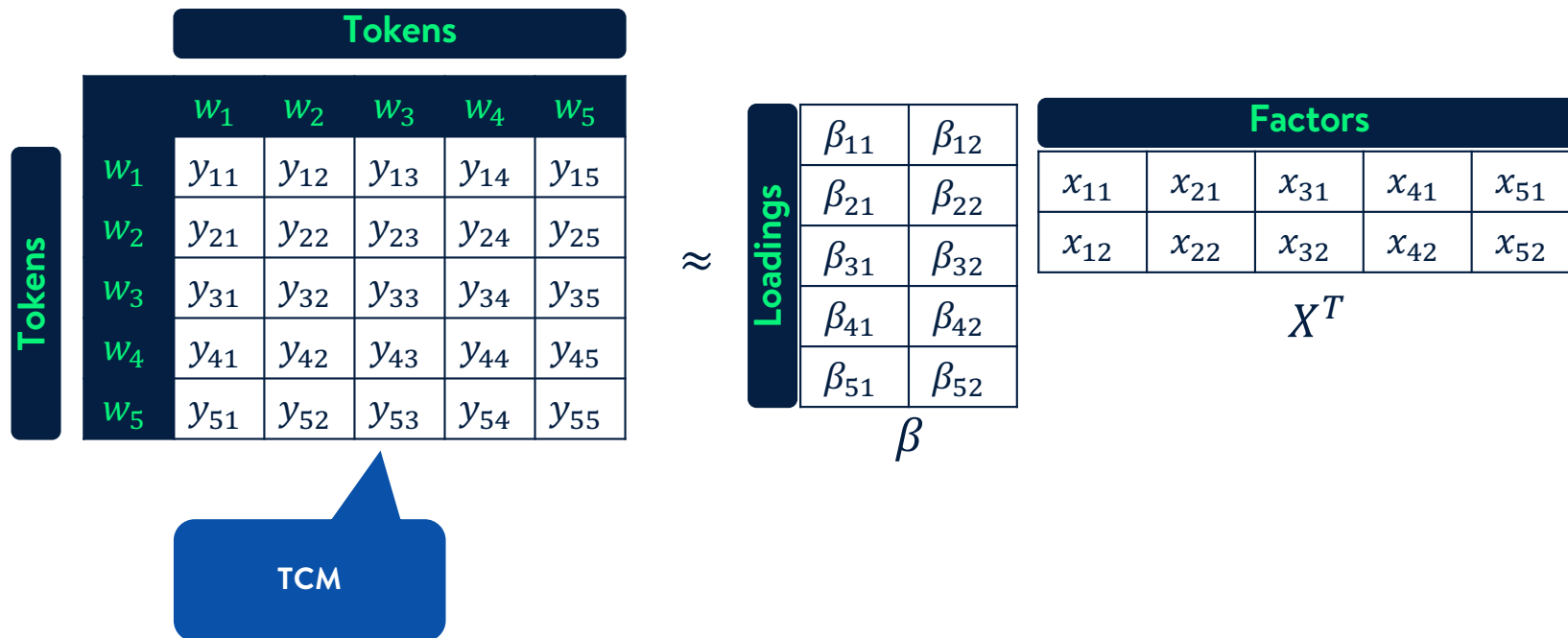
Context Window : ∞

	child	got	rash	diaper	not	fit
child	0	1	1	2	1	1
got	1	0	1	1	0	0
rash	1	1	0	1	0	0
diaper	2	1	1	0	1	1
not	1	0	0	1	0	1
fit	1	0	0	1	1	0

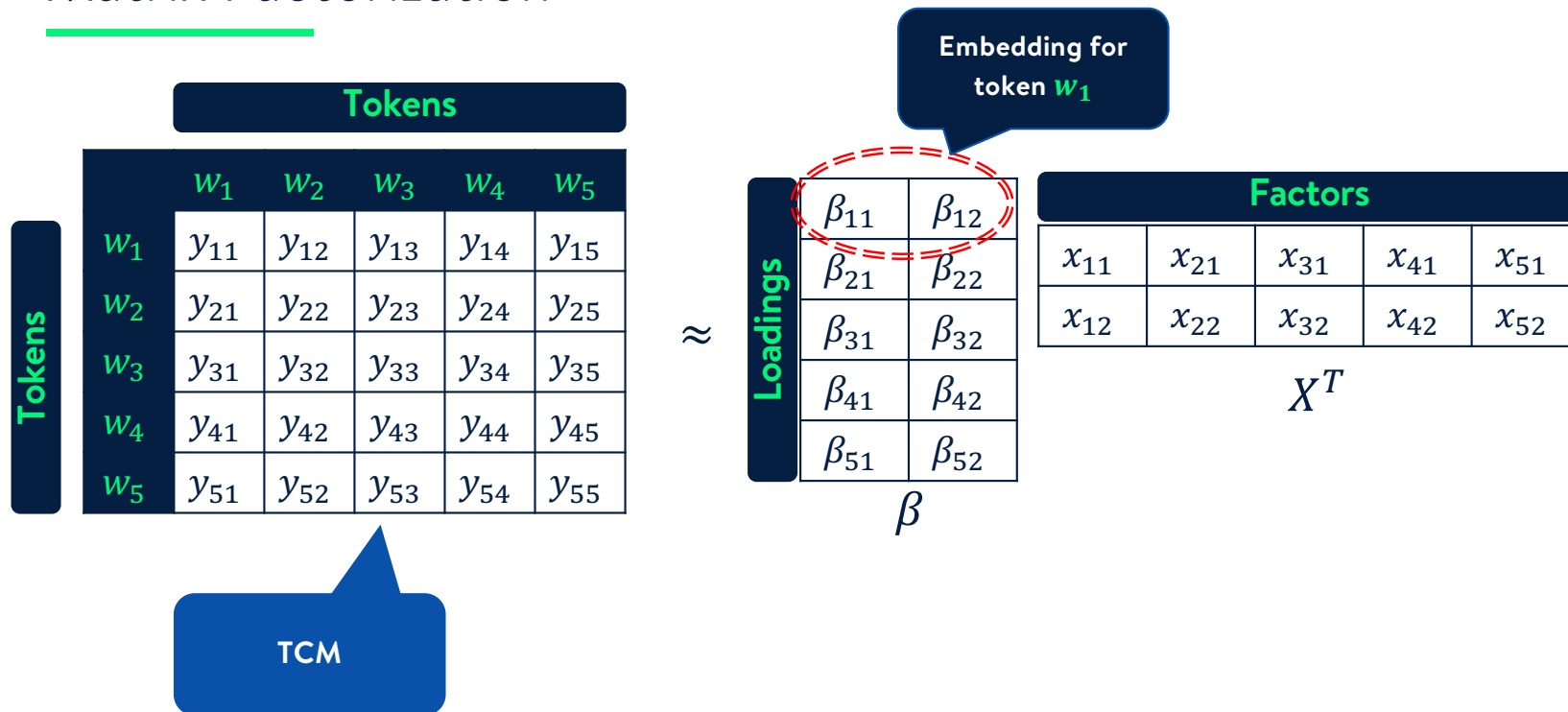
TCM



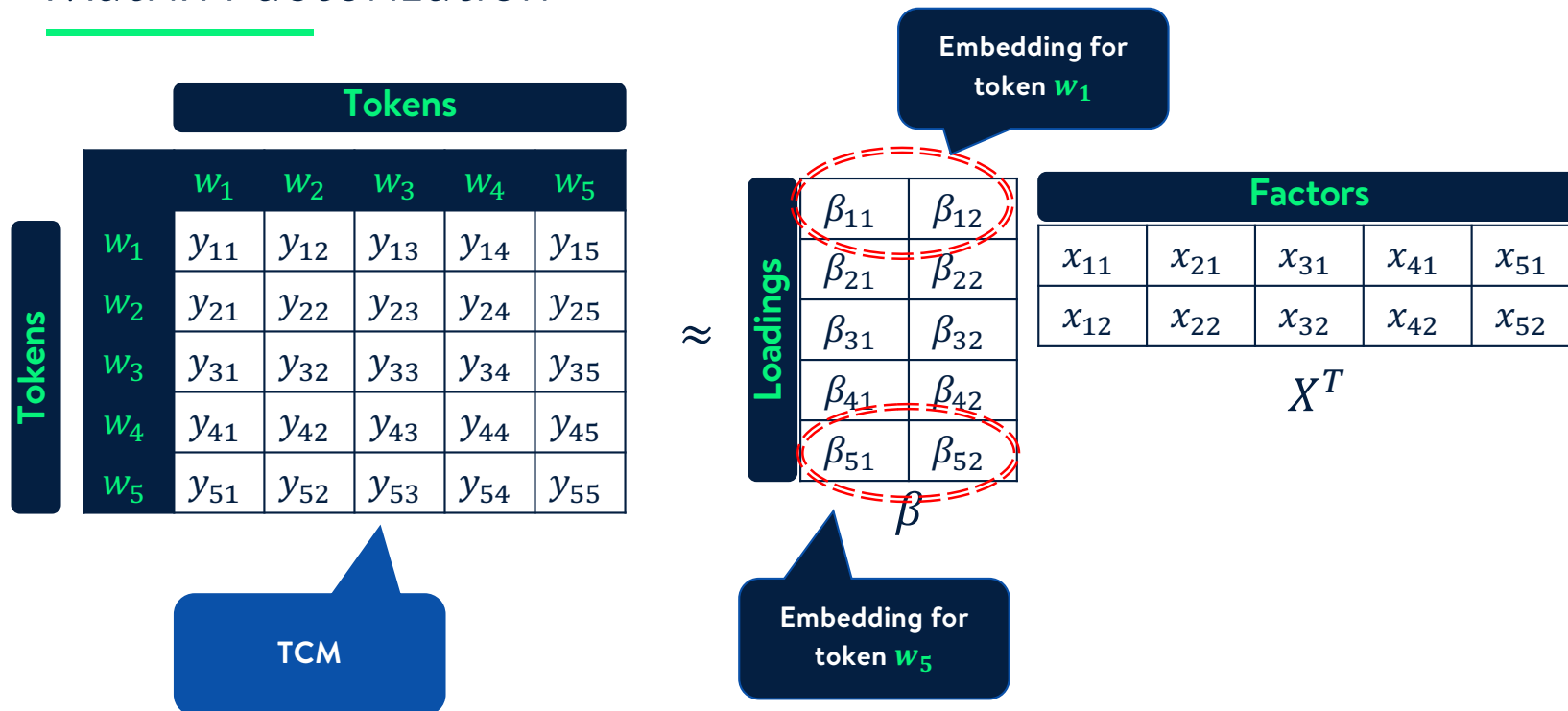
Matrix Factorization



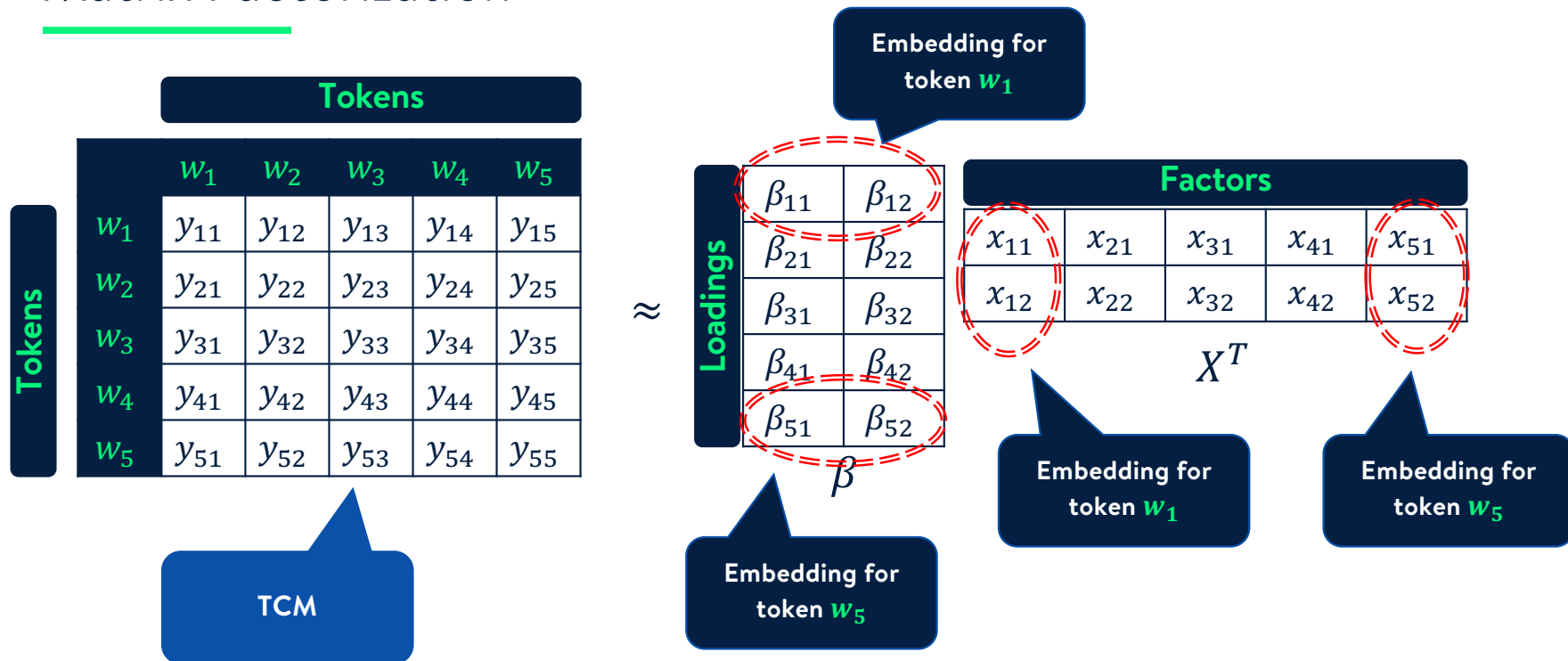
Matrix Factorization



Matrix Factorization



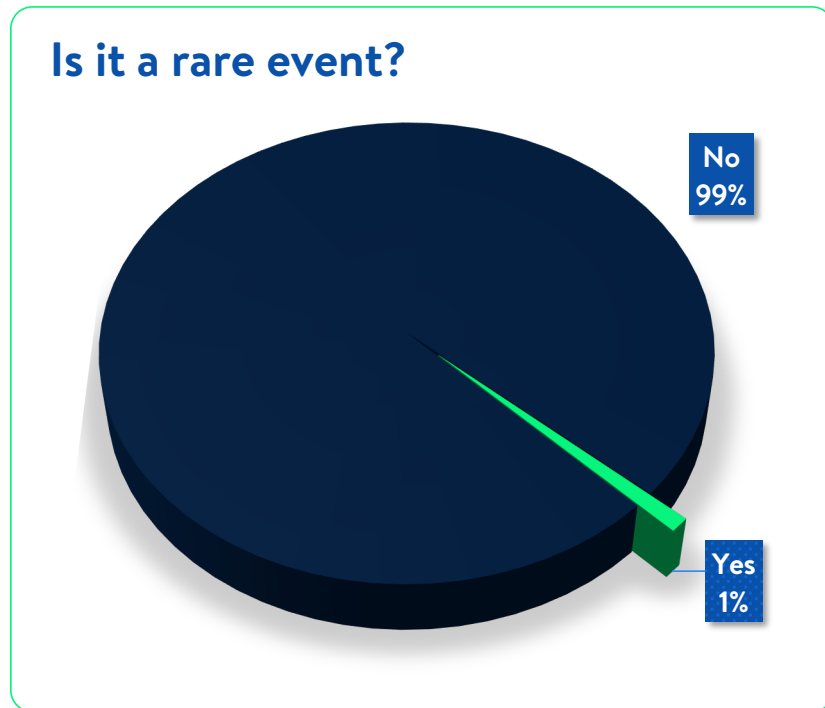
Matrix Factorization



Positive Unlabeled (PU) Learning



Data Quality Issues in Rare Event Extraction



Poor Quality and low quantity of Positives (Reportables)

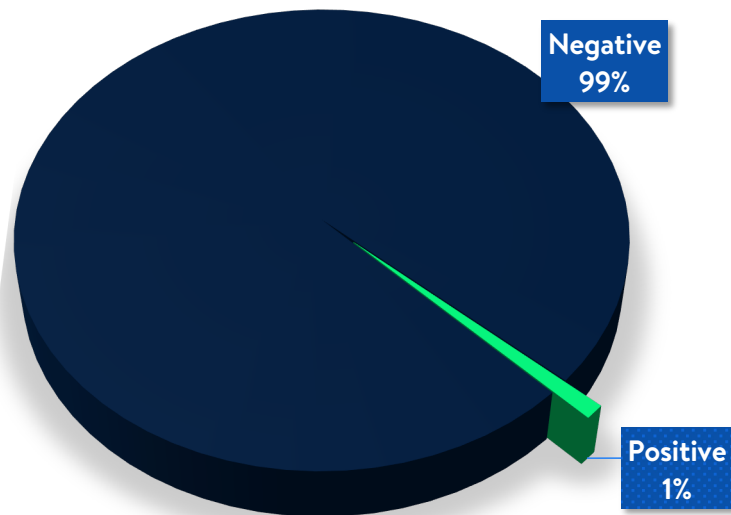
Poor Quality or, no record of Negatives (Non-Reportables)

Positive Unlabeled Learning



Rare Event Extraction as a Binary Classification technique

Is it a rare event of interest?



Rare Event of Interest: Consumer Safety



Negative

The text (review) is non-reportable
True > 99% of the times



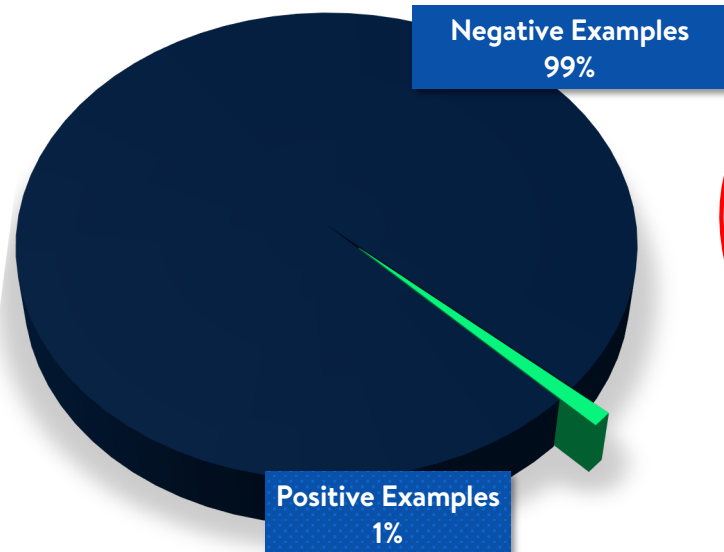
Positive

The text (review) is safety-reportable
True < 1% of the times



Training Data in an ideal set-up

Ideal distribution of training data



Rare Event of Interest: Consumer Safety



Negative

The text (review) is non-reportable
True > 99% of the times



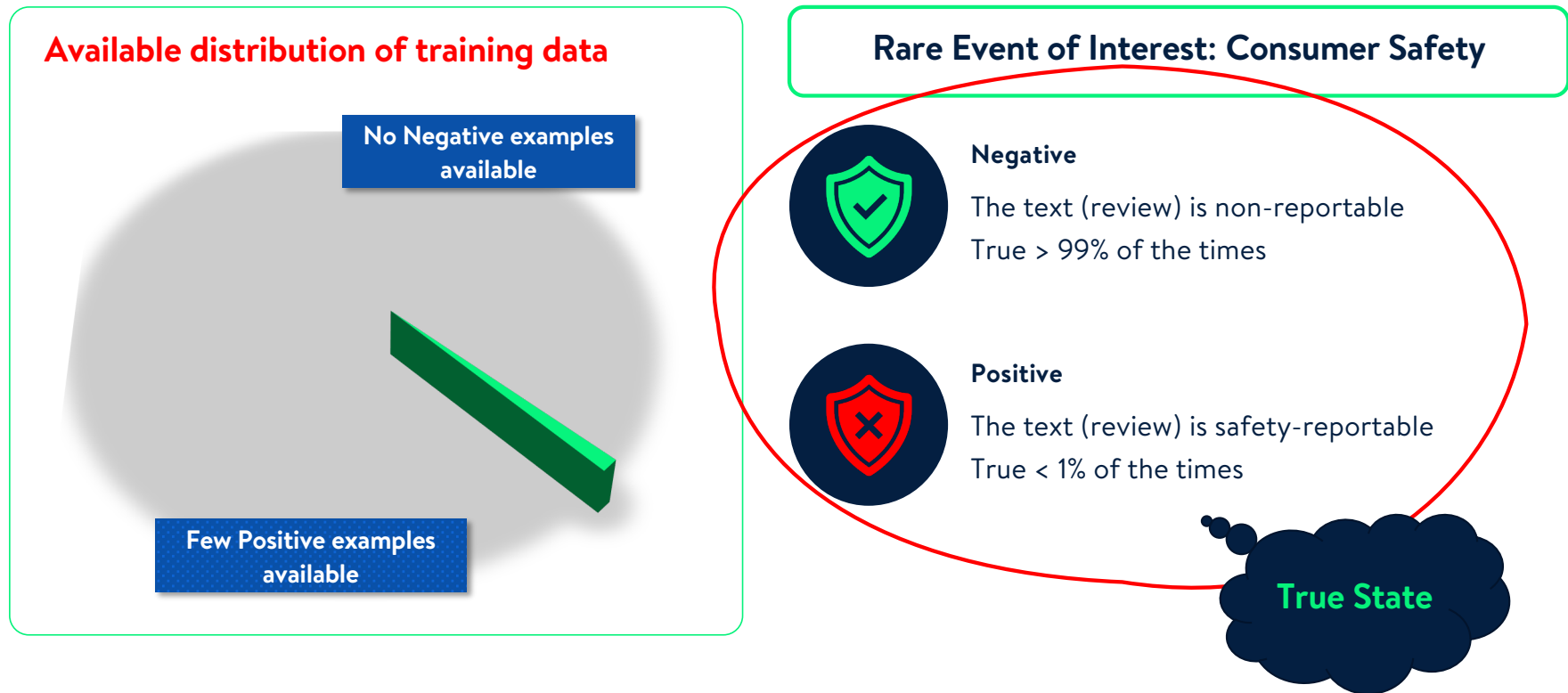
Positive

The text (review) is safety-reportable
True < 1% of the times

True State

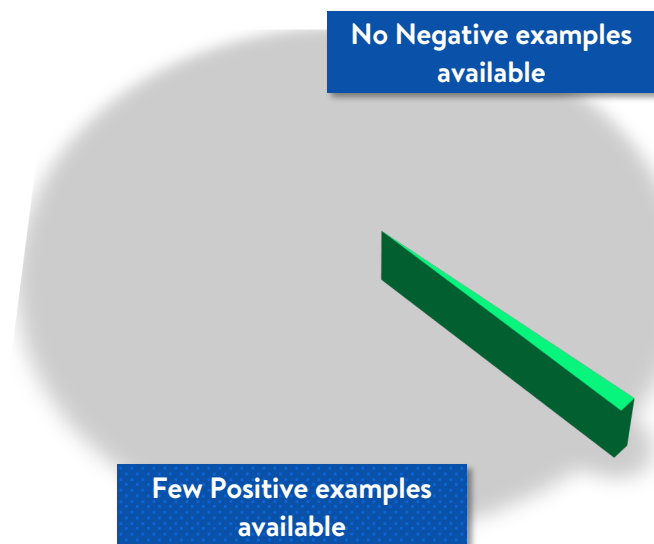


Training Data in a PU set-up (way more probable!)



Problem Formulation in an Imbalanced PU classifier

Available distribution of training data



- Binary Classification
- Two classes – Positive & Negative
- Class imbalance – Positive is a Rare Event
- Few examples available from Positive class
 - Data quality unknown
- No example available from Negative class
- Exact Rate of Imbalance unknown



Quick Break

Back in 10 mins



Semi Supervised Learners





When do we use semi-supervised learning?

- Limited Training Data
 - Usually labeled by experts
 - Hence, expensive
- Huge amount of unlabeled data
- Classes can be assumed to be distinctly separable
- Feature X is highly 'informative' about label Y



Is our use case eligible for SSL?

- **Limited Training Data** 
 - Incomplete Training Data
 - Only positive examples available
- **Huge amt of unlabeled data** 
 - Customer reviews database is huge






When do we use semi-supervised learning?

- Limited Training Data
 - Usually labeled by experts
 - Hence, expensive
- Huge amount of unlabeled data
- Classes can be assumed to be distinctly separable
- Feature X is highly 'informative' about label Y



Is our use case eligible for SSL?

- **Limited Training Data** 
 - Incomplete Training Data
 - Only positive examples available
- **Huge amt of unlabeled data** 
 - Customer reviews database is huge
- **We will assume** 
 - **Classes are distinctly separable**
 - Text is enough to identify positive cases i.e., **X is informative about Y**



Entropy Regularization in SSL

- Entropy is a measure of randomness (stability) in a random variable



Entropy Regularization in SSL

- Entropy is a measure of randomness (stability) in a random variable
- Entropy of our label Y *given* X :

$$H(Y|X) = E(-\ln P(Y|X))$$



Entropy Regularization in SSL

- Entropy is a measure of randomness (stability) in a random variable
- Entropy of our label Y *given* X :

$$H(Y|X) = E(-\ln P(Y|X))$$

Given the value of X , how stable is the value of Y



Entropy Regularization in SSL

- Entropy is a measure of randomness (stability) in a random variable

- Entropy of our label Y *given* X :

$$H(Y|X) = E(-\ln P(Y|X))$$

Given the value of X , how stable is the value of Y

Y can take only two values: **0 and 1**



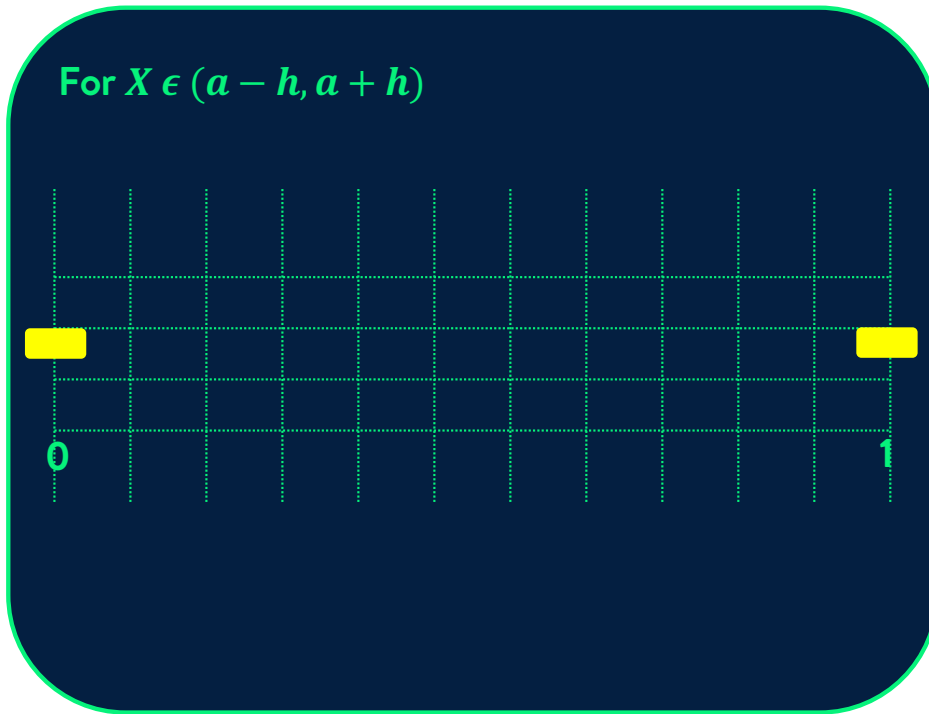
Why regularize on Entropy?

- **SSL claims X is informative about Y**
 - i.e., X alone can act as a good predictor for Y
- **High $H(Y|X)$ indicates Y can have high variance for fixed values of X**
- **SSL uses unlabeled data to make predictions on Y**
 - This makes sense only if X actually has enough information about Y
 - i.e., **$H(Y|X)$ is low**



Entropy Regularization in SSL

For $X \in (a - h, a + h)$



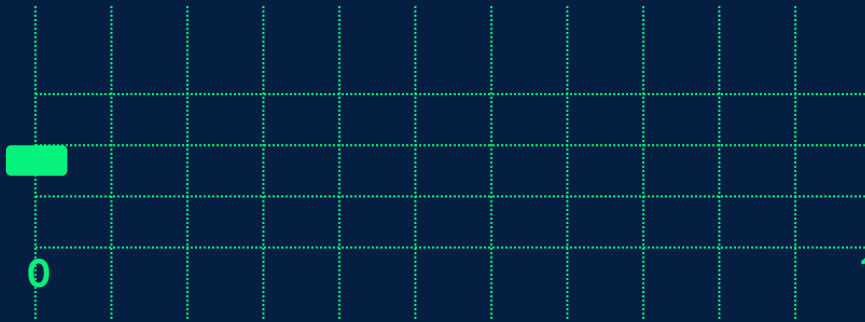
Why regularize on Entropy?

- **SSL claims X is informative about Y**
 - i.e., X alone can act as a good predictor for Y
- **High $H(Y|X)$ indicates Y can have high variance for fixed values of X**
- **SSL uses unlabeled data to make predictions on Y**
 - This makes sense only if X actually has enough information about Y
 - i.e., **$H(Y|X)$ is low**
 - **Y is stable for given ε -nbh of X**



Entropy Regularization in SSL

For $X \in (a - h, a + h)$



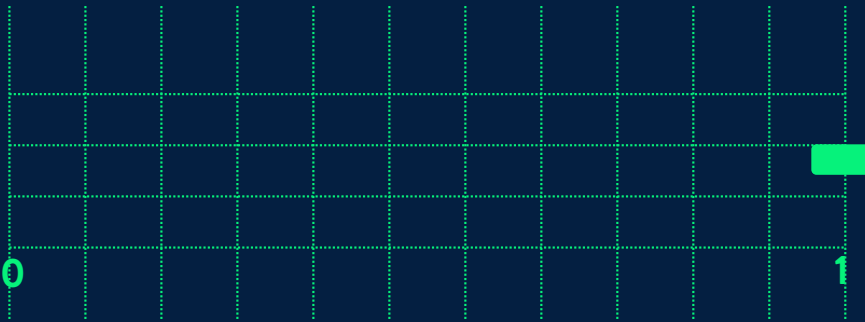
Why regularize on Entropy?

- SSL claims X is informative about Y
 - i.e., X alone can act as a good predictor for Y
- High $H(Y|X)$ indicates Y can have high variance for fixed values of X
- SSL uses unlabeled data to make predictions on Y
 - This makes sense only if X actually has enough information about Y
 - i.e., $H(Y|X)$ is low
 - Y is stable for given ε -nbh of X



Entropy Regularization in SSL

For $X \in (a - h, a + h)$



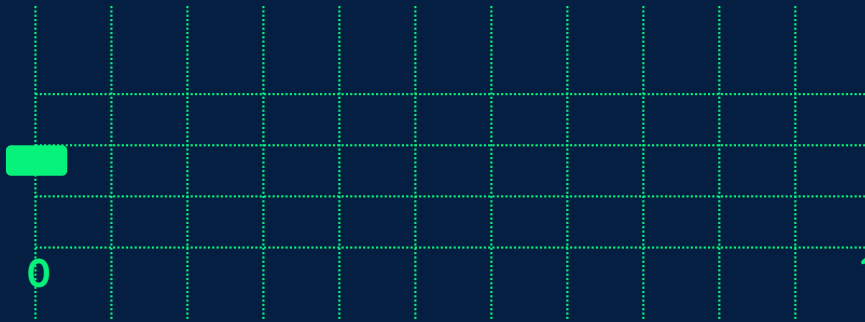
Why regularize on Entropy?

- **SSL claims X is informative about Y**
 - i.e., X alone can act as a good predictor for Y
- **High $H(Y|X)$ indicates Y can have high variance for fixed values of X**
- **SSL uses unlabeled data to make predictions on Y**
 - This makes sense only if X actually has enough information about Y
 - i.e., **$H(Y|X)$ is low**
 - **Y is stable for given ϵ -nbh of X**



Entropy Regularization in SSL

For $X \in (a - h, a + h)$



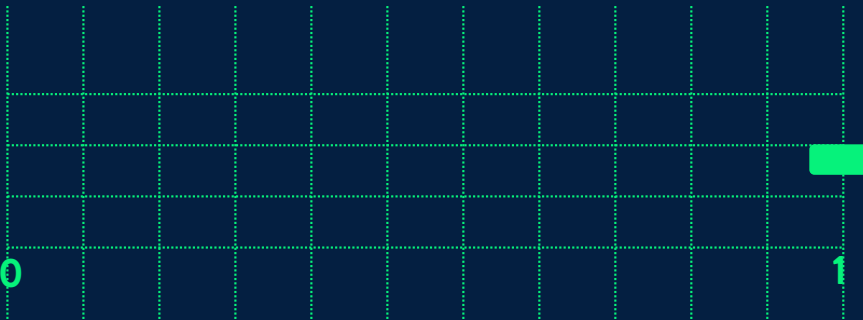
Why regularize on Entropy?

- SSL claims X is informative about Y
 - i.e., X alone can act as a good predictor for Y
- High $H(Y|X)$ indicates Y can have high variance for fixed values of X
- SSL uses unlabeled data to make predictions on Y
 - This makes sense only if X actually has enough information about Y
 - i.e., $H(Y|X)$ is low
 - Y is stable for given ϵ -nbh of X



Entropy Regularization in SSL

For $X \in (a - h, a + h)$



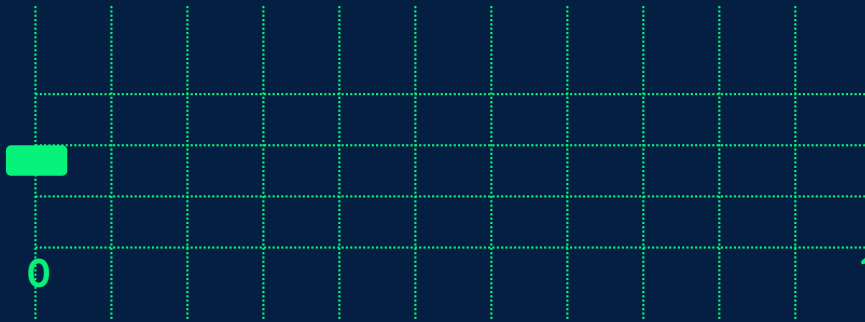
Why regularize on Entropy?

- SSL claims X is informative about Y
 - i.e., X alone can act as a good predictor for Y
- High $H(Y|X)$ indicates Y can have high variance for fixed values of X
- SSL uses unlabeled data to make predictions on Y
 - This makes sense only if X actually has enough information about Y
 - i.e., $H(Y|X)$ is low
 - Y is stable for given ϵ -nbh of X



Entropy Regularization in SSL

For $X \in (a - h, a + h)$



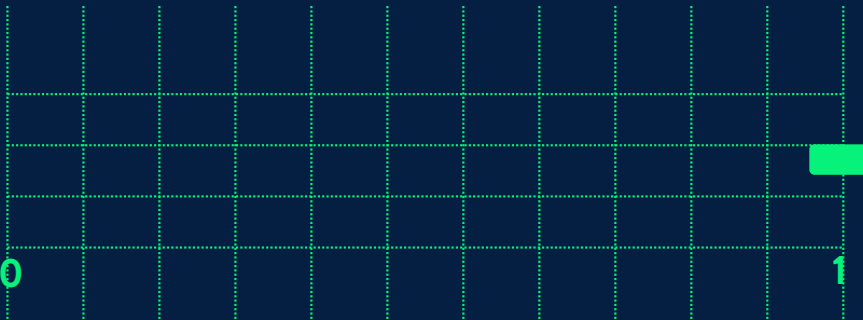
Why regularize on Entropy?

- SSL claims X is informative about Y
 - i.e., X alone can act as a good predictor for Y
- High $H(Y|X)$ indicates Y can have high variance for fixed values of X
- SSL uses unlabeled data to make predictions on Y
 - This makes sense only if X actually has enough information about Y
 - i.e., $H(Y|X)$ is low
 - Y is stable for given ϵ -nbh of X



Entropy Regularization in SSL

For $X \in (a - h, a + h)$



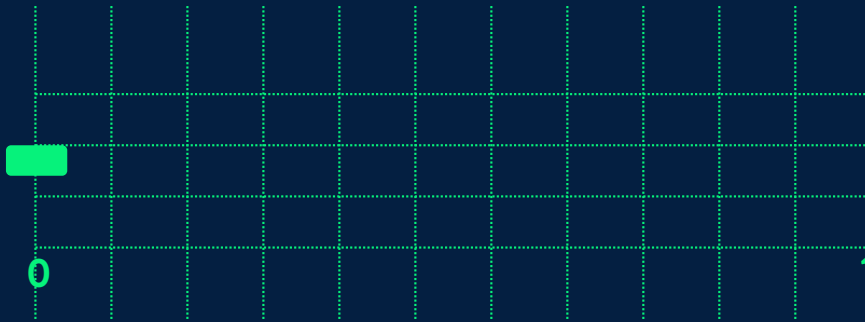
Why regularize on Entropy?

- SSL claims X is informative about Y
 - i.e., X alone can act as a good predictor for Y
- High $H(Y|X)$ indicates Y can have high variance for fixed values of X
- SSL uses unlabeled data to make predictions on Y
 - This makes sense only if X actually has enough information about Y
 - i.e., $H(Y|X)$ is low
 - Y is stable for given ϵ -nbh of X



Entropy Regularization in SSL

For $X \in (a - h, a + h)$



Why regularize on Entropy?

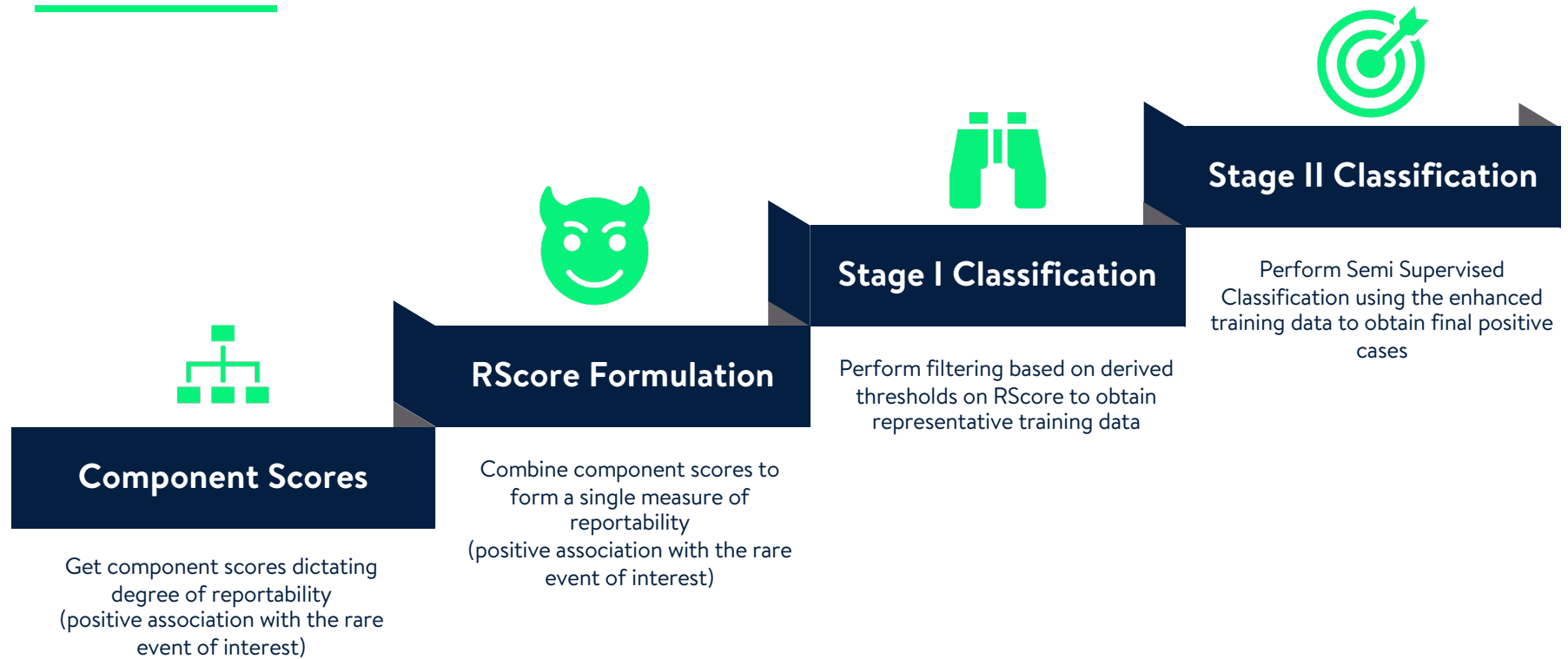
- SSL claims X is informative about Y
 - i.e., X alone can act as a good predictor for Y
- High $H(Y|X)$ indicates Y can have high variance for fixed values of X
- SSL uses unlabeled data to make predictions on Y
 - This makes sense only if X actually has enough information about Y
 - i.e., $H(Y|X)$ is low
 - Y is stable for given ε -nbh of X



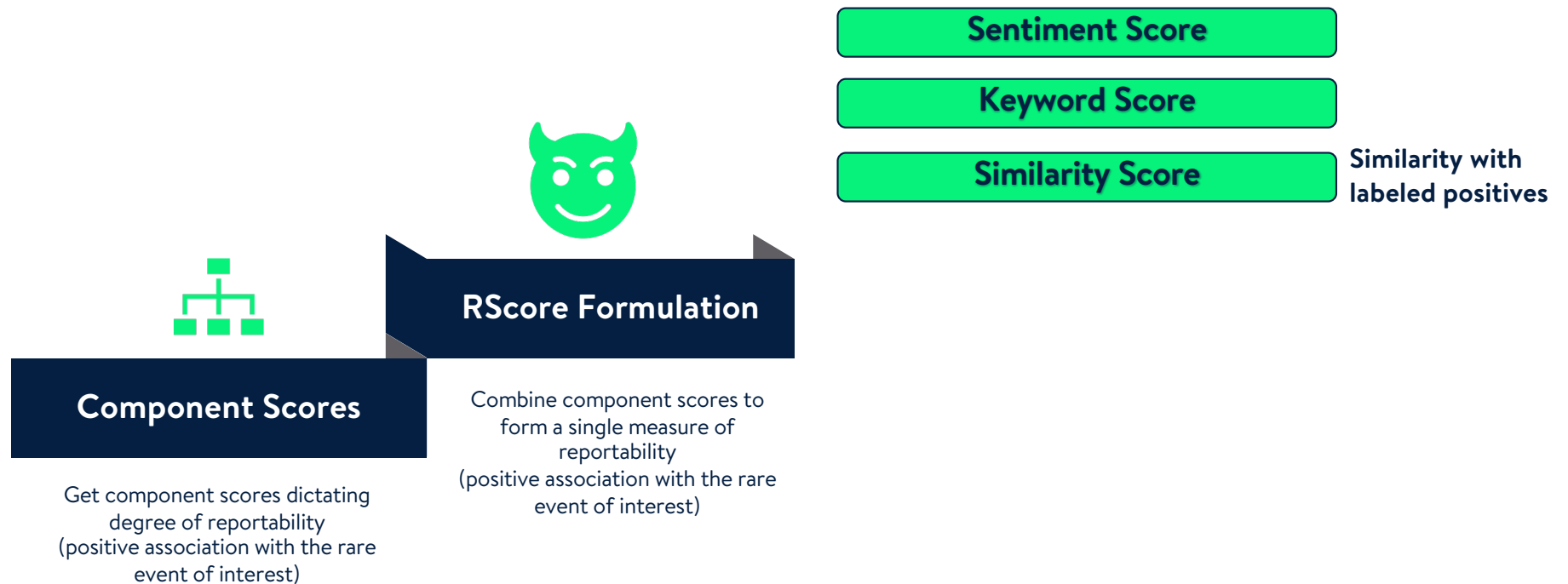
iCASSTLe



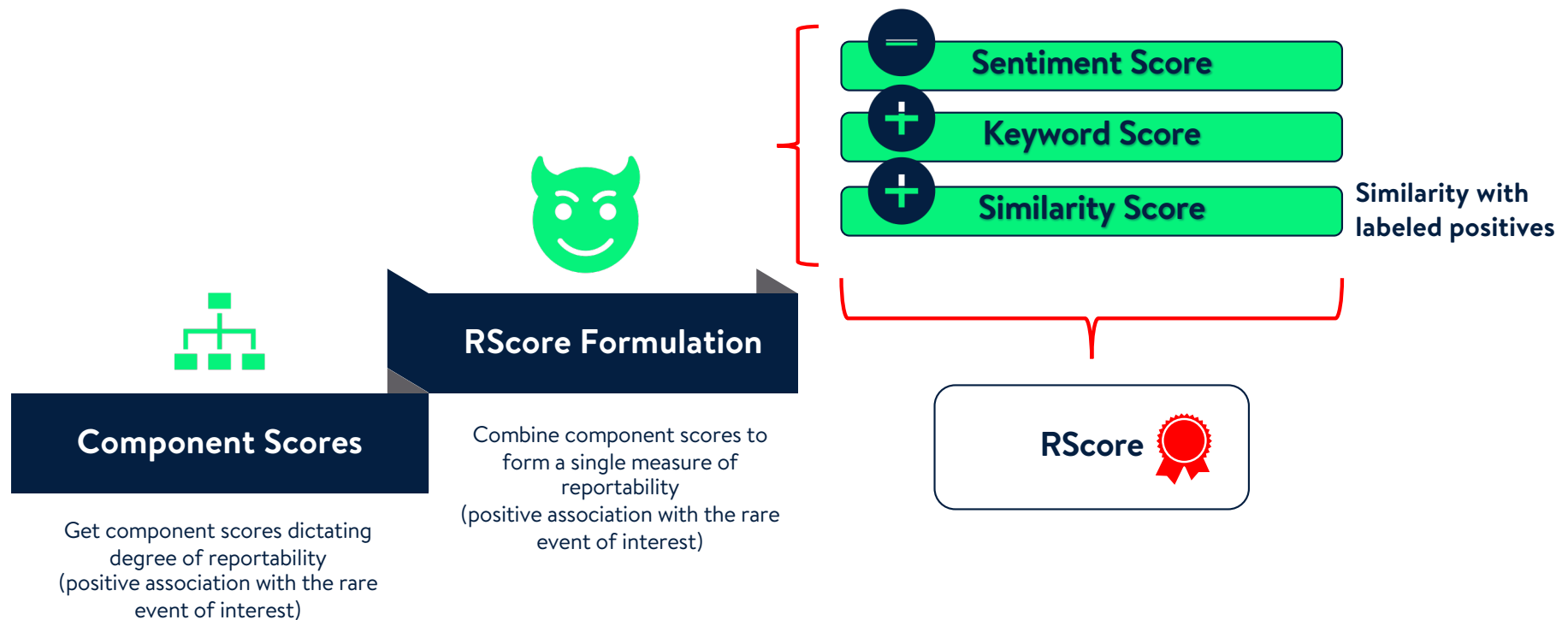
Algorithm Overview



Component Scores for Reportability & RScore Formulation



Component Scores for Reportability & RScore Formulation



Stage I Classification: Obtaining Training Negatives



Stage I Classification

Perform filtering based on derived thresholds on RScore to obtain representative training data



Global Threshold for RScore

Q_{LR}^{K1} : K_1 $_{th}$ quantile of the Rscore values for labeled positive examples



Local Threshold for RScore

Q_U^{K1} : K_2 $_{th}$ quantile of the Rscore values for unlabeled cases

Anything not classified as Stage I Positive is labeled as **Stage I Negative**

The j^{th} unlabeled case is classified as Stage I Positive iff
 $RScore_j > \min(Q_{LR}^{K1}, Q_U^{K2})$



Stage II Classification: SSL



Stage II Classification

Perform Semi Supervised Classification using the enhanced training data to obtain final positive cases



Labeled Data for Stage II (SSL)

- Original labeled Positives
- Top $K\%$ of positives + negatives obtained in **Stage I** (ranked by RScore)



Unlabeled Data for Stage II (SSL)

All original test cases except top $K\%$ of positives + negatives obtained in **Stage I** (ranked by RScore)

Quantity to Minimize: $L + \lambda H(Y|X)$

Entropy of Y
given X

Loss
Function

Regularization
Constant



Q&A

Please post on Event X Ai Platform

Channel: **wed-debanjana-banerjee-finding-rare-events-in-text**



Thank You!

debanjanabanerjee1993@gmail.com

github.com/debanjana-banerjee

linkedin.com/in/debanjana-banerjee



Appendix

