

1. Working with text
2. Working with unbalanced datasets
3. Hyperparameter tuning
4. If we have time, we will briefly discuss strategy

Feature Extraction:

The house plans to vote on the Senate's bipartisan bill but is scrapping an earlier plan to put Democrats' social policy bill first.

Convert text into a feature vector.

Documents as feature vectors

ouse leader said they

The actress is
drawing Oscar talk

The house plans to
vote on the Senate's
bipartisan bill but is
scrapping an earlier
plan to put
Democrats' social
policy bill first.

Bag of words

to the on is but
the bill on scrapping
Senate's policy to earlier
social bill vote Democrats' plans
an house bipartisan
is plan first put

| | |
|---|-----------|
| 0 | politics |
| 1 | house |
| 2 | bill |
| 0 | president |
| 0 | leader |
| 1 | social |
| ⋮ | ⋮ |

$\Phi(\mathbf{x})$

(NYTimes Nov. 5, 2021)

\mathbf{x}

Transform the feature vectors to emphasize more “relevant” words

Turning text into a feature vector

Document 1

The quick brown
fox jumped over
the lazy dog's
back.

Document 2

Now is the time
for all good men
to come to the
aid of their party.

- ❑ Document is natively text
- ❑ Must represent as a numeric vector
- ❑ Represent by word counts
 - Enumerate all words
 - Each document is count of frequencies
- ❑ Stopwords

Discussion Questions

❑ Is the absolute number of times a word appears the correct metric?

❑ What about the length of the document?

❑ What about the frequency of the word?

❑ What words “matter”?

the, for, a, in

convolutional, gradient

❑ Perhaps:

- if a word appears frequently, it is important (give it a high score)
- If a word appears in many documents, it is not important (give it a low score)

Ideas:

$$\begin{aligned} TF_{\text{"this"}, d_1} &= \frac{1}{5} = 0.2 & IDF_{\text{"this"}} &= \log\left(\frac{2}{2}\right) = 0 & TF_{\text{"example"}, d_1} &= \frac{0}{5} = 0 & IDF_{\text{"example"}} &= \log\left(\frac{2}{1}\right) = 1 \\ TF_{\text{"this"}, d_2} &= \frac{1}{7} \approx 0.14 & & & TF_{\text{"example"}, d_2} &= \frac{3}{7} \approx 0.429 & & \end{aligned}$$

Example modified from <https://en.wikipedia.org/wiki/Tf%E2%80%93idf>

❑ How can we categorize how important a word is in a document?

❑ Perhaps:

- if a word appears frequently, it is important (give it a high score) convolutional, gradient
- except if the word appears in many documents, it is not important (give it a low score)

❑ Steps:

- Count the frequency of every word in the document

Term frequency

$$TF_{i,n} = \frac{\text{num times word } i \text{ in doc } n}{\text{total num words in doc } n}$$

- Determine how much information a word provides: Inverse Document Frequency (IDF)

The more common a word is
the lower its IDF score

Inverse doc frequency

$$IDF_i = \log \left[\frac{\text{Total num docs in corpus}}{\text{Num docs with word } i} \right]$$

Document 1

| Term | Term Count |
|--------|------------|
| this | 1 |
| Is | 1 |
| a | 2 |
| sample | 1 |

Document 2

| Term | Term Count |
|---------|------------|
| this | 1 |
| Is | 1 |
| another | 2 |
| example | 3 |

Ideas:

$$\begin{aligned}
 TF_{\text{"this"}, d_1} &= \frac{1}{5} = 0.2 & IDF_{\text{"this"}} &= \log\left(\frac{2}{2}\right) = 0 & TF_{\text{"example"}, d_1} &= \frac{0}{5} = 0 & IDF_{\text{"example"}} &= \log\left(\frac{2}{1}\right) = 1 \\
 TF_{\text{"this"}, d_2} &= \frac{1}{7} \approx 0.14 & & & TF_{\text{"example"}, d_2} &= \frac{3}{7} \approx 0.429 & &
 \end{aligned}$$

Example modified from <https://en.wikipedia.org/wiki/Tf%E2%80%93idf>

Frequency is relative to the size of the document

How can we categorize how important a word is in a document?

Perhaps:

- if a word appears frequently, it is important (give it a high score) *convolutional, gradient*
- except if the word appears in many documents, it is not important (give it a low score)

Steps:

- Count the frequency of every word in the document

Term frequency

$$TF_{i,n} = \frac{\text{num times word } i \text{ in doc } n}{\text{total num words in doc } n}$$

- Determine how much information a word provides: Inverse Document Frequency (IDF)

The more common a word is the lower its IDF score

Inverse doc frequency

$$IDF_i = \log \left[\frac{\text{Total num docs in corpus}}{\text{Num docs with word } i} \right]$$

Document 1

| Term | Term Count |
|--------|------------|
| this | 1 |
| Is | 1 |
| a | 2 |
| sample | 1 |

Document 2

| Term | Term Count |
|---------|------------|
| this | 1 |
| Is | 1 |
| another | 2 |
| example | 3 |

Term Frequency - Inverse Document Frequency

Use TF-IDF weight for vectors:

$$X[n, i] = \text{Document weight vector} \times \text{Term frequency} \times \text{Inverse doc frequency}$$

$$TF_{i,n} = \frac{\text{num times word } i \text{ in doc } n}{\text{total num words in doc } n}$$

$$IDF_i = \log \left[\frac{\text{Total num docs in corpus}}{\text{Num docs with word } i} \right]$$

$$\begin{aligned} TF_{\text{"this"}, d_1} &= \frac{1}{5} = 0.2 & IDF_{\text{"this"}} &= \log \left(\frac{2}{2} \right) = 0 & TF_{\text{"example"}, d_1} &= \frac{0}{5} = 0 & IDF_{\text{"example"}} &= \log \left(\frac{2}{1} \right) = 1 \\ TF_{\text{"this"}, d_2} &= \frac{1}{7} \approx 0.14 & & & TF_{\text{"example"}, d_2} &= \frac{3}{7} \approx 0.429 & & \end{aligned}$$

$$\begin{aligned} TF-IDF_{\text{"this"}, d_1} &= 0.2 \times 0 = 0 \\ TF-IDF_{\text{"this"}, d_2} &= 0.14 \times 0 = 0 \end{aligned}$$

$$\begin{aligned} TF-IDF_{\text{"example"}, d_1, D} &= 0. \times 1 = 0 \\ TF-IDF_{\text{"example"}, d_2} &= 0.429 \times 1 = 0.429 \end{aligned}$$

Term Frequency - Inverse Document Frequency

Document 1

| Term | Term Count |
|--------|------------|
| this | 1 |
| is | 1 |
| a | 2 |
| sample | 1 |

Document 2

| Term | Term Count |
|---------|------------|
| this | 1 |
| is | 1 |
| another | 2 |
| example | 3 |

| Term | TF-IDF Example 1 | TF-IDF Example 2 |
|---------|------------------|------------------|
| this | 0 | 0 |
| is | 0 | 0 |
| a | 0.4 | 0 |
| sample | 0.2 | 0 |
| example | 0 | 0.429 |
| another | 0 | 0.286 |

$$\text{TF-IDF}_{\text{"is"}, d_1} = 0.2 \times 0 = 0$$

$$\text{TF-IDF}_{\text{"is"}, d_2} = 0.143 \times 0 = 0$$

$$\text{TF-IDF}_{\text{"a"}, d_1} = 0.4 \times 1 = 0.4$$

$$\text{TF-IDF}_{\text{"a"}, d_2} = 0 \times 1 = 0$$

$$\text{TF-IDF}_{\text{"sample"}, d_1} = 0.2 \times 1 = 0.2$$

$$\text{TF-IDF}_{\text{"sample"}, d_2} = 0 \times 1 = 0$$

$$\text{TF-IDF}_{\text{"another"}, d_1} = 0.286 \times 1 = 0.286$$

$$\text{TF-IDF}_{\text{"another"}, d_2} = 0 \times 1 = 0$$

This is a very small example. Obviously we would usually compute the TF-IDF for a large collection of documents

Learn more at:

https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.TfidfVectorizer.html

1. Working with text



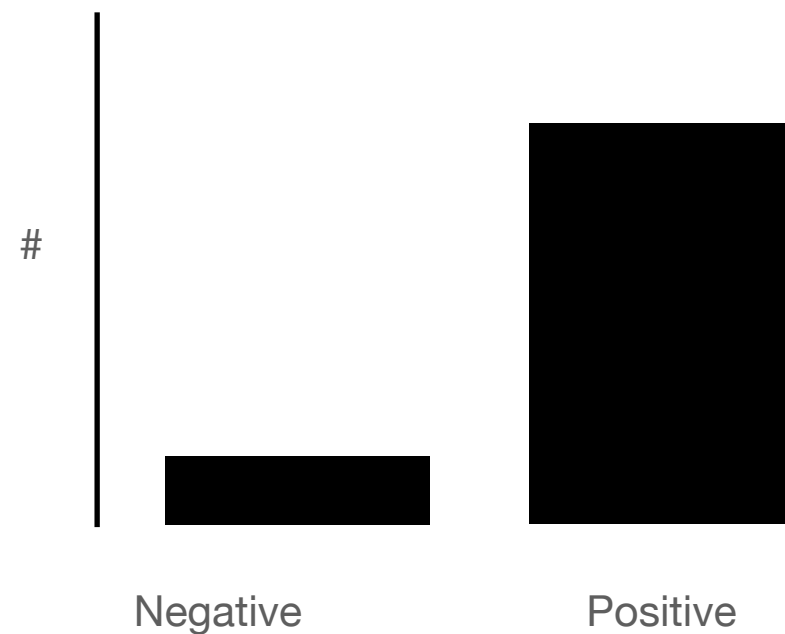
2. Working with unbalanced datasets

3. Hyperparameter tuning

4. If we have time, we will briefly discuss strategy

Imbalanced Data

Asymmetric data - Potentially not do well on the minority class



Mild: 20-40%
Moderate: 1-20%
Extreme: <1%

Imbalanced Data

Basic Approaches

- Do nothing. See if your hypothesis works.
- Under-sampling: Reduce the number of the majority class
- Over-sampling: Increase the number in the minority class
- Weighted Learners. Create a weighted objective function.

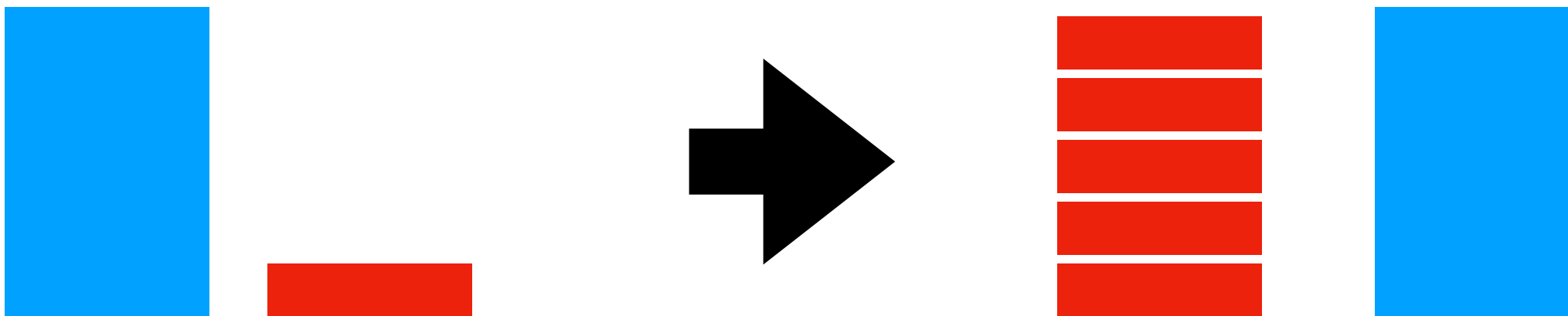


Oversampling

Increase the number of examples in the minority class

- Duplicate minority class examples
- Generate new synthetic data
 - e.g. SMOTE, ADASYN

Advantage: use all the original data

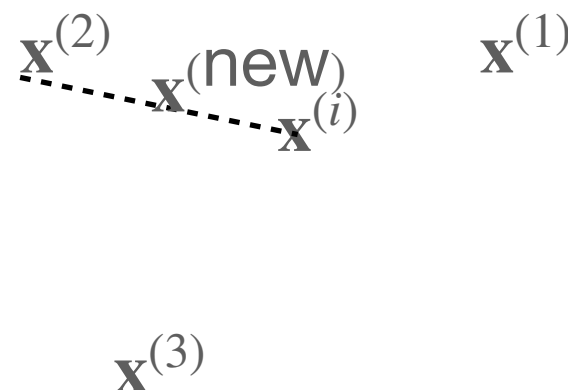


SMOTE

Create synthetic minority class training examples. There are many variations on the following basic algorithm (including ADASYN).

Creates new examples by:

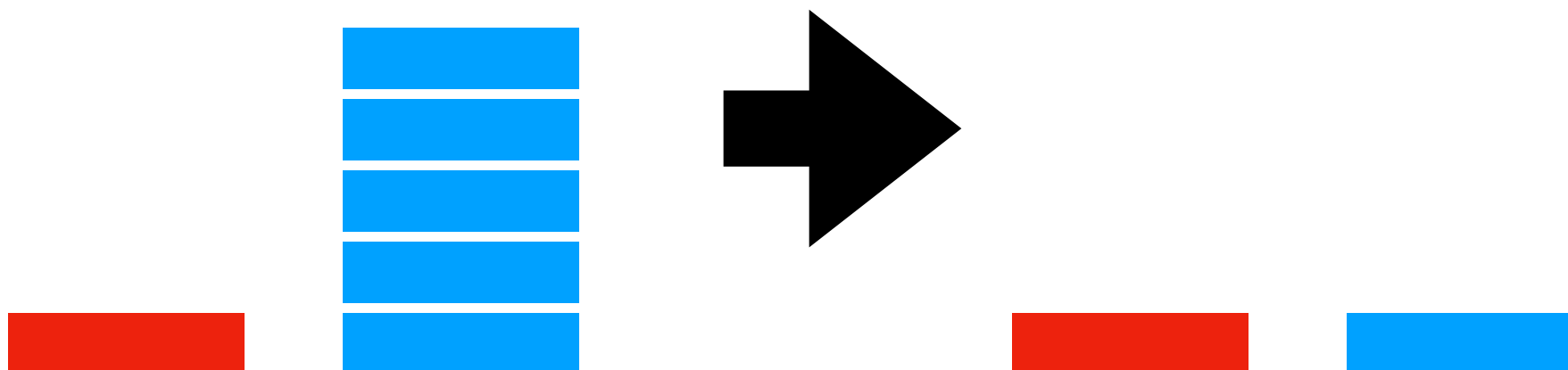
- Choosing a random example $\mathbf{x}^{(i)}$ from the minority class and finding the k closest minority training examples to it
- Randomly selects $\mathbf{x}^{(j)}$ one of the k closest minority examples and creates a new example: $\mathbf{x}^{(\text{new})} = \mathbf{x}^{(i)} + \lambda(\mathbf{x}^{(j)} - \mathbf{x}^{(i)})$ where λ is randomly chosen in the range $[0,1]$



Undersampling

If there are m minority-class elements, use only m of the majority-class examples. Some ways to select the majority-class examples

- Randomly
- Selectively. The goal is to remove borderline examples and redundant noisy examples. e.g., Tomek links and condensed nearest neighbor



Weighted Learners

Popular technique!!

- Incorporate different penalties for different classes in objective function
- Higher weight for the minority class, smaller weight for the majority class

Logistic regression

$$\ell(\mathbf{w}) = \sum_{i=1}^N \left[\text{weight}_1 y^{(i)} \ln \sigma(\mathbf{w}^T \mathbf{x}^{(i)}) + \text{weight}_0 (1 - y^{(i)}) \ln (1 - \sigma(\mathbf{w}^T \mathbf{x}^{(i)})) \right]$$

Support vector machine

$$\min_{w_0, \mathbf{w}, \{\xi^{(i)}\}_{i=1}^N} \|\mathbf{w}\|_2^2 + \text{weight}_1 C \sum_{j: y^{(j)} = +1} \xi^{(j)} + \text{weight}_0 C \sum_{j: y^{(j)} = -1} \xi^{(j)}$$

subject to $y^{(i)} (w_0 + \mathbf{w}^T \mathbf{x}^{(i)}) \geq 1 - \xi^{(i)}$ for all $i = 1, \dots, N$

$$\xi^{(i)} \geq 0$$

Learn more at:

[https://chrisalbon.com/code/machine_learning/logistic_regression/
handling_imbalanced_classes_in_logistic_regression/](https://chrisalbon.com/code/machine_learning/logistic_regression/handling_imbalanced_classes_in_logistic_regression/)

[https://chrisalbon.com/code/machine_learning/
support_vector_machines/imbalanced_classes_in_svm/](https://chrisalbon.com/code/machine_learning/support_vector_machines/imbalanced_classes_in_svm/)