## 1. Loading the NIPS papers

The NIPS conference (Neural Information Processing Systems) is one of the most prestigious yearly events in the machine learning community. At each NIPS conference, a large number of research papers are published. Over 50,000 PDF files were automatically

of NIPS (Neural Information Processing Systems) First, we will explore the CSV file to determine what type of data we can use for the analysis and how it is structured. A research paper typically consists of a title, an abstract and the main text. Other data such as figures and tables were not extracted from the PDF files. Each paper discusses a novel technique or improvement. In this analysis, we will focus on analyzing these papers with natural language

downloaded and processed to obtain a dataset on various machine learning techniques. These NIPS papers are stored in datasets/papers.csv. The CSV file contains information on the different NIPS papers that were published from 1987 until 2017 (30 years!). These papers discuss a wide variety of topics in machine learning, from neural networks to optimization methods and many more. The logo

warnings.filterwarnings('ignore')

# Read datasets/papers.csv into papers papers = pd.read csv("papers.csv")

# Print out the first rows of papers

Self-Organization of Associative

A Mean Field Theory of Layer IV of

Bayesian Query Construction for

Neural Network Ensembles, Cross

Storing Covariance by the

Associative Long-Ter...

Neural Network...

Database and ...

Visual Cort...

import warnings

papers.head()

id year

1 1987

10 1987

100 1988

1000 1994

# Importing modules import pandas as pd

In [20]:

Out[20]:

In [22]:

1

processing methods.

pdf\_name

database-an...

iv-of-visual-c...

associative-long...

cross-validation...

1000-bayesian-query-

1-self-organization-of-associative-

10-a-mean-field-theory-of-layer-

100-storing-covariance-by-the-

construction-for-neural-ne...

1001-neural-network-ensembles-

We will analyze this text data using natural language processing. Since the file contains some metadata such as id's and filenames, it is

abstract

By looking at the number of published papers per year, we can understand the extent of the machine learning 'revolution'! Typically, this

Let's now analyze the titles of the different papers to identify machine learning trends. First, we will perform some simple preprocessing on the titles in order to make them more amenable for analysis. We will use a regular expression to remove any punctuation in the title. Then

we will perform lowercasing. We'll then print the titles of the first rows before and after applying the modification.

papers['title processed'] = papers['title'].map(lambda x: re.sub('[,\.!?]', '', x))

papers['title processed'] = papers['title processed'].map(lambda x: x.lower())

significant increase in popularity is attributed to the large amounts of compute power, data and improvements in algorithms.

abstract

Abstract

Missing

Abstract

Missing

Abstract

Missing

Abstract

Missing

Abstract

Missing

683\n\nA MEAN FIELD THEORY OF LAYER IV OF VISU...

Bayesian Query Construction for Neural\nNetwor...

Neural Network Ensembles, Cross\nValidation, a...

394\n\nSTORING COVARIANCE BY THE ASSOCIATIVE\n...

paper\_text

767\n\nSELF-ORGANIZATION OF

683\n\nA MEAN FIELD THEORY OF

Bayesian Query Construction for

Neural Network Ensembles,

394\n\nSTORING COVARIANCE BY THE

paper\_text

ASSOCIATIVE DATABA...

LAYER IV OF VISU...

ASSOCIATIVE\n...

Neural\nNetwor...

Cross\nValidation, a...

title event\_type

NaN

NaN

NaN

NaN

1001 1994 NaN Validation, an... 2. Preparing the data for analysis For the analysis of the papers, we are only interested in the text data associated with the paper as well as the year the paper was published

necessary to remove all the columns that do not contain useful text information. In [21]: # Remove the columns papers.drop(columns=['id', "event type", "pdf name"], inplace=True) # Print out the first rows of papers papers.head()

Out[21]: year title Self-Organization of Associative Database and ... Abstract Missing 767\n\nSELF-ORGANIZATION OF ASSOCIATIVE DATABA... **0** 1987 **1** 1987 A Mean Field Theory of Layer IV of Visual Cort... Abstract Missing Storing Covariance by the Associative Long-Ter... Abstract Missing **2** 1988

1994 Bayesian Query Construction for Neural Network... Abstract Missing

Neural Network Ensembles, Cross Validation, an... Abstract Missing **4** 1994 3. Plotting how machine learning has evolved over time In order to understand how the machine learning field has recently exploded in popularity, we will begin by visualizing the number of publications per year.

# Group the papers by year groups = papers.groupby("year")

counts = groups.size()

%matplotlib inline counts.plot(kind='bar')

plt.xlabel("Year");

700

600

500

400

300

In [23]:

In [24]:

import re

# Determine the size of each group

# Visualise the counts as a bar plot import matplotlib.pyplot as plt

plt.title("Machine Learning Publications since 1987")

Machine Learning Publications since 1987

200

4. Preprocessing the text data

# Load the regular expression library

# Print the titles of the first rows

# Convert the titles to lowercase

print(papers['title'].head())

# Remove punctuation

## Self-Organization of Associative Database and ... A Mean Field Theory of Layer IV of Visual Cort... Storing Covariance by the Associative Long-Ter... Bayesian Query Construction for Neural Network...

Name: title processed, dtype: object

# Import the wordcloud library

# Create a WordCloud object

# Generate a word cloud

wordcloud.to image()

# Visualize the word cloud

wordcloud = wordcloud.WordCloud()

wordcloud.generate(long string)

# Join the different processed titles together. long string = " ".join(papers.title processed)

import wordcloud

Name: title, dtype: object

print(papers['title processed'].head())

# Print the processed titles of the first rows

Neural Network Ensembles, Cross Validation, an...

0 self-organization of associative database and ... a mean field theory of layer iv of visual cort... storing covariance by the associative long-ter... bayesian query construction for neural network... neural network ensembles cross validation and ...

Python has a massive number of open libraries! Instead of trying to develop a method to create word clouds ourselves, we'll use Andreas Mueller's wordcloud library.

5. A word cloud to visualize the preprocessed text data

addition, it allows us to verify whether we need additional preprocessing before further analyzing the text data.

In order to verify whether the preprocessing happened correctly, we can make a word cloud of the titles of the research papers. This will give us a visual representation of the most common words. Visualisation is key to understanding whether we are still on the right track! In

The main text analysis method that we will use is latent Dirichlet allocation (LDA). LDA is able to perform topic detection on large

LDA does not work directly on text data. First, it is necessary to convert the documents into a simple vector representation. This

vocabulary. For example, 'Analyzing machine learning trends with neural networks.' would be transformed into [1, 0, 1, ..., 1, 0].

document sets, determining what the main 'topics' are in a large unlabeled set of texts. A 'topic' is a collection of words that tend to cooccur often. The hypothesis is that LDA might be able to clarify what the different topics in the research titles are. These topics can then be

representation will then be used by LDA to determine the topics. Each entry of a 'document vector' will correspond with the number of times a word occurred in the document. In conclusion, we will convert a list of titles into a list of vectors, all with length equal to the

We'll then plot the 10 most common words based on the outcome of this operation (the list of document vectors). As a check, these words

Finally, the research titles will be analyzed using LDA. Note that in order to process a new set of documents (e.g. news articles), a similar set of steps will be required to preprocess the data. The flow that was constructed here can thus easily be exported for a new text dataset.

The only parameter we will tweak is the number of topics in the LDA algorithm. Typically, one would calculate the 'perplexity' metric to determine which number of topics is best and iterate over different amounts of topics until the lowest 'perplexity' is found. For now, let's play around with a different number of topics. From there, we can distinguish what each topic is about ('neural networks', 'reinforcement

Out[24]:

In [25]: # Load the library with the CountVectorizer method from sklearn.feature extraction.text import CountVectorizer import numpy as np # Helper function

import matplotlib.pyplot as plt

total counts = np.zeros(len(words))

total counts+=t.toarray()[0]

words = [w[0] for w in count dict] counts = [w[1] for w in count dict]

plt.bar(x pos, counts,align='center') plt.xticks(x pos, words, rotation=90)

plt.title('10 most common words')

# Fit and transform the processed titles

# Visualise the 10 most common words

x pos = np.arange(len(words))

count dict = (zip(words, total counts))

def plot 10 most common words(count data, count vectorizer):

# Initialise the count vectorizer with the English stop words

count data = count vectorizer.fit transform(papers['title processed'])

count vectorizer = CountVectorizer(stop words='english')

plot 10 most common words (count data, count vectorizer)

10 most common words

count dict = sorted(count dict, key=lambda x:x[1], reverse=True)[0:10]

words = count vectorizer.get feature names()

6. Prepare the text for LDA analysis

used as a starting point for further analysis.

should also occur in the word cloud.

for t in count data:

plt.xlabel('words') plt.ylabel('counts')

plt.show()

1600 1400

1200 1000

> 800 600

200

import warnings

# Helper function

number topics = 10 number words = 6

Topic #0:

Topic #1:

Topic #2:

Topic #4:

Topic #8:

In [26]:

7. Analysing trends with LDA

learning', 'kernel methods', 'gaussian processes', etc.).

warnings.simplefilter("ignore", DeprecationWarning)

def print topics(model, count vectorizer, n top words): words = count vectorizer.get feature names()

# Tweak the two parameters below (use int values below 15)

from sklearn.decomposition import LatentDirichletAllocation as LDA

networks

models

for topic\_idx, topic in enumerate(model.components\_): print("\nTopic #%d:" % topic idx) print(" ".join([words[i] for i in topic.argsort()[:-n top words - 1:-1]]))

# Load the LDA model from sk-learn

# Print the topics found by the LDA model print("Topics found via LDA:") Topics found via LDA:

networks neural learning using recurrent image

learning reinforcement bayesian multiple time using

visual learning clustering unsupervised deep information

learning classification large adaptive matrix sampling

models optimization stochastic markov gradient random

inference learning neural networks bayesian linear

learning state optimal representations sequential adaptation

Topic #5: sparse multi high learning using data Topic #6: recognition neural model vector speech support

lda.fit(count data) print topics(lda, count vectorizer, number words)

lda = LDA(n components=number topics)

# Create and fit the LDA model

Topic #7: learning self vlsi carlo monte analog

8. The future of machine learning Machine learning has become increasingly popular over the past years. The number of NIPS conference papers has risen exponentially, and

# The historical data indicates that:

people are continuously looking for ways on how they can incorporate machine learning into their products and services.

learning by going through blogs, tutorials, and courses.

Although this analysis focused on analyzing machine learning trends in research, a lot of these techniques are rapidly being adopted in industry. Following the latest machine learning trends is a critical skill for a data scientist, and it is recommended to continuously keep

In [27]: more papers published in 2018 = True