

6.2. Feature extraction

The [`sklearn.feature_extraction`](#) module can be used to extract features in a format supported by machine learning algorithms from datasets consisting of formats such as text and image.

Note: Feature extraction is very different from [Feature selection](#): the former consists in transforming arbitrary data, such as text or images, into numerical features usable for machine learning. The latter is a machine learning technique applied on these features.

6.2.1. Loading features from dicts

The class [`DictVectorizer`](#) can be used to convert feature arrays represented as lists of standard Python `dict` objects to the NumPy/SciPy representation used by scikit-learn estimators.

While not particularly fast to process, Python’s `dict` has the advantages of being convenient to use, being sparse (absent features need not be stored) and storing feature names in addition to values.

[`DictVectorizer`](#) implements what is called one-of-K or “one-hot” coding for categorical (aka nominal, discrete) features. Categorical features are “attribute-value” pairs where the value is restricted to a list of discrete possibilities without ordering (e.g. topic identifiers, types of objects, tags, names...).

In the following, “city” is a categorical attribute while “temperature” is a traditional numerical feature:

```
>>> measurements = [
...     {'city': 'Dubai', 'temperature': 33.},
...     {'city': 'London', 'temperature': 12.},
...     {'city': 'San Francisco', 'temperature': 18.},
... ]

>>> from sklearn.feature_extraction import DictVectorizer
>>> vec = DictVectorizer()

>>> vec.fit_transform(measurements).toarray()
array([[ 1.,  0.,  0., 33.],
       [ 0.,  1.,  0., 12.],
       [ 0.,  0.,  1., 18.]])

>>> vec.get_feature_names_out()
array(['city=Dubai', 'city=London', 'city=San Francisco', 'temperature'], ...)
```

[`DictVectorizer`](#) accepts multiple string values for one feature, like, e.g., multiple categories for a movie.

Assume a database classifies each movie using some categories (not mandatories) and its year of release.

```
>>> movie_entry = [{'category': ['thriller', 'drama'], 'year': 2003},
...                 {'category': ['animation', 'family'], 'year': 2011},
...                 {'year': 1974}]
>>> vec.fit_transform(movie_entry).toarray()
array([[0.000e+00, 1.000e+00, 0.000e+00, 1.000e+00, 2.003e+03],
       [1.000e+00, 0.000e+00, 1.000e+00, 0.000e+00, 2.011e+03],
       [0.000e+00, 0.000e+00, 0.000e+00, 0.000e+00, 1.974e+03]])
>>> vec.get_feature_names_out()
array(['category=animation', 'category=drama', 'category=family',
       'category=thriller', 'year'], ...)
>>> vec.transform({'category': ['thriller'],
...                 'unseen_feature': '3'}).toarray()
array([[0., 0., 0., 1., 0.]])
```

[`DictVectorizer`](#) is also a useful representation transformation for training sequence classifiers in Natural Language Processing models that typically work by extracting feature windows around a particular word of interest.

For example, suppose that we have a first algorithm that extracts Part of Speech (PoS) tags that we want to use as complementary tags for training a sequence classifier (e.g. a chunker). The following dict could be such a window of features extracted around the word ‘sat’ in the sentence ‘The cat sat on the mat.’:

```
>>> pos_window = [
...     {
...         'word-2': 'the',
...         'pos-2': 'DT',
...         'word-1': 'cat',
...         'pos-1': 'NN',
...         'word+1': 'on',
...         'pos+1': 'PP',
...     },
...     # in a real application one would extract many such dictionaries
... ]
```

This description can be vectorized into a sparse two-dimensional matrix suitable for feeding into a classifier (maybe after being piped into a [TfidfTransformer](#) for normalization):

```
>>> vec = DictVectorizer()
>>> pos_vectorized = vec.fit_transform(pos_window)
>>> pos_vectorized
<1x6 sparse matrix of type '<... 'numpy.float64'>'
  with 6 stored elements in Compressed Sparse ... format>
>>> pos_vectorized.toarray()
array([[1., 1., 1., 1., 1., 1.]])
>>> vec.get_feature_names_out()
array(['pos+1=PP', 'pos-1=NN', 'pos-2=DT', 'word+1=on', 'word-1=cat',
       'word-2=the'], ...)
```

As you can imagine, if one extracts such a context around each individual word of a corpus of documents the resulting matrix will be very wide (many one-hot-features) with most of them being valued to zero most of the time. So as to make the resulting data structure able to fit in memory the `DictVectorizer` class uses a `scipy.sparse` matrix by default instead of a `numpy.ndarray`.

6.2.2. Feature hashing

The class [FeatureHasher](#) is a high-speed, low-memory vectorizer that uses a technique known as [feature hashing](#), or the “hashing trick”. Instead of building a hash table of the features encountered in training, as the vectorizers do, instances of [FeatureHasher](#) apply a hash function to the features to determine their column index in sample matrices directly. The result is increased speed and reduced memory usage, at the expense of inspectability; the hasher does not remember what the input features looked like and has no `inverse_transform` method.

Since the hash function might cause collisions between (unrelated) features, a signed hash function is used and the sign of the hash value determines the sign of the value stored in the output matrix for a feature. This way, collisions are likely to cancel out rather than accumulate error, and the expected mean of any output feature’s value is zero. This mechanism is enabled by default with `alternate_sign=True` and is particularly useful for small hash table sizes (`n_features < 10000`). For large hash table sizes, it can be disabled, to allow the output to be passed to estimators like [MultinomialNB](#) or [chi2](#) feature selectors that expect non-negative inputs.

[FeatureHasher](#) accepts either mappings (like Python’s `dict` and its variants in the `collections` module), `(feature, value)` pairs, or strings, depending on the constructor parameter `input_type`. Mapping are treated as lists of `(feature, value)` pairs, while single strings have an implicit value of 1, so `['feat1', 'feat2', 'feat3']` is interpreted as `[('feat1', 1), ('feat2', 1), ('feat3', 1)]`. If a single feature occurs multiple times in a sample, the associated values will be summed (so `('feat', 2)` and `('feat', 3.5)` become `('feat', 5.5)`). The output from [FeatureHasher](#) is always a `scipy.sparse` matrix in the CSR format.

Feature hashing can be employed in document classification, but unlike [CountVectorizer](#), [FeatureHasher](#) does not do word splitting or any other preprocessing except Unicode-to-UTF-8 encoding; see [Vectorizing a large text corpus with the hashing trick](#), below, for a combined tokenizer/hasher.

As an example, consider a word-level natural language processing task that needs features extracted from `(token, part_of_speech)` pairs. One could use a Python generator function to extract features:

```
def token_features(token, part_of_speech):
    if token.isdigit():
        yield "numeric"
    else:
        yield "token={}".format(token.lower())
        yield "token,pos={},{}".format(token, part_of_speech)
    if token[0].isupper():
        yield "uppercase_initial"
    if token.isupper():
        yield "all_uppercase"
    yield "pos={}".format(part_of_speech)
```

Then, the `raw_x` to be fed to `FeatureHasher.transform` can be constructed using:

```
token_features(tok, pos_tagger(tok)) for tok in corpus)
```

Toggle Menu

and fed to a hasher with:

```
hasher = FeatureHasher(input_type='string')
X = hasher.transform(raw_X)
```

to get a `scipy.sparse` matrix `X`.

Note the use of a generator comprehension, which introduces laziness into the feature extraction: tokens are only processed on demand from the hasher.

6.2.2.1. Implementation details

[FeatureHasher](#) uses the signed 32-bit variant of MurmurHash3. As a result (and because of limitations in `scipy.sparse`), the maximum number of features supported is currently $2^{31} - 1$.

The original formulation of the hashing trick by Weinberger et al. used two separate hash functions h and ξ to determine the column index and sign of a feature, respectively. The present implementation works under the assumption that the sign bit of MurmurHash3 is independent of its other bits.

Since a simple modulo is used to transform the hash function to a column index, it is advisable to use a power of two as the `n_features` parameter; otherwise the features will not be mapped evenly to the columns.

References:

- Kilian Weinberger, Anirban Dasgupta, John Langford, Alex Smola and Josh Attenberg (2009). [Feature hashing for large scale multitask learning](#). Proc. ICML.
- [MurmurHash3](#).

6.2.3. Text feature extraction

6.2.3.1. The Bag of Words representation

Text Analysis is a major application field for machine learning algorithms. However the raw data, a sequence of symbols cannot be fed directly to the algorithms themselves as most of them expect numerical feature vectors with a fixed size rather than the raw text documents with variable length.

In order to address this, scikit-learn provides utilities for the most common ways to extract numerical features from text content, namely:

- **tokenizing** strings and giving an integer id for each possible token, for instance by using white-spaces and punctuation as token separators.
- **counting** the occurrences of tokens in each document.
- **normalizing** and weighting with diminishing importance tokens that occur in the majority of samples / documents.

In this scheme, features and samples are defined as follows:

- each **individual token occurrence frequency** (normalized or not) is treated as a **feature**.
- the vector of all the token frequencies for a given **document** is considered a multivariate **sample**.

A corpus of documents can thus be represented by a matrix with one row per document and one column per token (e.g. word) occurring in the corpus.

We call **vectorization** the general process of turning a collection of text documents into numerical feature vectors. This specific strategy (tokenization, counting and normalization) is called the **Bag of Words** or “Bag of n-grams” representation. Documents are described by word occurrences while completely ignoring the relative position information of the words in the document.

6.2.3.2. Sparsity

As most documents will typically use a very small subset of the words used in the corpus, the resulting matrix will have many feature values that are zeros (typically more than 99% of them).

For instance a collection of 10,000 short text documents (such as emails) will use a vocabulary with a size in the order of 100,000 unique words in total while each document will use 100 to 1000 unique words individually.

In order to be able to store such a matrix in memory but also to speed up algebraic operations matrix / vector, implementations will typically use a sparse representation such as the implementations available in the `scipy.sparse` package.

6.2.3.3. Common Vectorizer usage

[CountVectorizer](#) implements both tokenization and occurrence counting in a single class:

```
>>> from sklearn.feature_extraction.text import CountVectorizer
```

Toggle Menu

This model has many parameters, however the default values are quite reasonable (please see the [reference documentation](#) for the details):

```
>>> vectorizer = CountVectorizer()
>>> vectorizer
CountVectorizer()
```

Let's use it to tokenize and count the word occurrences of a minimalistic corpus of text documents:

```
>>> corpus = [
...     'This is the first document.',
...     'This is the second second document.',
...     'And the third one.',
...     'Is this the first document?',
... ]
>>> X = vectorizer.fit_transform(corpus)
>>> X
<4x9 sparse matrix of type '<... 'numpy.int64'>'
  with 19 stored elements in Compressed Sparse ... format>
```

The default configuration tokenizes the string by extracting words of at least 2 letters. The specific function that does this step can be requested explicitly:

```
>>> analyze = vectorizer.build_analyzer()
>>> analyze("This is a text document to analyze.") == (
...     ['this', 'is', 'text', 'document', 'to', 'analyze'])
True
```

Each term found by the analyzer during the fit is assigned a unique integer index corresponding to a column in the resulting matrix. This interpretation of the columns can be retrieved as follows:

```
>>> vectorizer.get_feature_names_out()
array(['and', 'document', 'first', 'is', 'one', 'second', 'the',
       'third', 'this'], ...)

>>> X.toarray()
array([[0, 1, 1, 1, 0, 0, 1, 0, 1],
       [0, 1, 0, 1, 0, 2, 1, 0, 1],
       [1, 0, 0, 0, 1, 0, 1, 1, 0],
       [0, 1, 1, 1, 0, 0, 1, 0, 1]]...)
```

The converse mapping from feature name to column index is stored in the `vocabulary_` attribute of the vectorizer:

```
>>> vectorizer.vocabulary_.get('document')
1
```

Hence words that were not seen in the training corpus will be completely ignored in future calls to the transform method:

```
>>> vectorizer.transform(['Something completely new.']).toarray()
array([[0, 0, 0, 0, 0, 0, 0, 0, 0]]...)
```

Note that in the previous corpus, the first and the last documents have exactly the same words hence are encoded in equal vectors. In particular we lose the information that the last document is an interrogative form. To preserve some of the local ordering information we can extract 2-grams of words in addition to the 1-grams (individual words):

```
>>> bigram_vectorizer = CountVectorizer(ngram_range=(1, 2),
...                                   token_pattern=r'\b\w+\b', min_df=1)
>>> analyze = bigram_vectorizer.build_analyzer()
>>> analyze('Bi-grams are cool!') == (
...     ['bi', 'grams', 'are', 'cool', 'bi grams', 'grams are', 'are cool'])
True
```

The vocabulary extracted by this vectorizer is hence much bigger and can now resolve ambiguities encoded in local positioning patterns:

```
>>> X_2 = bigram_vectorizer.fit_transform(corpus).toarray()
>>> X_2
array([[0, 0, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 1, 0],
       [0, 0, 1, 0, 0, 1, 1, 0, 0, 2, 1, 1, 1, 0, 1, 0, 0, 0, 1, 1, 0],
       [1, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 1, 1, 1, 0, 0, 0],
       [0, 0, 1, 1, 1, 1, 0, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1, 0, 1]]...)
```

In particular the interrogative form “Is this” is only present in the last document:

```
>>> feature_index = bigram_vectorizer.vocabulary_.get('is this')
>>> X_2[:, feature_index]
array([0, 0, 0, 1]...)
```

Using stop words

Stop words are words like “and”, “the”, “him”, which are presumed to be uninformative in representing the content of a text, and which may be removed to avoid them being construed as signal for prediction. Sometimes, however, similar words are useful for prediction, such as in classifying writing style or personality.

There are several known issues in our provided ‘english’ stop word list. It does not aim to be a general, ‘one-size-fits-all’ solution as some tasks may require a more custom solution. See [\[NQY18\]](#) for more details.

Please take care in choosing a stop word list. Popular stop word lists may include words that are highly informative to some tasks, such as *computer*.

You should also make sure that the stop word list has had the same preprocessing and tokenization applied as the one used in the vectorizer. The word *we’ve* is split into *we* and *ve* by CountVectorizer’s default tokenizer, so if *we’ve* is in `stop_words`, but *ve* is not, *ve* will be retained from *we’ve* in transformed text. Our vectorizers will try to identify and warn about some kinds of inconsistencies.

References

[\[NQY18\]](#)
J. Nothman, H. Qin and R. Yurchak (2018). [“Stop Word Lists in Free Open-source Software Packages”](#). In *Proc. Workshop for NLP Open Source Software*.

6.2.3.4. Tf–idf term weighting

In a large text corpus, some words will be very present (e.g. “the”, “a”, “is” in English) hence carrying very little meaningful information about the actual contents of the document. If we were to feed the direct count data directly to a classifier those very frequent terms would shadow the frequencies of rarer yet more interesting terms.

In order to re-weight the count features into floating point values suitable for usage by a classifier it is very common to use the tf–idf transform.

Tf means **term-frequency** while tf–idf means term-frequency times **inverse document-frequency**: $\text{tf-idf}(t,d) = \text{tf}(t,d) \times \text{idf}(t)$.

Using the `TfidfTransformer`’s default settings, `TfidfTransformer(norm='l2', use_idf=True, smooth_idf=True, sublinear_tf=False)` the term frequency, the number of times a term occurs in a given document, is multiplied with idf component, which is computed as

$$\text{idf}(t) = \log \frac{1+n}{1+\text{df}(t)} + 1,$$

where n is the total number of documents in the document set, and $\text{df}(t)$ is the number of documents in the document set that contain term t . The resulting tf-idf vectors are then normalized by the Euclidean norm:

$$v_{\text{norm}} = \frac{v}{\|v\|_2} = \frac{v}{\sqrt{v_1^2+v_2^2+\dots+v_n^2}}.$$

This was originally a term weighting scheme developed for information retrieval (as a ranking function for search engines results) that has also found good use in document classification and clustering.

The following sections contain further explanations and examples that illustrate how the tf-idfs are computed exactly and how the tf-idfs computed in scikit-learn’s [TfidfTransformer](#) and [TfidfVectorizer](#) differ slightly from the standard textbook notation that defines the idf as

$$\text{idf}(t) = \log \frac{n}{1+\text{df}(t)}.$$

In the [TfidfTransformer](#) and [TfidfVectorizer](#) with `smooth_idf=False`, the “1” count is added to the idf instead of the idf’s denominator:

$$\text{idf}(t) = \log \frac{n}{\text{df}(t)} + 1$$

This normalization is implemented by the [TfidfTransformer](#) class:

```
>>> from sklearn.feature_extraction.text import TfidfTransformer
>>> transformer = TfidfTransformer(smooth_idf=False)
>>> transformer
TfidfTransformer(smooth_idf=False)
```

Again please see the [reference documentation](#) for the details on all the parameters.

Let’s take an example with the following counts. The first term is present 100% of the time hence not very interesting. The two other features only in less than 50% of the time hence probably more representative of the content of the documents:

```
>>> counts = [[3, 0, 1],
...           [2, 0, 0],
...           [3, 0, 0],
...           [4, 0, 0],
...           [3, 2, 0],
...           [3, 0, 2]]
...
>>> tfidf = transformer.fit_transform(counts)
>>> tfidf
<6x3 sparse matrix of type '<... 'numpy.float64'>'
  with 9 stored elements in Compressed Sparse ... format>

>>> tfidf.toarray()
array([[0.81940995, 0.          , 0.57320793],
       [1.          , 0.          , 0.          ],
       [1.          , 0.          , 0.          ],
       [1.          , 0.          , 0.          ],
       [0.47330339, 0.88089948, 0.          ],
       [0.58149261, 0.          , 0.81355169]])
```

Each row is normalized to have unit Euclidean norm:

$$v_{norm} = \frac{v}{\|v\|_2} = \frac{v}{\sqrt{v_1^2+v_2^2+\dots+v_n^2}}$$

For example, we can compute the tf-idf of the first term in the first document in the `counts` array as follows:

$n = 6$

$df(t)_{term1} = 6$

$idf(t)_{term1} = \log \frac{n}{df(t)} + 1 = \log(1) + 1 = 1$

$tf-idf_{term1} = tf \times idf = 3 \times 1 = 3$

Now, if we repeat this computation for the remaining 2 terms in the document, we get

$tf-idf_{term2} = 0 \times (\log(6/1) + 1) = 0$

$tf-idf_{term3} = 1 \times (\log(6/2) + 1) \approx 2.0986$

and the vector of raw tf-idfs:

$tf-idf_{raw} = [3, 0, 2.0986].$

Then, applying the Euclidean (L2) norm, we obtain the following tf-idfs for document 1:

$\frac{[3,0,2.0986]}{\sqrt{(3^2+0^2+2.0986^2)}} = [0.819, 0, 0.573].$

Furthermore, the default parameter `smooth_idf=True` adds “1” to the numerator and denominator as if an extra document was seen containing every term in the collection exactly once, which prevents zero divisions:

$idf(t) = \log \frac{1+n}{1+df(t)} + 1$

Using this modification, the tf-idf of the third term in document 1 changes to 1.8473:

$tf-idf_{term3} = 1 \times \log(7/3) + 1 \approx 1.8473$

And the L2-normalized tf-idf changes to

$\frac{[3,0,1.8473]}{\sqrt{(3^2+0^2+1.8473^2)}} = [0.8515, 0, 0.5243]:$

```
>>> transformer = TfidfTransformer()
>>> transformer.fit_transform(counts).toarray()
array([[0.85151335, 0.          , 0.52433293],
       [1.          , 0.          , 0.          ],
       [1.          , 0.          , 0.          ],
       [1.          , 0.          , 0.          ],
       [0.55422893, 0.83236428, 0.          ],
       [0.63035731, 0.          , 0.77630514]])
```

The weights of each feature computed by the `fit` method call are stored in a model attribute:


```
>>> transformer.idf_
array([1. ..., 2.25..., 1.84...])
```

As tf-idf is very often used for text features, there is also another class called [TfidfVectorizer](#) that combines all the options of [CountVectorizer](#) and [TfidfTransformer](#) in a single model:

```
>>> from sklearn.feature_extraction.text import TfidfVectorizer
>>> vectorizer = TfidfVectorizer()
>>> vectorizer.fit_transform(corpus)
<4x9 sparse matrix of type '<... 'numpy.float64'>'
  with 19 stored elements in Compressed Sparse ... format>
```

While the tf-idf normalization is often very useful, there might be cases where the binary occurrence markers might offer better features. This can be achieved by using the `binary` parameter of [CountVectorizer](#). In particular, some estimators such as [Bernoulli Naive Bayes](#) explicitly model discrete boolean random variables. Also, very short texts are likely to have noisy tf-idf values while the binary occurrence info is more stable.

As usual the best way to adjust the feature extraction parameters is to use a cross-validated grid search, for instance by pipelining the feature extractor with a classifier:

- [Sample pipeline for text feature extraction and evaluation](#)

6.2.3.5. Decoding text files

Text is made of characters, but files are made of bytes. These bytes represent characters according to some *encoding*. To work with text files in Python, their bytes must be *decoded* to a character set called Unicode. Common encodings are ASCII, Latin-1 (Western Europe), KOI8-R (Russian) and the universal encodings UTF-8 and UTF-16. Many others exist.

Note: An encoding can also be called a ‘character set’, but this term is less accurate: several encodings can exist for a single character set.

The text feature extractors in scikit-learn know how to decode text files, but only if you tell them what encoding the files are in. The [CountVectorizer](#) takes an `encoding` parameter for this purpose. For modern text files, the correct encoding is probably UTF-8, which is therefore the default (`encoding="utf-8"`).

If the text you are loading is not actually encoded with UTF-8, however, you will get a `UnicodeDecodeError`. The vectorizers can be told to be silent about decoding errors by setting the `decode_error` parameter to either `"ignore"` or `"replace"`. See the documentation for the Python function `bytes.decode` for more details (type `help(bytes.decode)` at the Python prompt).

If you are having trouble decoding text, here are some things to try:

- Find out what the actual encoding of the text is. The file might come with a header or README that tells you the encoding, or there might be some standard encoding you can assume based on where the text comes from.
- You may be able to find out what kind of encoding it is in general using the UNIX command `file`. The Python `chardet` module comes with a script called `chardetect.py` that will guess the specific encoding, though you cannot rely on its guess being correct.
- You could try UTF-8 and disregard the errors. You can decode byte strings with `bytes.decode(errors='replace')` to replace all decoding errors with a meaningless character, or set `decode_error='replace'` in the vectorizer. This may damage the usefulness of your features.
- Real text may come from a variety of sources that may have used different encodings, or even be sloppily decoded in a different encoding than the one it was encoded with. This is common in text retrieved from the Web. The Python package [ftfy](#) can automatically sort out some classes of decoding errors, so you could try decoding the unknown text as `latin-1` and then using `ftfy` to fix errors.
- If the text is in a mish-mash of encodings that is simply too hard to sort out (which is the case for the 20 Newsgroups dataset), you can fall back on a simple single-byte encoding such as `latin-1`. Some text may display incorrectly, but at least the same sequence of bytes will always represent the same feature.

For example, the following snippet uses `chardet` (not shipped with scikit-learn, must be installed separately) to figure out the encoding of three texts. It then vectorizes the texts and prints the learned vocabulary. The output is not shown here.

```
>>> import chardet
>>> text1 = b"Sei mir gegr\xbc\x9ft mein Sauerkraut"
>>> text2 = b"holdselig sind deine Ger\xfcche"
>>> text3 = b"\xff\xfeA\x00u\x00f\x00 \x0F\x00l\x00\xfc\x00g\x00e\x00l\x00n\x00 \x0d\x00e\x00s\x00
\x0G\x00e\x00s\x00a\x00n\x00g\x00e\x00s\x00,\x00 \x0H\x00e\x00r\x00z\x00l\x00i\x00e\x00b\x00c\x00h\x00e\x00n\x00,\x00
\x0t\x00r\x00a\x00g\x00 \x0i\x00c\x00h\x00 \x0d\x00i\x00c\x00h\x00 \x0f\x00o\x00r\x00t\x00"
>>> decoded = [x.decode(chardet.detect(x)['encoding'])
...             for x in (text1, text2, text3)]
>>> v = CountVectorizer().fit(decoded).vocabulary_
>>> for term in v: print(v)
```

(Depending on the version of `chardet`, it might get the first one wrong.)

For an introduction to Unicode and character encodings in general, see Joel Spolsky’s [Absolute Minimum Every Software Developer Must Know About Unicode](#).

6.2.3.6. Applications and examples

The bag of words representation is quite simplistic but surprisingly useful in practice.

In particular in a **supervised setting** it can be successfully combined with fast and scalable linear models to train **document classifiers**, for instance:

- [Classification of text documents using sparse features](#)

In an **unsupervised setting** it can be used to group similar documents together by applying clustering algorithms such as [K-means](#):

- [Clustering text documents using k-means](#)

Finally it is possible to discover the main topics of a corpus by relaxing the hard assignment constraint of clustering, for instance by using [Non-negative matrix factorization \(NMF or NNMF\)](#):

- [Topic extraction with Non-negative Matrix Factorization and Latent Dirichlet Allocation](#)

6.2.3.7. Limitations of the Bag of Words representation

A collection of unigrams (what bag of words is) cannot capture phrases and multi-word expressions, effectively disregarding any word order dependence. Additionally, the bag of words model doesn’t account for potential misspellings or word derivations.

N-grams to the rescue! Instead of building a simple collection of unigrams (n=1), one might prefer a collection of bigrams (n=2), where occurrences of pairs of consecutive words are counted.

One might alternatively consider a collection of character n-grams, a representation resilient against misspellings and derivations.

For example, let’s say we’re dealing with a corpus of two documents: ['words' , 'wprds'] . The second document contains a misspelling of the word 'words'. A simple bag of words representation would consider these two as very distinct documents, differing in both of the two possible features. A character 2-gram representation, however, would find the documents matching in 4 out of 8 features, which may help the preferred classifier decide better:

```
>>> ngram_vectorizer = CountVectorizer(analyzer='char_wb', ngram_range=(2, 2))
>>> counts = ngram_vectorizer.fit_transform(['words', 'wprds'])
>>> ngram_vectorizer.get_feature_names_out()
array([' w', 'ds', 'or', 'pr', 'rd', 's ', 'wo', 'wp'], ...)
>>> counts.toarray().astype(int)
array([[1, 1, 1, 0, 1, 1, 1, 0],
       [1, 1, 0, 1, 1, 1, 0, 1]])
```

In the above example, char_wb analyzer is used, which creates n-grams only from characters inside word boundaries (padded with space on each side). The char analyzer, alternatively, creates n-grams that span across words:

```
>>> ngram_vectorizer = CountVectorizer(analyzer='char_wb', ngram_range=(5, 5))
>>> ngram_vectorizer.fit_transform(['jumpy fox'])
<1x4 sparse matrix of type '<... 'numpy.int64'>'
  with 4 stored elements in Compressed Sparse ... format>
>>> ngram_vectorizer.get_feature_names_out()
array([' fox ', ' jump', 'jumpy', 'umpy '], ...)

>>> ngram_vectorizer = CountVectorizer(analyzer='char', ngram_range=(5, 5))
>>> ngram_vectorizer.fit_transform(['jumpy fox'])
<1x5 sparse matrix of type '<... 'numpy.int64'>'
  with 5 stored elements in Compressed Sparse ... format>
>>> ngram_vectorizer.get_feature_names_out()
array(['jumpy', 'mpy f', 'py fo', 'umpy ', 'y fox'], ...)
```

The word boundaries-aware variant char_wb is especially interesting for languages that use white-spaces for word separation as it generates significantly less noisy features than the raw char variant in that case. For such languages it can increase both the predictive accuracy and convergence speed of classifiers trained using such features while retaining the robustness with regards to misspellings and word derivations.

While some local positioning information can be preserved by extracting n-grams instead of individual words, bag of words and bag of n-grams destroy most of the inner structure of the document and hence most of the meaning carried by that internal structure.

In order to address the wider task of Natural Language Understanding, the local structure of sentences and paragraphs should thus be taken into such models will thus be casted as “Structured output” problems which are currently outside of the scope of scikit-learn.

6.2.3.8. Vectorizing a large text corpus with the hashing trick

The above vectorization scheme is simple but the fact that it holds an **in- memory mapping from the string tokens to the integer feature indices** (the `vocabulary_` attribute) causes several **problems when dealing with large datasets**:

- the larger the corpus, the larger the vocabulary will grow and hence the memory use too,
- fitting requires the allocation of intermediate data structures of size proportional to that of the original dataset.
- building the word-mapping requires a full pass over the dataset hence it is not possible to fit text classifiers in a strictly online manner.
- pickling and un-pickling vectorizers with a large `vocabulary_` can be very slow (typically much slower than pickling / un-pickling flat data structures such as a NumPy array of the same size),
- it is not easily possible to split the vectorization work into concurrent sub tasks as the `vocabulary_` attribute would have to be a shared state with a fine grained synchronization barrier: the mapping from token string to feature index is dependent on ordering of the first occurrence of each token hence would have to be shared, potentially harming the concurrent workers' performance to the point of making them slower than the sequential variant.

It is possible to overcome those limitations by combining the “hashing trick” ([Feature hashing](#)) implemented by the [FeatureHasher](#) class and the text preprocessing and tokenization features of the [CountVectorizer](#).

This combination is implemented in [HashingVectorizer](#), a transformer class that is mostly API compatible with [CountVectorizer](#). [HashingVectorizer](#) is stateless, meaning that you don't have to call `fit` on it:

```
>>> from sklearn.feature_extraction.text import HashingVectorizer
>>> hv = HashingVectorizer(n_features=10)
>>> hv.transform(corpus)
<4x10 sparse matrix of type '<... 'numpy.float64'>'
  with 16 stored elements in Compressed Sparse ... format>
```

You can see that 16 non-zero feature tokens were extracted in the vector output: this is less than the 19 non-zeros extracted previously by the [CountVectorizer](#) on the same toy corpus. The discrepancy comes from hash function collisions because of the low value of the `n_features` parameter.

In a real world setting, the `n_features` parameter can be left to its default value of 2^{20} (roughly one million possible features). If memory or downstream models size is an issue selecting a lower value such as 2^{18} might help without introducing too many additional collisions on typical text classification tasks.

Note that the dimensionality does not affect the CPU training time of algorithms which operate on CSR matrices (`LinearSVC(dual=True)`, `Perceptron`, `SGDClassifier`, `PassiveAggressive`) but it does for algorithms that work with CSC matrices (`LinearSVC(dual=False)`, `Lasso()`, etc).

Let's try again with the default setting:

```
>>> hv = HashingVectorizer()
>>> hv.transform(corpus)
<4x1048576 sparse matrix of type '<... 'numpy.float64'>'
  with 19 stored elements in Compressed Sparse ... format>
```

We no longer get the collisions, but this comes at the expense of a much larger dimensionality of the output space. Of course, other terms than the 19 used here might still collide with each other.

The [HashingVectorizer](#) also comes with the following limitations:

- it is not possible to invert the model (no `inverse_transform` method), nor to access the original string representation of the features, because of the one-way nature of the hash function that performs the mapping.
- it does not provide IDF weighting as that would introduce statefulness in the model. A [TfidfTransformer](#) can be appended to it in a pipeline if required.

6.2.3.9. Performing out-of-core scaling with HashingVectorizer

An interesting development of using a [HashingVectorizer](#) is the ability to perform [out-of-core](#) scaling. This means that we can learn from data that does not fit into the computer's main memory.

A strategy to implement out-of-core scaling is to stream data to the estimator in mini-batches. Each mini-batch is vectorized using [HashingVectorizer](#) so as to guarantee that the input space of the estimator has always the same dimensionality. The amount of memory used at any time is thus bounded by the size of a mini-batch. Although there is no limit to the amount of data that can be ingested using such an approach, from a practical point of view the learning time is often limited by the CPU time one wants to spend on the task.

For a full-fledged example of out-of-core scaling in a text classification task see [Out-of-core classification of text documents](#).

6.2.3.10. Customizing the vectorizer classes

It is possible to customize the behavior by passing a callable to the vectorizer constructor:

```
>>> def my_tokenizer(s):
...     return s.split()
...
>>> vectorizer = CountVectorizer(tokenizer=my_tokenizer)
>>> vectorizer.build_analyzer()(u"Some... punctuation!") == (
...     ['some...', 'punctuation!'])
True
```

In particular we name:

- `preprocessor`: a callable that takes an entire document as input (as a single string), and returns a possibly transformed version of the document, still as an entire string. This can be used to remove HTML tags, lowercase the entire document, etc.
- `tokenizer`: a callable that takes the output from the preprocessor and splits it into tokens, then returns a list of these.
- `analyzer`: a callable that replaces the preprocessor and tokenizer. The default analyzers all call the preprocessor and tokenizer, but custom analyzers will skip this. N-gram extraction and stop word filtering take place at the analyzer level, so a custom analyzer may have to reproduce these steps.

(Lucene users might recognize these names, but be aware that scikit-learn concepts may not map one-to-one onto Lucene concepts.)

To make the preprocessor, tokenizer and analyzers aware of the model parameters it is possible to derive from the class and override the `build_preprocessor`, `build_tokenizer` and `build_analyzer` factory methods instead of passing custom functions.

Some tips and tricks:

- If documents are pre-tokenized by an external package, then store them in files (or strings) with the tokens separated by whitespace and pass `analyzer=str.split`
- Fancy token-level analysis such as stemming, lemmatizing, compound splitting, filtering based on part-of-speech, etc. are not included in the scikit-learn codebase, but can be added by customizing either the tokenizer or the analyzer. Here's a `CountVectorizer` with a tokenizer and lemmatizer using [NLTK](#):

```
>>> from nltk import word_tokenize
>>> from nltk.stem import WordNetLemmatizer
>>> class LemmaTokenizer:
...     def __init__(self):
...         self.wnl = WordNetLemmatizer()
...     def __call__(self, doc):
...         return [self.wnl.lemmatize(t) for t in word_tokenize(doc)]
...
>>> vect = CountVectorizer(tokenizer=LemmaTokenizer())
```

(Note that this will not filter out punctuation.)

The following example will, for instance, transform some British spelling to American spelling:

```
>>> import re
>>> def to_british(tokens):
...     for t in tokens:
...         t = re.sub(r"(...)our$", r"\1or", t)
...         t = re.sub(r"([bt])re$", r"\1er", t)
...         t = re.sub(r"([iy])s(e$|ing|ation)", r"\1z\2", t)
...         t = re.sub(r"ogue$", "og", t)
...         yield t
...
>>> class CustomVectorizer(CountVectorizer):
...     def build_tokenizer(self):
...         tokenize = super().build_tokenizer()
...         return lambda doc: list(to_british(tokenize(doc)))
...
>>> print(CustomVectorizer().build_analyzer()(u"color colour"))
[...'color', ...'color']
```

for other styles of preprocessing; examples include stemming, lemmatization, or normalizing numerical tokens, with the latter illustrated in:

- [Biclustering documents with the Spectral Co-clustering algorithm](#)

6.2.4. Image feature extraction

6.2.4.1. Patch extraction

The [extract_patches_2d](#) function extracts patches from an image stored as a two-dimensional array, or three-dimensional with color information along the third axis. For rebuilding an image from all its patches, use [reconstruct_from_patches_2d](#). For example let us generate a 4x4 pixel picture with 3 color channels (e.g. in RGB format):

```
>>> import numpy as np
>>> from sklearn.feature_extraction import image

>>> one_image = np.arange(4 * 4 * 3).reshape((4, 4, 3))
>>> one_image[:, :, 0] # R channel of a fake RGB picture
array([[ 0,  3,  6,  9],
       [12, 15, 18, 21],
       [24, 27, 30, 33],
       [36, 39, 42, 45]])

>>> patches = image.extract_patches_2d(one_image, (2, 2), max_patches=2,
...     random_state=0)
>>> patches.shape
(2, 2, 2, 3)
>>> patches[:, :, :, 0]
array([[[ 0,  3],
        [12, 15]],

       [[15, 18],
        [27, 30]]])
>>> patches = image.extract_patches_2d(one_image, (2, 2))
>>> patches.shape
(9, 2, 2, 3)
>>> patches[4, :, :, 0]
array([[15, 18],
       [27, 30]])
```

Let us now try to reconstruct the original image from the patches by averaging on overlapping areas:

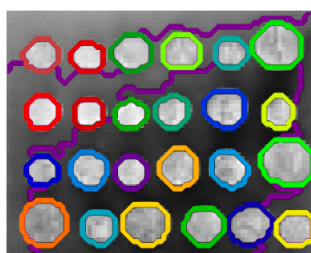
```
>>> reconstructed = image.reconstruct_from_patches_2d(patches, (4, 4, 3))
>>> np.testing.assert_array_equal(one_image, reconstructed)
```

The [PatchExtractor](#) class works in the same way as [extract_patches_2d](#), only it supports multiple images as input. It is implemented as an estimator, so it can be used in pipelines. See:

```
>>> five_images = np.arange(5 * 4 * 4 * 3).reshape(5, 4, 4, 3)
>>> patches = image.PatchExtractor(patch_size=(2, 2)).transform(five_images)
>>> patches.shape
(45, 2, 2, 3)
```

6.2.4.2. Connectivity graph of an image

Several estimators in the scikit-learn can use connectivity information between features or samples. For instance Ward clustering ([Hierarchical clustering](#)) can cluster together only neighboring pixels of an image, thus forming contiguous patches:



For this purpose, the estimators use a 'connectivity' matrix, giving which samples are connected.

The function [img_to_graph](#) returns such a matrix from a 2D or 3D image. Similarly, [grid_to_graph](#) build a connectivity matrix for images given the shape of these image.

These matrices can be used to impose connectivity in estimators that use connectivity information, such as Ward clustering ([Hierarchical clustering](#)), but also to build precomputed kernels, or similarity matrices.

- [A demo of structured Ward hierarchical clustering on an image of coins](#)
- [Spectral clustering for image segmentation](#)
- [Feature agglomeration vs. univariate selection](#)

© 2007 - 2023, scikit-learn developers (BSD License). [Show this page source](#)