## 3.4.2 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.

---

### 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 of 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 Fransisco', '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()
['city=Dubai', 'city=London', 'city=San Fransisco', 'temperature']
```

*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 *text.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()
['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`.

### 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.

If `non_negative=True` is passed to the constructor, the absolute value is taken. This undoes some of the collision handling, but allows the output to be passed to estimators like *sklearn.naive_bayes.MultinomialNB* or *sklearn.feature_selection.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 *text.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:

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

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.

### 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 $\xi$ 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.

### Text feature extraction

### 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.

## 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.

## Common Vectorizer usage

*CountVectorizer* implements both tokenization and occurrence counting in a single class:

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

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

```
>>> vectorizer = CountVectorizer(min_df=1)
>>> vectorizer
CountVectorizer(analyzer=...'word', binary=False, decode_error=...'strict',
        dtype=<... 'numpy.int64'>, encoding=...'utf-8', input=...'content',
        lowercase=True, max_df=1.0, max_features=None, min_df=1,
        ngram_range=(1, 1), preprocessor=None, stop_words=None,
        strip_accents=None, token_pattern=...'(?u)\\b\\w\\w+\\b',
        tokenizer=None, vocabulary=None)
```

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() == (
...     ['and', 'document', 'first', 'is', 'one',
...      'second', 'the', 'third', 'this'])
True

>>> 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]...)
```

### 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**: tf-idf(t,d) = tf(t,d) $\times$ idf(t).

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

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

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

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

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

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

In the *TfidfTransformer* and *TfidfVectorizer* with `smooth_idf=False`, the "1" count is added to the idf instead of the idf's denominator:

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

This normalization is implemented by the *TfidfTransformer* class:

```
>>> from sklearn.feature_extraction.text import TfidfTransformer
>>> transformer = TfidfTransformer(smooth_idf=False)
>>> transformer
TfidfTransformer(norm=...'l2', smooth_idf=False, sublinear_tf=False,
                 use_idf=True)
```

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:

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

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

$$n_{d,\text{term1}} = 6$$

$$\text{df}(d, t)_{\text{term1}} = 6$$

$$\text{idf}(d, t)_{\text{term1}} = log\frac{n_d}{\text{df}(d,t)} + 1 = log(1) + 1 = 1$$

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

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

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

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

and the vector of raw tf-idfs:

$$\text{tf-idf}_r aw = [3, 0, 2.0986].$$

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

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

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:

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

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

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

And the L2-normalized tf-idf changes to

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

```
>>> 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(min_df=1)
>>> 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*

### 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: