



Participate in the biggest free nature photography contest and compete for a photographic expedition! From the 1st to the 31st of July!

k-nearest neighbors algorithm

23 languages

Contents	hide
<div><div><div><div><div></div><div>(Top)</div></div><div><div></div><div>Statistical setting</div></div><div><div></div><div>Algorithm</div></div><div><div></div><div>Parameter selection</div></div><div><div></div><div>The 1-nearest neighbor classifier</div></div><div><div></div><div>The weighted nearest neighbour classifier</div></div><div><div></div><div>Properties</div></div><div><div></div><div>Error rates</div></div><div><div></div><div>Metric learning</div></div><div><div></div><div>Feature extraction</div></div><div><div></div><div>Dimension reduction</div></div><div><div></div><div>Decision boundary</div></div><div><div></div><div>Data reduction</div></div><div><div></div><div>Selection of class-outliers</div></div><div><div></div><div>Condensed Nearest Neighbor for data reduction</div></div><div><div></div><div>k-NN regression</div></div></div></div></div>	

From Wikipedia, the free encyclopedia

Article

Talk

Read

Edit

View history

Tools

Appearance

hide

Text

Small

Standard

Large

Width

Standard

Wide

Color (beta)

Automatic

Light

Dark

Report an issue with dark mode

*This article is about supervised classification/regression, but NOT a clustering algorithm. For algorithms for finding nearest neighbors, see *Nearest neighbor search*.*

*Not to be confused with *k-means clustering*.*

In *statistics*, the **k-nearest neighbors algorithm** (**k-NN**) is a non-parametric supervised learning method first developed by Evelyn Fix and Joseph Hodges in 1951,^[1] and later expanded by Thomas Cover.^[2] It is used for *classification* and *regression*. In both cases, the input consists of the *k* closest training examples in a *data set*. The output depends on whether k-NN is used for classification or regression:

- In *k-NN classification*, the output is a class membership. An object is classified by a plurality vote of its neighbors, with the object being assigned to the class most common among its *k* nearest neighbors (*k* is a positive *integer*, typically small). If *k* = 1, then the object is simply assigned to the class of that single nearest neighbor.
- In *k-NN regression*, the output is the property value for the object. This value is the average of the values of *k* nearest neighbors. If *k* = 1, then the output is simply assigned to the value of that single nearest neighbor.

k-NN is a type of *classification* where the function is only approximated locally and all computation is deferred until function evaluation. Since this algorithm relies on distance for classification, if the features represent different physical units or come in vastly different scales then *normalizing* the training data can improve its accuracy dramatically.^[3]

Both for classification and regression, a useful technique can be to assign weights to the contributions of the neighbors, so that the nearer neighbors contribute more to the average than the more distant ones. For example, a common weighting scheme consists in giving each neighbor a weight of *1/d*, where *d* is the distance to the neighbor.^[4]

The neighbors are taken from a set of objects for which the class (for k-NN classification) or the object property value (for k-NN regression) is known. This can be thought of as the training set for the algorithm, though no explicit training step is required.

A peculiarity of the k-NN algorithm is that it is sensitive to the local structure of the data.

Statistical setting

[edit]

Suppose we have pairs *(X*₁, *Y*₁), (*X*₂, *Y*₂), . . . , (*X*_{*n*}, *Y*_{*n*}) taking values in

R

d

×
{
1
,
2
}

{\displaystyle \mathbb {R} ^{d}\times \{1,2\}}

, where *Y* is the class label of *X*, so that *X*|*Y* = *r* ∼ *P*_{*r*} for *r* = 1, 2 (and probability distributions *P*_{*r*}). Given some norm

∥
⋅
∥

{\displaystyle \|\cdot \|}

 on

R

d

{\displaystyle \mathbb {R} ^{d}}

 and a point *x* ∈

R

d

{\displaystyle \mathbb {R} ^{d}}

, let (*X*₍₁₎, *Y*₍₁₎), . . . , (*X*_(*n*), *Y*_(*n*)) be a reordering of the training data such that

∥

X

(
1
)

−
x
∥
≤
⋅
⋅
≤
⋅
⋅
≤
∥

X

(
n
)

−
x
∥

{\displaystyle \|X_{(1)}-x\|\leq \cdots \leq \cdots \leq \|X_{(n)}-x\|}

.

Algorithm

[edit]

The training examples are vectors in a multidimensional feature space, each with a class label. The training phase of the algorithm consists only of storing the *feature vectors* and class labels of the training samples.

In the classification phase, *k* is a user-defined constant, and an unlabeled vector (a query or test point) is classified by assigning the label which is most frequent among the *k* training samples nearest to that query point.

A commonly used distance metric for *continuous variables* is *Euclidean distance*. For discrete variables, such as for text classification, another metric can be used, such as the **overlap metric** (or *Hamming distance*). In the context of gene expression microarray data, for example, k-NN has been employed with correlation coefficients, such as Pearson and Spearman, as a metric.^[5] Often, the classification accuracy of k-NN can be improved significantly if the distance metric is learned with specialized algorithms such as Large Margin Nearest Neighbor or Neighbourhood components analysis.

A drawback of the basic "majority voting" classification occurs when the class distribution is skewed. That is, examples of a more frequent class tend to dominate the prediction of the new example, because they tend to be common among the *k* nearest neighbors due to their large number.^[6]

One way to overcome this problem is to weight the classification, taking into account the distance from the test point to each of its *k* nearest neighbors. The class (or value, in regression problems) of each of the *k* nearest points is multiplied by a weight proportional to the inverse of the distance from that point to the test point. Another way to overcome skew is by abstraction in data representation. For example, in a *self-organizing map* (SOM), each node is a representative (a center) of a cluster of similar points, regardless of their density in the original training data. K-NN can then be applied to the SOM.

Parameter selection

[edit]

The best choice of *k* depends upon the data; generally, larger values of *k* reduces effect of the noise on the classification,^[7] but make boundaries between classes less distinct. A good *k* can be selected by various *heuristic* techniques (see *hyperparameter optimization*). The special case where the class is predicted to be the class of the closest training sample (i.e. when *k* = 1) is called the nearest neighbor algorithm.

The accuracy of the k-NN algorithm can be severely degraded by the presence of noisy or irrelevant features, or if the feature scales are not consistent with their importance. Much research effort has been put into *selecting* or *scaling* features to improve classification. A particularly popular^[*citation needed*] approach is the use of *evolutionary algorithms* to optimize feature scaling.^[8] Another popular approach is to scale features by the *mutual information* of the training data with the training classes.^[*citation needed*]

In binary (two class) classification problems, it is helpful to choose *k* to be an odd number as this avoids tied votes. One popular way of choosing the empirically optimal *k* in this setting is via bootstrap method.^[9]

The 1-nearest neighbor classifier

[edit]

The most intuitive nearest-neighbour type classifier is the one nearest neighbour classifier that assigns a point *x* to the class of its closest neighbour in the feature space, that is

C

1
n
n

n

(
x
)
=

Y

(

1
)

{\displaystyle C_{1nn}^{n}(x)=Y_{(1)}}

.

As the size of training data set approaches infinity, the one nearest neighbour classifier guarantees an error rate of no worse than twice the *Bayes error rate* (the minimum achievable error rate given the distribution of the data).

The weighted nearest neighbour classifier

[edit]

The *k*-nearest neighbour classifier can be viewed as assigning the *k* nearest neighbours a weight

1

/

k

{\displaystyle 1/k}

 and all others *0* weight. This can be generalised to weighted nearest neighbour classifiers. That is, where the *i*th nearest neighbour is assigned a weight

w

n
i

{\displaystyle w_{ni}}

 with

∑

i
=
1

n

w

n
i

=
1

{\displaystyle \sum _{i=1}^{n}w_{ni}=1}

. An analogous result on the consistency of weighted nearest neighbour classifiers also holds.^[10]

Let

C

n
n
n

{\displaystyle C_{nnn}^{n}}

 denote the weighted nearest classifier with weights

{

w

n
i

}

i
=
1

n

{\displaystyle \{w_{ni}\}_{i=1}^{n}}

. Subject to regularity conditions, which in asymptotic theory are conditional variables which require assumptions to differentiate among parameters with some criteria. On the class distributions the excess risk has the following asymptotic expansion^[11]

R

R

(

C

n
n
n

)
−

R

R

(

C

Bayes

)
=
(

B

1

s

n

2

+

B

2

t

n

2

)
{
1
+
o
(
1
)
}
,

{\displaystyle {\mathcal {R}}_{R}(C_{nnn}^{n})-{\mathcal {R}}_{R}(C^{\mathrm {Bayes} })=(B_{1}s_{n}^{2}+B_{2}t_{n}^{2})\{1+o(1)\},}

for constants *B*₁ and *B*₂ where

s

n

2

=

∑

i
=
1

n

w

n
i

2

{\displaystyle s_{n}^{2}=\sum _{i=1}^{n}w_{ni}^{2}}

 and

t

n

=

n

−2

/

d

∑

i
=
1

n

{

i

1
+
2

/

d

−
(
i
−
1

)

1
+
2

/

d

}

{\displaystyle t_{n}=n^{-2/d}\sum _{i=1}^{n}\{i^{1+2/d}-(i-1)^{1+2/d}\}}

.

The optimal weighting scheme

{

w

n
i

}

i
=
1

n

{\displaystyle \{w_{ni}\}_{i=1}^{n}}

, that balances the two terms in the display above, is given as follows: set

k

∗

=
⌊
B

n

2

/

d

+
4

⌋

{\displaystyle k^{*}=\lfloor Bn^{2/d+4}\rfloor }

.

w

n
i

∗

=

1

k

∗

[

1
+

d

2

−

d

2

k

∗

2

/

d

{

i

1
+
2

/

d

−
(
i
−
1

)

1
+
2

/

d

}

]

{\displaystyle w_{ni}^{*}={\frac {1}{k^{*}}}\left[1+{\frac {d}{2}}-{\frac {d}{2k^{*2/d}}}\{i^{1+2/d}-(i-1)^{1+2/d}\}\right]}

for *i* = 1, 2, . . . , *k*^{*} and

w

n
i

∗

=
0

{\displaystyle w_{ni}^{*}=0}

for *i* = *k*^{*} + 1, . . . , *n*.

With optimal weights the dominant term in the asymptotic expansion of the excess risk is

O

(

n

−4

/

d
+
4

)

{\displaystyle O(n^{-{\frac {4}{d+4}}})}

. Similar results are true when using a *bagged nearest neighbour classifier*.

Properties

[edit]

k-NN is a special case of a *variable-bandwidth*, *kernel density* "balloon" estimator with a uniform *kernel*.^{[12][13]}

The naive version of the algorithm is easy to implement by computing the distances from the test example to all stored examples, but it is computationally intensive for large training sets. Using an approximate *nearest neighbor search* algorithm makes k-NN computationally tractable even for large data sets. Many nearest neighbor search algorithms have been proposed over the years; these generally seek to reduce the number of distance evaluations actually performed.

k-NN has some strong *consistency* results. As the amount of data approaches infinity, the two-class k-NN algorithm is guaranteed to yield an error rate no worse than twice the *Bayes error rate* (the minimum achievable error rate given the distribution of the data).^[2] Various improvements to the k-NN speed are possible by using proximity graphs.^[14]

For multi-class k-NN classification, Cover and Hart (1967) prove an upper bound error rate of

R

∗

≤

R

k
N
N

≤

R

∗

(
2
−

M

R

∗

M
−
1

)

{\displaystyle R^{*}\leq R_{kNN}\leq R^{*}\left(2-{\frac {MR^{*}}{M-1}}\right)}

where *R*^{*} is the Bayes error rate (which is the minimal error rate possible), *R*_{kNN} is the asymptotic k-NN error rate, and *M* is the number of classes in the problem. This bound is tight in the sense that both the lower and upper bounds are achievable by some distribution.^[15] For *M* = 2 and as the Bayesian error rate *R*^{*} approaches zero, this limit reduces to "not more than twice the Bayesian error rate".

Error rates

[edit]

There are many results on the error rate of the *k* nearest neighbour classifiers.^[16] The *k*-nearest neighbour classifier is strongly (that is for any joint distribution on (*X*, *Y*)) *consistent* provided

k
∗

:=

k

n

{\displaystyle k^{*}:=k_{n}}

 diverges and

k

n

/

n

{\displaystyle k_{n}/n}

 converges to zero as *n* → ∞.

Let

C

k
n
n

{\displaystyle C_{knn}^{n}}

 denote the *k* nearest neighbour classifier based on a training set of size *n*. Under certain regularity conditions, the *excess risk* yields the following asymptotic expansion^[11]

R

R

(

C

k
n
n

)
−

R

R

(

C

Bayes

)
=
{

B

1

1

k

+

B

2

(

k

n

)

4

/

d

}
{
1
+
o
(
1
)
}
,

{\displaystyle {\mathcal {R}}_{R}(C_{knn}^{n})-{\mathcal {R}}_{R}(C^{\mathrm {Bayes} })=\left\{B_{1}{\frac {1}{k}}+B_{2}\left({\frac {k}{n}}\right)^{4/d}\right\}\{1+o(1)\},}

for some constants *B*₁ and *B*₂.

The choice

k

∗

=
⌊
B

n

4

/

d
+
4

⌋

{\displaystyle k^{*}=\lfloor Bn^{4/d+4}\rfloor }

 offers a trade off between the two terms in the above display, for which the *k*^{*}-

nearest neighbour error converges to the Bayes error at the optimal (minimax) rate

O

(

n

−4

/

d
+
4

)

{\displaystyle O\left(n^{-{\frac {4}{d+4}}}\right)}

.

Metric learning

[edit]

The K-nearest neighbor classification performance can often be significantly improved through (supervised) metric learning. Popular algorithms are *neighbourhood components analysis* and *large margin nearest neighbor*. Supervised metric learning algorithms use the label information to learn a new *metric* or *pseudo-metric*.

Feature extraction

[edit]

When the input data to an algorithm is too large to be processed and it is suspected to be redundant (e.g. the same measurement in both feet and meters) then the input data will be transformed into a reduced representation set of features (also named features vector). Transforming the input data into the set of features is called *feature extraction*. If the features extracted are carefully chosen it is expected that the features set will extract the relevant information from the input data in order to perform the desired task using this reduced representation instead of the full size input. Feature extraction is performed on raw data prior to applying k-NN algorithm on the transformed data in *feature space*.

An example of a typical *computer vision* computation pipeline for *face recognition* using k-NN including feature extraction and dimension reduction pre-processing steps (usually implemented with OpenCV):

- Haar face detection
- Mean-shift tracking analysis
- PCA or Fisher LDA projection into feature space, followed by k-NN classification

Dimension reduction

[edit]

For high-dimensional data (e.g., with number of dimensions more than 10) *dimension reduction* is usually performed prior to applying the k-NN algorithm in order to avoid the effects of the *curse of dimensionality*.^[17]

The *curse of dimensionality* in the k-NN context basically means that *Euclidean distance* is unhelpful in high dimensions because all vectors are almost equidistant to the search query vector (imagine multiple points lying more or less on a circle with the query point at the center; the distance from the query to all data points in the search space is almost the same).

Feature extraction and dimension reduction can be combined in one step using *principal component analysis* (PCA), *linear discriminant analysis* (LDA), or *canonical correlation analysis* (CCA) techniques as a pre-processing step, followed by clustering by k-NN on *feature vectors* in reduced-dimension space. This process is also called low-dimensional *embedding*.^[18]

For very-high-dimensional datasets (e.g. when performing a similarity search on live video streams, DNA data or high-dimensional *time series*) running a fast *approximate* k-NN search using *locally sensitive hashing*, "random projections",^[19] "sketches"^[20] or other high-dimensional similarity search techniques from the VLDB toolbox might be the only feasible option.

Decision boundary

[edit]

Nearest neighbor rules in effect implicitly compute the *decision boundary*. It is also possible to compute the decision boundary explicitly, and to do so efficiently, so that the computational complexity is a function of the boundary complexity.^[21]

Data reduction

[edit]

Data reduction is one of the most important problems for work with huge data sets. Usually, only some of the data points are needed for accurate classification. Those data are called the *prototypes* and can be found as follows:

- Select the *class-outliers*, that is, training data that are classified incorrectly by k-NN (for a given *k*)
- Separate the rest of the data into two sets: (i) the prototypes that are used for the k-NN classification decisions and (ii) the *absorbed points* that can be correctly classified by k-NN using prototypes. The absorbed points can then be removed from the training set.

Selection of class-outliers

[edit]

A training example surrounded by examples of other classes is called a class outlier. Caused of class outliers include:

- random error
- insufficient training examples of this class (an isolated example appears instead of a cluster)
- missing important features (the classes are separated in other dimensions which we don't know)
- too many training examples of other classes (unbalanced classes) that create a "hostile" background for the given small class

Class outliers with k-NN produce noise. They can be detected and separated for future analysis. Given two natural numbers, *k*>*r*>0, a training example is called a (*k*,*r*)-NN class-outlier if its *k* nearest neighbors include more than *r* examples of other classes.

Condensed Nearest Neighbor for data reduction

[edit]

Condensed nearest neighbor (CNN, the *Hart algorithm*) is an algorithm designed to reduce the data set for k-NN classification.^[22] It selects the set of prototypes *U* from the training data, such that 1NN with *U* can classify the examples almost as accurately as 1NN does with the whole data set.

Given a training set *X*, CNN works iteratively:

- Scan all elements of *X*, looking for an element *x* whose nearest prototype from *U* has a different label than *x*.
- Remove *x* from *X* and add it to *U*
- Repeat the scan until no more prototypes are added to *U*.

Use *U* instead of *X* for classification. The examples that are not prototypes are called "absorbed" points.

It is efficient to scan the training examples in order of decreasing border ratio.^[23] The border ratio of a training example *x* is defined as

a
(
x
)
=

|

x

∗

−
y
|

|

x
′
−
y
|

{\displaystyle a(x)={\frac {|x^{*}-y|}{|x'-y|}}}

where

|

x

∗

−
y
|

{\displaystyle |x^{*}-y|}

 is the distance to the closest example *y* having a different color than *x*, and

|

x
′
−
y
|

{\displaystyle |x'-y|}

 is the distance from *y* to its closest example *x*' with the same label as *x*.

The border ratio is in the interval [0,1] because

|

x

∗

−
y
|

{\displaystyle |x^{*}-y|}

 never exceeds

|

x
′
−
y
|

{\displaystyle |x'-y|}

. This ordering gives preference to the borders of the classes for inclusion in the set of prototypes *U*. A point of this different label than *x* is called external to *x*. The calculation of the border ratio is illustrated by the figure on the right. The data points are labeled by colors: the initial point is *x* and its label is red. External points are blue and green. The closest to *x* external point is *y*. The closest to *y* red point is *x*'. The border ratio

a
(
x
)
=
|

x

∗

−
y
|

/

|

x
′
−
y
|

{\displaystyle a(x)=|x^{*}-y|/|x'-y|}

 is the attribute of the initial point *x*.

Below is an illustration of CNN in a series of figures. There are three classes (red, green and blue). Fig. 1: initially there are 60 points in each class. Fig. 2 shows the 1NN classification map: each pixel is classified by 1NN using all the data. Fig. 3 shows the 5NN classification map. White areas correspond to the unclassified regions, where 5NN voting is tied (for example, if there are two green, two red and one blue points among 5 nearest neighbors). Fig. 4 shows the reduced data set. The crosses are the class-outliers selected by the (3,2)-NN rule (all the three nearest neighbors of these instances belong to other classes); the squares are the prototypes, and the empty circles are the absorbed points. The left bottom corner shows the numbers of the class-outliers, the prototype types and absorbed points for all three classes. The number of prototypes varies from 15% to 20% for different classes in this example. Fig. 5 shows that the 1NN classification map with the prototypes is very similar to that with the initial data set. The figures were produced using the Mirkes applet.^[24]

CNN model reduction for k-NN classifiers



Fig. 1. The dataset.



Fig. 2. The 1NN classification map.



Fig. 3. The 5NN classification map.



Fig. 4. The CNN reduced dataset.



Fig. 5. The 1NN classification map based on the CNN extracted prototypes.

k-NN regression

[edit]

In k-NN regression, the k-NN algorithm^[*citation needed*] is used for estimating continuous variables. One such algorithm uses a weighted average of the *k* nearest neighbors, weighted by the inverse of their distance. This algorithm works as follows:

- Compute the Euclidean or *Mahalanobis distance* from the query example to the labeled examples.
- Order the labeled examples by increasing distance.
- Find a heuristically optimal number *k* of nearest neighbors, based on *RMSE*. This is done using cross validation.
- Calculate an inverse distance weighted average with the *k*-nearest multivariate neighbors.

k-NN outlier

[edit]

The distance to the *k*th nearest neighbor can input data to their local density estimate and thus is also a popular outlier score in *anomaly detection*. The larger the distance to the *k*-NN, the lower the local density, the more likely the query point is an outlier.^[25] Although quite simple, this outlier model, along with another classic data mining method, *local outlier factor*, works quite well also in comparison to more recent and more complex approaches, according to a large scale experimental analysis^[25]

Validation of results

[edit]

A *confusion matrix* or "matching matrix" is often used as a tool to validate the accuracy