

ished in Towards Data Science

2 free member-only stories left this month.

[for Medium and get an extra one](#)

atsal

in 11, 2021 · 4 min read · ✨ Member-only · [Listen](#)



K Nearest Neighbours Explained

and the KNN algorithm and its implementation in Python
using the sklearn library



Image from : <https://unsplash.com/photos/IW25ZxpklN8>

In this article I will give a general overview, implementation, drawbacks and

is associated with the K Nearest Neighbours algorithm. Supervised is a subsection of machine learning generally associated with classification and regression based problems. Supervised learning implies we are training a model using a labelled dataset. K Nearest Neighbours falls under the supervised learning umbrella and is one of the core algorithms in machine learning. It's a highly used, simple yet efficient type of a non-parametric, lazy learner classification algorithm.

Lazy learner implies that it doesn't learn a discriminative function from training data but rather memorizes the training data instead

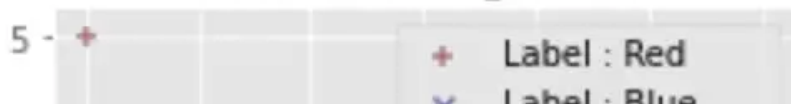
Non-parametric implies that the algorithm makes no assumptions about the distribution of the data.

The algorithm classifies unclassified data points based on their distance and similarity to other available data points. The underlying assumption this algorithm makes is that similar data points can be found close to another. It's commonly used to solve problems in various industries due to its ease of use, application to classification and regression problems, and the ease of interpretability of the results it generates.

Algorithm

It determines the number of nearest neighbours. When $K = 1$, the algorithm becomes the nearest neighbour algorithm. This is the simplest scenario. Given an unlabelled position X , the algorithm can predict its label by finding the closest labelled point to X and assigning that as the label.

Nearest Neighbour



Search Medium

Write

Sign up

Sign In



enario, the unlabelled black point X would be predicted to be a blue point based on the nearest proximity of labelled points (Image provided by Author)

rithm works as follows :

ose the number of K and a distance metric used to calculate
ximity
d the K nearest neighbours of the point we want to classify
ign the point a label by majority vote

Machine Learning Engineer

<https://www.linkedin.com/in/vatsal-pa57978149/>

Follow

More from Medium



Zach ... in Pipeline: A Data Engi...

3 Data Science Projects That Got Me 12 Interviews. And 1 That Got Me in Trouble.



Anil Tilbe in Level Up Coding

K-Nearest Neighbor (KNN): Why Do We Make It So Difficult? Simplified



Mark Schaefer

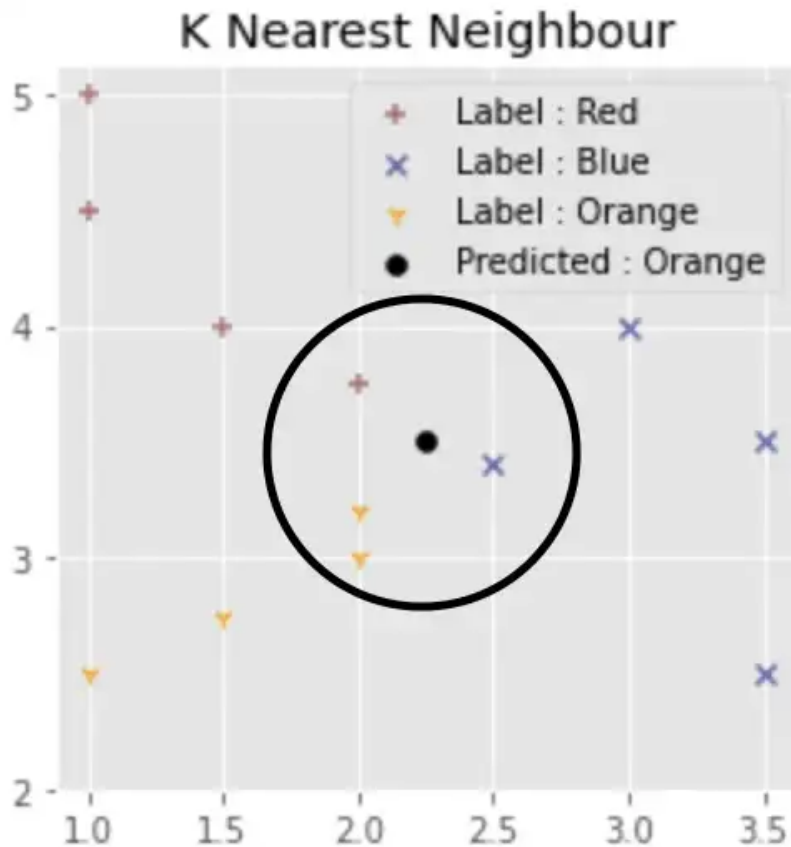
20 Entertaining Uses of ChatGPT You Never Knew Were Possible



Murtaza ... in Towards Data Scie...

10 Simple Things You Can Do to Improve Your Data Science Skills in 2023

[Help](#) [Status](#) [Writers](#) [Blog](#) [Careers](#) [Privacy](#) [Text to speech](#)



number of orange points is larger than the number of blue and red points in the proximity of the black point. Thus the algorithm will predict this to be an orange label (Image provided by Author)

Value of K

The value of K is crucial for the model, if chosen incorrectly it can cause the model to be over / under fit. A K value too small will cause noise in the model to have a high influence on the prediction, however a K value too large will make it computationally expensive.

The industry standard for choosing the optimal value of K is by taking the square root of N, where N is the total number of samples. Of course, take this with a grain of salt as it varies from problem to problem.

One experiment with various values of K and their associated accuracies. One practice to determine the accuracy of a KNN model is to use confusion matrices, cross validation or F1 scores.

Advantages & Disadvantages

ve listed some of the advantages and disadvantages of using the
orithm.

ges

le & intuitive — The algorithm is very easy to understand and
ement

ory based approach — Allows it to immediately adapt to new
ing data

ty of distance metrics — There is flexibility from the users side to
distance metric which is best suited for their application
idean, Minkowski, Manhattan distance etc.)

rtages

putational complexity — As your training data increases, the speed at
h calculations are made rapidly decrease

performance on imbalanced data — When majority of the data the
l is being trained on represents 1 label then that label will have a
likelihood of being predicted

nal value of K — If chosen incorrectly, the model will be under or
fitted to the data

entation

ry

ation, this article outlines that kNN is a lazy learner and non
ric algorithm. It works by assigning a label to an unlabelled point
the proximity of the unlabelled point to all the other nearest
points. It's main disadvantages are that it is quite computationally
nt and its difficult to pick the "correct" value of K. However, the
ges of this algorithm is that it is versatile to different calculations of
y, it's very intuitive and that it's a memory based approach.

ces

<https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html#sklearn.neighbors.KNeighborsClassifier>

enjoyed this read then check out my other works as well.

<div><div>Chain Explained</div><div>rticle I will explain and provide the python implementations ov chain. This article will not be a deep...</div><div>atascience.com</div></div>	
<div><div>Prediction Recommendation Engines with Node2Vec</div><div>ode Embeddings for Link Prediction</div><div>atascience.com</div></div>	
<div><div>Vec Explained</div><div>ng the Intuition of Word2Vec & Implementing it in Python</div><div>atascience.com</div></div>	
<div><div>Recommendation Systems Explained</div><div>ng & Implementing Content Based, Collaborative Filtering & Recommendation Systems in Python</div><div>atascience.com</div></div>	
<div><div>Monte Carlo Method Explained</div><div>ost I will introduce, explain and implement the Monte Carlo to you. This method of simulation is one of...</div><div>com</div></div>	



for The Variable

data Science

ay, the Variable delivers the very best of Towards Data Science: from hands-on tutorials and cutting-edge original features you don't want to miss. [Take a look.](#)

ou will create a Medium account if you don't already have one. Review [privacy policy](#) for more information about our privacy practices.



Get this newsletter