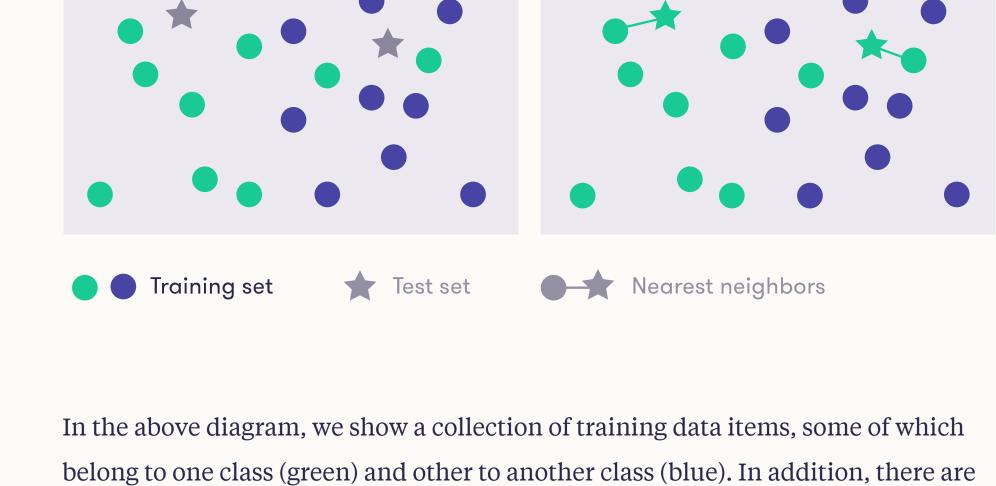
Menu <u></u>

Classifier The nearest neighbor classifier is among the simplest possible classifiers. When given an item to classify, it

II. The nearest neighbor

finds the training data item that is most similar to the new item, and outputs its label. An example is given in the following diagram.



The two test items are both classified in the "green" class because their nearest neighbors are both green (see diagram (b) above).

The position of the points in the plot represents in some way the properties of the items. Since we draw the diagram on a flat two-dimensional surface – you can

two test data items, the stars, which we are going to classify using the nearest

move in two independent directions: up-down or left-right – the items have two properties that we can use for comparison. Imagine for example representing patients at a clinic in terms of their age and blood-sugar level. But the above

diagram should be taken just as a visual tool to illustrate the general idea, which is to relate the class values to similarity or proximity (nearness). The general idea is by no means restricted to two dimensions and the nearest neighbor classifier can easily be applied to items that are characterized by many more properties than two.

What do we mean by nearest?

An interesting question related to (among other things) the nearest neighbor classifier is the definition of distance or similarity between instances. In the illustration above, we tacitly assumed that the standard geometric distance,

measure the distance between any two items by pulling a piece of thread straight

Note

Defining 'nearest'

Using the geometric distance to decide which is the nearest item may not always be

reasonable or even possible: the type of the input may, for example, be text, where it is

not clear how the items are drawn in a geometric representation and how distances

technically called the Euclidean distance, is used. This simply means that if the

points are drawn on a piece of paper (or displayed on your screen), you can

should be measured. You should therefore choose the distance metric on a case-by-case basis.

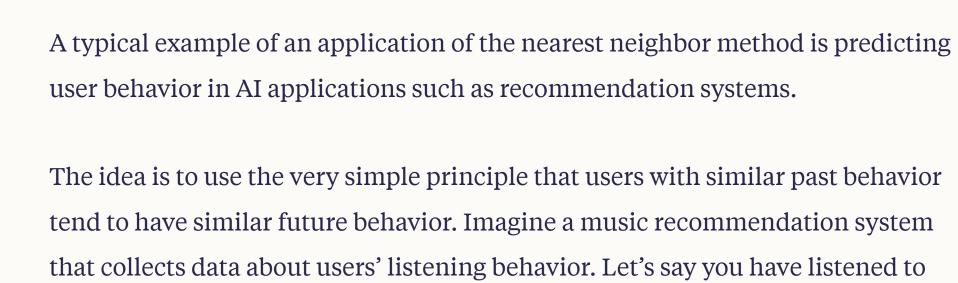
neighbor method.

In the MNIST digit recognition case, one common way to measure image similarity is to count pixel-by-pixel matches. In other words, we compare the pixels in the top-left corner of each image to one another and if the more similar color (shade of gray) they are, the more similar the two images are. We also compare the pixels in the bottom-right corner of each image, and all pixels inbetween. This technique is quite sensitive to shifting or scaling the images: if we take an image of a '1' and shift it ever so slightly either left or right, the outcome is that the two images (before and after the shift) are very different because the black pixels are in different positions

in the two images. Fortunately, the MNIST data has been preprocessed by

centering the images so that this problem is alleviated.

Music recommendations



scarce and coarse and it will only be able to give rough predictions.

Using nearest neighbors to predict user

behavior

The system now needs to predict whether you will like it or not. One way of doing this is to use information about the genre, the artist, and other metadata, entered by the good people of the service provider. However, this information is relatively

1980s disco music (just for the sake of argument). One day, the service provider gets

their hands on a hard-to-find 1980 disco classic, and adds it into the music library.

What current recommendation systems use instead of the manually entered metadata, is something called collaborative filtering. The collaborative aspect of it is that it uses other users' data to predict your preferences. The word "filter" refers to the fact that you will be only recommended content that passes through a filter: content that you are likely to enjoy will pass, other content will not (these kind of filters may lead to the so called filter bubbles, which we mentioned in Chapter 1. We will return to them later).

release and keep listening to it again and again. The system will identify the similar

past behavior that you and other 80s disco fanatics share, and since other users like

you enjoy the new release, the system will predict that you will too. Hence it will

show up at the top of your recommendation list. In an alternative reality, maybe

the added song is not so great and other users with similar past behavior as yours

don't really like it. In that case, the system wouldn't bother recommending it to you, or at least it wouldn't be at the top of the list of recommendations for you.

The following exercise will illustrate this idea.

Unanswered

Exercise 14: Customers who bought

history of four items and the item they bought after buying these four items:

In this exercise, we will build a simple recommendation system for an online shopping

headphones

Purchase

coffee

beans

socks

flip flops

Purchase

?

flip flops

sunglasses

boxing

gloves

application where the users' purchase history is recorded and used to predict which products the user is likely to buy next.

We have data from six users. For each user, we have recorded their recent shopping

Shopping History

boxing

gloves

2001: A

Space

(dvd)

Odyssey

the items have been purchased by both users.

Shopping History

green tea

User

Sanni

Henrik

sunscreen.

following products:

User

Travis

similar products

Jounit-shirtcoffee
beanscoffee
coffee makercoffee
beanscoffee
beansJaninasunglassessneakerst-shirtsneakers

Moby Dick

headphones

(novel)

Moby Dick Ville t-shirt flip flops sunglasses sunscreen (novel) 2001: A Moby Dick coffee Space coffee headphones Teemu Odyssey (novel) beans beans (dvd) The most recent purchase is the one in the rightmost column, so for example, after buying a t-shirt, flip flops, sunglasses, and Moby Dick (novel), Ville bought sunscreen.

Our hypothesis is that after buying similar items, other users are also likely to buy

To apply the nearest neighbor method, we need to define what we mean by nearest.

This can be done in many different ways, some of which work better than others. Let's

use the shopping history to define the similarity ("nearness") by counting how many of

For example, users Ville and Henrik have both bought a t-shirt, so their similarity is 1.

Note that flip flops doesn't count because we don't include the most recent purchase

Our task is to predict the next purchase of customer Travis who has bought the

when calculating the similarity — it is reserved for another purpose.

t-shirt

t-shirt

You can think of Travis being our test data, and the above six users make our training data.

sunglasses

Calculate the similarity of Travis relative to the six users in the training data (done by adding together the number of similar purchases by the users).
 Having calculated the similarities, identify the user who is most similar to Travis by selecting the largest of the calculated similarities.
 Predict what Travis is likely to purchase next by looking at the most recent

purchase (the rightmost column in the table) of the most similar user from the

What is the predicted purchase for Travis?

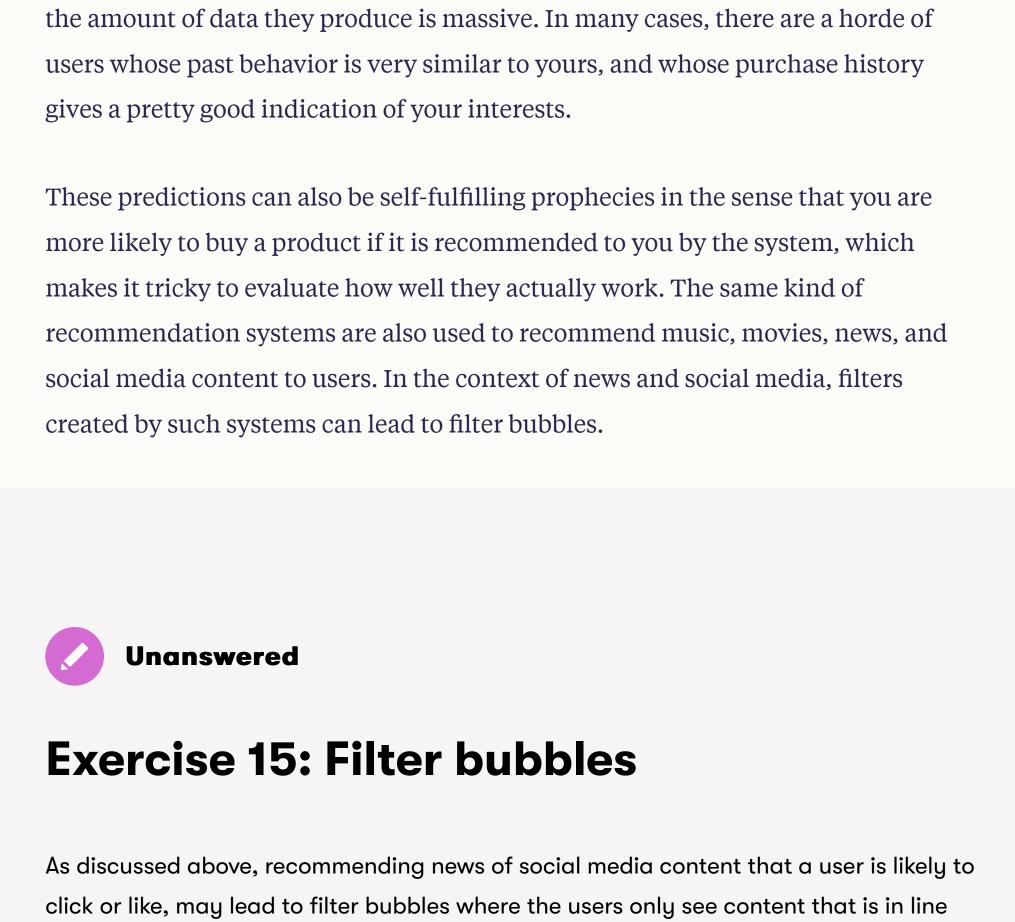
previous step.

Who is the user most similar to Travis?

 \rightarrow

Answer

Answer



In the above example, we only had six users' data and our prediction was probably

very unreliable. However, online shopping sites often have millions of users, and

may be associated with filter bubbles? Feel free to look for more information from other sources.2. Think of ways to avoid filter bubbles while still being able to recommend content to suit personal preferences. Come up with at least one suggestion. You can look

with their own values and views.

for ideas from other sources, but we'd like to hear your own ideas too!

Note: your answer should be at least a few sentences for each part.

1. Do you think that filter bubbles are harmful? After all, they are created by

recommending content that the user likes. What negative consequences, if any,

Your answer

Words: 0

Next section

III. Regression

My profile

HELSINGIN YLIOPISTO HELSINGFORS UNIVERSITET UNIVERSITY OF HELSINKI Reaktor