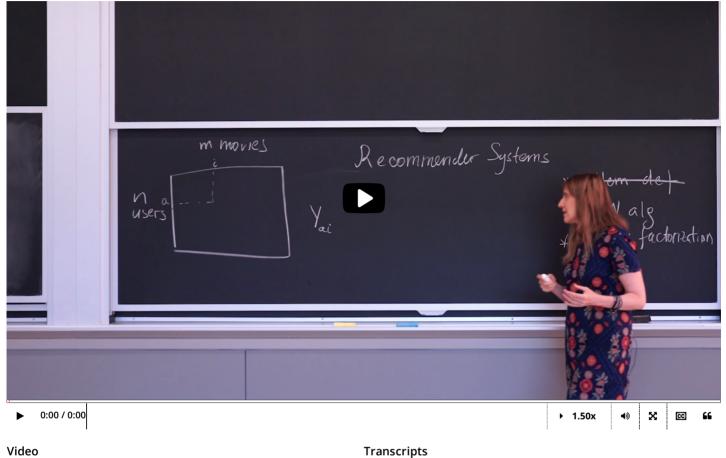


<u>Unit 2 Nonlinear Classification,</u> <u>Linear regression, Collaborative</u>

Course > Filtering (2 weeks)

> Lecture 7. Recommender Systems > 3. K-Nearest Neighbor Method

3. K-Nearest Neighbor Method K-Nearest Neighbor Method



Download video file

Transcripts

<u>Download SubRip (.srt) file</u>

<u>Download Text (.txt) file</u>

Our goal in the movie recommender system problem is to predict the movie ranking that a user would give on a movie that (s)he has not seen.

Let m be the number of movies and n the number of users. The ranking Y_{ai} of a movie $i \in \{1, \dots, m\}$ by a user $a \in \{1, \dots, n\}$ may already exist or not. Our goal is to predict Y_{ai} in the case when Y_{ai} does not exist.

K-Nearest Neighbour

The K-Nearest Neighbor method makes use of ratings by K other "similar" users when predicting Y_{ai} .

Let $\mathrm{KNN}\left(a\right)$ be the set of K users "similar to" user a, and let $\mathrm{sim}\left(a,b\right)$ be a **similarity measure** between users a and $b\in\mathrm{KNN}\left(a\right)$. The K-Nearest Neighbor method predicts a ranking Y_{ai} to be :

$$\widehat{Y}_{ai} = rac{\displaystyle\sum_{b \in ext{KNN}(a)} \sin\left(a,b
ight) Y_{bi}}{\displaystyle\sum_{b \in ext{KNN}(a)} \sin\left(a,b
ight)}.$$

The similarity measure $\sin{(a,b)}$ could be any distance function between the feature vectors x_a and x_b of users a and b, e.g. the euclidean distance $\|x_a - x_b\|$ and the cosine similarity $\cos{\theta} = \frac{x_a \cdot x_b}{\|x_a\| \|x_b\|}$. We will use these similarity measures again in *Unit 4 Unsupervised Learning*.

A drawback of this method is that the success of the K-Nearest Neighbor method depends heavily on the choice of the similarity measure. In the next section, we will discuss collaborative filtering, which will free us from the need to define a good similarity measure.

Topic: Unit 2 Nonlinear Classification, Linear regression, Collaborative Filtering (2 weeks):Lecture 7. Recommender Systems / 3. K-Nearest Neighbor Method

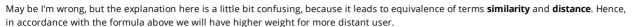
Add a Post

+

∢ All Posts

Similarity measure vs distance

discussion posted 8 months ago by Evgeny8



Probably, it worth to mention, that **similarity** and **distance** are mutually inverse measures, e.g. $\sin{(a,b)} = \frac{1}{\operatorname{dist}(a,b)+1}$

This post is visible to everyone.

Add a Response

1 response

+

•••

•••

•••

ptressel (Community TA)

8 months ago - endorsed 8 months ago by dkonomis (Staff)

That's a good point, Evgeny8!

Sometimes it may be easier to define a measure of how different two things are, than how alike they are, and then convert the difference value into a similarity value.

Hmm, one example of that might be "edit distance" -- how many one-character changes do you need to make to some text string, to convert it into another? (E.g. "cat" and "car" have an edit distance of 1, "cat" and "dam" have an edit distance of 2, "cat" and "chat" differ by 1, etc.) That's fairly straightforward, and probably easier than defining some similarity directly.

Another might be something like Euclidean distance, but when the scales of each axis don't mean the same thing, so they have to be adjusted. For instance, you might want to use the standard deviation of your data in each direction as the scale for that direction. In that case, you could use the <u>Mahalanobis distance</u>, that uses the inverse covariance matrix to rescale the axes. (Except for a constant, that's the square root of the expression in the exponent of a multivariate Gaussian.)

$$d\left(x_{1},x_{2}
ight)=\sqrt{\left(x_{1}-x_{2}
ight)^{ op}\Sigma^{-1}\left(x_{1}-x_{2}
ight)}$$

A formal, mathematical distance $d\left(x,y\right)$ has some requirements: It's non-negative: $d\left(x,y\right)\geq0$. It's zero only if and only if the two points are the same: $d\left(x,y\right)=0 \Leftrightarrow x=y$. It's symmetric: $d\left(x,y\right)=d\left(y,x\right)$. And it satisfies the triangle inequality: $d\left(x,z\right)\leq d\left(x,y\right)+d\left(y,z\right)$.

Not every measure of how different things are will have those properties. I see folks using the word "dissimilarity" instead of distance, so it doesn't exclude perfectly useful measures of difference that don't happen to be formal distances.

"A formal, mathematical distance"

It's called measure. See also similarity measure.

posted 8 months ago by Mark B2 (Community TA)

I meant distance...with those four properties. Not the same as measure. 📦

posted 8 months ago by ptressel (Community TA)

Just for fun, here's another definition of infinitesimal distance, which allows one to integrate to find the curve length: The <u>metric tensor</u>. It is frequently used in studying curved spaces such as the spacetime in general relativity, and also seems to have usage in <u>geometric deep learning</u>. I don't understand the detail as I lack knowledge in differential geometry though.

posted 8 months ago by FutureStar

Add a comment

Showing all responses

Add a response:

© All Rights Reserved