

# Pattern Recognition

SPR\_CHW2

K-Nearest Neighbors (KNN) For Classification

Mohammad Ali Basavad

40330475

## بیان مسئله

در این تمرین قصد داریم با داشتن مجموعه ای از داده های اعداد (عکس اعداد یک رقمی). با استفاده از دسته بند KNN ، یک classifier ایجاد کنیم. در این تمرین قصد داریم با روش kfold ابتدا k مناسب برای دسته بند KNN را پیدا کرده سپس میزان دقت مدل براساس این classifier را حساب کنیم

## دسته بند KNN

دسته بند KNN (k nearest neighbors) فاصله نمونه را با تمام داده های موجود محاسبه کرده و سپس به نمونه لیبل کلاسی را میزند که بیشترین نزدیکی با دیتا های آن را داشته باشد.

در واقع در این روش K داده نزدیک به نمونه پیدا شده و نمونه به کلاسی تعلق پیدا می کند که بیشترین میزان را در K داده نزدیک به نمونه داشته باشد.

برای محاسبه فاصله نمونه تا داده ها می توان از روش های مختلفی از جمله تابع cosine یا فاصله اقلیدسی استفاده کرد

$$\text{cosine} = \cos(\theta) = \frac{A \cdot B}{||A|| \cdot ||B||}$$

$$\text{Euclidean distance} = d(q, p) = \sqrt{\sum_{i=1}^n (q_i - p_i)^2}$$

## داده ها

در این تمرین از مجموعه داده openML mnist\_784، استفاده خواهیم کرد. این دیتاست شامل عکس های 784 پیکسلی از رقم های مختلف می باشد که هر پیکسل باینری می باشد.

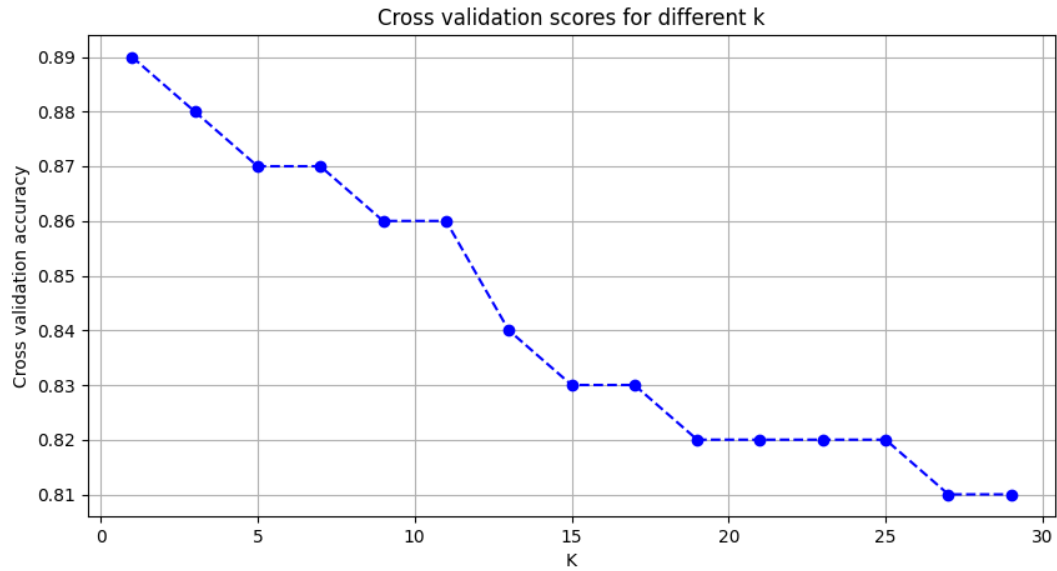
## جداسازی داده ها

ابتدا 80٪ از داده های ابتدای دیتاست را بعنوان داده های Train و باقی 20٪ را بعنوان داده های Test جداسازی کردیم.

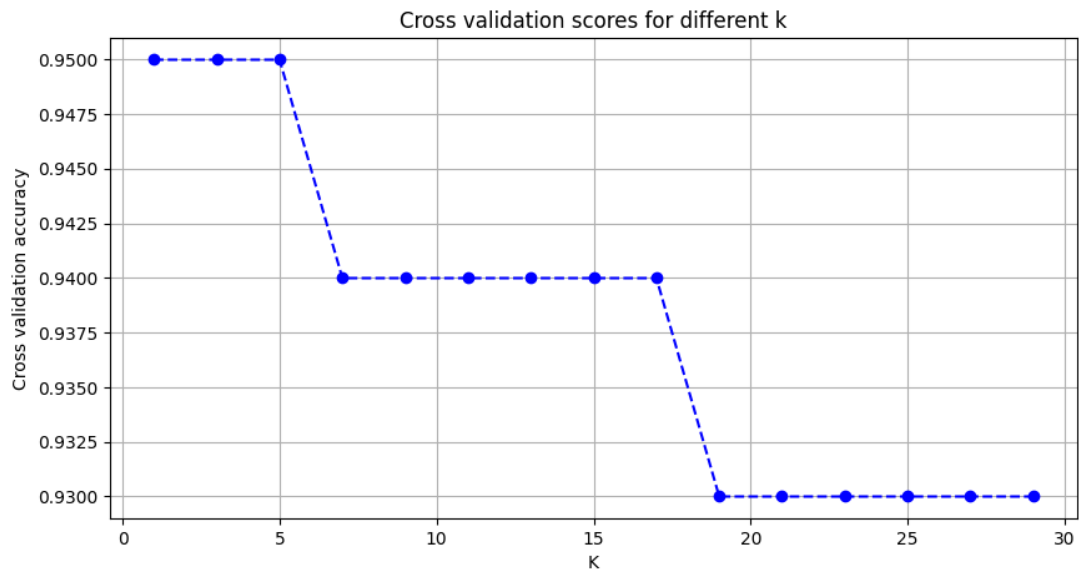
## بدست آوردن K بهینه

سپس برای پیدا کردن K بهینه با استفاده از روش kfold داده های آموزش را به 10 قسمت تبدیل کرده هر بار یک قسمت را بعنوان Test کنار گذاشته و با استفاده از باقی داده ها دسته بندی میکنیم و accuracy میانگین را برای 10 قسمت بدست آورده و اینکار را برای اعداد فرد بین 1 تا 29 بعنوان K، تکرار میکنیم و برای هر K یک accuracy پیدا می شود، در نتیجه ما K بهینه را از روی درصد دقت های بدست آمده پیدا می کنیم

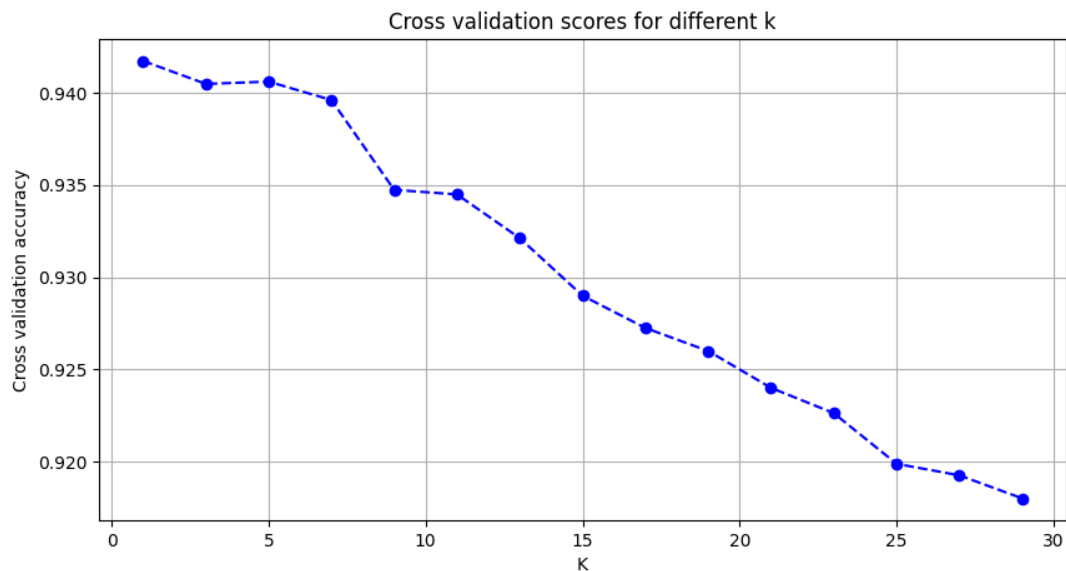
Cross validation using 1000 of datas and cosine metric



Cross validation using 10000 of datas and cosine metric



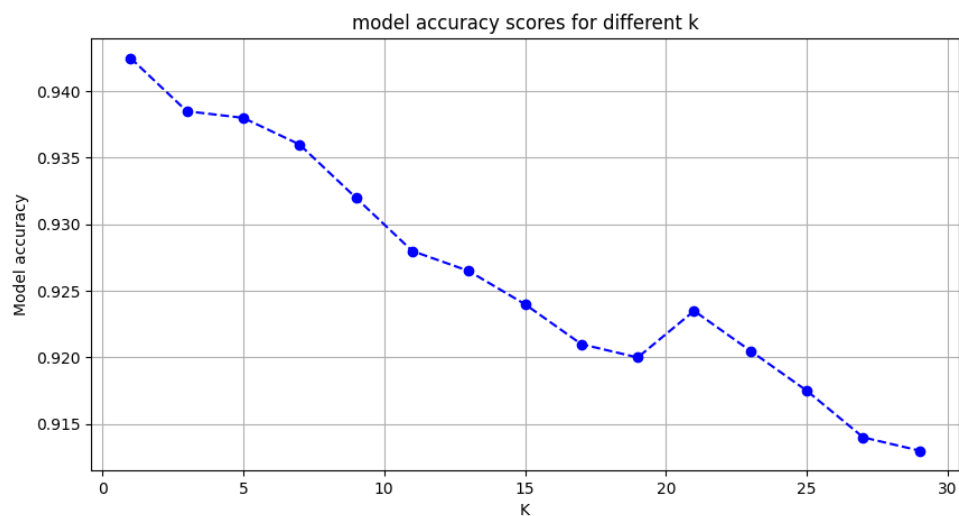
## Cross validation using 10000 of datas and Euclidean metric



## بدست آوردن دقت مدل

با توجه بدست آوردن K بهینه می توان از داده هایی که در ابتدا در دسته Test قرار دادیم برای بدست آوردن دقت نهایی مدل استفاده کرد

## Model accuracy by K



## نتیجه گیری

در نهایت می توان برای بهبود عملکرد زمانی مدل، میتوان با استفاده روش LeaveOneOut داده های Train را که به اشتباه دسته بندی می شوند را از مجموعه داده های خود حذف میکنیم. و بعد دوباره داده های Test خود را دسته بندی کنیم و میزان دقت مدل رو بسنجیم.

با استفاده از این روش زمان میانگین برای دسته بندی نمونه های جدید کاهش معنا داری پیدا می کند اما این روش مقداری از دقت الگوریتم کاهش میدهد

Accuracy in normal 1nn : `0.9425`

time in normal 1nn : `0.4150214195251465`

Accuracy in removeAmbiguousSamples 1nn : `0.939`

time in removeAmbiguousSamples 1nn : `0.33756589889526367`

After running each algorythm 50 times

Times in normal data :

[0.3538992404937744, 0.3870735168457031, 0.37909817695617676, 0.3584103584289551, 0.3645031452178955, 0.3682107925415039, 0.35208630561828613, 0.3610696792602539, 0.3524167537689209, 0.35890626907348633, 0.359588623046875, 0.3761718273162842, 0.3596780300140381, 0.35506725311279297, 0.35620665550231934, 0.349107027053833, 0.3656797409057617, 0.3680083751678467, 0.3544299602508545, 0.35216665267944336, 0.36009836196899414, 0.3586115837097168, 0.3549926280975342, 0.35550546646118164, 0.35799479484558105, 0.35588955879211426, 0.3580961227416992, 0.3530433177947998, 0.35185766220092773, 0.35704517364501953, 0.3809812068939209, 0.3900940418243408, 0.3591651916503906, 0.35320234298706055, 0.36304211616516113, 0.35268568992614746, 0.3656044006347656, 0.3523869514465332, 0.36023569107055664, 0.35240817070007324, 0.35839223861694336, 0.37230539321899414, 0.3569643497467041, 0.35579824447631836, 0.3570411205291748, 0.3600430488586426, 0.35428905487060547, 0.3559722900390625, 0.36656808853149414, 0.39223670959472656]

Time in removeAmbiguousSample :

[0.3433394432067871, 0.33911967277526855, 0.34980154037475586, 0.3541240692138672, 0.3412175178527832, 0.3474764823913574, 0.3972768783569336, 0.3671236038208008, 0.3542201519012451, 0.3452913761138916, 0.347916841506958, 0.3471405506134033, 0.36057567596435547, 0.34526586532592773, 0.35135960578918457, 0.3511793613433838, 0.349442720413208, 0.33513951301574707, 0.35129666328430176, 0.34374070167541504, 0.34174084663391113, 0.3500950336456299, 0.3433396816253662, 0.34102678298950195, 0.34658384323120117, 0.3479022979736328, 0.34037232398986816, 0.3432731628417969, 0.35001349449157715, 0.3440663814544678, 0.3547780513763428, 0.34546375274658203, 0.3423442840576172, 0.355391263961792, 0.34285712242126465, 0.340317964553833,

0.35012340545654297, 0.34327220916748047, 0.34969258308410645,  
0.34862852096557617, 0.3428654670715332, 0.3432917594909668,  
0.347675085067749, 0.34937429428100586, 0.35387516021728516,  
0.34657716751098633, 0.3445310592651367, 0.35310912132263184,  
0.3477611541748047, 0.3728015422821045]

$P\_value = 1.64 * 10^{-8}$

that shows there is meaningful difference between two methods  
and second method is much better