

En Yakın k-Komsu (k-Nearest Neighbor)

Yapay Öğrenim alanında örnek bazlı öğrenen algoritmalarından bilinen kNN, eğitim verinin kendisini sınıflama (classification) amaçlı olarak kullanır, yeni bir model ortaya çıkartmaz. Algoritma şöyle işler: etiketleri bilinen eğitim verisi alınır ve bir kenarda tutulur. Yeni bir veri noktası görüldüğünde bu veriye geri dönülür ve o noktaya “en yakın” k tane nokta bulunur. Daha sonra bu noktaların etiketlerine bakılır ve çoğunluğun etiketi ne ise, o etiket yeni noktanın etiketi olarak kabul edilir.

“En yakın” sözü bir koordinat sistemi anlamına geliyor, ve kNN, aynen k-Means ve diğer pek çok koordinatsal öğrenme yöntemi gibi eldeki çok boyutlu veri noktalarının elemanlarını bir koordinat sistemindeymiş gibi görür. Kiyasla mesela APriori gibi bir algoritma metin bazlı veriyle olduğu gibi çalışabilirdi.

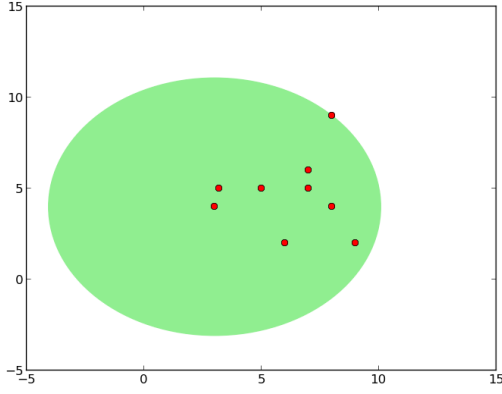
Peki arama bağlamında, bir veri obgesi içinden en yakın noktaları bulmanın en basit yolu nedir? Listeyi bastan sonra taramak (kaba kuvvet yöntemi -brute force-) ve her listedeki nokta ile yeni nokta arasındaki mesafeyi teker teker hesaplayıp en yakın k taneyi içinden seçmek bir yöntem. Bu basit algoritmanın yuku $O(N)$ 'dir. Eğer tek bir nokta arıyorsa olsaydık, bu kabul edilebilir olabilirdi. Fakat genellikle bir sınıflayıcı algoritmanın sürekli işlemesi, mesela bir online site için günde milyonlarca kez bazı kararları alması gerekebilir. Bu durumda N 'in çok büyük olduğu şartlarda, üstteki hız bile yeterli olmayabilir.

Arama işlemini daha hızlı yapmanın yolları var. Akıllı arama algoritmaları kullanarak eğitim verilerini bir ağac yapısı üzerinden tarayıp erişim hızını $O(\log N)$ 'e indirmek mümkündür.

Küre Ağaçları (Ball Tree, BT)

Bir noktanın diğer noktalara yakın olup olmadığının hesabında yapılması gereken en pahalı işlem nedir? Mesafe hesabıdır. BT algoritmasının puf noktası bu hesabi yapmadan, noktalara değil, noktaları kapsayan “kürelere” bakarak hız kazandırmasıdır. Noktaların her biri yerine o noktaları temsil eden kürenin mihenk noktasına (pivot -bu nokta küre içindeki noktaların ortalaması olarak merkezi de olabilir, herhangi bir başka nokta da-) bakılır, ve oraya olan mesafeye göre bir küre altındaki noktalara olabilecek en az ve en fazla uzaklık hemen anlaşılmış olur.

Mesela elimizde alttaki gibi noktalar var ve küreyi oluşturduk.

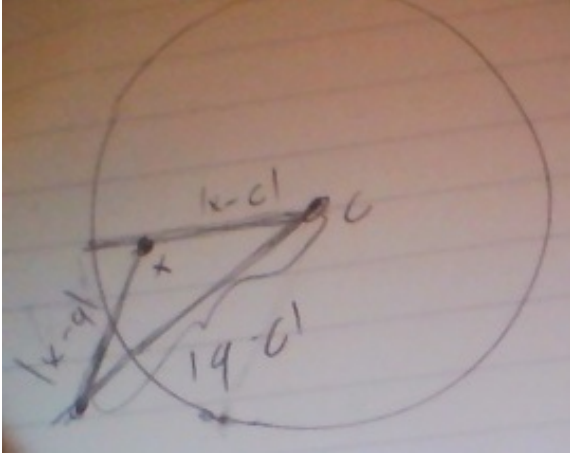


Bu kureyi kullanarak kure disindaki herhangi bir nokta q 'nun kuredeki “diger tum noktalar x 'e” olabilecegi en az mesafenin ne olacagini ucgensel esitsizlik ile anlayabiliriz.

Ucgensel esitsizlik

$$|x - y| \leq |x - z| + |z - y|$$

$||$ operatoru norm operatoru anlamina gelir ve uzaklik hesabının genelleştirilmiş halidir. Konu hakkında daha fazla detay için *Fonksinel Analiz* ders notlarımıza bakabilirsiniz. Kisaca söylenmek istenen iki nokta arasında direk gitmek yerine yolu uzatırsak, mesafe artacaktır. Tabii uzaklık, yol, nokta gibi kavramlar tamamen soyut matematiksel ortamda da işleyecek şekilde ayarlanmıştır. Mesela mesafe (norm) kavramını değiştirebiliriz, Oklitsel yerine Manhattan mesafesi kullanırız, fakat bu kavram bir norm olduğu ve belirttiğimiz uzayda geçerli olduğu için ucgensel esitsizlik üzerine kurulmuş tüm diğer kurallar geçerli olur.



Simdi diyelim ki disaridaki bir q noktasından bir kure icindeki diger tum x noktalarına olan mesafe hakkında bir seyler soylemek istiyoruz. Ustteki sekilden bir ucgensel esitsizlik cikartabiliriz,

$$|x - c| + |x - q| \geq |q - c|$$

Bunun dogru bir ifade oldugunu biliyoruz. Peki simdi yaricapi bu ise dahil edelim, cunku yaricap hesabi bir kere yapilip kure seviyesinde depolanacak ve bir daha hesaplanmasi gerekmeyecek, yani algoritmayi hizlandiracak bir sey olabilir bu, o zaman eger $|x - c|$ yerine yaricapi kullanirsak, esitsizlik hala gecerli olur, sol taraf zaten buyuktu, simdi daha da buyuk olacak,

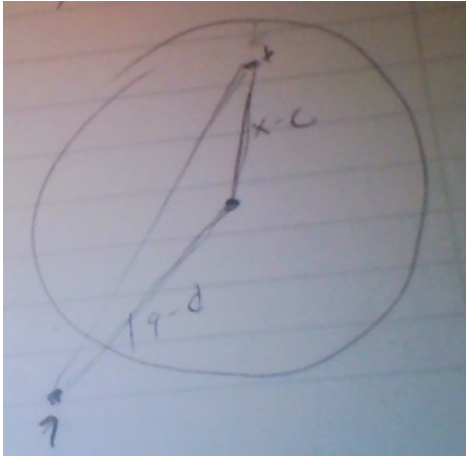
$$radius + |x - q| \geq |q - c|$$

Bunu nasil boyle kesin bilebiliyoruz? Cunku BT algoritmasi radius'u $|x - c|$ 'ten kesinlikle daha buyuk olacak sekilde secer). Simdi yaricapi saga gecirelim,

$$|x - q| \geq |q - c| - radius$$

Boylece guzel bir tanim elde ettik. Yeni noktanin kuredeki herhangi bir nokta x 'e olan uzakligi, yeni noktanin mihenke olan uzakliginin yaricapi cikartilmis halinden *muhakkak* fazladir. Yani bu cikartma isleminde ele gecen rakam yeni noktanin x 'e uzakligina bir “alt sinir (lower bound)” olarak kabul edilebilir. Diger tum mesafeler bu rakamdan daha buyuk olacaktir. Ne elde ettik? Sadece bir yeni nokta, mihenk ve yaricap kullanarak kuredeki “diger tum noktalar hakkında” bir irdeleme yapmamiz mumkun olacak. Bu noktalara teker teker bakmamiz gerekmeyecek. Bunun nasil ise yaradigini algoritma detaylarinda gorecegiz.

Benzer sekilde



Bu ne diyor?

$$|q - c| + |x - c| \geq |q - x|$$

$|x - c|$ yerine yaricap kullanirsak, sol taraf buyuyecegi icin buyukluk hala buyukluk olarak kalir,

$$|q - c| + radius \geq |q - x|$$

Ve yine daha genel ve hizli hesaplanan bir kural elde ettik (onceki ifadeye benzemesi icin yer duzenlemesi yapalim)

$$|q - x| \leq |q - c| + radius$$

Bu ifade ne diyor? Yeni noktanin mihenke olan uzakligina yaricap “eklenirse” bu

uzaklıktan, büyüklükten daha büyük bir yeni nokta / kure mesafesi olamaz, kuredeki hangi nokta olursa olsun. Bu esitsizlik te bize bir üst sınır (upper bound) vermiş oldu.

Algoritma

Kure Ağaçları (BT) metodu önce kureleri, ağaçları oluşturmaktadır. Bu kureler hiyerarşik şekilde planlanır, tüm noktaların içinde olduğu bir “en üst kure” vardır her kurenin iki tane çocuk kuresi olabilir. Belli bir (disaridan tanımlanan) minimum r_{min} veri noktasına gelinceye kadar sadece noktaları geometrik olarak kapsamakla ilgili kureler oluşturulur, kureler noktaları sahiplenmezler. Fakat bu r_{min} sayısına erişince (artık oldukça alttaki) kurelerin üzerine noktalar konacaktır.

Önce tek kurenin oluşturulmasına bakalım. Bir kure oluşumu için eldeki veri içinden herhangi bir tanesi mihenk olarak kabul edilebilir. Daha sonra bu mihenkten diğer tüm noktalara olan uzaklık ölçülür, ve en fazla, en büyük olan uzaklık yarıçap olarak kabul edilir (her şeyi kapsayabilmesi için).

Not: Bu arada “tüm diğer noktalara bakılması” dedik, bundan kaçınmaya çalışmıyor muyduk? Fakat dikkat, “kure oluşturulması” evresindeyiz, k tane yakın nokta arama evresinde değiliz. Yapmaya çalıştığımız aramaları hızlandırmak - eğitim / kure oluşturması bir kez yapılacak ve bu eğitilmiş kureler bir kenarda tutulacak ve sürekli aramalar için ardi ardına kullanılacaklar.

Kureyi oluşturma algoritması şöyledir: verilen noktalar içinde herhangi birisi mihenk olarak seçilir. Sonra bu noktadan en uzakta olan nokta f_1 , sonra f_1 'den en uzakta olan nokta f_2 seçilir. Sonra tüm noktalara teker teker bakılır ve f_1 'e yakın olanlar bir gruba, f_2 'ye yakın olanlar bir gruba ayrılır.

```
import pprint
import numpy as np
import dist

__rmin__ = 2

# node: [pivot, radius, points, [child1, child2]]
def new_node(): return [None, None, None, [None, None]]

def zero_if_neg(x):
    if x < 0: return 0
    else: return x

def form_tree(points, node):
    pivot = points[0]
    radius = np.max(dist.dist(points, pivot))
    node[0] = pivot
    node[1] = radius
    if len(points) <= __rmin__:
        node[2] = points
        return
    idx = np.argmax(dist.dist(points, pivot))
    furthest = points[idx, :]
    idx = np.argmax(dist.dist(points, furthest))
```

```

    furthest2 = points[idx,:]
    dist1=dist.dist(points,furthest)
    dist2=dist.dist(points,furthest2)
    diffs = dist1-dist2
    p1 = points[diffs <= 0]
    p2 = points[diffs > 0]
    node[3][0] = new_node() # left child
    node[3][1] = new_node() # right child
    form_tree(p1,node[3][0])
    form_tree(p2,node[3][1])

# knn: [min_so_far, [points]]
def search_tree(new_point, knn_matches, node, k):
    pivot = node[0]
    radius = node[1]
    node_points = node[2]
    children = node[3]

    # calculate min distance between new point and pivot
    # it is direct distance minus the radius
    min_dist_new_pt_node = dist.norm(pivot,new_point) - radius

    # if the new pt is inside the circle, its potential minimum
    # distance to a random point inside is zero (hence
    # zero_if_neg). we can only say so much without looking at all
    # points (and if we did, that would defeat the purpose of this
    # algorithm)
    min_dist_new_pt_node = zero_if_neg(min_dist_new_pt_node)

    knn_matches_out = None

    # min is greater than so far
    if min_dist_new_pt_node >= knn_matches[0]:
        # nothing to do
        return knn_matches
    elif node_points != None: # if node is a leaf
        print knn_matches_out
        knn_matches_out = knn_matches[:] # copy it
        for p in node_points: # linear scan
            if dist.norm(new_point,p) < radius:
                knn_matches_out[1].append([list(p)])
                if len(knn_matches_out[1]) == k+1:
                    tmp = [dist.norm(new_point,x) \
                           for x in knn_matches_out[1]]
                    del knn_matches_out[1][np.argmax(tmp)]
                    knn_matches_out[0] = np.min(tmp)
        else:
            dist_child_1 = dist.norm(children[0][0],new_point)
            dist_child_2 = dist.norm(children[1][0],new_point)
            node1 = None; node2 = None
            if dist_child_1 < dist_child_2:
                node1 = children[0]
                node2 = children[1]
            else:

```


caktır. Çocuk noktaları kararlaştırıldıktan sonra artık sonraki aşamaya geçilir, fonksiyon `form_tree` bu çocuk noktaları alarak, ayrı ayrı, her çocuk grubu için özyineli (recursive) olarak kendi kendini çağırır. Kendi kendini çağırarak `form_tree`, tekrar başladığında kendini yeni (bir) nokta grubu ve yeni bir düğüm objesi ile bas-basa bulur, ve hiçbir şeyden habersiz olarak işleme koyulur. Tabii her özyineli çağrı yeni düğüm objesini yaratırken bir referansı üstteki ebeveyn düğümüne koymayı unutmamıştır, böylece özyineli fonksiyon dünyadan habersiz olsa bile, ağacın en üstünden en altına kesintisiz bir bağlantı zinciri hep elimizde olur.

Not: `form_tree` içinde bir numara yaptık, tüm noktaların f_1 'e olan uzaklığı `dist1`, f_2 'e olan uzaklığı ise `dist2`. Sonra `diffs = dist1-dist2` ile bu iki uzaklığı birbirinden çıkartıyoruz ve mesela `points[diffs <= 0]` ile f_1 'e yakın olanları buluyoruz, çünkü bir tarafta f_1 'e yakınlık 4 diğer tarafta f_2 'ye yakınlık 6 ise, $4-6=-2$ ie o nokta f_1 'e yakın demektir. Ufak bir numara ile Numpy dilimleme (slicing) tekniğini kullanabilmiş olduk ve bu önemli çünkü böylece `for` döngüsü yazmıyoruz, Numpy'in arka planda C ile yazılmış hızlı rutinlerini kullanıyoruz.

Arama

Üstteki imaj içinde *BallKNN* olarak gösterilen ve bizim kodda `search_tree` olarak anılan fonksiyon arama fonksiyonudur. Aranan `new_point`'e olan k en yakın diğer veri noktaları. Disaridan verilen değışken `knn_matches` üzerinde fonksiyon özyineli bir şekilde arama yaparken “o ana kadar bulunmuş en yakın k nokta” ve o noktaların `new_point`'e olan en yakın mesafesi saklanır, arama işleyişi sırasında `knn_matches`, `knn_matches_out` sürekli verilip geri dondurulan değışkenlerdir, İngilizce koddaki P^{in}, P^{out} 'un karşılığıdır.

Arama algoritması şöyle işler: şimdi önceden oluşturulmuş küre hiyerarşisini üstten alta doğru gezmeye başlarız. Her basamakta yeni nokta ile o kürenin mihenkini, yarıçapını kullanarak bir “alt sınır mesafe hesabı” yaparız, bu mesafe hesabının arkasında yatan düşünceyi yazının başında anlatmıştık. Bu mesafe küre içindeki tüm noktalara olan bir en az mesafe idi, ve eğer eldeki `knn_matches` üzerindeki şimdiye kadar bulunmuş mesafelerin en azından daha az ise, o zaman bu küre “bakmaya değer” bir küredir, ve arama algoritması bu küreden işleme devam eder. Şimdiye kadar bulunmuş mesafelerin en azı `knn_matches` veri yapısı içine `min_so_far` olarak saklanıyor, İngilizce kodda D_{sofar} .

Bu irdeleme sonrası (yani vs küresinden yola devam kararı arkasından) işleme iki şekilde devam edilebilir, çünkü bir küre iki türden olabilir; ya nihai en alt kürelerden biridir ve üzerinde gerçek noktalar depolanmıştır, ya da ara kürelerden biridir (sona gelmedik ama doğru yoldayız, daha alta inmeye devam), o zaman fonksiyon yine özyineli bir şekilde bu kürenin çocuklarına bakacaktır - her çocuk için kendi kendini çağıracaktır. İkinci durumda, kürede noktalar depolanmıştır, artık basit lineer bir şekilde o tüm noktalara teker teker bakılır, eldekilerden daha yakın olanı alınır, eldeki liste sısmeye başlamışsa (k 'den daha fazla ise) en büyük noktalardan biri atılır, vs.

Daha alta inmemiz gereken birinci durumda yapılan iki çağrının bir özelliğine dikkat

cekmek isterim. Yeni noktanin bu cocuklara olan uzakligi da olculuyor, ve en once, en yakin olan cocuga dogru bir ozyineleme yapiliyor. Bu nokta cok onemli: niye boyle yapildi? Cunku icinde muhtemelen daha yakin noktalarin olabilecegi kurelere dogru gidersek, ozyineli cagrilarin teker teker bitip yukari dogru cikmaya baslamasi ve kaldiklari yerden bu sefer ikinci cocuk cagrilarini yapmaya baslamasi ardindan, elimizdeki `knn_matches` uzerinde en yakin noktalarini buyuk bir ihtimalle zaten bulmus olacagiz. Bu durumda ikinci cagri yapilsa bile tek bir alt sinir hesabi o kurede dikkate deger hicbir nokta olamayacagini ortaya cikaracak (cunku en iyiler zaten elimizde), ve ikinci cocuga olan cagrilar hic alta inmeden pat diye geri donecektir, hic asagi inilmeyecektir.

Bu muthis bir kazanimdir: zaten bu stratejiye liteturde “budamak (pruning)” adi veriliyor, bu da cok uygun bir kelime aslinda, cunku agaclarla ugrasiyoruz ve bir dugum (kure) ve onun altindaki hicbir alt kureye ugramaktan kurtularak o dallarin tamamini bir nevi “budamis” oluyoruz. Bir suru gereksiz islemden de kurtuluyoruz bu arada, ve aramayi hizlandiriyoruz.

Kaynaklar

[1] Liu, Moore, Gray, *New Algorithms for Efficient High Dimensional Non-parametric Classification*