

## En Yakın k-Komsu (k-Nearest Neighbor)

Yapay Ogrenim 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 sorgulunca bu veriye geri donulur 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. Mesela elde 1 kategorisi altında  $[2 \ 2]$ , 2 kategorisi altında  $[5 \ 5]$  var ise, yeni nokta  $[3, \ 3]$  için yakınlık açısından  $[2 \ 2]$  bulunmalı ve etiket olarak 1 sonucu döndürülmelidir.

Ustte tarif edilen basit bir ihtiyaç, yöntem gibi görülebilir. Fakat yapay öğrenim ve yapay zeka çok boyutlarda oruntu tanıma (pattern recognition) ile uğraşır, ve milyonlarca satırlık veri, onlarca boyut (ustteki örnekte 2, fakat çoğunlukla çok daha fazla boyut vardır) işler hakikaten zorlaşabilir. Mesela görüntü tanımada veri  $M \times N$  boyutundaki dijital imajlar (düzleştirilince  $M \cdot N$  boyutunda), ve onların içindeki resimlerin kime ait olduğu etiket bilgisi olabilir. kNN bu tür multimedia, çok boyutlu veri ortamında başarılı şekilde çalışabilmektedir. Ayrıca en yakın k komşunun içeriği tarifsel bilgi çıkarımı (knowledge extraction) amacıyla da kullanılabilir [2].

“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-) listedeki her nokta ile yeni nokta arasındaki mesafeyi teker teker hesaplayıp en yakın k taneyi içinden seçerdi, bu bir yöntemdir.. Bu basit algoritmanın yuku  $O(N)$ 'dir. Eğer tek bir nokta arıyor olsaydık, kabul edilebilir olabilirdi. Fakat genellikle bir sınıflayıcı (classifier) algoritmasının sürekli işlemesi, mesela bir online site için günde milyonlarca kez bazı kararları alması gerekebilir. Bu durumda ve  $N$ 'in çok büyük olduğu şartlarda, ustteki hız bile yeterli olmayacaktır.

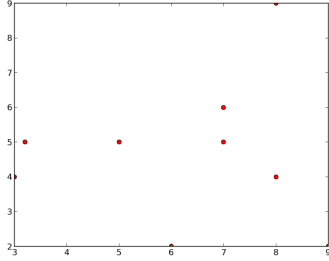
Arama işlemini daha hızlı yapmanın yolları var. Akıllı arama algoritmaları kullanarak eğitim verilerini bir ağac yapısı üzerinden tarayıp erişim hızını  $O(\log N)$ 'e indirmek mümkündür.

## 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 “kurelere” bakarak hız kazandırmasıdır. Noktaların her biri yerine o noktaları temsil eden kurenin mihenk noktasına (pivot -bu nokta kure içindeki noktaların ortalamasal olarak merkezi de olabilir, herhangi bir başka nokta da-) bakılır, ve oraya olan mesafeye göre bir kure altındaki noktalara olabilecek en az ve en fazla uzaklık hemen anlaşılmış olur.

Not: Kure kavrami uc boyutta anlamlı tabii ki, iki boyutta bir cemberden bahsetmek lazım, daha yüksek boyutlarda ise merkezi ve çapı olan bir “hiper yüzeyden” bahsetmek lazım. Tarifi kolaylaştırdığı için cember ve kure tanımlarını kullanıyoruz.

Mesela elimizde alttaki gibi noktalar var ve kureyi oluşturduk.

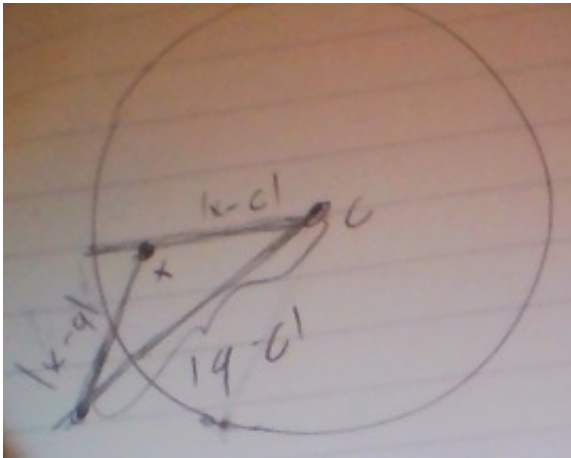


Bu kureyi kullanarak kure dışındaki herhangi bir nokta  $q$ 'nın kuredeki “diğer tüm noktalar  $x$ 'e” olabileceği en az mesafenin ne olacağını uçgensel eşitsizlik ile anlayabiliriz.

Uçgensel eşitsizlik

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

$\|$  operatörü norm operatörü anlamına gelir ve uzaklık hesabının genelleştirilmiş halidir. Konu hakkında daha fazla detay için *Fonksiyonel Analiz* ders notlarımıza bakabilirsiniz. Kısaca 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 uçgensel eşitsizlik üzerine kurulmuş tüm diğer kurallar geçerli olur.



Şimdi diyelim ki dışarıdaki bir  $q$  noktasından bir kure içindeki diğer tüm  $x$  noktalarına olan mesafe hakkında bir şeyler söylemek istiyoruz. Üstteki şekilde 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,

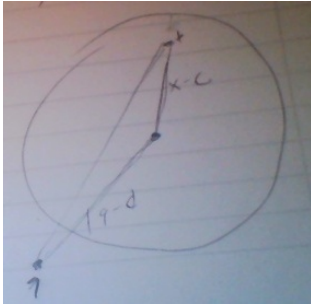
$$\text{radius} + |x - q| \geq |q - c|$$

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

$$|x - q| \geq |q - c| - \text{radius}$$

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 olacaktır. 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 gerek-meyecek. 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,

$$|q - c| + \text{radius} \geq |q - x|$$

---

```

ball_knn( $PS^{in}, node$ )
1   – Eger alttaki sart gecerli ise node icindeki bir noktanin daha once
2   – kesfedilmis k en yakin komsudan daha yakin olmasi imkansizdir
3   if  $D_{minp}^{node} \geq D_{sofar}$ 
4       return  $PS_{in}$  degismemis halde;
5   else if node bir cocuk noktasi ise
6        $PS_{out} = PS_{in}$ ;
7       for  $\forall x \in points(node)$ 
8           if  $(|x - q| < D_{sofar})$ ; – basit lineer arama yap
9            $x$ 'i  $PS_{out}$ 'a ekle;
10          if  $|PS^{out}| == k + 1$ ;
11              en uzak olan komsuyu  $PS^{out}$ 'tan cikart;
12               $D_{sofar}$ 'i guncelle;
13
14   – eger uc nokta degil ise iki cocuk dugumden daha yakin olanini
15   – incele, sonra daha uzakta olanina bak. buyuk bir ihtimalle
16   – arama devam ettirilirse bu arama kendiliginden kesilecektir
17   else
18        $node_1 = node$ 'un  $q$ 'ya en yakin cocugu;
19        $node_2 = node$ 'un  $q$ 'dan en uzak cocugu;
20        $PS^{temp} = ball\_knn(PS^{in}, node_1)$ ;
21        $PS^{out} = ball\_knn(PS^{temp}, node_2)$ ;

```

---

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

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

Bu ifade ne diyor? Yeni noktanin mihenke olan uzakligina yaricap “eklenirse” bu uzakliktan, buyuklukten daha buyuk bir yeni nokta / kure mesafesi olamaz, kuredeki hangi nokta olursa olsun. Bu esitsizlik te bize bir ust sinir (upper bound) vermis oldu.

Algoritma

```

Procedure BallKNN ( $PS^n, Node$ )
begin
  if ( $D_{\min}^{Node} \geq D_{\text{sofar}}$ ) then                                /* If this condition is satisfied, then impossible
    Return  $PS^n$  unchanged.                                     for a point in Node to be closer than the
                                                                previously discovered  $k^{th}$  nearest neighbor.*/
  else if (Node is a leaf)
     $PS^{out} = PS^n$ 
     $\forall x \in Points(Node)$ 
    if ( $\|x - q\| < D_{\text{sofar}}$ ) then                                /* If a leaf, do a naive linear scan */
      add  $x$  to  $PS^{out}$ 
      if ( $\|PS^{out}\| == k + 1$ ) then
        remove furthest neighbor from  $PS^{out}$ 
        update  $D_{\text{sofar}}$ 
  else                                                         /* If a non-leaf, explore the nearer of the two
     $node_1 = \text{child of Node closest to } q$                     child nodes, then the further. It is likely that
     $node_2 = \text{child of Node furthest from } q$                 further search will immediately prune itself.*/
     $PS^{temp} = \text{BallKNN}(PS^n, node_1)$ 
     $PS^{out} = \text{BallKNN}(PS^{temp}, node_2)$ 
end

```

Kure Agaclari (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 görevli 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 oluşturma 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 itertools

def dist(vect, x):
    return np.fromiter(itertools.imap
                        (np.linalg.norm, vect-x), dtype=np.float)

def norm(x, y): return np.linalg.norm(x-y)

points = np.array([[3., 3.], [2., 2.]])
q = [1., 1.]

```

```

print 'diff', points-q
print 'dist', dist(points,q)

diff [[ 2.  2.]
      [ 1.  1.]]
dist [ 2.82842712  1.41421356]

# k-nearest neighbor Ball Tree algorithm in Python
import pprint

__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,all_points,plot_tree=False):
    pivot = points[0]
    radius = np.max(dist(points,pivot))
    if plot_tree: plot_circles(pivot, radius, points, all_points)
    node[0] = pivot
    node[1] = radius
    if len(points) <= __rmin__:
        node[2] = points
        return
    idx = np.argmax(dist(points,pivot))
    furthest = points[idx,:]
    idx = np.argmax(dist(points,furthest))
    furthest2 = points[idx,:]
    dist1=dist(points,furthest)
    dist2=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],all_points)
    form_tree(p2,node[3][1],all_points)

# 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 = 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 norm(new_point,p) < radius:
            knn_matches_out[1].append([list(p)])
            if len(knn_matches_out[1]) == k+1:
                tmp = [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 = norm(children[0][0],new_point)
        dist_child_2 = norm(children[1][0],new_point)
        node1 = None; node2 = None
        if dist_child_1 < dist_child_2:
            node1 = children[0]
            node2 = children[1]
        else:
            node1 = children[1]
            node2 = children[0]

    knn_tmp = search_tree(new_point, knn_matches, node1, k)
    knn_matches_out = search_tree(new_point, knn_tmp, node2, k)

    return knn_matches_out

points = np.array([[3.,4.],[5.,5.],[9.,2.],[3.2,5.],[7.,5.],
                  [8.,9.],[7.,6.],[8,4],[6,2]])
tree = new_node()
form_tree(points,tree,all_points=points)
pp = pprint.PrettyPrinter(indent=4)
print "tree"
pp.pprint(tree)
newp = np.array([7.,7.])
dummysp = [np.Inf,np.Inf] # it should be removed immediately
res = search_tree(newp,[np.Inf, [dummysp]], tree, k=2)
print "done", res

tree
[ array([ 3.,  4.]),
  7.0710678118654755,
  None,
  [ [ array([ 8.,  9.]),
      3.1622776601683795,
```

```

        array([[ 8.,  9.],
[ 7.,  6.])),
        [None, None]],
        [
            array([ 3.,  4.]),
            6.324555320336759,
            None,
            [
                [
                    array([ 9.,  2.]),
                    3.6055512754639891,
                    None,
                    [
                        [
                            array([ 7.,  5.]),
                            1.4142135623730951,
                            array([[ 7.,  5.],
[ 8.,  4.])),
                            [None, None]],
                        [
                            array([ 9.,  2.]),
                            3.0,
                            array([[ 9.,  2.],
[ 6.,  2.])),
                            [None, None]]]]],
                        [
                            array([ 3.,  4.]),
                            2.2360679774997898,
                            None,
                            [
                                [
                                    array([ 5.,  5.]),
                                    0.0,
                                    array([[ 5.,  5.])),
                                    [None, None]],
                                [
                                    array([ 3.,  4.]),
                                    1.019803902718557,
                                    array([[ 3.,  4.],
[ 3.2,  5. ])),
                                    [None, None]]]]]]],
                                [None, None]]]]]]]]
None
done [1.0, [[[8.0, 9.0]], [[7.0, 6.0]]]]

```

Bu iki grup, o anda islemekte oldugumuz agac dugumun (node) iki cocuklari olacaktır. Çocuk noktaları kararlaştırıldıktan sonra artık sonraki asamaya gecilir, fonksiyon `form_tree` bu çocuk noktaları alarak, ayrı ayrı, her çocuk grubu için ozyineli (recursive) olarak kendi kendini çağırır. Kendi kendini çağırarak `form_tree`, tekrar başladığında kendini yeni (bir) nokta grubu ve yeni bir dugum objesi ile basbasa bulur, ve hiçbir seyden habersiz olarak isleme koyulur. Tabii her ozyineli çağrı yeni dugum objesini yaratırken bir referansi ustteki ebeveyn dugume koymayı unutmamıştır, böylece ozyineli 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) teknigini kullanabilmiş olduk ve bu önemli çünkü böylece `for` dongusu yazmıyoruz, Numpy'in arka planda C ile yazılmış hızlı rutinlerini kullanıyoruz.

Ek bazı bilgiler: kurelerin sinirlari kesisebilir.



## Arama

Ustte sozde program (pseudocode) BallKNN olarak gosterilen ve bizim kodda `search_tree` olarak anilan fonksiyon arama fonksiyonu. Aranan `new_point`'e olan  $k$  en yakin diger veri noktalar. Disaridan verilen degisken `knn_matches` uzerinde fonksiyon ozyineli bir sekilde arama yaparken "o ana kadar bulunmus en yakin  $k$  nokta" ve o noktalarin `new_point`'e olan en yakin mesafesi saklanir, arama isleyisi sirasinda `knn_matches`, `knn_matches_out` surekli verilip geri dondurulen degiskenlerdir, sozde programdaki  $P^{in}$ ,  $P^{out}$ 'un karsiligidirlar.

Arama algoritmasi soyle isler: simdi onceden olusturulmus kure hiyerarisisini ustten alta dogru gezmeye baslariz. Her basamakta yeni nokta ile o kurenin mi-henkini, yaricapini kullanarak bir "alt sinir mesafe hesabi" yapariz, bu mesafe hesabinin arkasinda yatan dusunceyi yazinin basinda anlatmistik. Bu mesafe kure icindeki tum noktalara olan bir en az mesafe idi, ve eger eldeki `knn_matches` uzerindeki simdiye kadar bulunmus mesafelerin en azindan daha az ise, o zaman bu kure "bakmaya deger" bir kuredir, ve arama algoritmasi bu kureden isleme devam eder. Simdiye kadar bulunmus mesafelerin en azi `knn_matches` veri yapisi icine `min_so_far` olarak saklaniyor, sozde programdaki  $D_{sofar}$ .

Bu irdeleme sonrasi (yani vs kuresinden yola devam karari arkasindan) isleme iki sekilde devam edilebilir, cunku bir kure iki turden olabilir; ya nihai en alt kurelerden biridir ve uzerinde gercek noktalar depolanmistir, ya da ara kurelerden biridir (sona gelmedik ama dogru yoldayiz, daha alta inmeye devam), o zaman fonksiyon yine ozyineli bir sekilde bu kurenin cocuklarina bakacaktır - her cocuk icin kendi kendini cagiracaktır. Ikinci durumda, kurede noktalar depolanmistir, artik basit lineer bir sekilde o tum noktalara teker teker bakilir, eldekilerden daha yakin olani alinir, eldeki liste sismeye baslamissa ( $k$ 'den daha fazla ise) en buyuk noktalardan biri atilir [3], vs.

Daha alta inmemiz gereken birinci durumda yapilan iki cagrinin bir ozelligine 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 noktalar 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.

## Mesafeler

Algoritmanın mesafeleri anlatan kısmında norm ve uzaylar gibi kavramlardan bahsettik. Yeni noktanın mihenk olan uzaklığının o küre içindeki tüm diğer noktalara olan uzaklığını temsil edebileceğini söyledik: peki niye bu kavramları direk bu şekilde anlatmadık, ve norm, üçgensel eşitsizlik gibi kavramlardan bahsettik? Çünkü 2 ve 3 boyut sonrası uzayları görsel olarak düşünmek mümkün değildir, istediğimiz kadar ellerimizi kollarımızı sallayalım, bu kavramları görsel olarak tarif edemeyiz, ve değişik bir norm (mesafe) ölçütü kullanmayı seçebiliriz. Bu her iki durumda da elimizde soyut matematik bağlamında sağlam bir temel olduğunu bilmek algoritmanın genelliğini, ve değişik şartlarda uygulanabilirliğini artırır. Mesela Oklit mesafesi yerine Manhattan mesafesi kullansam bile, bu mesafenin ölçütünün norm kurallarını uydugunu bildiğim için kNN yapısının geri kalanını olduğu gibi kullanabilirim, çünkü o yapının geçerliliğini normlar üzerinde geçerli üçgensel eşitsizlik üzerinde ispat ettim.

## Model

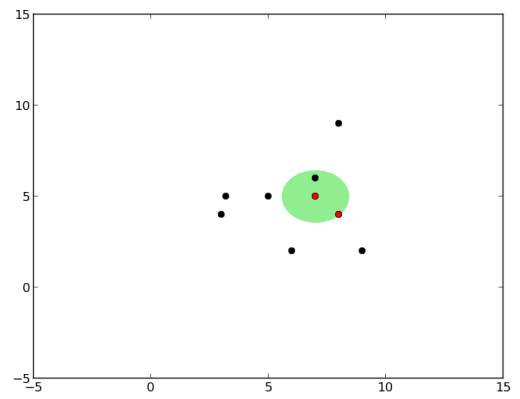
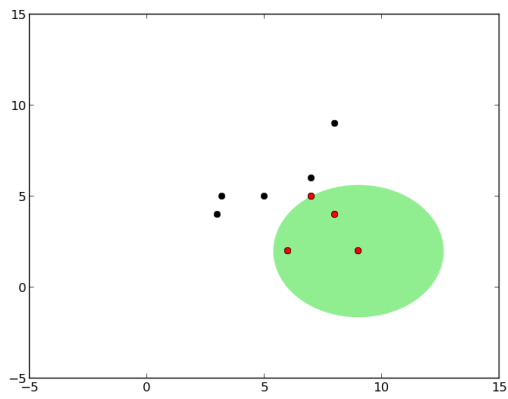
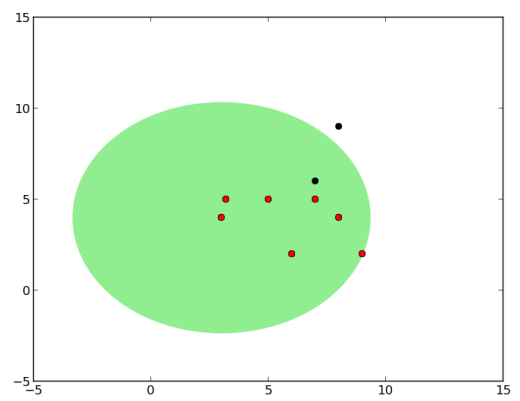
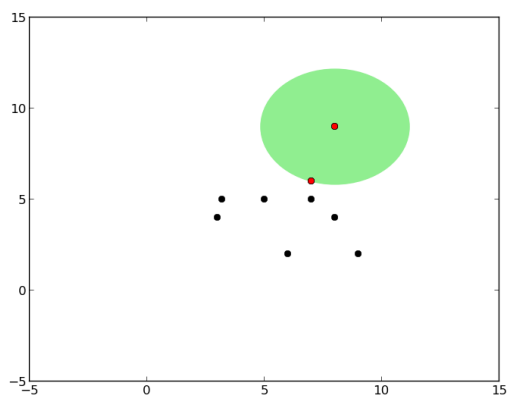
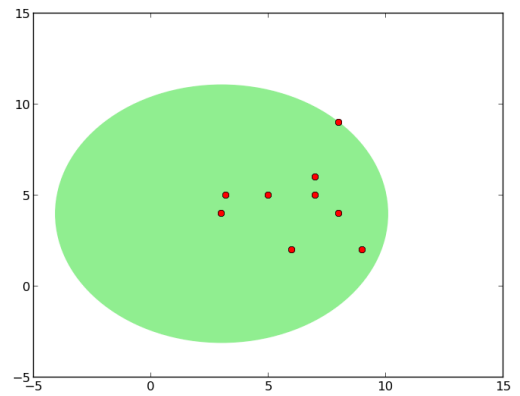
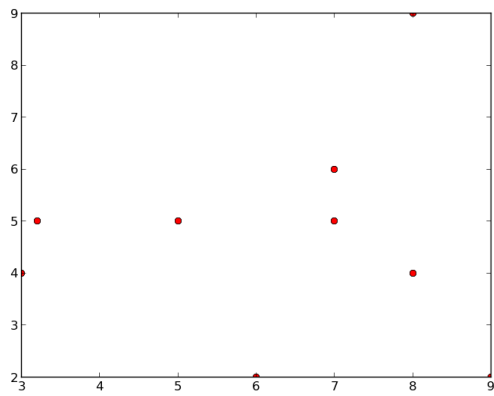
kNN'in model kullanmayan, model yerine verinin kendisini kullanan bir algoritma olarak tanıttık. Peki "eğitim" evresi sonrası ele geçen küreler ve ağac yapısı bir nevi model olarak görülebilir mi?

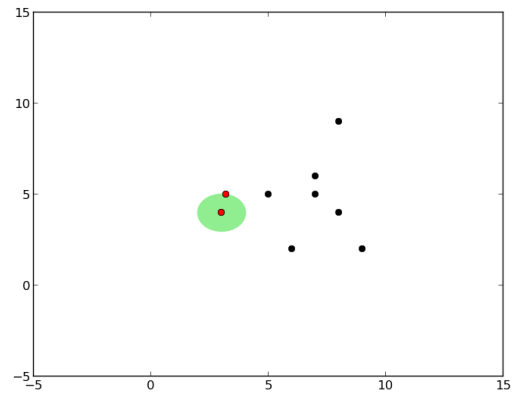
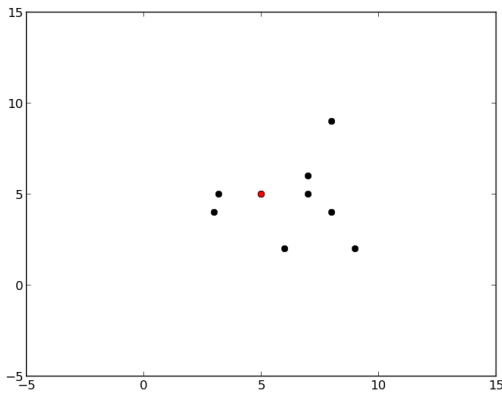
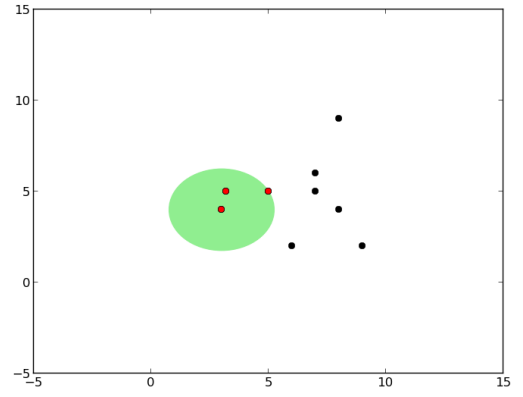
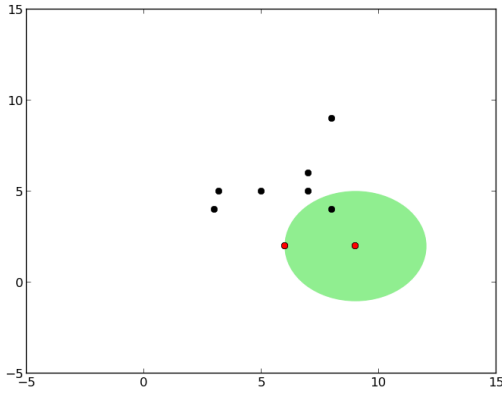
Bu önemli bir soru, ve bir bakıma, evet ağac yapısı sanki bir modelmiş gibi duruyor. Fakat, mesela istatistiksel, grafiksel, yapay sinir ağları (neural net) bağlamında bakılırsa bu yapıya tam bir model denemez. Model bazlı metodlarda model kurulunca veri atılır, ona bir daha bakılmaz. Fakat kNN, küre ve ağac yapısını hala eldeki veriye erismek için kullanmaktadır. Yani bir bakıma veriyi "indeksliyoruz", ona erişimi kolaylaştırıp hızlandırıyoruz, ama ondan model çıkartmıyoruz.

Not: Verilen Python kodu ve algoritma yakın noktaları hesaplıyor sadece, onların etiketlerinden hareketle yeni noktanın etiketini tahmin etme aşamasını gerçekleştiriyor. Fakat bu son aşama işin en basit tarafı, eğitim veri yapısına eklenecek bir etiket bilgisi ve sınıflama sonrası k noktanın ağırlıklı etiketinin hesabi ile basit şekilde gerçekleştirilebilir.

```
!python plot_circles.py
```

Ağacı oluşturma sırasında kürelerin grafiği alttadır.





## Kaynaklar, Notlar

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

[2] Alpaydın, Introduction to Machine Learning

[3] Silme islemi ornek kodumuzda Python `del` ile gercekleştirildi. Eger bu islem de hizlandirilmak istenirse, en alt kure seviyesindeki veriler bir oncelik kuyrug (priority queue) üzerinde tutulabilir, ve silme islemi hep en sondaki elemani siler, ekleme islemi ise yeni elemani (hep sirali olan) listede dogru yere koyar.