

## Felzenswalb Gruplamasi (Felzenswalb Clustering)

Minimum Kapsayan Agac (Minimum Spanning Tree -MST-) kavramini kullanan Felzenswalb kumelemesini gorecegiz. Literaturde bu metotun imaj gruplamasi icin verildigini gorebilirsiniz, biz imaj gruplamasi icinden veri kumelemesi yapan kısmi cikarttik ve bu algoritmayı ayri bir sekilde paylasiyoruz.

Felzenswalb'in matematiginde once imaj bolgelerinin (ya da veri kumeleri olarak dusunebiliriz) ikili karsilastirmaya bir olcut gerekir. Bu bolumde bir beyani D ortaya koyacagiz, ki bu beyan, imajdaki iki bilesen (ki imaj gruplamasinin dogru olarak bulmaya calisacagi bilesenler) arasinda bir sinir olup olmadigina dair kanitin olcusu olacak. Beyanin temeli sudur: iki bilesen arasindaki sinirin boyunda yer alan her iki tarafın ogelerinin farkliligina bak, ve onu her bilesenin kendi icindeki farkliliga gore oranla. Yani bu beyan, bir bilesenin ic farkliligini dis farkliligina kiyaslar, ve bu sebeple verinin yerel karakteristikleri gozetmis olur. Kiyaslama mesela, global, verinin her yerinde aynen gecerli olacak bir sabit esik degerine vs. bagli degildir.

### Tanim

Bir bilesen  $C \subseteq V$ , ki  $C$  bir bilesendir (component) ve  $V$  cizitin tum noktalaridir, *ic farkliligini*, o  $C$ 'nin minimum kapsayan agacinin, yani  $MST(C)$ 'sinin en buyuk kenar agirligi olarak aliyoruz. Bu ic farkliligi  $Int(C)$  olarak belirtirsek,

$$Int(C) = \max_{e \in MST(C,E)} w(e)$$

ki  $w((v_i, v_j))$  bir cizit  $G = (V, E)$ 'yi olusturan bir kenar  $(v_i, v_j) \in E$  agirligi olarak belirtilir.

### Tanim

İki bilesen  $C_1, C_2 \subseteq V$  arasindaki farki o iki bileseni birlestiren kenarlardan en ufagi olarak aliyoruz. İki bilesenin arasinda birden fazla baglanti olmasi mumkundur, tum bunlara bakiyoruz, ve en ufagini aliyoruz.

$$Dif(C_1, C_2) = \min_{v_i \in C_1, v_j \in C_2, (v_i, v_j) \in E} w((v_i, v_j))$$

Eger  $C_1, C_2$  arasinda bir kenar yok ise  $Dif(C_1, C_2) = \infty$  kabul ediliyor.

Prensip olarak iki bilesen arasindaki en minimal baglantinın problem cikartabilecegi dusunulebilirdi, niye en az, niye ortalama vs degil? Fakat pratikte bu olcutun cok iyi isledigini gorduk. Hatta iyi olmaktan ote, bu olcutu minimal yerine medyan, ya da diger ceyreksel (quantile) olcute degistirdigimiz zaman (ki bunu yaparak aykiri degerlere -outlier- karsi daha dayanikli olmasini istemistik), algoritma cetrefilligi NP-Zor haline geliyor. Yani gruplama kriterinde ufacak bir degisiklik problemin cozum zorlulugunda muthis bir degisim ortaya cikartiyor.

Simdi iki bilesenin karsilastirma beyani  $D$ 'nin tanimina geldik.  $D$  olcutu,  $Dif(C_1, C_2)$ 'nin

$\text{Int}(C_1)$  ya da  $\text{Int}(C_2)$ 'den herhangi birinden daha büyük olup olmadigina bakar. Ayrica bu karsilastirmayi bir esik degeri uzerinden pay ekleyerek yapar, eger irdeme olumlu ise, iki bilesen arasinda sinir vardır, yoksa yoktur.

$$D(C_1, C_2) = \begin{cases} \text{Dogru} & \text{Eger } \text{Dif}(C_1, C_2) > M\text{Int}(C_1, C_2) \text{ ise} \\ \text{Yanlis} & \text{Diger durumda} \end{cases}$$

Minimum ic fark  $M\text{Int}$  ise soyle tanimlidir,

$$M\text{Int}(C_1, C_2) = \min(\text{Int}(C_1) + \tau(C_1), \text{Int}(C_2) + \tau(C_2))$$

Esik fonksiyonu  $\tau$  ustteki irdeledigimiz fark hesaplarinin belli derecelerde disari-dan etkilemek icin koyulmustur. Eger bunu kullanmasaydik sadece  $\text{Int}$  fonksiy-onunu kullanmamiz gerekecekti, fakat bu olcut tek basina ufak bir bilesenin yerel karakteristiklerini gostermesi acisindan yeterli degildir. Asiri durumda mesela  $|C| = 1, \text{Int}(C) = 0$ , yani en kucuk  $C$  durumudur bu ( $|C|$  bilesenin icindeki oge sayisi), icinde tek oge vardır, ve hicbir kenar yoktur,  $\text{Int}(C) = 0$ .

Bu sebeple iyi bir  $\tau$  bilesenin buyuklugunu hesaba katarak, ona ters oranli bir rakam olusturursa iyi olur, mesela bir sabit  $k$  uzerinden,

$$\tau(C) = \frac{k}{|C|}$$

Bu demektir ki ufak bilesenler icin daha kuvvetli bir ispat ariyoruz, cunku ku-cuk  $|C|$ ,  $\tau$ 'yu buyutecektir, ve  $\text{Dif}$ 'in ondan buyuk olmasi daha zorlasacaktır. Tabii dikkat edelim,  $k$  bir "bilesen sayisi" degildir, yani fonksiyonuna dikkatli bakarsak, eger bilesenler arasinda yeterince buyuk bir fark var ise ufak bilesen-lere hala izin verilmistir.

Algoritma soyledir, girdi olarak  $G = (V, E)$  alır, ve  $V$ 'yi  $S$  bilesenlerine ayirir ki her  $S$  icinde ona ait olan kenarlar vardır, yani  $S = (C_1, \dots, C_r)$

---

$\text{felzenswalb}(G)$	
1	$E$ kenarlarini $\pi = (o_1, \dots, o_m)$ seklinde kucukten buyuge dogru sirala.
2	Ilk basta $S^0$ gruplamasini al. Bu durumda her kenar $v_i$ kendi bileseni icindedir.
3	for $q = 1, \dots, m$
4	$S^{q-1}$ gruplamasini baz alip $S^q$ gruplamasini soyle yarat; $q$ 'inci siradaki
5	kenarin birlestirdigi noktalar $v_i, v_j$ oldugunu farz edelim, yani $o_q = (v_i, v_j)$ .
6	Eger $v_i, v_j$ $S^{q-1}$ gruplamasi icinde farkli iki bilesen icindeyseler, ve $w(o_q)$ her
7	iki bilesenin icisel farkina kiyasla cok kucuk ise, bu iki bileseni birlestir,
8	yoksa hicbir sey yapma.
9	return $S = S^m$

---

Ustteki dongu icindeki en son irdelemede icisel farktan bahsediliyor, bu tabii ki  $M\text{Int}(C_1, C_2)$ . Daha formel sekilde  $M\text{Int}(C_1^{q-1}, C_2^{q-1})$  cunku bilesenlerin icerik-leri hangi adimda oldugumuza gore degisebilir,  $q$  adiminda bir onceki  $q - 1$ 'den

bize “miras kalan” gruplamalar ve bileşenler üzerinden iş yapıyoruz. Bir sonraki adıma ya birleşmiş, ya da birleşmemiş (aynı) gruplamaları aktarıyoruz.

```
import scipy.sparse as sps
import scipy.io as io
import itertools, numpy as np

def threshold(size, c): return c / size

S = {}

def find(C, u):
    if C[u] != u:
        C[u] = find(C, C[u])           # Path compression
    return C[u]

def union(C, R, u, v, S):
    u, v = find(C, u), find(C, v)
    if R[u] > R[v]:                     # Union by rank
        C[v] = u
        S[v] = S[u] + S[v]
    else:
        C[u] = v
        S[v] = S[u] + S[v]
    if R[u] == R[v]:                   # A tie: Move v up a level
        R[v] += 1

class Felzenswalb:
    def __init__(self, min_size, c):
        self.min_size_ = min_size
        self.c_ = c

    def fit(self, X):
        print X.shape
        G = {}
        for i in range(X.shape[0]): G[i] = {}
        for u,v,w in itertools.izip(X.row, X.col, X.data): G[u][v] = w
        E = [(G[u][v],u,v) for u in G for v in G[u]]
        E = sorted(E)
        T = set()
        C, R = {u:u for u in G}, {u:0 for u in G}   # Comp. reps and ranks
        S = {u:1 for u in range(len(G))}

        ts = {x:threshold(1,self.c_) for x in C}

        for w, u, v in E:
            if find(C, u) != find(C, v):
                if w <= ts[u] and w <= ts[v]:
                    T.add((u, v))
                    union(C, R, u, v, S)
                    ts[u] = w + threshold(S[u],self.c_)

        for _, u, v in E:
            if find(C, u) != find(C, v):
                if S[C[u]] < self.min_size_ or S[C[v]] < self.min_size_:
```

```

        union(C, R, u, v, S)

    self.labels_ = [np.nan for i in range(len(C))]
    for i in range(len(C)): self.labels_[i] = int(C[i])
    self.T_ = T

```

## Basit bir ornek

```

import scipy.sparse as sps, felz
import scipy.io as io
X = io.mmread('simple.mtx')
clf = felz.Felzenswalb(min_size=1,c=1.0)
clf.fit(X)
print clf.labels_

(5, 5)
[1, 1, 3, 3, 1]

import scipy.sparse as sps
import scipy.io as io, random
import pandas as pd, os, sys
syn = pd.read_csv("../kmeans/synthetic.txt", names=['a', 'b'], sep=" ")
data = np.array(syn)

from sklearn.metrics.pairwise import euclidean_distances
X = euclidean_distances(data, data)

X2 = X.copy()
# filter out large values / distances so matrix can be sparse
X2[X > 2000] = 0.0
X3 = sps.lil_matrix(X2)
X4 = sps.triu(X3)
print 'non-zero items', len(X4.nonzero()[0])
print X4.shape

non-zero items 87010
(3000, 3000)

import felz
clf = felz.Felzenswalb(min_size=20,c=800)
clf.fit(X4)

(3000, 3000)

syn['cluster'] = clf.labels_
print len(syn['cluster'].unique()), 'clusters found'
print syn[:5]

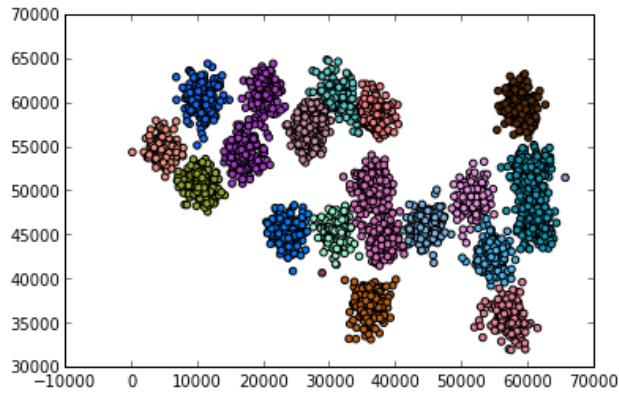
19 clusters found
      a      b cluster
0  54620  43523      120
1  52694  42750      120
2  53253  43024      120
3  54925  42624      120
4  54973  43980      120

```

```

import random
for clust in syn['cluster'].unique():
    tmp = np.array(syn[syn['cluster'] == clust][['a', 'b']])
    plt.scatter(tmp[:,0], tmp[:,1], c=np.random.rand(3,1))
plt.savefig('mstseg_01.png')

```



## Kaynaklar

[1] Pedro F. Felzenszwalb and Daniel P. Huttenlocher, *Efficient Graph-Based Image Segmentation*, [http://scikit-image.org/docs/dev/auto\\_examples/plot\\_segmentations.html](http://scikit-image.org/docs/dev/auto_examples/plot_segmentations.html)