

Produced by YE Qianying, GAO Ge, ZHAO Feier



# Your Nose Never Know

Mysterious Network Behind The World of Perfume





girl

- 01      Background**
- 02      Data Collection**
- 03      Data Cleaning**

- 04      Visualization**
- 05      Results & Conclusions**

# Background 01

Beauty  
is a  
dramatic  
power





## About Perfume

Perfume is closely linked to urbanites' daily life

- a. Influences our **sense of smell**.
- b. Influences our **mood state**.
- c. For consideration of **social needs**, people would like to wear perfumes in both formal and informal situations.
- d. Different fragrances show different **personalities**.
- e. **Fashion items** for women and men.



# Attitudes towards Perfume



**Chanel**

*"I wear nothing but a few drops of Chanel No.5."*



**Calvin Klein**

*"Be good, be bad, just be yourself."*



**Dior**

*"Gold is cold. Diamonds are dead. A Limousine is a car. Don't pretend. Feel what's real."*



**Burberry**

*"The good things in life never change."*



**Estee Lauder**

*"Is staying in touch the same as being in touch?"*

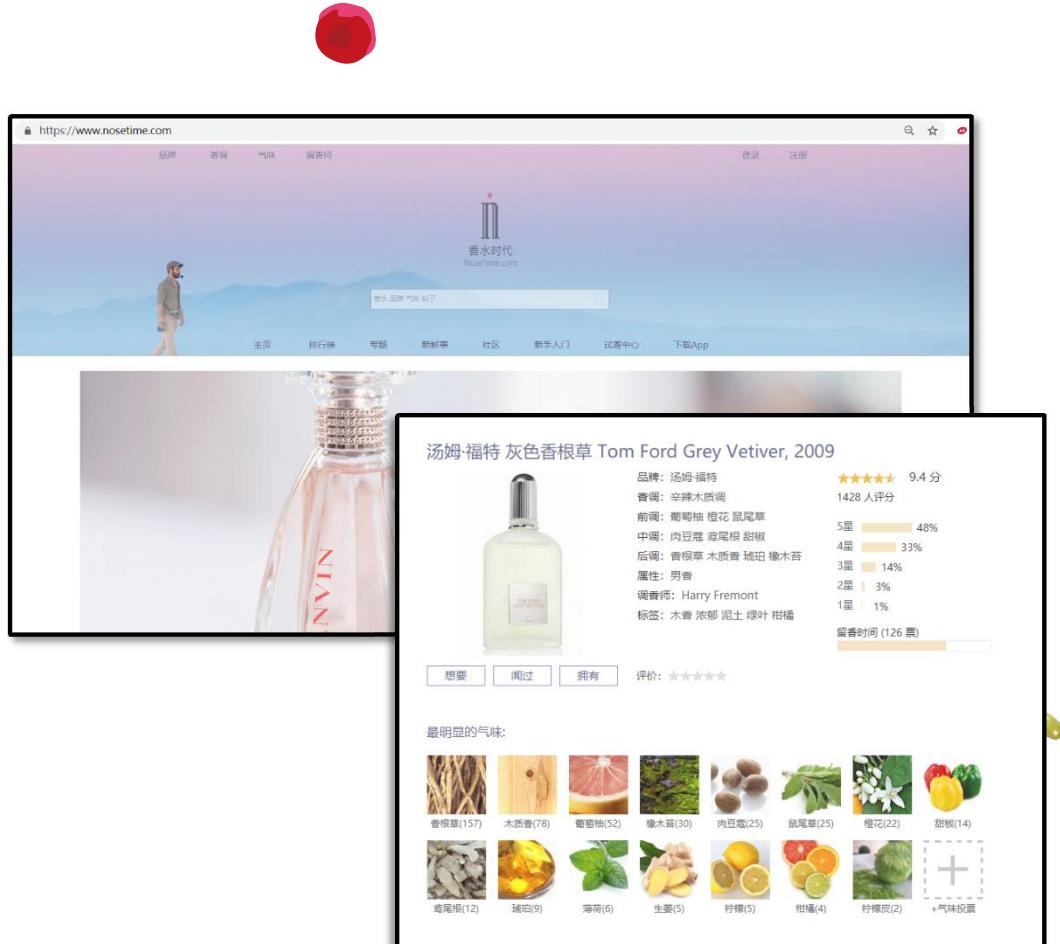


*Thousands of perfumes have been created by cosmetic manufacturers as “**name cards**” for consumers to choose, which also arouses our interests to explore more about **the world of perfume**...*





**Data Collection**



<https://www.nosetime.com>

# Data Source

## Nosetime 香水时刻

We used python to get the information of all brands' products.

- Got all information of **2040 brands**.
- Got **2226+1441=3667 urls**.
- Variables: ingredients, brand names, three notes (top notes, middle notes, base notes), years of perfumes, scores of perfumes, reply number of perfumes...



*Here is what we used.*



```
In [4]: from selenium import webdriver
import time
import re
from lxml import etree
from bs4 import BeautifulSoup
import bs4 as bs
from selenium.common.exceptions import NoSuchElementException
from selenium.webdriver.common.by import By
from selenium.webdriver.support.ui import WebDriverWait
from selenium.webdriver.support import expected_conditions as EC
import pandas as pd

URLS_brand = []
#browser = webdriver.Firefox()#Chrome('./chromedriver.exe')
whole = ['2-a','3-b','4-c','5-d','6-e','7-f','8-g','9-h','10-i','11-j','12-k','13-l','14-m','15-n','16-o','17-p',
         '18-q','19-r','20-s','21-t','22-u','23-v','24-w','25-x','26-y','27-z']

# For every page in the interval
for wh in whole:
    PATIENCE_TIME = 60
    driver = webdriver.Chrome()
    driver.get('https://www.nosetime.com/pinpai/' + str(wh) + '.html')
    driver.encoding = 'utf-8'
    driver.maximize_window()

    html = driver.page_source
    soup = BeautifulSoup(html, 'lxml')
    url_containers = soup.findall('div', class_='odorlist')

    for container in url_containers:
        #find all the date of the news
        try:
            for url_b in container.findall('a'):
                url_brand = url_b.get('href')
                URLs_brand.append(url_brand)
        except:
            url_brand = "no_url"
            URLs_brand.append(url_brand)
    driver.close()
driver.quit()
```

```
In [39]: #get the pure id of brand
Brand_id = []
for i in Cleaned_URLs_brand:
    brand_id = re.sub(r'[^0-9]', '', i)
    Brand_id.append(brand_id)
datapagel = [
    'Brand_id': Brand_id,
]
brandid = pd.DataFrame.from_dict(datapagel, orient='index').transpose()
print(brandid.info())
brandid.to_csv('Brand_id.csv', index=False, encoding='utf_8_sig')

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2040 entries, 0 to 2039
Data columns (total 1 columns):
Brand_id    2040 non-null object
dtypes: object(1)
memory usage: 16.0+ KB
None
```

*That is what we have got.*



**Got urls according  
to the brand names.**

```
In [83]: len(brand_id)
Out[83]: 2040

In [84]: set(brand_id)
len(brand_id)
Out[84]: 2040

In [75]: AllURLs = []
#browser = webdriver.Firefox()#Chrome('./chromedriver.exe')
import numpy as np
import random
import time
from time import sleep
#For each page in the interval
for bid in brand_id:
    for bid in brand_id:
        PATIENCE_TIME = 60
        driver = webdriver.Chrome()
        driver.get('https://www.nosetime.com/brand.php?id=' + str(bid))
        driver.encoding = 'utf-8'
        driver.maximize_window()

        html = driver.page_source
        soup = BeautifulSoup(html, 'lxml')
        all_url_container = soup.find_all('div', class_='next_news')

        for container in all_url_container:
            try:
                for url in container.find_all('a'):
                    AllURLs.append("https://www.nosetime.com/brand.php?id=" + str(bid) + "page" + url.text)

            except:
                allurls = 'no_allurls'
                AllURLs.append(allurls)
                driver.close()
                time.sleep(1)
                driver.quit()

In [130]: len(AllURLs)
Out[130]: 3424

In [135]: p_urls = []
error = []
for i in AllURLs:
    if not "下一页" in i:
        if not "尾页" in i:
            p_urls.append(i)
        else:
            error.append(i)
len(p_urls)
Out[135]: 2226

In [136]: print(error[10])
['http://www.nosetime.com/brand.php?id=10014016page=下一页', 'http://www.nosetime.com/brand.php?id=10054734page=下一页', 'http://www.nosetime.com/brand.php?id=1010461page=下一页', 'http://www.nosetime.com/brand.php?id=10061680page=下一页', 'http://www.nosetime.com/brand.php?id=1028489page=下一页', 'http://www.nosetime.com/brand.php?id=10098532page=下一页', 'http://www.nosetime.com/brand.php?id=1035496page=下一页', 'http://www.nosetime.com/brand.php?id=10030919page=下一页', 'http://www.nosetime.com/brand.php?id=10076007page=下一页', 'http://www.nosetime.com/brand.php?id=10077602page=下一页']

In [137]: datapage2 = [
    p_urls: p_urls,
]
p_urls = pd.DataFrame.from_dict(datapage2, orient = 'index').transpose()
print(p_urls.info())
p_urls.to_csv('part_urls.csv', index = False, encoding = 'utf_8_sig')

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2226 entries, 0 to 2225
Data columns (total 1 columns):
 p_urls: 2226 non-null object
dtypes: object(1)
memory usage: 17.5+ KB
None

In [142]: import csv
purl = []
with open('part_urls.csv', 'r') as f:
    reader = csv.reader(f)
    for line in f:
        purl = [row[0] for row in reader]

In [153]: import re
part_id = []
for i in purl:
    new1 = re.split('\d+', i) #将url后的部分
    new1 = re.sub(r'\d+', '', new1[0])
    part_id.append(new1)

In [157]: part_id = set(part_id)
len(part_id)
Out[157]: 599

In [161]: another_url = []
for i in brand_id:
    if i not in part_id: #find the singel page brand id
        another_url.append("https://www.nosetime.com/brand.php?id=" + str(i)) #append the single page url
len(another_url) #the number is correct
Out[161]: 1441
```

*That is what we have got.*



```
1 from selenium import webdriver
2 import time
3 import re
4 from lxml import etree
5 from bs4 import BeautifulSoup
6 import bs4 as bs
7 from selenium.common.exceptions import NoSuchElementException
8 from selenium.webdriver.common.by import By
9 from selenium.webdriver.support.ui import WebDriverWait
10 from selenium.webdriver.support import expected_conditions as EC
11 import pandas as pd
12
13 Names_perfume = []
14 Whole_info = []
15 Score = []
16 #browser = webdriver.Firefox()#Chrome('./chromedriver.exe')
17 for url in URLs:
18     PATIENCE_TIME = 60
19     driver = webdriver.Chrome()
20     driver.get(url)
21     driver.encoding = 'utf-8'
22     driver.maximize_window()
23
24     html = driver.page_source
25     soup = BeautifulSoup(html, 'lxml')
26     perfume_containers = soup.find_all('div', class_='item')
27
28     for container in perfume_containers:
29         #find all the date of the news
30         try:
31             for name in container.find_all('h2'):
32                 Names_perfume.append(name.text)
33         except:
34             name = "no_date"
35             Names_perfume.append(name)
36
37         #find all the ingredients of the perfumes
38         try:
39             for whole_info in container.find_all('div', class_='info'):
40                 Whole_info.append(whole_info.text)
41
42         except:
43             whole_info = "no_whole_info"
44             Whole_info.append(whole_info)
45
46         #find all the scores of the perfumes
47         try:
48             for score in container.find_all('div', class_='score'):
49                 Score.append(score.text)
50         except:
51             score = "no_score"
52             Score.append(score.text)
53
54     driver.close()
55 driver.quit()
```

# Some variables.

*That is what we have got.*



1	Names_perfume	Whole_info	Score
2	皇家香氛 Parfums Regence Eau De Regence	前调：中调：后调： 气味：董衣草	0 人评价
3	皇家香氛 Parfums Regence Lavendou	前调：中调：后调：	0 人评价
4	皇家香氛 Parfums Regence Lavendou	前调：中调：后调：	0 人评价
5	皇家香氛 Parfums Regence Vert Regence	前调：中调：后调：	0 人评价
6	皇家香氛 Parfums Regence Santal	前调：中调：后调：	0 人评价
7	皇家香氛 Parfums Regence Oakmos	前调：中调：后调： 气味：香柠檬, 董衣草, 广藿香, 白松香, 劳丹脂, 檀木苔	0 人评价
8	皇家香氛 Parfums Regence Kolnisch Juchten	前调：中调：后调：	0 人评价
9	皇家香氛 Parfums Regence Eau De Regence	前调：中调：后调： 气味：柑橘, 果香	0 人评价
10	皇家香氛 Parfums Regence Santal	前调：中调：后调： 气味：檀香木	0 人评价
11	皇家香氛 Parfums Regence Kolnisch Juchten	前调：中调：后调： 气味：乳香, 闭鞘姜, 皮革, 麻香	0 人评价
12	皇家香氛 Parfums Regence Oakmos	前调：中调：后调：	0 人评价
13	布鲁诺·法兹拉瑞 Bruno Fazzolari Jimmy, 2013	前调：柠檬, 依兰中调：玫瑰, 老鹳草后调：紫罗兰叶, 天芥菜, 檀香木	2 人评价
14	布鲁诺·法兹拉瑞 Bruno Fazzolari Lampblack, 2013	前调：胡椒, 檳子中调：葡萄柚后调：莎草, 安息香脂, 香根草	6 人评价
15	布鲁诺·法兹拉瑞 Bruno Fazzolari Room 237, 2015	前调：小雏菊, 乙烯基, 当归, 红没药, 闭鞘姜, 绿叶, 花香, 麻香	3 人评价
16	布鲁诺·法兹拉瑞 Bruno Fazzolari Au Delà - Narcisse des Montagnes Bruno Fazzolari Au Delà Narcisse, 2014	前调：香柠檬, 水仙花, 茉莉, 檀木苔, 檀花, 玻珀	2 人评价
17	布鲁诺·法兹拉瑞 Bruno Fazzolari Au Delà, 2013	前调：壳果中调：橙花油, 茉莉后调：橙花, 树脂, 檀木苔	2 人评价
18	布鲁诺·法兹拉瑞 Bruno Fazzolari Seyrig, 2015	前调：醛, 橘子中调：依兰, 玫瑰, 铃兰, 丁香花后调：檀木苔, 麻香	2 人评价
19	布鲁诺·法兹拉瑞 Bruno Fazzolari Feu Secret, 2017	前调：鸢尾根, 云杉, 尤加利, 姜黄, 粉红胡椒, 雪松, 檀木	1 人评价
20	布鲁诺·法兹拉瑞 Bruno Fazzolari Monserrat, 2013	前调：葡萄柚, 胡萝卜籽中调：杏, 绿叶, 茉莉后调：龙涎香, 麻香, 石膏	1 人评价
21	布鲁诺·法兹拉瑞 Bruno Fazzolari Five, 2013	前调：柠檬, 檳子中调：迷迭香, 舌瓣叶后调：木质香	1 人评价
22	布鲁诺·法兹拉瑞 Bruno Fazzolari Vetiverissimo, 2018	前调：辛香料, 香根草, 雪松, 树脂, 木质香	0 人评价
23	布鲁诺·法兹拉瑞 Bruno Fazzolari Vetiverissimo, 2018	前调：辛香料, 香根草, 雪松, 树脂, 木质香	0 人评价
24	布鲁诺·法兹拉瑞 Bruno Fazzolari Cadavre Exquis, 2016	前调：血橙, 檀脑, 依兰, 金盏花, 干果, 八角, 黑巧克力, 柚树, 安息香脂, 香草, 麻猫香	0 人评价
25	布鲁诺·法兹拉瑞 Bruno Fazzolari Fontevraud, 2018	前调：香柠檬, 番石榴, 梨, 玫瑰, 檀木苔, 红没药, 广藿香, 玻珀	0 人评价
26	布鲁诺·法兹拉瑞 Bruno Fazzolari Unsettled, 2017	前调：香柠檬, 茶叶, 快乐鼠尾草, 菠萝, 檀香木, 劳丹脂, 香草, 海水	0 人评价
27	伊露米琳 白色橘子花瓣 Illuminum White Gardenia Petals, 2011	前调：百合, 香柠檬, 醇栗叶中调：茉莉, 铃兰, 桔子花, 依兰后调：琥珀, 木质香	7.1 分 66 人评价
28	伊露米琳 Illuminum Wild Tobacco, 2011	前调：公丁香, 快乐鼠尾草中调：雪松, 烟草后调：劳丹脂, 海狸香, 雪松	6.5 分 11 人评价
29	伊露米琳 Illuminum Tribal Black Tea, 2011	前调：香柠檬, 小豆蔻, 柠檬中调：茱萸, 花香, 茉莉后调：琥珀, 木质香, 内豆蔻, 雪松	7 人评价
30	伊露米琳 Illuminum Moroccan Tuberose	前调：老鹳草, 玫瑰中调：依兰, 玫瑰后调：晚香玉, 雪松, 麻香	2 人评价
31	伊露米琳 Illuminum Ginger Pear, 2011	前调：柑橘, 梨, 香柠檬中调：茱萸, 花香, 姜黄, 小茴香后调：麝香, 檀香木	8 人评价
32	伊露米琳 Illuminum Piper Leather, 2011	前调：壳果, 胡椒中调：皮革, 茉莉, 胡萝卜籽后调：乳香, 麝猫香, 麻香	6 人评价
33	伊露米琳 Illuminum Rose Oud, 2011	前调：茉莉, 铃兰, 罗勒, 芦荟中调：玫瑰, 老鹳草后调：海狸香, 沉香(乌木), 广藿香	5 人评价
34	伊露米琳 Illuminum Taif Rose	前调：老鹳草, 玫瑰中调：塔夫玫瑰, 玫瑰后调：麝香, 檀香木	5 人评价
35	伊露米琳 Illuminum Scarlet Oud, 2011	前调：胡椒, 蜂蜜中调：沉香(乌木), 香根草, 树脂后调：麝猫香, 海狸香	3 人评价
36	伊露米琳 Illuminum Wild Berry Blossom, 2011	前调：醋栗叶, 仙客来中调：牡丹, 野花, 树莓花, 木兰后调：雪松, 麻香	2 人评价
37	伊露米琳 Illuminum Black OUD	前调：广藿香, 柑橘, 公丁香中调：沉香(乌木)后调：麝猫香	2 人评价
38	伊露米琳 Illuminum Saffron Amber, 2011	前调：香柠檬, 桃子, 芦荟中调：依兰, 藏红花, 茉莉后调：琥珀, 麻香	2 人评价
39	伊露米琳 Illuminum White on White	前调：肉桂, 公丁香, 小豆蔻中调：广藿香后调：安息香, 香豆素, 檀木苔	2 人评价
40	伊露米琳 Illuminum Hindi OUD	前调：桂花, 鸡蛋花, 桃子花中调：茉莉, 沉香(乌木), 黑檀红木后调：麝香	2 人评价
41	伊露米琳 Illuminum Skin Petals, 2011	前调：苹果花, 森林水果, 黑加仑中调：铃兰, 紫罗兰, 茉莉后调：琥珀, 麻香	1 人评价
42	伊露米琳 Illuminum White Saffron	前调：老鹳草, 茉莉中调：藏红花, 玫瑰, 香根草后调：愈创木, 檀香木, 麻香	1 人评价
43	伊露米琳 Illuminum Cashmere Musk, 2011	前调：椰子, 绿叶, 依兰, 仙客来中调：风信子, 雪松, 檀香木后调：开司米木, 麝香, 木质香	1 人评价

**Original data.**



# 03

## Data cleaning





*In the process of data cleaning, we chose the variables from the whole information we get. For example, split the English name, Chinese name and year from the column of "Names\_perfume", and split the three notes from the information of whole fragrances...*

# Step 1

Extracted the variables we need from the raw data and eliminated some unwanted information.

# For missing values, we filled in with "None" or "0".



```

1 import jieba
2 f1= open("cn_names.txt", "a", encoding="utf-8",newline="\n")
3 for i in Chinese_names:
4     i.encode("utf-8")
5     seg_list = jieba.cut(i.cut_all=False)
6     seg_list = ",".join(seg_list)
7     f1.write(seg_list)
8     f1.write("\n")

```

```

1 import re
2 Brand_names = []
3 Chinese_names = []
4 for i in Names_perfume:
5     i = i.split(" ")
6     cn = i[1]
7     ncn = re.sub(r'[a-zA-Z]', '', cn)
8     ncn1 = re.sub(u"\u20([.*])|\u20([.*])|\u20([.*])|\u20([.*])", "", ncn)
9     #i = re.sub(r'\ufe00-\ufeff', '', ncn1)
10
11     if ncn1 != '':
12         Chinese_names.append(ncn1)

```

```

1 Years = []
2 for i in Names_perfume:
3     if not "," in i:
4         Years.append("None")
5     else:
6         i = i.split(",")
7         i = i[-1]
8         i = re.sub(r'(\s+)', '', i)
9         Years.append(i)
10
11 print(Years[100:120])

```

```

['2015', 'None', 'None', 'None', 'None', 'None', '2016', 'None', '2013', 'None', 'None', 'None', '2013', '2010', '2010', 'None', 'None', '2010']

```

```

1 print(Chinese_names[100:200])

```

```

['爱神周一晚报', '爱恋', '金色露雅', '斯嘉丽天堂', '安妮浓缩版', '捕捉我之水', '之水', '诺亚天堂', '斯嘉丽', '爱神之水', '露雅之水2018', '趣味火烈鸟之水', '露露', '闪耀爱神', '露雅', '安妮', '趣味火烈鸟之水', '露露', '爱神之水', '诺亚之水', '爱神丘比特之水火烈鸟', '优雅小姐天堂', '优雅小姐之水', '伊甸园', '诺亚之水火烈鸟', '爱神丘比特', '爱神丘比特', '红衣女郎', '拜占庭', '红衣女郎', '水之男', '拜占庭', '诱惑之水', '胡须', '他', '原版光辉', '同名', '点水威金', '谜', '蓝海珍珠', '光辉', '渴望女士香水', '秘密玫瑰极致', '环球', '鸢尾瀑布', '马加茂', '神秘沉香', '沉思之水男士香水', '夏日偷闲', '绝对', '花之水', '流行', '摇滚缪斯', '阳光', '幸运之星', '拜占庭之旅', '幻彩柑橘', '渴望男士香水', '秘密', '极度晶红', '勇敢', '沉思之水', '玫瑰', '游艇', '浓情时光', '游艇', '时光', '鼠尾草与海盐', '蓝色风铃草', '英国梨与小苍兰', '橙花', '杏桃花与蜂蜜', '红玫瑰', '黑莓与月桂叶', '青柠罗勒与柑橘', '牡丹与胭脂玫瑰', '伯爵茶与黄瓜', '白茉莉与薄荷', '黑石榴', '浓古龙-乌木与佛手柑', '浓古龙-丝绒玫瑰与乌木', '英国橡树与榛果', '浓古龙-黑琥珀与姜百合', '含羞草与小豆蔻', '柚子', '罗勒与橙花', '英国橡树与红醋栗', '琥珀与薰衣草', '浓古龙-晚香玉与天使草', '普提花', '伦敦雨-黑雪松与杜松', '合欢花', '龙舌兰可可', '豆蔻辛姜', '浓古龙-茉莉与冬加豆', '樱花', '浓古龙-焚香与香橼', '星花木兰', '桂花', '夏之青梅']

```

# Step 1

**Extracted the variables we need from the raw data and eliminated some unwanted information.**

# For missing values, we filled in with "None" or "0".

```
1 Years = []
2 for i in Names_perfume:
3     if not "," in i:
4         Years.append("None")
5     else:
6         i = i.split(",")
7         i = i[-1]
8         i = re.sub(r'(\s+)', '', i)
9         Years.append(i)
10 print(Years[100:120])
```

```
1 Score_number=[]
2 for i in Score:
3     if "分" in i:
4         i1 = i.split("分")
5         i1 = i1[0]
6         Score_number.append(i1)
7     elif "(" in i:
8         i2 = i.split("(")
9         i2 = i2[0]
10        Score_number.append(i2)
11    else:
12        Score_number.append("None")
13 print(Score_number[100:120])
```

```
1 Reply_number = []
2 for i in Score:
3     if "分" in i:
4         i1 = i.strip("人评价")
5         i1 = i1.split("分")
6         i1 = i1[-1]
7         Reply_number.append(i1)
8     elif "分" in i:
9         i2 = i2.strip("评价")
10        i2 = i2.split("分")
11        i2 = i2[-1]
12        Reply_number.append(i2)
13    else:
14        Reply_number.append("0")
15 print(Reply_number[100:120])
```



# Step 1

**Extracted the variables we need from the raw data and eliminated some unwanted information.**

# For missing values, we filled in with "None" or "0".

```

1 Top = []
2 for i in Main_info:
3     i = i.strip('中调')
4     i = i.split('中调')
5     top = i[0]
6     if top != '':
7         Top.append(top)
8 print(Top[0:100])

```

「柠檬,依兰」,胡椒,橙子,「芙蓉」,醛,橘子,「南茜怕,胡萝卜籽」,「柠檬,橙子」,百合,香柠檬,甜茉莉,「公丁香,快乐鼠尾草」,「垂柳香,小麦草,小豆蔻,柠檬」,老鹳草,玫瑰,「柑橘,香柠檬」,芸叶,胡蘿,茉莉,芦荟,「迷迭香,长枝桔」,胡椒,蜂蜜,「薰草叶,仙人掌」,「薰衣草,柑橘,公丁香」,香柠檬,茉莉,芸香,「肉桂,公丁香」,小白苣,桔梗,鸡蛋花,瓶子花,「苹果花,荔枝果黑,果皮」,「薰衣草,茉莉」,椰子,绿叶,依兰,仙客来,「香柠檬,黑加仑,重瓣茉莉,橙子」,香柠檬,血桔,桔梗,甜桔,「柑橘,柠檬,香柠檬,小麦草」,「香柠檬,血桔,柠檬,甜桔」,「橙花,苦橙叶」,「玫瑰,藏红花」,铃兰,莲花,「黑加仑,黑加仑,苦茉莉」,桂花,「玫瑰,老鹳草」,「柠檬,香柠檬,橙花,小麦草」,「芙蓉,孜然」,「芙蓉,孜然」,「茉莉,孜然」,「香柠檬,去魅惑」,「去魅惑,玫瑰」,「苦柑叶,高丽,橙子」,香柠檬,依兰,「茉莉,玫瑰」,「薰衣草,百里香,广藿香」,「无花果,孜然」,「茉莉,桂花」,「香柠檬,去魅惑」,「去魅惑,玫瑰」,「苦柑叶,高丽,橙子」,

香柠檬,依兰,「茉莉,玫瑰」,「薰衣草,百里香,广藿香」,「无花果,孜然」,「茉莉,桂花」,「香柠檬,去魅惑」,「去魅惑,玫瑰」,「苦柑叶,高丽,橙子」,

薰衣草,茉莉怕,「南茜怕,芙蓉」,「依兰,玫瑰,铃兰,丁香花」,「苦,绿叶,茉莉」,「迷迭香,桂叶怕」,「茉莉,铃兰,栀子花,依兰」,「薰松,烟草草」,「茉莉,花苞,茉莉怕,「依兰,我怕」,「茶叶,花香,生姜,大黄,小茴香」,「皮草,茉莉,胡萝卜籽」,「孜然,老鹳草怕」,「塔夫玫瑰,玫瑰怕」,「沉木(木本),雅得草,胭脂」,「牡丹,白花,桂花,木兰」,「沉木(木本)」,「沉木,珠江花,茉莉」,「广藿香」,「茉莉,沉木(木本),黑檀木」,「铃兰,紫罗兰,茉莉」,「藏红花,玫瑰草」,「凤信子,雪莲,茉莉」,「薰衣草,番木鳖,香柠檬」,「茉莉,紫罗兰,日本柚子,百香果,桃子」,「桂花,油,百里香,胡椒」,「茉莉,茉莉,玫瑰,玫瑰红胡椒」,「雅得草,凤信子,雪莲,茉莉」,「薰衣草,番木鳖,香柠檬」,「茉莉,紫罗兰,日本柚子,百香果,桃子」,

```
1 Alone_info = []
2 Main_info = []
3 for i in Whole_info:
4     if re.findall(r"^\u00d7\u00d7", i):
5         alone = i.replace("\u00d7\u00d7", "")
6         #alone = alone.replace(", ", "\n")
7         Alone_info.append(alone)
8     else:
9         Main_info.append(i)
10 print(A lone_info[100:120])
11 print(Main_info[100:110])
```

[‘顿加豆，干草，可可果，蜂蜜，烟草’，‘威士忌，玫瑰，甜椒，雪松，蜡烛’，‘小麦，绿叶，葡萄柚，绣线菊’，‘橘子，鼠尾草，蜂蜡，烟草’，‘报春花，玉米穗，含羞草，香草，847’，‘睡莲，绿叶，依兰，麝香，香草’，‘劳丹脂，愈创木’，‘杏仁，红醋栗，檀香木’，‘风信子，老鹳草，香根草’，‘palo santo wood, siam wood, white champ flower’, ‘black pepper, 康乃馨, 肉桂, 老鹳草’, ‘青草, 茉莉, 檀香, wood’, ‘檀香, sampaguita, ylang ylang’, ‘茉莉, 柠檬, 玫瑰, 檀香木’, ‘citrus, 玫瑰, cedar wood fire tree’, ‘草桂花, 皮革, rosewood’, ‘苦艾, balsam, oud’, ‘香柠檬, 茉莉, 广藿香, 香草’, ‘小豆蔻, indonesian vetiver, woods’, ‘香柠檬, 茉莉, 广藿香, 香草’]

[前调：鼠尾草，柏树中调：红苹果，纸莎草，雪花莲，木质感后调：香根草，麝香，皮革，琥珀，木质香，'前调：橘子花，晚香玉，无花果叶，紫丁香，紫罗兰，大麻，鸦片，罗勒中调：老鹳草，猕猴桃，菠萝，葡萄柚，公丁香，晚香玉后调：檀香木，雪松，'前调：香根草，桃花，柠檬中调：薄荷，黑加仑后调：琥珀，麝香，樟木本，雪松，'前调：柚子，粉红胡椒中调：花香，果仁糖，树莓后调：柏树，雪松，'前调：柠檬，黑加仑，红醋栗中调：玫瑰，树莓，石榴后调：麝香，雪松，紫罗兰，'前调：荔枝，日本柚子，柠檬中调：木兰，小苍兰，玫瑰后调：琥珀，麝香，木质香，'前调：橙子，桃子，梨中调：茉莉，莲花，忍冬后调：雪松，琥珀，香草，'前调：橘子，粉红胡椒，芒果中调：竹子，玫瑰，牡丹后调：香根草，檀香木，雪草，'前调：火龙果，玫瑰，梨，红浆果中调：牡丹，小苍兰，木兰花后调：麝香，檀香木，木质香，香草，树莓，'前调：橘子，薰衣草，树莓，桂花中调：花香，紫露草，茉莉，薰衣草，铃兰，玫瑰后调：雪松，香豆，檀香木，琥珀，雪松]



# Step 2

## Selected 100 “core nodes” by frequency counting.

# To simplify the networks and easily find underlying patterns from the most important nodes.



```
In [14]: import csv
import pandas as pd
datapage1 = [
    'Top',
]
pd_perfume = pd.DataFrame.from_dict(datapage1, orient='index').transpose()
print(pd_perfume.info())
pd_perfume.to_csv('perfume_top.csv', 'a', index = False, encoding = 'utf_8_sig')

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 15702 entries, 0 to 15701
Data columns (total 1 columns):
Top      15702 non-null object
dtypes: object(1)
memory usage: 122.8+ KB
None
```

```
In [15]: import csv
import pandas as pd
datapage2 = [
    'Middle',
]
pd_perfume = pd.DataFrame.from_dict(datapage2, orient='index').transpose()
print(pd_perfume.info())
pd_perfume.to_csv('perfume_middle.csv', 'a', index = False, encoding = 'utf_8_sig')

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 15687 entries, 0 to 15686
Data columns (total 1 columns):
Middle   15687 non-null object
dtypes: object(1)
memory usage: 122.6+ KB
None
```

```
In [16]: import csv
import pandas as pd
datapage3 = [
    'Base',
]
pd_perfume = pd.DataFrame.from_dict(datapage3, orient='index').transpose()
print(pd_perfume.info())
pd_perfume.to_csv('perfume_Base.csv', 'a', index = False, encoding = 'utf_8_sig')

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 15691 entries, 0 to 15690
Data columns (total 1 columns):
Base     15691 non-null object
dtypes: object(1)
memory usage: 122.7+ KB
None
```

```
In [20]: topall = []
for i in Main_info:
    i = i.strip('前调：')
    i = i.split('中调')
    top = i[0]
    if top != '':
        top = top.replace('，', '/n')
        top = top.split('/')
        Topall.append(top)
Top_every = []
for i in range(len(Topall)):
    Top.every.extend(Topall[i])
print(Top.every[0:100])
```

```
In [23]: from collections import Counter
Mdict = Counter(Middle_every)
Bdict = Counter(Base_every)
Tdict = sorted(Tdict.items(),key=lambda item:item[1],reverse=True)
Mdict = sorted(Mdict.items(),key=lambda item:item[1],reverse=True)
Bdict = sorted(Bdict.items(),key=lambda item:item[1],reverse=True)
```

```
In [28]: print(Tdict[10])
print(Mdict[10])
print(Bdict[10])
```

```
In [24]: len(Tdict)
Out[24]: 664
```

```
In [10]: len(Mdict)
Out[10]: 755
```

```
In [11]: len(Bdict)
Out[11]: 559
```

```
In [20]: Tlist=[]
for i in Tdict:
    if int(i[1]) <100:
        Tlist.append(i[0])
Mlist = []
for i in Mdict:
    if int(i[1]) <100:
        Mlist.append(i[0])
Blist = []
for i in Bdict:
    if int(i[1]) <100:
        Blist.append(i[0])
```

```
In [113]: len(Tlist)
Out[113]: 574
```

```
In [114]: len(Mlist)
Out[114]: 666
```

```
In [115]: len(Blist)
Out[115]: 516
```

```
In [309]: NTop= []
for topi in Top:
    ni = topi.split(',')
    try:
        for element in Tlist:
            if element in ni:
                ni.remove(element)
        NTop.append(ni)
    except:
        continue
```

```
NMiddle=[]
for midi in Middle:
    ni = midi.split(',')
    try:
        for element in Mlist:
            if element in ni:
                ni.remove(element)
        NMiddle.append(ni)
    except:
        continue
```

```
NBase=[]
for bai in Base:
    ni = bai.split(',')
    try:
        for element in Blist:
            if element in ni:
                ni.remove(element)
        NBase.append(ni)
    except:
        continue
```

# Step 3

## Extracted the variables.

*By using the same method as above, we extracted the variables according to the year, score, and number of evaluations.*

*Among them, the year is divided into three periods: before 1900, 1900-2010, and 2019. The score is 9.0 and above. The reply number is more than 1,000.*



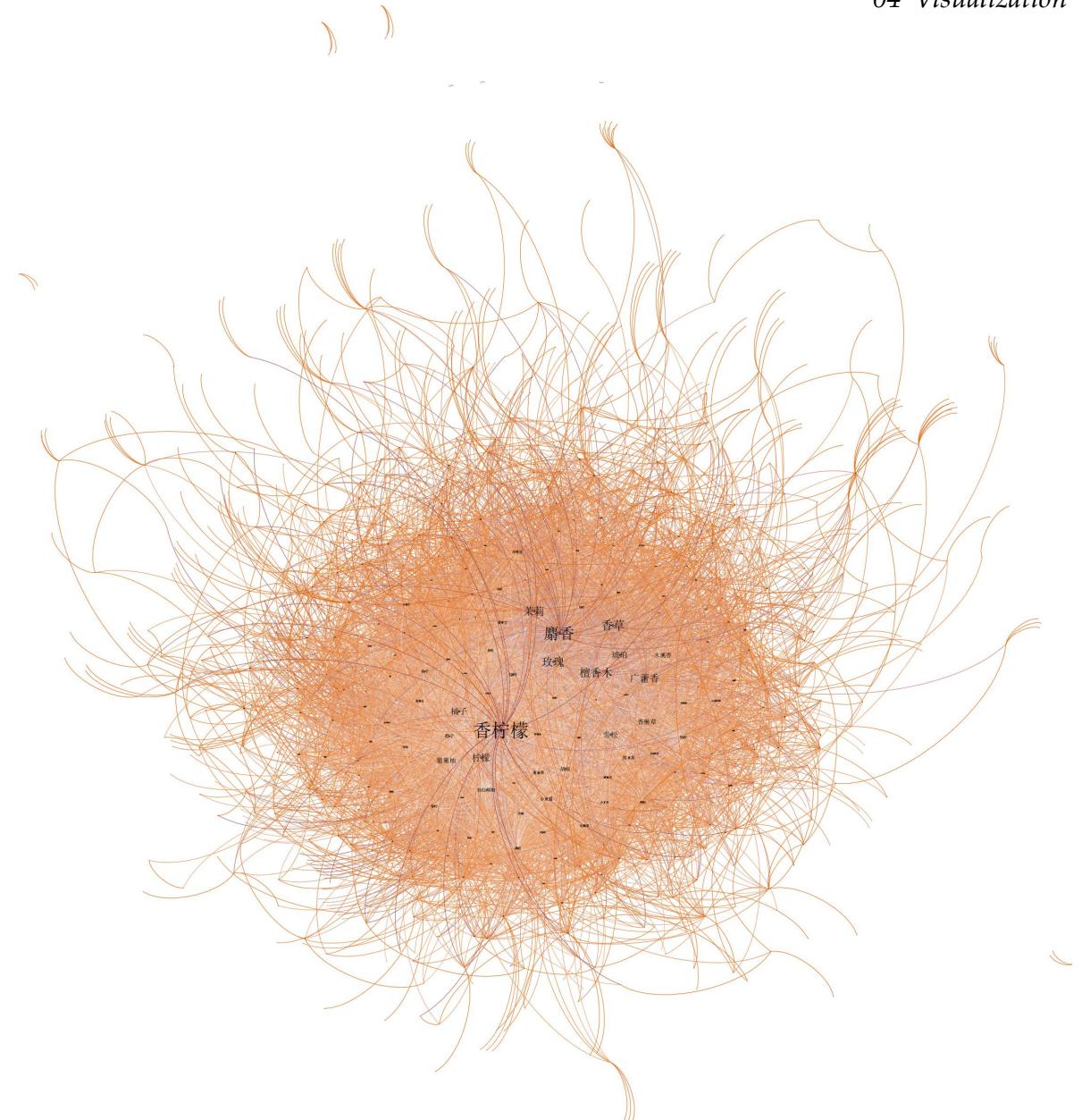
# 04 Visualization



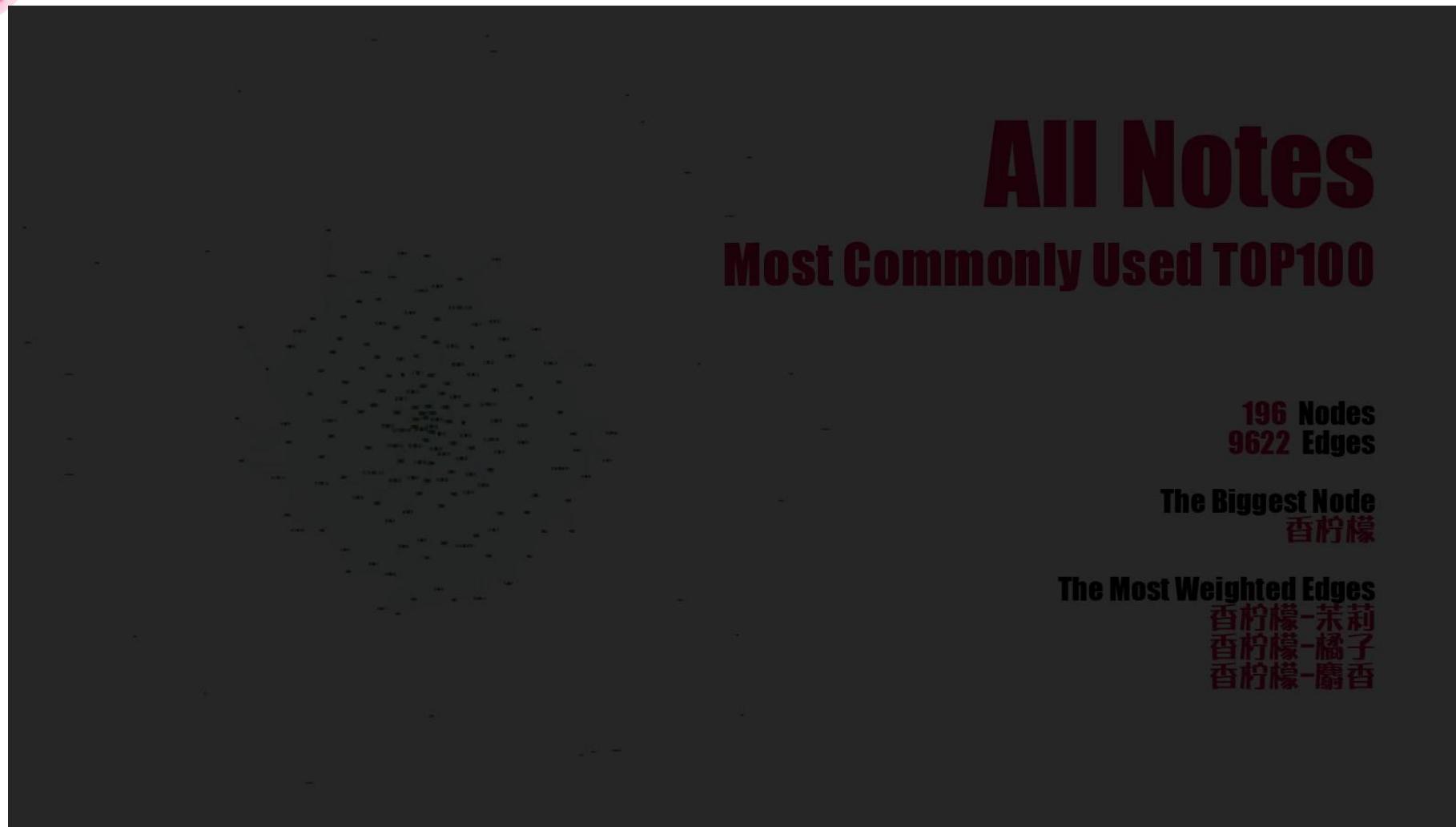


After data collection and analysis, we generated three clusters of the ingredients used in top notes, middle notes and base notes of perfumes. Then we drew several networks based on different variables.

We got 700+ nodes.



# Notes Commonly Used TOP100



# Notes Used Before 1900

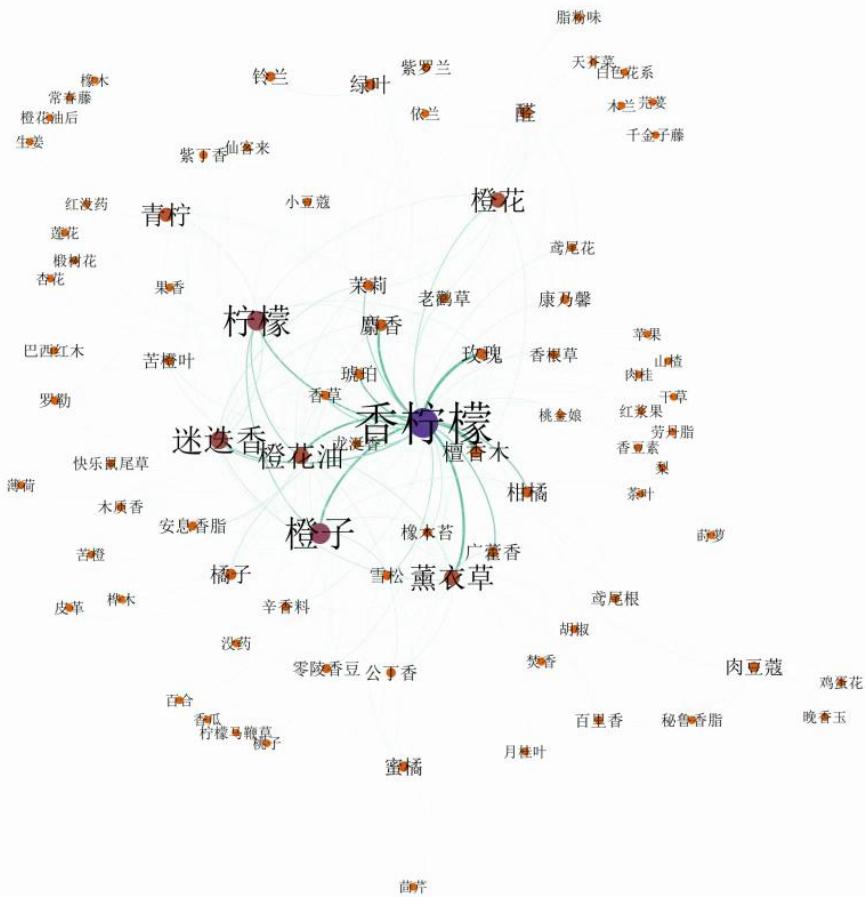
# All Notes Before 1900

**87 Nodes  
220 Edges**

## The Biggest Nodes

# The Most Weighted Edges

木香-玫瑰



# Notes Used After 1900



# Notes Used in 2019

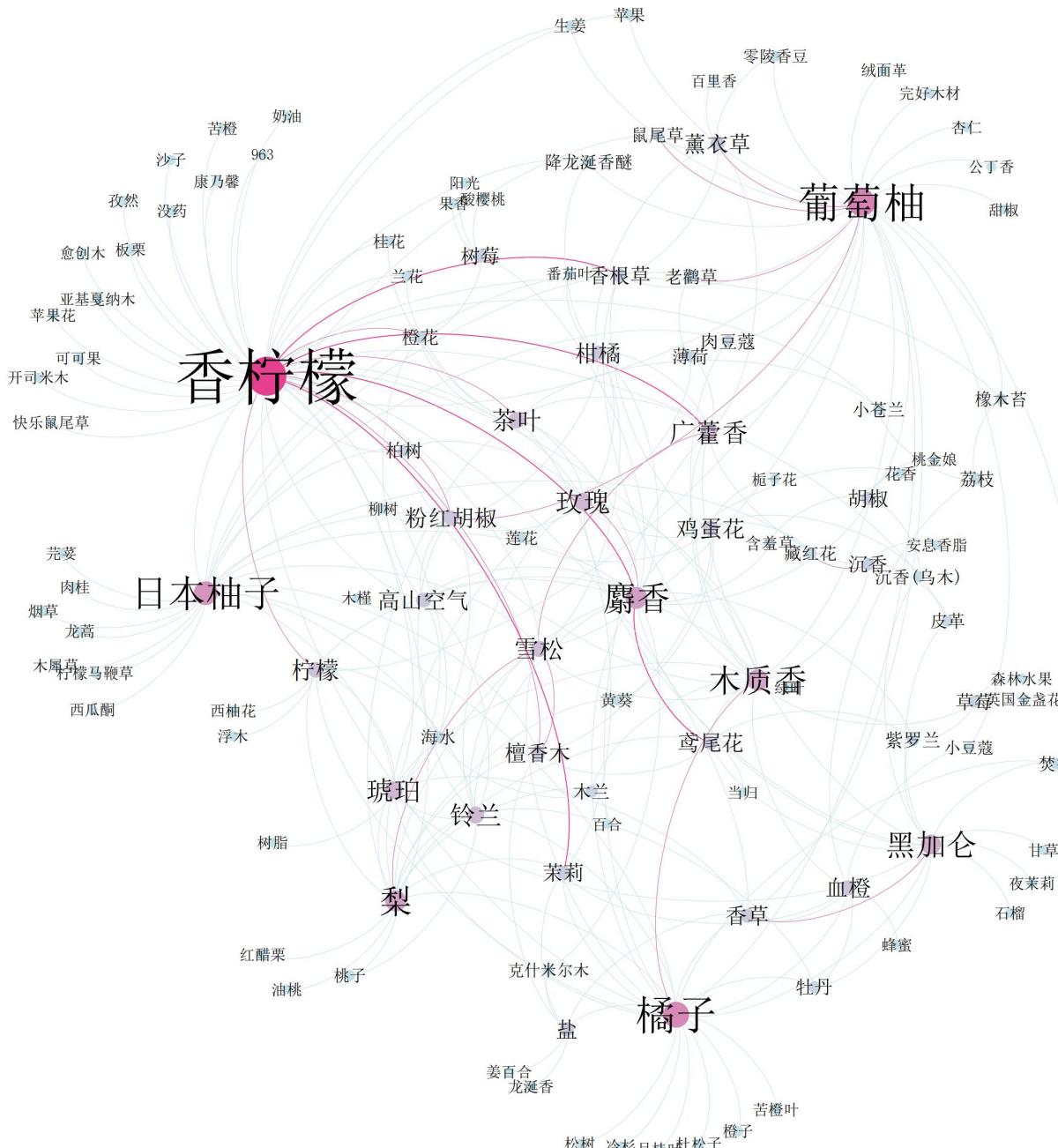
**119 Nodes  
253 Edges**

# The Biggest Nodes

## 香柠檬 葡萄柚

## The Most Weighted Edges

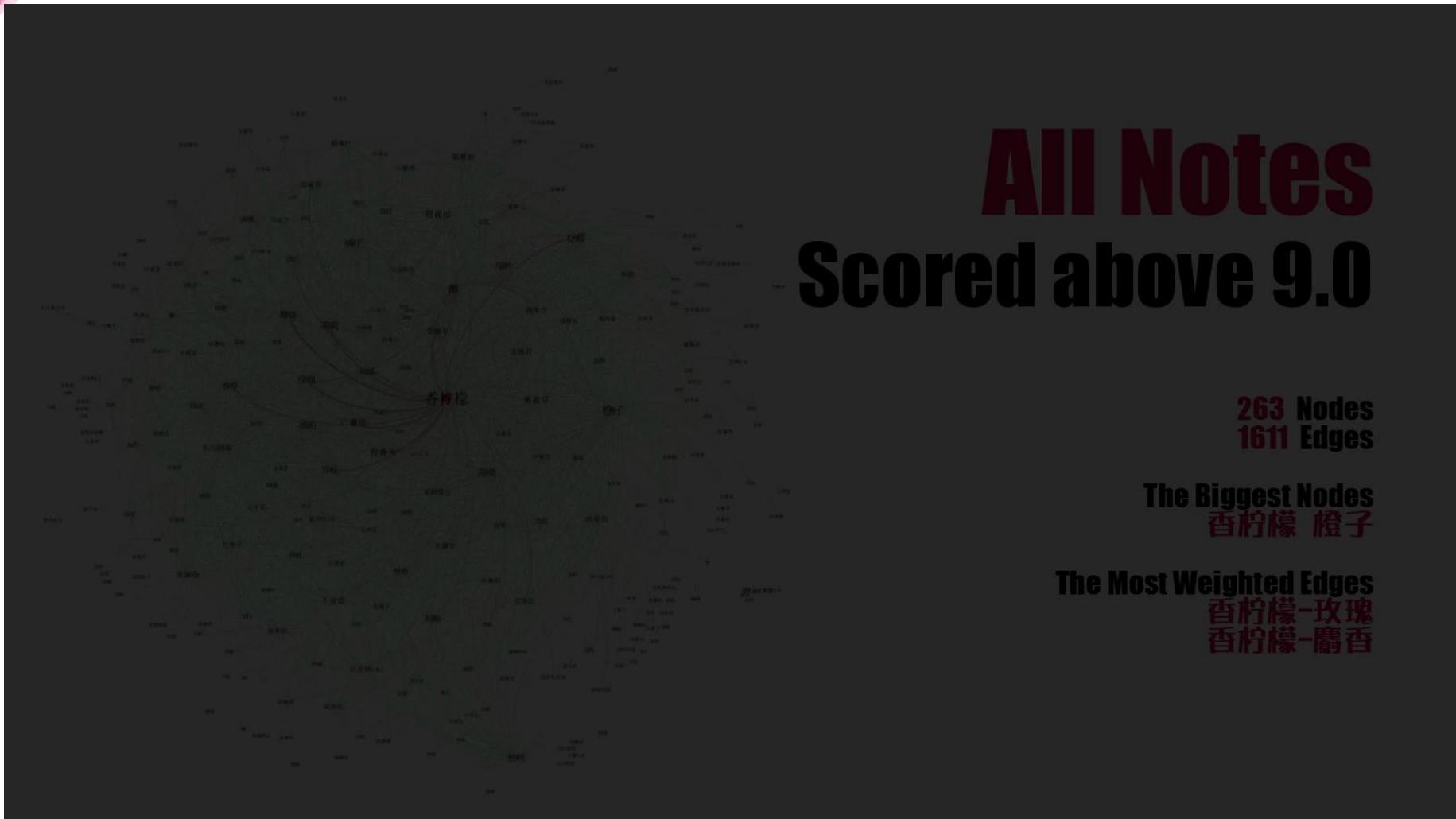
香柠檬-葡萄柚  
香柠檬-香根草  
麝香-鸢尾花



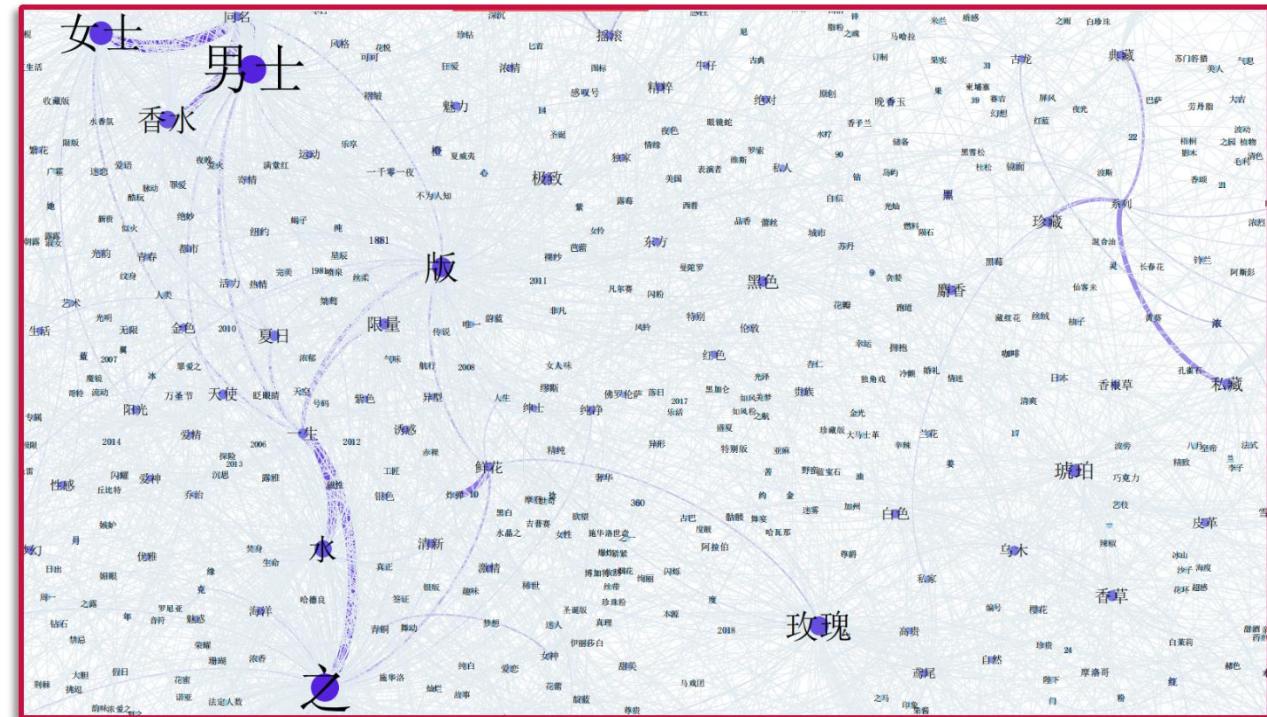
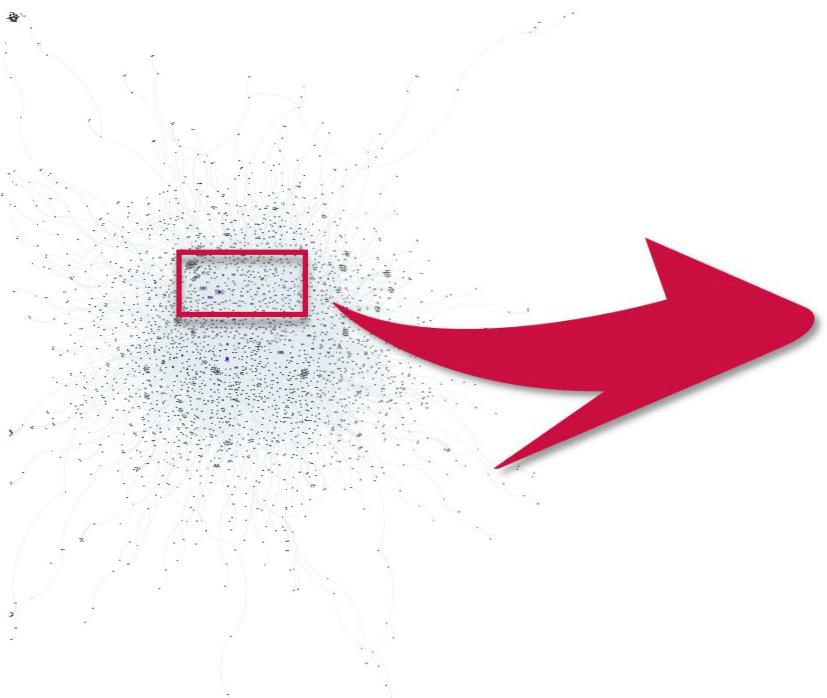
# Notes Most Commented



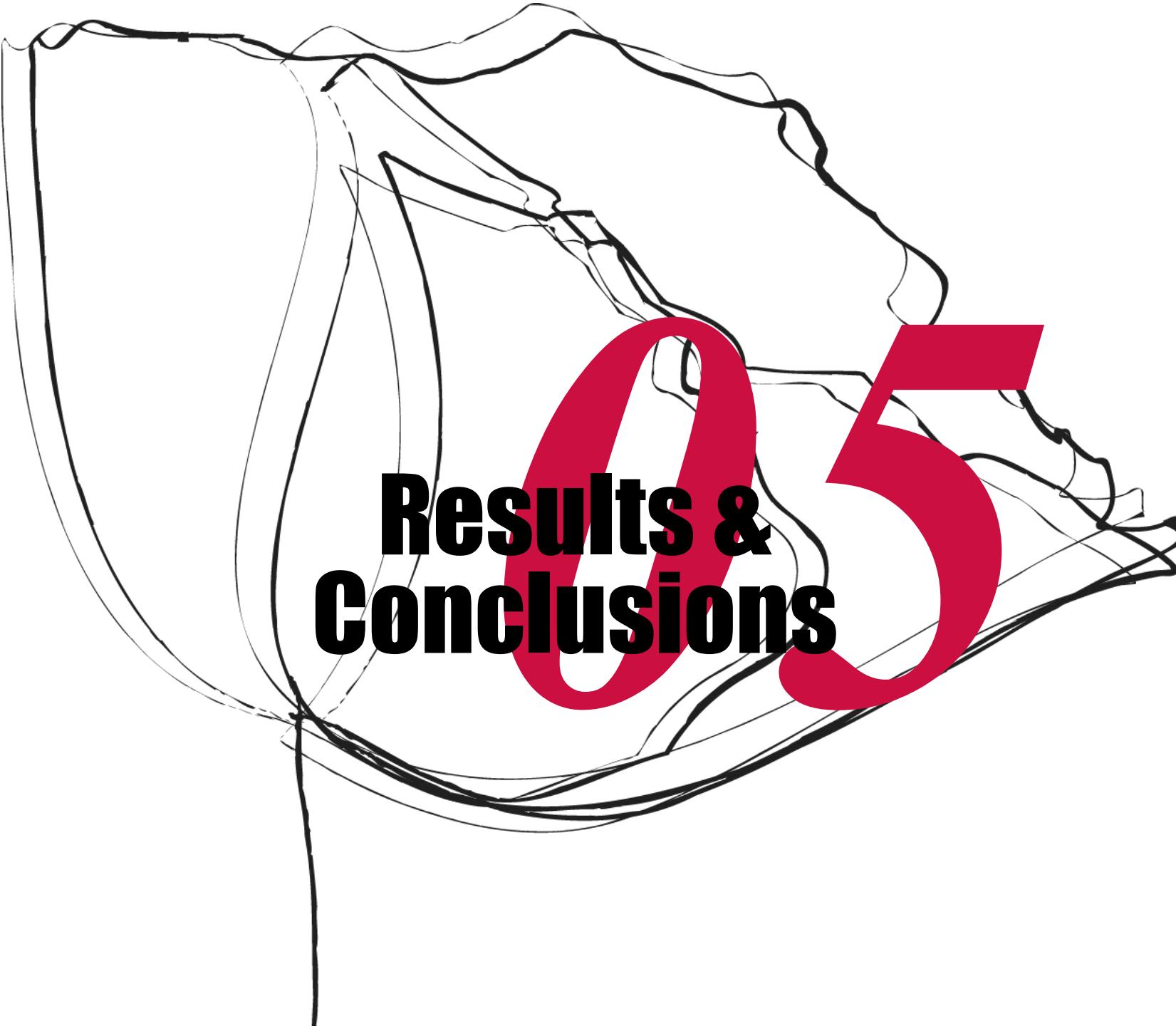
# Notes Scored above 9.0



# Brand Names in Chinese



**3586 Nodes  
6782 Edges**



# **Results & Conclusions**

The number '05' is overlaid on the bottom right of the text, with the '0' being white and the '5' being red.





*Although we got over 30 pictures about connected networks of perfumes, we still wonder that how and why any two nodes could be connected.*

*Another interesting thing is, some ingredients have been commonly used from ancient times to the present. Why are they so popular?*

*To answer these questions, we decided to **introduce a formula** to evaluate the performances of different ingredients.*

```
In [70]: 1 import csv
2 Escore = {}
3 wdct = {}
4 with open('all_score.csv', 'r', encoding="utf-8") as f:
5     for line in f:
6         line = line.strip("\n")
7         line = line.strip("\ufeff")
8         line = line.strip("\r\n")
9         line = line.split(",")
10        newline = []
11
12        for i in line:
13            if i!="":
14                newline.append(i)
15        wordcount = len(newline)-1
16        everyscore = float(newline[-1])/int(wordcount)
17
18        for i in newline[:-2]:
19            i = i.strip(",")
20            if i in wdct:
21                count = wdct[i]
22            else:
23                count=0
24            count = count + 1
25            wdct[i]=count
26
27            if i in Escore.keys():
28                wholescore = float(Escore[i]) + everyscore
29                Escore[i] = wholescore
30            else:
31                Escore[i] = everyscore
```

```
In [66]: 1 for k, v in Escore.items():
2     if wdct.get(k):
3         Escore[k]=round(float(Escore[k])/wdct[k],2)
4 print(Escore)
```

```
In [66]: 1 for k, v in Escore.items():
2     if wdct.get(k):
3         Escore[k]=round(float(Escore[k])/wdct[k],2)
4 print(Escore)
```

```
{'草莓': 0.91, '奶黄': 1.71, '橙子': 0.81, '苦橙': 0.93, '橘子': 0.81, '橙花': 0.82, '茉莉': 0.8, '桂花': 0.86, '牡丹': 0.89, '报春花': 1.34, '玫瑰': 0.83, '麝香': 0.79, '含羞草': 0.86, '康乃馨': 0.53, '晚香玉': 0.79, '醋栗叶': 0.7, '兰花': 0.73, '小苍兰': 0.82, '紫丁香': 0.79, '鸡蛋花': 0.88, '提亚蕾花': 0.86, '辛香料': 0.87, '琥珀': 0.77, '香草': 0.83, '快乐鼠尾草': 0.74, '焚香': 0.79, '蜂蜜': 0.79, '木质香': 1.0, '红没药': 0.66, '檀香脂': 0.77, '零陵香豆': 0.73, '安息香脂': 0.73, '黑巧克力': 1.0, '鸢尾花': 0.82, '芫荽': 0.65, '龙涎香': 0.82, '檀香木': 0.74, '广藿香': 0.78, '橡木苔': 0.64, '848': 2.4, '矢车菊': 2.0, '榛子': 1.06, '阿米香树花': 0.96, '鸢尾根': 0.61, '椰子': 0.8, '黄葵': 0.93, '阿魏': 0.51, '香柠檬': 0.75, '巴西红木': 0.6, '艾蒿': 0.66, '薄荷': 0.87, '绿叶': 0.85, '桃子': 0.68, '依兰': 0.67, '雪松': 0.79, '皮革': 0.77, '马黛茶': 0.94, '藏红花': 0.84, '肉桂': 0.7, '粉红胡椒': 0.87, '朗姆酒': 0.85, '莎草': 0.74, '沉香(乌木)': 0.93, '杜松子': 0.78, '紫罗兰叶': 0.81, '橙花油': 0.87, '铃兰': 0.71, '月桂叶': 0.67, '甜椒': 0.82, '没药': 0.79, '愈创木': 0.78, '葛缕子': 0.64, '白松香': 0.74, '罗勒': 0.7, '肉豆蔻': 0.7, '香根草': 0.78, '松树': 0.68, '老鹳草': 0.7, '安息香': 0.68, '葡萄柚': 0.85, '柠檬': 0.77, '仙客来': 0.68, '洋甘菊': 0.83, '胡椒': 0.8, '冷杉': 0.68, '黑加仑': 0.79, '杏仁': 0.9, '木兰': 0.81, '杜松': 0.77, '纸莎草': 0.78, '牛奶': 1.22, '百里香': 0.79, '薰衣草': 0.77, '咖啡': 0.88, '苹果': 0.81, '鼠尾草': 0.72, '梨': 0.89, '小豆蔻': 0.74, '生姜': 0.85, '苦橙叶': 0.88, '桦木': 0.75, '果香': 0.92, '树脂': 0.81, '焦糖': 0.91, '孜然': 0.74, '柏树': 0.77, '桃金娘': 0.78, '迷迭香': 0.71, '龙蒿': 0.57, '竹子': 0.8, '红浆果': 0.88, '八角': 0.85, '开司米木': 0.82, '柚子': 0.82, '树莓': 0.77, '可可果': 0.95, '百合': 0.86, '栀子花': 0.81, '巴西坚果': 3.13, '芭蕉叶': 0.61, '香蕉': 0.71, '白色花系': 1.19, '红醋栗': 0.97, '天芥菜': 0.75, '柠檬花': 0.84, '柑橘': 1.07, '紫罗兰': 0.78, '水仙花': 0.7, '忍冬': 0.82, '葫芦巴': 0.78, '闭鞘姜': 0.39, '乳香': 0.76, '风信子': 0.68, '金雀花': 0.98, '青草': 0.97, '木犀草': 0.72, '茴芹': 0.71, '干邑白兰地': 0.91, '蓬': 0.61, '莲花': 0.87, '椴树花': 0.79, '克里曼丁红橘': 0.83, '蜜橘': 0.78, '绒面革': 0.76, '干草': 1.05, '山楂': 0.75, '公丁香': 0.65, '无花果叶': 0.89, '蓝莓': 0.85, '黑莓': 0.84, '金属': 1.06, '剑兰': 1.33, '菊花': 0.91, '桃花心木': 0.74, '棉花': 1.18, '杏': 0.76, '香瓜': 0.71, '睡莲': 0.85, '香豌豆': 1.0, '洋槐': 0.73, '百香果': 0.88, '番石榴': 0.9, '芒果': 0.84, '西瓜': 0.85, '英国金盏花': 0.58, '野玫瑰果': 0.81, '西番莲': 0.93, '荔枝': 0.85, '红苹果': 0.95, '樱花': 1.0, '猕猴桃': 0.81, '桔梗': 0.84, '芍药': 0.84, '玫瑰': 0.84, '紫藤': 0.84, '紫罗兰': 0.84, '茉莉': 0.84}
```

$$\text{Score} = \frac{\sum (\text{perfume's score} / \text{the number of ingredients it contains})}{\text{the frequency of node of whole perfumes}}$$

1	巴西坚果	3.13
2	柳树	2.87
3	848	2.4
4	印蒿	2.33
5	矢车菊	2
6	苦艾酒	1.9
7	玉米穗	1.82
8	泥炭	1.76
9	顿加豆	1.74
10	奶黄	1.71
11	烈酒	1.7
12	麦芽	1.7
13	红色水果	1.68
14	云莓	1.68
15	乌木	1.66
16	番荔枝	1.64
17	面包	1.6
18	榛子可可酱	1.58
19	蛋白糖饼	1.51
20	鸦片	1.51
21	鱼子酱	1.47
22	树莓花	1.43
23	法式焦糖布丁	1.41
24	厚壳桂	1.39
25	油桃花	1.39
26	沙漠玫瑰	1.38
27	杜松子酒	1.37
28	牛肝菌	1.35
29	小麦	1.35
30	报春花	1.34



## High Socre But Less Used

Here are some possible reasons:

- a. Raw materials are **expensive**.
- b. **Complex** production process. (e.g., flavor components can easily volatilize )
- c. Some materials have strong smells, which leads to that they can **hardly find their "good partners"**.



370	橘子花	0.75
371	马天尼	0.75
372	蔓越莓	0.75
373	蔓长春花	0.75
374	木瓜	0.75
375	山楂	0.75
376	天芥菜	0.75
377	铁杉花	0.75
378	香柠檬	0.75
379	白松香	0.74
380	核桃	0.74
381	苦艾	0.74
382	快乐鼠尾草	0.74
383	李子	0.74
384	莎草	0.74
385	檀香木	0.74
386	桃花心木	0.74
387	小豆蔻	0.74
388	杏花	0.74
389	吲哚	0.74
390	柚木	0.74
391	芸香	0.74
392	孜然	0.74
393	安息香脂	0.73
394	菠萝	0.73
395	兰花	0.73
396	零陵香豆	0.73
397	玫瑰草	0.73
398	牛膝草	0.73

## Commonly Used But Relatively Low Score

*It might be due to:*

- a. Most of them are **common materials**.
- b. Light smell, usually serve as a **contrast or foil**. “甘作绿叶当陪衬”





## Limitations

- a. Did not make full use of all information we scraped from Nosetime.  
(Why we failed to create more complicated networks)
- b. The connections between brands could not be showed. (Time limitation)
- c. The features of nodes was not represented well.
- d. Most of our visualized outputs were static images.



*Thanks for your attention*

