

# Lecture 02: Clustering class demo

## Contents

- Hierarchical clustering on recipes dataset
- Let's cluster images!!



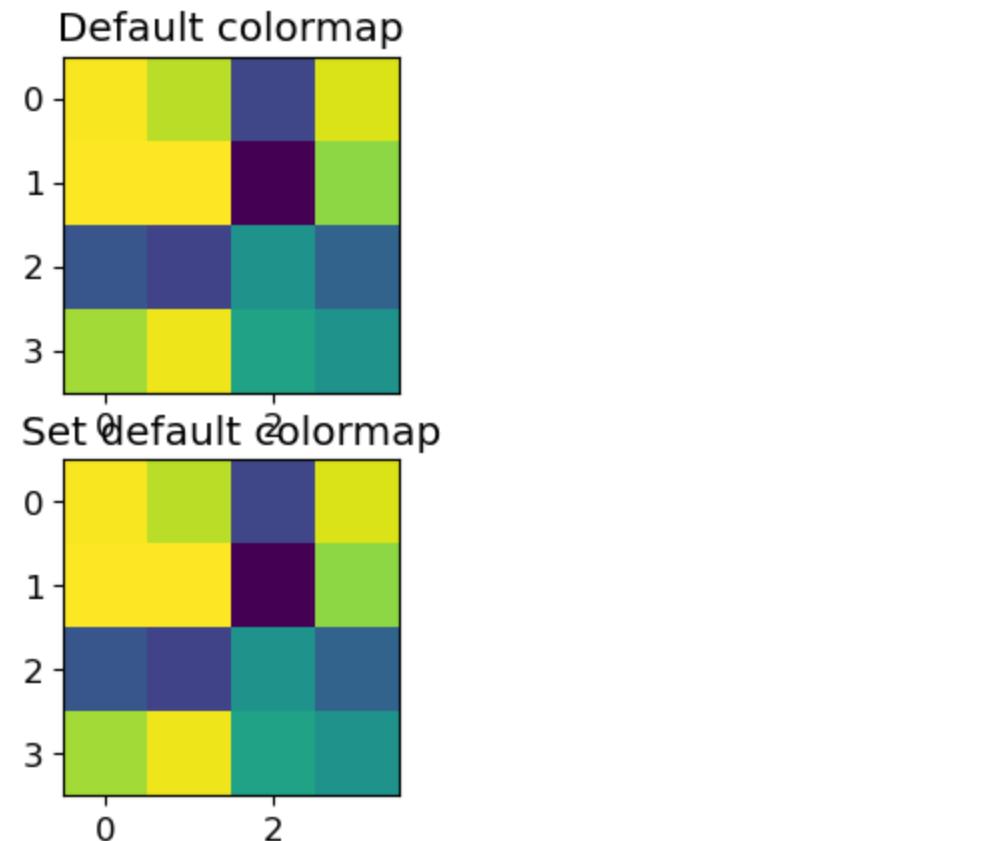
To run the code below, you need to install pytorch and torchvision in the course conda environment.

```
conda install pytorch torchvision -c pytorch
```

```

import numpy as np
import pandas as pd
import sys
import os
import torch, torchvision
from torchvision import datasets, models, transforms, utils
from PIL import Image
from torchvision import transforms
from torchvision.models import vgg16
import matplotlib.pyplot as plt
DATA_DIR = os.path.join(os.path.abspath("."), "data/")
sys.path.append(os.path.join(os.path.abspath("."), "code"))
from plotting_functions import *

```



## Hierarchical clustering on recipes dataset

Let's create a dendrogram on a more realistic dataset.

We'll use [Kaggle's Food.com recipes corpus](#). This corpus contains 180K+ recipes and 700K+ recipe reviews. In this lab, we'll only focus on recipes and **not** on reviews. The recipes are present in [RAW\\_recipes.csv](#). Our goal is to find categories or groupings of recipes from this corpus based on their names.

```
orig_recipes_df = pd.read_csv(DATA_DIR + "RAW_recipes.csv")
orig_recipes_df.shape
```

(231637, 12)

```
def get_recipes_sample(orig_recipes_df, n_tags=300, min_len=5):
    orig_recipes_df = orig_recipes_df.dropna() # Remove rows with NaNs.
    orig_recipes_df = orig_recipes_df.drop_duplicates(
        "name"
    ) # Remove rows with duplicate names.
    orig_recipes_df = orig_recipes_df[orig_recipes_df["name"].apply(len) >= min_len]
    first_n = orig_recipes_df["tags"].value_counts()[0:n_tags].index.tolist()
    recipes_df = orig_recipes_df[orig_recipes_df["tags"].isin(first_n)]
    return recipes_df
```

```
recipes_df = get_recipes_sample(orig_recipes_df)
recipes_df.shape
```

(9100, 12)

```
recipes_df["name"]
```

```
42          i yam what i yam  muffins
101         to your health  muffins
129         250 00 chocolate chip cookies
138                     lplermagronen
163         california roll  salad
...
231430      zucchini wheat germ cookies
231514      zucchini blueberry bread
231547      zucchini salsa burgers
231596      zuppa toscana
231629      zydeco salad
Name: name, Length: 9100, dtype: object
```

```
from sentence_transformers import SentenceTransformer
embedder = SentenceTransformer('all-MiniLM-L6-v2')
```



```
-----  
AttributeError Traceback (most recent call last)  
File ~/miniforge3/envs/jbook/lib/python3.12/site-packages/transformers/utils/integration/_utils.py:1862, in <module>  
     1862     try:  
--> 1863         return importlib.import_module("." + module_name, self.__name__)  
     1864     except Exception as e:  
  
File ~/miniforge3/envs/jbook/lib/python3.12/importlib/__init__.py:90, in import()  
     89     level += 1  
---> 90     return _bootstrap._gcd_import(name[level:], package, level)  
  
File <frozen importlib._bootstrap>:1387, in _gcd_import(name, package, level)  
  
File <frozen importlib._bootstrap>:1360, in _find_and_load(name, import_)  
  
File <frozen importlib._bootstrap>:1331, in _find_and_load_unlocked(name, import_)  
  
File <frozen importlib._bootstrap>:935, in _load_unlocked(spec)  
  
File <frozen importlib._bootstrap_external>:995, in exec_module(self, module)  
  
File <frozen importlib._bootstrap>:488, in _call_with_frames_removed(f, *args, **kwargs)  
  
File ~/miniforge3/envs/jbook/lib/python3.12/site-packages/transformers/integration/_utils.py:23, in <module>  
     23     import numpy as np  
---> 24     from tokenizers import Tokenizer, decoders, normalizers, pre_tokenizers  
     25     from tokenizers.models import BPE, Unigram  
  
File ~/miniforge3/envs/jbook/lib/python3.12/site-packages/tokenizers/__init__.py:78, in <module>  
     78     from .tokenizers import (  
     79         AddedToken,  
     80         Encoding,  
     81         ...  
     92         __version__,  
     93     )  
---> 94     from .implementations import (  
     95         BertWordPieceTokenizer,  
     96         ByteLevelBPETokenizer,  
     97         CharBPETokenizer,  
     98         SentencePieceBPETokenizer,  
     99         SentencePieceUnigramTokenizer,  
    100     )  
  
File ~/miniforge3/envs/jbook/lib/python3.12/site-packages/tokenizers/implementations/_base_tokenizer.py:1, in <module>  
----> 1     from .base_tokenizer import BaseTokenizer  
     2     from .bert_wordpiece import BertWordPieceTokenizer  
  
File ~/miniforge3/envs/jbook/lib/python3.12/site-packages/tokenizers/implementations/_decoders/_decoder.py:3, in <module>  
     3     from tokenizers import AddedToken, EncodeInput, Encoding, InputSequence  
----> 4     from tokenizers.decoders import Decoder  
     5     from tokenizers.models import Model  
  
File ~/miniforge3/envs/jbook/lib/python3.12/site-packages/tokenizers/decoders/_sequence.py:14, in <module>  
     14     Sequence = decoders.Sequence  
----> 15     DecodeStream = decoders.DecodeStream
```

```
AttributeError: module 'decoders' has no attribute 'DecodeStream'
```

The above exception was the direct cause of the following exception:

```
RuntimError                                     Traceback (most recent call last)
File ~/miniforge3/envs/jbook/lib/python3.12/site-packages/transformers/utils/ir
  1862     try:
-> 1863         return importlib.import_module("." + module_name, self.__name__)
  1864     except Exception as e:
...
File ~/miniforge3/envs/jbook/lib/python3.12/importlib/__init__.py:90, in import_
   89         level += 1
---> 90     return _bootstrap._gcd_import(name[level:], package, level)
...
File <frozen importlib._bootstrap>:1387, in _gcd_import(name, package, level)
...
File <frozen importlib._bootstrap>:1360, in _find_and_load(name, import_
...
File <frozen importlib._bootstrap>:1331, in _find_and_load_unlocked(name, impo
...
File <frozen importlib._bootstrap>:935, in _load_unlocked(spec)
...
File <frozen importlib._bootstrap_external>:995, in exec_module(self, module)
...
File <frozen importlib._bootstrap>:488, in _call_with_frames_removed(f, *args,
...
File ~/miniforge3/envs/jbook/lib/python3.12/site-packages/transformers/models/a
  22     from typing import List, Union
---> 24     from ...configuration_utils import PretrainedConfig
  25     from ...dynamic_module_utils import get_class_from_dynamic_module, reso
...
File ~/miniforge3/envs/jbook/lib/python3.12/site-packages/transformers/configu
  28     from .dynamic_module_utils import custom_object_save
---> 29     from .modeling_gguf_pytorch_utils import load_gguf_checkpoint
  30     from .utils import (
  31         CONFIG_NAME,
  32         PushToHubMixin,
  (...))
  41         logging,
  42     )
...
File ~/miniforge3/envs/jbook/lib/python3.12/site-packages/transformers/modeling
  21     from tqdm.auto import tqdm
---> 23     from .integrations import (
  24         GGUF_CONFIG_MAPPING,
  25         GGUF_TOKENIZER_MAPPING,
  26         _gguf_parse_value,
  27     )
  28     from .utils import is_torch_available
...
File <frozen importlib._bootstrap>:1412, in _handle_fromlist(module, fromlist,
...
File ~/miniforge3/envs/jbook/lib/python3.12/site-packages/transformers/utils/ir
  1850     elif name in self._class_to_module.keys():
-> 1851         module = self._get_module(self._class_to_module[name])
```

```

1852     value = getattr(module, name)

File ~/miniforge3/envs/jbook/lib/python3.12/site-packages/transformers/utils/imports.py:1864
1864 except Exception as e:
-> 1865     raise RuntimeError(
1866         f"Failed to import {self.__name__}.{module_name} because of the following"
1867         f" traceback):\n{e}"
1868 ) from e

```

`RuntimeError: Failed to import transformers.integrations.ggml because of the following`  
`module 'decoders' has no attribute 'DecodeStream'`

The above exception was the direct cause of the following exception:

```

RuntimeError                                     Traceback (most recent call last)
Cell In[6], line 1
----> 1 from sentence_transformers import SentenceTransformer
      2 embedder = SentenceTransformer('all-MiniLM-L6-v2')

File ~/miniforge3/envs/jbook/lib/python3.12/site-packages/sentence_transformers/12
      7 import os
      9 from sentence_transformers.backend import (
10     export_dynamic_quantized_onnx_model,
11     export_optimized_onnx_model,
12     export_static_quantized_openvino_model,
13 )
----> 14 from sentence_transformers.cross_encoder.CrossEncoder import CrossEncoder
15 from sentence_transformers.datasets import ParallelSentencesDataset, Seq2SeqDataset
16 from sentence_transformers.LoggingHandler import LoggingHandler

File ~/miniforge3/envs/jbook/lib/python3.12/site-packages/sentence_transformers/12
      1 from __future__ import annotations
----> 3 from .CrossEncoder import CrossEncoder
      5 __all__ = ["CrossEncoder"]

File ~/miniforge3/envs/jbook/lib/python3.12/site-packages/sentence_transformers/12
      12 from torch.utils.data import DataLoader
      13 from tqdm.autonotebook import tqdm, trange
----> 14 from transformers import AutoConfig, AutoModelForSequenceClassification
      15 from transformers.tokenization_utils_base import BatchEncoding
      16 from transformers.utils import PushToHubMixin

File <frozen importlib._bootstrap>:1412, in _handle_fromlist(module, fromlist,
      1
File ~/miniforge3/envs/jbook/lib/python3.12/site-packages/transformers/utils/imports.py:1850
      1850 elif name in self._class_to_module.keys():
      1851     module = self._get_module(self._class_to_module[name])
-> 1852     value = getattr(module, name)
      1853 elif name in self._modules:
      1854     value = self._get_module(name)

File ~/miniforge3/envs/jbook/lib/python3.12/site-packages/transformers/utils/imports.py:1849
      1849     value = Placeholder
      1850 elif name in self._class_to_module.keys():
-> 1851     module = self._get_module(self._class_to_module[name])
      1852     value = getattr(module, name)

```

```
1853 elif name in self._modules:
```

```
File ~/miniforge3/envs/jbook/lib/python3.12/site-packages/transformers/utils/imports.py:1863
1863     return importlib.import_module("." + module_name, self.__name__)
1864 except Exception as e:
-> 1865     raise RuntimeError(
1866         f"Failed to import {self.__name__}.{module_name} because of the"
1867         f" traceback):\n{e}"
1868     ) from e
```

`RuntimeError: Failed to import transformers.models.auto.configuration_auto because of the following error:`  
`Failed to import transformers.integrations.ggml because of the following error:`  
`module 'decoders' has no attribute 'DecodeStream'`

```
recipe_names = recipes_df["name"].tolist()
embeddings = embedder.encode(recipe_names)
recipe_names_embeddings = pd.DataFrame(
    embeddings,
    index=recipes_df.index,
)
recipe_names_embeddings
```

	0	1	2	3	4	5	
<b>42</b>	0.019592	-0.088336	0.072677	-0.034575	-0.048741	-0.049801	0
<b>101</b>	-0.000567	-0.011825	0.073199	0.058175	0.031688	-0.015428	0
<b>129</b>	-0.022604	0.065034	-0.033065	0.014450	-0.105039	-0.050559	0
<b>138</b>	-0.066915	0.025988	-0.087689	-0.006847	-0.012861	0.049035	0
<b>163</b>	-0.007068	-0.007308	-0.026629	-0.004153	-0.052810	0.011126	0.
...	...	...	...	...	...	...	...
<b>231430</b>	-0.036366	0.087173	-0.039641	0.002705	0.097142	-0.075385	0.
<b>231514</b>	-0.052718	0.008980	-0.046014	0.030194	0.005201	0.009964	-0.
<b>231547</b>	-0.080801	0.004295	-0.044325	0.038307	-0.030125	-0.063566	0.
<b>231596</b>	-0.060801	0.111672	-0.108755	0.052323	-0.099851	-0.027532	0.
<b>231629</b>	-0.085059	0.009065	-0.088442	-0.008907	0.002234	-0.055801	0.

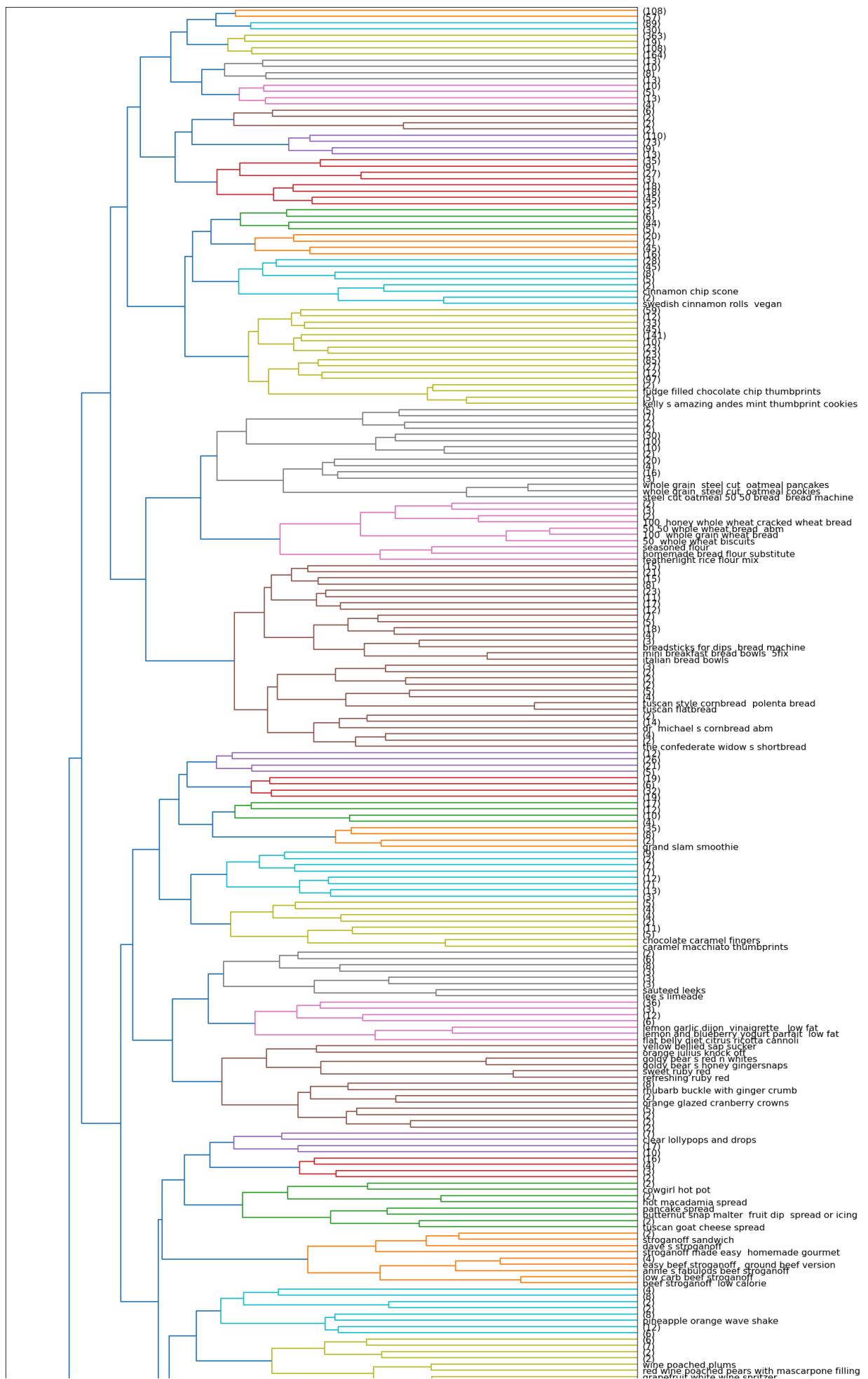
9100 rows × 384 columns

```
def plot_dendrogram(Z, labels, w=15, h=20, p=10):
    fig, ax = plt.subplots(figsize=(w, h))
    dendrogram(
        Z,
        p=p,
        truncate_mode="level",
        orientation="left",
        leaf_font_size=12,
        labels=labels,
        ax=ax,
    )
    ax = plt.gca()
    ax.set_ylabel("Cluster distances", fontsize=w)
```

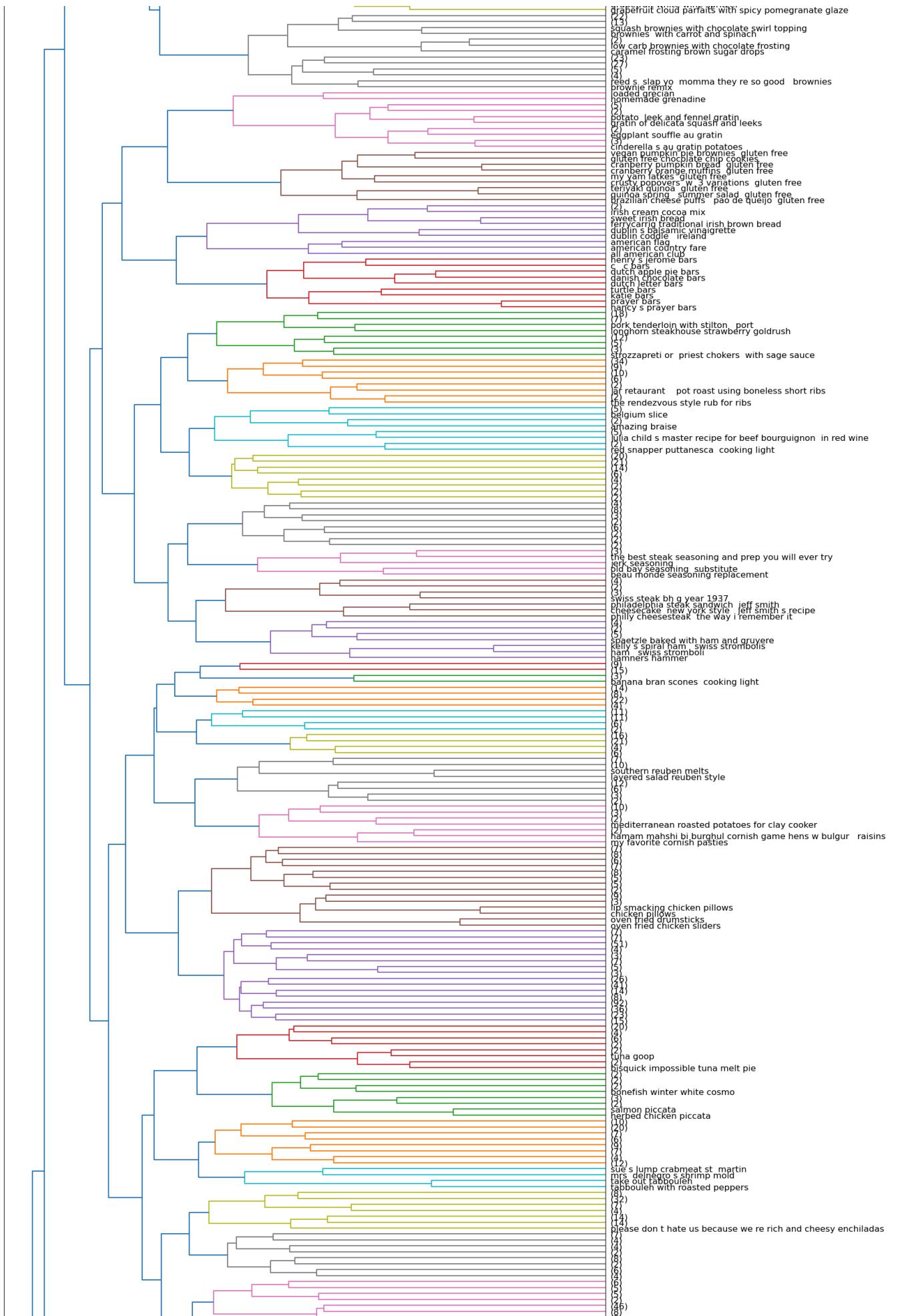
```
from scipy.cluster.hierarchy import dendrogram, linkage

recipe_names = recipes_df["name"].tolist()
Z = linkage(embeddings, method="complete", metric="cosine")
```

```
plot_dendrogram(Z, labels=recipe_names, w=15, h=230, p=10)
```



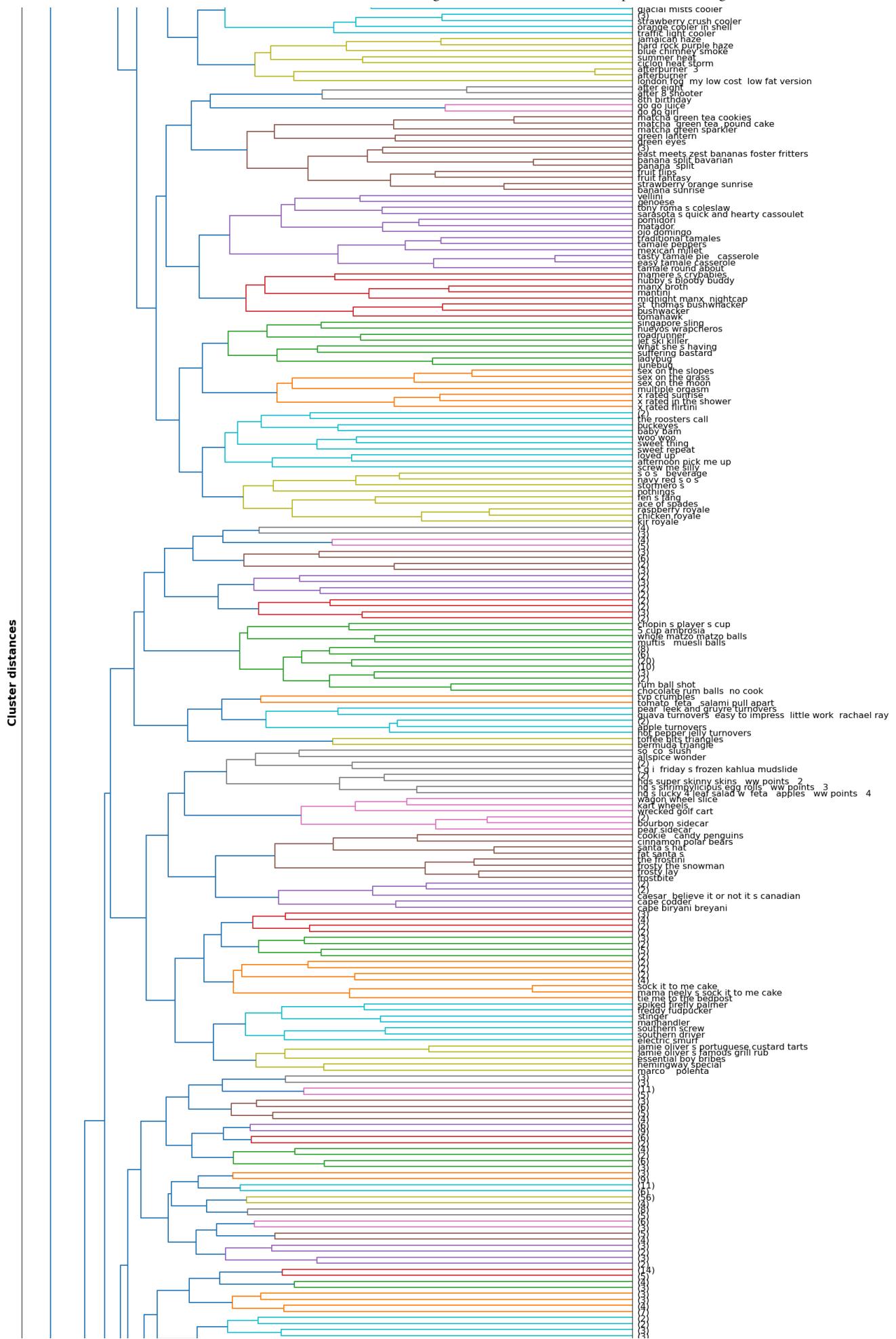
Lecture 02: Clustering class demo — DSCI 563 Unsupervised Learning

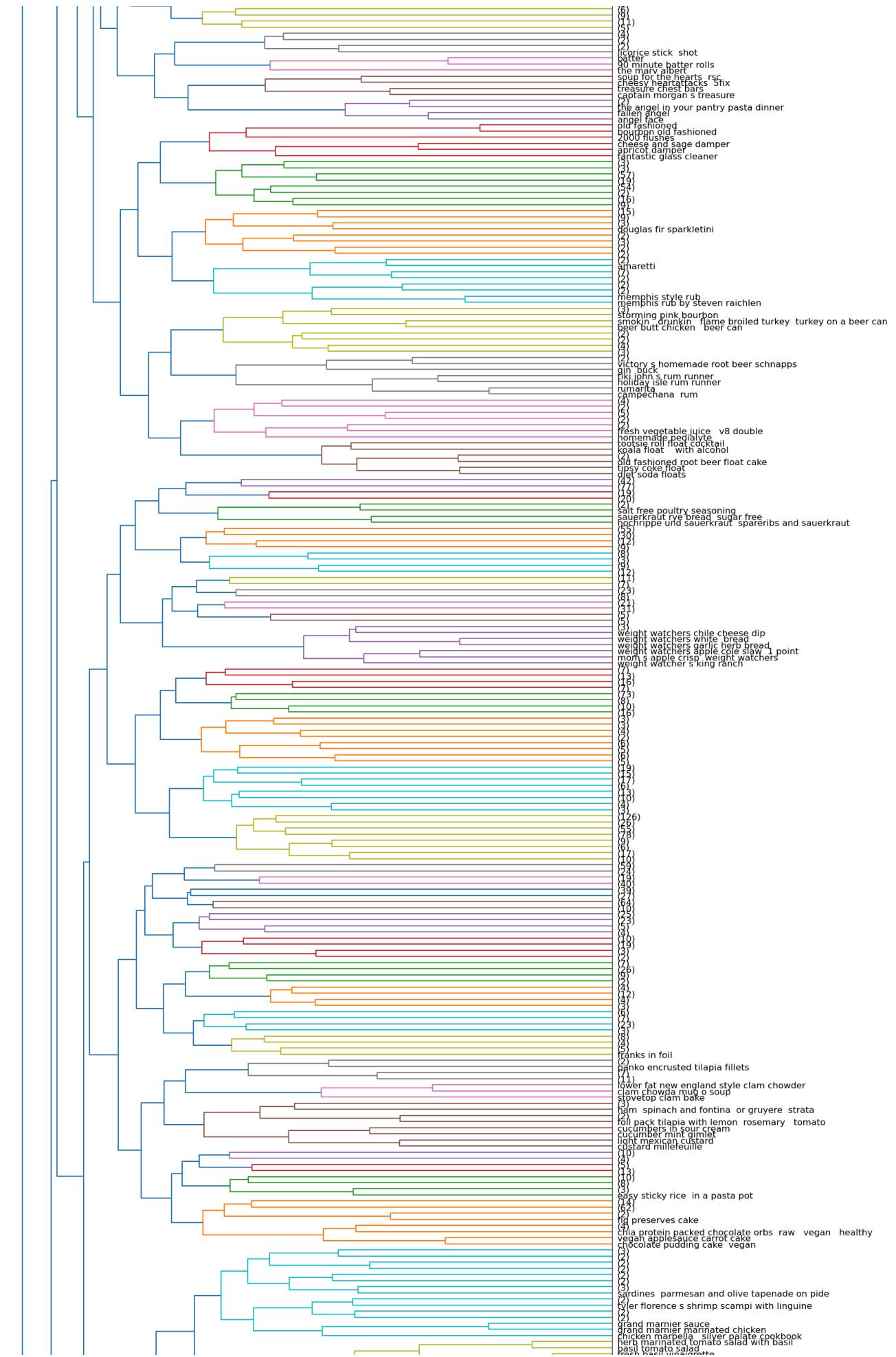


Lecture 02: Clustering class demo — DSCI 563 Unsupervised Learning



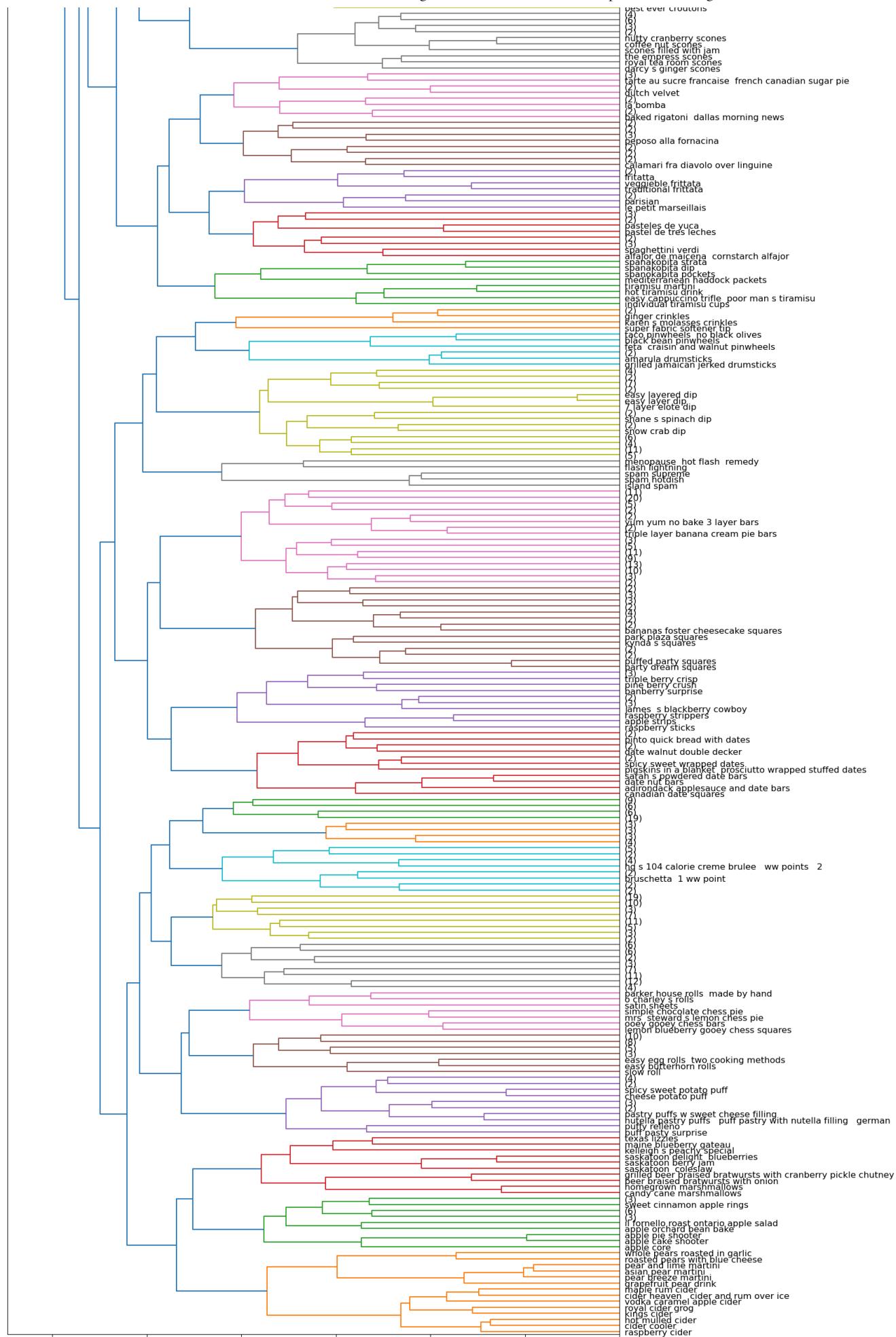
Lecture 02: Clustering class demo — DSCI 563 Unsupervised Learning





Lecture 02: Clustering class demo — DSCI 563 Unsupervised Learning





# Let's cluster images!!

For this demo, I'm going to use the following image dataset:

1. A tiny subset of [Food-101](#) from last lecture (available [here](#)).
2. A small subset of [Human Faces dataset](#) (available [here](#)).

Let's start with small subset of birds dataset. You can experiment with a bigger dataset if you like.

```
device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
```

```
import random
def set_seed(seed=42):
    torch.manual_seed(seed)
    np.random.seed(seed)
    random.seed(seed)
```

```
set_seed(seed=42)
```

```

import glob
IMAGE_SIZE = 224
def read_img_dataset(data_dir, BATCH_SIZE):
    data_transforms = transforms.Compose(
        [
            transforms.Resize((IMAGE_SIZE, IMAGE_SIZE)),
            transforms.ToTensor(),
            transforms.Normalize([0.5, 0.5, 0.5], [0.5, 0.5, 0.5]),
        ])

    image_dataset = datasets.ImageFolder(root=data_dir, transform=data_transfo
    dataloader = torch.utils.data.DataLoader(
        image_dataset, batch_size=BATCH_SIZE, shuffle=True, num_workers=0
    )
    dataset_size = len(image_dataset)
    class_names = image_dataset.classes
    inputs, classes = next(iter(dataloader))
    return inputs, classes

```

```

def plot_sample_imgs(inputs):
    plt.figure(figsize=(10, 70)); plt.axis("off"); plt.title("Sample Training")
    plt.imshow(np.transpose(utils.make_grid(inputs, padding=1, normalize=True))

```

```

def get_features(model, inputs):
    """Extract output of densenet model"""
    model.eval()
    with torch.no_grad(): # turn off computational graph stuff
        Z = model(inputs).detach().numpy()
    return Z

```

```

densenet = models.densenet121(weights="DenseNet121_Weights.IMGNET1K_V1")
densenet.classifier = torch.nn.Identity() # remove that last "classification"

```

```

data_dir = "../data/food"
file_names = [image_file for image_file in glob.glob(data_dir + "/*/*.jpg")]
n_images = len(file_names)
BATCH_SIZE = n_images # because our dataset is quite small
food_inputs, food_classes = read_img_dataset(data_dir, BATCH_SIZE)
n_images

```

350

```
X_food = food_inputs.numpy()
```

```
plot_sample_imgs(food_inputs[0:24,:,:,:])
```

Sample Training Images



```
Z_food = get_features(
    densenet, food_inputs,
)
```

Z\_food.shape

(350, 1024)

```
from sklearn.cluster import KMeans
k = 7
km = KMeans(n_clusters=k, n_init='auto', random_state=123)
km.fit(Z_food)
```

huggingface/tokenizers: The current process just got forked, after parallelism  
To disable this warning, you can either:

- Avoid using `tokenizers` before the fork if possible
- Explicitly set the environment variable TOKENIZERS\_PARALLELISM=(true



KMeans



KMeans(n\_clusters=7, random\_state=123)

```
for cluster in range(k):
    get_cluster_images(km, Z_food, X_food, cluster, n_img=6)
```

82

Image indices: [ 82 334 265 35 114 343]

## Cluster center 0



276

Image indices: [276 191 186 29 265 343]

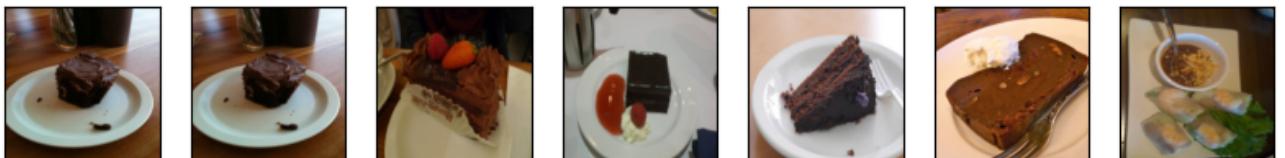
## Cluster center 1



161

Image indices: [161 102 190 241 56 76]

## Cluster center 2



76

Image indices: [ 76 39 259 295 313 138]

## Cluster center 3



124

Image indices: [124 253 25 27 223 280]

**Cluster center 4**

214

Image indices: [214 20 141 249 203 274]

**Cluster center 5**

48

Image indices: [ 48 76 61 5 89 138]

**Cluster center 6**

Let's try DBSCAN.

```
dbSCAN = DBSCAN()

labels = dbSCAN.fit_predict(Z_food)
print("Unique labels: {}".format(np.unique(labels)))
```

Unique labels: [-1]

It identified all points as noise points. Let's explore the distances between points.

```
from sklearn.metrics.pairwise import euclidean_distances

dists = euclidean_distances(Z_food)
np.fill_diagonal(dists, np.inf)
dists_df = pd.DataFrame(dists)
dists_df
```

	<b>0</b>	<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>	<b>5</b>	
<b>0</b>	inf	27.170071	23.031963	29.344770	27.464485	25.905249	28.6
<b>1</b>	27.170071	inf	22.623034	24.378839	25.491880	21.179726	25.2
<b>2</b>	23.031963	22.623034	inf	28.764124	25.544821	23.159750	26.1
<b>3</b>	29.344770	24.378839	28.764124	inf	28.520519	23.345842	28.1
<b>4</b>	27.464485	25.491880	25.544821	28.520519	inf	21.172096	26.4
...	...	...	...	...	...	...	...
<b>345</b>	26.458879	23.107458	25.684532	21.547804	25.304211	22.533875	26.8
<b>346</b>	28.636688	21.258495	26.029633	26.480219	25.343542	21.611742	26.2
<b>347</b>	23.513239	19.471357	20.417009	24.750483	23.945749	19.392572	23.9
<b>348</b>	26.565344	20.257633	21.786682	25.170492	23.584101	19.035461	22.5
<b>349</b>	27.202259	22.101439	22.903080	26.520855	25.460316	21.709145	24.9

350 rows × 350 columns

```
dists.min(), np.nanmax(dists[dists != np.inf]), np.mean(dists[dists != np.inf])
```

```
(np.float32(10.06717), np.float32(36.652683), np.float32(24.538565))
```

```
for eps in range(13, 20):
    print("\neps={}".format(eps))
    dbscan = DBSCAN(eps=eps, min_samples=3)
    labels = dbscan.fit_predict(Z_food)
    print("Number of clusters: {}".format(len(np.unique(labels))))
    print("Cluster sizes: {}".format(np.bincount(labels + 1)))
```

```

eps=13
Number of clusters: 2
Cluster sizes: [347 3]

eps=14
Number of clusters: 5
Cluster sizes: [334 3 6 4 3]

eps=15
Number of clusters: 4
Cluster sizes: [299 26 8 17]

eps=16
Number of clusters: 4
Cluster sizes: [248 86 3 13]

eps=17
Number of clusters: 2
Cluster sizes: [205 145]

eps=18
Number of clusters: 2
Cluster sizes: [160 190]

eps=19
Number of clusters: 2
Cluster sizes: [116 234]

```

```

dbscan = DBSCAN(eps=14, min_samples=3)
dbscan_labels = dbscan.fit_predict(Z_food)
print("Number of clusters: {}".format(len(np.unique(dbscan_labels))))
print("Cluster sizes: {}".format(np.bincount(dbscan_labels + 1)))
print("Unique labels: {}".format(np.unique(dbscan_labels)))

```

```

Number of clusters: 5
Cluster sizes: [334 3 6 4 3]
Unique labels: [-1 0 1 2 3]

```

```
print_dbscan_clusters(Z_food, food_inputs, dbscan_labels)
```

cluster 0



cluster 0



cluster 0



cluster 1



cluster 1



cluster 1



cluster 1



cluster 1



cluster 1



cluster 2



cluster 2



cluster 2



cluster 2



cluster 3



cluster 3



cluster 3



Let's examine noise points identified by DBSCAN.

```
print_dbSCAN_noise_images(Z_food, food_inputs, dbSCAN_labels)
```