
CLUSTERING

Unsupervised classification / clustering

Unsupervised classification

- ▶ **Input:** $x_1, x_2, \dots, x_n \in \mathbb{R}^d$, target cardinality $k \in \mathbb{N}$.
- ▶ **Output:** function $f: \mathbb{R}^d \rightarrow \{1, 2, \dots, k\} =: [k]$.
- ▶ **Typical semantics:** hidden subpopulation structure.

Unsupervised classification / clustering

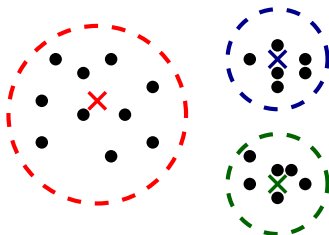
Unsupervised classification

- ▶ **Input:** $x_1, x_2, \dots, x_n \in \mathbb{R}^d$, target cardinality $k \in \mathbb{N}$.
- ▶ **Output:** function $f: \mathbb{R}^d \rightarrow \{1, 2, \dots, k\} =: [k]$.
- ▶ **Typical semantics:** hidden subpopulation structure.

Clustering

- ▶ **Input:** $x_1, x_2, \dots, x_n \in \mathbb{R}^d$, target cardinality $k \in \mathbb{N}$.
 - ▶ **Output:** partitioning of x_1, x_2, \dots, x_n into k groups.
 - ▶ Often done via unsupervised classification;
 \Rightarrow “clustering” often synonymous with “unsupervised classification”.
 - ▶ Sometimes also have a “representative” $c_j \in \mathbb{R}^d$ for each $j \in [k]$
(e.g., average of the x_i in j th group) \longrightarrow **quantization**.
-

Unsupervised classification / clustering



Clustering

- ▶ **Input:** $x_1, x_2, \dots, x_n \in \mathbb{R}^d$, target cardinality $k \in \mathbb{N}$.
 - ▶ **Output:** partitioning of x_1, x_2, \dots, x_n into k groups.
 - ▶ Often done via unsupervised classification;
 \Rightarrow “clustering” often synonymous with “unsupervised classification”.
 - ▶ Sometimes also have a “representative” $c_j \in \mathbb{R}^d$ for each $j \in [k]$
(e.g., average of the x_i in j th group) \rightarrow **quantization**.
-

Uses of clustering: feature representations

“One-hot” / “dummy variable” encoding of $f(\mathbf{x})$

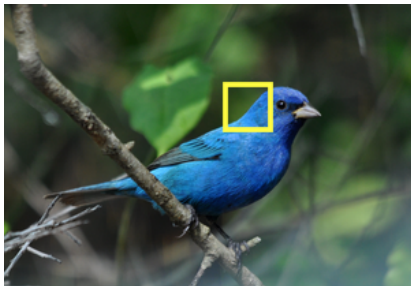
$$\phi(\mathbf{x}) = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \leftarrow f(\mathbf{x}) \text{ position}$$

(Often used together with other features.)

Uses of clustering: feature representations

Histogram representation

- ▶ Cut up each $\mathbf{x}_i \in \mathbb{R}^d$ into different parts $\mathbf{x}_{i,1}, \mathbf{x}_{i,2}, \dots, \mathbf{x}_{i,m} \in \mathbb{R}^p$ (e.g., small patches of an image) .
- ▶ Cluster all the parts $\mathbf{x}_{i,j}$: get k representatives $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k \in \mathbb{R}^p$.
- ▶ Represent \mathbf{x}_i by a histogram over $\{1, 2, \dots, k\}$ based on assignments of \mathbf{x}_i 's parts to representatives.



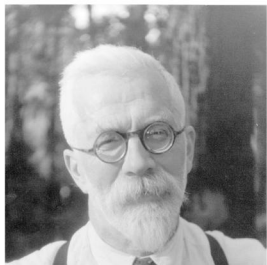
Uses of clustering: compression

Quantization

Replace each x_i with its representative

$$x_i \mapsto c_f(x_i).$$

Example: quantization at image patch level.



k -MEANS CLUSTERING

k -means clustering

Problem

- ▶ **Input:** $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{R}^d$, target cardinality $k \in \mathbb{N}$.
- ▶ **Output:** k representatives (“centers”, “means”) $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k \in \mathbb{R}^d$.
- ▶ **Objective:** choose $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k \in \mathbb{R}^d$ to minimize

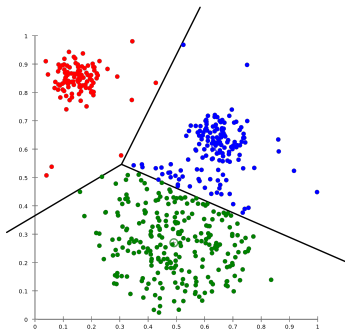
$$\sum_{i=1}^n \min_{j \in [k]} \|\mathbf{x}_i - \mathbf{c}_j\|_2^2.$$

k -means clustering

Problem

- ▶ **Input:** $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{R}^d$, target cardinality $k \in \mathbb{N}$.
- ▶ **Output:** k representatives (“centers”, “means”) $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k \in \mathbb{R}^d$.
- ▶ **Objective:** choose $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k \in \mathbb{R}^d$ to minimize

$$\sum_{i=1}^n \min_{j \in [k]} \|\mathbf{x}_i - \mathbf{c}_j\|_2^2.$$



Natural assignment function

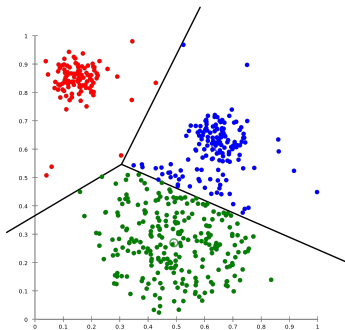
$$f(\mathbf{x}) := \arg \min_{j \in [k]} \|\mathbf{x} - \mathbf{c}_j\|_2^2.$$

k -means clustering

Problem

- ▶ **Input:** $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{R}^d$, target cardinality $k \in \mathbb{N}$.
- ▶ **Output:** k representatives (“centers”, “means”) $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k \in \mathbb{R}^d$.
- ▶ **Objective:** choose $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k \in \mathbb{R}^d$ to minimize

$$\sum_{i=1}^n \min_{j \in [k]} \|\mathbf{x}_i - \mathbf{c}_j\|_2^2.$$



Natural assignment function

$$f(\mathbf{x}) := \arg \min_{j \in [k]} \|\mathbf{x} - \mathbf{c}_j\|_2^2.$$

NP-hard, even if $k = 2$ or $d = 2$.

The easy cases

k -means clustering for $k = 1$

Problem: Pick $\mathbf{c} \in \mathbb{R}^d$ to minimize

$$\sum_{i=1}^n \|\mathbf{x}_i - \mathbf{c}\|_2^2.$$

The easy cases

k -means clustering for $k = 1$

Problem: Pick $\mathbf{c} \in \mathbb{R}^d$ to minimize

$$\sum_{i=1}^n \|\mathbf{x}_i - \mathbf{c}\|_2^2.$$

Solution: “bias/variance decomposition”

$$\frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{c}\|_2^2 = \|\boldsymbol{\mu} - \mathbf{c}\|_2^2 + \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - \boldsymbol{\mu}\|_2^2$$

where $\boldsymbol{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$.

The easy cases

k -means clustering for $k = 1$

Problem: Pick $\mathbf{c} \in \mathbb{R}^d$ to minimize

$$\sum_{i=1}^n \|\mathbf{x}_i - \mathbf{c}\|_2^2.$$

Solution: “bias/variance decomposition”

$$\frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{c}\|_2^2 = \|\boldsymbol{\mu} - \mathbf{c}\|_2^2 + \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - \boldsymbol{\mu}\|_2^2$$

where $\boldsymbol{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$.

Therefore, optimal choice for \mathbf{c} is $\boldsymbol{\mu}$.

The easy cases

k -means clustering for $k = 1$

Problem: Pick $\mathbf{c} \in \mathbb{R}^d$ to minimize

$$\sum_{i=1}^n \|\mathbf{x}_i - \mathbf{c}\|_2^2.$$

Solution: “bias/variance decomposition”

$$\frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{c}\|_2^2 = \|\boldsymbol{\mu} - \mathbf{c}\|_2^2 + \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - \boldsymbol{\mu}\|_2^2$$

where $\boldsymbol{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$.

Therefore, optimal choice for \mathbf{c} is $\boldsymbol{\mu}$.

k -means clustering for $d = 1$

Dynamic programming in time $O(n^2k)$.

Alternating optimization algorithm

Assignment variables

For each data point \mathbf{x}_i , let $\phi_i \in \{0, 1\}^k$ denote its “one-hot” representation:

$$\phi_{i,j} = \mathbb{1}\{\mathbf{x}_i \text{ is assigned to cluster } j\}.$$

Alternating optimization algorithm

Assignment variables

For each data point \mathbf{x}_i , let $\phi_i \in \{0, 1\}^k$ denote its “one-hot” representation:

$$\phi_{i,j} = \mathbb{1}\{\mathbf{x}_i \text{ is assigned to cluster } j\}.$$

Objective becomes (for optimal setting of ϕ_i s)

$$\sum_{i=1}^n \min_{j \in [k]} \|\mathbf{x}_i - \mathbf{c}_j\|_2^2 = \sum_{i=1}^n \left\{ \sum_{j=1}^k \phi_{i,j} \|\mathbf{x}_i - \mathbf{c}_j\|_2^2 \right\}.$$

Alternating optimization algorithm

Assignment variables

For each data point \mathbf{x}_i , let $\phi_i \in \{0, 1\}^k$ denote its “one-hot” representation:

$$\phi_{i,j} = \mathbb{1}\{\mathbf{x}_i \text{ is assigned to cluster } j\}.$$

Objective becomes (for optimal setting of ϕ_i s)

$$\sum_{i=1}^n \min_{j \in [k]} \|\mathbf{x}_i - \mathbf{c}_j\|_2^2 = \sum_{i=1}^n \left\{ \sum_{j=1}^k \phi_{i,j} \|\mathbf{x}_i - \mathbf{c}_j\|_2^2 \right\}.$$

Lloyd's algorithm (sometimes called *the* k -means algorithm)

Initialize $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k \in \mathbb{R}^d$ somehow. Then repeat until convergence:

- ▶ Holding $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k$ fixed, pick optimal $\phi_1, \phi_2, \dots, \phi_n$.
 - ▶ Holding $\phi_1, \phi_2, \dots, \phi_n$ fixed, pick optimal $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k$.
-

Alternating optimization algorithm

Assignment variables

For each data point \mathbf{x}_i , let $\phi_i \in \{0, 1\}^k$ denote its “one-hot” representation:

$$\phi_{i,j} = \mathbb{1}\{\mathbf{x}_i \text{ is assigned to cluster } j\}.$$

Objective becomes (for optimal setting of ϕ_i s)

$$\sum_{i=1}^n \min_{j \in [k]} \|\mathbf{x}_i - \mathbf{c}_j\|_2^2 = \sum_{i=1}^n \left\{ \sum_{j=1}^k \phi_{i,j} \|\mathbf{x}_i - \mathbf{c}_j\|_2^2 \right\}.$$

Lloyd's algorithm (sometimes called *the* k -means algorithm)

Initialize $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k \in \mathbb{R}^d$ somehow. Then repeat until convergence:

- Holding $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k$ fixed, pick optimal $\phi_1, \phi_2, \dots, \phi_n$.

Set ϕ_i so \mathbf{x}_i is assigned to closest \mathbf{c}_j .

- Holding $\phi_1, \phi_2, \dots, \phi_n$ fixed, pick optimal $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k$.
-

Alternating optimization algorithm

Assignment variables

For each data point \mathbf{x}_i , let $\phi_i \in \{0, 1\}^k$ denote its “one-hot” representation:

$$\phi_{i,j} = \mathbb{1}\{\mathbf{x}_i \text{ is assigned to cluster } j\}.$$

Objective becomes (for optimal setting of ϕ_i s)

$$\sum_{i=1}^n \min_{j \in [k]} \|\mathbf{x}_i - \mathbf{c}_j\|_2^2 = \sum_{i=1}^n \left\{ \sum_{j=1}^k \phi_{i,j} \|\mathbf{x}_i - \mathbf{c}_j\|_2^2 \right\}.$$

Lloyd's algorithm (sometimes called *the* k -means algorithm)

Initialize $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k \in \mathbb{R}^d$ somehow. Then repeat until convergence:

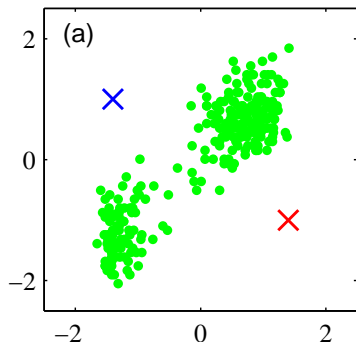
- ▶ Holding $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k$ fixed, pick optimal $\phi_1, \phi_2, \dots, \phi_n$.

Set ϕ_i so \mathbf{x}_i is assigned to closest \mathbf{c}_j .

- ▶ Holding $\phi_1, \phi_2, \dots, \phi_n$ fixed, pick optimal $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k$.

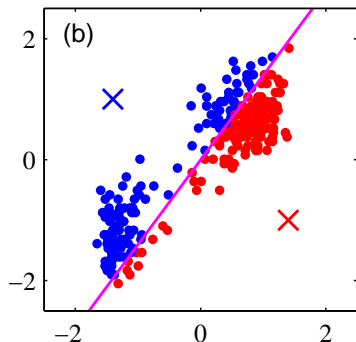
Set \mathbf{c}_j to be the average of the \mathbf{x}_i assigned to cluster j .

Sample run of Lloyd's algorithm



Arbitrary initialization of c_1 and c_2 .

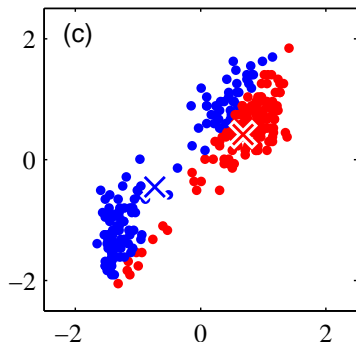
Sample run of Lloyd's algorithm



Iteration 1

Optimize assignments ϕ_i .

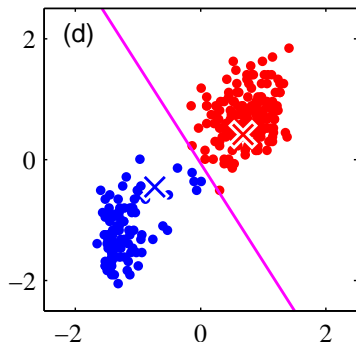
Sample run of Lloyd's algorithm



Iteration 1

Optimize representatives c_j .

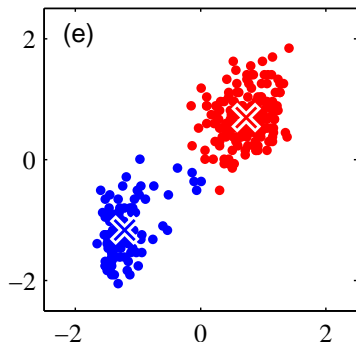
Sample run of Lloyd's algorithm



Iteration 2

Optimize assignments ϕ_i .

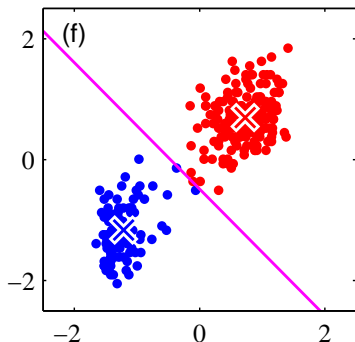
Sample run of Lloyd's algorithm



Iteration 2

Optimize representatives c_j .

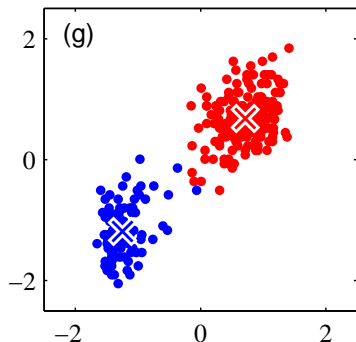
Sample run of Lloyd's algorithm



Iteration 3

Optimize assignments ϕ_i .

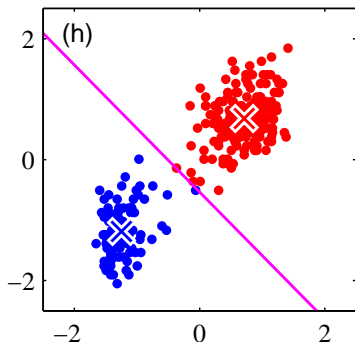
Sample run of Lloyd's algorithm



Iteration 3

Optimize representatives c_j .

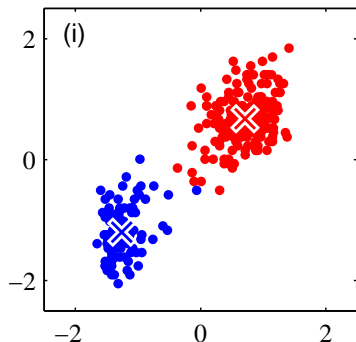
Sample run of Lloyd's algorithm



Iteration 4

Optimize assignments ϕ_i .

Sample run of Lloyd's algorithm



Iteration 4

Optimize representatives c_j .

Initializing Lloyd's algorithm

Basic idea: Choose initial centers to have good coverage of the data points.

Farthest-first traversal

For $j = 1, 2, \dots, k$:

- Pick $\mathbf{c}_j \in \mathbb{R}^d$ from among $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ farthest from previously chosen $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_{j-1}$.
(\mathbf{c}_1 chosen arbitrarily.)

Initializing Lloyd's algorithm

Basic idea: Choose initial centers to have good coverage of the data points.

Farthest-first traversal

For $j = 1, 2, \dots, k$:

- Pick $\mathbf{c}_j \in \mathbb{R}^d$ from among $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ farthest from previously chosen $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_{j-1}$.
(\mathbf{c}_1 chosen arbitrarily.)

But this can be thrown off by outliers...

Initializing Lloyd's algorithm

Basic idea: Choose initial centers to have good coverage of the data points.

Farthest-first traversal

For $j = 1, 2, \dots, k$:

- ▶ Pick $\mathbf{c}_j \in \mathbb{R}^d$ from among $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ farthest from previously chosen $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_{j-1}$.
(\mathbf{c}_1 chosen arbitrarily.)

But this can be thrown off by outliers...

A better idea:

D^2 sampling (a.k.a. “ k -means++”)

For $j = 1, 2, \dots, k$:

- ▶ Randomly pick $\mathbf{c}_j \in \mathbb{R}^d$ from among $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ according to distribution

$$P(\mathbf{x}_i) \propto \min_{\ell=1,2,\dots,j-1} \|\mathbf{x}_i - \mathbf{c}_\ell\|_2^2.$$

(Uniform distribution when $j = 1$.)

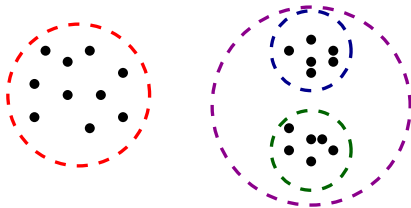
Choosing k

- ▶ Usually by hold-out validation / cross-validation on auxiliary task (e.g., supervised learning task).

Choosing k

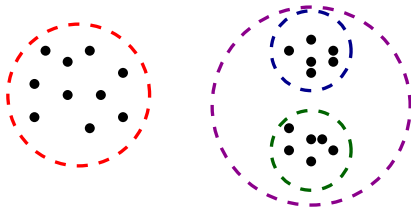
- ▶ Usually by hold-out validation / cross-validation on auxiliary task (e.g., supervised learning task).
- ▶ *Heuristic*: Find large gap between $(k - 1)$ -means cost and k -means cost.

Clustering at multiple scales



$k = 2$ or $k = 3$?

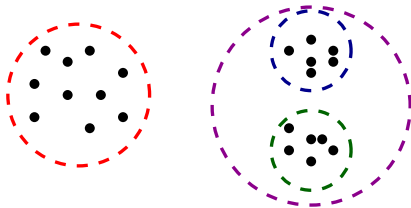
Clustering at multiple scales



$k = 2$ or $k = 3$?

Hierarchical clustering: encode clusterings for all values of k in a tree.

Clustering at multiple scales



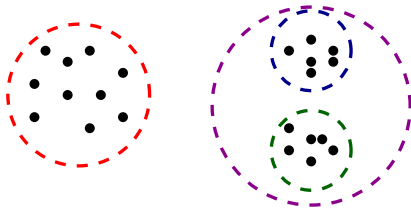
$k = 2$ or $k = 3$?

Hierarchical clustering: encode clusterings for all values of k in a tree.

Caveat: not always possible.



Clustering at multiple scales



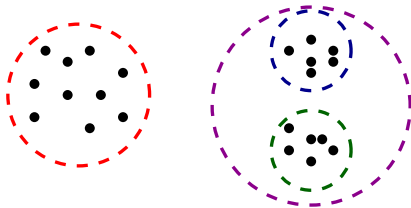
$k = 2$ or $k = 3$?

Hierarchical clustering: encode clusterings for all values of k in a tree.

Caveat: not always possible.



Clustering at multiple scales



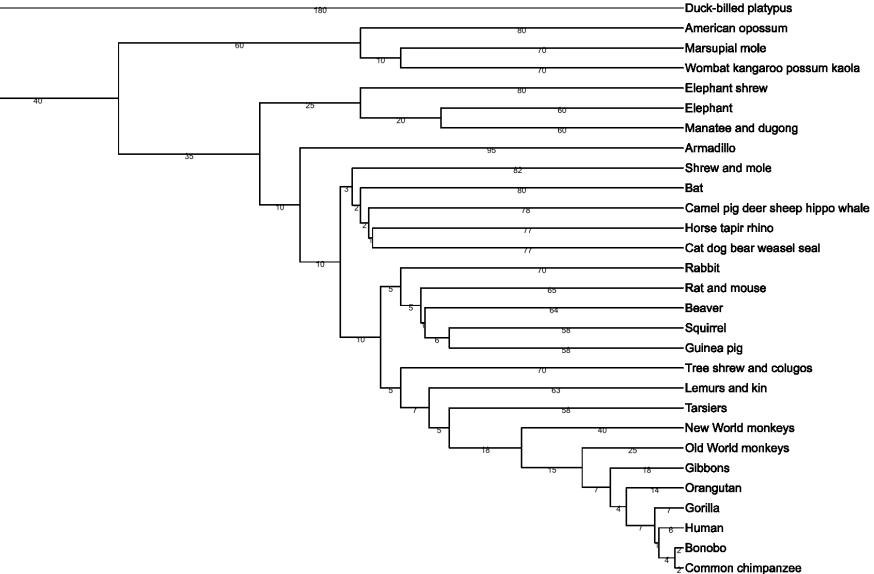
$k = 2$ or $k = 3$?

Hierarchical clustering: encode clusterings for all values of k in a tree.

Caveat: not always possible.



Example: phylogenetic tree



Hierarchical clustering

Divisive (top-down) clustering

- ▶ Partition data into two groups (e.g., via k -means clustering with $k = 2$).
- ▶ Recurse on each part.

Hierarchical clustering

Divisive (top-down) clustering

- ▶ Partition data into two groups (e.g., via k -means clustering with $k = 2$).
- ▶ Recurse on each part.

Agglomerative (bottom-up) clustering

- ▶ Start with every point x_i in its own cluster.
 - ▶ Repeatedly merge “closest” pair of clusters.
-

Hierarchical clustering

Divisive (top-down) clustering

- ▶ Partition data into two groups (e.g., via k -means clustering with $k = 2$).
- ▶ Recurse on each part.

Agglomerative (bottom-up) clustering

- ▶ Start with every point x_i in its own cluster.
- ▶ Repeatedly merge “closest” pair of clusters.

Example: *Ward's average linkage method*

$$\text{dist}(C, \tilde{C}) := \frac{|C| \cdot |\tilde{C}|}{|C| + |\tilde{C}|} \|\text{mean}(C) - \text{mean}(\tilde{C})\|_2^2$$

(the increase in k -means cost caused by merging C and \tilde{C}).

Recap

- ▶ Uses of clustering:
 - ▶ Unsupervised classification (“hidden subpopulations”).
 - ▶ Quantization
 - ▶ ...
- ▶ k -means clustering: popular objective for clustering and quantization.
- ▶ Lloyd's algorithm: alternating optimization, needs good initialization.
- ▶ Hierarchical clustering: clustering at multiple levels of granularity.