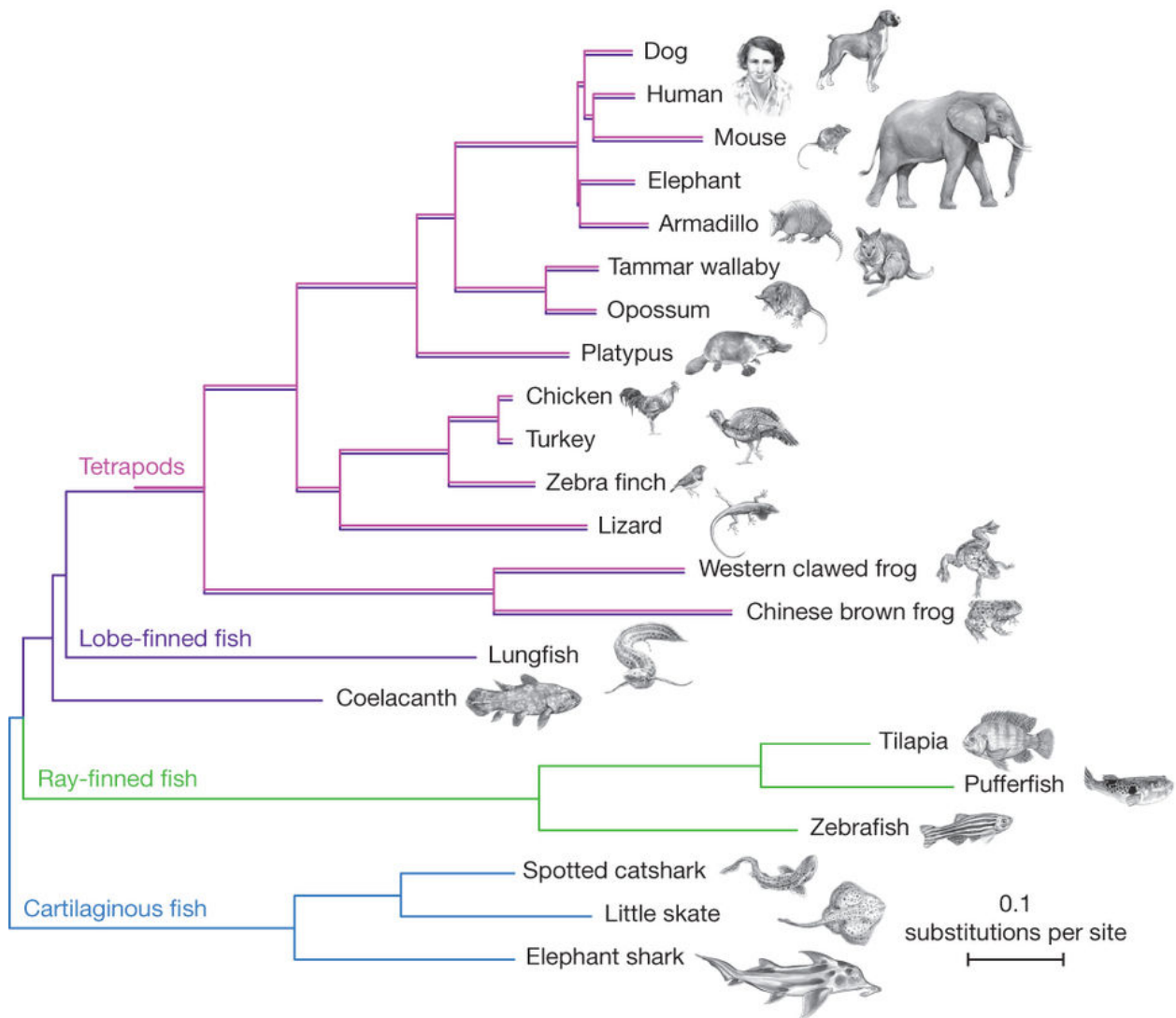


1 Hierarchical Clustering for Phylogenetic Trees

A phylogenetic tree (or “evolutionary tree”) is a way of representing the branching nature of evolution. Early branches represent major divergences in evolution (for example, modern vertebrates diverging from modern invertebrates), while later branches represent smaller branches in evolution (for example, modern humans diverging from modern monkeys). An example is shown below.



Creating phylogenetic trees is a popular problem in computational biology. We are going to combine what we know about clustering, decision trees, and unsupervised learning.

We start with all the samples (in this case, animals) in a single cluster and gradually divide this up, until each animal is in its own cluster. This should remind you of decision trees! After k steps, we have at most 2^k clusters. Since we do not have labels, we need to find some way of deciding how to split the samples (other than using entropy).

We will use the same objective as in k -means clustering to determine how good our proposed clustering is (here, $|x|$ denotes length and $\|x\|$ denotes norm):

$$\min L = \sum_{i=1}^k \sum_{X_j \in S_i} \|X_j - \mu_i\|_2^2, \text{ where}$$

$$\mu_i = \frac{1}{|S_i|} \sum_{X_j \in S_i} X_j, \quad i = 1, 2, \dots, k.$$

At each iteration, we will split each cluster with more than one element into two clusters, choosing the split that achieves the minimum resulting objective value L . The algorithm terminates when every sample point is in its own cluster, yielding an objective value $L = 0$.

(a) Consider the following six animals and their two features. Create the resulting decision tree.

Animal	Lifespan	Wings
Dog	12	0
Human	80	0
Mouse	2	0
Elephant	60	0
Chicken	8	2
Turkey	10	2

(b) Prove that an optimal clustering on $k + 1 < n$ clusters has an objective value that is at least as small as that of the optimal clustering on k clusters.

2 Kernel k-means

Suppose we have a dataset $\{x_i\}_{i=1}^N, x_i \in \mathbb{R}^n$ that we want to split into K clusters. Furthermore, suppose we know a priori that this data is best clustered in a large feature space \mathbb{R}^m , and that we have a feature map $\phi : \mathbb{R}^n \rightarrow \mathbb{R}^m$. How should we perform clustering in this space?

- (a) Write the objective for K-means clustering in the feature space (using the squared L_2 norm in the feature space). Do so by explicitly constructing cluster centers $\{\mu_k\}_{k=1}^K$ with all $\mu_k \in \mathbb{R}^m$.
- (b) Write an algorithm that minimizes the objective in (a).
- (c) Write an algorithm that minimizes the objective in (a) without explicitly constructing the cluster centers $\{\mu_k\}$. Assume you are given a kernel function $\kappa(x, y) = \phi(x) \cdot \phi(y)$.

3 Regularization and Kernel k-Means

Recall that in k -means clustering we attempt to minimize the objective

$$\min_{C_1, C_2, \dots, C_k} L = \sum_{i=1}^k \sum_{X_j \in C_i} \|X_j - \mu_i\|_2^2, \text{ where}$$

$$\mu_i = \operatorname{argmin}_{\mu_i \in \mathbb{R}^d} \sum_{X_j \in C_i} \|X_j - \mu_i\|_2^2 = \frac{1}{|C_i|} \sum_{X_j \in C_i} X_j, \quad i = 1, 2, \dots, k.$$

The sample points are $\{X_1, \dots, X_n\}$, where $X_j \in \mathbb{R}^d$. C_i is the set of sample points assigned to cluster i and $|C_i|$ is the cluster's cardinality. Each sample point is assigned to exactly one cluster.

- (a) What is the minimum value of the objective when $k = n$ (the number of clusters equals the number of sample points)?
- (b) (Regularized k -means) Suppose we add a regularization term to the above objective. The regularized k -means objective is now

$$\min_{\mu_i \in \mathbb{R}^d} \sum_{i=1}^k \left(\lambda \|\mu_i\|_2^2 + \sum_{X_j \in C_i} \|X_j - \mu_i\|_2^2 \right).$$

Show that the optimum of

$$\min_{\mu_i \in \mathbb{R}^d} \lambda \|\mu_i\|_2^2 + \sum_{X_j \in C_i} \|X_j - \mu_i\|_2^2$$

is obtained at

$$\mu_i = \frac{1}{|C_i| + \lambda} \sum_{X_j \in C_i} X_j.$$

- (c) (Kernel k -means) Suppose we have a dataset $\{X_i\}_{i=1}^n$, $X_i \in \mathbb{R}^\ell$ that we want to split into k clusters, i.e., finding the best k -means clustering (without regularization). Furthermore, suppose we know *a priori* that this data is best clustered in an impractically high-dimensional feature space \mathbb{R}^m with an appropriate metric. Fortunately, instead of having to deal with the (implicit) feature map $\Phi : \mathbb{R}^\ell \rightarrow \mathbb{R}^m$ and (implicit) distance metric¹, we have a kernel function $\kappa(X_1, X_2) = \Phi(X_1) \cdot \Phi(X_2)$ that we can compute quickly on the raw samples. How should we perform the kernelized counterpart of k -means clustering?

Derive the missing portion of this algorithm, and show your work in deriving it. The primary issue to consider is that although we define the means μ_i in the usual way, we can't ever compute Φ explicitly because it's way too big. Therefore, in the step where we determine which cluster each sample point is assigned to, we must use the kernel function κ to obtain the right result. Review the lecture on kernels if you don't remember how that's done.

¹Just as how the interpretation of kernels in kernelized ridge regression involves an implicit prior/regularizer as well as an implicit feature space, we can think of kernels as generally inducing an implicit distance metric as well. Think of how you would represent the squared distance between two points in terms of pairwise inner products and operations on them.

Algorithm 1: Kernel k -means

Require: Data matrix $X \in \mathbb{R}^{n \times d}$; number of clusters K ; kernel function $\kappa(X_1, X_2)$

Ensure: Cluster $\text{class}(j)$ assigned for each sample point X_j .

function KERNEL-K-MEANS(X, K)

Randomly initialize $\text{class}(j)$ to be an integer in $1, 2, \dots, K$ for each X_j .

while *not converged* **do**

for $i \leftarrow 1$ **to** K **do**

$S_i \leftarrow \{j \in \{1, 2, \dots, n\} : \text{class}(j) = i\}$.

for $j \leftarrow 1$ **to** n **do**

 Set $\text{class}(j) \leftarrow \arg \min_k$ (Action item: derive the missing portion here)

 Return S_i for $i \leftarrow 1, 2, \dots, K$.

end function

- (d) The expression you derived may have unnecessary terms or redundant kernel computations, especially when you compute $\text{class}(j)$ for every $j \in [1, n]$. Explain how to eliminate them; that is, how to perform the computation quickly without doing irrelevant computations or redoing computations already done.