

A Randomized Linear-Time Algorithm to Find Minimum Spanning Trees

David R. Karger, Stanford

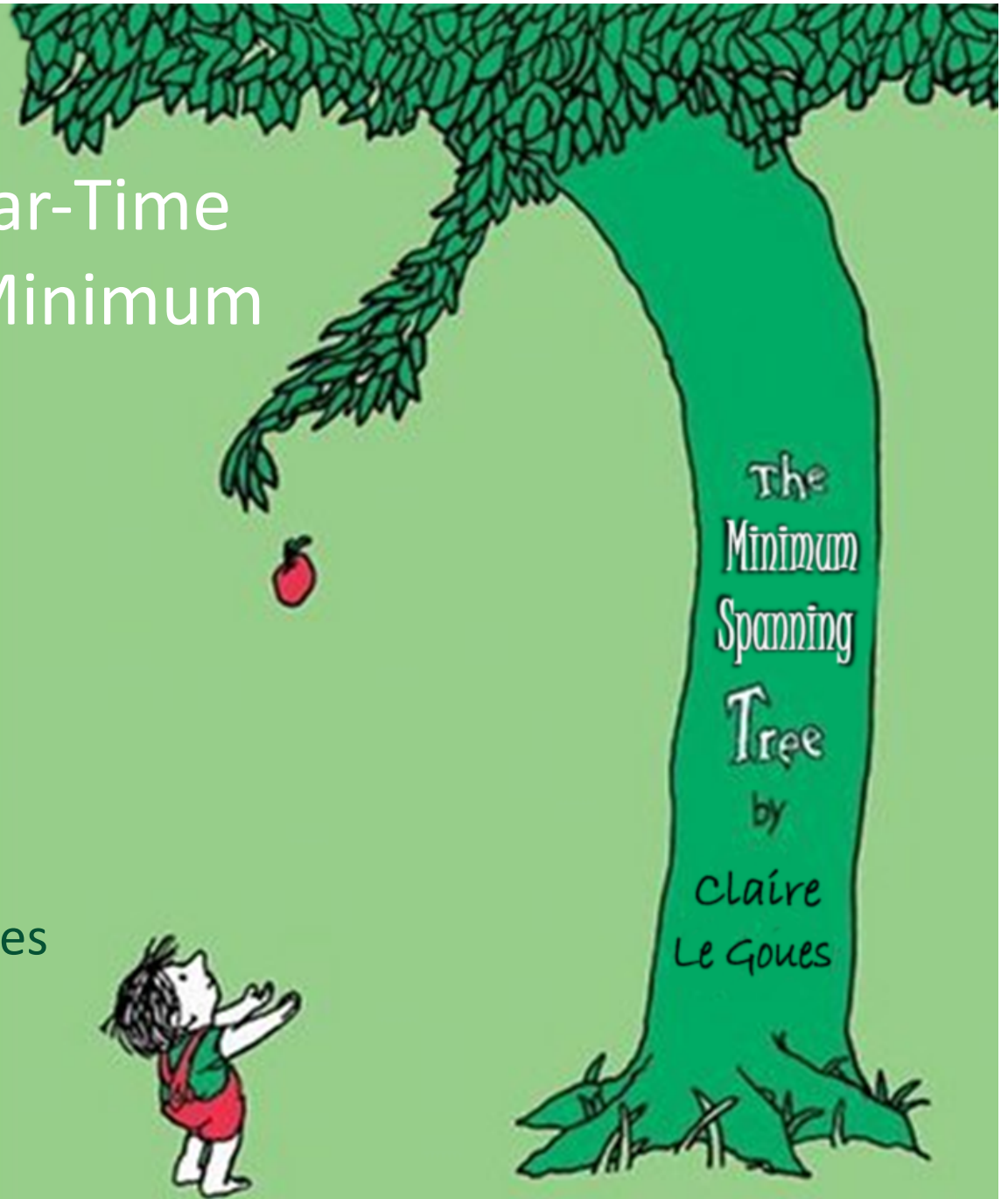
Philip N. Klein, Brown

Robert E. Tarjan, Princeton

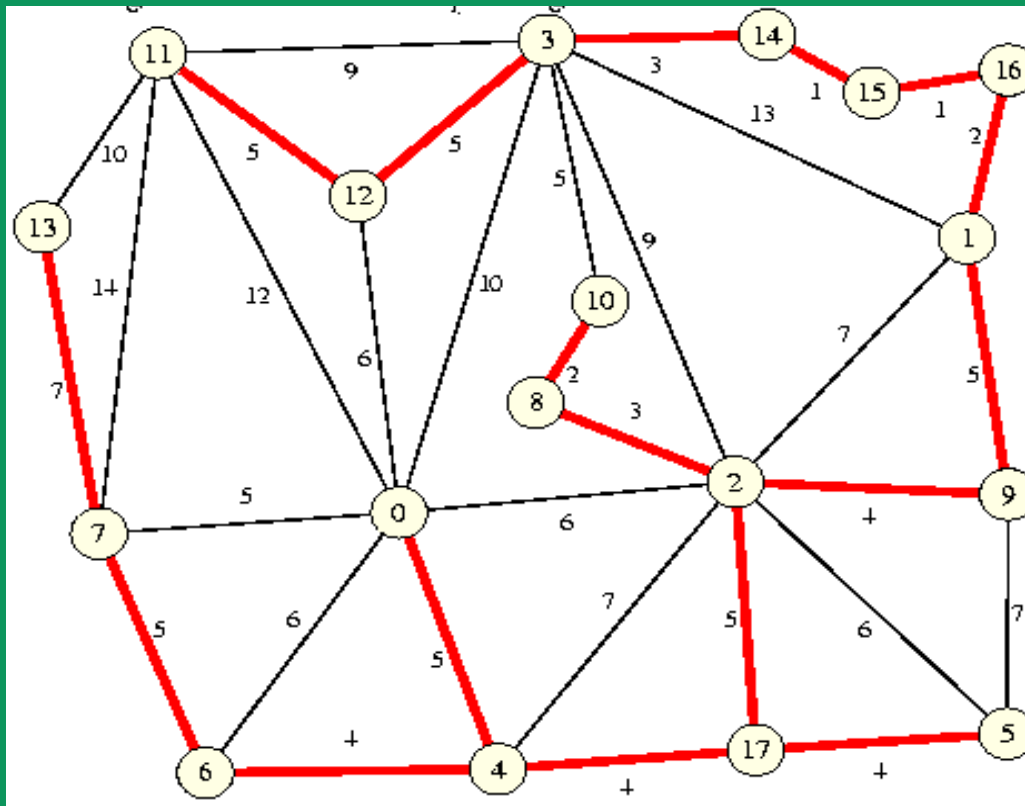
Presented By: Claire Le Goues

Theory lunch

August 14, 2008



Minimum Spanning Tree (MST)



(the weight of this tree is 65)

- **Definition:** a subset of edges in an undirected, weighted graph such that all the vertices are touched by at least one of the edges and the combined weight of all the edges is minimal



MST: Why?

- Cable/phone line layouts
- Airplane route scheduling
- Electricity grids
- ...and so on.
- Doing it quickly and efficiently is important!



Previous (Deterministic) Approaches

- Greedy algorithms, polynomial time:
 - Baruvka (1926)
 - Prim's and Kruskal's algorithms.
- Better algorithms, $O(n \log n)$:
 - Chazelle
 - Gabow et al.
 - etc

Can we do better?



SURE!

(probabilistically)

Introducing:
RANDOMNESS

A Sorting Digression

- Quicksort (vs. Mergesort, for example)
- Pick a random pivot point and recurse on parts.
- If we pick the point properly, it's much better!
- But we need to know how/what to pick, and how to combine the subparts
- KKT's algorithm is analogous.

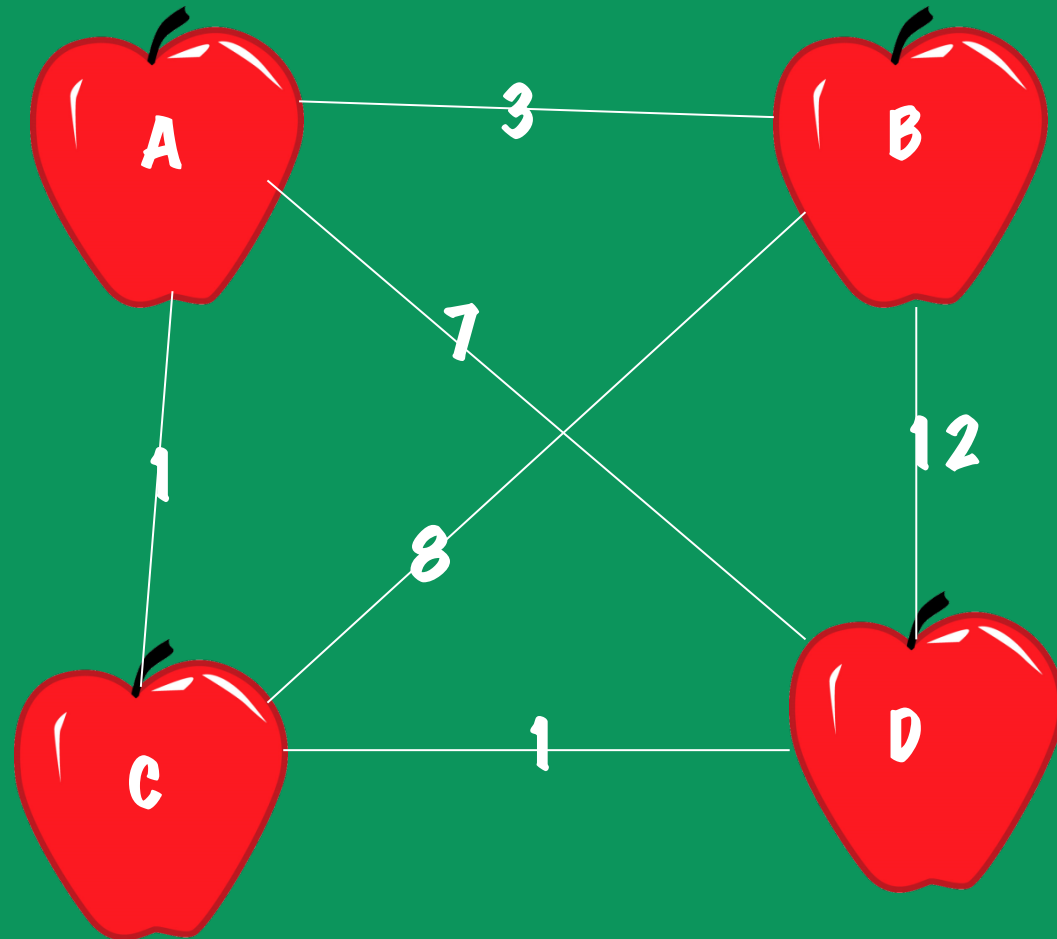
Preliminaries: Cut and Cycle Properties

- **Cut property:** For any proper nonempty subset X of the vertices, the lightest edge with exactly one endpoint in X belongs to the minimum spanning tree
- **Cycle property:** For any cycle C in a graph, the heaviest edge in C does not appear in the minimum spanning tree

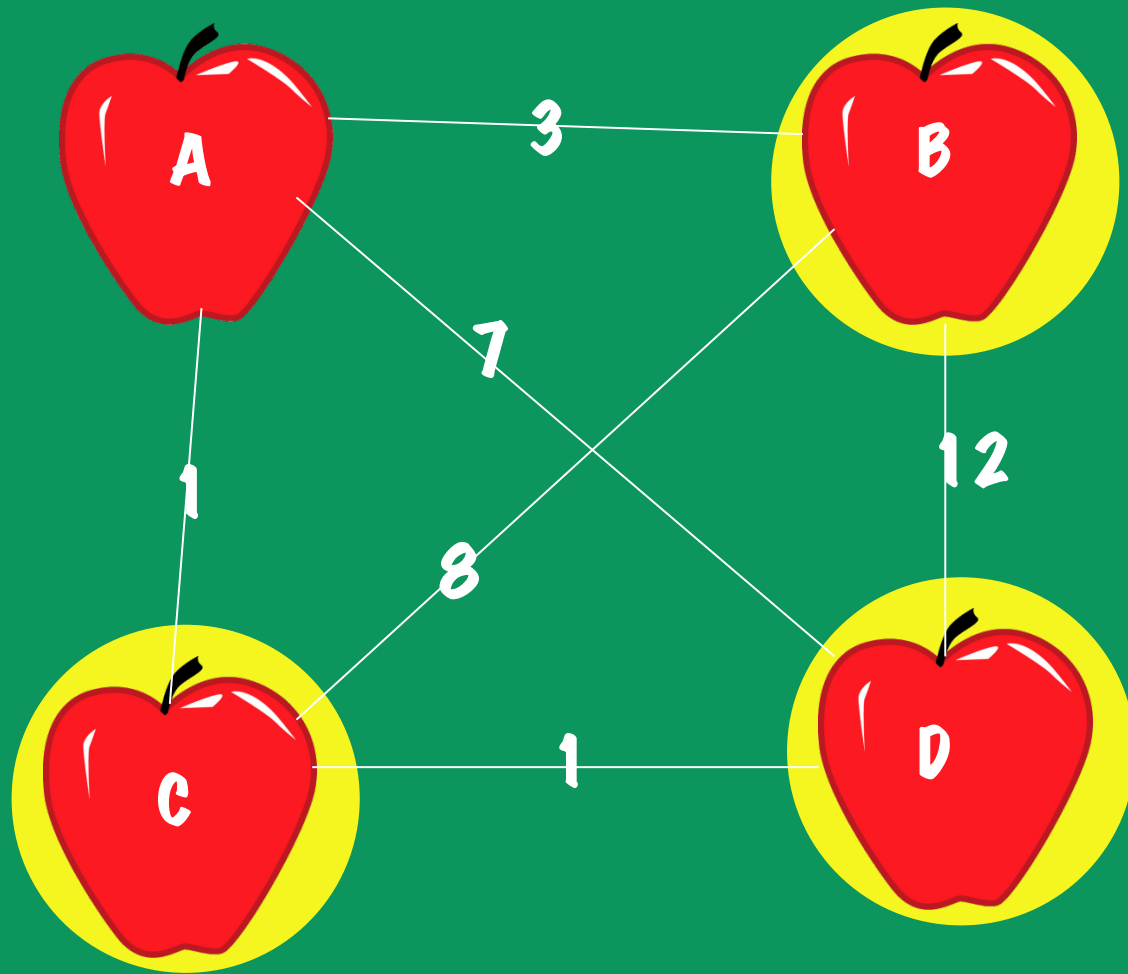
More preliminaries

- $w(x,y)$ is the weight of the edge $\{x, y\}$
- Where F is a forest in G , $F(x,y)$ is the path connecting x and y in F , and $w_F(x,y)$ is the maximum weight of an edge on $F(x,y)$
- An edge (x,y) is **F-heavy** if $w(x,y) > w_F(x,y)$ and **F-light** otherwise.
- *F-heavy edges don't belong in an MST!*

F-Heavy Edges

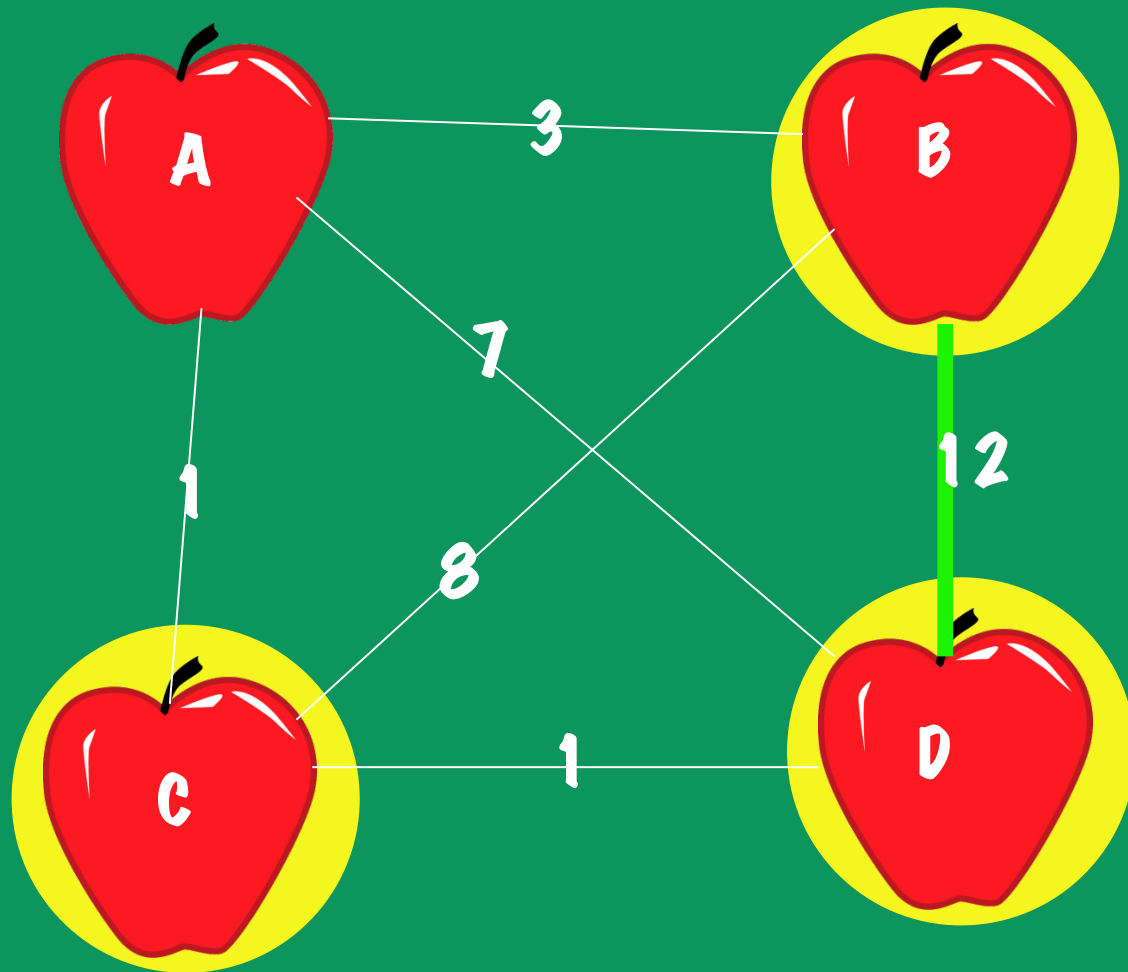


F-Heavy Edges

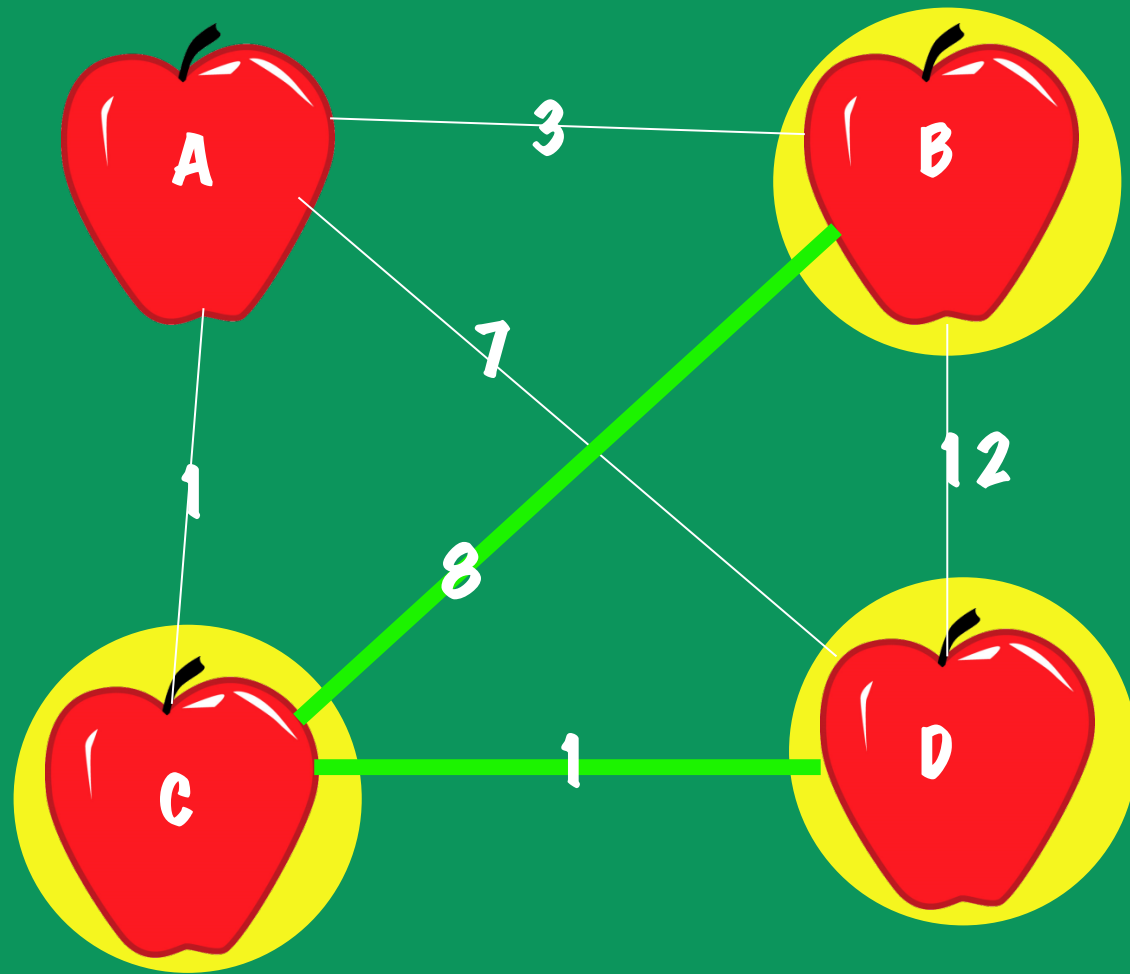


F-Heavy Edges

$$w(B,D) = 12$$

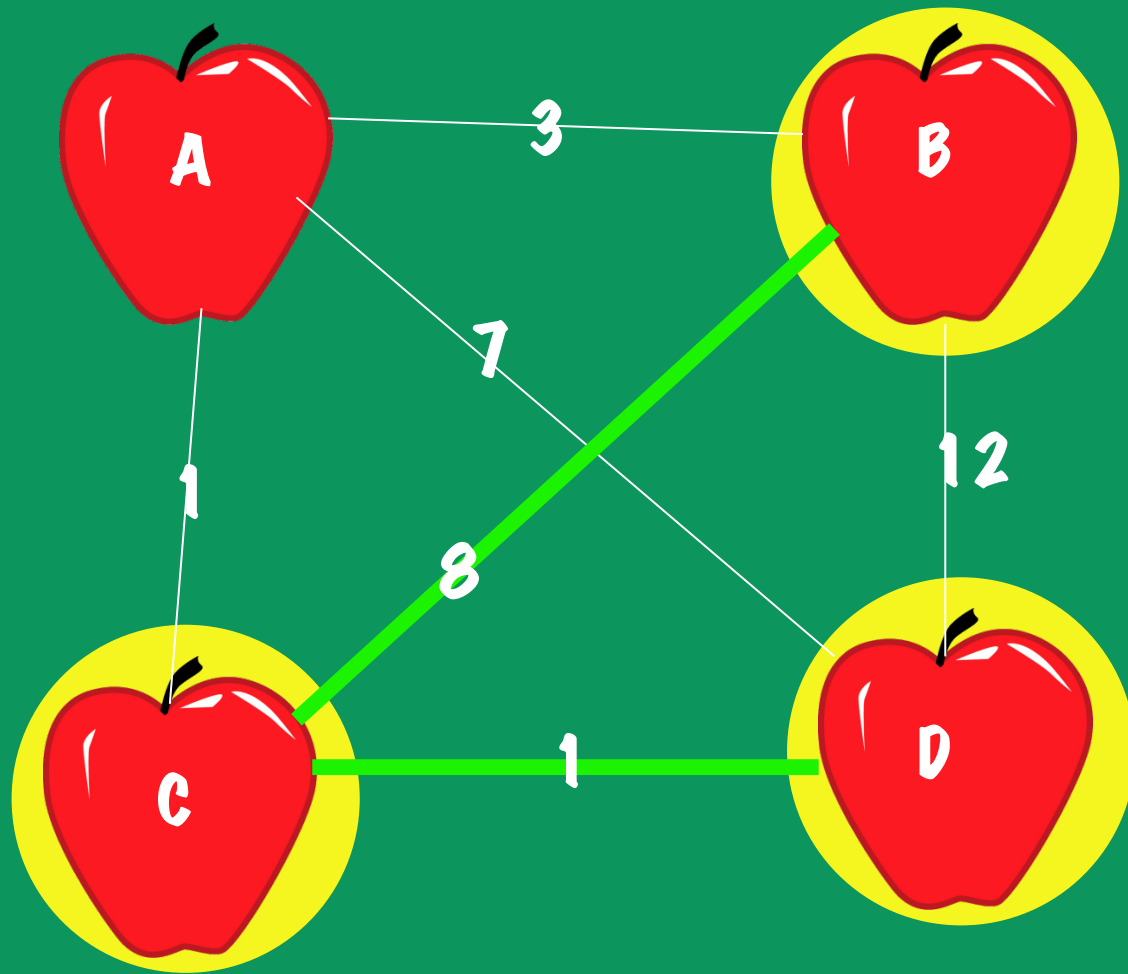


F-Heavy Edges



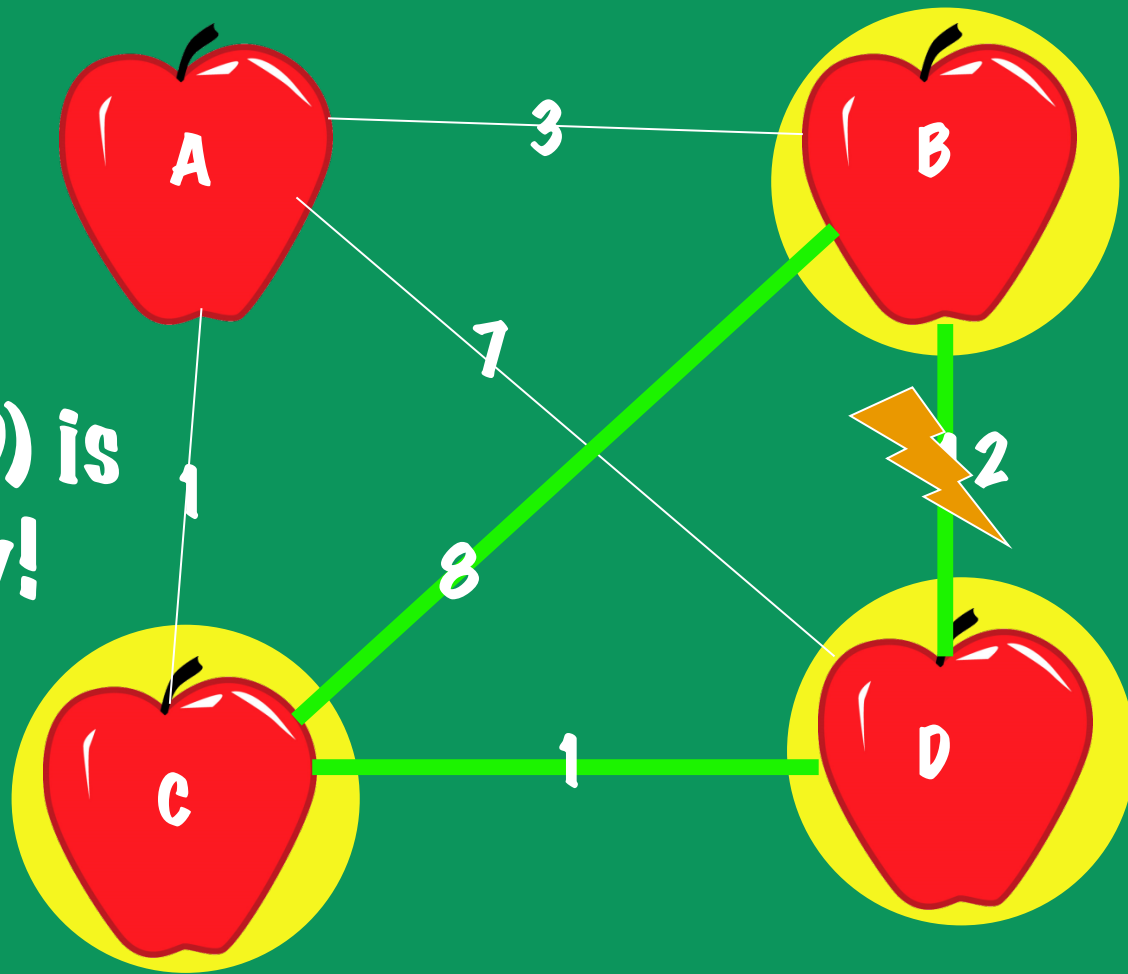
F-Heavy Edges

$$W_f(B,D) = 8$$

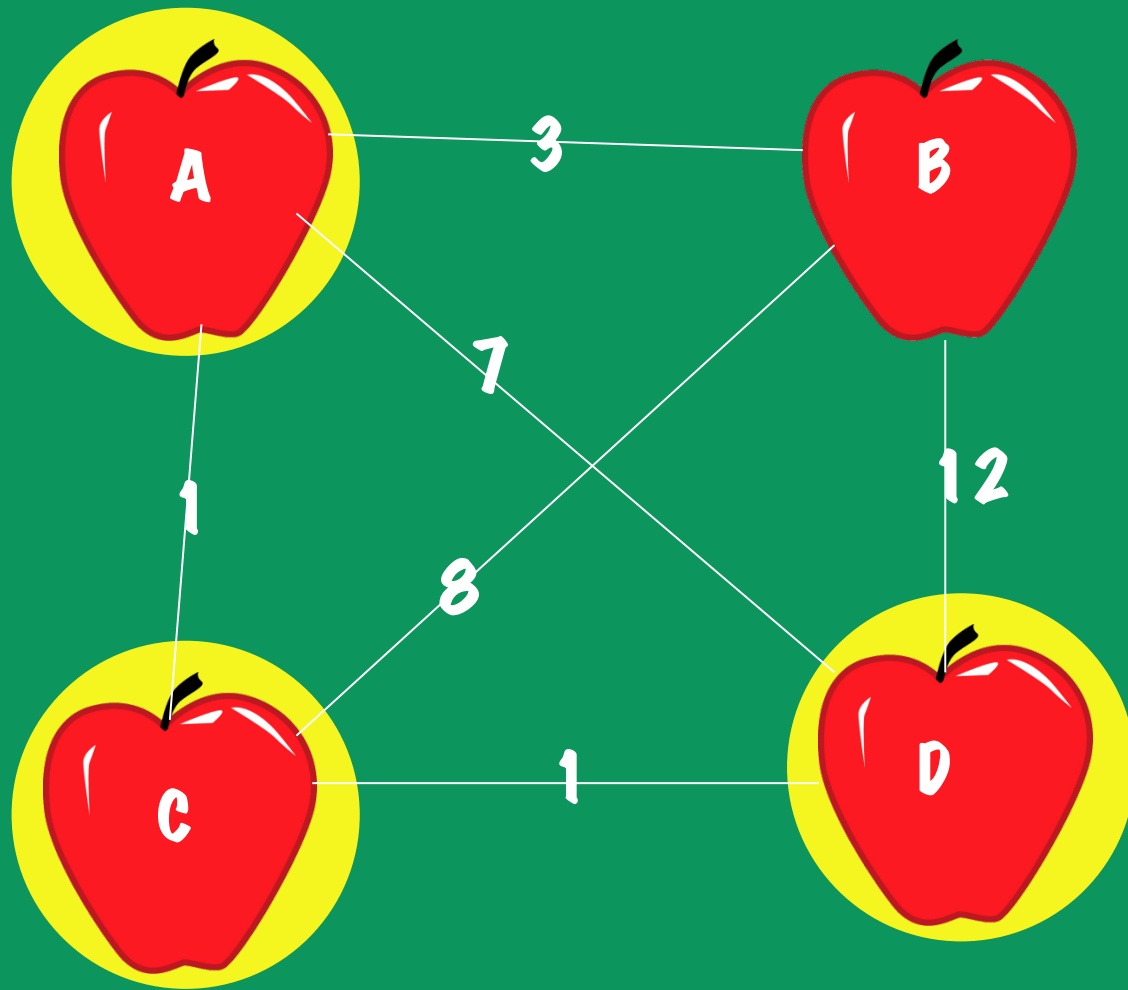


F-Heavy Edges

$12 > 8$,
So $w(B,D)$ is
F-Heavy!

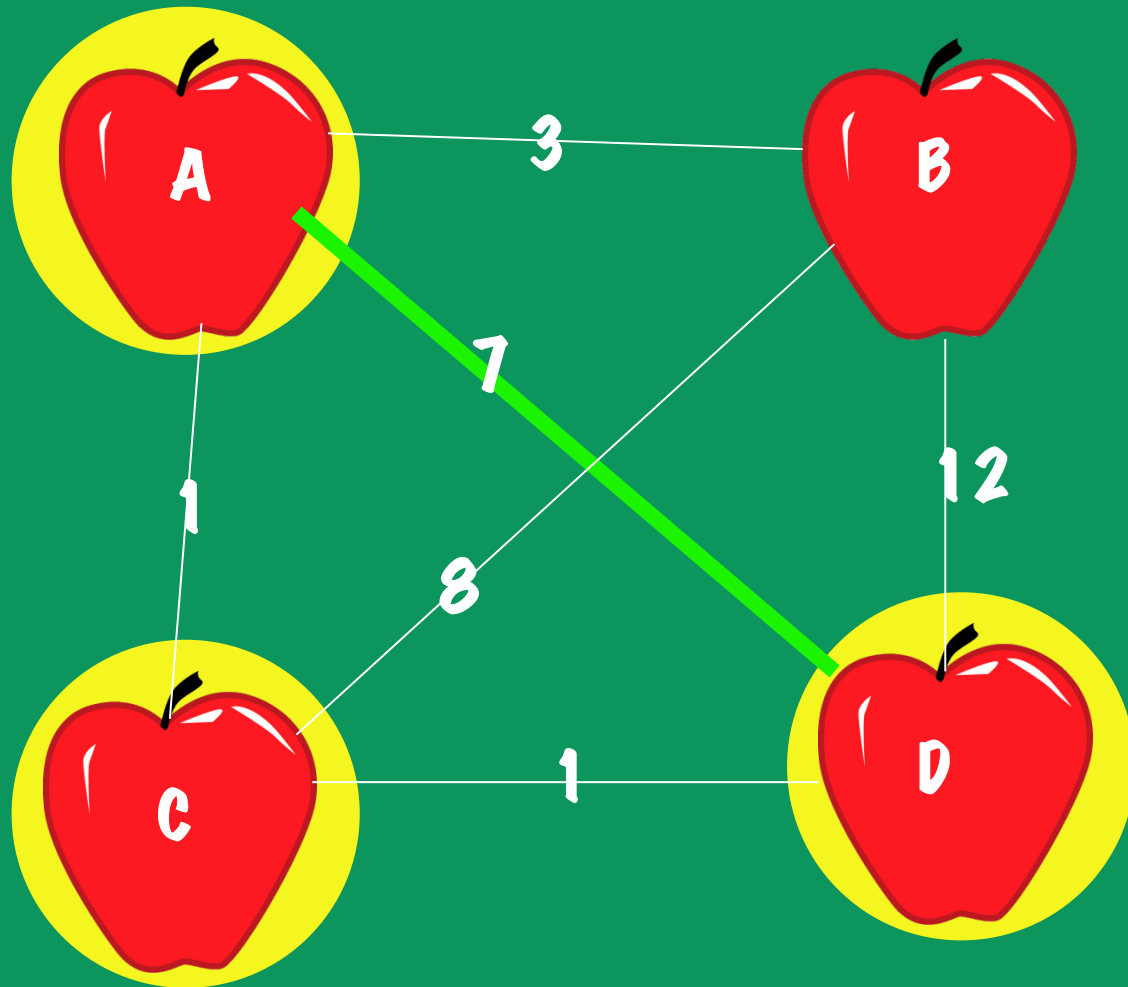


F-Heavy Edges



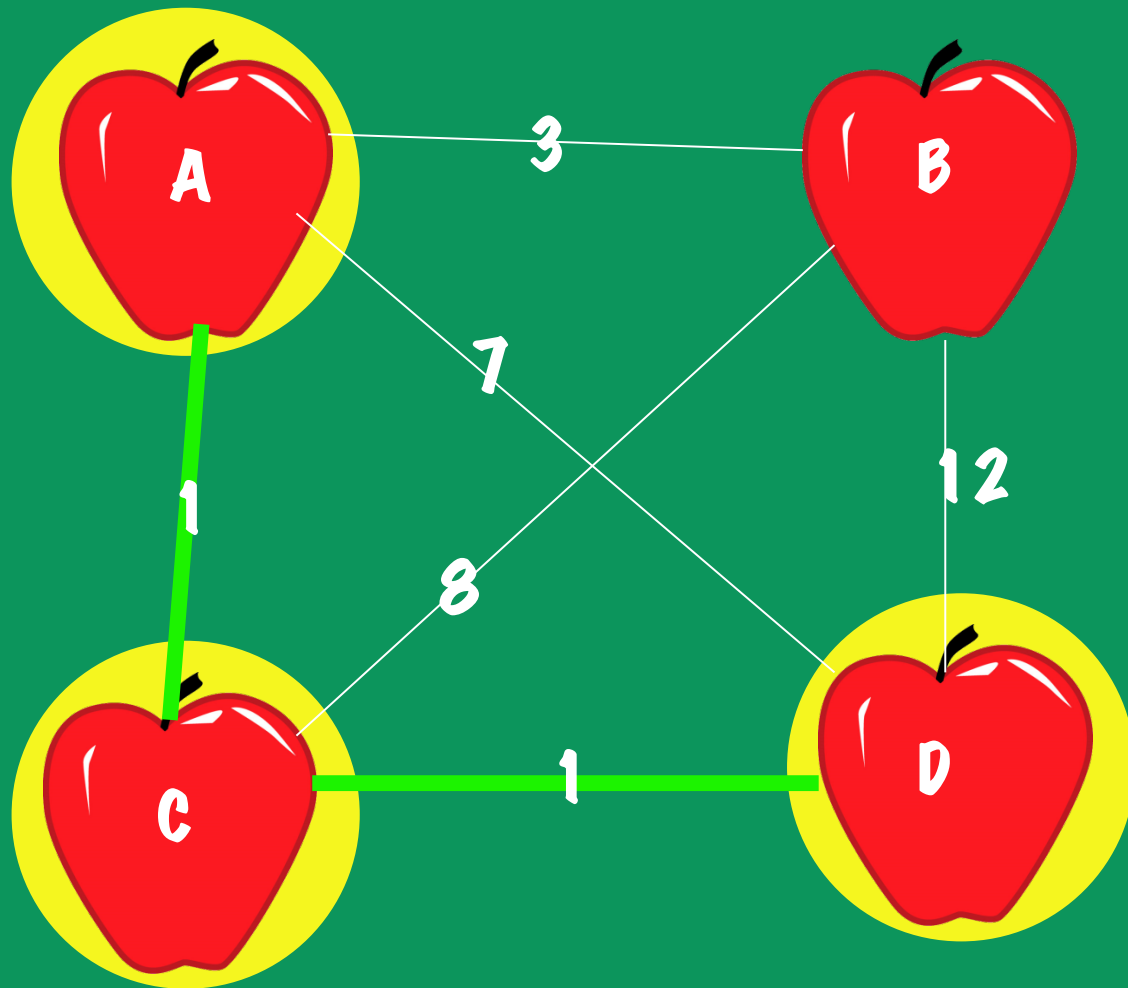
F-Heavy Edges

$$w(A,D) = 7$$



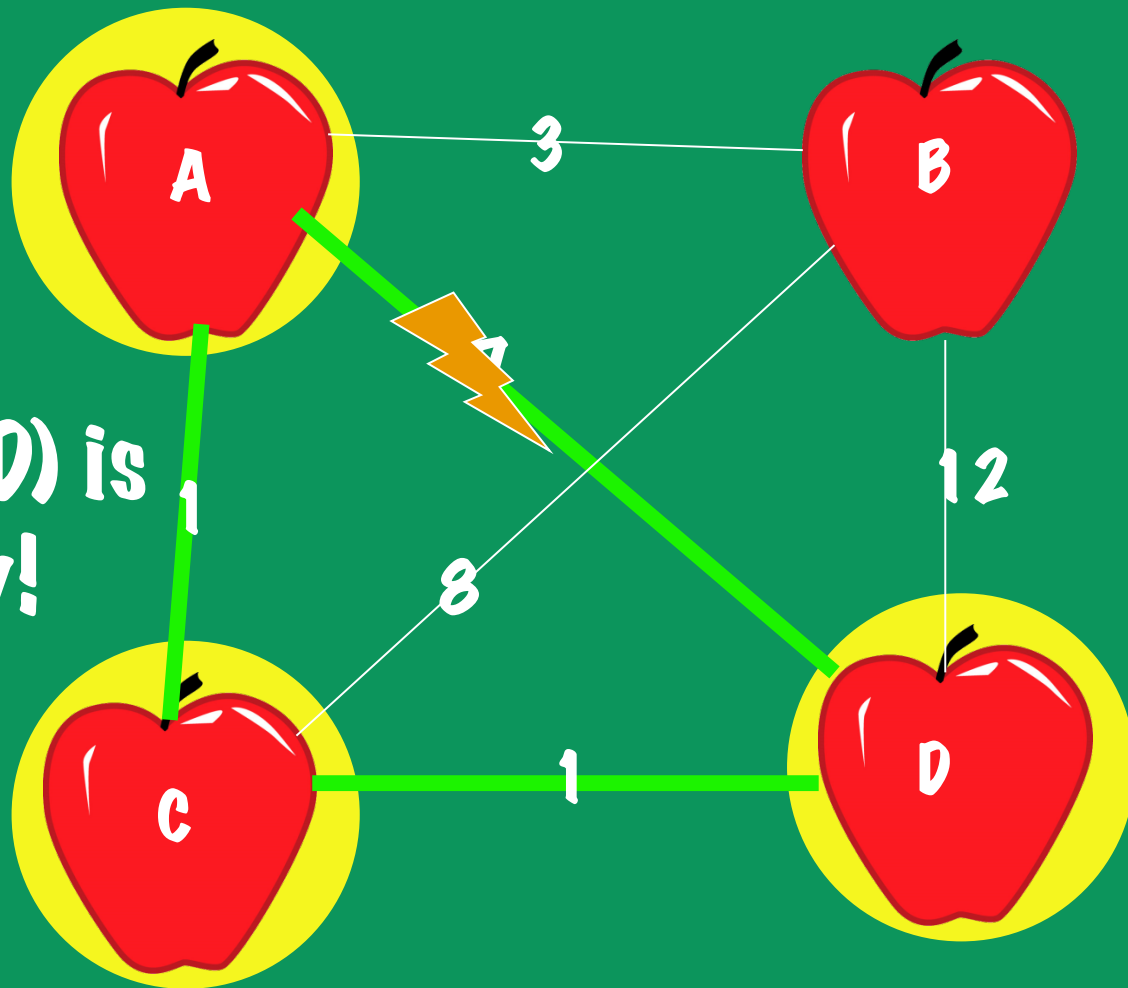
F-Heavy Edges

$$w_f(A,D) = 1$$

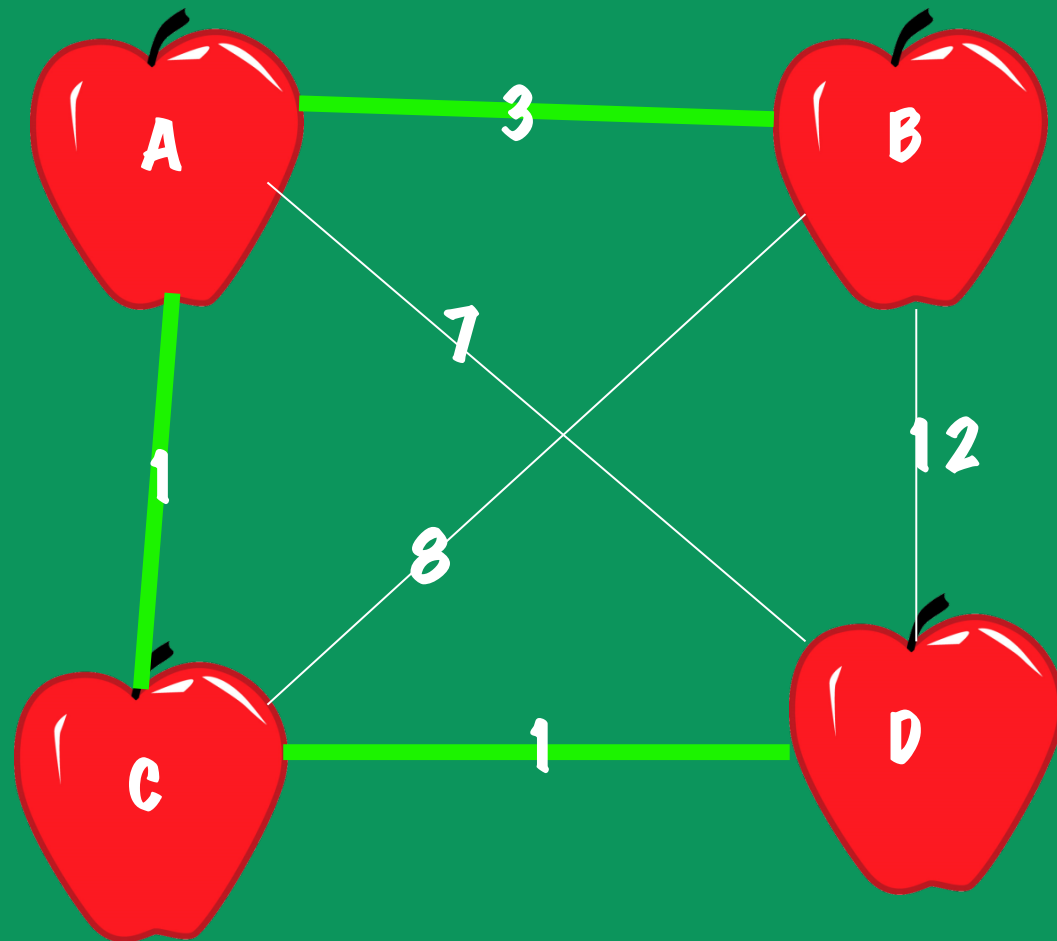


F-Heavy Edges

$7 > 1$,
So $w(A,D)$ is
F-Heavy!



F-Heavy Edges - The Final MST

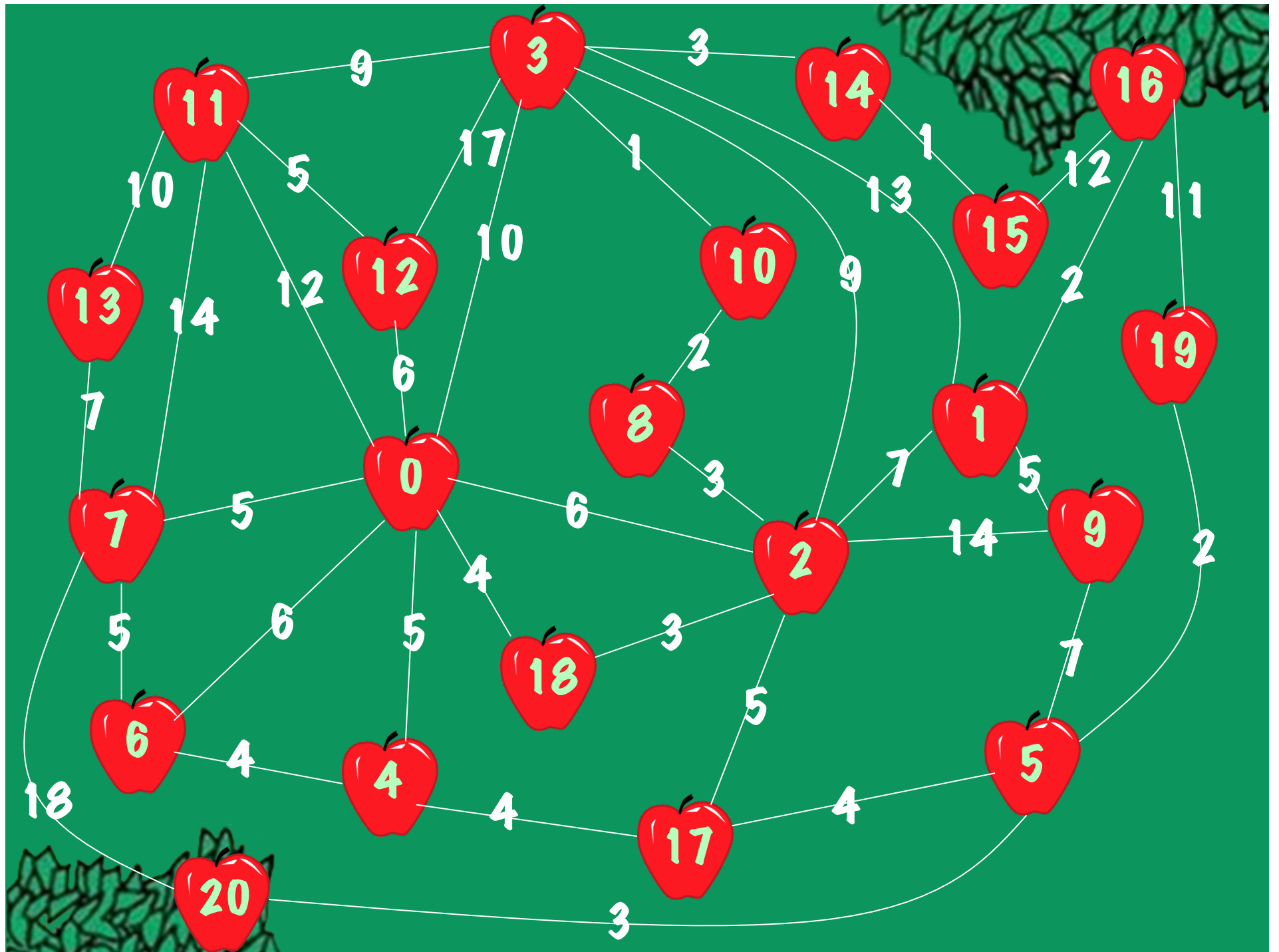


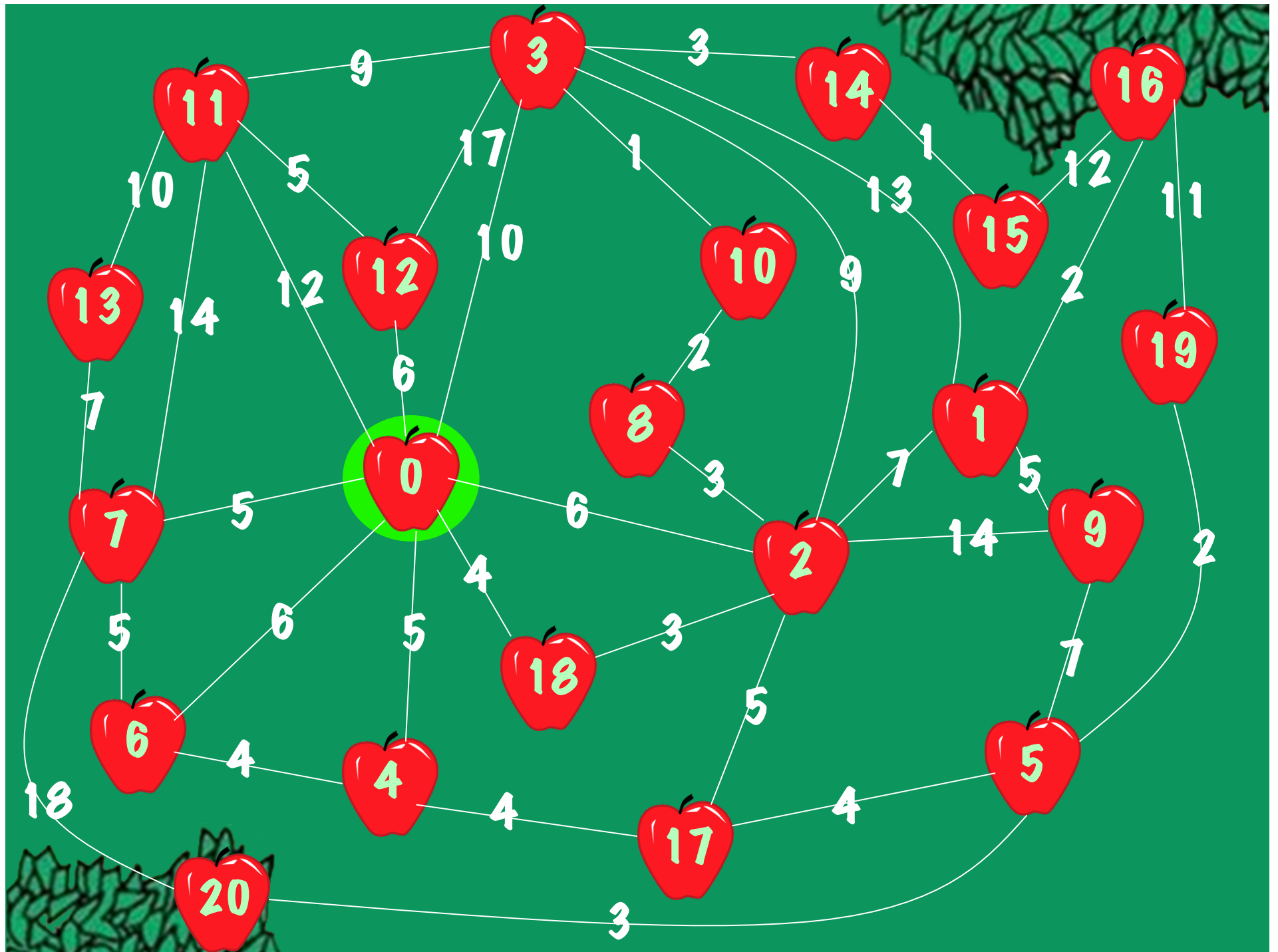
One more thing: Boruvka Step

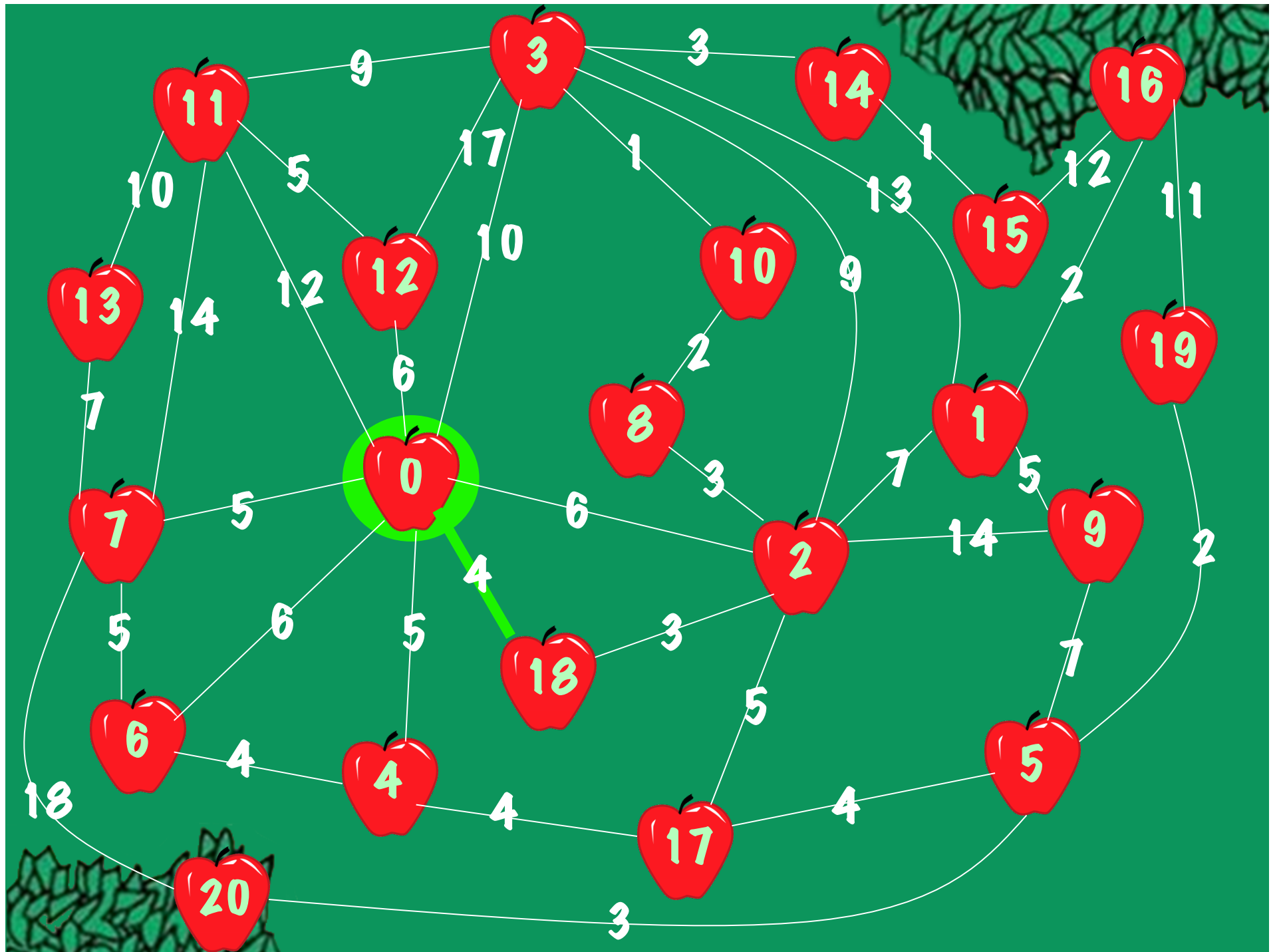
- For each vertex v , select the minimum-weight edge incident to v .
- The selected edges defined connected components; coalesce these components into single vertices.
- Delete all resulting loops, isolated vertices, and all but the lowest-weight edge among each set of multiple edges.

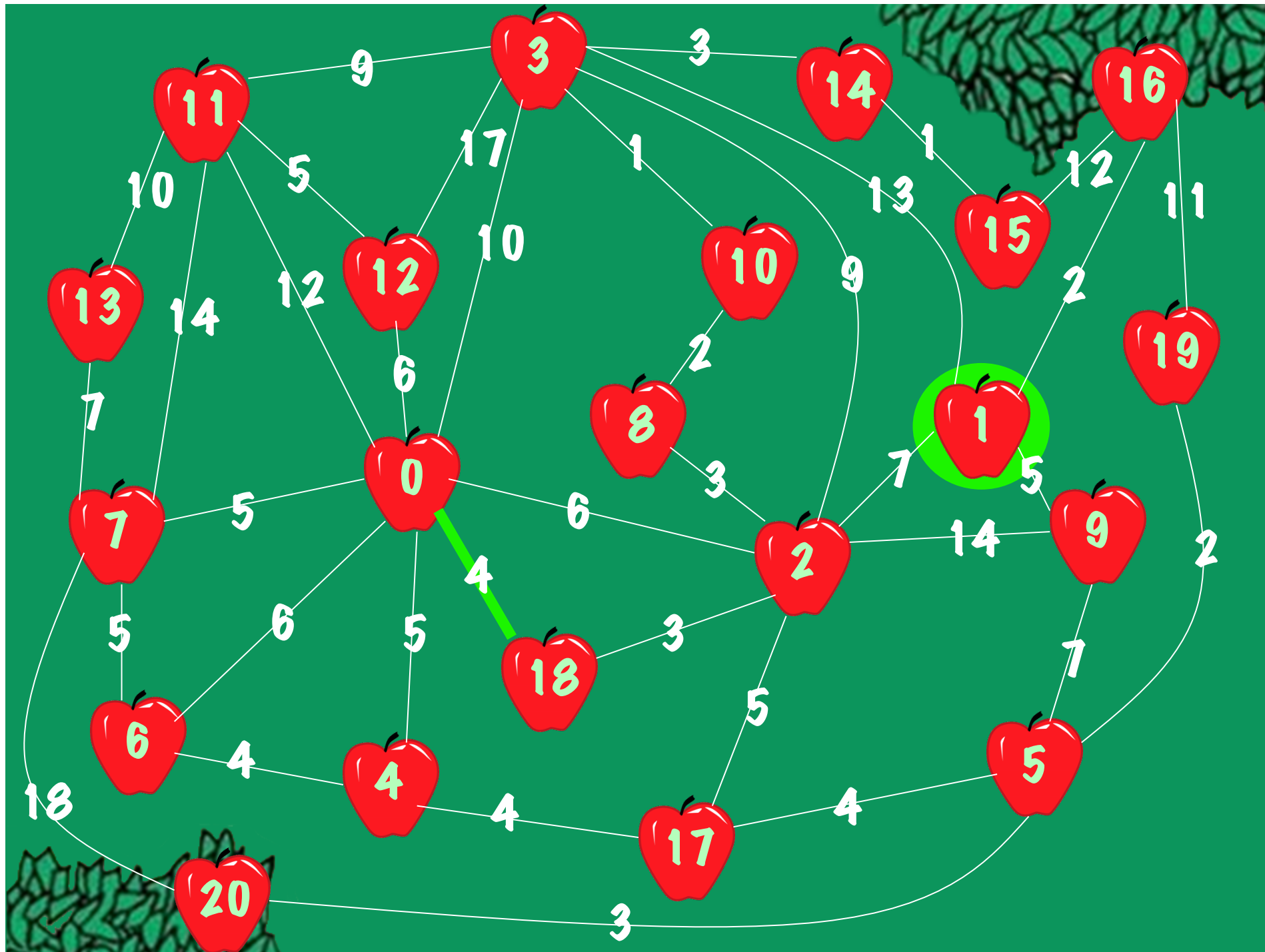
The Algorithm

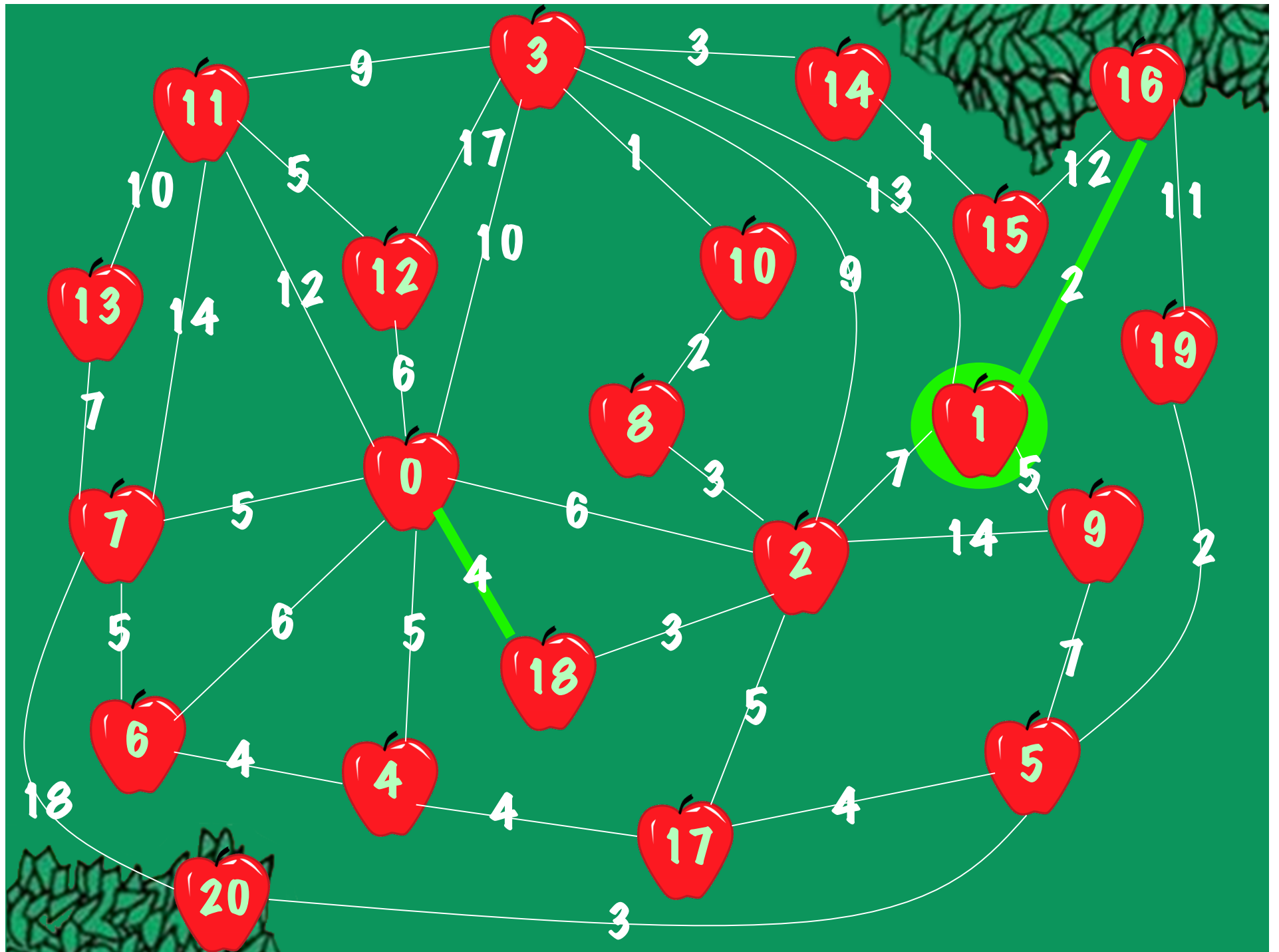
- STEP 1: Perform two Boruvka steps.
- STEP 2: Choose a subgraph \mathcal{H} by selecting each edge independently with probability $1/2$. Apply the algorithm to the subgraph, producing a minimum spanning forest \mathcal{F} of \mathcal{H} . Throw away all **F-heavy** edges in the entire graph.
- STEP 3: Recurse on the remaining graph to compute a spanning forest \mathcal{F}' . Return edges contracted in STEP 1 along with the edges of \mathcal{F}'

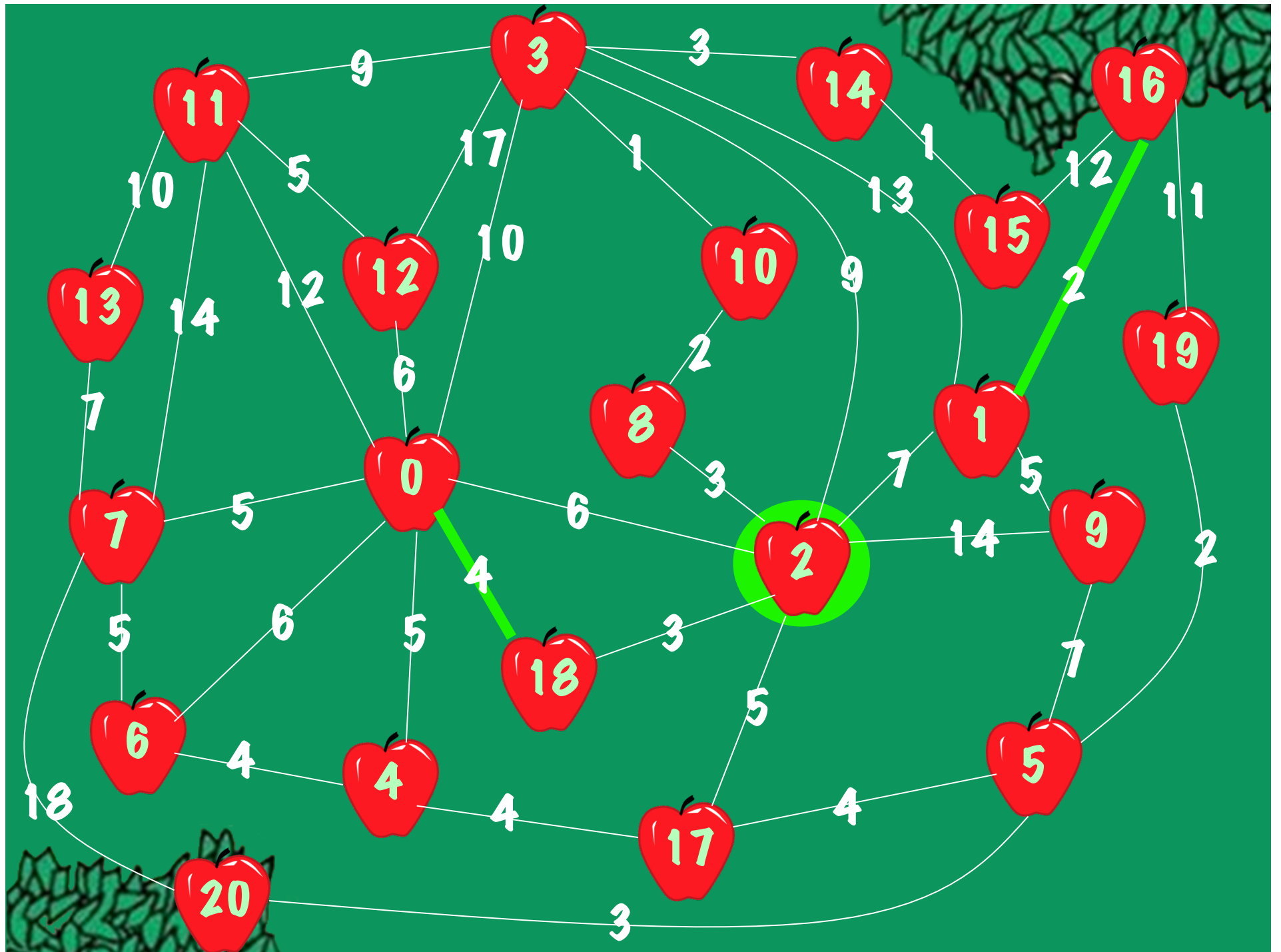


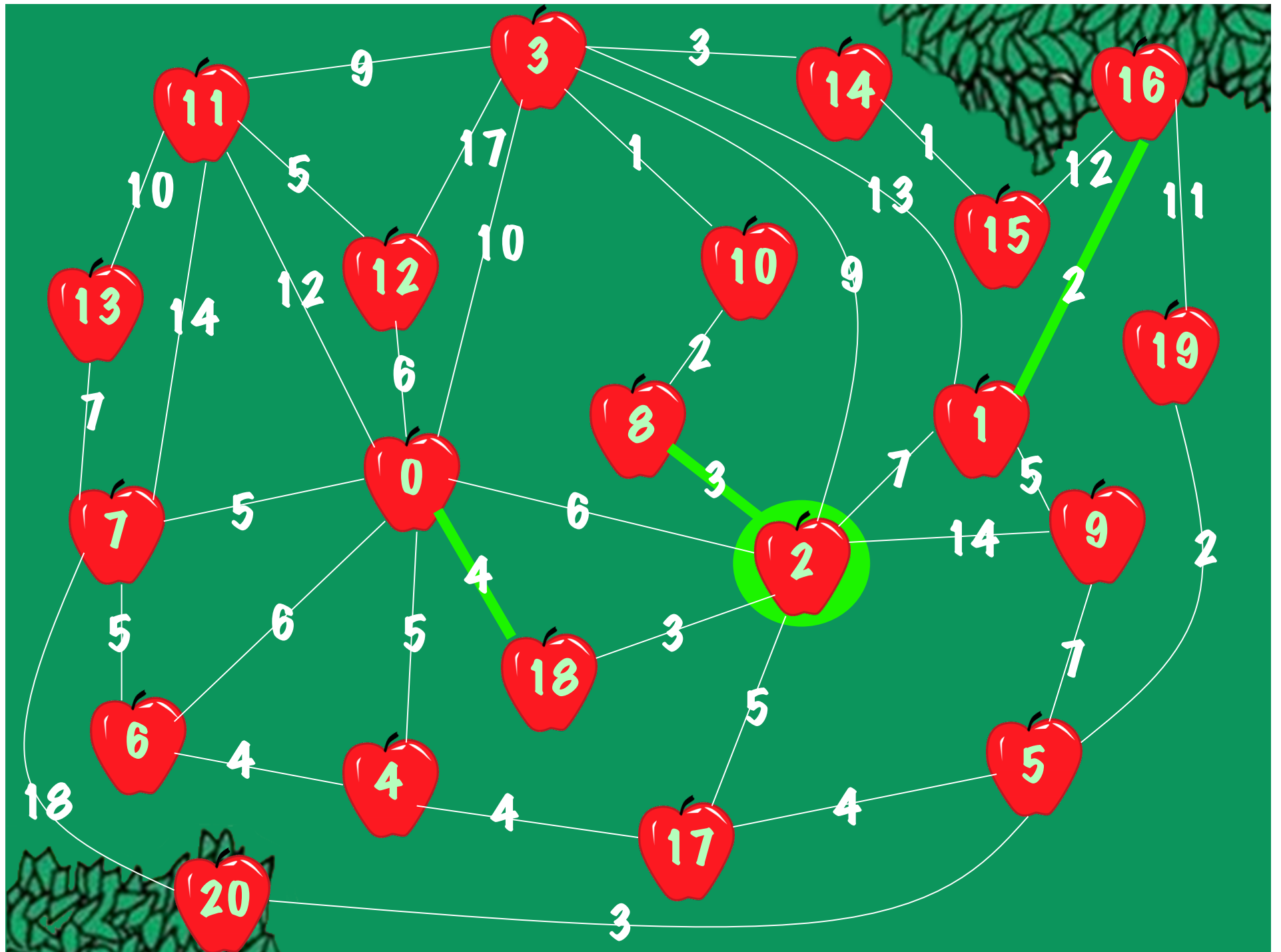


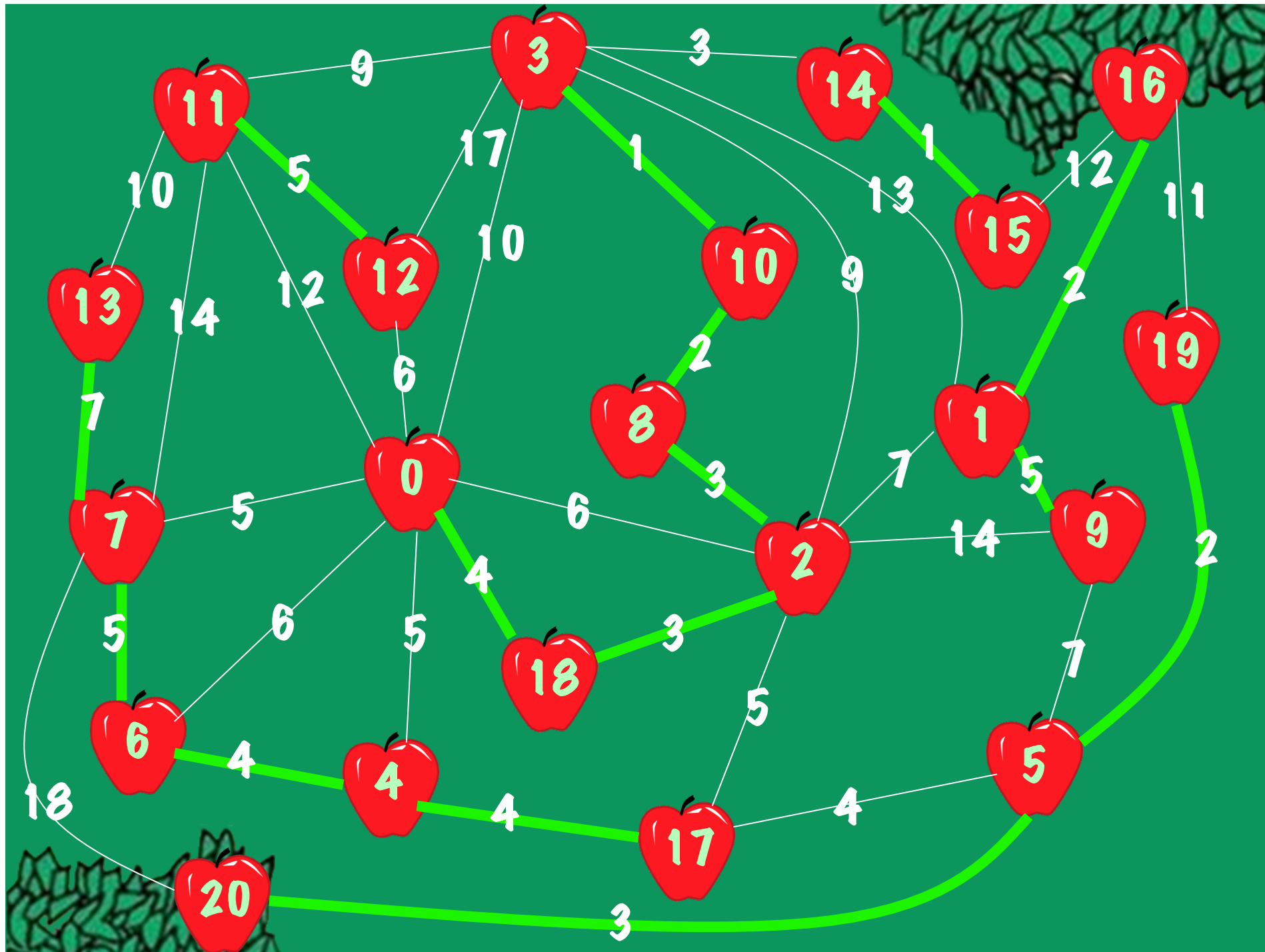


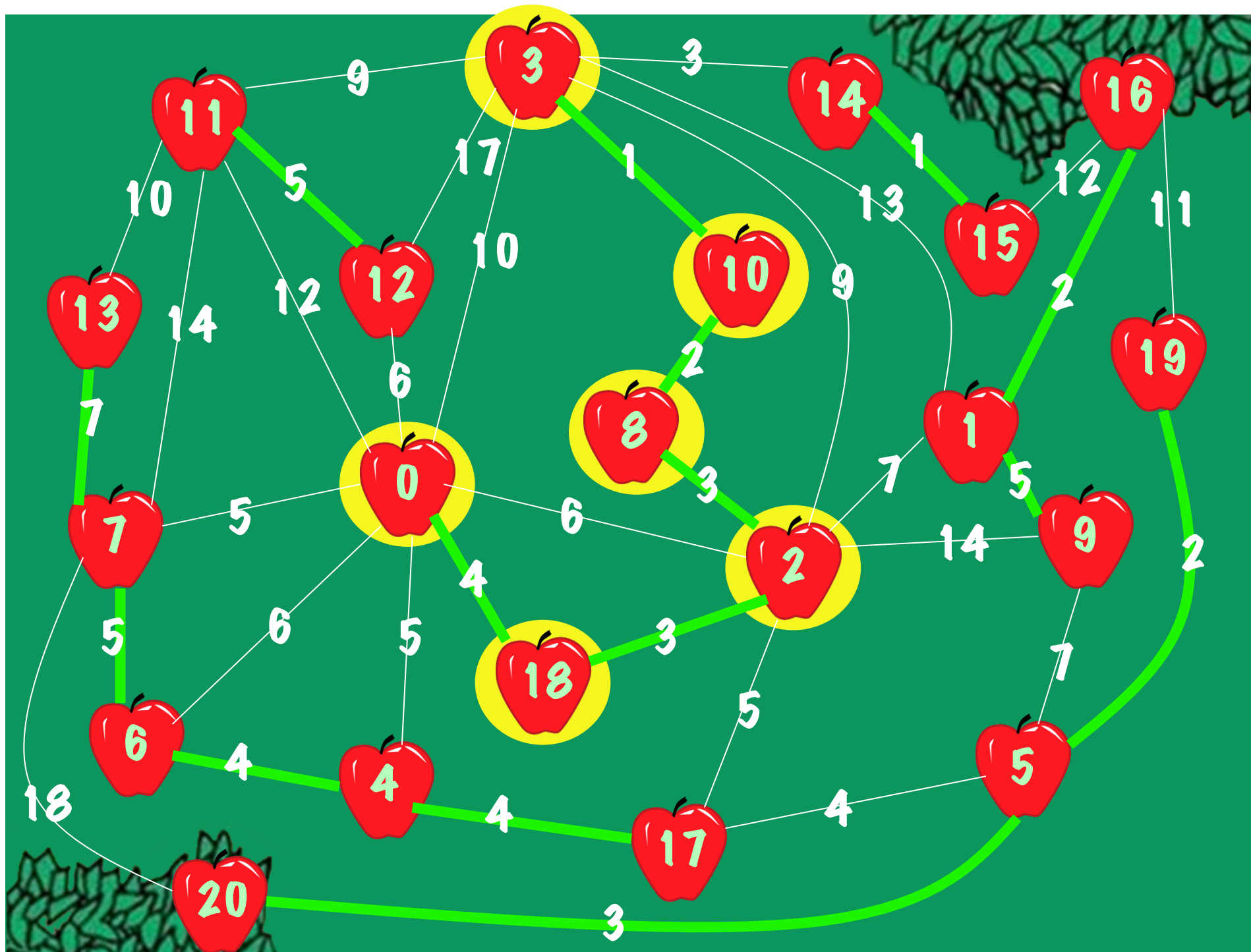


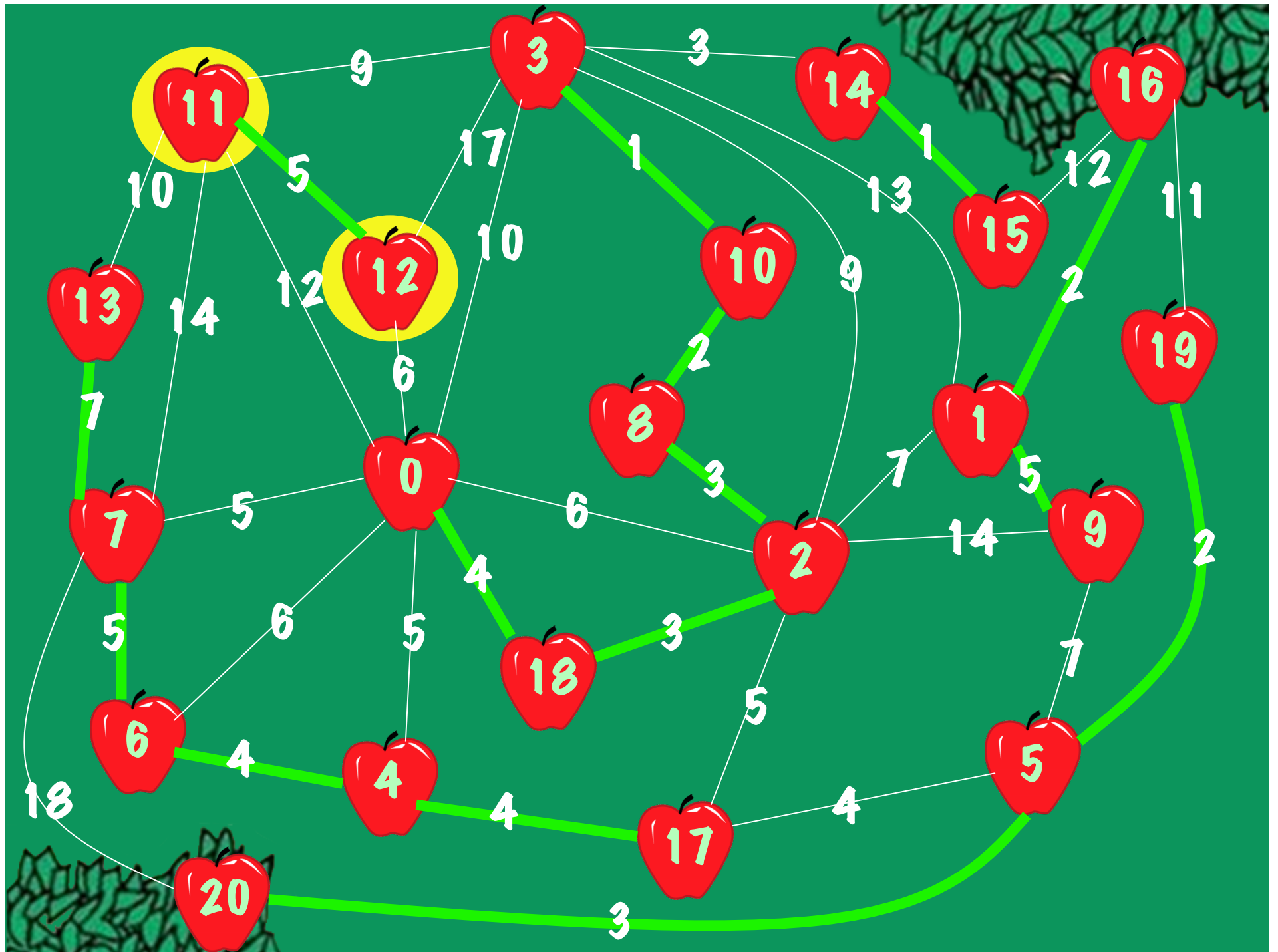


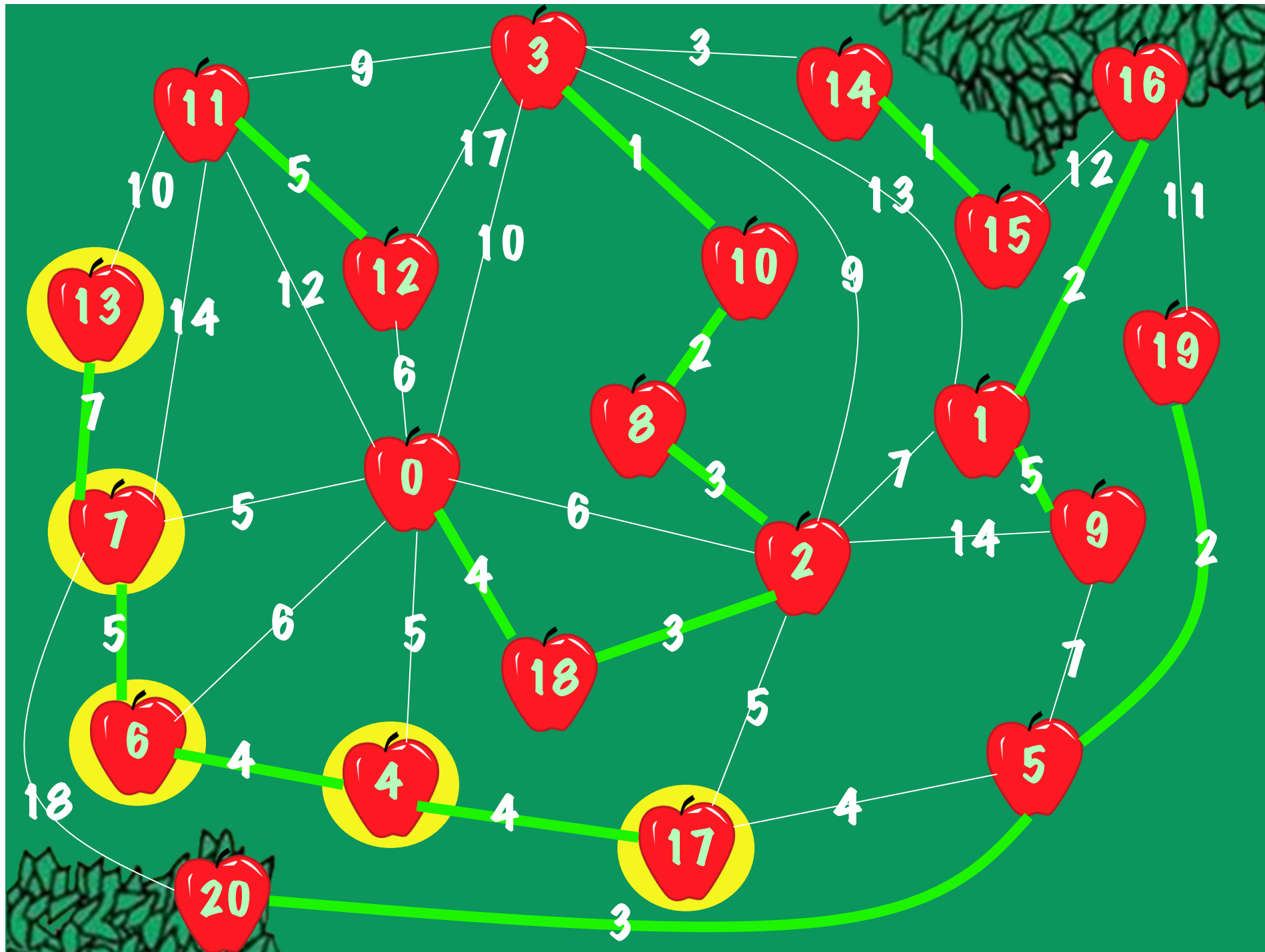


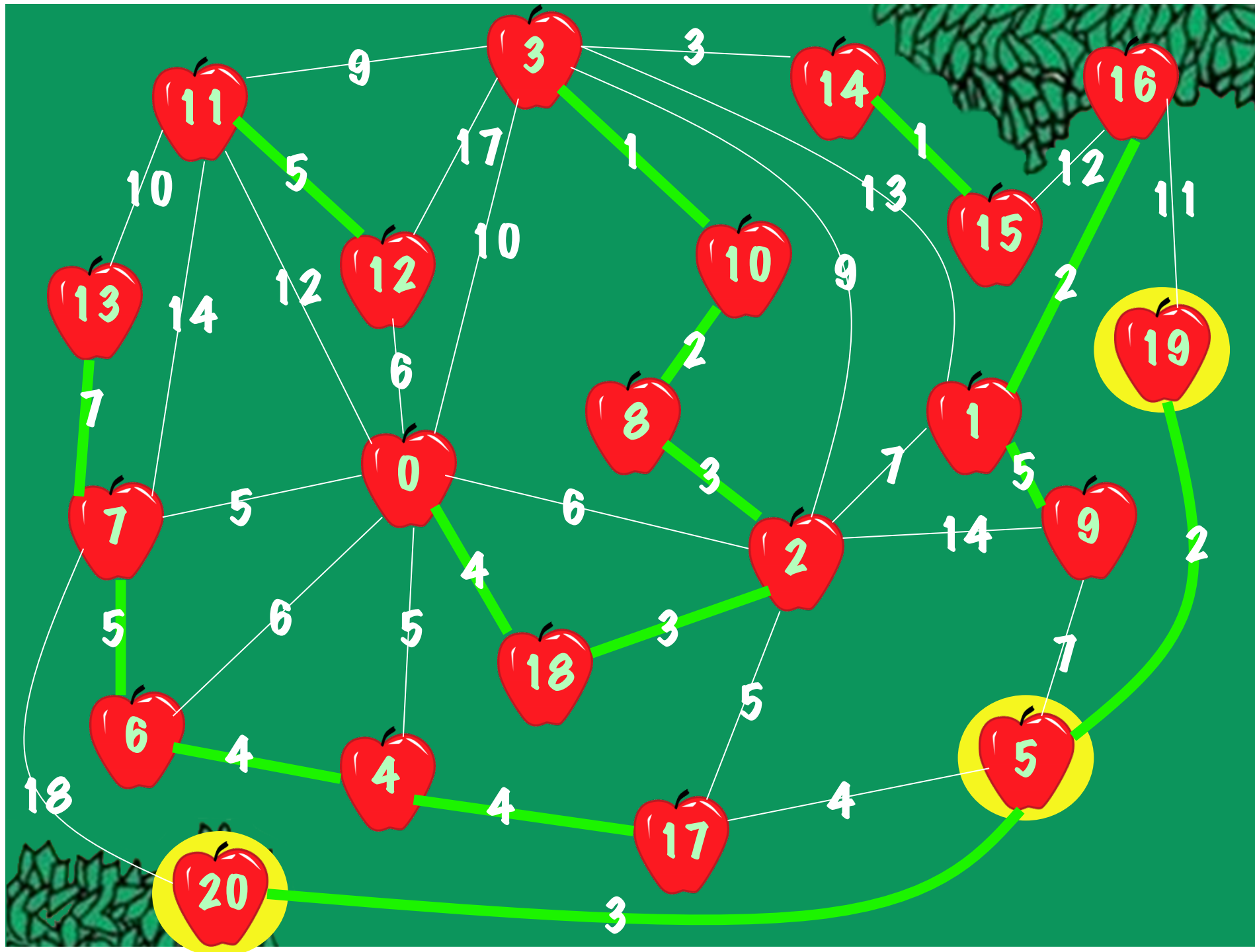


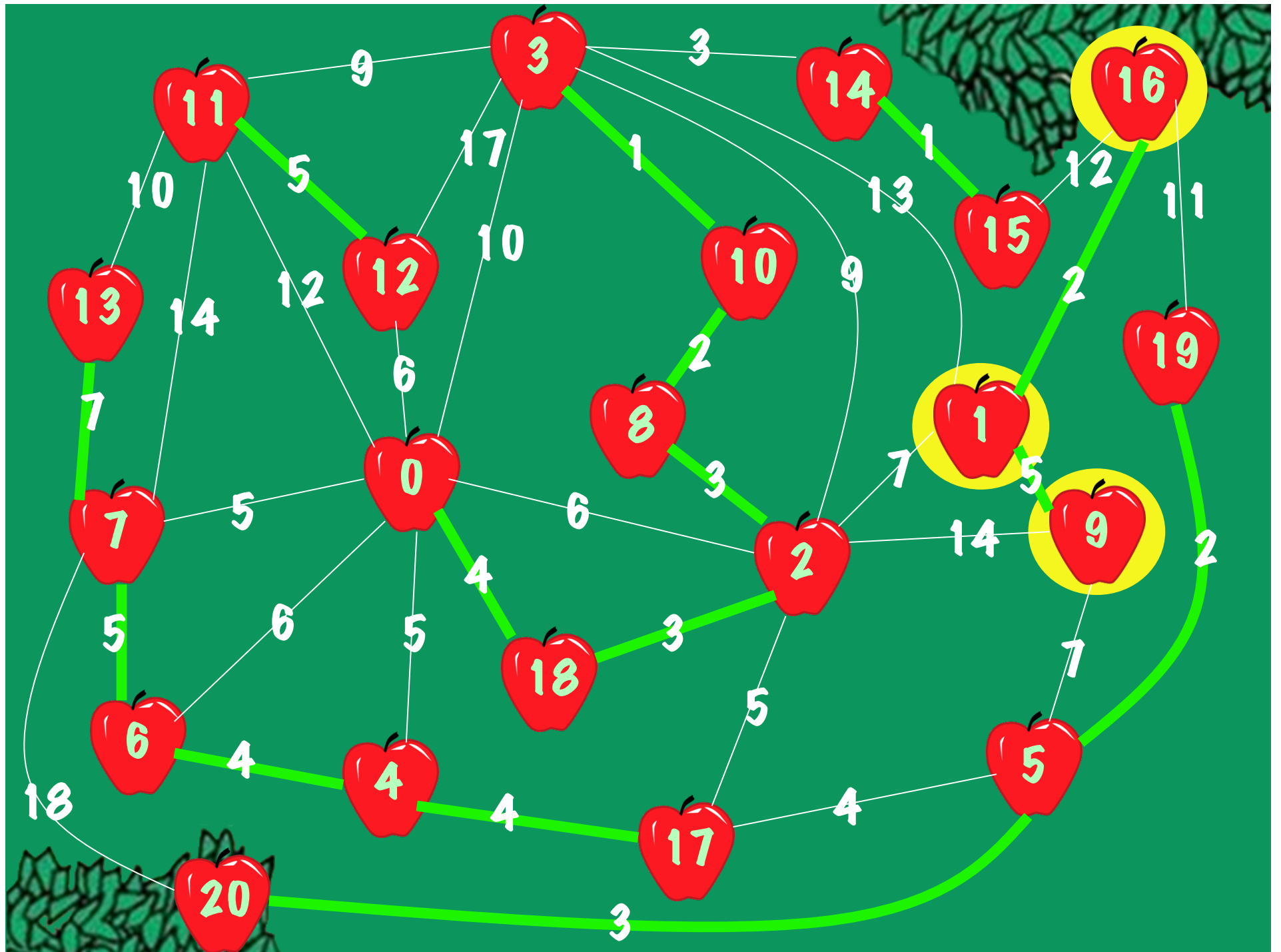


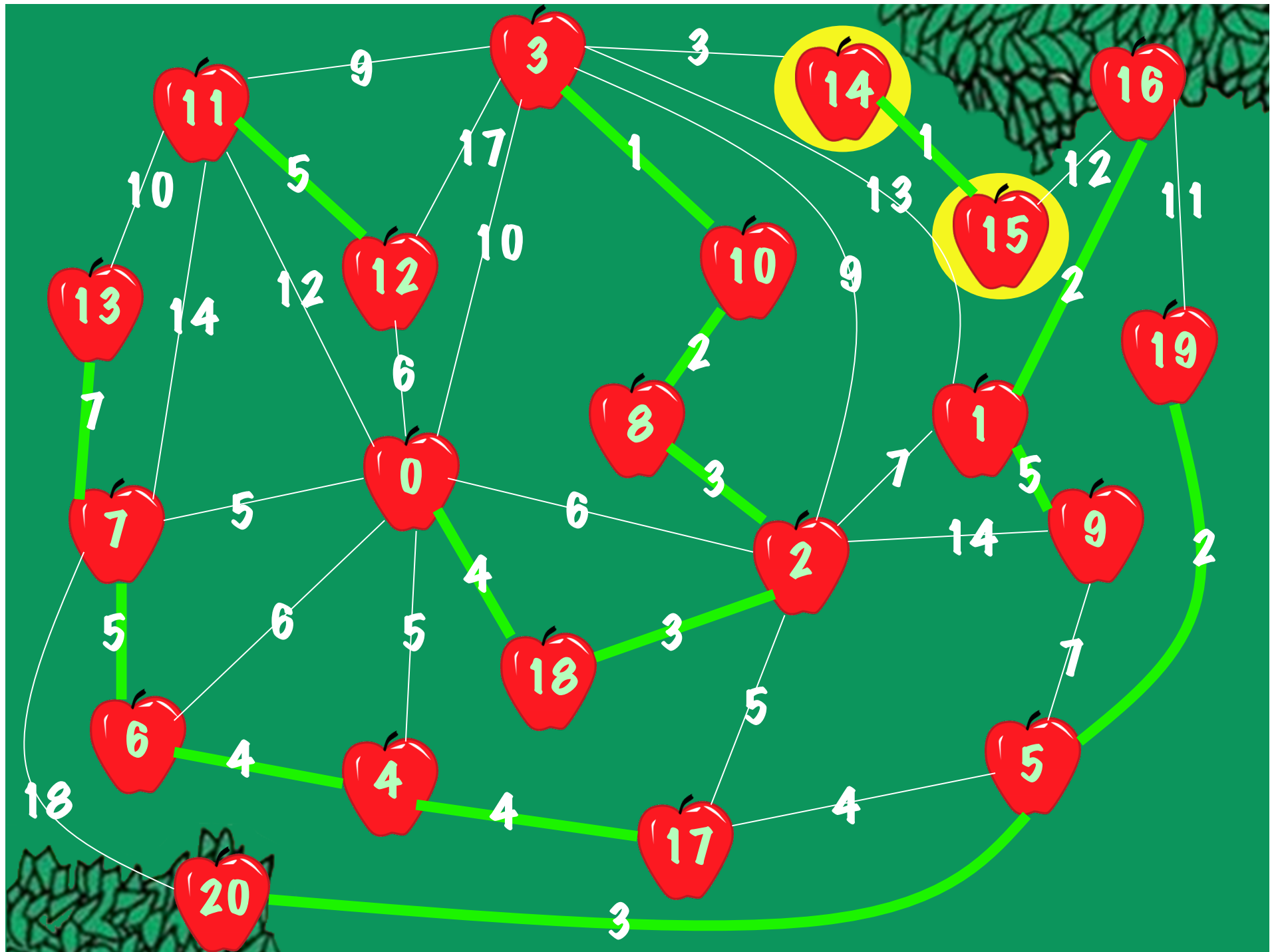


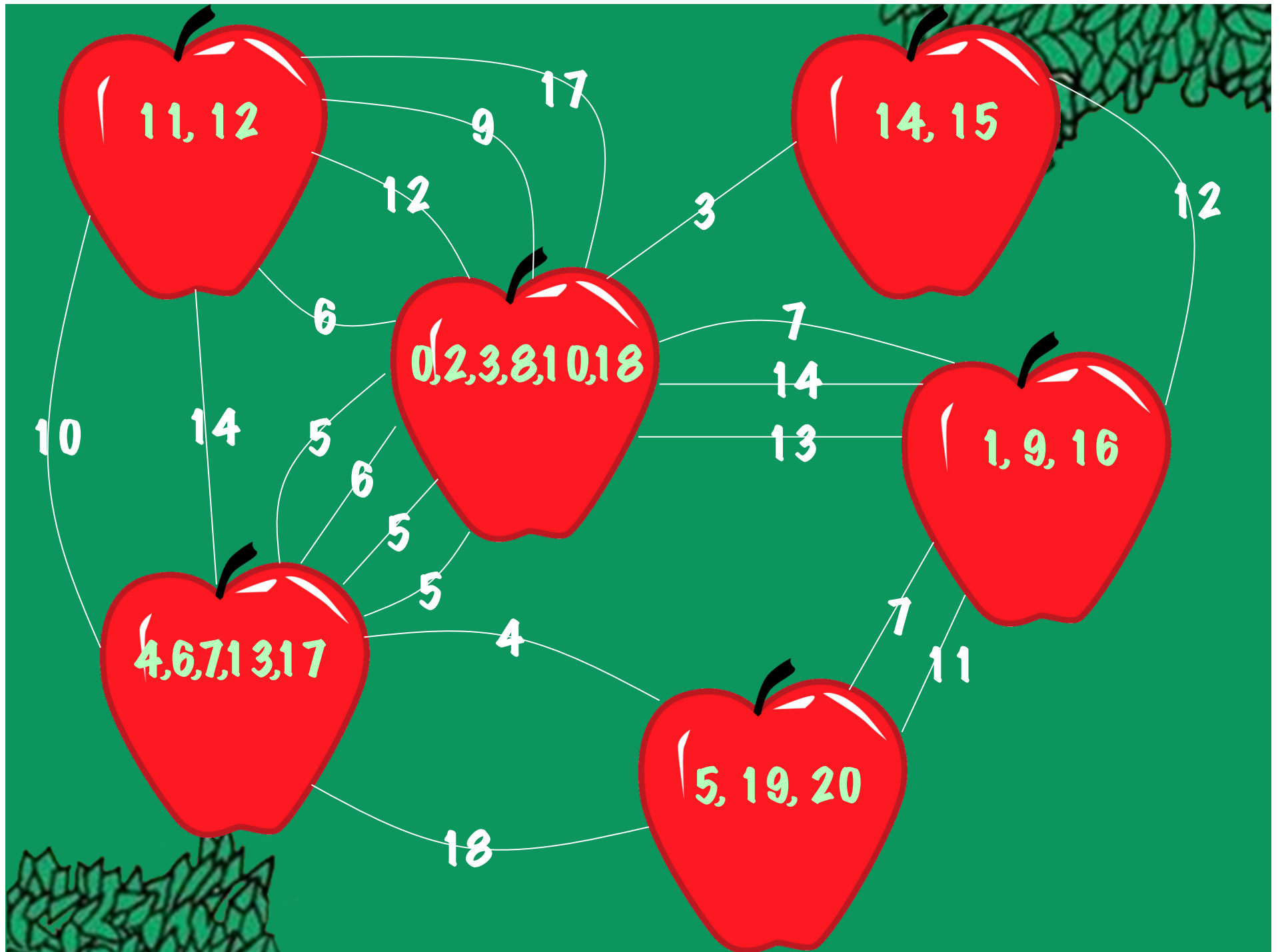


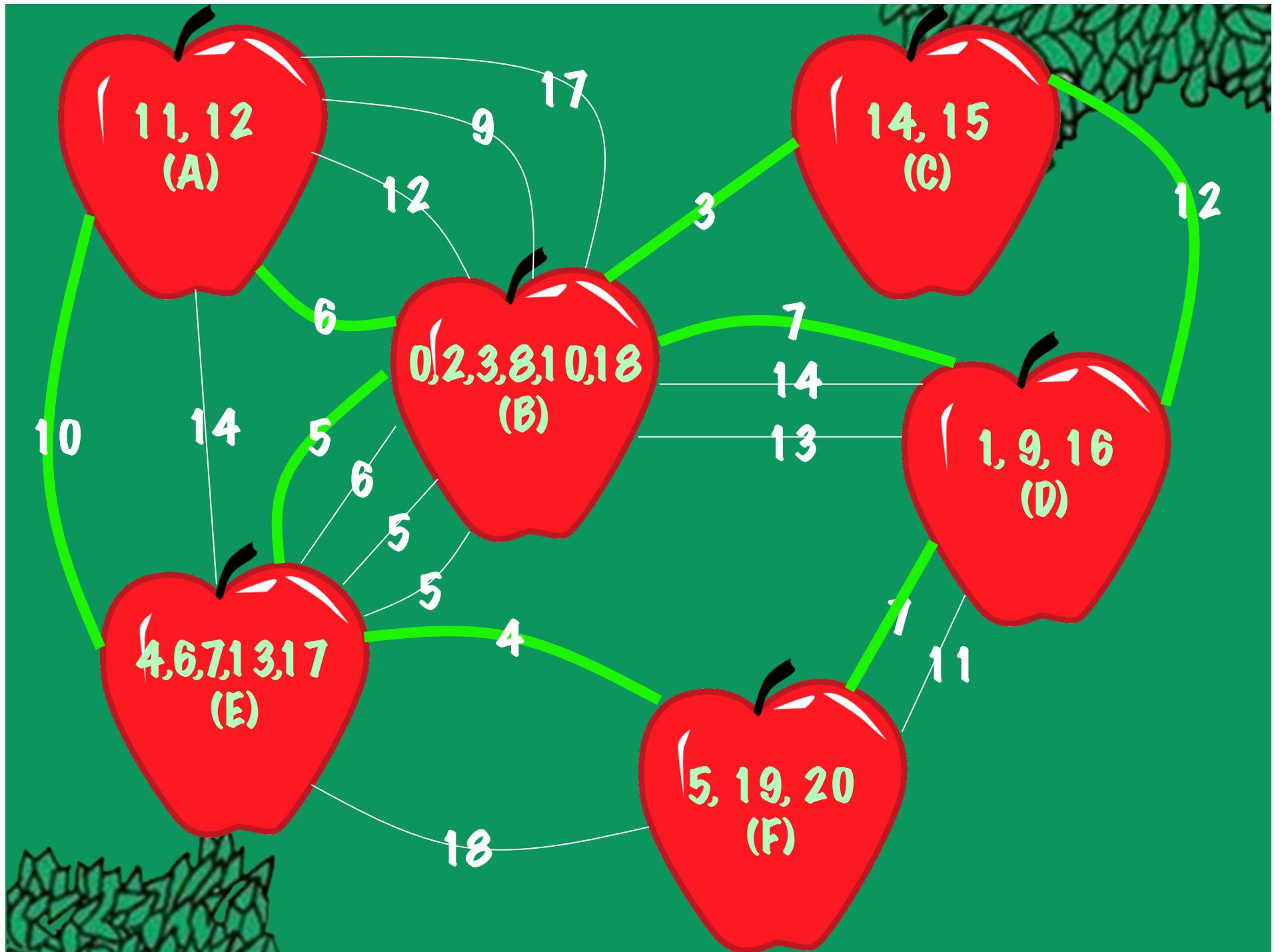


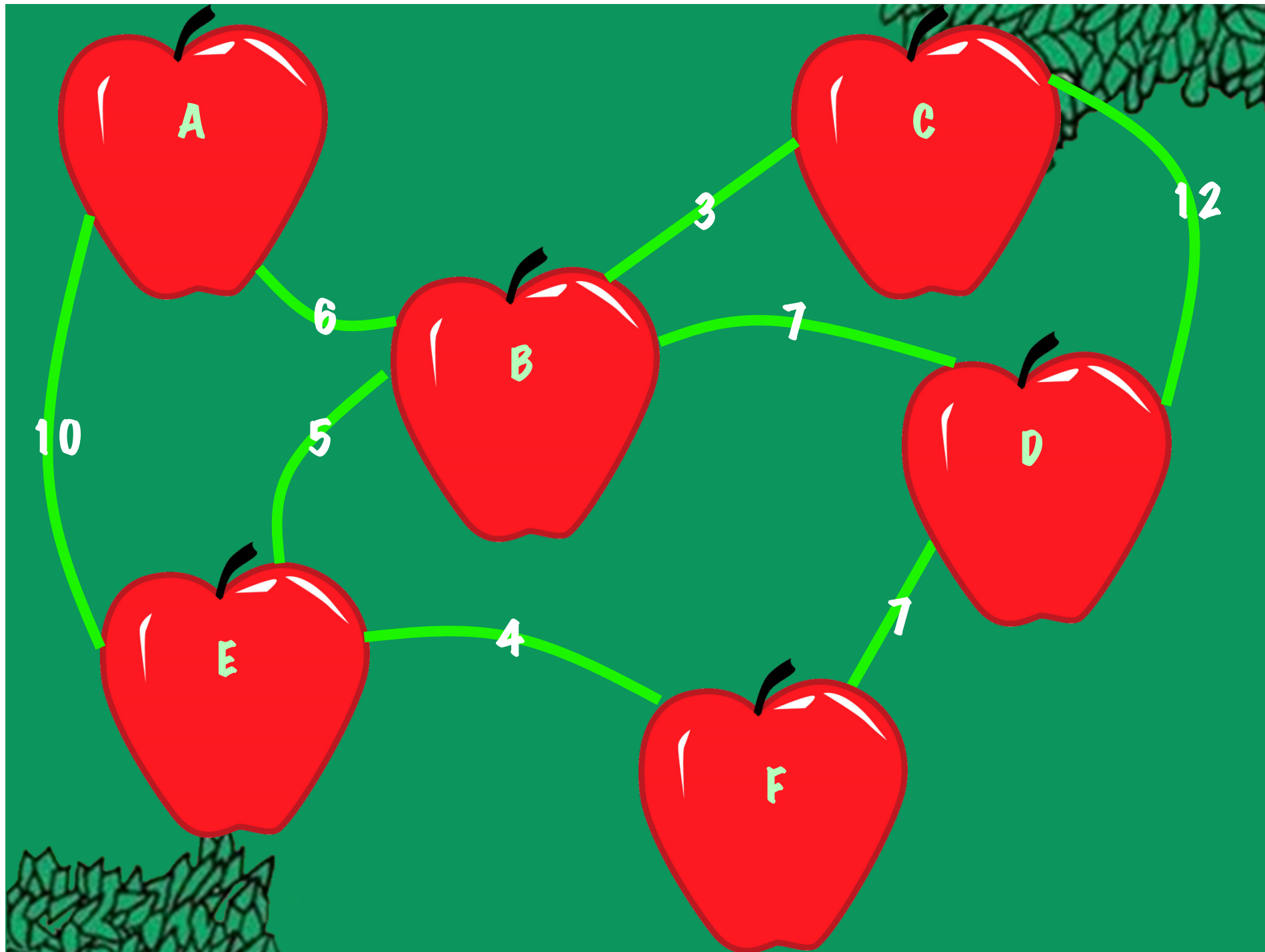


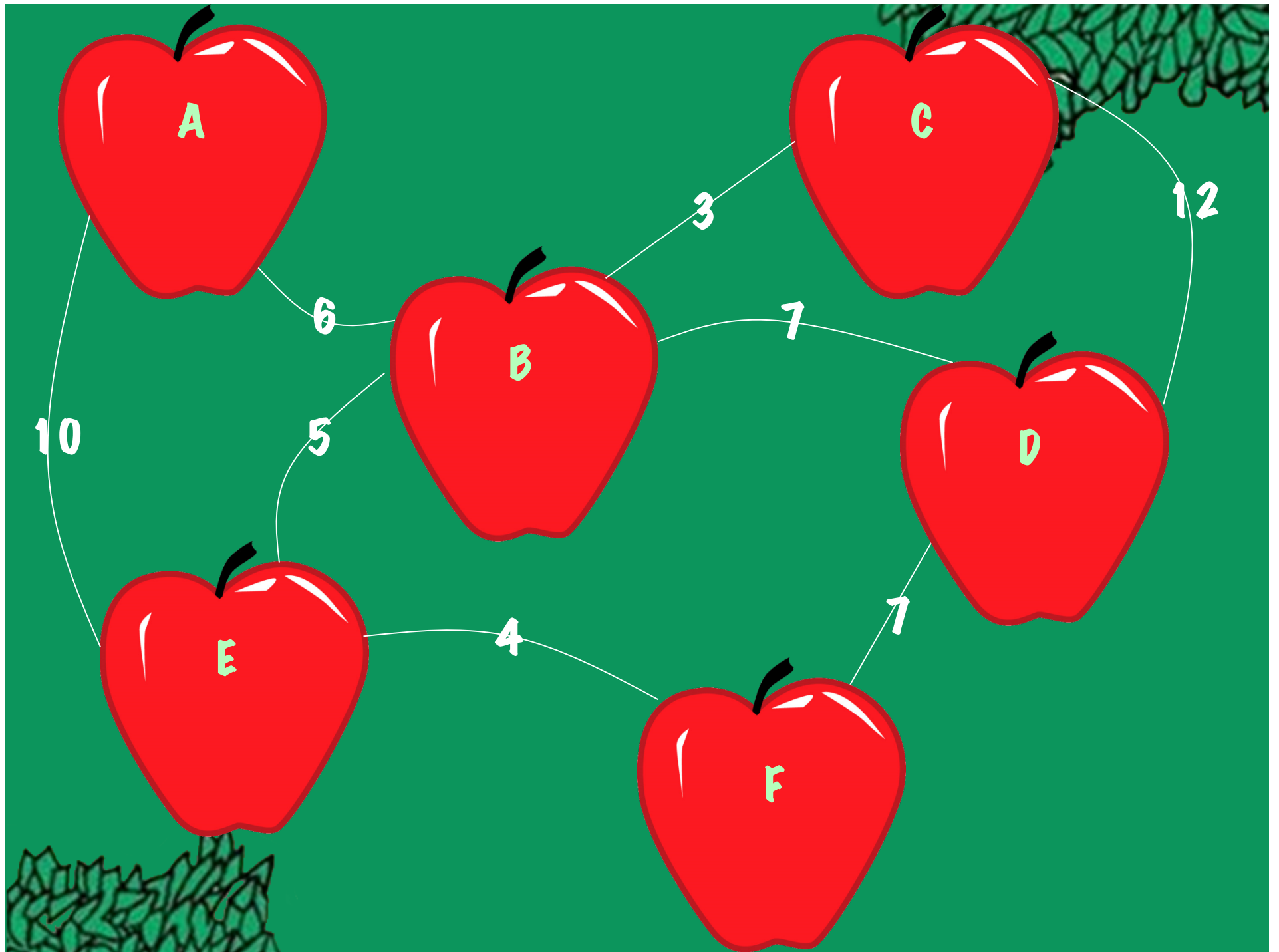


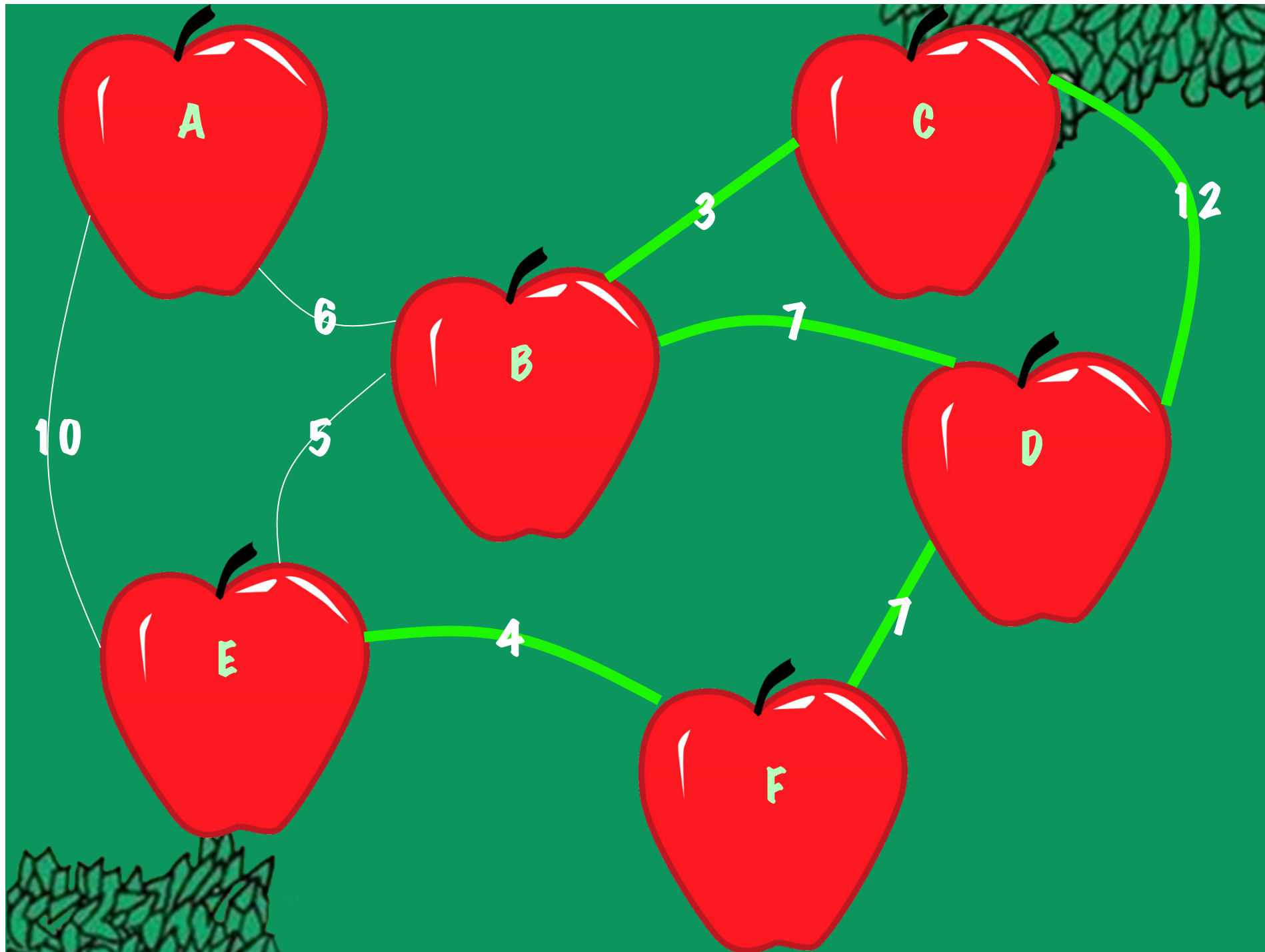


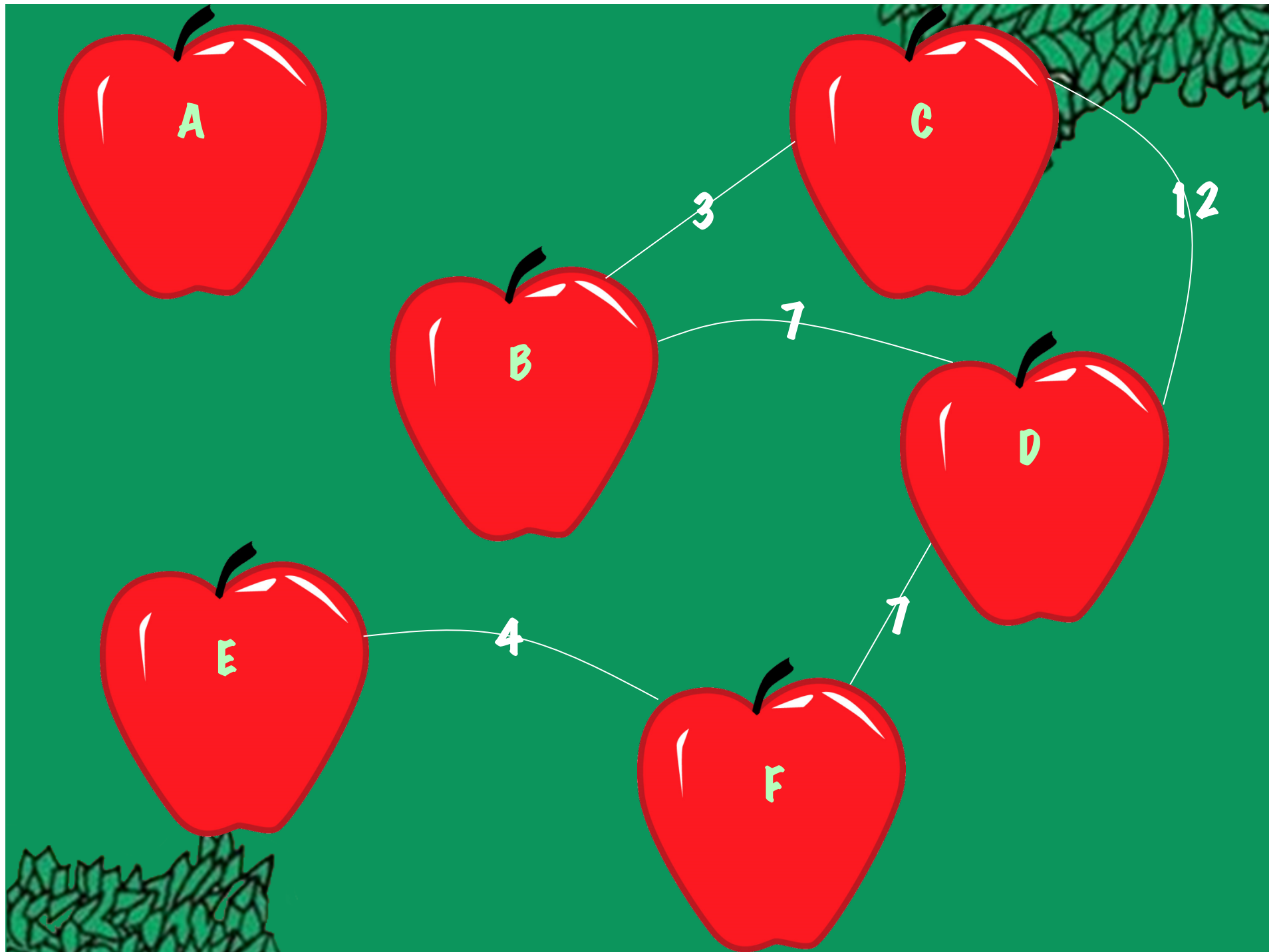


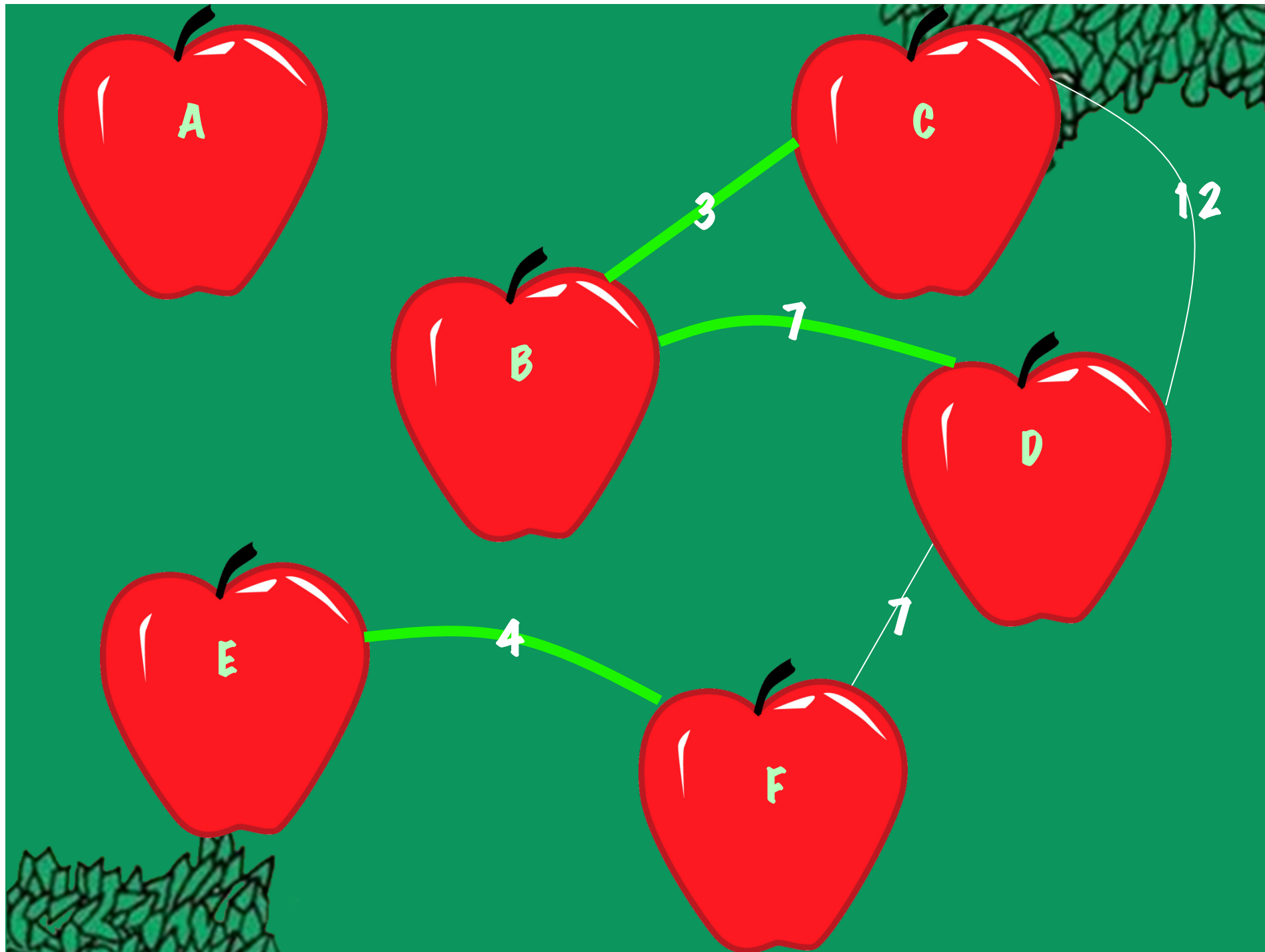


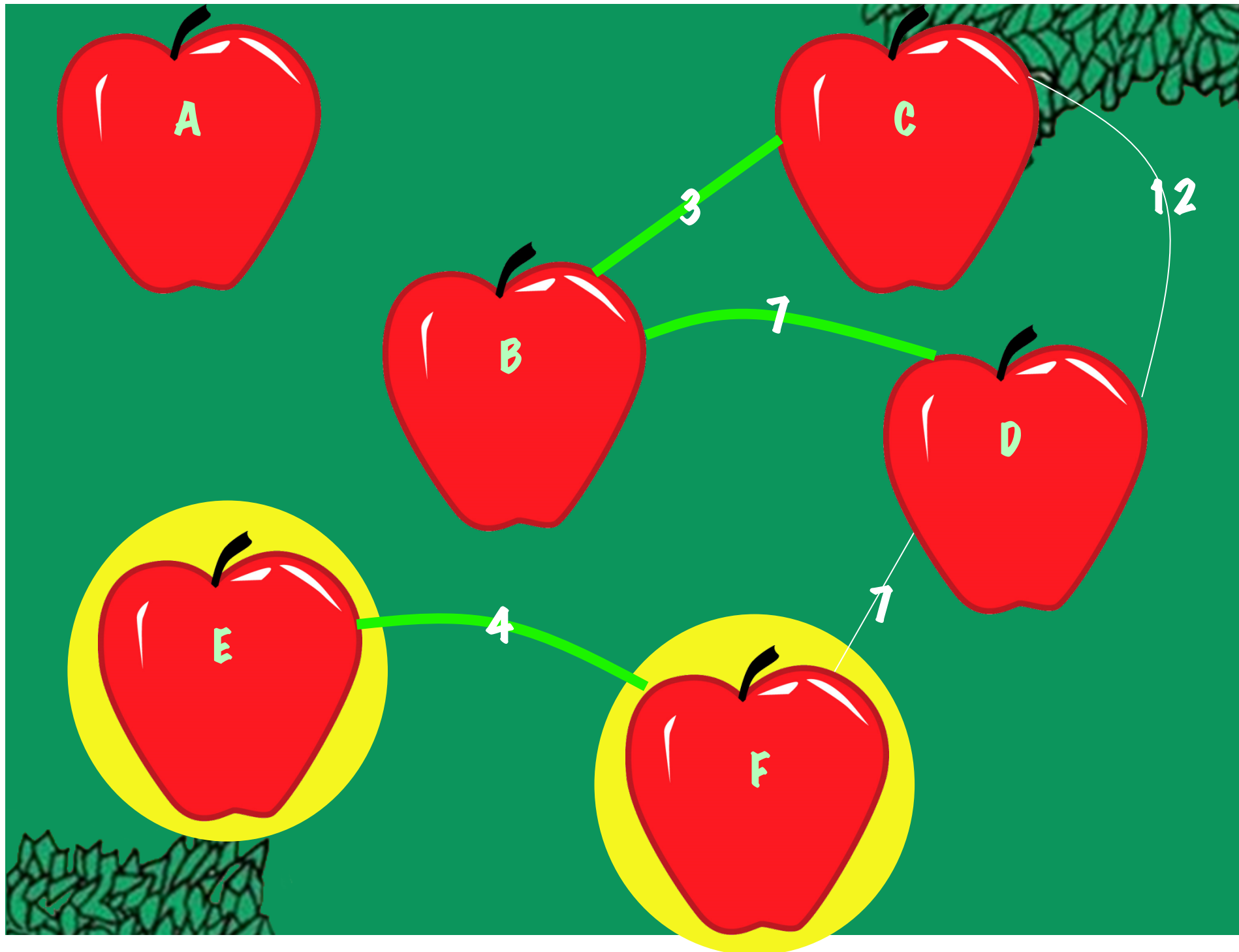


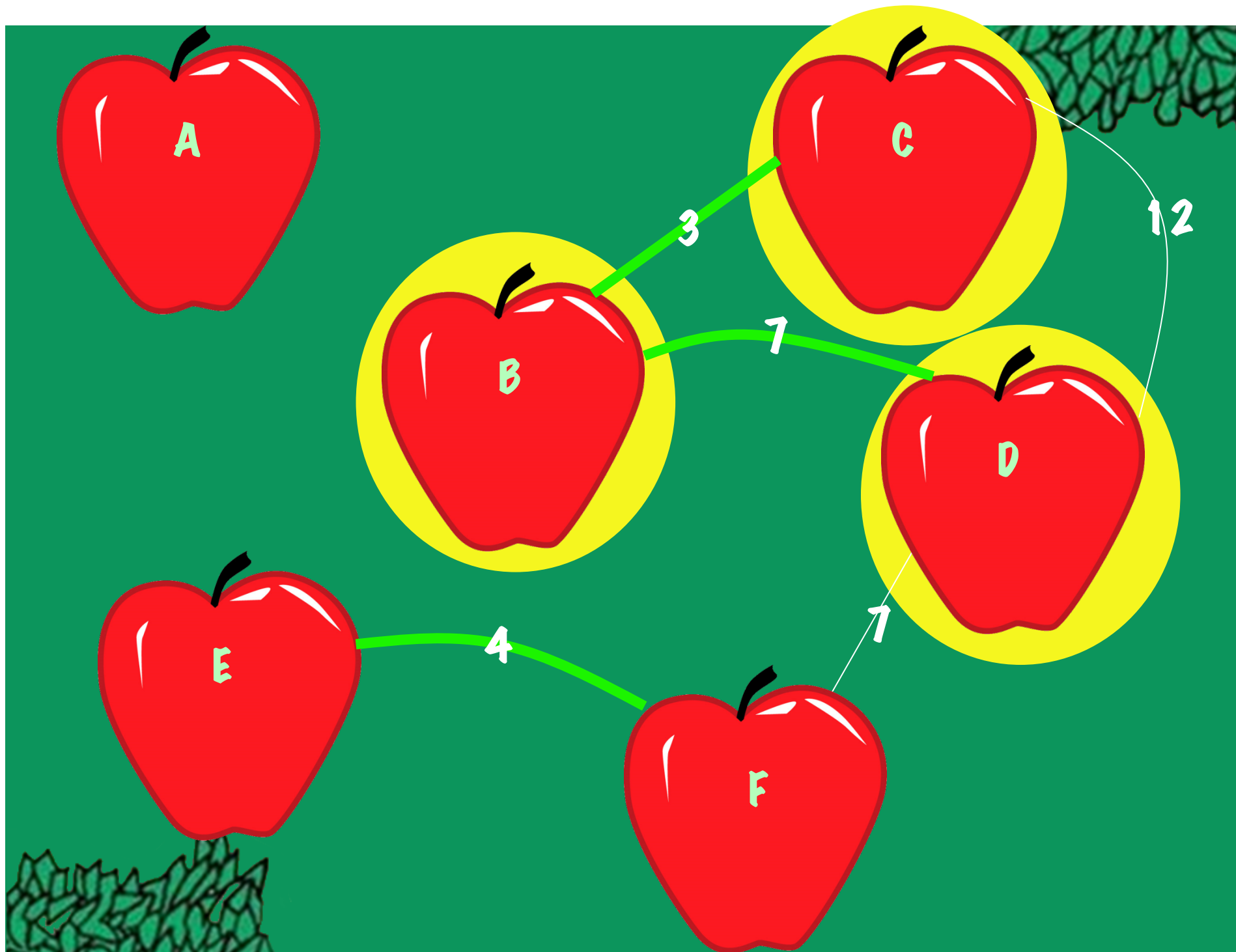


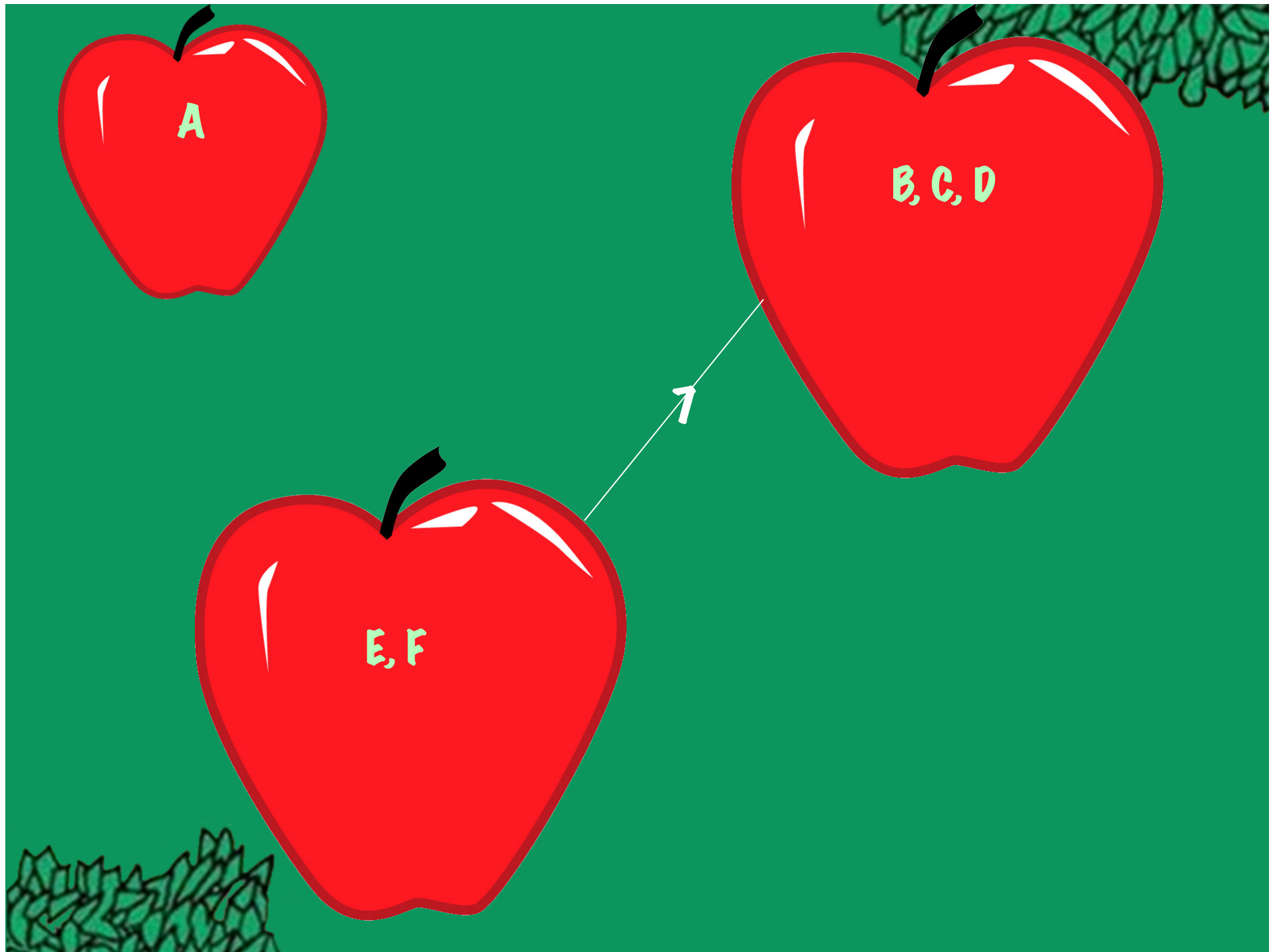


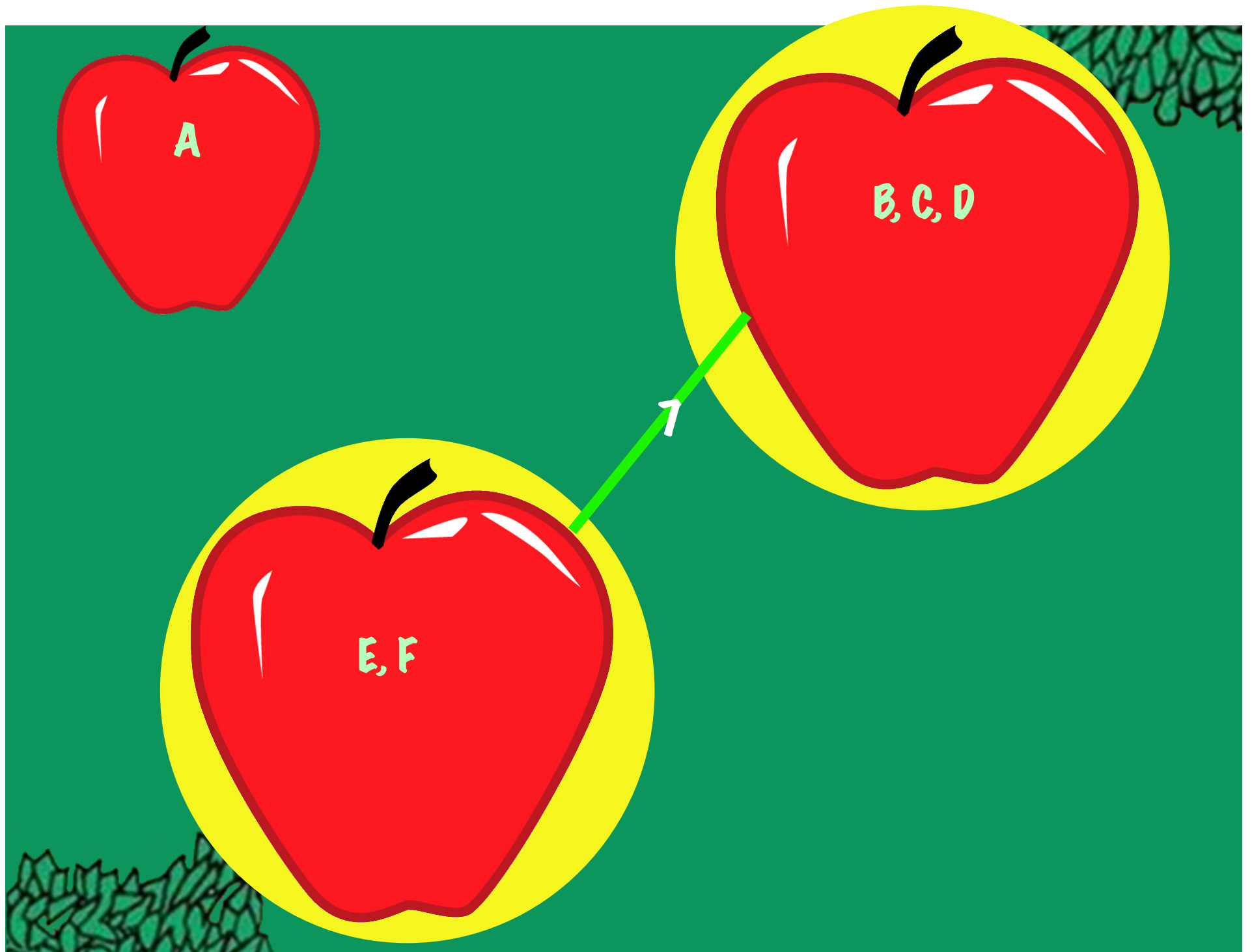




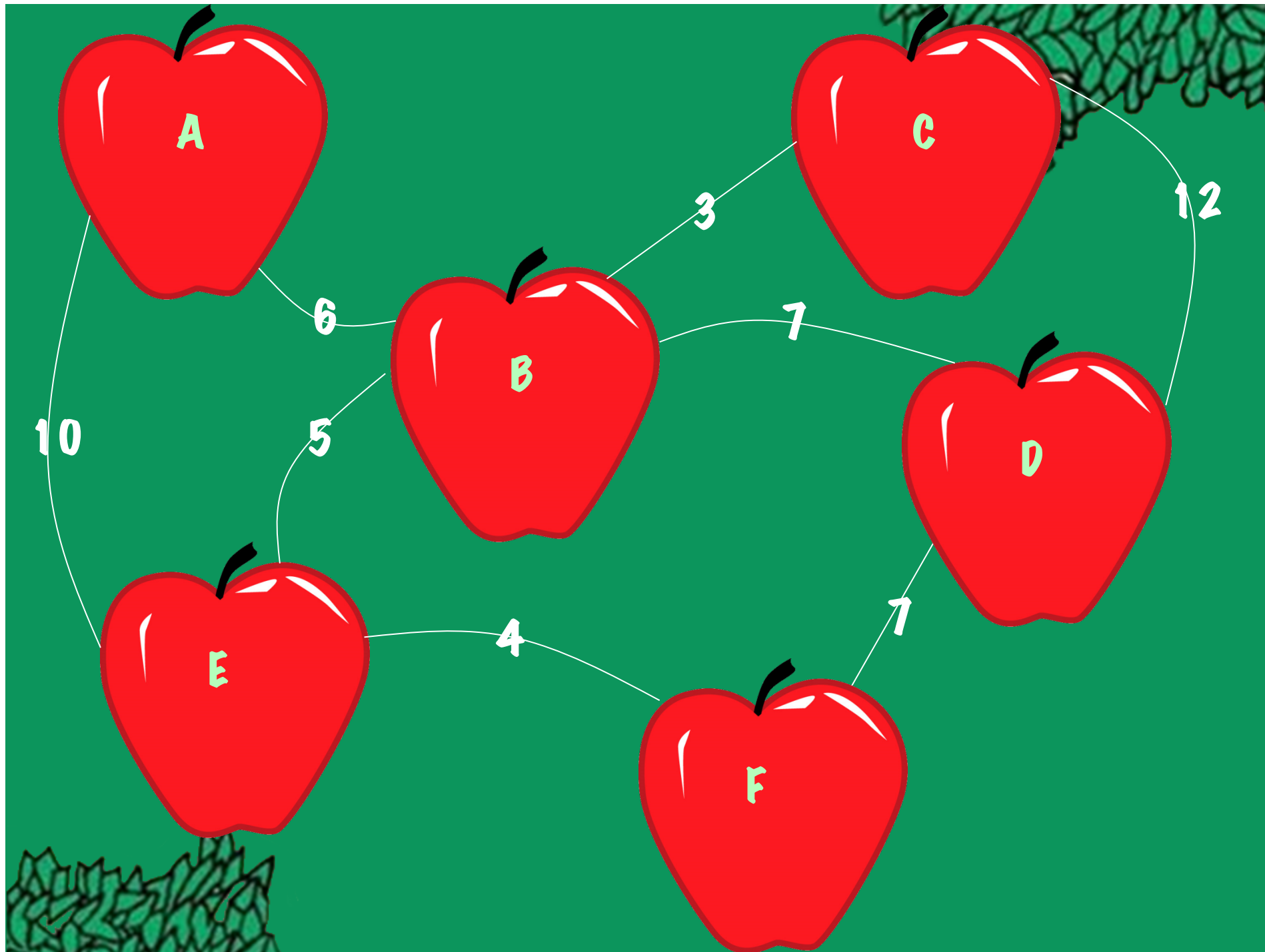


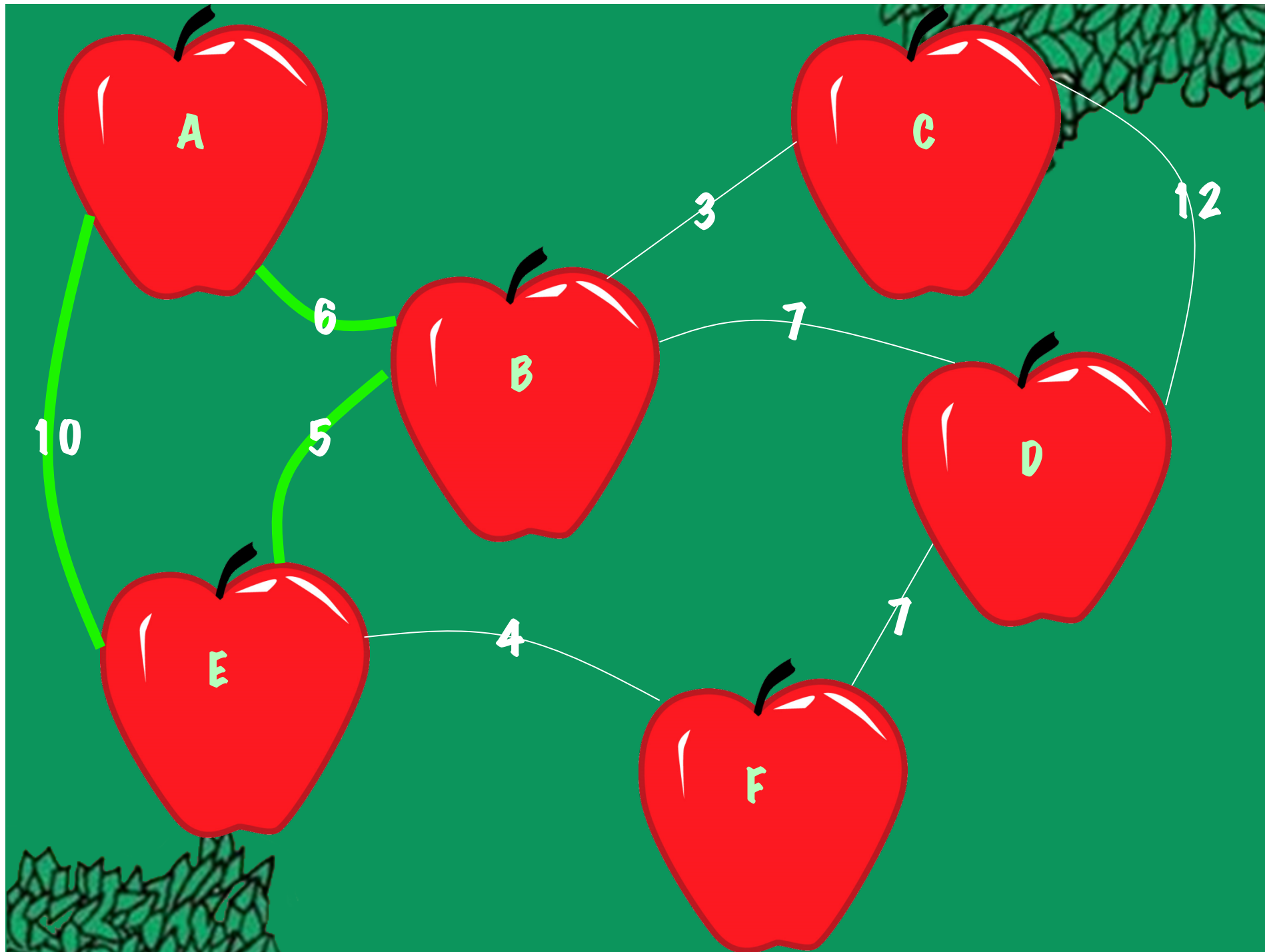


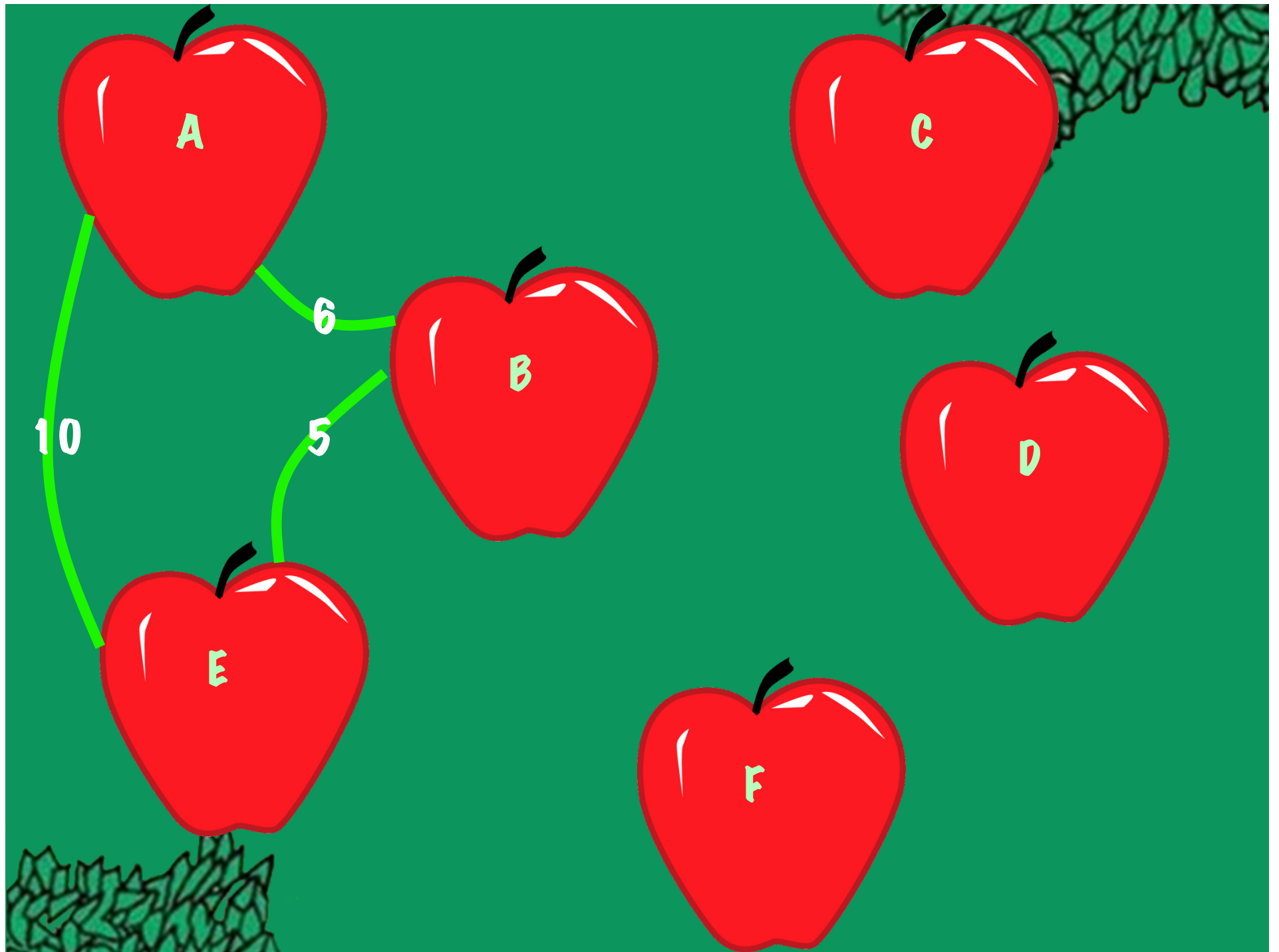


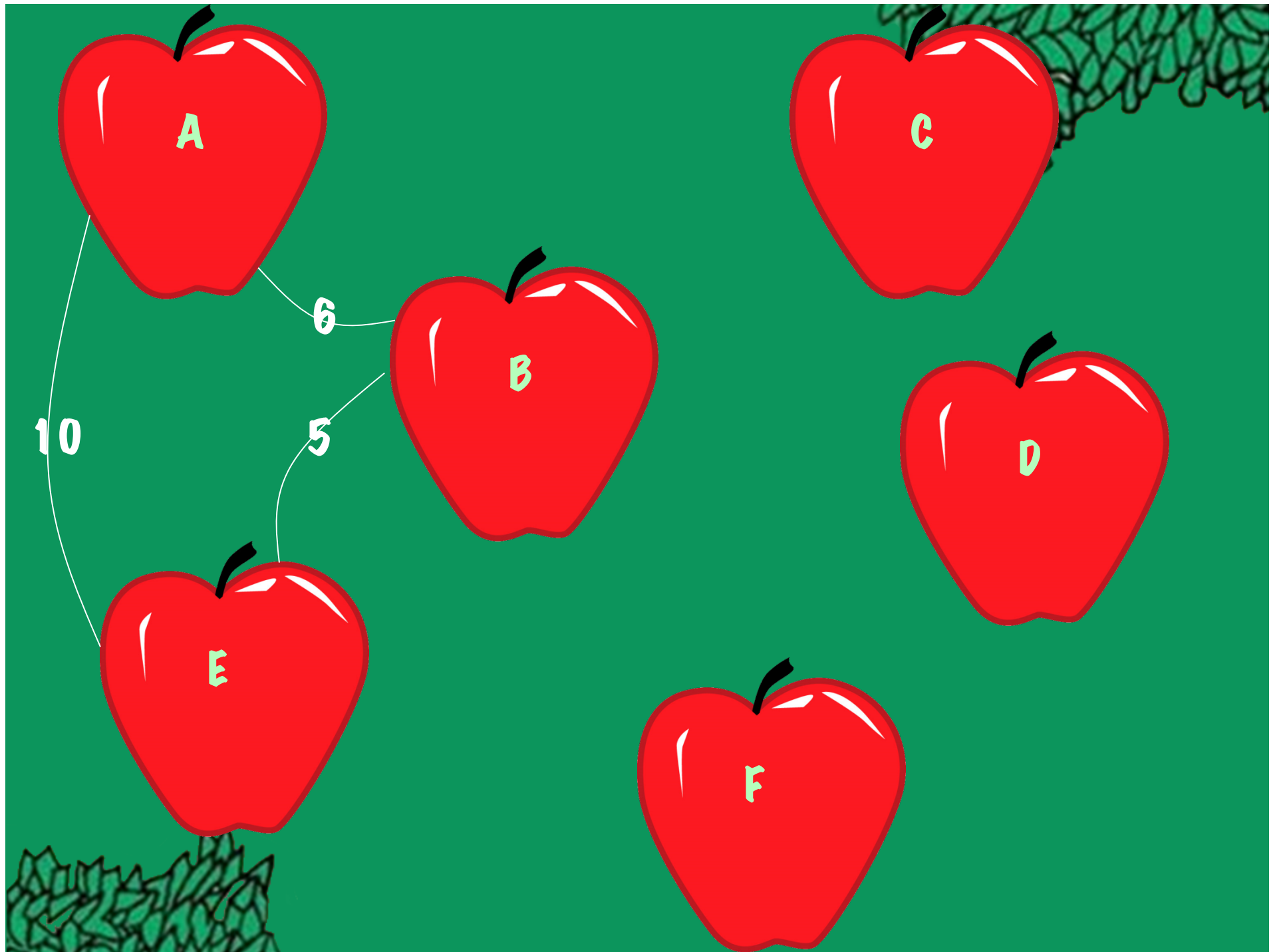


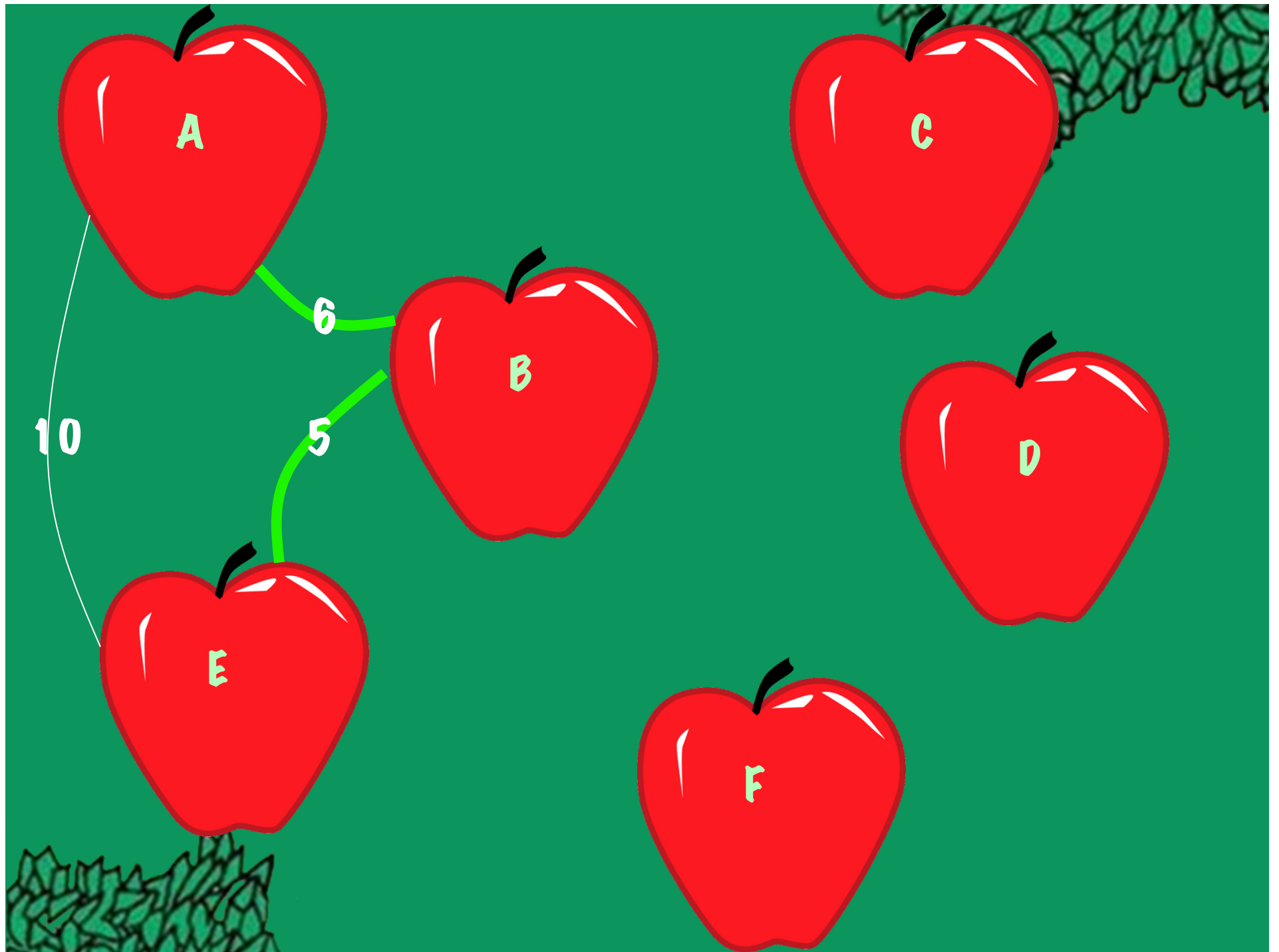


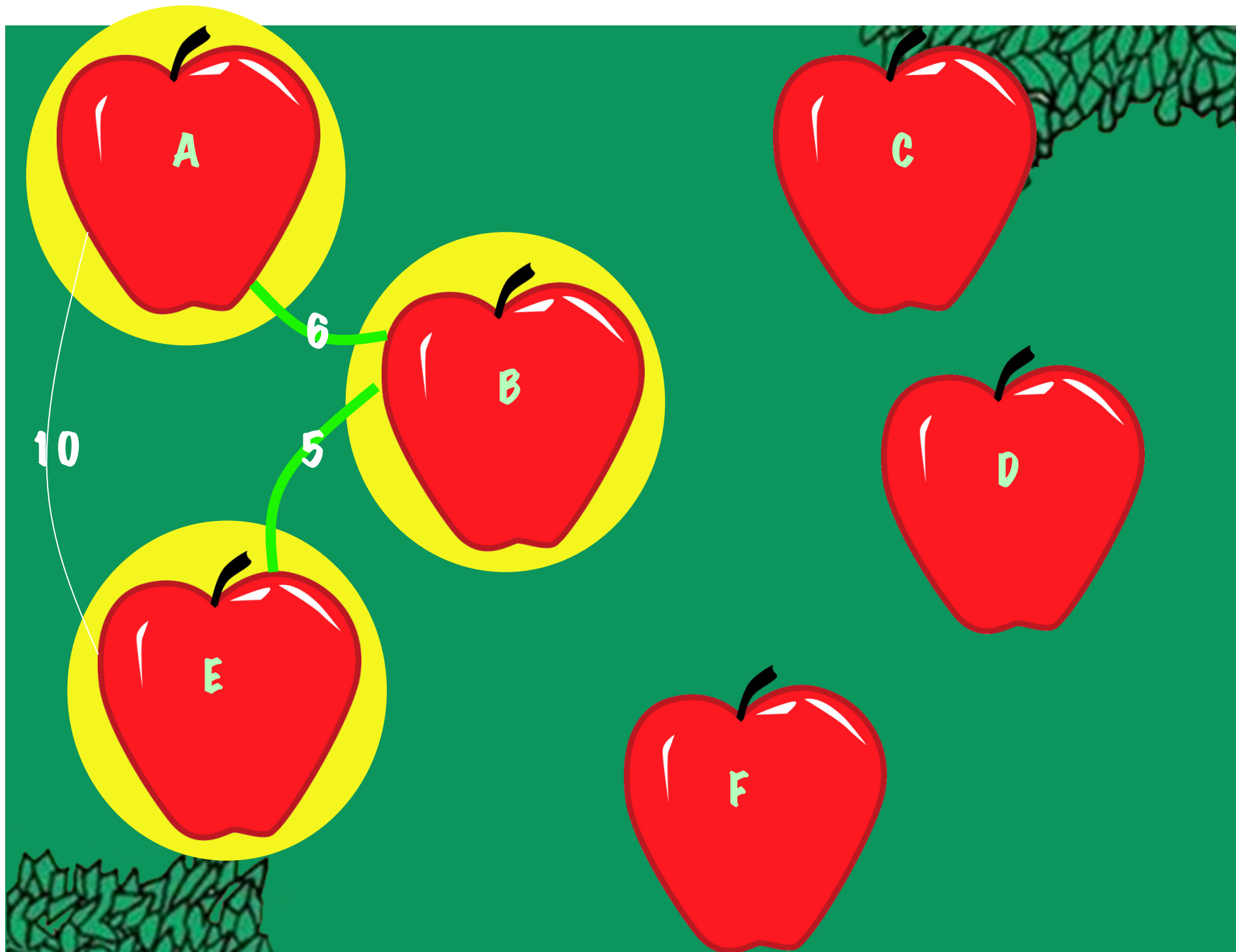


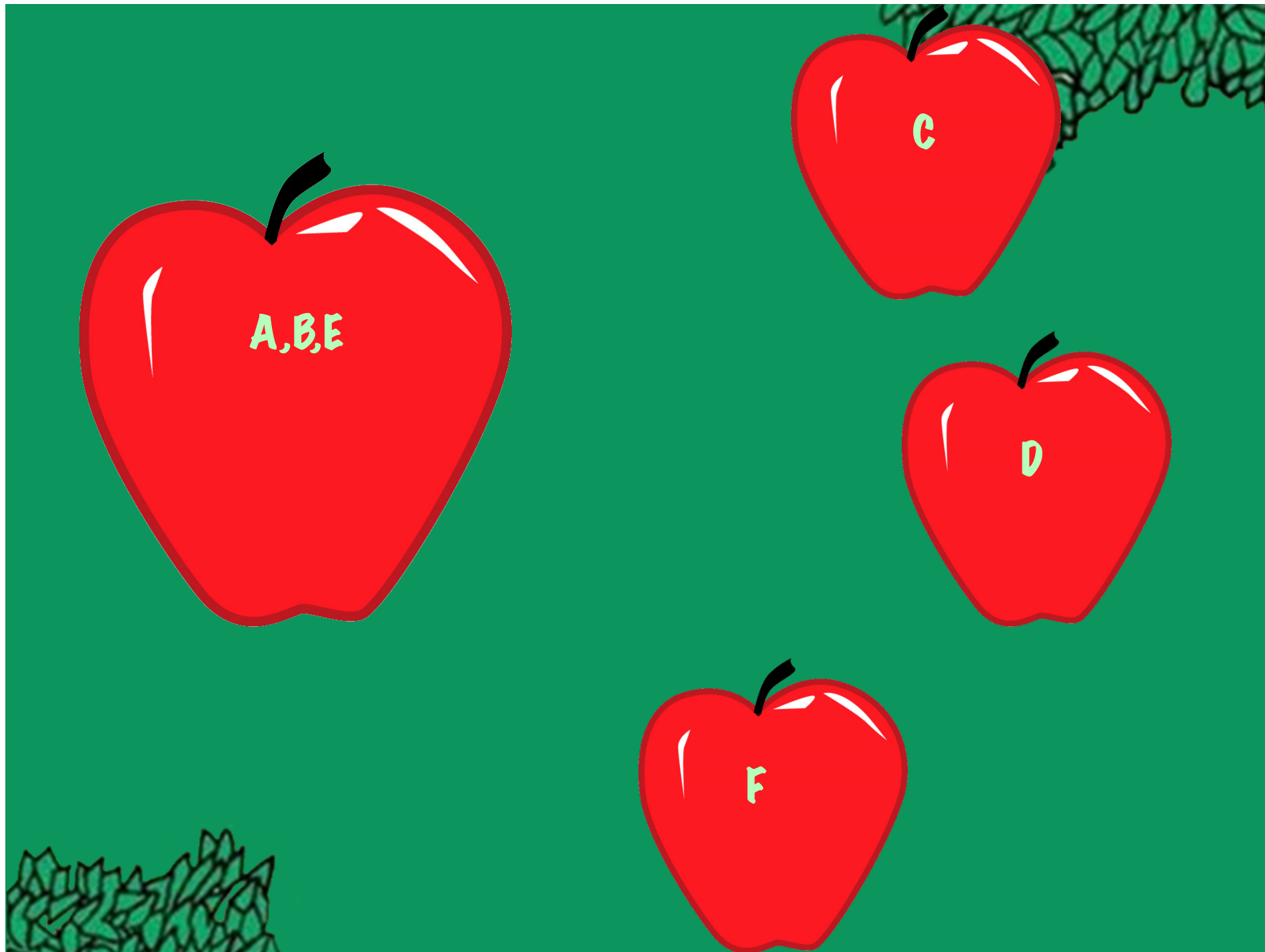


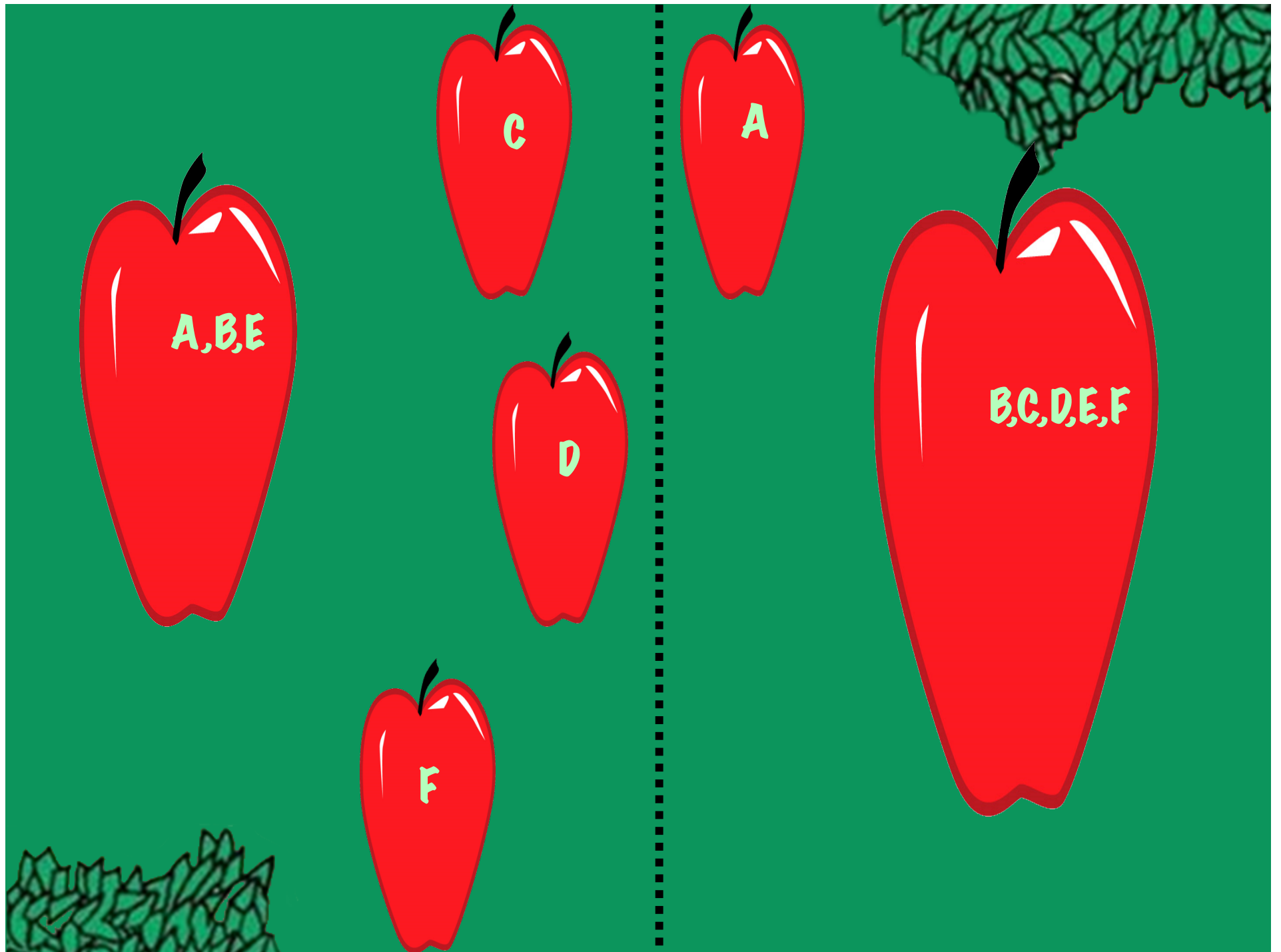


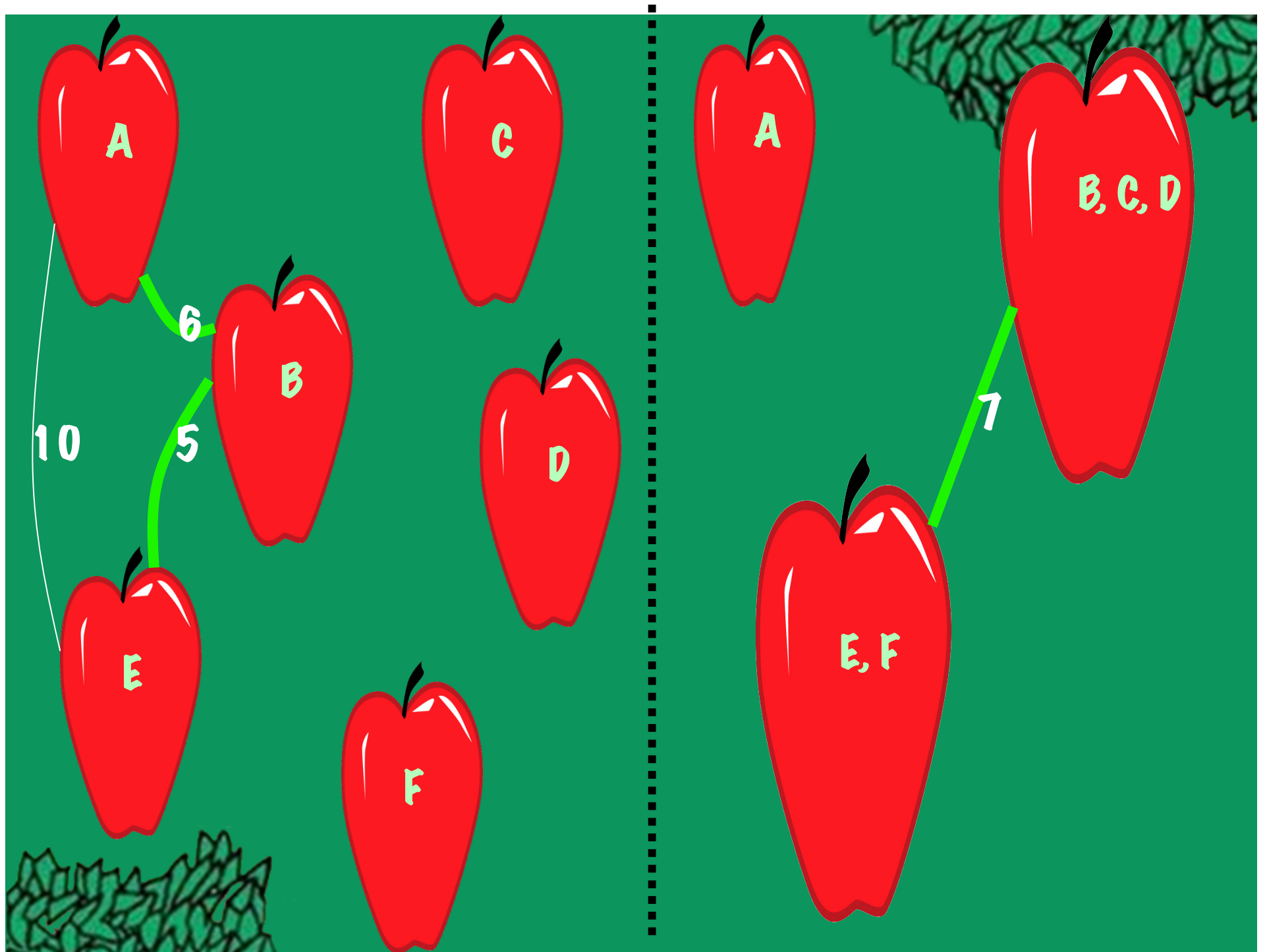


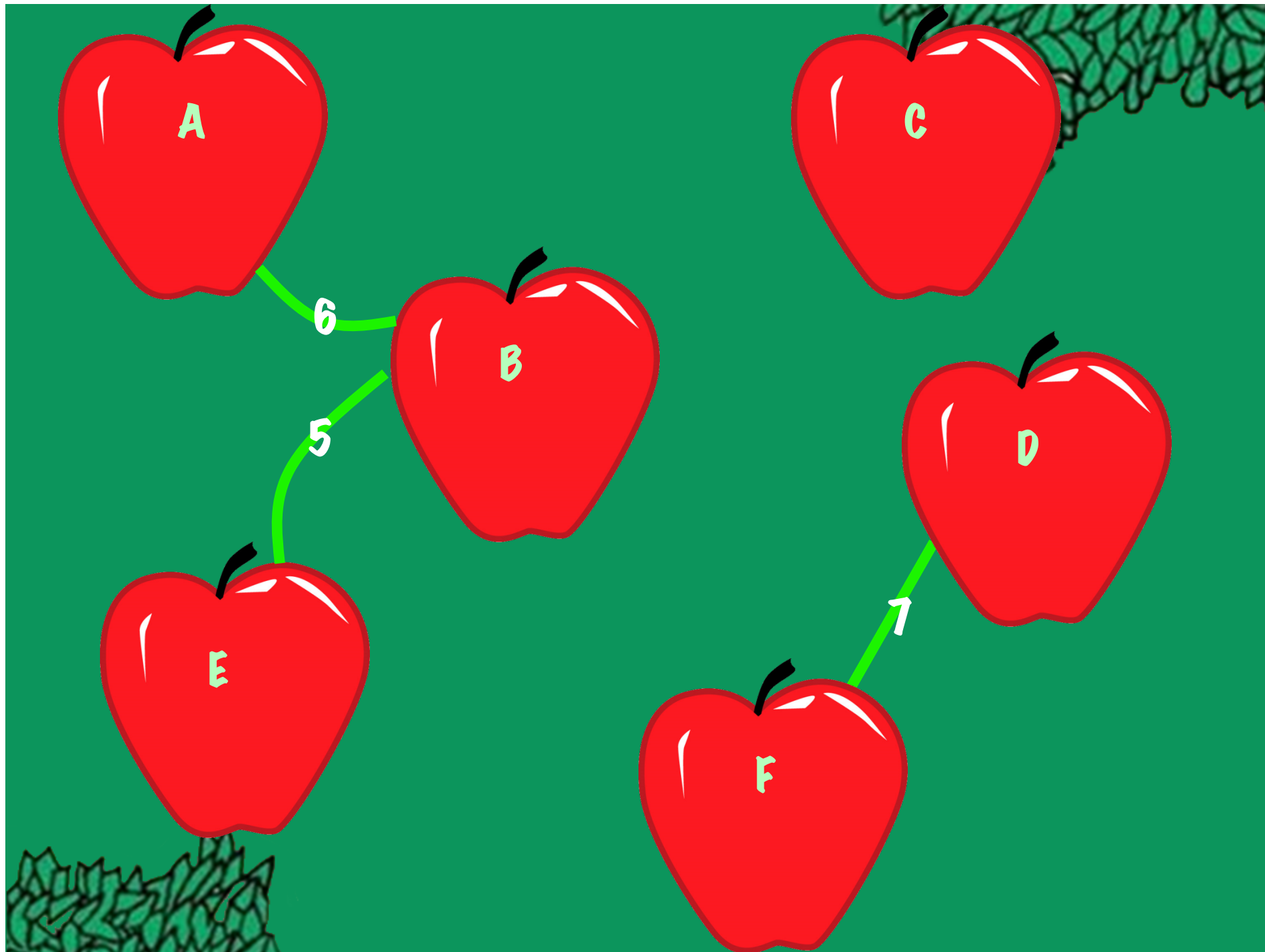


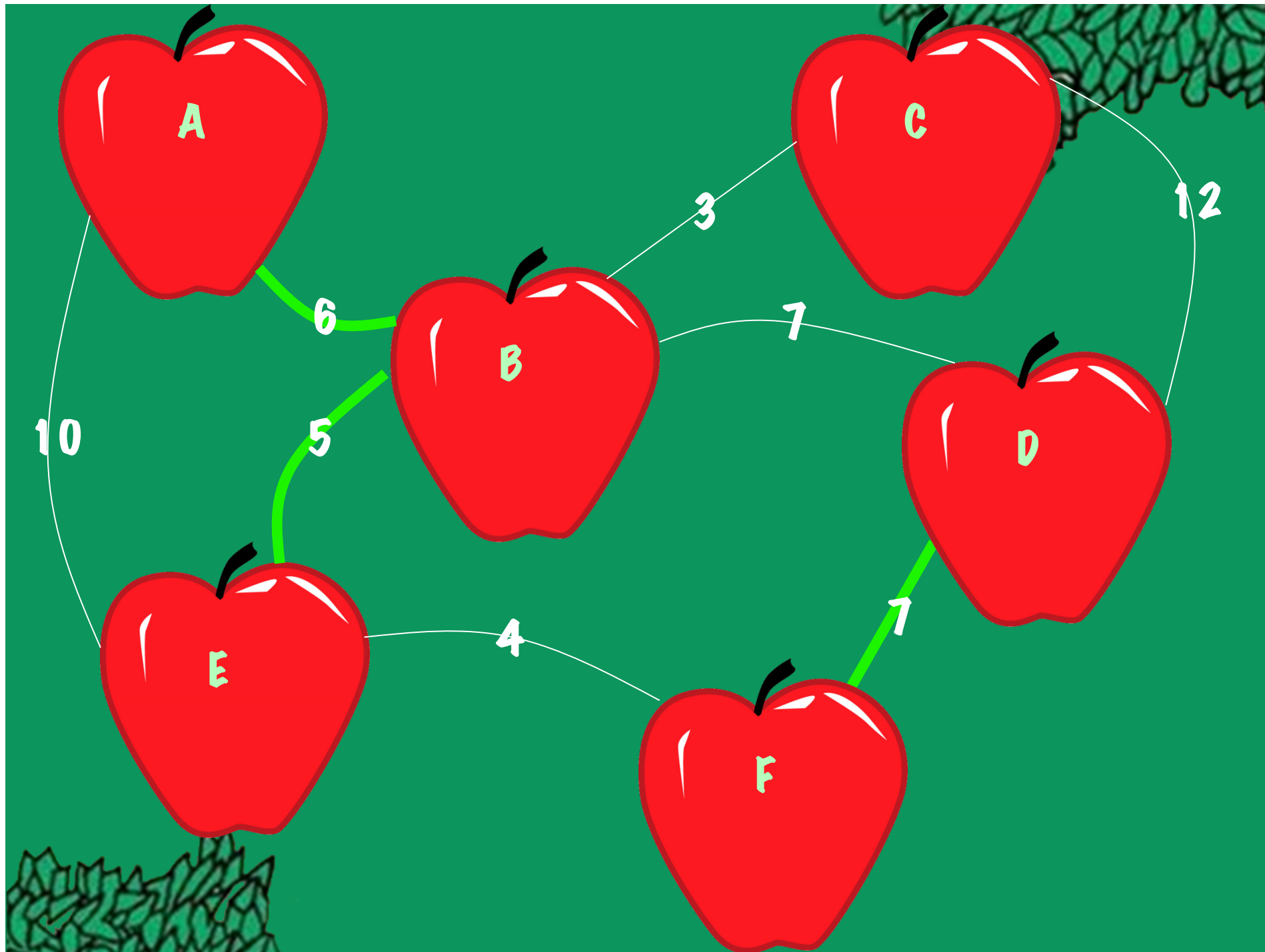


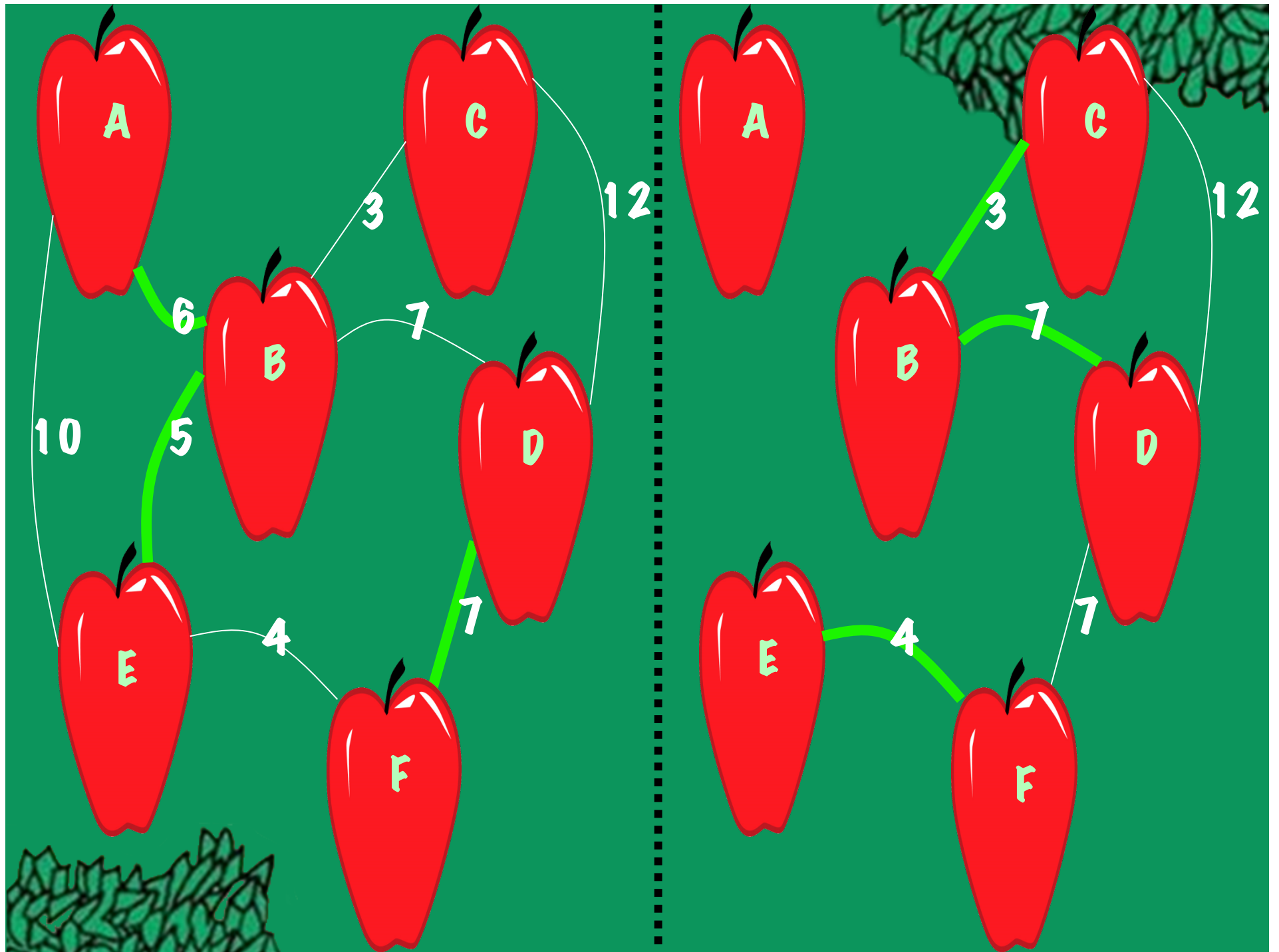


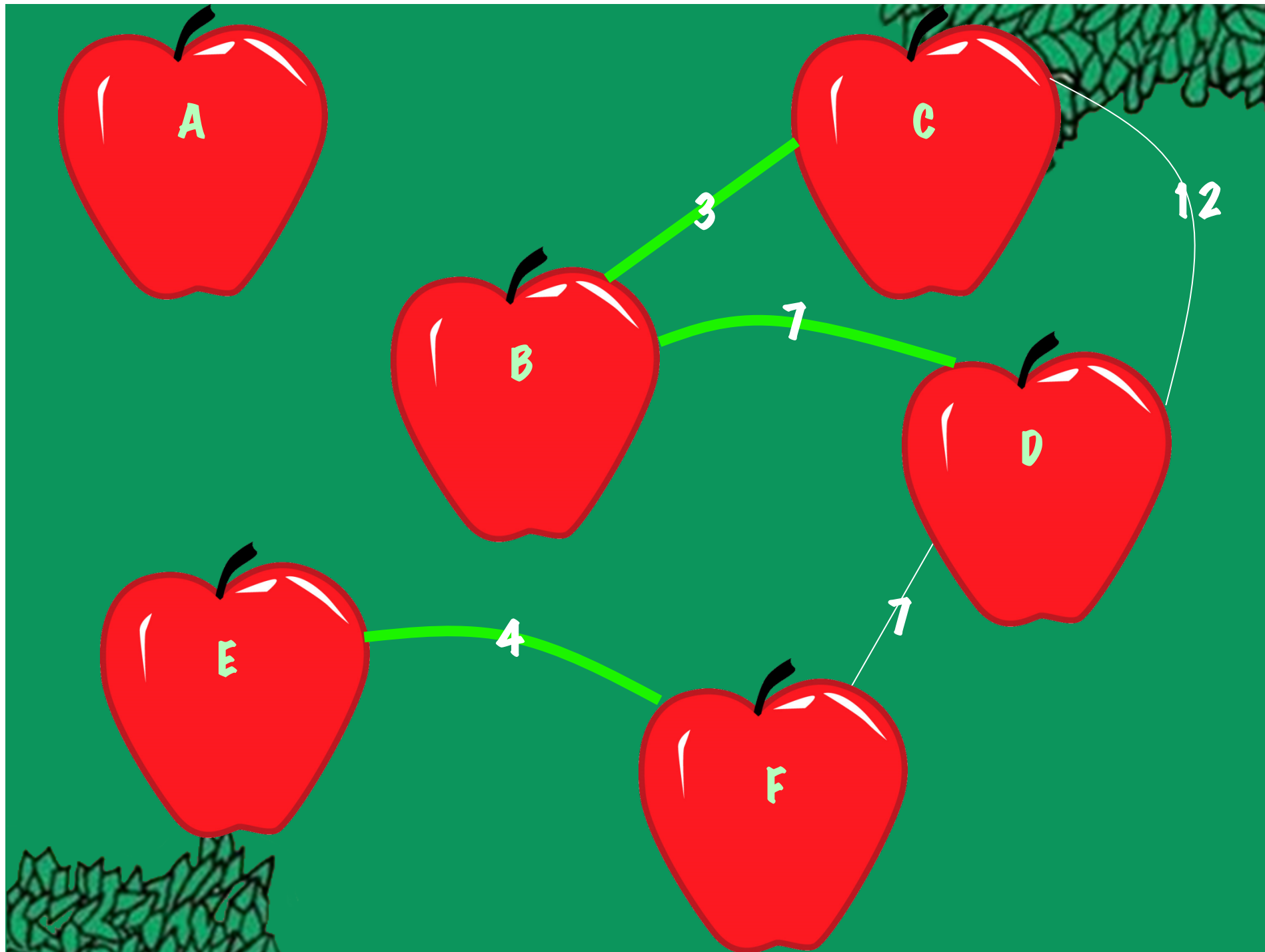


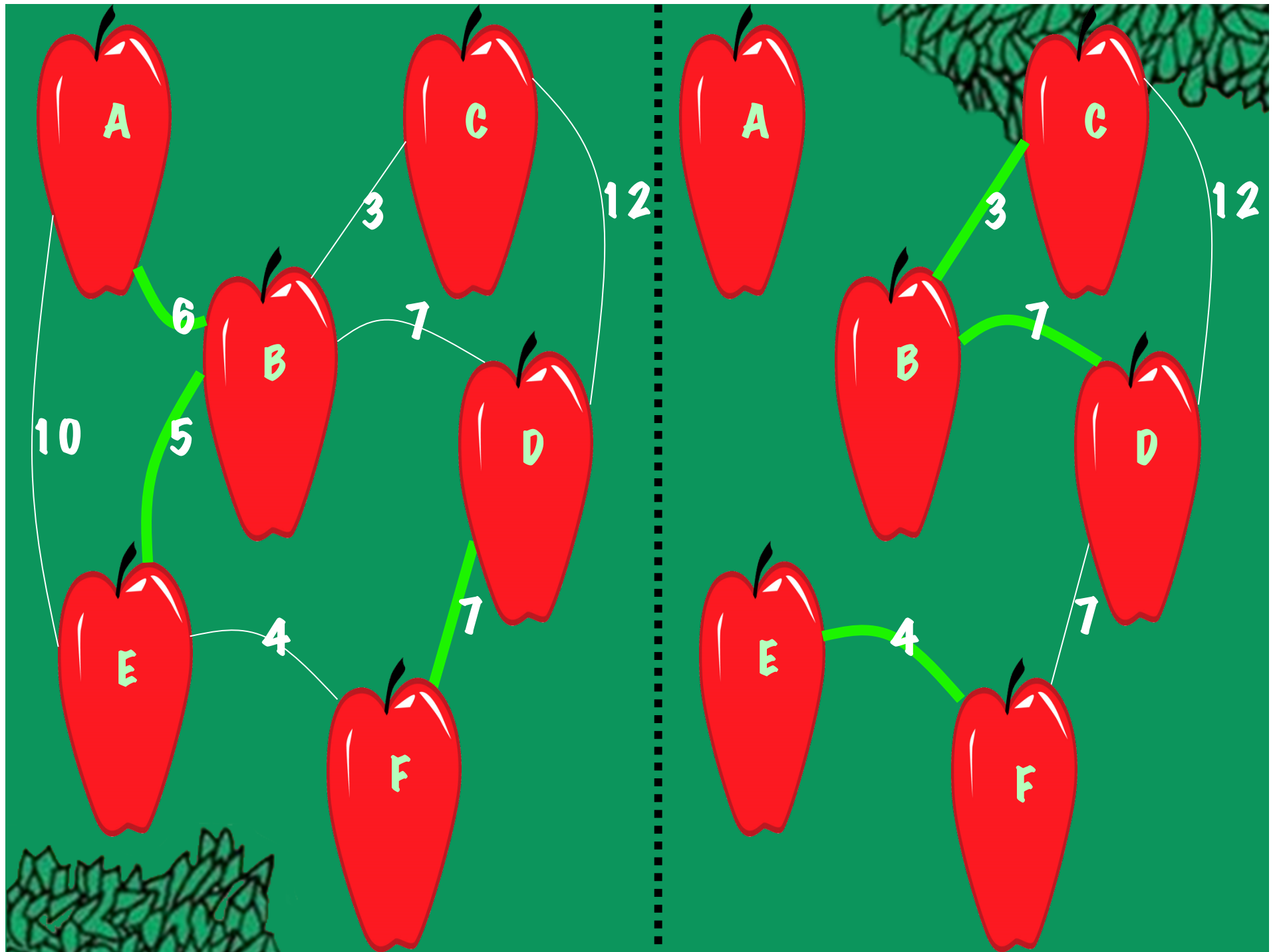


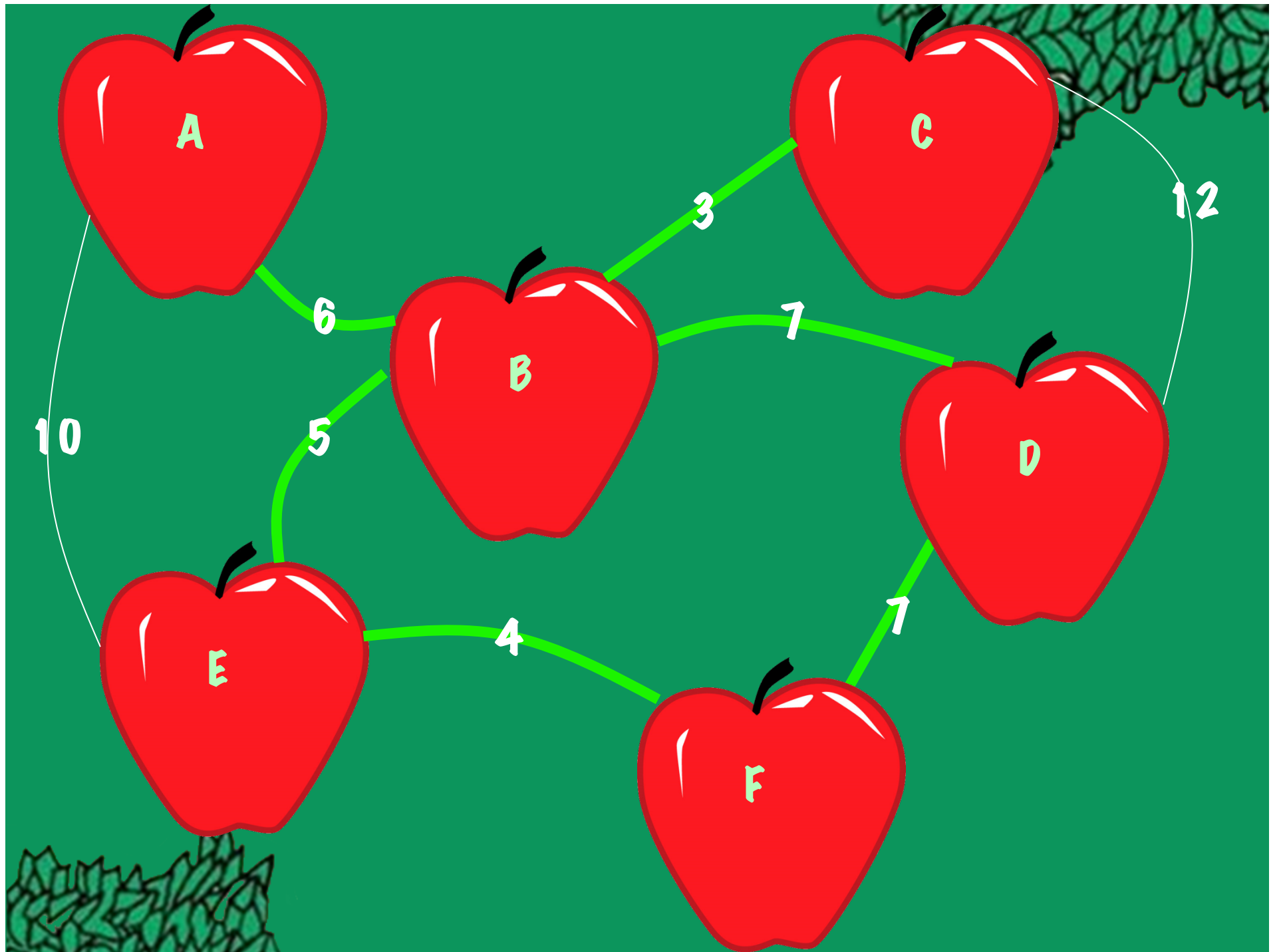


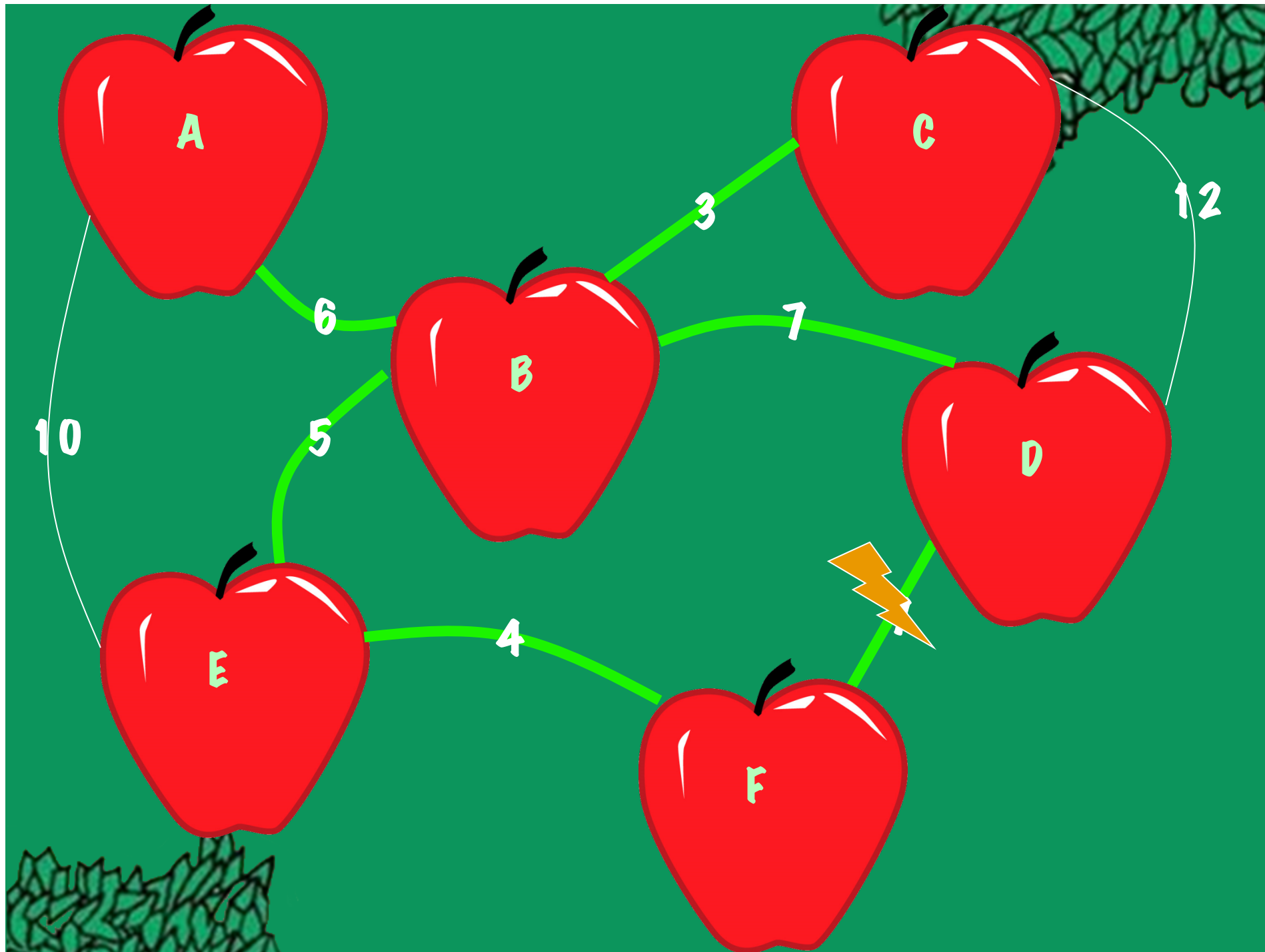


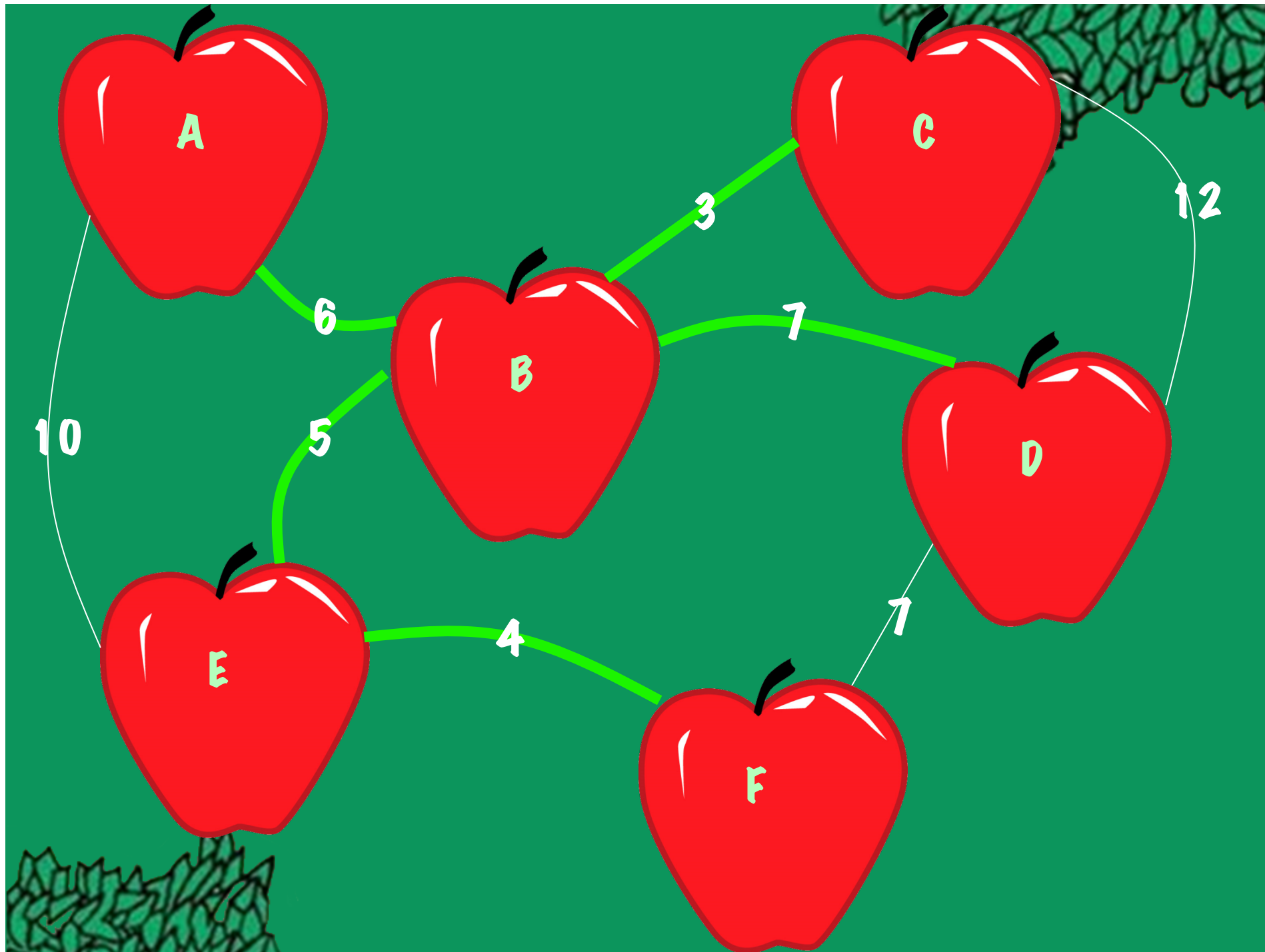


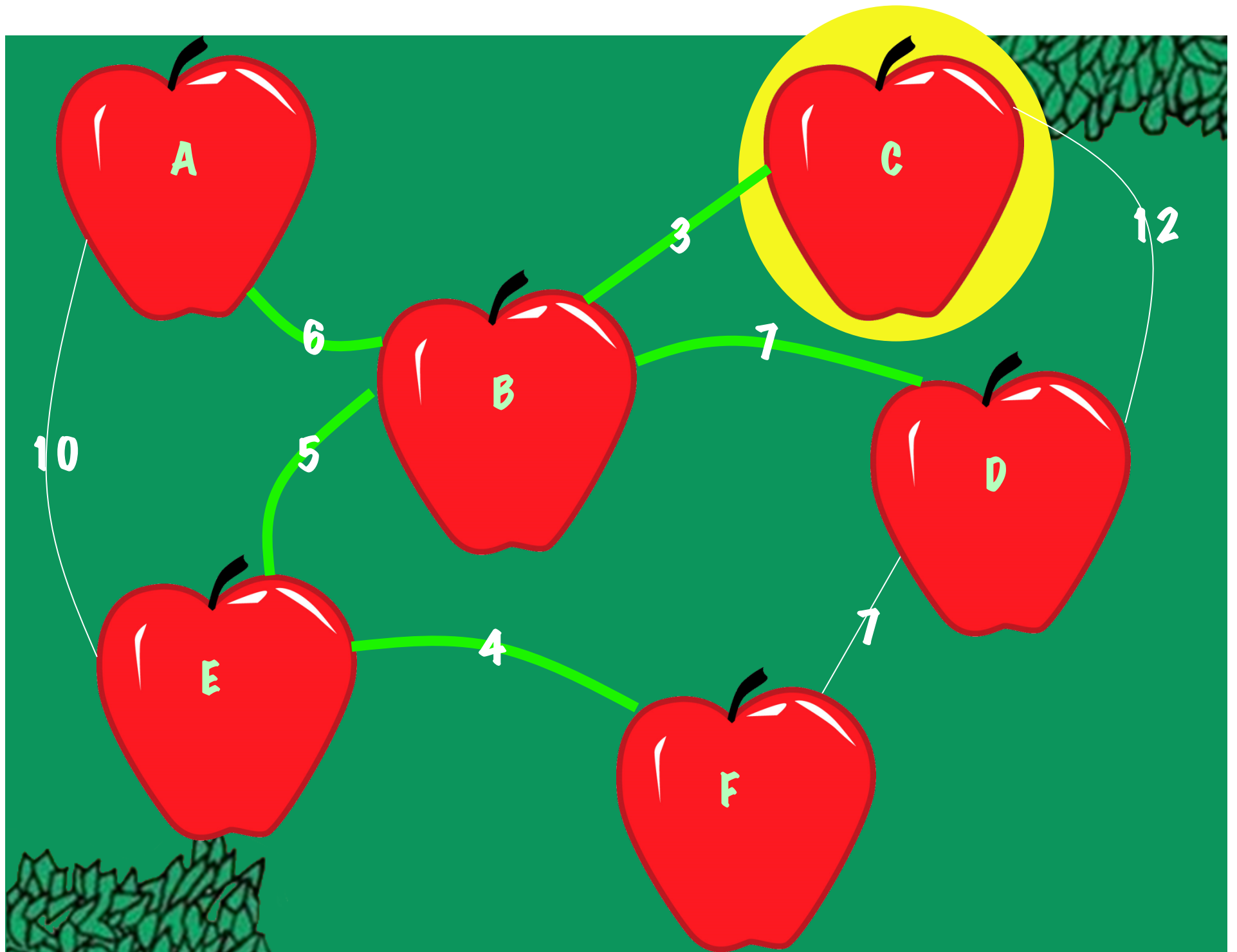


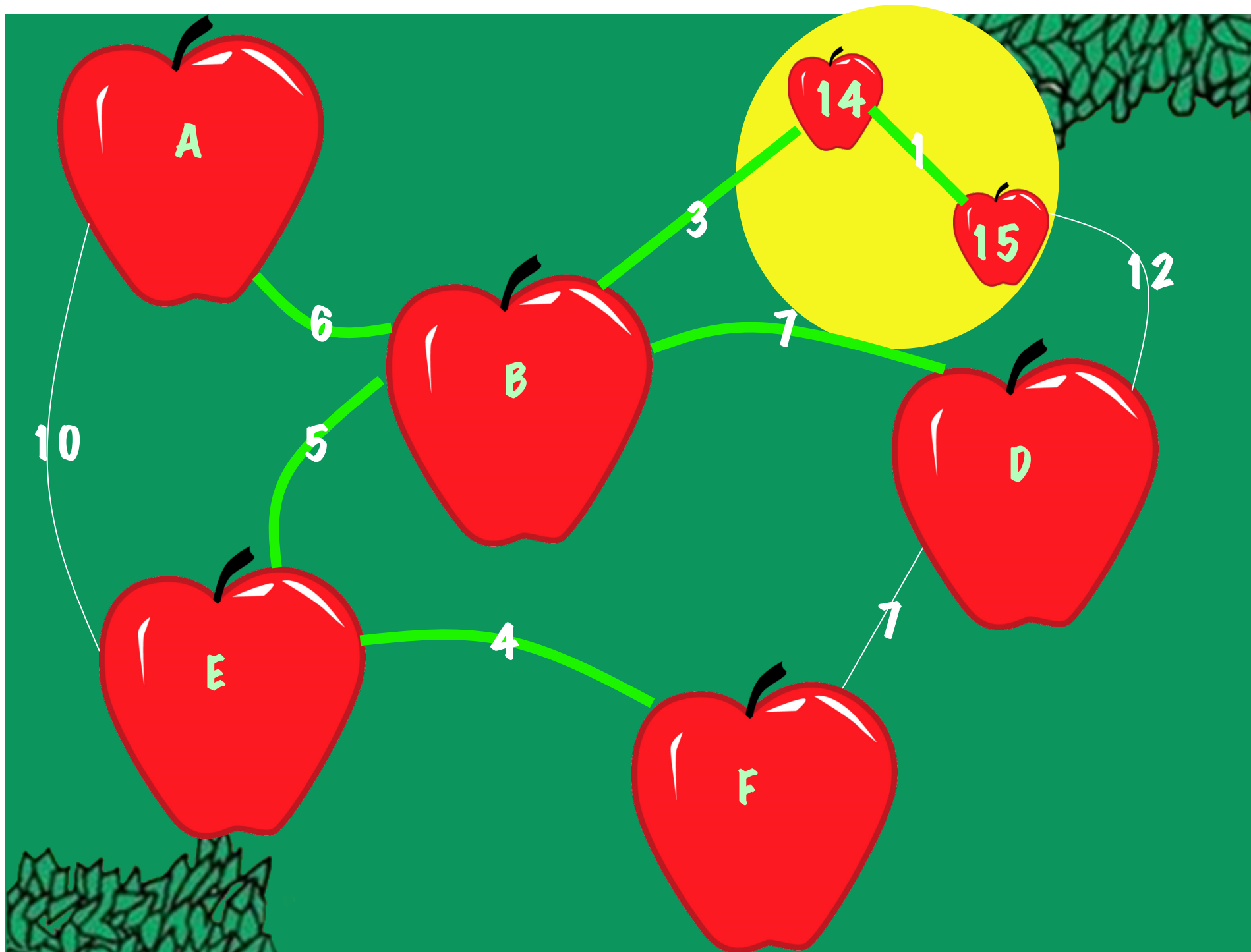


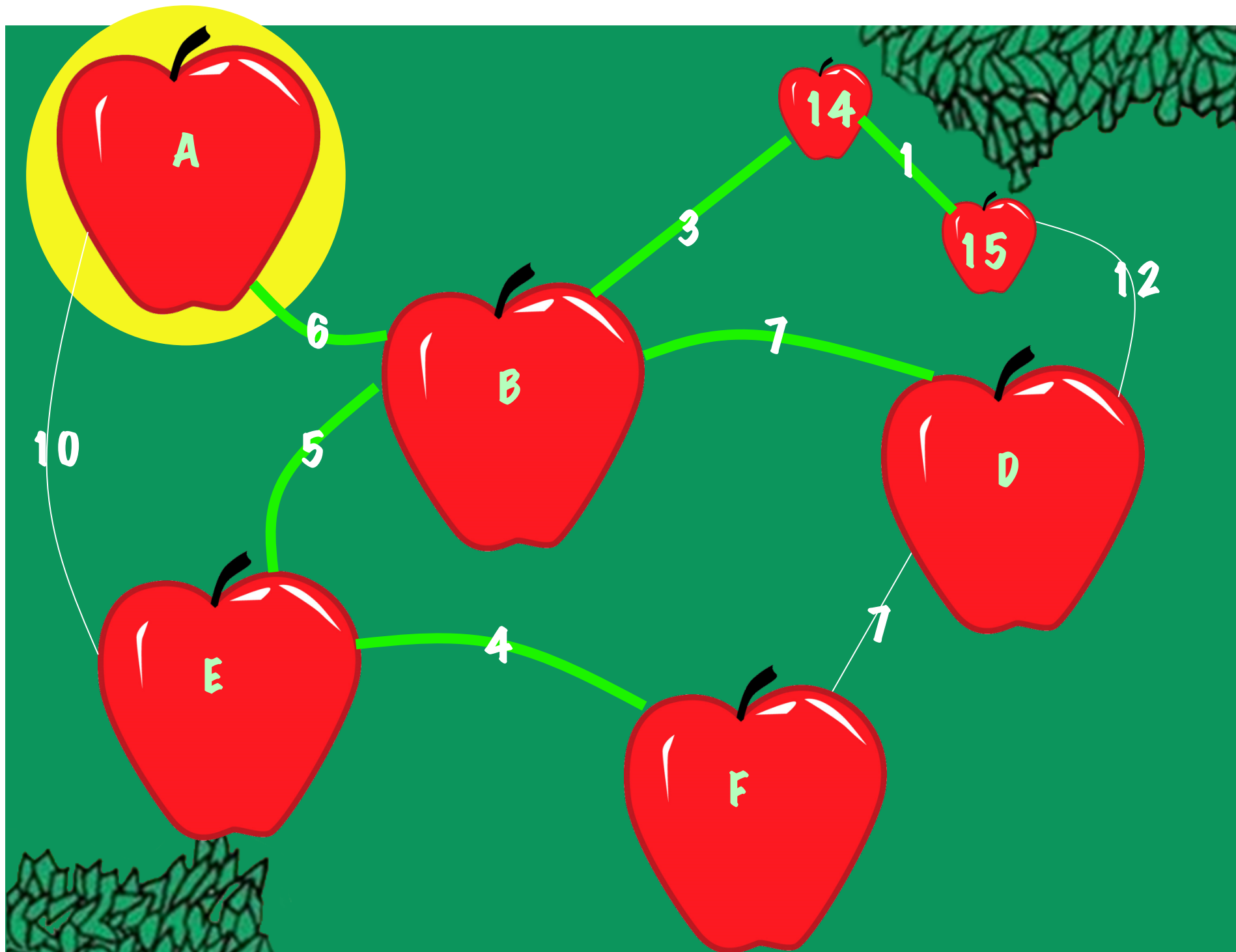


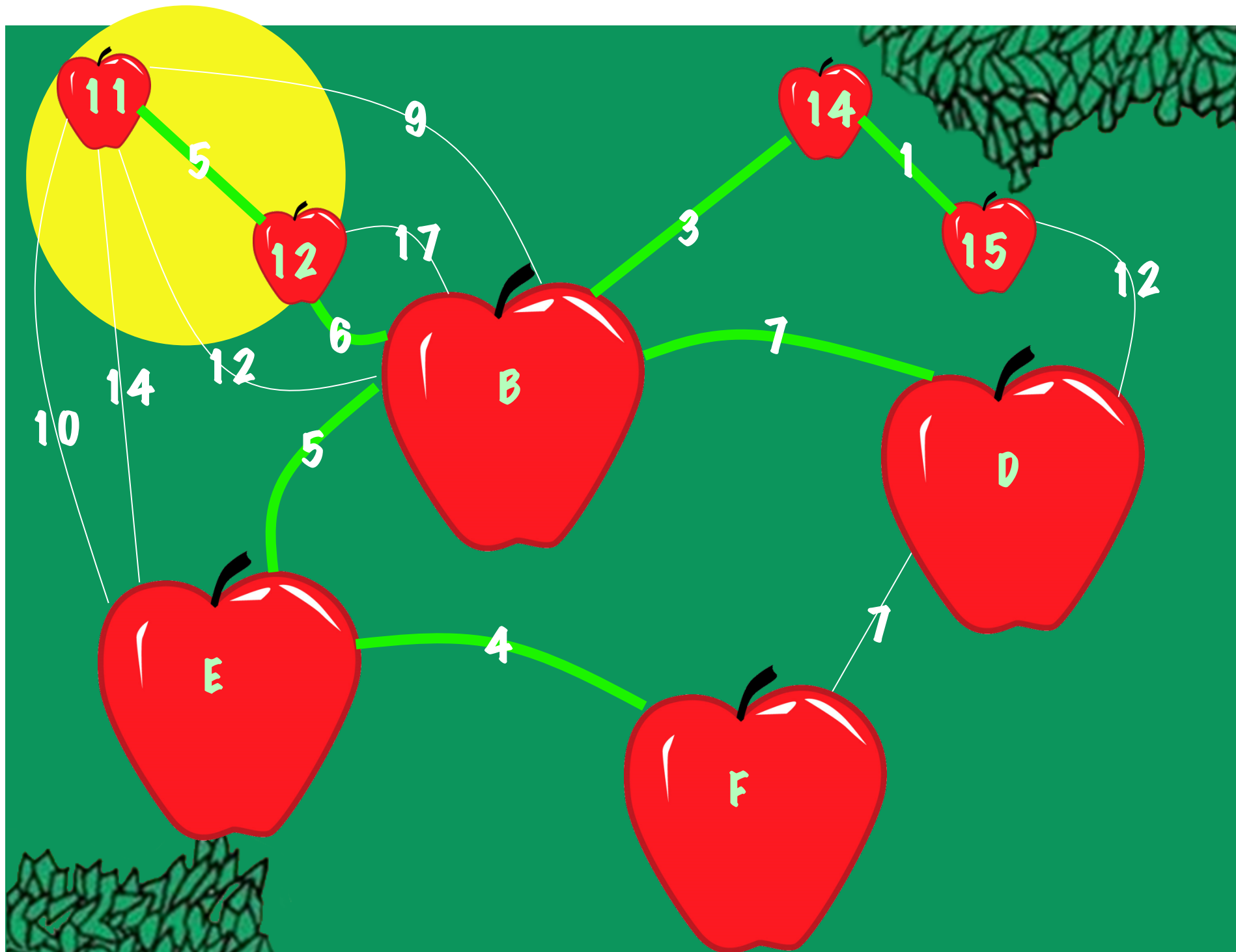


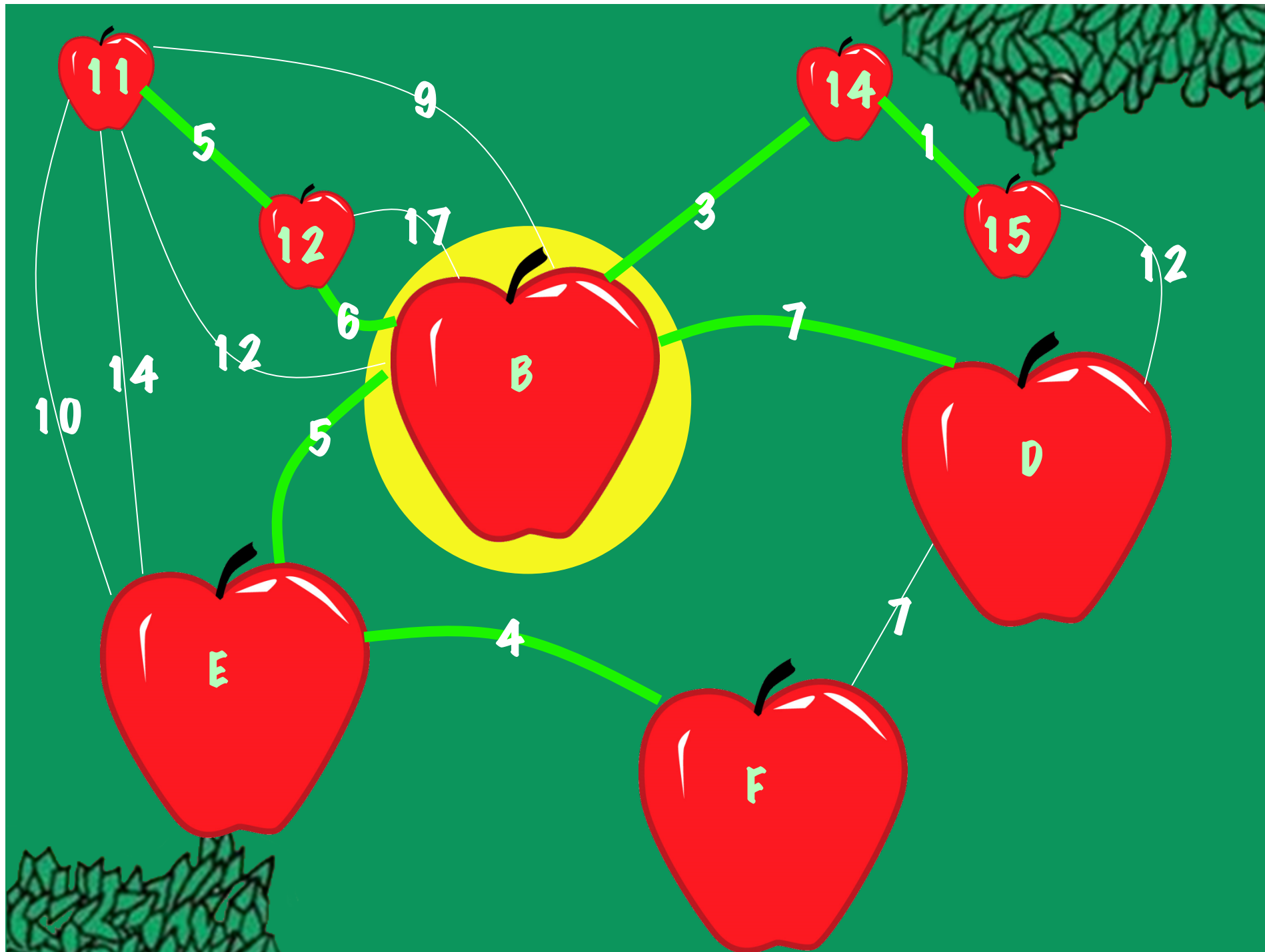


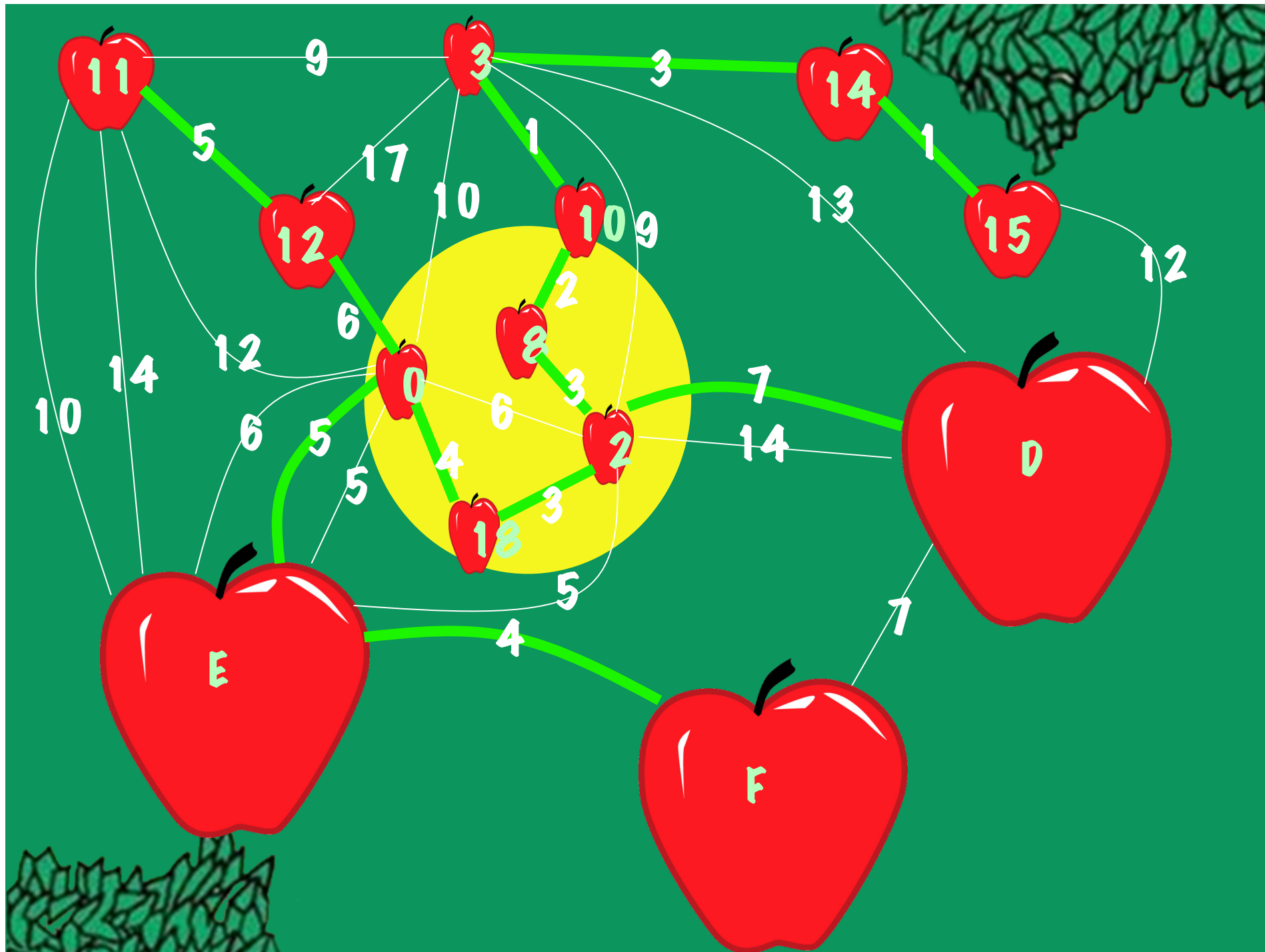


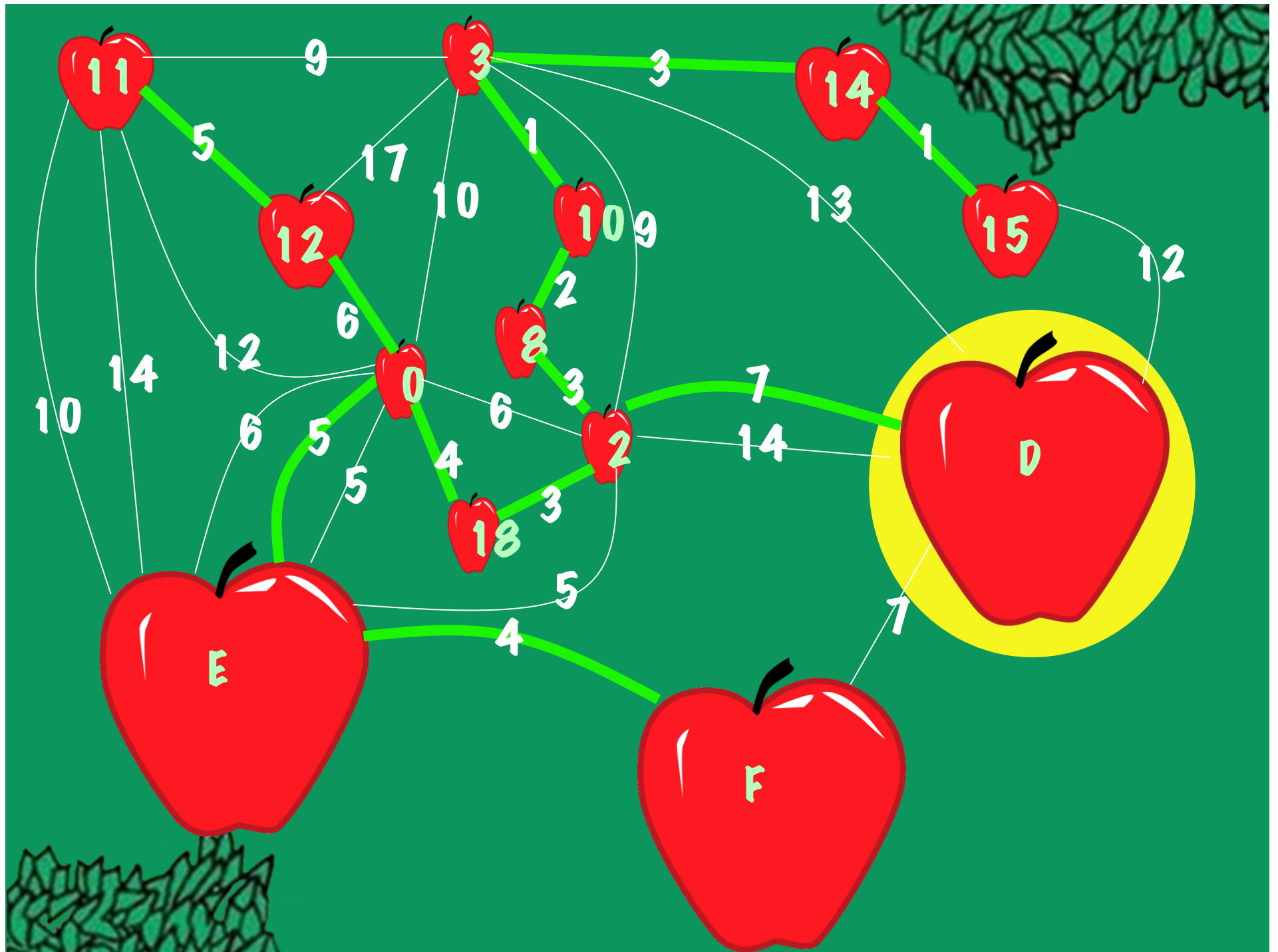


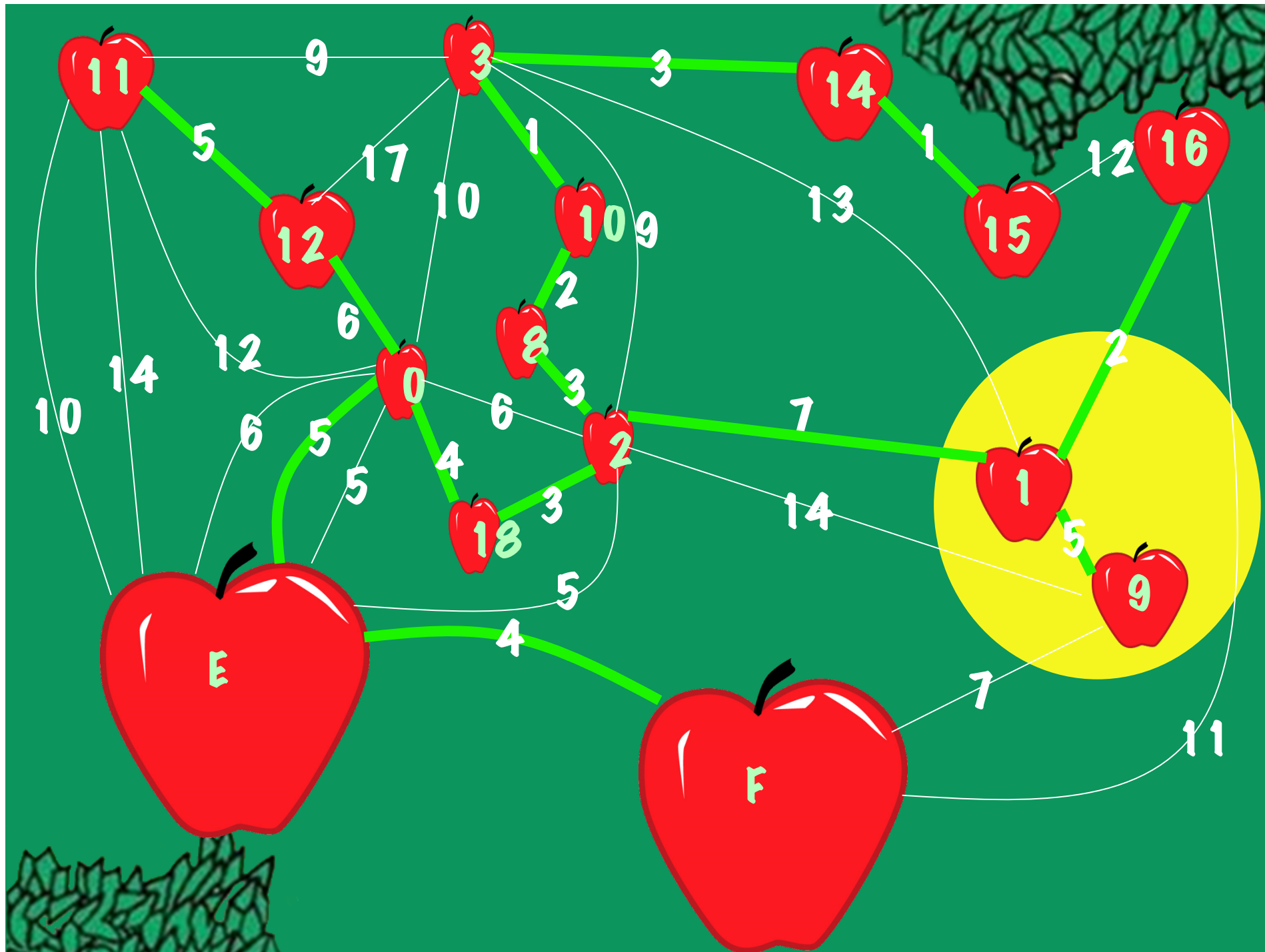


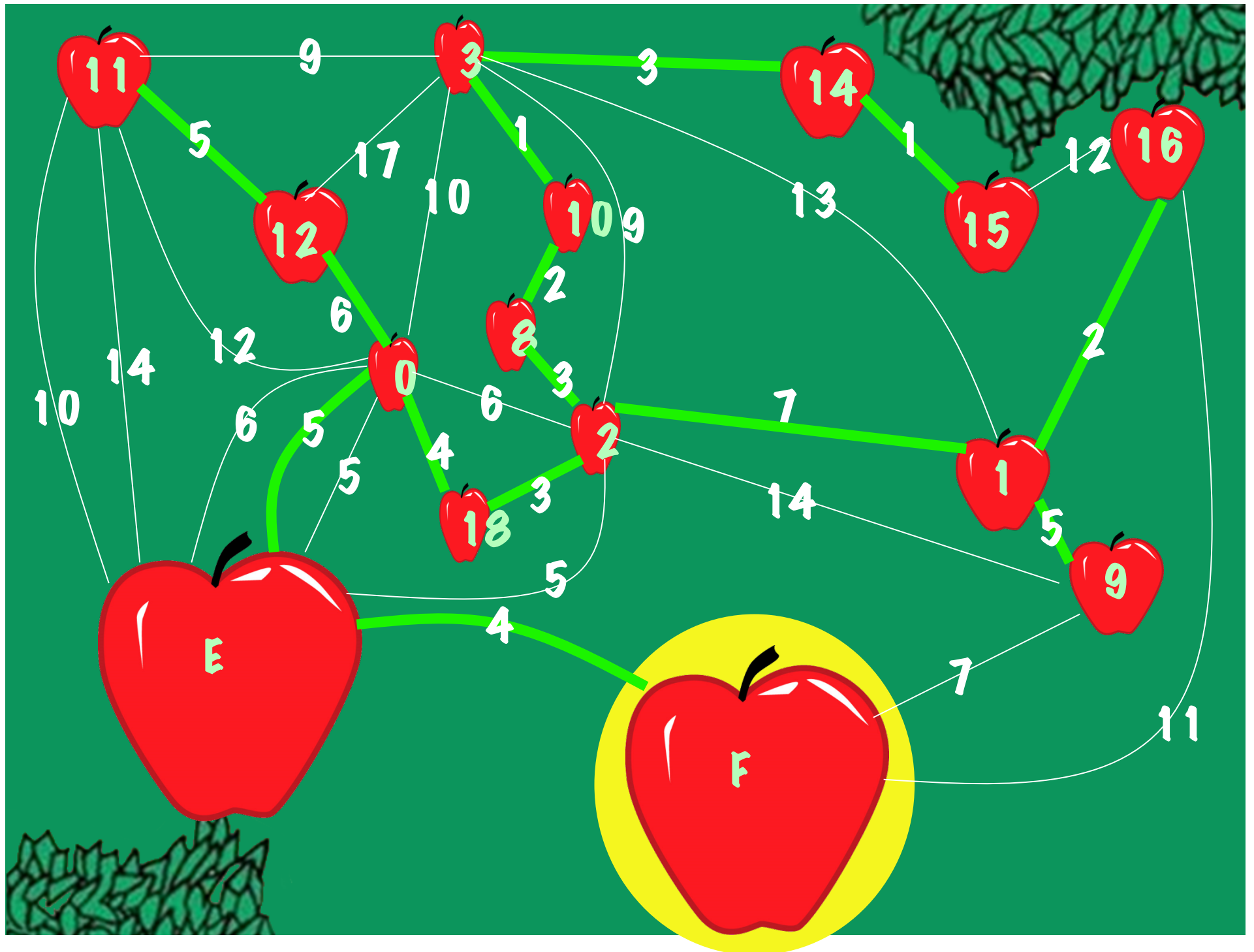


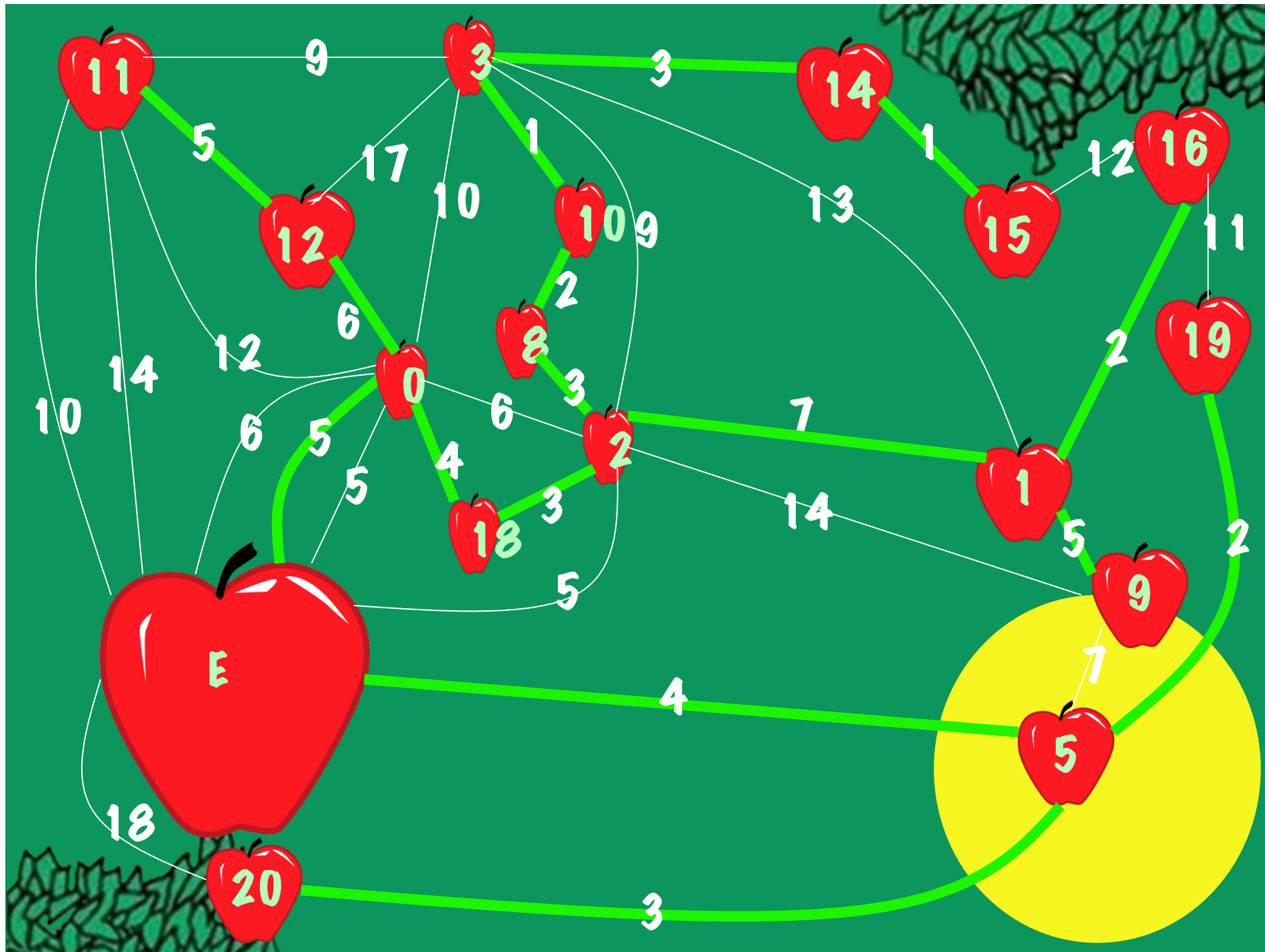


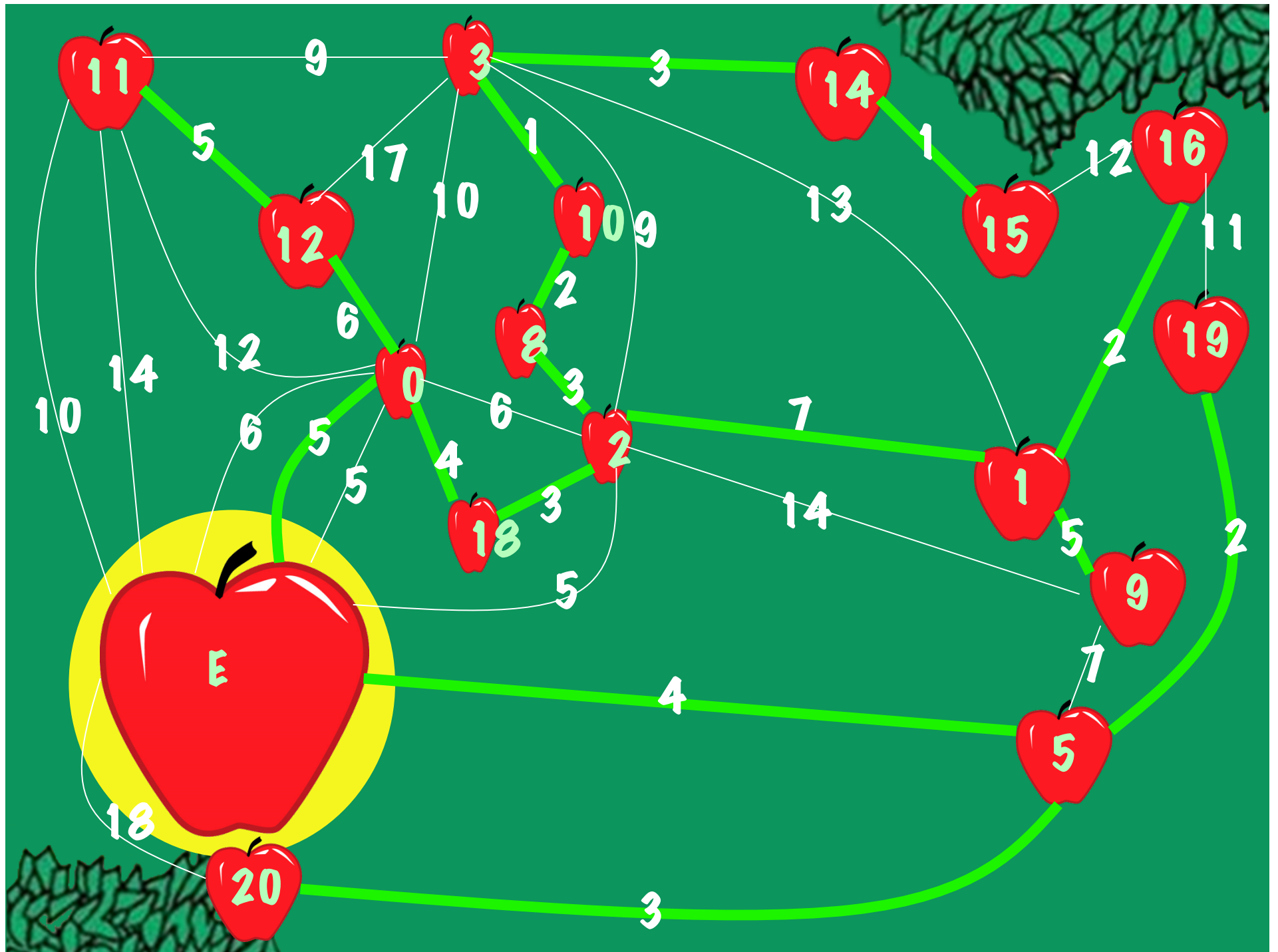


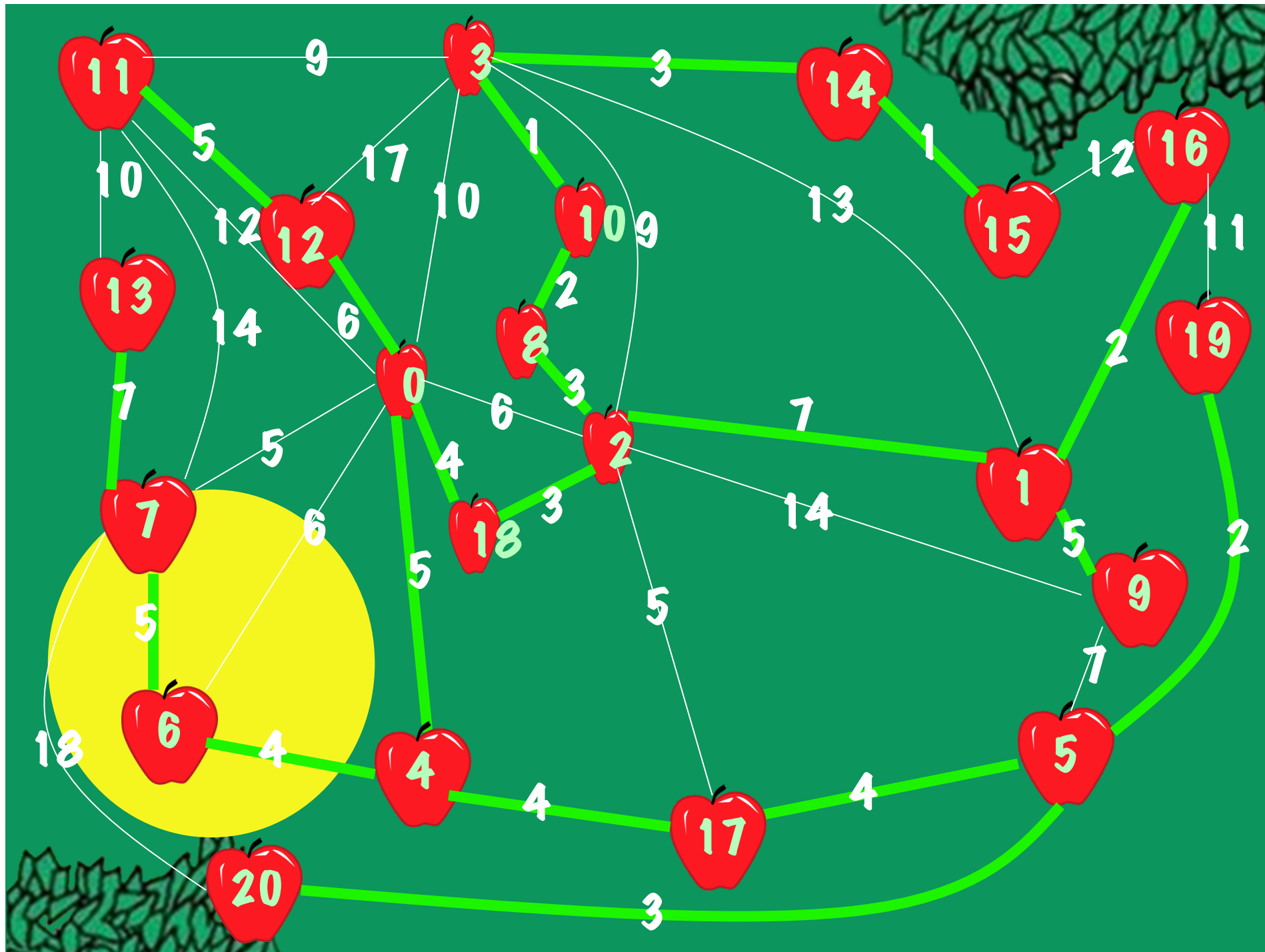


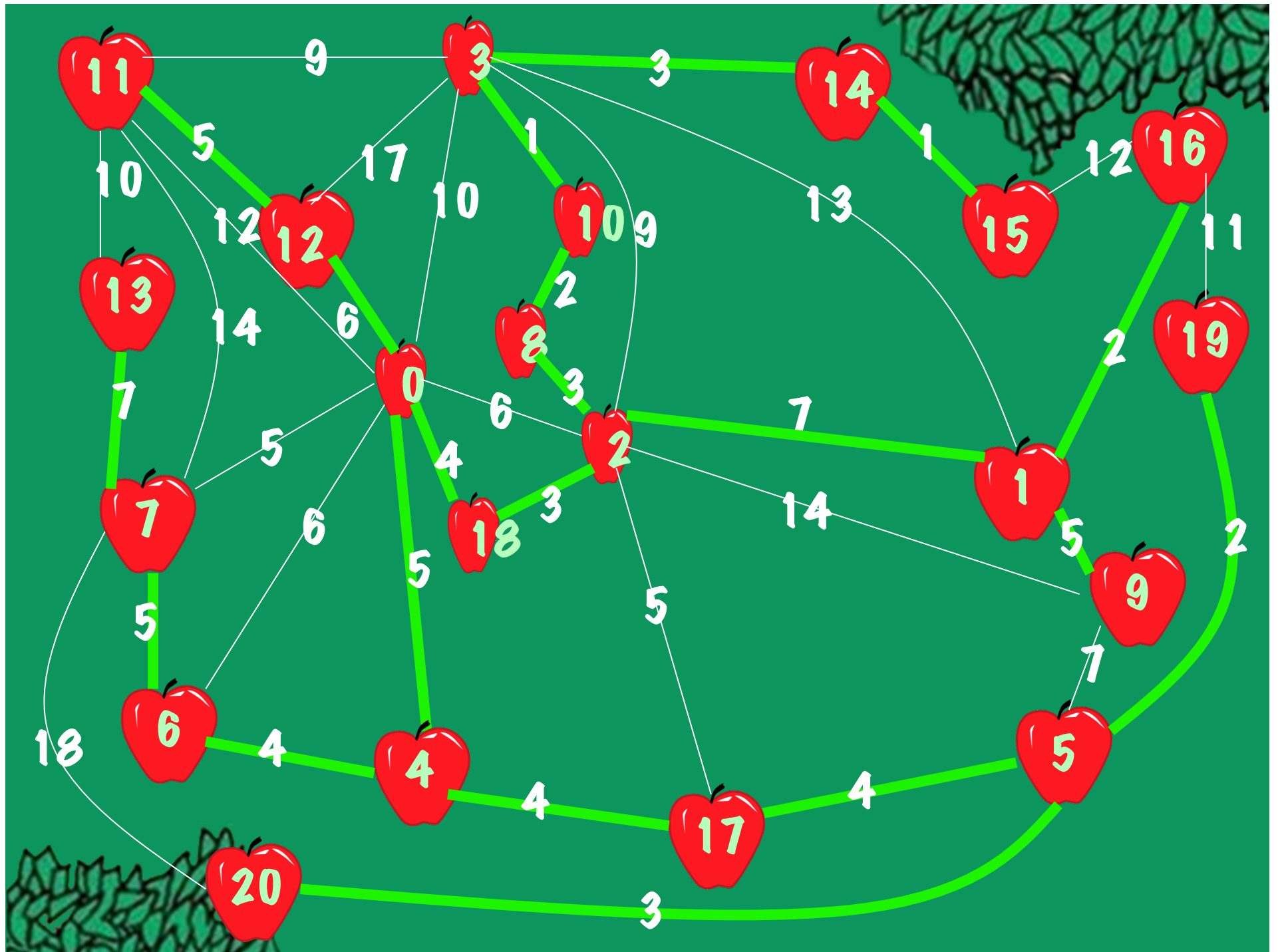


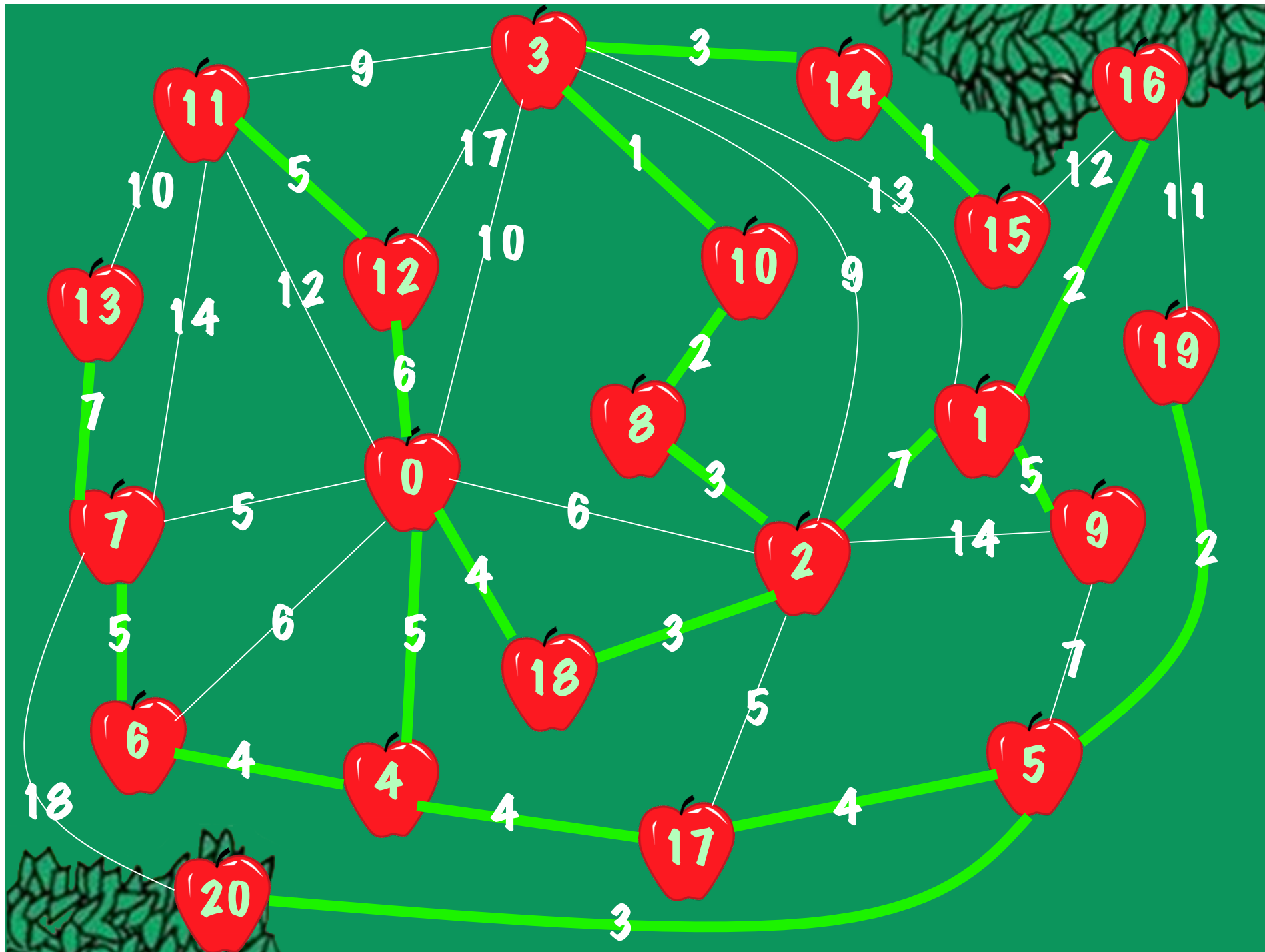












Why Does It Terminate?

- The Boruvka step reduces the size of the graph by a factor of 2
- Terminate when all the nodes are coalesced in all the subgraphs, then build back up

Why Does It Work?

- Proof by induction:
 - **Cut property:** edges contracted in STEP 1 are in the MST.
 - **Cycle property:** edges deleted in STEP 2 not in the MST.
 - Remaining edges of the *original graph's* MST form an MST of the *contracted graph*.
 - By the inductive hypothesis, MST of the remaining graph is correctly determined in recursive call of STEP 3.

How Long Does It Take?

- WORST CASE:
 - Total time spent in steps 1-3, ignoring recursion, is linear in the number of edges.
 - Step 1 is two steps of Buruvka, linear
 - Step 2 is linear using the modified Dixon-Rauch-Tarjan verification algorithm
 - Total running time is bounded by the total number of edges and in all recursive subproblems
 - So how many edges are there in all of the problems?

Worst Case, Continued

- Assume a graph with n vertices and m edges.
 $m \geq n/2$.
- Each invocation generates at most two recursive subproblems.
- Consider the binary tree of recursive subproblems, where the root is the initial problem.
- First subproblem is the left child (step 2), second is right child (step 3).

Worst Case, Continued

- At depth d , the tree of subproblems has at most 2^d nodes, each a problem on a graph of at most $n/4^d$ vertices.
- The the depth of the tree is at most $\log_v n$, and there are at most $2n$ vertices total in the original problem and all subproblems.
- A subproblem at depth d contains at most $(n/4^d)^2/2$ edges. The total number of edges in all subproblems at an recursive depth d is therefore, at most, m .

Worst Case, Finally!

- Since there are $O(\log n)$ depths in the tree, the total number of edges in all recursive subproblems is $O(m \log n)$
- This is the same as Baruvka's original algorithm.

Best Case

- The expected running time is $O(m)$
- The expected number of edges in the left subproblem is at most half the expected number of edges in its parent.
- $m + (m/2) + (m/4) \dots = 2m$
- The sum of expected number of edges in every subproblem in the left side of the tree is therefore at most $2m$.

Best Case, continued

- The total number of edges is bounded by $2m$ and the expected total number of edges in all right subproblems
- The number of edges in the right subproblem is at most twice the number of vertices in the subproblem.

Best Case, continued

- The total number of vertices in all right subproblems is at most $n/2$.
- So the expected number of edges in the original problem and all subproblems is at most $2m+n$
- The algorithm runs in $O(m)$ time with probability $1 - (\exp(-\Omega(m)))$



THE END

...questions?

