

R User Group Singapore

Meetup

Network Analysis with tidygraph and ggraph in R

Naman Agrawal
National University of Singapore

2024-11-29

Table of Contents

1. Introduction to Network Analysis and Graphs
2. Basics of tidygraph
3. Visualization with ggraph
4. An Example
5. Trees in tidygraph (if time)

Introduction to Network Analysis and Graphs

What is Network Analysis?

Network analysis is a method used to study and understand the relationships and interactions between entities within a system.

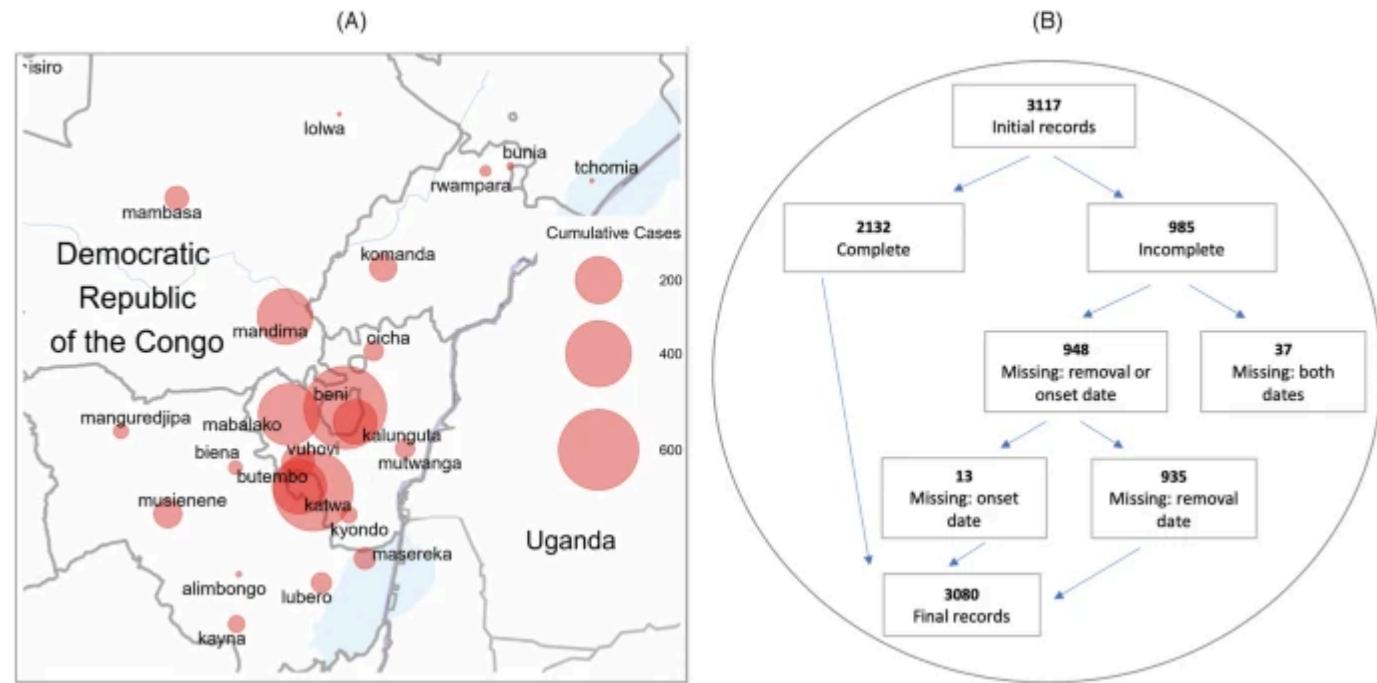
- **Examples:**
 - Understanding social networks (e.g., identifying influencers on social media).
 - Analyzing biological systems such as protein interactions.

How Network Analysis Differs from Tabular Data Analysis

Aspect	Tabular Data	Network Data
Focus	Observations	Relationships
Data Structure	Rows and Columns	Nodes and Edges
Visualization	Scatter plots, histograms	Graph diagrams (node-link)

Importance of Network Analysis?

- Helps uncover hidden patterns.
- Optimizes decision-making in complex systems.
- Identifies influential entities and critical connections.



<https://www.nature.com/articles/s41598-022-09564-4>

Introduction to Graphs

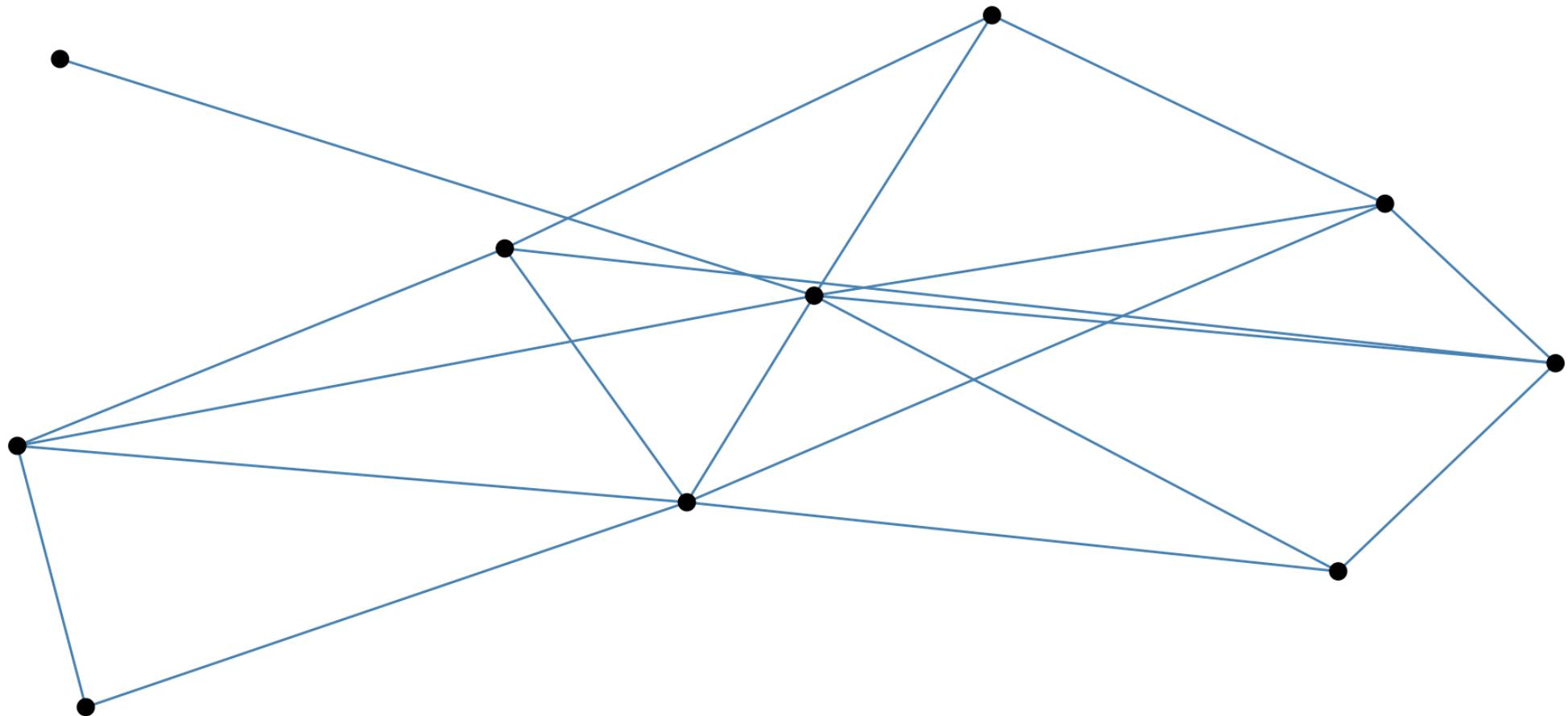
A graph is a mathematical structure used to model pairwise relationships between entities.

Components of a Graph

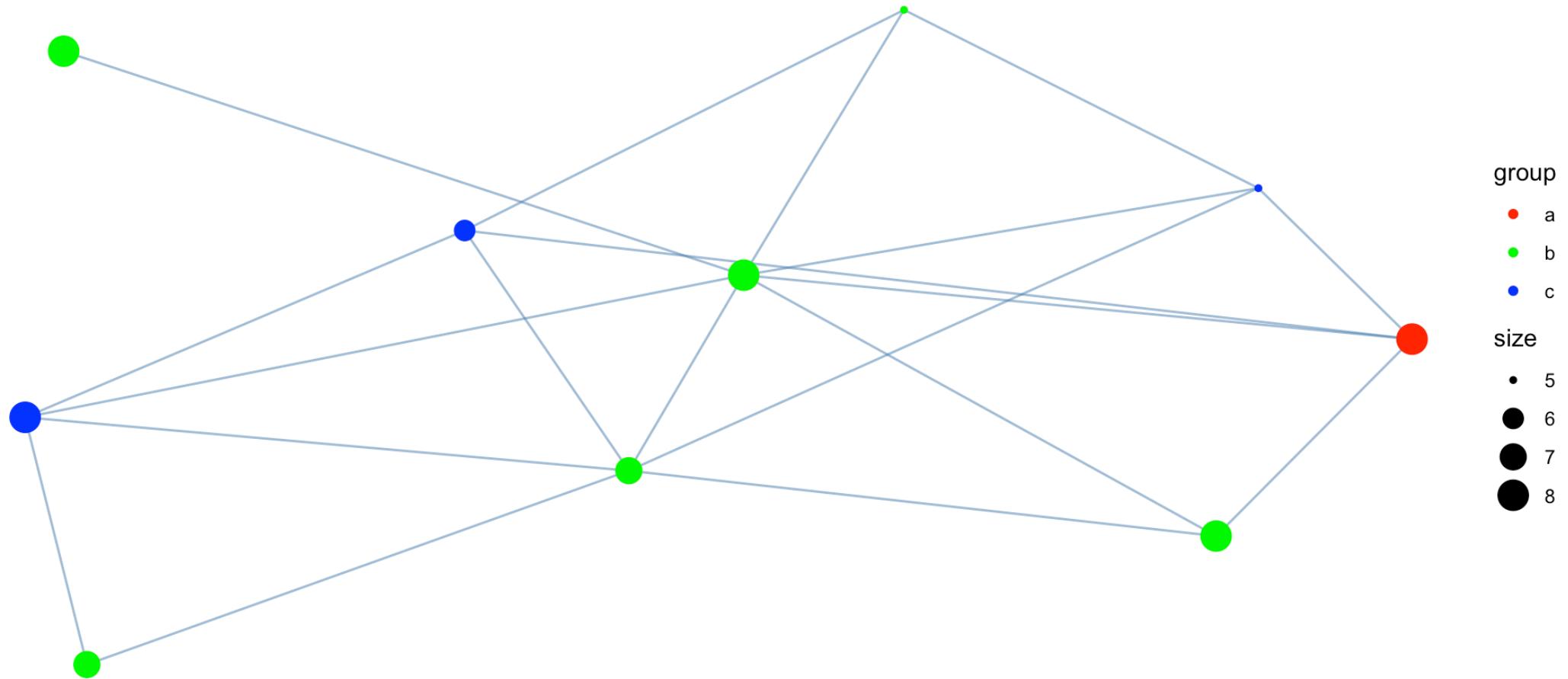
1. **Nodes (Vertices)**: Represent entities.
2. **Edges (Links)**: Represent relationships between nodes.

A graph may also have **attributes**, which are additional data about nodes or edges (e.g., weights, labels).

Introduction to Graphs



Introduction to Graphs



Introduction to Graphs

Graph Types:

- **Directed:** Relationships have a direction (e.g., $A \rightarrow B$).
- **Undirected:** Relationships are bidirectional.
- **Weighted:** Edges have weights (e.g., cost, distance).
- **Unweighted:** Edges do not carry weights.

Ways to Represent Graphs

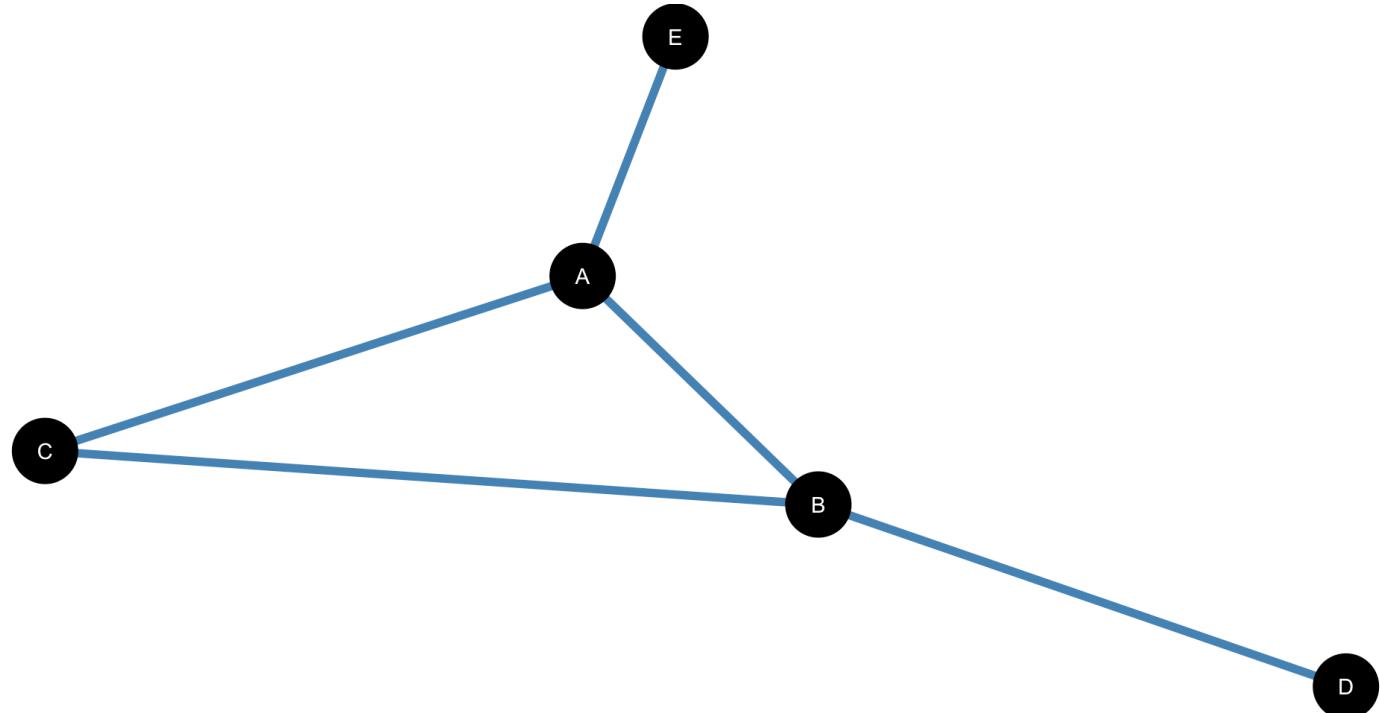
Graph Representations

1. **Edge List:** A list of node pairs connected by edges.
2. **Adjacency Matrix:** Rows and columns represent nodes; cell values indicate edges.
3. **Adjacency List:** A collection of lists where each list corresponds to a node and contains the nodes it is connected to via edges.

and many more ...

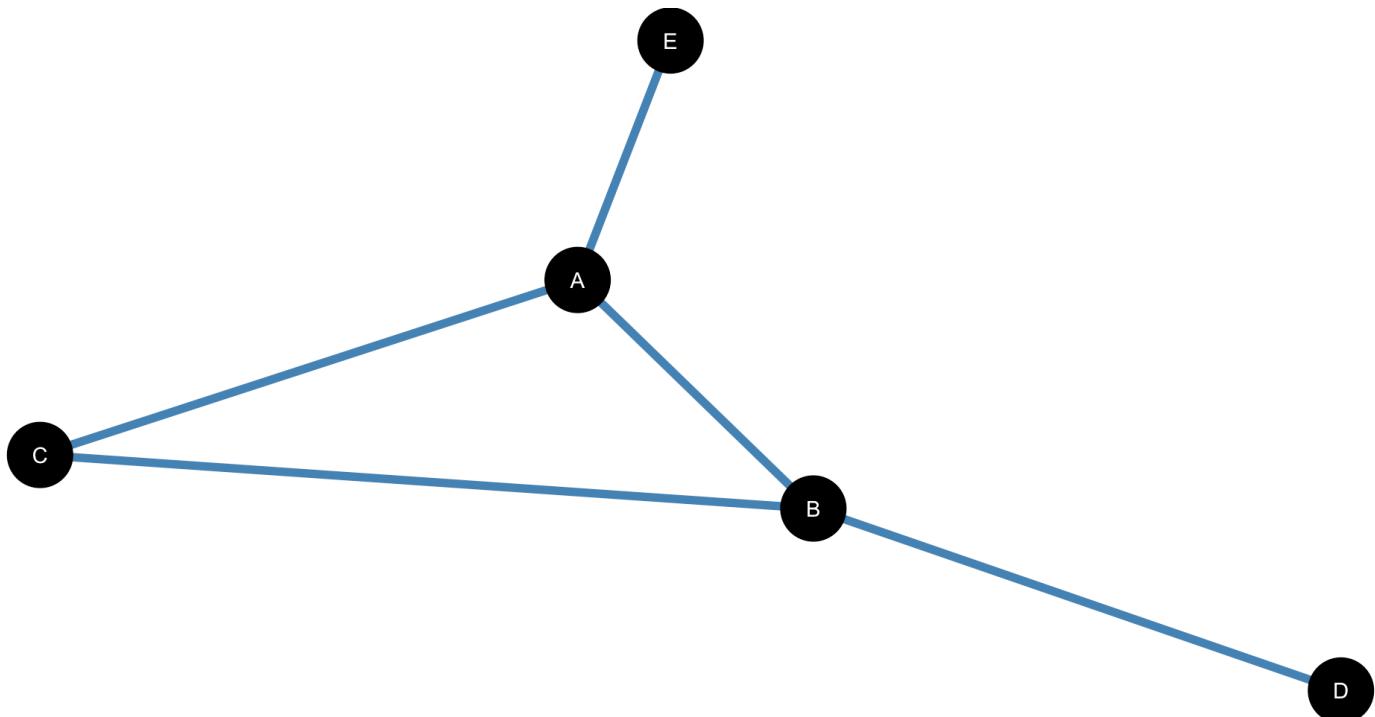
Representations: Edge List

- | | |
|---|--------|
| 1 | (A, B) |
| 2 | (A, C) |
| 3 | (B, D) |
| 4 | (B, C) |
| 5 | (A, E) |



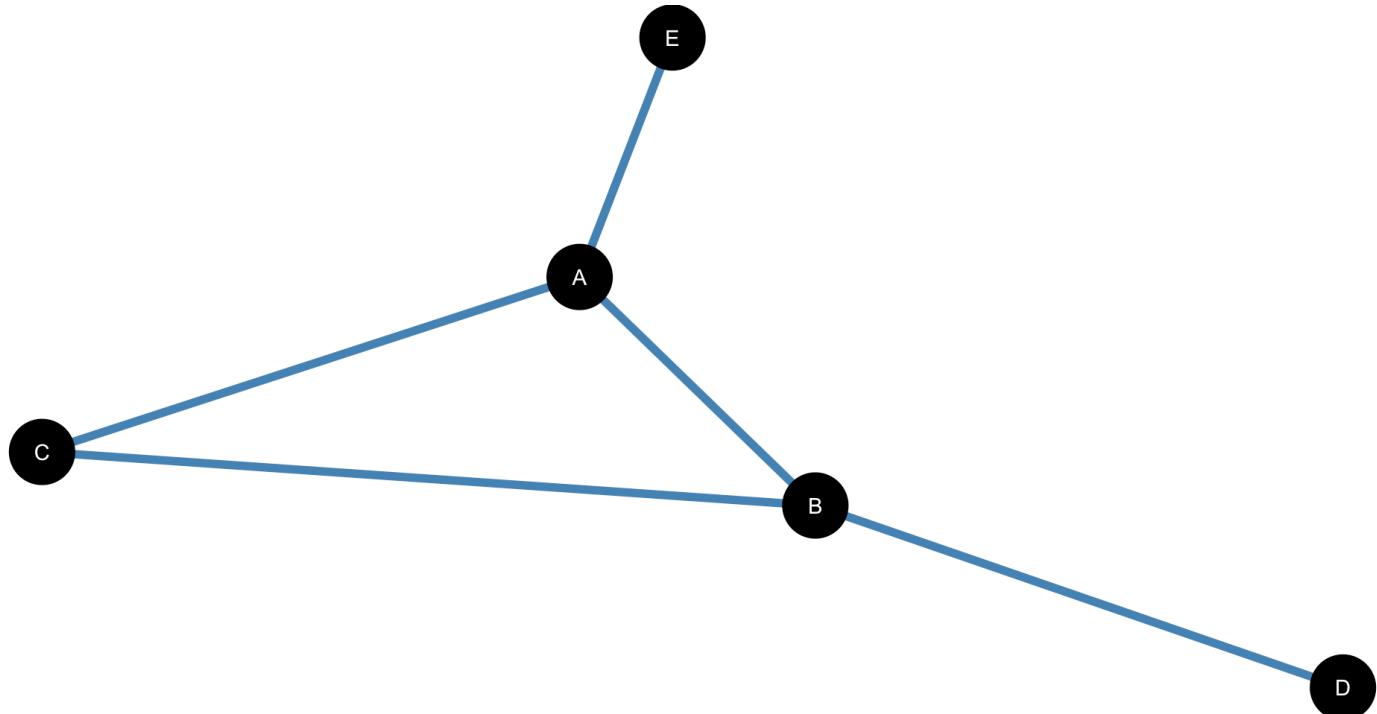
Representations: Adjacency Matrix

1	A	B	C	D	E	
2	A	0	1	1	0	1
3	B	1	0	1	1	0
4	C	1	1	0	0	0
5	D	0	1	0	0	0
6	E	1	0	0	0	0



Representations: Adjacency List

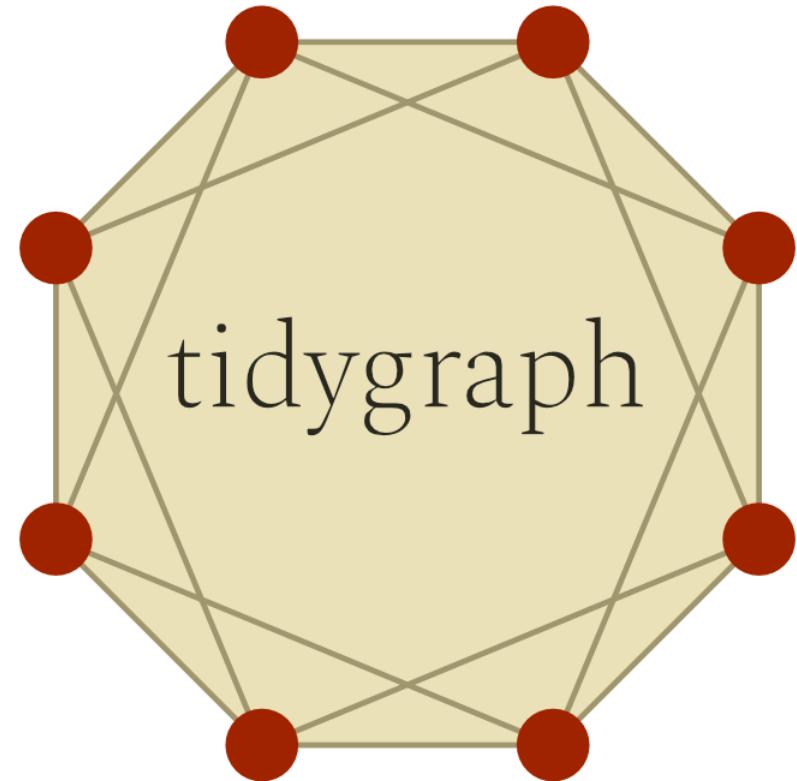
```
1 A: B, C, E  
2 B: A, D, C  
3 C: A, B  
4 D: B  
5 E: A
```



Basics of tidygraph

What are tidygraph and ggraph?

- tidygraph is used for graph manipulation (built upon igraph).
- ggraph is used for visualizing graphs created or manipulated using tidygraph.



<https://www.data-imaginist.com/>

How tidygraph works

tidygraph stores graph data using two tibbles:

- Nodes tibble: Contains node attributes.
- Edges tibble: Contains edge attributes.

```
1 library(tidygraph)
2 # Create a graph
3 graph <- tbl_graph(
4   nodes = tibble(name = c("A", "B", "C")),
5   edges = tibble(from = c(1, 2, 3), to = c(2, 3, 1))
6 )
7 print(graph)

# A tbl_graph: 3 nodes and 3 edges
#
# A directed simple graph with 1 component
#
# Node Data: 3 × 1 (active)
#   name
#   <chr>
1 A
2 B
3 C
#
# Edge Data: 3 × 2
#   from     to
#   <int> <int>
1     1     2
2     2     3
3     3     1
```

How to Create a Graph?

Two ways:

1. `tbl_graph(nodes = ..., edges = ...)`: used when you have explicit data frames for nodes and edges.

- Nodes: A data frame where each row represents a node, and columns may include attributes (e.g., name, type).
- Edges: A data frame where each row represents a connection between nodes, typically with from and to columns to indicate relationships.

```
1 library(tidygraph)
2 nodes <- data.frame(name = c("A", "B", "C"), type = c("1", "1", "2"))
3 edges <- data.frame(from = c(1, 2), to = c(2, 3))
4 graph <- tbl_graph(nodes = nodes, edges = edges)
```

How to Create a Graph?

Two ways:

2. `as_tbl_graph()`: from common graph structures like igraph objects, adjacency matrices, or edge lists. Very general, works with a lot of data types like data frames or matrices:

```
1 library(tidygraph)
2 adj_matrix <- matrix(
3   c(0, 1, 0,
4     1, 0, 1,
5     0, 1, 0), nrow = 3, byrow = TRUE,
6   dimnames = list(c("A", "B", "C"), c("A", "B", "C")))
7 )
8 graph <- as_tbl_graph(adj_matrix)
```

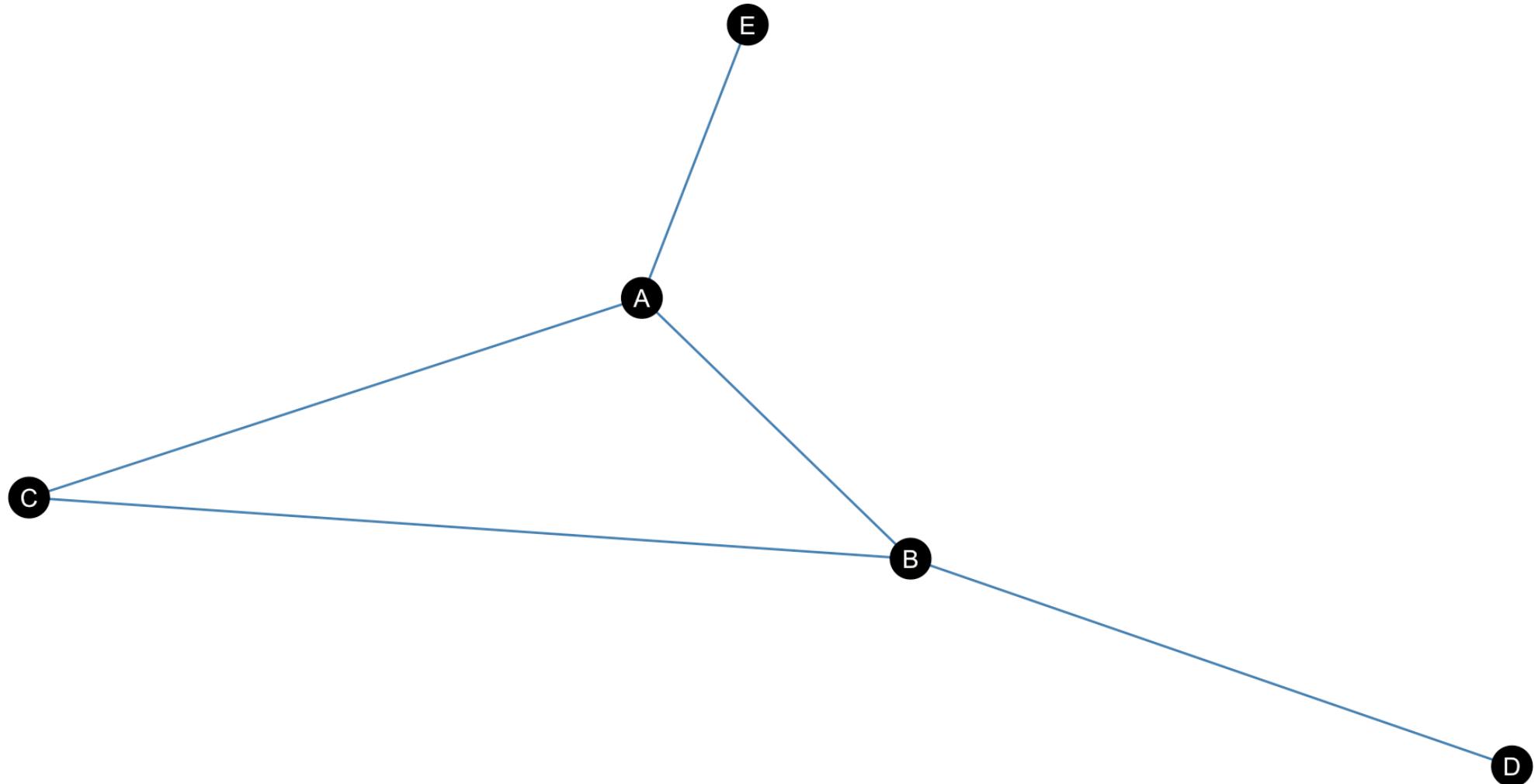
Key Operations in tidygraph

- `activate()`: Switch between node and edge context.
- `filter()`: Subset nodes or edges.
- `mutate()`: Add or modify attributes.

```
1 library(tidygraph)
2 edges <- data.frame(
3   from = c("A", "A", "B", "B", "E"),
4   to = c("B", "C", "D", "C", "A")
5 )
6 graph <- as_tbl_graph(edges, directed = FALSE)
7 print(graph)

# A tbl_graph: 5 nodes and 5 edges
#
# An undirected simple graph with 1 component
#
# Node Data: 5 × 1 (active)
#   name
#   <chr>
1 A
2 B
3 E
4 C
5 D
#
# Edge Data: 5 × 2
#   from     to
#   <int> <int>
1     1     2
2     1     4
3     2     5
# i 2 more rows
```

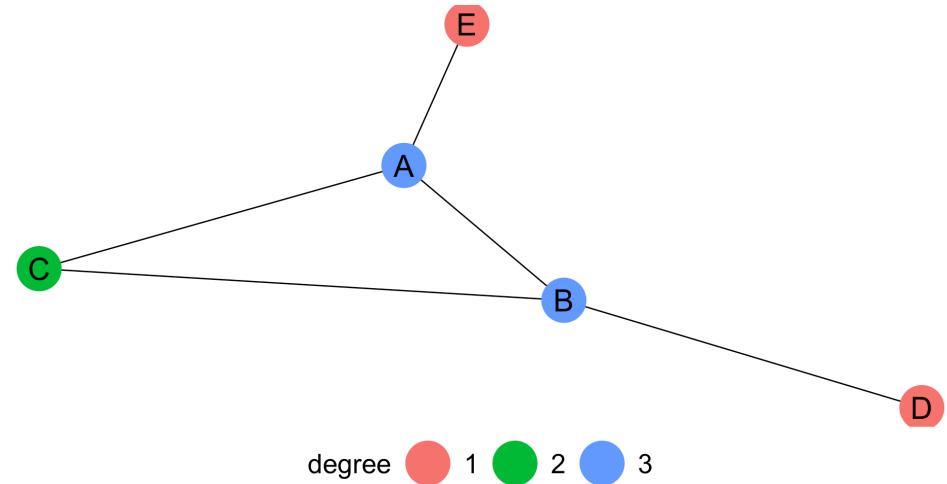
Key Operations in tidygraph



Key Operations in tidygraph

```
1 # Activate nodes
2 g = graph %>%
3   activate(nodes) %>%
4   mutate(degree = as.factor(centrality_degree())) .
5 print(g)

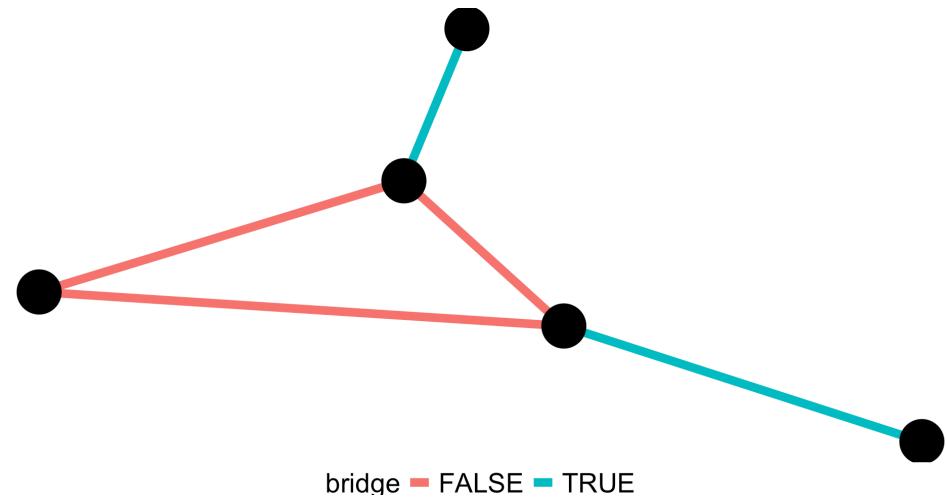
# A tbl_graph: 5 nodes and 5 edges
#
# An undirected simple graph with 1 component
#
# Node Data: 5 × 2 (active)
  name degree
  <chr> <fct>
1 A     3
2 B     3
3 E     1
4 C     2
5 D     1
#
# Edge Data: 5 × 2
  from    to
  <int> <int>
1     1     2
2     1     4
3     2     5
# i 2 more rows
```



Key Operations in tidygraph

```
1 # Add a new attribute to edges
2 g = graph %>%
3   activate(edges) %>%
4   mutate(bridge = tidygraph::edge_is_bridge())
5 print(g)

# A tbl_graph: 5 nodes and 5 edges
#
# An undirected simple graph with 1 component
#
# Edge Data: 5 × 3 (active)
#   from      to bridge
#   <int> <int> <lgl>
1     1      2 FALSE
2     1      4 FALSE
3     2      5 TRUE
4     2      4 FALSE
5     1      3 TRUE
#
# Node Data: 5 × 1
#   name
#   <chr>
1 A
2 B
3 E
# i 2 more rows
```



Key Operations in tidygraph

```
1 # Filter nodes with specific names
2 graph %>%
3   activate(nodes) %>%
4   filter(name != "B")
```

```
# A tbl_graph: 4 nodes and 2 edges
#
# An unrooted forest with 2 trees
#
# Node Data: 4 × 1 (active)
  name
  <chr>
1 A
2 E
3 C
4 D
#
# Edge Data: 2 × 2
  from    to
  <int> <int>
1     1      3
2     1      2
```

```
1 # Filter edges
2 graph %>%
3   activate(edges) %>%
4   filter(from == 1)
```

```
# A tbl_graph: 5 nodes and 3 edges
#
# An unrooted forest with 2 trees
#
# Edge Data: 3 × 2 (active)
  from    to
  <int> <int>
1     1      2
2     1      4
3     1      3
#
# Node Data: 5 × 1
  name
  <chr>
1 A
2 B
3 E
# i 2 more rows
```

Visualization with ggraph

What is ggraph?

ggraph provides an easy way to visualize graphs using intuitive functions and multiple layout options.

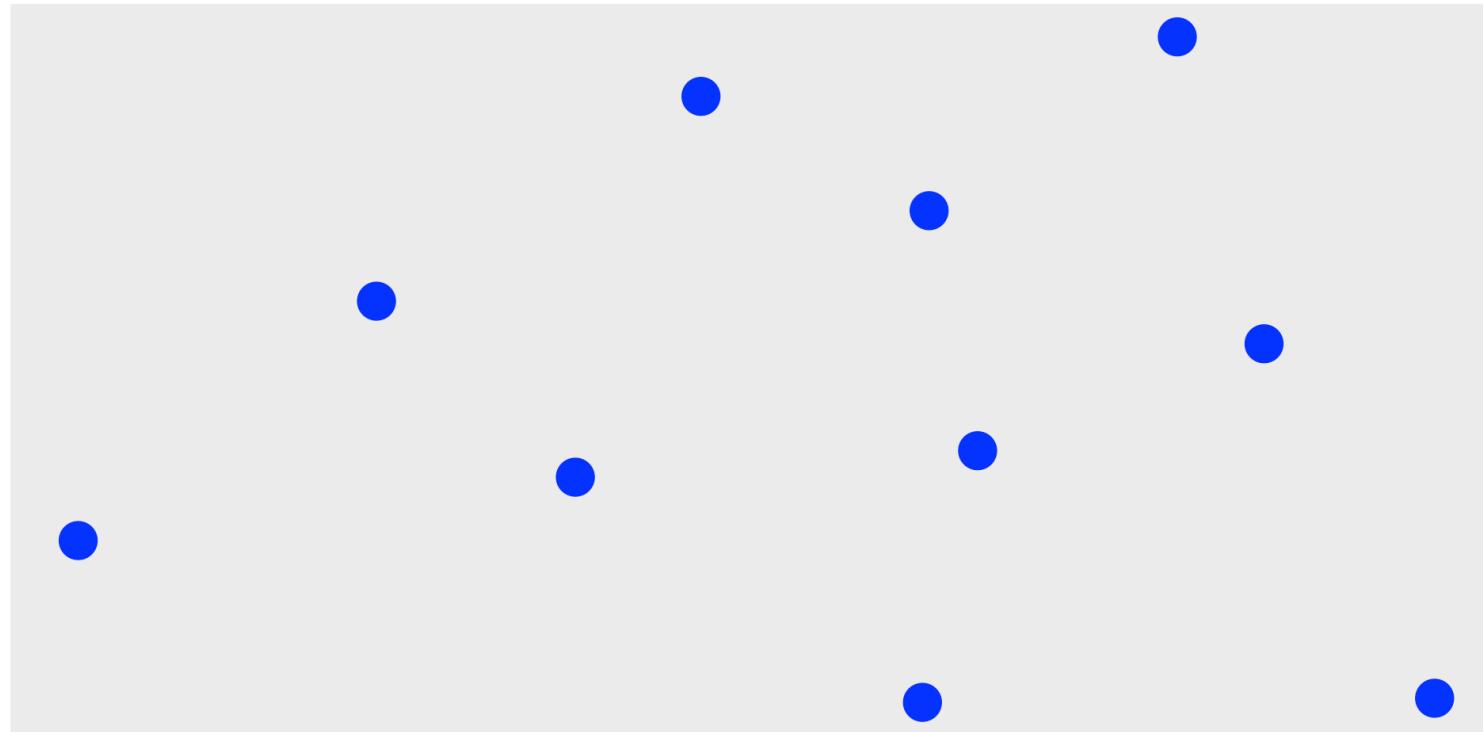
- Supports layouts like tree, circular, and force-directed.
- Seamlessly integrates with ggplot2.

Overview of Node Layers

- Node layers in ggraph are used to represent the entities of a graph visually.
- Common aesthetics include size (importance of nodes), color (group membership or categories), and shape (node type).
- The `geom_node_point()` function is the primary method to plot nodes.

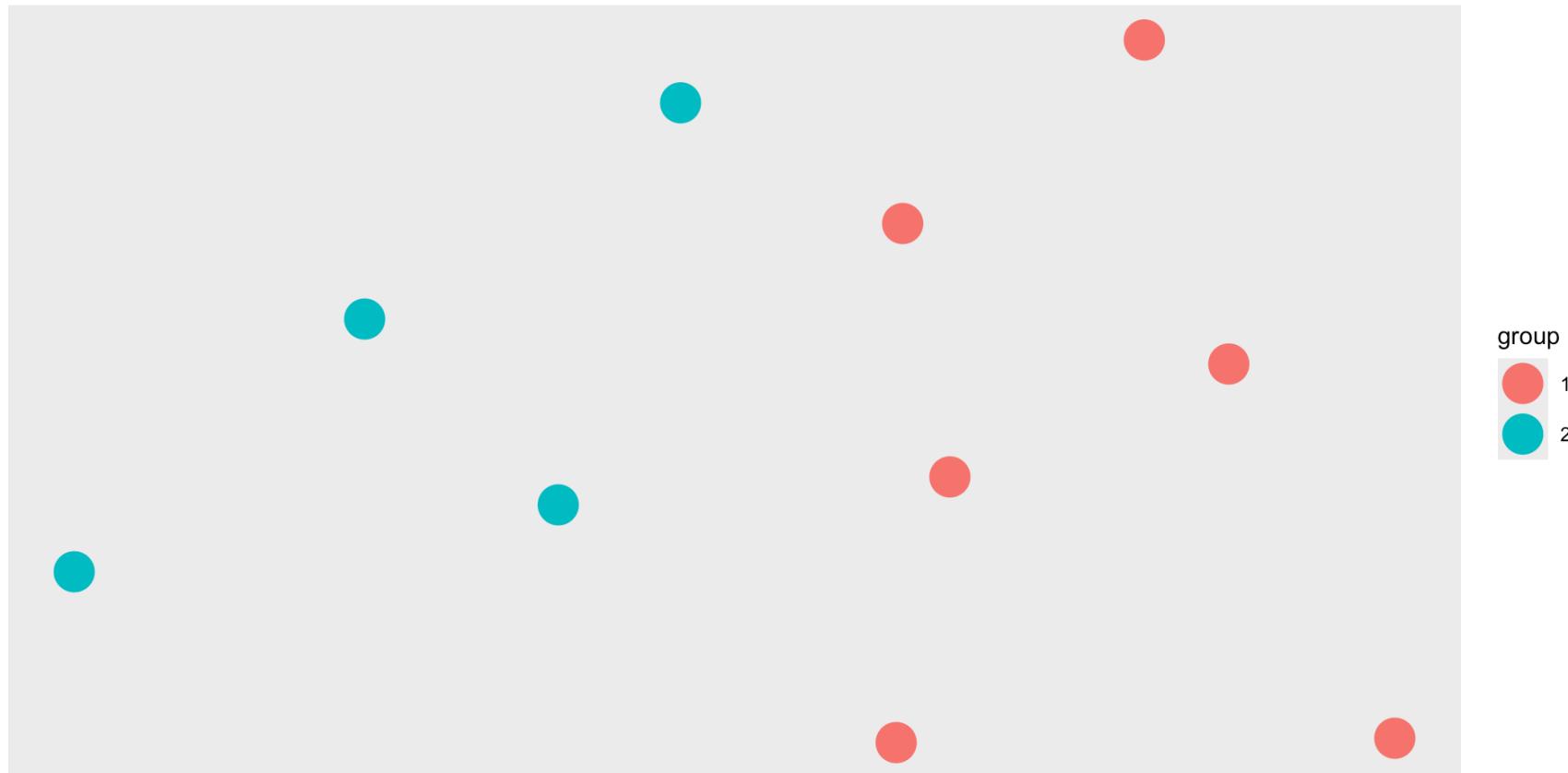
geom_node_point

```
1 library(ggraph)
2 set.seed(7)
3 graph <- sample_gnp(n = 10, p = 0.3) # Random graph
4 graph = as_tbl_graph(graph, directed = FALSE)
5 ggraph(graph) +
6   geom_node_point(size = 8, color = 'blue')
```



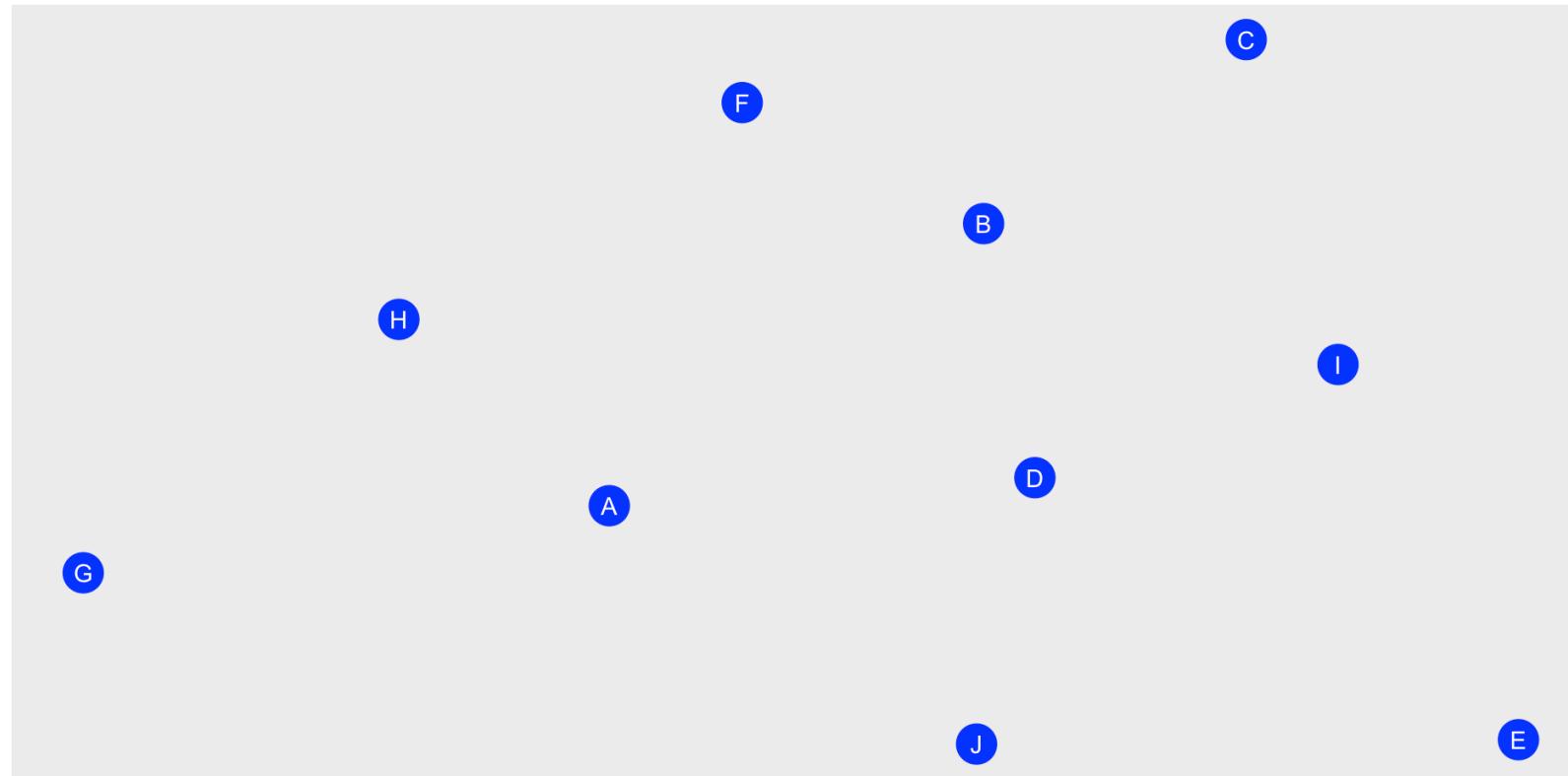
geom_node_point: aesthetics

```
1 gr <- graph %>% activate(nodes) %>%
2   mutate(group = as.factor(group_louvain()))
3 ggraph(gr) +
4   geom_node_point(aes(color = group), size = 8)
```



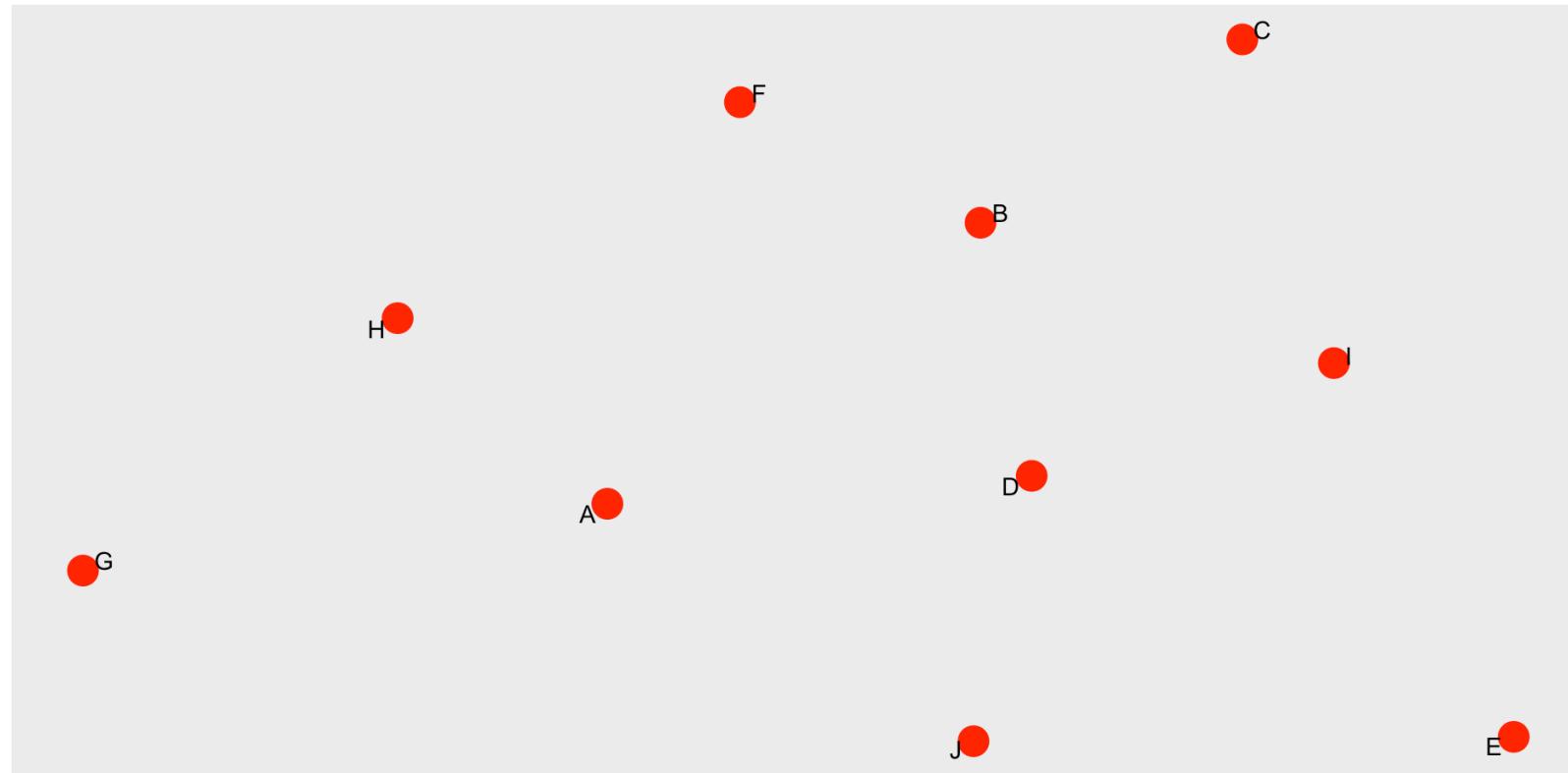
Adding Node Labels

```
1 gr = graph %>% activate(nodes) %>%
2   mutate(name = LETTERS[1:10])
3 ggraph(gr) +
4   geom_node_point(size = 8, color = 'blue') +
5   geom_node_text(aes(label = name), color = "white")
```



Adding Node Labels: repel

```
1 gr = graph %>% activate(nodes) %>%
2   mutate(name = LETTERS[1:10])
3 ggraph(gr) +
4   geom_node_point(size = 6, color = 'red') +
5   geom_node_text(aes(label = name), color = "black", repel = T)
```

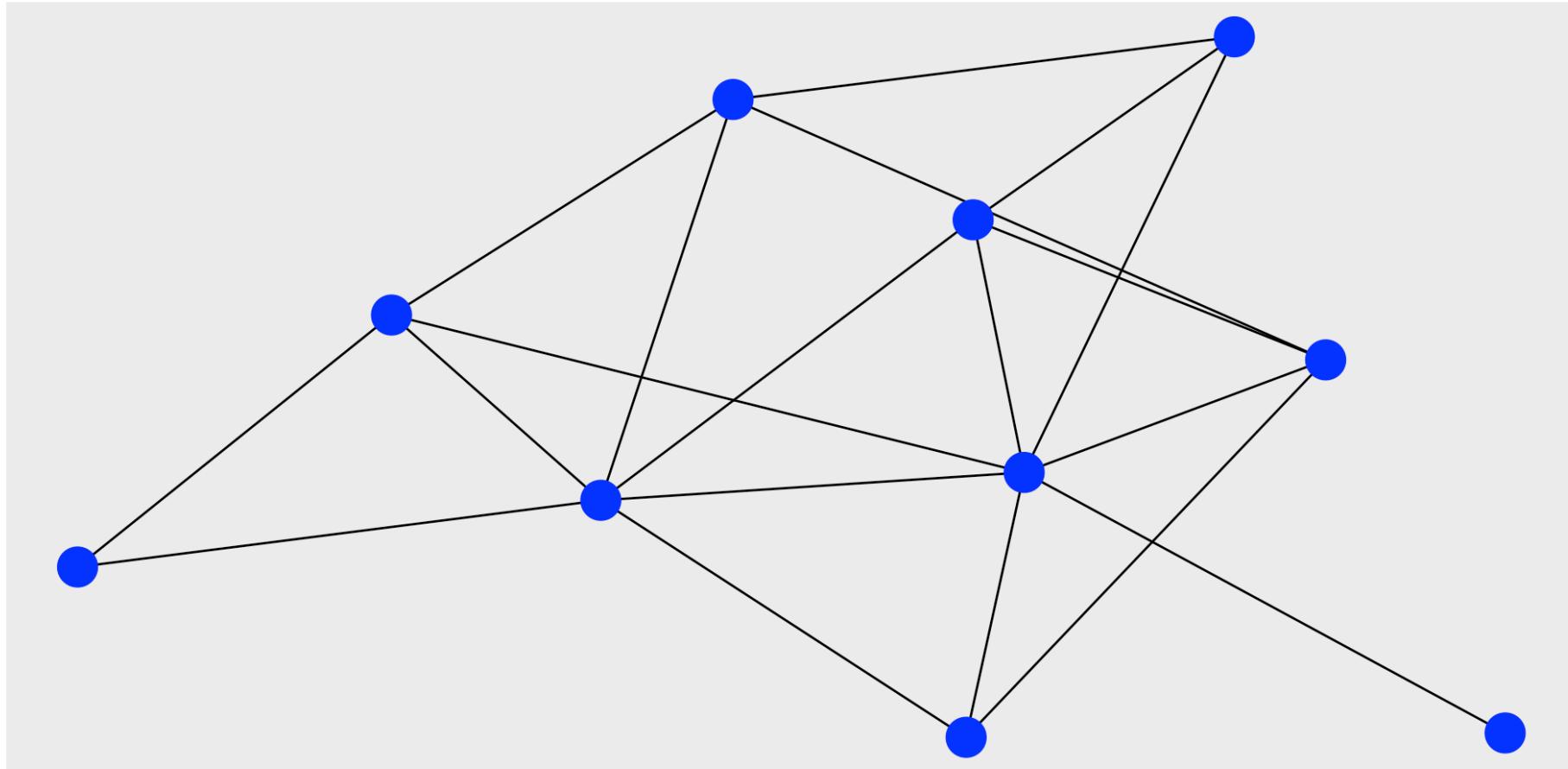


Overview of Edge Layers

- Edges represent the connections or relationships between nodes in a graph.
- Edge layers in ggraph can visualize various aspects of relationships, such as strength, direction, or type through aesthetics like width, color, and transparency.
- The `geom_edge_link()` function is the most commonly used layer for visualizing edges.

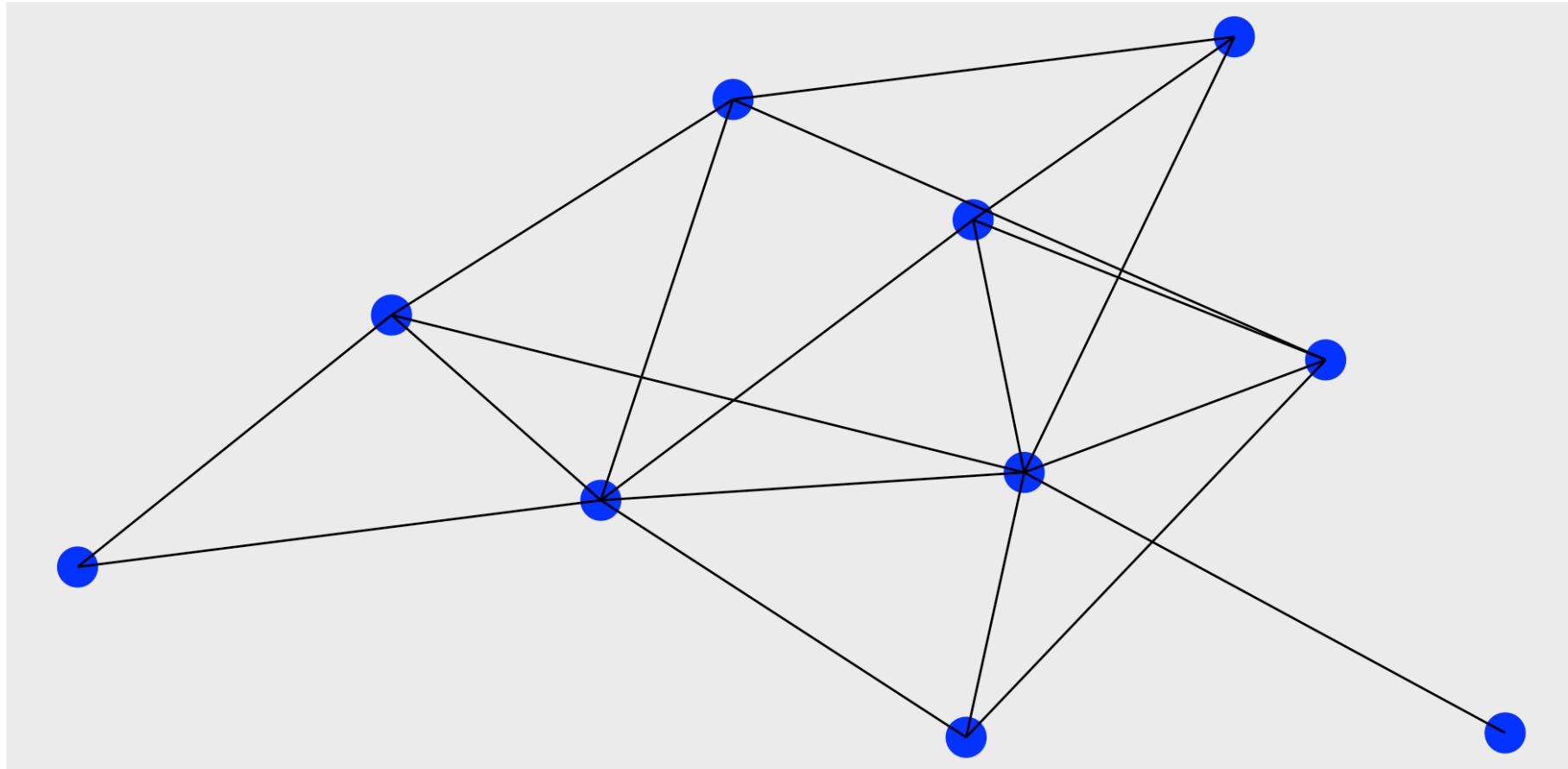
geom_edge_link

```
1 ggraph(graph) +  
2   geom_edge_link() +  
3   geom_node_point(size = 8, color = 'blue')
```



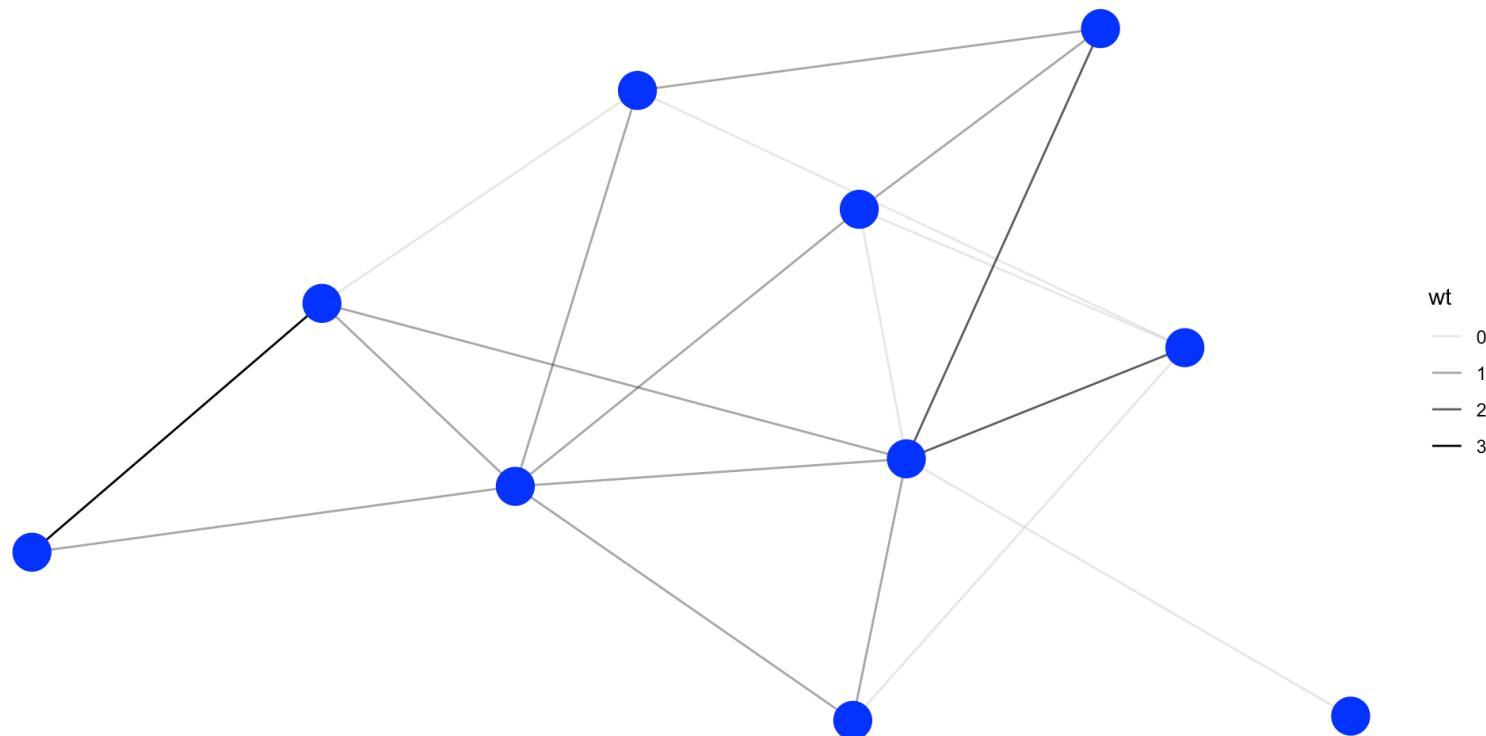
geom_edge_link: Order Matters

```
1 ggraph(graph) +  
2   geom_node_point(size = 8, color = 'blue') +  
3   geom_edge_link()
```



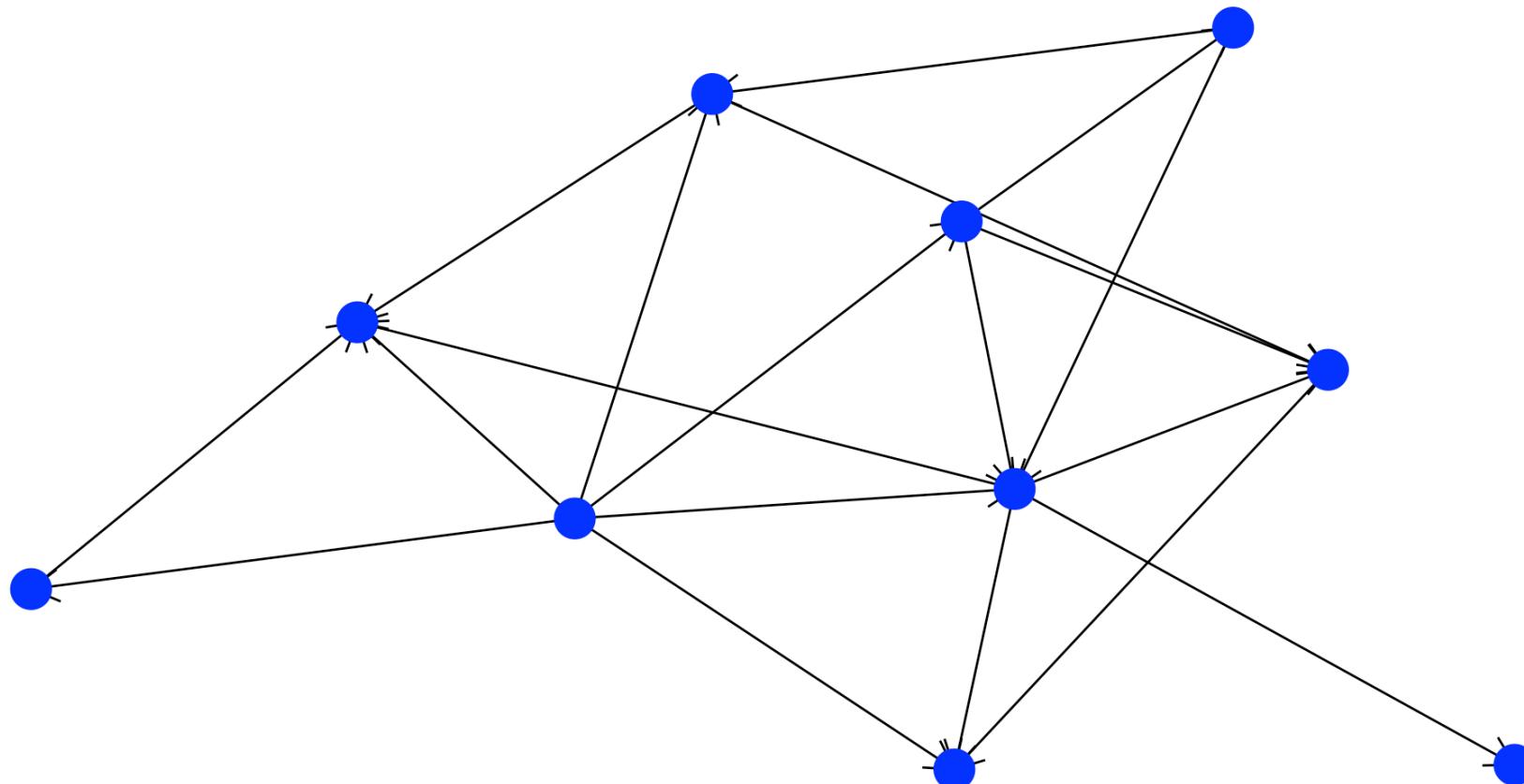
geom_edge_link: aesthetics

```
1 gr <- graph %>% activate(edges) %>%
2   mutate(wt = rpois(n(), lambda = 1))
3 ggraph(gr) +
4   geom_edge_link(aes(alpha = wt)) +
5   geom_node_point(size = 8, color = 'blue') +
6   theme_void()
```



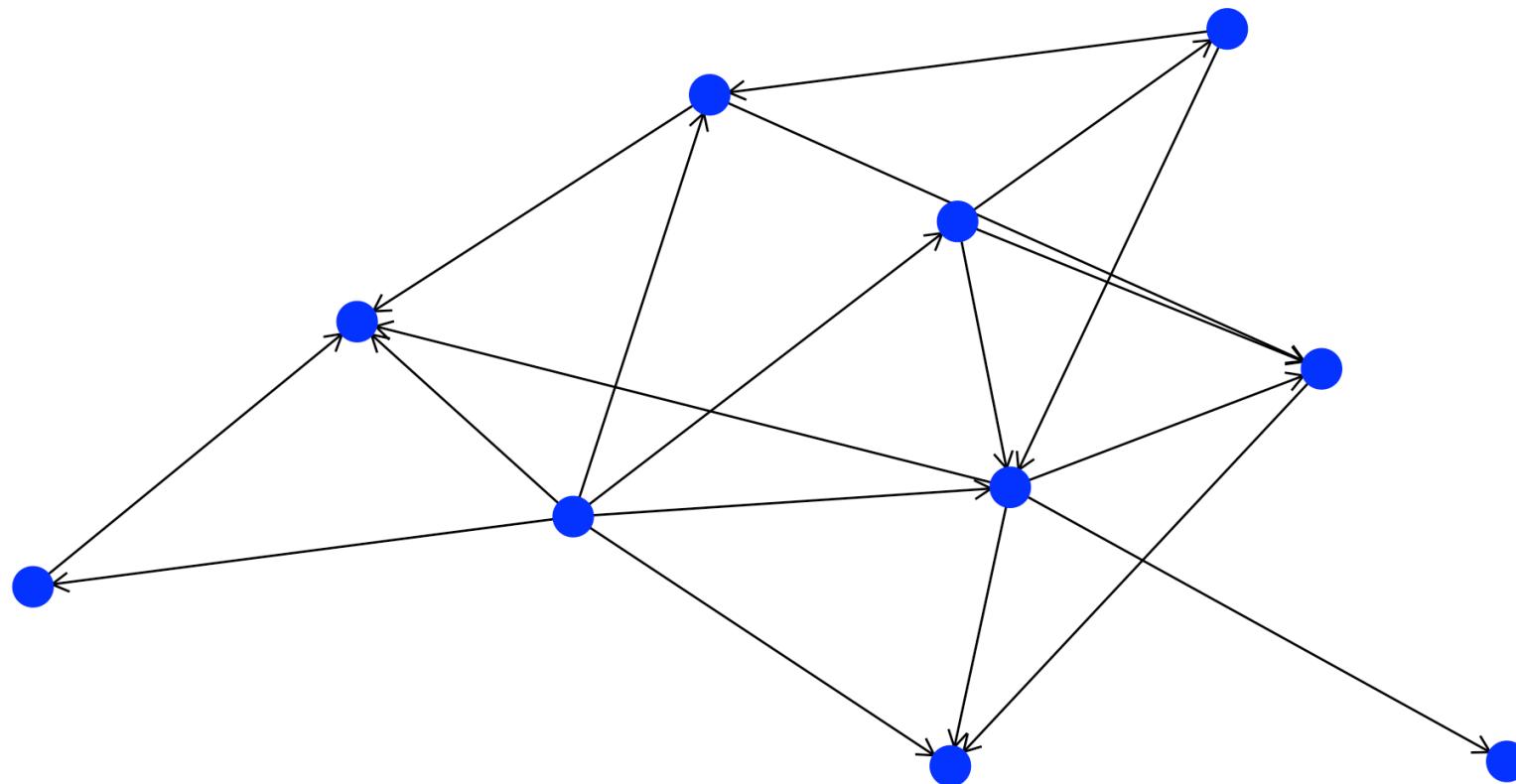
Edges in a Directed Graphs

```
1 ggraph(graph) +  
2   geom_edge_link(arrows = arrow(length = unit(5, 'mm')))) +  
3   geom_node_point(size = 8, color = 'blue') +  
4   theme_void()
```



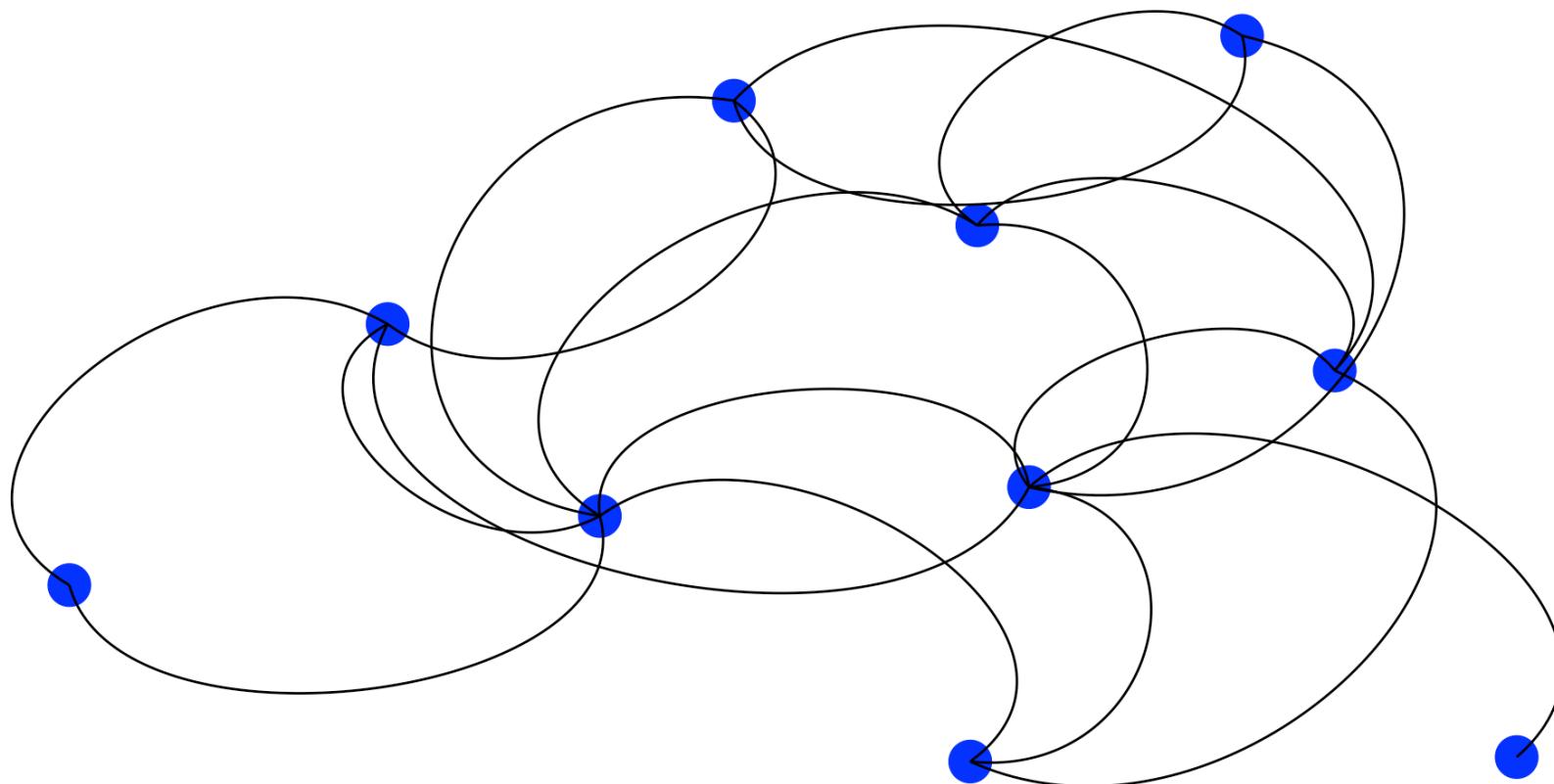
Edges in a Directed Graphs: end_cap

```
1 ggraph(graph) +  
2   geom_edge_link(arrows = arrow(length = unit(3, 'mm')),  
3                     end_cap = circle(3, 'mm')) +  
4   geom_node_point(size = 8, color = 'blue') +  
5   theme_void()
```



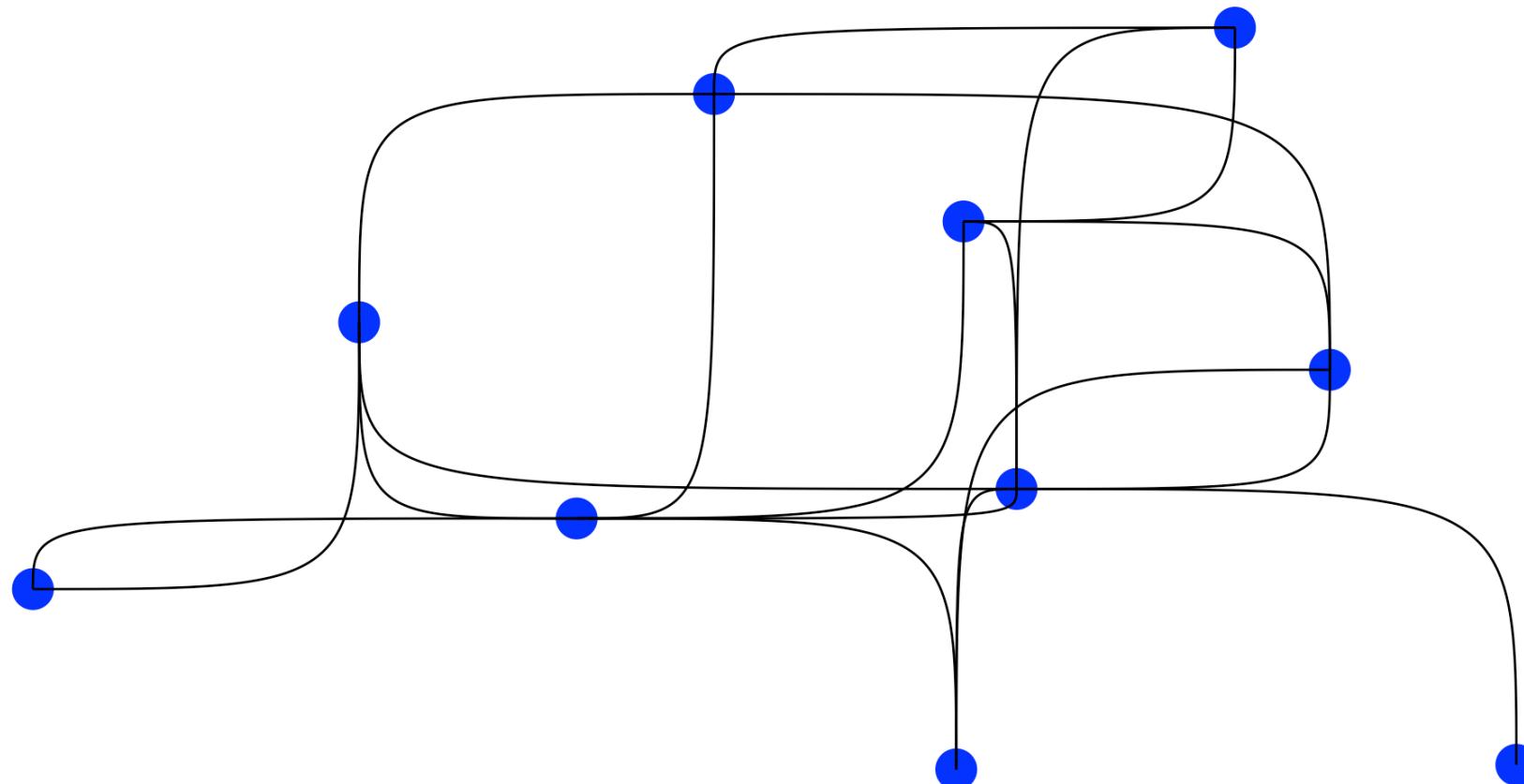
Other Kinds of Edges: arc

```
1 ggraph(graph) +  
2   geom_node_point(size = 8, color = 'blue') +  
3   geom_edge_arc() +  
4   theme_void()
```



Other Kinds of Edges: bend

```
1 ggraph(graph) +  
2   geom_node_point(size = 8, color = 'blue') +  
3   geom_edge_bend() +  
4   theme_void()
```



Overview of Layouts in ggraph

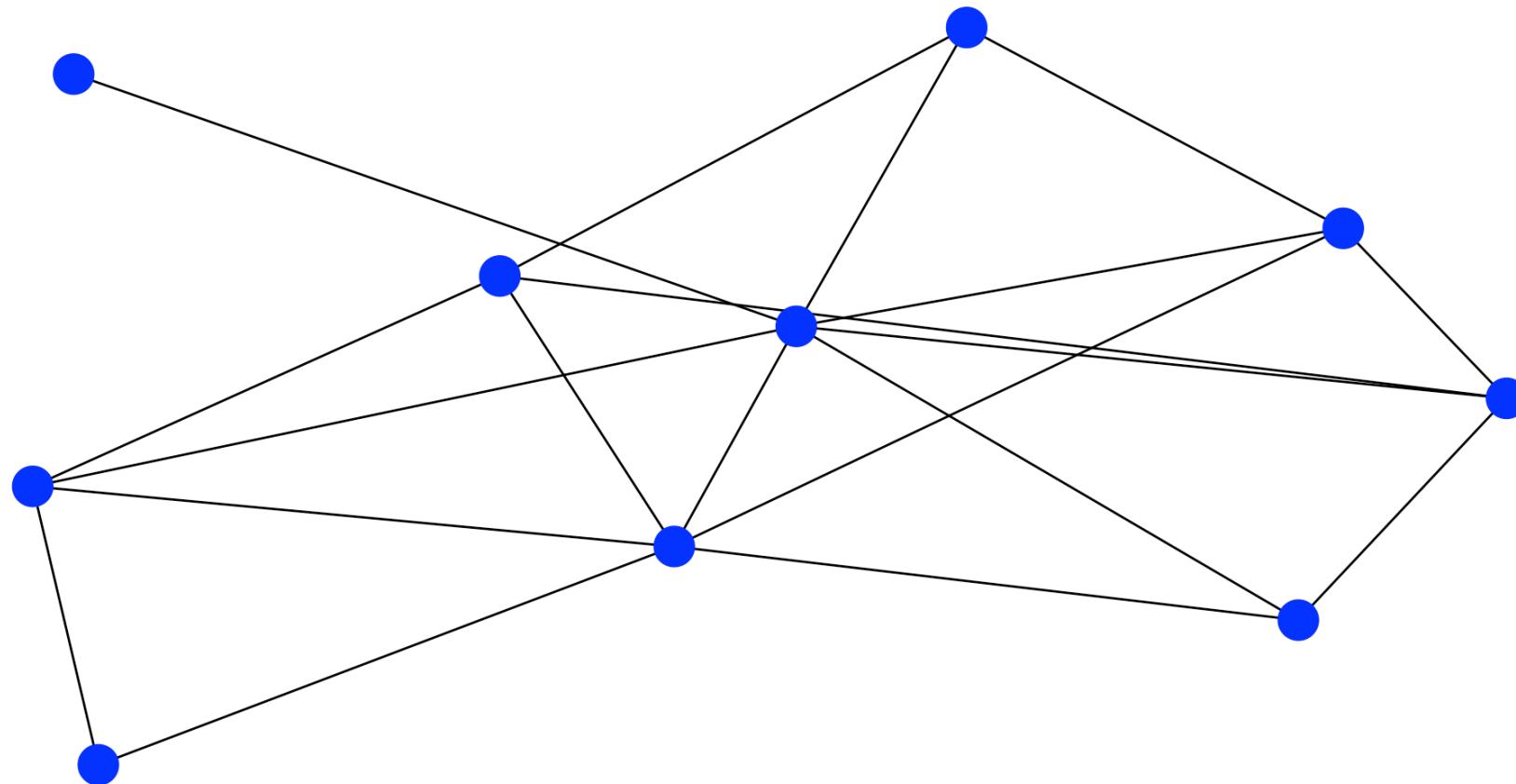
- Layouts determine the spatial arrangement of nodes and edges in a graph.
- Choosing the right layout is essential for clarity and aesthetics.
- ggraph supports a variety of layouts, including force-directed, circular, and hierarchical types.

Force-directed Layouts

- Force-directed layouts, such as ‘fr’ (Fruchterman-Reingold) and ‘kk’ (Kamada-Kawai), simulate physical forces to position nodes.
- They are ideal for visualizing relationships dynamically.

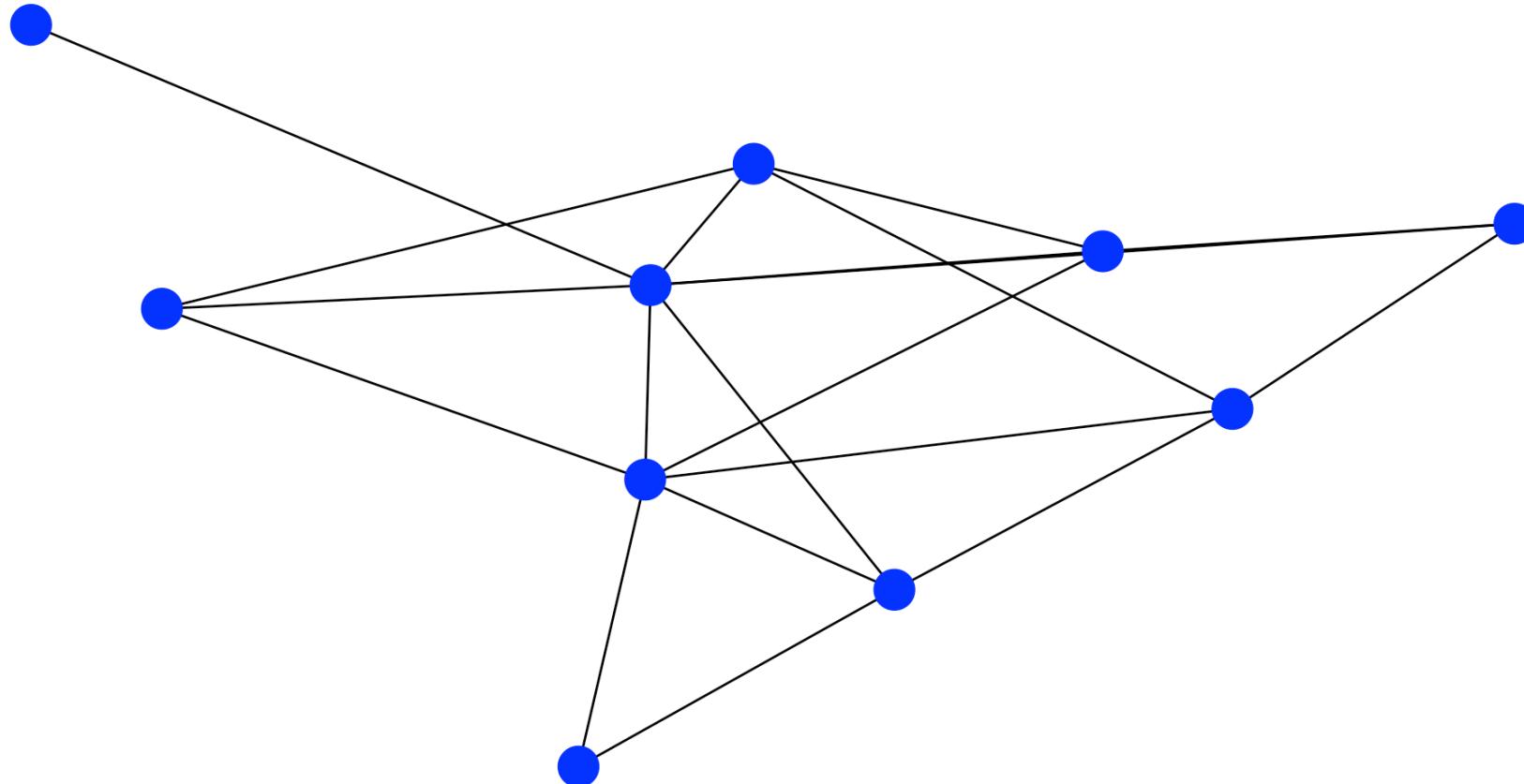
Kamada-Kawai Layout

```
1 ggraph(graph, layout = "kk") +  
2   geom_edge_link() +  
3   geom_node_point(size = 8, color = 'blue') +  
4   theme_void()
```



Fruchterman-Reingold Layout

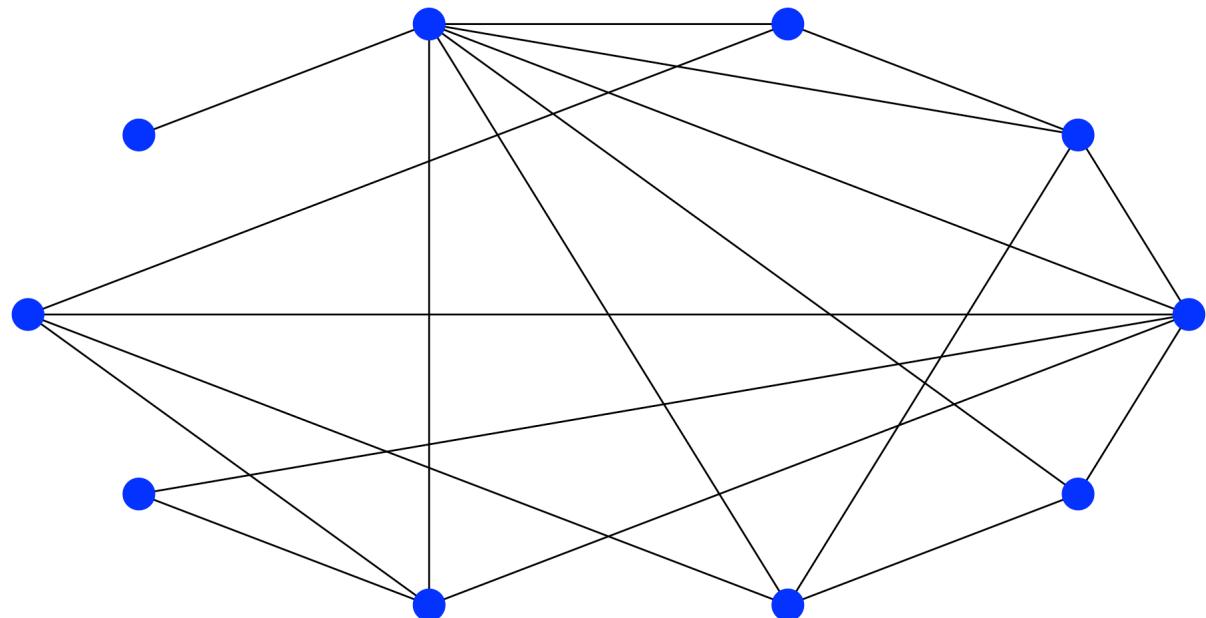
```
1 ggraph(graph, layout = "fr") +  
2   geom_edge_link() +  
3   geom_node_point(size = 8, color = 'blue') +  
4   theme_void()
```



Circular Layouts

Circular layouts position nodes around a circle, making them suitable for cyclical or hierarchical relationships.

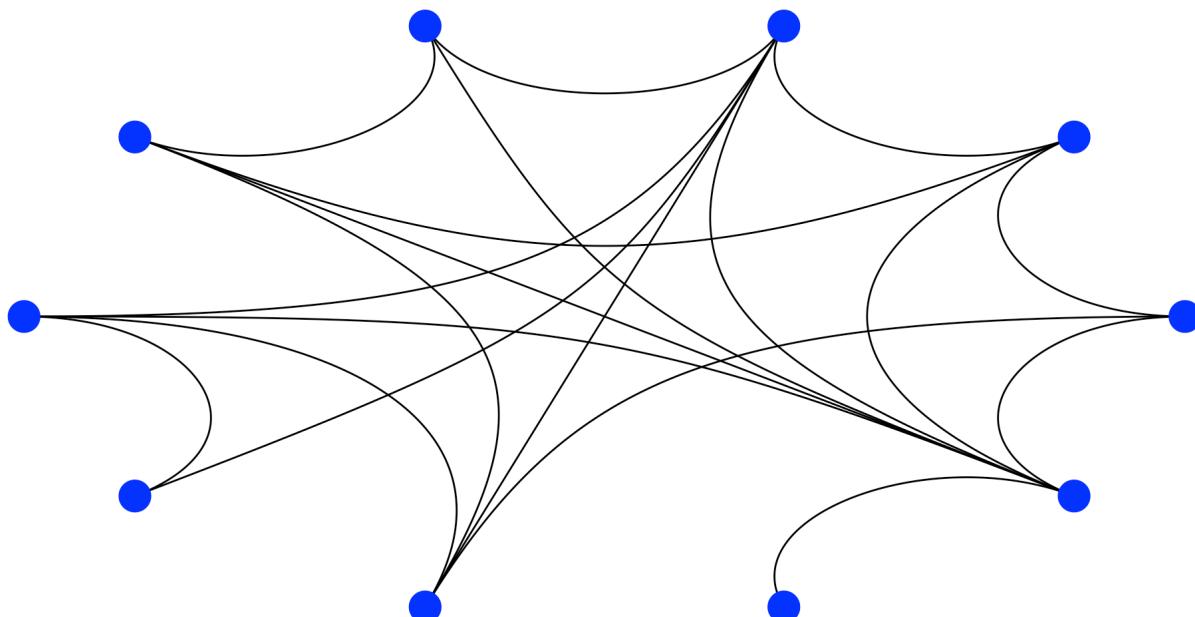
```
1 ggraph(graph, layout = "circle") +  
2   geom_edge_link() +  
3   geom_node_point(size = 8, color = 'blue') +  
4   theme_void()
```



Chord Diagrams

Used to compare similarities between groups of data or within a dataset.

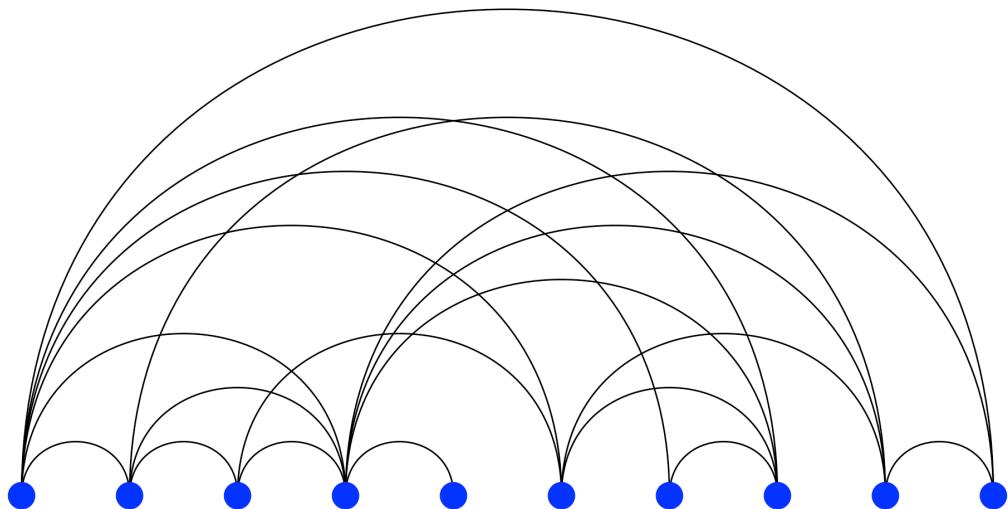
```
1 ggraph(graph, layout = "linear", circular = T) +  
2   geom_edge_arc() +  
3   geom_node_point(size = 8, color = 'blue') +  
4   theme_void()
```



Arc Diagrams

Nodes are placed along a single line (a one-dimensional axis) and arcs are used to show connections between those nodes, useful in finding the co-occurrence within the data.

```
1 ggraph(graph, layout = "linear") +  
2   geom_edge_arc() +  
3   geom_node_point(size = 8, color = 'blue') +  
4   theme_void()
```



An Example

nycflights Data

```
1 library(nycflights13)
2 df = flights
3 head(df)

# A tibble: 6 × 19
  year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
  <int> <int> <int>     <int>           <int>      <dbl>     <int>           <int>
1 2013     1     1      517            515        2         830           819
2 2013     1     1      533            529        4         850           830
3 2013     1     1      542            540        2         923           850
4 2013     1     1      544            545       -1        1004          1022
5 2013     1     1      554            600       -6         812           837
6 2013     1     1      554            558       -4         740           728
# i 11 more variables: arr_delay <dbl>, carrier <chr>, flight <int>,
# tailnum <chr>, origin <chr>, dest <chr>, air_time <dbl>, distance <dbl>,
# hour <dbl>, minute <dbl>, time_hour <dttm>
```

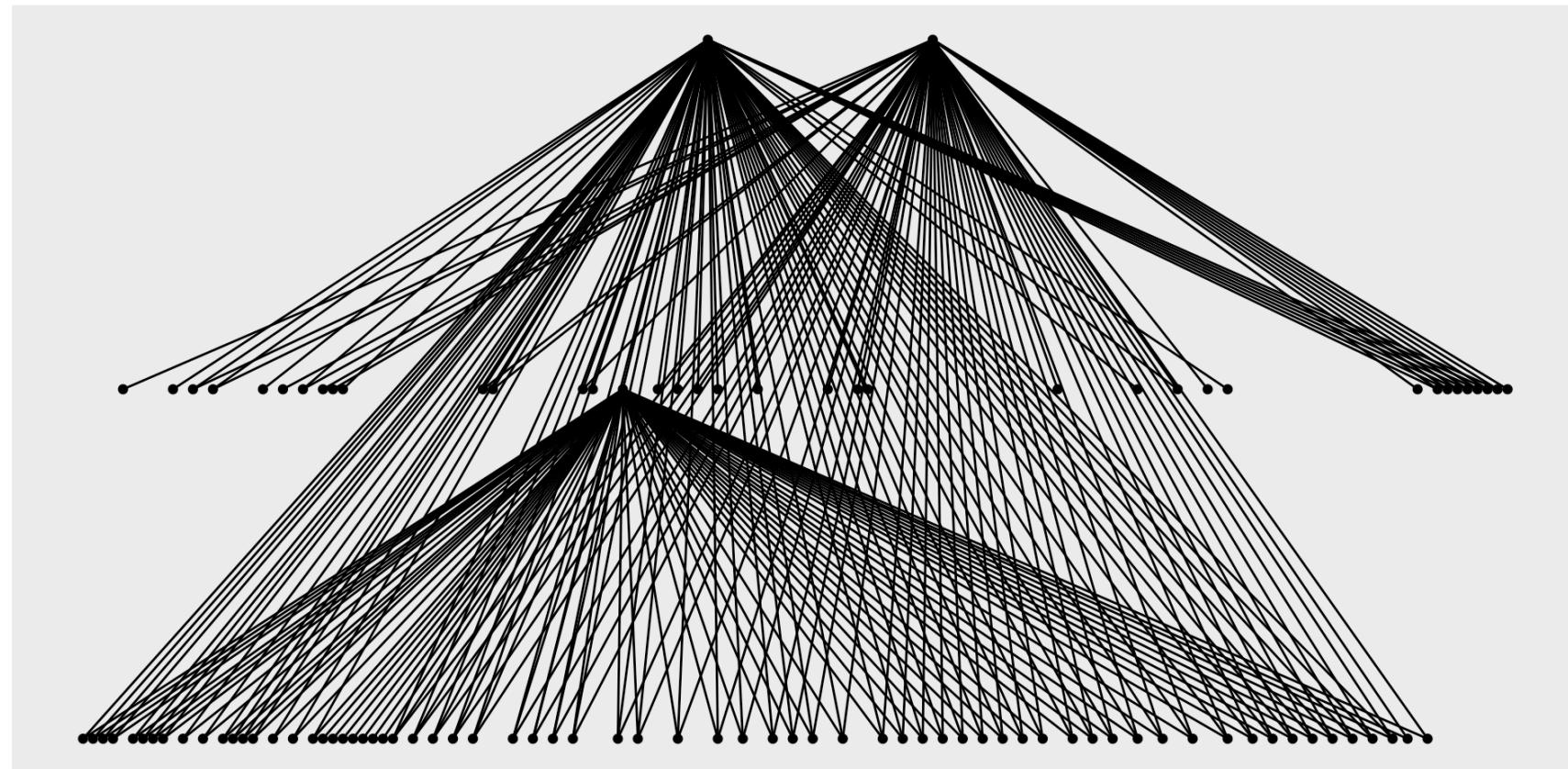
nycflights Data: Graph

```
1 gr = df %>% count(origin, dest) %>%
2   as_tbl_graph()
3 gr

# A tbl_graph: 107 nodes and 224 edges
#
# A directed acyclic simple graph with 1 component
#
# Node Data: 107 × 1 (active)
  name
  <chr>
1 EWR
2 JFK
3 LGA
4 ALB
5 ANC
6 ATL
7 AUS
8 AVL
~ ---
```

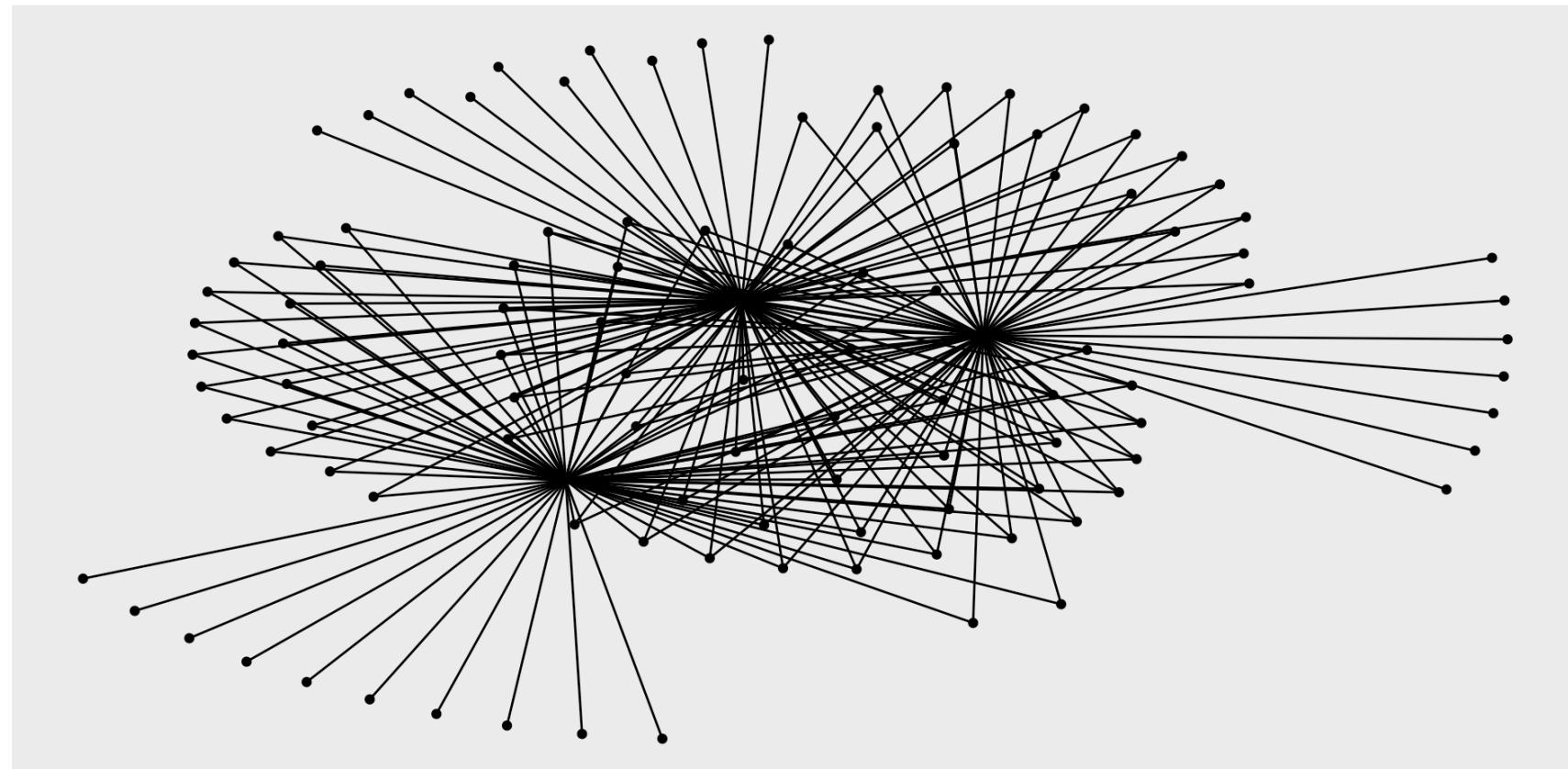
nycflights Data: Simple Visualization

```
1 # base graph  
2 gr %>% ggraph() +  
3   geom_edge_link() +  
4   geom_node_point()
```



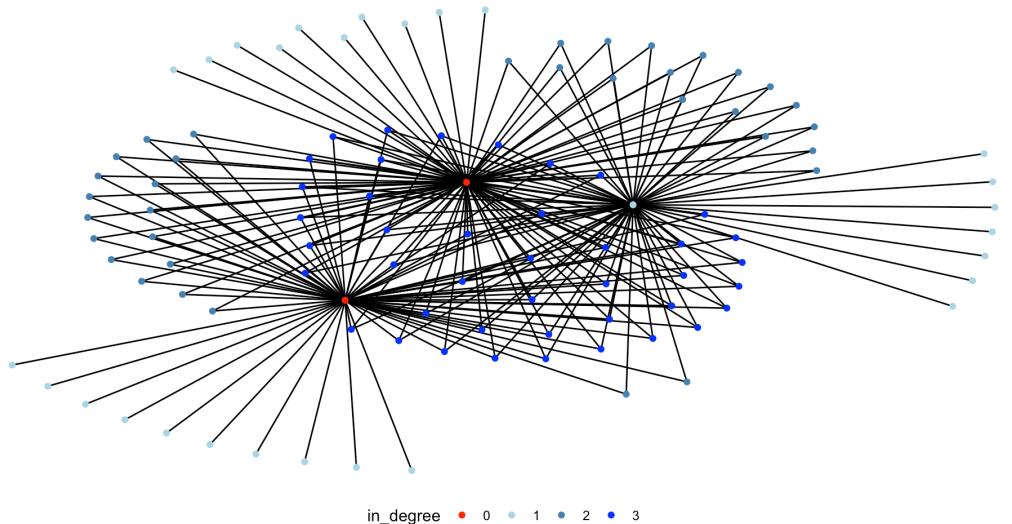
nycflights Data: Layout

```
1 # adjust layout
2 gr %>% ggraph(layout = "kk") +
3   geom_edge_link() +
4   geom_node_point()
```

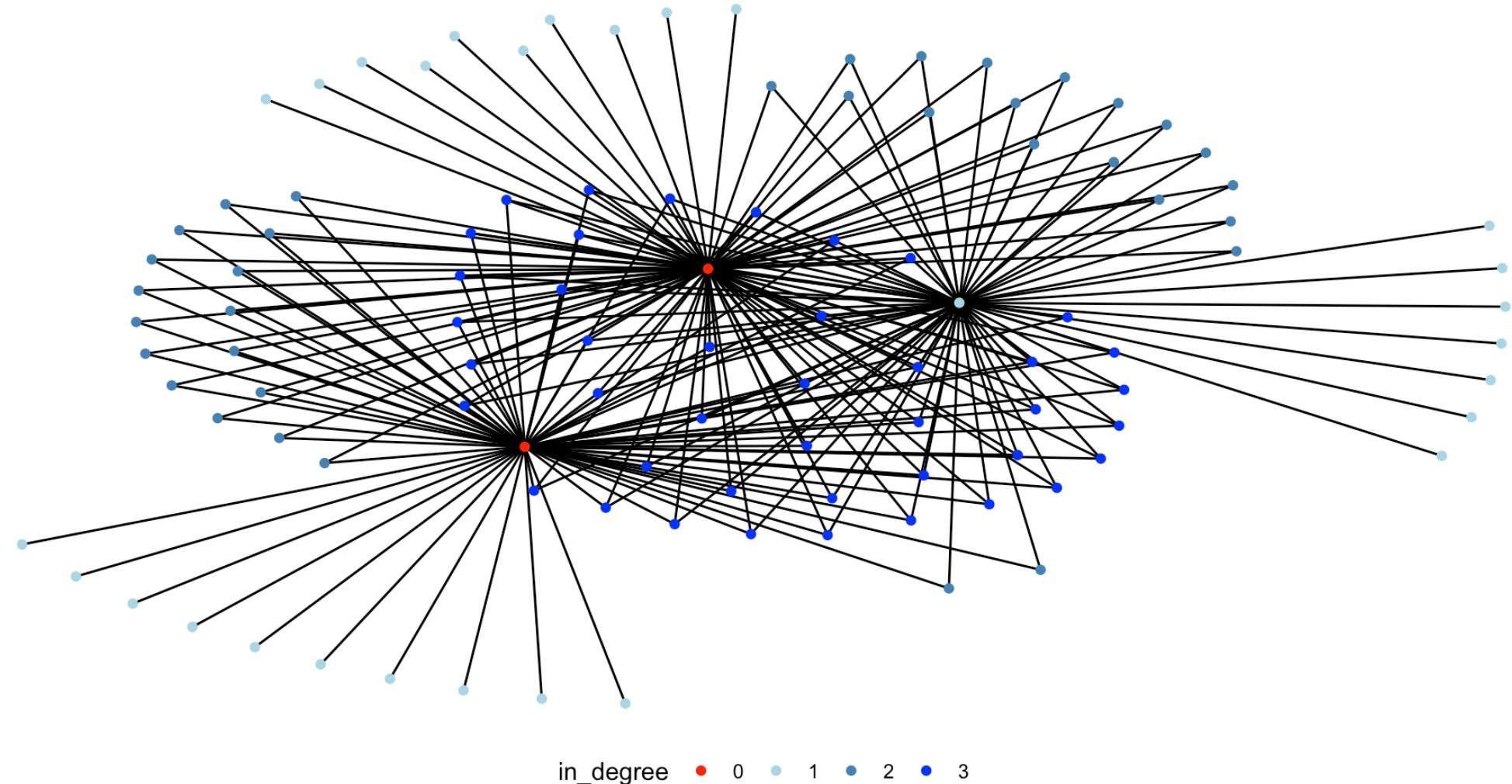


nycflights Data: In-Degree

```
1 # color nodes based on in-degree
2 gr %>%
3   activate(nodes) %>%
4   mutate(in_degree = as.factor(centrality_degree(mode = "in")))) %>%
5   ggraph(layout = "kk") +
6   geom_edge_link() +
7   geom_node_point(aes(color = in_degree)) +
8   scale_color_manual(values = c("0" = "red", "1" = "lightblue", "2" = "steelblue",
9   theme_void() +
10  theme(legend.position = "bottom")
```

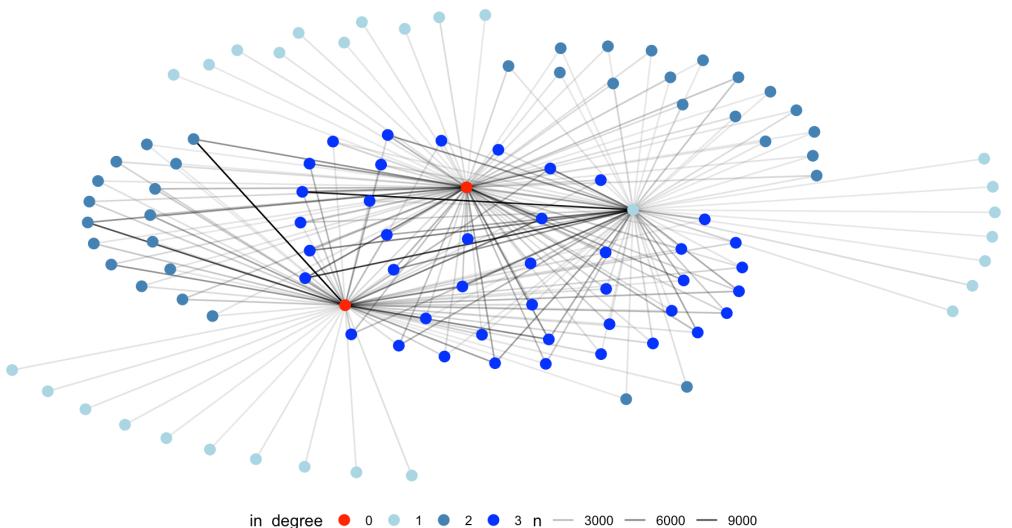


nycflights Data: In-Degree

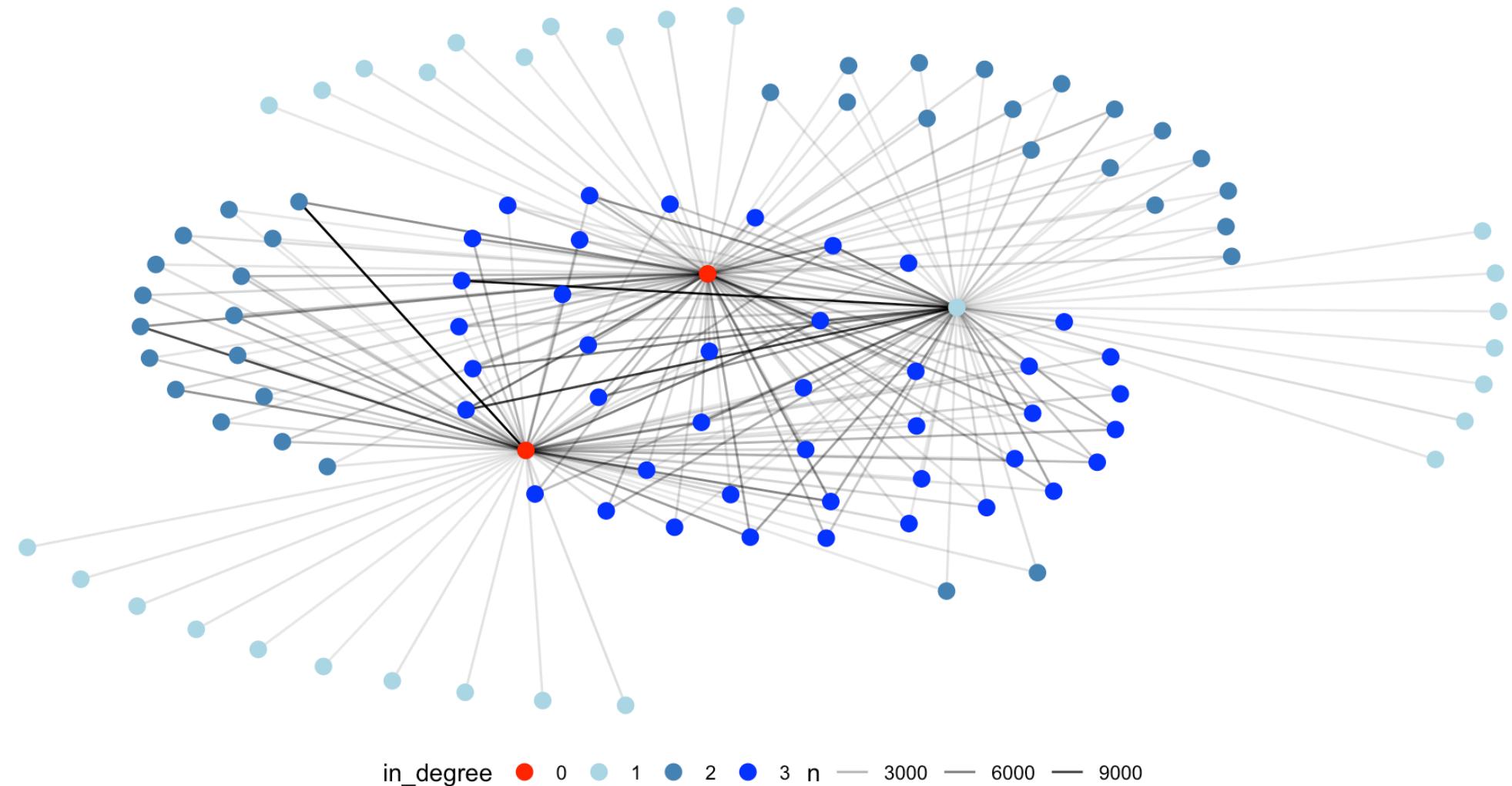


nycflights Data: Number of Flights

```
1 # add info on number of flights to edges + increase size
2 gr %>%
3   activate(nodes) %>%
4   mutate(in_degree = as.factor(centrality_degree(mode = "in")))) %>%
5   ggraph(layout = "kk") +
6   geom_edge_link(aes(alpha = n)) +
7   geom_node_point(aes(color = in_degree), size = 3) +
8   scale_color_manual(values = c("0" = "red", "1" = "lightblue", "2" = "steelblue",
9   theme_void() +
10  theme(legend.position = "bottom")
```

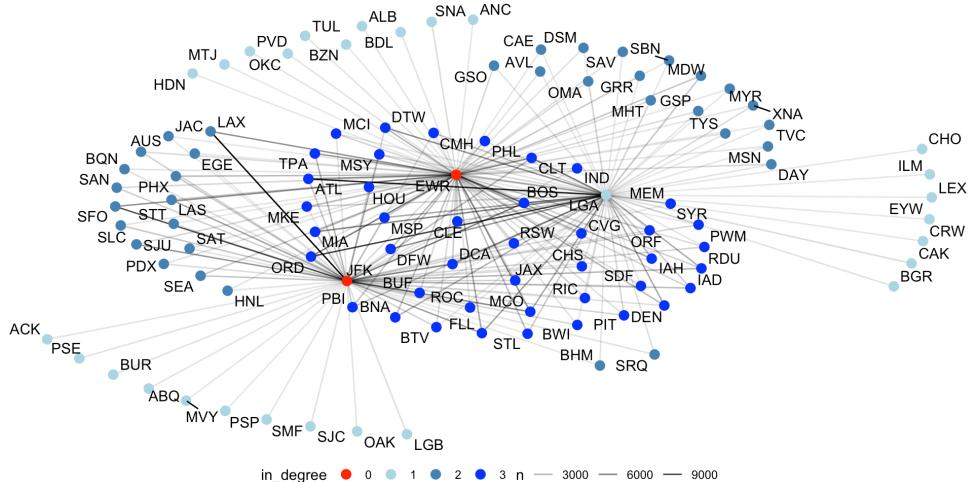


nycflights Data: Number of Flights

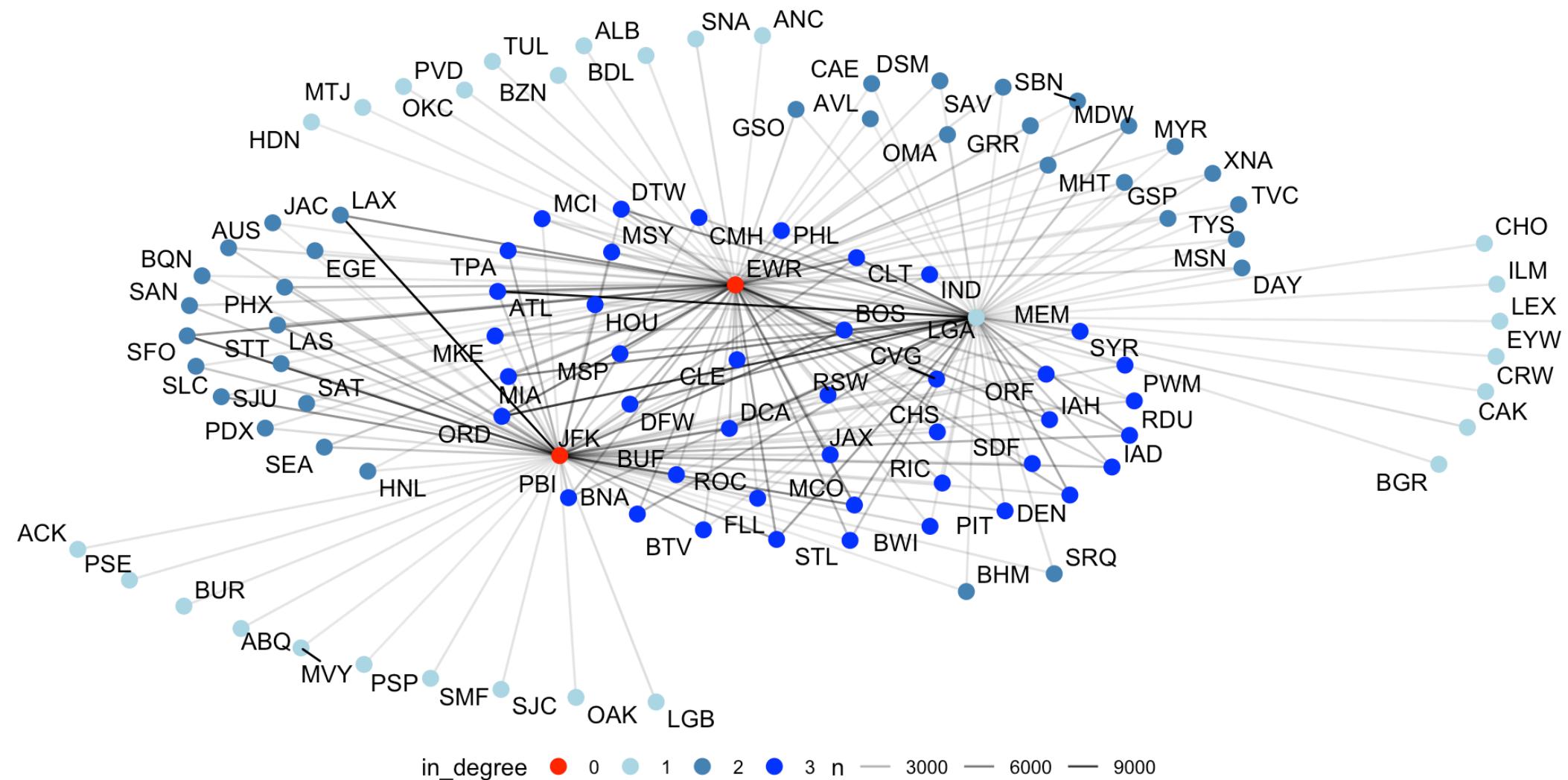


nycflights Data: Labels (1)

```
1 # add labels: repel
2 gr %>%
3   activate(nodes) %>%
4   mutate(in_degree = as.factor(centrality_degree(mode = "in")))) %>%
5   ggraph(layout = "kk") +
6   geom_edge_link(aes(alpha = n)) +
7   geom_node_point(aes(color = in_degree), size = 3) +
8   scale_color_manual(values = c("0" = "red", "1" = "lightblue", "2" = "ste
9 theme_void() +
10 theme(legend.position = "bottom") +
11 geom_node_text(aes(label = name), repel = T)
```

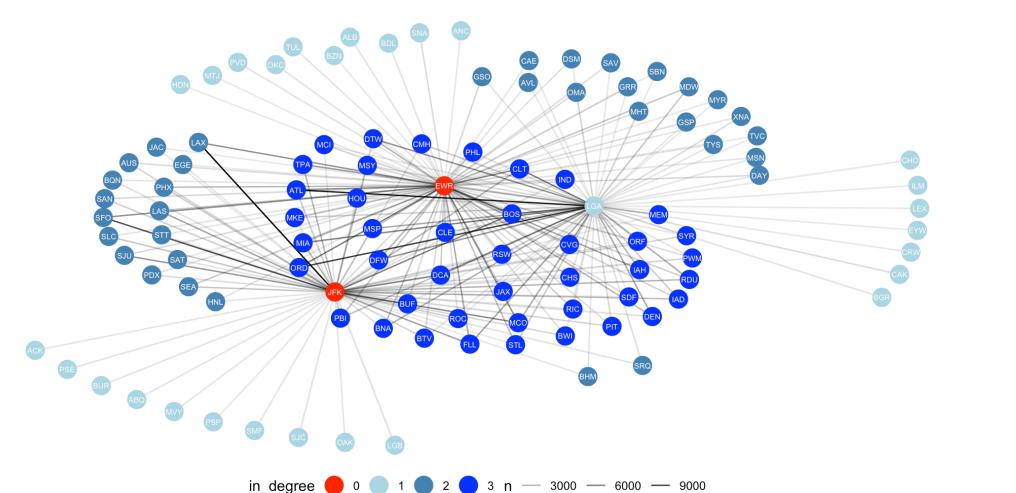


nycflights Data: Labels (1)

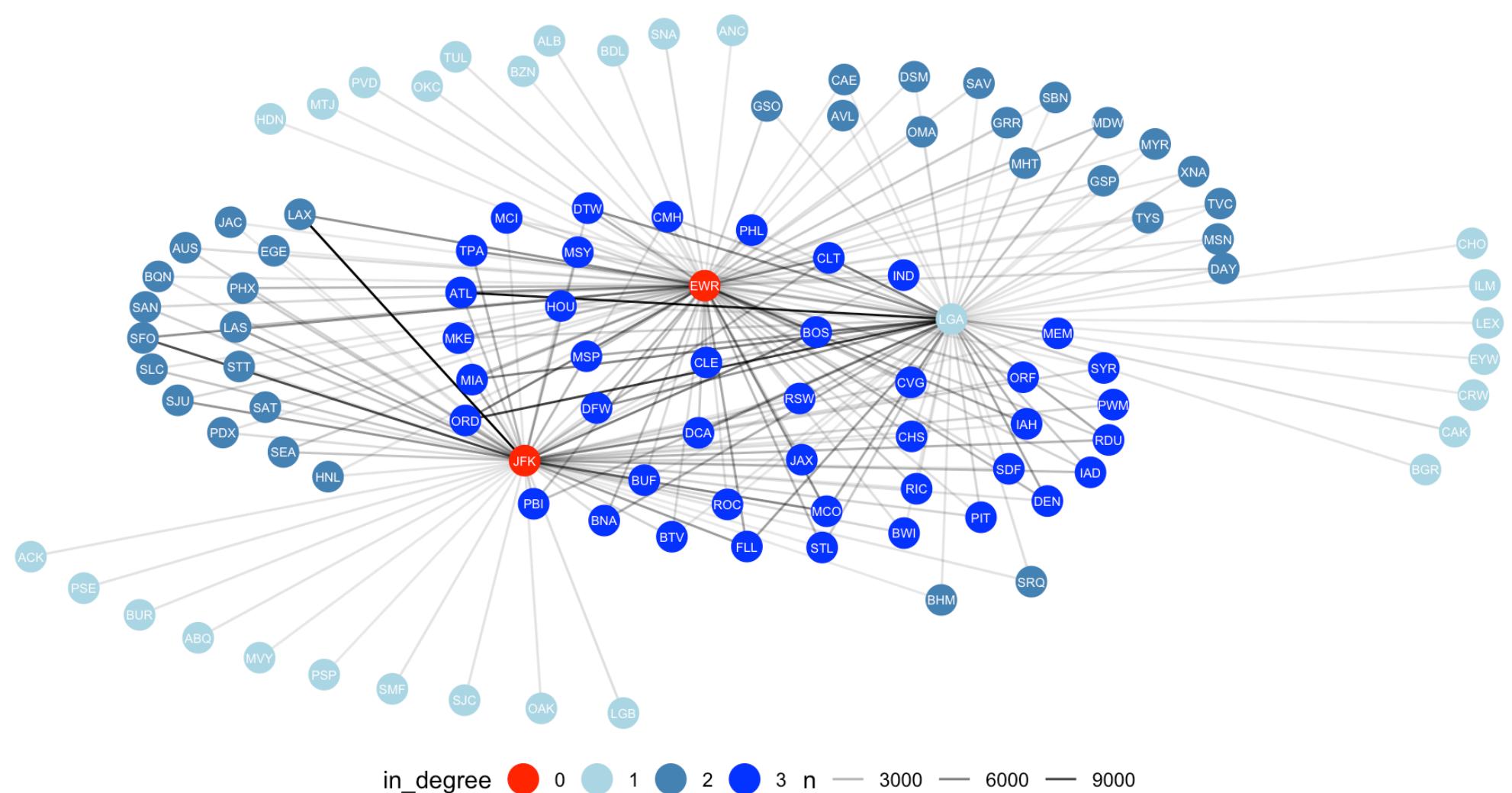


nycflights Data: Labels (2)

```
1 # add labels: inside points
2 gr %>%
3   activate(nodes) %>%
4   mutate(in_degree = as.factor(centrality_degree(mode = "in")))) %>%
5   ggraph(layout = "kk") +
6   geom_edge_link(aes(alpha = n)) +
7   geom_node_point(aes(color = in_degree), size = 6) +
8   scale_color_manual(values = c("0" = "red", "1" = "lightblue", "2" = "steelblue", "3" = "darkblue")) +
9   theme_void() +
10  theme(legend.position = "bottom") +
11  geom_node_text(aes(label = name), size = 2, color = "white")
```

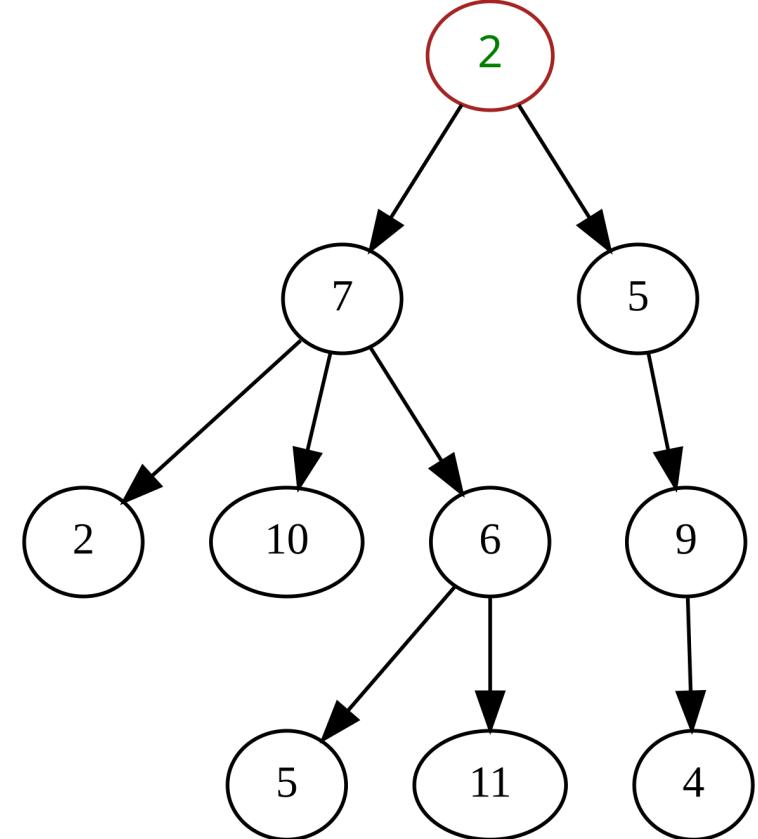


nycflights Data: Labels (2)



Trees in tidygraph

- A tree is a data structure used to represent hierarchical relationships.
- Nodes connected by edges, and each node in the tree contains data and is connected to other nodes in a parent-child relationship.



[https://en.wikipedia.org/wiki/Tree_\(graph_theory\)](https://en.wikipedia.org/wiki/Tree_(graph_theory))

Components of a Tree

- Node: Represents an element or data point in the tree.
- Edge: Represents the connection between two nodes.
- Root: The topmost node in a tree, with no parent. All other nodes are descendants of this root.
- Parent: A node that has one or more children.
- Child: A node that is a descendant of another node.
- Leaf: A node that does not have any children.
- Subtree: A tree formed by a node and its descendants.

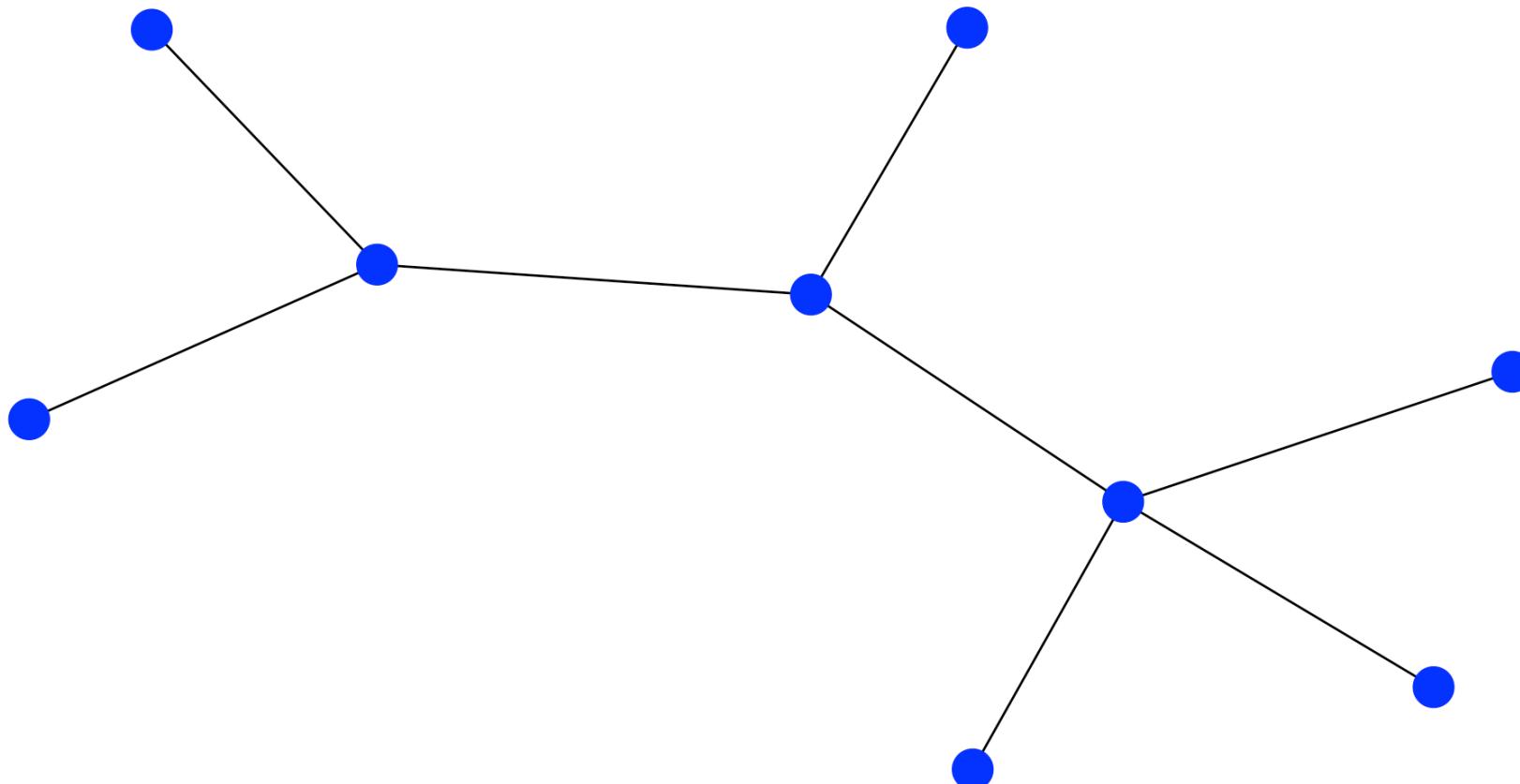
Trees in tidygraph

```
1 # Create a sample hierarchical graph
2 graph <- data.frame(from = c("A", "A", "A", "B", "B", "C", "C", "C"),
3                      to = c("B", "C", "I", "D", "E", "F", "G", "H"))
4 graph <- as_tbl_graph(graph)
5 print(graph)
```

```
# A tbl_graph: 9 nodes and 8 edges
#
# A rooted tree
#
# Node Data: 9 × 1 (active)
  name
  <chr>
1 A
2 B
3 C
4 I
5 D
6 E
7 F
8 G
9 --
```

Trees in tidygraph: kk layout

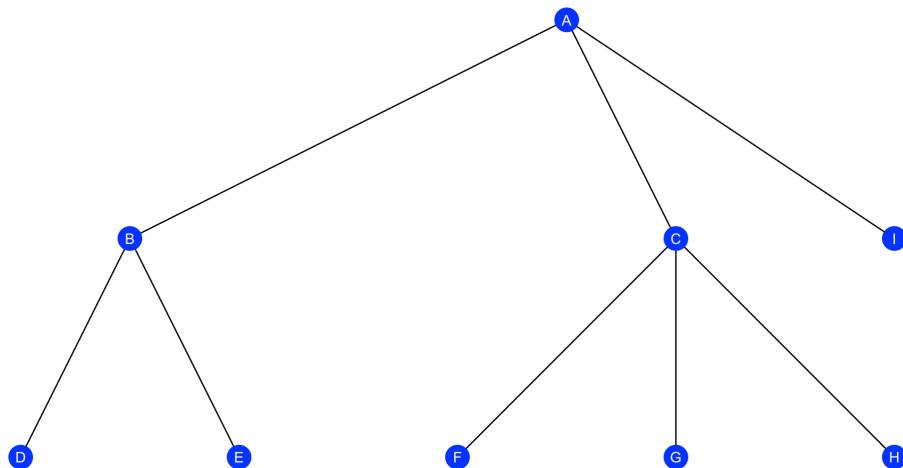
```
1 ggraph(graph, layout = "kk") +  
2   geom_edge_link() +  
3   geom_node_point(size = 8, color = 'blue') +  
4   theme_void()
```



Creating a Tree Layout

dendrogram layout arranges nodes and edges linearly in a hierarchical structure. The layout positions nodes based on their depth in the tree.

```
1 ggraph(graph, layout = 'tree') +  
2   geom_edge_link() +  
3   geom_node_point(size = 8, color = 'blue') +  
4   geom_node_text(aes(label = name), color = "white") +  
5   theme_void()
```



Creating a Dendrogram Layout

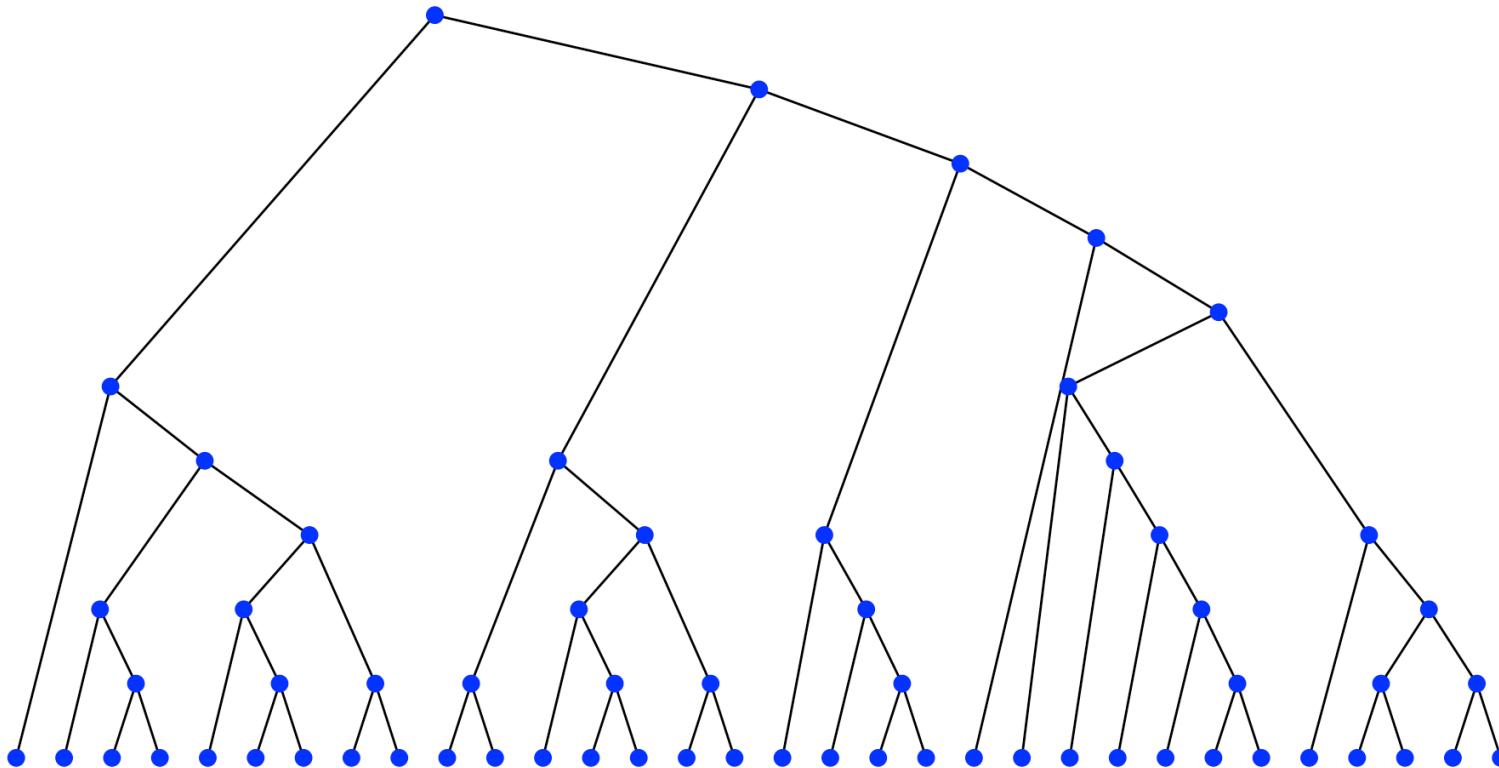
Dendrogram layouts are also very nice to show the output of hierarchical clustering

```
1 head(mtcars)
```

		mpg	cyl	disp	hp	drat	wt	qsec	vs	am	gear	carb
Mazda	RX4	21.0	6	160	110	3.90	2.620	16.46	0	1	4	4
Mazda	RX4 Wag	21.0	6	160	110	3.90	2.875	17.02	0	1	4	4
Datsun	710	22.8	4	108	93	3.85	2.320	18.61	1	1	4	1
Hornet	4 Drive	21.4	6	258	110	3.08	3.215	19.44	1	0	3	1
Hornet	Sportabout	18.7	8	360	175	3.15	3.440	17.02	0	0	3	2
Valiant		18.1	6	225	105	2.76	3.460	20.22	1	0	3	1

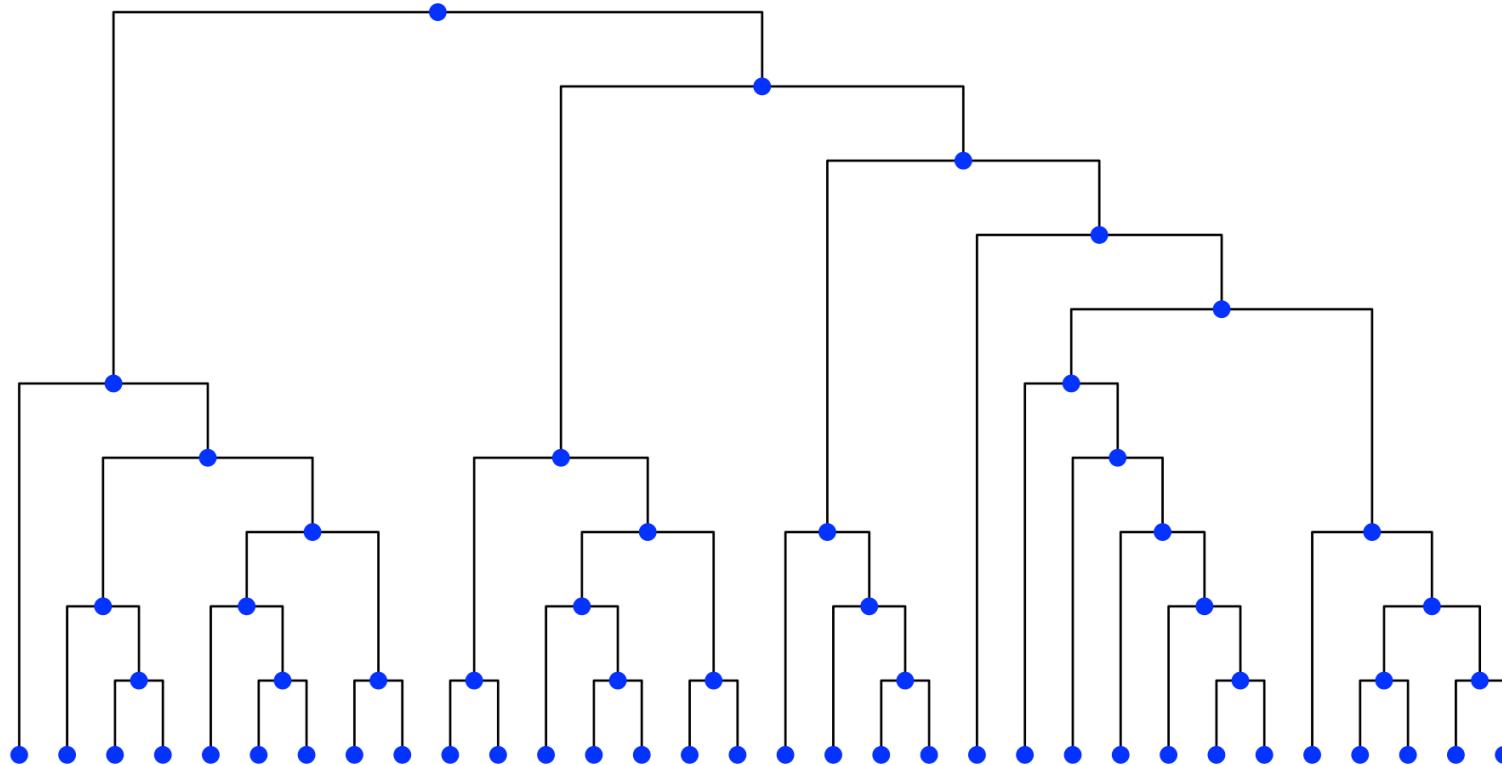
Creating a Dendrogram Layout

```
1 hclust_gr = as_tbl_graph(hclust(dist(mtcars)))
2 ggraph(hclust_gr, layout = 'dendrogram') +
3   geom_edge_link() +
4   geom_node_point(size = 3, color = 'blue') +
5   theme_void()
```



Creating a Dendrogram Layout: Edges

```
1 hclust_gr = as_tbl_graph(hclust(dist(mtcars)))
2 ggraph(hclust_gr, layout = 'dendrogram') +
3   geom_edge_elbow() +
4   geom_node_point(size = 3, color = 'blue') +
5   theme_void()
```



Icicle Plots

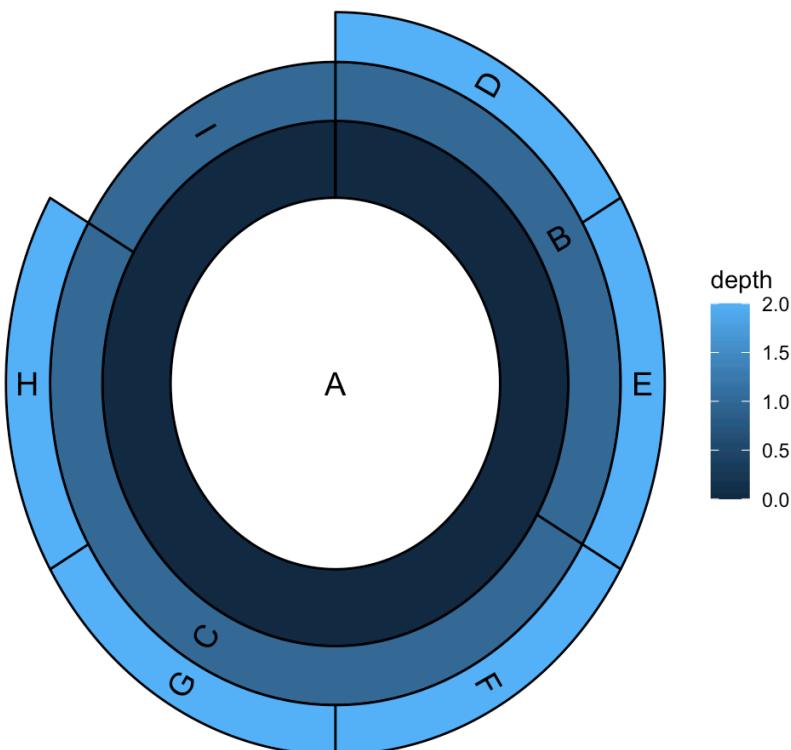
Nodes are represented as rectangles, with the area or position indicating relationships.

```
1 ggraph(graph, layout = 'partition') +  
2   geom_node_tile(aes(fill = depth)) +  
3   geom_node_text(aes(label = name), color = "white", size = 8) +  
4   theme_void()
```



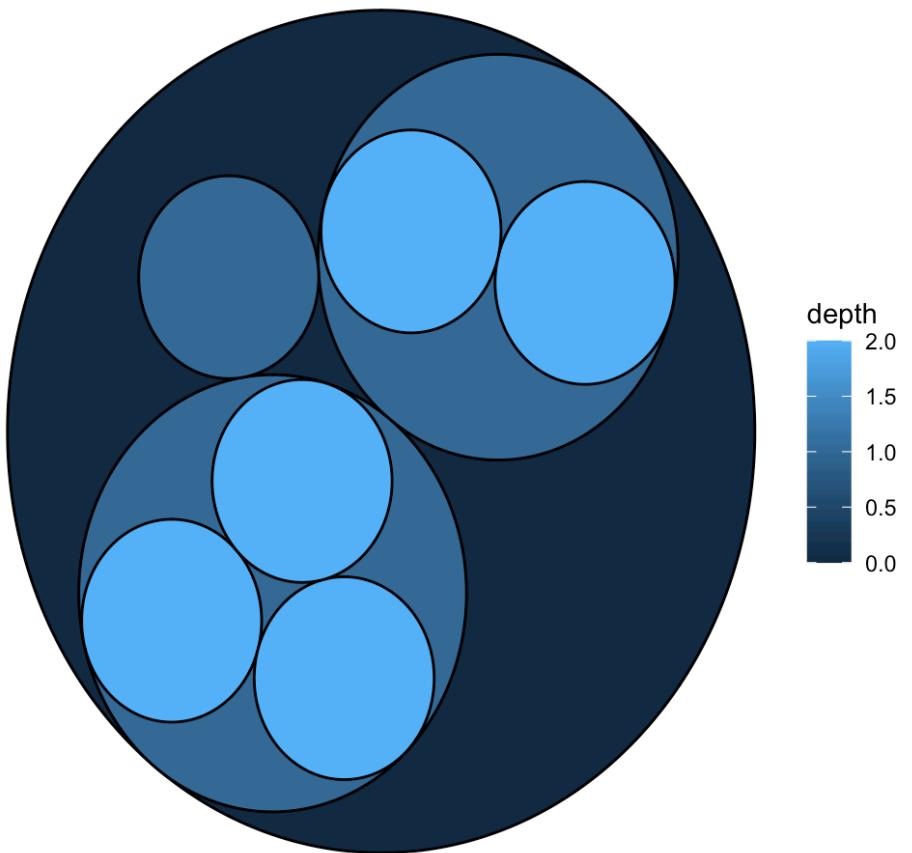
Sunburst Plots

```
1 ggraph(graph, layout = 'partition', circular = TRUE) +
2   geom_node_arc_bar(aes(fill = depth)) +
3   geom_node_text(aes(label = name, angle = node_angle(x, y)),
4                 size = 5) +
5   theme_void()
```



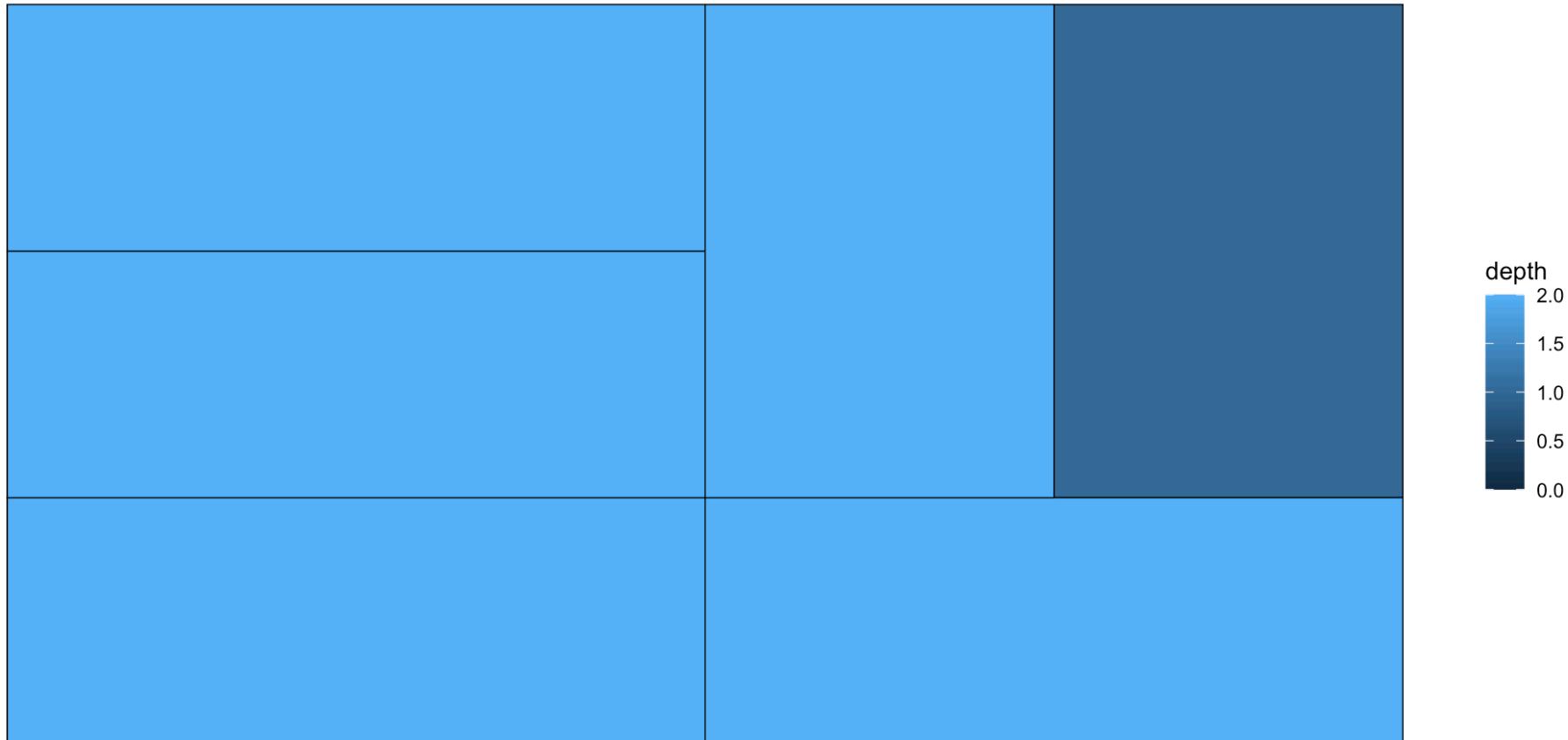
Circle Packs

```
1 ggraph(graph, layout = 'circlepack') +  
2   geom_node_circle(aes(fill = depth)) +  
3   theme_void()
```



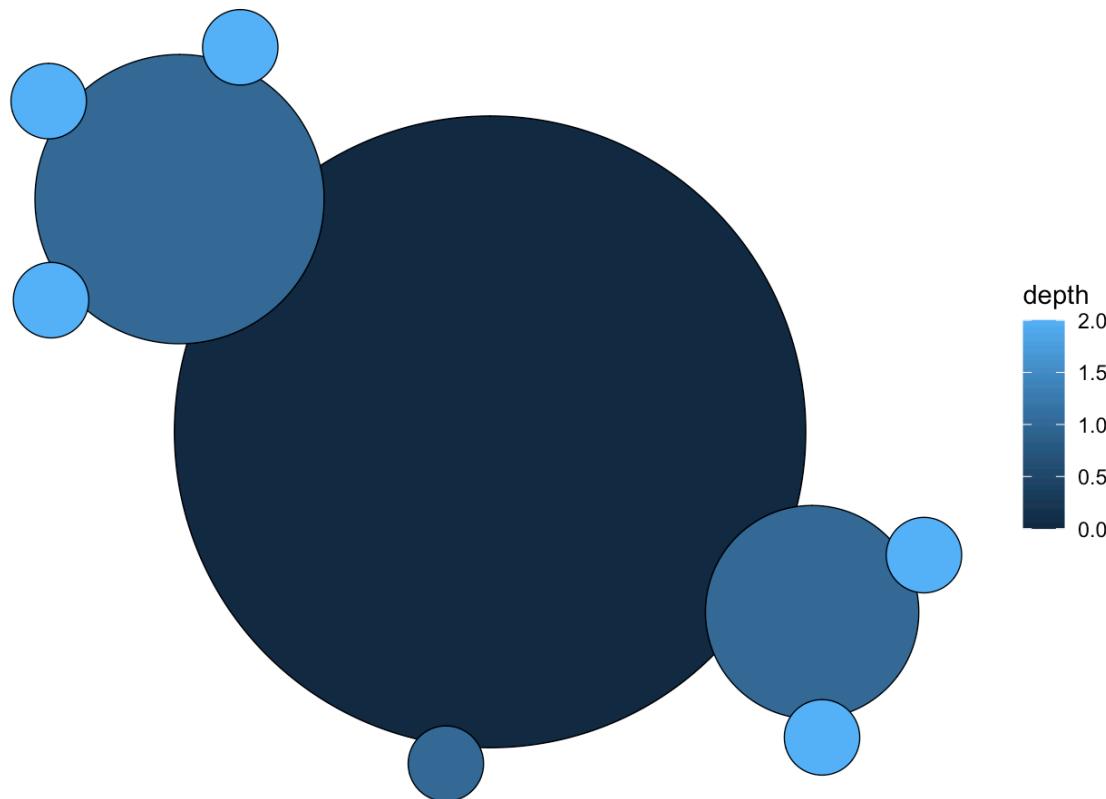
Tree Maps

```
1 ggraph(graph, layout = 'treemap') +  
2   geom_node_tile(aes(fill = depth), size = 0.25) +  
3   theme_void()
```

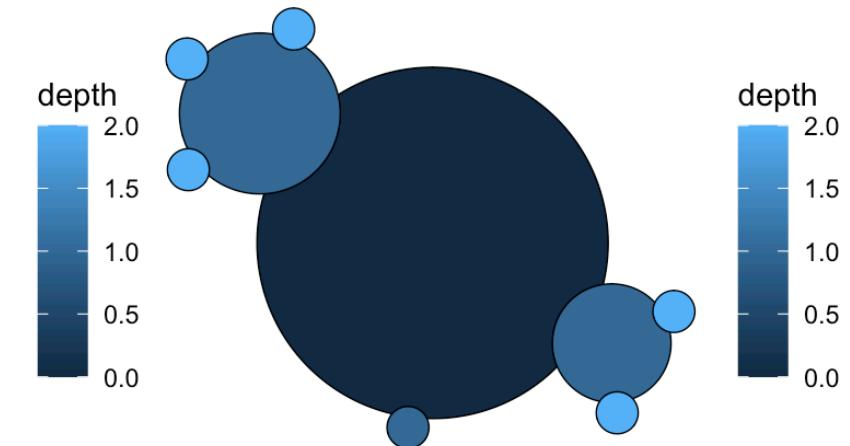
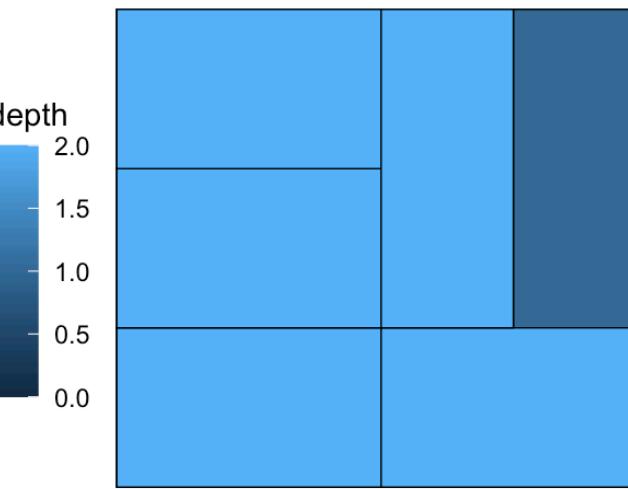
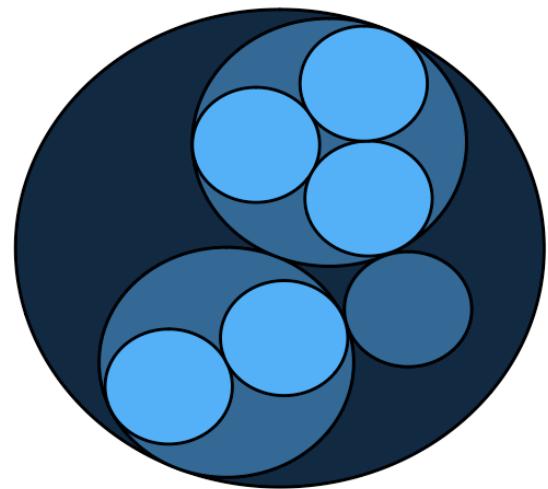
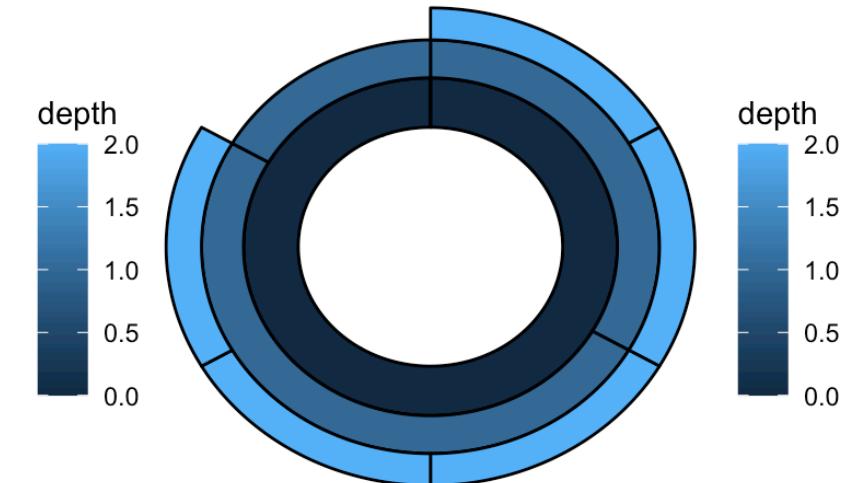
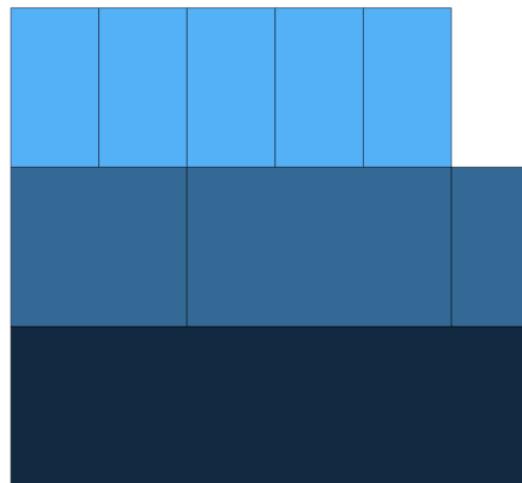
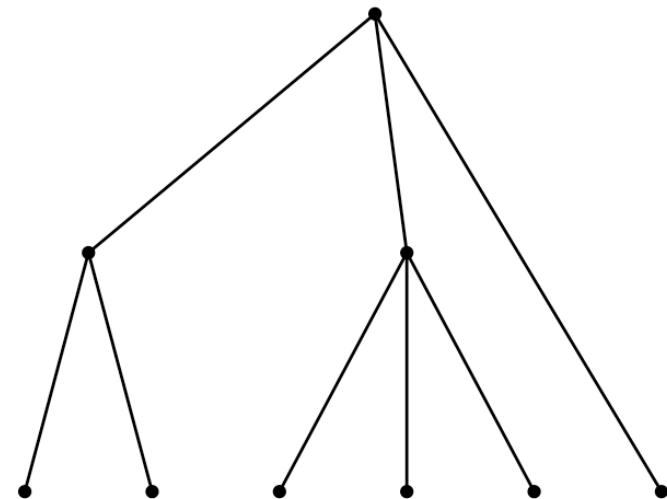


Cactustree

```
1 ggraph(graph, layout = 'cactustree') +  
2   geom_node_circle(aes(fill = depth), size = 0.25) +  
3   coord_fixed() +  
4   theme_void()
```



A Comparison



Thank You!

Email: naman.agr03@gmail.com

