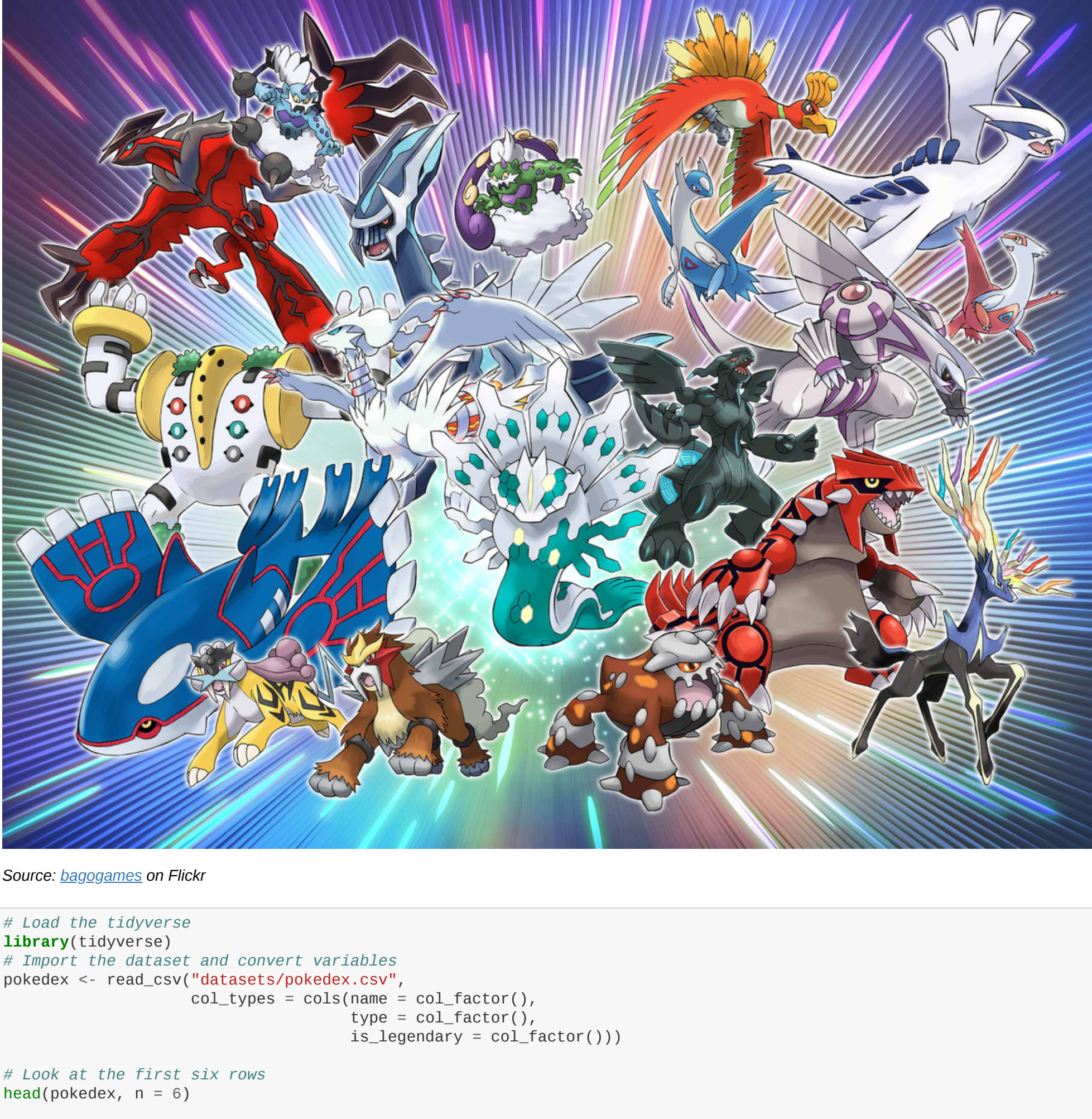


1. Introduction

In the world of Pokémon academia, one name towers above any other – Professor Samuel Oak. While his colleague Professor Elm specializes in Pokémon evolution, Oak has dedicated his career to understanding the relationship between Pokémon and their human trainers. A former trainer himself, the professor has first-hand experience of how obscure Pokémon can be – particularly when they hold legendary status.

For his latest research project, Professor Oak has decided to investigate the defining characteristics of legendary Pokémon to improve our understanding of their temperament. Hearing of our expertise in classification problems, he has enlisted us as the lead researchers.

Our journey begins at the professor's research lab in Pallet Town, Kanto. The first step is to open up the Pokédex, an encyclopaedic guide to 801 Pokémon from all seven generations.



Source: [@asapovitski](#) on Flickr

```
In [3]: # Load the tidyverse
library(tidyverse)

# Import the dataset and convert variables
pokedex <- read_csv("data/pokedex.csv")
col_types <- cols(name = col_factor(),
                  type = col_factor(),
                  is_legendary = col_factor())

# Look at the first six rows
head(pokedex, n = 6)

# Examine the structure
str(pokedex)

## Attaching packages: ..... tidyverse 1.3.0 ##
# vplot2 2.2.1    vpurrr 0.3.3
# tidyr 1.1.3    v dplyr 0.8.3
# lubridate 1.0.9 v stringr 1.4.0
# readr 1.3.1    v forcats 0.4.0
# Conflicts: ..... tidyverse_conflicts() ..
# dplyr::filter() masks stats::filter()
# dplyr::lag()   masks stats::lag()

pokedex_number | name      | attack | defense | height_m | hp | percentage_male | sp_attack | sp_defense | speed | type | weight_kg | generation | is_
1 | Bulbasaur | 49      | 49      | 0.7       | 45 | 86.1            | 65        | 65         | 45    | grass | 6.9       | 1          | 0
2 | Ivysaur  | 62      | 63      | 1.0       | 60 | 88.1            | 80        | 80         | 60    | grass | 13.0      | 1          | 0
3 | Venusaur | 100     | 123     | 2.0       | 80 | 88.1            | 122       | 120        | 80    | grass | 100.0     | 1          | 0
4 | Charmander | 52     | 43      | 0.6       | 39 | 86.1            | 60        | 50         | 60    | fire  | 8.5       | 1          | 0
5 | Charmeleon | 58    | 58      | 1.1       | 38 | 86.1            | 60        | 65         | 60    | fire  | 19.0     | 1          | 0
6 | Charizard | 104     | 78      | 1.7       | 78 | 86.1            | 150       | 115        | 100   | fire  | 90.5     | 1          | 0

Classes: 'spec_tbl_df', 'tbl_df', 'tbl' and 'data.frame':   801 obs. of  14 variables:
 $ pokedex_number: num  1 2 3 4 5 6 7 8 9 10 ...
 $ name          : chr  "Bulbasaur", "Ivysaur", "Venusaur", ...
 $ attack       : num  49 62 100 52 64 104 48 43 103 50 ...
 $ defense      : num  49 62 100 52 64 104 48 43 103 50 ...
 $ height_m     : num  0.7 1.0 2.0 0.6 1.1 1.7 0.5 1.1 0.3 ...
 $ hp           : num  45 60 80 39 68 78 44 59 79 45 ...
 $ percentage_male: num  86.1 88.1 88.1 86.1 88.1 88.1 88.1 88.1 88.1 88.1 59 ...
 $ sp_attack    : num  65 80 122 60 88 150 59 65 135 20 ...
 $ sp_defense   : num  65 80 120 50 65 135 64 60 135 20 ...
 $ speed        : num  45 60 80 65 80 100 43 58 78 45 ...
 $ type         : chr  "grass", "grass", "fire", ...
 $ weight_kg    : num  6.9 13 19 8.5 15 90.5 22.5 85.5 2.9 ...
 $ generation   : num  1 1 1 1 1 1 1 1 1 ...
 $ is_legendary : factor w/ 2 levels "0","1": 1 1 1 1 1 1 1 1 1 1 ...
 #> #> #>
.. cols:
.. pokedex_number = col_double(),
.. name = col_factor(levels = NULL, ordered = FALSE, include_na = FALSE),
.. attack = col_double(),
.. defense = col_double(),
.. height_m = col_double(),
.. hp = col_double(),
.. percentage_male = col_double(),
.. sp_attack = col_double(),
.. sp_defense = col_double(),
.. speed = col_double(),
.. type = col_factor(levels = NULL, ordered = FALSE, include_na = FALSE),
.. weight_kg = col_double(),
.. generation = col_double(),
.. is_legendary = col_factor(levels = NULL, ordered = FALSE, include_na = FALSE)
.. )
```

2. How many Pokémon are legendary?

After browsing the Pokédex, we can see several variables that could feasibly explain what makes a Pokémon legendary. We have a series of numerical (stat) stats – attack, defense, speed, etc. – as well as a categorization of Pokémon type (bug, dark, dragon, etc.). `is_legendary` is the binary classification variable we will eventually be predicting, tagged 1 if a Pokémon is legendary and 0 if it is not.

Before we explore these variables in any depth, let's first out how many Pokémon are legendary out of the 801 total, using the handy `count()` function from the `dplyr` package.

```
In [2]: # Prepare the data
legendary_pokemon <- pokedex %>%
  count(is_legendary) %>%
  mutate(prop = n / sum(n))

# Print the data frame
legendary_pokemon

#> #> #>
  is_legendary  n      prop
1             0 731 0.9126024
2             1   70  0.0873976
```

3. Legendary Pokémon by height and weight

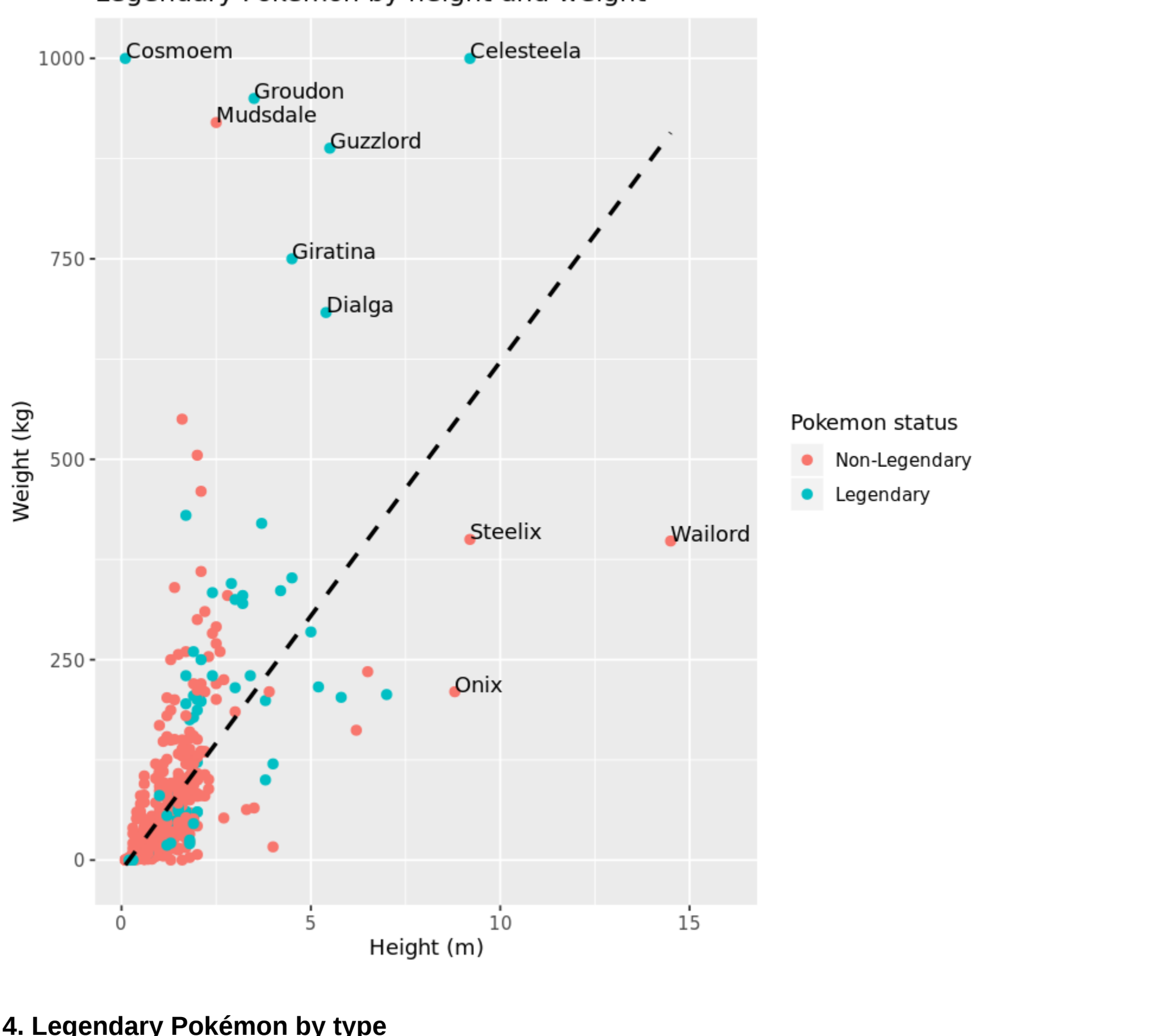
We now know that there are 70 legendary Pokémon – a scarce minority at 9% of the population. Let's start to explore some of their distinguishing characteristics.

First of all, we'll plot the relationship between `height_m` and `weight_kg` for all 801 Pokémon, highlighting those that are classified as legendary. We'll also add conditional labels to the plot, which will only print a Pokémon's name if it is taller than 7.5m or heavier than 600kg.

```
In [3]: # Prepare the plot
legend_by_heightweight_plot <- pokedex %>%
  geom_point(aes(color = is_legendary, size = 2)) +
  geom_text(aes(label = ifelse(height_m > 7.5 | weight_kg > 600, as.character(name), '')),
            vjust = 0, hjust = 8) +
  facet_wrap(~ is_legendary, col = "black", linetype = "dashed") +
  expand_limits(x = 10) +
  labs(title = "Legendary Pokemon by height and weight",
       x = "height (m)",
       y = "weight (kg)") +
  guides(color = guide_legend(title = "Pokemon status")) +
  scale_color_manual(labels = c("Non-Legendary", "legendary"),
                    values = c("grey20", "red2"))

# Print the plot
legend_by_heightweight_plot

#> #> #>
Warning message:
'Removed 28 rows containing non-finite values (stat.smooth).':Warning message:
'Removed 28 rows containing missing values (geom.point).':Warning message:
'Removed 28 rows containing missing values (geom.text).':
```



4. Legendary Pokémon by type

It seems that legendary Pokémon are generally heavier and taller, but with many exceptions. For example, Onix (Gen 1), Steelix (Gen 2) and Wailord (Gen 3) are all extremely tall, but none of them have legendary status. There must be other factors at play.

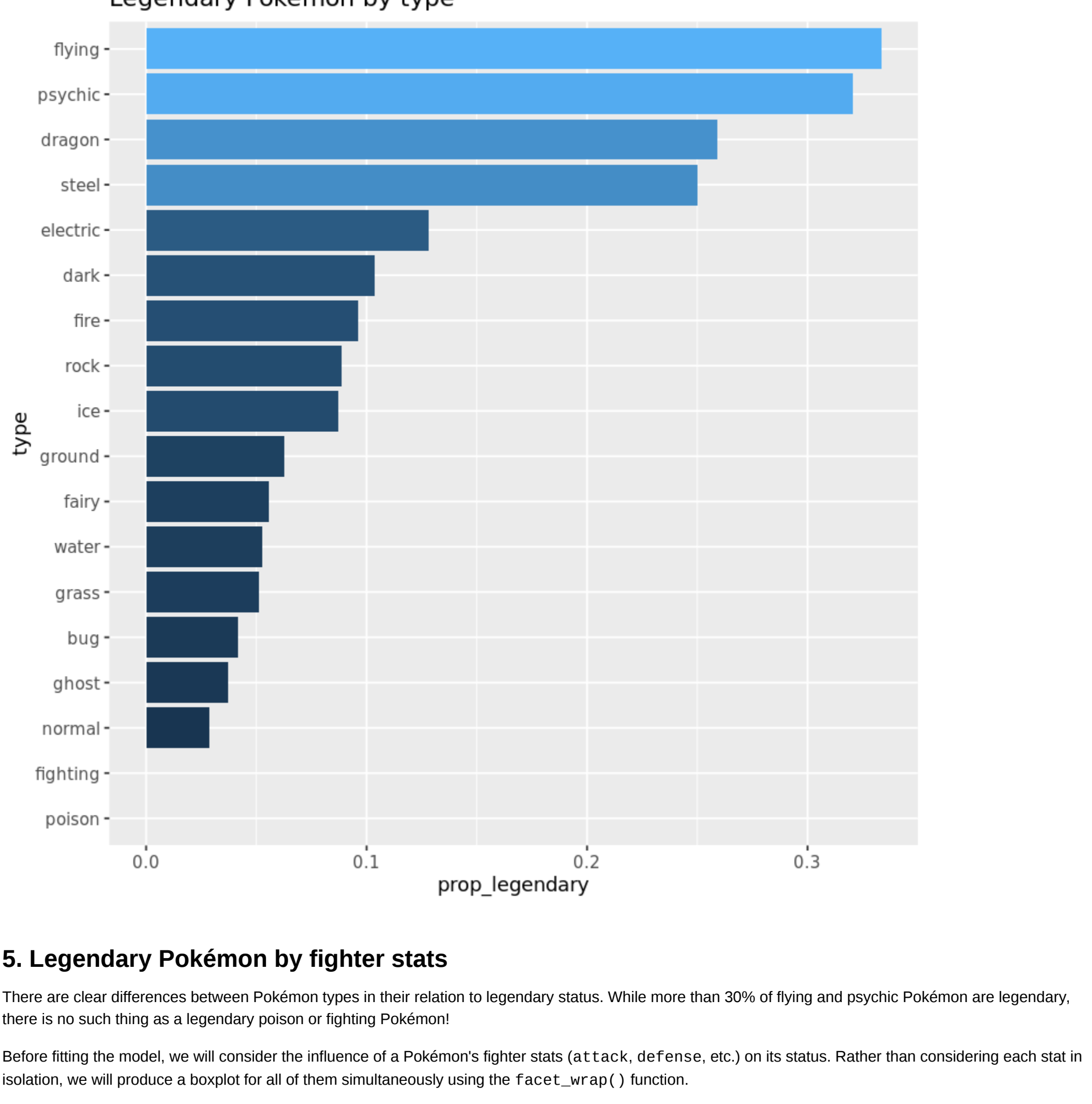
We will now look at the effect of a Pokémon's type on its legendarity/legendary classification. There are 18 possible types, ranging from the common (Grass / Normal / Water) to the rare (Fairy / Flying / Ice). We will calculate the proportion of legendary Pokémon within each category, and then plot these proportions using a simple bar chart.

```
In [4]: # Prepare the data
legend_by_type <- pokedex %>%
  group_by(type) %>%
  mutate(prop_legendary = as.numeric(is_legendary) > 0) %>%
  summarise(prop_legendary = mean(is_legendary) %>%
    ungroup()) %>%
  mutate(type = fct_reorder(type, prop_legendary))

# Prepare the plot
legend_by_type_plot <- legend_by_type %>%
  ggplot(aes(x = type, y = prop_legendary, fill = prop_legendary)) +
  geom_col() +
  labs(title = "Legendary Pokemon by type",
       coord_flip() = FALSE) +
  guides(fill = FALSE)

# Print the plot
legend_by_type_plot

#> #> #>
```



5. Legendary Pokémon by fighter stats

There are clear differences between Pokémon types in their relation to legendary status. While more than 30% of flying and psychic Pokémon are legendary, there is not such thing as a legendary poison or fighting Pokémon.

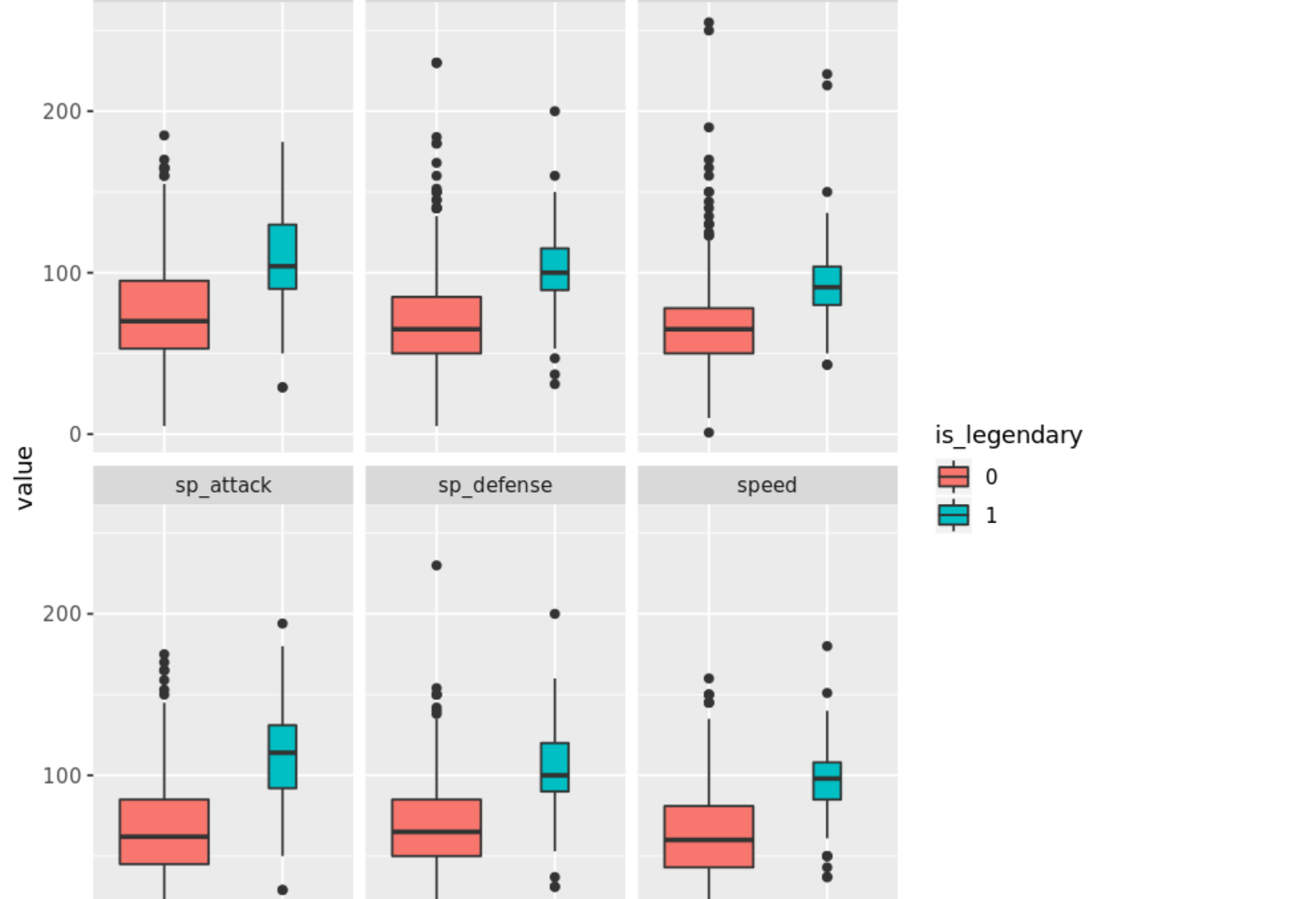
Before fitting the model, we will consider the influence of a Pokémon's fighter stats (attack, defense, etc.) on its status. Rather than considering each stat in isolation, we will produce a boxplot for all of them simultaneously using the `Facet_wrap()` function.

```
In [5]: # Prepare the data
legend_by_stats <- pokedex %>%
  select(is_legendary, attack, sp_attack, defense, sp_defense, hp, speed) %>%
  gather(key = fct_stat, value = value, is_legendary)

# Prepare the plot
legend_by_stats_plot <- legend_by_stats %>%
  ggplot(aes(x = is_legendary, y = value, fill = is_legendary)) +
  geom_boxplot(varwidth = TRUE) +
  facet_wrap(~ fct_stat) +
  labs(title = "Pokemon fighter statistics",
       x = "legendary status",
       guides(fill = "legend"))

# Print the plot
legend_by_stats_plot

#> #> #>
```



6. Create a training/test split

As we might expect, legendary Pokémon outshine their ordinary counterparts in all fighter stats. Although we haven't formally tested a difference in means, the boxplots suggest a significant difference with respect to all six variables. Nonetheless, there are a number of outliers in each case, meaning that some legendary Pokémon are exceptional.

We have now explored all of the predictor variables we will use to explain what makes a Pokémon legendary. Before fitting our model, we will split the pokedex into a training set (`pokedex_train`) and a test set (`pokedex_test`). This will allow us to test the model on unseen data.

```
In [6]: # Set seed for reproducibility
set.seed(1234)

# Save number of rows in dataset
n <- nrow(pokedex)

# Generate 68% sample of rows
sample_rows <- sample(n, n * 0.68)

# Create training set
pokedex_train <- pokedex %>%
  filter(row_number() %in% sample_rows)

# Create test set
pokedex_test <- pokedex %>%
  filter(row_number() %in% sample_rows)
```

7. Fit a decision tree

Now we have our training and test sets, we can go about building our classifier. But before we fit a random forest, we will fit a simple **classification decision tree**. This will give us a baseline to against which to compare the results of the random forest, as well as an informative graphical representation of the model.

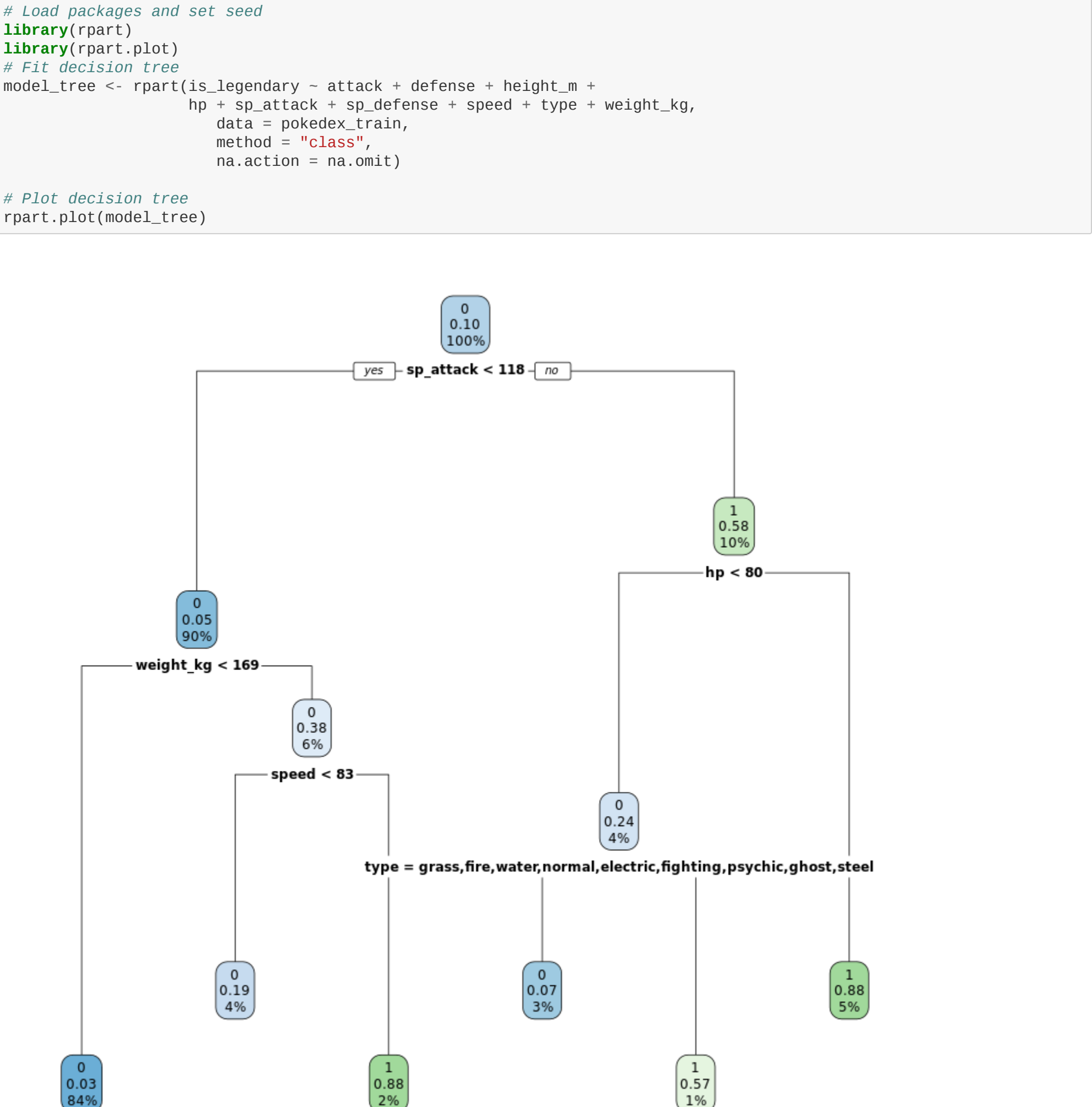
Here, and also in the random forest, we will omit incomplete observations by setting the `na.action` argument to `na.omit`. This will remove a few Pokémon with missing values for `height_m` and `weight_kg` from the training set.

```
In [7]: # Load packages and set seed
library(rpart)
library(rpart.plot)

# Fit decision tree
model_tree <- rpart(is_legendary ~ attack + defense + height_m +
                    hp + sp_attack + sp_defense + speed + type + weight_kg,
                    data = pokedex_train,
                    method = "class",
                    na.action = na.omit)

# Plot decision tree
rpart.plot(model_tree)

#> #> #>
```



8. Fit a random forest

Each node of the tree shows the predicted class, the probability of being legendary, and the percentage of Pokémon in that node. The bottom-left node, for example – for those with `sp_attack < 118` and `weight_kg < 169` – represents 84% of Pokémon in the training set, predicting that each only has a 9% chance of being legendary.

Decision trees place the most important variables at the top and exclude any they don't find to be useful. In this case, `sp_attack` occupies node 1 while `attack`, `defense`, `sp_defense` and `height_m` are all excluded.

However, decision trees are unstable and sensitive to small variations in the data. It therefore makes sense to fit a **random forest** – an ensemble method that averages over several decision trees all at once. This should give us a more robust model that classifies Pokémon with greater accuracy.

```
In [8]: # Load package and set seed
library(randomForest)
set.seed(1234)

# Fit random forest
model_forest <- randomForest(is_legendary ~ attack + defense + height_m +
                             hp + sp_attack + sp_defense + speed + type + weight_kg,
                             data = pokedex_train,
                             importance = TRUE,
                             na.action = na.omit)

# Print model output
model_forest

#> #> #>
randomForest 4.6-14
Type rfNew() to see new Features/changes/bug fixes.
Attaching package: 'randomForest'

The following object is masked from 'package:dplyr':
  combine

The following object is masked from 'package:ggplot2':
  margin

Call:
randomForest(formula = is_legendary ~ attack + defense + height_m + hp + sp_attack + sp_defense + speed + type + weight_kg, data = pokedex_train, importance = TRUE, na.action = na.omit)
Type of random forest: classification
Number of trees: 588
No. of variables tried at each split: 3

OOB estimate of error rate: 7.85%
Confusion matrix:
  0 1 class.error
0 411 9 8.6242857
1 24 24 0.9998989
```

9. Assess model fit

Looking at the model output, we can see that the random forest has an out-of-bag (OOB) error of 7.49%, which isn't bad by most accounts. However, since there are **24 true positives** and **24 false negatives**, the model only has a recall of 50%, which means that it struggles to successfully retrieve every legendary Pokémon in the dataset.

In order to allow direct comparison with the decision tree, we will plot the **ROC curves** for both models using the `ROC` package, which will visualize their true positive rate (TPR) and false positive rate (FPR) respectively. The closer the curve is to the top left of the plot, the higher the area under the curve (AUC) and the better the model.

```
In [9]: # Load the ROC package
library(ROC)

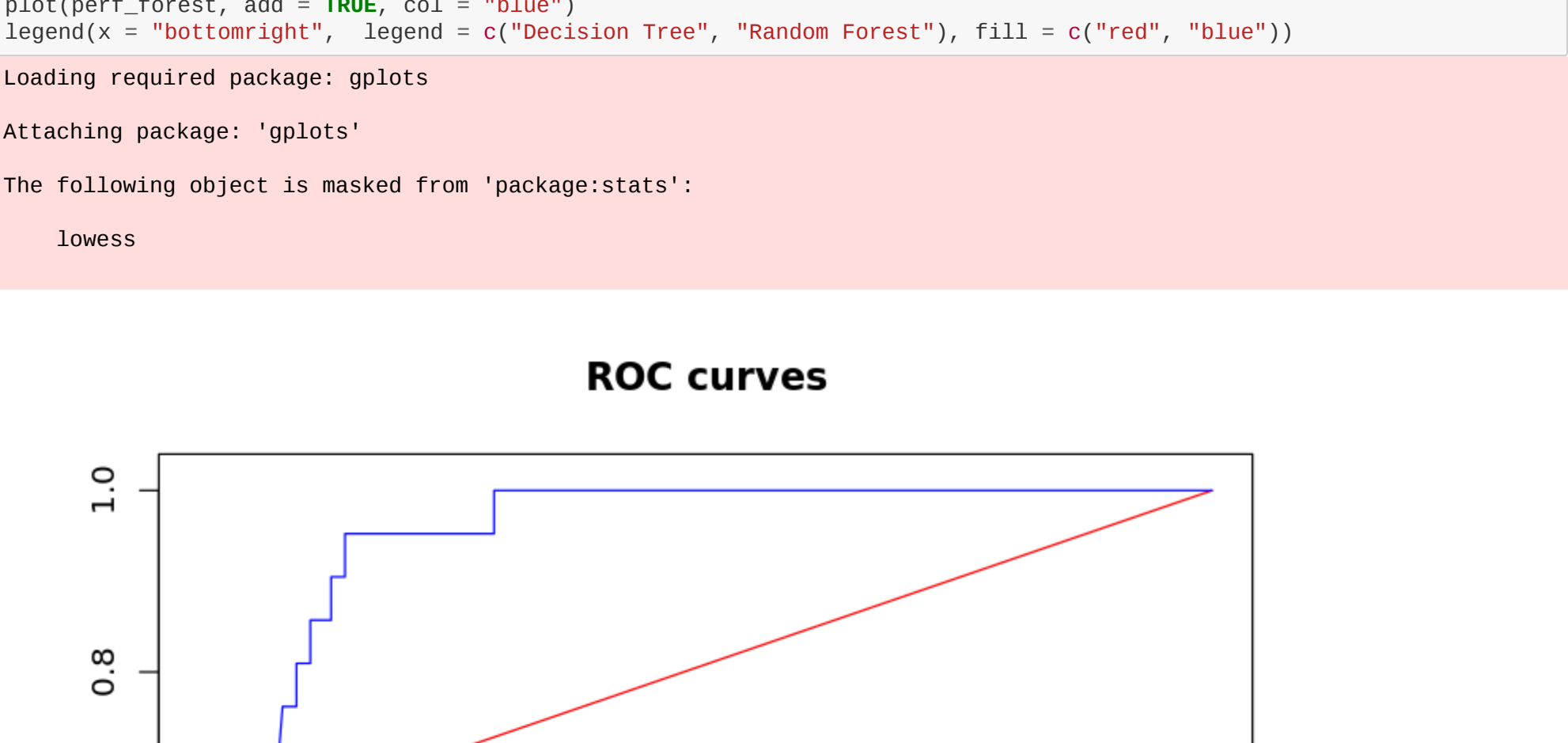
# Create prediction and performance objects for the decision tree
pred_tree <- predict(model_tree, pokedex_test, type = "prob")
perf_tree <- prediction(perf_tree[,2], pokedex_test$is_legendary)
perf_tree <- performance(pred_tree, "tpr", "fpr")

# Create prediction and performance objects for the random forest
pred_forest <- predict(model_forest, pokedex_test, type = "prob")
perf_forest <- prediction(perf_forest[,2], pokedex_test$is_legendary)
perf_forest <- performance(pred_forest, "tpr", "fpr")

# Plot the ROC curves: first for the decision tree, then for the random forest
plot(perf_tree, col = "red", main = "ROC curves")
plot(perf_forest, add = TRUE, col = "blue")
legend(x = "bottomright", legend = c("Decision Tree", "Random Forest"), fill = c("red", "blue"))

Loading required package: gplots
Attaching package: 'gplots'

The following object is masked from 'package:stats':
  lowess
```



10. Analyze variable importance

It's clear from the ROC curves that the random forest is a substantially better model, boasting an AUC (not calculated above) of 91% versus the decision tree's 78%. When calculating variable importance, it makes sense to do so with the best model available, so we'll use the random forest for the final part of our analysis.

Note that the random forest returns two measures of variable importance:

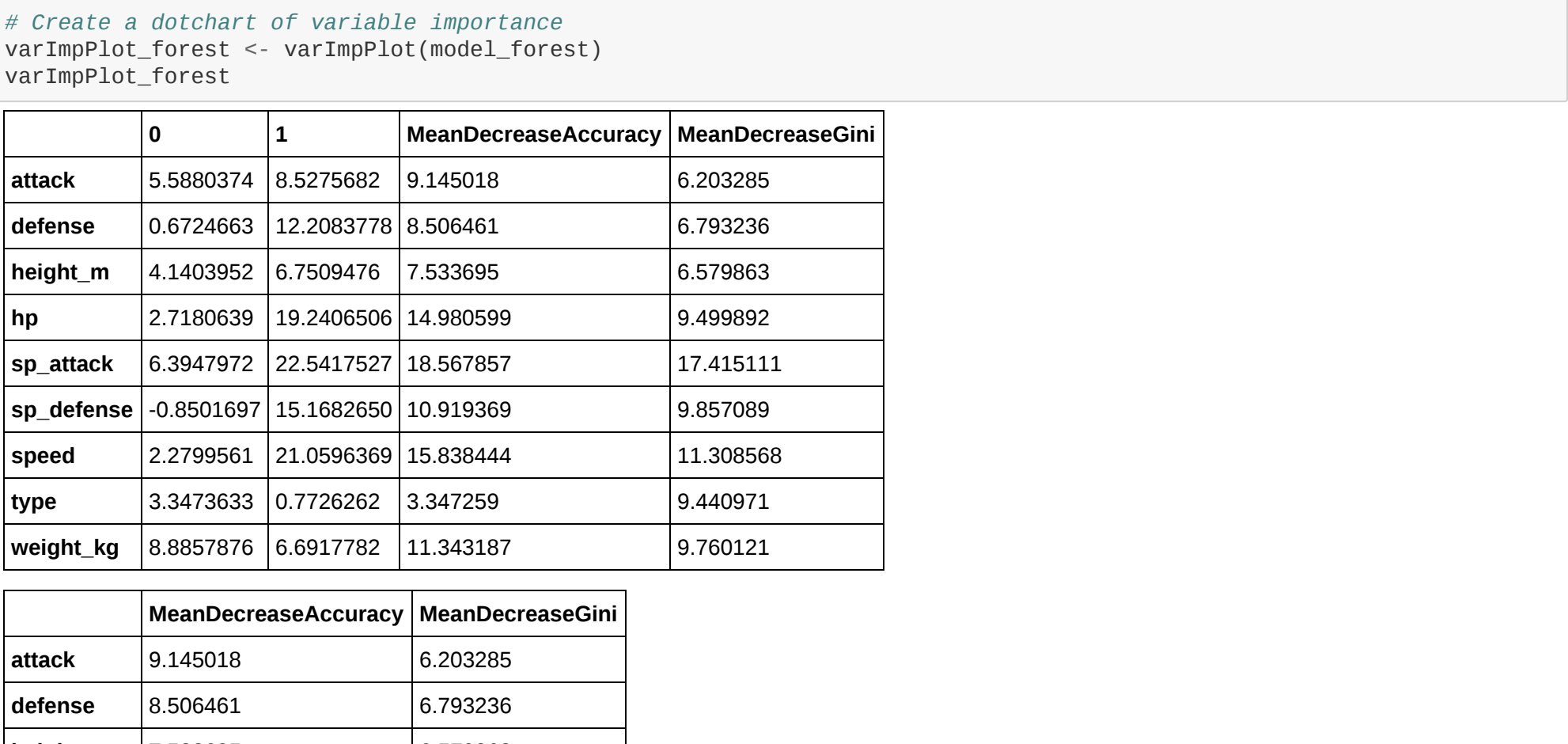
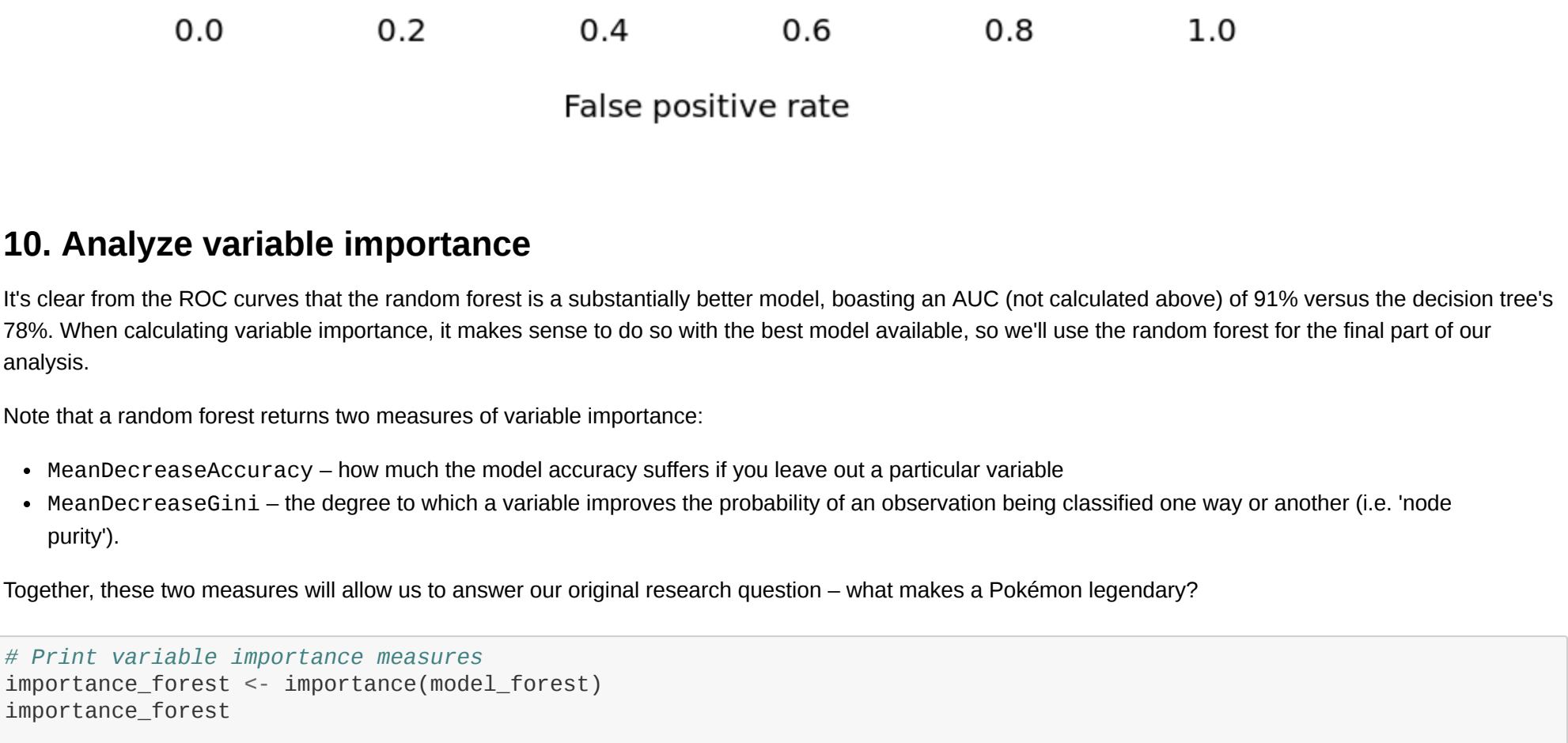
- `MeanDecreaseAccuracy` – how much the model accuracy suffers if you leave out a particular variable
- `MeanDecreaseGini` – the degree to which a variable improves the probability of an observation being classified one way or another (i.e. 'node purity').

Together, these two measures will allow us to answer our original research question – what makes a Pokémon legendary?

```
In [13]: # Print variable importance measures
importance_forest <- importance(model_forest)
importance_forest

# Create a bar chart of variable importance
varimp_plot <- varImpPlot(model_forest)
varimp_plot

#> #> #>
```



11. Conclusion

According to the variable importance plot, `sp_attack` is the most important factor in determining whether or not a Pokémon is legendary, followed by speed. The plot doesn't tell us whether the variables have a positive or a negative effect, but we know from our exploratory analysis that the relationship is generally positive. We therefore conclude that legendary Pokémon are characterized primarily by the power of their special attacks and secondarily by their speediness, while also exhibiting higher fighting abilities across the board.