

STOR 664 — Part 1 Combined EDA (Pokémon)

Team: Yinyu Yao, Kareena Legare, Samuel Moore, Irene Zhang,

2025-11-15

```
# 1) distribution statistics (with box plots)

#Missing values
core_wo_t2_gen <- select(pokemon, -type_2, -generation_id)
n_sparse <- sum(!complete.cases(core_wo_t2_gen))
cat("Number of rows with at least one N.A.(not counting type_2 or generation_id):", n_sparse, "\n")

## Number of rows with at least one N.A.(not counting type_2 or generation_id): 711

empty_type2 <- sum(is.na(pokemon$type_2))
empty_generation_id <- sum(is.na(pokemon$generation_id))
unique_generation_id <- sort(unique(pokemon$generation_id))
cat("Number of rows with missing type_2:", empty_type2, "\n")

## Number of rows with missing type_2: 439

cat("Number of rows with missing generation_id:", empty_generation_id, "\n")

## Number of rows with missing generation_id: 147

cat("Unique generation IDs:", paste(unique_generation_id, collapse = ", "), "\n")

## Unique generation IDs: 1, 2, 3, 4, 5, 6, 7

gen_tab <- as.data.frame(table(pokemon$generation_id, useNA = "ifany"))
names(gen_tab) <- c("generation_id", "count")
type_tab <- as.data.frame(table(pokemon$type_1, useNA = "ifany"))
names(type_tab) <- c("type_1", "count")
save_tbl(gen_tab, "tables/count_by_generation.csv")
```

generation_id	count
1	151
2	100
3	135
4	107
5	156

generation_id	count
6	72
7	81
NA	147

```
save_tbl(type_tab, "tables/count_by_type1.csv")
```

type_1	count
bug	79
dark	37
dragon	39
electric	61
fairy	19
fighting	31
fire	59
flying	4
ghost	40
grass	84
ground	36
ice	29
normal	111
poison	35
psychic	64
rock	65
steel	30
water	126

```
knitr::kable(gen_tab, caption = "Counts of Pokémon by generation_id")
```

Table 3: Counts of Pokémon by generation_id

generation_id	count
1	151
2	100
3	135
4	107
5	156
6	72
7	81
NA	147

```
knitr::kable(type_tab, caption = "Counts of Pokémon by type_1")
```

Table 4: Counts of Pokémon by type_1

type_1	count
bug	79
dark	37
dragon	39
electric	61
fairy	19
fighting	31
fire	59
flying	4
ghost	40
grass	84
ground	36
ice	29
normal	111
poison	35
psychic	64
rock	65
steel	30
water	126

```
# Conditional distribution  $P(\text{type}_2 \mid \text{type}_1) + X^2$  ----
pokemon_types <- pokemon |> filter(!is.na(type_2))

type_pair_counts <- pokemon_types |>
  count(type_1, type_2, name = "n") |>
  group_by(type_1) |>
  mutate(row_total = sum(n), prop = n / row_total) |>
  arrange(type_1, desc(prop)) |>
  ungroup()
save_tbl(type_pair_counts, "tables/type1_type2_joint_counts_proportions.csv")
```

type_1	type_2	n	row_total	prop
bug	flying	14	61	0.230
bug	poison	12	61	0.197
bug	steel	7	61	0.115
bug	grass	6	61	0.098
bug	electric	5	61	0.082
bug	fighting	4	61	0.066
bug	rock	3	61	0.049
bug	water	3	61	0.049
bug	fairy	2	61	0.033
bug	fire	2	61	0.033
bug	ground	2	61	0.033
bug	ghost	1	61	0.016
dark	flying	5	25	0.200
dark	dragon	4	25	0.160
dark	fire	3	25	0.120
dark	normal	3	25	0.120
dark	fighting	2	25	0.080

type_1	type_2	n	row_total	prop
dark	ghost	2	25	0.080
dark	ice	2	25	0.080
dark	psychic	2	25	0.080
dark	steel	2	25	0.080
dragon	ground	8	27	0.296
dragon	flying	6	27	0.222
dragon	psychic	4	27	0.148
dragon	fighting	3	27	0.111
dragon	ice	3	27	0.111
dragon	electric	1	27	0.037
dragon	fairy	1	27	0.037
dragon	fire	1	27	0.037
electric	flying	6	21	0.286
electric	steel	4	21	0.190
electric	fairy	2	21	0.095
electric	normal	2	21	0.095
electric	dragon	1	21	0.048
electric	fire	1	21	0.048
electric	ghost	1	21	0.048
electric	grass	1	21	0.048
electric	ice	1	21	0.048
electric	psychic	1	21	0.048
electric	water	1	21	0.048
fairy	flying	2	2	1.000
fighting	psychic	3	9	0.333
fighting	steel	2	9	0.222
fighting	dark	1	9	0.111
fighting	flying	1	9	0.111
fighting	ghost	1	9	0.111
fighting	ice	1	9	0.111
fire	fighting	7	28	0.250
fire	flying	7	28	0.250
fire	ground	3	28	0.107
fire	dragon	2	28	0.071
fire	normal	2	28	0.071
fire	psychic	2	28	0.071
fire	dark	1	28	0.036
fire	ghost	1	28	0.036
fire	rock	1	28	0.036
fire	steel	1	28	0.036
fire	water	1	28	0.036
flying	dragon	2	2	1.000
ghost	grass	11	30	0.367
ghost	fairy	4	30	0.133
ghost	poison	4	30	0.133
ghost	fire	3	30	0.100
ghost	flying	3	30	0.100
ghost	dragon	2	30	0.067
ghost	ground	2	30	0.067
ghost	dark	1	30	0.033
grass	poison	15	45	0.333
grass	flying	7	45	0.156

type_1	type_2	n	row_total	prop
grass	fairy	5	45	0.111
grass	dark	3	45	0.067
grass	fighting	3	45	0.067
grass	ice	3	45	0.067
grass	steel	3	45	0.067
grass	dragon	2	45	0.044
grass	psychic	2	45	0.044
grass	ghost	1	45	0.022
grass	ground	1	45	0.022
ground	flying	4	21	0.190
ground	dark	3	21	0.143
ground	rock	3	21	0.143
ground	steel	3	21	0.143
ground	dragon	2	21	0.095
ground	ghost	2	21	0.095
ground	psychic	2	21	0.095
ground	electric	1	21	0.048
ground	fire	1	21	0.048
ice	ground	3	14	0.214
ice	water	3	14	0.214
ice	flying	2	14	0.143
ice	psychic	2	14	0.143
ice	steel	2	14	0.143
ice	fairy	1	14	0.071
ice	ghost	1	14	0.071
normal	flying	27	44	0.614
normal	fairy	5	44	0.114
normal	fighting	4	44	0.091
normal	psychic	3	44	0.068
normal	grass	2	44	0.045
normal	dragon	1	44	0.023
normal	ground	1	44	0.023
normal	water	1	44	0.023
poison	dark	5	20	0.250
poison	fire	3	20	0.150
poison	flying	3	20	0.150
poison	water	3	20	0.150
poison	fighting	2	20	0.100
poison	ground	2	20	0.100
poison	bug	1	20	0.050
poison	dragon	1	20	0.050
psychic	fairy	7	23	0.304
psychic	flying	7	23	0.304
psychic	fighting	3	23	0.130
psychic	ghost	2	23	0.087
psychic	dark	1	23	0.043
psychic	fire	1	23	0.043
psychic	grass	1	23	0.043
psychic	steel	1	23	0.043
rock	flying	18	53	0.340
rock	ground	6	53	0.113
rock	water	6	53	0.113

type_1	type_2	n	row_total	prop
rock	electric	3	53	0.057
rock	fairy	3	53	0.057
rock	steel	3	53	0.057
rock	bug	2	53	0.038
rock	dark	2	53	0.038
rock	dragon	2	53	0.038
rock	grass	2	53	0.038
rock	ice	2	53	0.038
rock	psychic	2	53	0.038
rock	fighting	1	53	0.019
rock	poison	1	53	0.019
steel	psychic	7	25	0.280
steel	fairy	5	25	0.200
steel	ghost	4	25	0.160
steel	rock	3	25	0.120
steel	flying	2	25	0.080
steel	ground	2	25	0.080
steel	dragon	1	25	0.040
steel	fighting	1	25	0.040
water	ground	10	60	0.167
water	dark	8	60	0.133
water	flying	7	60	0.117
water	psychic	6	60	0.100
water	fairy	4	60	0.067
water	rock	4	60	0.067
water	fighting	3	60	0.050
water	grass	3	60	0.050
water	ice	3	60	0.050
water	poison	3	60	0.050
water	bug	2	60	0.033
water	dragon	2	60	0.033
water	electric	2	60	0.033
water	ghost	2	60	0.033
water	steel	1	60	0.017

```
knitr::kable(head(type_pair_counts, 30), digits = 3,
  caption = "P(type_2 | type_1): top 30 rows.") # as in your Table 3
```

Table 6: P(type_2 | type_1): top 30 rows.

type_1	type_2	n	row_total	prop
bug	flying	14	61	0.230
bug	poison	12	61	0.197
bug	steel	7	61	0.115
bug	grass	6	61	0.098
bug	electric	5	61	0.082
bug	fighting	4	61	0.066
bug	rock	3	61	0.049
bug	water	3	61	0.049
bug	fairy	2	61	0.033

type_1	type_2	n	row_total	prop
bug	fire	2	61	0.033
bug	ground	2	61	0.033
bug	ghost	1	61	0.016
dark	flying	5	25	0.200
dark	dragon	4	25	0.160
dark	fire	3	25	0.120
dark	normal	3	25	0.120
dark	fighting	2	25	0.080
dark	ghost	2	25	0.080
dark	ice	2	25	0.080
dark	psychic	2	25	0.080
dark	steel	2	25	0.080
dragon	ground	8	27	0.296
dragon	flying	6	27	0.222
dragon	psychic	4	27	0.148
dragon	fighting	3	27	0.111
dragon	ice	3	27	0.111
dragon	electric	1	27	0.037
dragon	fairy	1	27	0.037
dragon	fire	1	27	0.037
electric	flying	6	21	0.286

```

type2_pref <- type_pair_counts |>
  group_by(type_1) |>
  slice_max(prop, n = 1, with_ties = TRUE) |>
  arrange(type_1, desc(prop)) |>
  ungroup()
save_tbl(type2_pref, "tables/type2_preference_by_type1.csv")

```

type_1	type_2	n	row_total	prop
bug	flying	14	61	0.230
dark	flying	5	25	0.200
dragon	ground	8	27	0.296
electric	flying	6	21	0.286
fairy	flying	2	2	1.000
fighting	psychic	3	9	0.333
fire	fighting	7	28	0.250
fire	flying	7	28	0.250
flying	dragon	2	2	1.000
ghost	grass	11	30	0.367
grass	poison	15	45	0.333
ground	flying	4	21	0.190
ice	ground	3	14	0.214
ice	water	3	14	0.214
normal	flying	27	44	0.614
poison	dark	5	20	0.250
psychic	fairy	7	23	0.304
psychic	flying	7	23	0.304
rock	flying	18	53	0.340
steel	psychic	7	25	0.280

type_1	type_2	n	row_total	prop
water	ground	10	60	0.167

```
knitr::kable(type2_pref, digits = 3,
  caption = "Most frequent secondary type(s) for each primary type.")
```

Table 8: Most frequent secondary type(s) for each primary type.

type_1	type_2	n	row_total	prop
bug	flying	14	61	0.230
dark	flying	5	25	0.200
dragon	ground	8	27	0.296
electric	flying	6	21	0.286
fairy	flying	2	2	1.000
fighting	psychic	3	9	0.333
fire	fighting	7	28	0.250
fire	flying	7	28	0.250
flying	dragon	2	2	1.000
ghost	grass	11	30	0.367
grass	poison	15	45	0.333
ground	flying	4	21	0.190
ice	ground	3	14	0.214
ice	water	3	14	0.214
normal	flying	27	44	0.614
poison	dark	5	20	0.250
psychic	fairy	7	23	0.304
psychic	flying	7	23	0.304
rock	flying	18	53	0.340
steel	psychic	7	25	0.280
water	ground	10	60	0.167

```
tab_t1_t2 <- table(type_1 = pokemon_types$type_1, type_2 = pokemon_types$type_2)
chi_res <- chisq.test(tab_t1_t2) # you printed the warning; keep behavior
```

```
## Warning in stats::chisq.test(x, y, ...): Chi-squared approximation may be
## incorrect
```

```
chi_tbl <- tibble(statistic = as.numeric(chi_res$statistic),
  df = as.integer(chi_res$parameter),
  p_value = as.numeric(chi_res$p.value))
save_tbl(chi_tbl, "tables/chisq_type1_type2_independence.csv")
```

statistic	df	p_value
698.644	289	0

```
chi_res # prints Pearson X^2 (X^2=698.64, df=289, p<2.2e-16). :
```



```
##
## Pearson's Chi-squared test
##
## data:  tab_t1_t2
## X-squared = 698.64, df = 289, p-value < 2.2e-16

# By generation_id: distribution tables with boxplots
pokemon <- pokemon |>
  mutate(generation_id_f = forcats::fct_na_value_to_level(as.factor(generation_id), level = "Unknown"))
stat_vars <- c("hp", "attack", "defense", "special_attack", "special_defense", "speed")

gen_stats_all <- pokemon |>
  group_by(generation_id_f) |>
  summarise(
    n = n(),
    across(all_of(stat_vars),
      list(mean = ~mean(.x, na.rm = TRUE),
           sd = ~sd(.x, na.rm = TRUE)), .names = "{.col}_{.fn}")
  ) |>
  arrange(generation_id_f)
save_tbl(gen_stats_all, "tables/gen_stats_by_generation_all_stats.csv", digits = 1)
```

generation_id	hp	hp_mean	hp_sd	attack	attack_mean	attack_sd	defense	defense_mean	defense_sd	special_attack	special_attack_mean	special_attack_sd	special_defense	special_defense_mean	special_defense_sd	speed	speed_mean	speed_sd
1	151	64.2	28.6	72.9	26.8	68.2	26.9	67.1	28.5	66.1	24.2	69.1	27.0					
2	100	71.0	31.2	68.3	28.4	69.7	35.2	64.5	25.6	72.3	31.5	61.4	27.2					
3	135	65.7	25.2	73.1	30.4	69.0	31.1	67.9	28.3	66.5	28.5	61.6	26.9					
4	107	73.1	24.7	80.2	30.9	75.2	30.7	73.3	31.2	74.5	27.8	69.5	27.6					
5	156	70.3	21.6	81.0	29.4	71.2	23.0	69.2	29.8	67.3	21.9	66.6	28.2					
6	72	68.9	21.7	72.5	25.6	75.2	31.7	72.5	28.0	74.7	30.8	65.7	25.9					
7	81	70.7	28.2	83.2	32.6	77.0	29.9	73.5	33.4	74.8	29.3	64.5	29.1					
Unknown	147	70.1	25.1	98.9	37.5	87.7	34.4	92.0	42.0	84.6	26.2	87.3	31.4					

```
knitr::kable(gen_stats_all, digits = 1,
  caption = "Means and SDs of stats by generation (including Unknown).") # your Table 6
```

Table 11: Means and SDs of stats by generation (including Unknown).

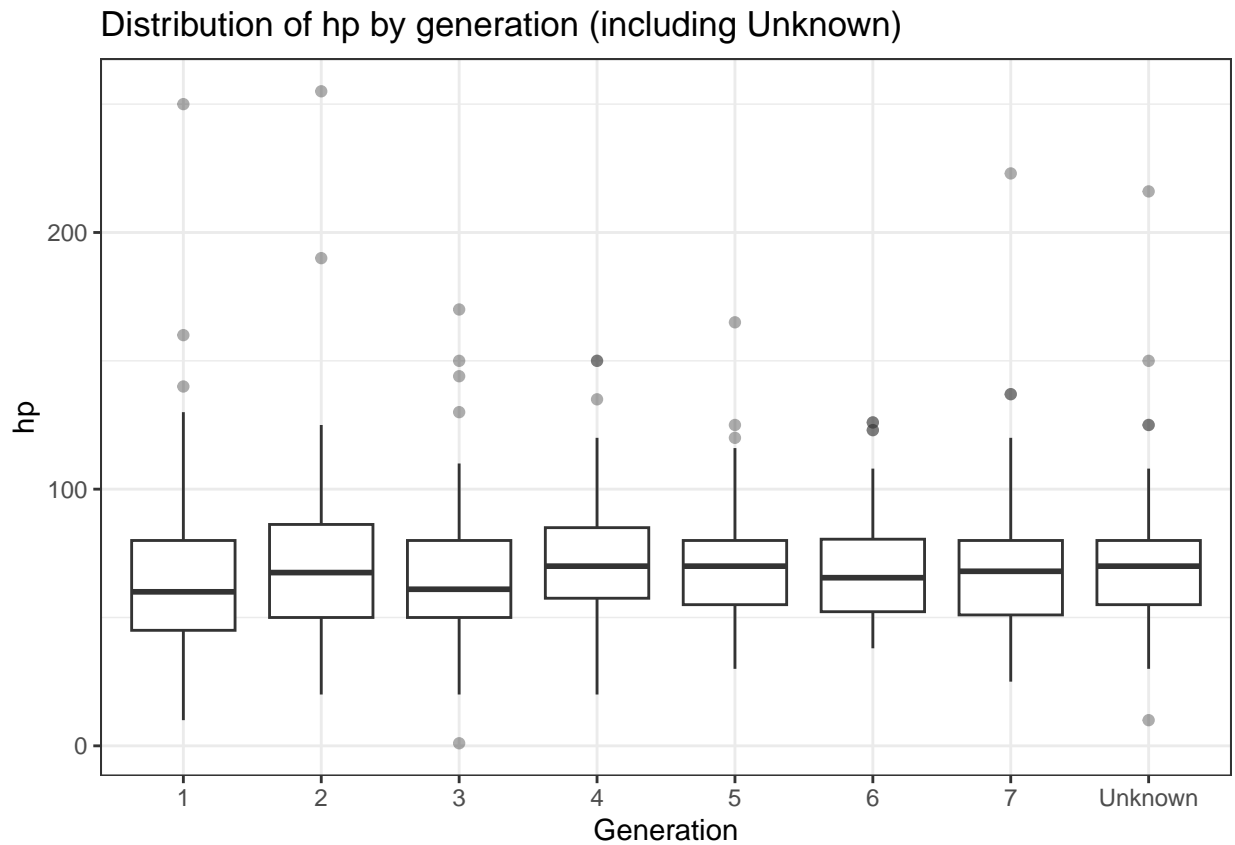
generation_id	hp	hp_mean	hp_sd	attack	attack_mean	attack_sd	defense	defense_mean	defense_sd	special_attack	special_attack_mean	special_attack_sd	special_defense	special_defense_mean	special_defense_sd	speed	speed_mean	speed_sd
1	151	64.2	28.6	72.9	26.8	68.2	26.9	67.1	28.5	66.1	24.2	69.1	27.0					
2	100	71.0	31.2	68.3	28.4	69.7	35.2	64.5	25.6	72.3	31.5	61.4	27.2					
3	135	65.7	25.2	73.1	30.4	69.0	31.1	67.9	28.3	66.5	28.5	61.6	26.9					
4	107	73.1	24.7	80.2	30.9	75.2	30.7	73.3	31.2	74.5	27.8	69.5	27.6					
5	156	70.3	21.6	81.0	29.4	71.2	23.0	69.2	29.8	67.3	21.9	66.6	28.2					
6	72	68.9	21.7	72.5	25.6	75.2	31.7	72.5	28.0	74.7	30.8	65.7	25.9					
7	81	70.7	28.2	83.2	32.6	77.0	29.9	73.5	33.4	74.8	29.3	64.5	29.1					
Unknown	147	70.1	25.1	98.9	37.5	87.7	34.4	92.0	42.0	84.6	26.2	87.3	31.4					

```

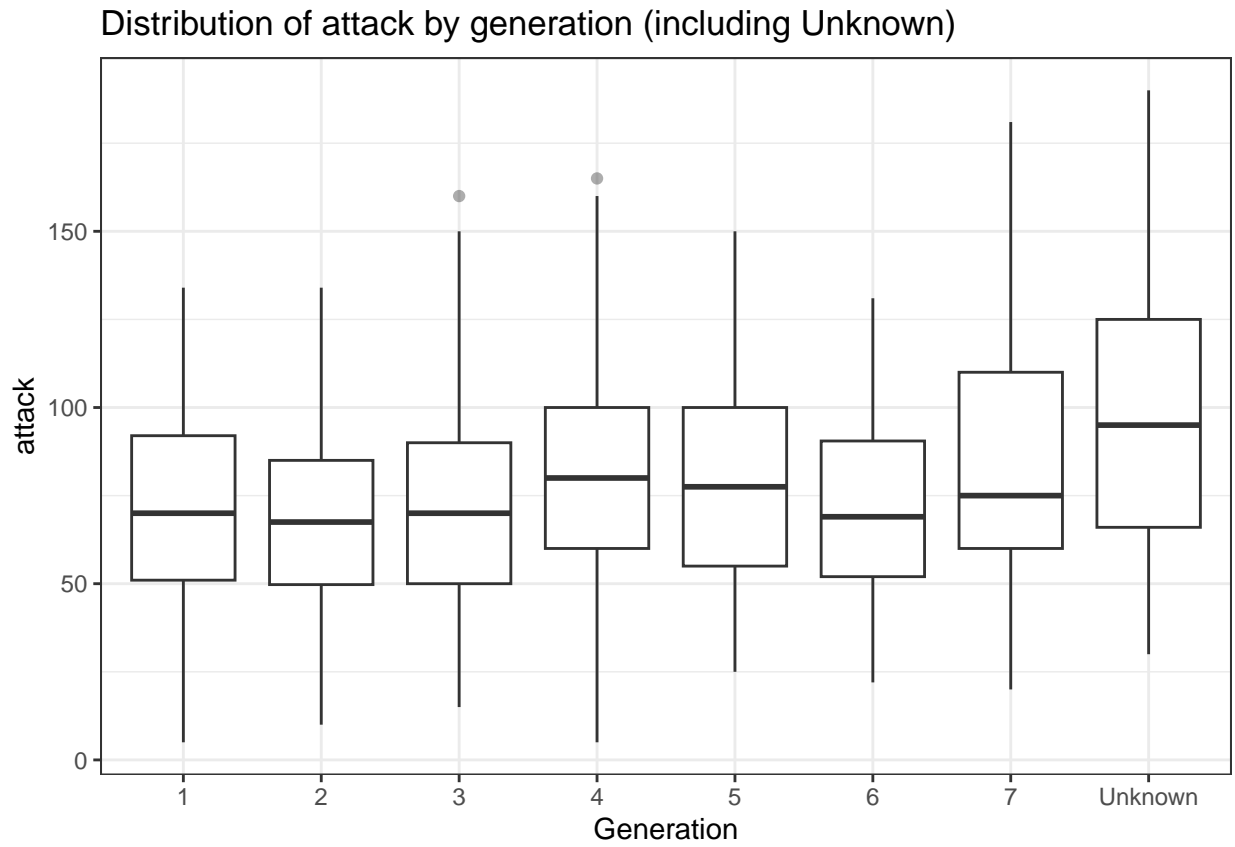
for (s in stat_vars) {
  p <- ggplot(pokemon, aes(x = generation_id_f, y = .data[[s]])) +
    geom_boxplot(outlier.alpha = 0.4) +
    labs(title = paste("Distribution of", s, "by generation (including Unknown)",
      x = "Generation", y = s) + theme_bw()
  save_plot(p, paste0("figures/box_", s, "_by_generation.png")) ; print(p)
}

```

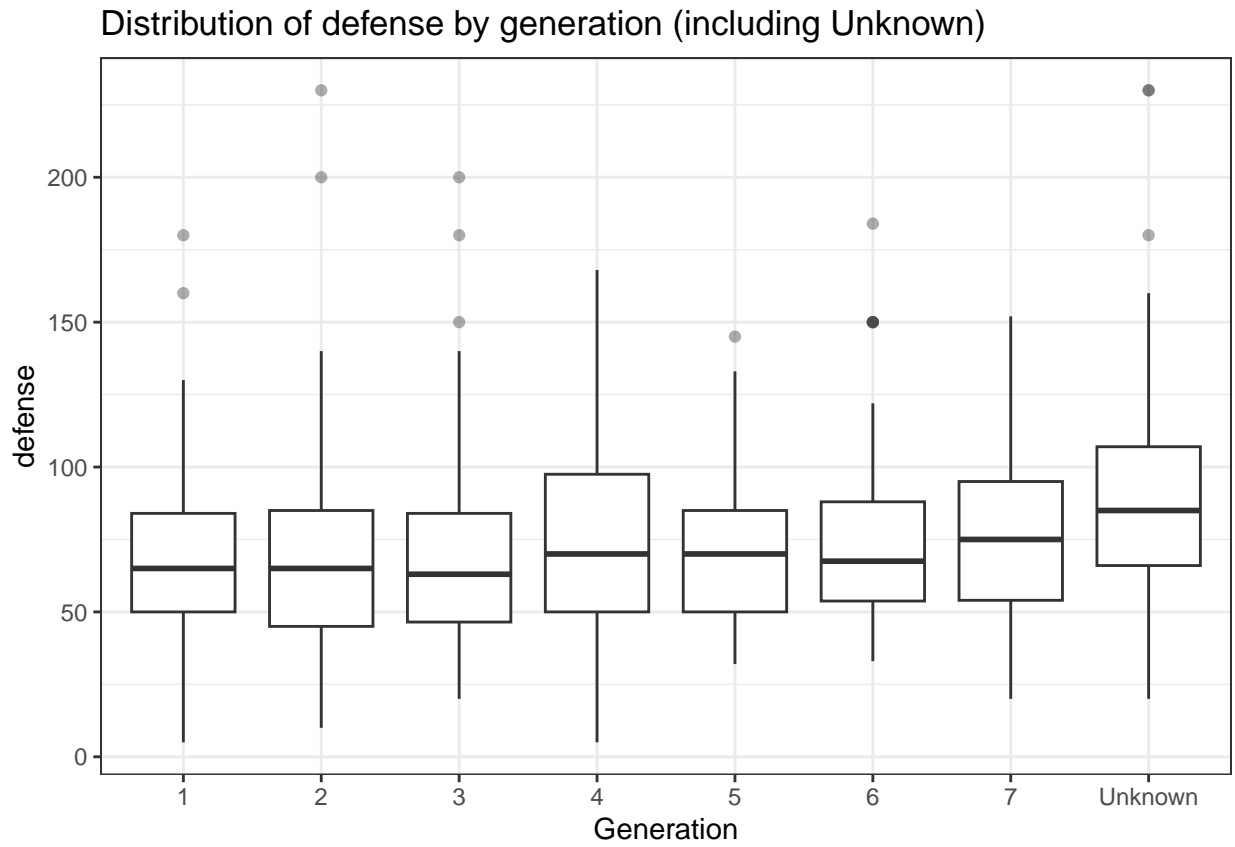
Saved: figures/box_hp_by_generation.png



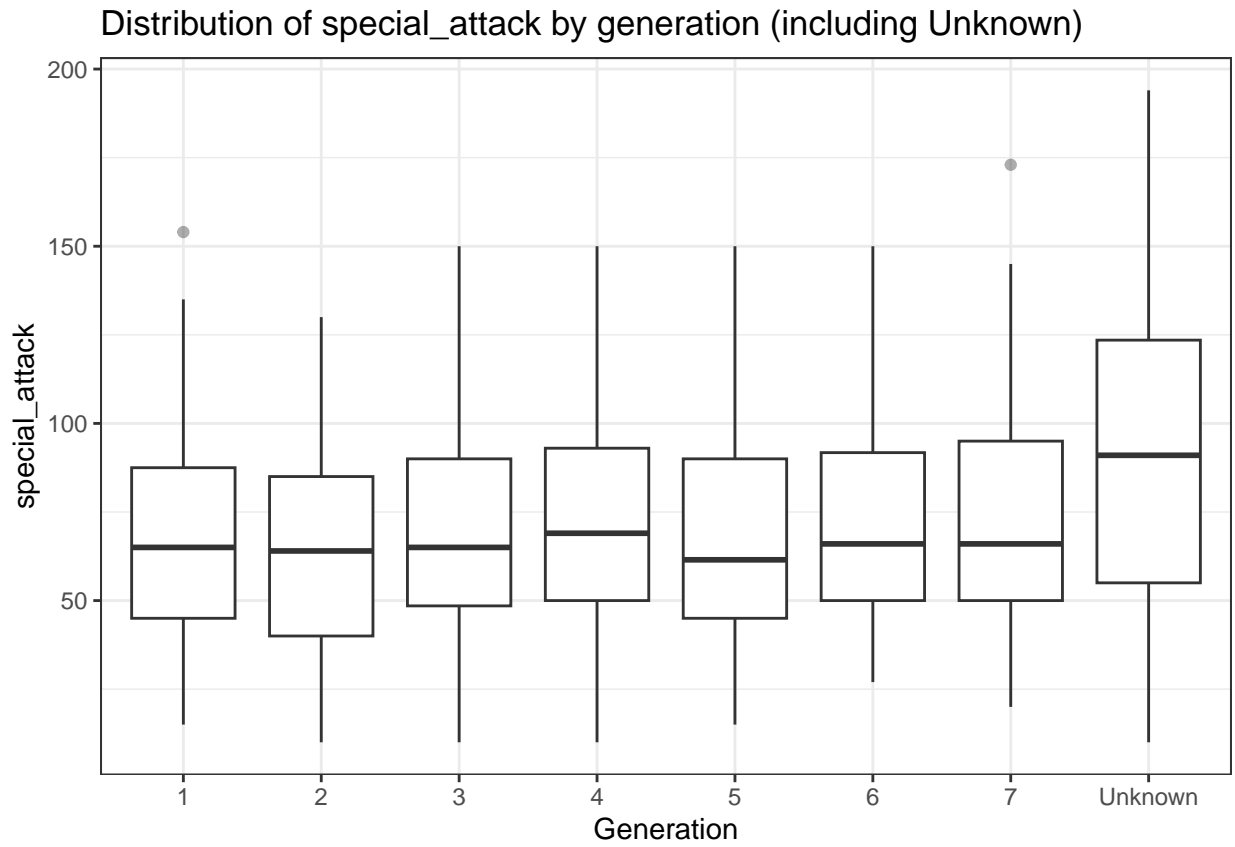
Saved: figures/box_attack_by_generation.png



Saved: figures/box_defense_by_generation.png

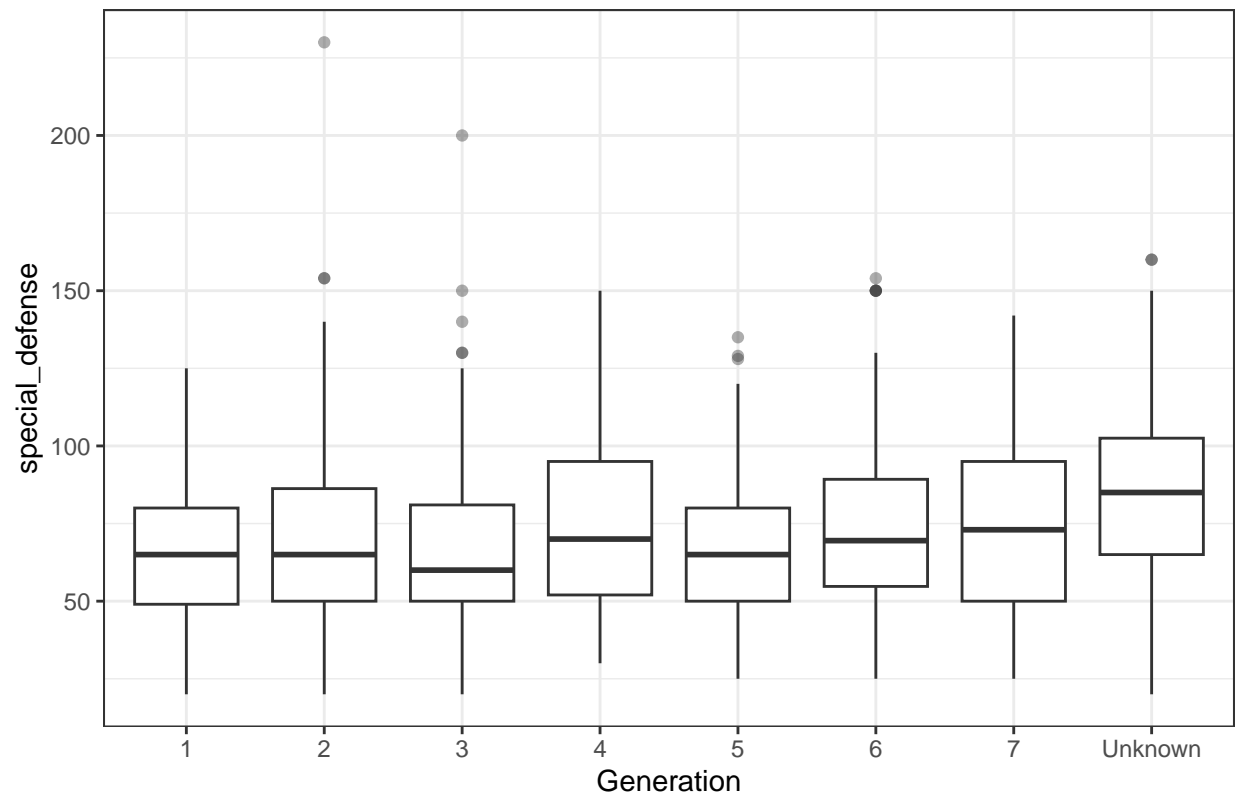


Saved: figures/box_special_attack_by_generation.png



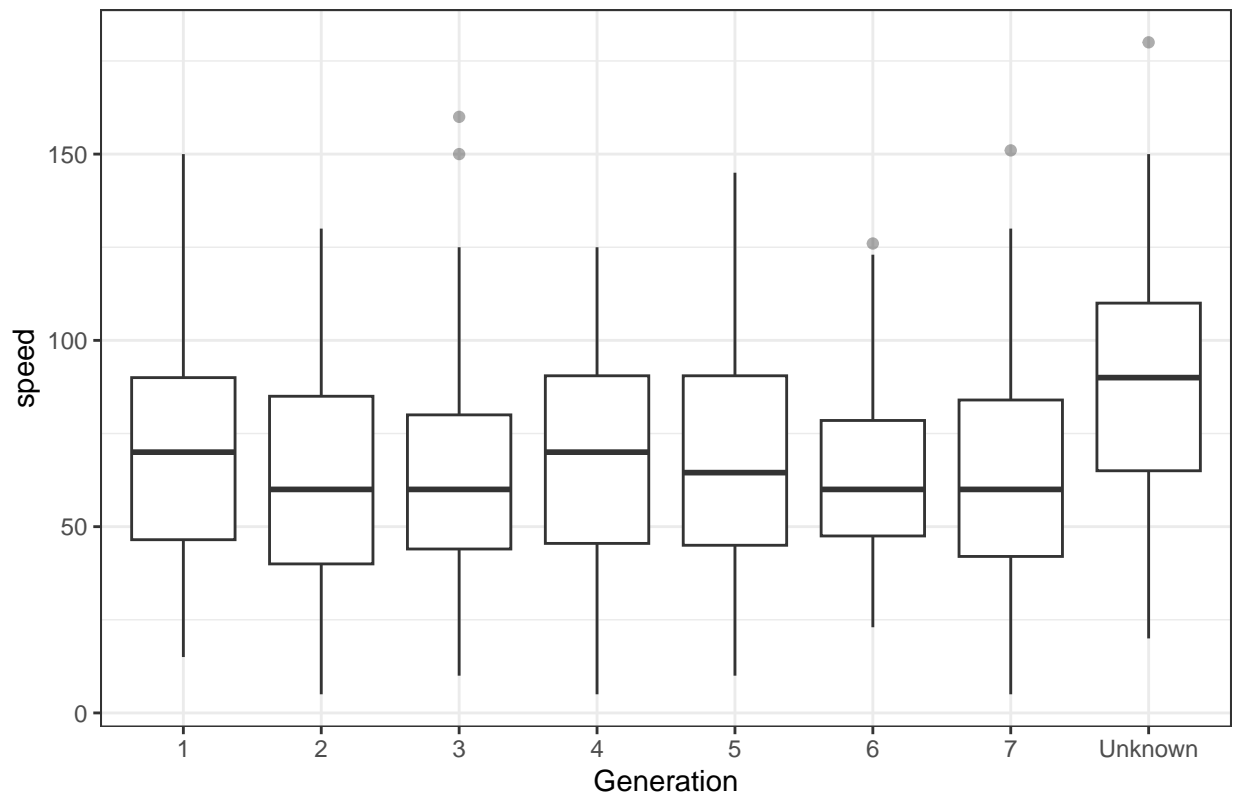
Saved: figures/box_special_defense_by_generation.png

Distribution of special_defense by generation (including Unknown)



Saved: figures/box_speed_by_generation.png

Distribution of speed by generation (including Unknown)



```
# By type1: distribution tables with boxplots
pokemon <- pokemon |>
  mutate(type_1_f = forcats::fct_infreq(as.factor(type_1)))

type1_stats <- pokemon |>
  group_by(type_1_f) |>
  summarise(
    n = n(),
    across(all_of(stat_vars),
      list(mean = ~mean(.x, na.rm = TRUE),
           sd   = ~sd(.x,   na.rm = TRUE)), .names = "{.col}_{.fn}")
  ) |>
  arrange(desc(n))
save_tbl(type1_stats, "tables/type1_stats_all_stats.csv", digits = 1)
```

type_1_f	hp	mp	hp_sd	attack	attack_sd	defense	defense_sd	special_attack	special_attack_sd	special_defense	special_defense_sd	speed	speed_sd
water	126	71.1	26.4	74.7	29.0	73.5	28.1	75.6	30.3	72.4	29.6	65.7	24.8
normal	111	77.3	34.8	75.9	30.0	60.5	23.5	57.5	25.1	64.4	25.2	70.3	27.7
grass	84	66.7	18.9	75.2	29.3	71.7	25.0	76.2	27.0	70.7	22.1	60.3	27.9
bug	79	57.7	17.3	72.5	37.2	72.1	34.1	57.7	31.0	64.4	31.0	63.3	34.0
rock	65	65.3	19.4	90.0	31.9	94.3	34.4	66.2	28.0	75.3	30.2	64.3	33.2
psychic	64	72.6	29.8	72.5	42.0	69.9	29.1	98.0	39.1	87.0	31.1	80.9	36.3
electric	61	55.8	18.3	68.1	22.4	61.2	23.7	83.0	32.6	69.1	21.4	87.2	24.0
fire	59	69.7	18.8	84.3	27.5	69.3	25.0	87.6	29.0	71.9	22.0	73.5	24.8

type	in_f	hp	mp	sp	attack	defense	defense	special	special	special	special	special	special	special
ghost	40	64.2	28.7	76.3	28.5	82.0	29.6	77.3	30.8	78.7	25.5	65.7	29.1	
dragon	39	84.5	32.0	108.7	32.4	89.3	25.0	93.4	39.8	88.3	28.4	82.9	22.8	
dark	37	69.7	33.2	84.7	26.3	67.6	24.5	71.4	32.8	67.8	24.1	76.6	26.8	
ground	36	71.6	27.5	96.2	32.2	82.6	33.5	55.3	26.8	62.9	20.7	64.6	27.9	
poison	35	67.2	19.6	73.5	19.8	69.3	24.4	62.5	21.7	65.9	23.6	64.2	25.7	
fighting	31	71.6	25.7	99.1	28.0	67.2	18.2	53.8	27.3	64.9	21.9	67.6	26.9	
steel	30	67.3	16.6	93.1	28.8	124.8	42.7	73.0	34.3	83.6	29.0	56.1	24.6	
ice	29	70.5	20.8	73.1	27.3	73.4	32.5	72.3	29.1	74.7	35.0	64.6	24.2	
fairy	19	72.9	22.9	61.2	28.1	67.1	18.7	81.2	28.9	88.3	30.2	53.6	26.6	
flying	4	70.8	20.7	78.8	37.5	66.2	21.4	94.2	34.8	72.5	22.2	102.5	32.1	

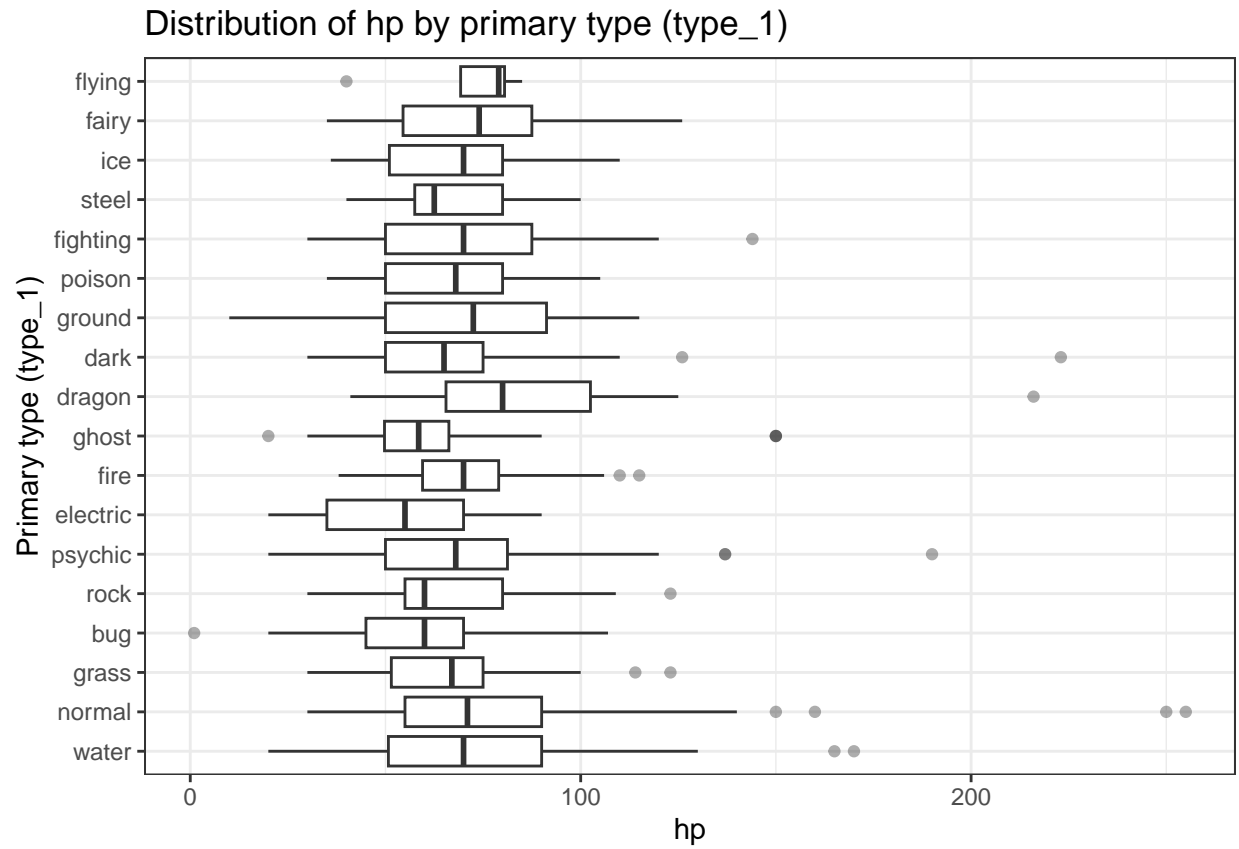
```
knitr::kable(type1_stats, digits = 1,
              caption = "Means and SDs of stats by primary type (type_1).") # your Table 7
```

Table 13: Means and SDs of stats by primary type (type_1).

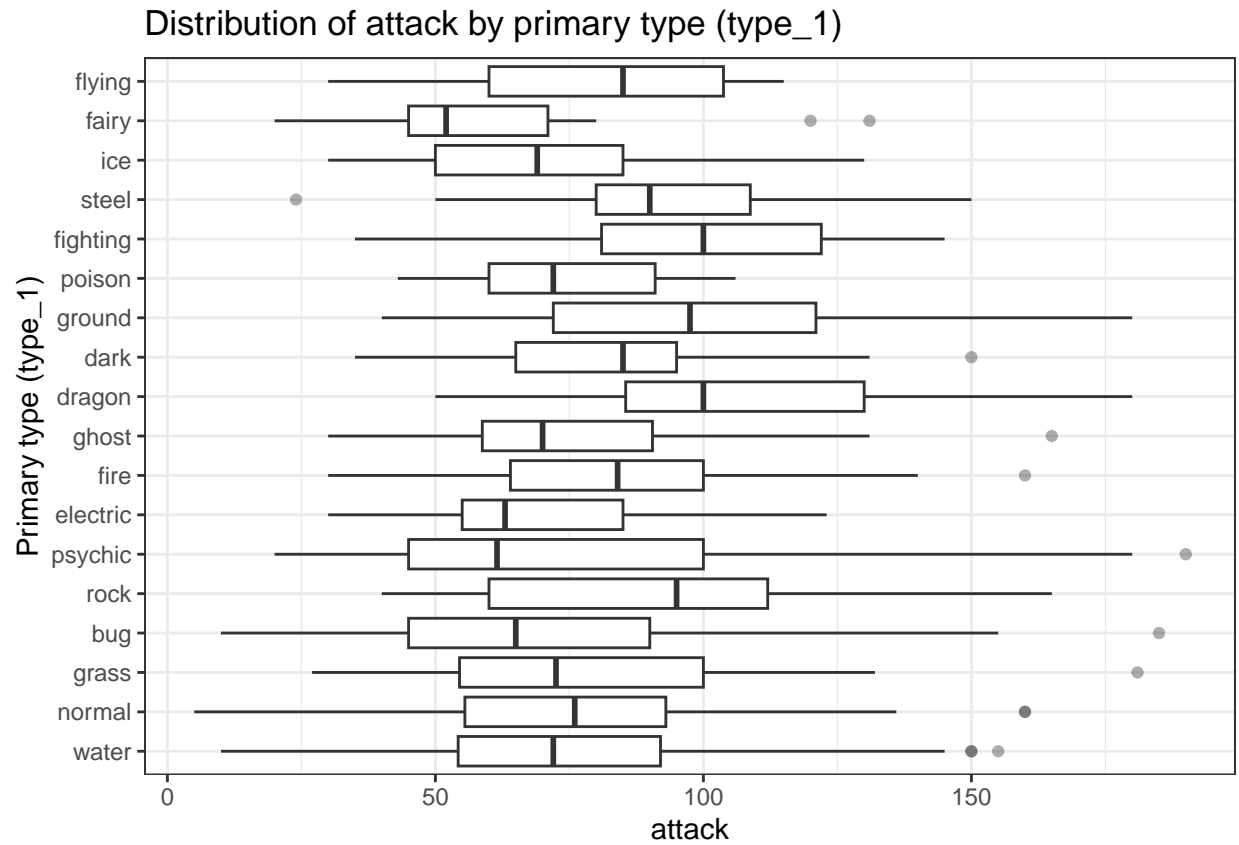
type	l	f	hp	mp	sp	attack	defense	special_attack	special_defense	speed	sd		
water	126	71.1	26.4	74.7	29.0	73.5	28.1	75.6	30.3	72.4	29.6	65.7	24.8
normal	11	77.3	34.8	75.9	30.0	60.5	23.5	57.5	25.1	64.4	25.2	70.3	27.7
grass	84	66.7	18.9	75.2	29.3	71.7	25.0	76.2	27.0	70.7	22.1	60.3	27.9
bug	79	57.7	17.3	72.5	37.2	72.1	34.1	57.7	31.0	64.4	31.0	63.3	34.0
rock	65	65.3	19.4	90.0	31.9	94.3	34.4	66.2	28.0	75.3	30.2	64.3	33.2
psychic	64	72.6	29.8	72.5	42.0	69.9	29.1	98.0	39.1	87.0	31.1	80.9	36.3
electric	61	55.8	18.3	68.1	22.4	61.2	23.7	83.0	32.6	69.1	21.4	87.2	24.0
fire	59	69.7	18.8	84.3	27.5	69.3	25.0	87.6	29.0	71.9	22.0	73.5	24.8
ghost	40	64.2	28.7	76.3	28.5	82.0	29.6	77.3	30.8	78.7	25.5	65.7	29.1
dragon	39	84.5	32.0	108.7	32.4	89.3	25.0	93.4	39.8	88.3	28.4	82.9	22.8
dark	37	69.7	33.2	84.7	26.3	67.6	24.5	71.4	32.8	67.8	24.1	76.6	26.8
ground	36	71.6	27.5	96.2	32.2	82.6	33.5	55.3	26.8	62.9	20.7	64.6	27.9
poison	35	67.2	19.6	73.5	19.8	69.3	24.4	62.5	21.7	65.9	23.6	64.2	25.7
fighting	31	71.6	25.7	99.1	28.0	67.2	18.2	53.8	27.3	64.9	21.9	67.6	26.9
steel	30	67.3	16.6	93.1	28.8	124.8	42.7	73.0	34.3	83.6	29.0	56.1	24.6
ice	29	70.5	20.8	73.1	27.3	73.4	32.5	72.3	29.1	74.7	35.0	64.6	24.2
fairy	19	72.9	22.9	61.2	28.1	67.1	18.7	81.2	28.9	88.3	30.2	53.6	26.6
flying	4	70.8	20.7	78.8	37.5	66.2	21.4	94.2	34.8	72.5	22.2	102.5	32.1

```
for (s in stat_vars) {
  p <- ggplot(pokemon, aes(x = type_1_f, y = .data[[s]])) +
    geom_boxplot(outlier.alpha = 0.4) +
    coord_flip() +
    labs(title = paste("Distribution of", s, "by primary type (type_1)",
                      x = "Primary type (type_1)", y = s) + theme_bw()
    save_plot(p, paste0("figures/box_", s, "_by_type1.png"), w = 7.5, h = 6) ; print(p)
}
```

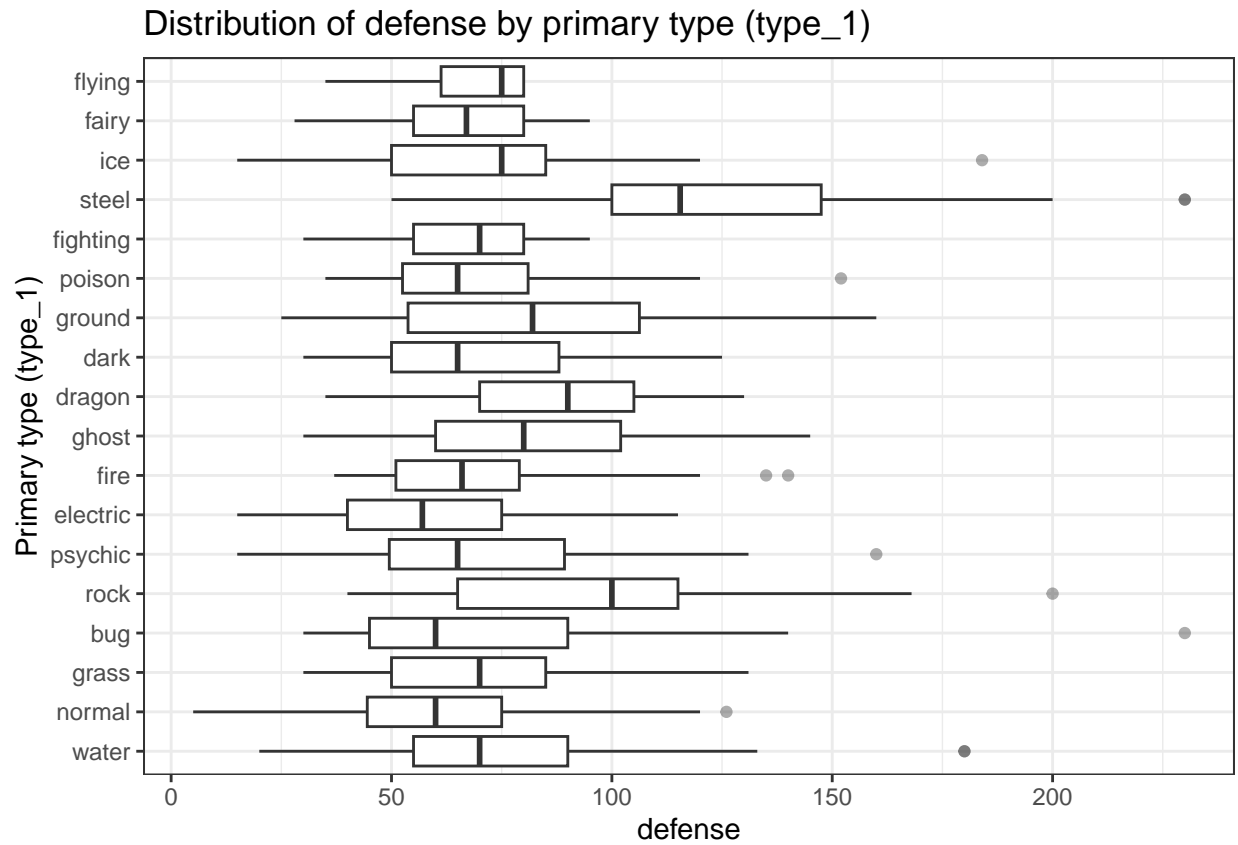
Saved: figures/box_hp_by_type1.png



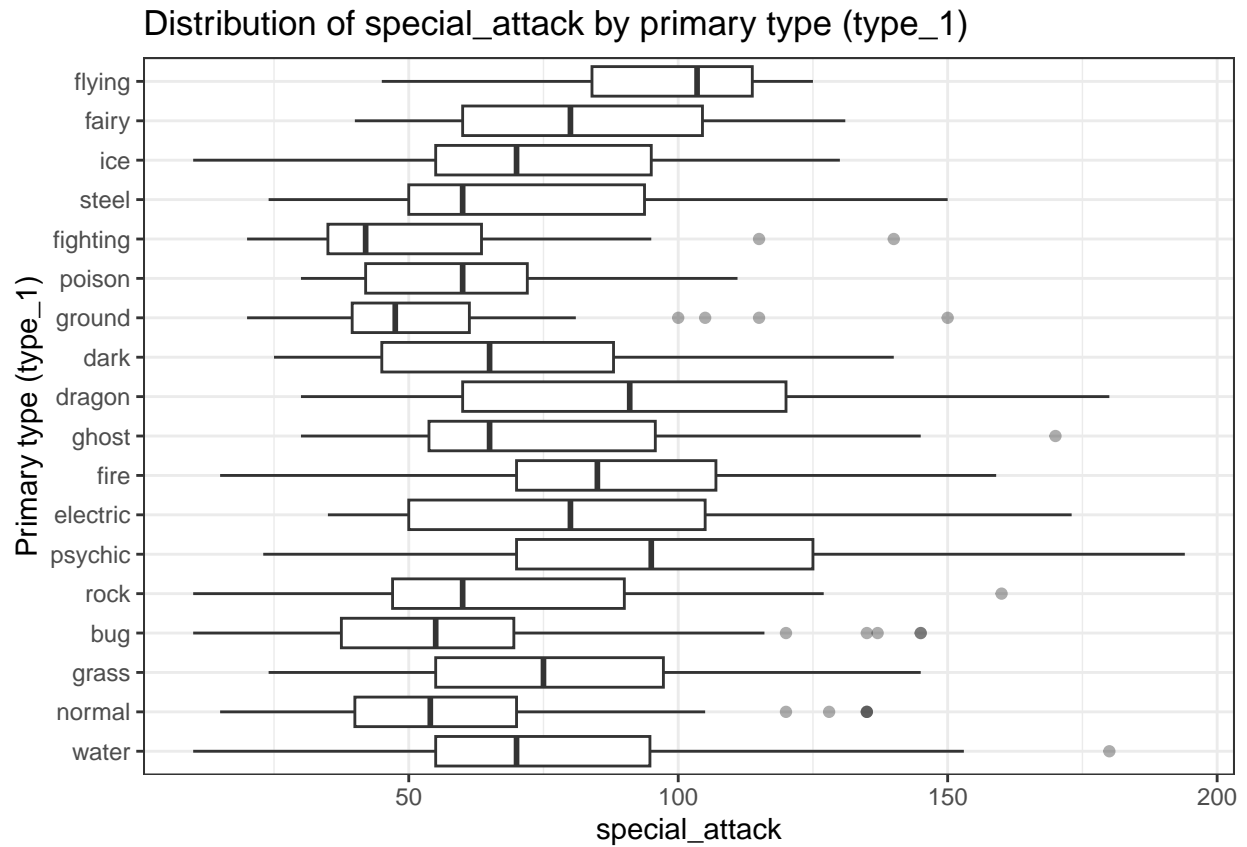
Saved: figures/box_attack_by_type1.png



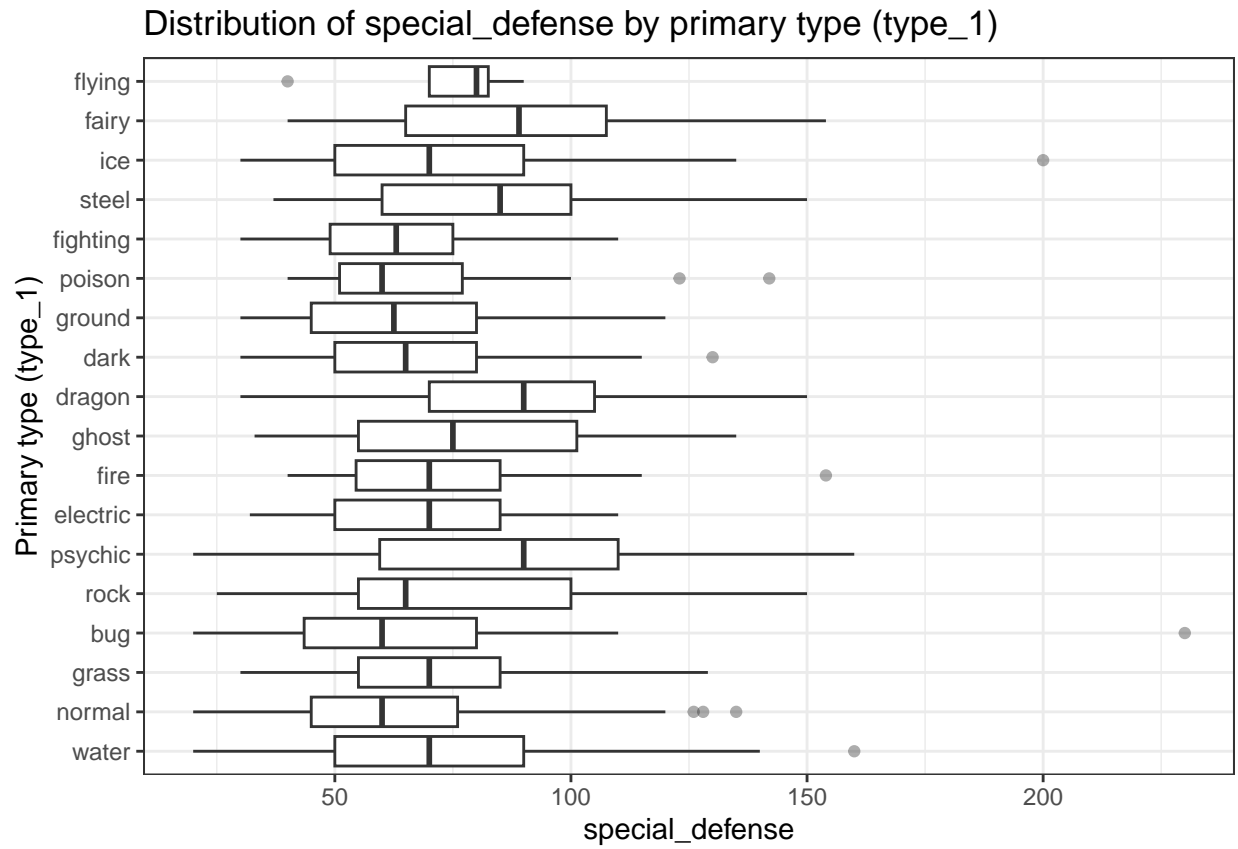
Saved: figures/box_defense_by_type1.png



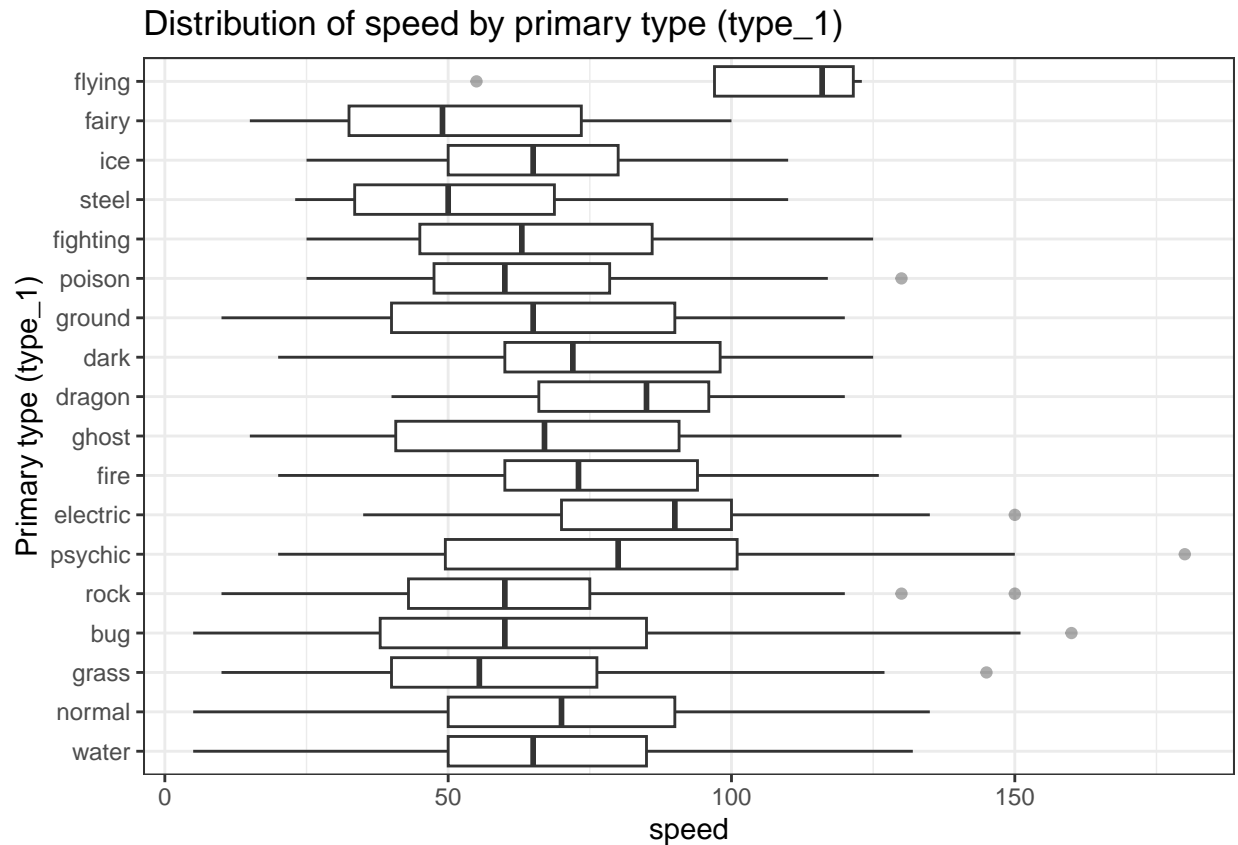
Saved: figures/box_special_attack_by_type1.png



Saved: figures/box_special_defense_by_type1.png



Saved: figures/box_speed_by_type1.png



```
# NA counts across ALL columns (1-row audit table)
na_counts <- pokemon |> summarise(across(everything(), ~sum(is.na(.))))
save_tbl(na_counts, "tables/na_counts_all_columns.csv") # see pp. 1-2. :contentReference[oa]
```

pokemon_id	height	weight	base_experience	type_1	type_2	attack	defense	special	special_defense	speed	egg_group	generation	is_legendary	is_banned	log_weight	log_height	log_experience	generation_id	type_id
0	0	0	0	0	0	439	0	0	0	0	0	0	711	147	0	0	0	0	0

```
# Type_1 counts
type1_counts <- pokemon |> count(type_1) |> arrange(desc(n))
save_tbl(type1_counts, "tables/type1_counts_kareena.csv")
```

type_1	n
water	126
normal	111
grass	84
bug	79
rock	65
psychic	64
electric	61
fire	59
ghost	40
dragon	39

type_1	n
dark	37
ground	36
poison	35
fighting	31
steel	30
ice	29
fairy	19
flying	4

```
knitr::kable(type1_counts, caption = "Counts of primary types (type_1).")
```

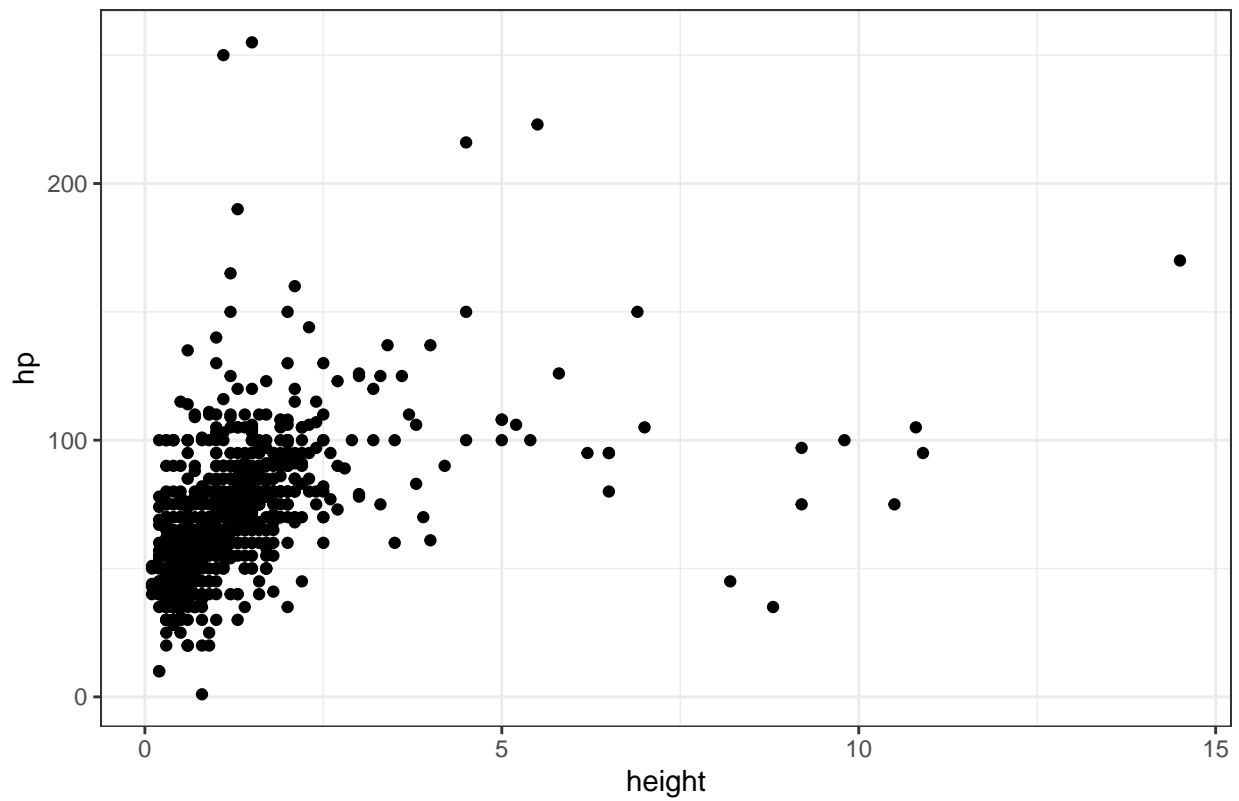
Table 16: Counts of primary types (type_1).

type_1	n
water	126
normal	111
grass	84
bug	79
rock	65
psychic	64
electric	61
fire	59
ghost	40
dragon	39
dark	37
ground	36
poison	35
fighting	31
steel	30
ice	29
fairy	19
flying	4

```
# Size → HP (raw vs log)
p_hp_h <- ggplot(pokemon, aes(height, hp)) + geom_point() + theme_bw() + labs(title="HP vs height")
p_hp_lh <- ggplot(pokemon, aes(log(height), hp)) + geom_point() + theme_bw() + labs(title="HP vs log(height)")
p_hp_w <- ggplot(pokemon, aes(weight, hp)) + geom_point() + theme_bw() + labs(title="HP vs weight")
p_hp_lw <- ggplot(pokemon, aes(log(weight), hp)) + geom_point() + theme_bw() + labs(title="HP vs log(weight)")
save_plot(p_hp_h, "figures/k_hp_vs_height.png")
```

```
## Saved: figures/k_hp_vs_height.png
```

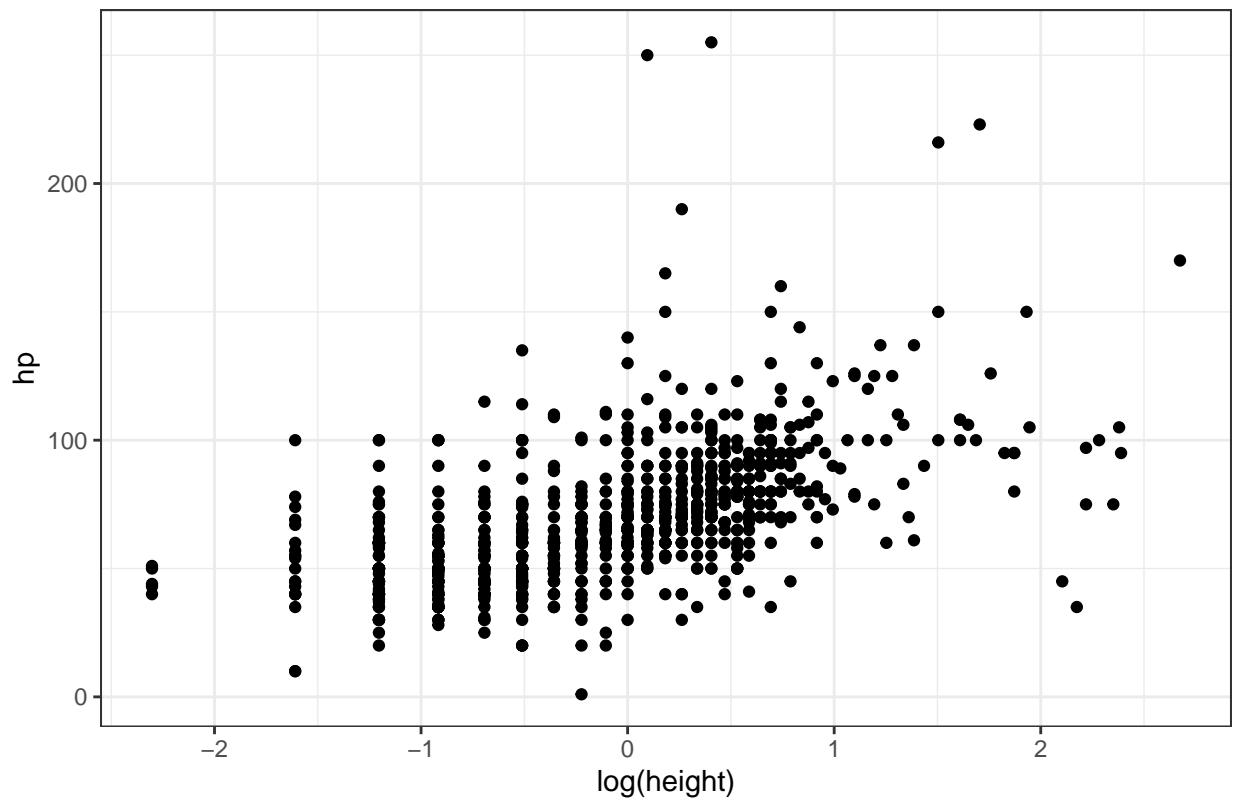
HP vs height



```
save_plot(p_hp_lh, "figures/k_hp_vs_logheight.png")
```

```
## Saved: figures/k_hp_vs_logheight.png
```

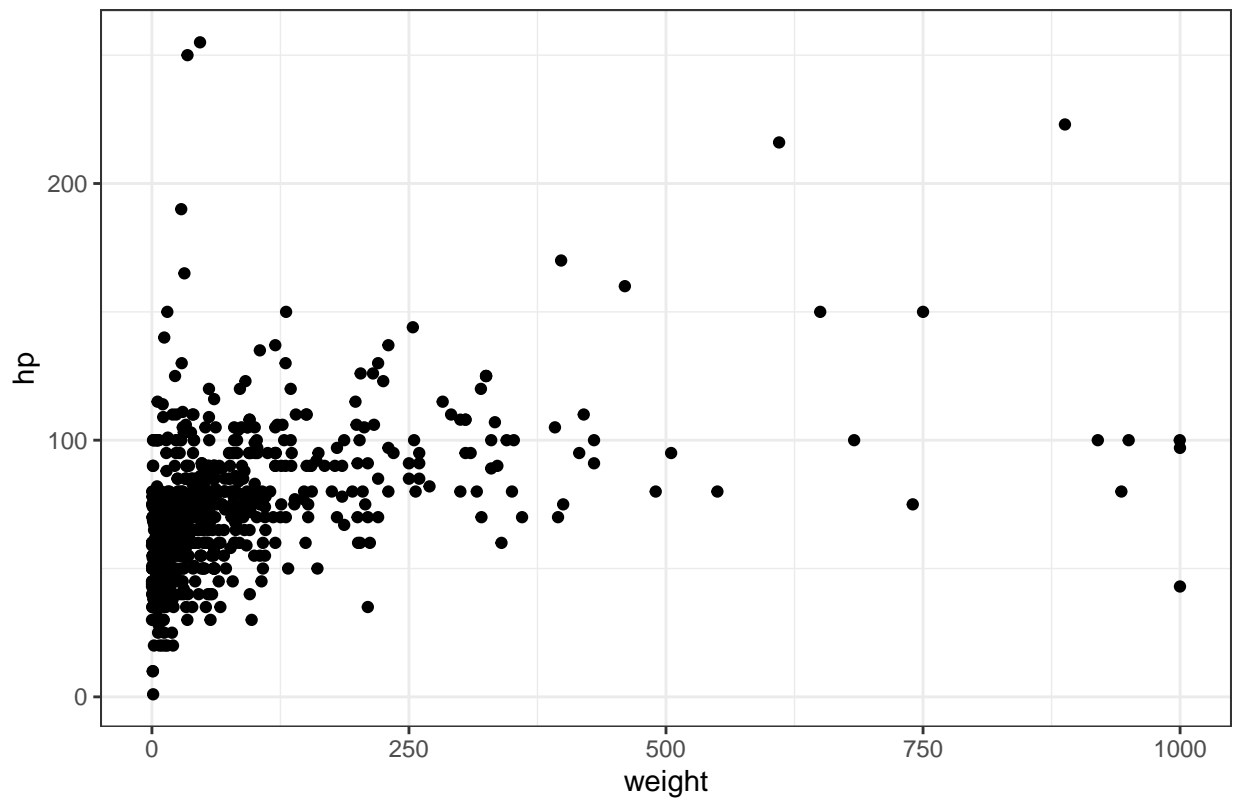

HP vs log(height)



```
save_plot(p_hp_w, "figures/k_hp_vs_weight.png")
```

```
## Saved: figures/k_hp_vs_weight.png
```

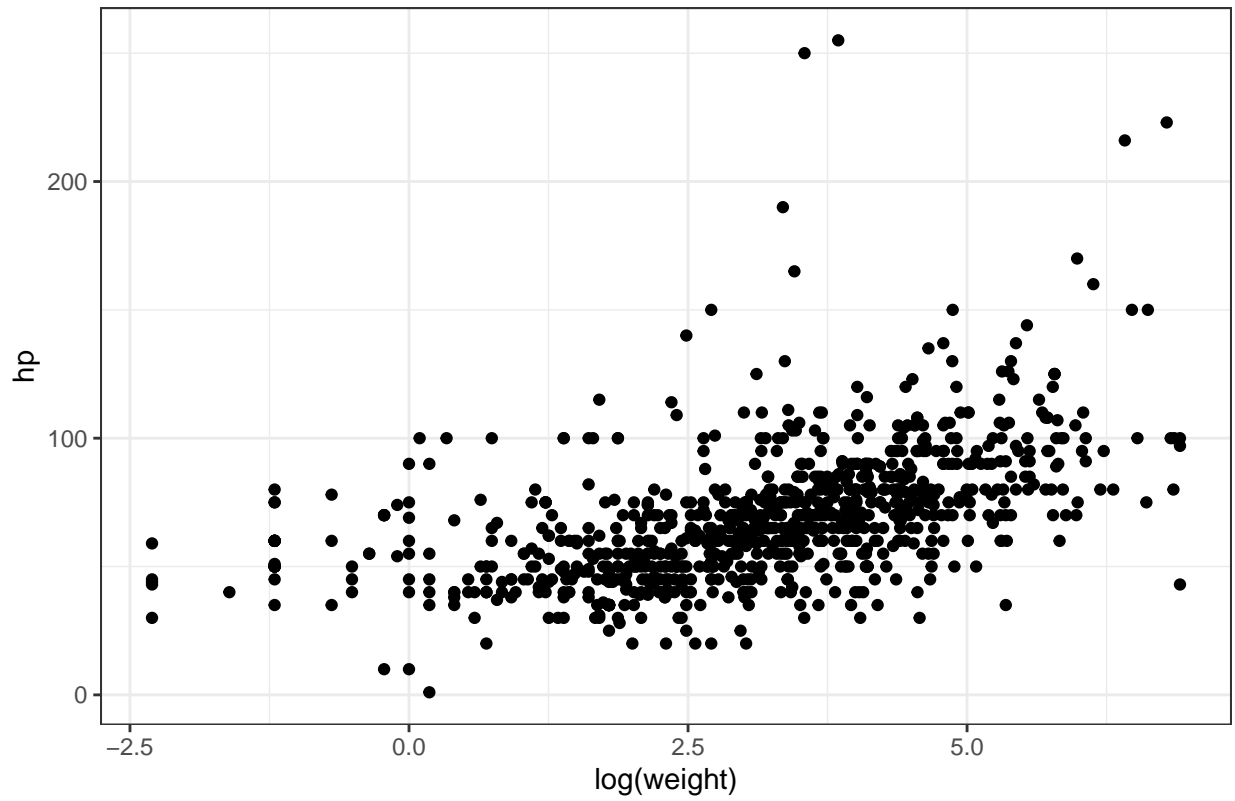
HP vs weight



```
save_plot(p_hp_lw, "figures/k_hp_vs_logweight.png")
```

```
## Saved: figures/k_hp_vs_logweight.png
```

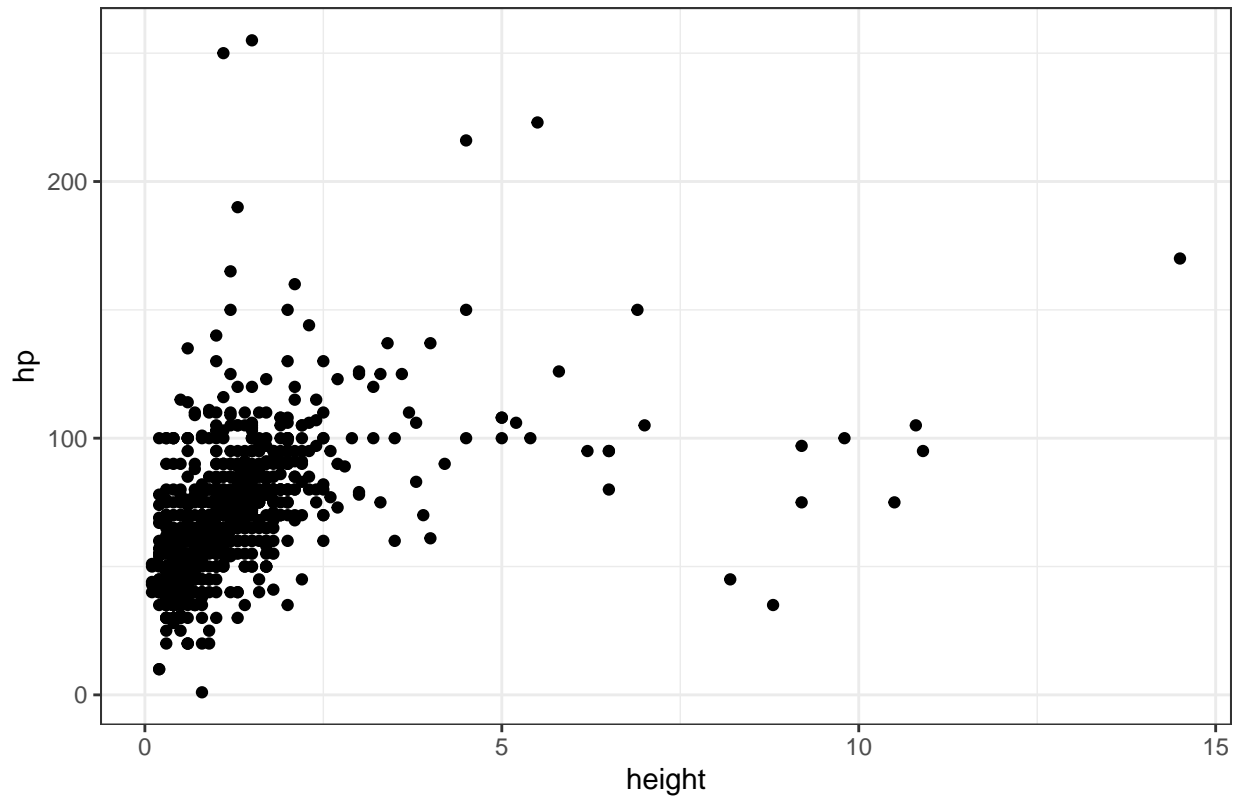
HP vs log(weight)



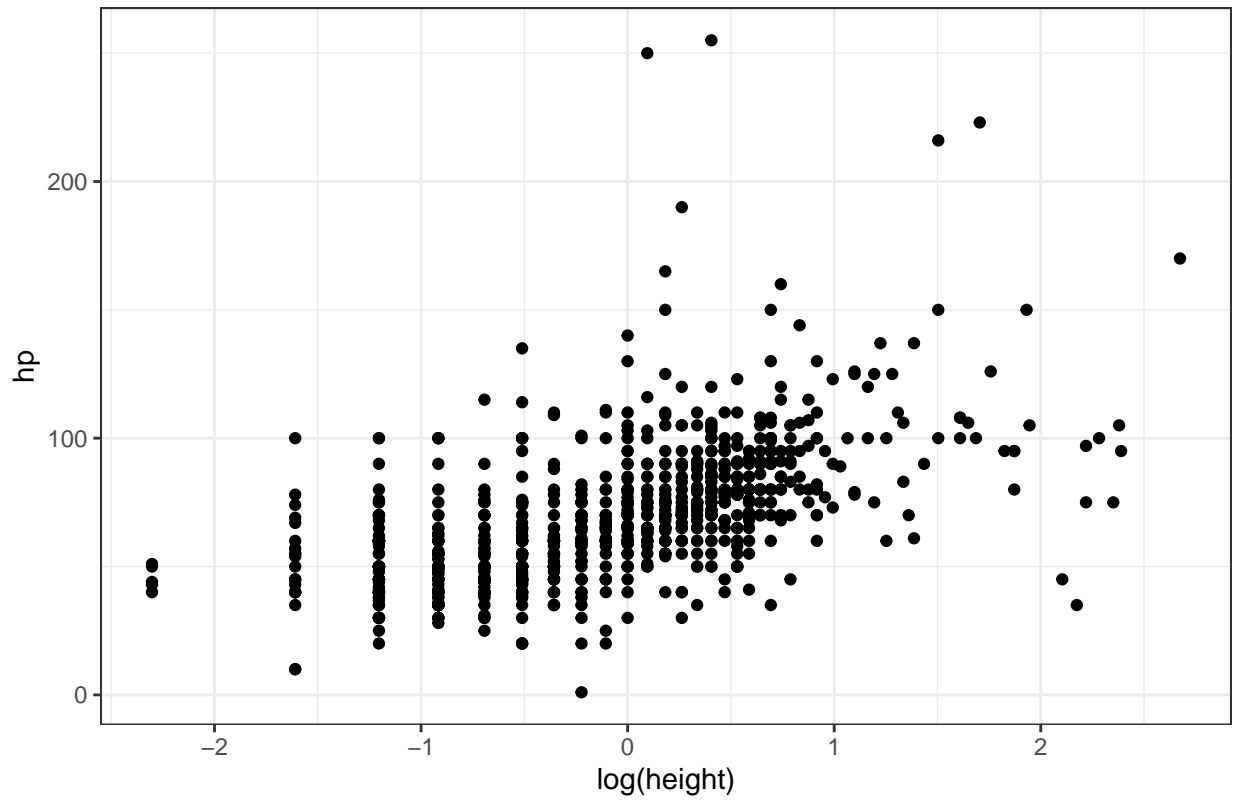
p_hp_h; p_hp_lh; p_hp_w; p_hp_lw

pp. 3-5. :contentReference[

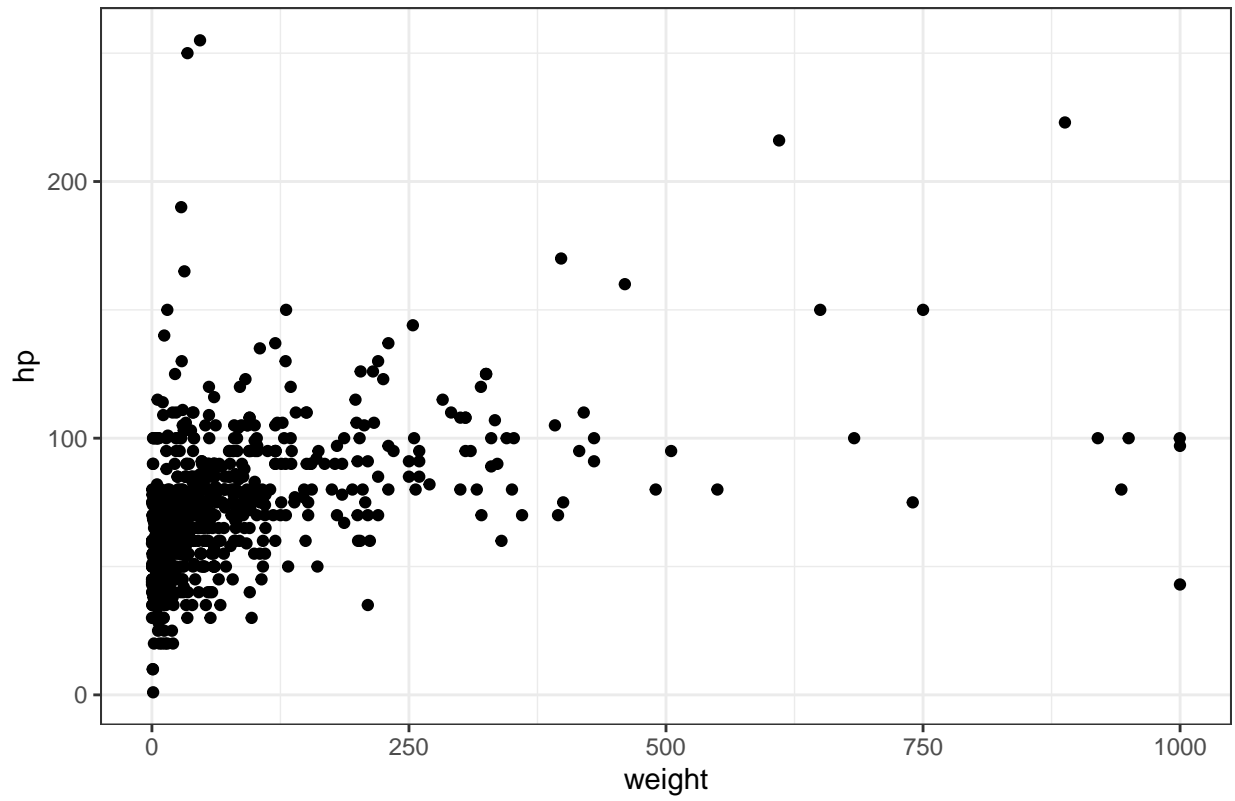
HP vs height



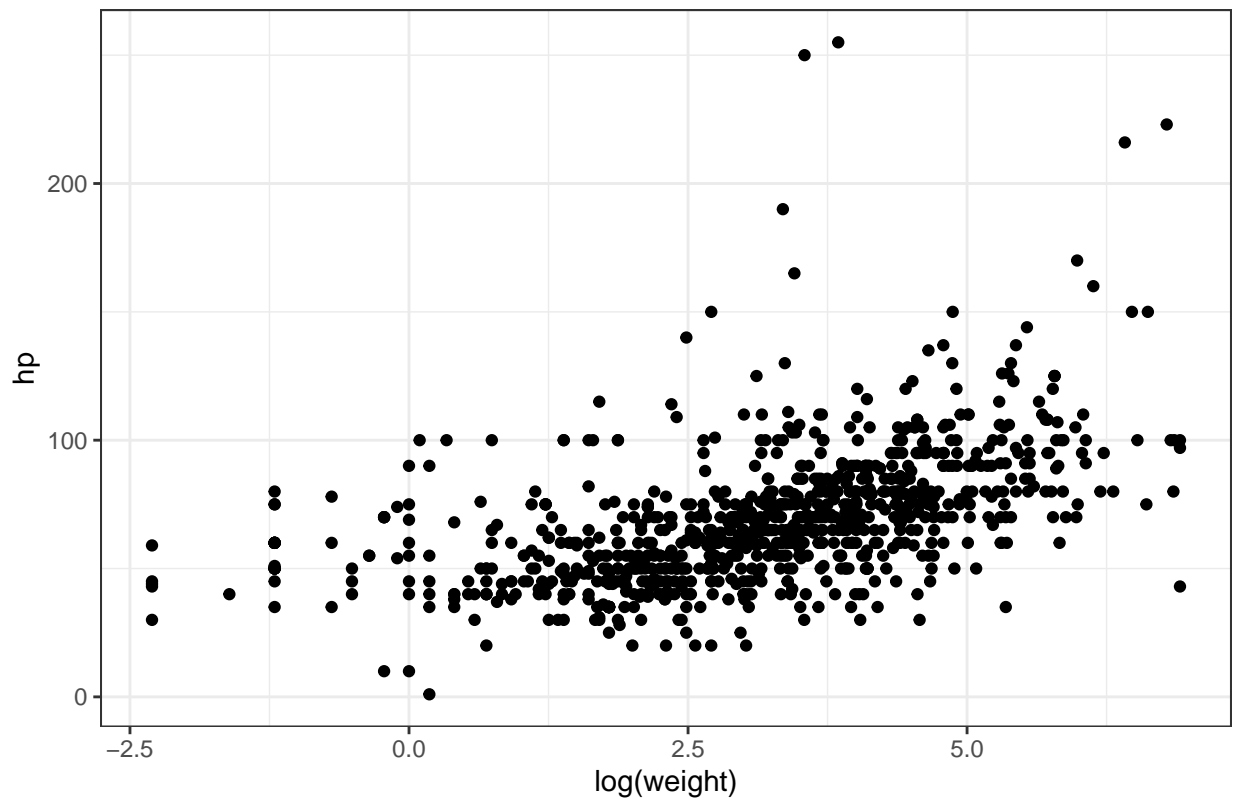
HP vs log(height)



HP vs weight



HP vs log(weight)



```
# Attack vs Special Attack; Defense vs Special Defense + identical-value counts
```

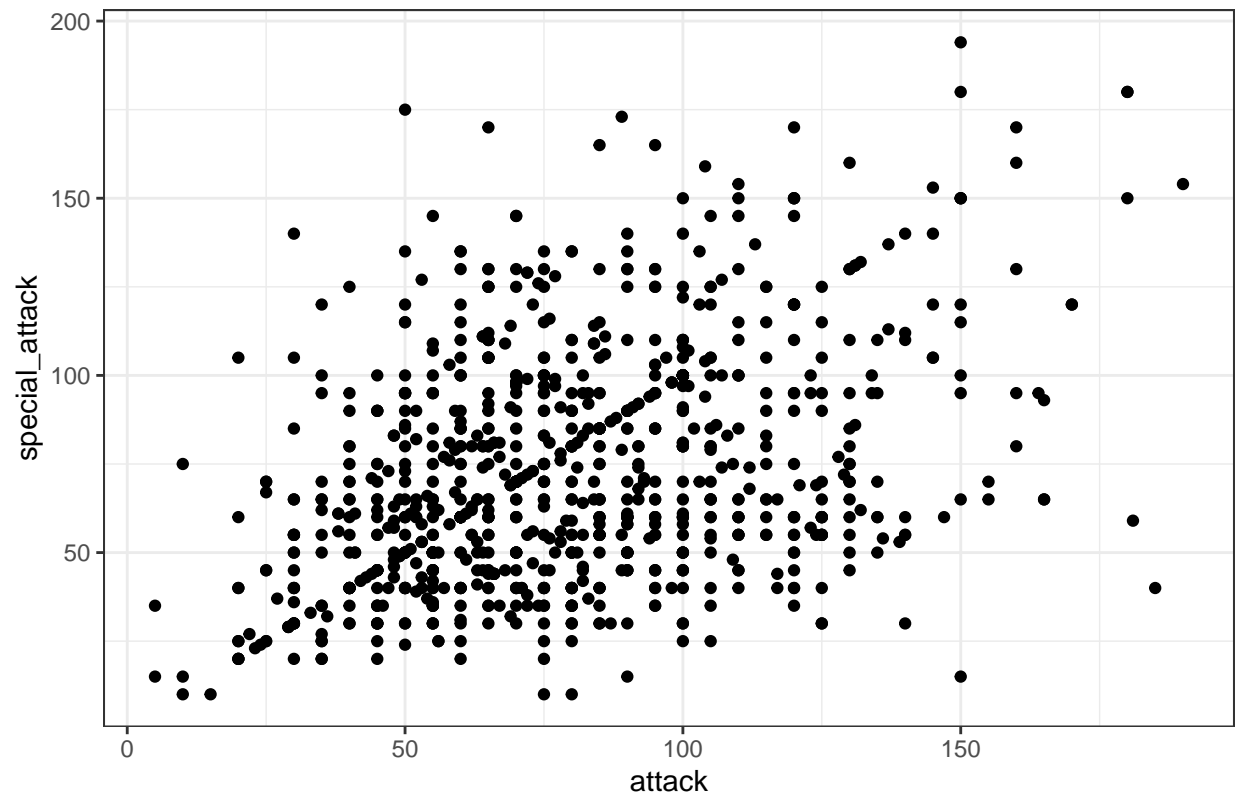
```
p_atk_spa <- ggplot(pokemon, aes(attack, special_attack)) + geom_point() + theme_bw() + labs(title="Sp
```

```
p_def_spd <- ggplot(pokemon, aes(defense, special_defense)) + geom_point() + theme_bw() + labs(title="S
```

```
save_plot(p_atk_spa, "figures/k_special_attack_vs_attack.png")
```

```
## Saved: figures/k_special_attack_vs_attack.png
```

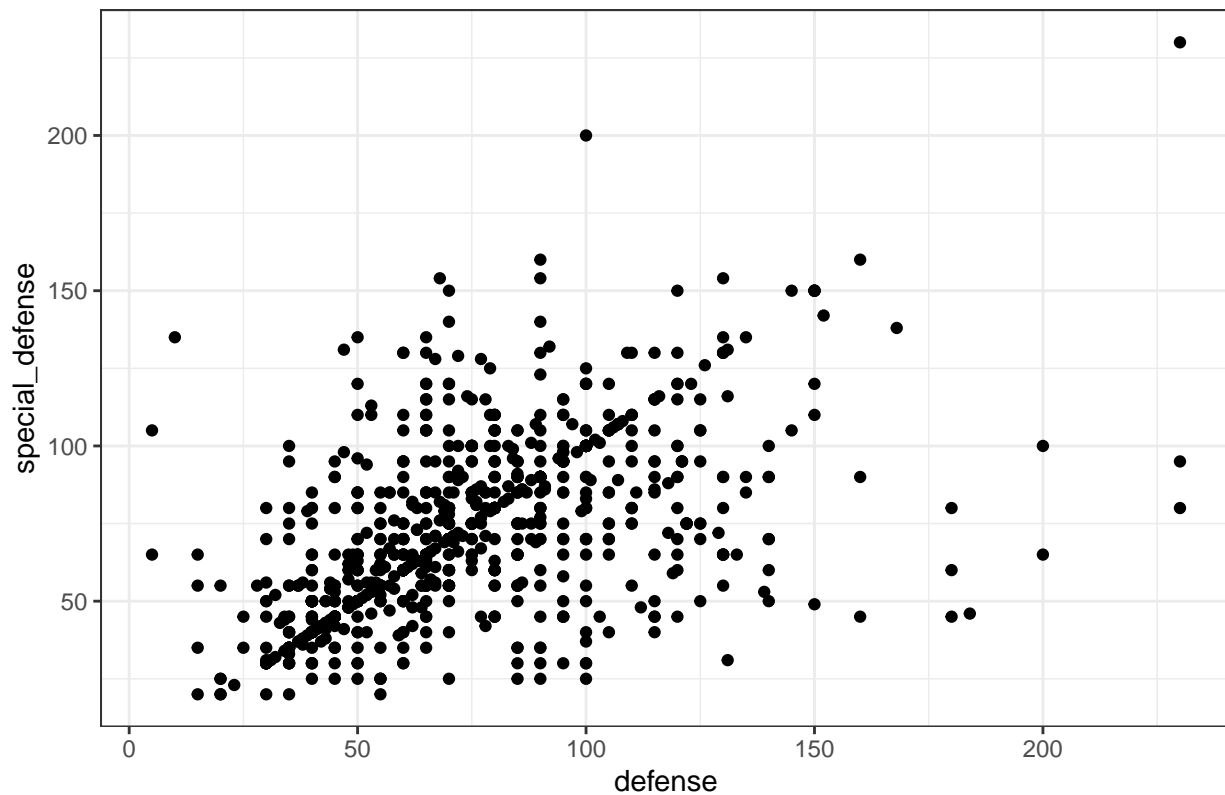
Special Attack vs Attack



```
save_plot(p_def_spd, "figures/k_special_defense_vs_defense.png")
```

```
## Saved: figures/k_special_defense_vs_defense.png
```


Special Defense vs Defense



```
eq_counts <- tibble(
  attack_eq_spatk = sum(pokemon$attack == pokemon$special_attack, na.rm = TRUE),
  defense_eq_spdef = sum(pokemon$defense == pokemon$special_defense, na.rm = TRUE),
  both_equal = sum((pokemon$attack == pokemon$special_attack) &
    (pokemon$defense == pokemon$special_defense), na.rm = TRUE)
)
save_tbl(eq_counts, "tables/k_equal_stat_counts.csv")
```

p. 6. :contentReference[0]

attack_eq_spatk	defense_eq_spdef	both_equal
157	321	116

Hypothesis triptych + simple correlations

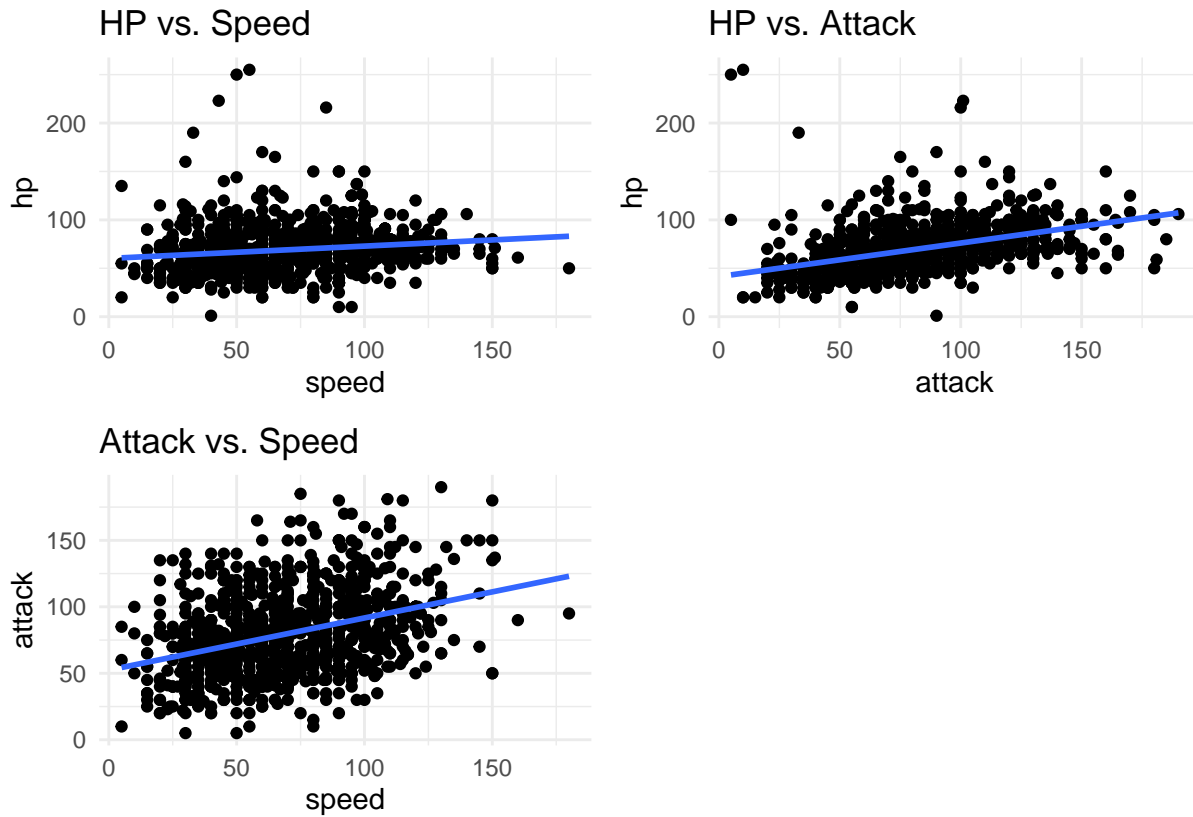
```
p1 <- ggplot(pokemon, aes(speed, hp)) + geom_point() + geom_smooth(method='lm', se=FALSE) + theme_m
p2 <- ggplot(pokemon, aes(attack, hp)) + geom_point() + geom_smooth(method='lm', se=FALSE) + theme_m
p3 <- ggplot(pokemon, aes(speed, attack)) + geom_point() + geom_smooth(method='lm', se=FALSE) + theme_m
g_triptych <- p1 + p2 + p3 + patchwork::plot_layout(ncol=2)
ggsave("figures/k_hypothesis_triptych.png", plot = g_triptych, width = 9, height = 6, dpi = 150)
```

```
## 'geom_smooth()' using formula = 'y ~ x'
## 'geom_smooth()' using formula = 'y ~ x'
## 'geom_smooth()' using formula = 'y ~ x'
```

g_triptych

pp. 13-14. :contentRe

```
## 'geom_smooth()' using formula = 'y ~ x'  
## 'geom_smooth()' using formula = 'y ~ x'  
## 'geom_smooth()' using formula = 'y ~ x'
```



```
corrs_simple <- pokemon |>  
  summarise(`cor(speed, hp)` = cor(speed, hp),  
            `cor(speed, attack)` = cor(speed, attack),  
            `cor(attack, hp)` = cor(attack, hp))  
save_tbl(corrs_simple, "tables/k_simple_correlations.csv")
```

p. 9. :contentReferenc

cor(speed, hp)	cor(speed, attack)	cor(attack, hp)
0.143	0.359	0.427

```
# egg_group_1 counts  
egg1_counts <- pokemon |> count(egg_group_1) |> arrange(desc(n))  
save_tbl(egg1_counts, "tables/k_egg_group1_counts.csv")
```

egg_group_1	n
ground	218
no-eggs	118
monster	94
water1	84
bug	80
mineral	67
indeterminate	64
flying	59
humanshape	40
fairy	38
plant	35
water2	21
water3	16
dragon	14
ditto	1

```
knitr::kable(egg1_counts, caption = "Counts of egg_group_1")
```

p. 9. :contentReferenc

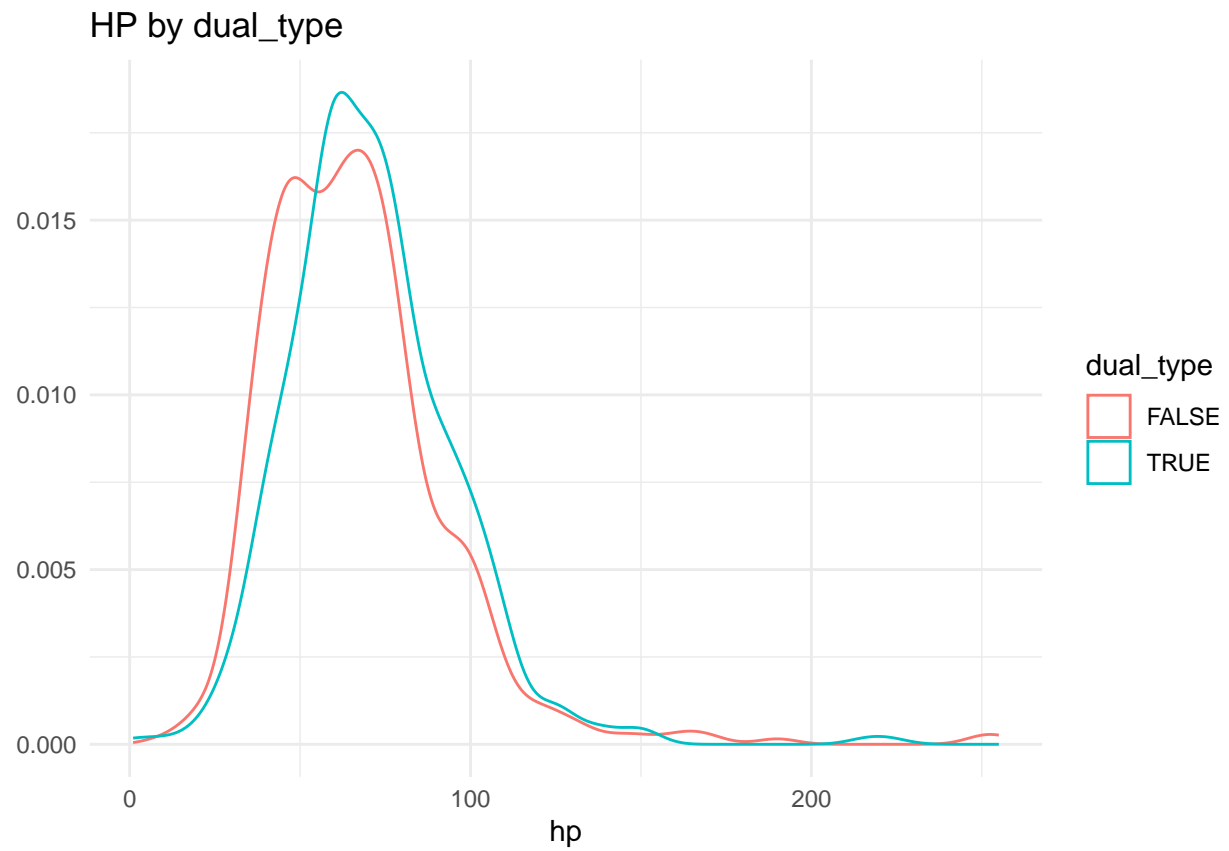
Table 20: Counts of egg_group_1

egg_group_1	n
ground	218
no-eggs	118
monster	94
water1	84
bug	80
mineral	67
indeterminate	64
flying	59
humanshape	40
fairy	38
plant	35
water2	21
water3	16
dragon	14
ditto	1

```
# Dual-type densities
```

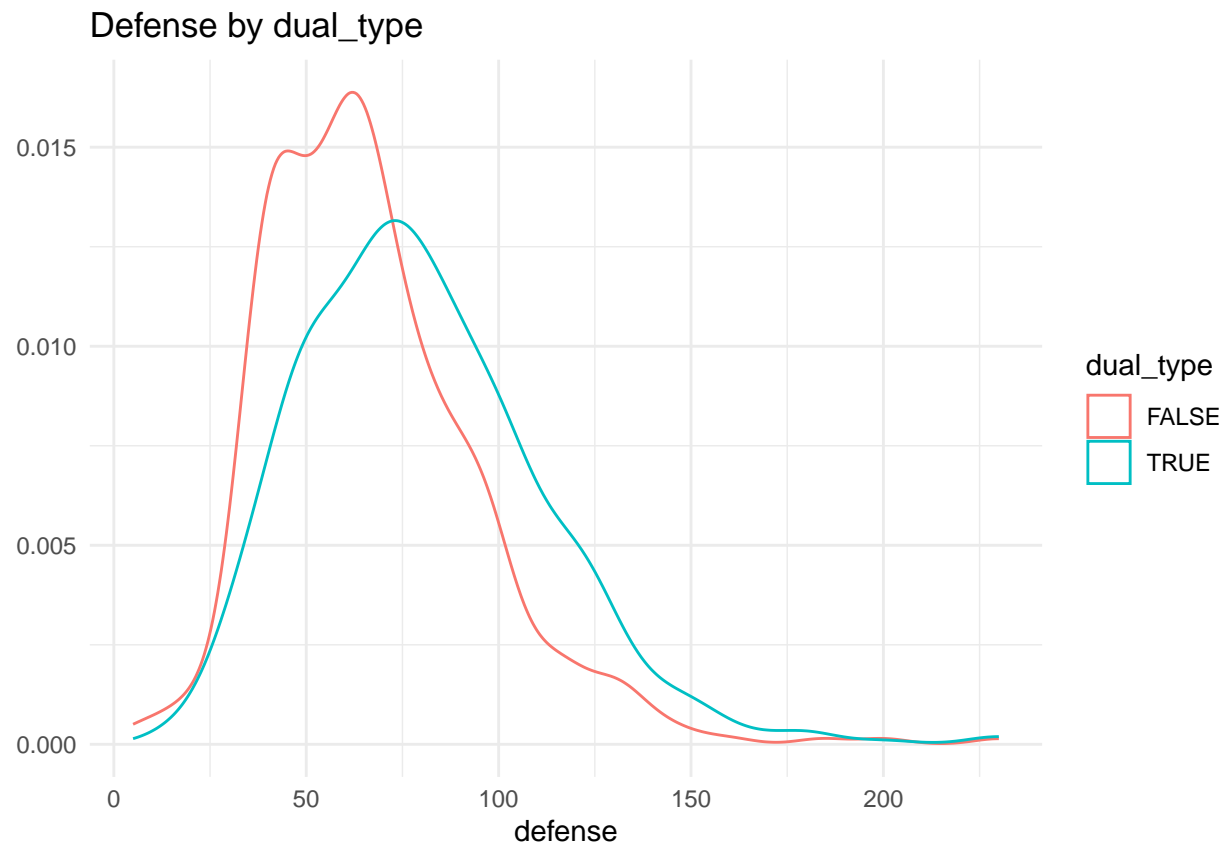
```
p_den_hp <- ggplot(pokemon, aes(hp, color = dual_type)) + geom_density() + theme_minimal()
p_den_def <- ggplot(pokemon, aes(defense, color = dual_type)) + geom_density() + theme_minimal()
p_den_exp <- ggplot(pokemon, aes(base_experience, color = dual_type)) + geom_density() + theme_minimal()
save_plot(p_den_hp, "figures/k_density_hp_by_dualtype.png")
```

```
## Saved: figures/k_density_hp_by_dualtype.png
```



```
save_plot(p_den_def, "figures/k_density_defense_by_dualtype.png")
```

```
## Saved: figures/k_density_defense_by_dualtype.png
```

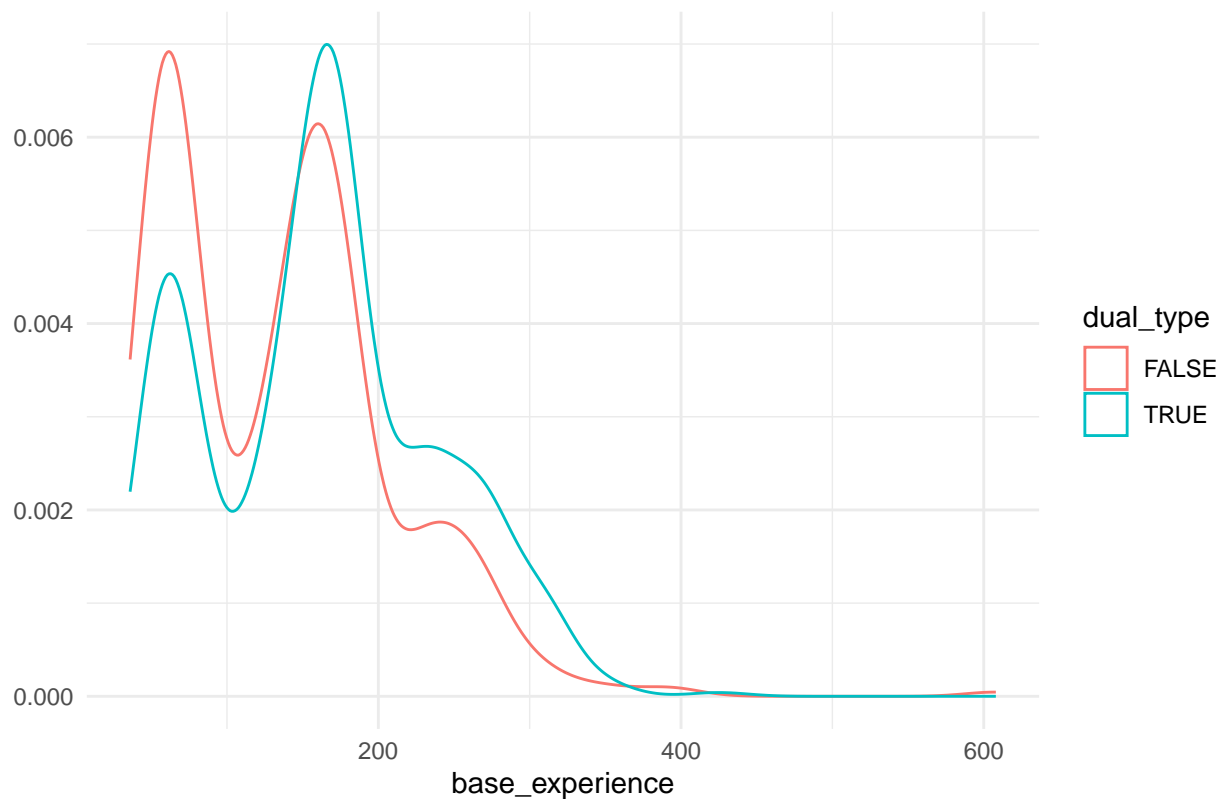


```
save_plot(p_den_exp, "figures/k_density_baseexp_by_dualtype.png")
```

pp. 10-11. :contentRef

```
## Saved: figures/k_density_baseexp_by_dualtype.png
```

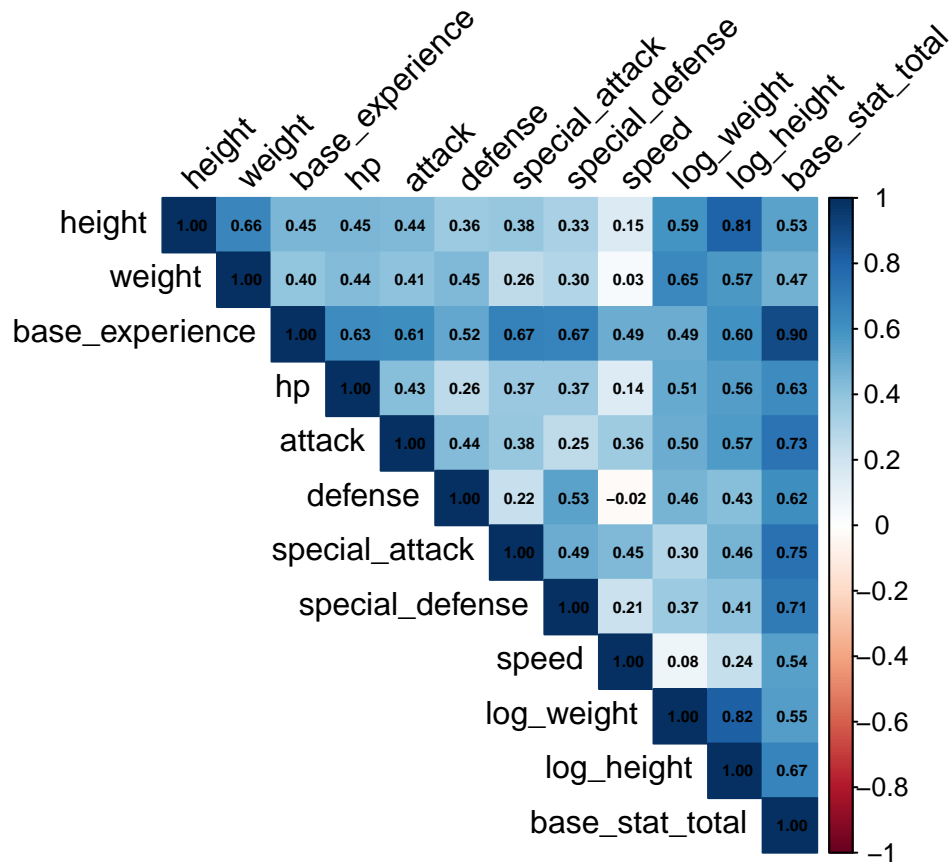
Base experience by dual_type



```
dual_means_atk <- pokemon |> group_by(dual_type) |> summarise(`mean(attack)` = mean(attack, na.rm=TRUE))
save_tbl(dual_means_atk, "tables/k_dual_type_mean_attack.csv")
```

dual_type	mean(attack)
FALSE	74.415
TRUE	83.825

```
# Correlation heatmap (Kareena used corrplot on numeric cols excl. generation_id)
num_keep <- pokemon |> select(where(is.numeric)) |> select(!generation_id)
M <- cor(num_keep, use = "pairwise.complete.obs")
corrplot::corrplot(M, method = "color", type = "upper", addCoef.col = "black", tl.col = "black",
  number.cex = 0.5, tl.srt = 45)
```



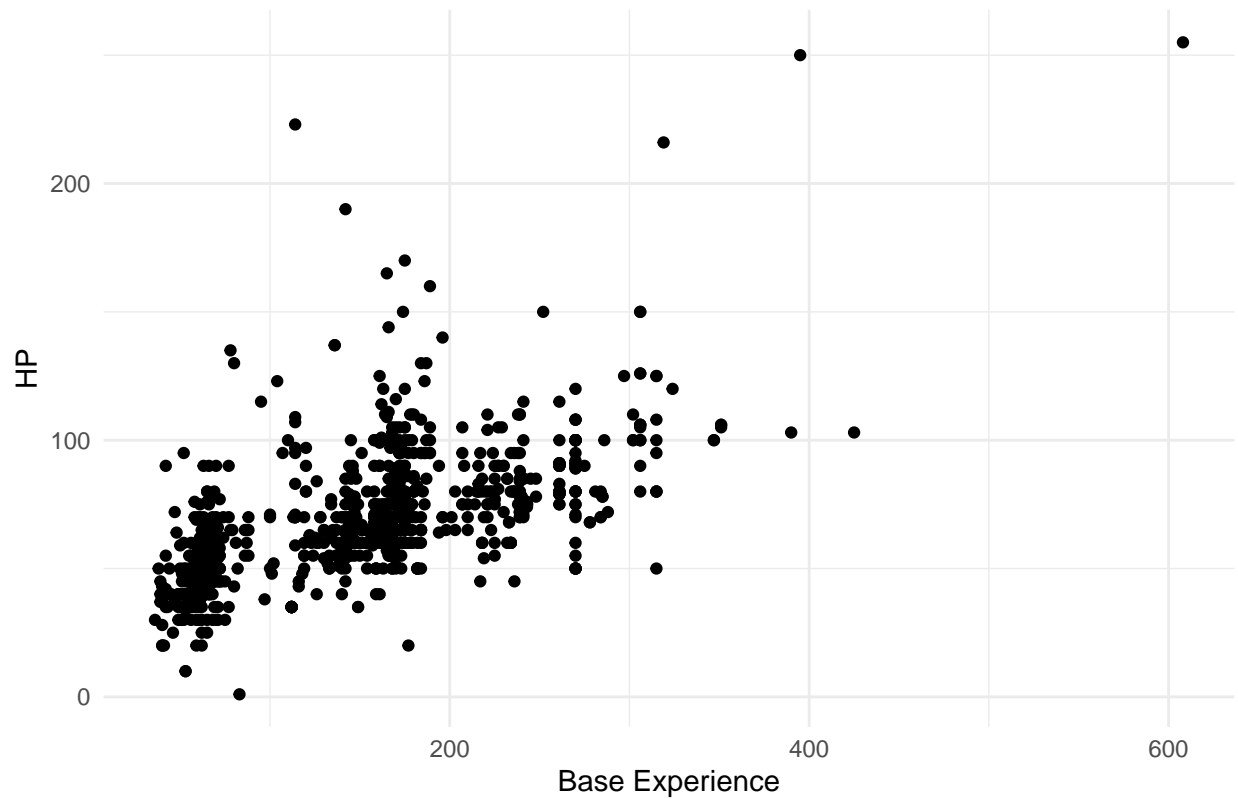
```
png("figures/k_corrplot_matrix.png", width = 1200, height = 1200, res = 140)
corrplot::corrplot(M, method = "color", type = "upper", addCoef.col = "black", tl.col = "black",
  number.cex = 0.5, tl.srt = 45)
dev.off()
```

```
## pdf
## 2
```

```
# HP vs Base Experience + identical-value highlights (Kareena)
p_hp_exp <- ggplot(pokemon, aes(base_experience, hp)) + geom_point() + theme_minimal() +
  labs(title = "HP vs. Base Experience", x = "Base Experience", y = "HP")
save_plot(p_hp_exp, "figures/k_hp_vs_base_experience.png")
```

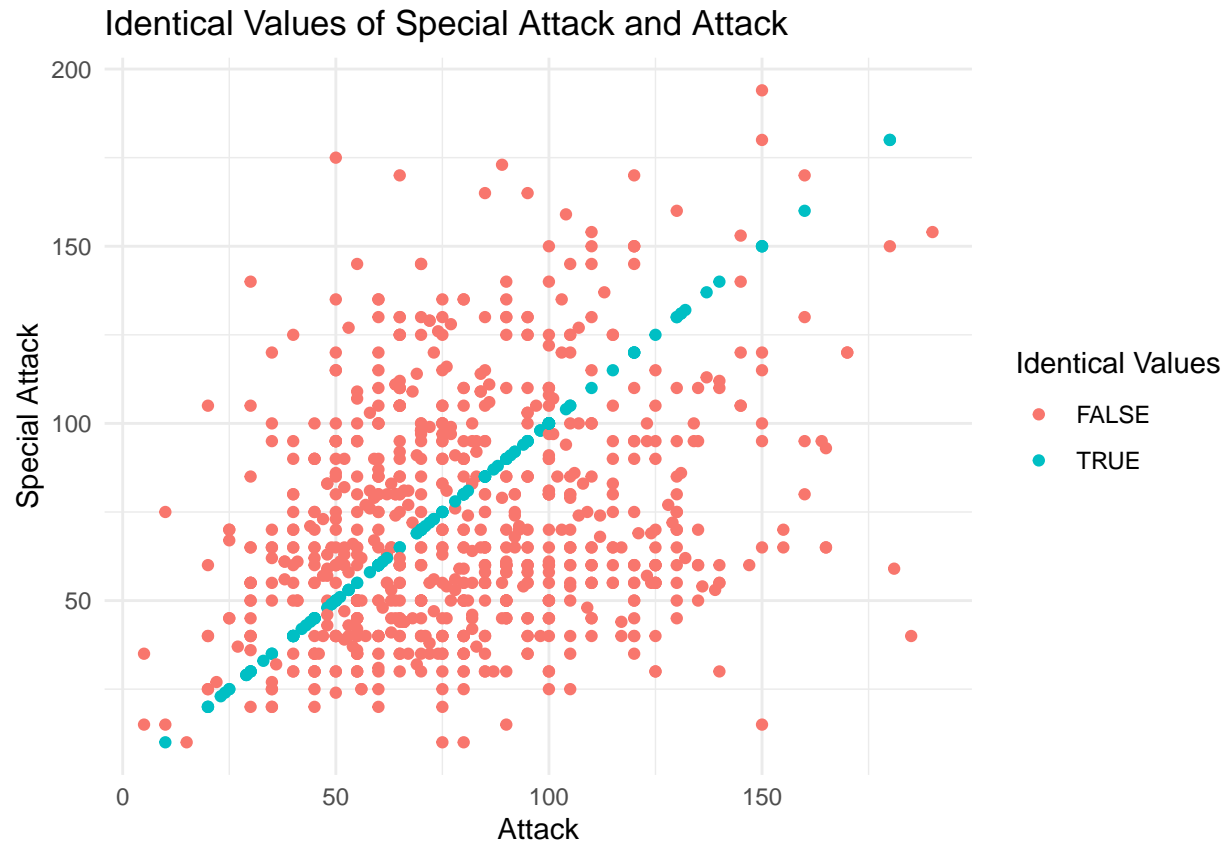
```
## Saved: figures/k_hp_vs_base_experience.png
```

HP vs. Base Experience



```
special_df <- pokemon |>
  mutate(same_attack = special_attack == attack,
         same_defense = special_defense == defense)
p_same_atk <- ggplot(special_df, aes(attack, special_attack, color = same_attack)) +
  geom_point() + theme_minimal() + labs(title="Identical Values of Special Attack and Attack", x="Attack")
p_same_def <- ggplot(special_df, aes(defense, special_defense, color = same_defense)) +
  geom_point() + theme_minimal() + labs(title="Identical Values of Special Defense and Defense", x="Defense")
save_plot(p_same_atk, "figures/k_identical_spatk_attack.png")
```

```
## Saved: figures/k_identical_spatk_attack.png
```

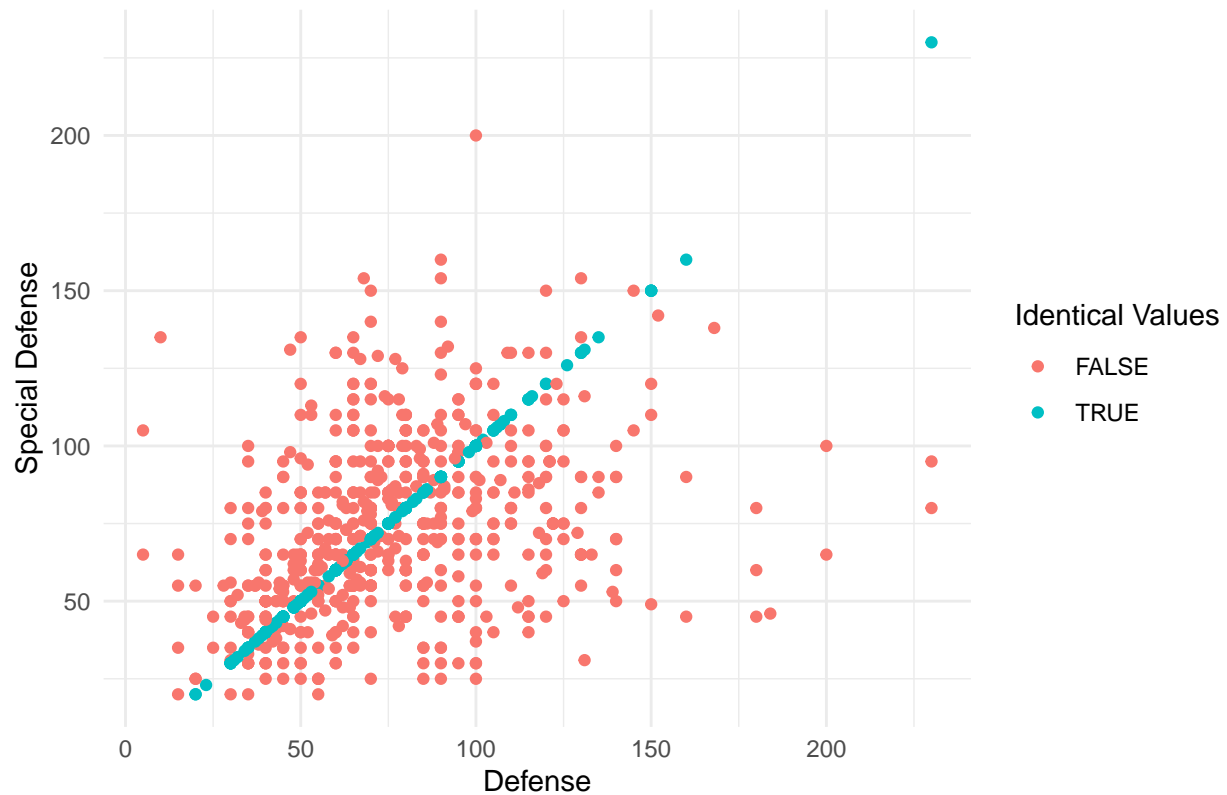



```
save_plot(p_same_def, "figures/k_identical_spdef_defense.png")
```

pp. 16-17. :contentRef

```
## Saved: figures/k_identical_spdef_defense.png
```

Identical Values of Special Defense and Defense

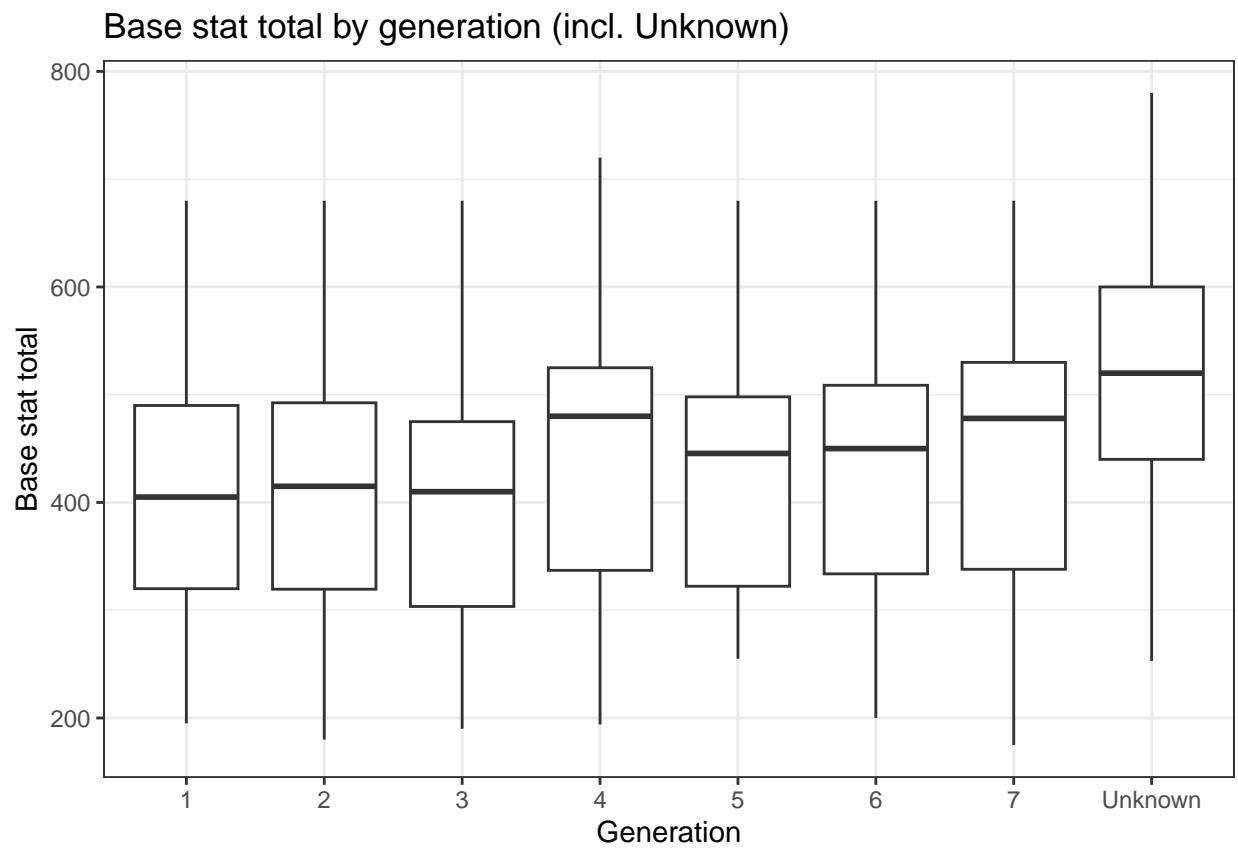


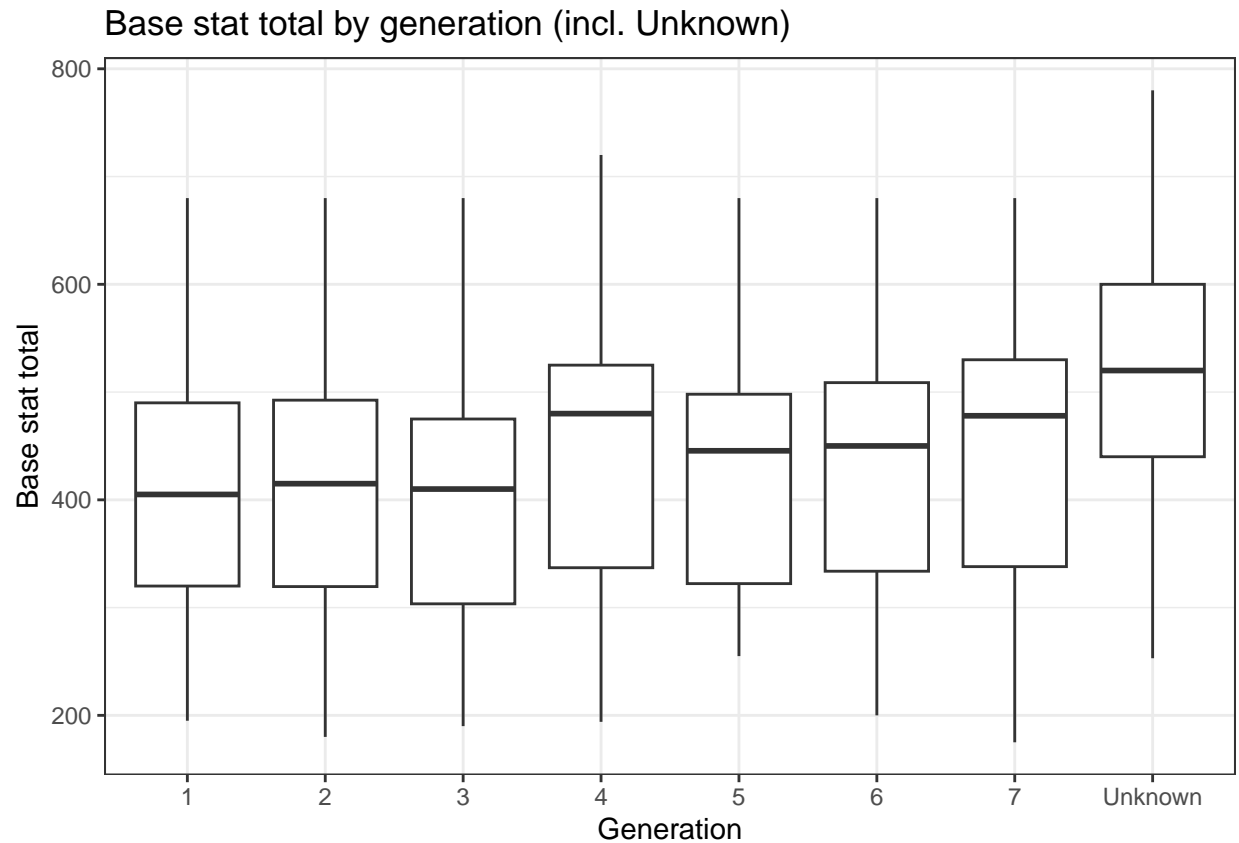
```
# Power creep: base_stat_total by generation (table + boxplot)
gen_total <- pokemon |>
  group_by(generation_id_f) |>
  summarise(n = n(),
            mean_total = mean(base_stat_total, na.rm=TRUE),
            sd_total = sd(base_stat_total, na.rm=TRUE)) |>
  arrange(generation_id_f)
save_tbl(gen_total, "tables/gen_base_stat_total_summary.csv", digits = 1)
```

generation_id_f	n	mean_total	sd_total
1	151	407.6	99.9
2	100	407.2	112.5
3	135	403.7	115.6
4	107	445.8	117.5
5	156	425.8	102.5
6	72	429.6	111.9
7	81	443.8	118.7
Unknown	147	520.5	122.9

```
p_gen_total <- ggplot(pokemon, aes(generation_id_f, base_stat_total)) +
  geom_boxplot(outlier.alpha=.4) + theme_bw() +
  labs(title="Base stat total by generation (incl. Unknown)", x="Generation", y="Base stat total")
save_plot(p_gen_total, "figures/box_basetotal_by_generation.png") ; p_gen_total
```

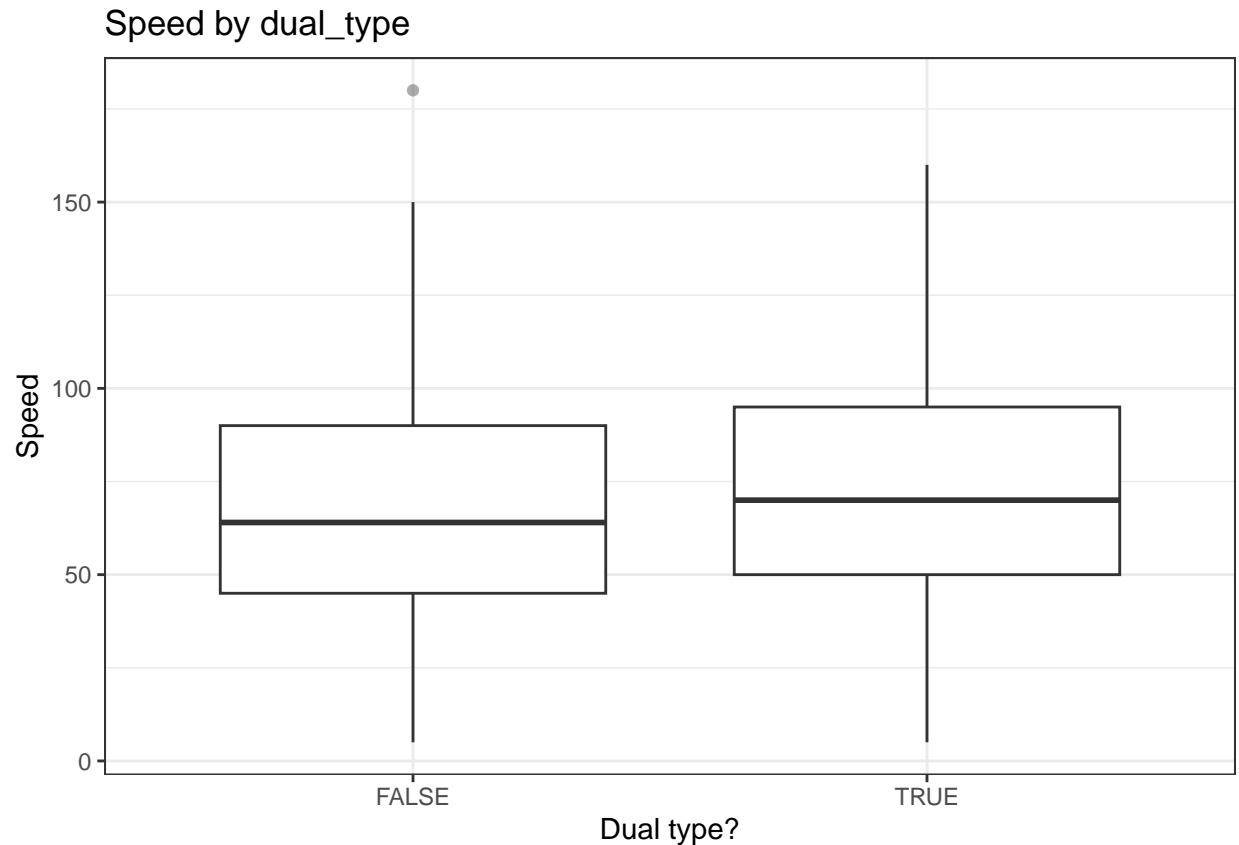
Saved: figures/box_basetotal_by_generation.png





```
# Speed by dual_type (boxplot + quick t-test)
p_speed_dual <- ggplot(pokemon, aes(x = dual_type, y = speed)) +
  geom_boxplot(outlier.alpha=.4) + theme_bw() +
  labs(title="Speed by dual_type", x="Dual type?", y="Speed")
save_plot(p_speed_dual, "figures/box_speed_by_dualtype.png")
```

```
## Saved: figures/box_speed_by_dualtype.png
```



```
tt_speed <- broom::tidy(t.test(speed ~ dual_type, data = pokemon))
save_tbl(tt_speed, "tables/t_test_speed_dualtype.csv")
```

estimate	estimate1	estimate2	statistic	p.value	parameter	conf.low	conf.high	method	alternative
-5.163	66.248	71.412	-2.726	0.007	927.629	-8.881	-1.446	Welch Two Sample t-test	two.sided

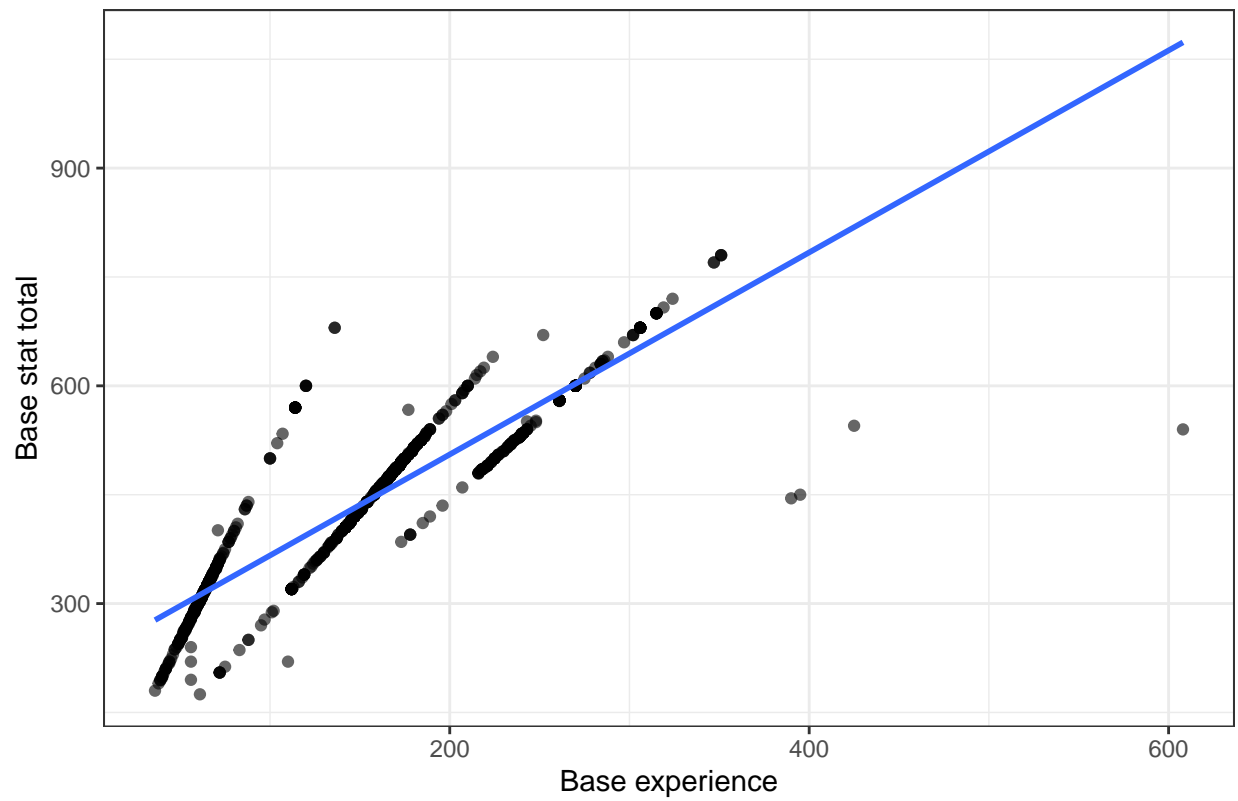
```
# Base experience vs base stat total (scatter + OLS line)
p_exp_total <- ggplot(pokemon, aes(base_experience, base_stat_total)) +
  geom_point(alpha=.6) + geom_smooth(method="lm", se=FALSE) + theme_bw() +
  labs(title="Base experience vs base stat total", x="Base experience", y="Base stat total")
save_plot(p_exp_total, "figures/baseexp_vs_basetotal.png") ; p_exp_total
```

```
## 'geom_smooth()' using formula = 'y ~ x'
```

```
## Saved: figures/baseexp_vs_basetotal.png
```

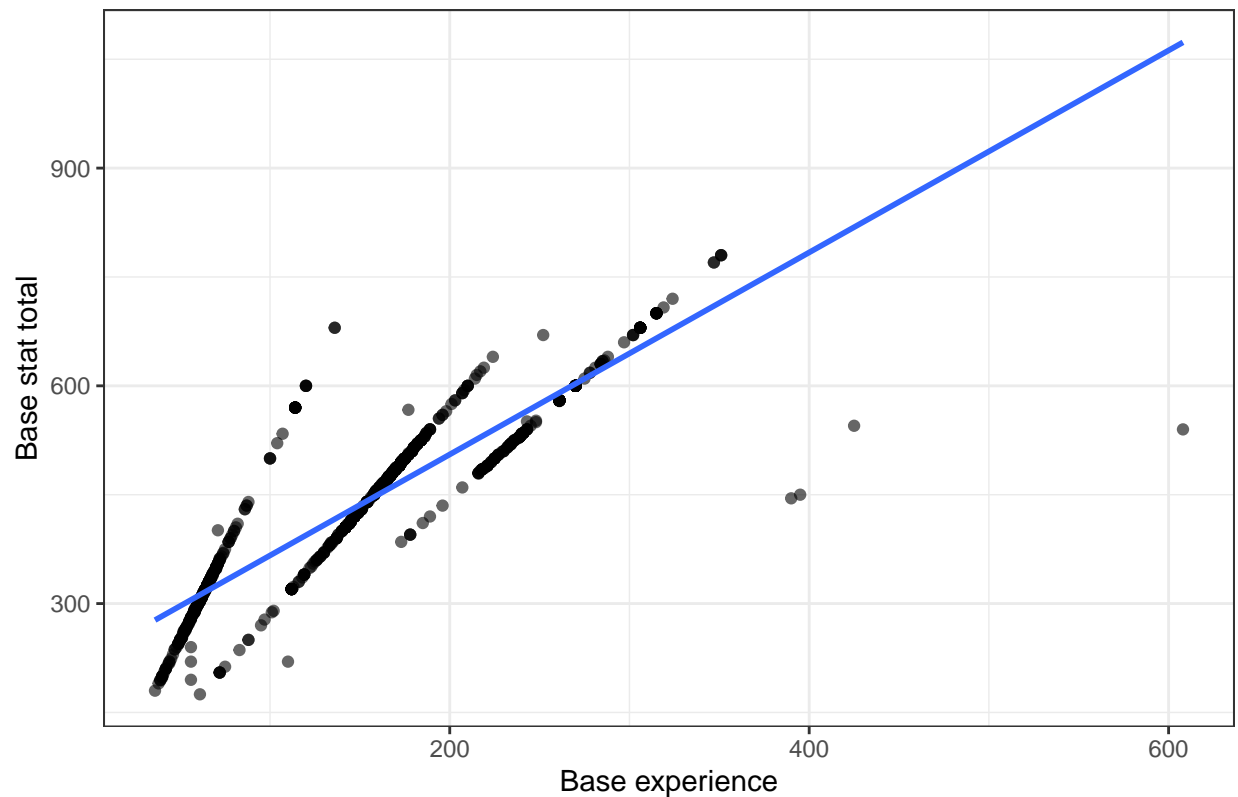
```
## 'geom_smooth()' using formula = 'y ~ x'
```

Base experience vs base stat total



```
## 'geom_smooth()' using formula = 'y ~ x'
```

Base experience vs base stat total



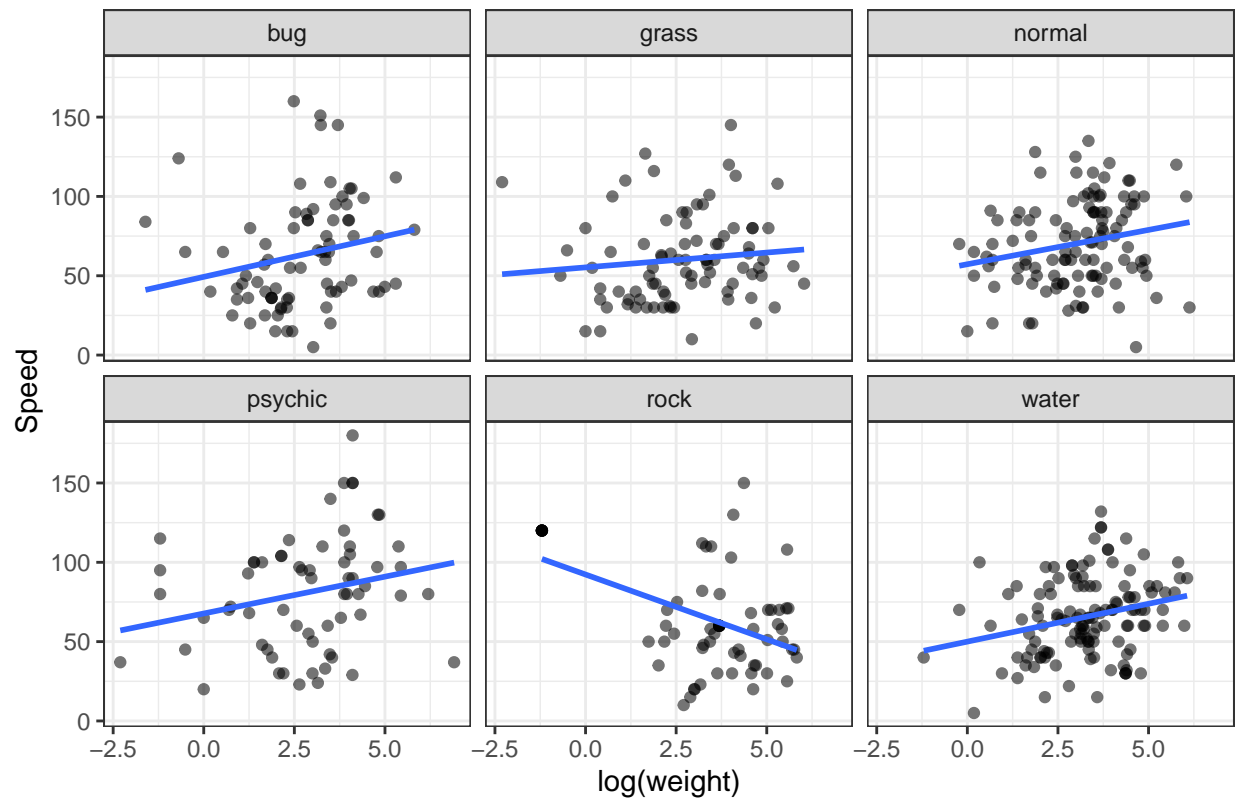
```
# Speed vs log(weight), faceted by top 6 primary types
top6 <- pokemon |> count(type_1) |> arrange(desc(n)) |> slice_head(n=6) |> pull(type_1)
p_speed_lw_facets <- pokemon |>
  filter(type_1 %in% top6, !is.na(log_weight)) |>
  ggplot(aes(log_weight, speed)) +
  geom_point(alpha=.55) + geom_smooth(method="lm", se=FALSE) +
  facet_wrap(~type_1, ncol=3) + theme_bw() +
  labs(title="Speed vs log(weight): top 6 primary types", x="log(weight)", y="Speed")
save_plot(p_speed_lw_facets, "figures/speed_vs_logweight_top6_types.png", w=9, h=6) ; p_speed_lw_facets
```

```
## 'geom_smooth()' using formula = 'y ~ x'
```

```
## Saved: figures/speed_vs_logweight_top6_types.png
```

```
## 'geom_smooth()' using formula = 'y ~ x'
```

Speed vs log(weight): top 6 primary types



```
## 'geom_smooth()' using formula = 'y ~ x'
```


Speed vs log(weight): top 6 primary types

