

Importance of open source

Overview

This script creates bar plots from question 2 on the survey, which is about the perceived importance of open source for different job categories and different tasks. My favorite plot is the one in the final section, “Percent more than moderately important”.

Import packages and utilities

```
project_root <- here::here() # requires that you be somewhere in the
# project directory (not above it)
# packages
suppressMessages(source(file.path(project_root, "scripts/packages.R")))
# functions and objects used across scripts
suppressMessages(source(file.path(project_root, "scripts/utils.R")))
```

Define functions

get_percent_more_than_moderate

- Arguments:
 - **df**: A data frame with a column named **importance_level**. Should contain only rows that you want to count. Extra columns are okay.
- Details:
 - A simple function to count what percent of rows in a data frame have an **importance_level** of more than Moderately important. Checks that there are no extra rows with an unexpected value in the **importance_level** column.

- Outputs:
 - A scalar value representing the percentage of respondents who selected an `importance_level` of more than Moderately important, rounded to 2 decimal places.

```
get_percent_more_than_moderate <- function(df) {
  # check that df has the required column
  if (!"importance_level" %in% names(df)) {
    stop("`df` must have a column named 'importance_level'.")
  }

  high_importance_cats <- c(
    "Very important",
    "Important"
  )
  low_importance_cats <- c(
    "Moderately important",
    "Slightly important",
    "Not at all important"
  )

  n_high_rows <- nrow(df %>% filter(importance_level %in% high_importance_cats))
  n_low_rows <- nrow(df %>% filter(importance_level %in% low_importance_cats))
  total <- nrow(df)
  if (n_high_rows + n_low_rows != total) {
    stop("`df` has at least one unexpected value in 'importance_level'.")
  }
  pct <- round(n_high_rows / total * 100, 2)
  return(pct)
}
```

Load data

```
data <- load_qualtrics_data("deidentified_no_qual.tsv")
```

Wrangle data

Let's extract the columns we care about.

```
importance_and_job <- data %>%
  select(
    starts_with("importance_opensrc") | starts_with("job_category")
  )
head(importance_and_job)
```

	importance_opensrc_1	importance_opensrc_2	importance_opensrc_3
1	Very important	Very important	Very important
2	Very important	Moderately important	Important
3	Very important	Very important	Very important
4	Very important	Slightly important	Important
5	Very important	Important	Very important
6	Very important	Non-applicable	Important

	importance_opensrc_4	importance_opensrc_5
1	Very important	Very important
2	Important	Non-applicable
3	Very important	Non-applicable
4	Important	Non-applicable
5	Very important	Non-applicable
6	Important	Non-applicable

	job_category
1	Faculty
2	Post-Doc
3	Other research staff (e.g., research scientist, research software engineer)
4	Faculty
5	Faculty
6	Other research staff (e.g., research scientist, research software engineer)

Let's reshape the data from wide to long format.

```
long_data <- importance_and_job %>%
  pivot_longer(
    cols = starts_with("importance_opensrc"),
    names_to = "importance_area",
    values_to = "importance_level"
  )

long_data <- long_data %>%
  mutate(
    importance_area = recode(
      importance_area,
```

```

      "importance_opensrc_1" = "Research",
      "importance_opensrc_2" = "Teaching",
      "importance_opensrc_3" = "Learning",
      "importance_opensrc_4" = "Professional Development",
      "importance_opensrc_5" = "Job"
    )
  )
long_data

```

```

# A tibble: 1,660 x 3
  job_category importance_area importance_level
  <chr>         <chr>         <chr>
1 Faculty      Research      Very important
2 Faculty      Teaching      Very important
3 Faculty      Learning      Very important
4 Faculty      Professional Development Very important
5 Faculty      Job           Very important
6 Post-Doc     Research      Very important
7 Post-Doc     Teaching      Moderately important
8 Post-Doc     Learning      Important
9 Post-Doc     Professional Development Important
10 Post-Doc     Job           Non-applicable
# i 1,650 more rows

```

STOP!!! At this point, I manually compared this data frame to the results table in Qualtrics to make sure the columns (e.g. `importance_opensrc_1`) correspond to the options I expect (e.g. “Research”). I had to use peoples’ email addresses to make sure I was comparing the same rows in each table. I assumed that the variables were ordered by their order on the survey, but you never know. In this case, my assumption was correct. I’ve commented out the code for this because it only needed to be done once.

```

# pii <- load_qualtrics_data("pii.tsv")
# emails <- pii %>%
#   select(starts_with("stay_in_touch_email"))

# t <- cbind(emails, importance_and_job)
# subset(t, startsWith(stay_in_touch_email, "PERSON_EMAIL_HERE"))

```

Back to data wrangling.

Here, I removed all rows that contain an empty string in any column. Since both questions were mandatory, I’m actually only removing people who never saw the demographic questions:

people who are not affiliated with UC (2) + people who are neither past nor future open source contributors (36). $(2+36)*5$ importance areas = 190 rows removed.

```
dim(long_data)
```

```
[1] 1660    3
```

```
long_data <- long_data %>%  
  filter(!if_any(everything(), ~ . == ""))
```

```
dim(long_data)
```

```
[1] 1470    3
```

Shorten this one long category name. Other research staff (e.g., research scientist, research software engineer) becomes simply Other.

```
long_data$job_category <- gsub(  
  "^Other.*",  
  "Research Staff",  
  long_data$job_category  
)
```

Reorder factor levels for plotting.

```
long_data$importance_level <- factor(  
  long_data$importance_level,  
  levels = c(  
    "Very important",  
    "Important",  
    "Moderately important",  
    "Slightly important",  
    "Not at all important",  
    "Non-applicable"  
  ),  
  ordered = TRUE  
)  
long_data
```

```
# A tibble: 1,470 x 3
  job_category importance_area importance_level
  <chr>         <chr>         <ord>
1 Faculty      Research      Very important
2 Faculty      Teaching      Very important
3 Faculty      Learning      Very important
4 Faculty      Professional Development Very important
5 Faculty      Job          Very important
6 Post-Doc     Research      Very important
7 Post-Doc     Teaching      Moderately important
8 Post-Doc     Learning      Important
9 Post-Doc     Professional Development Important
10 Post-Doc     Job          Non-applicable
# i 1,460 more rows
```

Bar plots

Simple bar plot for teachers

Now let's start making some bar plots. Let's start by making a bar plot showing how teachers rate the importance of open source for their teaching. Since we didn't ask people "Do you teach?", and since there was a "Non-applicable" option, we will simply assume that if they gave an answer for the "Teaching" option, they must be a teacher.

```
teaching <- long_data %>%
  filter(
    importance_area == "Teaching"
  ) %>%
  filter(
    importance_level != "Non-applicable"
  )

# For our bar plot, we only care about how many times each 'importance level' was selected.
teaching_to_plot <- teaching %>% select(-c(job_category, importance_area))

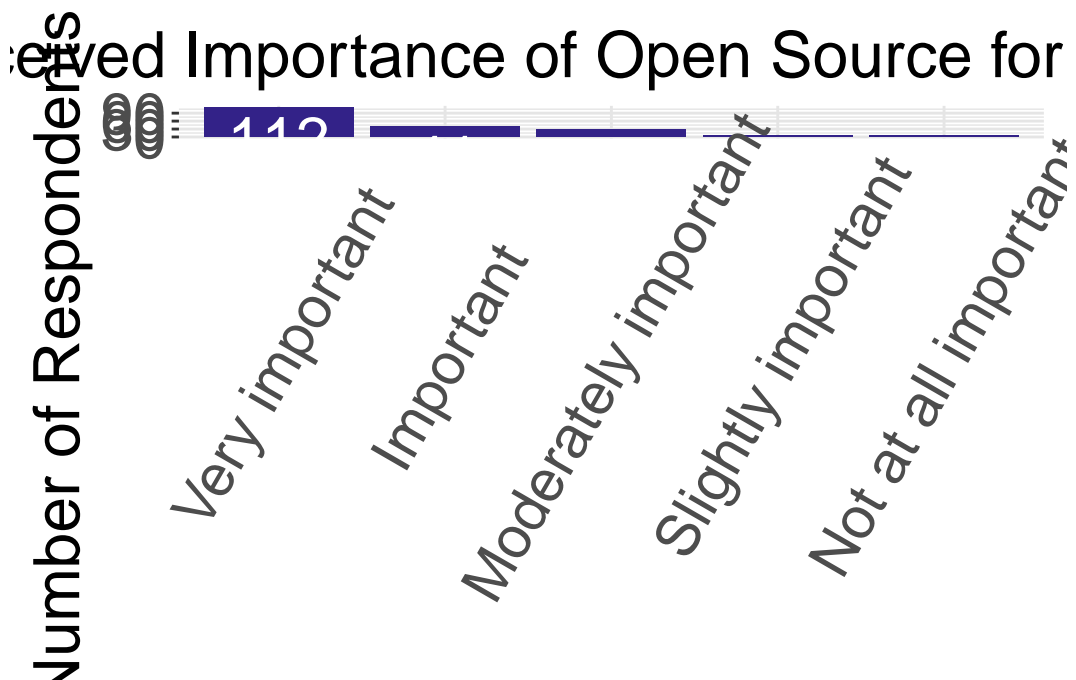
teaching_to_plot <- teaching_to_plot %>%
  count(importance_level, name = "Counts")

# By sheer luck, the columns are already ordered by response rates.
teaching_to_plot
```

```
# A tibble: 5 x 2
  importance_level Counts
  <ord>           <int>
1 Very important    112
2 Important         41
3 Moderately important 27
4 Slightly important  7
5 Not at all important 6
```

Now let's make that bar chart using a function that lives in my utilities script (`scripts/utls.R`).

```
basic_bar_chart(
  teaching_to_plot,
  x_var = "importance_level",
  y_var = "Counts",
  title = "Perceived Importance of Open Source for Teaching",
  ylabel = "Number of Respondents (Teachers Only)",
  show_bar_labels = TRUE
)
```



Save the plot using a function that lives in my utilities script (`scripts/utls.R`).

```
#save_plot("importance_teachers.tiff", 8, 5)
```

Grouped bar plot for researchers

Now let's look at researchers, and the importance categories that apply to all researchers. The importance categories again are:

Research

Teaching → Does not apply

Learning

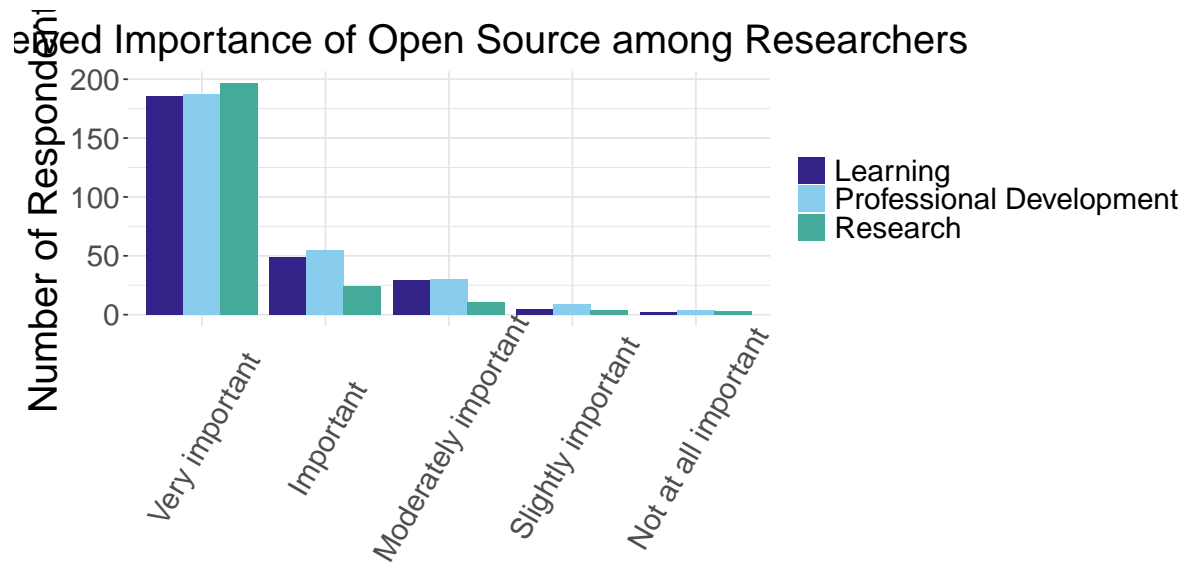
Professional Development

Job (For non-research staff) → Does not apply

So we'll make a bar plot with just those three categories that apply to all researchers. As with teachers above, we will assume that if they didn't select "Non-applicable", they must be a researcher.

The `grouped_bar_chart` function, like the `basic_bar_chart` function, lives in my utility script.

```
research_learning_pd <- long_data %>%  
  filter(  
    importance_area == "Research" |  
    importance_area == "Learning" |  
    importance_area == "Professional Development"  
  ) %>%  
  filter(importance_level != "Non-applicable")  
  
grouped_bar_chart(  
  df = research_learning_pd,  
  x_var = "importance_level",  
  fill_var = "importance_area",  
  title = "Perceived Importance of Open Source among Researchers"  
)
```

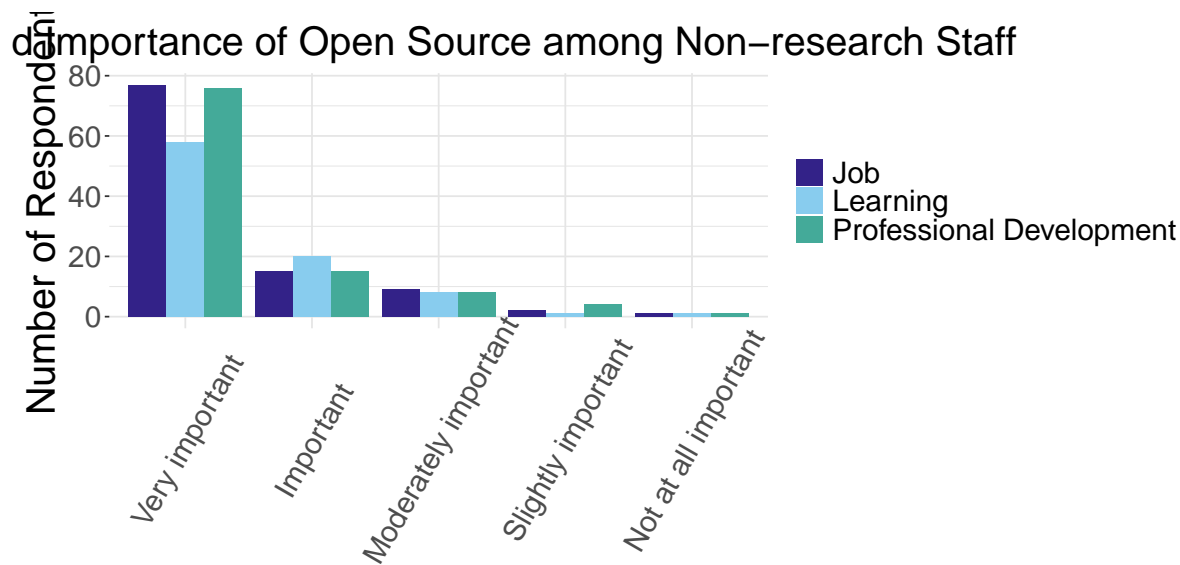
```
#save_plot("importance_researchers.tiff", 10, 5)
```

Grouped bar plot for non-research staff

This is very similar to what I did above, except the three applicable categories have changed.

```
job_learning_pd <- long_data %>%
  filter(
    importance_area == "Job" |
    importance_area == "Learning" |
    importance_area == "Professional Development"
  ) %>%
  filter(job_category == "Non-research Staff") %>%
  filter(importance_level != "Non-applicable")

grouped_bar_chart(
  df = job_learning_pd,
  x_var = "importance_level",
  fill_var = "importance_area",
  title = "Perceived Importance of Open Source among Non-research Staff"
)
```



```
#save_plot("importance_nrstaff.tiff", 10, 5)
```

Percent more than moderately important

Renata suggested I try to combine all these data into one figure that summarizes the question at a glance. Here's my attempt.

I think a useful “statistic” is the percent of a particular group that said OS is more than moderately important for a particular area of work. Let's make a dataframe with those percentages. I'd ultimately like to turn this into a bar plot where the color or design of the bars corresponds to the five job categories, and the x-axis shows four groups that I think are most relevant: teachers, researchers, non-research staff, and students. However, not all these groups were explicit survey categories, and not all 5 importance areas apply to all groups, so we'll need to do a fair amount of data wrangling.

To start, let's get the percent of teachers who said teaching was more than moderately important.

```
teaching <- long_data %>%
  filter(
    importance_area == "Teaching"
  ) %>%
  filter(
    importance_level != "Non-applicable"
```

```

)

more_than_mod <- data.frame(
  job_category = "Teachers",
  importance_area = "Teaching",
  pct = get_percent_more_than_moderate(teaching)
)

more_than_mod

```

```

  job_category importance_area  pct
1    Teachers      Teaching 79.27

```

The code is basically the same for researchers and researching: as with teaching, we will assume that anyone who gave an answer for research (i.e, didn't select "N/A") is a researcher.

```

research <- long_data %>%
  filter(
    importance_area == "Research"
  ) %>%
  filter(
    importance_level != "Non-applicable"
  )

more_than_mod <- rbind(
  more_than_mod,
  list(
    "Researchers",
    "Research",
    get_percent_more_than_moderate(research)
  )
)

```

The code for non-research staff and “Job” is slightly different. Our intention was that only non-research staff would answer this question, but there were some people who answered this but did not select “non-research staff” as their job category. So let’s just ensure that we’re only looking at responses from non-research staff by filtering for non-research staff using the `job_category` column.

```

nrstaff <- long_data %>%
  filter(
    job_category == "Non-research Staff"
  ) %>%
  filter(
    importance_area == "Job"
  ) %>%
  filter(
    importance_level != "Non-applicable"
  )

more_than_mod <- rbind(
  more_than_mod,
  list(
    "Non-research Staff",
    "Job",
    get_percent_more_than_moderate(nrstaff)
  )
)

```

Next, I want to look at two importance areas, Learning and Professional Development, for all four job categories: Teachers, Researchers, Non-research staff, and Students. We'll have to determine teachers and researchers based on who answered the teaching question and who answered the research question, respectively. Meanwhile, for "Student", we'll have to combine grad students and undergrads into one group.

We will have to go back to an earlier data frame and redo some of the data wrangling. (We want to filter for people who answered e.g. teaching, but look at their answers for the other questions. This information was lost when we rearranged from wide to long format.)

```

# Rename this one long job category
importance_and_job$job_category <- gsub(
  "^Other.*",
  "Research Staff",
  importance_and_job$job_category
)

#Rename columns for readability
importance_and_job <- importance_and_job %>%
  rename(
    Research = importance_opensrc_1,
    Teaching = importance_opensrc_2,

```

```

    Learning = importance_opensrc_3,
    `Professional Development` = importance_opensrc_4,
    Job = importance_opensrc_5
  )

# Remove rows that contain any empty strings
importance_and_job <- importance_and_job %>%
  filter(!if_any(everything(), ~ . == ""))

```

Let's keep rows from teachers, but keep columns for Learning and Professional Development. Then we change the job_category column to "Teacher".

```

teachers_learn_pd <- importance_and_job %>%
  filter(Teaching != "Non-applicable") %>%
  select(Learning, `Professional Development`, job_category)

teachers_learn_pd$job_category <- "Teacher"
head(teachers_learn_pd)

```

	Learning	Professional Development	job_category
1	Very important	Very important	Teacher
2	Important	Important	Teacher
3	Very important	Very important	Teacher
4	Important	Important	Teacher
5	Very important	Very important	Teacher
6	Moderately important	Moderately important	Teacher

Now we can add two more rows to more_than_mod.

```

teachers_learning <- teachers_learn_pd %>%
  select(Learning, job_category) %>%
  #unlikely but you never know
  filter(Learning != "Non-applicable") %>%
  # Change the column name because our function expects it
  rename(importance_level = Learning)

more_than_mod <- rbind(
  more_than_mod,
  list(
    "Teachers",
    "Learning",

```

```

    get_percent_more_than_moderate(teachers_learning)
  )
)

teachers_pd <- teachers_learn_pd %>%
  select(`Professional Development`, job_category) %>%
  #unlikely but you never know
  filter(`Professional Development` != "Non-applicable") %>%
  # Change the column name because our function expects it
  rename(importance_level = `Professional Development`)

more_than_mod <- rbind(
  more_than_mod,
  list(
    "Teachers",
    "Professional Development",
    get_percent_more_than_moderate(teachers_pd)
  )
)
more_than_mod

```

	job_category	importance_area	pct
1	Teachers	Teaching	79.27
2	Researchers	Research	92.47
3	Non-research Staff	Job	88.46
4	Teachers	Learning	84.29
5	Teachers Professional Development		81.48

And let's do the same for researchers.

```

researchers_learn_pd <- importance_and_job %>%
  filter(Research != "Non-applicable") %>%
  select(Learning, `Professional Development`, job_category)

researchers_learn_pd$job_category <- "Researcher"

researchers_learning <- researchers_learn_pd %>%
  select(Learning, job_category) %>%
  #unlikely but you never know
  filter(Learning != "Non-applicable") %>%
  # Change the column name because our function expects it

```

```

    rename(importance_level = Learning)

more_than_mod <- rbind(
  more_than_mod,
  list(
    "Researchers",
    "Learning",
    get_percent_more_than_moderate(researchers_learning)
  )
)

researchers_pd <- researchers_learn_pd %>%
  select(`Professional Development`, job_category) %>%
  #unlikely but you never know
  filter(`Professional Development` != "Non-applicable") %>%
  # Change the column name because our function expects it
  rename(importance_level = `Professional Development`)

more_than_mod <- rbind(
  more_than_mod,
  list(
    "Researchers",
    "Professional Development",
    get_percent_more_than_moderate(researchers_pd)
  )
)
more_than_mod

```

	job_category	importance_area	pct
1	Teachers	Teaching	79.27
2	Researchers	Research	92.47
3	Non-research Staff	Job	88.46
4	Teachers	Learning	84.29
5	Teachers Professional Development		81.48
6	Researchers	Learning	85.90
7	Researchers Professional Development		83.98

Now get percentages for non-research staff. This is straightforward since it was a survey category.

```

nrstaff_learn_pd <- importance_and_job %>%
  filter(job_category == "Non-research Staff") %>%
  select(Learning, `Professional Development`, job_category)

nrstaff_learning <- nrstaff_learn_pd %>%
  select(Learning, job_category) %>%
  #unlikely but you never know
  filter(Learning != "Non-applicable") %>%
  # Change the column name because our function expects it
  rename(importance_level = Learning)

more_than_mod <- rbind(
  more_than_mod,
  list(
    "Non-research Staff",
    "Learning",
    get_percent_more_than_moderate(nrstaff_learning)
  )
)

nrstaff_pd <- nrstaff_learn_pd %>%
  select(`Professional Development`, job_category) %>%
  #unlikely but you never know
  filter(`Professional Development` != "Non-applicable") %>%
  # Change the column name because our function expects it
  rename(importance_level = `Professional Development`)

more_than_mod <- rbind(
  more_than_mod,
  list(
    "Non-research Staff",
    "Professional Development",
    get_percent_more_than_moderate(nrstaff_pd)
  )
)

```

Finally, let's get students.

```

students_learn_pd <- importance_and_job %>%
  filter(job_category == "Undergraduate" | job_category == "Grad Student") %>%
  select(Learning, `Professional Development`, job_category)

```



```

students_learn_pd$job_category <- "Student"

students_learning <- students_learn_pd %>%
  select(Learning, job_category) %>%
  #unlikely but you never know
  filter(Learning != "Non-applicable") %>%
  # Change the column name because our function expects it
  rename(importance_level = Learning)

more_than_mod <- rbind(
  more_than_mod,
  list(
    "Students",
    "Learning",
    get_percent_more_than_moderate(students_learning)
  )
)

students_pd <- students_learn_pd %>%
  select(`Professional Development`, job_category) %>%
  #unlikely but you never know
  filter(`Professional Development` != "Non-applicable") %>%
  # Change the column name because our function expects it
  rename(importance_level = `Professional Development`)

more_than_mod <- rbind(
  more_than_mod,
  list(
    "Students",
    "Professional Development",
    get_percent_more_than_moderate(students_pd)
  )
)

```

FINALLY, let's plot it!

```

more_than_mod$job_category <- factor(
  more_than_mod$job_category,
  levels = c(
    "Teachers",
    "Researchers",
    "Non-research Staff",

```

```

    "Students"
  )
)

more_than_mod$importance_area <- factor(
  more_than_mod$importance_area,
  levels = c(
    "Learning",
    "Professional Development",
    "Teaching",
    "Research",
    "Job"
  )
)

```

I'm not using my `grouped_bar_chart` function in `scripts/utils.R` because I have pre-computed the bar heights, and that function counts rows. Since I'm currently only creating this kind of bar chart once, I'm not bothering to create a new function (or incorporate this option into the `grouped_bar_chart` function).

```

p <- ggplot(
  more_than_mod,
  aes(
    x = job_category,
    y = pct,
    fill = importance_area
  )
) +
  geom_col(position = position_dodge(width = 0.8), width = 0.7) +
  ggtitle("Perceived Importance of Open Source\nfor Different Kinds of Work") +
  labs(
    y = "Percent of Respondents Who Said OS\nIs More than Moderately Important"
  ) +
  ylim(0, 100) +
  theme(
    axis.title.x = element_blank(),
    axis.title.y = element_text(size = 14),
    axis.text.x = element_text(angle = 60, vjust = 0.6, size = 12),
    axis.ticks.x = element_blank(),
    legend.title = element_blank(),
    legend.text = element_text(size = 12),
    panel.background = element_blank(),
  )

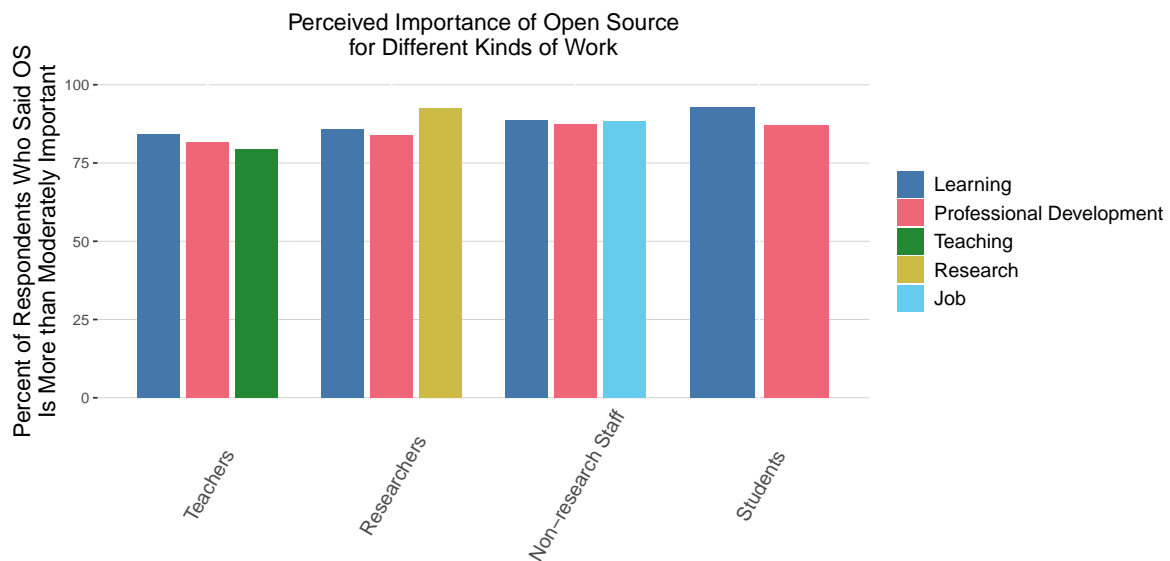
```

```

plot.title = element_text(hjust = 0.5, size = 14),
plot.margin = unit(c(0.3, 0.3, 0.3, 0.3), "cm"),
panel.grid.major.y = element_line(color = "lightgray", linewidth = 0.2)
) +
#https://sronpersonalpages.nl/~pault/
scale_fill_manual(
  values = c(
    '#4477AA',
    '#EE6677',
    '#228833',
    '#CCBB44',
    '#66CCEE'
  )
)
)

```

p



```

save_plot("importance_all_pct.tiff", 10, 5)

```

Session Info

```
sessionInfo()
```

```
R version 4.4.2 (2024-10-31)
Platform: aarch64-apple-darwin20
Running under: macOS Sequoia 15.4.1
```

```
Matrix products: default
```

```
BLAS: /Library/Frameworks/R.framework/Versions/4.4-arm64/Resources/lib/libRblas.0.dylib
```

```
LAPACK: /Library/Frameworks/R.framework/Versions/4.4-arm64/Resources/lib/libRlapack.dylib;
```

```
locale:
```

```
[1] en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8
```

```
time zone: America/Los_Angeles
```

```
tzcode source: internal
```

```
attached base packages:
```

```
[1] tools      stats      graphics  grDevices datasets  utils      methods
[8] base
```

```
other attached packages:
```

```
[1] treemap_2.4-4      tidyr_1.3.1        stringr_1.5.1
[4] scales_1.4.0       readr_2.1.5        pwr_1.3-0
[7] patchwork_1.3.0    mvabund_4.2.1      languageserver_0.3.16
[10] here_1.0.1         gtools_3.9.5       fpc_2.2-13
[13] forcats_1.0.0      factoextra_1.0.7   ggplot2_3.5.2
[16] dplyr_1.1.4        corrplot_0.95      cluster_2.1.8.1
```

```
loaded via a namespace (and not attached):
```

```
[1] gtable_0.3.6      xfun_0.52          ggrepel_0.9.6
[4] processx_3.8.6    lattice_0.22-6     callr_3.7.6
[7] tzdb_0.5.0        vctrs_0.6.5        ps_1.9.1
[10] generics_0.1.4    stats4_4.4.2       parallel_4.4.2
[13] flexmix_2.3-20    tibble_3.2.1       DEoptimR_1.1-3-1
[16] pkgconfig_2.0.3   data.table_1.17.6  RColorBrewer_1.1-3
[19] lifecycle_1.0.4   compiler_4.4.2     farver_2.1.2
[22] statmod_1.5.0     httpuv_1.6.16      htmltools_0.5.8.1
[25] class_7.3-22      yaml_2.3.10        later_1.4.2
[28] pillar_1.10.2     prabclus_2.3-4     MASS_7.3-61
[31] diptest_0.77-1    mclust_6.1.1       mime_0.13
[34] robustbase_0.99-4-1 tidyselect_1.2.1   digest_0.6.37
```

[37] stringi_1.8.7	purrr_1.0.4	kernlab_0.9-33
[40] labeling_0.4.3	rprojroot_2.0.4	fastmap_1.2.0
[43] grid_4.4.2	colorspace_2.1-1	cli_3.6.5
[46] magrittr_2.0.3	utf8_1.2.5	withr_3.0.2
[49] promises_1.3.3	tweedie_2.3.5	rmarkdown_2.29
[52] igraph_2.1.4	nnet_7.3-19	modeltools_0.2-24
[55] hms_1.1.3	shiny_1.11.0	evaluate_1.0.3
[58] knitr_1.50	rlang_1.1.6	Rcpp_1.0.14
[61] xtable_1.8-4	gridBase_0.4-7	glue_1.8.0
[64] xml2_1.3.8	renv_1.1.4	jsonlite_2.0.0
[67] R6_2.6.1		