

Pilot Analysis: Exploratory Visualization

This document is intended to give an overview of the response distributions from our pilot.

Data

Load Worker Responses from Pilot

The data is already anonymous and in a tidy format at this stage in the analysis pipeline. We just need to read it in and do some preprocessing.

```
# read in data
full_df <- read_csv("pilot-anonymous.csv")
```

```
## Parsed with column specification:
## cols(
##   .default = col_double(),
##   workerId = col_character(),
##   batch = col_integer(),
##   condition = col_character(),
##   start_means = col_character(),
##   numeracy = col_integer(),
##   gender = col_character(),
##   age = col_character(),
##   education = col_character(),
##   chart_use = col_character(),
##   intervene = col_integer(),
##   outcome = col_character(),
##   pSup = col_integer(),
##   trial = col_character(),
##   trialIdx = col_character()
## )
```

```
## See spec(...) for full column specifications.
```

```

# preprocessing
responses_df <- full_df %>%
  rename( # rename to convert away from camel case
    worker_id = workerId,
    account_value = accountValue,
    ground_truth = groundTruth,
    p_award_with = pAwardWith,
    p_award_without = pAwardWithout,
    p_superiority = pSup,
    start_time = startTime,
    resp_time = respTime,
    trial_dur = trialDur,
    trial_idx = trialIdx
  ) %>%
  filter(trial_idx != "practice", trial_idx != "mock") %>% # remove practice and mock trials from responses dataframe, leave in full version
  mutate( # mutate to jitter probability of superiority away from boundaries
    p_superiority = ifelse(p_superiority == 0, 0.5, p_superiority),           # avoid responses equal to zero
    p_superiority = ifelse(p_superiority == 100, 99.5, p_superiority)          # avoid responses equal to one-hundred
  ) %>%
  mutate( # mutate to rows where intervene == -1 for some reason
    intervene = if_else(intervene == -1,
      # repair
      if_else((payoff == (award_value - 1) | payoff == -1),
        1, # payed for intervention
        0), # didn't pay for intervention
      # don't repair
      as.numeric(intervene) # hack to avoid type error
    )
  ) %>%
  # add a variable to note whether the chart they viewed showed means
  mutate(means = as.factor((start_means == "True" & as.numeric(trial) < 16) | (start_means == "False" & as.numeric(trial) >= 16)))

head(responses_df)

```

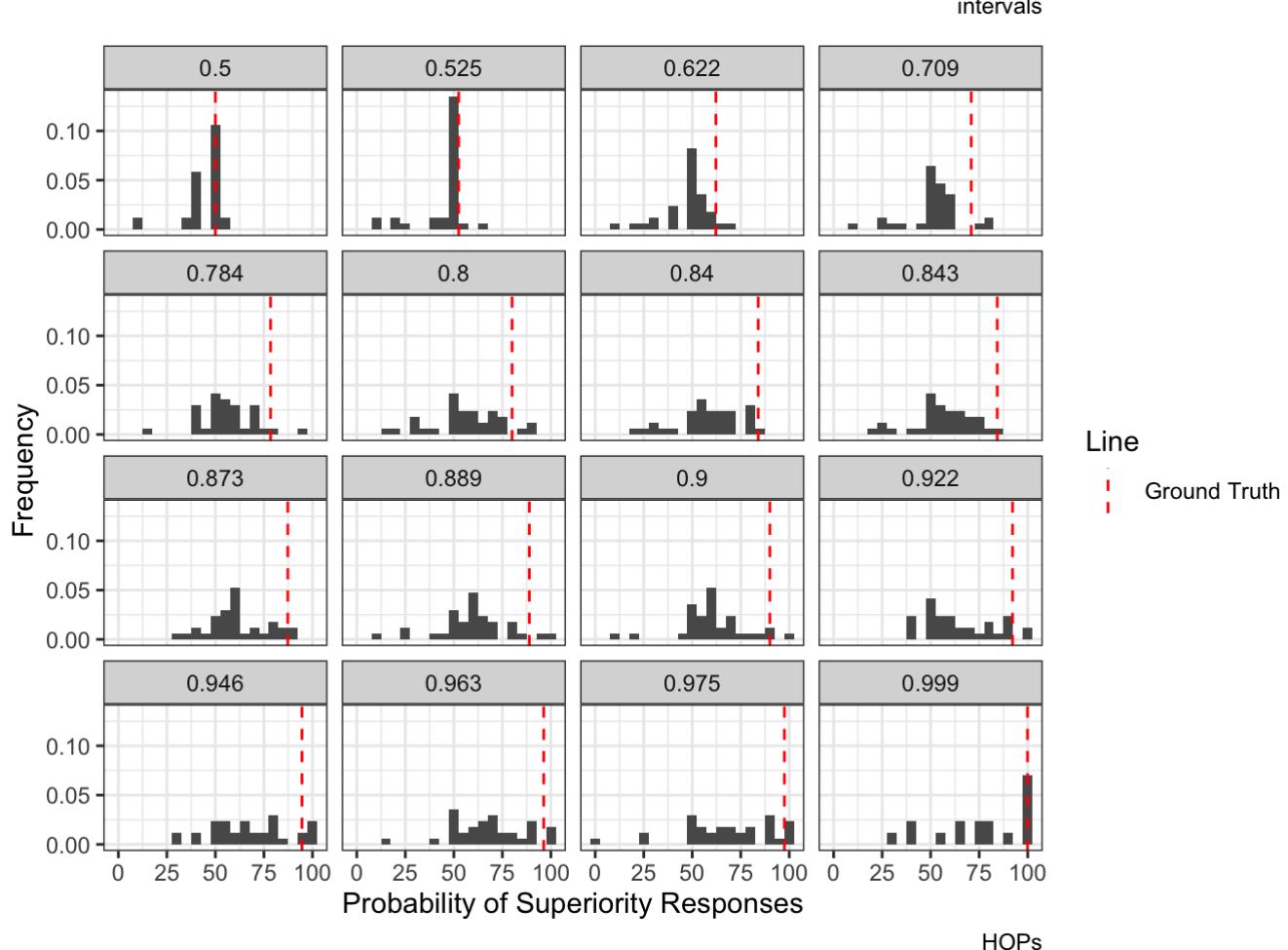
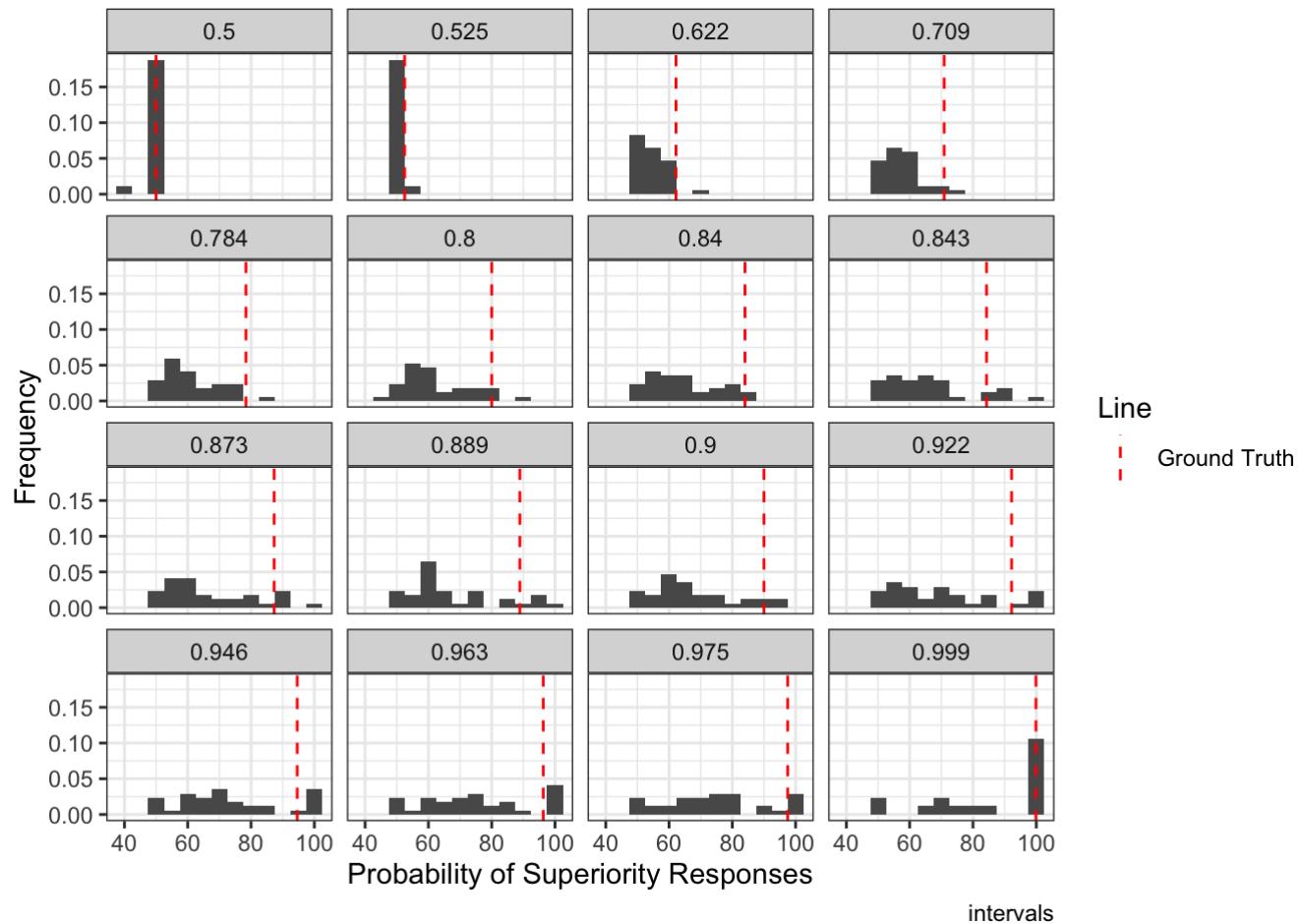
```
## # A tibble: 6 x 30
##   worker_id batch condition baseline award_value exchange start_means
##   <chr>      <int> <chr>        <dbl>      <dbl>    <dbl> <chr>
## 1 c3de4118     4 intervals     0.5       2.25     0.2 False
## 2 c3de4118     4 intervals     0.5       2.25     0.2 False
## 3 c3de4118     4 intervals     0.5       2.25     0.2 False
## 4 c3de4118     4 intervals     0.5       2.25     0.2 False
## 5 c3de4118     4 intervals     0.5       2.25     0.2 False
## 6 c3de4118     4 intervals     0.5       2.25     0.2 False
## # ... with 23 more variables: total_bonus <dbl>, duration <dbl>,
## #   numeracy <int>, gender <chr>, age <chr>, education <chr>,
## #   chart_use <chr>, account_value <dbl>, ground_truth <dbl>,
## #   intervene <dbl>, outcome <chr>, pAwardCurrent <dbl>, pAwardNew <dbl>,
## #   p_award_with <dbl>, p_award_without <dbl>, p_superiority <dbl>,
## #   payoff <dbl>, resp_time <dbl>, start_time <dbl>, trial <chr>,
## #   trial_dur <dbl>, trial_idx <chr>, means <fct>
```

Response Distributions

Probability of Superiority Judgments

Let's plot histograms of probability of superiority judgments at each level of the ground truth probability of superiority. We show the ground truth in red. This will give us an overview of bias and precision in judgments. We do this separately for each visualization condition to limit the number of faceted subplots in a single view.

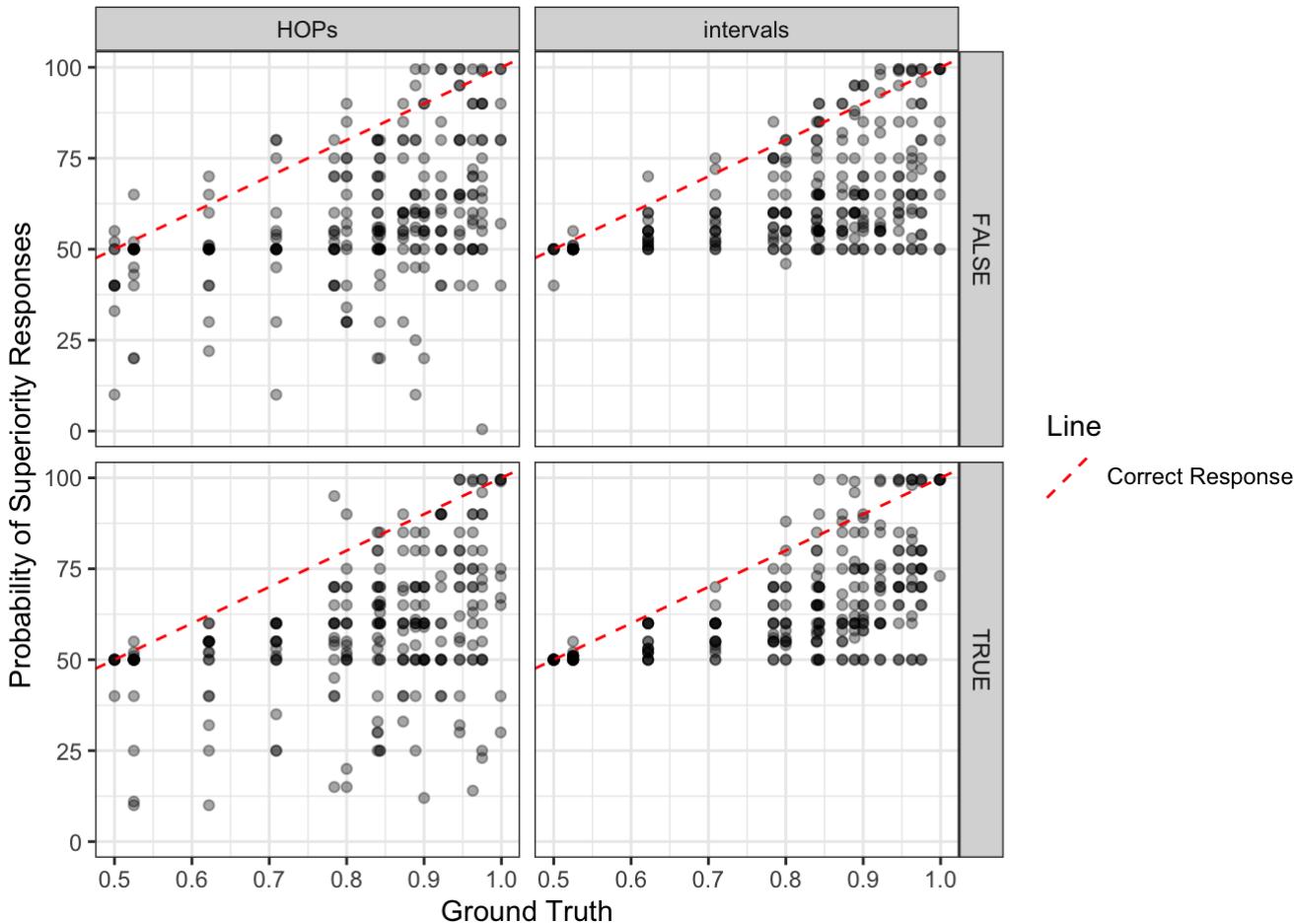
```
for (cond in unique(responses_df$condition)) {
  plt <- responses_df %>% filter(condition == cond) %>%
    ggplot(aes(x = p_superiority)) +
    geom_histogram(aes(y = ..density..), binwidth = 5) +
    geom_vline(aes(xintercept = ground_truth * 100, linetype = "Ground Truth"), color =
    "red") +
    scale_linetype_manual(name = "Line", values = c(2,1), guide=guide_legend(override.ae
    s = list(color = c("red")))) +
    theme_bw() +
    labs(
      caption=cond,
      x = "Probability of Superiority Responses",
      y = "Frequency"
    ) +
    facet_wrap(~ ground_truth)
  print(plt)
}
```



As we would expect based on a linear log odds representation of probability, probability of superiority judgments tend to be biased toward 50% relative to the ground truth.

Another more compact way of looking at the relationship between estimated probability of superiority and the ground truth is to just plot them against one another. Let's look at this even though its sort of a mess.

```
# plot estimated probability of superiority vs the ground truth
responses_df %>%
  ggplot(aes(x = ground_truth, y = p_superiority)) +
  geom_point(alpha = 0.35) +
  geom_abline(aes(intercept = 0, slope = 100, linetype = "Correct Response"), color = "red") +
  scale_linetype_manual(name = "Line", values = c(2,1), guide=guide_legend	override.aes
= list(color = c("red")))) +
  theme_bw() +
  labs(
    x = "Ground Truth",
    y = "Probability of Superiority Responses"
  ) +
  facet_grid(means ~ condition)
```



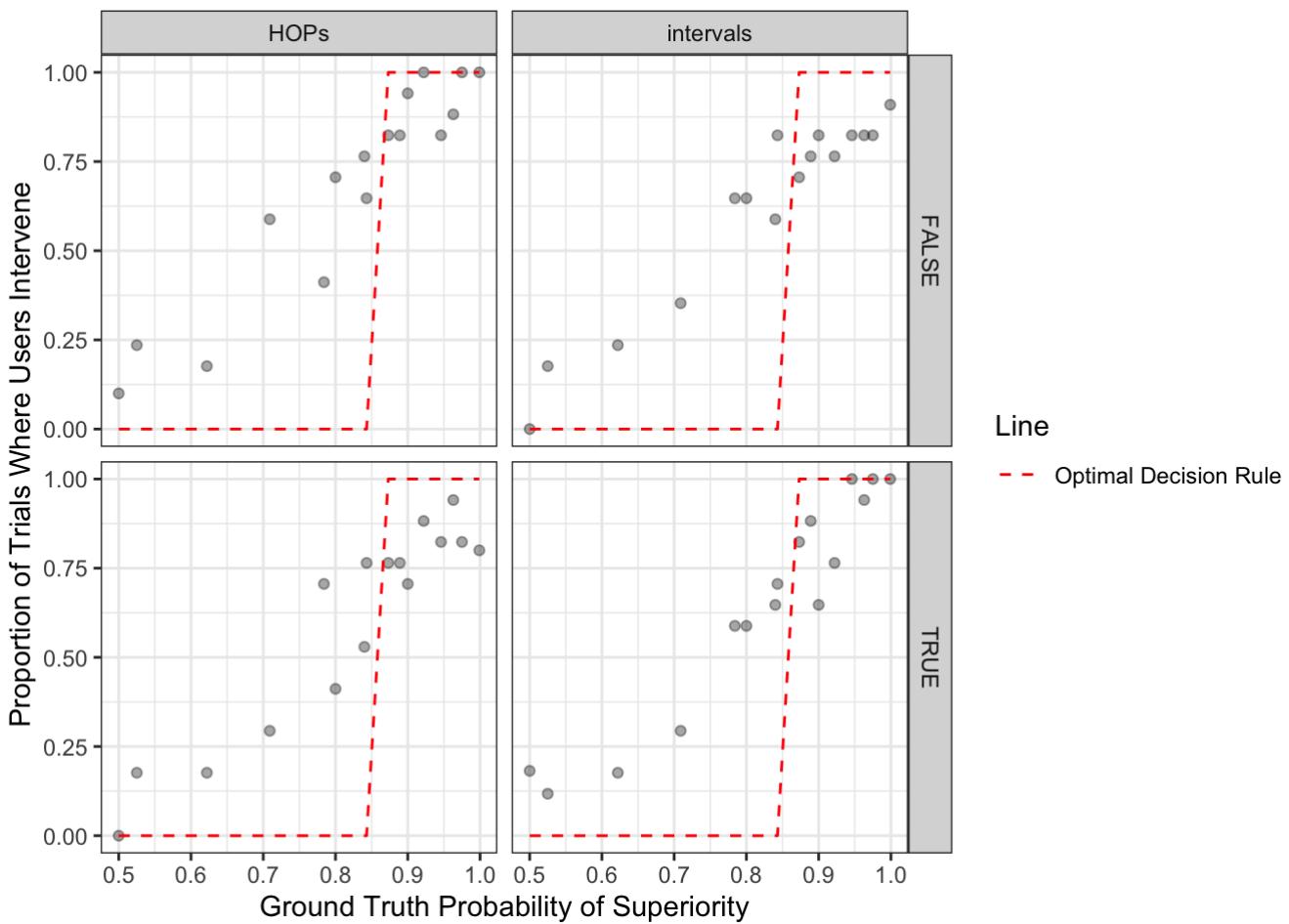
Decisions to Intervene

In order to see how people are doing on the decision task, we want to benchmark their performance against a utility optimal decision rule. The rule is different depending on whether the task is framed as a gain or a loss (i.e., whether the ground truth probability of superiority is greater than or less than 50%).

```
# determine whether or not intervention is utility optimal on each trial
responses_df <- responses_df %>%
  mutate(should_intervene = if_else(ground_truth > 0.5,
                                    (p_award_with - p_award_without) > 1 / award_value,
# gain framing decision rule
                                    ((1 - p_award_without) - (1 - p_award_with)) > 1 / award_value) # loss framing decision rule
  )
```

Let's plot the proportion of users who intervene at each level of ground truth probability of superiority in each visualization condition. People should intervene more often at extreme probabilities. We show the utility optimal decision threshold in red. This should give us an overview of decision quality.

```
# summarise the data as the overall proportion of trials where users intervene vs what they should do at each level of ground_truth * condition * baseline
responses_df %>%
  group_by(means, condition, ground_truth) %>%
  summarise(
    proportion_intervene = sum(intervene) / n(),
    optimal_decision = mean(should_intervene)
  ) %>%
  ggplot(aes(x = ground_truth, y = proportion_intervene)) +
  geom_point(alpha = 0.35) +
  geom_line(aes(y = optimal_decision, linetype="Optimal Decision Rule"), color="red") +
  scale_linetype_manual(name="Line", values = c(2,1), guide=guide_legend(override.aes=list(color=c("red")))) +
  theme_bw() +
  labs(
    x = "Ground Truth Probability of Superiority",
    y = "Proportion of Trials Where Users Intervene"
  ) +
  facet_grid(means ~ condition)
```

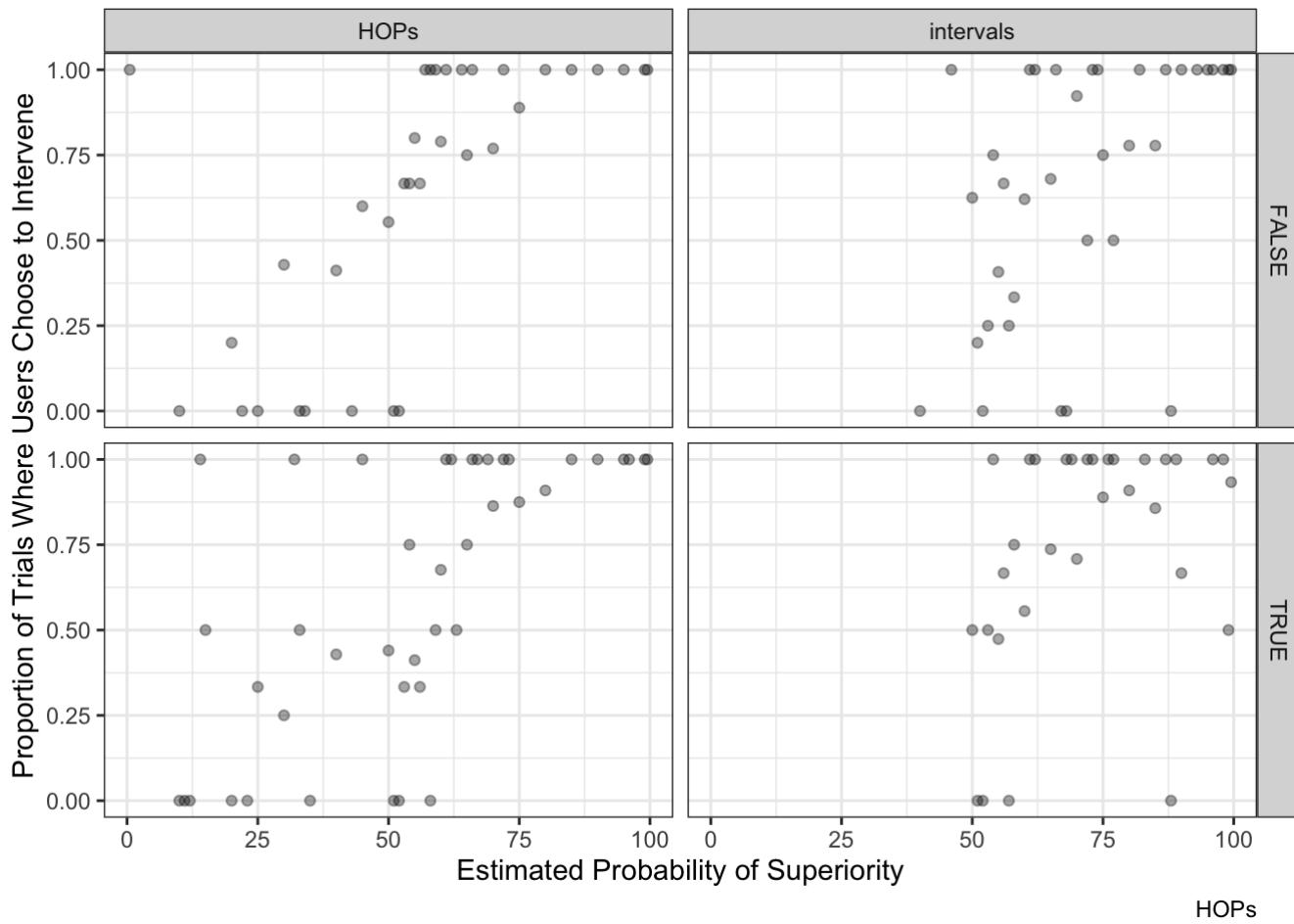


In the aggregate, differences between conditions are pretty subtle. It looks like there may be a slight discrepancy in performance between visualization conditions. We'll need to tease these effects out using statistical inference.

Probability of Superiority Judgments vs Decisions to Intervene

It might also be interesting to see how decisions correspond to probability of superiority judgments. We omit the ground truth and optimal decision rule from this chart.

```
# summarise the data as the overall proportion of trials where users choose to intervene
# at each level of condition * baseline * p_superiority
responses_df %>%
  group_by(means, condition, p_superiority) %>%
  summarise(proportion_intervene = sum(intervene) / n()) %>%
  ggplot(aes(x = p_superiority, y = proportion_intervene)) +
  geom_point(alpha = 0.35) +
  theme_bw() +
  labs(
    caption=cond,
    x = "Estimated Probability of Superiority",
    y = "Proportion of Trials Where Users Choose to Intervene"
  ) +
  facet_grid(means ~ condition)
```



People's probability of superiority judgments and decisions are correlated in the way that you would expect if they understood the task.

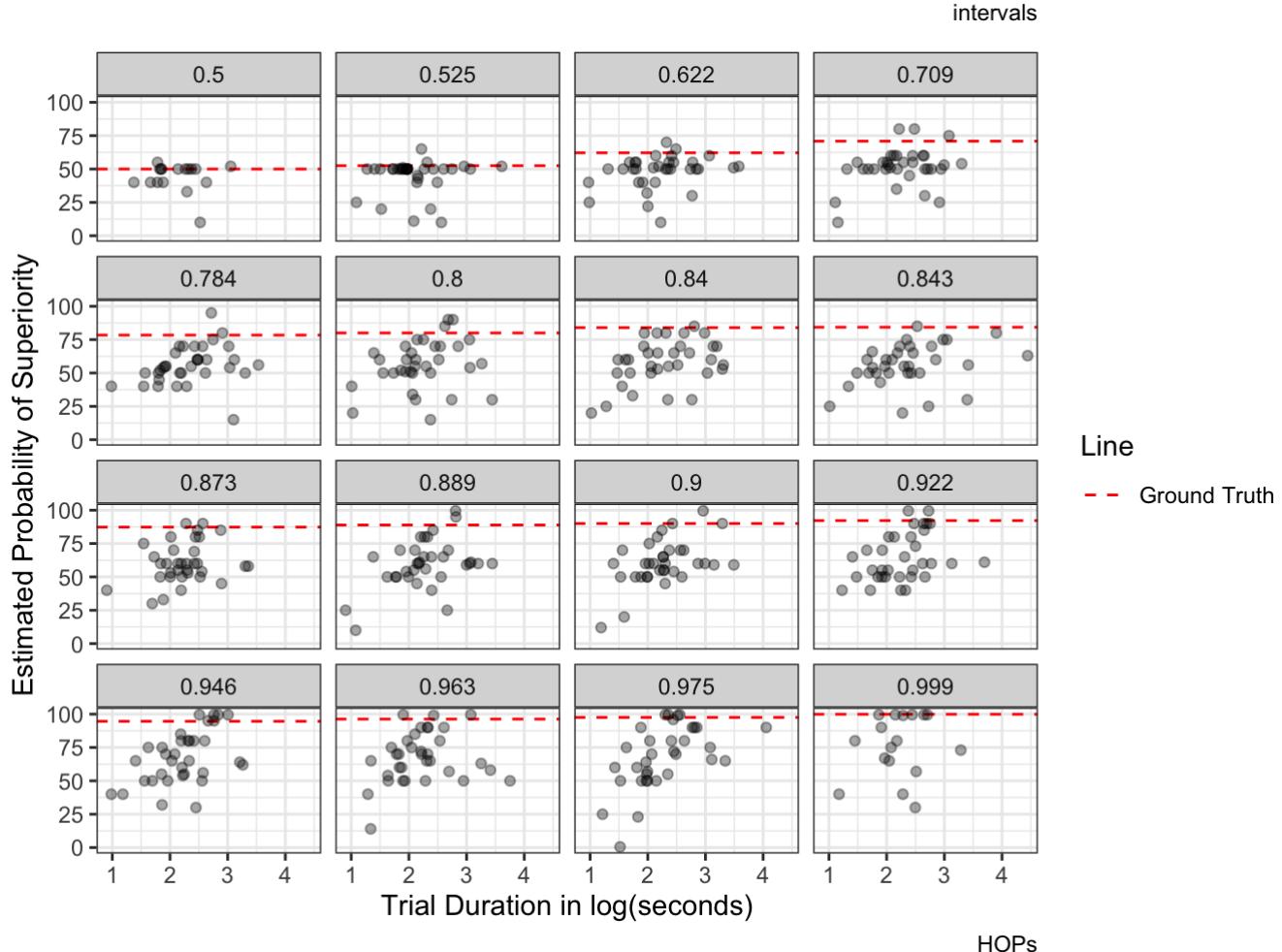
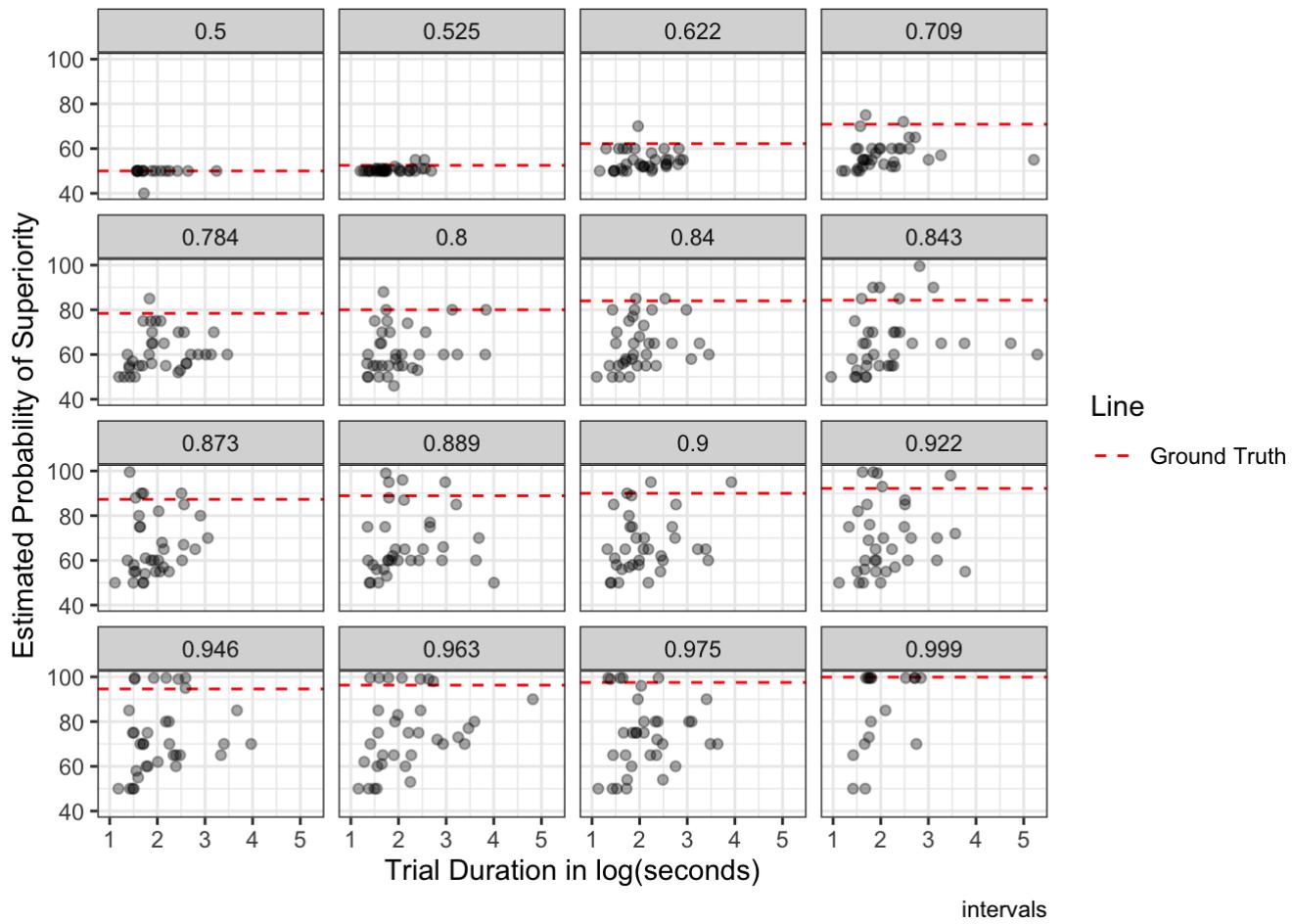
Relationships with Trial Duration

We want to know when, if at all, spending more time on a response results in improved performance.

Trial Duration vs Probability of Superiority Judgments

Let's look at probability of superiority estimates as a function of trial duration. As before, we show the ground truth in red and separate visualization conditions into different views to limit the number of faceted subplots in a single view.

```
for (cond in unique(responses_df$condition)) {  
  plt <- responses_df %>% filter(condition == cond) %>%  
    ggplot(aes(x = log(trial_dur), y = p_superiority)) +  
    geom_hline(aes(yintercept = ground_truth * 100, linetype = "Ground Truth"), color =  
    "red") +  
    scale_linetype_manual(name = "Line", values = c(2,1), guide=guide_legend(override.ae  
s = list(color = c("red")))) +  
    geom_point(alpha = 0.35) +  
    theme_bw() +  
    labs(  
      caption=cond,  
      x = "Trial Duration in log(seconds)",  
      y = "Estimated Probability of Superiority"  
    ) +  
    facet_wrap(~ ground_truth)  
  print(plt)  
}
```



HOPs

Probability of superiority judgments seem somewhat more accurate at longer response times, particularly in the HOPs condition.

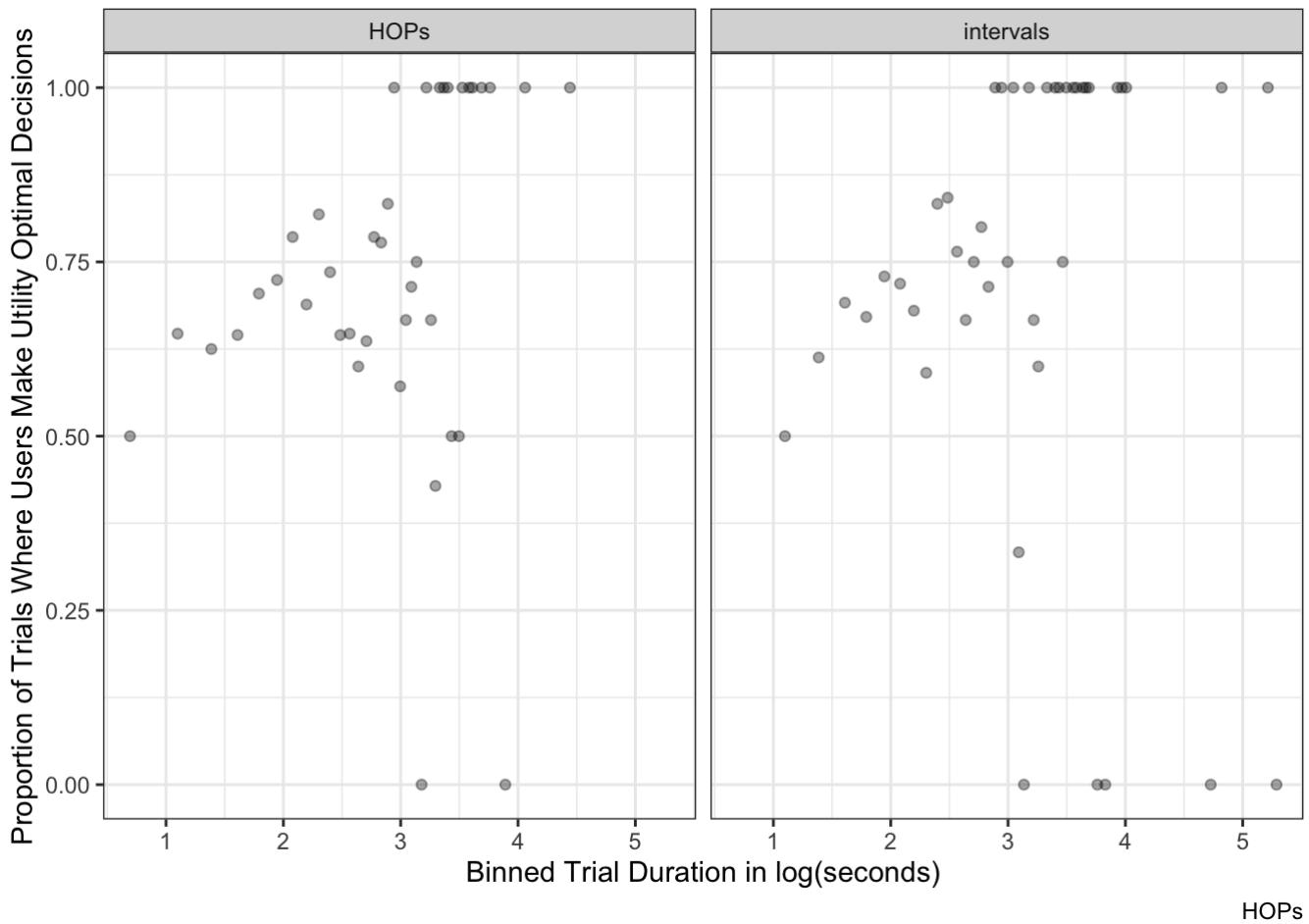
Trial Duration vs Decision Quality

A nice metric for decision quality is whether users responded “correctly” or in line with the normative utility optimal decision rule. We calculate whether the user was “correct” or not on each trial.

```
# determine whether response on each trial is utility optimal
responses_df <- responses_df %>%
  mutate(correct = intervene == should_intervene)
```

Let's look at the proportion correct as a function of trial duration, faceting visualization conditions as above.

```
# summarise the data as the overall proportion of trials where users make utility optimal decisions at each level of condition * baseline * trial_duration
responses_df %>%
  mutate(trial_dur_binned = round(trial_dur)) %>%
  group_by(condition, baseline, trial_dur_binned) %>%
  summarise(proportion_correct = sum(correct) / n()) %>%
  ggplot(aes(x = log(trial_dur_binned), y = proportion_correct)) +
  geom_point(alpha = 0.35) +
  theme_bw() +
  labs(
    caption=cond,
    x = "Binned Trial Duration in log(seconds)",
    y = "Proportion of Trials Where Users Make Utility Optimal Decisions"
  ) +
  facet_grid(. ~ condition)
```



Trial duration seems to have little to do with decision quality.

Error Analysis

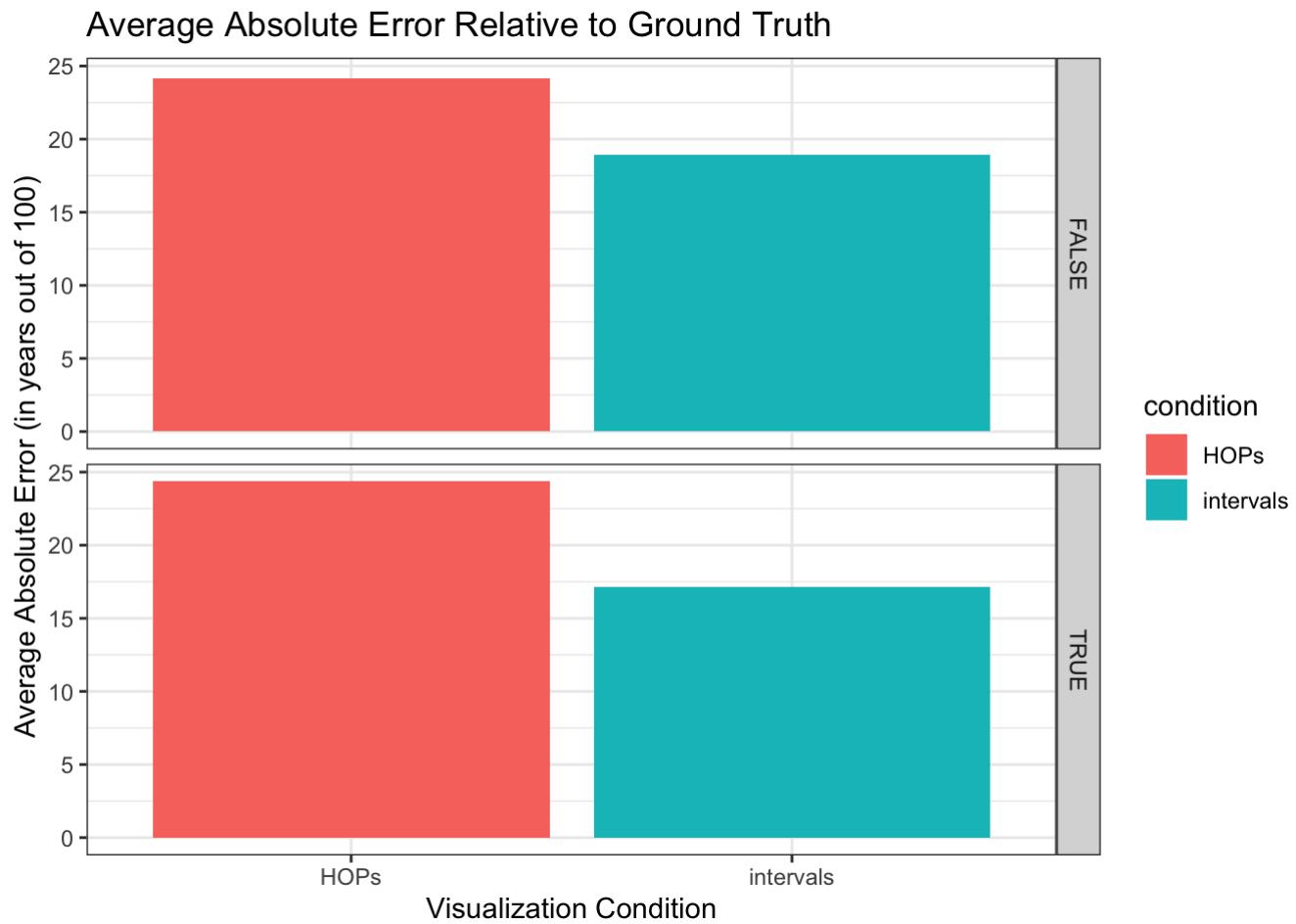
In this section, we look for patterns of interest in response errors. We'll start by adding error and absolute error in probability of superiority judgments to the dataframe. We already have a metric for correctness of decisions.

```
# add error and absolute error to df
responses_df <- responses_df %>%
  mutate(
    err_p_sup = ground_truth * 100 - p_superiority,
    abs_err_p_sup = abs(err_p_sup)
  )
```

Mean Absolute Error

Let's look at the average absolute error in probability of superiority judgments in each condition, regardless of the ground truth.

```
# avg absolute error per condition
responses_df %>%
  group_by(means, condition) %>%
  summarise(avg_abs_err_p_sup = mean(abs_err_p_sup)) %>%
  ggplot(aes(x = condition, y = avg_abs_err_p_sup, fill = condition)) +
  geom_bar(stat = "identity") +
  theme_bw() +
  labs(title = "Average Absolute Error Relative to Ground Truth",
       x = "Visualization Condition",
       y = "Average Absolute Error (in years out of 100)"
     ) +
  facet_grid(means ~ .)
```

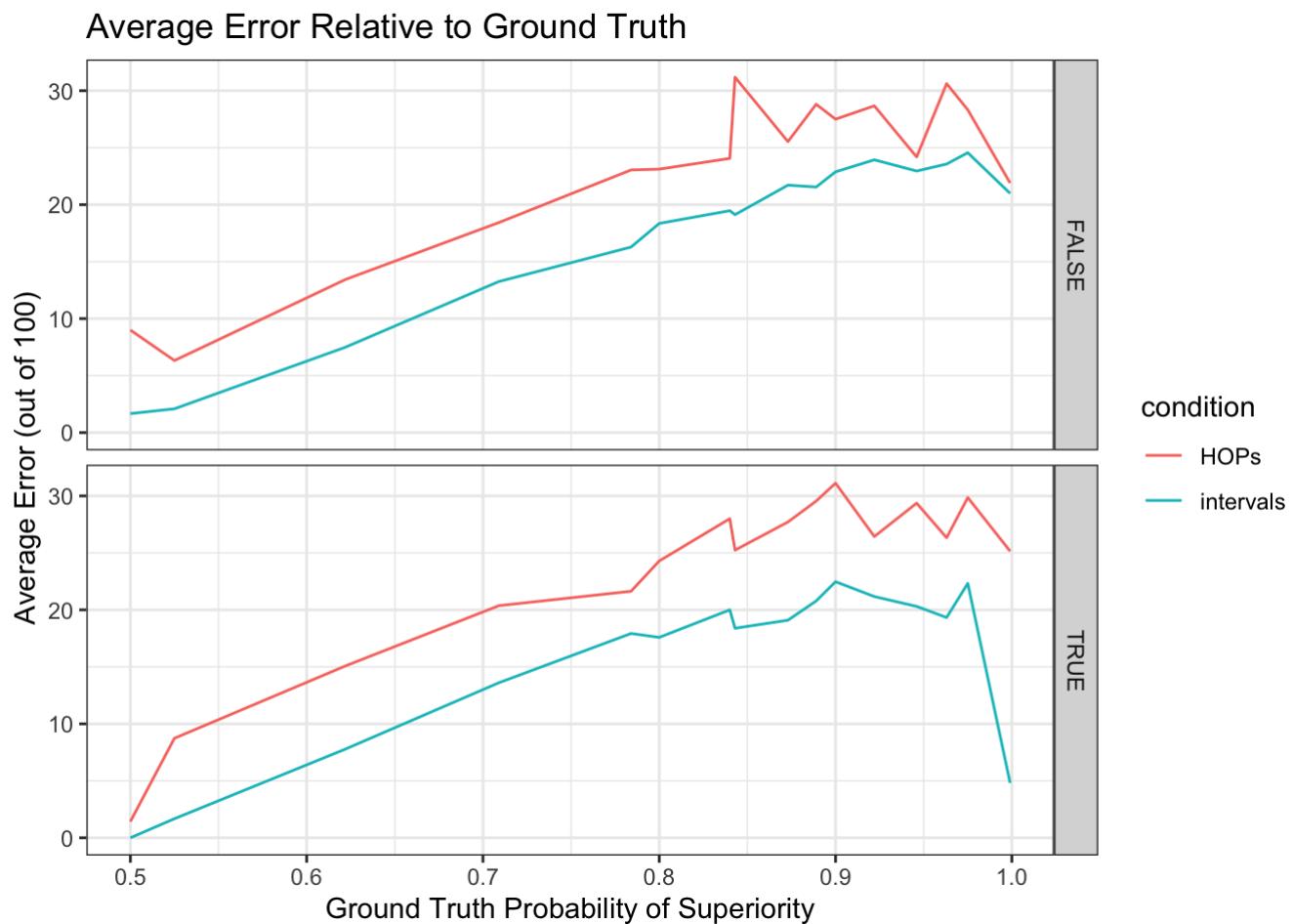


On average, errors in probability of superiority judgments are high across the board, with the average error equal to about a quarter the range of possible responses. Error rates are notably high in the HOPs condition.

Mean Error vs Ground Truth

Let's look at the average signed error in probability of superiority judgments. This time we'll plot error in each condition in relation to ground truth.

```
# error by ground truth, per condition
responses_df %>%
  group_by(ground_truth, means, condition) %>%
  summarise(avg_err_p_sup = mean(err_p_sup)) %>%
  ggplot(aes(x = ground_truth, y = avg_err_p_sup, color = condition)) +
  geom_line() +
  theme_bw() +
  labs(title = "Average Error Relative to Ground Truth",
       x = "Ground Truth Probability of Superiority",
       y = "Average Error (out of 100)")
) +
facet_grid(means ~ .)
```

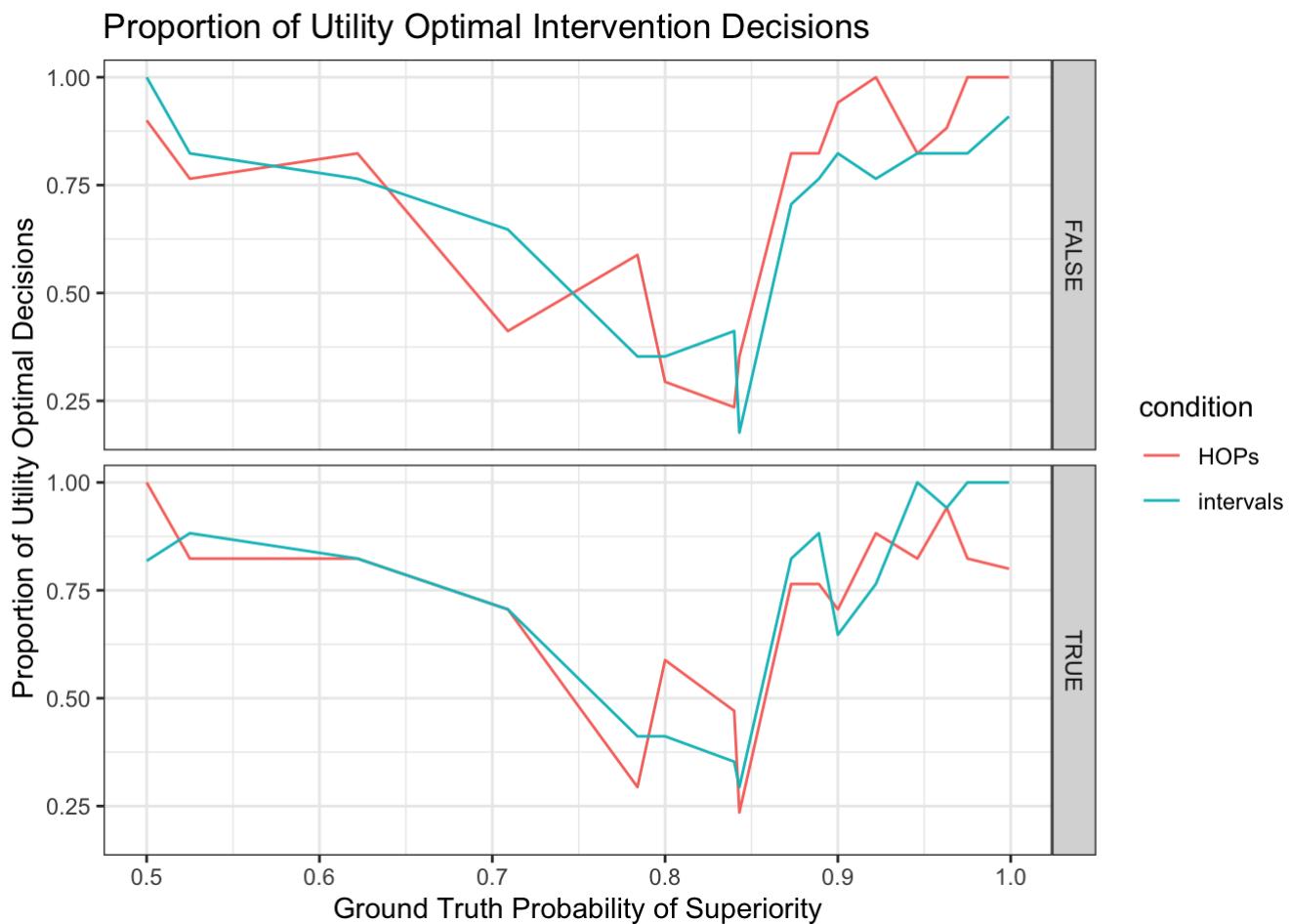


Again, we can see that errors are large on average, especially at the extreme end of the probability scale near 90%. Higher errors at the extremes of the probability scale are expected based on *the central tendency of judgments*. We can also see that people make judgments with consistently lower error when using intervals rather than HOPs.

Proportion of Utility Optimal Decisions

Let's take a similar approach to visualizing decisions by looking at the proportion of utility optimal decisions as a function of ground truth probability of superiority and condition.

```
# error by ground truth, per condition
responses_df %>%
  group_by(ground_truth, means, condition) %>%
  summarise(proportion_correct = sum(correct) / n()) %>%
  ggplot(aes(x = ground_truth, y = proportion_correct, color = condition)) +
  geom_line() +
  theme_bw() +
  labs(title = "Proportion of Utility Optimal Intervention Decisions",
       x = "Ground Truth Probability of Superiority",
       y = "Proportion of Utility Optimal Decisions"
     ) +
  facet_grid(means ~ .)
```



We can see that there are dips in performance near 83% probability of superiority. This is expected considering that this is at the decision threshold where the intervention decision is most ambiguous.

Individual Patterns of Behavior

As is often the case with judgments from visualizations, the data seem highly heterogenous. We try to get a sense of this by looking at individual patterns of responses in conjunction with individual characteristics such as gender, age, education, chart use, and numeracy.

We'll also use these individual views to develop some exclusion criteria. We use two attention checks, one in the middle of each block. If a user understands the decision problem, they should intervene when probability of superiority is 99.9% and should not intervene when probability of superiority is 50%. We also look for participants

who responded the same on every trial for either question since these folks were probably speeding. Last, we'll look for people who often respond that probability of superiority is less than 50% since this response reflects a misunderstanding of the response scale.

Below we create an overview of performance and individual characteristics for each participant separately.

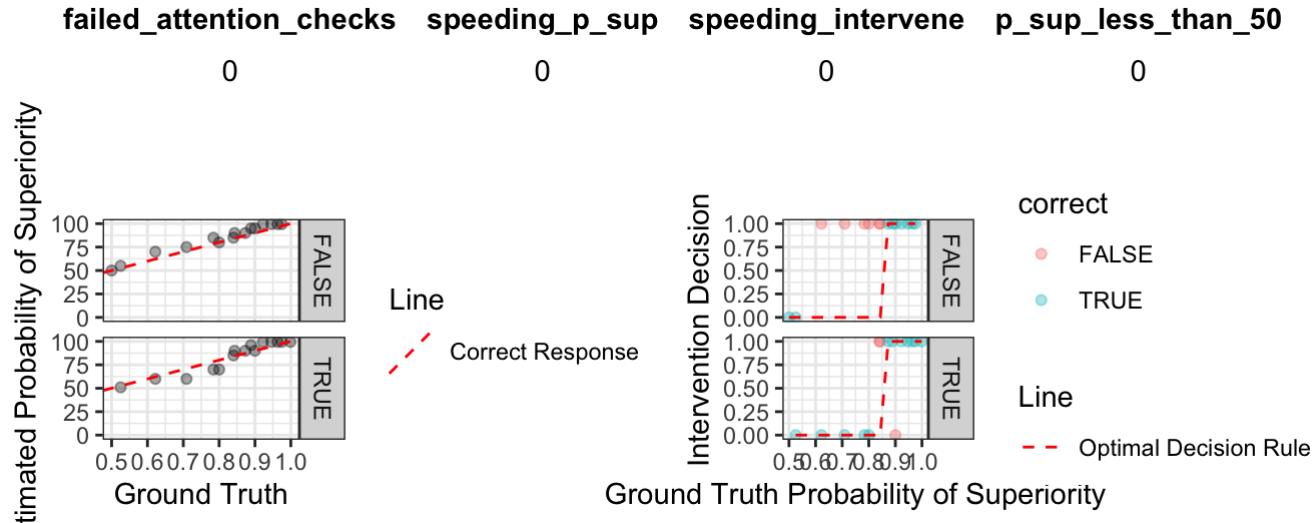
```

for (worker in unique(responses_df$worker_id)) {
  # get a df for just this worker
  worker_df <- responses_df %>% filter(worker_id == worker)
  # plot probability of superiority judgments vs ground truth
  p_sup_plt <- worker_df %>%
    ggplot(aes(x = ground_truth, y = p_superiority)) +
    geom_point(alpha = 0.35) +
    geom_abline(aes(intercept = 0, slope = 100, linetype = "Correct Response"), color =
"red") +
    scale_linetype_manual(name = "Line", values = c(2,1), guide=guide_legend	override.ae
s = list(color = c("red")))) +
    theme_bw() +
    ylim(0, 100) +
    labs(
      x = "Ground Truth",
      y = "Estimated Probability of Superiority"
    ) +
    facet_grid(means ~ .)
  # plot intervention decisions vs ground truth, noting which are in line with the utili
  ty optimal decision rule
  decision_plt <- worker_df %>%
    ggplot(aes(x = ground_truth, y = intervene, color = correct)) +
    geom_point(alpha = 0.35) +
    geom_line(aes(y = as.numeric(should_intervene), linetype="Optimal Decision Rule"), c
olor="red") +
    scale_linetype_manual(name="Line", values = c(2,1), guide=guide_legend	override.aes=
list(color=c("red")))) +
    theme_bw() +
    labs(
      x = "Ground Truth Probability of Superiority",
      y = "Intervention Decision"
    ) +
    facet_grid(means ~ .)
  # create a table summarizing this worker
  summary_table <- worker_df %>%
    group_by(worker_id) %>%
    summarise(
      condition = unique(condition),
      gender = unique(gender),
      age = unique(age),
      education = unique(education),
      chart_use = unique(chart_use),
      numeracy = unique(numeracy)
    ) %>%
    select(-worker_id) %>%
    ggttexttable(rows = NULL, theme = ttheme("blank"))
  # create table summarizing potential exclusion criteria
  exclusion_table <- worker_df %>%
    # attention check trials where ground truth = c(0.5, 0.999)
    mutate(failed_check = (ground_truth == 0.5 & intervene != 0) | (ground_truth == 0.99
9 & intervene != 1)) %>%
    group_by(worker_id) %>%
    summarise(

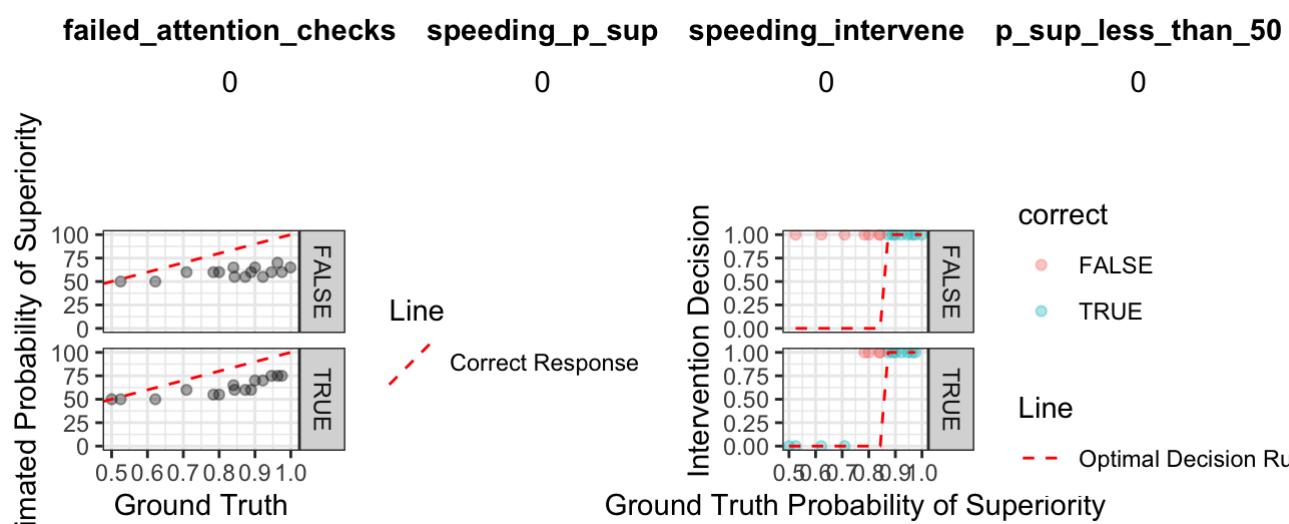
```

```
failed_attention_checks = sum(failed_check),
speeding_p_sup = as.logical(length(unique(p_superiority)) == 1),
speeding_intervene = as.logical(length(unique(intervene)) == 1),
p_sup_less_than_50 = sum(p_superiority < 50) / n()
) %>%
select(-worker_id) %>%
ggtexttable(rows = NULL, theme = ttheme("blank"))
# stitch together these three views
charts <- ggarrange(p_sup_plt, decision_plt, ncol = 2, nrow = 1)
figure <- ggarrange(summary_table, exclusion_table, charts, ncol = 1, nrow = 3)
print(figure)
}
```

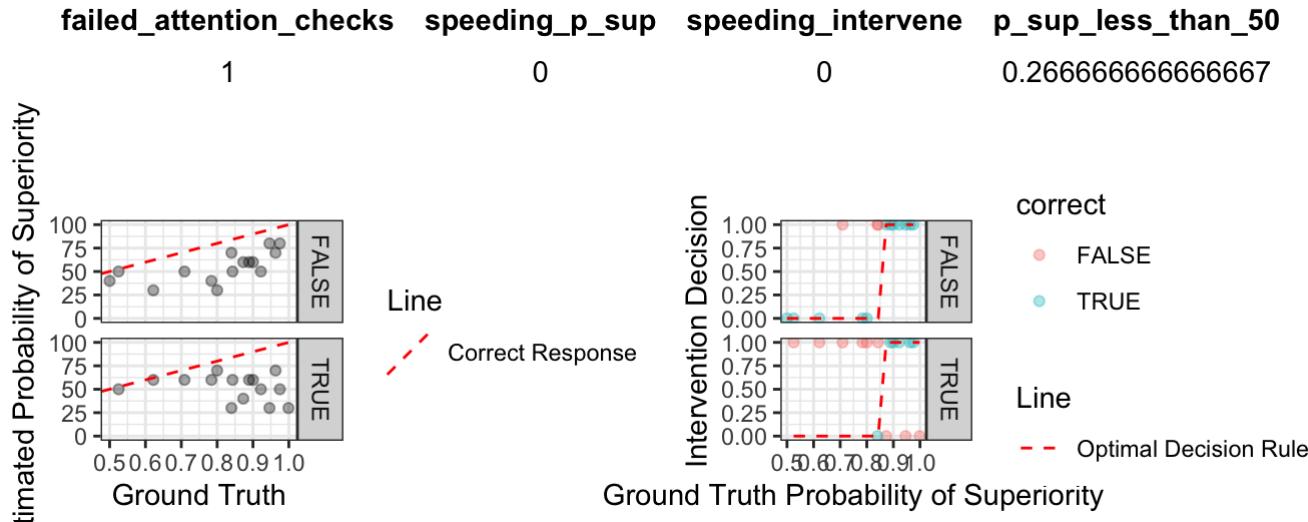
condition	gender	age	education	chart_use	numeracy
intervals	M	25-34	Some college, no degree	Daily	11



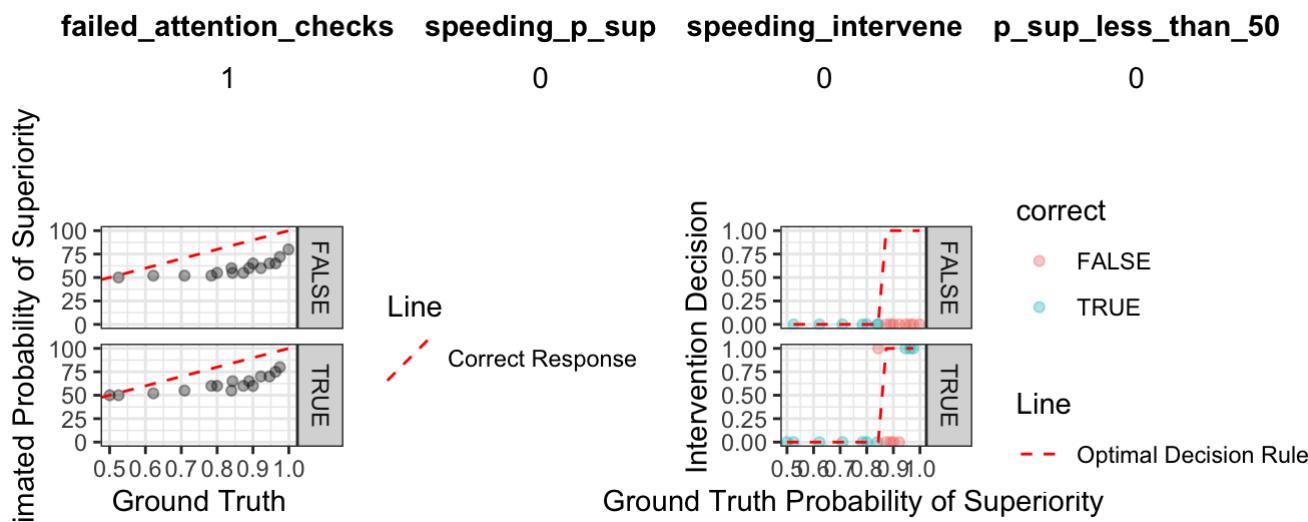
condition	gender	age	education	chart_use	numeracy
intervals	F	55-64	Bachelor's degree	Weekly	11



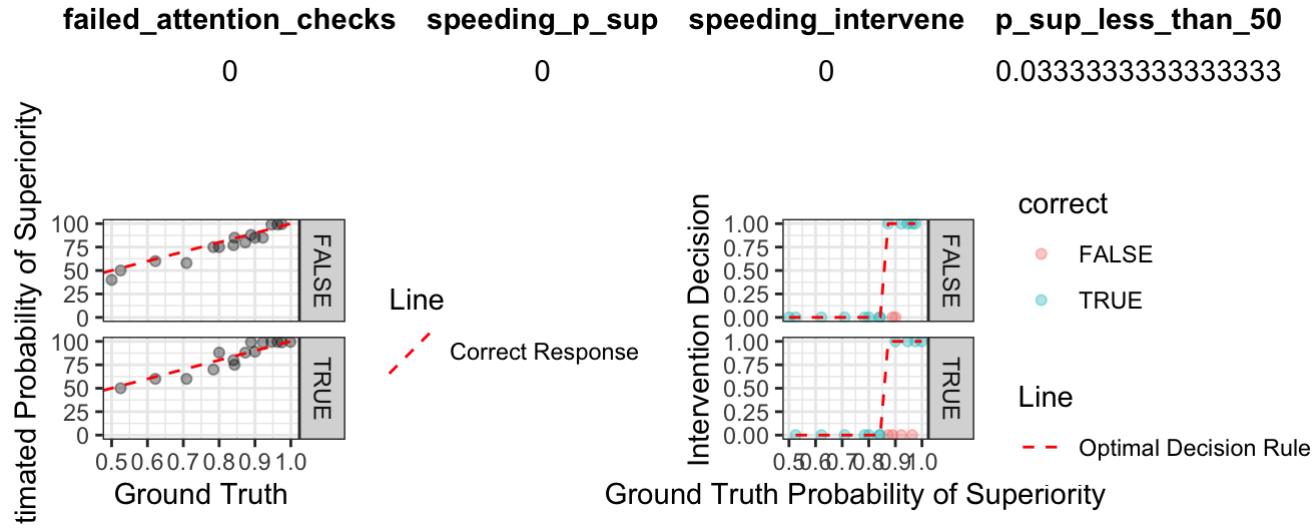
condition	gender	age	education	chart_use	numeracy
HOPs	F	35-44	Master's degree	Monthly or less	10



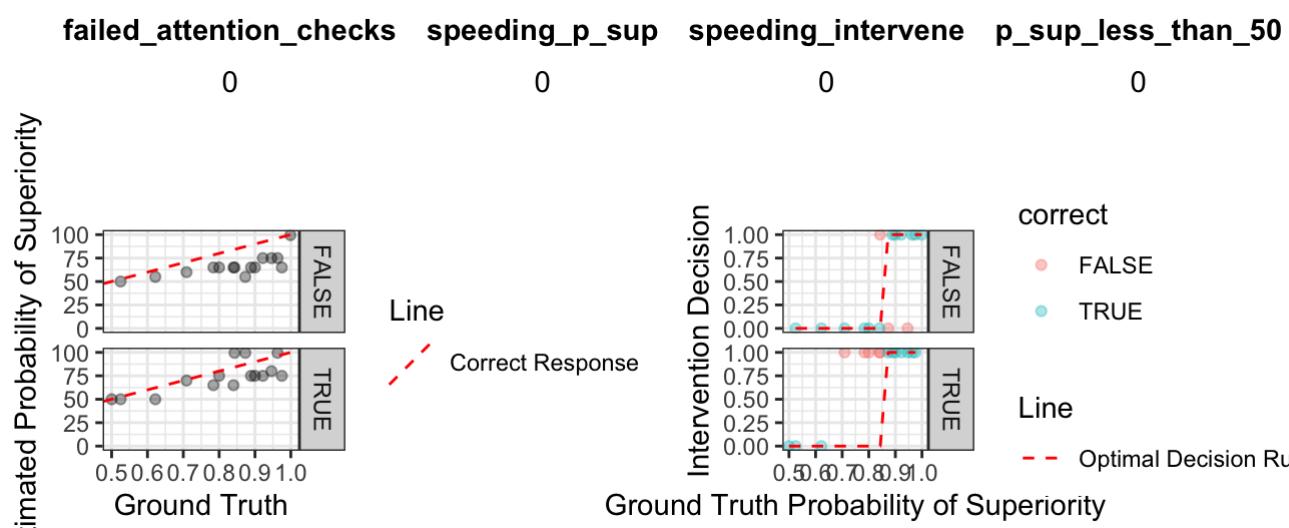
condition	gender	age	education	chart_use	numeracy
intervals	M	25-34	Bachelor's degree	Weekly	9



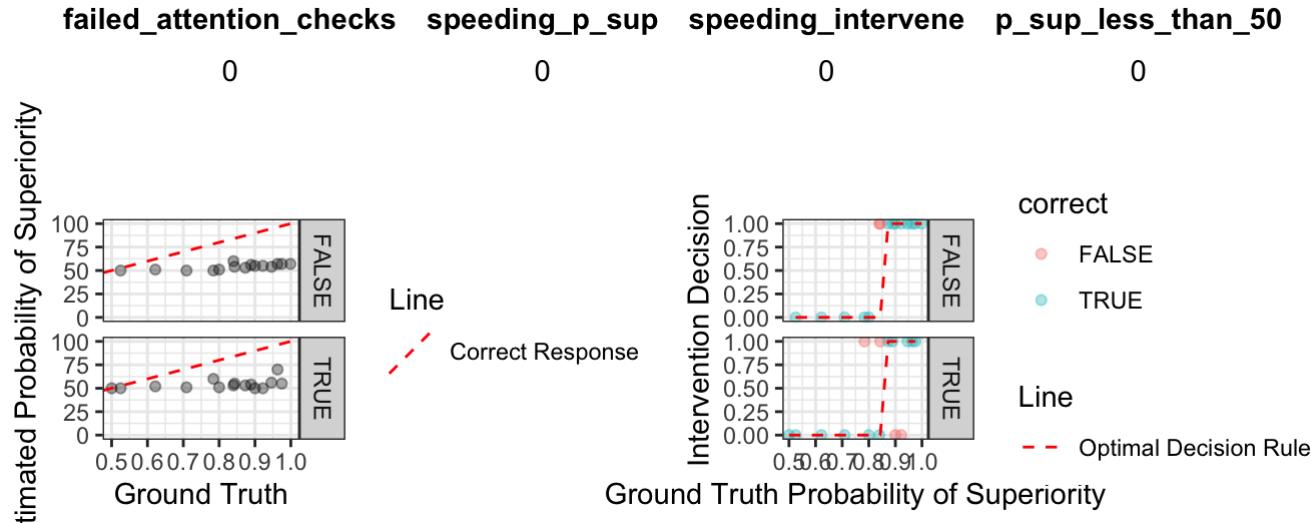
condition	gender	age	education	chart_use	numeracy
intervals	F	25-34	Some college, no degree	Weekly	10



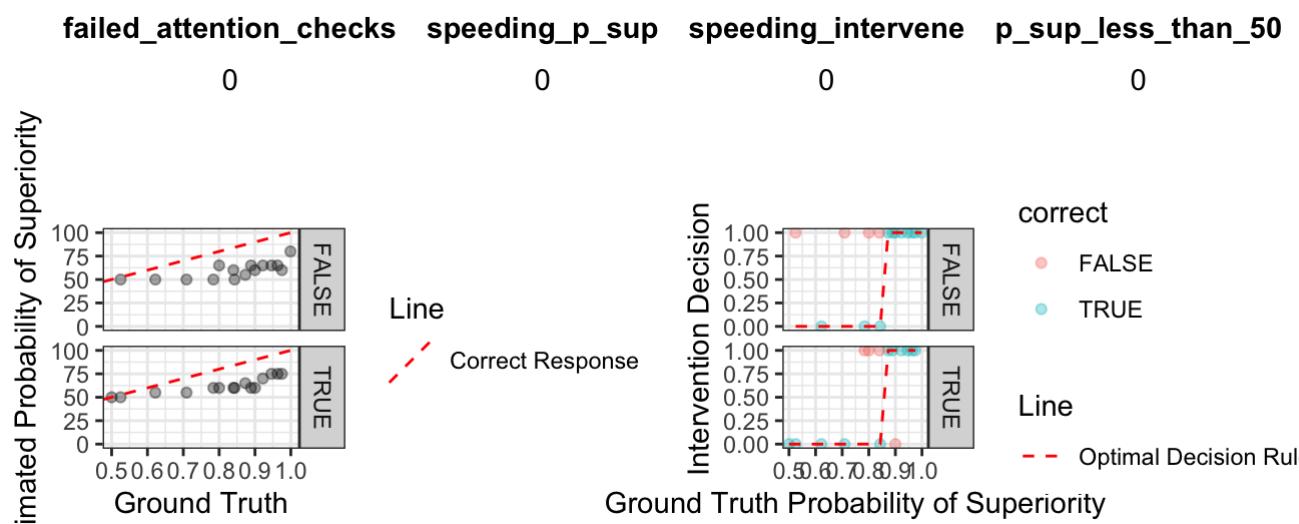
condition	gender	age	education	chart_use	numeracy
intervals	F	25-34	Some college, no degree	Monthly or less	7



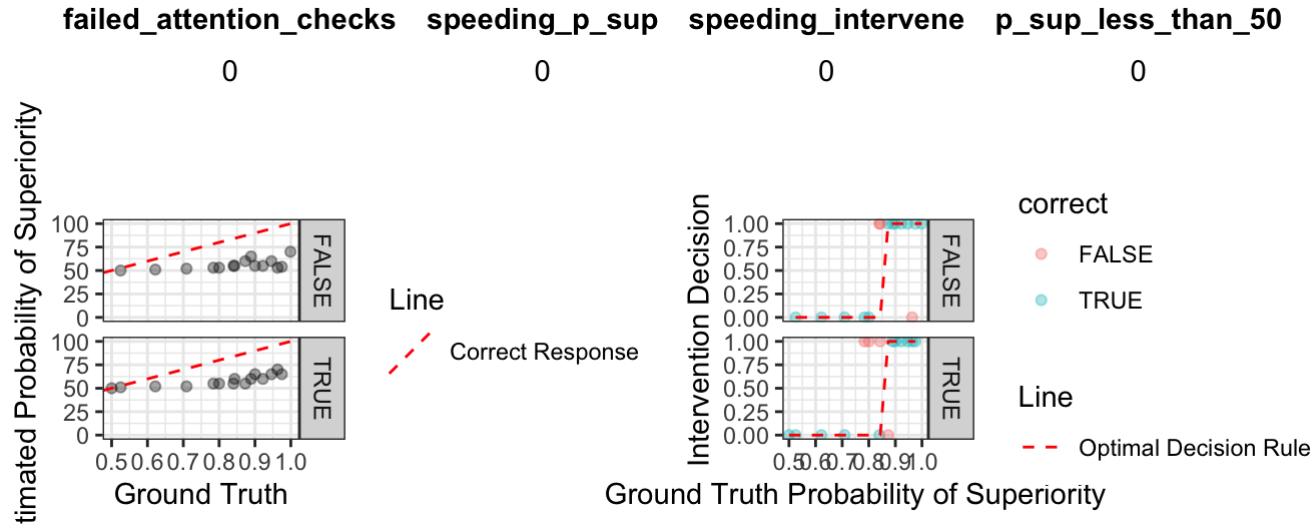
condition	gender	age	education	chart_use	nnumeracy
HOPs	F	45-54	Bachelor's degree	Weekly	11



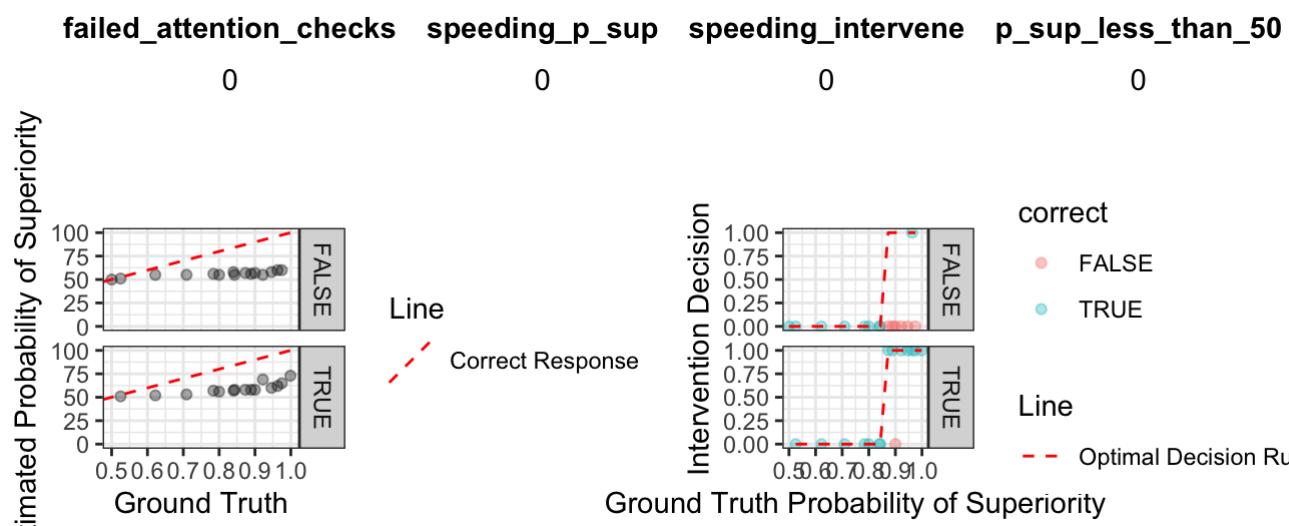
condition	gender	age	education	chart_use	nnumeracy
HOPs	M	25-34	Bachelor's degree	Daily	10



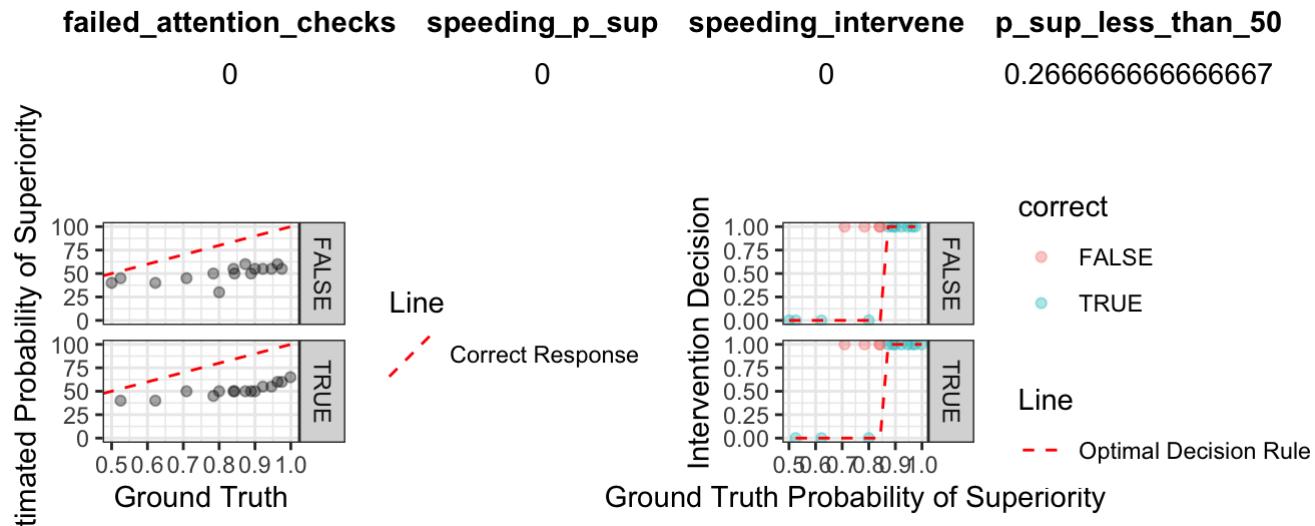
condition	gender	age	education	chart_use	numeracy
intervals	M	45-54	Bachelor's degree	Monthly or less	9



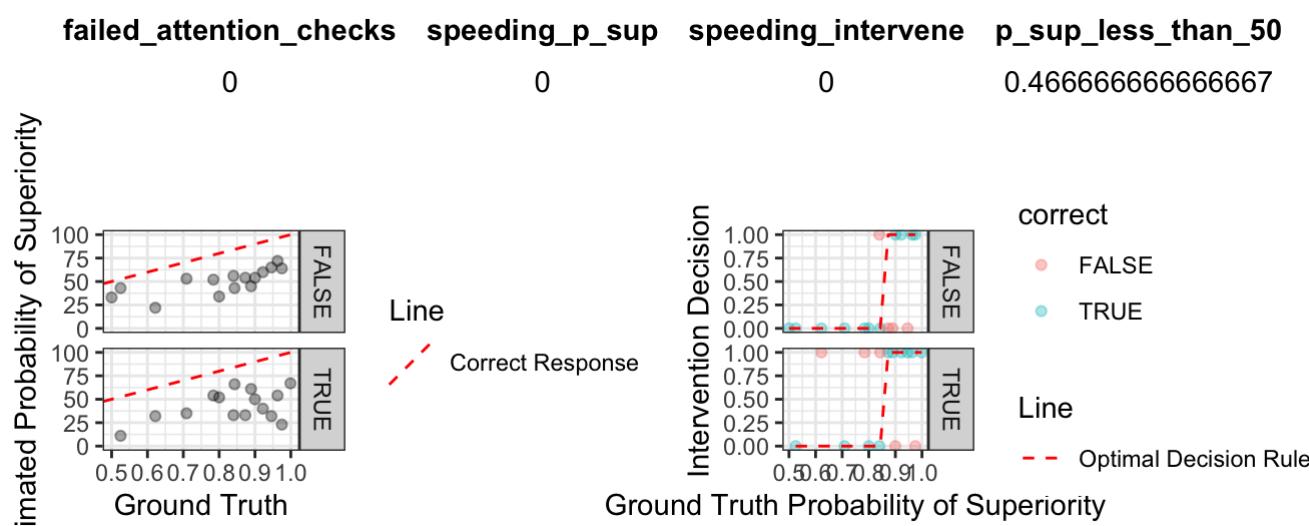
condition	gender	age	education	chart_use	numeracy
intervals	M	25-34	Some college, no degree	Monthly or less	10



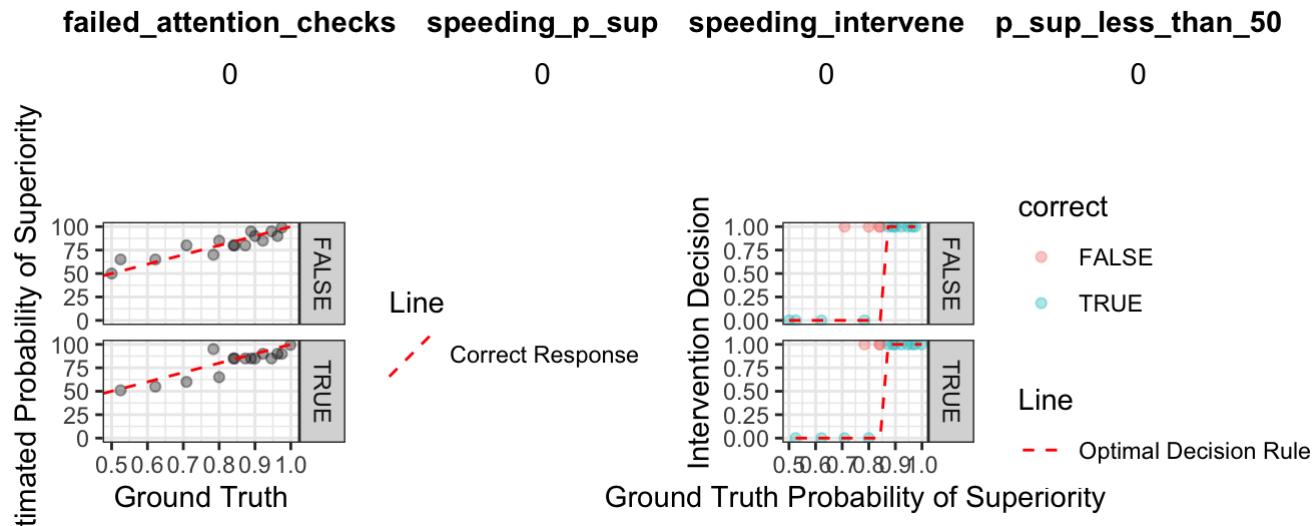
condition	gender	age	education	chart_use	numeracy
HOPs	M	35-44	Bachelor's degree	Daily	10



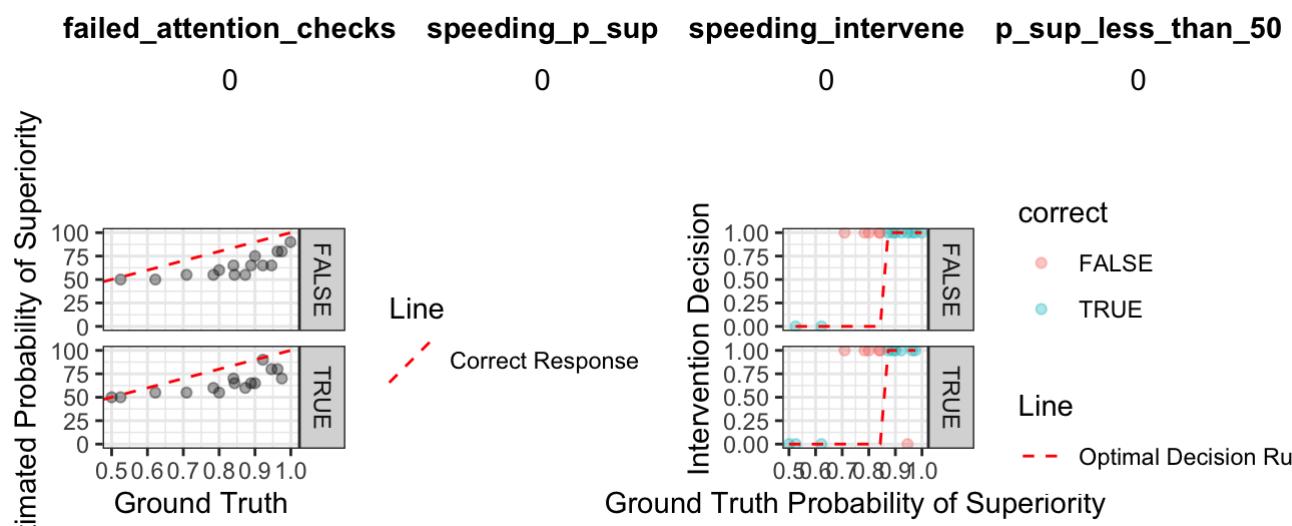
condition	gender	age	education	chart_use	numeracy
HOPs	F	55-64	High school diploma or GED	Monthly or less	8



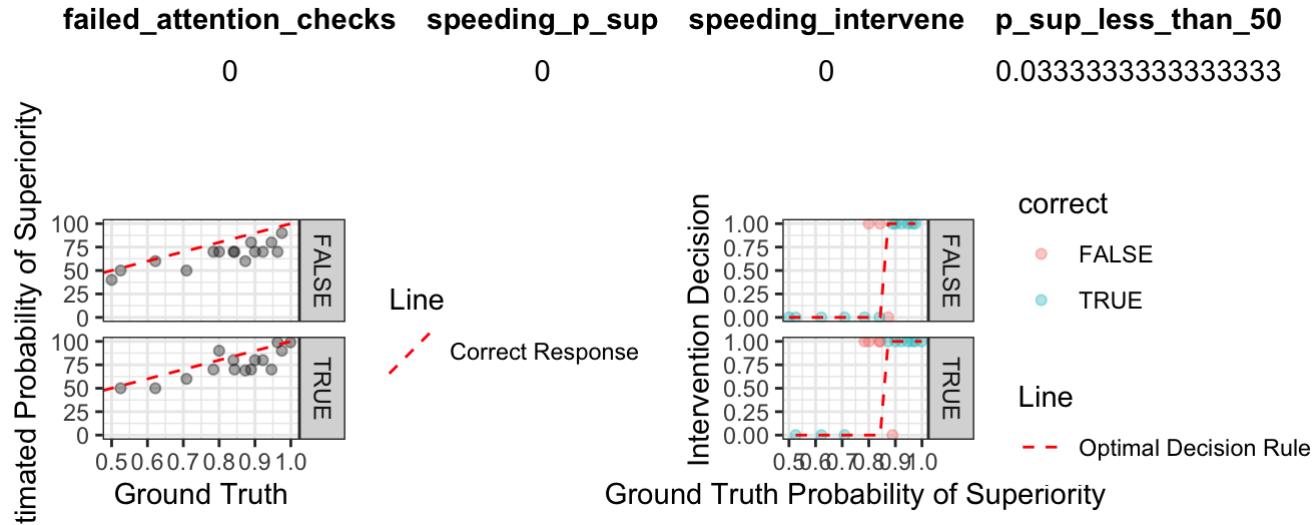
condition	gender	age	education	chart_use	nnumeracy
HOPs	M	35-44	Bachelor's degree	Daily	11



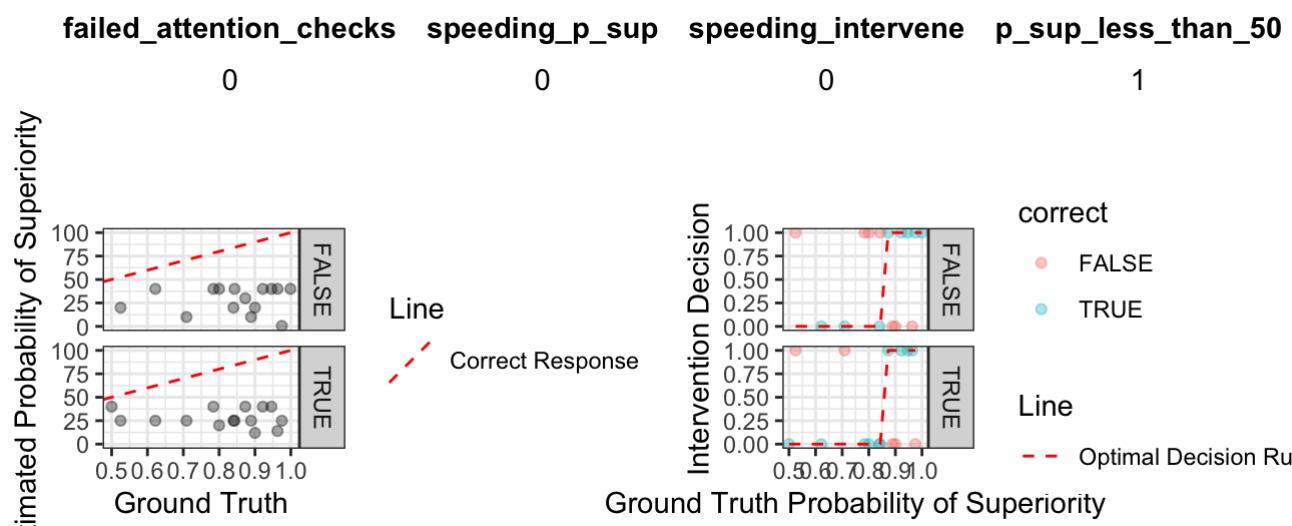
condition	gender	age	education	chart_use	nnumeracy
HOPs	M	35-44	Bachelor's degree	Daily	10



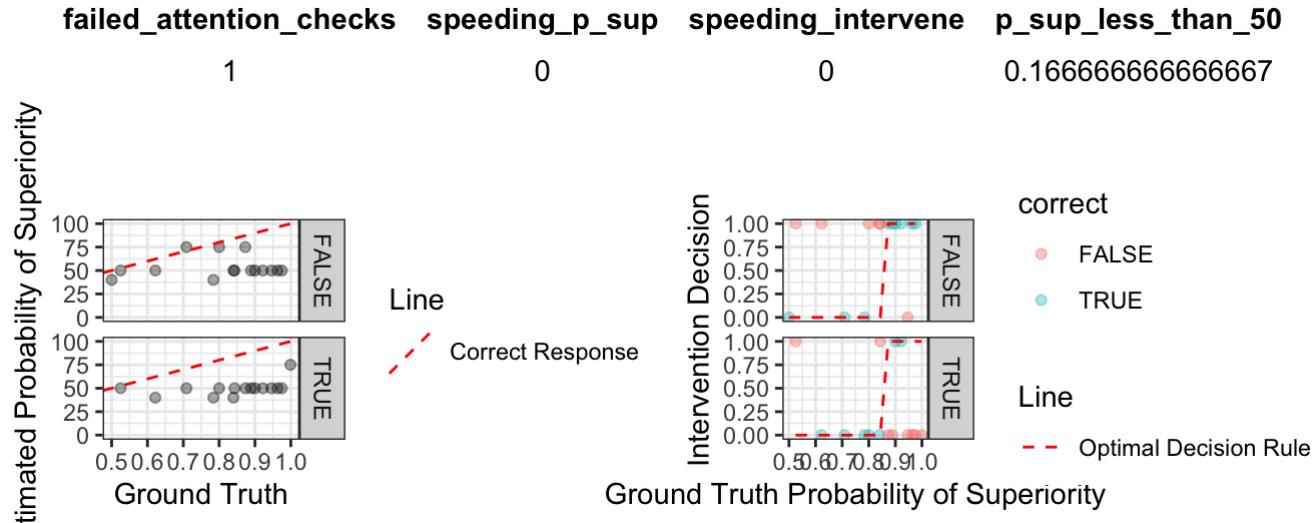
condition	gender	age	education	chart_use	numeracy
HOPs	M	25-34	Associate's degree	Monthly or less	9



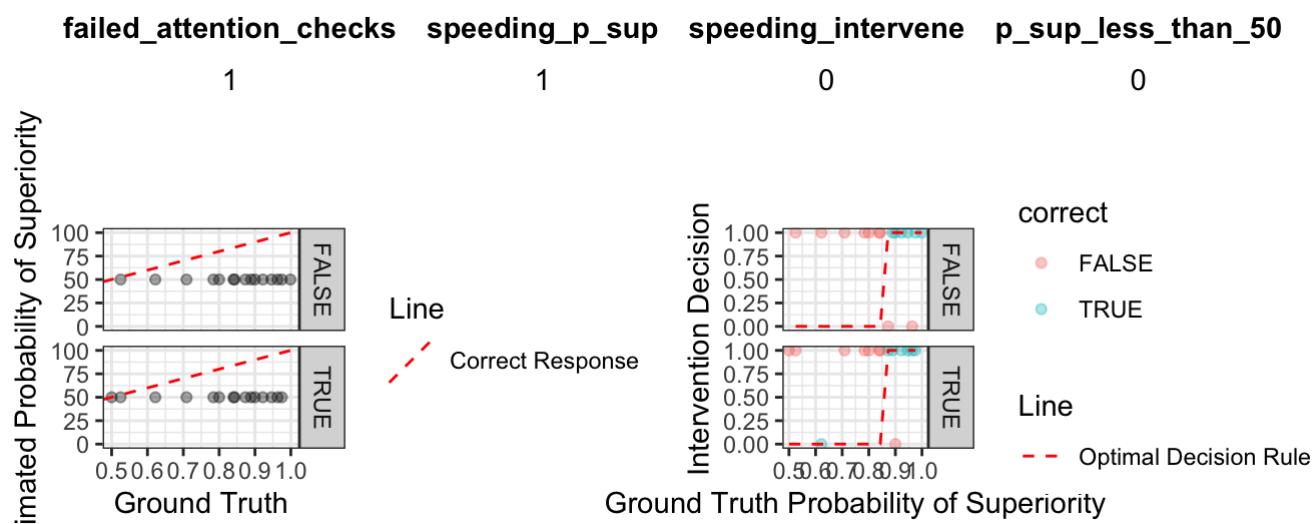
condition	gender	age	education	chart_use	numeracy
HOPs	M	25-34	Bachelor's degree	Monthly or less	7



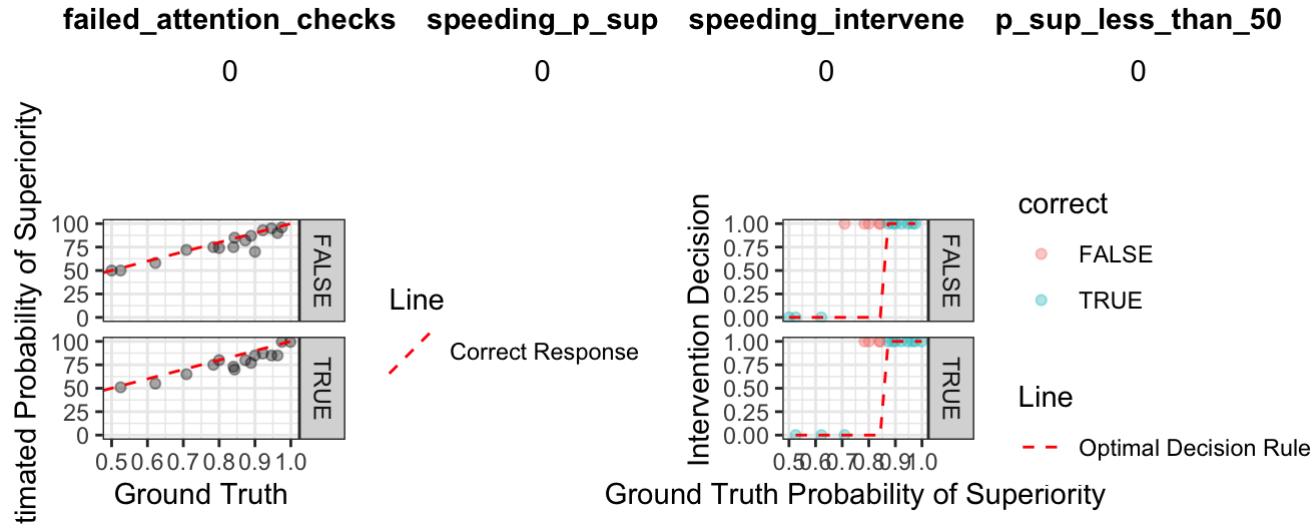
condition	gender	age	education	chart_use	numeracy
HOPs	F	25-34	Some college, no degree	Monthly or less	7



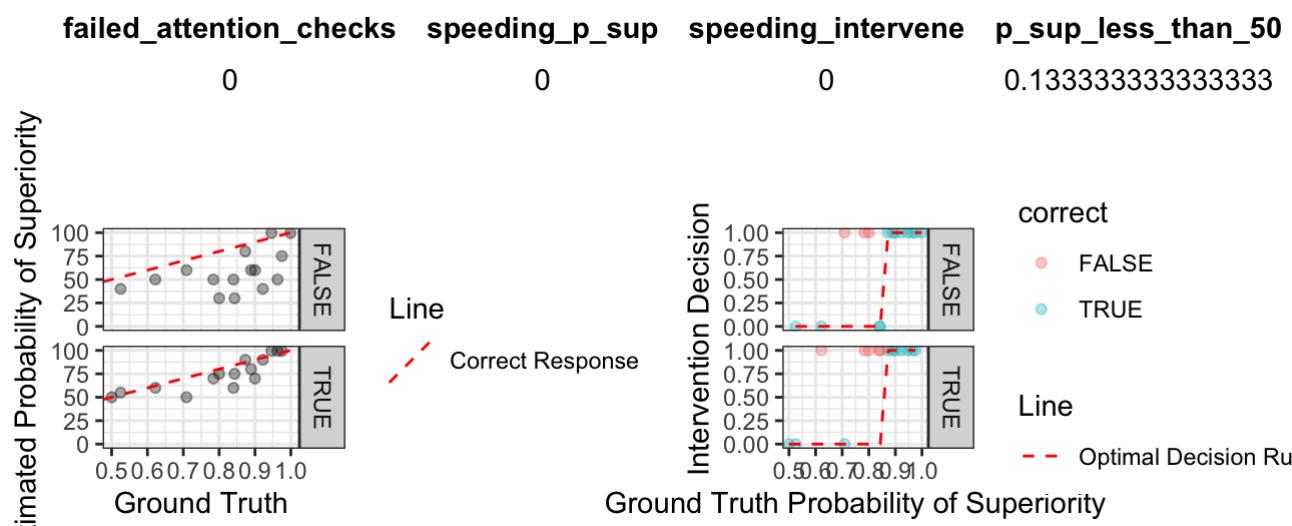
condition	gender	age	education	chart_use	numeracy
intervals	F	45-54	Bachelor's degree	Weekly	9



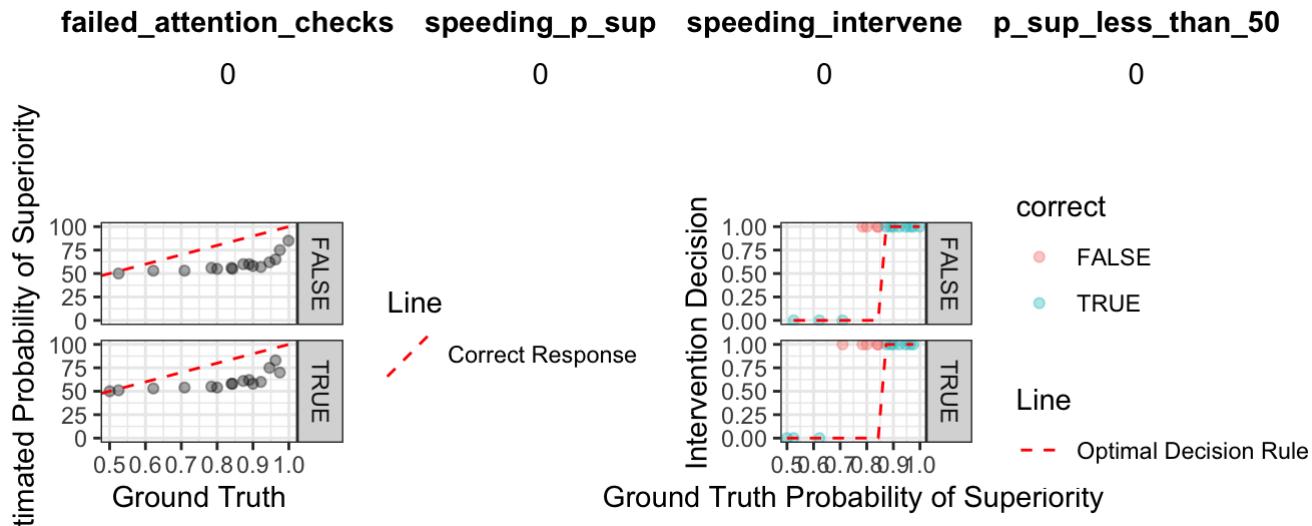
condition	gender	age	education	chart_use	nnumeracy
intervals	M	25-34	Bachelor's degree	Weekly	11



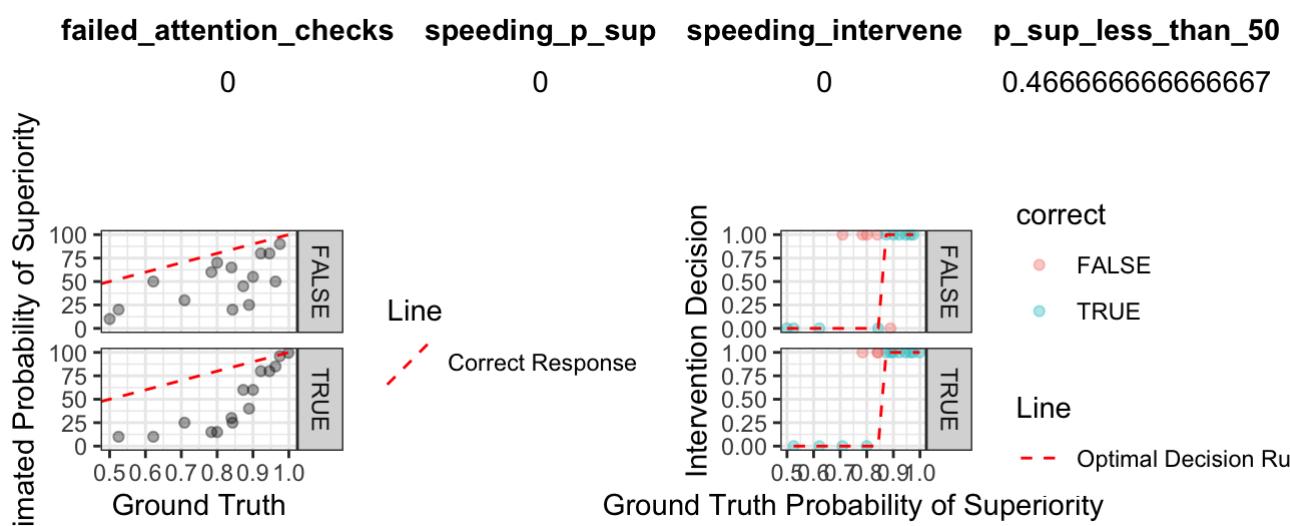
condition	gender	age	education	chart_use	nnumeracy
HOPs	F	35-44	Bachelor's degree	Monthly or less	11



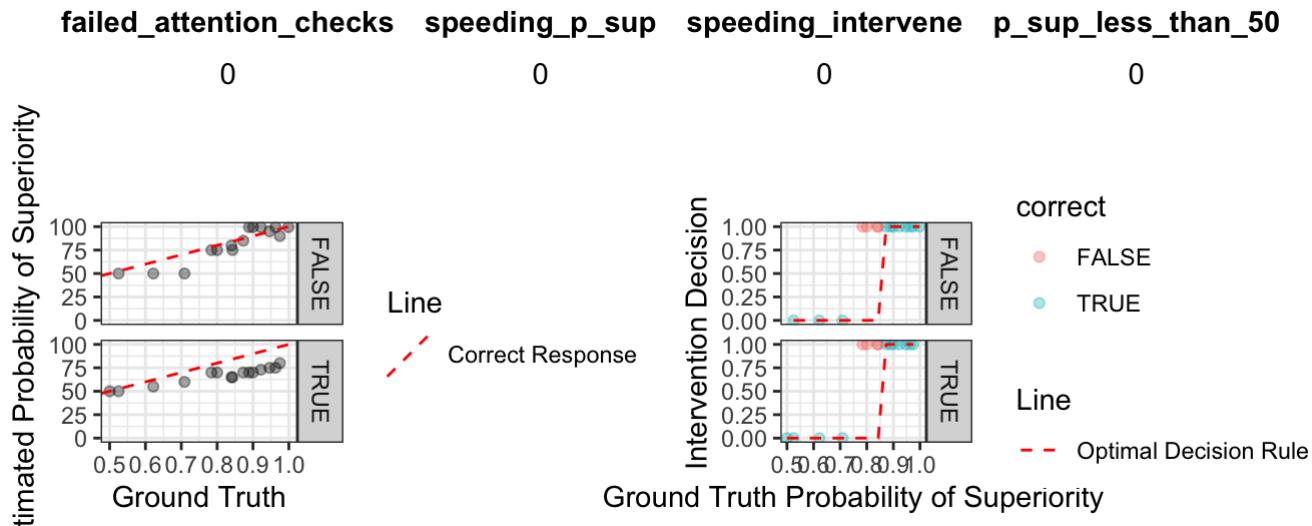
condition	gender	age	education	chart_use	numeracy
intervals	F	25-34	Associate's degree	Daily	10



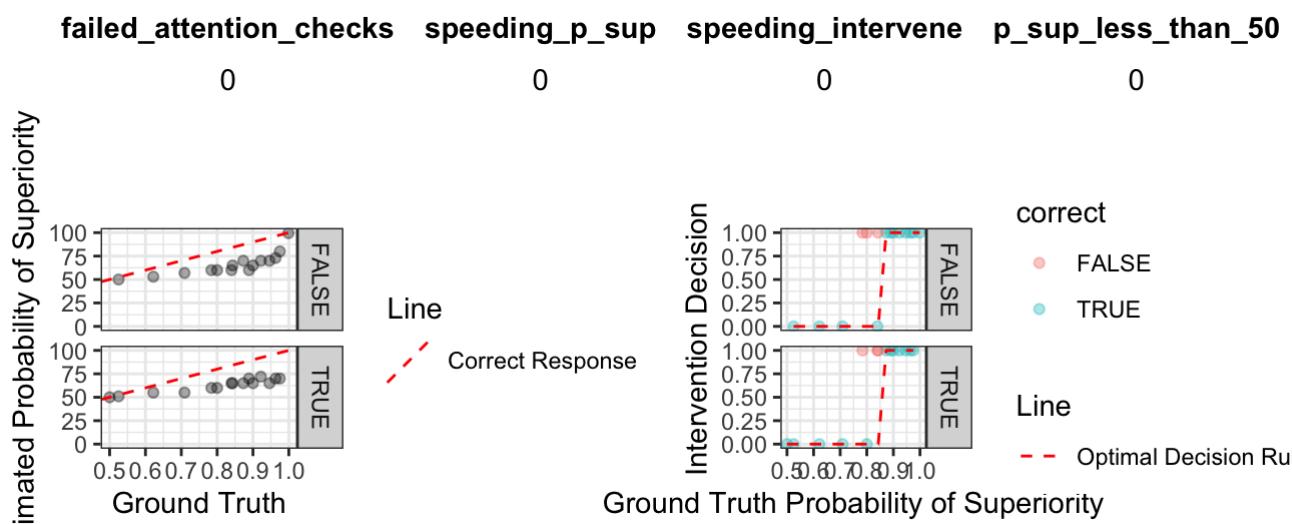
condition	gender	age	education	chart_use	numeracy
HOPs	F	35-44	Bachelor's degree	Weekly	10



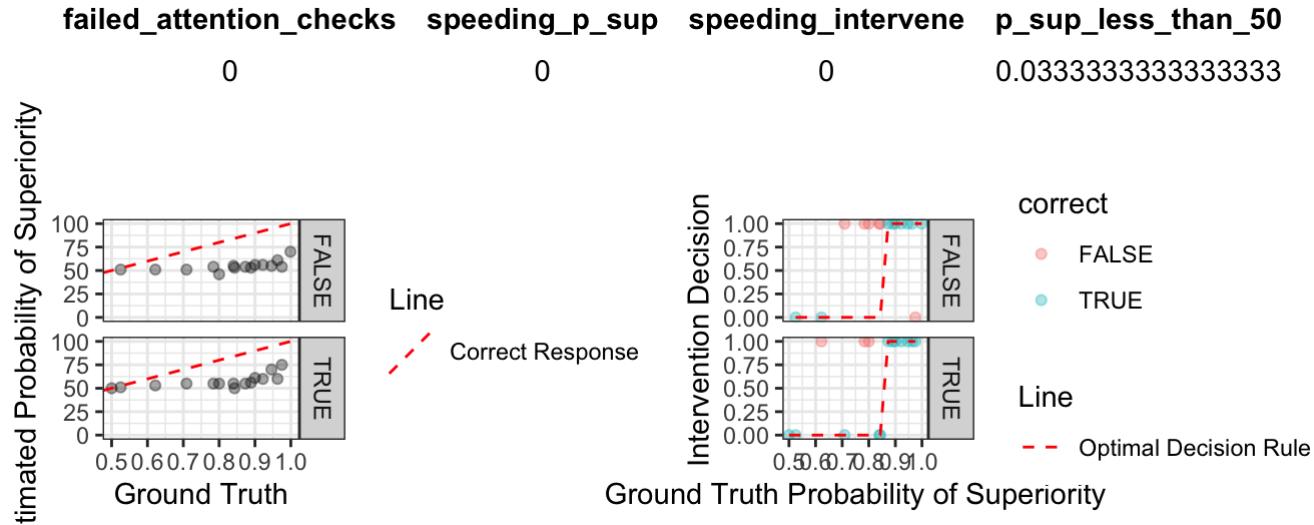
condition	gender	age	education	chart_use	numeracy
HOPs	F	25-34	Bachelor's degree	Weekly	11



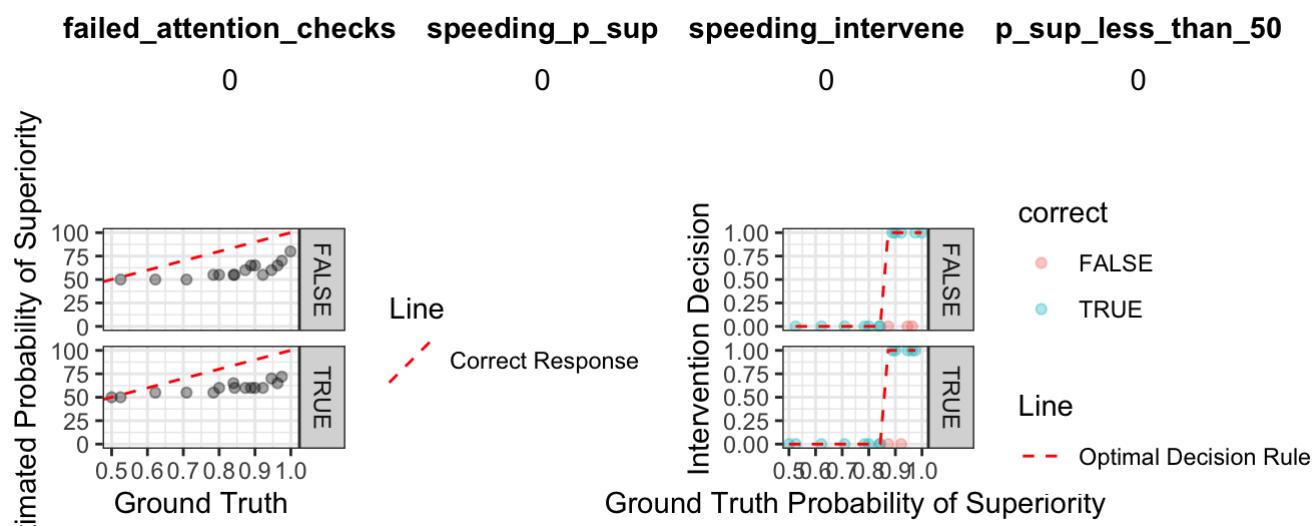
condition	gender	age	education	chart_use	numeracy
intervals	M	25-34	High school diploma or GED	Monthly or less	11



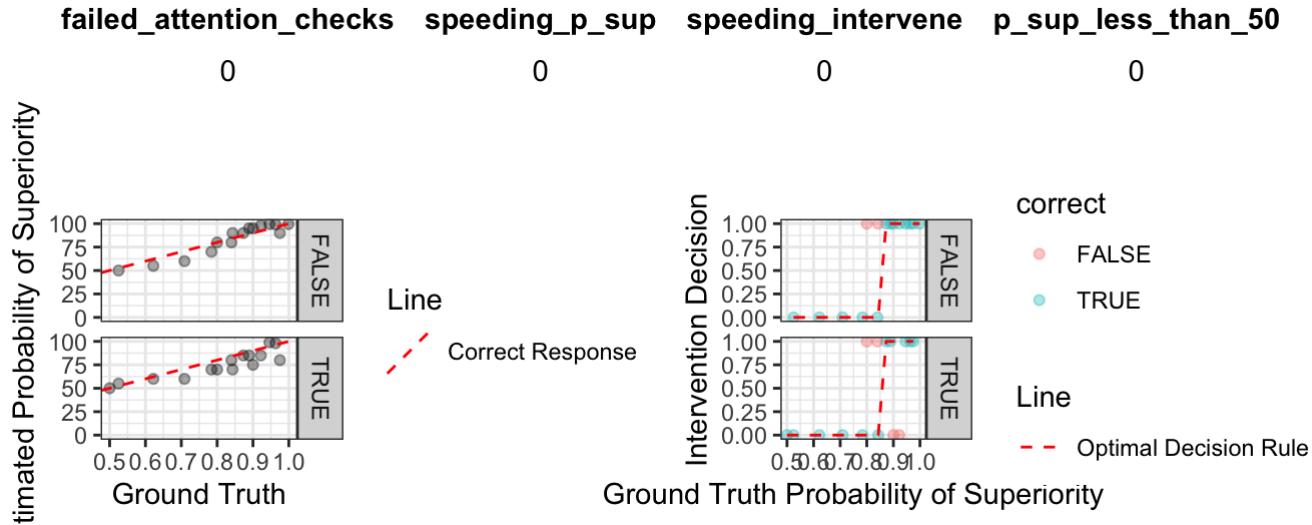
condition	gender	age	education	chart_use	numeracy
intervals	M	25-34	Bachelor's degree	Weekly	9



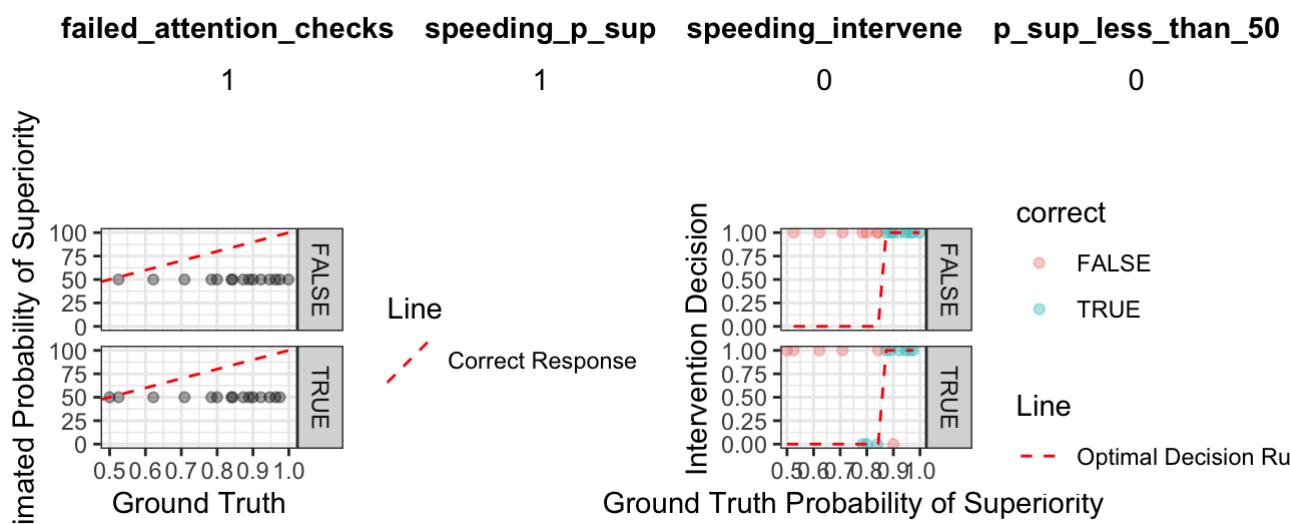
condition	gender	age	education	chart_use	numeracy
HOPs	F	25-34	High school diploma or GED	Daily	9



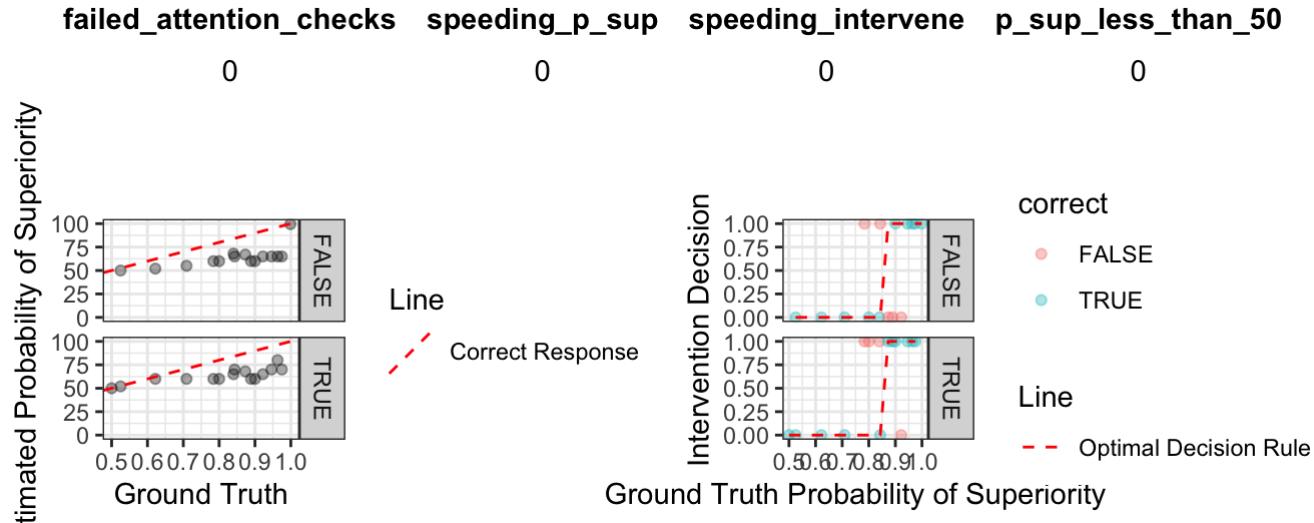
condition	gender	age	education	chart_use	numeracy
intervals	F	35-44	Bachelor's degree	Monthly or less	10



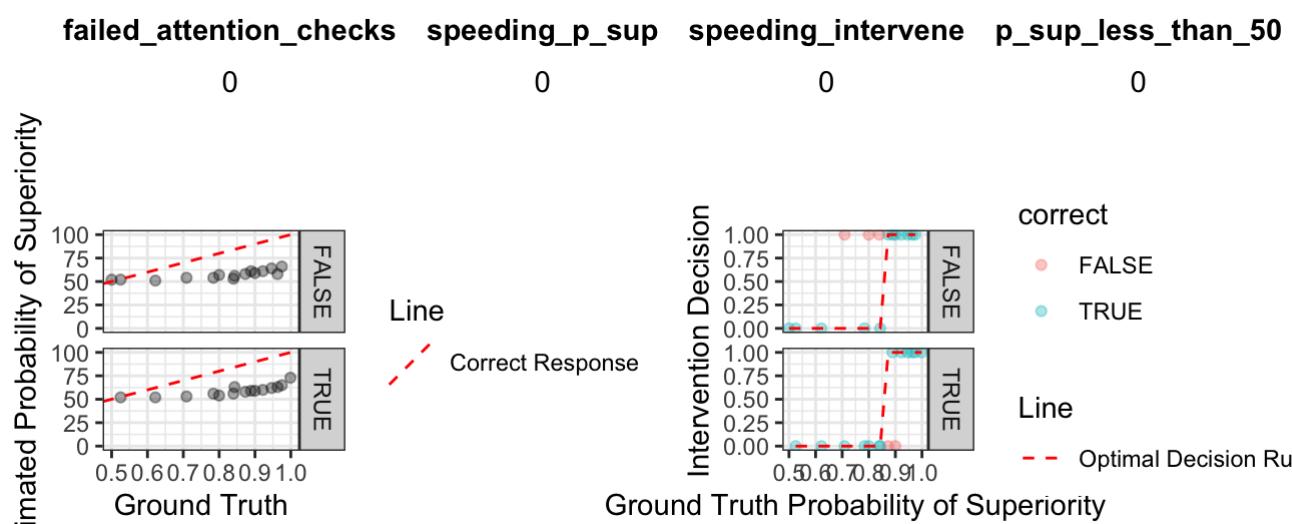
condition	gender	age	education	chart_use	numeracy
intervals	M	25-34	Some college, no degree	Weekly	4



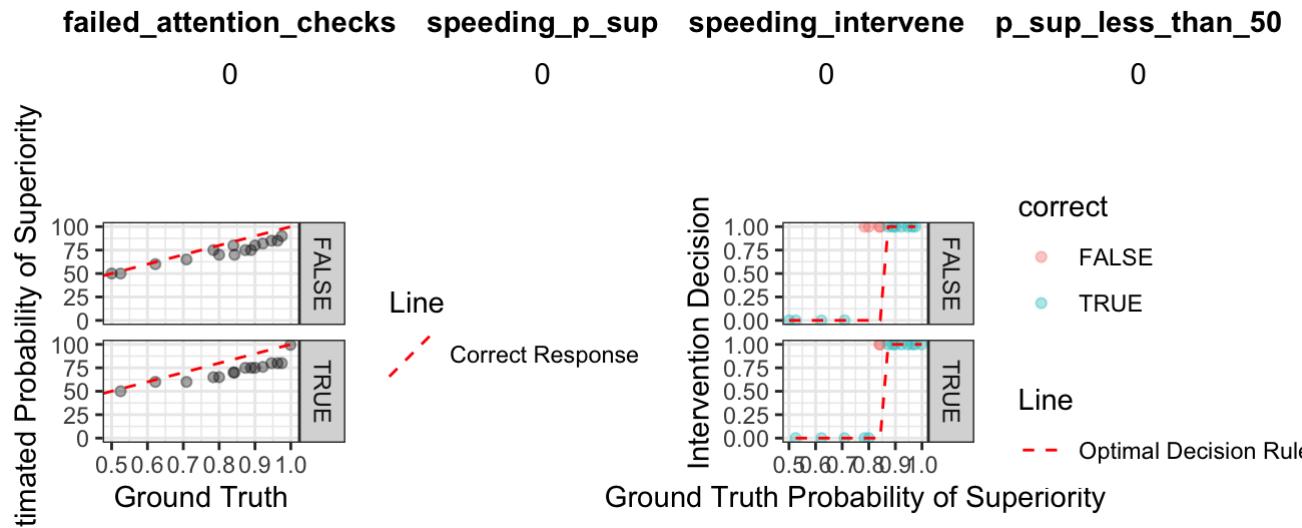
condition	gender	age	education	chart_use	numeracy
intervals	M	25-34	Some college, no degree	Monthly or less	10



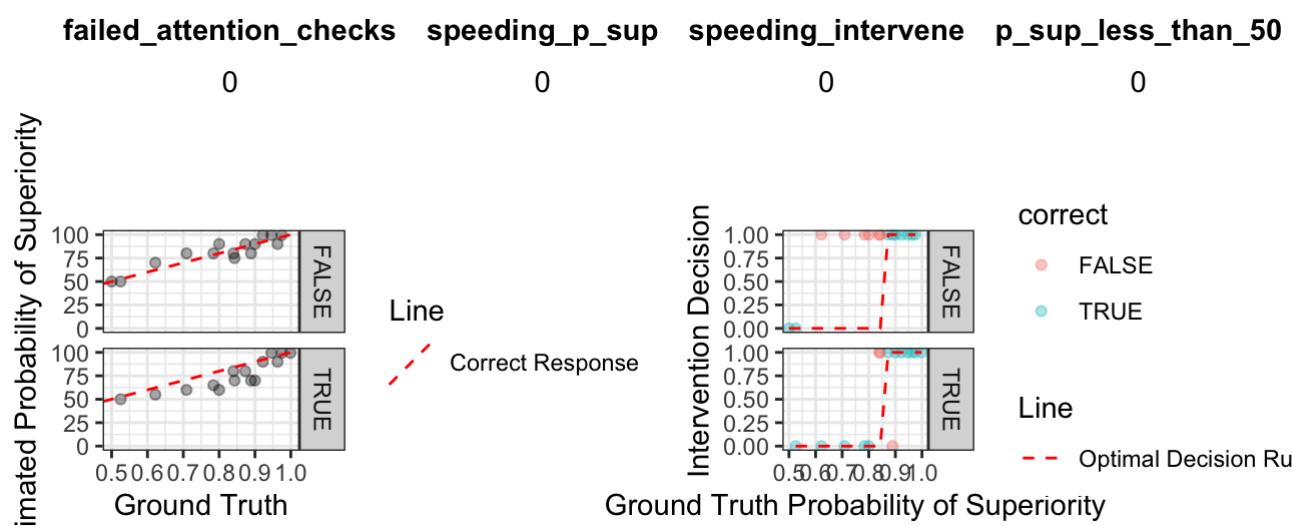
condition	gender	age	education	chart_use	numeracy
HOPs	M	25-34	Bachelor's degree	Daily	11



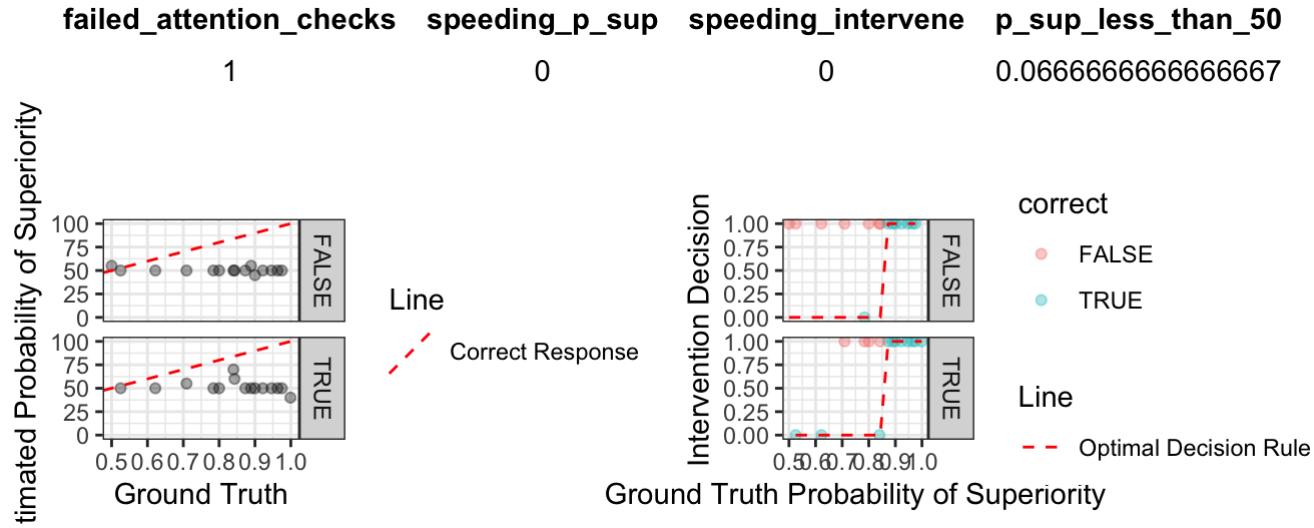
condition	gender	age	education	chart_use	numeracy
intervals	M	18-24	Bachelor's degree	Daily	11



condition	gender	age	education	chart_use	numeracy
HOPs	M	35-44	High school diploma or GED	Monthly or less	11



condition	gender	age	education	chart_use	numeracy
HOPs	M	55-64	Associate's degree	Monthly or less	10



condition	gender	age	education	chart_use	numeracy
intervals	M	35-44	Bachelor's degree	Monthly or less	10

