

Pilot 2

Alex Kale

5/14/2019

Recap of Pilot 1

In pilot 1 we asked participants to view distributions of past scores for two teams competing in a game and judge the probability that Team A would score higher than Team B in a future game (a.k.a., probability of superiority or common language effect size). We then asked them to bet between 1 and 1000 coins that team A would win. We structured the payoff so that the utility-optimal bet was a linear function of the probability that Team A would win. Each participant completed 20 trials at different combinations of effect size and standard deviation of scores using one of three visualizations: means and intervals, means only, or HOPs.

We found that user judgments of probability of superiority were substantially more biased toward 50% in the conditions that emphasized the mean. However, the results of the betting task were less conclusive. Bets seemed to cluster around 1, 500, and 1000 coins suggesting that participants did not use the full granularity of the betting scale. We tried categorizing bets as low, medium, and high but still were unable to make much of the results. One possible reason for this is that the betting data were very noisy, perhaps because the payoff scheme of the task required participants to trade off a tiered tax on winnings from their bet vs a constant tax on the amount they did not bet. Even though there are real world financial decisions that resemble these incentives, it is highly unlikely that users were able to form an intuition about what their decision rule should be.

Pilot 2

In pilot two, we keep the probability of superiority judgment that worked well in pilot 1, but we change the incentivized decision task.

Inspiration for New Betting Task

We draw inspiration from Joslyn's road salting paradigm where participants must decide whether to spend money to salt the roads based on weather forecasts about the probability of freezing temperatures. The key aspect of this study that we find compelling is that participants are deciding whether to make an intervention by trading off the cost of intervention with the cost and probability of potential damages. However, this task is slightly simplified relative to real world situations where people make decisions based on effect size information. First, in Joslyn's task applying salt to the road guarantees the prevention of damage entirely, whereas in most situations intervention will only reduce either the probability or the severity of potential damages. Second, in Joslyn's task participants should always intervene if probabilities of freezing temperatures exceed 17%, whereas in real world situations stakes are dynamic and decisions may be subject to non-linear distortions in the perception of probability. Third, Joslyn's paradigm only examines incentivized decisions which are framed as potential losses, whereas real world situations are sometimes framed as potential gains. Also, the use of the familiar context of weather introduces biases which may not generalize such as the tendency to mistake confidence intervals on forecasts for daily low and high temperatures. In our new task, we aim for a similarly intuitive incentivized decision about interventions, but we also modify the scenario and payoff scheme to make it more realistic and representative of the types of situations in which people make decisions based on effect size information.

New Betting Task: Scenario and Payoff Scheme

In our new task, users play the role of an owner of a mid-sized manufacturing business that makes widgets. They are presented with scenarios in which they must make a decision about whether to continue to use their existing machinery or to pay for a new machine based on the probability of gaining or losing a contract.

In the gain framing trials, users are told they will *pick up a new contract* worth \$X if they produce more than a certain *number of widgets* next year, but they can improve their chances of getting that contract if they pay \$C for a new machine. The optimal decision rule in this case maximizes expected gains.

$$X * p(\text{contract} | \sim \text{machine}) < X * p(\text{contract} | \text{machine}) - C$$

If we assume a constant ratio between the value of the contract and the cost of intervention $K = \frac{X}{C}$, this can be expressed as in terms of the difference in probability of gaining the contract with and without a new machine.

$$p(\text{contract} | \sim \text{machine}) + \frac{1}{K} < p(\text{contract} | \text{machine})$$

The decision rule can also be expressed in terms of the risk ratio of getting the contract with vs without intervention.

$$1 + \frac{1}{K * p(\text{contract} | \sim \text{machine})} < \frac{p(\text{contract} | \text{machine})}{p(\text{contract} | \sim \text{machine})}$$

The loss framing trials are similarly set up. Users are told they will *lose an existing contract* worth \$X if they produce more than a certain *number of defective widgets* next year, but they can improve their chances of keeping that contract if they pay \$C for a new machine. The optimal decision rule in this case minimizes expected losses.

$$X * p(\sim \text{contract} | \sim \text{machine}) > X * p(\sim \text{contract} | \text{machine}) + C$$

Again, if we assume a constant ratio between the value of the contract and the cost of intervention $K = \frac{X}{C}$, this can be expressed as in terms of the difference in probability of losing the contract with and without a new machine.

$$p(\sim \text{contract} | \sim \text{machine}) - \frac{1}{K} > p(\sim \text{contract} | \text{machine})$$

The decision rule can also be expressed in terms of the risk ratio of losing the contract with vs without intervention.

$$1 - \frac{1}{K * p(\sim \text{contract} | \sim \text{machine})} > \frac{p(\sim \text{contract} | \text{machine})}{p(\sim \text{contract} | \sim \text{machine})}$$

Each trial, users will receive feedback based on their decision and a simulated outcome regarding the contract in question. We will tally the dollar value of contracts gained/kept minus the cost of interventions across trials, and users will receive a proportional amount as a bonus through MTurk.

Visualization Conditions

Users will be shown distributions of the number of widgets or defective widgets produced in past years by both their current equipment and the new machine they might want to buy. Number of widgets will be visualized as *means only*, *means with intervals*, *densities*, *quantile dotplots*, and *HOPs*. These conditions will help us test the impact of the salience of the mean as well as discrete vs continuous presentations of effect size on probability of superiority judgments and incentivized decisions about intervention.

We might also manipulate how data is processed before it is visualized to test different ways of presenting effect size. In addition to showing historical distributions of numbers of widgets and defective widgets for each machine, we could also show single distributions of the difference between the number of widgets produced by the two machines or even a the ratio of the number of widgets produced by the two machines. In the single distribution conditions we would use supplementary text to describe the average widget production with the current machine as a baseline. These manipulations would test whether showing a single derived measure is actually more helpful to users than just showing two distributions. We will need to think about situations in which this comparison should be considered representative.

Data Conditions

We manipulate the probability of the new machine producing more widgets than the old machine ($p_{\text{superiority}}$), sampling at linear intervals in logodds units. When $p_{\text{superiority}}$ is greater than 0.5, the decision task is framed as a gain scenario where the user needs to manufacture at least 500 million widgets next year to get a new contract. When $p_{\text{superiority}}$ is less than 0.5, the decision task is framed as a loss scenario where the user needs to manufacture no more than 75 defective widgets per million next year to keep an existing contract.

We also manipulate the baseline probability of gaining/keeping the contract with the old machine. We sample two levels of this baseline probability: 0.5 where the old machine is as likely as a coin flip to result in the contract, and 0.15 where the old machine is fairly unlikely to result in the contract. These values are chosen in part to guarantee an equal number of trials where intervention is and is not the correct decision at a single level of K .

We control the proportion of the value of the contract over the cost of the new machine (K) setting it equal to 2.1. This is the only value for which there are an equal number of trials where users should vs shouldn't intervene at each level of baseline \times gain/loss framing.

As stated above, we set the threshold for gaining the new contract 500 million widgets and the threshold for keeping the old contract at 75 defective widgets per million.

We control the standard deviation of the distribution of the difference in widgets between the two machines (sd_{diff}) by setting it to 15. In the gain framing this is 15 million widgets. In the loss framing, this is 15 defective widgets per million. Since the value of sd_{diff} is relative to the threshold for gaining/keeping the contract, we can think of this variable as constant across trials.

We derive the mean difference in the number of widgets produced by the new minus the old machine ($mean_{\text{diff}}$) from sd_{diff} and $p_{\text{superiority}}$. We derive the standard deviation of the number of widgets produced by the machines from year to year (sd) from sd_{diff} , variance sum law, and the assumption that the machines have equal and independent variances. We derive the mean number of widgets produced by each machine ($mean$) from the threshold for gaining/keeping the contract, the sd of widgets for each machine, and the $mean_{\text{diff}}$ between the number of widgets for the new minus the old machine. We derive the probability of gaining/keeping the contract from the threshold, mean, and sd .

```

# linear sampling of log odds for ground truth probability of superiority for the new
machine
logodds <- seq(log(0.025/(1-0.025)), log(0.975/(1-0.975)), length.out = 8)
p_superiority <- 1 / (1 + exp(-logodds))

# baseline probability of gaining/keeping a contract with the old machine
baseline <- c(.5, .15)
# baseline <- c(.65, .5, .35)

# ratio (K) of value of contract (X) over cost of intervention (C)
K <- c(2.1)
# K <- seq(1.5, 5, .1)

# initialize data conditions dataframe
conds_df <- data.frame(
  "p_superiority" = rep(p_superiority, length(baseline) * length(K)),
  "baseline" = rep(sort(rep(baseline, length(p_superiority))), length(K)),
  "K" = sort(rep(K, length(p_superiority) * length(baseline)))

# label gain vs loss framing trials based on p_superiority and add contract threshold
s
conds_df <- conds_df %>%
  mutate(frame = if_else(p_superiority > .5, "gain", "loss"),
         threshold = if_else(frame=="gain",
                             500, # million widgets required to gain
                             75)) # defective widgets per million req
uired to keep contract

# add columns for the mean and standard deviation of the difference in the number of
widgets produced by the new vs old machine
# depending on the gain vs loss frame, these values represent millions of widgets vs
defective widgets per million
conds_df <- conds_df %>%
  mutate(sd_diff = 15, # std(new - old)
         mean_diff = sd_diff * qnorm(p_superiority)) # mean(new - old)

# double the length of the dataframe to add information per machine, creating a stimu
lus dataframe with a row per distribution to visualize
stim_df <- map_df(seq_len(2), ~conds_df)
stim_df$machine <- sort(rep(c("new", "old"), length(stim_df$p_superiority)/2))

# add columns for the mean and standard deviation of widgets for each machine and the
probability of gaining/keeping the contract
stim_df <- stim_df %>%
  mutate(sd = sqrt(stim_df$sd_diff ^ 2 / 2), # assume equal and independent variances
         in the number of widgets produced by each machine
         mean = if_else(machine=="old",
                        if_else(frame=="gain", # old machine is at baseline
                                threshold - sd * qnorm(1 - baseline),
                                threshold - sd * qnorm(baseline)),
                        if_else(frame=="gain", # new machine is at difference from bas
eline
                                threshold - sd * qnorm(1 - baseline) + mean_diff,
                                threshold - sd * qnorm(baseline) + mean_diff)),
         p_contract = if_else(frame=="gain", # probability of exceeding threshold to g
ain/keep contract

```

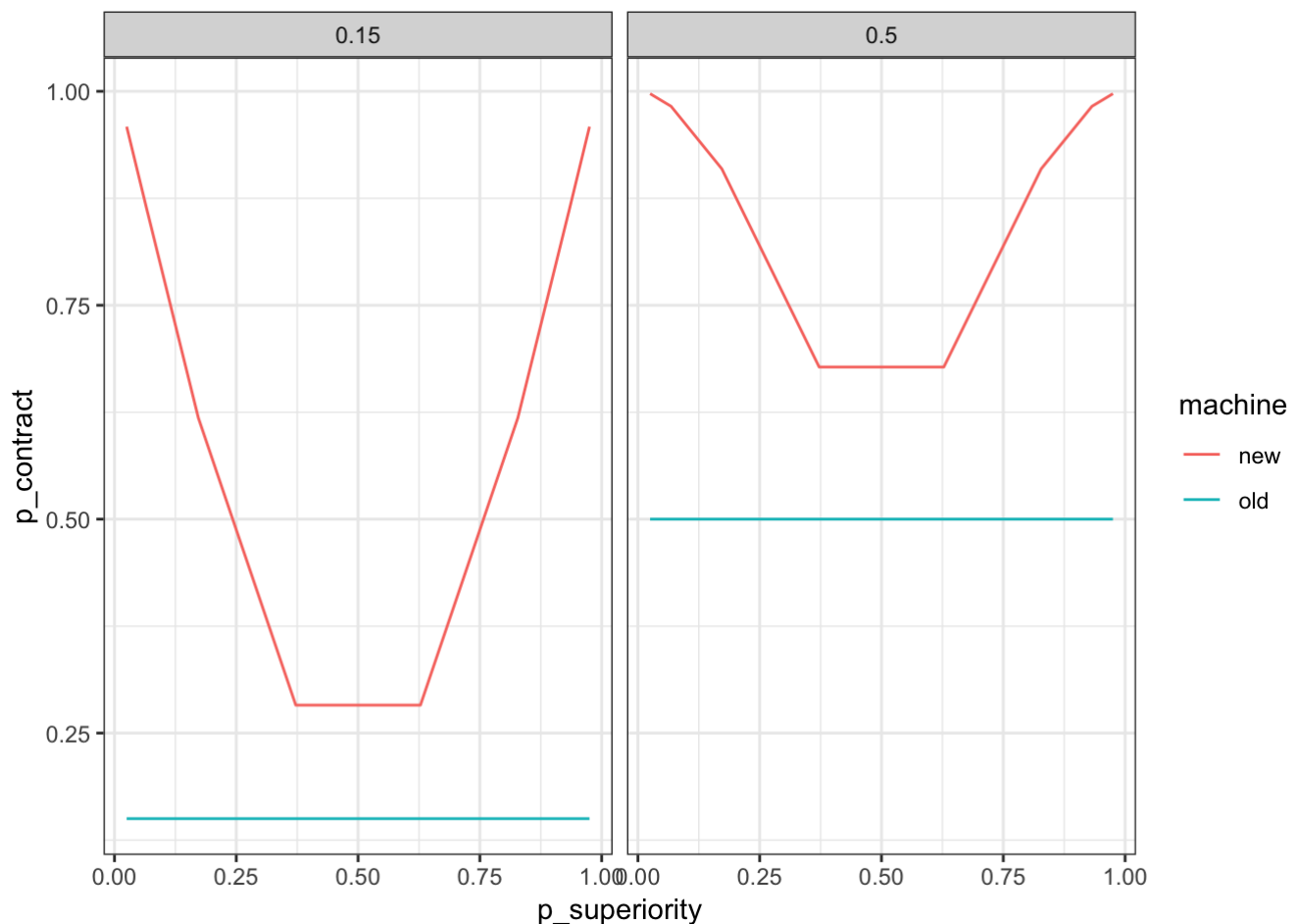
```

1 - pnorm((threshold - mean)/sd),
pnorm((threshold - mean)/sd)))

# spread values per machine across columns to get back to a conditions dataframe one
row per trial
conds_df <- stim_df %>% # explanation: https://kieranhealy.org/blog/archives/2018/11/06/spreading-multiple-values/
  gather(variable, value, -(p_superiority:machine)) %>%
  unite(temp, machine, variable) %>%
  spread(temp, value)

```

This results in an experimental design where the probability of gaining/keeping the contract for the new machine increases monotonically with $p_{\text{superiority}}$. This means that users should intervene only at extreme values of $p_{\text{superiority}}$. Even though the decision rule is not defined in terms of $p_{\text{superiority}}$, users can use effect size as a proxy for the decision task.



We want to check that we have an equal number of trials where intervening is and is not the optimal choice. We also want to make sure that this balance is maintained across all levels of baseline x K x framing.

```
# determine whether or not intervention is utility optimal on each trial
conds_df <- conds_df %>%
  mutate(should_intervene = if_else(frame=="gain",
                                     (old_p_contract + 1 / K) < new_p_contract, # gain
                                     ((1 - old_p_contract) - 1 / K) > (1 - new_p_contr
act))) # loss framing decision rule

conds_df %>%
  group_by(baseline, K, frame) %>%
  summarise(intervene = sum(should_intervene), n_trials = n())
```

```
## # A tibble: 4 x 5
## # Groups:   baseline, K [?]
##   baseline      K frame intervene n_trials
##   <dbl> <dbl> <chr>      <int>    <int>
## 1    0.15    2.1 gain         2         4
## 2    0.15    2.1 loss         2         4
## 3    0.5     2.1 gain         2         4
## 4    0.5     2.1 loss         2         4
```

Although it would be interesting to look at other values of K this would complicate and unbalance the the design of the study. Similarly, it would be interesting to study higher values of baseline such as 0.65 (i.e., Clinton's predicted chance of winning in 2016), where the old machine seems more sure to result in a contract. However, it turns out the intervention is not optimal when the baseline is much above 0.5 unless the intervention costs far less than the value of the contract. This would require different baselines to be tested at different levels of K, which would make it hard to disentangle the effects of baseline vs K.

I'm thinking we can use more extreme values of baseline and K as attention checks. For example, if the baseline probability of gaining/keeping the contract with the old machine is very high, the user should obviously not intervene. If the K is extremely high such that the intervention costs very little relative to the value of the contract, the users should obviously intervene. If users fail in these obvious cases, then we will know that they are not paying attention or do not understand the task. However, for the purpose of the experimental manipulations, I think we should focus on cases where the choice to intervene is more of a toss up in terms of the incentive structure of the situation, and the decision ultimately comes down to the effect size of the intervention (i.e., how much better is the new machine).

Copy

Here's how we will frame the task for participants within the experimental interface.

Instructions Page

In this HIT, you will use charts to inform bets on...

Practice

Here's an example of...

Q1.

Q2.

Task Page

Task

Trial n out of N

The chart below shows... Use this chart to answer the questions below.

Q1.

Q2.

Betting Prompt

(This should be displayed somewhere on the task screen as a reminder.)