

Team mTurk - Motivating Quality Work

Kevin Hanna, Kevin Stone, Changjing Zhao

Contents

Motivating Quality Work	1
What motivates crowdsourced workers to do quality work?	1
Our Datasets	1
1. Our First Pilot	2
2. Our Second Pilot	3
2.1 Score Summary Statistics	3
2.1 Power Test	4
2.2 Analysis	4
2.3 Learnings from our second experiment	5
Experiment 3: Treatment, incentive of future work	5
3.1 Simple regression analysis	5
3.2 Analysis with covariates	5
3.3 Power Test	7
Experiment 4, More data	7
4.1 Analysis	7
4.2 Power Test	7
4.3 More data	8
Experiment 5, threats don't work	8
Experiment 6, threats still don't work	8

Motivating Quality Work

What motivates crowdsourced workers to do quality work?

Our scoring metric measures the accuracy of the bounding box by calculating the euclidean distance of the Turkers bounds to the correct bounding. Therefor a **lower score is better**. When the treatment should cause a negative reaction, the score should increase if our hypothesis is correct.

Our Datasets

1. Bound 20 images with negative treatment (Government Surveillance)
2. Bound a single image with negative treatment (Government Surveillance)
3. Bound a single image with positive treatment (Potential future work)
4. Increase subjects for above dataset, 3
5. Bound a single image with negative treatment, reward 2 cents (Threat of not paying for poor performance)
6. Bound a single image with negative treatment, increased reward to 5 cents (Threat of not paying for poor performance)
7. Increase subjects for above datasets 3 & 4 above, smaller reward.

dataset_no	is_pilot	in_treatment	count	mean_score	std_dev
1	1	0	397	137.60402	295.29711
1	1	1	396	136.39470	296.29412
2	1	0	187	15.25847	36.79851
2	1	1	189	17.09362	42.27343

dataset_no	is_pilot	in_treatment	count	mean_score	std_dev
3	0	0	48	19.02776	23.08104
3	0	1	47	22.71446	36.73351
4	0	0	93	40.35981	139.47131
4	0	1	94	14.51884	25.36227
5	0	0	96	13.55187	23.01864
5	0	1	97	11.61424	11.00214
6	0	0	94	13.56319	20.70507
6	0	1	92	13.15357	17.04807
7	0	0	181	21.52917	98.68055
7	0	1	191	13.17927	16.12633

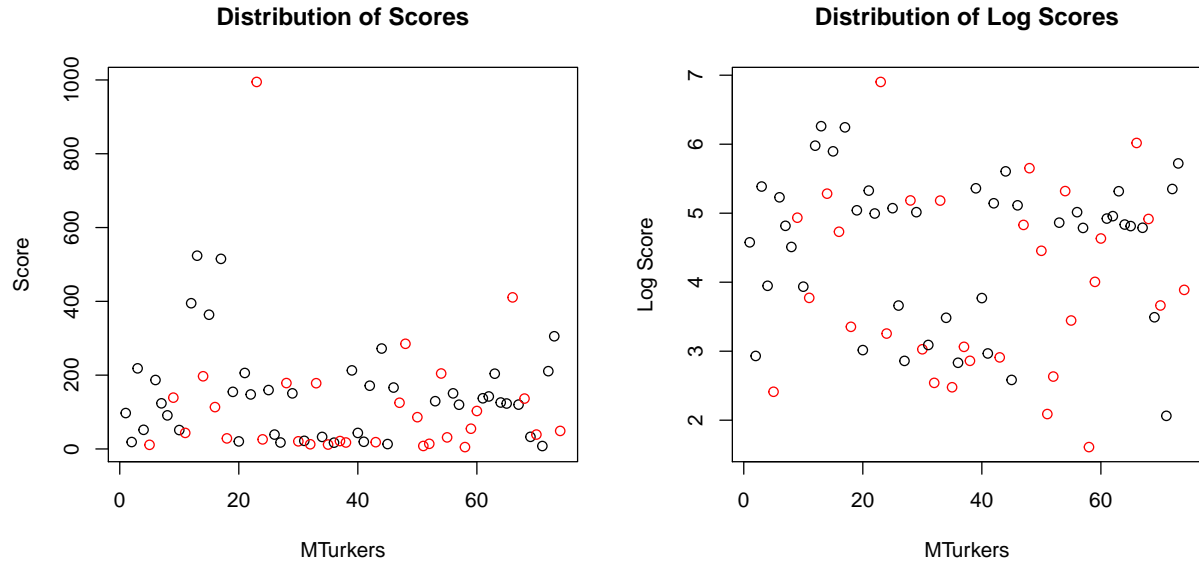
dataset_no	Mean Score	In Treatment	In Control	Total	Has is_mobile	Reward Amount	Standard Deviation
1	137.00089	396	397	793	0	\$0.02	295.60836
2	16.17850	189	187	376	0	\$0.02	39.59534
3	20.79096	47	48	95	93	\$0.20	30.26851
4	27.36948	94	93	187	182	\$0.20	100.54784
5	12.57798	97	96	193	0	\$0.02	17.98909
6	13.36058	92	94	186	0	\$0.05	18.93443
7	17.20553	191	181	372	360	\$0.05	69.52288

bounding_box_score
Min. : 1.000
1st Qu.: 6.681
Median : 12.414
Mean : 60.752
3rd Qu.: 27.917
Max. :1284.400
NA's :18

1. Our First Pilot

For our pilot, we gave the Turkers a negative treatment and asked that they draw a single bounding box on each of 20 images. We first collected some information about the subject through a survey and then randomly assigned those subjects to treatment and control. Our primary goal was to understand how our scoring scheme worked, gauge level of variance we should expect in future experiments and test if our covariates collected from our survey were helpful. We had high attrition and due to a misunderstanding of the Mechanical Turk platform, our assignments to treatment and control failed and we ended up with Turkers not in our experiment in our results, and many ended up in both treatment and control.

We were not able to trust any ATE, but we could at least see the variance, which was exceptionally high.



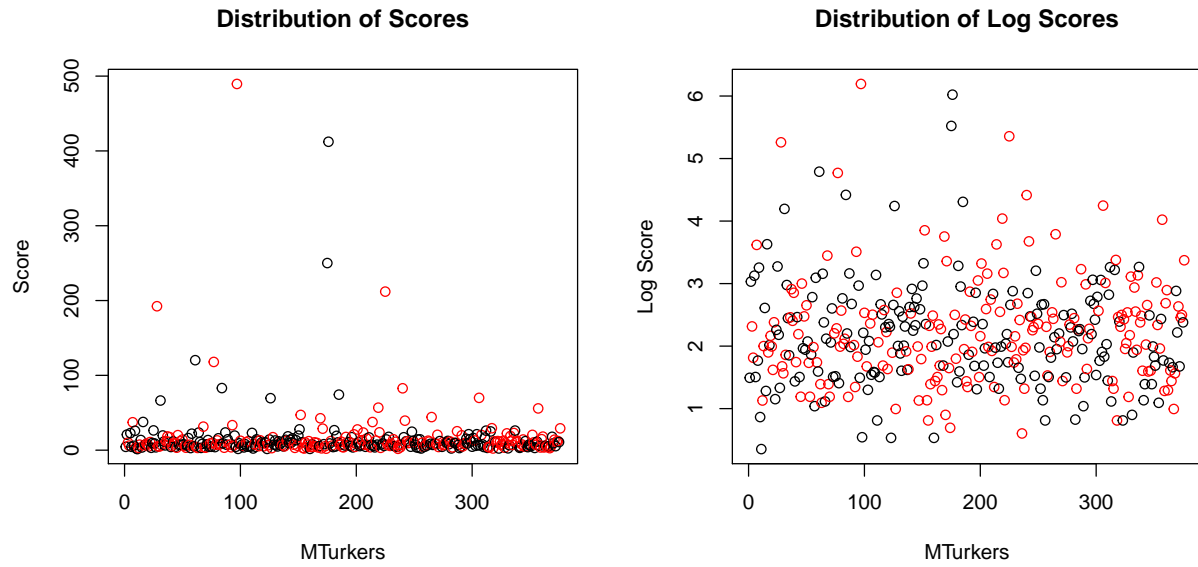
mean_worker_score
Min. : 5.003
1st Qu.: 25.938
Median :119.894
Mean :135.288
3rd Qu.:178.160
Max. :994.601
NA's :1

in_treatment	mean_score	std_dev
0	146.7838	125.4300
1	118.8101	190.9985

#TODO Gauge if effort decreases with more HITs

2. Our Second Pilot

With the first pilot behind us, we decided we needed to focus on increasing our statistical power and hypothesized that collecting the same number of bounding boxes but using more subjects & fewer experiments would provide more statistical power. Each subject was presented a single image and created a single bounding box.



2.1 Score Summary Statistics

Summary Statistics for Score

bounding_box_score
Min. : 1.423
1st Qu.: 5.054
Median : 8.320
Mean : 16.178
3rd Qu.: 13.498
Max. :489.540
NA's :3

in_treatment	mean_score	std_dev
0	15.25847	36.79851
1	17.09362	42.27343

2.1 Power Test

To achieve the statistical power of 0.8 at the 0.05 confidence-level with the variance we had in this experient, we would require nearly 5,800 subjects in both control and treatment.

```
##
##      Two-sample t test power calculation
##
##              n = 5756.986
##            delta = 1.835148
##              sd = 39.59534
##      sig.level = 0.05
##              power = 0.8
##      alternative = one.sided
```


 ## NOTE: n is number in *each* group

2.2 Analysis

The results of our regression failed to show any reliable affect of our treatment. The coefficient is negative, which for our scoring means there is a positive influence from the treatment. But with a p-value of 0.66 there no information can be gleaned from this with any confidence.

Table 8:

	<i>Dependent variable:</i>
	bounding_box_score
in_treatment	1.835 p = 0.656
Constant	15.258*** p = 0.00000
Data Subset	All
Observations	373
R ²	0.001
Adjusted R ²	-0.002
Residual Std. Error	39.638 (df = 371)
F Statistic	0.200 (df = 1; 371)
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01

With this pilot, the only covariate we had was the amount of time each Turker spent on the task. And working time doesn't seem to be affected by our treatment.

Table 9:

	<i>Dependent variable:</i>
	WorkTimeInSeconds
in_treatment	-7.720 p = 0.663
Constant	86.059*** p = 0.000
Data Subset	All
Observations	376
R ²	0.001
Adjusted R ²	-0.002
Residual Std. Error	171.347 (df = 374)
F Statistic	0.191 (df = 1; 374)
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01

The results suggest the negative treatment caused Turkers to spend less time on the task, but the p-value is far from statistically significant again.

2.3 Learnings from our second experiment

The estimated 11,600 subjects required to achieve the statistical power we needed was far too many, With a p-value of 0.389, even with the 11,600 subjects, we weren't likely to find a statistically significant ATE. We need to change our experiment and collect more covariates.

Experiment 3: Treatment, incentive of future work

In both of our pilots, we used a treatment which we hypothesized would cause the Turkers in treatment to work less hard, and the ATE was positive, which in our scoring means the bounding was less accurate. We also wanted to test if a positive treatment would have a larger ATE, so the Turkers in treatment were told we were looking for Turkers to perform some future work with the hypothesis that if the Turkers thought of the task as a test with the incentive of future work they would try harder. So we ran a small experiment to test this theory.

3.1 Simple regression analysis

At first look there doesn't seem to be any significant treatment affect, the last p-value had gone down from 0.66 in the previous experiment to 0.56 in this, but we only used a quarter the number of subjects.

Table 10:

	<i>Dependent variable:</i>	
	bounding_box_score Incentivized	Negative Treatment
	(1)	(2)
in_treatment	3.687 p = 0.563	1.835 p = 0.656
Constant	19.028*** p = 0.00004	15.258*** p = 0.00000
Data Subset	All	All
Observations	92	373
R ²	0.004	0.001
Adjusted R ²	-0.007	-0.002
Residual Std. Error	30.379 (df = 90)	39.638 (df = 371)
F Statistic	0.338 (df = 1; 90)	0.200 (df = 1; 371)

Note: *p<0.1; **p<0.05; ***p<0.01

3.2 Analysis with covariates

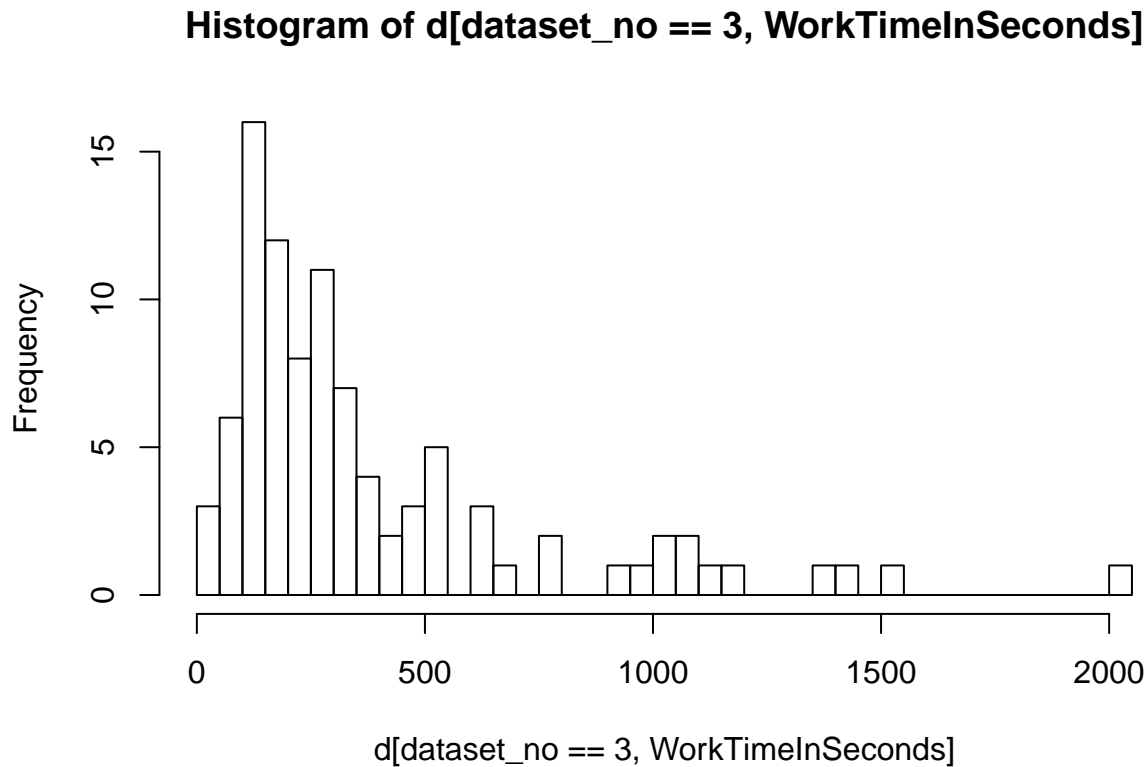
In this experiment we asked the Turkers to answer some questions about the device they were using, their experience doing these types of tasks and some demographic info.

The only covariate which seemed to act as any type of control was the education question, though it wasn't very significant. However, all of the coefficients for the screensize question were negative, and by a fairly significant amount. The baseline value was cellphone, which can be significantly smaller than all the other types of screens. So we tested that on its own.

If the subject is using a cellphone to do the task, their accuracy goes down (score increases), which is intuitive. Having cellphone as a control decreases the p-value from 0.56 to 0.077. With more data, this could be even lower.

As with the previous experiment, we also analyzed how the treatment affected the amount of time they spent on the task.

The regression shows those in treatment on average spent 23 seconds more time, this alone is concerning, as the task itself shouldn't take that much time.



There are a lot of values suggesting that Turkers are not concentrating on our task, it could be they are spawning multiple tabs. Regardless, working time is not helpful for our experiment.

3.3 Power Test

With a lot of speculation about whether our statistical significance would go up with more data, we tested that theory by doing a power calculation.

```
##
##      Two-sample t test power calculation
##
##              n = 834.1739
##              delta = 3.686703
##              sd = 30.26851
##      sig.level = 0.05
##              power = 0.8
##      alternative = one.sided
##
## NOTE: n is number in *each* group
```

The power calculation when using the negative treatment, telling those in treatment that they were doing work for a government surveillance system estimated we needed 5,800 subjects. Using an incentive of possible

future work as the treatment, the ATE has less variance, and estimated that we only need 835 subjects to get 0.80 statistical power.

Experiment 4, More data

To improve the statistical power from Experiment 3, we are adding more data and sending out another 100 control tasks to Turkers and 100 with the same treatment.

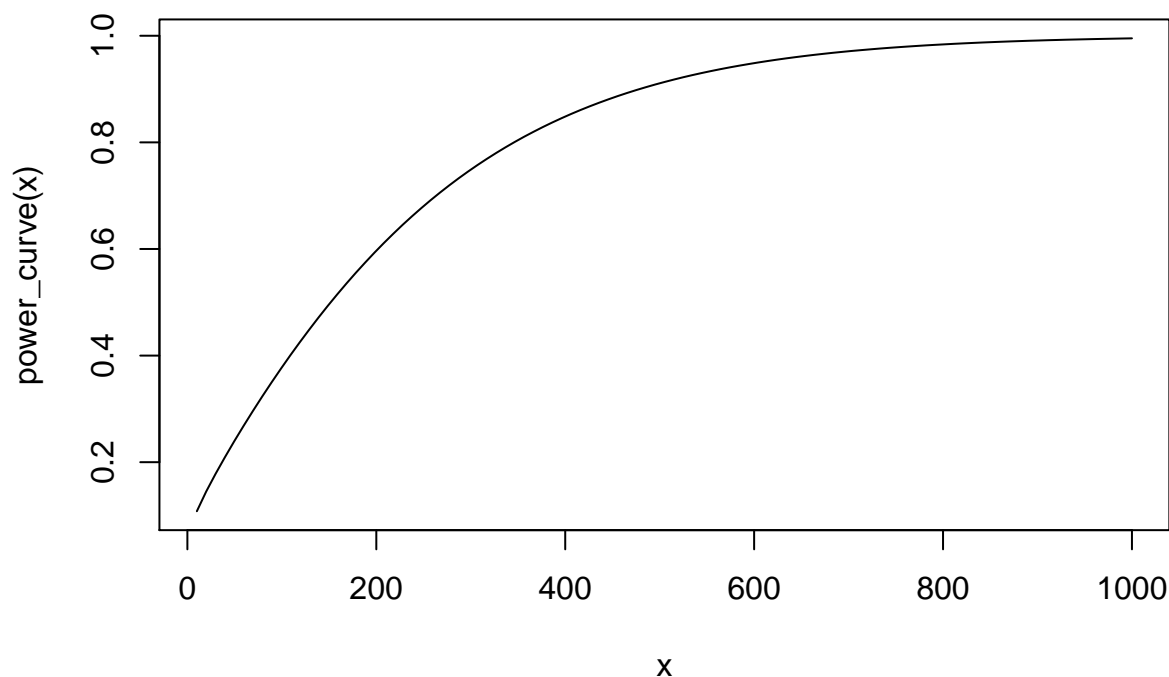
TODO covariate balance check, demographic info show how random it is.

4.1 Analysis

The results are much better, adding another 200 subjects helped decrease the p-value from 0.56 to 0.12, and our ATE is -15.9, a negative number means the bounding boxes from treatment are more accurate. Controlling using mobile devices as a control, we see a much of the variance is explained by the use of mobile devices, though our p-value decreased when we used this control.

4.2 Power Test

```
##
##      Two-sample t test power calculation
##
##              n = 345.7973
##             delta = 15.89495
##              sd = 83.97393
##             sig.level = 0.05
##              power = 0.8
##      alternative = one.sided
##
## NOTE: n is number in *each* group
```

4.3 More data

p-value 0.058

Just some dummy text

Experiment 5, threats don't work

Call: `lm(formula = bounding_box_score ~ in_treatment)`

Residuals: Min 1Q Median 3Q Max -12.005 -7.388 -4.123 0.854 193.311

Coefficients: Estimate Std. Error t value Pr(>|t|)

(Intercept) 13.552 1.848 7.334 6.38e-12 *** in_treatment -1.938 2.606 -0.743 0.458

— Signif. codes: 0 ‘’ **0.001** ’’ 0.01 ’’ 0.05 ‘.’ 0.1 ’’ 1

Residual standard error: 18.01 on 189 degrees of freedom (2 observations deleted due to missingness) Multiple R-squared: 0.002916, Adjusted R-squared: -0.00236 F-statistic: 0.5527 on 1 and 189 DF, p-value: 0.4582

Experiment 6, threats still don't work

Call: `lm(formula = bounding_box_score ~ in_treatment)`

Residuals: Min 1Q Median 3Q Max -12.02 -8.64 -6.39 -1.38 107.83

Coefficients: Estimate Std. Error t value Pr(>|t|)

(Intercept) 13.5632 1.9581 6.927 6.99e-11 *** in_treatment -0.4096 2.7842 -0.147 0.883

— Signif. codes: 0 ‘’ **0.001** ’’ 0.01 ’’ 0.05 ‘.’ 0.1 ’’ 1

Residual standard error: 18.98 on 184 degrees of freedom Multiple R-squared: 0.0001176, Adjusted R-squared: -0.005317 F-statistic: 0.02164 on 1 and 184 DF, p-value: 0.8832

Table 11:

	<i>Dependent variable:</i>			
	Target Alone	Monitor size	bounding_box_score Did task before	Age
	(1)	(2)	(3)	(4)
in_treatment	3.687 p = 0.563	7.794 p = 0.231	4.708 p = 0.470	2.681 p = 0.640
as.factor(monitor)largescreen		-66.462*** p = 0.0004		
as.factor(monitor)midsized		-60.451*** p = 0.0005		
as.factor(monitor)smalllaptop		-57.383*** p = 0.002		
as.factor(monitor)tablet		-35.061* p = 0.064		
as.factor(didbf)no			7.732 p = 0.612	
as.factor(didbf)yes			8.810 p = 0.535	
as.factor(age)31to40				-9.611 p = 0.372
as.factor(age)41to50				97.704*** p = 0.00001
as.factor(age)lto21				-12.327 p = 0.653
as.factor(educ)highschool				
as.factor(educ)masterorabove				
as.factor(educ)somecollege				
as.factor(income)gt30klt60k				
as.factor(income)gt60klt90k				
as.factor(income)gt90k				
as.factor(income)lt10k				
Constant	19.028*** p = 0.00004	72.889*** p = 0.00004	9.539 p = 0.474	18.250*** p = 0.00003

Table 12:

	<i>Dependent variable:</i>	
	bounding_box_score	
	Target Alone	Used Cellphone
	(1)	(2)
in_treatment	3.687 p = 0.563	11.196* p = 0.077
is_mobile		34.558*** p = 0.0002
Constant	19.028*** p = 0.00004	10.296** p = 0.031
Data Subset	All	All
Observations	92	90
R ²	0.004	0.160
Adjusted R ²	-0.007	0.141
Residual Std. Error	30.379 (df = 90)	28.327 (df = 87)
F Statistic	0.338 (df = 1; 90)	8.291*** (df = 2; 87)
<i>Note:</i> *p<0.1; **p<0.05; ***p<0.01		

Table 13:

	<i>Dependent variable:</i>	
	WorkTimeInSeconds	
	Incentivized	Negative Treatment
	(1)	(2)
in_treatment	22.983 p = 0.766	-7.720 p = 0.663
Constant	377.208*** p = 0.000	86.059*** p = 0.000
Data Subset	All	All
Observations	95	376
R ²	0.001	0.001
Adjusted R ²	-0.010	-0.002
Residual Std. Error	374.924 (df = 93)	171.347 (df = 374)
F Statistic	0.089 (df = 1; 93)	0.191 (df = 1; 374)
<i>Note:</i> *p<0.1; **p<0.05; ***p<0.01		

Table 14:

	<i>Dependent variable:</i>	
	bounding_box_score	
	n=100 + n=200	n=100
	(1)	(2)
in_treatment	-15.895 p = 0.116	3.687 p = 0.563
Constant	33.046*** p = 0.00001	19.028*** p = 0.00004
Data Subset	All	All
Observations	277	92
R ²	0.009	0.004
Adjusted R ²	0.005	-0.007
Residual Std. Error	83.748 (df = 275)	30.379 (df = 90)
F Statistic	2.494 (df = 1; 275)	0.338 (df = 1; 90)
<i>Note:</i> *p<0.1; **p<0.05; ***p<0.01		

Table 15:

	<i>Dependent variable:</i>	
	bounding_box_score	
	Target Alone	Controlling for Mobile
	(1)	(2)
in_treatment	-15.895 p = 0.116	-12.436 p = 0.217
is_mobile		51.273*** p = 0.003
Constant	33.046*** p = 0.00001	25.710*** p = 0.001
Data Subset	All	All
Observations	277	270
R ²	0.009	0.041
Adjusted R ²	0.005	0.034
Residual Std. Error	83.748 (df = 275)	82.362 (df = 267)
F Statistic	2.494 (df = 1; 275)	5.694*** (df = 2; 267)
<i>Note:</i> *p<0.1; **p<0.05; ***p<0.01		

Table 16:

	<i>Dependent variable:</i>	
	bounding_box_score	
	Target Alone	Controlling for Mobile
	(1)	(2)
in_treatment	-11.238* p = 0.058	-12.436 p = 0.217
is_mobile	81.852*** p = 0.000	51.273*** p = 0.003
Constant	21.251*** p = 0.00000	25.710*** p = 0.001
Data Subset	All	All
Observations	625	270
R ²	0.071	0.041
Adjusted R ²	0.068	0.034
Residual Std. Error	73.872 (df = 622)	82.362 (df = 267)
F Statistic	23.772*** (df = 2; 622)	5.694*** (df = 2; 267)
<i>Note:</i>		*p<0.1; **p<0.05; ***p<0.01

Table 17:

	<i>Dependent variable:</i>		
	bounding_box_score		
	Reward		
	(1)	(2)	(3)
0.05	0.773 p = 0.685	0.723 p = 0.901	
0.20		12.547** p = 0.019	
0.20 ^a			12.190*** p = 0.007
in_treatment	-1.184 p = 0.535	-7.414* p = 0.094	-7.418* p = 0.093
Constant	13.173*** p = 0.000	16.305*** p = 0.0005	16.663*** p = 0.00001
Data Subset	All	All	$x == 1$
Observations	377	654	654
R ²	0.001	0.016	0.016
Adjusted R ²	-0.004	0.011	0.013
Residual Std. Error	18.477 (df = 374)	56.365 (df = 650)	56.323 (df = 651)
F Statistic	0.278 (df = 2; 374)	3.451** (df = 3; 650)	5.177*** (df = 2; 651)

Note:

*p<0.1; **p<0.05; ***p<0.01