

# Evaluating Non-Expert Annotations on Mechanical Turk for Opinion Mining on Spanish Data

## Abstract

One of the major bottlenecks in the development of data-driven AI Systems is the lack of reliable human annotations. The recent advent of several crowdsourcing platforms such as Amazon's Mechanical Turk, allowing requesters the access to affordable and rapid results of a global workforce, greatly facilitates the creation of massive training sets. Most of the available studies on the effectiveness of crowdsourcing report on English data. We use Mechanical Turk annotations to train an Opinion Mining System to classify Spanish consumer comments. We design three different Human Intelligence Task (HIT) strategies and show high inter-annotator agreement between non-experts and expert annotators. We evaluate the advantages/drawbacks of each HIT design and report <results of the classifier>.

## 1 Introduction

<Intro.> <Some possible citations for remainder paper.> In (Snow et al., 2008), it is shown that ...

In (Sheng et al., 2008), it is shown that ...

In (Kittur et al., 2008), it is shown that ...

In (Su et al., 2007), it is shown that ...

## 2 Task Outline and Goals

We compare different HIT design strategies by evaluating the usefulness of resulting Mechanical Turk (AMT) annotations to train an Opinion Mining System on Spanish consumer data. More specifically,

we address the following research questions:

(i) Annotation quality: how do the different AMT annotations compare to expert annotations? We compare the inter-annotator agreement in expert annotations with the inter-annotator agreement in the different AMT annotations.

(ii) Annotation applicability: how does the performance of an Opinion Mining classifier vary after training on different (sub)sets of AMT and expert annotations? Given a simple classification technique, we evaluate the system performance by using AMT annotations, expert annotations and the combination of both as training data. The idea is not to evaluate the classification technique *per se*, but to measure the influence of the training material.

(iii) Return on Investment (ROI): how does the use of AMT annotations compare economically against the use of expert annotations? AMT offers the possibility of obtaining inexpensive annotations the quality of which tends to be worse than expert annotations <include ref>. We show that, for the task at hand, the ROI is positive.

(iv) Language barriers: x% of all AMT tasks are designed for English speakers <include ref>. How easy is it to get reliable AMT results for Spanish?

## 3 HIT Design

We selected a dataset of 1000 sentences containing user opinions on cars from the automotive section of [www.ciao.es](http://www.ciao.es) (Spanish). This website was chosen because it contains a large and varied pool of Spanish customer comments suitable to train an Opinion Mining System and because opinions in-

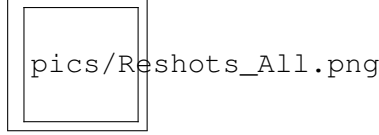


Figure 1: An example sentence (a) and the three HIT designs used in the experiments: (b) HIT1: a simple categorization scheme, (c) HIT2: a graded categorization scheme, and (d) HIT3: a continuous triangular categorization scheme containing both a horizontal positive-negative axis and a vertical subjective-objective axis.

clude simultaneously global numeric and specific ratings over particular attributes of the subject matter. Section 5.1 contains more detailed information about the selection of the dataset. An example of a sentence from the data set can be found in (1):

- (1) 'No te lo pienses más, cómpratelo!'  
(= 'Don't think twice, buy it!')

The sentences in the dataset were presented to the AMT workers in three different HIT designs. HIT1 is a simple categorization scheme in which workers are asked to classify each sentence as being either *positive*, *negative* or *neutral*, as is shown in Figure 1, section b. HIT2 is a graded categorization template in which workers had to assign a score between -5 (negative) and +5 (positive) to each example sentence, as is shown in Figure 1, section c. Finally, HIT3 is a continuous triangular categorization template that allows workers to use both a horizontal positive-negative axis and a vertical subjective-objective axis by placing the example sentence anywhere inside the triangle. The subjective-objective axis expresses the degree to which the sentence contains opinionated content and was earlier used by (Esuli and Sebastiani, 2006). For example, the sentence '*I think this is a wonderful car*' clearly marks an opinion and should be positioned towards the subjective end, while the sentence '*The car has six cilinders*' should be located towards the objective end. Figure 1, section d contains an example of HIT3. In order not to burden the workers with overly complex instructions, we did not mention this subjective-objective axis but asked them instead to place ambiguous sentences towards the center of the horizontal positive-negative axis and more objective, non-opinionated sentences towards the lower *neutral* tip of the triangle.

ID	Country	HIT1	HIT2	HIT3	Acc
A16MC82ITK70QZ	US	77/8.2			x%
A19835WFUL4B52	US	43/5.3			x%
A198YDDSSOBP8A	Mexico	794/11.0			x%
A1COK1GRYUJA1M	US	3/15.7			x%
A1F70TQGR00PTQ	US	980/			x%

Table 1: Statistics on AMT workers: (fictional) ID, Country, per HIT type nr. hits/average completion time, Accuracy.

<Explain three independent annotations. Total of 9000 annotations. Total cost.>

## 4 Annotation Task Results and Analysis

After designing the HITs, we uploaded 30 random samples for testing purposes. These HITs were completed in a matter of seconds, mostly by workers in India. After a brief inspection of the results, it was obvious that most answers corresponded to random clicks. Therefore, we decided to include a small competence test to ensure that future workers would possess the necessary linguistic skills to perform the task. The test consists of six simple categorisation questions of the type of HIT1 that a skilled worker would be able to perform in under a minute.

### 4.1 Annotation Statistics

These are the statistics:

<among non-experts and experts vs. non-experts. (Snow et al., 2008) is good for this.>

### 4.2 Annotation Quality

<This is a very useful reference: (Dawid and Skene, 1979)>  
<This is another useful reference: (Mason and Watts, 2009)>

The annotation quality of AMT workers can be measured by comparing them to expert annotations. This is usually done by calculating inter-annotator agreement (ITA) scores. Note that, since a single HIT can contain more than one assignment and each assignment is typically performed by more than one annotator, we can only calculate ITA scores between batches of assignments, rather than between individual workers. Therefore, we describe the ITA scores

in terms of batches. In Table 4.2, we present a comparison of standard kappa<sup>1</sup> calculations (?) between batches of assignments in HIT1 and expert annotations.

We found an inter-batch ITA score of 0.598, which indicates a moderate agreement due to fairly consistent annotations between workers. When comparing individual batches with expert annotations, we found similar ITA scores, in the range between 0.628 and 0.649. This increase with respect to the inter-batch score suggests a higher variability among AMT workers than between workers and experts. In order to filter out noise in worker annotations, we applied a simple majority voting procedure in which we selected, for each sentence in HIT1, the most voted category. This results in an additional batch of annotations. This batch, referred in Table 4.2 as *Majority*, produced a considerably higher ITA score of 0.716, which confirms the validity of the majority voting scheme to obtain better annotations.

In addition, we calculated ITA scores between three expert annotators on a separate, 500-sentence dataset, randomly selected from the same corpus as described at the start of Section 3. This collection was later used as test set in the experiments described in Section 5. The inter-expert ITA scores on this separate dataset contains values of 0.725 for  $\kappa_1$  and 0.729 for  $\kappa_2$ , only marginally higher than the *Majority* ITA scores. Although we are comparing results on different data sets, these results seem to indicate that multiple AMT annotations are able to produce a similar quality to expert annotations. This might suggest that a further increase in the number of HIT assignments would outperform expert ITA scores, as was previously reported in (Snow et al., 2008).

## 5 Incidence of annotations on supervised polarity classification

```
<Two experiments: (i) AMT
annotations vs. original
Ciao annotations and (ii)
AMT annotations vs. expert
annotations. >
```

<sup>1</sup>In reality, we found that fixed and free margin Kappa values were almost identical, which reflects the balanced distribution of the dataset.

	$\kappa_1$	$\kappa_2$
Inter-batch	0.598	0.598
Batch_1 vs. Expert	0.628	0.628
Batch_2 vs. Expert	0.649	0.649
Batch_3 vs. Expert	0.626	0.626
Majority vs. Expert	0.716	0.716
Experts <sup>2</sup>	0.725	0.729

Table 2: Interannotation Agreement as a measure of quality of the annotations in HIT1.  $\kappa_1$  = Fixed Margin Kappa.  $\kappa_2$  = Free Margin Kappa.

```
<Make it very clear in the intro
that only HIT1, and not HIT2 and
HIT3 are considered for classifier
training>
```

This section intends to evaluate the incidence of AMT-generated annotations on a polarity classification task. According to this, a comparative evaluation between two polarity classification systems is conducted. More specifically, baseline or reference classifiers trained with noisy available metadata are compared with contrastive classifiers trained with AMT generated annotations. Although more sophisticated classification schemas can be conceived for this task, a simple SVM-based binary supervised classification approach is considered here.

### 5.1 Description of datasets

As was mentioned in Section 3, all sentences were extracted from a corpus of user opinions on cars from the automotive section of [www.ciao.es](http://www.ciao.es) (Spanish). For conducting the experimental evaluation, three different datasets were considered:

1. Baseline: constitutes the dataset used for training the baseline or reference classifiers. Automatic annotation for this dataset was obtained by using the following naive approach: those sentences extracted from comments with ratings<sup>3</sup> equal to 5 were assigned to category ‘positive’, those extracted from comments with ratings equal to 3 were assigned to ‘neutral’, and those extracted from comments with ratings equal to 1 were assigned to ‘negative’. This

<sup>3</sup>The corpus at [www.ciao.es](http://www.ciao.es) contains consumer opinions marked with a score between 1 (negative) and 5 (positive).  
<Rafael, please correct.>

dataset contains a total of 5570 sentences, with a vocabulary coverage of 11797 words.

2. Annotated: constitutes the dataset that was manually annotated by AMT workers in HIT1. This dataset is used for training the contrastive classifiers which are to be compared with baseline system. The three independent annotations generated by AMT workers for each sentence within this dataset were consolidated into one unique annotation by majority voting: if the three provided annotations happened to be different<sup>4</sup>, the sentence was assigned to category ‘neutral’; otherwise, the sentence was assigned to the category with at least two annotation agreements. This dataset contains a total of 1000 sentences, with a vocabulary coverage of 3022 words.
3. Evaluation: constitutes the gold standard used for evaluating the performance of classifiers. This dataset was manually annotated by three experts in an independent manner. The gold standard annotation was consolidated by using the same criterion used in the case of the previous dataset<sup>5</sup>. This dataset contains a total of 500 sentences, with a vocabulary coverage of 2004 words.

These three datasets were constructed by randomly extracting sample sentences from an original corpus of over 25000 user comments containing more than 1000000 sentences in total. The sampling was conducted with the following constraints in mind: (i) the three resulting datasets should not overlap, (ii) only sentences containing more than 3 tokens are considered, and (iii) each resulting dataset must be balanced, as much as possible, in terms of the amount of sentences per category. Table 3 presents the distribution of sentences per category for each of the three considered datasets.

## 5.2 Experimental settings

As mentioned above, a simple SVM-based supervised classification approach was considered for the

<sup>4</sup>This kind of total disagreement among annotators occurred only in 13 sentences out of 1000.

<sup>5</sup>In this case, annotator inter-agreement was above 80%, and total disagreement among annotators occurred only in 1 sentence out of 500

	Baseline	Annotated	Evaluation
Positive	1882	341	200
Negative	1876	323	137
Neutral	1812	336	161
Totals	5570	1000	500

Table 3: Sentence-per-category distributions for baseline, annotated and evaluation datasets.

polarity detection task under consideration. According to this, two different groups of classifiers were considered: a baseline or reference group, and a contrastive group. Classifiers within these two groups were trained with data samples extracted from the baseline and annotated datasets, respectively. Within each group of classifiers, three different binary classification subtasks were considered: positive/not\_positive, negative/not\_negative and neutral/not\_neutral. All trained binary classifiers were evaluated by computing precision and recall for each considered category, as well as overall classification accuracy, over the evaluation dataset.

A feature space model representation of the data was constructed by considering the standard bag-of-words approach. In this way, a sparse vector was obtained for each sentence in the datasets. Stop-word removal was not conducted before computing vector models, and standard normalization and TF-IDF weighting schemes were used.

Multiple-fold cross-validation was used in all conducted experiments to tackle with statistical variability of the data. In this sense, twenty independent realizations were actually conducted for each experiment presented and, instead of individual output results, mean values and standard deviations of evaluation metrics are reported.

Each binary classifier realization was trained with a random subsample set of 600 sentences extracted from the training dataset corresponding to the classifier group, i.e. baseline dataset for reference systems, and annotated dataset for contrastive systems. Training subsample sets were always balanced with respect to the original three categories: ‘positive’, ‘negative’ and ‘neutral’.

## 5.3 Results and discussion

Table 4 presents the resulting average values of precision and recall for each considered category in

classifier	baseline	annotated
positive/not_positive	59.63 (3.04)	69.53 (1.70)
negative/not_negative	60.09 (2.90)	63.73 (1.60)
neutral/not_neutral	51.27 (2.49)	62.57 (2.08)

Table 5: Average accuracy (with standard deviations provided in parenthesis) for each classification subtasks trained with either the baseline or the annotated dataset.

classifiers trained with either the baseline or the annotated dataset. As observed in the table, with the exception of recall for category ‘negative’ and precision for category ‘not\_negative’, both metrics are substantially improved when the annotated dataset is used for training the classifiers. The most impressive improvements are observed for ‘neutral’ precision and recall.

Table 5 presents the resulting average values of accuracy for each considered subtask in classifiers trained with either the baseline or the annotated dataset. As observed in the table, all subtasks benefit from using the annotated dataset for training the classifiers; however, it is important to mention that while similar absolute gains are observed for the ‘positive/not\_positive’ and ‘neutral/not\_neutral’ subtasks, this is not the case for the subtask ‘negative/not\_negative’, which actually gains much less than the other two subtasks.

After considering all evaluation metrics, the benefit provided by human-annotated data availability for categories ‘neutral’ and ‘positive’ is evident. However, in the case of category ‘negative’, although some gain is also observed, the benefit of human-annotated data does not seem to be as much as for the two other categories. This, along with the fact that the ‘negative/not\_negative’ subtask is actually the best performing one (in terms of accuracy) when baseline training data is used, might suggest that low rating comments contains a better representation of sentences belonging to category ‘negative’ than medium and high rating comments do with respect to classes ‘neutral’ and ‘positive’.

In any case, this experimental work only verifies the feasibility of constructing training datasets for opinionated content analysis, as well as it provides an approximated idea of costs involved in the generation of this type of resources, by using AMT.

Apart from the xxx, ... explain.

System	Experts	Batch_1	Batch_2	Batch_3	Majority	All
Maxent	59.2	55.8	57.6	54.0	57.6	58.6
C45	42.2	33.6	42.0	41.2	41.6	45.0
Winnow	44.2	43.6	40.4	47.6	46.2	50.6
SVM	57.6	53.0	55.4	54.0	57.2	52.8

Table 6: Accuracy figures of four different classifiers (Maxent, C45, Winnow and SVM) trained on six different datasets: Experts=training set annotated by experts, Batch\_1=first batch of HIT1, Batch\_2=second batch of HIT1, Batch\_3=third batch of HIT1, Majority=batch obtained by majority voting between Batch\_1, Batch\_2 and Batch\_3, All=batch obtained by aggregating Batch\_1, Batch\_2 and Batch\_3.

<extend on cost/benefit.>

## 5.4 Economic impact

In order to decide to use the Amazon’s ”mechanical turk” it’s important to measure the quality of the results, but, of course, the cost of performing the annotation using turkers instead of in-house experts and the time needed to perform the annotation must also be taken into account. To compute the in-House costs, the expert or engineer is paid at a rate of 70\$ an hour, including salary and all other structural costs. The in-house annotation tasks of 1000 sentences would last for three hours, giving a cost of 210\$ and a marginal cost of preparing the data. By means of Amazon we spent 75\$ to annotate the same sentences three times, to get the same quality of results. But we must add the costs of designing the Hits and the qualifying test and the task of upload the data into Amazon’s servers. These tasks highly depend on the experience of the engineers developing it, and can range from a couple of hours to a couple of days (for the first times). As the volume of data to annotate increases, the economical benefit of the AMT is more evident, because the design time becomes a small portion of the budget. For annotating 30,000 sentences the differences rises up to 1350\$ that are more than enough to design the hit. But the main important economic impact may not be the cost but the time to perform the task. In-house massive annotation may take a lot of time, and can

	baseline	baseline	annotated	annotated
category	precision	recall	precision	recall
positive	50.10 (3.79)	62.00 (7.47)	60.21 (2.07)	71.00 (2.18)
not_positive	69.64 (2.70)	58.05 (7.54)	77.95 (1.32)	68.54 (2.75)
negative	35.25 (2.63)	53.46 (10.55)	39.07 (1.78)	55.52 (3.26)
not_negative	78.04 (2.19)	62.62 (6.76)	79.73 (1.10)	66.87 (2.31)
neutral	32.51 (3.02)	48.03 (7.33)	44.72 (2.00)	67.12 (2.96)
not_neutral	68.17 (2.65)	52.81 (3.84)	79.41 (1.58)	60.40 (2.96)

Table 4: Average precision and average recall (with standard deviations provided in parenthesis) for each considered category in classifiers trained with either the baseline or the annotated dataset.

become a hard task, as a few users need to annotate a lot of sentences. As this task can become tiresome and hard it needs to be elapsed in time. In AMT the a hit may be distributed into many volunteers that 24 hours a day are ready to work in parallel to produce the results in a few hours.

## 6 Conclusions

In this paper we have analyzed the use of Amazon’s ”mechanical turk” to annotate polarity of sentences. We have observed that even that the quality of single annotators may be lower than using experts, the price of each annotation is so cheap that allows the use of multiple annotators to reach similar quality still at lower prices. From the different designs of the hits that we did, we can conclude that the complexity of the task was simpler than we thought, (and maybe two cents was a high price) and that what a priory seemed simpler the turkers took more time than the multiple answers task that includes richest data. From a first trial Hit we observed some users answering at random, also the task was in Spanish, and somehow we had to be sure that the turkers would understand the language. We designed a qualifying test (from which we don’t have any information about how many users tried it and failed, we think that could be interesting to have this information). The design of the qualifying test was done using sentences difficult to translate using automatic tools (like Google translate).

Future work: HIT Design: what is the optimal design of the annotation task for our purpose and what is the effect of a suboptimal design on the system

scores? We present AMT workers with three different HIT designs and evaluate the impact of each of them on the overall system performance. Train classifiers with results from HIT2/HIT3. two ways of measuring quality: kappa & performance

## Acknowledgments

We thank Amazon for generously sponsoring this Shared Task.

## References

- A. P. Dawid and A. M. Skene. 1979. Maximum likelihood estimation of observer error-rates using the EM algorithm. *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, 28(1):20–28.
- A. Esuli and F. Sebastiani. 2006. SentiWordNet: a publicly available lexical resource for opinion mining. In *Proceedings of LREC*, volume 6.
- A. Kittur, E. H Chi, and B. Suh. 2008. Crowdsourcing user studies with mechanical turk.
- W. Mason and D. J Watts. 2009. Financial incentives and the performance of crowds. In *Proceedings of the ACM SIGKDD Workshop on Human Computation*, pages 77–85.
- V. S Sheng, F. Provost, and P. G Ipeirotis. 2008. Get another label? improving data quality and data mining using multiple, noisy labelers. In *Proceeding of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 614–622.
- R. Snow, B. O’Connor, D. Jurafsky, and A. Y Ng. 2008. Cheap and fastbut is it good?: evaluating non-expert annotations for natural language tasks. In *Proceedings*

*of the Conference on Empirical Methods in Natural Language Processing*, pages 254–263.

- Q. Su, D. Pavlov, J. H Chow, and W. C Baker. 2007. Internet-scale collection of human-reviewed data. In *Proceedings of the 16th international conference on World Wide Web*, pages 231–240.