
Towards a Shared Rubric for Dataset Annotation

Andrew Marc Greene
Adobe
agreene@adobe.com

Abstract

When arranging for third-party data annotation, it can be hard to compare how well the competing providers apply best practices to create high-quality datasets. This leads to a “race to the bottom,” where competition based solely on price makes it hard for vendors to charge for high-quality annotation. We propose a voluntary rubric which can be used (a) as a scorecard to compare vendors’ offerings, (b) to communicate our expectations of the vendors more clearly and consistently than today, (c) to justify the expense of choosing someone other than the lowest bidder, and (d) to encourage annotation providers to improve their practices.

1 Why we need a rubric

When evaluating data annotation services, the comparison points are often velocity and price. To compete on those terms, some vendors take shortcuts that diminish the value of the data, or achieve lower prices through unethical treatment of the human beings doing the actual annotation. Discussions of the vendor’s processes and their adherence to best practices sometimes appear on their websites or in conversation, but it can be difficult to objectively compare providers. Having a widely used rubric shared by the data community can formalize these definitions, remind us of the hidden layers in data sourcing, enable apples-to-apples comparison, and reduce surprises when the dataset is delivered.

Furthermore, when a project team has selected a vendor who isn’t the lowest bid, they often need to justify their choice to a procurement team. A scorecard that objectively measures each vendor against a set of standard criteria should make it easier to explain why spending additional money is appropriate, in pursuit of essential high-quality data.

Vendors will also benefit from such a rubric, which can help them more effectively explain to new clients the value that they provide. It will reduce the pressure to compete solely on time and money, allowing a richer ecosystem to thrive with different providers offering different tradeoffs among price, schedule, and quality.

This rubric is also useful in assessing one’s own annotation and curation practices.

2 Interpreting the rubric

We propose 6 major categories (ethical treatment of annotators, preparation of an ontology and guidelines, assessing annotation quality, merging of individual annotations, preparing the data for annotation, and data delivery), several of which have subcategories, for a total of 15 specific areas.

Each category presents a summary of practices we have encountered with various vendors, assigned to four levels. (An empty level in a category means we have not yet encountered a practice that we would assign to that level.)


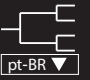

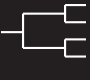
- **Excellent:** This is a best practice and exceeds expectations.
- **Good:** This is adequate for most cases and is the usual expectation.


- **Poor:** Below expectations; a warning sign that the provider may deliver poor-quality data.
- **Unacceptable:** A major deficiency; even one of these usually disqualifies a provider.


This rubric, while shared, is not to be mindlessly applied across all projects. Your team might adjust the “grade” for some items to emphasize what you consider important. We welcome your feedback on what should be added or judged differently; please send email to agreene@adobe.com and check datacuration.org/rubric for the latest version.


Those interested in a detailed explanation of many of these topics will find Monarch (2021) useful.


3 The rubric


 Taxonomy/ Ontology/ Annotation Guide: Versioning	Excel.	Uses semantic versioning for the instructions. (See semvar.org)
	Good	Uses timestamp or other linear version number for the instructions. Maintains a change log.
	Poor	Lacks a formal process for tracking changes and ensuring client agreement on changes to the taxonomy/ontology/instructions.
	Unacc.	Makes unilateral changes; makes it difficult to keep client and provider versions “in sync”.
 Taxonomy/ Ontology/ Annotation Guide: Language	Excel.	Manually translates instructions into the annotators’ native language, when appropriate. The client is given the opportunity to review the translation.
	Good	When annotators are not fluent in the language in which the client has written guidelines, ensures that the instructions are easy to understand, with client collaboration and approval.
	Poor	Does not review instructions for readability by non-fluent speakers of the language in which the guidelines are written, even when that is needed.
	Unacc.	Relies on machine translation of instructions without manual verification.
 Taxonomy/ Ontology/ Annotation Guide: Questions	Excel.	Has a UI in the annotation tool for annotator questions and client responses. Considers annotator feedback important and includes that time in their pay.
	Good	Can collect annotators’ questions in the UI, but responses delivered externally. Provides client with opportunity to test-drive the annotators’ experience.
	Poor	Uses an out-of-context system (e.g., a shared spreadsheet) for annotator queries.
	Unacc.	Lacks any way for annotators to ask questions.
 Taxonomy/ Ontology/ Annotation Guide: Testing and Refinement	Excel.	Has experts in annotation techniques and HCI review the guidelines and suggest improvements to increase accuracy and reduce cognitive load. Coaches underperforming annotators; firing them only as a last resort.
	Good	Tests the instructions for clarity and completeness with small group of annotators before scaling up. Annotators have option to respond “can’t answer”. Provides annotators sufficient time to understand the task. Distinguishes between training material and the taxonomy.
	Poor	Insists that guidelines must cover every eventuality, even though that adds cognitive costs while coverage only asymptotically approaches being complete.
	Unacc.	Rushes to scale up annotation without first collecting data that confirms that the guidelines are clear and are being consistently applied.


 <p>Ethical treatment of annotators: Payment</p>	Excel.	Full-time, paid salary, with benefits.
	Good	Part-time, paid hourly, including time spent in training and on breaks.
	Poor	Paid per annotation, works out to a living wage in the annotator's location.
	Unacc.	Paid per annotation, works out to below living wage in the annotator's location.


 <p>Ethical treatment of annotators: Work conditions</p>	Excel.	Annotation interface is accessible. Coaches underperforming annotators, firing them only as a last resort.
	Good	Provides annotators with appropriate equipment (e.g., large monitors). Provides breaks and variety to avoid fatigue.
	Poor	Neglects ergonomic needs of annotators. Sets unrealistic quotas.
	Unacc.	Lacks appropriate pandemic safety precautions.


 <p>✓ 97±0.4% Assessing annotation quality</p>	Excel.	Provider's quality team manually reviews a random sample of data frequently. Provides dashboard (updated daily) for client to monitor detailed metrics. In multi-stage workflows, annotators can flag bad results from previous stages. Builds models to predict annotation errors based on metadata. Reports accuracy estimates for each task with confidence/credibility intervals.
	Good	Identifies high-risk items for additional manual review. Sets and meets higher quality goals for test data. Provides UI for "acceptance testing" by client on a rapid cadence. Reports accuracy estimates for each task using a statistically valid approach.
	Poor	Uses only simplistic inter-annotator agreement to monitor dataset consistency. Fails to revisit earliest annotations once annotators have gained experience.
	Unacc.	Flawed data is discarded (which can introduce bias) instead of being corrected. Task-level accuracy not reported or lacks an explanation of how it is computed.


 <p>✓ 97% Assessing annotator reliability</p>	Excel.	Uses statistical tests to identify outlier annotators for each question.
	Good	Regularly adds items whose correct answer is known, to monitor annotator quality throughout the project. Uses statistical tests to monitor/identify questions with high disagreement. Uses statistical tests and tracking of IP addresses to identify bots or collusion between supposedly independent annotators.
	Poor	Uses simplistic IAA to identify outlier annotators. Rejects minority opinions out of hand (instead of trying to understand the cause for the disagreement).
	Unacc.	Lacks IAA or statistical monitoring. Fails to exclude or review data previously obtained from annotators who turn out to be unreliable.


 <p>Merging/ adjudication of individual annotations</p>	Excel.	Merging strategy accounts for individual annotators' previous accuracy. (E.g., reweighting of individuals or modeling correlations among annotators.) Annotators can review/visit their work before it is finalized.
	Good	Client can specify merging strategy (e.g., median, or priority voting, etc.) Number of annotators per task clearly specified and statistically justified.
	Poor	Decisions depend only on majority vote of annotators.
	Unacc.	Some decisions are a single annotator's opinion. (Note: When annotator has demonstrated high reliability or is a designated SME, this may be defensible.)

	Data Delivery	Excel.	Provides detailed raw data including annotator ID, date+time of annotation, elapsed time for annotation, annotator's location and/or locale (for modeling sources of error such as unanticipated cultural bias or a poor translation of the instructions), previous versions for this item from this annotator, version number of the instructions under which each datum was collected and annotated (plus merged annotation data).
		Good	Provides merged data and individual responses but with incomplete metadata.
		Poor	Provides merged data plus individual responses but with minimal metadata.
		Unacc.	Provides merged data only. Does not exclude or revisit previously gathered data from annotators who turn out to be unreliable.

	Acquiring Unlabeled Data: Ethics	Excel.	Obtains informed consent from data providers. Offers compensation to content creators when appropriate.
		Good	Relies on sources in the public domain and CC-style licenses.
		Poor	Scrapes the web for publicly visible data, relying on Fair Use carve-outs in copyright law in the relevant geographies.
		Unacc.	Violates copyright law in the relevant geographies. Does not comply with privacy laws such as GDPR.

	Acquiring Unlabeled Data: Bias and Domain Shift	Excel.	Uses ML approaches such as clustering to monitor distribution for societal bias and domain mismatch, and resamples as needed.
		Good	Uses heuristics to monitor distribution for societal bias and domain mismatch, and resamples as needed.
		Poor	Only monitors for domain shift.
		Unacc.	Does not monitor for bias or domain shift.

	Selection Function / Prioritizing Items	Excel.	Provides methods for (a) distribution sampling using provider's embeddings, and (b) uncertainty sampling using provider's off-the-shelf model.
		Good	Provides API to allow client to prioritize data (distribution, uncertainty based on model under development). Provides control over ratio of sampling methods (random, distribution, uncertainty) and allows for that to be updated over time.
		Poor	Does not empower client to control prioritization or selection within the queue.
		Unacc.	Uses filenames or database IDs to control the order of annotation (because these may encode metadata, such as when items are in chronological order, and this may prime or bias the annotators).

	Starting with Seeded Data	Excel.	Compares seeded and unseeded (control-group) annotation tasks to measure impact of anchoring bias. Can seed data via in-house computational approaches if client desires.
		Good	Can seed annotations using data provided by client.
		Poor	
		Unacc.	Automates annotation or seeds data via heuristics, models, or algorithms without client's knowledge and assent.

Rubric version: 0.3.2, last saved 2021-09-30 19:53

References

[1] Monarch, Robert. (2021) *Human-in-the-Loop Machine Learning* Shelter Island: Manning.