

Assignment 6

Title : Perform proper data labelling operation on data set

Theory:

Data labelling

Explore the uses & benefits of data labelling, including different approaches & best practices.

What is data labelling?

Data labelling, or data annotation, is part of preprocessing stage when developing a machine learning (ML) Model it requires the identification of raw data (i.e. images, text files, videos) & then the addition of one or more labels to that data to specify its context for the models, allowing the machine learning model to make accurate predictions.

Data labeling underpins different machine learning & deep learning use cases, including computer vision & natural language processing (NLP).

How does data labelling work?

Companies integrate software, processes & data annotators to clean, structure & label data. This training data becomes the foundation for machine learning models. These labels allows analysts to isolate variables within datasets & this, enables the selection of optimal data predictors for ML models. The labels identify the appropriate data vectors to be pulled in for model training, where the model, then learns to make the best predictions. Along with machine assistance, data labelling tasks require "human-in-the-loop (HITL)" participation HITL leverages the judgment of human "data labelers" toward creating, training, fine tuning & testing ML Models They help guide the data labeling process by feeding the models datasets that are most applicable to a given project.

Labeled data vs unlabeled data

- Labeled data is used in supervised learning, whereas unlabeled data is used in unsupervised learning.
- Labeled data is more difficult to acquire & store (i.e. time consuming & expensive), whereas unlabeled data is easier to acquire & store.
- Labeled data can be used to determine actionable insights, whereas unlabeled data is more limited its usefulness. Unsupervised learning method can help discover new clusters of data, allowing for new categorizations when labelling.

Data labeling approaches:

Data labeling is a critical step in developing a high-performance ML model. Though labeling appears simple, it's not always easy to implement. As a result, companies must consider multiple factors & methods to determine the best approach to labeling. Since each data labeling method has its pros & cons a detailed assessment of task complexity, as well as the size, scope & duration of the project is advised.

Some paths to labelling your data:

- Internal labelling - Using in-house data science experts simplifies tracking, provides greater accuracy & increases quality. However, this approach typically requires more time & favors large companies with extensive resources.
- Synthetic labeling - This approach generates new project data from pre-existing datasets, which enhances data quality & time efficiency. However, synthetic labeling requires extensive computing power, which can increase pricing.
- Programmatic labeling - This automated data labeling process uses scripts to reduce time ~~cost~~ consumption & the need for human annotation. However, the possibility of technical problems requires HITL to remain a part of the quality assurance (QA) process.
- Outsourcing - This can be an optimal choice for high-level temporary projects, but developing & managing a Freelance-oriented workflow can also be time-consuming. Though Freelance platforms provide comprehensive candidate information to ease the

vetting process, hiring managed data labeling teams - provides prevetted staff & pre-built data labeling tools.

- Crowdsourcing - This approach is quicker & more cost-effective due to its micro-tasking capability and web-based distribution. However, worker quality, QA & project management vary across crowdsourcing platforms. One of the most famous examples of crowdsourced data labeling is Recaptcha. This project was twofold in the it controlled for ~~lots~~ lots while simultaneously improving data annotation of images.

Benefits & challenges of data labeling:

The general tradeoff of data labelings is that while it can decrease a business time to scale, it tends to come at a cost. More accurate data generally improves model predictions, so despite its high cost, the value that provides is usually well worth the investment. Since data annotation provides more content to datasets, it enhances the performance of exploratory data analysis as well as machine learning (ML) & artificial intelligence (AI) application. For ex, data labeling produces more relevant search result across search engine platforms & better product recommendations on e-commerce platform lets delve deeper into other key benefits & challenges.

Benefits

Data labeling provides users, teams & companies with greater context, quality & usability.

- More Precise Predictions: accurate data labeling ensures better quality assurance within machine learning algorithms, allowing model to train & yield the expected output. Otherwise, as the old saying goes, "garbage in, garbage out", Properly labeled data provide the "ground truth" (i.e. how labels reflect the 'real world' scenarios) from testing & iterating subsequent model.

Challenges:

Data labeling is not without its challenge. In particular, some of the most common challenge are:

- **Expensive & time consuming** - While data labeling is critical for machine learning models. It can be costly from both a resource & time perspective. If a business takes a more automated approach, engineering teams will still need to set up data pipelines prior to data processing & manual labeling will almost always be expensive & time consuming.
- **Prone to Human - Error**: These labeling approaches are also subject to human-error, which can decrease the quality of data. This, in turn, leads to inaccurate data processing & modeling. Quality assurance checks are essential to maintaining data quality.

Data labeling best practices:

- **Intuitive & streamlined task interfaces** minimize cognitive load and context switching for human labelers.
- **Consensus**: Measures the rate of agreement betⁿ multiple labelers. A consensus score is calculated by dividing the sum of agreeing labels by the total no of labels per annotator.
- **Label auditing**: Verifies the accuracy of labels & updates them as needed.
- **Transfer learning**: Takes one or more pretrained models from one dataset & applies them to another. This can include multi-task learning, in which multiple tasks are learned in tandem.
- **Active learning**: A category of ML algorithms & subset of semi-supervised learning that helps human identify the most appropriate datasets. Active learning approaches include:
 - **Membership query synthesis** - Generates a synthetic instance & require a label for it.
 - **Pool based sampling** - Ranks all unlabeled instances according to informativeness measurement & selects the best queries to annotate.
 - **Stream-based selective sampling** - Selects unlabeled instances one by one & labels or ignores them depending on their informativeness or uncertainty.
- **Natural language processing (NLP)**: A branch of AI that combines computational linguistics with statistical, machine learning & deep learning models to identify & tag important sections of text that generate training data for sentiment analysis, entity, name recognition & optional character recognition. NLP is increasingly being used in enterprise solⁿ like spam detection, machine translation, speech recognition.

text summarization, virtual assistants and chat bots & voice-operated GPS systems.

IBM & data labeling:

- IBM cloud Annotations (link resides outside IBM) -
A collaborative open source image annotation tool that uses AI models to help developers create fully labelled datasets of images, in real time, without manually drawing the labels.
- IBM cloud object storage - encrypted at-rest & accessible from anywhere it stores sensitive data & safeguards data integrity, availability & confidentiality via information Dispersal Algorithm (IDA) & All-or-Nothing Transform (AONT).
- IBM Watson - AI perform with NLP driven tools & services that enables organization to optimize employees time, automate complex business processes & gain critical business insights to predict future outcomes

Conclusion:

Thus, we have studied data labeling operation on dataset.