# Classification of Product Manager Types to Understand the Job Market

**Lena Corredor[1]**    **Saliia Asanova[1]***    **Zikai Liu[1]***
[1]University of California, Berkeley
{lena_corredor, saliia.asanova, zikailiu}@berkeley.edu

## 1  Introduction

Planning for one's career and figuring out the viability of a path is difficult. This is not only the case for students and recent grads, but also the case for experienced and talented individuals. Recently, this is complicated by factors such as the tech recession and a negative macrotrend quoted by McKinsey (2023): The average lifespan of a company listed on the S&P 500 was 61 years in 1958 and is now 18 years. This creates a lot of volatility in the market and thus impacts companies and individuals. Natural Language Processing can play a significant role in helping candidates understand the job market. One task that helps professionals is the classification of jobs into different categories. Since we have knowledge of the product manager (PM) role, our goal is to use NLP techniques to classify PM job postings into technology categories. This helps students and professionals plan for their careers. For example: with this understanding, some may choose to upskill into specific technologies such as AI product manager while others may want to focus on moving into a people manager role. A second task is to calculate the distance between the different types and visualize them in 2D. This helps a working professional understand the relative effort it takes to transition between them and how similar/different between two types.

## 2  Literature Review

Through the applied NLP course, we learnt of important considerations when using NLP models for a specific task. Elements that affect the results are: 1) the data the model is trained on (the type of data, the size of the data), 2) the model selection, and the 3) evaluation process. The task here is to classify a product manager job posting into different categories. We conducted literature review in areas of job classification, taxonomy, prior studies on product management and NLP methods relating to

classification and methods we used. Our contribution is to increase accuracy with annotation of data and evaluate the models selected for this task.

In terms of the NLP method, we tested it with sentenceBERT. The sentence BERT Reimers and Gurevych (2019) is a masked method that uses the Stanford Natural Language Inference Dataset Bowman et al. (2015) and the Multi-Genre NLI dataset Williams et al. (2018). The dataset was manually labeled 'entailment', 'contradiction', or 'neutral'. The researchers Bowman et al. also used GloVe vectors and common crawl. The sentenceBERT derives the semantic meanings using sentence embedding that can be compared using cosine-similarity. To create the labels for product management roles, we leverage our domain knowledge in this area with web search results.

Here are literature reviews relating to the method and domain we are interested in. In the article: "Deep Learning Approaches for Big Data-Driven Metadata Extraction in Online Job Postings Skondras et al. (2023)", the authors used TF-IDF on job descriptions. Then used cosine similarity to measure the differences between job postings. Their model accurately captures both the overarching and detailed nuances of the professional landscape. This is between roles like cook and receptionist whereas we are focusing on categories of product roles. In another study - "An Improved BERT model for precise Job Title Classification Using Job Descriptions Maglyas et al. (2013)", they used BERT for multi-label classification of job titles and achieved a notable performance with an accuracy of 0.842. This paper - Employing Natural Language Processing Techniques for Online Job Vacancies Classification Varelas et al. (2022)" indicated that it's promising. In the article - Machine Learning and Job Posting Classification: A Comparative Study Nasser and Alzaanin (2020), the authors used TF-IDF to extract important words for their model, this contributed to high evalua-

tion metrics 'Accuracy', 'Recall', 'Precision' and F-measure. This prompted us to use TF-IDF as one of our methods. We used logistic regression to weight and do a weight and bias matrix map from the existing embedding to the 9 classes we labeled. This is similar to the methods "Classification and regression tree analysis in public health: Methodological review and comparison with logistic regression Lemon et al. (2003)". We learned from paper "A Comparative analysis of logistic regression, random forest and KNN models for the text classification Shah et al. (2020)" that logistic regression obtained precision of 0.94 compared to K-Nearest Neighbor at 0.96 in the business section, indicating it can be a suitable candidate for classification.

We also learnt about the role in the 2013 paper - What are the roles of software product managers Maglyas et al. (2013). Although this paper is a while back, it does outline the core responsibilities of a PM: 1) identifying features that prove significant value for the customer; 2) work with engineering and other teams to build the product; 3) gain authority and influence for the product. This is considered as a general product manager and it's categorized as "Other Type of Product Management" within the technology dimension. One paper used ISCO taxonomy for job classification. Because of annotating resource limitations, we have to consider other approaches to get a bigger training dataset. In this article "Deep Learning Approaches for Big Data-Driven Metadata Extraction in Online Job Postings Skondras et al. (2023)", the author substantiated the benefits of using synthetic data to enhance job posting classification. This research confirms that AI-generated datasets can enhance the efficacy of NLP algorithms, especially in the domain of multi-class classification job postings. We also need balanced data to increase model effectiveness (Dealing with Imbalanced Data for Fraudulent Job Postings - An Application of BERT and Sampling Techniques".

## 3 Data Collection and Preparation

We obtained the job posting dataset from Indeed using a web scraper tool with the search term 'product manager' as of October 2024. This search also picks up non-PM related roles, but has the terms product and manager in the title. From the initial corpus of 5,000 job postings extracted from Indeed, we manually annotated a subset of 782 postings to ensure precise identification of legitimate Product Manager positions. Instead of manually filtering non-PM roles during preprocessing, we dedicated one classification category called: Non-PM, which is for any role that is not computer product manager type. To enhance data quality, we implemented a heuristic filter to retain only postings with more than 200 words.

We first created the labeling guidelines for 9 categories of product manager by technology. The 9 categories are listed below. Their labeling guidelines are in this spreadsheet due to their length.

- **AI/ML Product Manager**

- **Platform Product Manager**

- **Cloud Product Manager**

- **Other Types of Product Manager**

- **Fintech Product Manager**

- **Data Product Manager**

- **API/Integrations Product Manager**

- **Trust & Safety Product Manager**

- **Non-PM**

To evaluate annotation reliability, we conducted an Inter-Annotator Agreement (IAA) study among three annotators. The analysis was performed on a randomly selected subset of 15 job postings, with each annotator independently classifying the positions into the 9 categories according to our defined taxonomy. With this, we calculated the inter annotator agreement values using Cohen's Kappa. The Cohen's Kappa values we obtained are:

| Annotator Pair | Cohen's Kappa |
|---|---|
| Annotator 1 vs 2 | 0.4599 |
| Annotator 1 vs 3 | 0.6847 |
| Annotator 2 vs 3 | 0.4529 |

Table 1: Cohen's Kappa for Annotator Agreement

Given the moderate Cohen's kappa values obtained from initial annotations, we aimed to move annotator understanding to substantial. Three annotators discussed the differences in understanding and made changes to the labeling guidelines so that we agree on the new classification.

| Cohen's Kappa Range | Agreement |
| --- | --- |
| 0.41 - 0.60 | Moderate |
| 0.61 - 0.80 | Substantial |
| 0.81 - 1.00 | Almost Perfect |

Table 2: Common Interpretations of Cohen's Kappa

Following the refinement of annotation guidelines, we classified a total of 782 job postings across nine categories. The resulting distribution revealed "Other type of PM" and "Non-PM" as the two most prevalent categories, comprising 283 (36.2%) and 231 (29.5%) samples respectively. In contrast, Trust & Safety (12 samples) and API (10 samples) emerged as the least represented ones.

We constructed the training dataset by integrating job titles and descriptions, paired with our defined nine categories as classification labels, forming the basis for our supervised learning task.

## 4 Method

### 4.1 Feature Engineering & Model Selection

For this study, we explored textual feature engineering techniques—specifically TF-IDF and BERT embeddings—and additionally introduced a chain-of-thought (CoT) prompting approach using GPT-4 for comparison and benchmarking.

TF-IDF (Term Frequency-Inverse Document Frequency) Sammut and Webb (2010) weights words based on their frequency within a document (in this case, job descriptions) and rarity across the dataset, enabling identification of role-specific terms (e.g., "learning" for *AI/ML* PMs, "financial" for *Payments* PMs). Our TF-IDF vectorizer implementation extracts 768 features to match BERT's dimensionality, incorporating unigrams and bigrams while filtering terms based on document frequency (appearing in 0.3-90% of documents) and removing English stop words. This feature representation serves as input for a Logistic Regression classifier.

By maintaining consistent dimensionality (768) across both BERT and TF-IDF representations, we establish a controlled framework for comparative analysis, though the fundamental nature of these embeddings differs substantially. While BERT generates dense vectors where most dimensions carry semantic information, TF-IDF produces sparse representations where non-zero values directly correspond to term importance scores Sammut and

Webb (2010). This architectural difference—dense versus sparse embeddings—provides an interesting lens through which to examine their respective strengths in categorizing PM roles. Given TF-IDF's inherent sparsity, we opted against regularization (C=1.0) in the logistic regression classifier to preserve its representation capacity.

To capture deeper contextual nuances, we employed BERT (Bidirectional Encoder Representations from Transformers) Embeddings (BERT-base-uncased) Devlin et al. (2019), a pre-trained transformer model renowned for its ability to generate context-aware embeddings for textual input. In this work, we fine-tuned BERT to encode both job titles and descriptions as a single input sequence for each classification task, leveraging its bidirectional architecture to understand the words in context. The fine-tuned model was then connected to a custom fully connected layer, optimized for multi-class classification. Our empirical analysis demonstrated that [CLS] token representation achieves superior performance compared to mean pooling, aligning with established research on the effectiveness of [CLS] tokens in capturing hierarchical contextual relationships Devlin et al. (2019).

To adapt BERT for our application and maximize its performance, we configured the model to process input sequences with a maximum token length of 512, a batch size of 16 during training, and an epoch size of 7, following the recommendations for small sample sizes as outlined in Sun et al. (2019). During training, we implemented early stopping to preserve the model weights that achieved the highest validation accuracy. This choice balances GPU memory efficiency with the ability to capture meaningful long-term dependencies in job descriptions. However, as many job descriptions exceeded the 512-token limit, truncation of text inevitably occurred, which might have led to the loss of critical information. To address this, we tested two preprocessing techniques: a head-only BERT-Based structure to focus on the title and introductory details, and a dynamic regular expression match model BERT-RegEx to extract the job responsibilities section. Additionally, a dropout layer, inspired by its implementation in Sun et al. (2019), was added during fine-tuning to mitigate overfitting. The effectiveness of these configurations in product management role categorization is empirically evaluated in Section 5.

In addition to the TF-IDF and BERT approaches, we integrated a state-of-the-art Large Language

Model (LLM) with a chain-of-thought (CoT) prompt, specifically employing the GPT-4 CoT Wei et al. (2023). By incorporating an advanced LLM into our experimental framework, we aimed to benchmark our BERT approaches against a State-of-the-Art model thereby gain insights into the relative difficulty of the task. Due to computational constraints, however, we limited the GPT-4 CoT evaluation to a curated subset of 211 samples.

To ensure a robust evaluation of our model's performance, we incorporated two complementary validation techniques: detailed misclassification analysis and bootstrap confidence interval estimation. Misclassification analysis aims to identify instances where the model's predictions deviate from the true labels, providing qualitative insights into common error patterns. To quantitatively assess model performance with statistical rigor, we utilized bootstrap confidence interval (CI) estimation Bestgen (2022) which involves repeatedly resampling the validation dataset with replacement to generate a distribution of accuracy scores (Figure 1). Together, these techniques facilitate a deeper understanding of model behavior and aiding in the comparison of different model configurations.

Further discussions on model limitations, particularly concerning token truncation and category imbalance, are presented in Section 5. Our findings emphasize the importance of tailoring pre-trained transformer models for domain-specific tasks.

We employed stratified sampling to split the data into 80% training and 20% testing sets, ensuring that the class distribution was preserved across both subsets. We utilized a single train-test split supplemented with bootstrap confidence intervals to maintain consistent assessment protocols across all models. Four models were chosen for evaluation:

- **Majority Class Baseline**: This baseline model consistently predicts the most frequent class, the category 'Other type of PM', in the data set which comprises 36.91% of the data.

- **TF-IDF with Logistic Regression**: This approach achieved a training accuracy of 71.36%, but the validation accuracy decreased to 52.23% (95% Confidence interval: [44.59%, 59.87%]).

- **BERT-base Classifier with Initial 512 Tokens**: Utilizing the first 512 tokens (including both the title and truncated job description),

this BERT-Based model attained a training accuracy of 85.81% and a validation accuracy of 61.78% (95% CI: [54.78%, 68.81%]).

- **BERT-RegEx Classifier Enhanced with Regular Expression**: This variant incorporates a regex-based extraction of the "About the Job" section before feeding the data into the BERT model. It achieved both the highest training accuracy 89.92%, and the highest validation accuracy of 64.97% (95% CI: [57.32%, 71.97%]).

- **GPT-4 Chain-of-Thought (CoT) Model**: The state-of-the-art large language model with chain-of-thought reasoning achieved a classification accuracy of 67.30% on a curated subset of 211 samples.

## 4.2 Clustering

Given BERT's demonstrated effectiveness in capturing PM role semantics through our classification results, we leveraged the base BERT model (bert-base-uncased) without fine-tuning to analyze relationships between job categories. For each posting, we obtained embeddings from the [CLS] token using the pre-trained BERT model, concatenating the job title with RegEx-filtered description text (truncated to 512 tokens). Category-level representations were then constructed by averaging these embeddings within each job type (e.g., AI/ML, Cloud, Platform). To facilitate analysis, we reduced the 768-dimensional embeddings to 2D using PCA for visualization, while maintaining the original dimensionality for computing pairwise cosine distances between categories.

## 5 Results

### 5.1 Classification Results

In this study, we employed three distinct models to classify Product Management (PM) job postings into nine predefined categories based on technology requirements: BERT-Based, BERT-RegEx (utilizing regex for token matching), and TF-IDF. Additionally, we benchmarked these models against a GPT-4 CoT (Chain-of-Thought) model. Classification performance was evaluated using accuracy, precision, recall, and F1-score metrics.

The BERT-Based model achieved an overall accuracy of 61.78%, representing a significant improvement over the TF-IDF model's 52.23% accuracy. Specifically, BERT models demonstrated

superior performance in categories such as *Cloud* (F1-score: 0.667 vs. 0.182) and *Payments* (F1-score: 0.588 vs. 0.250), highlighting their ability to better generalize to clearly defined and distinct classes. However, the improvement in accuracy is not statistically significant when considering confidence interval overlap, even with the BERT-RegEx.

The BERT-RegEx model, which integrates regex-based token matching to refine feature extraction, exhibited an overall accuracy of 64.97%, outperforming the BERT-Based model by approximately 3.19%. Notably, the BERT-RegEx model showed substantial category-wise improvements in *AI/ML* (F1-score: 0.800 vs. 0.700) and *Payments* (F1-score: 0.727 vs. 0.588) compared to BERT-Based. These results suggest that for technical domains such as AI and financial services, job description text contains highly discriminative features that are effectively captured through the integration of regex-based pattern matching.

Both BERT variants and TF-IDF model exhibit signs of overfitting, as evidenced by the diverging training and validation loss curves after epoch 3 (Figure 3) and the substantial disparity between training and validation accuracy in TF-IDF.

A notable observation across all models is the zero F1-score for the *API* and *Trust & Safety* categories in Figure 2, indicating a complete failure in identifying these classes. Several hypothesis may explain this. First, the support for these classes is extremely low (only 2 instances each), which likely leads to inadequate learning during training. In addition, the definitions of API may overlap with other categories or lack distinct linguistic markers, as seen in the TF-IDF's feature weights for each category in Table 3, making it challenging for models to differentiate them accurately.

The BERT-RegEx model achieves a validation accuracy of 64.97%, approaching the performance of our state-of-art GPT-4 CoT benchmark (67.30%) for this classification task.

To further understand the performance across models, we conducted a qualitative analysis by examining sample misclassifications for each method. For all models, common misclassification patterns included predicting *Other Type of PM* as *Non-PM* and vice versa (Table 4). This error may arise from the presence of marketing-related (*Non-PM*) project manager roles within the dataset, where job descriptions heavily utilize common product manager keywords. This overlap, although being accurately identified by human labelers, adversely

| Class | Top Features |
|---|---|
| AI/ML | technologies, talent, systems, content, consumer, platform, experiences, world, learning, ai |
| API | location, development, software development, usa, benefits, agile, partners, partner, digital, google |
| Cloud | architecture, network, background, microsoft, technical, aws, security, infrastructure, google, cloud |
| Data | owner, product owner, problems, roadmap, product vision, reporting, capital, aws, analytics, data |
| Non-PM | issues, brand, manufacturing, required, quality, project, support, product marketing, sales, marketing |
| Other Type of PM | people, product owner, owner, ad, market, ideas, lines, mobile, digital, product manager |
| Payments | services, small, businesses, risk, capital, banking, financial, credit, bank, payments |
| Platform | gather, engineering, applications, risk, mission, features, infrastructure, technical, enterprise, platform |
| Trust & Safety | *NA* |

Table 3: TF-IDF Feature Contribution

affects model performance. Effectively mimicking human "attention" mechanisms within language models remains a key area for future work to enhance classification accuracy.

We also tested all 3 methods along with GPT-4 CoT benchmark (Table 5) with a sample Gemini Cloud Assist Senior Product Manager at Google responsible for developing the Gemini (Generative AI) Assist. Human labelers agree as *AI/ML* by identifying critical terms such as "Gemini" and "Generative AI" within job descriptions. In contrast, simple models like TF-IDF and BERT-Based (with a 512-token limit) misclassify these roles as *Cloud* due to the frequent occurrence of the term "Cloud" in the job posting titled "Gemini Cloud Assist Senior Product Manager, Google Cloud,". However, the BERT-RegEx model, utilizing regex to extract the "About the Job" section, successfully captured the essential responsibilities, such as "Develop a Gemini Cloud Assist," and correctly classified the role as *AI/ML*. Similarly, the GPT-4 CoT model employed chain-of-thought reasoning to discern

that despite the cloud-related terminology, the primary focus on AI/ML technologies justified the correct classification. This demonstrates that incorporating domain-specific feature extraction and advanced reasoning mechanisms significantly enhances model performance in accurately classifying nuanced job roles.

| Model | Top Misclassification Patterns |
|---|---|
| BERT-Based | 1. Platform → Other type of PM (11)<br>2. Other type of PM → Non-PM(10)<br>3. Non-PM → Other type of PM (10) |
| BERT-RegEx | 1. Non-PM → Other type of PM (18)<br>2. Platform → Other type of PM (6)<br>3. Other type of PM → Non-PM (3) |
| TFIDF | 1. Non-PM → Other type of PM (14)<br>2. Other type of PM → Non-PM (13)<br>3. Platform → Other type of PM (10) |

Table 4: Top 3 misclassification patterns for each model (True → Predicted). Numbers in parentheses indicate the count of misclassified instances in validation

Future work should focus on addressing data imbalance through augmentation strategies, refining *API* and *Trust & Safety* class definitions, and exploring hybrid models combining rule-based and transformer architectures. Larger batch sizes and increased regularization techniques could mitigate the observed overfitting issues.

| Model | Rationale and Job Snippet |
|---|---|
| BERT-Based (Cloud) | Focused on the keyword "Cloud" present in the job title: *"Gemini Cloud Assist Senior Product Manager, Google Cloud"* |
| BERT-RegEx (AI/ML) | Utilized regex to extract job responsibilities: *"Develop a Gemini Cloud Assist"* |
| TF-IDF (Cloud) | Misled by frequent "Cloud" term without capturing context-specific keywords like "Generative AI" |
| GPT-4-CoT (AI/ML) | Applied chain-of-thought reasoning to identify key terms: *"Generative AI or Large Language Models"* |

Table 5: Comparative Qualitative Analysis of Model Classifications for an AI/ML Role

## 5.2 Clustering Results

The PCA dimensionality reduction effectively visualizes the embedding space in 2D, capturing 62.90% of the total variance despite some information loss. Our clustering analysis in Figure 4 and Table 6 reveals that technical PM roles show

notable differences, with *AI/ML* positions being particularly distinct from other categories (cosine distance of 0.013-0.015 from the non-PM roles ). This distinctiveness aligns with our model's superior classification performance for *AI/ML* positions. In contrast, traditional PM roles show high semantic similarity, as evidenced by the minimal distances between *Other type of PM*, *Data*, and *Platform* categories (cosine distances of 0.002-0.004). These findings suggest that while core PM competencies remain consistent across most domains, *AI/ML* product management has evolved into a specialized field with unique requirements. For job seekers, this implies that transitioning between traditional PM roles may require less role-specific training compared to moving into *AI/ML* product management.

## 6 Conclusion

In this study, we developed and evaluated several approaches for classifying product manager job postings across nine technology categories, with our BERT-RegEx model achieving comparable performance to the GPT-4 CoT benchmark. Our analysis revealed that while traditional PM roles share substantial semantic overlap, AI/ML product management positions exhibit distinct characteristics, suggesting specialized skill requirements for these roles. While our models showed promising results for well-defined categories like AI/ML and Cloud, they struggled with under-represented categories such as API and Trust & Safety, highlighting opportunities for future improvements through data augmentation and refined class definitions. These findings provide valuable insights for both job seekers planning career transitions and organizations designing PM role specifications.

| Category Pair | Distance |
|---|---|
| *Most Dissimilar Pairs* | |
| Non-PM vs Cloud | 0.017 |
| Cloud vs Payments | 0.015 |
| Non-PM vs API | 0.015 |
| Non-PM vs AI/ML | 0.013 |
| AI/ML vs Payments | 0.012 |
| *Most Similar Pairs* | |
| Other type of PM vs Data | 0.002 |
| Non-PM vs Other type of PM | 0.002 |
| Other type of PM vs Platform | 0.004 |
| Platform vs Data | 0.004 |
| Other type of PM vs Payments | 0.004 |

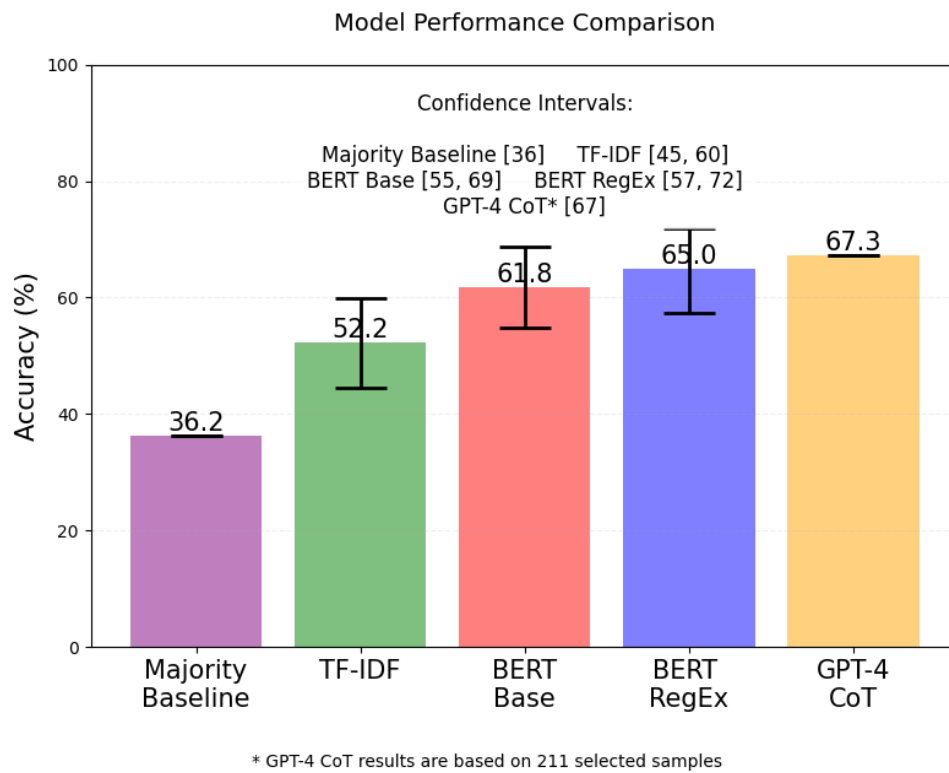Table 6: Pairwise Cosine Distances Between Category BERT Embeddings

Figure 1: Model Results
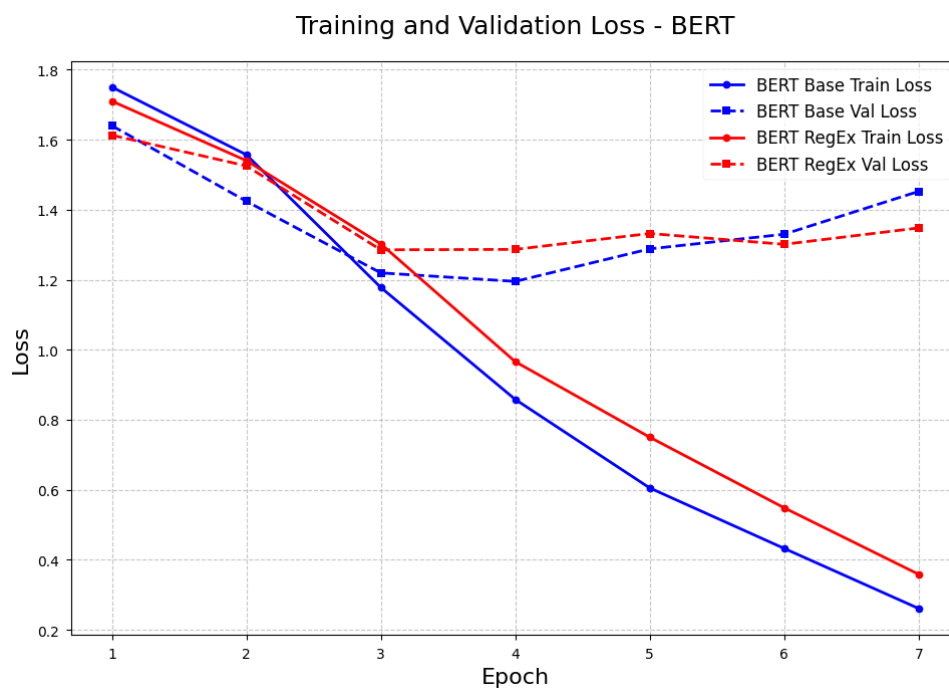


Figure 2: Model Results by Category
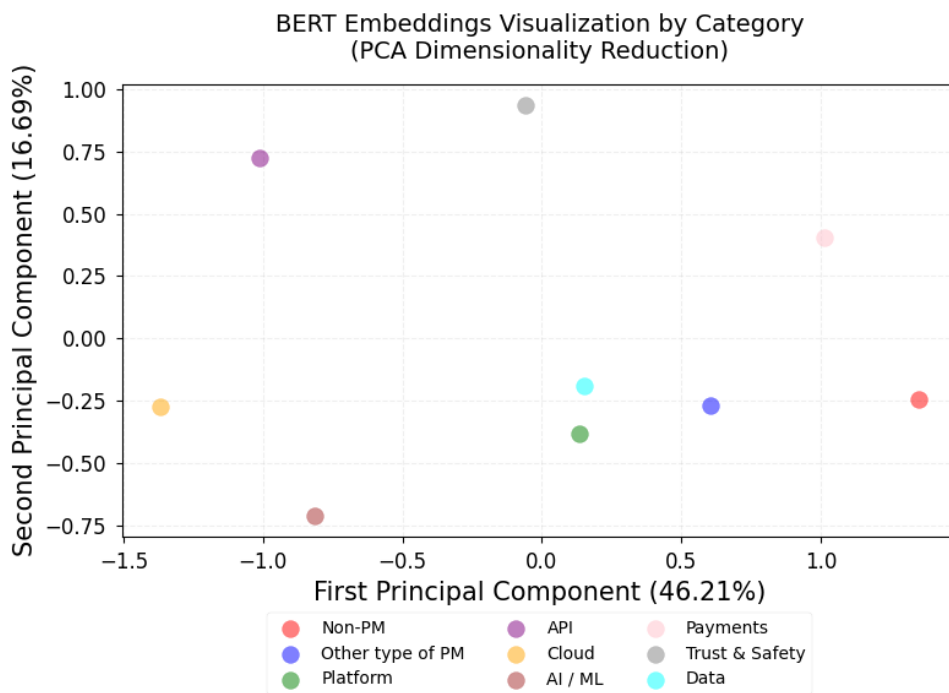
Figure 3: BERT Models Loss Curve



Figure 4: Embeddings by Category in 2D

# References

Yves Bestgen. 2022. Please, don't forget the difference and the confidence interval when seeking for the state-of-the-art status. *arXiv preprint arXiv:2205.11134*.

Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. A large annotated corpus for learning natural language inference. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 632–642, Lisbon, Portugal. Association for Computational Linguistics.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. *Preprint*, arXiv:1810.04805.

Stephen C. Lemon, Jason Roy, Melissa A. Clark, Peter D. Friedmann, and William Rakowski. 2003. Classification and regression tree analysis in public health: Methodological review and comparison with logistic regression. *Annals of Behavioral Medicine*, 26(3):172–181.

Andrey Maglyas, Uolevi Nikula, and Kari Smolander. 2013. What are the roles of software product managers? an empirical investigation. *Journal of Systems and Software*, 86(12):3071–3090.

McKinsey. 2023. Mckinsey on finance. *McKinsey on Finance*. Available online at McKinsey.com.

Ibrahim Nasser and Amjad H. Alzaanin. 2020. Machine learning and job posting classification: A comparative study. *International Journal of Engineering and Information Systems (IJEAIS)*, 4(9):6–14. Available at SSRN: `https://ssrn.com/abstract=3723037`.

Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. *Preprint*, arXiv:1908.10084.

Claude Sammut and Geoffrey I. Webb, editors. 2010. *TF–IDF*, pages 986–987. Springer US, Boston, MA.

Krunal Shah, Harshil Patel, Dhaval Sanghvi, and Maitri Shah. 2020. A comparative analysis of logistic regression, random forest and knn models for the text classification. *Augmented Human Research*, 5:12.

Panagiotis Skondras, Nikos Zotos, Dimitris Lagios, Panagiotis Zervas, Konstantinos C. Giotopoulos, and Giannis Tzimas. 2023. Deep learning approaches for big data-driven metadata extraction in online job postings. *Information*, 14(11).

Chi Sun, Xipeng Qiu, Yige Xu, and Xuanjing Huang. 2019. How to fine-tune BERT for text classification? *CoRR*, abs/1905.05583.

George Varelas, Dimitris Lagios, Spyros Ntouroukis, Panagiotis Zervas, Kenia Parsons, and Giannis Tzimas. 2022. Employing natural language processing techniques for online job vacancies classification. In *Artificial Intelligence Applications and Innovations. AIAI 2022 IFIP WG 12.5 International Workshops*, pages 333–344, Cham. Springer International Publishing.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. Chain-of-thought prompting elicits reasoning in large language models.

Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1112–1122, New Orleans, Louisiana. Association for Computational Linguistics.