

SAMAR SRIVASTAVA

☎ +91-70600 04225 • Github • ✉ samarsrivastava44@gmail.com • LinkedIn • Gurugram, India

EDUCATION

Dr. APJ Abdul Kalam Technical University • B.Tech
Computer Science

August 2015 – May 2019

TECHNICAL SKILLS

Programming Languages • Tools

Python, Elastic Search, Kibana, Apache AirFlow, OOP, Bash, Git

Packages • Frameworks

sklearn, pandas, folium, seaborn, nltk, spacy, Flask, FastAPI, Selenium, BeautifulSoup, spacy, nltk

Cloud • DevOps

AWS, Docker, Azure

Domain Expertise • Domain Knowledge

NLP, Machine Learning, Data Scraping, Text mining, Text analytics, data cleaning, data pre-processing, data visualizations, Statistics, ChatGPT, Rest APIs

WORK EXPERIENCE

Data Scientist – WiseStep
Gurugram, India

November 2020 – Present

- Building LLM based features using ChatGPT to help recruiters in parse and screen resumes, & generate job descriptions.
- Created an industry classification model to improve job-candidate matching by categorizing job descriptions and candidate work experiences.
- Stabilized in-house resume parser which resulted in a 140% increase in CV processing speeds from 50 CVs per minute to 120 CVs per minute. CV parser now runs at a speed of 1 million+ CVs monthly. It predicts and extracts 70+ attributes from a resume like the industry of the candidate, education gaps, experience gaps, finding names from text via NER, pretrained language models etc.
- Increased name & location identification accuracy by 35% using zero shot classification, & NER models.
- Orchestrated strategy and logic for filtering out unique and duplicate active consultants across sources like emails, and job boards, marketed by different recruiters, using Levenshtein distance across multiple fields. This resulted in reducing duplicate candidates on the platform and achieving a higher degree of uniqueness, as a result, a recruiter can reach out to more unique candidates.
- Responsible for developing job description parser (JD Parser) using nltk to scan through large amount of raw documents (emails) and extract information pertaining to job title, company name, compensations, locations, industry. JD Parser helped end users by attenuating job creation process by 70%.
- Day to day task involves data collection to Elasticsearch indices, EDA for job titles, company names normalizing them and performing pre-processing to build knowledge base for data products.
- Developed Airflow DAGs to perform repetitive task like index cleaning, transportation pipelines to move data from ES to redis. This ended up boosting teams productivity by a significant factor.
- Generate and manage data dashboards for the Customer Success team and stakeholders.
- **Tech Stack** - python, flask, elasticsearch, kibana, nltk, spacy, transformers, regex, airflow, AWS, chatGPT, LLMs

Machine Learning Engineer – Scanta Inc.
Gurugram, India

April 2019 – October 2020

- Worked on data dashboard generation by analysing virtual assistants requests and response to detect anomaly in conversations and report malicious events.
- Setup the pipeline for basic NLP preprocessing like text cleaning, tokenization, generating bag of words, evaluating n-grams. Evaluate similarity between text using cosine similarity, and much more.
- Used T-SQL Server as primary database and redis as in secondary database for super-fast data fetching.
- Experimented with Uber AI's Plug & Play Language Model to induce personalities in text.
- Responsible for development of a paraphrasing tool using Transformers.
- Deployed 3 products on AWS using various services like EC2, API gateway, and AWS Lambda.
- Responsible for end to end engineering on NLP products pipelines from data mining, data cleaning, to modelling and deployment on cloud.
- **Major tech stack** - python3, huggingface tokenizers, nltk, sklearn, Docker, T-SQL, HTML, CSS, JS, spacy

PERSONAL PROJECTS (AVAILABLE ON GITHUB)

Classification of Business Licence Status Kaggle

2020

Source Code | Tech Stack - python, data science, nltk, pandas, regression, ML

Predicting Stack Overflow Tags Kaggle

2020

Source Code | Tech Stack - python, data science, nltk, pandas, regression, ML