

m.8.Exercises.an.templates → Navigating Data Employment

Overview	Use Apache Spark and machine learning to determine sentence authorship labels.
Data	https://www.kaggle.com/competitions/spooky-author-identification/code

Exercise -

Jupyter Not <decriptio

Dataset Description: The spooky author

Data Professional Skills by Skill Domain Skills Ontology Self-Assessment Template

ID	Focus & Medium	YesNo (Yn)
1	Data Professional Skills by Skill Domain	dfdf
2	└ Education	
3	└ Advanced degree in a quantitative discipline	
4	└ Mathematics, Linguistics, Computer Science	
5	└ Enrolled in an M.S./Ph.D. program in Comp. Science or Elect. Engineer	
6		
7	└ Experience	
8	└ Industry or academic experience in applied NLP - 2+ years	
9	└ Research experience in fields such as machine learning, languages	
10	└ program synthesis, software eng., or human-computer interaction	
11	└ Research or practical experience in applying deep learning	
12	└ on large-scale and real-world data - 3+years	
13		
14	└ Programming and Technical Skills	
15	└ Familiarity with OCR libraries like Tesseract, PyOCR, OpenCV, .NET, SDK	
16	└ Extracting, cleaning, and preprocessing data sets using NumPy and Pandas	
17	└ Knowledge of supervised and unsupervised machine learning techniques	
18	└ regression models, decision tree models, clustering, deep learning	
19	└ with tools like Scikit-learn, Tensorflow, Keras, or PyTorch	
20	└ Data visualization skills using tools such as Matplotlib, Tableau, etc	
21	└ Familiarity with rule-based NLP like CFG, constituency, and parsing	
22	└ and related libraries including NLTK, spaCy, Stanford NLP	
23	└ Specialization in OCR and familiarity with Transformers, ELMo, and BERT	
24	└ Experience with Python NLP packages like Spacy, NLTK, and	
25	└ Statistical packages familiarity like R, Python, SPSS, SAS, STATA	
26	└ Experience with deep learning techniques and publishing in related	

27		└ conferences (ICML, CVPR, NeurIPS)	
28		└ Handling and analyzing data at scale w Hadoop, Dask, Spark, MapReduce	
29		└ Working knowledge of data store tools like SQL, Elasticsearch	
30			
31		└ Analytical and Problem-Solving Skills	
32		└ Proficiency in quantitative and qualitative analytical techniques rooted	
33		└ in business, economic, and statistical analysis	
34		└ Ability to perform business analysis of market competitiveness,	
35		└ financial analysis, social media monitoring	
36		└ Expertise in statistical analysis (linear regression, logistic regression,	
37		└ nonparametric statistics, probabilistic modeling, spatial modeling	
38		└ Ability to tell stories using data	
39		└ Strong problem-solving abilities	
40			
41		└ Additional Skills and Preferences	
42		└ Knowledge of healthcare industry practices and medical coding (a plus)	
43		└ Experience with computational imaging, cyber security, dist systems,	
44		└ logistics, next-generation networking, quantum information processing,	
45		└ sensor systems, speech and language processing, etc.	
46		└ Security Clearance (for specific positions)	
47		└ Experience managing, coding, and analyzing qualitative data using	
48		└ content analysis software	
49		└ Time series analysis expertise (Prophet, ARIMA, LSTMs)	
50		└ Writing maintainable, testable, production-grade Python code	
51		└ Understanding of different machine learning and deep learning algorithm	
52		└ families and their tradeoffs	
53		└ Experience with Selenium and SeleniumGrid	
54		└ Data analytics, data mining, or other data science skills	
55		└ Database experience, preferably working with Mongo databases	
56		└ Experience working with data in Information Security, Cybersecurity,	
57		└ or Threat Intelligence	
58		└ Experience working with bulletin boards and forums	
59			
60			

Dataset Description: The spooky author identification dataset contains text from works of fiction written by spooky authors of the public domain: Edgar Allan Poe, HP Lovecraft and Mary

TASK SUMMARY

Jupyter Not <decscription here>

- **a.:** Jupyter Notebook
- **B:** Jupyter Notebo