

Data dictionary (confidential)

Prepared by Daubry Byagogo, Black Saber Software

Current employees dataset

Path: "data/black-saber-current-employees.csv"

This dataset contains data on all *current* employees for the whole duration of their employment. This dataset has been anonymised. Each row represents salary, demonstrated leadership and productivity for an employee in a given quarter. The le

Variable	Description
<code>employee_id</code>	5 digit unique identifier for each employee
<code>gender</code>	Gender of employee: 'Man', 'Woman', 'Prefer not to say'
<code>team</code>	Which one of the 8 teams the employee works for
<code>financial_q</code>	Financial quarter salary, leadership and productivity
<code>role_seniority</code>	Least senior to most senior: "Entry-level", "Junior I", "Junior II", "Senior I", "Senior II", "Senior III", "Manager", "Director", "Vice president"
<code>leadership_for_level</code>	Quality of demonstrated leadership, taking into account role level (i.e. "Appropriate for level" requires much less for entry-level employees than for a manager)
<code>productivity</code>	work output in relation to job description, rated on a 0-100 scale with 50 being satisfactory and above 50 indicating better than expected productivity
<code>salary</code>	Salary at at the given financial quarter (note: these are effective yearly values for the current wage, but don't take in to account previous salary steps in the same year, etc.)

Hiring data sets

Black Saber has been trialling a new AI recruitment pipeline manager for the Data and Software teams. There are three phases, outlined below, each narrowing down the field of applicants. Based on advice from our legal team we are not able to provide the original application data, be we can provide these anonymised indicators/ratings from each phase. `applicant_id` is consistent across phases.

		Data collected
Phase 1	Initial application	Team applied for, Cover letter, CV, GPA, Gender, Extracurriculars, Internship experience,
Phase 2	Technical task, writing sample, pre-recorded video	Technical skills, Writing skills, Leadership presence, Speaking skills
Phase 3	Final interview	Interviewer 1 rating Interviewer 2 rating

Phase 1

Path: "data/phase1-new-grad-applicants-2020.csv"

In the first phase of the hiring pipeline applicants complete a form and are asked to submit a CV and cover letter. Extracurriculars and internship experience are autorated based on the descriptions applicants provide in the application form.

Variable	Description
<code>applicant_id</code>	A unique ID assigned to applicants in Phase 1
<code>team_applied_for</code>	Software or Data
<code>cover_letter</code>	0 if absent, 1 if present
<code>cv</code>	0 if absent, 1 if present
<code>gpa</code>	0.0 to 4.0
<code>gender</code>	Gender of employee: 'Man', 'Woman', 'Prefer not to say' only options provided
<code>extracurriculars</code>	The description of extracurricular involvement is assessed against a proprietary key term and phrase bank and given a 0, 1 or 2 for where 2 indicates several high relevance and/or skills building extracurriculars, 1 indicates some relevant and/or skills building extracurriculars and 0 indicates no extracurriculars describes or that those describe were not rated as high relevance or high skills building
<code>work_experience</code>	Similar to <code>extracurriculars</code> , the description applicants provided is assessed against a proprietary key term and phrase bank, that also consideres company names and reputations, to give a 0, 1 or 2 score, with 2 being the best, 0 the worst

Phase 2

Path: "data/phase2-new-grad-applicants-2020.csv"

We don't know exactly how these are being assessed by the AI, the algorithm is obviously commercially sensitive but their demonstrations of the system were impressive.

Variable	Description
<code>applicant_id</code>	A unique ID assigned to applicants in Phase 1
<code>technical_skills</code>	Score from 0 to 100 on a timed technical task, AI autograded
<code>writing_skills</code>	Score from 0 to 100 on a timed writing task, AI autograded
<code>speaking_skills</code>	A rating of speaking ability based on pre-recorded video, AI autograded
<code>leadership_presence</code>	A rating of 'leadership presence' based on pre-recorded video, AI autograded

Phase 3

Path: "data/phase3-new-grad-applicants-2020.csv"

This is the interview phase. Being listed as 'first' or 'second' interviewer is arbitrary and who the interviewers were is not available from our tracking system.

Variable	Description
<code>applicant_id</code>	A unique ID assigned to applicants in Phase 1
<code>interviewer_rating_1</code>	The overall rating of job fit given by the first interviewer on a scale of 0 to 100
<code>interviewer_rating_2</code>	The overall rating of job fit given by the second interviewer on a scale of 0 to 100

Final hires

Path: "data/final-hires-newgrad_2020.csv"

This data set contains the applicant IDs of everyone who was sent an offer letter. In this cohort, everyone accepted.

Variable	Description
<code>applicant_id</code>	A unique ID assigned to applicants in Phase 1