

HDSC Summer'22 Capstone Project Presentation: Gender Statistics – Earnings of Male and Female Employees

A Project by Team Data Warehouse

INTRODUCTION

The Australian labor market is highly gender-segregated by industry and occupation, a pattern that has persisted over the past two decades. The gender wage gap is an indicator of women's earnings compared to men and it is derived by dividing the average annual earnings for women by the average annual earnings for men. According to the international labor organization (ILO), on average, women globally are paid 20% less than men. The ILO further attributes discrimination based on gender as the largest contributory factor to the pay gap. Data Science techniques can be used in an organization to identify wage disparities between male and female employees, hence serving as an inference for companies to examine their hiring practices and policies and resolve gender pay gaps.

PROBLEM STATEMENT

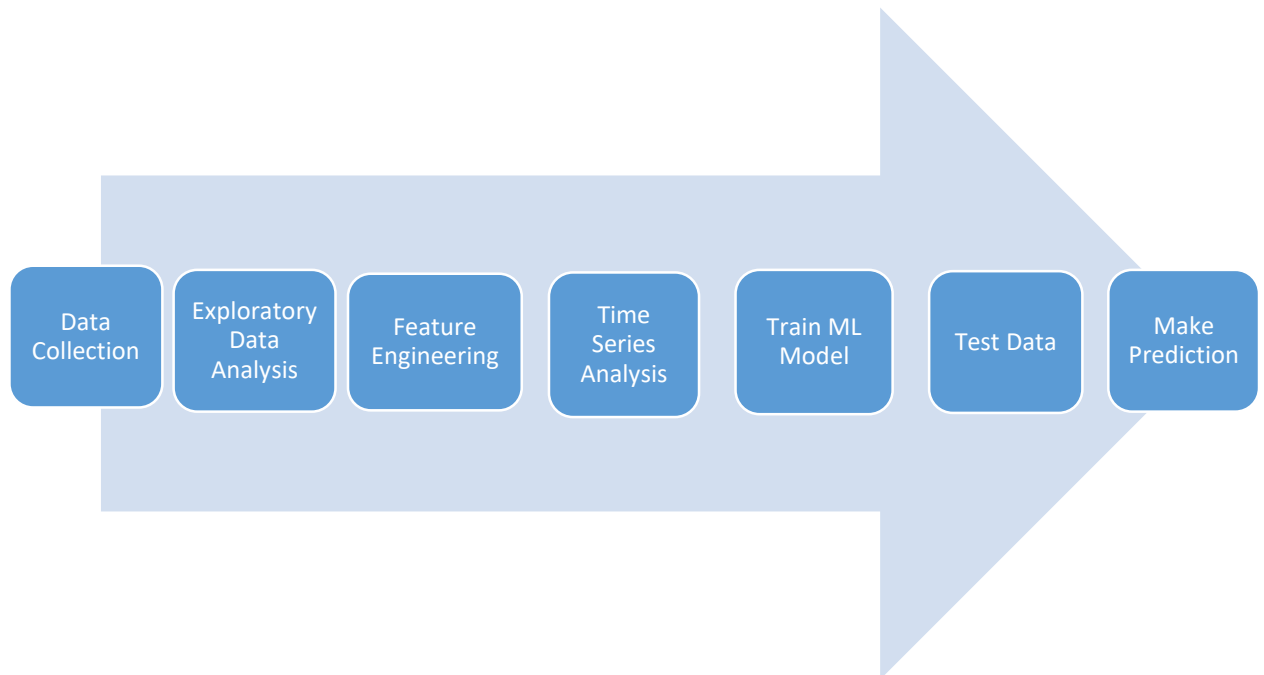
The Workplace Gender Equality Agency (WGEA) estimates that Australian women earn, on average, 15.3% less than men. That is \$1,387 for women per week compared with \$1,638 for men. To bridge this gap, the Australian government has committed huge funds to sponsor more women across different fields in recent years, but despite this sponsorship, female average earnings are still significantly less than their male counterparts. The Data Warehouse team of the Hamoye 2022 Internship seeks to explore and gain significant insights from this disparity.

AIMS AND OBJECTIVES

The aim of this project is to

1. Explore and analyze the gender pay disparity of male and female employees in different job roles between 2004 to 2017 in Australia
2. Build a machine learning model that predicts the future wage gap and its impact on the female workforce over the next 4 years.

FLOW PROCESS



DATA COLLECTION

The data used in this analysis was collected and downloaded from Kaggle at

<https://www.kaggle.com/datasets/mpwolke/cusersmarildownloadsearningscsv>

DATA PRE-PROCESSING

The data set used in the analysis is a CSV file containing the average hourly earnings of females, males and person's, with 14 rows and 28 columns representing each data from 2004 to 2017.

```
(In [9]): df.head()
```

```
(Out[9]):
```

	year	Gender	FemaleNonmanager	FemaleProfessional	FemaleService	FemaleTotal
0	2004	Female	25.14	29.03	17.90	16.95
1	2005	Female	25.10	30.00	18.00	16.20
2	2006	Female	25.50	30.50	18.84	16.94
3	2007	Female	26.43	31.33	20.00	19.85
4	2008	Female	32.00	32.03	20.30	20.97

5 rows x 7 columns

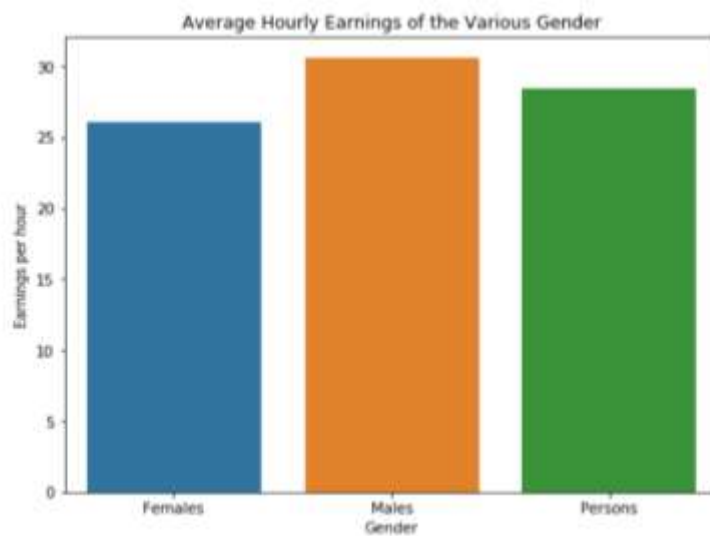
The data was cleaned and manipulated by renaming its columns with a proper naming convention and was found to have no missing nor duplicated values.

Out[12]:

	job_title	females	female_managers	female_professionals	female_technicians_and_trades_workers	female_community_and_personal_service_workers	female_
0	2004 Females		26.14	29.02	17.99		18.01
1	2005 Females		28.10	30.00	18.86		18.20
2	2006 Females		28.60	30.56	19.64		18.84
3	2007 Females		36.43	31.93	20.98		19.85
4	2008 Females		32.68	32.93	20.30		20.97
5	2009 Females		33.30	34.72	20.19		21.00
6	2010 Females		34.41	35.52	22.86		21.97
7	2011 Females		35.22	36.48	23.19		23.25
8	2012 Females		36.24	37.72	23.95		23.52
9	2013 Females		37.26	38.95	23.90		24.50
10	2014 Females		42.63	44.13	27.86		25.01
11	2015 Females		42.12	44.75	27.50		26.87
12	2016 Females		44.35	45.17	26.28		26.54
13	2017 Females		46.02	47.04	28.30		27.84

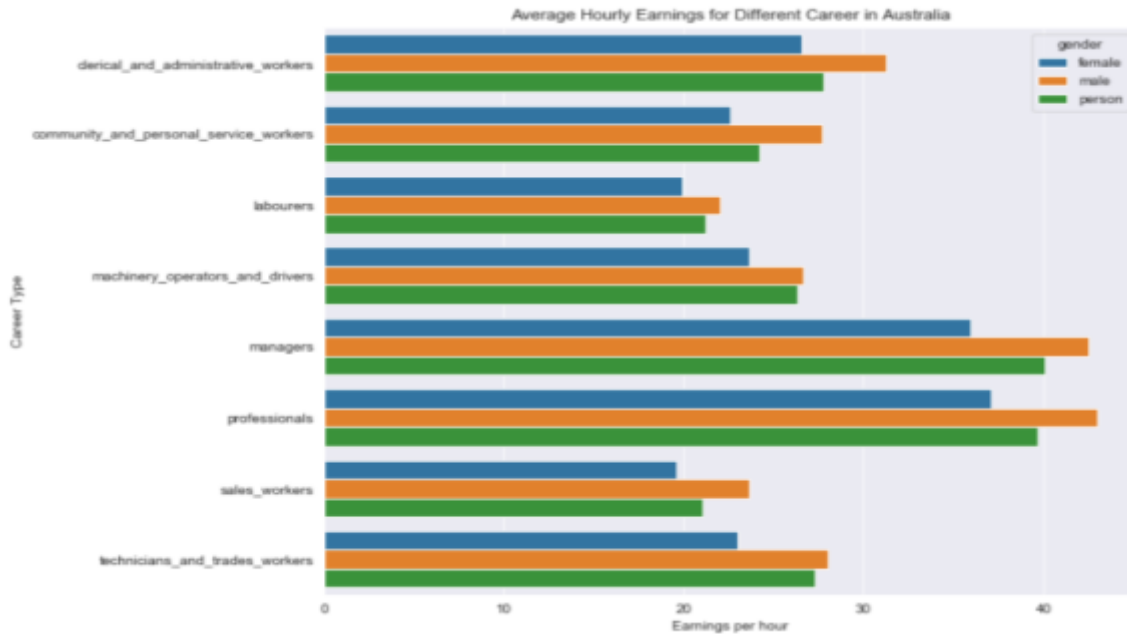
EXPLORATORY DATA ANALYSIS

Exploratory data analysis was carried out to understand the trend of female and male earnings over the years and reveal the job positions with the highest male-to-female earning disparities.



A bar chart of Hourly earnings vs Gender

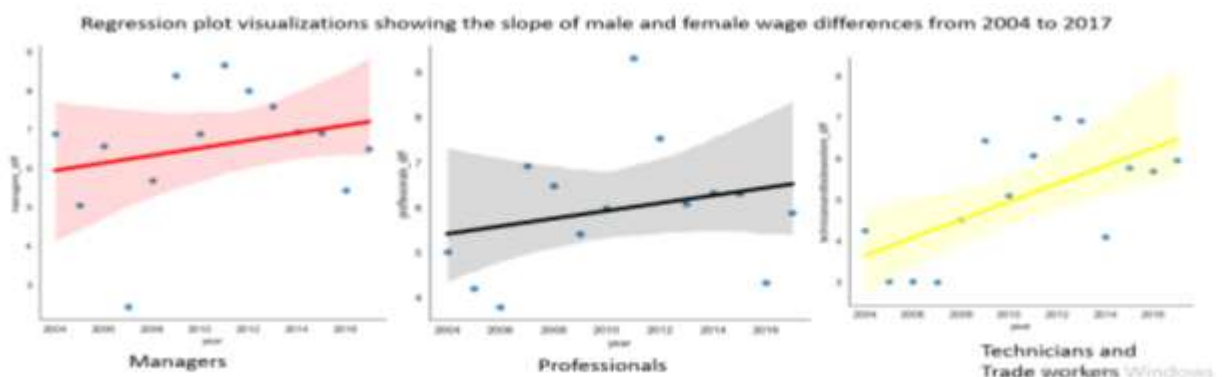
The figure above shows the average hourly earnings of the eight cadres over the total period of analysis (2004 - 2017) for the various gender. The male gender receives the highest hourly wage of 30.59 AUD while the female gender receives the least of 26.05 AUD. Generally, hourly earnings increased regardless of gender or the type of job over the years.



The Plot above shows how hourly wages differ across the various careers and it further reveals that the male gender has the highest hourly pay while the female has the least. Also, the difference between male and the other genders tend to increase for careers with more hourly pay. It is observed that managers and professionals are the most paid job titles.

REGRESSION ANALYSIS

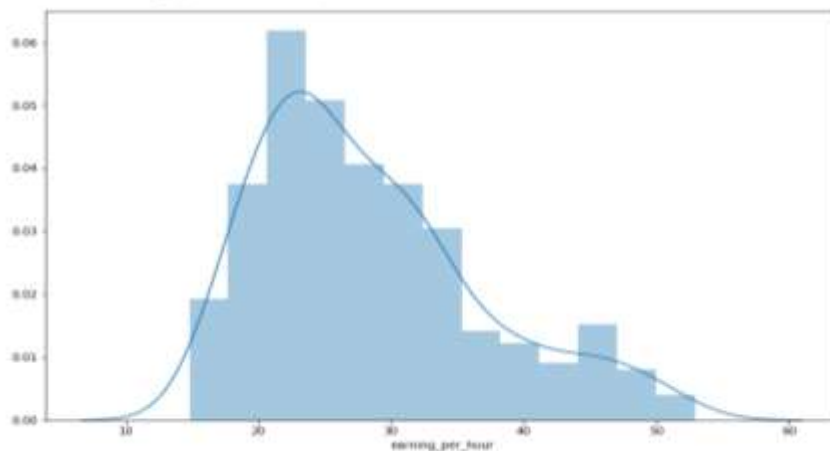
The male and female data were extracted from the data to avoid dropping the 'persons' data. The person's data stems from people who did not specify their gender: including them will be irrelevant to our analysis. The average hourly earnings per year for both male and female employees regardless of their job cadre was used to form two new columns: "male_hourly_data" and "female_hourly_data", in addition to the "year" column. Our target variable (wage_gap) is the difference between the male and female hourly earnings per year.



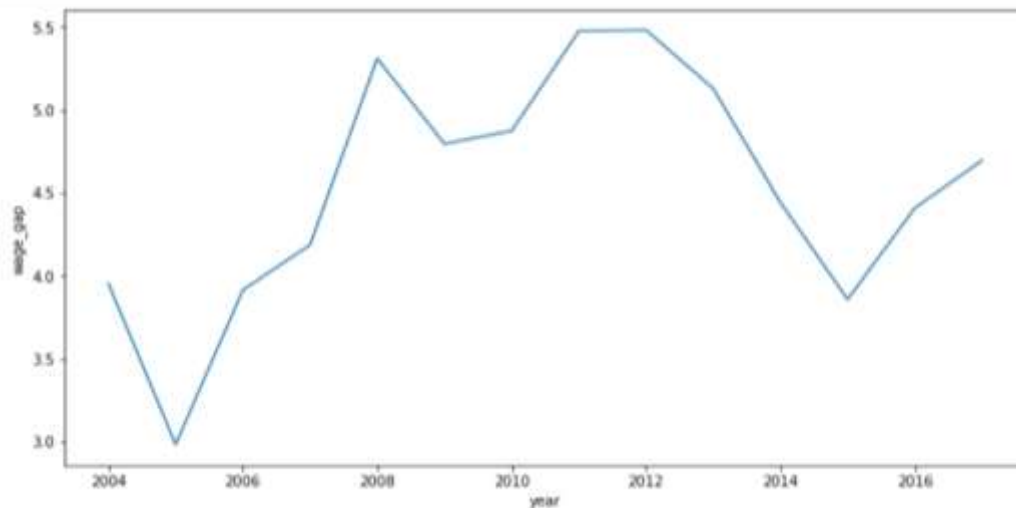
The plots above show the relationship between male and female earnings with respect to job titles. In the relationship, the difference for most jobs is positive meaning that the male earns more than the female counterpart.

FEATURE ENGINEERING

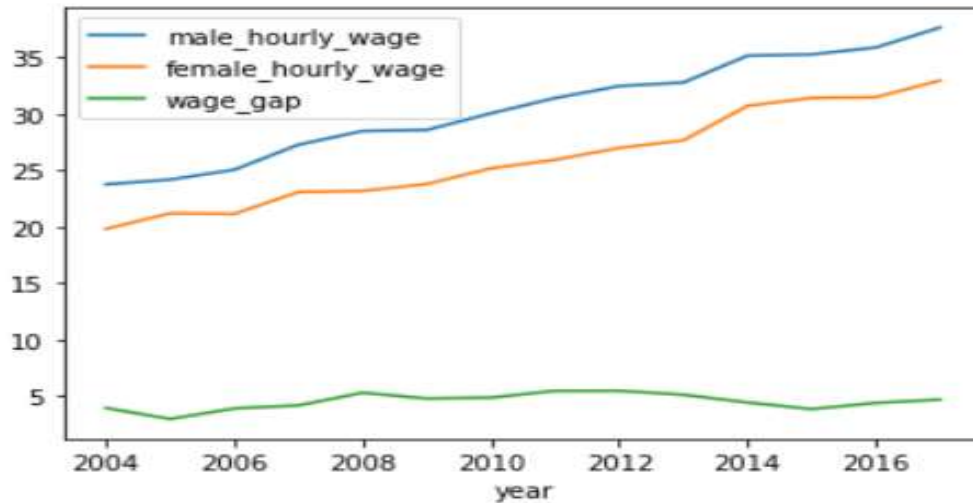
The target variable (wage gap) was created from the difference between the male and female earnings for each year in the dataset, using feature engineering. This was done by querying out the male and female data and then creating an average hourly wage each year for both males and females. Next, the hourly wage gap between males and females was created. The wage gap was then plotted against each year to achieve the histogram below.



From the visualization above, it can be seen that overall, irrespective of cadre, most employees in Australia earn an average hourly wage of about 25 AUD.



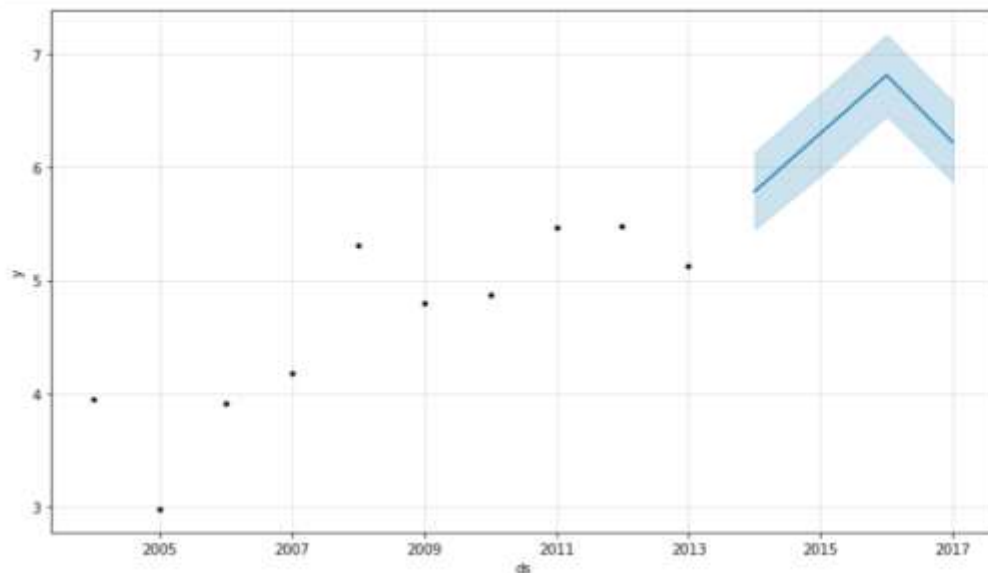
A line graph of wage gap versus year.

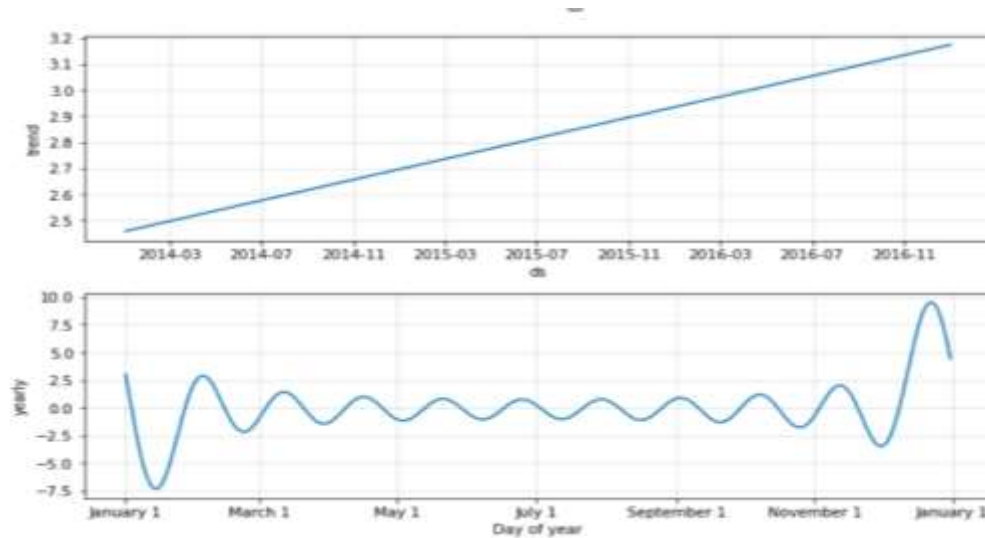


We realize the wage gap has been on a steady trend around 5 dollars in hourly wage. We can easily infer that this data doesn't have seasonality. The wage gap is independent of time.

MODEL TRAINING AND EVALUATION

The model used a time Series analysis and modeling (The Box - Jenkins Method) to predict the future wage gap. The analysis conducted through the depiction of the ACF and PACF plots, respectively, resulted in the selection of the ARIMA model order, i.e. q and p.





CONCLUSION AND RECOMMENDATIONS

The dataset had limited features and was modeled on time series. The Arima model was able to determine the future wage gap, further confirming the inequality in earnings. The evidence from this project revealed that the average hourly wages for men across all job titles were greater than those for women in Australia. However, in the years 2005 and 2015, there was an exceptional event where women who worked as drivers and machinery operators made more money. This could be attributed to the strong employment growth during these periods according to the [Australian Government Treasury Report](#) and [Employment in Australia](#).

One of the challenges faced while modeling the data was the insufficiency of independent features (variables) to reliably forecast our target variable without bias, and probably make recommendations on factors that can be put into consideration to reduce the wage gap disparity.

We recommend that:

- To automate and predict the gap across all industries, data on the main causes of the wage gap should be made public.
- Formal policies and/or strategies that support gender equality should be in place inside organizations.
- Top cadres have a lower proportion of women, hence more women ought to be inspired to enroll in university professional programs.
- To determine whether and where discrepancies may exist, the employer should always do a pay gap analysis for all the roles, and perhaps the reason for the salary disparity and how to close it.
- Consistent efforts should be made to lessen the impact of the major causes of the pay disparity.