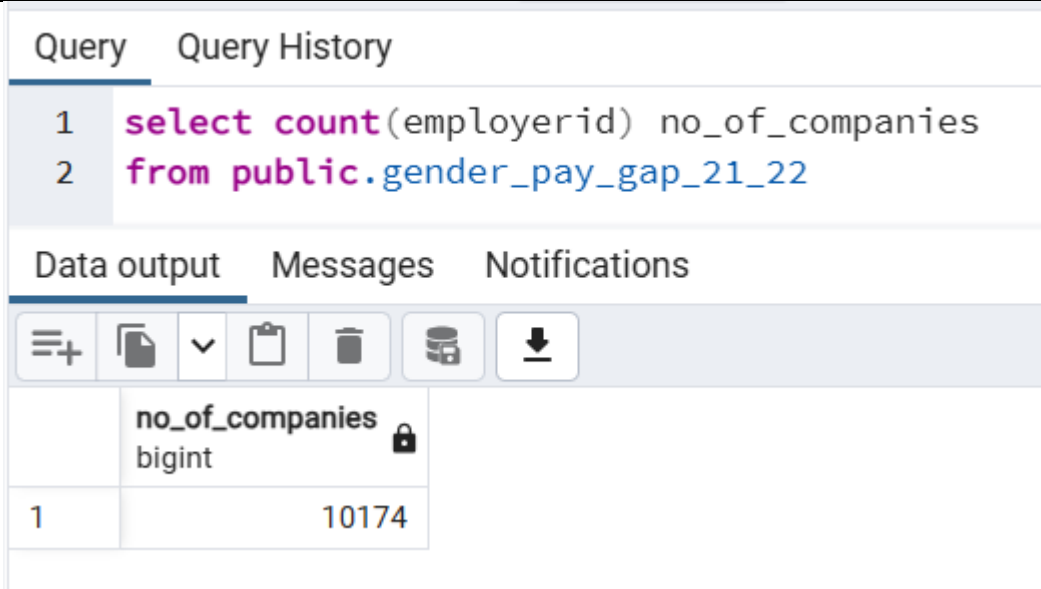
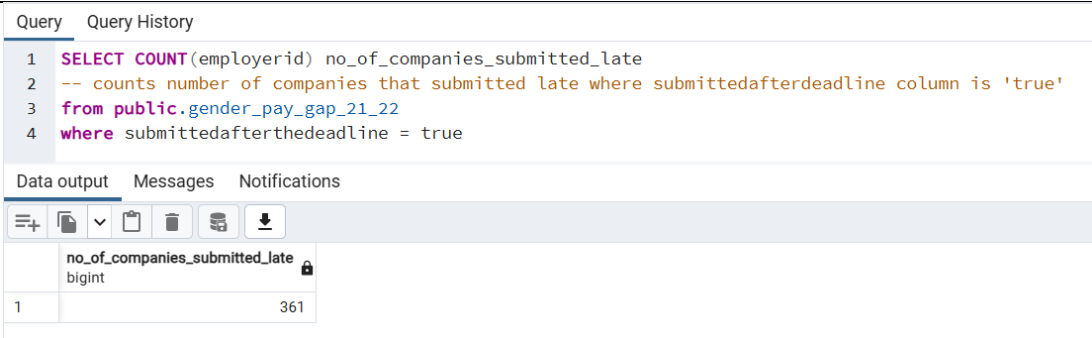



## Gender Pay Gap Analysis – SQL Project

Q/N	Question & Answer				
1	<p>How many companies are in the data set?</p>  <pre> 1  select count(employerid) no_of_companies 2  from public.gender_pay_gap_21_22 </pre> <table border="1"> <thead> <tr> <th colspan="2">no_of_companies bigint</th> </tr> </thead> <tbody> <tr> <td>1</td> <td>10174</td> </tr> </tbody> </table>	no_of_companies bigint		1	10174
no_of_companies bigint					
1	10174				
2	<p>How many of them submitted their data after the reporting deadline?</p>  <pre> 1  SELECT COUNT(employerid) no_of_companies_submitted_late 2  -- counts number of companies that submitted late where submittedafterdeadline column is 'true' 3  from public.gender_pay_gap_21_22 4  where submittedafterthedeathline = true </pre> <table border="1"> <thead> <tr> <th colspan="2">no_of_companies_submitted_late bigint</th> </tr> </thead> <tbody> <tr> <td>1</td> <td>361</td> </tr> </tbody> </table>	no_of_companies_submitted_late bigint		1	361
no_of_companies_submitted_late bigint					
1	361				
3	<p>How many companies have not provided a URL?</p>  <pre> 1  SELECT COUNT(employerid) companies_without_url FROM public.gender_pay_gap_21_22 2  -- counts number of employers that does not have 'http' in the companylinktogpginfo column 3  WHERE companylinktogpginfo NOT ILIKE 'http%' 4  -- ILIKE is used as a few URLs have capitalised 'HTTP' instead </pre> <table border="1"> <thead> <tr> <th colspan="2">companies_without_url bigint</th> </tr> </thead> <tbody> <tr> <td>1</td> <td>3700</td> </tr> </tbody> </table>	companies_without_url bigint		1	3700
companies_without_url bigint					
1	3700				
4	<p>Which measures of pay gap contain too much missing data, and should not be used in our analysis?</p> <pre> SELECT (SELECT COUNT(*) no_of_missing_dmhp FROM public.gender_pay_gap_21_22 WHERE diffmeanhourlypercent::text = '0'), (SELECT COUNT(*) no_of_missing_dmdhp </pre>				

```

FROM public.gender_pay_gap_21_22
WHERE diffmedianhourlypercent::text = '0'),
(SELECT COUNT(*) no_of_missing_dmbp
FROM public.gender_pay_gap_21_22
WHERE diffmeanbonuspercent::text = '0'),
(SELECT COUNT(*) no_of_missing_dmdbp
FROM public.gender_pay_gap_21_22
WHERE diffmedianbonuspercent::text = '0'),
(SELECT COUNT(*) no_of_missing_mbp
FROM public.gender_pay_gap_21_22
WHERE malebonuspercent::text = '0'),
(SELECT COUNT(*) no_of_missing_fbp
FROM public.gender_pay_gap_21_22
WHERE femalebonuspercent::text = '0'),
(SELECT COUNT(*) no_of_missing_mfq
FROM public.gender_pay_gap_21_22
WHERE malelowerquartile::text = '0'),
(SELECT COUNT(*) no_of_missing_fmql
FROM public.gender_pay_gap_21_22
WHERE femalelowerquartile::text = '0'),
(SELECT COUNT(*) no_of_missing_mlmql
FROM public.gender_pay_gap_21_22
WHERE malelowermiddlequartile::text = '0'),
(SELECT COUNT(*) no_of_missing_fmqlmql
FROM public.gender_pay_gap_21_22
WHERE femalelowermiddlequartile::text = '0'),
(SELECT COUNT(*) no_of_missing_mumql
FROM public.gender_pay_gap_21_22
WHERE maleuppermiddlequartile::text = '0'),
(SELECT COUNT(*) no_of_missing_fmumql
FROM public.gender_pay_gap_21_22
WHERE femaleuppermiddlequartile::text = '0'),
(SELECT COUNT(*) no_of_missing_mtql
FROM public.gender_pay_gap_21_22
WHERE maletopquartile::text = '0'),
(SELECT COUNT(*) no_of_missing_ftql
FROM public.gender_pay_gap_21_22
WHERE femaletopquartile::text = '0')

```

Data output Messages Notifications

no_of_missing_dmbp bigint	no_of_missing_dmbp bigint	no_of_missing_dmbp bigint	no_of_missing_dmbp bigint	no_of_missing_mbp bigint	no_of_missing_fbp bigint	no_of_missing_mfq bigint	no_of_missing_fmql bigint
1	49	404	2787	3381	1356	1374	212
no_of_missing_mlmql bigint	no_of_missing_fmql bigint	no_of_missing_mumql bigint	no_of_missing_fmumql bigint	no_of_missing_mtql bigint	no_of_missing_ftql bigint		
203	215	197	222	190	223		

Assumption: value of 0.0 is an input provided whereas 0 is no input provided/missing data. To only count 0 and not 0.0, I used cast to convert data type to TEXT, and then only counted '0'.

	<p>Diffmeanbonuspercent and diffmedianbonuspercent were found to have the most missing data – 2787 (&gt;27%) and 3381 (&gt;33%) out of 10174 rows respectively. Hence, they should not be used in the analysis.</p> <p>Additionally, malebonuspercent and femalebonuspercent were found to have significant missing data – 1356 (&gt;13%) and 1374 (&gt;13%) respectively. Hence they should also not be used in the analysis.</p>																
5	<p>Choose which column you will use to calculate the pay gap. Will you use DiffMeanHourlyPercent or DiffMedianHourlyPercent? Can you justify your choice?</p> <pre>12 SELECT 13 PERCENTILE_CONT(0.75) WITHIN GROUP (ORDER BY diffmedianhourlypercent ASC) + 1.5* 14 (PERCENTILE_CONT(0.75) WITHIN GROUP (ORDER BY diffmedianhourlypercent ASC) - 15 PERCENTILE_CONT(0.25) WITHIN GROUP (ORDER BY diffmedianhourlypercent ASC)) AS Q3_plus1_5_IQR, 16 PERCENTILE_CONT(0.25) WITHIN GROUP (ORDER BY diffmedianhourlypercent ASC) - 1.5* 17 (PERCENTILE_CONT(0.75) WITHIN GROUP (ORDER BY diffmedianhourlypercent ASC) - 18 PERCENTILE_CONT(0.25) WITHIN GROUP (ORDER BY diffmedianhourlypercent ASC)) Q1minus1_5_IQR 19 FROM public.gender_pay_gap_21_22</pre> <p>Data output Messages Notifications</p> <table><tr><th></th><th>q3_plus1_5_iqr double precision</th><th>q1minus1_5_iqr double precision</th></tr><tr><td>1</td><td>52.75</td><td>-30.049999999999997</td></tr></table> <pre>21 SELECT 22 COUNT(*) no_of_outliers_dmdhp 23 FROM public.gender_pay_gap_21_22 24 WHERE diffmedianhourlypercent &gt; 52.75 25 OR diffmedianhourlypercent &lt; -30.05</pre> <p>Data output Messages Notifications</p> <table><tr><th></th><th>no_of_outliers_dmdhp bigint</th></tr><tr><td>1</td><td>199</td></tr></table> <pre>12 SELECT 13 PERCENTILE_CONT(0.75) WITHIN GROUP (ORDER BY diffmeanhourlypercent ASC) + 1.5* 14 (PERCENTILE_CONT(0.75) WITHIN GROUP (ORDER BY diffmeanhourlypercent ASC) - 15 PERCENTILE_CONT(0.25) WITHIN GROUP (ORDER BY diffmeanhourlypercent ASC)) AS Q3_plus1_5_IQR, 16 PERCENTILE_CONT(0.25) WITHIN GROUP (ORDER BY diffmeanhourlypercent ASC) - 1.5* 17 (PERCENTILE_CONT(0.75) WITHIN GROUP (ORDER BY diffmeanhourlypercent ASC) - 18 PERCENTILE_CONT(0.25) WITHIN GROUP (ORDER BY diffmeanhourlypercent ASC)) Q1minus1_5_IQR 19 FROM public.gender_pay_gap_21_22</pre> <p>Data output Messages Notifications</p> <table><tr><th></th><th>q3_plus1_5_iqr double precision</th><th>q1minus1_5_iqr double precision</th></tr><tr><td>1</td><td>48.587500000000006</td><td>-22.112500000000004</td></tr></table>		q3_plus1_5_iqr double precision	q1minus1_5_iqr double precision	1	52.75	-30.049999999999997		no_of_outliers_dmdhp bigint	1	199		q3_plus1_5_iqr double precision	q1minus1_5_iqr double precision	1	48.587500000000006	-22.112500000000004
	q3_plus1_5_iqr double precision	q1minus1_5_iqr double precision															
1	52.75	-30.049999999999997															
	no_of_outliers_dmdhp bigint																
1	199																
	q3_plus1_5_iqr double precision	q1minus1_5_iqr double precision															
1	48.587500000000006	-22.112500000000004															

</

	<p>i. Difference in dollar value for pay gap is not captured in the data. With only percentage difference being given, this may not be a fair/accurate representation of the actual pay gap.</p> <p>ii. Percentage may be a good way of understanding the situation at a specific company as other variables like no. of staff is known. However, aggregating percentage across companies with varying sizes will likely lead to an outcome that is inaccurate from the actual pay gap situation (compared to aggregating using dollar pay gap, treating the whole nation as a big 'company'), as company size is not constant across companies and reported across a range</p>
--	---

8

What are the 10 companies with the largest pay gaps skewed towards men?

Query Query History

```
1 SELECT employername, diffmedianhourlypercent FROM
2 public.gender_pay_gap_21_22
3 ORDER BY diffmedianhourlypercent DESC
4 LIMIT 11
5 -- used LIMIT 11 to look for top 10 to investigate if there is a tie
```

Data output Messages Notifications

	employername character varying	diffmedianhourlypercent numeric
1	M. ANDERSON CONSTRUCTION LIMITED	100.0
2	ATFC LIMITED	100
3	HPI UK HOLDING LTD.	100
4	PSJ FABRICATIONS LTD	100
5	HULL COLLABORATIVE ACADEMY TRUST	93
6	SERVICE INNOVATION GROUP-UK LIMITED	90.4
7	BRAND ENERGY & INFRASTRUCTURE SERVICES UK, L...	89.0
8	ROBINSON WEBSTER (HOLDINGS) LIMITED	85.6
9	THE LEARNING FOR LIFE PARTNERSHIP	82.6
10	GREENBROOK HEALTHCARE (HOUNSLOW) LIMITED	77.1
11	THE FALLIBROOME TRUST	74.4

9

What do you notice about the results? Are these well-known companies?

Query Query History

```
1 SELECT employersize, siccodes, employername, diffmedianhourlypercent FROM
2 public.gender_pay_gap_21_22
3 ORDER BY diffmedianhourlypercent DESC
4 LIMIT 10
5
```

Data output Messages Notifications



	employersize character varying	siccodes character varying	employername character varying	diffmedianhourlypercent numeric
1	250 to 499	55100	HPI UK HOLDING LTD.	100
2	250 to 499	41100	M. ANDERSON CONSTRUCTION LIMITED	100.0
3	Less than 250	25110	PSJ FABRICATIONS LTD	100
4	250 to 499	56101	ATFC LIMITED	100
5	Not Provided	85200, 85310	HULL COLLABORATIVE ACADEMY TRUST	93
6	250 to 499	82990	SERVICE INNOVATION GROUP-UK LIMITED	90.4
7	1000 to 4999	96090	BRAND ENERGY & INFRASTRUCTURE SERVICES UK, L..	89.0
8	250 to 499	47710, 47910	ROBINSON WEBSTER (HOLDINGS) LIMITED	85.6
9	Not Provided	85200	THE LEARNING FOR LIFE PARTNERSHIP	82.6
10	500 to 999	86210	GREENBROOK HEALTHCARE (HOUNSLOW) LIMITED	77.1

Query Query History

```
1 SELECT employersize, COUNT(*) no_of_companies FROM
2 public.gender_pay_gap_21_22
3 WHERE diffmedianhourlypercent >= 77.1
4 GROUP BY employersize
5
6
```

Data output Messages Notifications



	employersize character varying	no_of_companies bigint
1	1000 to 4999	1
2	250 to 499	5
3	500 to 999	1
4	Less than 250	1
5	Not Provided	2

They are not well-known companies, and majority (6 out of 10) are of relatively smaller companies (499 or less employees). This could be that with less employees there will be less salary data. Companies with 100% diffmedianhourlypercent are also assumed to not have any female employees, hence  $(M-W)/M * 100\%$  would be 100%.

10 Apply some additional filtering to pick out the most significant companies with large pay gaps.

Query

Query History

1

2

3

4

5

6

SELECT employersize, employername, diffmedianhourlypercent FROM

public.gender\_pay\_gap\_21\_22

WHERE employersize IN ('5000 to 19,999', '20,000 or more')

ORDER BY diffmedianhourlypercent DESC

LIMIT 10

Data output

Messages

Notifications

	employersize character varying	employername character varying	diffmedianhourlypercent numeric
1	5000 to 19,999	EASYJET AIRLINE COMPANY LIMITED	63.6
2	5000 to 19,999	INDEPENDENT VETCARE LIMITED	49.4
3	5000 to 19,999	CVS (UK) LIMITED	42.9
4	5000 to 19,999	H&M HENNES & MAURITZ UK LIMITED	42.6
5	5000 to 19,999	SAVILLS (UK) LIMITED	41.1
6	5000 to 19,999	VETPARTNERS PRACTICES LIMITED	41
7	20,000 or more	LLOYDS BANK PLC	40.9
8	5000 to 19,999	STONEGATE PUB COMPANY LIMITED	39.2
9	5000 to 19,999	CIVICA UK LIMITED	36.5
10	5000 to 19,999	EDF ENERGY LIMITED	36.3

11 How would you report on the results? Can we say that these companies are engaging in unlawful pay discrimination?

As these companies are large (5000+ employees), there should be sufficient data points in the salary data to see an accurate representation of the pay gap. With that in mind, before drawing any conclusions in terms of unlawful pay discrimination, I would investigate further into the individual companies as to what are the reasons the pay gap is this large. For example, in airline companies where typically Pilots tend to be male and cabin crew tend to be female, it could explain the large pay gap.

12 What's the average pay gap in London versus outside London?

Query

Query History

1

SELECT

2

(SELECT PERCENTILE\_CONT(0.50) WITHIN GROUP (ORDER BY diffmedianhourlypercent ASC) median\_paygap\_percent\_London

3

FROM public.gender\_pay\_gap\_21\_22

4

WHERE address ILIKE ('%London%')

5

AND diffmedianhourlypercent::text != '0'),

6

-- filters away values that are deemed to be missing

7

(SELECT PERCENTILE\_CONT(0.50) WITHIN GROUP (ORDER BY diffmedianhourlypercent ASC) median\_paygap\_percent\_outside\_London

8

FROM public.gender\_pay\_gap\_21\_22

9

WHERE address NOT ILIKE ('%London%')

10

AND diffmedianhourlypercent::text != '0')

11

-- filters away values that are deemed to be missing

Data output

Messages

Notifications

median\_paygap\_percent\_london

double precision

median\_paygap\_percent\_outside\_london

double precision

1

12.9

10

13 What's the average pay gap in London versus Birmingham?

Query Query History

```
1 SELECT
2 (SELECT PERCENTILE_CONT(0.50) WITHIN GROUP (ORDER BY diffmedianhourlypercent ASC) median_paygap_percent_London
3 FROM public.gender_pay_gap_21_22
4 WHERE address ILIKE ('%London%'))
5 AND diffmedianhourlypercent::text != '0'),
6 -- filters away values that are deemed to be missing
7 (SELECT PERCENTILE_CONT(0.50) WITHIN GROUP (ORDER BY diffmedianhourlypercent ASC) median_paygap_percent_Birmingham
8 FROM public.gender_pay_gap_21_22
9 WHERE address ILIKE ('%Birmingham%'))
10 AND diffmedianhourlypercent::text != '0')
11 -- filters away values that are deemed to be missing
```

Data output Messages Notifications

	median_paygap_percent_london double precision	median_paygap_percent_birmingham double precision
1	12.9	8.8

## 14 What is the average pay gap within schools?

	<div> <div>Section P</div> <div>Education</div> </div> <table> <tr><td>85100</td><td>Pre-primary education</td></tr> <tr><td>85200</td><td>Primary education</td></tr> <tr><td>85310</td><td>General secondary education</td></tr> <tr><td>85320</td><td>Technical and vocational secondary education</td></tr> <tr><td>85410</td><td>Post-secondary non-tertiary education</td></tr> <tr><td>85421</td><td>First-degree level higher education</td></tr> <tr><td>85422</td><td>Post-graduate level higher education</td></tr> <tr><td>85510</td><td>Sports and recreation education</td></tr> <tr><td>85520</td><td>Cultural education</td></tr> <tr><td>85530</td><td>Driving school activities</td></tr> <tr><td>85590</td><td>Other education n.e.c.</td></tr> <tr><td>85600</td><td>Educational support services</td></tr> </table> <p>My assumption is Schools refer to any / group of education institutes. To filter these companies, I used the uk gov website for siccodes (<a href="https://resources.companieshouse.gov.uk/sic/">https://resources.companieshouse.gov.uk/sic/</a>) under Section P – Education.</p> <div>Query Query History</div> <pre> 1 SELECT PERCENTILE_CONT(0.50) WITHIN GROUP (ORDER BY diffmedianhourlypercent ASC) median_paygap_percent_schools 2 FROM public.gender_pay_gap_21_22 3 WHERE siccodes LIKE '%85100%' 4 OR siccodes LIKE '%85200%' 5 OR siccodes LIKE '%85310%' 6 OR siccodes LIKE '%85320%' 7 OR siccodes LIKE '%85410%' 8 OR siccodes LIKE '%85421%' 9 OR siccodes LIKE '%85422%' 10 OR siccodes LIKE '%85510%' 11 OR siccodes LIKE '%85520%' 12 OR siccodes LIKE '%85530%' 13 OR siccodes LIKE '%85590%' 14 OR siccodes LIKE '%85600%' </pre> <div>Data output Messages Notifications</div> <table> <tr> <td></td><td>median_paygap_percent_schools double precision</td></tr> <tr> <td>1</td><td>21</td></tr> </table>	85100	Pre-primary education	85200	Primary education	85310	General secondary education	85320	Technical and vocational secondary education	85410	Post-secondary non-tertiary education	85421	First-degree level higher education	85422	Post-graduate level higher education	85510	Sports and recreation education	85520	Cultural education	85530	Driving school activities	85590	Other education n.e.c.	85600	Educational support services		median_paygap_percent_schools double precision	1	21
85100	Pre-primary education																												
85200	Primary education																												
85310	General secondary education																												
85320	Technical and vocational secondary education																												
85410	Post-secondary non-tertiary education																												
85421	First-degree level higher education																												
85422	Post-graduate level higher education																												
85510	Sports and recreation education																												
85520	Cultural education																												
85530	Driving school activities																												
85590	Other education n.e.c.																												
85600	Educational support services																												
	median_paygap_percent_schools double precision																												
1	21																												
15	What is the average pay gap within banks?																												



Section K Financial and insurance activities

64110 Central banking

64191 Banks

My assumption is Banks refer to the 2 above siccodes  
(<https://resources.companieshouse.gov.uk/sic/>)

QueryQuery History

1SELECT PERCENTILE\_CONT(0.50) WITHIN GROUP (ORDER BY diffmedianhourlypercent ASC) median\_paygap\_percent\_banks

2FROM public.gender\_pay\_gap\_21\_22

3WHERE siccodes LIKE '%64110%'

4OR siccodes LIKE '%64191%'

Data outputMessagesNotifications

median\_paygap\_percent\_banks

double precision

1

32.4

16 Is there a relationship between the number of employees at a company and the average pay gap?

QueryQuery History

1SELECT employersize, PERCENTILE\_CONT(0.50) WITHIN GROUP (ORDER BY diffmedianhourlypercent ASC) median\_paygap\_percent

2FROM public.gender\_pay\_gap\_21\_22

3WHERE diffmedianhourlypercent::text != '0'

4GROUP BY employersize

5HAVING employersize != 'Not Provided'

6ORDER BY median\_paygap\_percent DESC

Data outputMessagesNotifications

employersize

character varying

median\_paygap\_percent

double precision

1

250 to 499

11.5

2

500 to 999

10.350000000000001

3

Less than 250

10.2

4

1000 to 4999

9.8

5

5000 to 19,999

9.5

6

20,000 or more

7.55

Query Query History

```
1 SELECT employersize, PERCENTILE_CONT(0.50) WITHIN GROUP (ORDER BY diffmedianhourlypercent ASC) median_payg
2 FROM public.gender_pay_gap_21_22
3 WHERE diffmedianhourlypercent::text != '0'
4 GROUP BY employersize
5 HAVING employersize != 'Not Provided'
6 ORDER BY median_paygap_percent DESC
```

Data output Messages Notifications



	employersize character varying	median_paygap_percent double precision
1	250 to 499	11.5
2	500 to 999	10.3500000000000001
3	Less than 250	10.2
4	1000 to 4999	9.8
5	5000 to 19,999	9.5
6	20,000 or more	7.55