



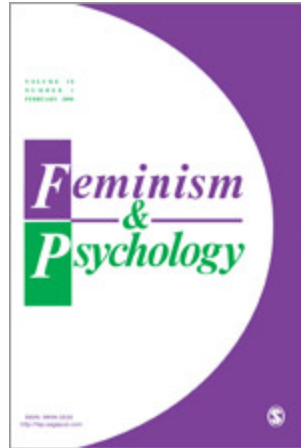
BIROn - Birkbeck Institutional Research Online

Whiley, Lilith and Walasek, L. and Juanchich, M. (2023) Contributions to reducing online gender harassment: social re-norming and appealing to empathy as tried-and-failed techniques. *Feminism & Psychology* 33 (1), pp. 83-104. ISSN 0959-3535.

Downloaded from: <https://eprints.bbk.ac.uk/id/eprint/48269/>

Usage Guidelines:

Please refer to usage guidelines at <https://eprints.bbk.ac.uk/policies.html> or alternatively contact lib-eprints@bbk.ac.uk.



Contributions to reducing online gender harassment: Social re-norming and appealing to empathy as tried-and-failed techniques.

Journal:	<i>Feminism & Psychology</i>
Manuscript ID	Draft
Manuscript Type:	Article
Keywords:	online gender harassment, sexism, misogyny, social norms, empathy, social media
Abstract:	Inspired by similar methods that have been shown to be effective in reducing online racist harassment, we designed two tweets aimed at reducing online gender harassment. Our interventions were based on the principles of social re-norming and appealing to harassers' empathy. In a sample of 666 Twitter users, we found that our intervention tweets were not successful at reducing the number of sexist slurs or sexist users either 7 days or 31 days after being sent. Our attempts also did not affect the valence, nor the arousal, of the subsequent tweets posted by our sample of Twitter users. We discuss the conceptual, methodological, and ethical challenges associated with activist research aimed at reducing online gender harassment.

SCHOLARONE™
Manuscripts

1
2
3 **Contributions to reducing online gender harassment:**
4 **Social re-norming and appealing to empathy as tried-and-failed techniques.**
5
6
7

8 **ABSTRACT**
9

10 Inspired by similar methods that have been shown to be effective in reducing online
11 racist harassment, we designed two tweets aimed at reducing online gender harassment. Our
12 interventions were based on the principles of social re-norming and appealing to harassers'
13 empathy. In a sample of 666 Twitter users, we found that our intervention tweets were not
14 successful at reducing the number of sexist slurs or sexist users either 7 days or 31 days after
15 being sent. Our attempts also did not affect the valence, nor the arousal, of the subsequent
16 tweets posted by our sample of Twitter users. We discuss the conceptual, methodological,
17 and ethical challenges associated with activist research aimed at reducing online gender
18 harassment.
19
20
21
22
23
24
25
26

27 **KEYWORDS:**
28

29 Online gender harassment; sexism; misogyny; social norms; empathy.
30
31
32

33 **INTRODUCTION**
34

35 Despite being unlawful under the Equality Act (2010), 71% of women in the United
36 Kingdom (UK) still experience gender harassment in shared spaces (UN Women UK, 2021).
37 Fifty four percent of women hear men wolf-whistling at them and 39% are called names;
38 23% are even groped (Action Aid, 2016). Out of 510 000 annual reported cases of sexual
39 assault (Office of National Statistics, 2018), only 1758 proceed to prosecution (HM Crown
40 Prosecution Service Inspectorate, 2019). One domestic abuse call is received every minute by
41 the UK police (Amnesty International, 2020). One hundred and eighty eight women were
42 murdered last year in the UK alone (Office of National Statistics, 2021). Online behaviour
43 mirrors these lived experiences – up to 45% of gender harassment is done virtually (Rights of
44 Women, 2021). Indeed, online gender harassment reflects the wider misogynist treatment of
45 women; it is “firmly grounded in the material realities of women’s everyday experiences of
46 sexism in patriarchal society” (Megarry, 2014, p. 49). The online space is a highly gendered
47 space (Locke et al., 2018). It is a heteronormative and hegemonically masculine space
48 (Drakett et al., 2018) - Han (2018) even deems it a space of *toxic* masculinity, “of
49 technological privilege where the masculine elite dominates the archetypal passive
50
51
52
53
54
55
56
57
58
59
60

1
2
3 sexualised woman” (Lock et al., 2018, p.7). To illustrate, tweets that blame-the-victim and
4 slut-shame rape survivors have more followers and retweets than those who support the
5 women (Stubbs-Richardson et al., 2018). Over 400 000 sexist slurs are posted on Twitter...
6 every day (Felmlee et al., 2020). Offenders typically attack women’s physical appearance
7 (e.g., ‘ugly cunt’), their intelligence (e.g., ‘stupid slut’), and their age (e.g., ‘old bitch’). The
8 harassment is often based on accusing women of ‘failing’ to meet patriarchal norms of
9 femininity (i.e., hegemonic standards of thin, young, innocent, and passive beauty). Women
10 receive death threats and calls for rape (Chen et al., 2020); comments can be particularly vile,
11 for example: “She gave great blowjobs before her fall, now imagine the pleasure she will
12 bring without her front teeth” (Jane, 2014, p. 561). Women are harassed on online dating sites
13 (Thompson, 2018) and are re-victimised in so-called ‘revenge porn’ networks – while abusers
14 hide behind the protective anonymity of online spaces (Uhl et al., 2018). The fact that these
15 attacks (and they *are* attacks) are being perpetuated online creates a false sense of triviality –
16 a misconception that because they are not physical attacks they should simply be ignored by
17 women who are, surely, over-reacting (again) (Chadha et al., 2020).
18
19
20
21
22
23
24
25
26
27
28
29
30

31 Yet, sixty one percent of women who are sexually harassed online have trouble
32 sleeping afterwards and 55% experience anxiety (Amnesty International, 2017). Women who
33 experience online abuse also experience fear and depression (Lindsay et al., 2016). The
34 COVID-19 pandemic has further exacerbated women’s suffering because online gender
35 harassment has increased due to working from home; for example, one woman shared how
36 offenders are now in one’s home and bedroom: “I feel my privacy has been invaded and
37 nowhere is safe” (Rights of Women, 2021). Indeed, the purpose of sexual harassment is to
38 violate someone’s dignity, to intimidate, degrade, and humiliate them, and create a hostile
39 environment (Citizens Advice, 2021). Up to 67% of women who experience online gendered
40 harassment feel apprehensive about using social media again (Amnesty International, 2017).
41 Women are more cautious about what they post in order to “keep quiet so as to reduce abuse”
42 (Adams, 2017, p. 7), actively avoid voicing their opinions in online discussions (Chadha et
43 al., 2020), and “[watch] over [their] shoulder in cyberspace” (Chen et al., 2020, p. 887). For
44 these reasons, some women decide to leave social networking altogether (Citron, 2014).
45 Online gender harassment thus limits women’s equal participation in online communities and
46 social networks (Megarry, 2014). Constraints and limitations are imposed on women’s
47 freedom in both the physical world and the online one (Vera-Gray, 2017). They can also have
48 a profound impact on women’s livelihood in what Jane (2018) terms ‘economic vandalism’
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 by either directly or indirectly impacting on women's professional lives. Online gender
4 harassment is insidious and proliferates in almost every aspect of women's lives: "it was
5 actually the destruction of my person" (Chen et al., 2020, p. 884). In these ways, online
6 gender harassment becomes another means by which women's behaviour is monitored,
7 policed, and contained – especially when they are perceived to be breaching patriarchal
8 hegemonic social norms.
9
10
11
12
13
14

15 What channels exist for dealing with online harassment? Online gender harassment
16 can be reported to the police as either 'harassment' or 'malicious communications' (Met
17 Police, 2021). It can also be reported directly to the social media platform, but despite
18 Facebook, Twitter, and YouTube agreeing a Code of Conduct on Countering Illegal Hate
19 Speech Online with the European Commission, 43% of women in the UK still think that the
20 responses from social media giants are inadequate in addressing online gender harassment
21 (Amnesty International, 2020). Just this year, Slack disabled its private DM function over its
22 'potential' to facilitate harassment. Women explain how, despite several complaints, few if
23 any posts are deleted; responses also range from automated emails to speedy investigations
24 exonerating the offenders. Certainly, social media giants lament how difficult it is to regulate
25 'hate speech' while citing 'freedom of speech' as one reason for their limited intervention
26 (House of Commons, 2017). Interestingly, male internet users think that 'censorship' is their
27 greatest threat, whereas women believe it to be 'privacy' (Herring, 2003). Many of the
28 recommended courses of action such as 'unfriend the person', 'block the person', and 'don't
29 retaliate' (House of Commons, 2017), do little in the way of giving women resources to
30 *actually* respond to offenders. Indeed, the advice to ignore the problem is harmful; Mallett et
31 al. (2019) found that when women did not confront instances of harassment, it desensitised
32 them and increased their tolerance for future abuse.
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

48 Against this backdrop, we wanted to contribute by giving women resources to stand
49 up against online gender harassment. Jane (2019, p. 1) said – and we agree, that "shutting
50 down sexual harassment shouldn't make you shake in fear or feel like your stomach just fell
51 10 stories. We don't just need to be empowered, but released from the burden of protecting
52 men's comfort at the expense of ourselves". We are by no means alone in our endeavour;
53 TrollBusters are an "online pest control" fighting against the online harassment of women
54 (Ferrier & Garud-Patkar, 2018, p. 316), the campaign #OutThem actively denounces male
55 harassers, and #MenCallMeThings gives women voice by revealing and re-tweeting sexist
56
57
58
59
60

1
2
3 comments. Indeed, online spaces can *also* be ‘safe spaces’ to learn about feminism and
4 connect with other feminists (Jackson, 2018). We therefore set out to create simple messages
5 that women could actively post in response to online gender harassment. Although in doing
6 so we are both, engaging in activism to ‘shout back’ at online misogyny via feminist action
7 (e.g., Turley & Fisher, 2018), and responding to academic calls to design interventions to
8 strengthen women’s voices in online spaces (e.g., Jane, 2014), we have nevertheless received
9 a significant amount of criticism from colleagues for delivering our interventions from a
10 *female* Twitter account. In fact, two editors and reviewers from elsewhere urged us to
11 replicate our study from a male Twitter account. We declined; in our view, women should
12 feel perfectly capable of standing up for themselves. We agree that online harassment is
13 disproportionately done *by men to* women, and the greatest responsibility lies with men to
14 change *their* behaviour. At the same, women have the right to speak up and should not have
15 to suffer online gender harassment in silence; indeed, women do have a range of strategies in
16 their repertoire for responding to online abuse, for example in gaming (Cote, 2017), and
17 online spaces can be spaces for feminist activism too (e.g., #MeToo, #FreeTheNipple). How,
18 therefore, can women respond to online gender harassment?
19
20
21
22
23
24
25
26
27
28
29
30
31

32
33 There is certainly an extensive amount of research on the importance of social norms
34 in promoting behavioural change (see the review by Paluck & Green, 2009). The idea is to
35 encourage people to change their behaviour without any external incentives by simply
36 communicating information about ‘what is commonly done’ (Schultz et al., 2018). People
37 begin to realise that others do not engage in the same behaviour as much and would
38 disapprove of them. Social norm campaigns have been successful in reducing alcohol
39 consumption (Perkins & Craig, 2006) and smoking (Hancock and Henry, 2003). They have
40 also had some success online, for example, in a community group with 13 million
41 subscribers, Matias (2019) found that announcing socially normative expectations of
42 members’ behaviours increased compliance and reduced harassment. Given that one purpose
43 of online abuse is to harass women into conforming to patriarchal social norms (e.g., Felmlee
44 et al. 2020), what would happen if women attempted to ‘re’-norm offenders’ beliefs?
45 Accordingly, we reasoned that if we informed misogynist offenders that most people
46 disapprove of their sexist language, that this could reduce the number and frequency of their
47 sexist Tweets. Indeed, most men over-estimate others’ sexism and educating them about this
48 could be the first step (Kilmartin et al., 2008).
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Another way to tackle online harassment could be to appeal to offenders' emotions and invite them to take the perspective of those that are discriminating against (Dovidio et al., 2004). Interventions that encourage people to focus on the feelings of another person have been shown to arouse feelings of empathy and reduce prejudice towards members of an outgroup (Batson et al., 2002; Galinsky & Moskowitz, 2000). Studies show that taking the perspective of a stigmatised group in particular can improve attitudes towards that group (Vescio et al., 2003), reduce in-group favouritism (Galinsky and Moskowitz, 2000), and increase helping behaviours (Mallett et al., 2008). One way of tackling online gender harassment could therefore be to encourage perpetrators to take the victims' perspective and to inform them about the negative emotional consequences of their misogynist tweets. Doing so might encourage harassers to think about the impact of their language, appeal to their empathy, and prompt a reduction of subsequent sexist tweets.

Method

Overview

Our methodological approach was inspired by Munger (2016), who showed that targeted tweets can be effective in reducing online racist harassment. More precisely, in a sample of 242 Twitter users, Munger (2016) found that participants who were 'told off' by a (male) Twitter user significantly reduced the number of racist slurs in their future tweets. We designed a similar approach to tackling online gender harassment. We identified a sample of misogynist Twitter users who frequently tweeted sexist slurs and posted two tweets using the @function. One message aimed to socially re-norm sexist users and the other called for empathy. We then reviewed the pre-and-post streams of tweets to assess if our interventions had any effect. To foreshadow our results, they did not. We transparently share our methodological approach and decision making below.

Step 1: Identify misogynist Twitter users

Step 1 was not that difficult given the large population group (!), but we still needed to identify a sample of users with whom we could try our interventions. Given that tweets are only 280 characters long, we needed a very concise and precise way to identify online gender harassment, and to do so, we operationalised it via the presence of either one of two sexist slurs: "fucking bitch" and "fucking cunt". The most popular derogatory term on Twitter is "fuck", which accounts for 34.73% of all curse word occurrences (Wang et al., 2014). "Bitch", "cunt", and "slut" are gendered slurs that target women specifically and are

1
2
3 commonly posted on Twitter (Felmlee et al., 2020). Initially, we double barreled our slurs
4 and combined “fuck” with all three terms, but “fucking slut” brought up mostly pornographic
5 contents. While we acknowledge that some porn can be misogynist, this was beyond the
6 scope of our stud, so we selected the sexist slurs “fucking bitch” and “fucking cunt”. Our first
7 objective was to obtain a sample of tweets that featured these sexist slurs. To do so, we used
8 the StreamR package in R (Barbera, 2014) to connect to Twitter’s official application
9 programming interfaces (API) and collected tweets (i.e., scraped) over a period of six days.
10 Our initial sample consisted of whopping 89,939 tweets.
11
12
13
14
15
16
17
18

19 We proceeded to automatically remove non-alphanumeric symbols, links, excessive
20 white space, numbers, and usernames (e.g., “@username” inside the tweet’s body). We
21 screened out all retweets and removed duplicates. The initial filtering process left us with
22 6,024 tweets (out of the initial 89,939) that featured at least one of our two sexist slurs (i.e.,
23 “fucking bitch” and “fucking cunt”). We still found that a large proportion of these tweets
24 were pornographic content (e.g., advertisements), and for this reason, we decided to
25 strengthen our exclusion criteria. To do so, we removed tweets that included more than three
26 hashtags (i.e., #word) and tweets from users who tweeted most often (upper quartile of
27 average activity in our sample [75th to 100th] - since these users turned out to be
28 predominantly advertisers). The remaining sample included 2,970 misogynistic tweets that
29 featured at least one of our two sexist slurs; these were posted from 2,844 Twitter users.
30
31
32
33
34
35
36
37
38
39

40 We then proceeded to manually code each tweet to confirm that it was indeed aimed
41 at harassing women. This process was arduous, time consuming, and shocking – we were
42 disappointed and saddened at the vehement violence that was directed at women on Twitter.
43 We also experienced several methodological challenges. For example, we had to identify and
44 wean out tweets where the sexist slur was negated or those where the slur was used in a
45 power affirming way (e.g., “Well done you fucking bitch! You nailed it!”), but the intent was
46 not always easy to decipher. Our coding framework (figure 1) emerged iteratively by toing-
47 and-froing between the tweets and discussions between authors to assess their relevance.
48
49
50
51
52
53
54

55 **Code as relevant:**

- 56
57 1. Tweets that use a sexist slur in an unambiguously derogatory way.
58 2. The slur is made against/about/in reference to women/a particular woman.
59
60

- The tweet should NOT be about a man.
3. The slur is NOT a report of someone else being derogatory (e.g., “someone screamed fucking bitch while I was driving. Okay cool”)
 4. The slur is NOT used in an endearing/empowering way (e.g., “my fucking bitch”).
 5. The slur is NOT self-deprecating (e.g., I am a fucking bitch I know...)
 6. The slur is not associated with joking/laughing (e.g., “fucking bitch gave me a fright lol”)
 7. The slur is NOT negated (e.g., “my mum is happy I am **not** a fucking bitch”).
 8. The tweet does NOT come from porn companies (e.g., porn tweet: “slim east Asian, fucking bitch with her long dildo”).

Figure 1. Coding framework for manually identifying sexist slurs.

We then assessed inter-rater reliability amongst the three authors (highest Cronbach’s alpha = .634) and selected the tweets that were above the chance threshold with an agreement rate of 65% or more; in this way, we arrived at a sample of 1,000 tweets containing sexist slurs harassing women. Given that manually coding such a large sample of tweets was time consuming, by the time we had accomplished our goal, only 847 offending users were still active on Twitter. The sample attrition may have occurred because some users might have closed their accounts, changed their privacy settings, or changed their username handle.

Step 2: Designing and delivering our interventions

We create two tweets based on re-norming and encouraging empathy, respectively:

@_____ *Most people believe that some of your tweets against women are simply unacceptable.*

@_____ *Women are hurt by some of your tweets. Take a minute to think about how they feel.*

To assess their suitability for our purpose and determine whether people would indeed interpret these statements as communicating social norms and appealing to empathy, we presented both tweets to an independent and unrelated sample of 272 participants (136 participants evaluated each tweet). We asked whether these tweets were believable (yes/no), if their presumed goal was to stop online gender harassment (yes/no), and might they allude to social norms or empathy. For both interventions, we also asked a “check” question that stated an erroneous goal (that the tweet was aimed at encouraging people to recycle) to avoid capturing acquiescence as evidence of understanding. The results showed that participants

1
2
3 believed that both our tweets were realistic (90% and 85% for the re-norming and empathy
4 interventions respectively) and that their aim was to stop the recipient from harassing women
5 online (93% and 79%). Participants also correctly identified the re-norming tweet as
6 communicating social disapproval from most people (89%) and the empathy tweet as
7 appealing to the recipient's emotions (82%).
8
9
10
11
12

13 We then proceeded to randomly allocate our sample of 847 users into two
14 experimental groups (n=282 in the re-norming tweet condition and 282 in the empathy one)
15 and a control group (n = 283) to whom we did not send any intervention tweet. We sent the
16 intervention tweets at regular intervals to abide by Twitter's rules and regulations concerning
17 the limited number of tweets that can be sent to other users in any given hour. All tweets
18 were sent from a research account that we had named "Lizzy _____" belonging to a
19 fictional woman named Elizabeth _____. We addressed each unique user specifically via the
20 "@username [intervention message]" format. Two users replied: "Hiya Lizzy he just dmed
21 me telling you to lick his bald head" and "#Balded". Someone retweeted our intervention
22 tweet and someone liked our intervention tweet.
23
24
25
26
27
28
29
30
31
32

33 Tweets were continuously monitored, and data were collected for a period of 62 days
34 - 31 days before and after the intervention Tweets. Afterwards, we individually tweeted the
35 messages below to our sample to debrief them and give them the opportunity to withdraw
36 their data. No one requested to withdraw their data.
37
38
39

40 @_____ You have been part of a study on online behaviour towards women.
41 We are interested in finding solutions to reduce poor online behaviour such as being
42 derogatory against women.
43
44

45 @_____ We hope you value our interest in improving girls and women's
46 lives. If you would like to withdraw your participation from our study, please let us know by
47 emailing: withdrawresearch@gmail.com with your Twitter username.
48
49
50

51 Data analysis

52
53 For each user, we extracted their Twitter activity exactly 31 days prior and 31 days
54 after the intervention tweets. For the control group, we used a 62-day window of activity that
55 we split in two 31-day periods: pre and post *non*-intervention to make the number of tweets
56 comparable across conditions. There was a further sample attrition because some people did
57 not tweet at all or tweeted very rarely during this time frame (accounts that tweeted fewer
58
59
60

than five times either before or after the intervention were excluded). Our final data sample included 487,659 tweets from 666 users; 218 were in the re-norming condition, 214 were in the empathy condition, and 234 were in the control condition. Table 1 shows descriptive statistics of our final sample.

Table 1. Total and daily tweeting frequency and follower counts across experimental groups for Twitter users in our studies.

	Condition			
	Re-norming	Empathy	Control	Total
Number of users	218	214	234	666
Number of tweets over 62 days	164,621	145,335	177,703	487,659
Median number of tweets per day	7	8	9	8
Median followers count	775	529	573	619

We assessed the frequency with which users posted sexist slurs by developing a list of expressions derogating women from urbandictionary.com (e.g., “ballbuster”, “cocktease” – see appendix I for a full list). This approach allowed us to form a ‘big’ picture overview and to assess changes in discourse more generally; at the same time, there were two challenges that needed to be addressed. Firstly, some of the sexist slurs above could be used in non-derogatory ways (e.g., ‘tart’ to refer to a pie). Secondly, specific slurs are limited in capturing other forms of online gender harassment that are also misogynistic (e.g., “this woman was so fucking stupid that it was actually fun to see her fail”) or threatening (e.g., “I would like to kill this woman”). We therefore complemented the focus on frequencies of specific sexist slurs by assessing the valence and arousal of the words composing the tweets across condition. Our reasoning is based on the premise that words carry and evoke emotions in people (e.g., happy, unhappy etc.) (Warriner, Kuperman, and Brysbaert, 2013). Words can thus be understood in terms of the valence of that emotion (i.e., positive or negative) and arousal (i.e., low or high intensity). Some words can have both a positive valence and high arousal (e.g., “excited”) and others can have a neutral valence and low arousal (e.g., “table”). Offensive words have both high negative valence and high arousal (e.g., “bitch”). They imply very negative feelings and high levels of intensity. Using Warriner et al.’s (2013) coding of 13,915 words and matching them to our sample of tweets, we were also able to

1
2
3 explore if there were any changes in the valence and arousal of users' tweets following our
4 intervention tweets.
5
6
7

8 RESULTS

11 Effects of the intervention on sexist tweets and users

12 To compare a user's propensity to tweet a sexist slur, we focused on the normalised
13 variables: (a) the frequency of tweets featuring a sexist slur out of the total number of tweets
14 sent by a given user, and (b) the number of users who tweeted a sexist slur (at least once) out
15 of the total number of users in a given condition. We also checked the transience of our
16 interventions on both the short-term (i.e., 7 days after our intervention) (see table 2) and
17 longer-term (i.e., 31 days after our intervention) (see table 3). To compare the rate of sexist
18 tweets and sexist users before and after the intervention, we computed: (a) the number of
19 sexist tweets after our intervention and deducted from this the number of sexist tweets that
20 were being posted before our intervention (scores ranged from -9.09% to +14.29%), and (b)
21 the number of sexist users after our intervention minus the number of sexist users before our
22 intervention (ranges from -1 to 1). A difference score of 0 meant that the intervention did not
23 have an effect, whereas a positive difference meant that the rate of sexist tweets and sexist
24 users increased after the intervention, and finally, a negative difference meant a decrease in
25 the rate of sexist tweets and sexist users.
26
27
28
29
30
31
32
33
34
35
36
37
38

39 As shown in Table 2, the rate of sexist slurs and sexist users did not vary greatly
40 before and after the intervention (see rows in bold). In the social re-norming condition, there
41 was an increase in the number of sexist slurs while the number of Twitter users who tweeted
42 a sexist slur remained stable. In the empathy condition, we noticed both an increased trend in
43 the number of sexist slurs and an increase in the number of users who tweeted a sexist slur.
44 However, the most important increase in sexist slurs and users occurred in the control
45 condition. To assess significance, we used a non-parametric Kruskal-Wallis test, which
46 showed that the effects were not statistically significant in either the short (7 days) or the long
47 term (31 days), $Kruskal-Wallis(2) = 2.98, p = .225$ and $Kruskal-Wallis(2) = 1.15, p = .564$.
48 A chi square comparing the change in proportion of sexist users across conditions was not
49 statistically significant either, whether we considered the short term effect (7 days) or the
50 longer one (31 days), $\chi^2(4, N = 578) = 7.53, p = .110$, $Cramer's V = .08$ and $\chi^2(666) = 2.56, p$
51 $= .634$, $Cramer's V = .05$.
52
53
54
55
56
57
58
59
60

Table 1. Percentage of sexist tweets and sexist users 7 days after our intervention

7 days after our intervention tweets			
Condition of sample		% of sexist tweets	% of sexist users
Social re-norming	Before	0.42%	20%
	After	0.63%	21%
	Difference	+0.24%	+1%
Empathy	Before	0.73%	27%
	After	0.69%	24%
	Difference	+0.03%	-3%
Control	Before	0.57%	22%
	After	1.11%	26%
	Difference	+0.54%	+4%
Total	Before	0.57%	23%
	After	0.82%	24%

Table 3. Percentage of sexist tweets and users 31 days after our intervention tweets

31 before/after			
Condition of sample		% sexist tweets	% of sexist users
Social re-norming	Before	0.57%	46%
	After	0.60%	46%
	Difference	+0.11%	-/+0%
Empathy	Before	0.48%	51%
	After	0.56%	53%
	Difference	+0.06%	+2%
Control	Before	0.68%	50%
	After	0.70%	54%
	Difference	-0.03%	+4%
Total	Before	0.58%	49%
	After	0.62%	51%

Effect of the intervention of the valence and arousal of tweets

Figure 2 illustrates the valence and arousal averages for all tweets across the 62-day research window of our study. To establish that tweets that included one of the sexist slurs would be more negative and more arousing than tweets that did not, we evaluated the valence and arousal of tweets for tweets that included a sexist slur and those that did not. We found a difference, with tweets that included the slurs were markedly less positive and generated stronger arousal. However, as is clear from the flat pattern over time, our intervention tweets did not have any effect on the valence and arousal of the words being used. Users' tweets following our interventions were neither less negative nor less emotionally loaded.

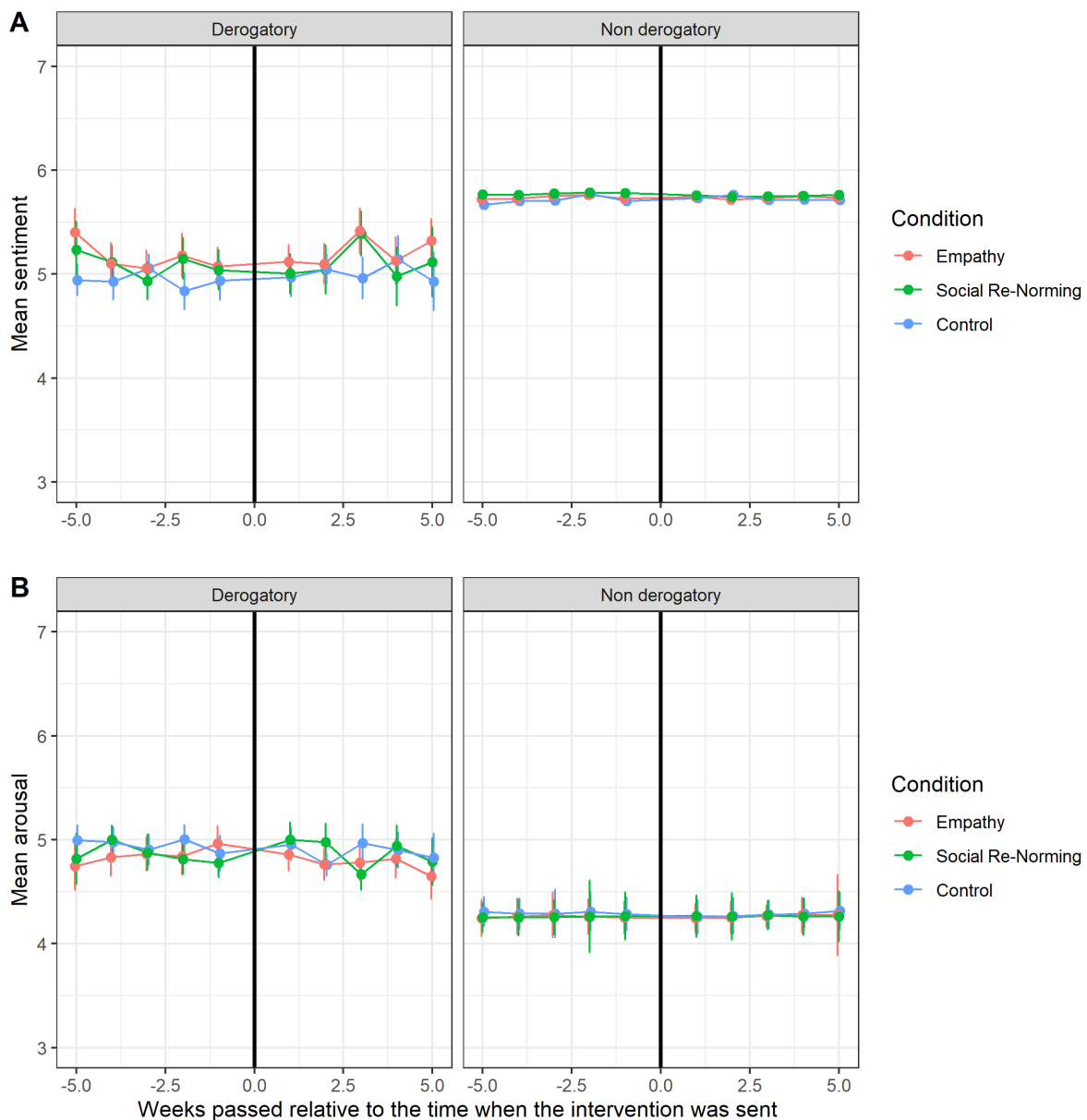


Figure 2. Panel A: Valence of tweets over the course of the study in weeks (ranging from 1: completely unhappy to 9: completely happy). Panel: B: Arousal of tweets over the course of

1
2
3 the study in weeks (ranging from 1: completely calm to 9: completely aroused). In both Panel
4 A and B, left panel shows tweets that include one of the sexist slurs; right panel: the
5 remaining tweets. Error bars represent 2 standard errors of the means (they are too small to be
6 clearly visible for non-derogatory tweets).
7
8
9

10 11 12 **GENERAL DISCUSSION**

13
14
15 Responding to calls by Turley and Fishers (2018) and Jane (2014) to empower women in
16 online spaces, we designed two straightforward responses that women could tweet in response
17 to online gender harassment. Our preconceptions were that (1) social re-norming, and (2)
18 appealing to empathy could decrease sexist slurs – in the same that these types of messages
19 were found to reduce online racist harassment (Munger, 2016). We also tested to see if the
20 valence and arousal of tweets posted before and after our interventions changed. Regrettably,
21 our interventions did not reduce the frequency of sexist tweets nor the number of sexist users
22 either 7 days or 31 days after. We did not observe a change in the valence or arousal of users’
23 tweets, nor a reduction in the overall rate that users tweeted (with or without sexist slurs).
24 Although these findings are disappointing, it is important to reflect on the possible reasons for
25 our interventions’ lack of effect and discuss the conceptual, technical, and ethical challenges
26 associated with reducing online gender harassment.
27
28
29
30
31
32
33
34
35
36
37

38 **Conceptual Challenges**

39
40 Our first tweet attempted to socially ‘re’-norm offenders by reminding them that most
41 people found their tweets against women unacceptable. Communicating social norms is an
42 effective way to nudge behaviour change (Paluck & Green, 2008). For example, research
43 shows that social norm interventions are successful in reducing excessive towel use in hotels
44 (Goldstein, Cialdini, & Griskevicius, 2008), enhancing compliance with community rules
45 (Matias, 2019), limiting alcohol consumption (Perkins & Craig, 2006) – and even reducing
46 intentions to harass in Facebook groups (Van Royen et al., 2017). Following this logic, our
47 tweet to socially re-norm offending users should have had some effect on their subsequent
48 tweets, but this was not the case.
49
50
51
52
53
54
55
56

57 Our second intervention was based on highlighting the affective consequences of using
58 misogynistic language and appealing to users’ empathy. Research shows that encouraging
59 people to take the perspective of others and develop empathy can decrease prejudice (Batson
60

1
2
3 et al., 2002; Galinsky and Moskowitz, 2000; Vescio et al., 2003), and intervention messages
4 that highlight the negative consequences of online harassment (e.g., “This comment may be
5 hurtful for the receiver. Are you sure to post it?”) were found to be successful in reducing the
6 intention to harass on Facebook (Van Royen et al., 2017). Yet again, this was not the case in
7 our study, and we did not find that our intervention tweets had an effect on the frequency of
8 sexist slurs tweeted or the number of users tweeting derogatory material.
9
10
11
12
13
14

15 There are several reasons why our interventions might not have reduced the use of sexist
16 slurs. First, it is possible that this result is a type II error: the effect exists, but we were not
17 able to statistically capture it in this sample. Our study focused on a sample of 666 Twitter
18 users who posted before and after our interventions, and they were split across three
19 conditions: social re-norming, empathy, and control. When comparing one of our two
20 experimental conditions to the control condition, we had a 90% power (with a 5% alpha) to
21 detect a small to medium between-subject mean difference in the number of tweets including
22 a slur (Cohen’s $d = .29$). We could argue that even a difference of 0.5% could actually be
23 meaningful and represent a large number of tweets (we found 89,939 tweets featuring
24 “fucking cunt” or “fucking bitch” in only 6 days). Alternatively, it may be that single tweets
25 are simply not powerful enough to prompt misogynist behaviour change. We are exposed to
26 such an inane amount of content on social media that a single tweet may have been drowned
27 in masses of other emotion-rich contents, and it may be that a greater number of intervention
28 tweets could actually have an impact - for example, by sending multiple similarly worded
29 messages from several different accounts, but this is also problematic from an ethical
30 perspective for this would constitute ‘harassing the harassers’. Notwithstanding, we do know
31 that at least some offenders in our sample did receive and noted our messages because we
32 received a few reactions to our tweets (e.g., likes, retweets, replies). Yet, nevertheless, a tweet
33 is only a micro-intervention in a macro-level system of entrenched in sexism.
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49

50 Further, sexism is so deeply ingrained in our society (e.g., #MeToo, Time’s Up) and online
51 gender harassment is so normalised on the internet (e.g., Felmelee et al., 2020), that we were
52 perhaps overtly optimistic in attempting to reduce it via a couple of tweets – despite this
53 approach being shown to be successful in other online studies (e.g., Munger, 2016; Pennycook
54 et al., 2021). Racism is believed to be more offensive than sexism (Woodzicka et al., 2015) and
55 individuals who are called up on using racist slurs might feel more embarrassed at being so
56 openly confronted than people using sexist slurs ...after all, sexist attitudes are very common
57
58
59
60

(Georgeac et al., 2019). Felmelee et al. (2020) found over 2.9 million tweets in just one week that contained sexist slurs. This shocking rate maintains the online gender harassment cycle because these tweets reinforce the idea that ‘everybody does it’. Certainly, sexist harassers might feel less chastised in online spaces than they might in real-life, especially given the protection of anonymity. Disclosing one’s true identity can reduce the use of offensive words (Cho and Acquisti, 2013); for example, Lapidot-Lefler and Barak (2012) found that participants assigned to the eye-contact condition via webcam were twice less likely to engage in flaming behaviours than those assigned to the no-eye-contact condition. In our case, although some users did display demographic data, many did not, and it was impossible to tell whether those who did used their real information. It could therefore be that our tweet did not threaten to expose them in any meaningful way – like campaigns such as #OutThem do so successfully.

Yet another reason, grounded in patriarchy and misogyny, might be that our intervention tweet was posted by someone who clearly appeared to be a woman: Elizabeth_____. Women who confront sexism are often denigrated as being hysterical ‘whiners’ (Doyle, 2011) and ‘over-reactors’ (Czopp et al., 2006), thereby enabling their views to be more easily discounted. Men, of course, are taken more seriously than women when they confront sexism (Drury and Kaiser, 2014). Although we tried to mitigate this by using the gender neutral ‘most people’ as our reference group in the re-norming condition, we nevertheless recognise that the confronter’s apparent gender might play a role in the intervention’s lack of effect – a male confronter (similar to Munger’s 2016 design on online racist harassment) might yield different results. Although we do support campaigns that engage men in the eradication of gender harassment (e.g., He4She), women also need resources to speak up and call out harassment. Our aim is, specifically, to empower *women’s* voices in online spaces, and to those who criticise our decision to use a female confronter and urge us to replicate our study with a male confronter, we, respectfully, ask you to replicate yours with a female confronter.

Technical and Ethical Considerations

Efficiently identifying online gender harassment for research purposes is difficult on social media because despite using stringent filtering criteria on raw tweets, we nevertheless had to resort to manual coding. Our initial sample of tweets contained an overwhelming amount of pornography. We managed to exclude a substantial amount of those by filtering out tweets that featured web links and more than three hashtags; nonetheless, we found a significant number of pornographic tweets while manually inspecting our data. This is not

1
2
3 just an obstacle for research, but a major social concern – why are social media sites
4 permitted to host pornographic contents? To put this into perspective, we would not normally
5 expect to find pornographic contents on trains or in the park because these are communally
6 shared spaces, frequented by both adults and children, yet, online, communal shared spaces
7 are somehow exempt from these standards. Anyone can see and access porn on Twitter. It is a
8 different matter altogether to host pornographic contents on sites specifically designated for
9 adults such as PornHub. More generally, we ask, should websites take more responsibility
10 and actions for policing offensive contents? Indeed, it appears that they must when it comes
11 to copyrighted material (e.g., Copyright Directive 2019) so why not harassing content? Yet
12 despite several high-profile cases and activism by groups such as Amnesty International,
13 social media giants are largely only meekly ‘policing’ themselves – with little to no impact
14 on harassed women’s actual lived experiences (e.g., Chadha et al., 2020; Amnesty
15 International, 2020). Moreover, the Home Affairs Committee (2017: 31) in the UK has
16 criticised social media companies’ reliance on users to report abuse as “outsourcing the vast
17 bulk of their safeguarding responsibilities at zero expense”. This is simply one example of a
18 larger issue around social media companies failing to adequately address hate speech and
19 misinformation on their platforms.
20
21
22
23
24
25
26
27
28
29
30
31
32
33

34 A second reason why efficiently identifying online gender harassment is challenging for
35 research purposes is because it is not possible to automatically detect slurs that are used in an
36 empowering way. For example, marginalised groups often ‘take ownership’ of derogatory
37 words that have been historically used against them (Galinsky et al., 2013) (e.g., the adoption
38 of the word ‘queer’ by gender non-conforming persons), but manually coding such a large
39 dataset is resource intensive (see Schwartz and Ungar, 2015 for further guidance on how to
40 review social media posts). The creation of algorithms to automatically detect a range of
41 negative content online is currently a pressing topic to tackle all forms of harassment
42 including hate speech (Schmidt & Wiegand, 2017) and cyberbullying (Van Hee et al., 2018)
43 – see Zimmerman et al. (2018) for discussions on how to improve detection. Other avenues
44 for research include how women might take ownership of sexist discourse in online spaces in
45 an empowering way, and how it is precisely the *femininity* in sexist slurs that is perceived to
46 be offensive (see Hoskin's 2019 work on femmephobia), for example, by ‘insulting’ a male
47 footballer in saying that he plays like a ‘bitch’.
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 Our involvement in this study also brought up interesting debates about conducting
4 ethical research online. We were engaging with publicly available data and did, of course,
5 acquire ethical approval from our university. It is however necessary that users (yes, even the
6 sexist ones) whose behaviours are being monitored, are made aware that they are
7 participating in a research project, but in some instances, it is not feasible to do so before the
8 study because its premise relies on being covert. In those cases – as in ours, we felt it was
9 important to debrief participants after the study and give them the opportunity to withdraw
10 their data. Yet, this also brings up another uncomfortable dilemma for researchers. Do we
11 want to open ourselves to being harassed by people that are clearly prone to harassing? Vera-
12 Gray (2017) has already documented the noted dangers of women academics being trolled for
13 simply doing research online. Despite the time-consuming nature of the activity, we manually
14 sent individual debrief @tweets to everyone in our sample, explaining that we were doing
15 research and giving users the opportunity to withdraw their data, but we chose not to disclose
16 our identity and directed participants to an anonymous email research account. For those
17 interested in ethical internet-based research, see the BPS Ethics Guidelines for Internet-
18 Mediated Research, 2017, or the AoIR Internet Research: Ethical Guidelines 3.0 (2019) for a
19 guide.
20
21
22
23
24
25
26
27
28
29
30
31
32
33

34 CONCLUSION

35
36 Online gender harassment is an extension of the violence that is done to women
37 (Lindsay et al., 2016). Several scholars have called for activism to tackle this problem (e.g.,
38 Turley & Fisher, 2018; Jane, 2014). We therefore designed two straightforward tweets based
39 on principles of social re-norming and empathy and tested these on a sample of 666 Twitter
40 users. Our intervention tweets did not, regrettably, reduce the number of sexist slurs or sexist
41 users in our sample. They also did not affect the valence or arousal of subsequent tweets.
42 Disappointing, but perhaps not altogether surprisingly given how prolific sexist slurs are on
43 social media and how normalised online gender harassment has become. We add our voices to
44 calls for further activism.
45
46
47
48
49
50
51
52

53 REFERENCES

54
55
56 Action Aid. (2016). *Three in four women experience harassment and violence in UK and*
57 *global cities*. [https://www.actionaid.org.uk/latest-news/three-in-four-women-experience-](https://www.actionaid.org.uk/latest-news/three-in-four-women-experience-harassment-and-violence)
58 [harassment-and-violence](https://www.actionaid.org.uk/latest-news/three-in-four-women-experience-harassment-and-violence)
59
60

- 1
2
3 Amnesty International. (2017). *More than a quarter of UK women experiencing online abuse*
4 *and harassment receive threats of physical or sexual assault - new research.*
5
6 [https://www.amnesty.org.uk/press-releases/more-quarter-uk-women-experiencing-](https://www.amnesty.org.uk/press-releases/more-quarter-uk-women-experiencing-online-abuse-and-harassment-receive-threats)
7 [online-abuse-and-harassment-receive-threats](https://www.amnesty.org.uk/press-releases/more-quarter-uk-women-experiencing-online-abuse-and-harassment-receive-threats)
8
9
10 Amnesty International. (2020). *Violence against women.*
11
12 <https://www.amnesty.org.uk/violence-against-women>
13
14 Barbera, P. (2014). *streamR: Access to Twitter Streaming API via R. R package version 0.2.1.*
15
16 <https://cran.r-project.org/package=streamR>
17
18 Batson, C. D., Chang, J., Orr, R., & Rowland, J. (2002). Empathy, attitudes, and action: Can
19 feeling for a member of a stigmatized group motivate one to help the group? *Personality*
20 *and Social Psychology Bulletin*, 28(12), 1656–1666.
21
22 <https://doi.org/10.1177/014616702237647>
23
24 Chadha, K., Steiner, L., Vitak, J., & Ashktorab, Z. (2020). Women’s Responses to Online
25 Harassment. *International Journal of Communication*, 14(0), 19.
26
27 Chen, G. M., Pain, P., Chen, V. Y., Mekelburg, M., Springer, N., & Troger, F. (2020). ‘You
28 really have to have a thick skin’: A cross-cultural perspective on how online harassment
29 influences female journalists. *Journalism*, 21(7), 877–895.
30
31 <https://doi.org/10.1177/1464884918768500>
32
33
34 Cho, D., & Acquisti, A. (2013). The More Social Cues , The Less Trolling? An Empirical
35 Study of Online Commenting Behavior. *The Twelfth Workshop on the Economics of*
36 *Information Security, Weis.*
37
38
39 Citizens Advice. (2021). *Sexual Harassment.* [https://www.citizensadvice.org.uk/law-and-](https://www.citizensadvice.org.uk/law-and-courts/discrimination/what-are-the-different-types-of-discrimination/sexual-harassment/)
40 [courts/discrimination/what-are-the-different-types-of-discrimination/sexual-harassment/](https://www.citizensadvice.org.uk/law-and-courts/discrimination/what-are-the-different-types-of-discrimination/sexual-harassment/)
41
42
43 Citron, D. (2014). *Hate crimes in cyberspace.* Harvard University Press.
44
45 Cote, A. C. (2017). “I Can Defend Myself”: Women’s Strategies for Coping with Harassment
46 while Gaming Online. *Games and Culture*, 12(2), 136–155.
47
48 <https://doi.org/10.1177/1555412015587603>
49
50 Dovidio, J. F., Ten Vergert, M., Stewart, T. L., Gaertner, S. L., Johnson, J. D., Esses, V. M.,
51 Riek, B. M., & Pearson, A. R. (2004). Perspective and prejudice: Antecedents and
52 mediating mechanisms. *Personality and Social Psychology Bulletin*, 30(12), 1537–1549.
53
54 <https://doi.org/10.1177/0146167204271177>
55
56
57 Drakett, J., Rickett, B., Day, K., & Milnes, K. (2018). Old jokes, new media -online sexism
58 and constructions of gender in internet memes. *Feminism and Psychology*, 28(1), 109–
59 127. <https://doi.org/10.1177/0959353517727560>
60

- 1
2
3 Felmlee, D., Inara Rodis, P., & Zhang, A. (2020). Sexist Slurs: Reinforcing Feminine
4 Stereotypes Online. *Sex Roles*, 83(1–2), 16–28. [https://doi.org/10.1007/s11199-019-](https://doi.org/10.1007/s11199-019-01095-z)
5
6 01095-z
7
8 Ferrier, M., & Garud-Patkar, N. (2018). TrollBusters: Fighting Online Harassment of Women
9 Journalists. In J. Vickery & T. Everback (Eds.), *Mediating Misogyny* (Issue February,
10 pp. 311–332). Palgrave Macmillan. <https://doi.org/10.1007/978-3-319-72917-6>
11
12 Galinsky, A., & Moskowitz, G. (2000). Perspective-Taking: Decreasing Stereotype
13 Expression, Stereotype Accessibility, and In-Group Favoritism. *Journal of Personality*
14 and *Social Psychology*, 78(4), 708–724. <https://doi.org/10.1037//0022-3514.78.4.708>
15
16 Goldstein, N. J., Cialdini, R. B., & Griskevicius, V. (2008). A Room with a Viewpoint: Using
17 Social Norms to Motivate Environmental Conservation in Hotels. *Journal of Consumer*
18 *Research*, 35(3), 472–482. <https://doi.org/10.1086/586910>
19
20 Han, X. (2018). Searching for an online space for feminism? The Chinese feminist group
21 Gender Watch Women’s Voice and its changing approaches to online misogyny.
22 *Feminist Media Studies*, 18(4), 734–749.
23
24 <https://doi.org/10.1080/14680777.2018.1447430>
25
26 Herring, S. (2003). Gender and power in online communication. In *The Handbook of*
27 *Language and Gender* (pp. 202–228).
28
29 HM Crown Prosecution Service Inspectorate. (2019). *2019 Rape Inspection*. December.
30
31 Hoskin, R. A. (2019). Femmephobia: The Role of Anti-Femininity and Gender Policing in
32 LGBTQ+ People’s Experiences of Discrimination. *Sex Roles*, 81(11–12), 686–703.
33
34 <https://doi.org/10.1007/s11199-019-01021-3>
35
36 House of Commons. (2017). *Online harassment and cyber bullying* (Issue 07967).
37
38 Jackson, S. (2018). Young feminists, feminism and digital media. *Feminism & Psychology*,
39 28(1), 32–49.
40
41 Jane, E. (2014). Back to the kitchen, cunt: Speaking the unspeakable about online misogyny.
42 In *Continuum* (Vol. 28, Issue 4, pp. 558–570). Taylor & Francis.
43
44 <https://doi.org/10.1080/10304312.2014.924479>
45
46 Jane, E. (2018). Gendered cyberhate as workplace harassment and economic vandalism.
47 *Feminist Media Studies*, 18(4), 575–591.
48
49 <https://doi.org/10.1080/14680777.2018.1447344>
50
51 Jane, T. (2019). *Creepy men slide into women’s DMs all the time, but they can be shut down*.
52 The Guardian. [https://www.theguardian.com/commentisfree/2019/may/07/creepy-men-](https://www.theguardian.com/commentisfree/2019/may/07/creepy-men-dm-online-harassment)
53
54 dm-online-harassment
55
56
57
58
59
60

- 1
2
3 Kilmartin, C., Smith, T., Green, A., Heinzen, H., Kuchler, M., & Kolar, D. (2008). A real
4 time social norms intervention to reduce male sexism. *Sex Roles*, *59*(3–4), 264–273.
5 <https://doi.org/10.1007/s11199-008-9446-y>
6
7
8 Lapidot-Lefler, N., & Barak, A. (2012). Effects of anonymity, invisibility, and lack of eye-
9 contact on toxic online disinhibition. *Computers in Human Behavior*, *28*(2), 434–443.
10 <https://doi.org/10.1016/j.chb.2011.10.014>
11
12
13 Lindsay, M., Booth, J. M., Messing, J. T., & Thaller, J. (2016). Experiences of Online
14 Harassment Among Emerging Adults: Emotional Reactions and the Mediating Role of
15 Fear. *Journal of Interpersonal Violence*, *31*(19), 3174–3195.
16 <https://doi.org/10.1177/0886260515584344>
17
18
19 Locke, A., Lawthom, R., & Lyons, A. (2018). Social media platforms as complex and
20 contradictory spaces for feminisms: Visibility, opportunity, power, resistance and
21 activism. *Feminism and Psychology*, *28*(1), 3–10.
22 <https://doi.org/10.1177/0959353517753973>
23
24
25 Mallett, K. A., Bachrach, R. L., & Turrisi, R. (2008). Are all negative consequences truly
26 negative? Assessing variations among college students' perceptions of alcohol related
27 consequences. *Addictive Behaviors*, *33*(10), 1375–1381.
28
29
30 Megarry, J. (2014). Online incivility or sexual harassment? Conceptualising women's
31 experiences in the digital age. *Women's Studies International Forum*, *47*(PA), 46–55.
32 <https://doi.org/10.1016/j.wsif.2014.07.012>
33
34
35 Met Police. (2021). *I'm being harassed by someone on social media. What can I do?*
36 [https://www.met.police.uk/advice/advice-and-information/har/harassment-on-social-
37 media/#:~:text=You can report either harassment,by calling us on 101.](https://www.met.police.uk/advice/advice-and-information/har/harassment-on-social-media/#:~:text=You can report either harassment,by calling us on 101.)
38
39
40
41
42
43 Munger, K. (2016). Tweetment Effects on the Tweeted: Experimentally Reducing Racist
44 Harassment. *Political Behavior*, 1–21. <https://doi.org/10.1007/s11109-016-9373-5>
45
46
47 Matias, N. (2019). Preventing harassment and increasing group participation through social
48 norms in 2,190 online science discussions. *Proceedings of the National Academy of
49 Sciences of the United States of America*, *116*(20), 9785–9789.
50 <https://doi.org/10.1073/pnas.1813486116>
51
52
53 Office of National Statistics. (2018). *Sexual offences in England and Wales: year ending
54 March 2017.*
55 [https://www.ons.gov.uk/peoplepopulationandcommunity/crimeandjustice/articles/sexual
56 offencesinenglandandwales/yearendingmarch2017](https://www.ons.gov.uk/peoplepopulationandcommunity/crimeandjustice/articles/sexual-offencesinenglandandwales/yearendingmarch2017)
57
58
59 Office of National Statistics. (2021). *Homicide in England and Wales: year ending March*
60

2020.
<https://www.ons.gov.uk/peoplepopulationandcommunity/crimeandjustice/articles/homicideinenglandandwales/yearendingmarch2020>
- Paluck, E. L., & Green, D. P. (2009). Prejudice reduction: what works? A review and assessment of research and practice. *Annual Review of Psychology*, *60*, 339–367.
<https://doi.org/10.1146/annurev.psych.60.110707.163607>
- Pennycook, G., Epstein, Z., Mosleh, M. Arechar, A., Eckles, D., & Rand, D. (2021) Shifting attention to accuracy can reduce misinformation online. *Nature* *592*, 590–595.
<https://doi.org/10.1038/s41586-021-03344-2>
- Rights of Women. (2021). *Rights of Women survey reveals online sexual harassment has increased, as women continue to suffer sexual harassment whilst working through the Covid-19 pandemic*. <https://rightsofwomen.org.uk/news/rights-of-women-survey-reveals-online-sexual-harassment-has-increased-as-women-continue-to-suffer-sexual-harassment-whilst-working-through-the-covid-19-pandemic/>
- Schmidt, A., & Wiegand, M. (2017). *A Survey on Hate Speech Detection using Natural Language Processing*. *2012*, 1–10. <https://doi.org/10.18653/v1/w17-1101>
- Schultz, P. W., Nolan, J. M., Cialdini, R. B., Goldstein, N. J., & Griskevicius, V. (2018). The Constructive, Destructive, and Reconstructive Power of Social Norms: Reprise. *Perspectives on Psychological Science*, *13*(2), 249–254.
<https://doi.org/10.1177/1745691617693325>
- Schwartz, H. A., & Ungar, L. H. (2015). Data-Driven Content Analysis of Social Media: A Systematic Overview of Automated Methods. *Annals of the American Academy of Political and Social Science*, *659*(1), 78–94. <https://doi.org/10.1177/0002716215569197>
- Stubbs-Richardson, M. S., Rader, N. E., & Cosby, A. G. (2018). Tweeting rape culture: Examining portrayals of victim blaming in discussions of sexual assault cases on Twitter. *Feminism & Psychology*, *28*(1), 90–108.
- Thompson, L. (2018). ‘I can be your Tinder nightmare’: Harassment and misogyny in the online sexual marketplace. *Feminism & Psychology*, *28*(1), 69–89.
- Turley, E., & Fisher, J. (2018). Tweeting back while shouting back: Social media and feminist activism. *Feminism & Psychology*, *28*(1), 128–132.
- Uhl, C., Rhyner, K., & Lugo, N. (2018). An examination of nonconsensual pornography websites. *Feminism & Psychology*, *28*(1), 50–68.
- UN Women UK. (2021). *Prevalence and reporting of sexual harassment in UK public spaces - A report by the APPG for UN Women*. March, 1–28.

- 1
2
3 Van Hee, C., Jacobs, G., Emmery, C., Desmet, B., Lefever, E., Verhoeven, B., De Pauw, G.,
4 Daelemans, W., & Hoste, V. (2018). Automatic detection of cyberbullying in social
5 media text. *ArXiv*, 1–22. <https://doi.org/10.17605/OSF.IO/RGQW8>.
6
7
8 Vera-Gray, F. (2017). Talk about a cunt with too much idle time’: Trolling feminist research.
9 *Feminist Review*, 115(1), 61–78. <https://doi.org/10.1057/s41305-017-0038-y>
10
11
12 Vescio, T. K., Sechrist, G. B., & Paolucci, M. P. (2003). Perspective taking and prejudice
13 reduction: The mediational role of empathy arousal and situational attributions.
14 *European Journal of Social Psychology*, 33(4), 455–472.
15
16 <https://doi.org/10.1002/ejsp.163>
17
18
19 Wang, W., Chen, L., Thirunarayan, K., & Sheth, A. P. (2014). Cursing in English on twitter.
20 *Proceedings of the 17th ACM Conference on Computer Supported Cooperative Work &*
21 *Social Computing - CSCW '14*, 415–425. <https://doi.org/10.1145/2531602.2531734>
22
23
24 Warriner, A. B., Kuperman, V., & Brysbaert, M. (2013). Norms of valence, arousal, and
25 dominance for 13,915 English lemmas. *Behavior Research Methods*, 45(4), 1191-1207.
26
27
28 Zimmerman, S., Kruschwitz, U., & Fox, C. (2018, May). Improving hate speech detection
29 with deep learning ensembles. In Proceedings of the Eleventh International Conference
30 on Language Resources and Evaluation (LREC 2018).
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Appendix I – List of commonly used sexist slurs

TERM	INCLUDE
arm candy	1
asking for it/asked for it	1
ball-breaker	1
ballbuster	1
battle axe	1
bimbo	1
bimbo	1
bint	1
bitch	0
bridezilla	1
bunny boiler	1
butch	1
butterface	1
catfight	1
chavette/girl chav?	1
cock tease	1
cougar	0
crank whore	1
crocadillapig	1
crone	1
cunt	0
daft bimbo	1
daft bitch	1
daft cow	1
daft cunt	1
damaged good	1
ditz	1
dizty	1
essex girl	1
fag hag	1

TERM	INCLUDE
frigid bitch	1
frump	1
frumpy	1
fucking bimbo	1
fucking bitch	1
fucking cunt	1
gagging for it	1
ghetto bird	1
ghetto ho	1
gold digger	1
harridan	1
hoe	0
hooch	1
hoochie	1
hussy	1
huzzie	1
milf	0
MILF	0
minger	1
moll	1
moose	1
mousey	1
old bag	1
pass around pussy	1
poon	1
poontang	1
prostitute	1
prude	1
pussy	0
sausage jockey	1

feminazi	1
flange	1
flipper	1
floozy	1
floozy	1
frigid	1
stupid bimbo	1
tart	1
town bike	1
tramp	1
troglodyte	1
trollop	1
vamp	1
village bicycle	1
what's-her-face	1
whatshername	1
whore	0

shrew	1
skank	1
skeezy ho	1
slag	1
slapper	1
sleaze	1

1 = include; 0 = do not include.