**Supplementary Table 1**. *Sensitive Information Categories and Relevant Regular Expressions*

| Sensitive Information | Subcategory | Regular expression patterns | Explanations |
|---|---|---|---|
| **(A) Names** | Firstname, lastame | ^[A-Z]\'?[-a-zA-Z]+$ | match names, A'Bsfs, Absssfs, A-Bsfsfs |
| | Title | (Dr\.\|Mr\.\|Mrs\.\|Ms\.\|Miss\|Sir\|Madam)\s (([A-Z]\'?[A-Z]?[\-a-z]+(\s[A-Z]\'?[A-Z]?[\-a-z]+)*) ) | match salutation |
| | Middle name | \*\*PHI\*\*,? (([A-CE-LN-Z][Rr]?\|[DM])\.?) \| (([A-CE-LN-Z][Rr]?\|[DM])\.?),? \*\*PHI\*\* | match middle initial |
| **(B) Geographic subdivisions** | Postal code | \b(\d{5}(-\d{4})?)\b | match postal code |
| | Address | address_indictor = ['street', 'avenue', 'road', 'boulevard', 'drive', 'trail', 'way', 'lane', 'ave', 'blvd', 'st', 'rd', 'trl', 'wy', 'ln', 'court', 'ct', 'place', 'plc', 'terrace', 'ter', 'highway', 'freeway', 'autoroute', 'autobahn', 'expressway', 'autostrasse', 'autostrada', 'byway', 'auto-estrada', 'motorway', 'avenue', 'boulevard', 'road', 'street', 'alley', 'bay', 'drive', 'gardens', 'gate', 'grove', 'heights', 'highlands', 'lane', 'mews', 'pathway', 'terrace', 'trail', 'vale', 'view', 'walk', 'way', 'close', 'court', 'place', 'cove', 'circle', 'crescent', 'square', 'loop', 'hill', 'causeway', 'canyon', 'parkway', 'esplanade', 'approach', 'parade', 'park', 'plaza', 'promenade', 'quay', 'bypass' 'dr.', 'ave.', 'rd.', 'st.', 'blvd.','pkwy.','city'] | address keywords |
| **(C) Dates** | General date | \b(.*?(?=\b(\d{1,2}[-./\s]\d{1,2}[-./\s]\d{2}\|\d{1,2}[-./\s]\d{1,2}[-./\s]\d{4}\|\d{2}[-./\s]\d{1,2}[-./\s]\d{1,2}\|\d{4}[-./\s]\d{1,2}[-./\s]\d{1,2})\b))\b<br><br>month_name = Jan(uary)?\|Feb(ruary)?\|Mar(ch)?\|Apr(il)?\|May\|Jun(e)?\|Jul(y)?\|Aug(ust)?\|Sep(tember)?\|Oct(ober)?\|Nov(ember)?\|Dec(ember)?<br><br>\b( | Date in various formats. YYYY/MM-YYYY/MM MM/YYYY-MM/YYYY MM/YY-MM/YY MM/YY-MM/YY MM/YYYY-MM/YYYY MM/DD-MM/DD DD/MM-DD/MM DD/MM-DD/MM MM/DD/YY MM/DD/YYYY |

| | | | |
|---|---|---|---|
| | | \d{4}[\-/](0?[1-9]\|1[0-2]\|"""+month_name+r""")\-\d{4}[\-/](0?[1-9]\|1[0-2]\|"""+month_name+r""") # YYYY/MM-YYYY/MM<br>\|(0?[1-9]\|1[0-2]\|"""+month_name+r""")[\-/]\d{4}\-(0?[1-9]\|1[0-2]\|"""+month_name+r""")[\-/]\d{4} # MM/YYYY-MM/YYYY<br>\|(0?[1-9]\|1[0-2]\|"""+month_name+r""")/\d{2}\-(0?[1-9]\|1[0-2]\|"""+month_name+r""")/\d{2} # MM/YY-MM/YY<br>\|(0?[1-9]\|1[0-2]\|"""+month_name+r""")/\d{2}\-(0?[1-9]\|1[0-2]\|"""+month_name+r""")/\d{4} # MM/YYYY-MM/YYYY<br>\|(0?[1-9]\|1[0-2]\|"""+month_name+r""")/([1-2][0-9]\|3[0-1]\|0?[1-9])\-(0?[1-9]\|1[0-2]\|"""+month_name+r""")/([1-2][0-9]\|3[0-1]\|0?[1-9]) #MM/DD-MM/DD<br>\|([1-2][0-9]\|3[0-1]\|0?[1-9])/(0?[1-9]\|1[0-2]\|"""+month_name+r""")\-([1-2][0-9]\|3[0-1]\|0?[1-9])/(0?[1-9]\|1[0-2]\|"""+month_name+r""") #DD/MM-DD/MM<br>\|(0?[1-9]\|1[0-2]\|"""+month_name+r""")[\-/\s]([1-2][0-9]\|3[0-1]\|0?[1-9])[\-/\s]\d{2} # MM/DD/YY<br>\|(0?[1-9]\|1[0-2]\|"""+month_name+r""")[\-/\s]([1-2][0-9]\|3[0-1]\|0?[1-9])[\-/\s]\d{4} # MM/DD/YYYY<br>\|([1-2][0-9]\|3[0-1]\|0?[1-9])[\-/\s](0?[1-9]\|1[0-2]\|"""+month_name+r""")[\-/\s]\d{2} # DD/MM/YY<br>\|([1-2][0-9]\|3[0-1]\|0?[1-9])[\-/\s](0?[1-9]\|1[0-2]\|"""+month_name+r""")[\-/\s]\d{4} # DD/MM/YYYY<br>\|\d{2}[\-./\s](0?[1-9]\|1[0-2]\|"""+month_name+r""")[\-\./\s]([1-2][0-9]\|3[0-1]\|0?[1-9]) # YY/MM/DD<br>\|\d{4}[\-./\s](0?[1-9]\|1[0-2]\|"""+month_name+r""")[\-\./\s]([1-2][0-9]\|3[0-1]\|0?[1-9]) # YYYY/MM/DD<br>\|\d{4}[\-/](0?[1-9]\|1[0-2]\|"""+month_name+r""") # YYYY/MM<br>\|(0?[1-9]\|1[0-2]\|"""+month_name+r""")[\-/]\d{4} # MM/YYYY | DD/MM/YY<br>DD/MM/YYYY<br>YY/MM/DD<br>YYYY/MM/DD<br>MM/YYYY<br>MM/YY<br>MM/YYYY<br>MM/DD<br>DD/MM |

| | | | |
|---|---|---|---|
| | | \|(0?[1-9]\|1[0-2]\|"""+month_name+r""")/\d{2}  # MM/YY<br>\|(0?[1-9]\|1[0-2]\|"""+month_name+r""")/\d{2}  # MM/YYYY<br>\|(0?[1-9]\|1[0-2]\|"""+month_name+r""")/([1-2][0-9]\|3[0-1]\|0?[1-9])  #MM/DD<br>\|([1-2][0-9]\|3[0-1]\|0?[1-9])/(0?[1-9]\|1[0-2]\|"""+month_name+r""")  #DD/MM<br>)\b | |
| | DOB | \b(.*?(?=\b(\d{1,2}[-./\s]\d{1,2}[-./\s]\d{2}  # X/X/XX<br>\|\d{1,2}[-./\s]\d{1,2}[-./\s]\d{4}         # XX/XX/XXXX<br>\|\d{2}[-./\s]\d{1,2}[-./\s]\d{1,2}         # xx/xx/xx<br>\|\d{4}[-./\s]\d{1,2}[-./\s]\d{1,2}         # xxxx/xx/xx<br>)\b)<br>)\b | Match date of birth |
| | Age | \b(<br>age\|year[s-]?\s?old\|y.o[.]?<br>)\b | Match age |
| **(D)Telephone/ FAX numbers, (I) Health plan beneficiary numbers, Account numbers, Medical record numbers** | | \b((\d[\(\)\-\']?\s?){6}([\(\)\-\']?\d)+<br>\|(\d[\(\)\-.\']?){7}([\(\)\-.\']?\d)+  # test<br>)\b<br>\b(\d{5}[A-Z0-9]*)\b<br>\b([A-Z0-9\-/]{6}[A-Z0-9\-/]*)\b | Phone/fax/account number/MRN, etc. in the format of SSN/PHONE/FAX XXX-XX-XXXX, XXX-XXX-XXXX, XXX-XXXXXXXX, etc. |
| **(F)  Email addresses** | | \b([a-zA-Z0-9_.+-@\"]+@[a-zA-Z0-9-\:\]\[]+[a-zA-Z0-9-.]*)\b | e.g., xxx.xxx@xxxx.xxx |
| **(G)  Social security numbers** | | \b((\d[\(\)\-\']?\s?){6}([\(\)\-\']?\d)+\|(\d[\(\)\-.\']?){7}([\(\)\-.\']?\d)+)\b | e.g., xxx-xx-xxxx, etc. |
| **(N)  Web URLs** | | \b((http[s]?://)?([a-zA-Z0-9$-_@.&+:!\*\(\),])*[\.\/]([a-zA-Z0-9$-_@.&+:\!\*\(\),])*)\b | e.g., http://xxx.xxx.xxx/xxx.xxx |

| (O) Internet Protocol (IP) addresses | | ^(?:(?:25[0-5]|2[0-4][0-9]|[01]?[0-9][0-9]?)\.){3}(?:25[0-5]|2[0-4][0-9]|[01]?[0-9][0-9]?)$ | e.g., xxx.xxx.xxx.xxx |
|---|---|---|---|
| (R) Other unique identifiers | Race | "(White/Caucasian|White|Caucasian|American|Indian|Alaska|Native|Asian|Black|African|Native|Hawaiian|Pacific|Islander)" | Match race |
| | Gender | "(Male|Female)" | Match gender |
| | Ethnicity | "(Hispanic|Latino)" | Match ethnicity |