



OULUN YLIOPISTO
UNIVERSITY of OULU

Roskapostin tunnistaminen koneoppimisen avulla sosiaalisessa mediassa

Oulun yliopisto
Tieto- ja sähkötekniikan tiedekunta
LuK-tutkielma
Nuutti Kinnunen
30.5.2022

Tiivistelmä

Tämä työn tarkoitus on selvittää, miten koneoppimista hyödynnetään suodattamaan roskapostia sosiaalisesta mediasta. Tämän lisäksi koneoppimista vertaillaan muihin tapoihin suodattaa roskapostia. Aihe on tärkeä, koska lähivuosina roskapostista on tullut suuri ongelma sosiaalisen median alustoille. Roskapostin tunnistamiseen manuaalisesti liittyy kuitenkin haittoja, joita ovat suurella sosiaalisen median alustalla suuret kulut sekä epäkäytännöllisyys suuren viestimäärän tarkastamiseen.

Työ suoritettiin kirjallisuuskatsauksena. Aiempien tutkimusten perusteella koneoppimista voidaan hyödyntää tämän ongelman lieventämiseen. Koneoppimisen avulla roskapostia pystytään suodattamaan automatisoidusti ilman että tarvitsee tehdä monimutkaisia käsin kirjoitettuja sääntöjä. Koneoppimisalgoritmeina voidaan käyttää esimerkiksi Naive Bayesia ja neuroverkkoja. Työssä käsitellyn aiemman tutkimuksen mukaan Naive Bayes suoriutuu roskapostin suodattamisesta kokonaisuudessa neuroverkkoja paremmin. Työ tarjoaa yleisen katsauksen aiheeseen.

Avainsanat

Tekoäly, koneoppiminen, roskaposti

Ohjaaja

Väitöskirjatutkija Leevi Rantala

Sisällysluettelo

Tiivistelmä.....	2
Sisällysluettelo	3
1. Johdanto.....	4
2. Keskeiset käsitteet	5
2.1 Tekoäly.....	5
2.2 Koneoppiminen.....	5
2.3 Roskaposti.....	5
3. Tutkimusmenetelmät	6
4. Sosiaalinen media	7
4.1 Erilaisia tapoja roskapostin suodattamiseen	7
4.2 Automatisoitu roskapostin suodattaminen	8
4.3 Automatisoidun roskapostin suodattamisen haasteita	8
5. Koneoppimistekniikat.....	10
5.1 Naive Bayes luokittelu	10
5.2 Neuroverkot	11
6. Koneoppimisen käyttäminen roskapostin tunnistamiseen sosiaalisessa mediassa	12
6.1 Naive Bayes luokittelun soveltaminen käytännössä	12
6.2 Naive Bayesin ja neuroverkkojen suorituskyky	14
7. Pohdinta	16
8. Yhteenveto.....	18
Lähteet.....	20

1. Johdanto

Roskaposti on vaivannut sähköpostia jo internetin alkuaajoista lähtien. Ensimmäinen digitaalisesti tehty roskapostittaminen tapahtui vuonna 1978 ARPANET verkossa, kun Digital Equipment Corporation mainosti uutta tietokonettaan yli 400 verkon jäsenelle. Roskapostittaminen yleistyi 1990-luvulla sähköpostin yleistymisen myötä. (Ferrara, 2019)

Roskapostia lähetetään pääasiassa kahden eri syyn takia, joita ovat mainostaminen ja huijaaminen. Mainostaminen voi olla esimerkiksi jonkin tuotteen tai palvelun mainostamista, ja huijaaminen voi olla esimerkiksi tietojenkalastelua. (Ferrara, 2019)

Roskaposti on yksi internetin isoimmista ongelmista, ja se aiheuttaa taloudellista vahinkoa yrityksille sekä ärsyttää monia internetin käyttäjiä. Roskaposti aiheuttaa monia vakavia ongelmia, joista monet johtavat suoraan taloudellisiin menetyksiin. (Blanzieri & Bryl, 2008)

Lähivuosina sosiaalisen median suosio on kasvanut, ja tämä on houkuttanut myös roskapostin lähettäjiä sosiaaliseen mediaan. Roskapostin lähettäminen sosiaalisessa mediassa voi olla esimerkiksi linkkien lähettämistä johonkin kaupalliseen tuotteeseen. (Chakraborty ym., 2016) Ongelma pahenee vuosi vuodelta, koska roskapostittajat kehittävät aina vain tehokkaampia tapoja roskapostin lähettämiseen (Ferrara, 2019). Tämän takia internetissä toimivien sosiaalisen median palveluiden ylläpitäjillä on tarvetta kehittää myös tehokkaampia menetelmiä roskapostin tunnistamiseen.

Tutkimusongelmana tässä työssä on *roskapostin tunnistaminen koneoppimisen avulla sosiaalisessa mediassa*. Tähän vastataan tutkimuskysymyksen avulla, joka on ”*Miten roskapostia voi tunnistaa koneoppimisen avulla sosiaalisessa mediassa, ja mitä etuja ja haittoja sillä on verrattuna muihin menetelmiin?*”. Tutkimusmenetelmänä toimii kirjallisuuskatsaus, eli omaa empiiristä tutkimusta ei tehdä, vaan tässä työssä perehdytään aiempaan tutkimukseen. Tämä työ perustuu vuonna 2021 tekemääni JTT-tutkimukseen.

2. Keskeiset käsitteet

Tässä luvussa esittelen tämän tutkimuksen kannalta keskeiset käsitteet. Keskeisiä käsitteitä tässä tutkimuksessa ovat tekoäly, koneoppiminen ja roskaposti.

2.1 Tekoäly

Tekoäly tarkoittaa ohjelmaa, jossa katsotaan olevan älykkyyttä. Tekoälyä on vaikea määritellä, koska älykkyydelle ei ole täysin selvää määritelmää, eikä tekoäly suurelta osin muistuta ihmisen älykkyyttä (Kaplan, 2016). Koneoppiminen on osa tekoälyn käsitettä (Bi ym., 2019). Yksi esimerkki tekoälyn algoritmeista on neuroverkot (Kukreja ym., 2016).

2.2 Koneoppiminen

Koneoppiminen on tekoälyn osa-alue. Sillä tarkoitetaan ohjelmia, jotka pystyvät oppimaan vastaanottamastaan tiedosta ja kehittymään ajan kuluessa. Koneoppimisalgoritmit pystyvät muuttamaan käyttäytymistään parantaakseen suorituskykyään jossain tehtävässä. (Bi ym., 2019) Koneoppimisalgoritmeja ovat esimerkiksi Naive Bayes ja neuroverkot (Jain ym., 2019).

2.3 Roskaposti

Roskaposti tarkoittaa ilman vastaanottajien lupaa lähetettyä massapostitusta. Roskapostia voidaan lähettää esimerkiksi sähköpostissa tai sosiaalisessa mediassa. (Castillo, 2006) Roskapostia voidaan lähettää myös pikaviestimissä, hakukoneissa, wiki -palveluissa ja arvostelusivustoilla. Roskapostin tarkoitus on yleensä mainostaminen tai huijaaminen. Mainostaminen voi olla jonkin kaupallisen palvelun tai tuotteen mainostamista, ja huijaaminen voi olla esimerkiksi tietojenkalastelua. (Ferrara, 2019)

3. Tutkimusmenetelmät

Tässä tutkielmassa tutkimusmenetelmänä toimii kirjallisuuskatsaus, eli omaa empiiristä tutkimusta ei toteuteta. Tutkielmassa perehdytään aiempaan tutkimukseen aiheesta, ja kerätään niistä tutkimustuloksia. Olen valinnut vain lähteitä, jotka liittyvät koneoppimiseen tai roskapostiin sosiaalisessa mediassa.

Hakukoneena käytin suurimmaksi osaksi Google Scholar-hakukonetta, jonka lisäksi käytin Scopus tietokantaa. Hakulausekkeissa käytin pääosin termejä ”social media”, ”spam”, ”naive bayes”, ”bag of words”, ”artificial neural network”, ”moderation” ja ”bayes theorem”. Osa lähteistä on löydetty hakusanoilla löytämäni artikkelien lähdeluettelosta.

Analysoin Google Scholarista ja Scopuksesta löytämäni aineiston seuraavalla tavalla. Ensin kävin läpi artikkelien otsikkoja ja arvioin niiden sopivuutta suhteessa tutkimuskysymykseeni. Relevanteiksi tunnistamani artikkelit etenivät seuraavaan vaiheeseen, jossa tutkin niiden tiivistelmät. Tämän jälkeen päätin, sopiiko kyseinen artikkeli lähteeksi tutkimukseeni. Artikkelien lukemisessa keskityin tutkielman kannalta oleellisimpiin osiin.

4. Sosiaalinen media

Sosiaalinen media viittaa keskustelulliseen ja hajautettuun sisällön tuottamiseen, levittämiseen ja kommunikointiin yhteisöjen välillä. Sosiaalinen media erottuu perinteisestä mediasta siten, että kirjoittajan ja lukijan välistä eroa ei ole. (Zeng ym., 2010) Sosiaalisella medialla on nykyään paljon käyttäjiä. Facebookilla oli vuoden 2021 lopussa yli 1,9 miljardia päivittäistä käyttäjää (Hamilton, 2022). Twitterillä taas oli 217 miljoonaa aktiivista päivittäistä käyttäjää vuoden 2021 lopussa (Statista, 2022). Suuri osa ihmisistä tekee päätöksiä sosiaalisen median avulla esimerkiksi tuotteiden arvostelujen ja palautteen perusteella. Kuka tahansa pystyy kirjoittamaan sosiaalisessa mediassa arvosteluja ja antamaan palautetta, mikä tekee siitä hyvän alustan roskapostin lähettämiseen. (Shehnepoor ym., 2017)

Sosiaalisessa mediassa arviolta yksi jokaisesta 200 viestistä ja yksi jokaisesta 21 tweetistä on roskapostia (Inuwa-Dutse ym., 2018). Kaksi pääsyytä minkä takia roskapostia lähetetään, mainostaminen ja huijaaminen pätee myös sosiaaliseen mediaan. Roskapostin lähettäjät käyttävät yleensä linkkejä ihmisten huijaamiseksi esimerkiksi tietojenkalastelu sivulle (Cao ym., 2015). Noin 25 % sosiaalisen median viesteistä sisältää linkkejä, joka tekee haitallisten linkkien löytämisen niiden joukosta vaikeaa (Cao ym., 2015).

Sosiaalisen median erikoispiirteenä on poliittinen vaikuttaminen roskapostin avulla. Roskapostilla sosiaalisessa mediassa voidaan vaikuttaa ihmisten poliittisiin mielipiteisiin ja vaaleihin. Esimerkiksi vuonna 2010 Massachusetts osavaltion vaalien alla, roskapostittajat tekivät 9 käyttäjää Twitteriin ja tavoittivat niillä yli 60 000 käyttäjää päivässä, käytännössä ilmaiseksi (Metaxas ym., 2012).

Roskapostin tunnistaminen sosiaalisessa mediassa eroaa hieman sähköpostista. Tämä johtuu siitä, että esimerkiksi Twitterissä viestit ovat yleensä lyhyitä ja ne vaihtelevat paljon (Jain ym., 2019). Sosiaalisessa mediassa käyttäjien luoma sisältö voi sisältää esimerkiksi tekstiä, kuvia tai videoita (Almadhoor, 2021). Kaikki näistä voivat sisältää roskapostia, joka tekee roskapostin suodattamisesta haastavampaa (Almadhoor, 2021). Seuraavaksi esittelen erilaisia tapoja roskapostin suodattamiseen, jonka jälkeen tarkastelemme automatisoitua roskapostin suodatusta ja siihen liittyviä haasteita.

4.1 Erilaisia tapoja roskapostin suodattamiseen

Sosiaalisen median alustoilla on roskapostin suodattamiseen neljä erilaista vaihtoehtoa. Näistä ensimmäinen on, että käyttäjien luoma sisältö tarkistetaan ennen kuin se tulee muiden käyttäjien nähtäville. Toinen tapa on, että käyttäjien luoma sisältö tarkistetaan sen jälkeen, kun se on tullut käyttäjien nähtäville. Kolmas tapa on hajautettu käyttäjien luoman sisällön tarkastaminen ja neljäs tapa on automatisoitu käyttäjien luoman sisällön tarkastaminen (Veglis, 2014).

Sisällön tarkastaminen ennen kuin se tulee käyttäjien nähtäville tarkoittaa sitä, että palvelun moderaattori tarkastaa sisällön, ennen kuin se tulee käyttäjien nähtäville. Tämä tapa tarjoaa parhaan kontrollin sisällöstä. Tätä tapaa kannattaa käyttää, jos on tärkeää, että roskapostia ei tule ollenkaan käyttäjien nähtäville. Tämän tavan käyttö kuitenkin johtaa siihen, että käyttäjien luoman sisällön määrä tippuu 40–50 %. Käyttäjä ei

myöskään saa välitöntä tyydytystä viestinsä lähettämisestä, vaan jää odottamaan moderaattorin hyväksyntää. Tämän tavan käytöstä voi koitua myös suuria kustannuksia, jos käyttäjien luomaa sisältöä tulee paljon. (Veglis, 2014)

Sisällön tarkastaminen sen jälkeen, kun se tulee käyttäjien nähtäville tarkoittaa sitä, että sisältö tulee toisille käyttäjille näkyviin heti kun se lähetetään ja se tarkastetaan seuraavan 24 tunnin kuluessa. Kaikki käyttäjien luoma sisältö lisätään jonoon, josta palvelun moderaattorit joko hyväksyvät sen tai poistavat sen palvelusta. Hyvä puoli tässä tavassa on se, että käyttäjien ei tarvitse jäädä odottamaan sisällön hyväksymistä. Huono puoli tässä on se, että kaikki lähetetty roskaposti tulee aina osan käyttäjistä nähtäville. (Veglis, 2014)

Hajautettu käyttäjien luoman sisällön tarkastaminen tarkoittaa sitä, että myös käyttäjät osallistuvat sisällön tarkastamiseen. Tämä toimii parhaiten sivuilla, joilla on suuri käyttäjämäärä. (Veglis, 2014)

Automatisoitu käyttäjien luoman sisällön tarkastaminen eroaa muista tavoista siten että se ei vaadi ihmisiä valvomaan sisältöä. Automaattiseen sisällön tarkastamiseen voidaan käyttää erilaisia tekniikkoja, joista yksi on koneoppiminen. Tässä tavassa yleensä on alkukustannuksia, mutta ei käyttökustannuksia. (Veglis, 2014)

4.2 Automatisoitu roskapostin suodattaminen

Sosiaalisessa mediassa on mahdollista suodattaa roskapostia tekemällä tiettyjä sääntöjä, joiden täytyessä viesti katsotaan roskapostiksi. Yleensä viesti saa tietyn määrän pisteitä jokaisesta säännöstä ja kun jokin ennalta määritelty määrä pisteitä täyttyy, niin viesti katsotaan roskapostiksi. (Jain ym., 2019)

Ennalta määritetyillä säännöillä ei kuitenkaan saada kaikkea roskapostia suodatettua. Tämä johtuu siitä, että roskapostiviesti voidaan tehdä niin että se saa kierrettyä säännöt. Parempi lopputulos saadaan, kun käytetään näitä menetelmiä yhdessä koneoppimisen kanssa. (Jain ym., 2019)

Naive Bayes on yksi suosituimmista koneoppimiseen perustuvista luokittelualgoritmeista sosiaalisen median roskapostin tunnistamiseen. Muita ovat SVM, Random Forests ja neuroverkot. (Jain ym., 2019)

4.3 Automatisoidun roskapostin suodattamisen haasteita

Luonnollisen kielen käsittely automatisoidusti tuo monia haasteita. Ihmisten käyttämä kieli on joskus epäselvää (esimerkiksi englannin kielen lauseet ”I ate pizza with friends” ja ”I ate pizza with olives”). Luonnollinen kieli on myös vaihtelevaa (lauseen ”I ate pizza with friends” ydinsanomana voi välittää myös sanomalla ”I met my friends and shared some pizza”). Luonnollinen kieli myös muuttuu ja kehittyy jatkuvasti. (Goldberg, 2017)

Ihmiset ovat hyviä ymmärtämään ja tuottamaan luonnollista kieltä, ja pystyvät kehittämään ja tulkitsemaan todella monimutkaisia lauserakenteita. Vaikka ihmisinä pystymmekin käyttämään luonnollista kieltä, niin olemme tosi huonoja ymmärtämään ja kuvailemaan sen sääntöjä. Kuitenkin että luonnollista kieltä pystytään käsittelemään

automatisoidusti, tietokoneelle pitää kertoa sen säännöt ja tämän takia luonnollisen kielen käsittely tietokoneita hyödyntäen on todella haasteellista. (Goldberg, 2017)

Esimerkiksi jos on tehtävänä luokitella dokumentti yhteen neljästä eri luokasta, jotka ovat urheilu, politiikka, juoru ja talous, niin ihminen pystyy yleensä päättämään mihin luokkaan dokumentti saattaa kuulua. Lukija pystyy siis helposti määrittelemään tämän dokumentin luokan, mutta vaikeampaa on yrittää kirjoittaa ylös sääntöjä, joiden perusteella dokumentti luokitellaan. (Goldberg, 2017)

Tässä kohtaa avuksi voi ottaa ohjatut koneoppimisalgoritmit. Algoritmi tarvitsee avukseen muutama sataa ihmisen lajittelemaa dokumenttia, ja näiden perusteella se pystyy kehittämään mallin, jolla se pystyy lajittelemaan dokumentin luokan. Koneoppimisalgoritmit toimivat erinomaisesti ongelmassa, joissa sääntöjä ongelman ratkaisemiseen on hankala määrittää, mutta ongelman ratkaisu ihmiselle on kuitenkin suhteellisen yksinkertaista. (Goldberg, 2017)

5. Koneoppimistekniikat

Oppiminen tarkoittaa uuden tiedon, käyttäytymisen, arvojen, taitojen tai mieltymysten hankkimista tai muuttamista (Alzubi ym., 2018). He mainitsevat myös, että ihmiset oppivat uutta kokemusten perusteella, mutta tietokoneet oppivat sen sijasta datan perusteella. Ihmisten ja tietokoneiden oppiminen eroaa siis huomattavasti toisistaan. Koneoppiminen on tekoälyn osa-alue, joka mahdollistaa tietokoneiden ajattelemisen ja oppimisen ilman ihmisen apua (Alzubi ym., 2018). Tekijät kirjoittavat myös, että koneoppiminen tarkoittaa, että tietokoneet pystyvät muuttamaan toimintaansa siten että niiden tarkkuus paranee. Tarkkuus tarkoittaa tässä sitä, että kuinka monta kertaa niiden toiminta johtaa oikeaan lopputulokseen.

Koneoppimista voidaan hyödyntää monessa tietojenkäsittelyn osa-alueessa (Alzubi ym., 2018). Sen avulla voidaan esimerkiksi suodattaa roskapostia sähköpostista, tunnistaa huijauksia sosiaalisesta mediasta, käydä osakekauppaa verkossa, tunnistaa ihmisen kasvoja tai muita muotoja, tehdä lääketieteellisiä diagnooseja, ennustaa liikennettä, tunnistaa merkkejä painetulta tai kirjoitetulta paperiarkilta ja antaa tuotesuositteluja yksittäisille käyttäjille mieltymyksiensä mukaan (Alzubi ym., 2018). Googlen itseohjautuvat autot, Netflixin esittely elokuvista ja televisiosarjoista, joista käyttäjä saattaa pitää ja muut verkossa toimivat suosittelualgoritmit kuten Facebookin kaverisuositukset ovat kaikki esimerkkejä tosielämän koneoppimisen sovelluksista (Alzubi ym., 2018). Seuraavaksi esittelen kaksi koneoppimisen luokittelualgoritmia, jotka ovat Naive Bayes ja neuroverkot. Valitsin nämä algoritmit, koska niitä käytetään yleisesti luonnollisen kielen käsittelyyn.

5.1 Naive Bayes luokittelu

Naive Bayes luokittelualgoritmit ovat algoritmeja, jotka perustuvat Bayesin teoreemaan. Kaikista luokittelualgoritmeista varsinkin nämä ovat suosittuja sähköpostin kaupallisissa sekä avoimen lähdekoodin roskapostisuodattimissa (Rusland ym., 2017).

Bayesin teoreeman avulla voidaan laskea tapahtuman A todennäköisyys ehdolla B. Bayesin teoreema kuvaa kahden ehdollisen todennäköisyyden välillä olevaa suhdetta. (Efron, 2013)

Bayesin teoreeman kaava:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} \quad (1)$$

missä

- $P(A|B)$ on A:n todennäköisyys, kun B on totta
- $P(B|A)$ on B:n todennäköisyys, kun A on totta
- $P(A)$ on A:n todennäköisyys ilman tietoa B:stä
- $P(B)$ on B:n todennäköisyys ilman tietoa A:sta

Naive Bayes algoritmi on yksinkertainen luokittelija, joka laskee datan todennäköisyyden kuulua johonkin luokkaan (Rusland ym., 2017). Se tekee tämän laskemalla erilaisten arvojen lukumäärän ja kombinaation datassa (Rusland ym., 2017). Naive Bayes algoritmi

on suosittu tapa lajitella tekstiä, koska on huomattu, että se toimii yllättävän hyvin monella osa-alueella (Schneider, 2003). Tämän lisäksi se on yksinkertainen verrattuna muihin luokittelijoihin (Schneider, 2003).

5.2 Neuroverkot

Artificial Neural Network eli neuroverkot koostuvat suuresta joukosta neuroneita. Guptan (2013) mukaan ne kehitettiin alun perin jäljittelemään ihmisaivojen toimintoja. Ihmisaivoissa on yli 10 miljardia toisiinsa yhdistettyä neuronua, ja ne ottavat vastaan, prosessoivat sekä lähettävät informaatiota eteenpäin. Keinotekoiset neuroverkot ovat matemaattisia malleja, jotka jäljittävät näitä ihmisaivojen toimintoja. (Zhang, 2018)

Neuroverkot koostuvat toisiinsa yhdistetyistä keinotekoisista neuroneista. Jokainen neuroni koostuu synapseista, joille on määritelty synaptinen paino, summaajasta, joka summaa syötesignaalit sekä aktivaatiofunktiosta, joka rajoittaa neuronin ulostulon määrän johonkin rajalliseen arvoon. (Dongare ym., 2012)

Neuroverkko koostuu yleensä monista kerroksista. Ensimmäisenä on syötekerros, jonka neuronit ottavat vastaan neuroverkolle annetun syöteen. Tämän jälkeen voi olla yksi tai useampi piilokerros, jossa neuronit saavat syötteensä toisilta neuroneilta ja antavat tulosteensa toisille neuroneille. Viimeisenä on tulostekerros, joka antaa neuroverkon tulosteen. (Kukreja ym., 2016)

Neuroverkkoa pitää kouluttaa, että se voi tehdä tehtävänsä mahdollisimman hyvin. Kouluttaminen voidaan jakaa kahteen eri kategoriaan, jotka ovat valvottu ja valvomaton koulutus. Valvotussa kouluttamisessa algoritmille annetaan syöteen lisäksi myös haluttu ulostulo. Valvomattomassa koulutuksessa taas algoritmille ei anneta syöteen haluttua ulostuloa, vaan sen pitää tunnistaa yhtäläisyyksiä datan elementtien välillä. (Dreyfus, 2005)

6. Koneoppimisen käyttäminen roskapostin tunnistamiseen sosiaalisessa mediassa

Koneoppimistekniikoita on mahdollista soveltaa sosiaalisen median roskapostin suodattamiseen. Edellisessä luvussa esittelin Naive Bayesin ja neuroverkot. Tässä luvussa tutustumme tarkemmin, miten näiden avulla pystytään suodattamaan sosiaalisen median roskapostia.

Naive Bayes algoritmin toimivuus roskapostien tunnistamisessa perustuu siihen, että joitakin tiettyjä sanoja esiintyy todennäköisemmin roskapostissa tai halutussa sähköpostissa (Rusland ym., 2017). Esimerkiksi jos tietäisimme että sana "ilmainen" ei voi ikinä ilmestyä halutussa sähköpostissa, vaan ainoastaan roskapostissa ja näemme että kyseinen sana ilmestyy sähköpostissa, niin voisimme luokitella sen roskapostiksi.

Naive Bayes algoritmia käyttävät roskapostia tunnistavat algoritmit oppivat yleensä sen, että viesti on suurella todennäköisyydellä roskapostia, jos siinä esiintyy sanoja kuten "ilmainen". Jos siinä taas on esimerkiksi ystävien tai perheenjäsenten nimiä, niin se on todennäköisesti haluttu viesti. (Rusland ym., 2017)

Neuroverkot ovat tehokas työkalu, jonka käyttö voi olla hyödyllistä, kun käsitellään luonnollista kieltä. Vaikka neuroverkot ovat tehokas työkalu, niiden käyttöönotto voi olla hankalaa. (Goldberg, 2017)

Bag of words datamallia käytetään yleisesti, kun käsitellään luonnollista kieltä (Ma ym., 2018). Bag of words datamallia käytetään yleisesti luonnollisen kielen käsittelyssä, sekä Naive Bayesin että neuroverkkojen kanssa (Rusland ym., 2017; Goldberg, 2017). Bag of words datamalli sisältää jokaisen tekstissä ilmestyvän sanan, sekä tiedon siitä kuinka monta kertaa se esiintyy (Boulis & Ostendorf, 2005; Rusland ym., 2017). Datamallissa jätetään huomioimatta kieliooppi ja sanojen järjestys (Sharmin & Zaman, 2017). Monet koneoppimistekniikat tarkastelevat sanojen esiintymiskertoja tekstissä, minkä takia tätä mallia käytetään yleisesti (Boulis & Ostendorf, 2005).

Seuraavaksi esittelen miten Naive Bayesia pystyy soveltamaan käytännössä roskapostien suodattamiseen sosiaalisessa mediassa. Tämän jälkeen vertailemme Naive Bayesin ja neuroverkkojen suorituskykyä, kun suodatetaan roskapostia sosiaalisesta mediasta.

6.1 Naive Bayes luokittelun soveltaminen käytännössä

Sharmin & Zaman (2017) testasivat tutkimuksessaan Naive Bayes, K Nearest Neighbour (KNN), Bagging ja Support Vector Machine (SVM) algoritmien suorituskykyä tunnistamaan roskapostia Youtube kommentista. Tutkimuksessaan he käyttivät julkista Youtube-Spam-Collection Youtube kommenttien joukkoa. Siinä on viisi datajoukkoa, jotka koostuvat 1956 oikeasta viestistä. Viestit on kerätty viidestä videosta, jotka olivat kymmenen katsotuimman videon joukossa, kun data kerättiin. Eri datajoukkojen nimet ovat Psy, KatyPerry, LMFAO, Eminem ja Shakira. Kaikissa datajoukon näytteissä on kommentin kirjoittajan nimi, julkaisuajankohta ja kommentin metadata.

Ensiksi Sharmin & Zaman (2017) esikäsittelivät datan. Datan esikäsittelyssä he poistivat datajoukosta kaiken metadatan, ja jättivät vain kommenttien tekstit. Seuraavaksi he muodostivat tekstistä bag of words datamallin, jossa tekstiä edustaa laukku sanoja. Tämän jälkeen he tekivät Term Frequency-Inverse Document Frequency (TF-IDF) tilaston. TF-IDF mittaa lukumäärän sijasta sitä, kuinka relevantti jokainen sana on.

Esikäsittelyn jälkeen Sharmin & Zaman (2017) suorittivat piirteenalinnan. Piirteenalinnan tarkoituksena on parantaa luokittelun tehokkuutta, laskennallista tehokkuutta tai molempia (Sharmin & Zaman, 2017). Aggressiivinen piirteiden vähentäminen on monissa tapauksissa johtanut pieneen tarkkuuden vähenemiseen, ja suorituskyky parannuksiin (Sharmin & Zaman, 2017). Sharmin & Zaman (2017) käyttivät hukkasanojen poistoon perustuvia piirteenalinta menetelmiä. Hukkasanojen poistossa dokumentista poistetaan yleiset sanat, joiden perusteella dokumenttia ei voida luokitella (Sharmin & Zaman, 2017).

Sharmin & Zaman (2017) suorittivat luokittelun WEKA koneoppimistyökalun avulla. WEKA on avoimen lähdekoodin ohjelma, jossa on koneoppimisalgoritmeja tiedonlouhintaa varten (Sharmin & Zaman, 2017). Sharmin & Zaman (2017) rakensivat mallin algoritmeille työkalun avulla.

Viimeisenä Sharmin & Zaman (2017) suorittivat arvioinnin. Sharmin & Zaman (2017) käyttivät Naive Bayes algoritmin suorituskyvyn mittaamiseen monia mittareita. Näitä olivat ulkoinen tarkkuus (eng. accuracy), luokitteluvirhe, sisäinen tarkkuus (eng. precision), saanti (eng. recall), Matthews Correlation Coefficient (MCC) ja F-measure. Ulkoinen tarkkuus tarkoittaa mallin kykyä luokitella sille uusi data oikein. Sisäinen tarkkuus on oikeiden positiivisten tulosten lukumäärä jaettuna oikeiden positiivisten ja väärin positiivisten summalla. Saanti on oikeiden positiivisten tulosten lukumäärä jaettuna oikeiden positiivisten ja väärin negatiivisten summalla. F-measure lasketaan testin sisäisestä tarkkuudesta ja saannista. MCC ottaa huomioon oikeat ja väärät positiiviset ja negatiiviset tulokset, ja sen katsotaan olevan tasapainoinen menetelmä. Sitä voidaan käyttää, vaikka luokat olisivat todella eri kokoisia. MCC palauttaa luvun -1 ja +1 väliltä, jossa +1 tarkoittaa täydellistä tulosta, 0 vastaa satunnaista arvausta ja -1 tarkoittaa täydellistä ennustuksen ja havainnon päinvastaisuutta. (Sharmin & Zaman, 2017) Tutkimuksen algoritmien ulkoiset tarkkuudet löytyvät taulukosta 1 ja MCC arvot taulukosta 2.

Taulukko 1. Sharmin & Zaman (2017) tutkimuksen algoritmien ulkoinen tarkkuus

Algoritmi	KatyPerry	Shakira	Psy	Eminem	LMFAO
1-KNN	84,00 %	90,52 %	94,29 %	93, 82 %	90, 18 %
3-KNN	78,28 %	85, 14 %	92, 29 %	92, 27 %	89,04 %
Naive Bayes	92,00 %	92,16 %	94, 57 %	92, 49 %	89,47 %
Bagging	93,00 %	91,08 %	91,43 %	94,92 %	90,64 %
SVM	57,43 %	70, 27 %	93, 71 %	92, 05 %	91, 09 %

Taulukko 2. Sharming & Zaman (2017) tutkimuksen algoritmien MCC arvot

Algoritmi	KatyPerry	Shakira	Psy	Eminem	LMFAO
1-KNN	0,718	0,813	0,88	0,883	0,817
3-KNN	0,628	0,714	0,85	0,857	0,799
Naive Bayes	0,841	0,843	0,894	0,85	0,801
Bagging	0,866	0,83	0,83	0,902	0,826
SVM	0,265	0,48	0,88	0,851	0,834

Naive Bayes ja Bagging algoritmit antoivat kokonaisuudessa parhaan ulkoisen tarkkuuden ja MCC arvon Sharmin & Zaman (2017) tutkimuksessa. Tästä oli kuitenkin myös poikkeuksia. Esimerkiksi SVM algoritmi antoi tutkimuksessa parhaan ulkoisen tarkkuuden LMFAO datajoukolle. Ottaen huomioon myös sisäisen tarkkuuden ja saannin Naive Bayes ja Bagging algoritmit suoriutuivat Sharmin & Zaman (2017) tutkimuksessa parhaiten.

6.2 Naive Bayesin ja neuroverkkojen suorituskyky

Wang (2010) testasi tutkimuksessaan päätöspuu, neuroverkot, SVM ja Naive Bayes luokittelualgoritmien suorituskykyä tunnistamaan roskapostia Twitter viesteistä. Tutkimuksessa käytetty Twitter datajoukko kerättiin Twitterin API:n ja hakurobotin avulla. Datajoukko käsitti 25 847 käyttäjää ja noin 500 000 twiittiä.

Wang (2010) tutkimuksessa 500 käyttäjätiliä luokiteltiin kahteen eri luokkaan algoritmien testaamiseksi. Nämä luokat olivat roskaposti ja ei roskaposti. Jokainen käyttäjätili tarkistettiin manuaalisesti lukemalla 20 viimeistä twiittiä ja tarkistamalla käyttäjätilin kaverit ja seuraajat. Tuloksien mukaan datajoukon käyttäjätileistä 1 % oli roskapostia. Wang (2010) mukaan todellisuudessa 3 % Twitterin käyttäjätileistä on roskapostia, joten datajoukkoon lisättiin roskapostiksi luokiteltavia käyttäjätilejä. Tämän jälkeen myös datajoukossa 3 % kaikista käyttäjätileistä oli roskapostia. Suorituskyvyn mittaamiseksi Wang (2010) tutkimuksessa mittareina toimi sisäinen tarkkuus, saanti ja F-measure. Taulukossa 3 on esillä tutkimuksen tulokset.

Taulukko 3. Wang (2010) tutkimuksen tulokset

Luokittelija	Sisäinen tarkkuus	Saanti	F-measure
Decision Tree	0,667	0,333	0,444
Neuroverkot	1	0,417	0,588
SVM	1	0,25	0,4
Naive Bayes	0,917	0,917	0,917

Kokonaisuudessa Naive Bayes suoriutui Wang (2010) tutkimuksesta parhaiten, joka nähdään sen korkeimmista F-measure pisteistä. Neuroverkot suoriutuivat tutkimuksessa kokonaisuudessa toiseksi parhaiten, joka nähdään myös sen F-measure pisteistä. Tutkimuksessa huomattavaa on, että neuroverkkojen saanti oli vain 0,417. Naive Bayesin saanti oli taas paljon parempi 0,917. Naive Bayes siis tunnisti huomattavasti suuremman osan roskapostista kuin neuroverkot, tai mikään muista tutkimuksessa mukana olleista algoritmeista. Tutkimuksessa Naive Bayesin sisäinen tarkkuus oli hieman huonompi kuin neuroverkkojen tai SVM:n.

7. Pohdinta

Automatisoidulla roskapostin suodattamisella on monia etuja manuaaliseen suodattamiseen verrattuna. Jos käytetään manuaalista suodattamista, niin pitää valita, että annetaanko käyttäjän odottaa sisällön julkaisemista, kunnes palvelun moderaattori on tarkistanut sen vai tarkistetaanko sisältö vasta jälkeempään, joka johtaa siihen, että suurempi määrä roskapostista menee palvelun käyttäjien nähtäville (Veglis, 2014). Automatisoidussa suodattamisessa tätä valintaa ei tarvitse tehdä, ja automatisoitu suodattaminen voi monessa tapauksessa olla halvempaa, kuin se että palkataan henkilökuntaa tarkastamaan sisältöä (Veglis, 2014).

Varsinkin sosiaalisen median palveluissa on tyypillistä, että käyttäjät haluavat viestinsä heti muiden käyttäjien näkyville, eivätkä halua ensin odottaa moderaattorin tarkistusta (Veglis, 2014). Toisaalta taas, jos suuri määrä roskapostia tulee palvelun käyttäjien nähtäville, niin siitä voi aiheutua yritykselle taloudellista vahinkoa ja mainehaittaa (Blanzieri & Bryl, 2008).

Koneoppimista kannattaa käyttää roskapostin tunnistamiseen, kun roskapostia halutaan suodattaa automatisoidusti. Koneoppimista hyödyntämällä on mahdollista rakentaa roskapostia suodattava järjestelmä, johon ei tarvitse tehdä sääntöjä, joita pitää itse päivittää jatkuvasti, kun roskapostin lähettäjät oppivat kiertämään ne (Goldberg, 2017). Koneoppimisen avulla pystytään ratkaisemaan monia luonnollisen kielen käsittelyn ongelmia, ja automatisoimaan siinä prosesseja, mikä ei perinteisellä ohjelmoinnilla olisi mahdollista (Goldberg, 2017).

Tutkielmassa käydään läpi menetelmiä, joita voidaan käyttää roskapostin suodattamiseen, ja niistä esimerkkeinä olivat Naive Bayes algoritmit ja neuroverkot, joita molempia voidaan käyttää luokittelijana. Naive Bayes algoritmit ovat yksinkertaisempia ja helpompia toteuttaa kuin neuroverkot (Rusland ym., 2017). Neuroverkkojen suunnittelu voi olla monimutkaisempaa, ja siihen voi liittyä monia päätöksiä esimerkiksi liittyen piilokerrokseen (Goldberg, 2017). Ne ovat kumpikin erilaisia tapoja ratkaista sama ongelma.

Neuroverkot suoriutuivat Wang (2010) tutkimuksessa F-measure arvon mukaan Naive Bayesia huonommin. Huomattavaa kuitenkin on, että tutkimuksessa neuroverkoilla oli täydellinen sisäinen tarkkuus, kun taas Naive Bayes algoritmilla sisäinen tarkkuus oli 0,917. Neuroverkot tunnistivat tutkimuksessa oikein 100 % roskaposteista, mutta ne löysivät vain 41,7 % kaikista roskaposteista. Naive Bayes algoritmi taas löysi 91,7 % kaikista roskaposteista, mutta se lajitteli osan käyttäjätileistä väärin roskaposteiksi. Tämän tutkimuksen perusteella Naive Bayesia voidaan pitää neuroverkkoja parempana roskapostin suodattamiseen. Naive Bayes löysi yli 90 % kaikista roskaposteista, ja luokitteli vain vähän viestejä väärin roskapostiksi. Neuroverkot taas eivät luokitelleet yhtään viestiä väärin roskapostiksi, mutta ne pääsivät suurimman osan oikeasta roskapostista lävitse.

Koneoppimisen käyttämistä sosiaalisessa mediassa käsittelevässä kappaleessa huomasimme, että jos algoritmi tietää jonkin sanan, joka esiintyy aina roskapostissa, niin viesti voidaan luokitella 100 % varmasti roskapostiksi. Todellisuudessa ei kuitenkaan ole

olemassa sanoja, joita esiintyy ainoastaan roskapostissa tai halutussa viestissä, ja tämän takia roskapostin tunnistaminen perustuu todennäköisyyteen. Jos roskaposteissa yleisiä sanoja ilmestyy viestissä paljon, niin viesti voidaan mahdollisesti luokitella roskapostiksi (Rusland ym., 2017). Jos viestissä taas ilmestyy paljon sanoja, jotka ovat yleisiä halutuissa viesteissä, niin voidaan katsoa, että se ei ole roskaposti (Rusland ym., 2017). Koska ei ikinä voida sanoa 100 % varmuudella viestin olevan joko haluttu viesti tai roskaposti, niin haluttuja viestejä voidaan luokitella myös väärin roskapostiksi ja toisin päin. Sosiaalisen median alustalle viestin luokittelu väärin voi tuottaa taloudellista haittaa, ja tämä takia on tärkeätä, että nämä tapaukset minimoidaan.

Sosiaalisen median alustojen pitää valita kuinka varma algoritmin pitää olla, että viesti on roskaposti, ennen kuin mitään toimenpiteitä tehdään, jos ne käyttävät automatisoitua roskapostin suodattamista. Tässä niiden pitää käytännössä valita kahden pahan väliltä. Joko roskapostia pääsee palveluun enemmän läpi ja haluttuja viestejä ei poisteta turhaan palvelusta niin paljon, tai roskapostia ei pääse palveluun läpi niin paljon mutta haluttuja viestejä poistetaan enemmän. Kummatkin näistä vaihtoehdoista voivat aiheuttaa merkittävää mainehaittaa yritykselle. Nykyään varsinkin tieto siitä, jos sosiaalisen median alusta poistaa viestejä, ja mahdollisesti estää käyttäjiä lähettämästä uusia viestejä turhaan leviää nopeasti ja voi aiheuttaa merkittävää mainehaittaa yritykselle (Horn ym., 2015).

Tämän työn tuloksista voi olla hyötyä henkilöille, jotka haluavat tietää enemmän roskapostin suodattuksesta sosiaalisen median alustoilla ja heille, jotka suunnittelevat koneoppimisen käyttöä roskapostin tunnistamiseen. Tutkimusta voi käyttää hyväkseen, kun haluaa saada yleiskäsityksen aiheesta, mutta pelkästään sen perusteella ei ole mahdollista esimerkiksi rakentaa omaa koneoppimisalgoritmia.

8. Yhteenveto

Tämän työn tarkoituksena oli tutkia, miten koneoppimista voi käyttää hyväksi roskapostin suodattamisessa, sekä mitä eroja sillä on muihin tapoihin verrattuna. Tutkimuskysymyksenä työssä oli ”*Miten roskapostia voi tunnistaa koneoppimisen avulla sosiaalisessa mediassa, ja mitä etuja ja haittoja sillä on verrattuna muihin menetelmiin?*”. Aiheeseen liittyen on tehty paljon tutkimusta ja tämän työn tarkoituksena oli kerätä kyseisistä töistä oleellisimpia tutkimustuloksia.

Roskaposti on yleistä sosiaalisessa mediassa ja se aiheuttaa harmia internetin käyttäjille. Sosiaalisessa mediassa yksi jokaisesta 200 viestistä on roskapostia. Sosiaalisen median yrityksille roskaposti voi aiheuttaa taloudellista vahinkoa sekä mainehaittaa.

Automatisoidulla roskapostin suodattamisella on monia etuja verrattuna siihen, että roskaposti suodatettaisiin manuaalisesti. Etuja ovat ne, että automatisoitu suodattaminen voi tulla yritykselle halvemmaksi kuin manuaalinen suodattaminen ja roskaposti voidaan suodattaa reaaliajassa ilman että viestin lähettävän käyttäjän tarvitsee odottaa moderaattorin tarkistusta. Yksi haitta automatisoidulla suodattamisella on se, että siinä tapahtuu yleensä enemmän virheitä kuin manuaalisessa.

Koneoppimisella on monia etuja, kun luonnollista kieltä käsitellään automatisoidusti. Näitä ovat esimerkiksi se, että ohjelmaan ei tarvitse kirjoittaa monimutkaisia sääntöjä roskapostin tunnistamiseksi.

Kun halutaan hyödyntää koneoppimista, niin luokittelijana voidaan käyttää esimerkiksi Naive Bayes algoritmeja tai neuroverkkoja. Naive Bayes algoritmit ovat Bayesin teoreemaan perustuvia luokittelualgoritmeja. Neuroverkot ovat suuresta määrästä neuroneita koostuvia verkkoja, jotka jäljittävät ihmisaivojen toimintoja. Bag of words datamallia käytetään yleisesti näiden luokittelijoiden kanssa.

Naive Bayes on yksinkertaisempi toteuttaa kuin neuroverkot ja tutkielmassa käsiteltyjen tutkimuksien mukaan se suodattaa roskapostia tehokkaammin kuin neuroverkot. Wang (2010) tutkimuksessa neuroverkot lajittelivat kaikki roskapostit oikein, mutta ne löysivät vain murto-osan kaikista roskaposteista päästäen loput läpi. Samassa tutkimuksessa Naive Bayes taas löysi yli 90 % roskapostista ja lajitteli vain vähän viestejä väärin roskapostiksi.

Lisätutkimuksissa asiaa voisi tutkia myös empiirisesti. Yksi tapa olisi rakentaa oma koneoppimista hyödyntävä roskapostisuodatin. Tätä varten voisi implementoida esimerkiksi Naive Bayes algoritmin itse tai käyttää siihen valmista kirjastoa. Ensin ohjelmisto on pitää toteuttaa ja kouluttaa se antamalla sille tarvittava määrä roskapostia ja toivottuja viestejä. Tämän jälkeen voidaan mitata, kuinka suuren osan viesteistä se luokittelee oikein.

Kokonaisuutena voidaan siis todeta, että roskapostia voidaan tunnistaa koneoppimisen avulla sosiaalisessa mediassa hyödyntämällä esimerkiksi Naive Bayes algoritmeja tai neuroverkkoja. Etuja koneoppimisella verrattuna muihin menetelmiin on se, että kun koneoppimista hyödynnetään, ei ole tarvetta tehdä käsin kirjoitettuja sääntöjä roskapostin suodattamiseen. Huono puoli näissä säännöissä on, että niitä pitää päivittää jatkuvasti, kun roskapostin lähettäjät oppivat kiertämään ne. Yksi huono puoli koneoppimisella verrattuna muihin menetelmiin on, että koneoppimisalgoritmit lajittelevat osan viesteistä väärin. Samaa ongelmaa ei välttämättä ole roskapostin manuaalisella suodattamisella.

Tutkimuksessa käsiteltävien aiempien tutkimusten mukaan algoritmin valinnalla on tähän merkitystä. Käsitellyissä tutkimuksissa Naive Bayes suoriutui selvästi neuroverkkoja paremmin. Hyvin valitulla algoritmilla voidaan siis vähentää väärin luokiteltujen viestien määrää.

Lähteet

- Almadhoor, L. (2021). Social media and cybercrimes. *Turkish Journal of Computer and Mathematics Education (TURCOMAT)*, 12(10), 2972-2981.
- Alzubi, J., Nayar, A., & Kumar, A. (2018, November). Machine learning from theory to algorithms: an overview. In *Journal of physics: conference series* (Vol. 1142, No. 1, p. 012012). IOP Publishing.
- Bi, Q., Goodman, K. E., Kaminsky, J., & Lessler, J. (2019). What is machine learning? A primer for the epidemiologist. *American journal of epidemiology*, 188(12), 2222-2239.
- Blanzieri, E., & Bryl, A. (2008). A survey of learning-based techniques of email spam filtering. *Artificial Intelligence Review*, 29(1), 63-92.
- Boulis, C., & Ostendorf, M. (2005, April). Text classification by augmenting the bag-of-words representation with redundancy-compensated bigrams. In *Proc. of the International Workshop in Feature Selection in Data Mining* (pp. 9-16). Citeseer.
- Cao, C., & Caverlee, J. (2015, March). Detecting spam urls in social media via behavioral analysis. In *European conference on information retrieval* (pp. 703-714). Springer, Cham.
- Castillo, C., Donato, D., Becchetti, L., Boldi, P., Leonardi, S., Santini, M., & Vigna, S. (2006, December). A reference collection for web spam. In *ACM Sigir Forum* (Vol. 40, No. 2, pp. 11-24). New York, NY, USA: ACM.
- Chakraborty, M., Pal, S., Pramanik, R., & Chowdary, C. R. (2016). Recent developments in social spam detection and combating techniques: A survey. *Information Processing & Management*, 52(6), 1053-1073.
- Dreyfus, G. (2005). *Neural networks: methodology and applications*. Springer Science & Business Media.
- Dongare, A. D., Kharde, R. R., & Kachare, A. D. (2012). Introduction to artificial neural network. *International Journal of Engineering and Innovative Technology (IJEIT)*, 2(1), 189-194.
- Efron, B. (2013). Bayes' theorem in the 21st century. *Science*, 340(6137), 1177-1178.
- Ferrara, E. (2019). The history of digital spam. *Communications of the ACM*, 62(8), 82-91.
- Goldberg, Y. (2017). Neural network methods for natural language processing. *Synthesis lectures on human language technologies*, 10(1), 1-309.
- Gupta, N. (2013). Artificial neural network. *Network and Complex Systems*, 3(1), 24-28.

- Hamilton, I. A. (2022, helmikuu 3). Facebook's user numbers shrunk for the first time in its history. Business Insider. Viitattu 7.5.2022, saatavilla: <https://www.businessinsider.com/meta-facebook-user-numbers-shrink-first-time-ever-2022-2>
- Horn, I. S., Taros, T., Dirkes, S., Hüer, L., Rose, M., Tietmeyer, R., & Constantinides, E. (2015). Business reputation and social media: A primer on threats and responses. *Journal of Direct, Data and Digital Marketing Practice*, 16(3), 193-208.
- Inuwa-Dutse, I., Liptrott, M., & Korkontzelos, I. (2018). Detection of spam-posting accounts on Twitter. *Neurocomputing*, 315, 496-511.
- Jain, G., Sharma, M., & Agarwal, B. (2019). Spam detection in social media using convolutional and long short term memory neural network. *Annals of Mathematics and Artificial Intelligence*, 85(1), 21-44.
- Kaplan, J. (2016). *Artificial Intelligence: What Everyone Needs to Know*. Oxford University Press.
- Kukreja, H., Bharath, N., Siddesh, C. S., & Kuldeep, S. (2016). An introduction to artificial neural network. *Int J Adv Res Innov Ideas Educ*, 1, 27-30.
- Ma, S., Sun, X., Wang, Y., & Lin, J. (2018, July). Bag-of-Words as Target for Neural Machine Translation. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)* (pp. 332-338).
- Metaxas, P. T., & Mustafaraj, E. (2012). Social media and the elections. *Science*, 338(6106), 472-473.
- Rusland, N. F., Wahid, N., Kasim, S., & Hafit, H. (2017, August). Analysis of Naïve Bayes algorithm for email spam filtering across multiple datasets. In *IOP conference series: materials science and engineering* (Vol. 226, No. 1, p. 012091). IOP Publishing.
- Schneider, K. M. (2003, April). A comparison of event models for naive bayes anti-spam e-mail filtering. In *10th Conference of the European Chapter of the Association for Computational Linguistics*.
- Sharmin, S., & Zaman, Z. (2017, December). Spam detection in social media employing machine learning tool for text mining. In *2017 13th International Conference on Signal-Image Technology & Internet-Based Systems (SITIS)* (pp. 137-142). IEEE.
- Shehnepoor, S., Salehi, M., Farahbakhsh, R., & Crespi, N. (2017). NetSpam: A network-based spam detection framework for reviews in online social media. *IEEE Transactions on Information Forensics and Security*, 12(7), 1585-1595.
- Twitter global mdau 2021 (2022, toukokuu 5). Statista. Viitattu 7.5.2022, saatavilla: <https://www.statista.com/statistics/970920/monetizable-daily-active-twitter-users-worldwide/>

- Veglis, A. (2014, June). Moderation techniques for social media content. In *International conference on social computing and social media* (pp. 137-148). Springer, Cham.
- Wang, A. H. (2010, July). Don't follow me: Spam detection in twitter. In *2010 international conference on security and cryptography (SECRYPT)* (pp. 1-10). IEEE.
- Zeng, D., Chen, H., Lusch, R., & Li, S. H. (2010). Social media analytics and intelligence. *IEEE Intelligent Systems*, 25(6), 13-16.
- Zhang, Z. (2018). Artificial neural network. In *Multivariate time series analysis in climate and environmental research* (pp. 1-35). Springer, Cham.