



Kandidatutkielma

Tietojenkäsittelytieteen kandiohjelma

NLP-hyökkäysten käyttökohteet

Akira Taguchi

18.3.2022

MATEMAATTIS-LUONNONTIETEELLINEN TIEDEKUNTA
HELSINGIN YLIOPISTO

Ohjaaja(t)

FT Nikolaj Tatti

Yhteystiedot

PL 68 (Pietari Kalmin katu 5)
00014 Helsingin yliopisto

Sähköpostiosoite: info@cs.helsinki.fi

URL: <http://www.cs.helsinki.fi/>

Tiedekunta — Fakultet — Faculty		Koulutusohjelma — Utbildningsprogram — Study programme	
Matemaattis-luonnontieteellinen tiedekunta		Tietojenkäsittelytieteen kandiohjelma	
Tekijä — Författare — Author			
Akira Taguchi			
Työn nimi — Arbetets titel — Title			
NLP-hyökkäysten käyttökohteet			
Ohjaajat — Handledare — Supervisors			
FT Nikolaj Tatti			
Työn laji — Arbetets art — Level	Aika — Datum — Month and year	Sivumäärä — Sidoantal — Number of pages	
Kandidutkielma	18.3.2022	7 sivua	
Tiivistelmä — Referat — Abstract			
<p>Luonnollisen kielen käsittely on kätevä työkalu käsittelemään ihmisten puhumaa kieltä tietokoneaailmassa.</p> <p>Luonnollisen kielen käsittely on kuitenkin sellaisenaan haavoittuvainen erillaisille tekstipohjaisille hyökkäyksille.</p> <p>Tässä tutkielmassa tutustutaan näiden luonnollisen kielen käsittely historiaan, aikaisemmin mainittujen hyökkäysten mahdollistajiin sekä näiden torjuntametodeihin.</p>			
<p>ACM Computing Classification System (CCS)</p> <p>Security and privacy</p> <p>Computing methodologies → Artificial Intelligence → Natural language processing</p>			
Avainsanat — Nyckelord — Keywords			
nlp, unicode, nlp attack, machine learning, Natural Language Processing, cyber security			
Säilytyspaikka — Förvaringsställe — Where deposited			
Helsingin yliopiston kirjasto			
Muita tietoja — övriga uppgifter — Additional information			

Sisältö

1	Johdanto	1
2	Hyökkäystaksonomia	2
2.1	Näkymättömät merkit	2
2.2	Homoglyfit	2
2.3	Uudelleenjärjestelyt	3
2.4	Poistatukset	3
3	Puolustusmetodit	4
3.1	OCR-puolustus	4
3.2	Suorituskykykeskeinen puolustus	4
4	Skaalautuvuus	5
4.1	Tulevaisuus luonnollisen kielen käsittelylle	5
4.2	Tulevaisuus NLP -hyökkäyksille	5
5	Yhteenveto	6
	Lähteet	7

1 Johdanto

Koneoppimisen käyttötarkoitusten määrä kasvaa vuosi vuodelta suuremmaksi. Tätä teknologiaa voidaan hyödyntää muun muassa ihmisten puhuman kielen käsittely. Luonnollisen kielen käsittely (eng. Natural Language Processing, NLP) on alati kasvavassa kuluttajakäytössä johtuen laskentatehon kasvusta, suurien tietomäärien saatavuudesta, onnistuneiden koneoppimismetodien kehittämisestä sekä laajemmasta ihmiskielen ymmärryksestä ja sen käytöstä eri konteksteissa (Hirschberg ja Manning, 2015).

Luonnollisen kielen käsittely on kohdennetun mainonnan keskiössä. Viesti ystävälle mainoskohdennetussa viestipalvelussa antaa työstettävän datan NLP-mallille: “Mikä elokuva meidän pitäisi katsoa viikonloppuna?” NLP-mallin avulla automaattinen mainostaja ymmärtää mainostaa kyseiselle käyttäjälle miltei välittömästi sarjalippuja mainostavasta elokuvateatterista, suoratoistopalvelua tai mainostavaa aktiviteettikeskusta kyseiselle viikonlopulle.

Tässä tutkielmassa tarkastellaan NLP-hyökkäysten käyttökohteita. Tähän kuuluu hyökkäystaksonomia, puolustusmenetelmät sekä NLP-mallien sekä niihin kohdistuvien hyökkäysten tulevaisuus. Koska NLP-mallit ovat eksponentiaalisessa nousussa kuluttajakäytössä, on tärkeää spekuloida mahdollisia kehityksiä koneoppimisessa sekä tästä syntyviä haavoittuvuuksia. (Yang et al., 2021). (Boucher et al., 2021)

2 Hyökkäystaksonomia

Käymme seuraavaksi läpi neljä erilaista huomaamatonta hyökkäysmetodia. Nämä hyökkäykset eivät siis näy visuaalisesti ihmiskäyttäjälle näyttöpäätteellä erilaisina verrattuna viattomaan tekstiin. Tarkemmin keskitymme Unicoden ja muiden enkoodausmenetelmien hyväksikäyttämiseen NLP-malleja vastaan.

2.1 Näkymättömät merkit

Näkymättömät merkit vaikuttavat tietokoneen NLP-mallin ymmärtämään kontekstiin. Esimerkki tällaisesta on nollatilavuuden välilyönti -merkki, jonka Unicode merkintä on U+200B. Tällä merkillä voimme esimerkiksi vaikuttaa pelichattiin lähetettävän toksissuodatettavan merkkijonoon "olet huono"niin, että merkkijono menisi NLP-mallin läpi chätistä. Merkkijono `oletU+200Bhuono` saattaisi mennä läpi chatin suodattimesta, mutta vastapuolelle viesti olisi edelleen `olet huono`.

Kontekstin poistamisen lisäksi näkymättömillä merkeillä voidaan myös tuoda ja syrjäyttää konteksteja toisilla. Mikä pyhäinhäväistyksen rakennus! Miten onnistuit tekemään tämän näin laiskasti? -tekstin negatiivisuus voidaan syrjäyttää positiivisuudella syöttämällä NLP-mallille sen sijaan teksti `Mikä pyU+200BhäinhävU+200BäistyU+200BksenU+200Brakennus! Miten onnistuit tekemään tämän U+200BnäU+200Bin laU+200BiskasU+200Bti?.`

2.2 Homoglyfit

Homoglyfihyökkäykset NLP-malleja vastaan pohjautuvat pahantahtoisten merkkien virallisten esitysmuotojen näyttävän hyväntahtoisten merkkien virallisilta esityksiltä. Joissain kielissä tekstin merkitys muuttuu täysin yhden merkin vaihtuessa. Esimerkkinä homoglyyfistä on $A \rightarrow A$, missä viimeinen kirjain on todellisuudessa kyrillinen kirjain А. Näkymättömien merkkien lailla homoglyfihyökkäyksen toteutus riippuu ympäristön fontista.

2.3 Uudelleenjärjestelyt

Uudelleenjärjestelyhyökkäys pohjautuu näennäisen tekstin uudelleenjärjestämiseen pahan-tahtoisesti. Pankkitilinumeron 1234567 pystyy esimerkiksi vaihtamaan kaksisuuntaisella-algoritmillä tilinumeroksi 7654321 maksajan huomaamatta mitään.

2.4 Poistatukset

Poistatushyökkäyksiä on vaikeampi toteuttaa aikaisempiin metodeihin verrattuna. Tämä johtuu useimpien käyttöjärjestelmien estosta kopioida poistatusta sisältävää tekstiä leike-pöydälle.

3 Puolustusmetodit

3.1 OCR-puolustus

NLP-hyökkäykset voidaan estää alhaisemmalla tasolla overheadilla sekä korkeammalla tasolla edistyneen teknologian turvin. Näytöltäluvun (eng. OCR, On-Screen-Reading) avulla epäselvyydet tekstin aidosta luonteesta voidaan uudelleenrenderöidä tulkitsemalla aineisto uudestaan visuaalisesti. Tämä metodi lisää overheadia huomattavasti riippuen käyttötarcoituksesta, mutta poistaa pahantahtoiset merkit ilman NLP-mallin uudelleenkoulutusta.

3.2 Suorituskykykeskeinen puolustus

Keskitymme seuraavaksi näkymättömiin merkkeihin, -homoglyyfeihin, uudelleenjärjestelyihin -ja poistatuksiin perustuvien hyökkäyksien puolustamiseen.

Tietyt näkymättömät merkit voidaan poistaa suoraan syötteestä. Mikäli applikaatiossa näitä merkkejä ei voida poistaa, voidaan ne korvata non-`<unk>` upotuksilla.

Homoglyyfihyökkäysten torjuminen OCR-metodilla on ymmärrettävästi vaikeampaa verrattuina muihin merkkeihin. Paras keino torjua tällaisia hyökkäyksiä olisi mapata osa homoglyyfeistä niiden yleisemmin tunnettuihin vastineisiin. NLP-mallin ylläpitäjä joutuu tekemään tässä siis suurimman jalkatyön.

Uudelleenjärjestelyhyökkäykset voidaan torjua riisumalla bidi-ohjausmerkit syötteestä, varoittamalla käyttäjää bidi-ohjausmerkkien ilmestyessä syötteeseen tai käyttämällä bidi-algoritmia halutun syötteen selvittämiseen.

Poistatukset yleensä havaitaan NLP-mallien ulkopuolella syötteenannon ensivaiheessa. NLP-mallin tasolla tähän tarvitsee harvemmin puuttua ja käyttäjälle voidaan pahimmassa tapauksessa lähettää varoitus poistatusmerkkien olemassaolosta syötteessä.

4 Skaalautuvuus

4.1 Tulevaisuus luonnollisen kielen käsittelylle

4.2 Tulevaisuus NLP -hyökkäyksille

5 Yhteenveto

Luonnollisen kielen käsittely on kätevä työkalu käsittelemään ihmisten puhumaa kieltä tietokonemaailmassa.

Luonnollisen kielen käsittely on kuitenkin sellaisenaan haavoittuvainen erillaisille tekstipohjaisille hyökkäyksille.

Tässä tutkielmassa tutustuimme näiden luonnollisen kielen prosessoinnin historiaan, aikaisemmin mainittujen hyökkäysten mahdollistajiin sekä näiden torjuntametodeihin.

Lähteet

- Boucher, N., Shumailov, I., Anderson, R. ja Papernot, N. (2021). *Bad Characters: Imperceptible NLP Attacks*. arXiv: [2106.09898](https://arxiv.org/abs/2106.09898) [cs.CL].
- Hirschberg, J. ja Manning, C. D. (2015). "Advances in natural language processing". *Science* 349.6245, s. 261–266. DOI: [10.1126/science.aaa8685](https://doi.org/10.1126/science.aaa8685). eprint: <https://www.science.org/doi/pdf/10.1126/science.aaa8685>. URL: <https://www.science.org/doi/abs/10.1126/science.aaa8685>.
- Yang, W., Lin, Y., Li, P., Zhou, J. ja Sun, X. (elokuu 2021). "Rethinking Stealthiness of Backdoor Attack against NLP Models". Teoksessa: *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. Online: Association for Computational Linguistics, s. 5543–5557. DOI: [10.18653/v1/2021.acl-long.431](https://doi.org/10.18653/v1/2021.acl-long.431). URL: <https://aclanthology.org/2021.acl-long.431>.

