

Koneoppimisen hyödyntäminen avointen lähteiden tiedustelussa

Tietojenkäsittelytiede
Tietotekniikan laitos, Teknillinen tiedekunta
Kandidaatintutkielma

Laatija:
Eetu Torppa

Toukokuu 2025

Kandidaatintutkielma
Tietotekniikan laitos, Teknillinen tiedekunta
Turun yliopisto

Tutkinto-ohjelma: Tietojenkäsittelytiede

Tekijä: Eetu Torppa

Otsikko: Koneoppimisen hyödyntäminen avointen lähteiden tiedustelussa

Sivumäärä: 33 sivua

Päivämäärä: Toukokuu 2025

Internetin ja sosiaalisen median kasvu on lisännyt saatavilla olevan datan määrää. Avointen lähteiden tiedustelu on muodostunut merkittäväksi tiedustelun haaraksi muutosten myötä, mutta saatavilla olevan tiedon määrä on hankaloittanut perinteistä manuaalisesti tehtävää avointen lähteiden tiedustelua. Koneoppiminen on kehittynyt viimeisen kymmenen vuoden aikana merkittävästi ja tarjoaa menetelmiä suurten datamäärien läpikäymiseen. Tässä tutkielmassa käsitellään sitä, miten koneoppimista voidaan hyödyntää avointen lähteiden tiedustelussa. Tarkastelussa ovat avointen lähteiden tiedustelussa yleisimmin käytetyt koneoppimisen menetelmät sekä tulevaisuuden näkymät näiden kahden pääkäsitteen osalta.

Tutkielmassa havaittiin, että koneoppiminen on kasvavassa määrin tärkeä osa avointen lähteiden tiedustelua. Koneoppiminen on tällä hetkellä tärkeä osa tiedon keräämistä, käsittelyä ja analysointia. Yleisimmin käytettyjä menetelmiä ovat keinotekoisiiin neuroverkkoihin perustuvat konvoluutioverkot ja takaisinkytketyt verkot sekä tukivektorikone ja K-means-klusterointi. Tulevaisuudessa koneoppimista voidaan hyödyntää mahdollisesti enemmän myös muissa tiedustelun vaiheissa.

Asiasanat: Koneoppiminen, avointen lähteiden tiedustelu, neuroverkot, tiedustelu

Sisällysluettelo

| | | |
|----------|---|-----------|
| 1 | Johdanto | 1 |
| 2 | Avointen lähteiden tiedustelu | 4 |
| 2.1 | Tiedustelun vaiheet | 6 |
| 2.2 | Avointen lähteiden tiedustelun tiedonlähteet | 8 |
| 3 | Koneoppiminen | 10 |
| 3.1 | Koneoppimisen eri tyypit | 10 |
| 3.2 | Koneoppimisen menetelmiä | 12 |
| 4 | Koneoppimisen käyttö avointen lähteiden tiedustelussa | 18 |
| 4.1 | Koneoppimisen käyttö avointen lähteiden tiedustelun eri vaiheissa | 18 |
| 4.2 | Ukrainan ja Venäjän sota sovellusesimerkkinä | 20 |
| 5 | Yhteenveto | 23 |
| | Lähteet | 24 |

Lyhenteet

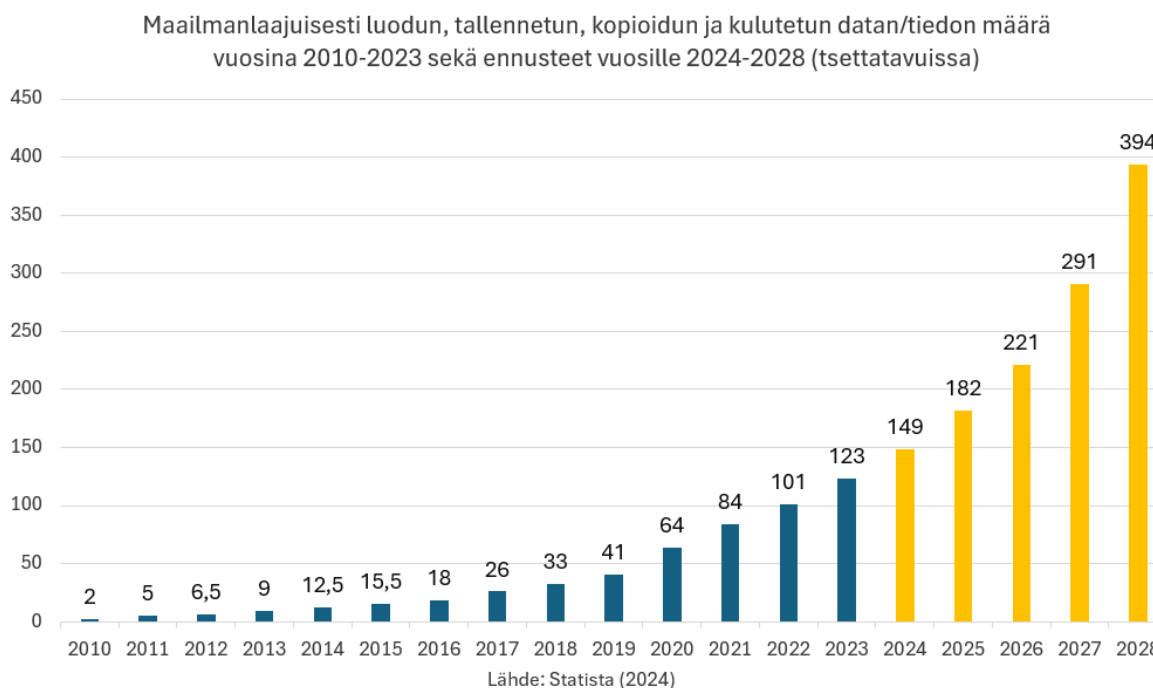
| | |
|----------------|--|
| AI | Tekoäly |
| ANN | Keinotekoiset neuroverkot |
| CNN | Konvoluutioverkko |
| GRU | Venäjän asevoimien sotilastiedusteluosasto |
| CIA | Central Intelligence Agency |
| IC | Intelligence Community |
| LSTM | Pitkäkestoinen lyhytkestomuisti |
| ML | Koneoppiminen |
| NOSINT | Non-Open Source Intelligence |
| OSS | Yhdysvaltain keskustiedustelupalvelu |
| OSINT | Avointen lähteiden tiedustelu |
| RNN | Takaisinkytketty neuroverkko |
| SOCMINT | Sosiaalisen median tiedustelu |
| SVM | Tukivektorikone |
| YOLOv5 | You Only Look Once v5 |
| ZB | Tsettattavu |

1 Johdanto

Tiedustelu mielletään usein valtiollisten toimijoiden, kuten sotilaallisten organisaatioiden toimialueeksi. Rislakin (2024) mukaan ”Tiedustelu ei ole vain valtioiden ja sotilaiden monopoli”. Avoimiin lähteisiin perustuva tiedustelu on kaikkien ulottuvilla, sillä kaikki tarvittava tiedustelumateriaali on vapaasti saatavilla esimerkiksi internetissä (U.S. Director of National Intelligence, 2006).

Ukrainan sodassa avointen lähteiden tiedustelua ja kansalaistiedustelua on hyödynnetty tehokkaasti (Rislakki, 2024). Venäjän joukkojen liikkeitä on voitu seurata sosiaaliseen mediaan kuvatuista videoista, puhelintiedoista ja satelliittikuvista sodan ensimmäisistä hetkistä lähtien. Rislakki antaa esimerkkinä, kuinka aamuyöllä 24. helmikuuta 2022 Googlen karttapalvelusta (Google Maps) pystyttiin näkemään liikenteen ruuhkautuminen Venäjän puolella Ukrainan rajan läheisyydessä. Kolme tuntia myöhemmin Venäjä aloitti hyökkäyksen Ukrainaan.

Avoimista lähteistä saatavan tiedon määrä on kasvanut merkittävästi internetin kasvun myötä. Joka vuosi internetissä olevan tiedon määrä kasvaa kiihtyvällä tahdilla (Kuva 1). Kun vuonna 2010 liikkuvan datan määrä oli vain 2 tsettatavua (ZB), niin viisi vuotta myöhemmin vuonna 2015 se oli jo 15,5 ZB. Vuonna 2020 tiedon määrä kasvoi COVID-19-pandemian takia nopeammin kuin oli odotettu, kun ihmiset tekivät töitä ja opiskelivat etänä. Kokonaisuudessaan luodun, tallennetun, kopioidun ja kulutetun datan määrä vuonna 2020 oli 64 ZB. Kolme vuotta myöhemmin vuonna 2023 oli määrä kasvanut jo 123 ZB:hen. Ennusteiden mukaan datan määrä jatkaa kasvuaan ja vuonna 2028 sen on ennustettu kasvavan jo 394 ZB:hen (Taylor, 2024).



Kuva 1 Datan määrä internetissä, muokattu lähteestä (Taylor, 2024)

Samalla kun internetissä olevan tiedon määrä on kasvanut, myös tiedon käsittelyyn tarvittava aika on kasvanut (Withorne, 2022). On mahdotonta käsitellä ja analysoida kaikkea saatavilla olevaa relevanttia tietoa manuaalisesti, kun tietoa on saatavilla niin paljon. Siksi nykyään kehitetään työkaluja tiedonkäsittelyn automatisoimiseksi. Koneoppiminen on yksi työkaluista, jota hyödynnetään tietojenkäsittelyn automaatioon. Koneoppimisella tarkoitetaan tekoälyn osaluuetta, jossa tietokone opetetaan löytämään malleja ja tekemään tulkintoja datan perusteella ilman erillistä ohjelmointia. Koneoppiminen mahdollistaa suurten tietomäärien prosessoimisen ja sen vuoksi sitä pidetään tehokkaana työkaluna tietojenkäsittelyn automatisoimiseksi. Kuten entinen tiedustelupäällikkö Pekka Toveri on todennut, ”Kunkku on se, joka pystyy valtavasta datamassasta analysoimaan olennaisen” (Rislakki, 2024)

Tutkielman pääkäsitteitä, avointen lähteiden tiedustelua ja koneoppimista on tutkittu akateemisessa kirjallisuudessa kattavasti. Näitä käsitteitä ei ole kuitenkaan tutkittu yhdessä laajasti, sillä koneoppimista on alettu hyödyntämään avointen lähteiden tiedustelussa enemmän vasta viime vuosina koneoppimisen kehityksen mukana. Tutkielman tarkoituksena on selvittää koneoppimisen käyttöä avointen lähteiden tiedustelussa.

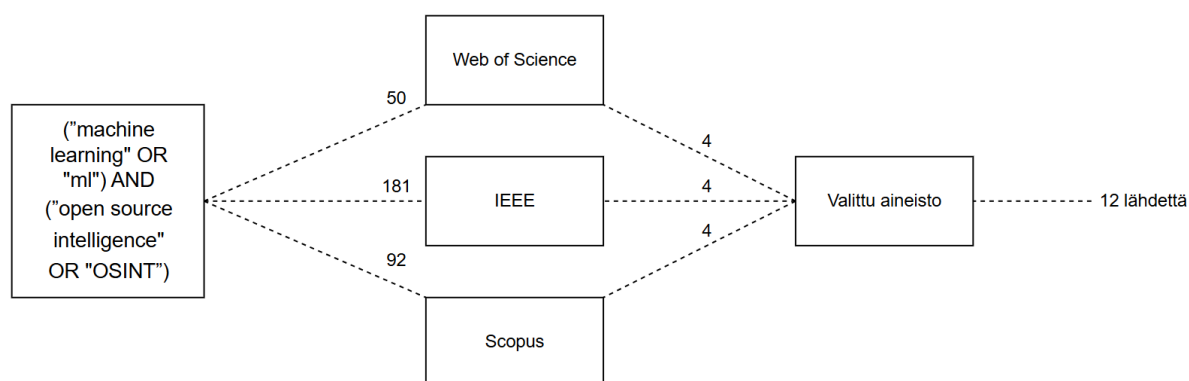
Tutkielman tutkimuskysymykset ovat:

Tk1: Mitä koneoppimisen menetelmiä nykyään käytetään avointen lähteiden tiedustelussa?

Tk2: Mitä koneoppimisen menetelmiä voidaan tulevaisuudessa käyttää avointen lähteiden tiedustelussa?

Tk3: Millaisissa tilanteissa koneoppimista voidaan hyödyntää avointen lähteiden tiedustelussa?

Tietokantahaku suoritettiin seuraavasti. Tietokantoina toimivat Web of science, IEEE, ja Scopus. Hakulausekkeena kaikissa tietokannoissa käytettiin ("machine learning" OR "ml") AND ("open source intelligence" OR "OSINT"). Haku tehtiin 6.3.2025 ja tuloksia löytyi tietokannoista yhteensä 323 (Kuva 2). Web of Sciencestä löytyi 50 tulosta, IEEE:stä 181 ja Scopusesta 92. Kaikista tuloksista otsikon perusteella valittiin 30 tarkempaan tarkasteluun. Näistä lähteistä valikoitui 12 lähdetä tiivistelmän ja muun sisällön perusteella. Lisäksi käytettiin kahta kirjaa, jotka valikoituivat sisällön sopivuuden perusteella.



Kuva 2 Aineiston haku

Seuraavassa luvussa käsitellään avointen lähteiden tiedustelua. Luvussa käydään läpi tarkemman määritelmän lisäksi tiedustelun vaiheet sekä tiedonlähteet. Kolmannessa luvussa käsitellään tutkimuksen toista pääkäsitettä, koneoppimista. Siinä kerrotaan koneoppimisen eri muodoista ja menetelmistä. Neljäs luku käsittelee näitä molempia, yhdistäen koneoppimisen ja avointen lähteiden tiedustelun. Luvussa tarkastellaan koneoppimisen käyttöä avointen lähteiden tiedustelun eri vaiheissa. Neljännen luvun lopussa tarkastellaan teemaa konkreettisemmin Ukrainan sodan kontekstissa. Yhteenvedossa vastataan aiempien lukujen pohjalta tutkimuskysymyksiin.

2 Avointen lähteiden tiedustelu

Avointen lähteiden tiedustelu (engl. Open Source INTelligence, OSINT) on julkisesti saatavilla olevan tiedon järjestelmällistä keräämistä, hyödyntämistä ja analysointia, siten että saadaan tuotettua ajankohtaista sekä kohdennettua tiedustelutietoa. Avoimilla lähteillä tarkoitetaan mitä tahansa julkisesti saatavilla olevaa tietoa, jota kuka tahansa voi laillisesti hankkia. (U.S. Director of National Intelligence, 2006). Esimerkiksi avoimesti saatavilla olevien satelliittikuvien tai sosiaalisen median kuvamateriaalin analysointi vihollisjoukkojen liikkeiden havaitsemiseksi on avointen lähteiden tiedustelua.

Avointen lähteiden tiedustelun historia juontaa juurensa toiseen maailmansotaan ja Yhdysvaltain keskustiedustelupalvelun (Office of Strategic Servicesin, OSS) perustamiseen (Colquhoun, 2016). OSS, jota pidetään OSINTin edelläkävijänä, loi ensimmäisenä osaston, jonka tehtävänä oli analysoida avoimista lähteistä saatavaa tietoa. Myöhemmin OSS:n seuraajaksi perustettiin Central Intelligence Agency (CIA), joka jatkoi ja kehitti edelleen avointen lähteiden tiedustelutoimintaa. OSINT on yksi tiedustelutoiminnan keskeisistä osa-alueista, ja sen merkitys on kasvanut erityisesti kylmän sodan päättymisen jälkeen. Arvioita on erilaisia, mutta joidenkin arvioiden mukaan jopa 70-80 % Yhdysvaltain tiedustelutiedoista on tullut OSINTin kautta (Hulnick, 2002). On todennäköistä, että OSINTin merkitys on kasvanut aiemmasta arviosta internetin kasvun ja avoimen tiedon saatavuuden lisääntyä.

Avointen lähteiden tiedustelua voidaan hyödyntää monilla eri aloilla, kuten lainvalvonnassa valtionhallinnossa (Yadav ym., 2023). Lainvalvonnassa poliisit ja viranomaiset voivat hyödyntää OSINTia rikostutkinnassa, kun selvitetään esimerkiksi ihmiskauppa- ja rahanpesurikoksia. Valtionhallinnossa erityisesti sotilas- ja turvallisuusviranomaiset voivat hyödyntää OSINTia vastatiedusteluun tai terroriuhkien havaitsemiseen.

OSINTia hyödyntävät myös yksityishenkilöt ja ei-valtiolliset toimijat, kuten Bellingcat (Rislakki, 2024). Eliot Higginsin aloittama Bellingcat-projekti on kansainvälinen avoimia lähteitä hyödyntävä tutkivan journalismin verkosto. Bellingcat on tullut tunnetuksi selvitettyään lennon MH17 pudonneen Venäjän asevoimille kuuluneen ohjuksen osuttua koneeseen. Maaliskuussa 2018 entinen Venäjän asevoimien sotilastiedustelun (GRU) upseeri Sergei Skripal myrkytettiin novitshok-hermomyrkyllä Etelä-Englannin Salisburyn kaupungissa. Oman raporttinsa mukaan Bellingcat selvitti Skripalin myrkyttäneiden GRU-upseereiden henkilöllisyydet (Bellingcat Investigation Team, 2018).

Toisen Bellingcatin raportin mukaan he onnistuivat myös selvittämään niiden GRU-upseereiden henkilöllisyydet, jotka elokuussa 2020 myrkyttivät novitshok-hermomyrkyllä Venäjän oppositiojohtajan Aleksei Navalnyin (Bellingcat Investigation Team, 2020).

Avointen lähteiden tiedustelulla on hyviä ja huonoja puolia, jotka vaikuttavat sen hyödyllisyyteen ja soveltamiseen tiedustelutoiminnassa. Yksi hyvistä puolista on tiedon keräämisen helppous (Ivanjko & Dokman, 2019). Tiedon kerääminen avoimista lähteistä on nopeaa ja digitaalista tietoa voidaan kerätä helposti suuria määriä. OSINTin käyttöön ei liity samoja fyysisiä tai poliittisia riskejä kuin esimerkiksi henkilötiedusteluun, jossa tiedonhankinta perustuu salaisten agenttien tai rekrytoitujen kontaktien toimintaan. Tällaisissa tilanteissa tietolähteen paljastuminen voi johtaa vakaviin poliittisiin seurauksiin tai jopa hengenvaaraan.

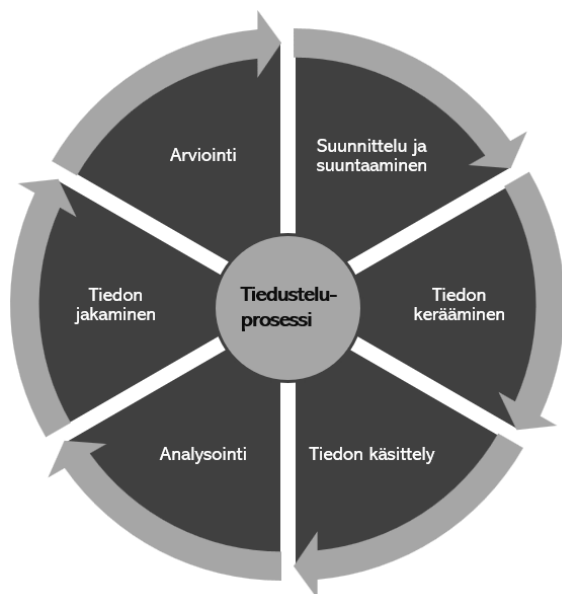
OSINT on myös kustannustehokas tiedustelumenetelmä verrattuna esimerkiksi henkilötiedusteluun tai signaalitiedusteluun (Ivanjko & Dokman, 2019). Jälkimmäisissä menetelmissä tiedonhankinta edellyttää henkilöstön koulutusta ja tietolähteiden värväämistä sekä kalliita teknisiä laitteita (Hwang ym., 2022). OSINT hyödyntää pääosin avoimia ja usein maksuttomia tietolähteitä, kuten uutismedioita, sosiaalista mediaa ja verkkosivustoja. Avointen lähteiden tiedustelua voidaan lisäksi tehdä yksinkertaisella laitteistolla, kuten älypuhelimella tai tavallisella tietokoneella. Lisäksi avoimista lähteistä kerätty tiedustelutieto on usein ajantasaista ja jopa lähes reaaliaikaista. Esimerkiksi sosiaalisesta mediasta voidaan saada tilannepäivityksiä nopeammin kuin henkilötiedustelun kautta.

Hyvien puolien lisäksi avointen lähteiden tiedusteluun liittyy myös haasteita, jotka vaikuttavat sen luotettavuuteen ja käytettävyyteen. Yksi merkittävistä ongelmista on tiedon luotettavuus (Ivanjko & Dokman, 2019). Koska kuka tahansa voi tuottaa ja julkaista sisältöä internetissä voi avoimista lähteistä kerätty tieto olla keskenään ristiriitaista tai harhaanjohtavaa. Disinformaatio ja valeuutiset vaikuttavat OSINTin tehokkuuteen ja tiedustelutiedon luotettavuuteen. Lähteiden kriittinen arviointi ja tietojen varmistaminen muista lähteistä ovat olennaisia lopputuloksen kannalta. Vaikka tiedon kerääminen avoimista lähteistä on nopeaa, voi tietojen käsittelyyn kulua paljon aikaa (Hwang ym., 2022). Kerätty tieto voi olla hajanaista, epäselvää tai vaikeasti tulkittavissa. Erityisenä haasteena on suurien datamäärien analysointi, kun seassa voi olla disinformaatiota ja tarkoituksellisesti vääristeltyä dataa.

2.1 Tiedustelun vaiheet

Tiedustelutoiminta perustuu järjestelmälliseen prosessiin, jonka avulla saadaan tuotettua tiedustelutietoa. Tiedustelun vaiheiden mallintamiseen on kehitetty erilaisia lähestymistapoja, siksi tiedustelun vaiheet vaihtelevat mallin mukaan. Eräs yleisesti käytetty tiedustelun vaiheita kuvaava malli on Intelligence Communityn (IC) luoma tiedusteluprosessi. Tiedusteluprosessin tarkoituksena on tiedustelulähteistä saadun tiedon prosessointi käyttökelpoiseksi tiedustelutiedoksi (Office of the Director of National Intelligence, 2011) (katso Laatikko 1).

Gibson (2016) on hyödyntänyt IC:n kehittämää tiedusteluprosessin mallia tutkimuksessaan ja luonut visuaalisen esityksen tiedusteluprosessin vaiheista. Kuva 3 havainnollistaa tiedusteluprosessin kuutta eri vaihetta, jotka ovat **suunnittelu ja suuntaaminen, tiedon kerääminen, tiedon käsittely, analysointi, tiedon jakaminen ja arviointi**.



Kuva 3 Tiedusteluprosessin vaiheet (Gibson, 2016), muokattu lähteestä (Gibson, 2016)

Laatikko 1 Tiedusteluprosessin vaiheet (Office of the Director of National Intelligence, 2011)

1. **Suunnittelu ja suuntaaminen:** Taustoitetaan, minkälaista tiedustelutietoa halutaan ja mihin tarkoitukseen. Kun on tiedossa mitä tietoa halutaan, täytyy suunnitella mistä ja miten kyseistä tietoa hankitaan. Suunnitteluvaiheessa päätetään myös minkälaista tiedustelutietoa halutaan tuottaa. Suunnitelman pohjalta lähdetään jatkamaan tiedusteluprosessia.
2. **Tiedon kerääminen:** Kun tiedustelun vaatimukset ja suunnitelma ovat tiedossa, kerätään tietoa tiedustelulähteistä, jotka soveltuvat parhaiten suunnitelman toteuttamiseen. Viisi perinteistä tiedustelutiedon lähdettä ovat geospaatialinen tiedustelu (GEOINT), henkilötiedustelu (HUMINT), mittaus- ja tunnusmerkkitiedustelu (MASINT), avointen lähteiden tiedustelu (OSINT) ja signaalitiedustelu (SIGINT).
3. **Tiedon käsittely:** Tiedustelulähteistä kerätty tieto muunnetaan ymmärrettävään ja käyttökelpoiseen muotoon. Tiedon käsittelyyn käytetään erikoistunutta henkilöstöä ja teknisiä laitteita, joilla tieto saadaan analysoitavaan muotoon. Tähän vaiheeseen kuuluvat esimerkiksi kielenkäännökset, salauksien purkaminen sekä kuvamateriaalin tulkitseminen.
4. **Analysointi:** Edellisessä vaiheessa käsitelty tieto valmistellaan lopulliseen muotoon, esimerkiksi tiedusteluraportiksi. Joissain tapauksissa tämä vaihe voidaan ohittaa, jos tiedustelutietoa ei tarvitse analysoida. Esimerkiksi lokakuussa 1962, Kuuban ohjuskriisin aikana Yhdysvaltojen presidentti John F. Kennedy tarvitsi vain Kuubassa sijaitsevan Neuvostoliiton kaluston määrän, muttei tiedon pohjalta tehtyä analyysiä.
5. **Tiedon jakaminen:** Prosessoitu ja analysoitu tieto toimitetaan tiedon tilaajalle sekä mahdollisesti muille tarvittaessa. Kun valmis tiedustelutieto on jaettu, siinä saatetaan havaita aukkoja, jolloin tiedusteluprosessi alkaa alusta.
6. **Arviointi:** Tiedusteluprosessin aikana kerätään palautetta prosessin eri vaiheista, joilla toimintaa voidaan kehittää ja tehostaa. Valmistaa tiedustelutietoa voidaan arvioida sen perusteella kuinka hyödyllistä sen on tai jäikö siihen aukkoja. Arviointivaihe on tärkeä osa tiedusteluprosessia, koska sen avulla tiedustelutiedon tuottamista voidaan jatkuvasti parantaa muuttuvissa tilanteissa.

Tiedusteluprosessia hyödynnetään kaikissa tiedustelumenetelmissä käyttökelpoisen tiedustelutiedon tuottamiseksi. Vaikka prosessin vaiheet ovat kaikille tiedustelumenetelmille samat, niiden tarkka sisältö vaihtelee käytetyn menetelmän mukaan. OSINTissa tiedusteluprosessin kulkuun vaikuttaa se, että käytetty tieto on julkisesti saatavilla. Suunnittelu ja suuntaus -vaiheessa korostuu se, mitä etsitään ja mitä tunnisteita tai lähteitä käytetään (Böhm & Lolagar, 2021). Tämä eroaa esimerkiksi henkilötiedustelusta, jossa tärkeässä roolissa ovat luottamukselliset kontaktit ja henkilöverkot.

Tiedon keräämisvaiheessa OSINT eroaa muista menetelmistä siinä, että kaikki tieto hankitaan avoimista lähteistä (Ivanjko & Dokman, 2019). Esimerkiksi signaalitiedustelussa tietoa kerätään sieppaamalla viestiliikennettä usein salassa, ja joskus jopa laittomasti.

Tiedon käsittely -vaiheessa OSINTin tarpeet tiedon käsittelemiselle on usein kevyempiä kuin esimerkiksi signaalitiedustelussa. Signaalitiedustelussa kerätty data voi olla salattua tai vaatia paljon muuta käsittelyä ja muotoilua ennen kuin se on analysoitavassa muodossa. OSINTissa tiedon käsittelyyn voi kuulua esimerkiksi satelliittikuvien asettelua helpommin tulkittavaan muotoon, kuten ennen ja jälkeen -vertailuun.

Analyysivaiheessa tieto asetetaan kontekstiin ja sen luotettavuutta arvioidaan (Böhm & Lolagar, 2021). Tämä vaihe ei poikkea merkittävästi menetelmän mukaan, vaikka käytetyt analyysitekniikat voivat vaihdella tiedon tyyppin mukaan. Esimerkiksi kuvadatan analyysi eroaa suullisten lausuntojen tulkinnasta. Tiedon jakamisvaiheessa ainoa suuri ero muihin menetelmiin on se, että OSINT-menetelmällä luodut tiedusteluraportit ovat useammin julkisia tai vapaasti saatavilla muille toimijoille. Tämä johtuu siitä, että OSINTia tekevät valtiollisten toimijoiden lisäksi monet yksityiset tahot ja organisaatiot, jotka haluavat jakaa tietoa eteenpäin. Arviointivaiheessa OSINT ei eroa suuresti muista tiedustelumenetelmistä.

2.2 Avointen lähteiden tiedustelun tiedonlähteet

Avointen lähteiden tiedustelun tiedonlähteet ovat muuttuneet vuosien saatossa. Aluksi tiedustelutietoa saatiin sanomalehdistä, artikkeleista, valokuvista ja radiolähetysistä (Colquhoun, 2016). Myöhemmin internetin kehittymisen ja kasvun myötä tietolähteitä on tullut lisää. OSINT-tietoa voidaan kerätä nykyään monista digitaalisista lähteistä, kuten sosiaalisesta mediasta, videonjakosivustoilta, hakukoneista ja verkkosivuilta (Hassan & Hijazi, 2018).

OSINT kattaa kaikki julkisesti saatavilla olevat tiedonlähteet, ja ne voidaan jakaa kahteen pääluokkaan: verkon ulkopuolisiin (engl. offline) ja verkkopohjaisiin (engl. online) lähteisiin (Browne ym., 2024). Verkon ulkopuolisia OSINT-lähteitä ovat esimerkiksi radio, sanomalehdet, osoitehakemistot ja televisio.

Verkkopohjaisia lähteitä ovat puolestaan verkkosivustot, sosiaalinen media, verkkotietokannat, avoimet tietoaaineistot, akateemiset kokoelmat, pimeä verkko, verkkotunnusrekisterit, valtiolliset verkkosivut, keskustelufoorumit, hakukoneet ja ohjelmointirajapinnat (API:t) (Browne ym., 2024). Toisin sanoen verkkopohjaisia lähteitä ovat kaikki internetissä julkisesti saatavilla olevat tiedonlähteet.

Korkean resoluution satelliittikuvia tarjoavia kaupallisia toimijoita on useita. Hassan & Hijazi (2018) mainitsevat suosituimpina palveluntarjoajina European Space Imagingin ja Digital Globen. Näiden satelliittipalveluiden avulla yksityiset toimijat, yritykset ja tutkijat voivat seurata reaaliajassa maailman tapahtumia, kuten sotilaallisia liikkeitä ja luonnonkatastrofien vaikutuksia ympäri maailmaa (Hassan & Hijazi, 2018).

Sosiaalisen median kasvu on tuonut mukanaan uuden tiedustelumenetelmän. Sosiaalisen median tiedustelu (engl. SOCial Media INTelligence, SOCMINT) tarkoittaa tiedon keräämistä ja analysointia sosiaalisen median alustoilta, kuten Facebookista, Instagramista, YouTubesta ja TikTokista. Sosiaalisiksi mediaksi määritellään kaikki internetissä olevat sovellukset ja sivustot, joissa käyttäjät ovat vuorovaikutuksessa keskenään ja voivat jakaa tietoa toisilleen (Omand, 2017). SOCMINTia voidaan pitää OSINTin alle kuuluvana käsitteenä, sillä sosiaalinen media kattaa suuren osan avoimista lähteistä joita OSINTissa käytetään (Böhm & Lolagar, 2021).

Vaikka OSINTiin kuuluu kaikki julkisesti saatavilla oleva tieto, voi tiedonhaun aikana löytyä myös huonosti suojattua tietoa ja salassa pidettäväksi tarkoitettua tietoa. Tällaisia ovat esimerkiksi Wikileaksin julkaisemat materiaalit tai pelifoorumeilla vuodetut asiakirjat (Allison, 2024; Hassan & Hijazi, 2018). Tätä kaikkien saatavilla olevaa tietoa, joka ei ole virallisesti julkista kutsutaan NOSINTiksi (Non-Open Source INTelligence).

3 Koneoppiminen

Koneoppiminen (engl. Machine Learning, ML) on tekoälyn osa-alue. Se keskittyy menetelmiin, joiden avulla tietokonejärjestelmien on mahdollista oppia ja tehdä päätöksiä ilman erikseen ohjelmoituja ohjeita (Gerard, 2021). Koneoppimisalgoritmit mukautuvat syötteenä saadun datan perusteella ja muodostavat itsenäisesti päätelmiä siitä. Tekoäly (engl. Artificial Intelligence, AI) on kattokäsite, joka viittaa tietokoneiden kykyyn suorittaa tehtäviä, jotka vaativat älykkyyttä. Tällaisia tehtäviä ovat esimerkiksi puheentunnistus, itseohjautuvien ajoneuvojen ajaminen ja strategisten pelien, kuten shakin, pelaaminen. Koneoppiminen on yksi tapa toteuttaa tekoälyä, ja sen avulla tietokonejärjestelmät voivat löytää datasta säännönmukaisuuksia ja tehdä ennusteita uusista tilanteista. Tässä luvussa käsitellään yleisimpiä avointen lähteiden tiedustelussa käytettyjä koneoppimisen menetelmiä.

3.1 Koneoppimisen eri tyypit

Koneoppimisessa on useita eri oppimismenetelmiä. Kolme keskeisintä oppimismenetelmää ovat ohjattu oppiminen, ohjaamaton oppiminen ja vahvistusoppiminen (Hiran ym., 2021). Oppimismenetelmän valinta riippuu käyttökohteesta ja saatavilla olevasta datasta.

Ohjattu oppiminen (engl. supervised learning) hyödyntää etukäteen valmistettua opetusdataa, jota käytetään koneoppimisalgoritmin kouluttamiseen (Hiran ym., 2021). Käytetty opetusdata on merkittävä eli se sisältää syötteet ja niihin liittyvät oikeat tulokset. Koulutusprosessin aikana koneoppimisalgoritmi etsii opetusdatasta piirteitä ja yhdistää niitä haluttuihin tuloksiin. Ohjatussa oppimisessa algoritmi koulutetaan oikeiksi tiedettyjen esimerkkien avulla.

Ohjattua oppimista voidaan hyödyntää, jos halutaan kehittää koneoppimismalli, joka kykenee esimerkiksi tunnistamaan ajoneuvoja satelliittikuvista. Ensin koneoppimismallille syötetään suuri määrä satelliittikuvia, joissa on erilaisia ajoneuvoja, kuten panssarivaunuja, kuorma-autoja ja siviiliajoneuvoja. Kuvat on luokiteltu sen mukaan mikä ajoneuvo niissä näkyy. Näin algoritmi oppii mitä eri piirteitä ajoneuvotyypeillä on. Tämän jälkeen koulutettu koneoppimisalgoritmi pystyy arvioimaan uusista satelliittikuvista mikä ajoneuvo on kyseessä.

Ohjattu oppiminen jakautuu kahteen pääryhmään: luokittelu ja regressio (Kan, 2017). Jos aineisto koostuu diskreeteistä luokista, eli rajallisesta määrästä selkeästi erilaisia vaihtoehtoja, kyseessä on luokittelu. Jos taas aineisto on jatkuvaa, eli arvoja on loputon määrä, käytetään regressiota.

Ohjaamaton oppiminen (engl. unsupervised learning) on koneoppimisen oppimismenetelmä, jossa käytetään merkitsemätöntä dataa (Kan, 2017). Koneoppimisalgoritmi oppii löytämään säännönmukaisuuksia aineistosta, jossa ei ole luokkia tai valmiiksi annettuja vastauksia (Hiran ym., 2021). Ohjaamattomassa oppimisessä ei käytetä opetusdataa, vaan algoritmille annetaan syötedataa, jota halutaan analysoida.

Klusterointi (engl. clustering) on tärkeä ohjaamattoman oppimisen menetelmä. Klusterointialgoritmi jakaa merkitsemättömän datan eri ryhmiin samankaltaisuuden perusteella. Jokainen ryhmä sisältää ne objektit, jotka muistuttavat eniten toisiaan ja eroavat piirteiltään muiden ryhmien objekteista. Esimerkiksi algoritmille voidaan antaa kuvia eri ajoneuvoista, ilman tietoa siitä, mitä ne esittävät. Algoritmi etsii samankaltaisuuksia kuvista ja muodostaa niistä ryhmiä piirteiden perusteella. Koska kyseessä on ohjaamaton oppiminen, algoritmi ei osaa nimetä ryhmiä. Ryhmä 1 voi sisältää henkilöautot ja ryhmä 2 kuorma-autot.

Vahvistusoppiminen (engl. reinforcement learning) on oppimismenetelmä, jossa koneoppimisalgoritmi eli agentti oppii vuorovaikutuksesta ympäristön kanssa (Hiran ym., 2021). Oppiminen tapahtuu palkkioiden ja rangaistusten kautta. Kun agentti suorittaa toimintoja ympäristössään, se saa palkkion tai rangaistuksen sen perusteella oliko toiminta haluttua vai ei. Vahvistusoppimisessä oppiminen tapahtuu palautteen ja kokemuksen kautta ilman merkittävää dataa.

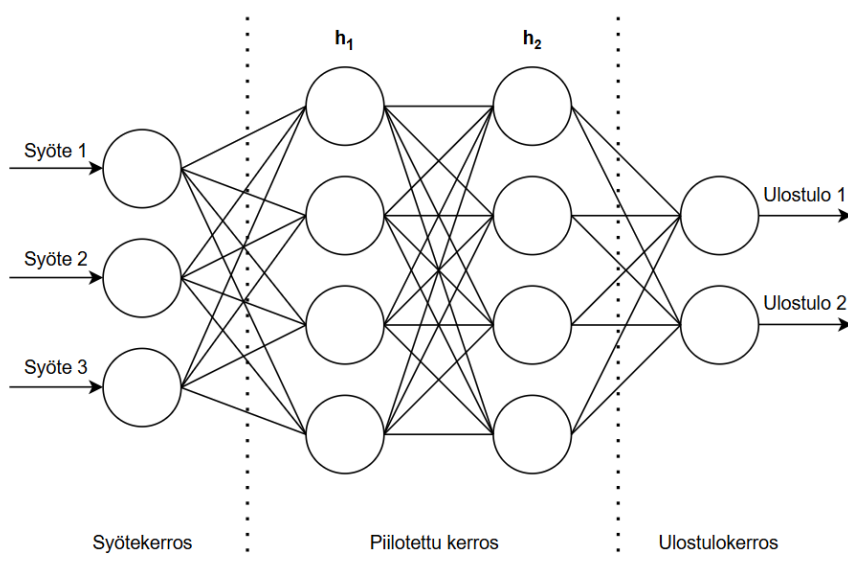
Konkreettinen esimerkki vahvistusoppimisesta on koiran kouluttaminen. Koira ei ymmärrä ihmiskieltä, joten sille ei voida antaa suoria ohjeita mitä pitää tehdä. Sen sijaan voimme palkita koiran oikeanlaisesta käyttäytymisestä. Ajan myötä koira oppii, mitkä toiminnot tuottavat palkkion. Samalla tavoin vahvistusoppimisessä agentin tavoitteena on maksimoida saamansa palkkiot, joten se oppii oikeat toiminnot. Esimerkiksi agenttina voi toimia drone, joka opetetaan keräämään tiedustelumateriaalia itsenäisesti. Arranz ym. (2023) esittelevät tutkimuksessaan droneparven, joka voidaan opettaa syvän vahvistusoppimisen avulla valvomaan alueita ja seuraamaan kohteita ilman ihmisen ohjausta. Tehtyjen simulaatioiden perusteella todettiin, että tämän menetelmän avulla voidaan opettaa droneparvi suorittamaan valvontatehtäviä ilman jatkuvaa ohjausta.

3.2 Koneoppimisen menetelmiä

Koneoppimisessa on useita erilaisia menetelmiä ja algoritmeja, joita käytetään eri käyttötarkoituksiin. Oikean menetelmän valinta on tärkeää lopputuloksen kannalta. Eri Koneoppimismenetelmät on kehitetty tietyn tyyppisen datan käsittelemiseen. Osa menetelmistä soveltuu hyvin kuvamuotoisen datan käsittelyyn, kun taas toiset toimivat paremmin muunlaisen datan, kuten tekstin kanssa (Gerard, 2021).

Keinotekoiset neuroverkot (engl. Artificial Neural Networks, ANN) kuuluvat syväoppimiseen (engl. deep learning), joka on koneoppimisen osa-alue (Gerard, 2021). Keinotekoiset neuroverkot on suunniteltu jäljittelemään aivojen hermoverkkojen toimintaa. Ne koostuvat toisiinsa kytketyistä neuroneista, jotka välittävät signaaleja eteenpäin muille neuroneille.

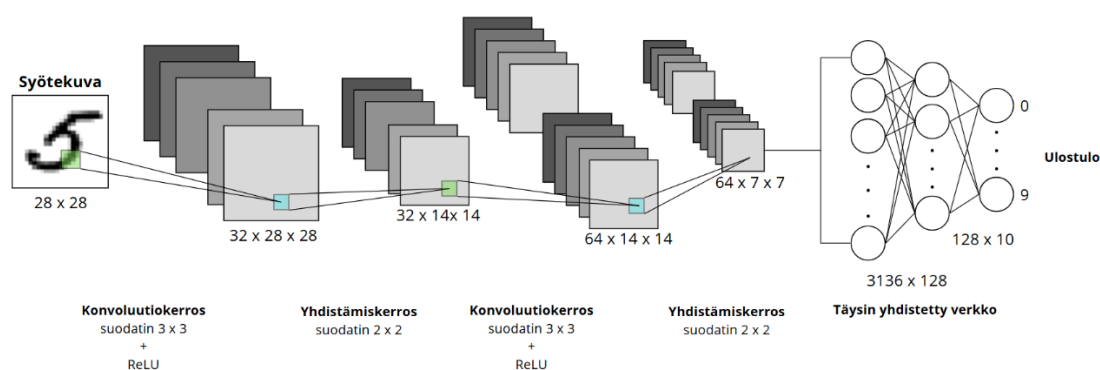
Keinotekoinen neuroverkko (Kuva 4) koostuu syötekerroksesta, yhdestä tai useammasta piilokerroksesta sekä ulostulokerroksesta (Bre ym., 2017). Kyseessä on eteenpäin kytketty neuroverkko, jossa tieto kulkee yksisuuntaisesti ilman takaisinkytkentää. Kerroksissa olevat keinotekoiset neuronit ovat yhteydessä toisiinsa painokertoimien kautta. Painokertoimia voidaan muuttaa koulutusprosessin aikana syötteen ja halutun ulostulon mukaan. Painokertoimet määrittävät kuinka suuri vaikutus kullakin syötteellä on neuroverkon lopulliseen tulokseen. Perinteisen keinotekoisin neuroverkon lisäksi on olemassa erilaisia neuroverkkorakenteita, joita on kehitetty eri tarkoituksiin. Konvoluutioverkko ja takaisinkytketty neuroverkko ovat erikoistuneita neuroverkkorakenteita.



Kuva 4 Neuroverkko. Mallinnettu lähteestä (Bre ym., 2017)

Konvoluutioverkko (engl. Convolutional Neural Network, CNN) on eteenpäin kytketty neuroverkko, joka hyödyntää konvoluutio-, yhdistämis- ja ReLU-kerroksia (engl. convolution, pooling ja Rectified Linear Unit) (Emmert-Streib ym., 2020). Konvoluutiokerrokset erottavat syötteestä erilaisia piirteitä (Hijazi ym., 2015). Ensimmäinen konvoluutiokerros erottaa matalan tason piirteitä, kuten reunaviivoja tai kulmia. Kerrokset siirtyvät korkeamman tason tunnistukseen, kunnes ne pystyvät erottamaan monimutkaisia rakenteita, kuten kasvojen osia tai tekstuureja. Konvoluutiokerroksessa kuvan yli liikutetaan pienempää matriisisuodatinta (engl. filter tai kernel). Suodatin on joukko painoja, jotka määritetään koulutusvaiheessa. Suodatin käy läpi kuvan eri kohdat ja tulokset tallennetaan piirrekuvaan (engl. feature map). Yhdistämiskerros on konvoluutioverkon osa, jonka tehtävänä on pienentää konvoluutiokerroksien luomien piirrekuvien kokoa säilyttäen tärkeät yksityiskohdat, joita kuva sisältää. Tämän avulla voidaan vähentää kuvassa esiintyvän kohinan ja vääristymien vaikutusta sekä pienentää konvoluutioverkon vaatimaa laskentatehoa (Guo ym., 2020). Yhdistämiskerros toimii siten että sille annettu syöte kuten piirrekuva jaetaan esimerkiksi 2×2 -kokoisiin ei-päällekkäisiin alueisiin ja kunkin alueen arvo tiivistetään yhdeksi arvoksi.

Perinteisissä keinotekoisissa neuroverkoissa jokainen neuroni on yhteydessä kaikkiin seuraavan kerroksen neuroneihin, mikä voi johtaa suureen määrään painokertoimia (Emmert-Streib ym., 2020). Konvoluutioverkossa jokainen neuroni on yhteydessä vain lähellä oleviin neuroneihin seuraavassa kerroksessa, joka voi vähentää painokertoimien määrää ja nopeuttaa algoritmin toimintaa. Perinteisissä konvoluutioverkoissa käytetään täysin yhdistettyä kerrosta (engl. fully connected layer) ennen viimeistä ulostulokerrosta (Kuva 5). Täysin yhdistetty kerros yhdistää kaikki edellisen kerroksen neuronit seuraavan kerroksen neuroneihin, joka mahdollistaa ei-lineaaristen suhteiden mallintamisen syötteessä.



Kuva 5 Perinteinen konvoluutioverkko. Mallinnettu lähteestä (Gerard, 2021)

Nykyään kuitenkin käytetään enemmän vaihtoehtoisia ratkaisuja, jotka eivät lisää verkkoon suurta määrää painokertoimia. Täysin yhdistetyn kerroksen sijaan voidaan käyttää esimerkiksi Global Average Pooling (GAP) -kerrosta (Guo ym., 2020). GAP-kerros pystyy tiivistämään kunkin piirrekartan tiedot yhdeksi arvoksi, ilman että ne muutetaan yksiulotteiseen muotoon. Tämä vähentää laskennallista kuormaa ja tekee mallista tehokkaamman. Konvoluutioverkot soveltuvat erityisen hyvin kuvien luokitteluun ja niitä voidaan käyttää esimerkiksi esineiden tunnistukseen tai kasvojentunnistukseen (Gerard, 2021).

Takaisinkytketty neuroverkko (engl. Recurrent Neural Network, RNN) on neuroverkkoarkkitehtuuri, jossa eri kerrosten välillä on takaisinkytkentöjä. Ne soveltuvat erityisen hyvin jaksolliseen dataan, kuten tekstin, puheen tai aikasarjojen käsittelyyn ja analysoimiseen (Al-Selwi ym., 2024). Poiketen tavallisista eteenpäin suunnatuista neuroverkoista, RNN voi käyttää sisäistä muistia syötteiden käsittelemiseen. Tämä muistirakenne mahdollistaa aikaisempien syötteiden vaikutuksen huomioimisen nykyhetken päätöksiin.

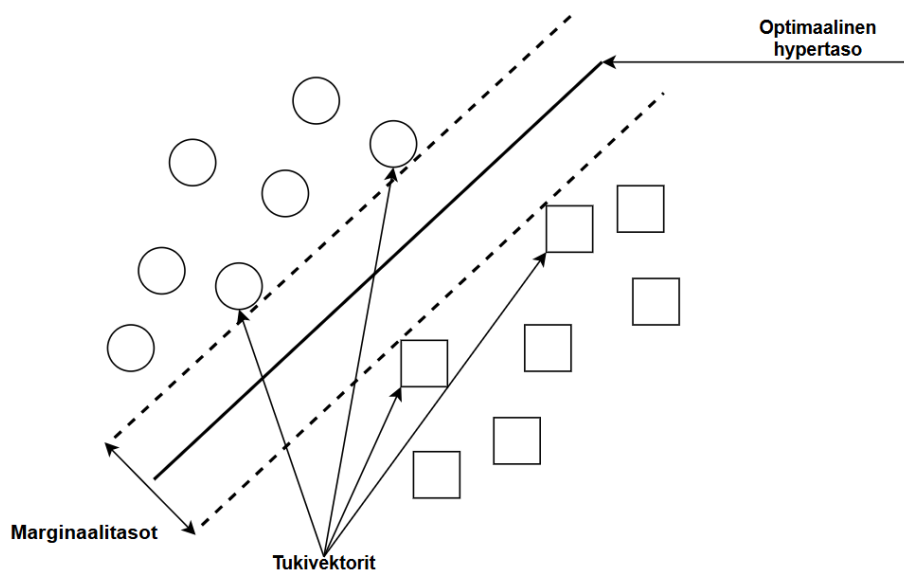
RNN-mallien haasteena on kuitenkin muistin ylläpitäminen, joka aiheuttaa niin sanotun katoavan gradientin ongelman (Sen & Mehtab, 2022). Siinä neuroverkkoa opetettaessa varhaisempien ja myöhempien kerrosten gradientit eroavat toisistaan suuresti. Katoavan gradientin ongelmaan on kehitetty ratkaisuna pitkäkestoinen lyhytkestomuisti (engl. Long Short Term Memory, LSTM). Se on muokattu RNN-malli, joka mahdollistaa pitkien syy-seurausketjujen oppimisen. LSTM-verkko sisältää muistiporotteja, jotka säätelevät mitä tietoa säilytetään tai poistetaan. LSTM on osoittautunut erityisen tehokkaaksi jaksollisen datan käsittelyssä ja on laajimmin käytetty RNN-malli (Al-Selwi ym., 2024).

RNN:t soveltuvat myös konenäön sovelluksiin, kuten ihmisen asennon tunnistamiseen, poikkeavan käyttäytymisen havaitsemiseen väkijoukosta, sekä kuvatekstien tuottamiseen tai videoiden tekstiselosteiden generoimiseen (Al-Selwi ym., 2024). Tällaisissa tapauksissa LSTM-malleja käytetään usein yhdessä konvoluutioverkkojen kanssa. Konvoluutioverkko tunnistaa visuaaliset piirteet, jotka syötetään LSTM-verkolle analysoitaviksi. Esimerkiksi Xue ym. (2019) esittelivät ST-HConvLSTM-mallin, joka on suunniteltu videoiden tapahtumien tunnistamiseen. Mallissa yhdistetään ajallis-paikallinen huomiointimoduuli ja hierarkkinen konvoluutio-LSTM-arkkitehtuuri (HConvLSTM).

Huomiointimoduuli auttaa mallia keskittymään merkityksellisiin kohtiin, kuten videolla esiintyvään henkilöön tai tiettyihin ajankohtiin liikkeen aikana (Xue ym., 2019). HConvLSTM puolestaan kykenee mallintamaan toiminnan rakenteen sekä ajassa että tilassa, mikä parantaa mallin kykyä ymmärtää monimutkaisia tapahtumaketjuja videoissa.

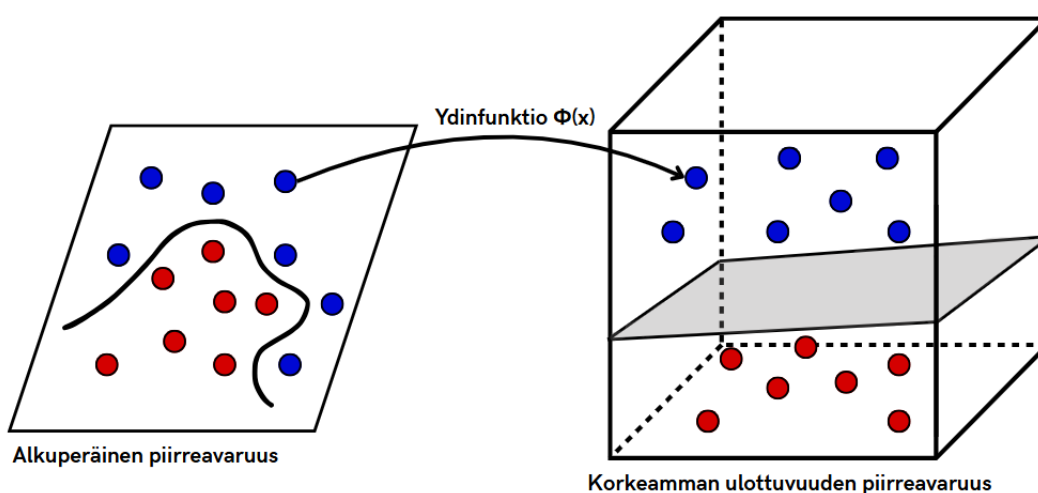
Tukivektorikone (engl. support-vector machine, SVM) on ohjatun oppimisen koneoppimismalli, jota käytetään luokittelu- ja regressioitehtäviin (Pisner & Schnyer, 2020). Tukivektorikoneen avulla voidaan löytää optimaalinen taso eli hyperpinta, joka erottaa eri luokkien datapisteet. Optimaalinen hyperpinta löydetään marginaalitason avulla. Marginaali, eli etäisyys lähimpiin datapisteisiin kustakin luokasta pyritään maksimoimaan. Hyperpintaa lähimmät datapisteet ovat nimeltään tukivektoreita. Ne määrittävät hyperpinnan asennon ja suunnan. Tukivektorikoneetta voidaan käyttää esimerkiksi kasvojentunnistukseen sekä kuvien luokitteluun (Chandra & Bedi, 2021; Pisner & Schnyer, 2020).

SVM voi olla lineaarinen tai epälineaarinen. Useimmiten luokitteluongelmat ovat epälineaarisia (Pisner & Schnyer, 2020). Hyperpinnan muoto lineaarisessa tapauksessa on suora ja oletuksena on, että luokkien piirteet ovat lineaarisesti erotettavissa (Kuva 6). Datassa olevat poikkeamat vaikeuttavat lineaarisen hypertason muodostamista. Pehmeä marginaali tukivektorikoneessa sallii virheitä käyttämällä lieventäviä muuttujia (engl. slack variables). Pehmeän marginaalin avulla yksittäisten poikkeamien vaikutusta hyperpinnan asentoon saadaan vähennettyä.



Kuva 6 Tukivektorikone mallinnettu lähteestä (Pisner & Schnyer, 2020)

Monissa tapauksissa luokkien datapisteet eivät ole erotettavissa yksinkertaisen suoran tai tason avulla, mikä tekee niistä epälineaarisia. Tällöin voidaan hyödyntää niin sanottua ytimen tempua (engl. kernel trick), joka muuntaa alkuperäisen datan korkeampaan ulottuvuuteen, jossa se voidaan erottaa lineaarisesti (Kuva 7) (Valkenborg ym., 2023). Datan muuntamiseen korkeampaan ulottuvuuteen käytetään ydinfunktioita, jotka laskevat yhtäläisyyksiä dataparien välillä. Tämä mahdollistaa monimutkaisen luokittelun epälineaarisesta datasta ilman laskennallisesti vaativampaa korkeampiulotteista laskentaa.



Kuva 7 Epälineaarisesti eroteltavien datapisteiden luokittelu ydintempun avulla

K-means-klusterointi on ohjaamattoman oppimisen menetelmä, joka ryhmittelee sille annetun datan alijoukkoihin eli klustereihin (Sinaga & Yang, 2020). Ryhmittely tapahtuu yksittäisten datapisteiden samankaltaisuuden perusteella. K-means-algoritmille annetaan haluttu klustereiden määrä, minkä jälkeen algoritmi etsii saman määrän keskuksia, joiden ympärille muodostaa klusterit. Klustereiden keskuksat päivittyvät vaiheittain siten että datapisteiden etäisyys omaan klusterikeskukseen on mahdollisimman pieni.

K-means-klusterointimenetelmän vahvuuksia ovat sen yksinkertaisuus ja nopeus (Xu & Lange, 2019). Se soveltuu hyvin erimuotoisen datan, kuten kuvien, sijaintitietojen tai tekstin piirteiden klusterointiin. K-means-algoritmin avulla voidaan havaita esimerkiksi kuvioita tai poikkeavuuksia datassa. Sinagan & Yangin (2020) mukaan K-means-algoritmin haasteina ovat klustereiden määrän valinta ja keskuksien alkusijainti, jotka vaikuttavat lopputulokseen. He esittävät tähän ongelmaan ratkaisuksi U-k-means algoritmia, joka löytää optimaalisen klusterimäärän automaattisesti.

Klusterointimenetelmiä, kuten K-means-klusterointia voidaan hyödyntää esimerkiksi sosiaalisen median käyttäjien ryhmittelymiseen (Yadav ym., 2023). Käyttäjistä voidaan kerätä tietoja, kuten seuraajamäärä, julkaisujen määrä, käytetyt sanat ja maantieteellinen sijainti. Näiden tietojen avulla voidaan ryhmitellä käyttäjät esimerkiksi seuraaviin klustereihin: tavalliset käyttäjät, botit, poliittisesti aktiiviset käyttäjät ja vaikuttajat tai organisaatiot.

4 Koneoppimisen käyttö avointen lähteiden tiedustelussa

Avointen lähteiden tiedustelun avulla voidaan kerätä suuria määriä ajantasaista ja monimuotoista tiedustelutietoa eri lähteistä. Internetin ja sosiaalisen median kautta saadaan dataa, mikä on ainutlaatuista ja tärkeää kattavan tilannekuvan muodostamiseksi. Samaan aikaan saatavilla olevan tiedon määrä on kasvanut valtavasti, joka muodostaa haasteen tiedon keräämiselle ja analysoimiselle. Manuaalisesti dataa ei siis voida analysoida tehokkaasti.

Koneoppiminen tarjoaa ratkaisuja tiedon määrään liittyviin ongelmiin. Sen vahvuus on suurien tietomäärien käsittely lyhyessä ajassa. Tehokkaiden koneoppimismallien kouluttamiseen tarvitaan paljon koulutusaineistoa, jota OSINT tarjoaa. Avoimista lähteistä saatua dataa voidaan käyttää koneoppimismallien koulutukseen ja koneoppiminen mahdollistaa tämän datan muuttamisen käyttökelpoiseksi tiedustelutiedoksi. Tässä luvussa käydään tarkemmin läpi, miten koneoppimista voidaan hyödyntää avointen lähteiden tiedustelussa.

4.1 Koneoppimisen käyttö avointen lähteiden tiedustelun eri vaiheissa

Koneoppimisen hyödyntämisellä voidaan tehostaa avointen lähteiden tiedustelua useassa eri vaiheessa (Browne ym., 2024). Tässä tutkielmassa käytetyistä tiedusteluprosessin vaiheista (Laatikko 1), sitä on hyödynnetty tähän mennessä eniten tiedon käsittelyyn ja analysointiin, mutta jonkin verran myös keräämiseen ja jakamiseen. Suunnittelu ja suuntaamis- sekä arviointivaiheissa koneoppimista ei vielä hyödynnetä laajasti.

Avointen lähteiden tiedustelutiedon keräämistä voidaan automatisoida tehokkaasti koneoppimisen avulla (Ghioni ym., 2024; Withorne, 2022). Siinä missä jokainen tietolähde jouduttiin ennen etsimään ja keräämään manuaalisesti, voi koneoppiminen nykyään hoitaa sen automaattisesti (Browne ym., 2024). Tämä nopeuttaa ja helpottaa tiedusteluprosessia, sillä tiedon keräämiseen ei enää kulu niin paljon henkilöresursseja. Koneoppimista käytetään pääosin tiedon keräämiseen verkkopohjaisista lähteistä, kuten uutisista, nettisivustoilta, ja pimeästä verkosta.

Sosiaalisen median alustat, kuten X (ent. Twitter), ovat nykyään yleinen tiedustelutiedon lähde (Vadapalli ym., 2019). Nämä alustat sisältävät valtavasti tietoa, jonka systemaattinen läpikäyminen manuaalisesti olisi erittäin hankalaa ja resursseja kuluttavaa. Eiverkkopohjaisista lähteistä tiedon keräämiseen koneoppiminen ei sovellu laajasti, sillä ne ovat usein fyysisiä.

Koneoppimista voidaan hyödyntää tiedon käsittelyvaiheessa, jossa kerätty data muunnetaan analysoitavaan muotoon (Layton & Watters, 2016). Koneoppimisalgoritmit voivat luokitella kerättyä tietoa kulloisenkin tarpeen pohjalta, jolloin tiedon analysointi manuaalisesti on huomattavasti nopeampaa. Koneoppimisalgoritmi voidaan esimerkiksi laittaa luokittelemaan avoimista lähteistä hankittuja kyberturvallisuuteen liittyviä raportteja teemoittain, jolloin tiettyyn uhkaan liittyvien raporttien löytäminen ja analysointi on nopeampaa (Yang & Lam, 2020). Koneoppimisen avulla voidaan myös suodattaa kerätystä datasta pois epäolennaiset ja merkityksettömät tiedot. Esimerkiksi satelliittikuvista voidaan koneoppimisen avulla valita manuaalisesti analysoitaviksi vain ne, joissa esiintyy tiedustelun kannalta merkittävää tietoa. Tämä vähentää manuaaliseen analysoimiseen käytettäviä henkilöresursseja.

Koneoppimisen menetelmien laaja käyttö juuri tiedon käsittelyn vaiheessa johtuu niiden algoritmien kyvystä luokitella dataa ja skaalautua suurille datamäärille. Esimerkiksi YOLOv5 (You Only Look Once) on konvoluutioverkkoon perustuva objektintunnistusmalli, joka pystyy tunnistamaan kuvissa tai videoissa esiintyviä esineitä ja asioita (Ballinger, 2023). Malli voidaan kouluttaa tiettyyn käyttötarkoitukseen, kuten sotilaallisten ajoneuvojen tai laivojen tunnistamiseen ja luokitteluun satelliittikuvista. Hyvin koulutettu YOLOv5 kykenee automaattisesti tunnistamaan ja luokittelemaan koulutusdataa vastaavia kohteita valtavasta määrästä kuvadataa. Sen avulla voidaan siis automatisoida merkityksellisten kohteiden havaitseminen ja datan luokittelu, mikä tekee siitä hyvän työkalun tiedon käsittelyyn.

Analyysivaiheessa koneoppimista voidaan käyttää datasta tehtävien tulkintojen helpottamiseksi. Kun data on analysoitavassa muodossa, voidaan tehdä siitä alustavia tulkintoja (Withorne, 2022). Tähän voidaan käyttää neuroverkkoja, tukivektorikonetta ja K-means-klusterointia. Esimerkiksi koneoppimisen avulla voitiin seurata Pohjois-Korean uraanikaivosten kasvua vuosien varrella (Park ym., 2021). Satelliittikuvista tunnistettiin automaattisesti keskeiset maastonpiirteet, kuten rakennukset, tiet, metsät, ruohikot, ja kaivoksen laajentumiseen viittaavat kohdat luokiteltiin ”muuksi”. Algoritmi pystyi havaitsemaan maastonmuutoksia, muttei sitä johtuivatko muutokset kaivostoiminnasta vai

luonnollisista ilmiöistä (Withorne, 2022). Koneoppimisen avulla tehdyt havainnot annettiin geologien analysoitaviksi. Koneoppimisen avulla voidaan siis tehdä alustava analyysi, mutta lisäksi tarvitaan asiantuntijoita ja tutkijoita tarkistamaan tulokset ja tekemään analyysi loppuun.

Tiedustelutiedon jakamisessa koneoppimista hyödynnetään toistaiseksi varsin vähän. Koneoppimista kuitenkin voidaan käyttää tiedon jakamiseen, esimerkiksi automaattisten varoitusten luonnin kautta (Ekwunife, 2020). Koneoppimista apuna käyttäen voidaan luoda järjestelmiä, jotka luovat hälytyksiä ja jakavat ne asianomaiselle yleisölle, kuten viranomaisistahoille, toimivallan ja luokituksen perusteella. Tiedustelutiedon jakamisvaiheen automatisointi koneoppimisen tai tekoälyn avulla voisi mahdollistaa reaaliaikaisemman tiedon jakamisen ja täysin automatisoidun tiedusteluprosessin luomisen (Browne ym., 2024).

Tekoälyn ja koneoppimisen soveltaminen avointen lähteiden tiedustelussa on yleistynyt vasta vuoden 2016 jälkeen (Evangelista ym., 2021). Kiinnostus aiheeseen on kasvanut tasaisesti ja tulevaisuudessa koneoppiminen tulee olemaan suuremmassa roolissa avointen lähteiden tiedustelussa. Mahdollisia kehityssuuntia ja tapoja, joilla koneoppimista voitaisiin tulevaisuudessa hyödyntää on useita (Browne ym., 2024). Tulevaisuuden konkreettisia käyttökohteita ovat muun muassa disinformaation tunnistus ja lähteiden luotettavuuden arviointi. Esimerkiksi konvoluutioverkkojen avulla voidaan tunnistaa muokattuja kuvia ja videoita. Muokatuista kuvista voidaan havaita pikselitasoisia epäjohdonmukaisuuksia, kuten valaistuksen tai tekstuurin välisiä eroja (Westerlund, 2019).

Tällä hetkellä koneoppimista hyödynnetään eniten tiedon keräämiseen, käsittelyyn ja analysoimiseen. Tulevaisuudessa koneoppimista voitaisiin hyödyntää enemmän myös suunnittelu- ja suuntausvaiheissa, esimerkiksi tiedonlähteiden valinnassa. Tietolähteitä voidaan luokitella sen mukaan, millaista tietoa niistä on aikaisemmin saatu (Iashvili, 2022). Koneoppimismallit voisivat myös suositella potentiaalisia tietolähteitä annettujen vaatimusten perusteella.

4.2 Ukrainan ja Venäjän sota sovellusesimerkinä

Ukrainan ja Venäjän välinen konflikti alkoi vuonna 2014 ja muuttui Venäjän täysimittaiseksi hyökkäykseksi Ukrainaan helmikuussa 2022. Konfliktiä on voitu seurata ensimmäisistä päivistä lähtien sosiaalisen median kautta ja päivityksiä on saatu lähes reaaliajassa (Winter ym., 2023). Sosiaaliseen mediaan julkaistut kuvat sotilaskalustosta ja niiden liikkeistä ovat

tarjonneet paljon dataa, jota on voitu hyödyntää avointen lähteiden tiedustelussa (Hockenhuil, 2022). Ukrainan asevoimat ovat omien sanojensa mukaan seuranneet esimerkiksi Tšetšeenijoukkojen liikkeitä sosiaalisen median alustojen, kuten Instagramin ja TikTokin avulla (Winter ym., 2023).

Kaupallisilla satelliittikuvilla on ollut merkittävä rooli Ukrainan asevoimille tilannekuvan muodostamisessa. Tekoälyä on käytetty monin tavoin nopeuttamaan tiedusteluprosessia, sillä nopeus antaa etulyöntiaseman vastapuoleen nähden (Hockenhuil, 2022). Satelliittikuvien analysointi on ollut yleinen tekoälyn ja koneoppimisen käyttökohde Ukrainan sodassa (Kamminga & Fontes, 2023). Satelliittikuvien analysointiin on käytetty esimerkiksi keinotekoisii neuroverkkoihin perustuvia menetelmiä. Niitä voidaan hyödyntää yhdistämään kuvamateriaalia monista eri lähteistä, jolloin saadaan tuotettua kattavampaa tiedustelutietoa.

Palantir Technologies on yhdysvaltalainen teknologiayhtiö, joka on kehittänyt useita kehittyneitä data-analytiikan työkaluja eri toimialojen käyttöön. Yksi tunnetuimmista näistä on Palantir Gotham, joka on suunnattu erityisesti valtiollisille toimijoille, kuten tiedustelupalveluille ja kansallisen turvallisuuden sektorille (Kosoy, 2025). Kesällä 2022 Palantir Technologies solmi yhteistyösopimuksen Ukrainan hallituksen kanssa, jossa Palantir tarjosi teknologiaratkaisujaan ilmaiseksi (Bergengruen, 2024). Palantirin kehittämät ohjelmistot pystyvät analysoimaan suuria tietomääriä ja tunnistamaan keskeisiä uhkia. Näistä järjestelmistä on tullut tärkeä osa Ukrainan puolustustoimia, ja niitä on hyödynnetty aktiivisesti muun muassa vihollisjoukkojen tunnistamisessa ja iskujen kohdistamisessa. Palantir Technologiesin toimitusjohtaja Alex Karp on muun muassa todennut ohjelmistojen olevan vastuussa suurimmasta osasta iskujen kohdistamisesta Ukrainassa, viitaten kohteiden paikantamiseen.

Palantirin ohjelmistot hyödyntävät tekoälyä ja koneoppimista yhdistääkseen ja analysoidakseen eri lähteistä saatavaa dataa. Näihin lähteisiin kuuluvat muun muassa satelliittikuvat, dronien tuottama kuvamateriaali, sosiaalisen median sisältö sekä kentältä saadut raportit. Eri lähteistä saadut tiedot yhdistetään kokonaisvaltaiseksi tilanneraportiksi (Mazarchuk, 2024).

Toinen tekoälyä ja koneoppimista hyödyntävä teknologia, jota Ukraina käyttää on Clearview AI, joka on kasvojentunnistusjärjestelmä. Toiminta perustuu avoimista lähteistä saatavaan tietoon, kuten sosiaalisesta mediasta kerättyihin kasvotietoihin, jotka on kerätty valtavaan tietokantaan (Bergengruen, 2023). Clearview AI muuntaa kerätyt kuvat kasvoista numeerisiksi vektoreiksi koneoppimisen ja syväoppimisen algoritmien avulla.

Ukraina otti käyttöön Clearview AI:n kasvojentunnistusjärjestelmän maaliskuussa 2022. (Dave & Dastin, 2022). Clearview tarjosi Ukrainalle ilmaisen pääsyn yli 10 miljardin julkisen valokuvan tietokantaansa, joista yli 2 miljardia oli peräisin Venäjän sosiaalisen median palvelusta VKontaktista. Tämä mahdollisti avoimista lähteistä saadun kuvamateriaalin käyttämisen venäläisten sotilaiden tunnistamiseen heidän kasvojensa perusteella (Bergengruen, 2023). Kasvojentunnistusjärjestelmän avulla on esimerkiksi pystytty tunnistamaan sotarikoksiin syyllistyneitä henkilöitä. Näitä henkilöitä on voitu tunnistaa sosiaalisen median julkaisujen avulla.

5 Yhteenveto

Tässä tutkielmassa tarkasteltiin koneoppimisen roolia avointen lähteiden tiedustelussa. Koneoppimisella on jo tällä hetkellä merkittävä rooli avointen lähteiden tiedustelussa, mutta sen merkityksen voidaan olettaa kasvavan tulevaisuudessa. Koneoppimisen menetelmien käyttö tehostaa avointen lähteiden tiedustelun prosesseja automatisoimalla suurten tietomassojen käsittelyn.

Yleisimmät nykyään käytetyt koneoppimisen menetelmät avointen lähteiden tiedusteluun ovat erilaiset neuroverkkoihin perustuvat menetelmät, kuten konvoluutioverkot tai takaisinkytketyt verkot. Lisäksi datan luokitteluun ja ryhmittelyyn käytetään usein tukivektorikonetta tai K-means-klusterointia.

Lähitulevaisuudessa näköpiirissä ei ole täysin uusia koneoppimisen menetelmiä, joita voitaisiin hyödyntää avointen lähteiden tiedustelussa. Sen sijaan jo käytössä olevia koneoppimisen menetelmiä voidaan käyttää uusiin käyttötarkoituksiin. Jo olemassa olevia menetelmiä voidaan myös kehittää soveltumaan paremmin uusiin tarkoituksiin.

Koneoppimisella on monia eri käyttökohteita avointen lähteiden tiedustelussa. Tällä hetkellä koneoppimista hyödynnetään useimmin tilanteissa, joissa joko datan kerääminen, käsittely tai analysointi manuaalisesti on haastavaa tai täysin mahdotonta. Koneoppiminen tarjoaa työkalut avointen lähteiden tarjoamien valtavien tietomäärien käsittelyyn.

Koneoppimisen hyödyntäminen avointen lähteiden tiedustelussa on edelleen verrattain uusi tutkimusalue. Sekä koneoppiminen että avointen lähteiden tiedustelu ovat kehittyneet nopeasti viime vuosien aikana. Avointen lähteiden tiedusteluun on syntynyt uusia tiedonlähteitä ja koneoppimisessa uusia tehokkaita menetelmiä. Tässä tutkielmassa tarkastellut menetelmät edustavat yleisimpiä ja vakiintuneimpia ratkaisuja.

Lähteet

- Allison, G. (2024). *Classified fighter jet specs leaked on War Thunder – again*.
<https://ukdefencejournal.org.uk/classified-fighter-jet-specs-leaked-on-war-thunder-again/>,
 viitattu 14.3.2025
- Al-Selwi, S. M., Hassan, M. F., Abdulkadir, S. J., Muneer, A., Sumiea, E. H., Alqushaibi, A., &
 Ragab, M. G. (2024). RNN-LSTM: From applications to modeling techniques and beyond—
 Systematic review. *Journal of King Saud University - Computer and Information Sciences*,
 36(5), 102068. <https://doi.org/10.1016/j.jksuci.2024.102068>
- Arranz, R., Carramiñana, D., Miguel, G. de, Besada, J. A., & Bernardos, A. M. (2023). Application of
 Deep Reinforcement Learning to UAV Swarming for Ground Surveillance. *Sensors*, 23(21),
 Article 21. <https://doi.org/10.3390/s23218766>
- Ballinger, O. (2023). *Remote Sensing for OSINT - Object Detection*.
https://bellingcat.github.io/RS4OSINT/C5_Object_Detection.html, viitattu 8.5.2025
- Bellingcat Investigation Team. (2018). *Full report: Skripal Poisoning Suspect Dr. Alexander Mishkin, Hero of Russia*. Bellingcat. <https://www.bellingcat.com/news/uk-and-europe/2018/10/09/full-report-skripal-poisoning-suspect-dr-alexander-mishkin-hero-russia/>, viitattu 12.3.2025
- Bellingcat Investigation Team. (2020). *Hunting the Hunters: How We Identified Navalny's FSB Stalkers*. Bellingcat. <https://www.bellingcat.com/resources/2020/12/14/navalny-fsb-methodology/>, viitattu 12.3.2025
- Bergengruen, V. (2023). *Ukraine's "Secret Weapon" Is a Controversial Tech Company*. TIME.
<https://time.com/6334176/ukraine-clearview-ai-russia/>, viitattu 13.5.2025
- Bergengruen, V. (2024). *How Tech Giants Turned Ukraine Into an AI War Lab*. TIME.
<https://time.com/6691662/ai-ukraine-war-palantir/>, viitattu 13.5.2025
- Bre, F., Gimenez, J., & Fachinotti, V. (2017). Prediction of wind pressure coefficients on building
 surfaces using Artificial Neural Networks. *Energy and Buildings*, 158.
<https://doi.org/10.1016/j.enbuild.2017.11.045>

- Browne, T. O., Abedin, M., & Chowdhury, M. J. M. (2024). A systematic review on research utilising artificial intelligence for open source intelligence (OSINT) applications. *International Journal of Information Security*, 23(4), 2911–2938. <https://doi.org/10.1007/s10207-024-00868-2>
- Böhm, I., & Lolagar, S. (2021). Open source intelligence. *International Cybersecurity Law Review*, 2(2), 317–337. <https://doi.org/10.1365/s43439-021-00042-7>
- Chandra, M. A., & Bedi, S. S. (2021). Survey on SVM and their application in imageclassification. *International Journal of Information Technology*, 13(5), 1–11. <https://doi.org/10.1007/s41870-017-0080-1>
- Colquhoun, C. (2016). *A Brief History of Open Source Intelligence*. Bellingcat. <https://www.bellingcat.com/resources/articles/2016/07/14/a-brief-history-of-open-source-intelligence/>, viitattu 14.3.2025
- Dave, P., & Dastin, J. (2022). Exclusive: Ukraine has started using Clearview AI’s facial recognition during war. *Reuters*. <https://www.reuters.com/technology/exclusive-ukraine-has-started-using-clearview-ais-facial-recognition-during-war-2022-03-13/>, viitattu 13.5.2025
- Ekwunife, N. (2020). National Security Intelligence through Social Network Data Mining. *2020 IEEE International Conference on Big Data (Big Data)*, 2270–2273. <https://doi.org/10.1109/BigData50022.2020.9377940>
- Emmert-Streib, F., Yang, Z., Feng, H., Tripathi, S., & Dehmer, M. (2020). An Introductory Review of Deep Learning for Prediction Models With Big Data. *Frontiers in Artificial Intelligence*, 3. <https://doi.org/10.3389/frai.2020.00004>
- Evangelista, J. R. G., Sassi J. R., Romero M., & Napolitano D. (2021). Systematic Literature Review to Investigate the Application of Open Source Intelligence (OSINT) with Artificial Intelligence. *Journal of Applied Security Research*, 16(3), 345–369. <https://doi.org/10.1080/19361610.2020.1761737>
- Gerard, C. (2021). The Basics of Machine Learning. Teoksessa *Practical Machine Learning in JavaScript* (ss. 1–1). Apress, an imprint of Springer Nature. https://doi.org/10.1007/978-1-4842-6418-8_1

- Ghioni, R., Taddeo, M., & Floridi, L. (2024). Open source intelligence and AI: A systematic review of the GELSI literature. *AI & SOCIETY*, 39(4), 1827–1842. <https://doi.org/10.1007/s00146-023-01628-x>
- Gibson, H. (2016). Acquisition and Preparation of Data for OSINT Investigations. Teoksessa B. Akhgar, P. S. Bayerl, & F. Sampson (Toim.), *Open Source Intelligence Investigation: From Strategy to Implementation* (s. 69–93). Springer International Publishing. https://doi.org/10.1007/978-3-319-47671-1_6
- Guo, Y., Xia, Y., Wang, J., Yu, H., & Chen, R.-C. (2020). Real-Time Facial Affective Computing on Mobile Devices. *Sensors*, 20(3), 870. <https://doi.org/10.3390/s20030870>
- Hassan, N. A., & Hijazi, R. (2018). *Open Source Intelligence Methods and Tools*. Apress. <https://doi.org/10.1007/978-1-4842-3213-2>
- Hijazi, S., Kumar, R., & Rowen, C. (2015). *Using Convolutional Neural Networks for Image Recognition*. 2015.
- Hiran, K. K., Jain, R. K., Lakhwani, D. K., & Doshi, D. R. (2021). *Machine Learning: Master Supervised and Unsupervised Learning Algorithms with Real Examples (English Edition)*. BPB Publications.
- Hockenhuil. (2022). *How open-source intelligence has shaped the Russia-Ukraine war*. GOV.UK. <https://www.gov.uk/government/speeches/how-open-source-intelligence-has-shaped-the-russia-ukraine-war>, viitattu 13.5.2025
- Hulnick, A. S. (2002). The Downside of Open Source Intelligence. *International Journal of Intelligence and CounterIntelligence*, 15(4), 565–579. <https://doi.org/10.1080/08850600290101767>
- Hwang, Y.-W., Lee, I.-Y., Kim, H., Lee, H., & Kim, D. (2022). Current Status and Security Trend of OSINT. *Wireless Communications and Mobile Computing*, 2022(1), 1290129. <https://doi.org/10.1155/2022/1290129>
- Iashvili, G. (2022). *Machine Learning Use In OSINT*. https://deepsec.net/docs/Slides/2022/Machine_Learning_Use_In_OSINT_Giorgi_Iashvili.pdf, viitattu 12.5.2025

- Ivanjko, T., & Dokman, T. (2019). *Open Source Intelligence (OSINT): Issues and trends*. 191–196.
<https://doi.org/10.17234/INFUTURE.2019.23>
- Kamminga, J., & Fontes, R. (2023). *Ukraine A Living Lab for AI Warfare*.
<https://www.nationaldefensemagazine.org/articles/2023/3/24/ukraine-a-living-lab-for-ai-warfare>, viitattu 13.5.2025
- Kan, A. (2017). Machine learning applications in cell image analysis. *Immunology and Cell Biology*, 95(6), 525–530. <https://doi.org/10.1038/icb.2017.16>
- Kosoy, D. (2025). *Palantir, the Secretive Tech Giant Shaping Ukraine's War Effort*. UNITED24 Media. <https://united24media.com/war-in-ukraine/palantir-the-secretive-tech-giant-shaping-ukraines-war-effort-5519>, viitattu 13.5.2025
- Layton, R., & Watters, P. A. (2016). *Automating open source intelligence: Algorithms for OSINT*. Elsevier/Syngress.
- Mazarchuk, A. (2024). *VIEWPOINT: AI for War and Peacetime: A Ukrainian Perspective*.
<https://www.nationaldefensemagazine.org/articles/2024/11/1/viewpoint-ai-for-war-and-peacetime-a-ukrainian-perspective>, viitattu 13.5.2025
- Office of the Director of National Intelligence. (2011). *U.S. National Intelligence—An Overview 2011*. Office of the Director of National Intelligence.
https://www.dni.gov/files/documents/IC_Consumers_Guide_2011.pdf, viitattu 3.3.2025
- Omand, D. (2017). Social Media Intelligence (SOCMINT). Teoksessa R. Dover, H. Dylan, & M. S. Goodman (Toim.), *The Palgrave Handbook of Security, Risk and Intelligence* (ss. 355–371). Palgrave Macmillan UK. https://doi.org/10.1057/978-1-137-53675-4_20
- Park S., McNulty T., Puccioni A. & Ewing R. C. (2021). Assessing Uranium Ore Processing Activities Using Satellite Imagery at Pyongsan in the Democratic People's Republic of Korea. *Science & Global Security*, 29(3), 111–144. <https://doi.org/10.1080/08929882.2021.1988258>
- Pisner, D. A., & Schnyer, D. M. (2020). Chapter 6—Support vector machine. Teoksessa A. Mechelli & S. Vieira (Toim.), *Machine Learning* (ss. 101–121). Academic Press.
<https://doi.org/10.1016/B978-0-12-815739-8.00006-7>
- Rislakki, J. (2024). *Tiedustelu ja vakoilu: Opit, operaatiot, agentit*. Docendo.

- Sen, J., & Mehtab, S. (2022). Long-and-Short-Term Memory (LSTM) Networks Architectures and Applications in Stock Price Prediction. Teoksessa *Emerging Computing Paradigms* (ss. 143–160). John Wiley & Sons, Ltd. <https://doi.org/10.1002/9781119813439.ch8>
- Sinaga, K. P., & Yang, M.-S. (2020). Unsupervised K-Means Clustering Algorithm. *IEEE Access*, 8, 80716–80727. <https://doi.org/10.1109/ACCESS.2020.2988796>
- Taylor, P. (2024). *Volume of data/information created, captured, copied, and consumed worldwide from 2010 to 2023, with forecasts from 2024 to 2028 (in zettabytes)*. Statista. <https://www.statista.com/statistics/871513/worldwide-data-created/>, viitattu 18.2.2025
- U.S. Director of National Intelligence. (2006). *Intelligence Community Directive Number 301: National Open Source Enterprise*. Office of the Director of National Intelligence. <https://irp.fas.org/dni/icd/icd-301.pdf>, viitattu 19.2.2025
- Vadapalli, S. R., Hsich, G., & Nauer K. (2019). *Twitter OSINT: Automated Cybersecurity Threat - ProQuest*. <https://www.proquest.com/docview/2153621548?fromopenview=true&pq-origsite=gscholar&sourcetype=Conference%20Papers%20&%20Proceedings>
- Valkenborg, D., Rousseau, A.-J., Geubbelmans, M., & Burzykowski, T. (2023). Support vector machines. *American Journal of Orthodontics and Dentofacial Orthopedics*, 164(5), 754–757. <https://doi.org/10.1016/j.ajodo.2023.08.003>
- Westerlund, M. (2019). The Emergence of Deepfake Technology: A Review. *Technology Innovation Management Review*, 9(11), 40–53. <https://doi.org/10.22215/timreview/1282>
- Winter, D. C., Gallacher, D. J., & Harris, A. (2023). *Artificial Intelligence, OSINT and Russia's Information Landscape*. <https://cetas.turing.ac.uk/publications/artificial-intelligence-osint-and-russias-information-landscape>
- Withorne, J. (2022). Assessing the Relationship between Machine Learning and Open Source Research in International Security. Teoksessa *Open Source Investigations in the Age of Google: Vsk. Volume 4* (ss. 302–317). WORLD SCIENTIFIC (EUROPE). https://doi.org/10.1142/9781800614079_0016
- Xu, J., & Lange, K. (2019). Power k-Means Clustering. *Proceedings of the 36th International Conference on Machine Learning*, 6921–6931. <https://proceedings.mlr.press/v97/xu19a.html>

- Xue, F., Ji, H., Zhang, W., & Cao, Y. (2019). Attention-based spatial–temporal hierarchical ConvLSTM network for action recognition in videos. *IET Computer Vision*, *13*(8), 708–718. <https://doi.org/10.1049/iet-cvi.2018.5830>
- Yadav, A., Kumar, A., & Singh, V. (2023). Open-source intelligence: A comprehensive review of the current state, applications and future perspectives in cyber security. *Artificial Intelligence Review*, *56*(11), 12407–12438. <https://doi.org/10.1007/s10462-023-10454-y>
- Yang, W., & Lam, K.-Y. (2020). Automated Cyber Threat Intelligence Reports Classification for Early Warning of Cyber Attacks in Next Generation SOC. Teoksessa J. Zhou, X. Luo, Q. Shen, & Z. Xu (Toim.), *Information and Communications Security* (ss. 145–164). Springer International Publishing. https://doi.org/10.1007/978-3-030-41579-2_9