

Jyväskylän yliopisto

# Tekoälyn hyödyntäminen ja potentiaali biopankkisovelluksissa

AI hub Keski-Suomi hanke

Harri Juutilainen, Liisa Petäinen, Timo Ojala, Sami Äyrämö  
ja Pekka Neittaanmäki  
11.10.2022

## Sisällysluettelo

1 Johdanto.....	2
2 Koneoppiminen ja tekoäly .....	3
2.1 Oppiminen koneiden näkökulmasta .....	3
2.2. Koneoppiminen .....	4
2.3 Tekoäly .....	4
3 Patologia ja biopankit .....	7
3.1 Odotukset tekoälyn käytöstä patologiassa .....	7
3.2 Tekoäly patologian tukena .....	7
3.3 Käyttöönotto .....	8
3.4 Biopankkitoiminta käytännössä .....	8
4 Esimerkkejä tekoälyn potentiaalista biopankkitoiminnassa .....	9
4.1 Hyperplasian tunnistus .....	9
4.2 Sydän- ja verisuonisairausriskin ennustaminen biopankkidatasta .....	9
4.3 Kognitiivinen konenäkökehys syövän arviointiin .....	9
4.4 Tekoälypohjainen rintasyövän etäpesäkkeiden havainnointi .....	10
4.5 Radiologian ja patologian tekoälysovellukset – eturauhassyöpä .....	10
4.6 Tekoäly mammografian ja digitaalisen kerroskuvauksen apuna .....	10
4.7 Ihmisten geneettisten piirteiden ennustaminen .....	11
5 Pohdinta .....	12
5.1 Tärkeimmät kehityskohteet .....	12
5.2 Datan saatavuus .....	13
5.3 Riskitekijät .....	14
Lähteitä .....	16

# 1 Johdanto

*”Biopankki on lääketieteellinen kokoelma, joka sisältää ihmisperäisiä näytteitä sekä niihin liittyviä tietoja.”*

<https://fi.wikipedia.org/wiki/Biopankki>

Tämän raportin tavoitteena on selvittää, minkälaiset tekoälyn sovellutukset ovat suomalaisille biopankeille hyödyllisiä. Tarkasteluissa yritetään ottaa huomioon myös biopankkien resurssit ja se mihin koneoppimiseen perustuviin tekoälytehtäviin on nykytilanteessa olemassa riittävästi dataa. Samalla tuodaan esille olemassa olevissa menetelmissä käytettyjä vaatimuksia ja tavoitteita: minkälaista dataa käytetään sekä mitä tavoitteita menetelmille asetettiin.

Raportissa esitettävän sisällön taustalla on sen kirjoittajien käymiä keskusteluja eri sidosryhmien asiantuntijoiden kanssa, tutkimusartikkeleiden sisältöjä, konferenssiesityksiä jne., minkä pohjalta kirjoittajat ovat muodostaneet oman näkemyksensä tekoälyn nykytilasta ja potentiaalista biopankkitoiminnassa ja siihen liittyvistä vaatimuksista ja rajoitteista.

Yleisesti ottaen biopankkitoiminnalla nähdään valtava potentiaali niin kansalaisten hyvinvoinnin ja terveyden edistämisessä kuin kaupallisten innovaatioiden luomisessa:

*”Näytekokoelmien mahdollisimman kattava käyttö lääketieteelliseen tutkimukseen on monella tavalla hyvä asia. Kansalaisia ei tarvitse kutsua uudelleen osallistumaan tutkimuksiin, eikä heistä tarvitse kerätä uusia näytteitä. Myös suurin kustannuksin kootut näytekokoelmat saadaan tehokkaampaan tutkimuskäyttöön ja niiden avulla voidaan kehittää parempia hoitoja ja lääkkeitä.”*

- [www.biopankki.fi](http://www.biopankki.fi)

Vuodesta 2020 alkaen Suomen biopankkitoimintaa ohjaa ja valtakunnallista biopankkirekisteriä ylläpitää Fimea. Valtakunnallisesti rekisteröityjä biopankkeja ovat lokakuussa 2022:

- Auria Biopankki
- THL Biopankki
- Suomen Hematologinen Rekisteri ja Biopankki FHRB
- Helsingin Biopankki
- Pohjois-Suomen Biopankki Borealis
- Tampereen Biopankki
- Itä-Suomen Biopankki
- Keski-Suomen Biopankki
- Veripalvelun Biopankki
- Suomen Terveystalon Biopankki
- Arctic Biopankki – Oulun yliopisto

## 2 Koneoppiminen ja tekoäly

### 2.1 Oppiminen koneiden näkökulmasta

Oppiminen on prosessi, jossa kokemusta muutetaan tiedoksi. Tämä pätee myös koneelliseen oppimiseen. Koneille kokemus on niissä toimiville algoritmeille syötetty opetusdata.

Yksinkertaisen tehtävien opettaminen koneelle voi onnistua ohjelmoinnin keinoin. Ohjelmointi kuitenkin edellyttää, että tunnemme ja hahmotamme kaikki ohjelmoitavan ilmiön taustalla olevat tekijät ja niiden väliset riippuvuudet. Esimerkiksi yksinkertaisen pelin sääntömaailma voi olla riittävän yksinkertainen mallinnettavaksi ohjelmoitavan logiikan keinoin, mutta esimerkiksi ihmiskehon toiminta on niin moniulotteinen kokonaisuus valtavan monimutkaisine riippuvaisuuksineen, että sen määrittely ohjelmoitavaa logiikkaa varten on käytännössä mahdotonta.

Perinteisen ohjelmoinnin sijaan tekoälyn kehittäminen perustuu nykyään lähes poikkeuksetta koneoppimiseen, jossa geneerisillä algoritmeilla sovitetaan kiinnostuksen kohteena olevasta ilmiöstä kerättyyn dataan laskennallinen malli. Tällöin tavoitteena voi olla opettaa kone vastaamaan automaattisesti ilmiöön liittyvään merkitykselliseen kysymykseen (esim. mikä on syöpäsairaalan potilaan N.N. elinajanodote, jos häntä hoidetaan lääkkeellä Y) tai tuottaa uutta tietämystä datan taustalla olevasta ilmiöstä (esim. mitä ovat tulehduksellisten suolistosairauksien riskitekijät). Mallin dataan sovittamisen lisäksi keskeistä on analysoida mallin yleistymiskyky, ts. arvioida tilastollisesti kuinka tarkka ja luotettava malli on kun sitä sovelletaan uusiin yksilöihin jotka eivät sisällyneet koneen opettamisessa käytettyyn dataan.

Shalev-Shwartz & Ben-David, 2009, määrittelivät oppimisparadigmoja [Luku 1.3: Types of learning], joiden pohjalta (kone)oppimistehtäviä voidaan luokitella:

- Ohjattu ja ohjaamaton oppiminen

Esimerkkinä roskapostin lajittelu. Ohjatussa oppimisessa voidaan käyttää sähköposteja, jotka on merkitty kahteen luokkaan (asiallinen/roskaposti) ja näiden perusteella luodaan luokittelumalli lajittelua varten. Ohjaamaton oppimisessa luokitustietoa ei ole käytössä vaan pelkästään merkitsemättömiä sähköposteja, joista etsitään poikkeamia ja tunnistamiseen kykenevää mallia.

- Aktiivinen tai passiivinen oppiminen

Yllä oleva sähköpostiesimerkki toimii esimerkkinä passiivisesta oppimisesta. Eli tekoälysovellus ei toimi vuorovaikutuksessa käyttäjän kanssa. Aktiivinen oppiminen on vuorovaikutuksessa toimintaympäristönsä kanssa parantaakseen luokittelutarkkuutta, esimerkiksi pyytämällä käyttäjiä merkitsemään roskasähköposteja.

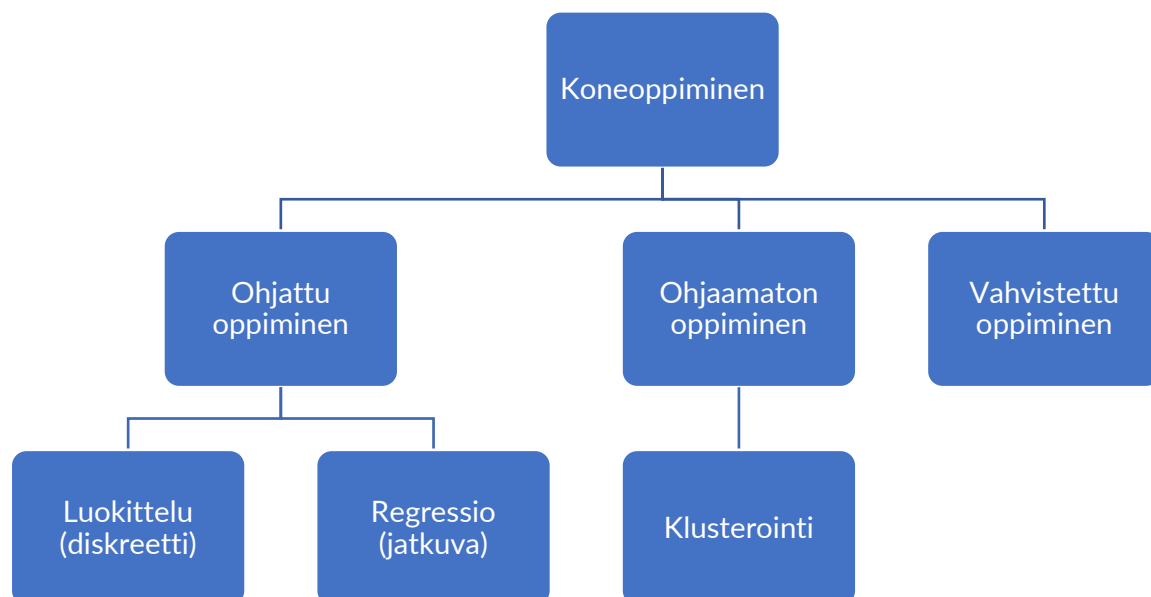
- Opetusdatan satunnaisuus

Jotta koneoppimismallista saadaan tarkempi, kannattaa opetusdatassa olla satunnaistettuja elementtejä. Mallista tulee vahvempi ja paremmin poikkeuksia kestävä, kun se opetetaan "vaikeammalla" ja enemmän vaihtelua sisältävällä opetusdatalla. Esimerkiksi digipatologisia kuvia voidaan kiertää tai muokata niiden väriavaruutta satunnaisesti ja lisätä muokatut kuvat opetusaineistoon.

## 2.2. Koneoppiminen

Koneoppiminen on tarpeellista datan kasvavan määrän vuoksi. Sitä voidaan Shalev-Shwartz & Ben-David, 2009, mukaan hyödyntää, jos jollekin kohteelle tunnetaan riittävän suuri joukko syöte- ja vastearvoja, mutta ei algoritmia joka annetulla syötteellä tuottaa tietyn vasteen. Muina hyödyntämismahdollisuuksina he pitävät ongelmia, jotka ovat liian monimutkaisia ohjelmoitaviksi, tai ongelmia, jotka vaativat jatkuvaa mukautumista ympäristön muutoksiin. Perinteisesti tieteelliset hypoteesit ovat perustuneet ihmisen havaintoihin ja kokemukseen, mutta koneoppimista voidaan käyttää myös uusien tieteellisten hypoteesien generoimiseen datasta.

Käytännön koneoppiminen on algoritmien sekä datapohjaisten mallien ja niiden arkkitehtuurien suunnittelua, kehitystä ja suorituskyvyn arviointia. Keskeisiä käyttökohteita ovat mm. piilevien muuttujien tunnistamisen datasta (esimerkiksi uusi syövä riskitekijä) ja päätöksenteon tukeminen luokittelu-/ennustusmallien avulla (esimerkiksi kasvaimen histopatologisten ominaisuuksien määrittäminen näytteestä). (Regitnig et al., 2020)



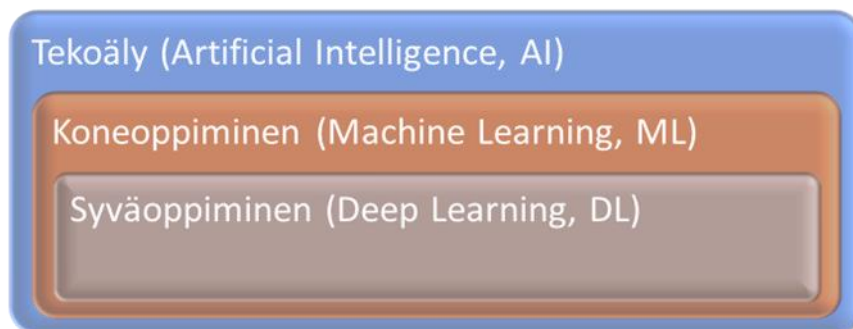
KUVA 1. KONEOPPIMISEN OSA-ALUEITA

## 2.3 Tekoäly

*“Älykkään käytöksen automatisointia, jonka tavoitteena on saavuttaa yleinen älykkyys.”*

Regitnig et al., 2020 määrittelevät tekoälyä yllä olevalla lainauksella. He tiivistävät syväoppimisen, koneoppimisen ja tekoälyn suhteet joukko-opin mukaan seuraavasti:

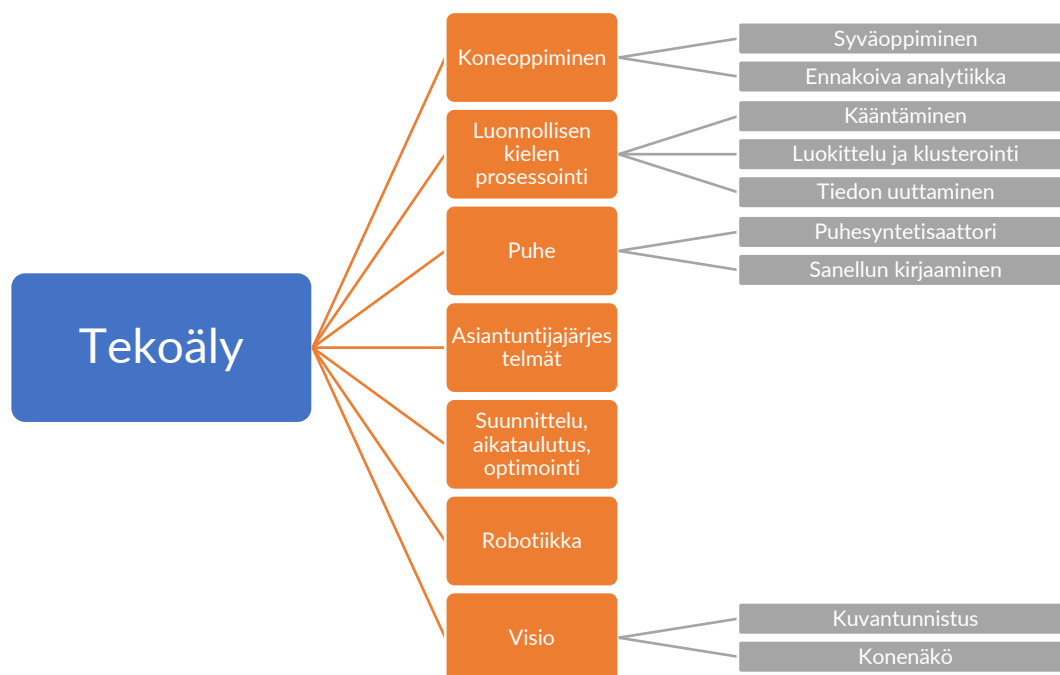
$$DL \subseteq ML \subseteq AI$$



KUVA 2. SYVÄOPPIMISEN, KONEOPPIMISEN JA TEKOÄLYN SUHDE

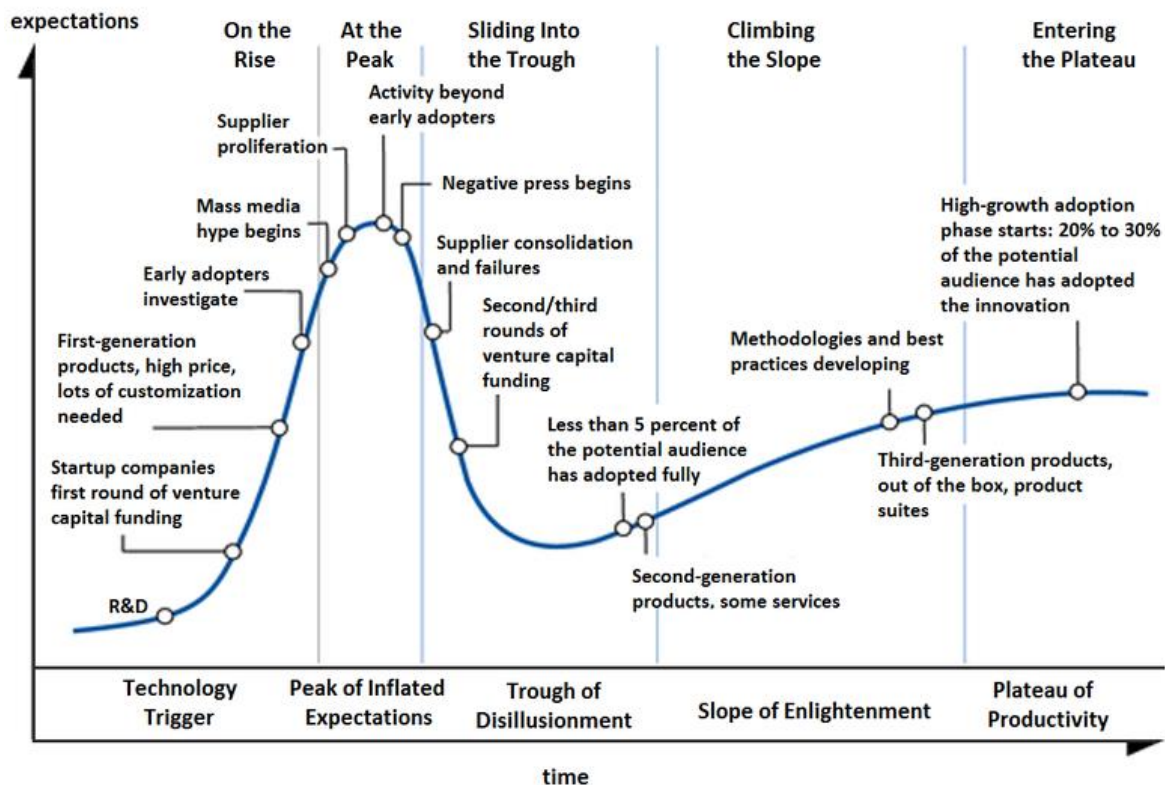
Koneoppiminen keskittyy enemmän algoritmisen ja tilastollisen oppimisen tekniseen puoleen, kun taas tekoäly soveltaa koneoppimista, mutta ottaa huomioon muun muassa eettiset ja filosofiset-, sekä päätöksenteon ongelmat. (Regitnig et al., 2020)

Modernit tekoälysovellukset perustuvatkin pääosin datasta oppimiseen. Ilman dataakin tietokoneet ovat tehokkaita ratkaisemaan monimutkaisia loogisia päättely- ja laskutehtäviä, mutta tosielämän ilmiöiden mallintamiseen riittävän informaation määrittely ja ohjelmoiminen on useimmiten vaikeaa. Ihmiselle helppojen mutta vaikeasti formaalissa muodossa esitettävien tehtävien ratkominen onkin tietokoneelle ja klassiselle tekoälylle haastavaa.



KUVA 3. RAKENTEELLINEN KUVAAUS TEKOÄLYSTÄ

Kuvan 4 hypesykli kuvastaa sitä, miten nousevat teknologiat pääsevät ns. "tuottavuuden tasangolle" (plateau of productivity). Teknologialla on tasangolle päästyään sekä paras taloudellinen tuottavuus että suurin mahdollinen merkitys. Biopankkien näkökulmasta tämä voitaisiin nähdä siten, että tasangolla tekoälystä saadaan irti suurin kliininen, tutkimuksellinen ja taloudellinen hyöty.



KUVA 4. NOUSEVIEN TEKNOLOGIOIDEN HYPESYKLI<sup>1</sup>

Digitaalisen terveyden ja älykkään terveydenhuollon tarkoituksena on tarkentaa diagnosointia ja tehostaa hoitoprosesseja, säästää kustannuksissa, sekä parantaa työ- ja asiakastyytyväisyyttä sairaaloissa -ja hoitolaitoksissa. Tekoäly on kyennyt jo auttamaan lääkäreitä diagnosoinnissa ja lääkekehityksen edistämässä. Sitä on hyödynnetty esimerkiksi unihäiriöiden hoidossa, onkologiassa, kardiologiassa sekä keuhkosairauksien tutkimuksessa ja hoidossa. Tekoälyn avulla diagnooseista voi tulla tarkempia ja nopeampia. Se myös tulee laajentamaan hoitohenkilökunnan toimintakykyä ja tarjoamaan mahdollisuuksia uudelle yhteistyölle eri toimijoiden välillä.

Terveydenhuollossa tekoälyä on kehitetty lisäksi stressin seurantaan, lääkärin dokumentaatiotyön helpottamiseen, aivojen kuvantamisen ja kuvien rekonstruoinnin kehittämiseen, terveystiedon louhintaan ja yksilön itsenäiseen terveyden seurantaan sekä hallintaan.

Jotta tekoäly kykenisi ihmismäiseen toimintaan, sen pitäisi pystyä oppimaan asioita esimerkiksi kuvista, puheesta, muista äänistä, teksteistä ja tapahtumaketjuista. Nykyiset tekoälyalgoritmit eivät pysty lisäämään tai luomaan oppimaansa malliin sellaista tietoa, jota niiden opettamiseen käytettävässä datassa ei ole. Tekoälyn laatu ja kattavuus on siis käytännössä lähes täysin sen opettamiseen saatavilla olevan datan määrästä ja laadusta riippuvaista.

<sup>1</sup> <https://commons.wikimedia.org/wiki/File:Hype-Cycle-General.png>; Olga Tarkovskiy, CC BY-SA 3.0 <<https://creativecommons.org/licenses/by-sa/3.0>>, via Wikimedia Commons

## 3 Patologia ja biopankit

Histopatologia määritellään englannin kielisessä wikipediassa seuraavasti: ”...the microscopic examination of tissue in order to study the manifestations of disease.”

Eli histopatologia tutkii mikroskooppinäkyvien avulla tautien ilmenemistä kudoksissa.

Digipatologialla taas tarkoitetaan kudoksista otettujen digitoitujen koepalojen hyödyntämistä ja datan hallinnointia.

Tutkittavia näytteitä histopatologiassa saadaan muun muassa koepaloista kirurgisen toimenpiteen tai ruumiinavauksen yhteydessä. Näytteitä tutkitaan mikroskoopilla ja niiden arviointi perustuu kuolleen tai muunnellun kudoksen ja terveän kudoksen vertailuun. (Slaoui & Fiette, 2010)

Histopatologiassa on ollut kaksi selvää murrosta: immunokemia ja myöhemmin geenilääkitys. Kolmas murros on Salto-Tellez et al., 2019 mukaan tekoäly. Tekoäly mahdollistaisi työläisten vaiheiden, kuten kudoksenäytteiden mikroskooppisen mittaamisen ja tulkinnan, automatisoinnin. Onkin mahdollista, että tekoäly korvaa osan patologin töistä. Tämä voi mahdollisesti parantaa työvoimapulaa sekä toisaalta vähentää patologien työkuormaa jolloin he voivat keskittyä aikaa vieviin monimutkaisempiin tapauksiin. (Salto-Tellez et al., 2019)

### 3.1 Odotukset tekoällyn käytöstä patologiassa

Patologiassa on vielä matkaa laajaan tekoällyn päivittäiseen käyttöön, mutta on nähtävissä että joidenkin tehtävien suorittamista on mahdollista edistää tekoällyn avulla. Tällaisia tehtäviä voivat olla esimerkiksi syövän havaitseminen ja arviointi, pienten kasvainten etsiminen imusolmukkeesta ja tekstin tulkinta kasvainkoodauksen virheiden vähentämiseksi.

Ehdottoman tärkeää patologian tekoälyjärjestelmissä on niiden luotettavuus, sillä väärä diagnoosi voi johtaa vakaviin hoitovirheisiin. Väärä positiivinen diagnoosi voi käytännössä tarkoittaa, että syöpäsoluja tappavaa kemoterapiaa annetaan henkilölle, jolla ei olekaan syöpää.

Yllä mainittujen tehtävien lisäksi parannuskohteeksi esitettiin kommunikaatiota ja yhteistyötä tekoälytutkijoiden ja patologioiden välillä. Kehittäjillä ei yleensä ole tarpeeksi lääketieteellistä tietoa, jotta kehityksen kohteena olevista järjestelmistä saataisiin lopulta kliinistä hyötyä (clinical value) ja toisaalta, patologeille ei ole selvää mitä mahdollisuuksia tekoäly käytännössä tuo. (Regitnig et al., 2020)

### 3.2 Tekoäly patologian tukena

Tekoällyn tuominen käytännön diagnostiikkaan edellyttää, että patologit luottavat tekoälymallien antamiin suosituksiin, diagnooseihin, ennusteisiin jne. Tämä vaatii patologioiden ja tekoälytutkijoiden



yhteistyötä. Tekoälytutkijat ovat kehittäneet tekoälymenetelmien evaluointia, jotta patologisista raporteista saadaan läpinäkyvämpiä. (Colling et al., 2019)

Tekoäly tehostaa lääkärin työtä, sillä esimerkiksi sairauksien diagnosointi ja hoidon suunnittelu käyvät helpommin, nopeammin ja luotettavammin. Parhaimmillaan ainakin tutkimusasetelmissa tekoälypohjaiseen kuva-analyysiin perustuvat sovellukset tarjoavat nopeita ja täsmällisyydessään asiantuntijoihin rinnastettavia vastauksia radiologian, patologian ja dermatologian alueilla. Nimenomaan koneen ja ihmisasiantuntijan yhteistyö voi tehostaa merkittävästi diagnosointiprosessia.

Onkologiaan liittyvissä tekoälytutkimuksissa on erityisesti painotettu syväoppimiseen ja kuvantunnistukseen perustuvien sovellusten kehittämistä. Usein nämä sovellukset voivat auttaa paitsi sairauden diagnosoinnissa, myös hoidon suunnittelussa. (TIEP1000)

### 3.3 Käyttöönotto

Histopatologian työtehtävät ovat säilyneet vuosikymmeniä isossa kuvassa samanlaisina. Digitaalisen patologian (perinteisen mikroskopian korvaaminen digitaalisten kuvien analyysillä) teknologioiden omaksuminen käyttöön on ollut hidasta. Uuden teknologian ja testien tuominen käytännön työhön vaikuttaa sekä työnkulkuun että henkilöstön koulutukseen. On myös mahdollista, että patologit alkavat luottaa liikaa tekoälyyn, joka voi johtaa heidän diagnostiikkakykynsä heikentymiseen. (Colling et al., 2019)

Tekoälyn potentiaali digitaalisessa patologiassa on laaja, mutta painopiste on nykyään kuva-analyysissa. Tekstianalytiikka on paljon kiinnostusta herättävä, mutta erityisesti pienten kielten kohdalla haastava sovellusalue.

Työkalujen arvioinnissa kiinnitetään huomiota tarkkuuden ja luotettavuuden lisäksi kliiniseen hyödyllisyyteen. Kliinisellä hyödyllisyydellä tarkoitetaan lääketieteellisen intervention asiaankuuluvuutta ja hyödyllisyyttä. Kliininen hyödyllisyys on usein sidosryhmien antama arvio, joka tehdään todistusaineiston näkökulmasta. Kliinisen hyödyllisyyden arviointi saattaa tästä syystä olla vaikeaa, koska se riippuu testeistä, lääkkeitä, sekä kontekstista. (Lesko et al., 2010)

### 3.4 Biopankkitoiminta käytännössä

Biopankkiin keräävät potilaista näytteitä ja tietoa erilaisia tulevaisuuden tutkimustarpeita varten. Näytteeseen liittyy myös muuta siihen liittyvää dataa.

biopankki.fi sivustolla määritellään: *”Biopankki eroaa käsitteenä perinteisistä tutkimusnäytekokoelmista siten, että biopankkiin ei kerätä näytteitä vain yhtä tiettyä tutkimusta varten, vaan myös erilaisiin tuleviin tutkimustarpeisiin.”*

Suomessa toimii yksitoista biopankkia sairaanhoitopiirien, yliopistojen sekä Terveystieteiden ja hyvinvoinnin laitoksen (THL), Veripalvelun ja Terveystalon yhteydessä. Biopankki itsessään ei ole itsenäinen rekisteri eikä sitä ole olemassa y-tunnuksena.

## 4 Esimerkkejä tekoälyn potentiaalista biopankkitoiminnassa

### 4.1 Hyperplasian tunnistus

Bukhari, 2020 sovelsivat koneoppimisen ja konenäön menetelmiä eturauhasen hyperplasian tunnistamiseen. Heidän tavoitteenansa oli tunnistaa sekä hyperplasia että alue, jossa se todennäköisesti esiintyy. Tuloksena saatiin malli, jolle todettiin 93 % tarkkuus riippumattomalla testidatalla.

Bukhari, 2020 käyttämä data koostui 59 kuvasta, joissa näkyi 169 aluetta, jossa esiintyi hyperplasiaa. Kuvat olivat erikokoisia, joten ne piti ensin normalisoida samankokoisiksi (datan esikäsittely). Lisäksi kuviin lisättiin variaatiota käyttämällä 33 % todennäköisyyttä jolla kuvat voivat peilaantua joko ylös, alas, vasemmalle tai oikealle.

70 % datasta käytettiin mallin opettamiseen, 10 % mallin valintaan ja 20 % yleistymiskyvyn testaukseen.

### 4.2 Sydän- ja verisuonisairausriskin ennustaminen biopankkidatasta

Dataan pohjautuvien menetelmien on mahdollista parantaa riskien ennustamista sydän- ja verisuonisairauksien osalta. Alaa et al., 2019, tekivät koneoppimiseen pohjautuvan mallin, jonka kehittämisessä hyödynnettiin AutoPrognosis -sovelluskehystä. AutoPrognosis käyttää päätöksentekoon kolmea eri koneoppimisen putkea, jotka kaikki hyödyntävät eri algoritmeja ja menetelmiä.

Alaa et al., 2019 käyttämä data oli haettu englantilaisesta biopankista. 423 604 henkilön dataa käytettiin ja sydän- ja verisuonisairauksien riski perustui 473 tekijään.

He käyttivät mallissaan tekijöitä, joita ei aiemmissa malleissa yleensä huomioitu, kuten koehenkilöiden kävelynopeutta tai koehenkilöiden itse ilmoittamaa arviota terveydestään. Lisäksi heidän mallinsa otti huomioon esimerkiksi henkilön historian diabeteksen kannalta.

### 4.3 Kognitiivinen konenäkökehys syövän arviointiin

Racoceanu & Capron, 2015, esittivät avainkonseptina semanttisuutta ohjaamaan päätöksentekoprosessia kuvien analysointiprotokollassa. Kaikki päätökset viedään semanttiseen ja formaaliin maailmaan, jossa niistä johdetaan olemassa olevan tiedon avulla päätös. Formalisoinnin avulla on mahdollista saada toistettavuutta, sekä jäljitettävyyttä päätöksille.

Kehyksen pohjalta kehitettiin prototyyppi, joka keskittyi syövän arviointiin ja erityisesti mitotoittisen lukeman mittaamiseen (*mitotic count*). Prototyyppi koostuu kahdesta osasta: semanttinen ydin ja kuvantunnistusalgoritmit. Ytimeen kuuluu semanttinen päätöksentekokone ja ontologia. Päätöksentekokone on vastuussa järjestelmän tekemistä päätöksistä ja ontologiassa on sitä avustavaa dataa sekä kuvankäsittelystä että histopatologiasta.

## 4.4 Tekoälypohjainen rintasyövän etäpesäkkeiden havainnointi

Ensisijaisen kasvaimen etäpesäkkeet vaikuttavat monen tapauksen kohdalla potilaan ennusteeseen ja hoitopäätöksiin. Imusolmukkeiden histologinen tunnistus on työlästä ja altista virheille. Tekoälyn on mahdollista käydä jokainen kudosesiintymä dialla läpi patologia tarkemmin. Tutkijat ovatkin kehittäneet LYNA (LYmph Node Assistant) - tekoälysovelluksen, joka tunnistaa koepaloista rintasyövän etäpesäkkeitä.

LYNA:aa käytettiin valmiiseen Camelyon16 - datasettiin, joka koostuu 399 kuvasta. Sen suoritusta arvioitiin kahdella tavalla: koko dian sekä kudosesiintymäkohtaisesti ROI (Region of Interest) mukaan. LYNA:aan eivät vaikuttaneet kohina kuten kuvan ilmakuplat tai huono tarkennus.

Kuvien väri vaihteli paljon, koska eri laboratoriot käyttivät erilaisia väriaineita, joten kuvien väriavaruus oli tärkeää normalisoida. Tähän käytettiin yksinkertaistettua versiota Bejnordi et al., 2016 kehittämästä algoritmista.

## 4.5 Radiologian ja patologian tekoälysovellukset – eturauhassyöpä

Eturauhassyöpä on maailmanlaajuisesti kuudenneksi yleisin syöpä. Multiparametrinen magneettikuvaus (mpMRI) on yleinen tapa tunnistaa eturauhassyöpää. Myös patologialla on suuri merkitys eturauhassyövän riskien arvioinnissa ja hoidon määrittämisessä.

Radiologian luokitus saattaa erota patologisesta luokituksesta. Harmon et al., mukaan tiheyden kartoittamisen parantaminen on yksi mahdollisuus, jolla saadaan täsmennettyä patologiasta arviota.

Tällä hetkellä sekä radiologia, että patologia ovat rajoittuneita eturauhassyövän luokittelussa. Suuri osa nykyisistä koneoppivista malleista on rakennettu heikosti luokitellulla datalla, joka ei ota huomioon kudoksen heterogeenisyyttä.

Päälöydökset:

- Tekoälyllä on potentiaalia auttaa eturauhassyövän tunnistamista ja luokittelua
- Nykyinen tekoälyn opetusdata ei vastaa patologiassa havaittua kudosten heterogeenisyyttä
- Patologiapohjaiset tekoälyratkaisut mahdollistavat korrelaatiota radiologian ja patologian välillä ja niitä voidaan hyödyntää radiologiapohjaisten tekoälyalgoritmien parantamiseen.

(Harmon et al., 2019)

## 4.6 Tekoäly mammografian ja digitaalisen kerroskuvauksen apuna

Perinteiset CAD-ohjelmat (*computer-aided diagnosis*) eivät ole johtaneet tarkempiin diagnooseihin, mutta ne ovat silti yleisesti käytössä. Koska syväoppivat (*deep learning*) neuroverkkoalgoritmit ovat kehittyneet pitkälle, on niitä mahdollista käyttää apuna mammografiassa ja rintojen digitaalisessa

kerroskuvauksessa (DBT). Tuloksien kliininen arviointi on kuitenkin vielä puutteellista eikä ole selvää miten syväoppimista tulisi käytännössä hyödyntää.

Tekoälyn on todettu arvioivan mammografiakuvia onnistuneesti. Joillakin syväoppivilla järjestelmillä on päästy jo radiologien tasolle (88%). Mallit jotka käyttivät yhdistelmää eri luokitusten määrittämisessä (luokka- ja pikselikohtainen luokittelu) olivat tarkimpia. Ongelmaksi mammografiakuissa osoittautui työläs datankeruu ja se että suorituskyky riippui lääkärin pohjatotuustiedon laadusta.

Toinen haaste mammografiakuissa on validointi, sillä eri valmistajat prosessoivat kuvia eri tavoin ja raakadataa ei yleensä ole saatavilla.

Sechopoulos, Teuwen, & Mann, 2020 kokosivat seuraavat tulokset:

- Neuroverkot pystyvät tekemään abstraktimpia esityksiä datasta ennen kuin kuva edes luokitellaan.
- Kuvissa saattaa olla eroja riippuen eri valmistajien kuvantamislaitteista ja tämä tulee ottaa huomioon ennen kuin otetaan syväoppivia malleja.
- Aiempia mammografiakuvia joissa on erilaisia modaliteetteja voidaan käyttää neuroverkon ennustuksen laadun parantamiseen.
- Syväoppivien järjestelmien suorituskyky on parempi kuin perinteisillä CAD- järjestelmillä.

## 4.7 Ihmisten geneettisten piirteiden ennustaminen

Ihmisen luontaisten piirteiden ennustusta syväoppivien järjestelmien avulla ei ole vielä täysin tutkittu. Bellot, Campos, & Pérez-Enciso, 2018 testasivat kahdella eri neuroverkkomenetelmällä UK biopankin tarjoamaa dataa. Käytettyjä menetelmiä olivat muun muassa monikerros- (MLP) ja konvoluutioneuroverkot (CNN).

Menetelmillä tarkasteltiin ”snippejä” (SNP<sup>2</sup>) eli populaatiossa esiintyviä pistemutaation aiheuttamia eroja DNA- ketjussa. Snippejä oli datassa noin 500 000. Käytetty data oli UK Biopankin aineistoa, josta valikoitiin valkoihoiset etäisesti toisilleen sukua ovat. Lopullisessa aineistossa oli noin 100 000 potilasta, joista noin 80 % käytettiin opetusdatana ja 20 % testidatana.

Testatut menetelmät toimivat lähes yhtä tehokkaasti, kun tarkasteltiin pituutta, joka on vahvasti periytyvä piirre. Muita piirteitä mitattaessa menetelmän tehokkuus riippui paljon snipistä, sekä piirteestä jota mitattiin. Kaiken kaikkiaan CNN oli kilpailukykyinen lineaarisen mallin kanssa, mutta ei löydetty menetelmää joka olisi toisaalta selkeästi lineaarista mallia parempi.

Käytetty CNN menetelmä pohjautui kuvantunnistukseen. Tutkijat totesivat, että tulokset saattaisivat parantua, jos CNN mukautettaisiin geenitutkimusta varten.

<https://www.genetics.org/content/210/3/809>

<sup>2</sup> [https://fi.wikipedia.org/wiki/Yhden\\_em%C3%A4ksen\\_monimuotoisuus](https://fi.wikipedia.org/wiki/Yhden_em%C3%A4ksen_monimuotoisuus)

## 5 Pohdinta

Pohdinta perustuu asiantuntijoiden kanssa käytyyn keskusteluun ja koottujen artikkelien tiivistelmiin ja mm. Lontoossa järjestetyn digipatologian konferenssin havaintoihin.

### Lainsäädännön vaikutus biopankkeihin ja tutkimukseen

Suomessa on lainsäädännön puitteissa kaksi vaihtoehtoa tutkimukselle: biopankki- ja rekisteritutkimus. Ensimmäiseen vaikuttaa biopankkilaki<sup>3</sup>, toiseen toissijaisen käytön laki<sup>4</sup>. Biopankkilaki on uudistumassa vastaamaan EU:n asettamia tietosuojavaatimuksia.

*“Laki määrää näytteenluovuttajan itsemääräämisoikeuden ja yksityisyyden suojan lisäksi monia biopankkien organisaatioon ja toimintaan liittyviä yksityiskohtia. Se säätää sekä aikaisemmin koottujen näytteiden siirtämisestä biopankkiin että tulevien kokoelmien keräämisestä.”, Biopankki.fi*

Rekisteritutkimus ei tarvitse (toistaiseksi) potilaan suostumusta ja kattaa kaikki rekisterissä olevat potilaat. Biopankkitutkimus taas kattaisi vain suostumuksen antaneet potilaat.

Uudistuksen myötä voi käydä niin, että biopankit nähdään lain silmin jatkossa näytevarastoina, joissa on vain itse näytteeseen suoraan liittyvää data, kuten näytteenottoaika ja pakastuslämpötila.

Sosiaali- ja terveysministeriön mukaan biopankkilain uudistus mahdollistaa näytteiden kokoamisen tutkimukseen, jonka yksityiskohdat eivät ole täysin tiedossa.

Lainsäädännön mukainen tietosuoja on kiristymässä 1.5.2021 alkaen. Tämä tarkoittaa muun muassa, että datan käyttöympäristöjen tulee olla sertifioituja ja auditoituja. Ennen tätä alkaneet tutkimukset voi jatkaa vanhalla lainsäädännöllä ja laki ei tule takautuvasti voimaan.

Eräs tutkimusta mahdollisesti hankaloittava tekijä on se, että perusterveydenhuollon tietoihin ei ole oikeuksia. Data koostuu näin ollen lähes yksinomaan erikoisterveydenhuollon datasta. Biopankkilainsäädännöllä ei ole tällä hetkellä mahdollista ylittää rekisterinpitäjää vaan data pyydetään Findatan kautta.

### 5.1 Tärkeimmät kehityskohteet

Hankkeen aikana asiantuntijoiden kanssa käydyissä keskusteluissa on noussut esille tekoälyyn liittyviä kehityskohteita, jotka olisivat relevantteja kotimaisesta näkökulmasta.

Tekstidatan hyödyntäminen tekoälyn keinoin on yksi tärkeimmistä suunnista, joissa edistystä toivotaan. Tekstissä on paljon asioita, joita ei ole rakenteisena olemassa ja tekoälyn hyödyntäminen tässä on mahdollista.

Auriassa on pyritty tunnistamaan tekoälyn keinoin esim. metastaasista kärsiviä potilaita tekstidatasta, koska diagnoosikoodin käyttö ei ole ollut kattavaa.

<sup>3</sup> <https://www.finlex.fi/fi/laki/alkup/2012/20120688>

<sup>4</sup> <https://www.finlex.fi/fi/laki/alkup/2019/20190552>

Histopatologiassa tekoälysovelluksilla voitaisiin tehostaa patologioiden työtä. Biopankeissa on paljon näytteitä ja niihin liittyvää historiallista dataa. Tällä datalla saataisiin tehtyä ja huomioitua paljon enemmän vaikutuksia kuin esim. avoimesti saatavalla rajoittuneilla aineistoilla.

Asiantuntijat toivovat selkeästi Suomeen myös kattavampaa digipatologiaa. Digipatologian avulla pystytään kaivamaan datasta esille asioita, joita patologi ei paljaalla silmällä näe. Jos onnistutaan keräämään riittävän paljon johdonmukaista dataa niin tekoälypohjaisten menetelmien toistettavuus voi olla myös parempi määritettäessä eri parameterien arvoja histopatologisista kuvista. Toisaalta tarkoitus ei ole korvata patologeja vaan tukea diagnoosien ja päätösten tekoa.

Kaiken kaikkiaan voidaan havaita selkeä tarve tekoälypohjaiselle päätöksenteon tueksi ja rakenteettoman tekstidatan tulkinneksi. Kuvien pitäisi liittää esimerkiksi potilaskertomuksista taustatietoa. Ideaalitulanteessa biopankkien käytössä voisi olla avoimella lisenssillä jaettava tekstinluokittelumenetelmä. Tämä olisi yksittäisiä pilotteja kestävämpi ratkaisu ja sillä olisi parempi validiteetti. On ehdotettu myös semanttisiin verkkoihin pohjautuvaa luokittelua.

Diagnostiikassa tekoäly voi avustaa mm. magneetti- ja röntgenkuvien tulkinneksi. Kehittäminen on mahdollista potilaskertomustiedon toisiokäytöllä. Näin voidaan saavuttaa nopeaa ja tehokasta diagnostiikkaa ja laadunvarmistusta rakenteellisen kirjaamisen avulla. Tekstinlouhinnalle ja avoimen lähdekoodin sovelluksille on kysyntää myös kuvantamissovelluksissa jotta taustatietoa saadaan tulkituksiin mukaan.

Tekoäly voidaan soveltaa myös palveluketjujen toiminnan kehittämiseen. Varhaisista ennusteista voidaan saada tukea esim. relapsiin, diabeteksen tai huonon haavanhoidon tunnistamiseen ja vanhuksien kanssa työskentelyyn. Tällä hetkellä tällaisia ei käytännössä ole mahdollista toteuttaa, koska riittävää dataa ei ole saatavilla.

## 5.2 Datan saatavuus

Valmiita aineistoja tekoälyn kehittämiseen on jo olemassa mm. biopankeissa, Kelalla ja THL:llä ja sitä on saatavilla tutkimuspyynnön kautta. On kuitenkin otettava huomioon se, että saatavilla oleva data on vain erikoisterveydenhuollosta.

Esimerkiksi diabeteksen tutkimus on vaikeaa ilman perusterveydenhuollon dataa. Haaste onkin, että perusterveydenhuollon tietojärjestelmissä ei ole tällä hetkellä valmiutta tehdä hyödyllisiä datapaketteja.

Asiantuntijat on esittäneet arvioita, että tekstidata on saatavilla jopa noin päivän viiveellä kun taas kuvantamisdata tulee arkistoista pyynnöstä.

Esimerkkinä voidaan esittää että Auriassa on noin 200 000 potilasta biopankkidatassa. Biopankkien määrät menevät yleensä suhteessa alueiden väestöpohjaan. Potilaista on käytössä käytännössä kaikki sairaalassa kirjattu data, mutta edelleen vain erikoissairanhoidon puolelta.

Kaikki data on saatavilla tutkimukseen tutkimuspyynnön kautta ja esimerkiksi Auriassa suurin osa vuodesta 2004 alkaen on jo valmiiksi digitoitua. Kuvantaminen tapahtuu käytännössä kokonaan tietokoneiden avulla, patologiakin siirtyy koko ajan enemmän digitaaliseen datan käsittelyyn. Tällä hetkellä kudokseteläsejä skannataan projektikohtaisesti digitaaliseen muotoon. Data päivittyy kerran yössä ja on käytettävissä heti.

## 5.3 Riskitekijät

Tekoälyn käyttöön liittyy luonnollisesti myös riskejä. Kehitys on vielä sillä tasolla, ettei tekoäly pysty päättämään milloin tiettyyn ongelmaan ei ole ratkaisua. Myös kykenemättömyys selittää logiikkaprosessia ja päätöksentekoa sekä luovuuden puute voivat vaikuttaa negatiivisesti tekoälyn käyttöön. (T. Siukonen, P. Neittaanmäki. "Mitä tulisi tietää tekoälystä")

Koska moderni tekoälyn kehittäminen perustuu lähes täysin dataan on tietosuojakysymyksistä tullut äärimmäisen tärkeitä. Toisiokäyttölain voimaan tultua terveystietojen käyttöympäristöjä koskevat vaatimukset ovat tiukentuneet. Toisiokäyttölain alaista terveystietoa voi jatkossa käsitellä vain Liikenne- ja viestintävirasto Traficomien hyväksymän tietoturvallisuuden arviointilaitoksen auditoimissa ympäristöissä. Toisiokäyttöympäristöjen rekisteriä ylläpitää Valvira ja sitä päivitetään sivulla: <https://www.valvira.fi/terveydenhuolto/toisilain-mukaiset-tietoturvalliset-kayttoymparistot/toisiokayttoymparistojen-rekisteri>

Biopankeissa tietoturva hoidetaan sairaaloiden infrastruktuuriratkaisujen avulla ja tiedot pysyvät palomuurien suojassa. Dataa myös käsitellään vain suojatuilla analyysikoneilla.

Tietosuojan edistäminen ja lainsäädäntö luonnollisesti myös monimutkaistaa monia asioita. Useammasta rekisteristä haettavien terveystietojen yhdistäminen pitää tapahtua Findatan (<https://findata.fi/>) toimesta.

Findata ohjeistaa tekemään tietolupa- tai muutoshakemuksen silloin, kun datan käyttö koskee:

- *usean julkisen sosiaali- ja terveysalan rekisterinpitäjän tietoja*
- *yhden tai usean yksityisen sosiaali- tai terveydenhuollon palvelunjärjestäjän rekisteritietoja, tai*
- *Kanta-palveluihin tallennettuja asiakastietoja.*

Esimerkiksi tilastokeskuksen kuolinsyydatan yhdistäminen sairaalan diagnoosidataan tapahtuu Findatan toimesta. Kehittäjille haasteita ovat kustannukset ja lupaprosesseihin kuluva aika. Yhtenä ratkaisuna prosessien nopeuttamiseen on esitetty synteettistä dataa. Mm. synteettisten nivelrikkoröntgenkuvien generointiin on kehitetty AI hub hankkeessa menetelmää tavoitteena tuottaa avoimesti jaettava nivelrikkoaineisto.

Kasvavan datamäärän myötä digipatologia on selkeästi tulevaisuuden ratkaisu datan käsittelyyn ja analysointiin. Tämä on havaittavissa kansainvälisissä konferensseissa yritysten esittelemissä tulevaisuuden sekä ohjelmistotuotteissa että tutkimuskirjallisuudessa.

Käytettävyys on keskeinen ominaisuus myös digipatologian sovelluksissa. Tärkeitä asioita ovat graafiset käyttöliittymät, reaaliaikainen vaste, tulkittavuus ja integroitavuus. Workflow ja käytettävyys tulisi olla isossa roolissa jotta työkalujen hyödyllisyys näkyisi selvästi käyttöönotossa. Tulevaisuudessa patologioiden saattaa olla välttämätöntä opiskella myös skriptien laatimisessa käytettävää Pythonia jolloin avoimia työkaluja voidaan kehittää tehokkaammin. Esimerkiksi HistoQC<sup>5</sup> on avoimen lähdekoodin työkalu histopatologisten näytelasien laaduntarkkailuun.

Tekoälyjärjestelmien toiminnan validointi on tärkeää riippumatta järjestelmän suoriutumisen tavoitetasosta. Pilvipalveluiden yleistyessä entistä tärkeämmäksi asiaksi nousee potilaiden

<sup>5</sup> <https://github.com/choosehappy/HistoQC>

tietoturva. Suomalaisen terveysdatan toisiokäsittelyn tietosuojavaatimukset koskevat myös ulkomaisia toimijoita.



## Lähteitä

1. Alaa, A. M., Bolton, T., Angelantonio, E. D., Rudd, J. H. F., & Schaar, M. V. D. (2019). Cardiovascular disease risk prediction using automated machine learning: A prospective study of 423,604 UK Biobank participants. *Plos One*, 14(5).  
<https://doi.org/10.1371/journal.pone.0213653>
2. Alpaydın, E. (2020). *Introduction to machine learning*. MIT Press.
3. Bejnordi, B. E., Litjens, G., Timofeeva, N., Otte-Holler, I., Homeyer, A., Karssemeijer, N., & Laak, J. A. V. D. (2016). Stain Specific Standardization of Whole-Slide Histopathological Images. *IEEE Transactions on Medical Imaging*, 35(2), 404–415.  
<https://doi.org/10.1109/tmi.2015.2476509>
4. Bukhari, S. usama K. (2020). Machine learning approaches for the histopathological diagnosis of prostatic hyperplasia. *Annals of Clinical and Analytical Medicine*, 11(9).  
<https://doi.org/10.4328/acam.20105>
5. Colling, R., Pitman, H., Oien, K., Rajpoot, N., & Macklin, P. (2019). Artificial intelligence in digital pathology: a roadmap to routine use in clinical practice. *The Journal of Pathology*, 249(2). <https://doi.org/https://doi.org/10.1002/path.5310>
6. Kayser, K., Gășrtler, J., Bogovac, M., Bogovac, A., Goldmann, T., Vollmer, E., & Kayser, G. (2010). AI (artificial intelligence) in histopathology--from image analysis to automated diagnosis. *Folia Histochemica Et Cytobiologica*, 47(3). <https://doi.org/10.2478/v10042-009-0087-y>
7. Lesko, L. J., Zineh, I., & Huang, S.-M. (2010). What Is Clinical Utility and Why Should We Care? *Clinical Pharmacology & Therapeutics*, 88(6), 729–733.  
<https://doi.org/10.1038/clpt.2010.229>
8. Liu, Y., Kohlberger, T., Norouzi, M., Dahl, G. E., Smith, J. L., Mohtashamian, A., ... Stumpe, M. C. (2018). Artificial Intelligence–Based Breast Cancer Nodal Metastasis Detection: Insights Into the Black Box for Pathologists. *Archives of Pathology & Laboratory Medicine*, 143(7), 859–868. <https://doi.org/10.5858/arpa.2018-0147-0a>
9. Racoceanu, D., & Capron, F. (2015). Towards semantic-driven high-content image analysis: An operational instantiation for mitosis detection in digital histopathology. *Computerized Medical Imaging and Graphics*, 42, 2–15.  
<https://doi.org/10.1016/j.compmedimag.2014.09.004>
10. Regitnig, P., Müller, H., & Holzinger, A. (2020). Expectations of Artificial Intelligence for Pathology. *Artificial Intelligence and Machine Learning for Digital Pathology Lecture Notes in Computer Science*, 1–15. [https://doi.org/10.1007/978-3-030-50402-1\\_1](https://doi.org/10.1007/978-3-030-50402-1_1)
11. Salto-Tellez, M., Maxwell, P., & Hamilton, P. (2019). Artificial intelligence-the third revolution in pathology. *Histopathology*, 74(3), 372–376. <https://doi.org/10.1111/his.13760>
12. Shalev-Shwartz, S., & Ben-David, S. (2009). *Understanding Machine Learning*.  
<https://doi.org/10.1017/cbo9781107298019>
13. Slaoui, M., & Fiette, L. (2010). Histopathology Procedures: From Tissue Sampling to Histopathological Evaluation. *Methods in Molecular Biology Drug Safety Evaluation*, 69–82.  
[https://doi.org/10.1007/978-1-60761-849-2\\_4](https://doi.org/10.1007/978-1-60761-849-2_4)
14. Harmon, S. A., Tuncer, S., Sanford, T., Choyke, P. L., & Turkbey, B. (2019). Artificial intelligence at the intersection of pathology and radiology in prostate cancer. *Diagnostic and Interventional Radiology*, 25(3), 183–188. <https://doi.org/10.5152/dir.2019.19125>

15. Sechopoulos, I., Teuwen, J., & Mann, R. (2020). Artificial intelligence for breast cancer detection in mammography and digital breast tomosynthesis: State of the art. *Seminars in Cancer Biology*. doi:10.1016/j.semcancer.2020.06.002
16. Bellot, P., Campos, G., & Pérez-Enciso, M. (2018, November 01). Can Deep Learning Improve Genomic Prediction of Complex Human Traits? Retrieved November 25, 2020, from <https://www.genetics.org/content/210/3/809>