



Ikä Tappaa Tieliikenteessä! - Osa 2 – Datan käsittely

Jatkamme tämän artikkelin edellisessä osassa mainitun Suomen Tieliikenneonnettomuudet 2005-2017 datan pyörittelyä, niin että voimme valituista muuttujista tehdä ennustemallin joka pyrkii ennustamaan todennäköisintä onnettomuusluokkaa.

[Väylävirasto](#) on julkaissut tietoaaineiston [Tieliikenneonnettomuudet](#) lisenssillä [Creative Commons Attribution 4.0 International License](#). Tietoaaineistossa on kullekin vuodelle kolme toisiinsa linkittyvää tiedostoa. Onnettomuudet-tiedostossa on onnettomuudet, joilla kullakin yksilöivä tunnus, ja massiivinen määrä taustatietoa itse onnettomuustilanteen tekijöistä. Henkilöt-tiedostossa on henkilötietoa, joilla kullakin yksilöivä tunnus ja se linkittyy onnettomuustunnuksella onnettomuudet-tiedostoon. Osalliset-tiedosto linkittyy sekä onnettomuustunnuksen että henkilötunnuksen kautta edellisiin tiedostoihin, tämä tiedosto kertoo lähinnä kunkin osallisen ajoneuvotiedon, sekä loukkaantumis/kuolemis tiedon, joka löytyy myös henkilöt-tiedostosta. Tietoaaineiston tehokas käyttö vaatii jonkin verran relaatiotietokanta tietämystä, että voi yhdistää kunkin tiedoston sisältämän datan toisiinsa.

Koska data on vuosikohtaisissa tiedostoissa, aloittakaamme niiden yhdistely yhteen tiedostoon. Lisäksi data sisältää relaatioita ja rivejä on satojatuhansia, niin käytetään tässä tapauksessa ehkä vähiten käytettyä Microsoft Office paketin ohjelmistoa nimeltä Access, joka on pieni relaatiotietokantaohjelma. Käytössä oli Microsoft Office 365 versio 1912. Excel tukee maksimissaan miljoonan rivin taulua, mutta osoittaa hyytymisen merkkejä jo muutaman sadantuhannen rivin kohdalla, jos aikoo käsitellä ja muokata dataa isompana massana, eikä se sovellu relaatioiden käsittelyyn kovinkaan luontevasti. Accessin rajana on 2 GB Access-tiedoston koko, jota voi tarvittaessa kiertää linkittämällä useita tiedostoja ulkoisina datalähteinä.

Huom! Kyseisen tietoaaineiston tiedostojen yhdistely onnistuu millä tahansa tekstinkäsittelyohjelmallakin, sillä jokaisen vuoden tiedostossa on samat sarakkeet ja sarakenimet. Seuraava harjoite on lähinnä kertomaan kuinka Accessiin voi tuoda ulkoista dataa.

Seuraava on yleistasoinen ohje Access:in käyttöön, josta selviää suoritettut askeleet muttei yksityiskohtaisia ohjeita ohjelmiston käyttöön, jottei artikkelista paisuisi 100 sivuista opusta. Vaikka et olisi koskaan käyttänyt Access:ia, suosittelen kokeilemaan, siihen tutustuminen pitäisi onnistua ihan alla olevilla ohjeillakin.

Luodaan ensin tyhjä Access-tietokanta nimeltään "Tieliikenneonnettomuudet-staging.accdb" datakansion juureen, jonne tietoaaineiston vuosikansiot ovat purettuna. Tämän jälkeen tuodaan kukin tiedosto Accessiin, External Data – New Data Source – From File – Text File. Avautuvasta dialogista valitaan "Import the source data into a new table in the current database", koska tässä tapauksessa data on staattista ja haluamme tarkistaa datan indeksoinnin (yksilöivät tunnuks) eheyden. Tuodaan data Accessin ehdottamalla datatyypeillä, koska haluamme vain yhdistää datan raakana yhteen tiedostoon ilman muutoksia. Näin vältämme sen ettei tässä vaiheessa karsita tietuita joiden sarakkeen datatyyppi vaihtelee soluttain, esim. numeroita ja kirjaimia.

Link Text Wizard

What delimiter separates your fields? Select the appropriate delimiter and see how it looks in the preview window.

Choose the delimiter that separates your fields:

☐ Tab ☒ Semicolon ☐ Comma ☐ Space ☐ Other:

☒ First Row Contains Field Names Text Qualifier:

Onnett_id	Osall_id	Henkilo_id	Osnro	Kulj_matk	Ikä
6551782	11677943	11520814	1	KU	20
6551782	11677944	11520815	2	KU	24
6580617	11677952	11520824	1	KU	78
6580617	11677953	11520825	2	MA	59
6580617	11677953	11520826	2	KU	55
6571962	11677972	11521122	1	KU	18
6571962	11677973	11521123	2	KU	
6673527	11677994	11521146	1	KU	44
6673527	11677995	11521147	2	KU	
6565101	11678000	11521152	1	MA	9
6565101	11678000	11521153	1	MA	7
6565101	11678000	11521154	1	KU	40
6565101	11678001	11521155	2	KU	58
6490335	11678803	11521180	1	KU	49

Advanced...

Cancel < Back Next > Finish

Tieliikenneonnettomuudet_2005_hlo Link Specification

File Format: ☒ Delimited ☐ Fixed Width Field Delimiter: ; Text Qualifier: {none}

Language: Finnish Code Page: Western European (Windows)

Dates, Times, and Numbers

Date Order: DMY ☒ Four Digit Years Date Delimiter: . ☐ Leading Zeros in Dates Time Delimiter: . Decimal Symbol: ,

Field Information:

Field Name	Data Type	Skip
Onnett_id	Long Integer	<input type="checkbox"/>
Osall_id	Long Integer	<input type="checkbox"/>
Henkilo_id	Long Integer	<input type="checkbox"/>
Osnro	Long Integer	<input type="checkbox"/>
Kulj_matk	Short Text	<input type="checkbox"/>
Ikä	Short Text	<input type="checkbox"/>
Sukupuoli	Short Text	<input type="checkbox"/>
Kortti	Short Text	<input type="checkbox"/>
Seuraus	Short Text	<input type="checkbox"/>
Seurausl	Short Text	<input type="checkbox"/>

Koska tietoaaineistossa kussakin tiedostossa on kullakin tietueella yksilöivä tunnus (ID) jota voidaan käyttää indeksinä, aseta oma Primary Key. Onnettomuudet-tiedoston kohdalla valitse "Onnett_id", Henkilöt-tiedoston kohdalla "Henkilo_id" ja Osalliset-tiedoston kohdalla "Osall_id". Näin voimme varmistua että indeksointi on eheä ja ettemme vahingossa tuo samaa dataa kahdesti, sillä se ei onnistu jos Primary Key sarakkeesta löytyy samoja arvoja.

Import Text Wizard

Microsoft Access recommends that you define a primary key for your new table. A primary key is used to uniquely identify each record in your table. It allows you to retrieve data more quickly.

☐ Let Access add primary key.

☒ Choose my own primary key. Onnett_id

☐ No primary key.

Onnett_id	Tienpit	Tienpitsel	Tie	Aosa	Aet	Ajr	Vuosi	Kk	Päivä	Kuolleet	Loukkaa
6639049	1	Liikennevirasto	1	3	0	1	2005	1	22.01.2005	0	0
6672801	1	Liikennevirasto	1	3	30	1	2005	1	11.01.2005	0	0
6780477	1	Liikennevirasto	1	3	200	1	2005	8	27.08.2005	0	1
6488786	1	Liikennevirasto	1	3	500	1	2005	2	24.02.2005	0	0
6456756	1	Liikennevirasto	1	3	1000	1	2005	10	17.10.2005	0	0
6675084	1	Liikennevirasto	1	3	1000	1	2005	12	23.12.2005	0	0
6852453	1	Liikennevirasto	1	3	2200	1	2005	6	11.06.2005	0	0
6632689	1	Liikennevirasto	1	3	2300	1	2005	1	21.01.2005	0	0
6729285	1	Liikennevirasto	1	3	3300	1	2005	11	13.11.2005	0	0
6551868	1	Liikennevirasto	1	4	0	1	2005	3	14.03.2005	0	0
6804235	1	Liikennevirasto	1	4	0	1	2005	5	27.05.2005	0	0
6895816	1	Liikennevirasto	1	4	0	1	2005	6	19.06.2005	0	0
6640531	1	Liikennevirasto	1	4	200	1	2005	1	10.01.2005	0	0
6497817	1	Liikennevirasto	1	4	800	1	2005	12	21.12.2005	0	0

Advanced...

Cancel < Back Next > Finish

Seuraavan vuoden tiedostoa tuodessa valitaan "Append a copy of the records to the table" ja lisätään seuraavan vuoden data edellä luotuun tauluun. Toistetaan niin kauan että kunkin vuoden tiedostot ovat vastaavissa tauluissaan, jonka jälkeen Accessissa tulisi olla kolme erillistä tietokanta taulua Onnettomuuksille, Henkilöille ja Osallisille.

Joidenkin vuosien kohdalla Onnettomuus-tiedosto (2009,2010, 2012 ja 2016) antoi ilmoituksen tuontivirheestä muutaman rivin osalta, samalla luoden ImportErrors taulun, johtuen epäonnistuneesta tyyppikonversiosta joidenkin sarakkeiden kohdalla, josta ei kannata tässä tapauksessa välittää. Muutaman rivin puuttumisella ei käyttöskenaariossamme ole merkitystä.

Kun data on saatu tuotua tauluihin, onkin aika viedä ne tekstitiedostoiksi, joka tapahtuu kutakin taulua oikeaklikkaamalla – Export – Text File.

Nyt kun kaikki data on kolmessa tiedostossa, voimme analysoida datan ominaisuuksia. Tässä tapauksessa käytämme siihen [Azure Machine Learning työtilasta](#) (Preview:ssa kirjoituksen hetkellä) löytyvää tietoaaineiston profilointi ominaisuutta (Datasets – Profile). Lisää Azure portaalissa tilaukseesi Azure Machine Learning Workspace, josta löytyy suora linkki työtilaportaaliisi. Datan profilointia varten Azure Machine Learning työtilassa tarvitsemme ensin palvelinklusterin (Compute – Training Cluster - New), kevyimmän mahdollisen, mieluummin kolme noodia kun kolme tiedostoa niin voi profiloida kaikki yhtä aikaa ja jota voimme käyttää jatkossa koneoppimismallien kouluttamiseen. Palvelinklusteri osaa sammuttaa noodit jos niitä ei käytetä, joten käyttökustannukset jäävät vähäisiksi vaikka klusteria ei poistaisikaan joka käytön jälkeen. Toki kannattaa seuralla käyttökustannuksia, sillä Preview-vaiheen alkupuolella klusteri aiheutti yllättäviä kuluja, ei ilmeisesti osannut sammutella noodeja.

New Training Cluster

Compute name * ?

small-3node-comp

① Machine Learning Compute is a managed training environment consisting of one or more nodes. [Learn more.](#)

Region * ?

westeurope

Virtual Machine size * ?

Standard_DS2_v2

Virtual Machine priority * ?

Dedicated

Low Priority

Minimum number of nodes * ?

0

Maximum number of nodes * ?

3

Idle seconds before scale down * ?

120

> Advanced settings

Kun klusteri luotu, voidaan ladata tietoaaineistot (Datasets – Create Dataset), valitse Confirm Details -sivulla "Profile this dataset after creation".

Create dataset from local files

✓ Basic info

✓ Datastore and file selection

✓ Settings and preview

✓ Schema

● Confirm details

Confirm details

Basic info

Name

Tieliikenneonnettomuudet_hlo

Dataset version

1

Dataset type

Tabular

File settings

File format

Delimited

Delimiter

Semicolon

Encoding

ISO-8859-1

Column headers

Use headers from the first file

Skip rows

None

☒ Profile this dataset after creation

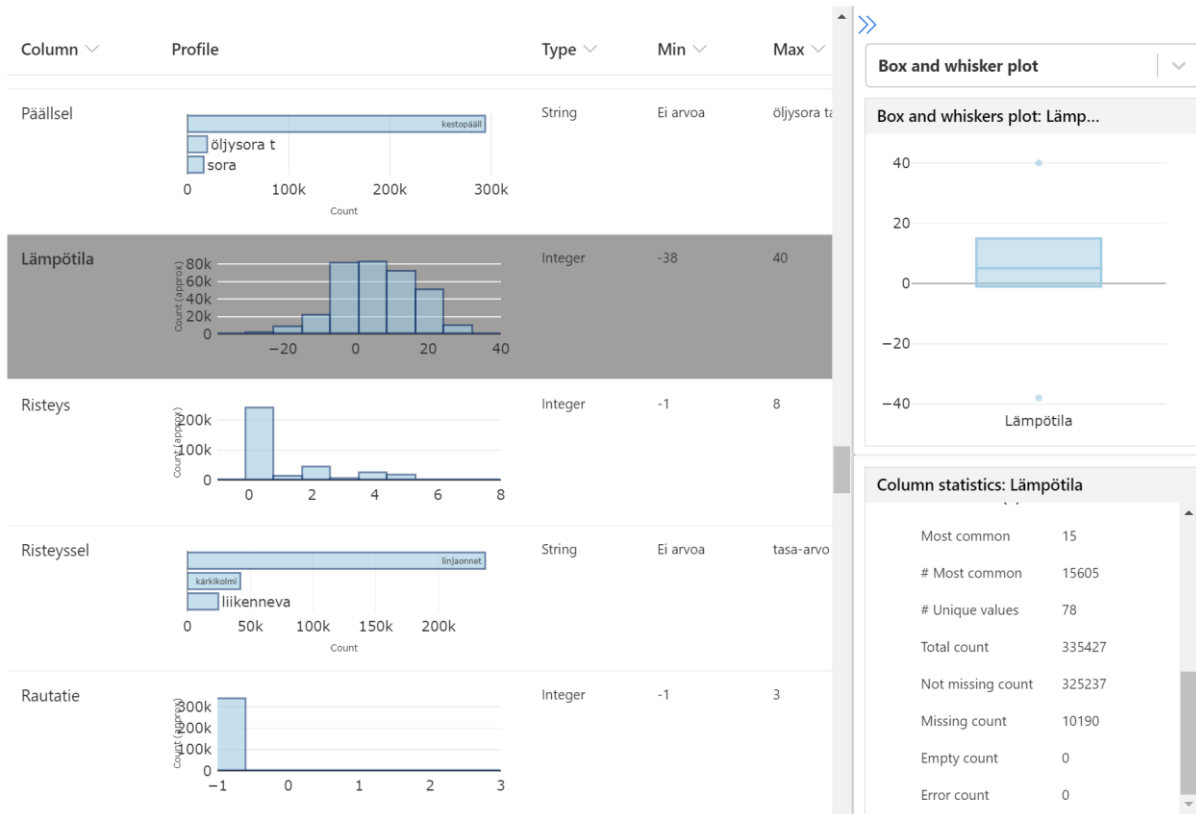
Select compute for profiling

small-3node-comp

Kun tietoaaineistot profiloitu, voi niitä tarkastella näkymästä Datasets – Explore – Profile. Profiili antaa monen näköistä статистиikkaa ja hyvän lähtökohdan arvioida mitkä muuttujat (feature) kannattaa ottaa lähempään tarkasteluun, jo siltä pohjalta ettei puuttuvia tai tyhjiä arvoja olisi suurin osa. Ko. tietoaaineistossa käytetty puuttuvina arvoina mm. "Ei arvoa", "-1", "-", "NA", "XX", "Tieto puuttuu" ja "Ei tiedossa" tyyppisiä syötteitä.

Preview **Profile**

Number of columns: 102 Number of rows: 335427



Muuttujia etsiessä kannattaa kiinnittää huomiota datan määrän lisäksi siihen että se on tasaisesti jakautunutta eri arvojen välillä eikä suuria vääristymiä olisi, ts. yksi arvo ylipainottunut muiden suhteen. Muuttujia valitessa on myös pidettävä mielessä mitä niistä on ennustemallin käyttäjälle helposti saatavilla, kuten kuukausi tai lämpötila, ja mitä on melko mahdoton kyseisen käyttäjän selvittää, kuten esimerkiksi tien ylläpitäjä tai tien näkyvän osuuden pituus.

Seuraavaksi luomme uuden Access tietokannan, "Tieliikenneonnettomuudet.accdb", jonne tuomme yhdistetyistä tietoaaineistoista haluamamme sarakkeet. Tällä kertaa linkitämme sen ulkoiseen datalähteeseen (Link to the data source by creating a linked table) eli edellä luomiimme tekstitiedostoihin, siltä varalta että saisimme lisää tuoretta vuosidataa (kyseistä tietoaaineistoa ei tosin enää päivitetä, dataa ei ole saatavilla ainakaan samassa muodossa ja linkitys ulkoiseen datalähteeseen luo myöhemmin selviävän haasteen datan jatkokäytölle). Toinen vaihtoehto olisi linkittää se suoraan edellä luotuun Access-tietokantaan (staging), jos sen datatyytit olisivat kunnossa, mitä ne eivät ole. Tällä kertaa kiinnitämme erityishuomion käytettyihin datatyypeihin emmekä tuo kaikkia sarakkeita.

Henkilöt-tiedostosta tuomme vain seuraavat sarakkeet, kaikki muut jätetään tuomatta (Skip):

- Onnett_id
- Osall_id
- Henkilö_id
- Osnro
- Kulj_matk
- Ikä

- Sukupuoli
- Seuraus
- Seuraussel

Samalla vaihdamme "Ikä" sarakkeen datatyyppiin numeroksi (Long Integer) ja tarkastamme että muissakin numerosarakkeissa datatyyppi on oikein.

Onnettomuus-tiedostosta tuomme vain seuraavat sarakkeet, kaikki muut jätetään tuomatta, datatyyppiin muutos suluissa:

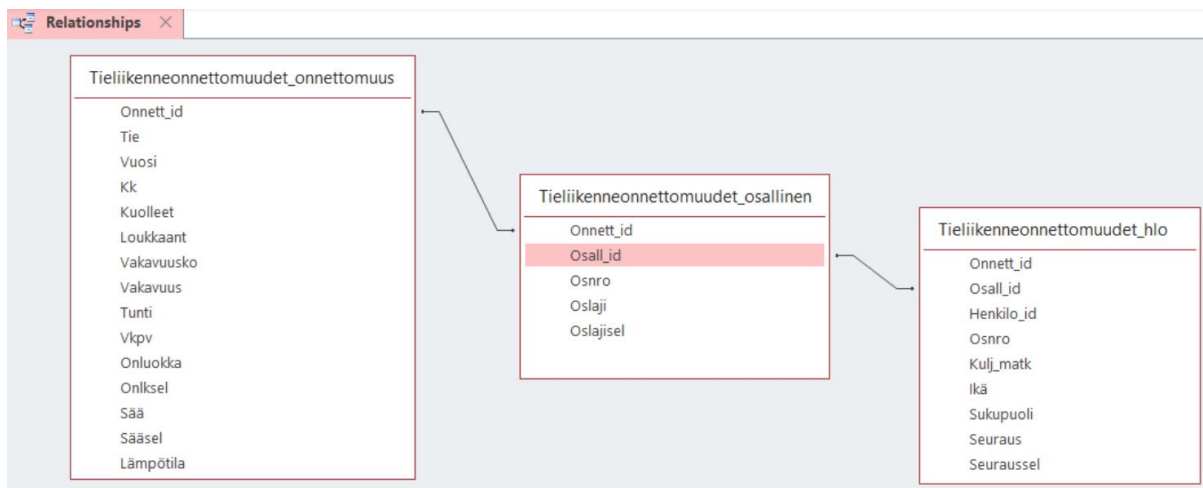
- Onnett_id
- Tie (Long Integer)
- Vuosi
- Kk
- Kuolleet
- Loukkaant
- Vakavuusko
- Vakavuus
- Tunti
- Vkp
- Onluokka
- Onlksel
- Sää
- Sääsel
- Lämpötila (Long Integer)

Osallinen-tiedostosta tuomme vain seuraavat sarakkeet, kaikki muut jätetään tuomatta:

- Onnett_id
- Osall_id
- Osnro
- Oslaji
- Oslajisel

Nyt kun ulkoiset datalähteet on linkitetty, voimme luoda niiden väliset suhteet (Database Tools – Relationships), tuo relaatiotietokantojen mystisin ominaisuus jossa voi mennä niin monella tavoin mönkään. Luomalla vääränlaisia suhteita voimme lisätä tai rajata tietueiden määrää vaikka se ei olisikaan tarkoitus. Linkitys tapahtuu Relationships-näkymässä lisäämällä tahdotut taulut ja raahaamalla linkittävistä (yksilöivästä, yleensä ID) sarakkeesta toiseen, sekä määrittelemällä linkityssuhde (Join Type).

Linkitetään Onnettomuus-tilu Osallinen-tiluun, niin että kummankin yhteneväsiet "Onnett_id" tietueet näytetään (Inner Join). Sen lisäksi linkitetään Osallinen-tilu Henkilöt-tiluun, niin että kummankin yhteneväsiet "Osall_id" tietueet näytetään. Näin siksi että haluamme vain semmoista dataa joka löytyy kaikista tauluista ja vähennämme jatkossa tarveta poistaa tietueita joille ei löydykään arvoa kaikista sarakkeista. Tämän jälkeen kun luodaan kyselyitä (Query) on linkitys taulujen välillä oletuksena luotu eikä sitä tarvitse enää tehdä, muuttaa vain niissä tapauksissa kun halutaan jotain muuta lopputulemaa kuin oletuksena on tarjolla, esim. kaikki tietueet Onnettomuus-tilusta vaikka niille ei löytyisikään vastinetta muista tauluista (Left Join).



Seuraavaksi voimmekin alkaa tekemään tarkoituksenmukaisia kyselyitä (Query) tauluihin, Create – Query Wizard – Simple Query Wizard. Luodaan ensin yleisnäkymä kaikkeen ja nimitetään se "masterQuery":ksi, eli lisää Onnettomuus-taulusta kaikki sarakkeet, Osallinen-taulusta kaikki paitsi "Onnett_id"-sarake ja Henkilöt-taulusta kaikki paitsi "Onnett_id", "Osall_id" ja "Osnro"-sarakeet.

Seuraavat kyselyt tehdään "masterQuery":ä vasten. Tämä siksi että mahdollisten taulumuutosten tapahtuessa on helpompi korjata yksi kysely kuin kymmentä. Esimerkiksi itse huomasin myöhemmin, kun olin luonut jo useita kyselyitä, että Henkilöt-taulun "Ikä"-sarake ei ollutkaan numeerinen vaan tekstiä, jonka lajitteluominaisuudet ovat tyystin erilaiset kuin numeroilla (1,10,11,12,...2,20,21), eikä datatyyppin vaihto suoraan taulun ominaisuuksista onnistunutkaan, joten jouduin tuomaan taulun datan uusiksi. Varovaisena toin taulun eri nimellä, että näin datatyyppin vaihdoksen vaikutukset, ja loin vuorovaikutussuhteet uuteen tauluun kuten aikaisemmin. Tämän jälkeen vaihdoin vain master-kyselyssä kyseltävän taulun uuteen, koska sarakeotsikot pysyivät samana oli vaihdos hyvin yksinkertainen. Muihin kyselyihin ei tarvinnut tehdä muutoksia ja ne toimivat samoin tein. Ketjutetut kyselyt voivat vaikuttaa suorituskyykyyn, mutta pienympäristössä erot ovat merkityksettömiä.

Mahtavaa! Olet tehnyt pohjatyöt ja Access-kantasi sekä "masterQuery" näyttää jotakuinkin tältä.

Onnett_id	Tie	Vuosi	Kk	Kuolleet	Loukkaant	Vakavuusko	Vakavuus	Tunti	Vkp	Onluokka	Onlksel	Sää
6422406		2011	4	0	0	0	0	15	Torstai	1	Yksittäisonnetti	4
6422407		2010	8	0	0	0	0	2	Torstai	1	Yksittäisonnetti	-1
6422408		2010	11	0	0	0	0	8	Perjantai	1	Yksittäisonnetti	-1
6422410		2010	1	0	0	0	0	12	Tiistai	2	Käntymisönn	1
6422410		2010	1	0	0	0	0	12	Tiistai	2	Käntymisönn	1
6422411		2009	10	0	0	0	0	15	Lauantai	6	peräajajo-onne	3
6422411		2009	10	0	0	0	0	15	Lauantai	6	peräajajo-onne	3
6422412		2009	10	0	0	0	0	2	Perjantai	1	Yksittäisonnetti	2
6422413		2010	9	0	0	0	0	17	Torstai	1	Yksittäisonnetti	-1
6422414		2010	9	0	0	0	0	15	Torstai	13	Muu onnettomu	-1
6422415		2009	1	0	0	0	0	21	Sunnuntai	2	Käntymisönn	2

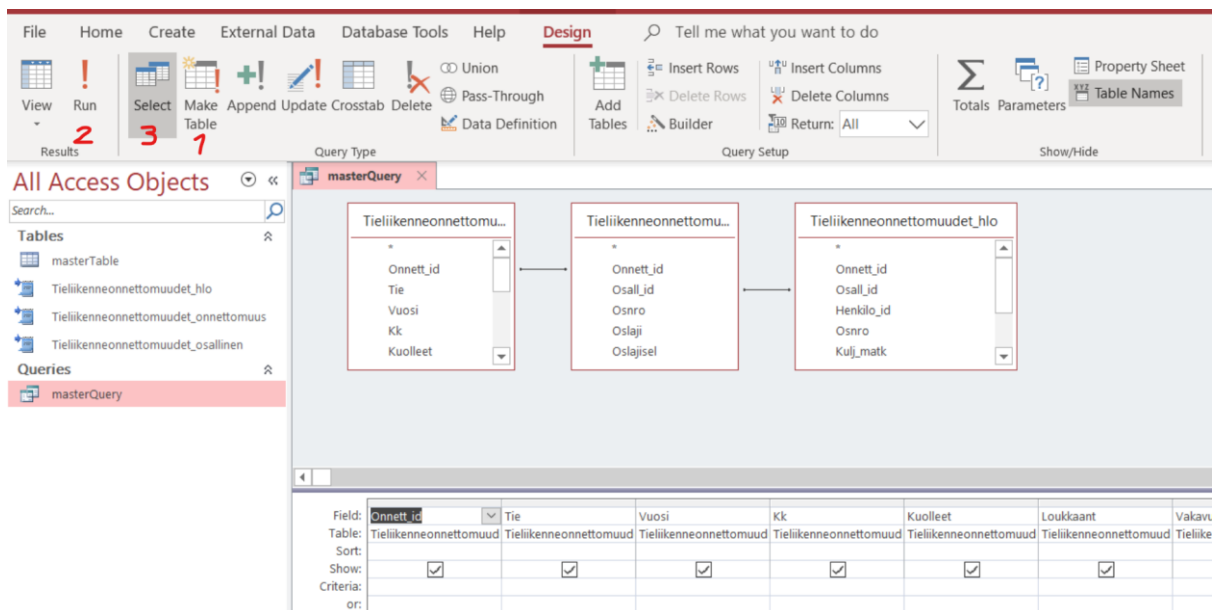
Jos haluat vilkaista miltä konepellin alla näyttää, vaihda näkymiä Datasheet View, SQL View ja Design View välillä esim. oikean alakulman pikanapeilla. Jos teet muutoksia ja haluat ne säilyttää, tallenna muutokset vasemman yläkulman Save-pikanapilla.

"MasterQuery":n SQL lauseke näyttää tältä.

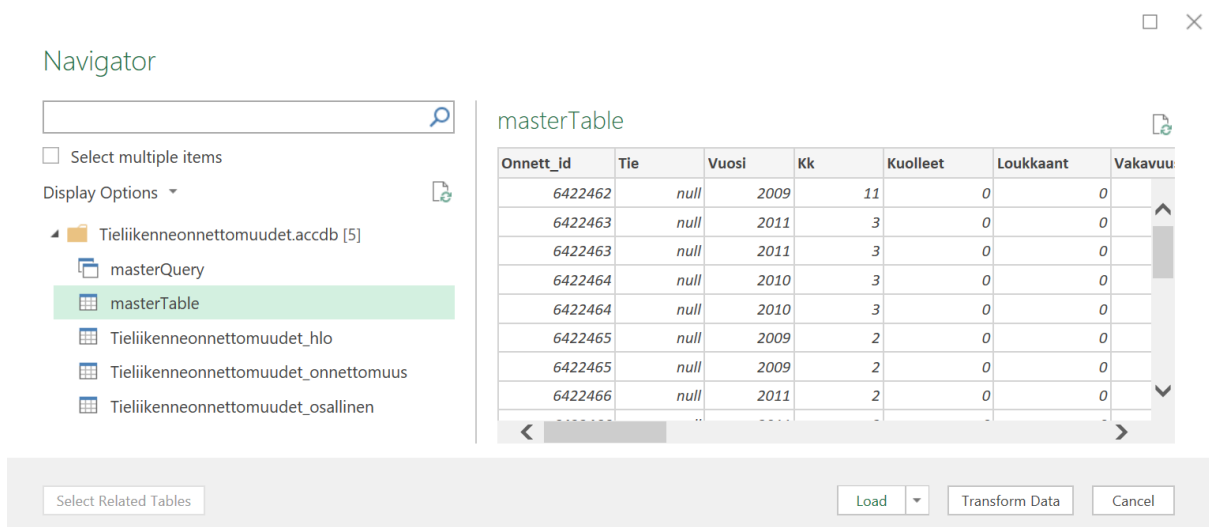
```
SELECT Tieliikenneonnettomuudet_onnettomuus.Onnett_id,  
Tieliikenneonnettomuudet_onnettomuus.Tie,  
Tieliikenneonnettomuudet_onnettomuus.Vuosi,  
Tieliikenneonnettomuudet_onnettomuus.Kk,  
Tieliikenneonnettomuudet_onnettomuus.Kuolleet,  
Tieliikenneonnettomuudet_onnettomuus.Loukkaant,  
Tieliikenneonnettomuudet_onnettomuus.Vakavuusko,  
Tieliikenneonnettomuudet_onnettomuus.Vakavuus,  
Tieliikenneonnettomuudet_onnettomuus.Tunti,  
Tieliikenneonnettomuudet_onnettomuus.Vkpv,  
Tieliikenneonnettomuudet_onnettomuus.Onluokka,  
Tieliikenneonnettomuudet_onnettomuus.Onlksel,  
Tieliikenneonnettomuudet_onnettomuus.Sää,  
Tieliikenneonnettomuudet_onnettomuus.Sääsel,  
Tieliikenneonnettomuudet_onnettomuus.Lämpötila,  
Tieliikenneonnettomuudet_osallinen.Osall_id,  
Tieliikenneonnettomuudet_osallinen.Osnro,  
Tieliikenneonnettomuudet_osallinen.Oslaji,  
Tieliikenneonnettomuudet_osallinen.Oslajisel,  
Tieliikenneonnettomuudet_hlo.Henkilo_id,  
Tieliikenneonnettomuudet_hlo.Kulj_matk, Tieliikenneonnettomuudet_hlo.Ikä,  
Tieliikenneonnettomuudet_hlo.Sukupuoli,  
Tieliikenneonnettomuudet_hlo.Seuraus,  
Tieliikenneonnettomuudet_hlo.Seuraussel  
  
FROM (Tieliikenneonnettomuudet_onnettomuus INNER JOIN  
Tieliikenneonnettomuudet_osallinen ON  
Tieliikenneonnettomuudet_onnettomuus.[Onnett_id] =  
Tieliikenneonnettomuudet_osallinen.[Onnett_id]) INNER JOIN  
Tieliikenneonnettomuudet_hlo ON Tieliikenneonnettomuudet_osallinen.[Osall_id]  
= Tieliikenneonnettomuudet_hlo.[Osall_id];
```

Koska Access:ista puuttuu datan visualisointi, niin graafien tekemiseen Excel toimii vallan mainiosti isommillakin datamäärillä. Voisimme käyttää ulkoisena datalähteenä juuri luomaamme "masterQuery":ä Access-tietokannassamme, jos se taasen EI käyttäisi ulkoiseen datalähteeseen linkitettyjä tauluja vaan data olisi tuota Access-kantaan. Jäljelle jää muutama vaihtoehto, "masterQuery":n datan vieminen ulos tekstitiedostoksi tai Make Table Query. Valitaan jälkimmäinen.

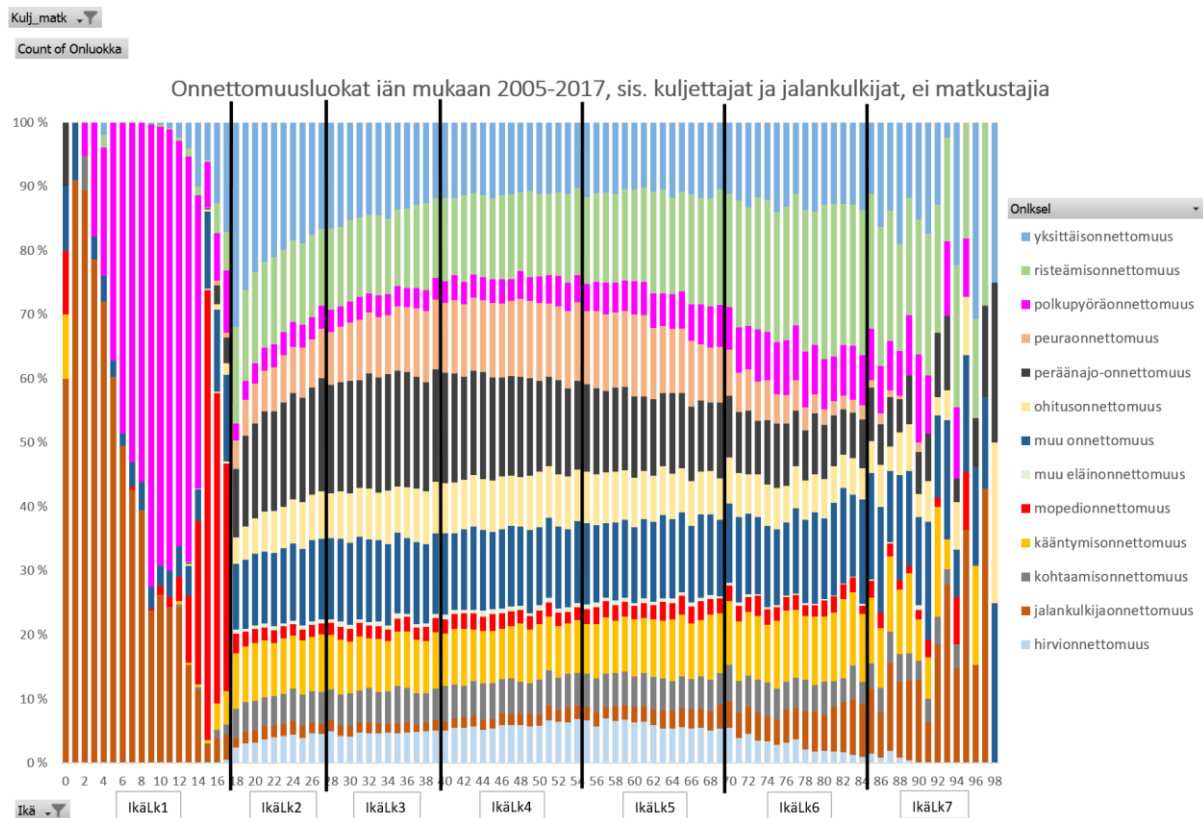
Mene "masterQuery":n Design näkymään, jos kysely ei ole avoinna esim. oikeaklikkaamalla kyselyn nimeä – Design View - paina Make Table nappia (1) - anna taululle nimi "masterTable" – OK – paina Run nappia (2) ja odota, kun uusi taulu on muodostunut - vaihda "masterQuery":n tyyppiä takaisin select query painamalla Select nappia (3).



Luo tyhjä Excel tiedosto "Tieliikenneonnettomuudet.xlsx", linkitä se Access-kannan "masterTable" tauluun, Data – Get Data – From Database – From Microsoft Access Database – etsi ja valitse Tieliikenneonnettomuudet Access-tiedosto – valitse "masterTable" - Load.



Luo näyttäviä graafeja, esimerkkejä voit katsella tämän artikkelisarjan ensimmäisestä osasta, ja tutki niiden kautta eri muuttujien vaikutusta onnettomuuksiin, jotta löydät vaihtoehtoisimmat valittavaksi ennustemalliin.



Kun käsitys muodostunut mitä muuttujia kannattaa ottaa mukaan ensimmäisiin ennustemallikokeiluihin, niin luodaan sen pohjalta kysely, jonka data viedään tekstitiedostoon. Käydään lävitse yhden mahdollisen muuttujavalikoiman luonti ja tällä kertaa ilman wizardia.

Access-tietokannassa, luo uusi Query, Create – Query Design – Select table dialogista vaihda Query näkymä ja valitse "masterQuery" - Add. Design näkymässä, kaksoisklikkaa lisäystä "masterQuery":stä "Ikä"-sarake – poista valinta Show laatikosta – lisää Criteria riville ehto: >0.

Yllä olevasta kuvasta huomaat että iät on jaettu luokkiin tiettyjen onnettomuusluokkakuvien muutosten perusteella, varsinkin nuoruusvuodet ja vanhuus, väliuodet pyritti jakamaan tasaisesti. Luokaamme kyseiset ikäluokat kyselyyn suoraan "Ikä" sarakkeen viereen.

Kopioi seuraava Field riville "Ikä"-sarakkeen viereen, voit painaa Builder nappia kun oikea sarake Field riviltä valittuna niin helpompi nähdä kokonaisuus: IkäLk: If([Ikä] Between 0 And 17;1;If([Ikä] Between 18 And 27;2;If([Ikä] Between 28 And 39;3;If([Ikä] Between 40 And 54;4;If([Ikä] Between 55 And 69;5;If([Ikä] Between 70 And 84;6;If([Ikä]>84;7))))))

Kaksoisklikkaa "Sää"-sarake – lisää Criteria riville ehto: >0. Kaksoisklikkaa "Kk"-sarake – lisää Criteria riville ehto: >0. Kaksoisklikkaa "Tunti"-sarake – lisää Criteria riville ehto: >0. Kaksoisklikkaa "Onluokka"-sarake – lisää Criteria riville ehto: >0 And <13.

Koska onnettomuusluokat sisältävät samankaltaisuuksia ja niitä on syytä yhdistellä ennustetarkkuuden parantamiseksi, luomme valmiiksi "UusiOnLk"-sarakkeen perustellulla yhdistelysäännöllä (onnettomuusprofiilit samankaltaisia), jota voi myöhemmin hienosäätää esim. uuden onnettomuusluokkanumeroinnin suhteen (numerointi kasvanut isoksi useamman iterointikierroksen jälkeen).

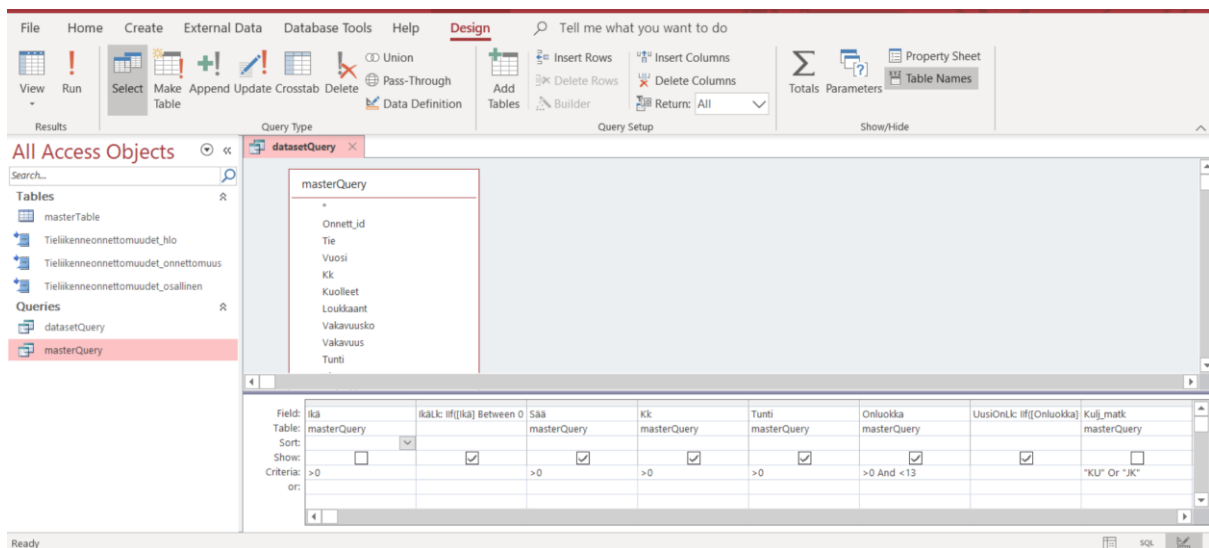
Kopioi seuraava Field riville "Onluokka"-sarakkeen viereen, muista käyttää Builder näkymää: UusiOnLk:

```
IIf([Onluokka]=1;1;IIf([Onluokka]=2;44;IIf([Onluokka]=3;43;IIf([Onluokka]=4;44;IIf([Onluokka]=5;43;
IIf([Onluokka]=6;6;IIf([Onluokka]=7;27;IIf([Onluokka]=8;27;IIf([Onluokka]=9;9;IIf([Onluokka]=10;30;
IIf([Onluokka]=11;30;IIf([Onluokka]=12;30;IIf([Onluokka]=13;13;[Onluokka])))))))))))
```

Kaksoisklikkaa "Kulj_matk"-saraketta, poista valinta Show laatikosta - lisää Criteria riville ehto: "KU" Or "JK". Näin valitsemme pois onnettomuuteen useimmiten syyttömät matkustajat, vain kuljettajat ja jalankulkijat ovat mukana.

Paina Save nappia, anna nimeksi "datasetQuery". Kun kysely tallennettu, paina Run nappia.

"DatasetQuery" näyttää Desing näkymässä tältä.



"DatasetQuery":n SQL lauseke näyttää tältä.

```
SELECT IIf([Ikä] Between 0 And 17,1,IIf([Ikä] Between 18 And 27,2,IIf([Ikä]
Between 28 And 39,3,IIf([Ikä] Between 40 And 54,4,IIf([Ikä] Between 55 And
69,5,IIf([Ikä] Between 70 And 84,6,IIf([Ikä]>84,7)))))) AS IkäLk, masterQuery.Sää,
masterQuery.Kk, masterQuery.Tunti, masterQuery.Onluokka,
IIf([Onluokka]=1,1,IIf([Onluokka]=2,44,IIf([Onluokka]=3,43,IIf([Onluokka]=4,44,IIf(
[Onluokka]=5,43,IIf([Onluokka]=6,6,IIf([Onluokka]=7,27,IIf([Onluokka]=8,27,IIf([O
nluokka]=9,9,IIf([Onluokka]=10,30,IIf([Onluokka]=11,30,IIf([Onluokka]=12,30,IIf([
Onluokka]=13,13,[Onluokka]))))))))))) AS UusiOnLk
```

```
FROM masterQuery
```

```
WHERE (((masterQuery.Ikä)>0) AND ((masterQuery.Sää)>0) AND
((masterQuery.Kk)>0) AND ((masterQuery.Tunti)>0) AND
((masterQuery.Onluokka)>0 And (masterQuery.Onluokka)<13) AND
((masterQuery.Kulj_matk)="KU" Or (masterQuery.Kulj_matk)="JK"));
```

Nyt ensimmäinen datasettisi on tekstitiedostoon vientiä vaille valmis, oikeaklikkaa "datasetQuery" – Export – Text File. Suosittelen että viet vain jomman kumman "Onluokka"- ja "UusiOnLk"-sarakkeista, poista Show valinta Design näkymässä toiselta – Save - Run. Ne ovat arvoja jota ennustemalli pyrkii ennustamaan ja on havaittu niiden ollessa yhdessä datasetissä aiheuttavan ongelmia ennustemallin pyytämiin muuttujiin vaikka kuinka olisi laitettu "Ignore" valintaa sarakkeeseen ennustemallia koulutettaessa.

Seuraavassa osassa pyrimmekin paneutumaan itse ennustemallin kouluttamiseen, miksi edellä kerrottu "UusiOnLk" oikein luotiin ja muita yksityiskohtia koneoppimismallista ja sen hyödyntämisestä.

Tämän LinkedIn artikkelisarjan osat:

Osa1: <https://www.linkedin.com/pulse/ikä-tappaa-tieliikenteessä-osa-1-raakadata-marko-peltojoki>

Osa2: <https://www.linkedin.com/pulse/ikä-tappaa-tieliikenteessä-osa-2-datan-käsittely-marko-peltojoki>

Osa3: <https://www.linkedin.com/pulse/ikä-tappaa-tieliikenteessä-osa-3-ennustemalli-marko-peltojoki>