

## Sukijan asennus- ja käyttöohje

Sukija on Javalla kirjoitettu Solr:n liitännäinen suomenkielisten tekstien indeksointiin.

Sukija analysoi sanat morfologisesti, muuttaa sanat perusmuotoon (joka on sanakirjoissa) ja indeksoi perusmuodot, jotta sanan kaikki taivutusmuodot löytyvät vain perusmuotoa etsimällä.

Sukija tallentaa perusmuodot Solr:n tietokantaan, josta niitä voi etsiä Solr:n käyttöliittymän kautta.

Sukija osaa indeksoida kaikkia niitä tiedostomuotoja, joita Apache Tika (<http://tika.apache.org/>) osaa lukea.

### Mitä tarvitaan ja mistä ne saa?

- Sukija: <https://github.com/ahomansikka/sukija> Koska luet tätä tekstiä, olet jo imuroinut tämän. (-:
- Corevoikko: <https://github.com/voikko/corevoikko> Tarvitaan versiot libvoikko.so.1.14.5 ja libvoikko-4.1.1.jar
- Solr 7.1.0: <http://lucene.apache.org/solr/> Tässä dokumentissa Solr:sta käytetään nimeä solr-x.y.z, missä x.y.z tarkoittaa version numeroa.

Lisäksi Sukija tarvitsee erinäisiä jar-tiedostoja, mutta Maven imuroi ne verkosta automaagisesti.

Tämä asennusohje olettaa, että libvoikko on asennettu hakemistoon `/usr/local/lib` ja että Solr on hakemistossa `$HOME/Lataukset/solr/solr-x.y.z`

Jos ne ovat jossain muualla, tiedostossa `sukija.properties` olevia Libvoikon tietoja pitää muuttaa vastaavasti.

### Solr:n asentaminen

Solr:ia ei tarvitse asentaa, vaan sitä voi käyttää suoraan hakemistosta `$HOME/Lataukset/solr/solr-x.y.z`.

Tuotantokäyttöä varten Solr on parasta asentaa palveluna (service). Tämä ohje olettaa, että Solr on asennettu palveluna.

Solr asennetaan palveluna skriptillä `install_solr_service.sh` (katso Apache Solr Reference Guide kohta Service Installation Script).

Jos et malta katsoa, se tehdään hakemistossa `$HOME/Lataukset/solr/` tällä tavalla

```
sudo bash solr-x.y.z/bin/install_solr_service.sh solr-x.y.z.tgz -u username
```

missä `username` on oma käyttäjätunnuksesi. Jos kohdan `-u username` jättää pois, Solr asennetaan tunnuksen `solr` alle.

Skripti asentaa Solr:n hakemistoihin `/opt/solr` ja `/var/solr`.

## Sukijan kääntäminen ja asentaminen

Sukija käännetään ja asennetaan komennolla

```
mvn install
```

Komento imuroi netistä tarvitsemansa Javan jar-paketit eli ensimmäinen kääntäminen saattaa kestää kauan. Erityisen kauan se kestää, jos et ole aiemmin käyttänyt mavenia.

Sukijan jar-tiedosto asennetaan maven-hakemistoon `${HOME}/.m2`

## Sukijan konfigurointi

Sukijan parametrit asetetaan tiedostossa `sukija.properties`, jossa on myös tarvittavat ohjeet.

Sen jälkeen komennolla `make` tiedostot tehdään Solr:n asetustiedostot `conf/data-config.xml` ja `conf/schema.xml`. Esimerkit näistä tiedostoista ovat hakemistossa `src/main/resources/`.

Seuraavaksi pitää tehdä Solr:iin Sukijan tarvitsemat hakemistot ja kopioida niihin Solr:n tarvitsemat tiedostot. Sitä varten on kaksi skriptiä

```
src/main/resources/install.sh
src/main/resources/install-sudo.sh
```

Jälkimmäistä käytetään, jos Solr on asennettu tunnuksella `solr`.

Jos asentamisen jälkeen konfiguroi Sukijaa, muutokset voi viedä Solr:iin jommallakummalla komennoista

```
cp -r conf/* /var/solr/data/sukija/conf/  
sudo -u solr cp -r conf/* /var/solr/data/sukija/conf/
```

Seuraavaksi mennään selaimella osoitteeseen

<http://localhost:8983/solr/> Vasemmalla paneelin alalaidassa lukee

```
No cores available  
Go and create one
```

Mennään sinne ja kirjoitetaan kohtaan `name` ja `instanceDir` ”sukija” (mutta ilman lainausmerkkejä) ja napsautetaan kohtaa `Add Core`.

Nyt asemalle paneeliin alalaitaan pitäisi tulla mahdollisuus valita indeksi (Solr käyttää siitä nimeä `core`) sukija. Sen jälkeen voi ruveta indeksoimaan. Katso sivu 3.

## Solr:n loki

Solr:n lokitulostus (<http://wiki.apache.org/solr/SolrLogging>) konfiguroidaan tiedostossa `/var/solr/log4j.properties`.

Mahdollisimman suuren lokitulostuksen saa lisäämällä tämän tiedoston loppuun tiedoston `src/main/resources/log4j.properties.sukija`.

Tällöin tulostus on paljon suurempi kuin indeksoitavat tiedostot (-:, mutta kaikkia ei tietenkään tarvitse lisätä.

Lokitulostuksen eri tasot (`ALL`, jne) voi katsua luokan `org.apache.log4j.Level` dokumentoinnista.

Lokitulostus menee kansioon `/var/solr/logs/`

## Indeksointi

Tiedostot indeksoidaan menemällä osoitteeseen

<http://localhost:8983/solr/sukija/dataimport?command=full-import>

Indeksoinnin voi aloittaa myös Solr:n käyttöliittymästä.

Enemmän tai vähemmän pitkän ajan kuluttua indeksoinnin lopputuloksen voi katsoa osoitteesta <http://localhost:8983/solr/sukija/dataimport> tai Solr:n käyttöliittymästä.

## Tietojen etsiminen

Sanoja etsitään menemällä osoitteeseen

<http://localhost:8983/solr/sukija/browse>

Etsittävien sanojen tulee olla perusmuodossa. Etsittäessä sanoja ei muuteta perusmuotoon siksi, että yhden sanan perusmuoto voi olla toisen sanan taivutusmuoto. Paras esimerkki tästä on "alusta", joka on sanojen "alusta", "alustaa", "alku", "alunen" ja "alus" taivutusmuoto. Tällöin herää kysymys, mitä sanaa pitää etsiä, vai etsitäänkö kaikkia?

Eri tavalla muotoillun tulostuksen saa osoitteesta

<http://localhost:8983/solr/sukija/select>

Esimerkiksi sanaa **sana** etsitään näin:

<http://localhost:8983/solr/sukija/select?q=sana>

Tämän tulostuksen ulkonäköä voi muuttaa muuttamalla Sukijan mukana tulevaa tiedostoa `conf/xslt/sukija.xsl`.

Tietoja voi etsiä myös Solr:n käyttöliittymällä.

## Tiedoston suggestions.xml konfigurointi

Tässä vaiheessa kaikki indeksoitavista tiedostoista luetut sanat on muutettu pieniksi kirjaimiksi eli tiedostossa `suggestions.xml` olevat erisnimetkin pitää kirjoittaa pienellä alkukirjaimella.

Konfiguroititiedoston formaatti on määritelty tiedostossa `/src/main/xsd/SuggestionInput.xsd`.

Konfigurointitiedostossa olevien säännöllisten lausekkeiden syntaksi on sama kuin Javan luokassa `java.util.regex.Pattern`

Konfiguroinnille ei ole käyttöliittymää, koska tiedoston `suggestions.xml` muokkaamiseen voi käyttää mitä tahansa XML-editoria.

`compoundWordEnd` tunnistaa yhdyssanan, jos se loppuu tiettyyn sanaan. Esimerkiksi

```
<compoundWordEnd>
  <input>joki joki</input>
  <input>joke joki</input>
  <input>joe joki</input>
</compoundWordEnd>
```

tunnistaa esimerkiksi merkkijonon ”aatsajoelle”. Tällä tavalla voidaan tunnistaa paikannimiä, jotka eivät ole sanastossa.

Jokaisessa input-lauseessa on kaksi osaa. Ensimmäinen on jokin merkkijono ja toinen jonkin sanan perusmuoto. Tunnistettaessa merkkijono katkaistaan siitä kohdasta, josta ensimmäinen merkkijono alkaa, ja jos merkkijonon loppuosan perusmuoto on argumentin toinen osa, perusmuotona palauteaan merkkijonon alkuosa + argumentin toinen osa.

Esimeriksi ”aatsajoelle” jaetaan kahtia osiin ”aatsa” ja ”joelle”, ja koska merkkijonon ”joelle” perusmuoto on ”joki”, merkkijonon ”aatsajoelle” perusmuodoksi tulee ”aatsajoki”.

Input-lauseita voi olla mielivaltainen määrä.

Ensimmäinen merkkijono ei ole säännöllinen lauseke, koska ne ovat hitaampia kuin merkkijonot varsinkin, kun merkkijonojen tunnistuksessa käytetään Aho-Corasick -algoritmia.<sup>1</sup> Algoritmissa sanan (ensimmäisen merkkijonon) etsimiseen käytetty aika ei riipu etsittävien sanojen määrästä, sillä se etsii kaikkia sanoja samanaikaisesti.

`prefix` tunnistaa etuliitteettömän sanan (”etuliite” voi olla mikä tahansa merkkijono). Esimerkiksi

```
<prefix>
  <prefix>abcdefg</prefix>
  <savePrefix>true</savePrefix>
  <saveWord>true</saveWord>
</prefix>
```

---

<sup>1</sup>Katso Alfred V. Aho and Margaret J. Corasick: *Efficient String Matching: An Aid to Bibliographic Search*. Communications of the ACM. June 1975 Volume 18 Number 6. Löytyy netistä googlaamalla.

poistaa sanan alusta merkkijohon "abcdefg" ja yrittää tunnistaa jäljelle jääneen sanan ("abcdefgsuomalaiselle" => "suomalaiselle") ja tallentaa sen perusmuodon ("suomalainen"). Jos `savePrefix` on true, tallentaa myös etuliitteen ("abcdefg") ja jos `saveWord` on true, tallentaa myös koko sanan perusmuodon ("'abcdefgsuomalainen").

Prefix-lauseita voi olla mielivaltainen määrä.

`char` muuttaa sanassa olevat merkit toiseksi. Tämä vastaa Unixin komentoa `tr`. Esimerkiksi

```
<char>
  <from>gbdkptvw</from>
  <to>kptgbdwv</to>
</char>
```

muuttaa g:n k:ksi, b:n p:ksi jne. Ohjelma testaa muutettavien kirjainten kaikki mahdolliset kombinaatiot. Esimerkiksi jos tiedostosta luettu sana on "piolokia", komento yrittää tunnistaa sanat "piolokia", "biolokia", "piologia" ja "biologia" (mutta ei välttämättä tässä järjestyksessä :-).

`regex` muuttaa säännöllisen lausekkeen. Esimerkiksi

```
<regex>
  <input>(ai)(j)([eou])  $1$3</input>
  <input>^([0-9]+)</input>
  <tryAll>true</tryAll>
</regex>
```

Ensimmäinen input-lause poistaa j-kirjaimen muun muassa sanoista "aijemmin", "aijomme" ja "kaijutin", ja toinen poistaa numerot sanan alusta.

Säännöllisessä lausekkeessa voi käyttää seuraavia lyhenteitä:

```
%A  [ää]
%C  [bcd fghjklmnpqrsšt vwxzž]
%O  [oö]
%U  [uy]
%V  [aeiouyäö]
%%  %
```

Esimerkiksi `<input>(%V)(h)(%V) $1hd$3</input>` muuttaa h:n hd:ksi esimerkiksi sanassa ”puhistus”.

Jos `tryAll` on `true`, ohjelma kokeilee kaikkia säännöllisiä lausekkeita, jos se on `false`, ohjelma lopettaa ensimmäisen tunnistetun sanan jälkeen.

Input-lauseen ensimmäinen osa on säännöllinen lauseke ja toinen merkkijono, miksi se muutetaan. Sen syntaksi on sama kuin Javan luokassa `java.util.regex.Matcher`.

Jos toista osaa ei ole, säännöllisen lausekkeen tunnistama merkkijono poistetaan. Input-lauseita voi olla mielivaltainen määrä.

`start` käy läpi kaikki sanan alut pisimmästä (`maxLength`) alkaen lyhimpään (`minLength`) asti ja lopettaa, kun löytyy ensimmäinen tunnistettu sana

```
<start>
  <minLength>4</minLength>
  <maxLength>10</maxLength>
  <baseFormOnly>true</baseFormOnly>
</start>
```

Jos `baseFormOnly` on `true` palautetaan sana vain, jos se on perusmuodossa (esim. ”autowerwww” palauttaa ”auto”), muuten palautetaan tunnistettu sana muutettuna perusmuotoon (”kuudenwww” palauttaa ”kuusi”, mutta jos `baseFormOnly` on `false`, ei palauteta mitään).

`apostrophe` poistaa sanasta heittomerkin ja yrittää tunnistaa sanan sen jälkeen. Jos tunnistaminen ei onnistu, poistaa sanasta heittomerkin ja kaikki sen jälkeiset merkit ja palauttaa jäljelle jääneet merkit sanan perusmuotona. Esimerkiksi yrittää tunnistaa merkkijonon `centime'in` muodossa `centimein`. Jos tunnistaminen ei onnistu, palauttaa merkkijonon `centime`.

```
<apostrophe/>
```

Tällä komennolla ei ole parametreja.