

Evaluation

Bastian Bolender

Pubmed&Medline

Medline ist eine öffentliche(engl.) Datenbank mit Texten aus dem Bereich Medizin

Pubmed ist der öffentliche Zugang zu Medline

z. T. von Fachleuten manuell annotiert

....unsere Referenzmenge....

Eingangsformat

```
<doc id="15971344">
<mwt name="Research Support" textPos="33 36">
<nr>H01.770.727</nr>
<nr>N03.219.483.645</nr>
</mwt>
</doc>
```

```
<doc id="15971340">
<mwt name="Facility Design and Construction" textPos="7 4">
<nr>J01.086.339</nr>
<nr>N02.278.200</nr>
</mwt>
</doc>
```

```
<Pagination>
  <MedlinePgn>115-7</MedlinePgn>
</Pagination>
<Abstract>
  <AbstractText>We report on a 30-year-old man with metastatic non seminomatous
  germ cell tumor of the left testicle to the abdomen and the lungs, who suddenly developed a
  bilateral spontaneous pneumothorax after the first course of salvage chemotherapy. Rapid
  destruction and lysis of lung nodules by chemotherapy seem to be the main mechanism of
  pneumothorax development. According to our case report and to the literature, the onset of acute
  dyspnea after chemotherapy for lung metastatic germ cell tumor should alert to the possibility
  of spontaneous pneumothorax.</AbstractText>
</Abstract>
  <Affiliation>Department of Hematology-Oncology, H.√¥tel-Dieu de France, University
  Hospital, Beirut, Lebanon.</Affiliation>
  <AuthorList CompleteYN="Y">
    <Author ValidYN="Y">
      <LastName>Loutfi</LastName>
      <ForeName>Rania</ForeName>
      <Initials>R</Initials>
    </Author>
  </AuthorList>
  <Language>eng</Language>
  <PublicationTypeList>
    <PublicationType>Case Reports</PublicationType>
    <PublicationType>Journal Article</PublicationType>
  </PublicationTypeList>
</Article>
<CitationSubset>IM</CitationSubset>
<MeshHeadingList>
  <MeshHeading>
    <DescriptorName MajorTopicYN="N">Adult</DescriptorName>
  </MeshHeading>
  <MeshHeading>
    <DescriptorName MajorTopicYN="N">Antineoplastic Combined Chemotherapy
  Protocols</DescriptorName>
    <QualifierName MajorTopicYN="Y">adverse effects</QualifierName>
  </MeshHeading>
```

Tools: Format

Bringt die Ergebnisse

- Bei Jannik und Sebastian "mwt"
- bei Temis "Radiation s=23 l=9",
- bei MA "<DescriptorName" und "<QualifierName"

in das Standardformat

15884692, Testicular Neoplasms, MA.

15884692, Tomography, X-Ray Computed, MA.

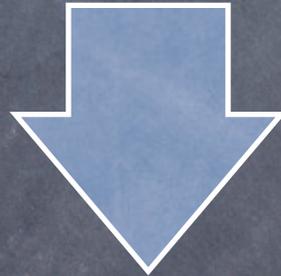
15884691, Amphotericin B, MA.

15884691, therapeutic use, MA.

Tools: Double

Doppelte Einträge werden aus Liste entfernt

15948488, Time Management, TE.
15941170, Fractures, Closed, TE.
15941022, Disability Evaluation, TE.
15941022, Disability Evaluation, TE.
15941161, Hip Dislocation, Congenital, TE.



15948488, Time Management, TE.
15941170, Fractures, Closed, TE.
15941022, Disability Evaluation, TE.
15941161, Hip Dislocation, Congenital, TE.

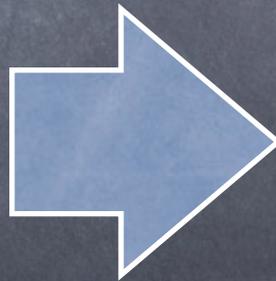
Tool: Reduce

Notwendig, da nicht alle Dateien aus einem Jahr mit dem gleichen Abstract anfangen.

Stellt dies fest und entfernt eine bestimmte Anzahl von Zeilen am Anfang von MA, damit bei allen Dateien der gleiche Abstract vorne steht.

16206981, Adolescent, MA.
16206982, Bone Neoplasms, MA.
16206982, secondary, MA.
16206983, Dopamine Agonists, MA.

16206982, Bone Neoplasms, JA.
16206982, secondary, JA.
16206983, Dopamine Agonists, JA.



16206982, Bone Neoplasms, MA.
16206982, secondary, MA.
16206983, Dopamine Agonists, MA.

16206982, Bone Neoplasms, JA.
16206982, secondary, JA.
16206983, Dopamine Agonists, JA.

Tool: Eval

Sortiert die X-on Headings von Jannik und Sebastian in fünf Dateien

Artikel in denen:

- nur Jannik etwas gefunden hat
- nur Sebastian etwas gefunden hat
- beide die gleichen Entryterms gefunden haben
- oder beide verschiedene Entryterms gefunden haben (2)

Schwierige Evaluierung I

Die Übereinstimmung unter manuellen
Annotierern liegt unter 50%

Richtiger Treffer für uns nur wenn:
mwt im Text gefunden wird, das so nicht in Mesh
steht und auch im Text nichtmehr so vorkommt,
dass Mesh es finden kann

Schwierige Evaluierung II

MeSH Heading	Lung Neoplasms
Tree Number	C04.588.894.797.520
Tree Number	C08.381.540
Tree Number	C08.785.520
Annotation	coord IM with histol type of neopl (IM)
Scope Note	Tumors or cancer of the LUNG .
Entry Term	Cancer of Lung
Entry Term	Lung Cancer
Entry Term	Pulmonary Cancer
Entry Term	Pulmonary Neoplasms
Entry Term	Cancer of the Lung
Entry Term	Neoplasms, Lung
Entry Term	Neoplasms, Pulmonary
See Also	Carcinoma, Non-Small-Cell Lung
See Also	Carcinoma, Small Cell
Allowable Qualifiers	BL BS CF CH CI CL CN CO DH DI DT
Entry Version	LUNG NEOPL
Unique ID	D008175

Tools: Vergleich I

(Das wichtigste)

Liest die Ergebnisse in Hashes ein und gleicht sie untereinander ab

Berechnet Precision & Recall

Der Output sieht dann etwa so aus:

Zahl der untersuchten Abstracts: 10000

Jannik findet 358 Treffer, die MA auch finden aber Temis nicht

Sebastian findet 38170 Treffer die Temis findet

Jannik, Sebastian&Temis finden 35970 Treffer die auch MA finden

MA finden 136830 Treffer die Jannik&Temis nicht finden

MA finden 130179 Treffer die Sebastian&Temis nicht finden

Tools: Vergleich II

(Precision & Recall)

Im Information Retrieval am häufigsten verwendete Maße zur Beschreibung der Güte von Suchergebnissen

Precision

Beschreibt die Genauigkeit eines Suchergebnisses

Definiert als der Anteil der gefundenen relevanten Dokumente von allen bei einer Suche gefundenen Dokumenten

Recall

Beschreibt die Vollständigkeit eines Suchergebnisses

Definiert als der Anteil der bei einer Suche gefundenen relevanten Dokumente (bzw. Datensätze) an den relevanten Dokumenten der Grundgesamtheit

Tools: Vergleich III

(Precision & Recall: Formeln)

Precision = (richtige Treffer) / (richtige Treffer + false positives)

Recall = (richtige Treffer) / (richtige Treffer + false negatives)

Die Referenzmenge wird vom Output der manuellen Annotierer gebildet

Precision = $(TE=MA) / (TE=MA + [TE-(TE=MA)])$

Recall = $(TE=MA) / (TE=MA) + [MA-(TE=MA)]$

Für gemeinsame Werte wird TE durch TEJA ersetzt

Maximaler Wert: 1 (Übereinstimmung mit der Referenzmenge)

Tools: Vergleich IV

(Precision&Recall: Berechnung)

Für 10.000 Abstracts aus dem Jahr 2004

Temis: Precision: $28.846 / (28.846 + 158.030) = 0,154359$

Recall: $28.846 / (28.846 + 137.176) = 0,173748$

Temis&Jannik: Precision: $29.204 / (29.204 + 158.165) = 0,155864$

Recall: $29.204 / (29.204 + 136.818) = 0,175904$

Die Dateien von Temis und MA beinhalten alle Treffer, unsere nur mwts

Mit anderen Worten:

Jannik findet 2.261 Treffer

Temis findet insgesamt 186.876 Treffer

MA finden insgesamt 166.022 Treffer

Jannik findet 590 Treffer (mwts), die auch MA finden

Jannik findet 358 Treffer (mwts), die MA auch finden aber TE nicht

Erwartungen für statistische Methode:

Recall steigt an, Precision wird geringer

Ergebnisse I

für 10.000 Abstracts

Treffer in 10.000 Abstracts: Jannik: 2261, Sebastian: 60.827, Temis: 186.803, MA: 166.055

Temis findet 28851 Treffer die auch MA finden

Jannik findet 590 Treffer die auch MA finden

Jannik findet 358 Treffer, die MA auch finden aber Temis nicht

Jannik findet 1768 Treffer die Temis findet

Sebastian findet 15176 Treffer die auch MA finden

Sebastian findet 7009 Treffer, die MA auch finden aber Temis nicht

Sebastian findet 38170 Treffer die Temis findet

Jannik&Temis finden 29209 Treffer die auch MA finden

Sebastian&Temis finden 35860 Treffer die auch MA finden

Jannik, Sebastian&Temis finden 35970 Treffer die auch MA finden

MA finden 136830 Treffer die Jannik&Temis nicht finden

MA finden 130179 Treffer die Sebastian&Temis nicht finden

Ergebnisse II

für 10.000 Abstracts

Precision Temis: 0.154385795928851

Recall Temis: 0.173743639155701

Precision Jannik&Temis: 0.154896564158858

Recall Jannik&Temis: 0.175899551353467

Precision Sebastian&Temis: 0.159396905406427

Recall Sebastian&Temis: 0.215952545843245

Precision Jannik, Sebastian&Temis: 0.159373670778392

Recall Jannik, Sebastian&Temis: 0.216614976965463

Danke fürs Zuhören