

HEPMED: data mining in **High Energy Physics and MED**incine



GIORNATA DI INCONTRO BORSE DI STUDIO GARR "ORIO CARLINI" ROMA

Roma, 27 Giugno2019

Borsisti Day 2019





Sviluppo di un modo efficiente di navigare grandi quantità di dati sfruttando tecniche di data warehousing, con applicazione diretta su database di fisica delle alte energie e di ambito medico

Abbiamo a disposizione dati di misure fatte dagli anni '70' ad oggi potrebbe esserci sfuggito qualcosa?



Correlando dati di esami medici diversi, potrei forse fare diagnosi più sicure?/







Creare uno strumento per rispondere alla domanda:



I dati sono consistenti con la teoria XXX?



Di cosa abbiamo bisogno?:

- 1)Un modo efficace per navigare i dati
 - Organizzare i dati già disponibili ed estrarre automaticamente quante più informazioni possibili tramite un interfaccia di supporto

Quali sono le misure rilevanti per la teoria da testare?

Vorrei confrontare dati di analisi diverse su un unico grafico...

- 2)Un tool per ottenere previsioni teoriche delle quantità misurate
- 3)Uno **strumento statistico** per quantificare l'accordo **tra** dati e teoria

Consortium

CONSOR

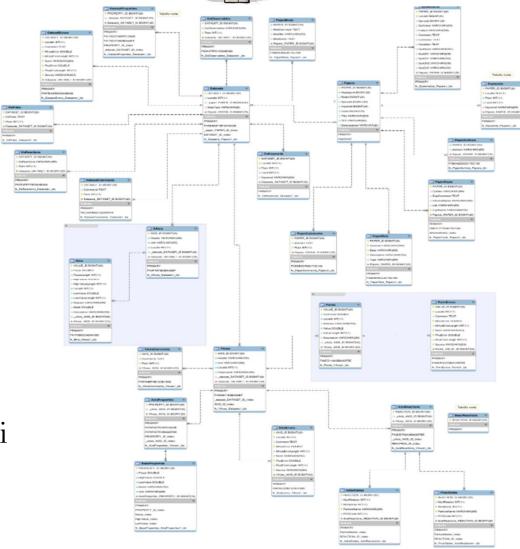


Start-point 1: database hepdata

- Database relazionale con struttura complessa
- Elemento base della ricerca è l'articolo scientifico

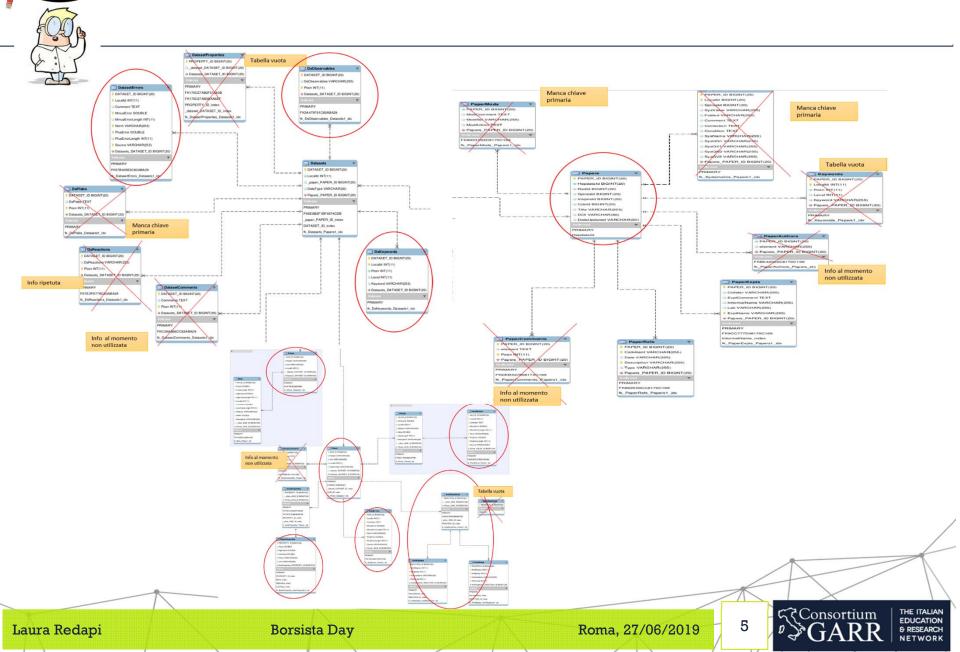
Ricerco tutti quei papers che...

- Possibilità di navigare i dati all'interno di una singola analisi selezionata



Laura Redapi



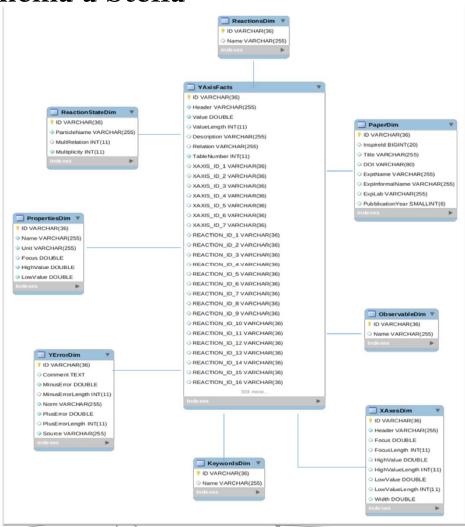






Il Nuovo Database MineHep

Schema a Stella



YAxisFact → dato numerico della singola misura

Le **dimension**i identificate sono :

- Properties
 - YError
 - XAxes
- Observable
 - Paper
- Keywords
- Reaction

Ogni dimensione è una tabella con tante colonne quante sono le caratteristiche necessarie ad identificarla.

Consortium | THE ITALIAN EDUCATION & RESEARCH NETWORK



-- RESULT: OK

GIORNATA DI INCONTRO BORSE DI STUDIO GARR "ORIO CARLINI" BORSISTI DAY 2019

PT(E)

> 35 GEV



Test di controllo



-- TEST: prendo tutte le misurazioni dela tabella 1 che sono comprese tra [0.8, 2.0] \$|\eta|\$ nella colonna HERAPDF select distinct Value from YAxisFacts yaf inner join PaperDim ppd on ppd.ID = yaf.PAPER_ID inner join XAxesDim xad on xad.ID = yaf.XAXIS_ID_1 where InspireId = 1118047 and TableNumber = 1 and xad.Header = '\$|\\eta|\$' and yaf.Header = "HERAPDF" and xad.LowValue >= 0.8 and xad. HighValue <= 2.0; # Value # 196 # 181 # 153 # 140 # 132 -- RESULT: OK

-- TEST: prendo tutte le misurazioni dela tabella 1 della colonna \$\mathcal{A}\$\$ che hanno un errore SYS +- 6

-- TEST: prendo tutte le misurazioni dela tabella 1 della deselect yaf.Value
from YAxisFacts yaf
inner join PaperDim ppd on ppd.ID = yaf.PAPER_ID
inner join YErrorDim yed on yed.ID = yaf.POINT_ERROR_ID_1
where InspireId = 1118047
and TableNumber = 1
and Source = 'SYS'
and Header= 'S\\mathcal{A}\$'
and (MinusError=-6 and PlusError=6);
Value
156
136

RE	P P> W+- < E+- NUE > X 7000.0 GeV				
SQRT(S)					
$ \eta $	\mathcal{A}	CT10	HERAPDF	MSTW	NNPDF
0.0 - 0.2	102.0 ±3.0 stat ±5.0 sys	109.0 ±5.0	106.0 +4.0	87.0 +3.0	107.0 ±5.0
0.2 - 0.4	111.0 ±3.0 stat ±5.0 sys	114.0 ±5.0	110.0 +4.0	89.0 +3.0	110.0 ±5.0
0.4 - 0.6	116.0 ±3.0 stat ±5.0 sys	119.0 ±5.0	115.0 +4.0	98.0 +3.0	116.0 ±5.0
0.6 - 0.8	123.0 ±3.0 stat ±5.0 sys	126.0 ±5.0	122.0 +4.0	103.0 +3.0	123.0 ±5.0
0.8 - 1.0	133.0 ±3.0 stat ±5.0 sys	138.0 +5.0	132.0 +4.0	115.0 +4.0	134.0 ±5.0
1.0 - 1.2	136.0 ±3.0 stat ±6.0 sys	146.0 ±6.0	140.0 +5.0	128.0 +4.0	145.0 ±5.0
1.2 - 1.4	156.0 ±3.0 stat ±6.0 sys	164.0 +6.0	153.0 *5.0	144.0 ±5.0	158.0 ±5.0
1.6 - 1.8	166.0 ±3.0 stat ±10.0 sys	195.0 *8.0	181.0 ±5.0	179.0 ±5.0	190.0 ±4.0
1.8 - 2.0	197.0 ±3.0 stat ±9.0 sys	207.0 +8.0	196.0 +4.0	200.0 +6.0	206.0 ±4.0
2.0 - 2.2	224.0 ±3.0 stat ±11.0 sys	224.0 +8.0	211.0 +5.0	213.0 +6.0	219.0 ±4.0
2.2 - 2.4	210.0 ±4.0 stat ±13.0 sys	241.0 +8.0	225.0 +9.0	231.0 +6.0	231.0 ±5.0

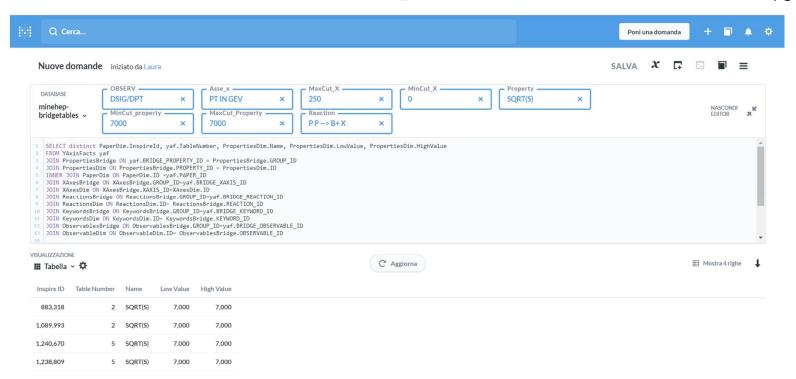
Consortium | THE ITALIAN EDUCATION & RESEARCH NETWORK





State of the art: MINEHEP + METABASE

Metabase, strumento opensource, per interrogare ed estrarre conoscenza dai dati in modo semplice tramite un'interfaccia SQL



Consortium THE ITALIAN EDUCATION & RESEARCH NETWORK

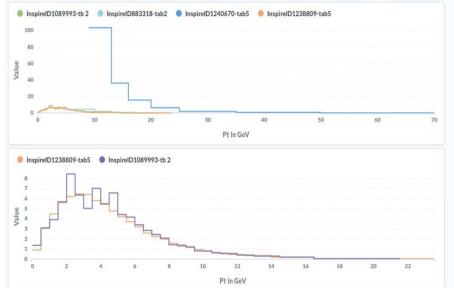




MINEHEP + METABASE







Consortium THE ITALIAN EDUCATION 8 RESEARCH NETWORK





Risultati ottenuti:



- Elemento principale della ricerca è il dato della singola analisi
- Possibilità di fare query articolate

Cercare tutti i dati che rispondono a certe keywords e a tagli su diverse proprietà scelte dall'utente

- Possibilità di navigare su più analisi contemporaneamente:

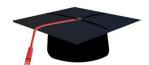
Osservo tutte le tabelle di interesse contemporaneamente e decido quali dati, anche di esperimenti diversi, riportare in uno o più grafici









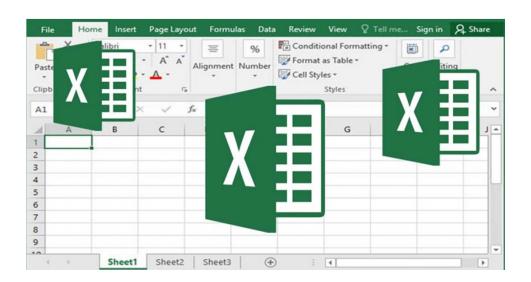




Start-point 2: database fisica medica



- Verifiche pretrattamento della radioterapia dell'Azienda Ospedaliera Universitaria di Careggi
- Risultati analisi gamma
- Fogli Excel





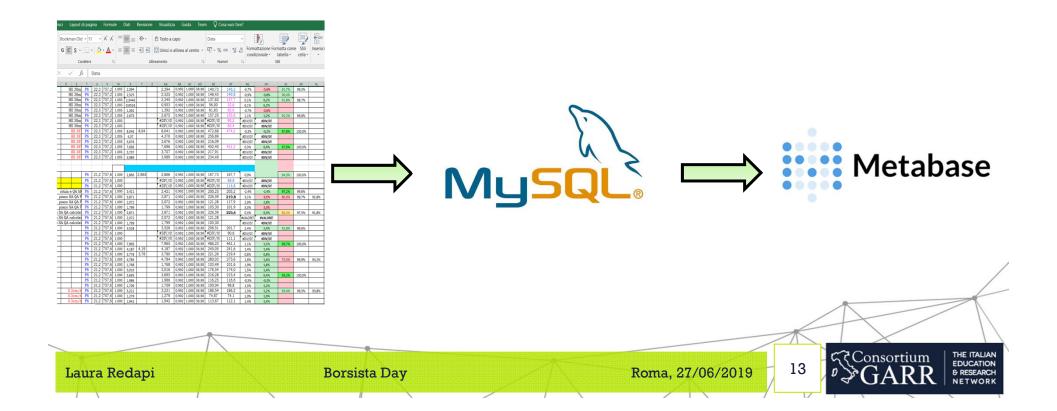




Work in progress



- Riorganizzazione del database in modo che sia interfacciabile al software di data mining METABASE







Risultati ottenuti:



- Raccolta e riorganizzazione dati in un unico foglio Excel
 - Progettazione diagramma a stella
 - Creazione database di test in MySQL

Consortium

GARR

THE ITALIAN EDUCATION 9 RESEARCH NETWORK







- Caricamento nuovo database su Metabase
 - Creazione query di interesse



Consortium | THE ITALIAI EDUCATION & RESEARCH NETWORK

15





Vitaliano Ciulli Pier Giulio Lenzi Andrea Ceccarelli Andrea Cioni Maria Vittoria Garzelli



Fisica Medica della Radioterapia di Careggi





Consortium THE ITALIAN EDUCATION & RESEARCH NETWORK