



Universitat Autònoma de Barcelona

FACULTAT DE CIÈNCIES

PRÀCTICA MONGODB

Bases de dades no relacionals

Sergi Cantón Simó - 1569251

Bernat Espinet Torrecassana - 1564342

Marc Llopart Enajas - 1569054

Gerard Vinyes Sanchez - 1563545

1/4/2022

Índex

1	Creació de les col·leccions	2
2	Creació del fitxer de python per generar la base de dades	2
3	Joc de proves	3
3.1	<i>Escàners diferents que hi ha a la BD. Mostra el device</i>	3
3.2	<i>Número total de nòduls que s'han utilitzat per l'entrenament (train=1) de l'experiment 1 del mètode "Method2"</i>	3
3.3	<i>Valor màxim, mínim i mitjà de BenignPrec agrupat per classificador (classifier). Mostra ID del mètode, MaxBenignPrec, MinBenignPrec, AvgBenignPrec.</i>	3
3.4	<i>Número total d'homes i dones. Mostra sexe i número total</i>	3
3.5	<i>Pacients amb més de dos nòduls. Mostra ID del Pacient, sexe, edat, diagnòstic del Pacient</i>	3
3.6	<i>Mostrar els 4 mètodes amb més repeticions de l'experiment. Mostra el ID del Mètode i número de repeticions de l'experiment</i>	3
3.7	<i>Per cada pacient els escàners (CTs) que s'ha fet. Mostra el ID del Pacient, device i la data del CT</i>	3
3.8	<i>Mostrar els pacients que tenen tots els seus nòduls amb diagnosis = "Benign" i el seu recompte</i>	3
3.9	<i>Modificar la ResolutionTV augmentant-la un 20 % dels escàners que es van realitzar amb DataCT = 18/11/2018</i>	4
4	Repartiment de tasques	4

1 Creació de les col·leccions

Per a la nostra base de dades, vam considerar 5 col·leccions: Nodule, Patient, CTScanner, Experiment i Relation.

Partint de la classe Nodule, inicialment vam relacionar-la amb Patient afegint l'atribut PatientID dins de Nodule. A continuació, vam també referenciar CTScanner dins de Nodule col·locant l'atribut CTID dins de la col·lecció Nodule. A més, però, vam afegir també l'atribut Diameter, de la relació entre Nodule i CTScanner, i els atributs Device i DataCT de la classe CTScanner, dins de Nodule. Vam prendre aquesta decisió perquè vam observar que aquests tres atributs s'han de fer servir sovint a les consultes que posteriorment se'ns demanen. La classe Method i Experiment les vam relacionar encastant tots els atributs de Method dins d'Experiment perquè Method no és una col·lecció massa gran i, basant-nos en les consultes, ens era el més òptim organitzar-ho així. Seguidament, Nodule i Experiment els vam relacionar referenciant Nodule dins d'Experiment afegint l'atribut NoduleID dins d'Experiment. També, vam afegir els atributs de la relació entre Nodule i Experiment, Diagnosis i Train, dins d'Experiment. Finalment, Relation ens relaciona les claus primàries de Patient, Nodule, Method i, també, inclou els atributs ExperimentRepetition, Train i RadiomicsDiagnosis. Relation és una col·lecció artificial que ens facilita dur a terme les consultes.

2 Creació del fitxer de python per generar la base de dades

Dins l'arxiu de python, vam establir la connexió amb el servidor local, vam crear la base de dades amb les seves col·leccions i vam llegir i insertar-hi les dades.

La base de dades que vam crear s'anomena *MongoDB-BDNR*. Després de generar-la, ens va fer falta crear concretament les col·leccions que hem especificat a l'apartat 1. Per crear cada col·lecció, es fa un tractament de l'error per si la col·lecció ja està creada. En cas afirmatiu, s'esborra la col·lecció i es torna a crear. Això ens permet evitar que s'insereixin valors duplicats dins la base de dades si hi ha qualsevol error d'execució, la base de dades es queda a mitja creació i s'ha de tornar a executar el fitxer de python.

A continuació, es genera un objecte de la classe *Options* que ens aporta les opcions necessàries per llegir el fitxer de dades, esborrar el contingut de les nostres bases de dades, etc.

Posteriorment, dins del fitxer de python es llegeixen les dades des de Excel fent servir una funció de la llibreria *Pandas* i s'afegeixen aquestes dades dins de les col·leccions que correspongui. En els casos on vam decidir referenciar, únicament ens va caler posar l'identificador d'una col·lecció dins de l'altre. Pels casos on vam decidir encastar, inicialment vam haver de crear el diccionari amb els atributs a encastar per, posteriorment, afegir aquest diccionari dins del diccionari corresponent a la col·lecció on volíem encastar.

Per acabar, vam modificar el fitxer *Options* per afegir la opció que esborri el contingut de les bases de dades al crear un objecte de la classe *Options*.

3 Joc de proves

L'última part de la pràctica va consistir en dur a terme una sèrie de consultes de la base de dades.

3.1 *Escàners diferents que hi ha a la BD. Mostra el device*

Resolem la consulta fent un *distinct* dels *Devices* de la base de dades a la col·lecció CTScanner.

3.2 *Número total de nòduls que s'han utilitzat per l'entrenament (train=1) de l'experiment 1 del mètode "Method2"*

A la col·lecció Relation, apliquem una pipeline (*aggregate*), dins de la qual seleccionem els atributs que ens interessin i fem un *count* del nombre de nòduls .

3.3 *Valor màxim, mínim i mitjà de BenignPrec agrupat per classificador (classifier). Mostra ID del mètode, MaxBenignPrec, MinBenignPrec, AvgBenignPrec.*

A la col·lecció Experiment, apliquem una pipeline (*aggregate*), dins de la qual separem l'array on guardem les dades del mètode.

3.4 *Número total d'homes i dones. Mostra sexe i número total*

A la col·lecció Patient, apliquem una pipeline (*aggregate*), dins de la qual agrupem per gènere i fem el recompte.

3.5 *Pacients amb més de dos nòduls. Mostra ID del Pacient, sexe, edat, diagnòstic del Pacient*

A la col·lecció Nodule, apliquem una pipeline (*aggregate*), dins de la qual fem un left join agrupant per la id del pacient i projectem allò que se'ns demana. A continuació, separem pel camp *info_pacient*, que ens hem creat nosaltres, i projectem per obtenir els atributs que ens interessin.

3.6 *Mostrar els 4 mètodes amb més repeticions de l'experiment. Mostra el ID del Mètode i número de repeticions de l'experiment*

A la col·lecció Experiment, apliquem una pipeline (*aggregate*), dins de la qual separem el camp *Method*, agrupem per id i comptem quants n'hi ha. Finalment, mostrem només aquells 4 que es repeteixen més cops.

3.7 *Per cada pacient els escàners (CTs) que s'ha fet. Mostra el ID del Pacient, device i la data del CT*

A la col·lecció Nodule, apliquem una pipeline (*aggregate*), dins de la qual fem un left join de Nodule amb CTScanner, per tenir les dades del scanner als documents dels nòduls. A continuació, separem els arrays de la informació que acabem d'afegir amb el join i agrupem per aquella informació que se'ns demana.

3.8 *Mostrar els pacients que tenen tots els seus nòduls amb diagnosis = "Benign" i el seu recompte*

A la col·lecció Nodule, apliquem una pipeline (*aggregate*), dins de la qual filtrem per la condició que el nòdul sigui benigne i comptem quants n'hi ha.

3.9 *Modificar la ResolutionTV augmentant-la un 20 % dels escàners que es van realitzar amb DataCT = 18/11/2018*

A la col·lecció CTScanner, fem un *update* per actualitzar la resolució multiplicant-la per 1.2 únicament en els escàners amb la data que se'ns demana.

4 Repartiment de tasques

A l'iniciar la pràctica, vam repartir les tasques equitativament. El script de python el vam fer principalment el Bernat i el Sergi i les consultes de MongoDB les vam fer majoritàriament el Marc i el Gerard. Respecte l'informe, cadascú va escriure allò referent a la part de la pràctica que havia dut a terme.

No obstant, tot i haver repartit les parts de la pràctica, tots hem estat al cas també de les parts que inicialment no ens corresponien i hem ajudat en allò que ha calgut, independentment de si se'ns havia assignat o no.