

CORSO PYTHON

LinkedIn Learning



DATA SCIENCE FOUNDATIONS:

PYTHON SCIENTIFIC STACK

- OBJETIVOS:
- JUPITER - ? ANDAS - foliogr
 - Numpy - Conda - GEO
 - ML con Scikit-learn
 - Plotting con Matplotlib y Bokeh
 - Numba, Cython, Deep learning and NLP.
 - VIRTUAL ENV. Kivy. Selenium
- NumPy → Matrices. Cálculo Numérico
- SciPy: Colección con FFT y Science capabilities.
- Matplotlib → Visualizaciones.
- Pandas → Data Set.
- Scikit-Learn → ML Algorithms
- Jupyter → ? → Ayuda
- ? ? → Ver código fuente
- ? - Pandas → Ver directory

% time → mide el
tiempo de ejecución del

!dir !ls . & Python -v comando
calcul corregido en Jupyter.

Puedes usar Markdown Cells q se ven de la
siguiente forma. En vez de usar
comandos terminal de Python puedes
usar `%%tex`, `%%py` y más.

Numpy

- NO ES LA MEJOR OPCIÓN CUANDO A
USO DE MEMORIA Y PERFORMANCE

`array = np.array([1, 2, 3])`

`array.dtype` → muestra el tipo
de variable de tu array

Puedes usar operadores lógicos y

Slicing Adentro de los arrays!

`array[1: ; 2:]` , and.

`array((array>7) & (array<9))`

array ([1, 2, 3], [4, 5, 6]) → [[1, 2, 3], [4, 5, 6]]

STARG

VIDEK

Bubbleo indexing (lo de fin peg
arriba) es mucho más rápido y
eficiente q usar for loops!

df (vector) → me muestra todo los elementos
q tiene el dicho.

Numpy espera valores en 2dicas p/ trig.

②. np·vectorje (funcion)

decorador de numpy

Pandas

Read-csv (file_name, parse_dates=True)

df['NOMBRE_COLUMNA'] → ACCEDER Cols.

df['LAT'][0]

COL Fila.

df.loc[0] o df.loc[2:7]

devuelve toda la fila de
LA o A LA 7.

df. [['AF', 'long'] [2 : 7]]

OBSERVAR filas 2:7 pero solo

de esas dos col.

df. loc [~~label~~ Name]

↳ Puedes ver las sp.

	Col 1	Col N
index 1	data		data
:			
index N	data		data

funko next steps ⇒ Proyecto de DATA

Personal PI seguir a prendiendo

(↳ keep)

finance

IBII - DATA SCIENCE

- PDF. CERTIFICATE.

CURSO 1 = WHAT IS DATA SCIENCE?

CURSO 1

- SERVIA, 1

5 V's OF Big DATA

- 1) Velocity 3) Variety 5) Value
- 2) Volume 4) VERACITY

- 1) A la que data se acuñaba. muy Rápido) sin paros.
- 2) la escala de los datos) el Acumulado
la cantidad de datos Almacenados.
- 3) diversidad. Estructurados) NO ESTRUCTURADOS
de trags, personas, procesos) ETC.
- 4) Quality and Origin .
- 5) Habilidades) Necesidad de conversión
data en valor. (No solo \$).

② 80% de los datos son No Estructurados.

→ HADOOP: ¿Qué es? → cloud. colm. lot
como un Apache, AWS, etc, para data.

→ **DATA MINING** → Autonomically and programmatically selecting and processing DATA Revealing previously hidden info

1) SET UP GOALS FOR THE EXERCISE.

Identify key questions that need to be ANSWERED. Consider the costs / Beneficios del ejercicio.

determine in advance the expected level of ACCURACY and DEFINIVENESS of the results. La plata es una restricción.

2) SELECTING DATA A veces la data no está disponible al alcance. en estos casos es necesario identificar otras fuentes y/o crear nuevas fuentes de recolección. El output de todo el proceso, su accuracy y etc depende en gran medida de la calidad de las mismas datos.

3) PREPROCESSING. DATA Eliminar data q no esté completa o q esté mal cargada. Desarrollar un formal method of dealing w/ missing data and determine if it is randomly or systematically missing.

- 4) Transforming Data determine the appropriate forever to store data.
 Reduzir el N° de vars. Al tratar de los ciclos se necesita el explicar el fenómeno. Principal component Analysis
- 5) Storing Data must be stored in a forever this gives unrestricted and unbounded read/write privileges to the data scientist. Servers & Backup and Store. + Security and Privacy.
- 6) Data Mining data mining algo, ML, data viz., etc.
- 7) Evaluating Mining Results formal evaluation. Testing predictive capabilities, "In-sample forecasts" & dealing w/ Stake Holders. ¡I REESTE!

→ ML vs DEEP LEARNING!

- Machine Learning = A subset of AI that uses composed algorithms to analyse data and make intelligent decisions based on rules it has learned, without being explicitly programmed. Trained on large sets of data, learn from examples, not follow rules-based algorithms.
- Deep Learning : A specialized subset of machine learning that uses layered neural networks to simulate human decision making. Can detect and recognize info and identify patterns. Continuously learns on the job.
- Artificial Neural Networks take inspiration from biological NN. Collection of small computing units called neurons, they take input and decide. usual neurons consider deep learning.

→ **AI vs DATA Science**

DS usa AI pero es en esencia un
diferimiento q 'incluye todo el proceso
de procesamiento' de datos.
Mientras q AI incluye lo q hace q
los PCs procedan aprendizaje y tomen
decisiones intelectuales.

⇒ **Regression** (No vimos nada acá...)

⇒ **CURSO 1 - Semana 3**

- How is DS solving lives? ML y DL para predicciones
económicas en Argentina??

"If you're unable to measure something, you
are unable to improve it."

⇒ Companies Should Be capturing DATA !!

"DS will change the way Companies compete
and operate"

→ "THE FINAL DELIVERABLE"

~ 1000 - 7000 words. * 4 communicating findings to stakeholders. * w/ academic. In business world, the FD ~ 1500 words w/ many figs or maybe several hundred pages as well.

Much more powerful, because of its narrative, than a Power Point presentation.

Before Analysts start their report and presentation, they should (thus) have discussed the scope of the FD. Not doing so could likely to result in a poor quality FD. (decide scope, then look for DATA AND Analytics).

→ COPPERS AND RECRUITING FOR DS

=> HOW TO BECOME A DS? High end DS are mostly PhD's. Come out of Physics, Maths, have to have CS background, MATH BG, DB's Stats Probability.

- FOR A DEAN => • KNOW HOW TO PROGRAM
- Algebra, Calculus, Prob, Stats. • DB's

- As you go up in the field, you really need to know a lot about CS theory, stats prob.
 - One of the ways you learn things is: you do them! ☺
- ⇒ Recruiting for DS: given the pool of applicants you have, who has the most passion with your firm's DNA? Who's passionate? * you can teach Analytics skills. Curiosity. Do they have a sense of humor? lol. At least, mechanical skills. * Story telling. Someone who is reliable. He / she will need to relate with many departments.
- ⇒ The Report Structure, it depends on the length of the document. This varies depending on its purpose.
- | | | |
|---------------------|-----------------------|--------------|
| - COVER PAGE | - ABSTRACT | - APPENDIX |
| - TABLE OF CONTENTS | - INTRODUCTORY Review | |
| - METHODOLOGY | - RESULTS | - DISCUSSION |
| - CONCLUSION | | |

Preggs. Placed your
final report

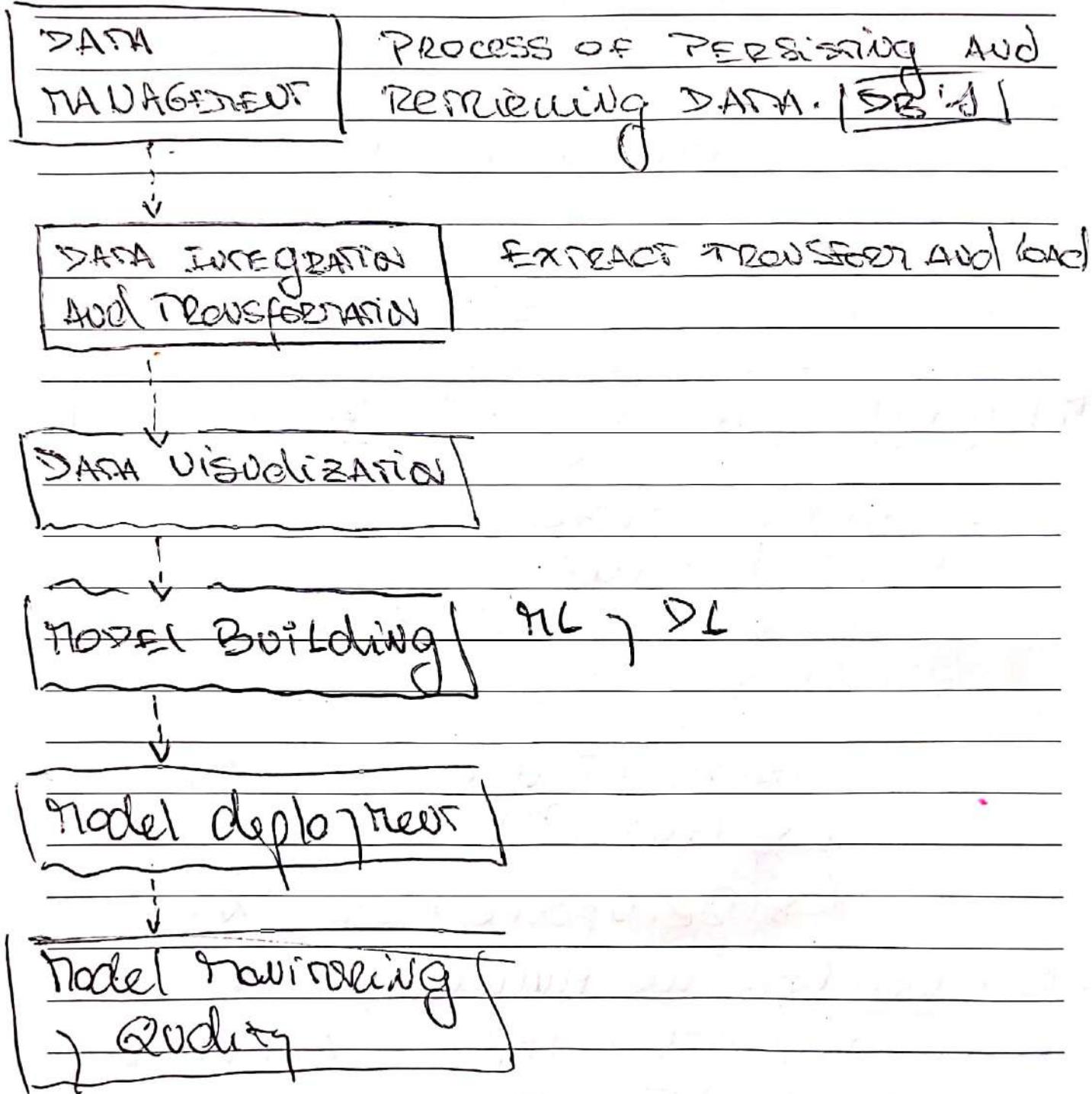
- 1 - Have you told readers at the outset, what they might get by reading your paper?
- 2 - Have you made the aim of your work clear?
- 3 - Have you explained the significance of your contribution?
- 4 - Have you set your work in the appropriate context & giving sufficient background (refs) to your work?
- 5 - Practicality? usefulness of your work?
- 6 - Future development that might result from your work?
- 7 - Have you summarized your paper in a clear and logical fashion?

Fin curso I
± 94 %

CURSO 2: TOOLS FOR DATA SCIENCE



• Python • R • SQL • OTROS.



- Execution Environments

- Data Asset Management

- Code "..."

- DEV. Environments (IDE)

- fully integrated visual tools

→ Sharing Surprise DAM → DAX = DATA ASSET Exchange, TOBL de IBM.

→ MACHINE LEARNING AND DL MODELS

- Identify Patterns.

- must be learned

- Create Predict, decide, etc.

- 3 main types

↳ Supervised learning

↳ Unsupervised learning.

↳ Reinforced learning.

• Supervisado = in humans provee los resultados correctos, labeled data → correct output

Regressión

SE usa para

clasiificación.

- Regressión → - Predecir Real Nums. Vals

- Classif → Clasify things into categories

• NO Supervisado)

LA DATA NO ES USOLED. EL MODELO TIENE Q APRENDER A CLASIFICAR DATOS SOLO CON LA INFO DE UNA DATA EN SI.

↳ Custering (dividir en subgrupos parecidos)

↳ Anomaly Detection (Anomalías)

• Reinforced

Parecido a como APRENDE LOS SÓLOS VIVOS. PUEDE J JERRO, PREMIOS. CHESS, GO, ETC.

⇒ [DEEP LEARNING]

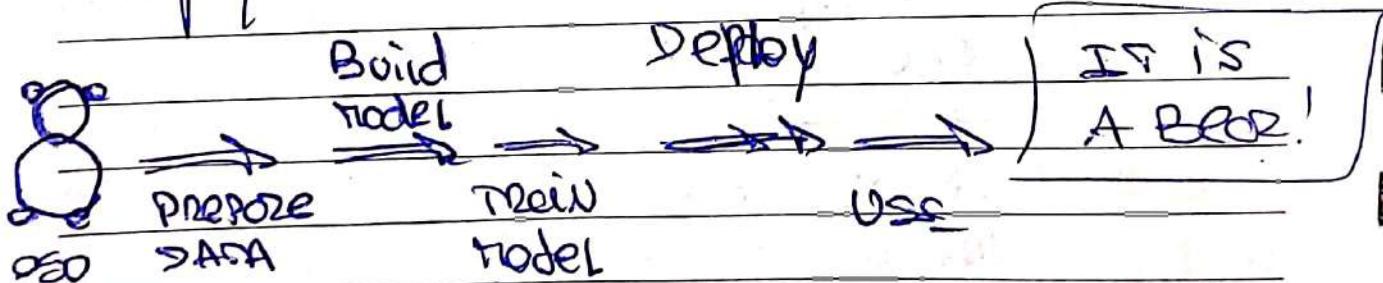
intenta reproducir la forma en la que los
Humanos aprendemos.

↳ NLP, Image processing, Audio y
Video Analysis, Time series forecasting,

↳ Necesita de gran cantidad de
of labeled DATA o training. Es
un proceso computacionalmente intenso.

↳ Build from Scratch or Public Repos.

- TensorFlow - Keras.
- PyTorch



INTERESTING.

↳ LONG TIME TO VALUE

↳ consume muchas veces usar Pre-

→ Podes usar Jupyter Notebook (Model Asset Exchange) de IBM.

■ SEMANA 2: Jupyter Notebooks

Allows to combine ~~real~~ code, tables, graphs, output y more.

1. Jupyter Lab → Es como el R studio

de Jupyter 

~~Actualmente~~

• EDA = Exploratory Data Analysis. Grafis, understanding the dataset.

Se hicieron SALTOS de Jupyter a R studio. No te di mucha Boleta pero ya lo use,) para Aprenderlo bien en estadística.

• Git y GitHub

→ Short Glossary

- SSH Protocol = Method for remote secure log in from one pc to another.
- Repository = Folders of the project set up for version control
- Fork = A copy of the Repo.
- Pull Request = The process that you use to request that someone reviews and approves your changes before they become final.

Commands

- INIT
 - ADD
 - STATUS
 - COMMIT
 - PULL
 - BRANCH
 - CHECKOUT
 - MERGE
- RESET → undo changes
- LOG → Browse previous changes.
- BRANCH → CREATE AN ISOLATED ENV. 4 CHANGES
- CHECKOUT → see y change branches.
- MERGE → put EVERYTHING back together

Pull req = lo haces cuando haces cambios
y querés q los demás los vean p/
sumercer al mejor. Esas propuestas
tus cambios.

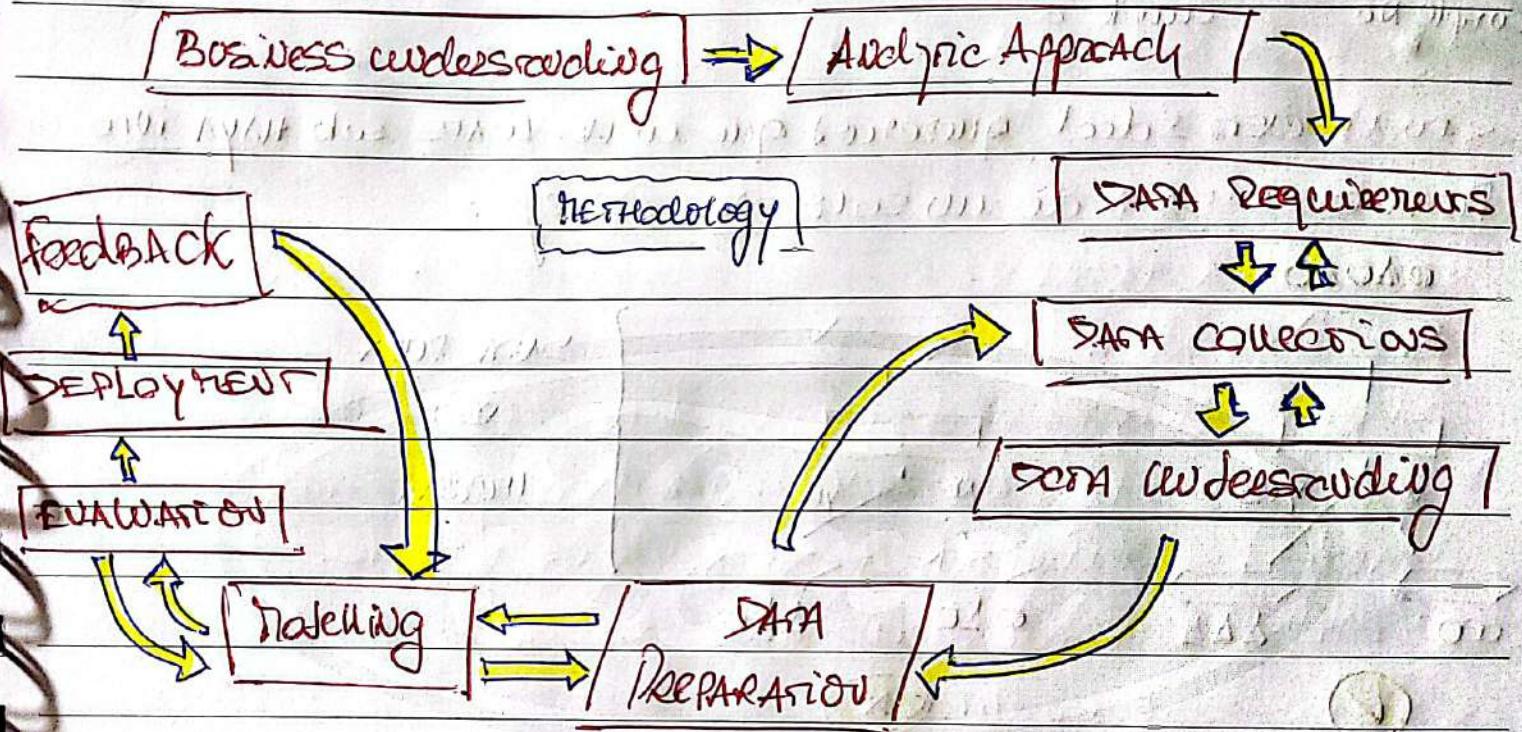
⇒ Sentencia 3 = IBM Watson Studio

La seguí con los cursos de Data Science
en uno cuaderno (Ford Grande).

→ ISTI - DATA SCIENCE METODOLOGÍA - ISTI DS P. CERIF.

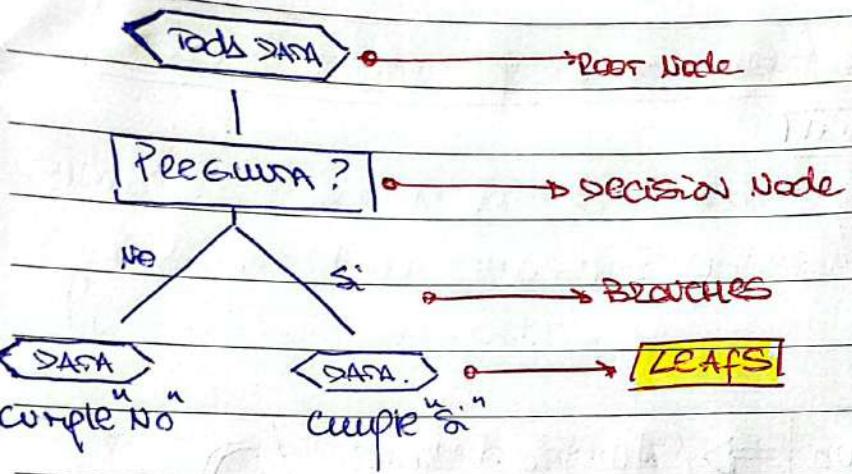
FROM PROBLEM TO APPROACH

- **CRISP-DM** = metodología iterativa para aplicar DS. Mucha énfasis en la importancia del Business Understanding y en hacer preguntas.

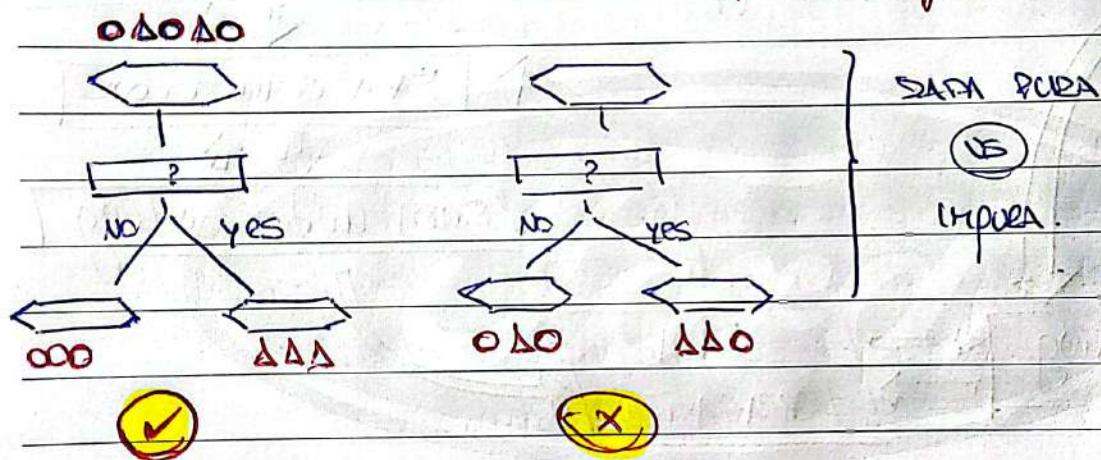


DECISION TREES

- Built using recursive Partitioning (such as subgroups) to classify DATA
- When partitioning, decision trees use the most predictive feature to split the DATA (to reduce a metric en general?)
- Predictiveness is based on decrease in entropy - gain in information, or impurity.



→ EN EL CASO IDEAL, queremos que en las LEAVES solo haya INFO PURA. OSEA, DATA de UN SOLO TIPO. Por ej:



A TREE STOPS GROWING AT A NODE WHEN

- Pure or nearly pure
- No remaining vars. to subset the data.
- It has grown passed a preselected size limit

From Requirements to Collection

- Allgoal q eu ta coociar, los ingredientes (data, tipos) serán obte por lo Bono Resumido
- Es importante saber que datos son la Necesaria para resolver y como hacerlos con esos datos.

From Understanding to Preparation

- Is the data that you collected representative of the problem you need to solve?

• Descriptive statistics

UNIVARIATE STATISTICS

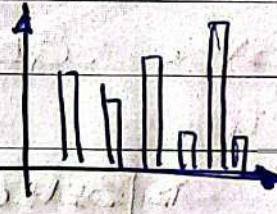
POLYWISE ~~STAT~~ CORRELATIONS

TESSERATURE

$$f(a) + \sum_{k=1}^N \frac{1}{k!} \frac{d^k}{dt^k} f(u(t)) \Big|_{t=0} = \int_0^1 f(u(t)) + \int_1^N \frac{1}{(t-1)!} \frac{d^{N+1}}{dt^{N+1}} f(u(t)) dt$$

$$f_{xy}(x,y) / f_{xy}(x,y) = f_x(x) f_y(y)$$

| Pobl | Freqüencia!



• Iterative data collection and

understanding

[ERC]

From Modelling to Evaluation

- Plantear un modelo de statisticas
- Evaluarlo
- Tunarlo. (ie cambiando parámetros pl q responde en resultados)

From Deployment to feedback

Una vez q el modelo sea desplegado lo medirás + pruebas y explicarás a stakeholders + feedback. Despues del 1er año se revisa el modelo.

DATABASES AND SQL FOR DATA SCIENCE

- BASICS.
- RELATIONAL DATABASES

"Structured Query Language"

WHAT IS A DATABASE?

- Repository of data
- Adding, modifying and querying data
- different kinds of databases store data in + formats.
- TABLE \leftrightarrow RELATIONAL-DB \leftrightarrow Rows + Cols.
- DBMS = DataBase Management System. (RDBMS = Relational)

BASIC SQL COMMANDS = CREATE, INSERT, SELECT, UPDATE, DELETE

SELECT STATEMENT

- `SELECT * FROM <TABLENAME>`

que pasa si quiero filtrar? \rightarrow WHERE

- `SELECT book_id, title FROM Books`

WHERE <predicate> \rightarrow True/false/unknown

$\downarrow \text{ref}$

`book_id = "B1"`

OP	SUMX	
Equal	=	Los Statement's q son
Greater	>	DML, son los q permiten
Lesser	<	modificar la data
g/eq	>=	"DATA Manipulation Lang."
l/eq	<=	#
Abreq	<>	[DDL] "Data def. lang"
		Define, change drop...

→ Select Statement's Expressions

- COUNT() = Retrieves N° of Rows matching the query criteria.
 - SELECT COUNT(*) FROM TABLE
 - SELECT COUNT(COUNTRY) FROM MEDALS
WHERE COUNTRY = "canada"
- DISTINCT() = Remove duplicate values from Query;
Retrieves unique values.
 - SELECT DISTINCT columnname FROM TABLENAME
- LIMIT() = N° max de Rows en el Query
 - SELECT * FROM TABLE
WHERE YEAR = 2018 LIMIT = 5

→ INSERT Statement

Agregar nuevos filas a la DB.

- INSERT INTO <TABLE> <(columnName, ...)> Values (value, ...)
 - Tengo q proveer igual cantidad de cols.

→ Se puede Agregar de A \oplus de una Row a la vez!

INSERT INTO AUTHOR (Id, Name, Email)
VALUES

(A1, Chrig, C1)
(A2, Jow, C2)

UPDATE and DELETE Statements

→ UPDATE TABLE

SET [ColName]=[Value]

Where [condición]

Puede ser \oplus de una \exists

Logos sin pones el

Where combinan todos
las filas!!! ⚡

→ DELETE FROM TABLE

Where [condición]

Sentencia 2

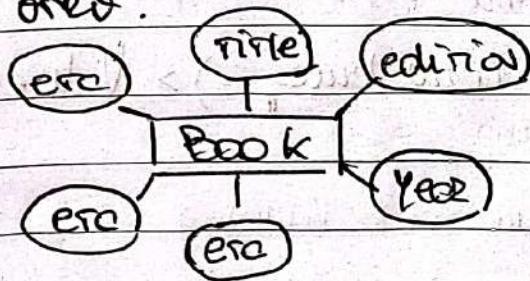
Relational DB Concepts

Information Model And DATA Models.

- Relational model → Allows for data independence

→ Entity Relationship Model = thinking a db as a collection
of entities. (u designing DB's). Objects q existen indep.

Entidades

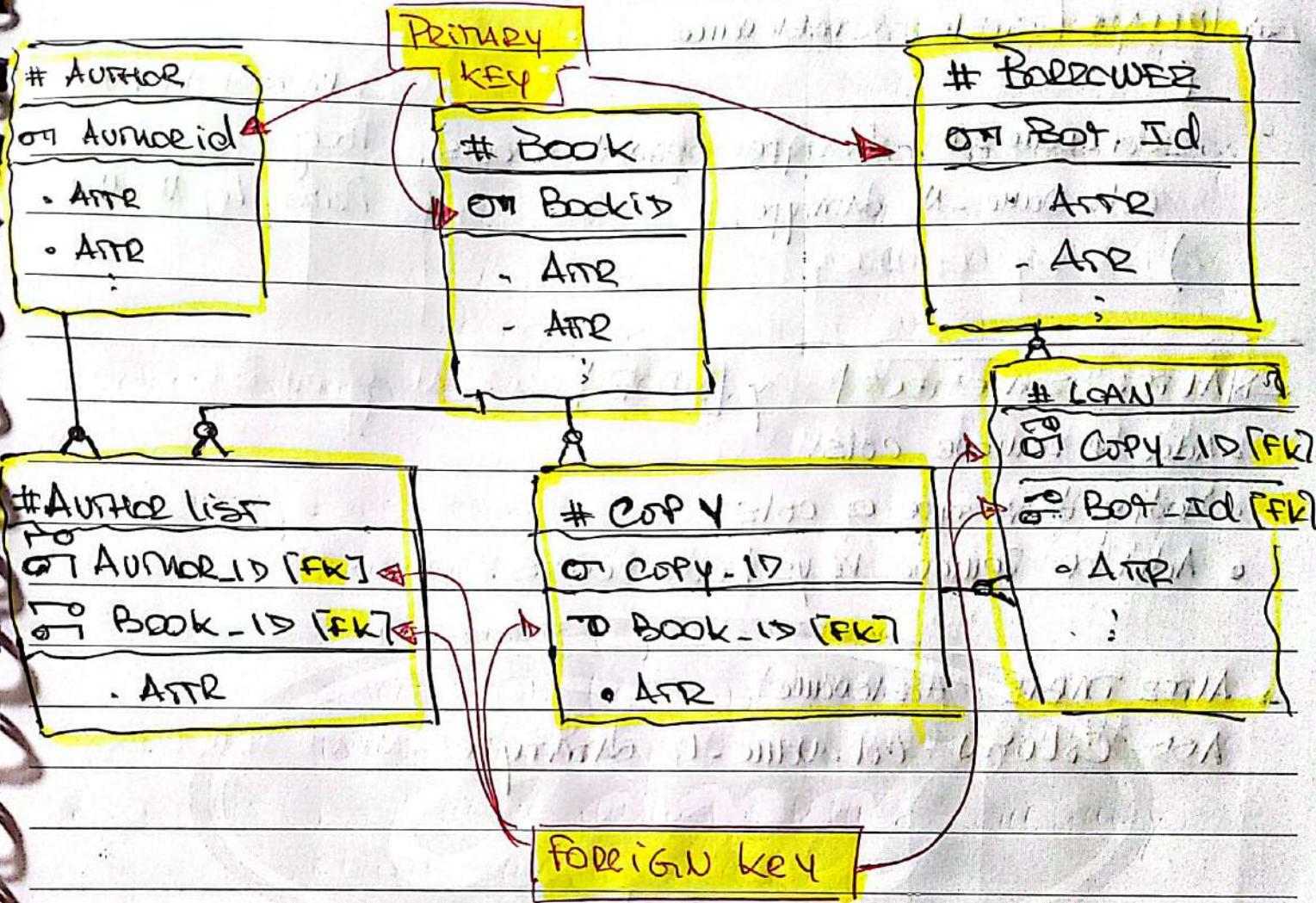


○ = ATTRIBUTES

□ = ENTITY

THE Book TABLE	
on Book id	ENTITY → Rows
• title	ATTRS → cols
• edition	
• Year	

Para nuestro ejemplo de la DB de libros, podría verse así:



- Primary key = Key una TABLA. Asigna valores únicos a cada fila.
- foreign key = Primary keys definidos en OTROS TABLAS q definen un link (relación) entre las tablas implicadas.

DDL = defn. def. larg. SISTEMAS (+ DML)

↳ define, change or drop data

↳ CREATE, ALTER, TRUNCATE, DROP.

ELIMINAR; ELIMINAR

Por lo de TABLA.

DBMS

→ **CREATE TABLE** tablename
(
Col-Name - 1 datatype optional - Pairs*, Prefs = Visos + Adelantas.
Col-Name - N datatype,
) Primary key Not Null.

→ **ALTER SCHEMA** | **DROP**

- Add or Remove cols.
- Modify data type of cols.
- Add or Remove keys / constraints.

→ **ALTER TABLE** tablename

ADD COLUMN col-name - 1 datatype

ADD COLUMN col-name - N datatype;

→ **ALTER TABLE** tablename

ALTER COLUMN col-name **SET DATA TYPE**

<datatype>; # OJO, si ya tiene otro ST, tira error

→ **ALTER TABLE** tablename

DROP COLUMN telephone - null;

→ **DROP TABLE** tableName;

→ **PRIVATE TABLE <Name>**

INTERMEDIATE: # BORRAR TODOS LOS DATOS (FIJAS). NO SE PUEDE DESHACER

⇒ **INTERMEDIATE SQL**

- Refining Results: String Patterns, Aggregating, Grouping
- Functions, Multiple Tables, Sub-queries.

↳ **USING STRING PATTERNS AND RANGES**

¿Qué pasa si no se que valor poner en el WHERE? Puede, se que de el tamaño empieza con "R" pero no me acuerdo exacto. ¿Cómo hacerlo?

→ WHERE column LIKE **<String Pattern>**

% = Representa missing entries. Si pongo "R%" estoy diciendo: Los q empiecen con R.

Ranges

→ select title, pages from Book

where pages **BETWEEN 290 AND 300**

IN

→ select name, country from Author
where country **IN ("AU", "BR")**

"set of values"

↳ **SORTING RESULT SETS: ORDER BY.** (Def = Ascending)

select title from Book

ORDER BY title DESC;

Tamb. puedes poner "order by 2" donde 2 = Segunda col del select statement.

GROUPING Result sets

SELECT COUNTRY, COUNT(COUNTRY)
FROM AUTHOR

COUNTRY	COUNT
ARG	100
BRA	0

GROUP BY COUNTRY:

TAMB SE LE PUEDEN AGREGAR CONDICIONES AL GROUP BY
CON LA KEYWORD "HAVING"

SELECT COUNTRY, COUNT(COUNTRY) AS COUNT
FROM AUTHOR

GROUP BY COUNTRY

HAVING COUNT(COUNTRY) > 4; BASICAMENTE 10 FILAS

→ Built-in Functions

Si bien puedes sacar la DATA, trabajos, tareas
etc ya viene con fns. y es mucho mas rapido
usar ESAS.

Tambien es posible definir NUESTRAS Propias
funciones.

Entre ellas hay: SUM, LENGTH, UCASE, DATETIME
etc.

→ Sub Queries and Nested Selects

SELECT column1 from TABLE

where column2 = (SELECT MAX(column2) from
TABLE)

SUBQUERY

Supongamos que quiero obtener una lista de
empleados que ganan mas que el salario
promedio

SELECT * from Employees

where salary > AVG(salary)

ERRORES!

EN CASO

(no se las vs Agg Functions pueden

crecer en el where, por
eso usamos N.Select)

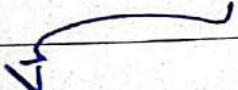


select * from employees

where salary <

(select Avg(salary) from employees);

→ Los N.Select tambié pueden ir en
otras partes. por ejemplo en columns.
A esto se le nombra



• Column Expressions

select empid, salary,

(select Avg(salary) from employees)

AS AVG-Salary

from Employees;

• DERIVED TABLES OR TABLE EXP.

N.Select.
(w from)

Pedj. Si queremos una tabla sin DATA sensible:

select * from

(select empid, f_name - l_name from employees)

AS EMP4ALL;

Eso es un ejemplo trivial q podríamos
querer pedirlo de uno, pero dice que
es falso.

Avg $\Delta t \approx 3.19 \text{ ms}$.
 → Tiempo total inicial $\approx 4-5\text{s}$.
 → 3.7K veces de llenado. \rightarrow ¿dónde ref?

⇒ Working w/ more than ONE TABLE

- ↳ • Job-queries
- Implicit join
- Join Operators (inner, outer, etc)

Select * from employees

Where DEP-ID in

(Select SEPF-ID-DEP from departments
 Where Loc-ID = "L002")

→ Un join implícito o full join es el que toma todo de los 2 TABLAS (o mas)

Select * from employees, departments;

TABLA 1

TABLA 2

⇒ Python + DB using

DBAPI: → Puede usar Python si INGRESA SQL en
 cualquier DB con el mismo código.

Conceptos clave

→ ~~BASES DE DATOS~~

- CONNECTIONS OFICIOS

→ Connections Objects

- DATABASE CONNECTIONS
- Manage Transactions.

→ CURSOR OBJECT

- DATABASE Queries
- Scroll through result set.
- Retrieve results.

from database import connect.

CONN = connect('db', 'user', 'pwd')

CUR = CONN.cursor()

CUR.execute('select * from myTable')

RESULTS = CUR.fetchall()

CUR.close()

CONN.close()

→ Creating tables, loading data and querying

SQL magic (fuzzy)

↳ nos dejó correr SQL "Adentro" de Python
de una forma 2 (x-x) Toma todo la celda.

↳ SQL SELECT * FROM TABLA_NOME

drawback = No podes elegir columna
CERRAR LA CONEXIÓN PI EJECUTAR
RECURSOS.

"Magic module"

→ Views / forma alternativa de traer la data
de una TABLA

- Pueden incluir columnas de múltiples TABLAS o VIEWS.
- Una vez creadas se puede hacer queries como una TABLA
- Solo la DEFINICIÓN de una VIEW se guarda
NO la DATA.

Views con

- SHOW PART OF THE DATA (SELECTIVE, ETC)
- COMBINE TWO OR MORE TABLES IN ONE VIEW FULL WAY.
- SIMPLIFY ACCESS TO DATA (HAS PERIODIC ALTER VIEW AND NO ALTER TABLES)

→

```
CREATE VIEW <Name> (<COLALIAS 1>,
<COLALIAS 2>, ... <COLALIAS .N>)
AS SELECT <COL1>, <COL2>, .. <COL N>
FROM <TABLE NAME>
WHERE <predicates>;
```

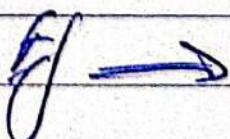
→ **STORED PROCEDURES** SER DE SQL Y GUARDADA EN EL SERVIDOR. EN CASO DE REQUERIR TODO EL TIEMPO DE UNA, MEJOR UN SP QUE COGRE TODOS LOS GUARDADOS POR OSE NOMBRE.

SE PUEDEN CODIFICAR EN TANTOS LARGOS.

VENTAJAS

- REUSE OF CODE EN VARIOUS APPS
- + SECURITY
- - DIFFICULTY
- + PERFORMANCE

(DATA | LOGIC | VIEW | JUNK)



CREATE Procedure update_sal (IN empNum CHAR(6),
IN Rating SMALLINT)

Language SQL

BEGIN

IF Rating = 1 THEN

UPDATE employees

SET salary = salary * 1.10

WHERE emp-id = emp-num;

ELSE

UPDATE employee

SET salary = salary * 1.05

WHERE emp-id = emp-num;

END IF;

END

→ CALL update-SAL ('E100', 1)

→ **ACID TRANSACTIONS**: MINVISIBLE UNIT OF WORK. CONSISTS OF 1 OR MORE SQLS. SUCCESS → SUCCEEDS OR NOT WITH 0 OR 1

Acid ↴

- Atomic = All changes will be performed successfully or not at all
- Consistent = Data must be in a consistent state before and after the transaction
- Isolated = No other process can change the data while the transaction is running.
- Durable = The changes made by the transaction must persist.
- PARA CODER UN ACID

Begin → ES IMPLICITO, NO SE ESCRIBEN
 UPDATE --- f , PROBLEMA EXPLICIT
 - - - f
 UPDATE - - . f 2
 - -

ROLLBACK . ↴ N