

Sample dataset (loan applications)

Loan ID	Gender	Married	Number of Dependents	Education	Self Employed ?	Applicant Income	Co-applicant Income	Loan Amount	Loan Term	Credit History	Property Area
LP001032	Male	No	0	Graduate	No	4950	0	125	360	1	Urban
LP001824	Male	Yes	1	Graduate	No	2882	1843	123	480	1	Semiurban
LP002928	Male	Yes	0	Graduate	No	3000	3416	56	180	1	Semiurban
LP001814	Male	Yes	2	Graduate	No	9703	0	112	360	1	Urban
LP002244	Male	Yes	0	Graduate	No	2333	2417	136	360	1	Urban
LP001854	Male	Yes	3+	Graduate	No	5250	0	94	360	1	Urban
...	Male	Yes	0	Graduate	No	3500	1667	114	360	1	Semiurban
LP001647	Male	Yes	0	Graduate	No	9328	0	188	180	1	Rural
LP001871	Female	No	0	Graduate	No	7200	0	120	360	1	Rural
LP001379	Male	Yes	2	Graduate	No	3800	3600	216	360	0	Urban
LP002789	Male	Yes	0	Graduate	No	3593	4266	132	180	0	Rural
LP001578	Male	Yes	0	Graduate	No	2439	3333	129	360	1	Rural
LP001318	Male	Yes	2	Graduate	No	6250	5654	188	180	1	Semiurban
LP001259	Male	Yes	1	Graduate	Yes	1000	3022	110	360	1	Urban
LP002804	Female	Yes	0	Graduate	No	4180	2306	182	360	1	Semiurban

Independent variables

Dependent variable ?

SAMPLE DATASET (automobiles)

KMs per liter	cylinders	displacement	horsepower	weight	acceleration	year	origin	Model name
18	8	307	130	3504	12	70	1	chevrolet chevelle malibu
15	8	350	165	3693	11.5	70	1	buick skylark 320
18	8	318	150	3436	11	70	1	plymouth satellite
16	8	304	150	3433	12	70	1	amc rebel sst
17	8	302	140	3449	10.5	70	1	ford torino
15	8	429	198	4341	10	70	1	ford galaxie 500
...
...
24	4	113	95	2372	15	70	3	toyota corona mark ii
22	6	198	95	2833	15.5	70	1	plymouth duster
18	6	199	97	2774	15.5	70	1	amc hornet
21	6	200	85	2587	16	70	1	ford maverick
27	4	97	88	2130	14.5	70	3	datsum pl510

Dependent variable ?

Sample data - jobs

Month	Total Filled Jobs
2004M07	1795610
2004M08	1792770
2004M09	1809590
2004M10	1815580
2004M11	1856360
2005M04	1871630
2005M05	1867870
2005M06	1857260
2005M07	1858360
2005M08	1856320
2005M09	1876270
2005M10	1866920
2011M10	1903630
2011M11	1940200
2011M12	1983070
2012M01	1865540
2012M02	1932380

Independent variables ?

Dependent variable ?

Stock prices

Date/time	\$\$	
8/7/2024 10am	10	
	11	
	12	
	10	

Stock price -> both dependent and independent!

Price at a moment is dependent on few other prices in the series (sequence)

11am price -> {9am, 10am, 10:30am}

TIME SERIES data

Examples : weather, all economic data, sensors, IOT, medical devices

Handling -->

- Traditional statistical methods (ARIMA, AR, MA, ARMA Variations of ARIMA)
- ML methods
 - o Classical methods
 - Sequence models

- Advanced
 - Transformer

- ML/DL methods are better than classical statistical method (ARIMA)

Sample data - text

Tweet id	Airline sentiment	Retweet count	Text	Tweet location
570306133677760000	neutral	0	~@VirginAmerica What @dhepburn said.	
570301130888120000	positive	0	@VirginAmerica plus you've added commercials to the experience... tacky.	
570301083672813000	neutral	0	@VirginAmerica I didn't today... Must mean I need to take another trip!	Let's Play
570301051407624000	negative	0	@VirginAmerica it's really aggressive to blast obnoxious "entertainment" in your guests' faces & amp; they have little recourse	
...
570300767074181000	negative	0	@VirginAmerica seriously would pay \$30 a flight for seats that didn't have this playing. it's really the only bad thing about flying VA	
570300616901320000	positive	0	@VirginAmerica yes, nearly every time I fly VX this scare worm is is go away :)	San Francisco CA
570300248553490000	neutral	0	@VirginAmerica Really missed a prime opportunity for Men Without Hats parody, there. https://t.co/mlWpG7grEZF	Los Angeles
570285904809598000	positive	0	@VirginAmerica Thanks!	San Francisco, CA
570282469121007000	negative	0	~@VirginAmerica SFO-FDX schedule is still MIA.	palo alto, ca
570277724385734000	positive	0	@VirginAmerica so excited for my first cross country flight LAX to MCO I've heard nothing but great things about Virgin America. #20DaysToGo	west covina
570276917301137000	negative	0	@VirginAmerica I flew from NYC to SFO last week and couldn't fully sit in my seat due to two large gentlemen on either side of me. HELP!	this place called NYC
57027684613923000	positive	0	I am flying @VirginAmerica. I'm in the air	Somewhere celebrating life.
570267956648792000	positive	0	@VirginAmerica you know what would be amazingly awesome? BOS-FLL PLEASE!!!!!! I want to fly with only you.	Boston Waltham
570258583513384000	negative	0	~@VirginAmerica why are your first fares in May over three times more than other carriers when all seats are available to select???	

Independent variables ?

Dependent variable ?

Processing text data (challenges)

- Unstructured
- Mixed cases
- Volume
- Language
- Typos, errors, grammatical errors
- Context of the text --> Semantics

"My flight is delayed. Brilliant!" -> Happy/Sad

Feature types

1. Numeric

a. Plain numeric

i. Float

- 1) Ratio (we will encounter 95% of the time)
- 2) Interval
 - a) Lacks a true zero
 - i) Temp in C
 One. 50 C and 55C

b. Discrete

i. Count

- 1) Number of children
- 2) Number of cars
- 3) Number of claims
- 4) Score?
- 5) No. of customers

2. Categorical (Qualitative variables)

a. Nominal

- i. No order
- ii. Can we apply any math op?
 - 1) No

b. Ordinal

- i. Which has order
 - 1) SML
 - 2) Rating scale
 - 3) Score
- ii. Can we apply any math op?
 - 1) Limited (sorting)

c. Binary

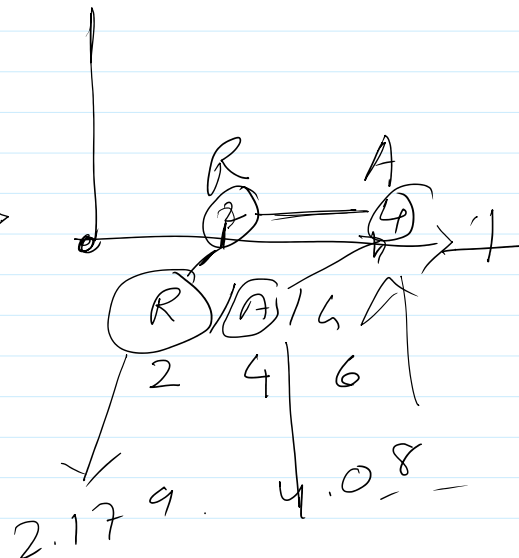
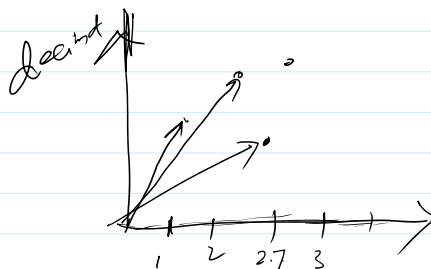
- i. Two categories

d. Cat data are BIG problem in AI

- i. ML/DL - are math wise - work well with decimals or cont data
 - 1) Assumption is that there is origin within the data (0,0)

3. Text

4. Date & time



5. Boolean
6. Spatial features (geo , locations)
7. Image/ videos
8. Audio

10. Derived form of features

Business scenario (from slide # 5) Lesson 3

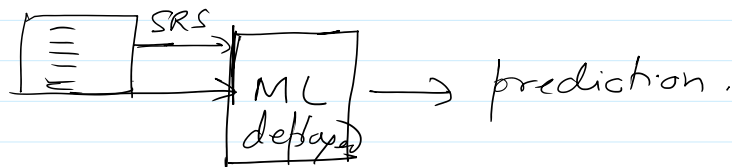
Healthcare org

- Manage
 - o Hospitals
 - o Colleges
 - o Clinics
 - o Labs
 - o Pharmacy
 - o Medicines
 - o staff
- Vast amount of data and huge number of data types
 - o Volume
 - o Variety
 - o Format
 - o dynamic
- 3 buckets for the data
 - o Purpose
 - o Relations
 - o Storage
- Dat sources
 - o Patient records
 - o Clinical systems
 - o Financial mgmt
 - o Operational/administrative
 - o Billing
 - o Compliance system
- **Bucket 1 (Patient data)**
 - o **Purpose**
 - Demographics (age ...
 - History of illness
 - Insurance info
 - Diagnosis
 - Lab results
 - Treatment/ regimen
 - Prescriptions
 - ...
 - ...
 - ...
 - EHRs , FHIR
 - Sensitivity
 - PHI
 - PII
 - Data secured
 - ◆ Access controls
 - ◆ Encryption
 - ◆ Backup and rollback ...(failsafe)
 - o Relations
 - Clinical info
 - Lab
 - Medication
 - Orders
 - Billing
 - ...
 - ...
 - o Storage methods
 - Centralized location

- Bucket 2 (clinical records)
 - o Decision making
 - o Accurate diagnosis
 - o Treatment plans
 - o Monitoring the patients
 - o Legal
 - o Fin
 - o Billing
 - o Insu

Population and samples

- Sampling techniques [...,...]
- o Simple random sampling (SRS)
 - Every element has equal chance of being picked
 - [10, 13, 15, 16, 18, 19.... 1000 numbers]
 - Implement
 - Length of array = len(arr) = 1000
 - **Rand**(len=1000) - rand - function in Python/R
 - ◆ Return any number from 1 ... 1000
 - ◇ 77
 - Arr[77]
 - Use cases
 - Large dataset
 - Prelim analysis or understanding
 - selecting a list of transactions to review for financial audit

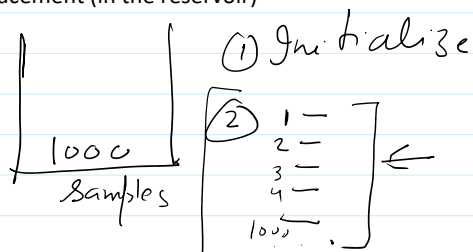


- Ensure FAIRNESS of the evaluation
 - ◆ SRS
 - If the dataset is HOMOGENOUS
 - ◆ No groups in the data (not heterogenous)
 - ◇ SRS
 - **Flip side**
 - Outliers
 - Chance based -> selections may not fully representative
 - Samples being small
- o Stratified random sampling
 - Subgroups in the data
 - Randomly pick samples from each of the strata (subgroups)
 - **Variability** = subgroups (ensured)
 - Use case
 - Clinical trial (new drug)
 - ◆ Effect on age groups, genders, cultures...
 - AI - all ML and DL model (evaluation)
- o Systematic sampling
 - Select every **nth** sample
 - Sorted, starting point
 - Maintain randomness in picking samples ???
 - May introduce some bias??? (periodicity)
 - Careful starting point, the interval (n) determination
 - Use case
 - Time series data

- Cluster sampling
 - Dividing the population into clusters or groups
 - Understanding cluster - Easy or difficult
 - ◆ Defining cluster
 - ◇ Using 1 col or multiple columns
 - 15 clusters
 - Subset of these 15 clusters
 - 5 of them
 - Use cases
 - Healthcare
 - ◆ Hospital patient records
 - ◆ Group of hospitals by location
 - ◆ Randomly select a subset of hospitals
 - ◇ Collect all patients from these subsets
 - Efficient
 - ◆ Large datasets
 - ◆ Spatial characteristics
 - ◆ CLUSTERS SHOULD BE NATURAL GROUPINGS
 - Flip side
 - Biased estimates
 - ◆ Some group
 - ◇ More variability than others
 - Miss out on other groups (completely)
 - Cluster size sufficient

- Reservoir sampling

- Randomized algo
- Pick FIXED number of samples --> model/analysis
- Useful -
 - Large datasets
 - Perform random sampling
 - Data too large to fit into memory
 - When we deal with streaming data (continuous data)
- Reservoir = bucket = fixed number of samples
- Process the incoming stream
- Random replacement (in the reservoir)



③ 1005th sample

3:1

Any random pos in the reservoir (1, 1000)

1005th → 517

3:2 replacing the existing reservoir (recency)

- Flip side
 - Representativeness
 - ◆ No

Measures of central tendencies

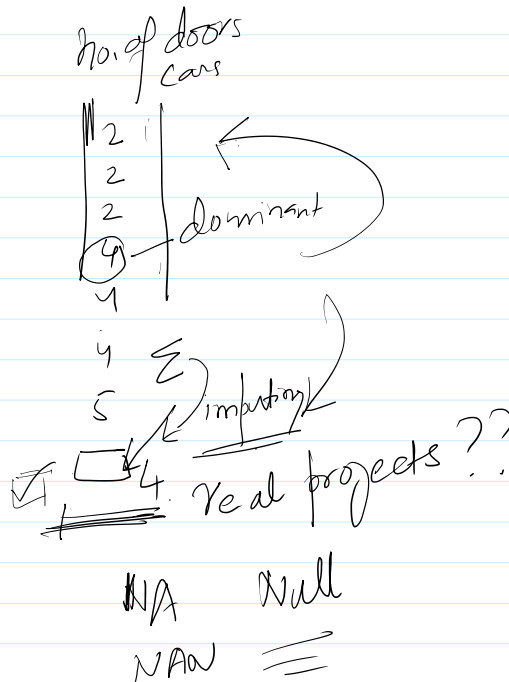
- Why we need understanding of CT
 - o Know the typical range of values for features

- Mean

- o Avg
- o Sum/n
- o Decimal (discrete)
- o Measure of the central value
- o **Affected by outliers**

o Use cases

- Aggregating
- **Imputing**
 - Fill in the missing spots



Handwritten notes:

- u₁ 5✓
- u₂ X 5✓

Geometric mean

- Column

- o Values are multiplicative in nature
 - Return rate 5%, 10%, 15%
 - Growth rate
 - Year 1 = initial value * (1 + int rate)
 - Year 2 = year 1 * (1 + in rate)
 - Year 3 = year 2 * (1 + int rate)

Handwritten note: exponential.



$$\sqrt[3]{x_1 \cdot x_2 \cdot x_3}$$

How do we find out if a feature exhibits exp nature?

- Plot it !!

Harmonic mean

- Features

- o Rates, ratios, speed

$$\text{Harmonic Mean} = n / (1/x_1 + 1/x_2 + 1/x_3 + \dots + 1/x_n)$$

• Where:

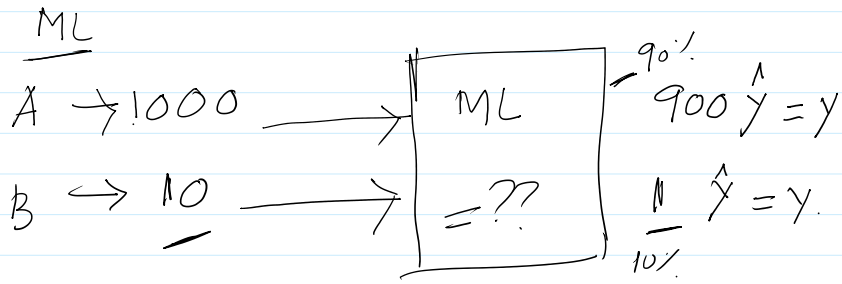
- $x_1, x_2, x_3, \dots, x_n$ are the individual values in the dataset.
- n is the total number of values in the dataset.

The harmonic mean places more weight on smaller values in the dataset.

Math intuition behind this formula

- Makes the small values more weighty

Use case



$$\begin{array}{r} .9 \\ .1 \quad 2 \\ \hline \frac{1}{.9} + \frac{1}{.1} \end{array}$$

$$\frac{900+1}{1000+10} = \frac{901}{1010} \approx 9\%$$