#### Contents

- What is data?
- Machine learning databases
- Terminology
  - Data preparation, Data modality, Data fidelity
- Data formats and sources
  - Experiments, Imaging, Sensing, Modeling and simulation
- Examples
  - Diamond data for feature-based pricing
  - Data collection from indentation testing

- Key input for mechanistic data science
- Where does the data come from? It can come from many sources and in many formats. → multimodal data collection and generation
  - Physical observation: very costly and difficult to control independent variables
  - Modern computer HW and SW: simulate the physical experiments and generate further complimentary data
- Efficient data collection and management through a database
  - Expedite the problem-solving timeline  $\rightarrow$  Help in rapid decision making
- Goal of mechanistic data science (MDS)
  - Mining the data intelligently to extract the science
  - Combining data and mechanisms for decision making

#### Data is the central piece for science

- Question: how are forces transmitted by structural members?
- Galileo's approach:
  - Data collection: performed many experiments on how <u>size and shape</u> of structural members affects their ability to carry and transmit loads
  - Observations : as length of a beam increases, its strength decreases, unless you increase the thickness and breadth at an even greater rate
  - Science: This led Galileo to recognize what we now call the <u>scaling</u> <u>problem</u>, there are limits to how big nature can make a tree, or an animal, for beyond a certain limit, the branches of the tree or the limbs of the animal, will break under their own weight.
  - *deflction*:  $\delta = \frac{FL^3}{3EI}$  This formula to calculate deflection of cantilever beams works for macroscopic beams made with all materials, size, shape and loads



#### Data is the central piece for science

- Evolution of scientific discovery: from data to empiricism to mechanism
  - Astronomical data: Observations of planetary orbits



Mechanistic data science is the hidden link between data to science

#### Kepler's Law (Data to Mechanism)

- First law (law of orbits): Each planet revolves around the sun in an elliptical orbit with the sun situated at one of the two foci.
- Second law (law of areas): The real velocity of a planet around the sun remains constant, OR, The radius vector drawn from the sun to the planet sweeps out equal areas in equal intervals of time.
- Third law (law of periods): The square of the time period period(T) of revolution of a planet around the sun is proportional to the cube of the semi major axis (r) of its elliptical orbit.

the p



### Mechanism to Science: Discovery of gravity

- Newtons universal law of gravitation
  - Every point mass attracts every single other point mass by a force acting along the line intersecting both points. The force is proportional to the product of the two masses and inversely proportional to the square of the distance between them.
- Force on a falling object (apple from a tree) in earth due to gravity is given by *F=mg*



### Discovery of gravitation from planetary motion data



uata Generation and Collection - 7

### What is Data?

Data: collection of information (numbers, words, measurements, observations) or descriptions of things



- Text, numbers, images, graphs, and signals are all common forms of data
- Data represents all industries and problems: finance, climate, transportation, etc.





**Discrete** data can only take certain values (ex. only whole numbers)

7 data points = 7 people. You can't measure height for "7.3" people.

#### Measured data

**Continuous** data can take any value (within a range)

Temperature can take any value within Earth's range. It changes continuously, forming an infinite curve.

### Common Databases for Machine Learning Applications

• Database: organized collection of data, generally stored and accessed electronically from a computer system



**Top Sources For Machine Learning Datasets (2019)** 

**Top Free Machine Learning Datasets to Use in 2025** 

### **Data Preparation for Analysis**

- Raw data: collected from the source directly
- Data wrangling: mapping and transforming raw data to another format for machine interpretation (ex. map Yes/No to 1/0)
- Data Formatting: formatting data for consistency (ex. formatting text data with labels)
- Data Cleaning: providing attributes to missing values and removing unwanted characters from the data
- Database preparation: adding data from 1+ sources to build your own database
- Feature Extraction
  - Identification of important features in the data
  - Determined with human expertise



#### Data Wrangling/Munging/Janitor Work



Mechanisti

### Dataset for Machine Learning (1)



- $f(X_i^N)$  maps input  $X_i^N$  to output  $y_j^N$ .
- Machine Learning Goal = find the functional form  $y_j^N = f(X_i^N)$ 
  - For Kaggle diamond dataset:
    - *i*=1...9 (Carat, Cut, etc.)
    - *j*=1 (Price)
    - *N*=53,940 (number of diamonds sampled)

### Dataset for Machine Learning (2)



- The dataset is divided into training (70%), validation (15%), and testing (15%) sets to find the functional relationship and confirm it is the **best possible fit**
- This process is iterative. The model is repeatedly trained, validated, and tested
- Final performance on the testing set is evaluated when function error is minimal



- □ Training set: Inputs and outputs fit to mapping function  $f(X_i^N)$
- Validation set: Evaluate function (frequently after each training step)
- □ Testing set: Evaluate the final function  $f(X_i^N)$

# Example: How can we identify a high quality diamond at a reasonable price?

1. The Pink Star



Image Source: Cosmopolitan Italia **Price:** \$71 million Sold: April 2017 at Sotheby's Auction **Carat Weight:** 59.6 carats **Color:** Pink

#### 4. The Princie Pink Diamond



Image Source: DailyMail.co.uk **Price:** \$39.3 million Sold: April 2013 **Carat Weight:** 36.45 carats **Color:** Pink

Mechanistic Data Science

2. Oppenheimer Blue Diamond



Image Source: Christie's **Price:** \$57.5 million Sold: May 2016 **Carat Weight:** 14.62 carats **Color:** Blue

#### 5. The Orange



Image Source: NY Post **Price:** \$35.54 million Sold: November 2013 **Carat Weight:** 14.82 carats **Color:** Orange

3. Graff Vivid Pink Diamond



Image Source: Diamondhistorygirl **Price:** \$46 million Sold: November 2010 **Carat Weight:** 24.78 carats **Color:** Pink

#### 6. The Largest Diamond Ever Sold



Image Source: CNBC **Price:** \$30.6 million Sold: Christie's in 2013 **Carat Weight:** 118.28 carat **Color:** Colorless

Data Generation and Collection - 14

#### Mohs Scale of Hardness

Earth'	s c	GEM SELECT	
hardes t	<sup>s</sup> Mohs	Hardness	Scale
materi	Name	Scale Number	Common Object
al	Diamond	10	
	Corundum	9	Masonry Drill Bit / 8.5
E	Topaz	8	(100000
S	Quartz	7	Steel Nail / 6.5
	Orthoclase	6	Knife / 5.5
۲	Apatite	5	and particular and a
	Fluorite	4	Penny (Copper) / 3.5
	Calcite	3	
	Gypsum	2	Fingernail / 2.5
	Talc	1	

https://www.gemselect.com/gem-info/gem-hardness-info.php Mechanistic Data Science

- German mineralogist Frederick Mohs (1773-1839)
- How to Perform the MOHS Test?
  - Scratch it!



https://geology.com/minerals/mohs-hardness-scale.shtm

#### Features used to Characterize Diamonds



SIZE (carats) 2.

- CLARITY 3.
- COLOR 4.
- 5. CUT
- 6. BRIGHTNESS
- FIRE (dispersion) 7.
- 8 SPARKLE
- 9. POLISH
- 10. SYMMETRY
- **11. FLUORESCENCE**
- **12. DURABILITY**
- 13. LUSTER



HEART

OVAL



Total internal

reflection



PRINCESS

MARQUISE



PEAR

"One carat" (100 points) equals the weight of 1/5 of a gram





#### **Data Science focuses on quantifiable features**

Mechanistic Data Science

#### Different features can represent the same problem



#### Example: Diamond Data for Feature-based Pricing

- Kaggle (Datasets/Diamonds)
  - 53,940 diamonds with 10 features

	carat	cut	color	clarity	depth	table	price	х	у	z
1	0.23	Ideal	Е	SI2	61.5	55	326	3.95	3.98	2.43
2	0.21	Premium	Е	SI1	59.8	61	326	3.89	3.84	2.31
3	0.23	Good	Е	VS1	56.9	65	327	4.05	4.07	2.31
4	0.29	Premium		VS2	62.4	58	334	4.2	4.23	2.63
5	0.31	Good	J	SI2	63.3	58	335	4.34	4.35	2.75
6	0.24	Very Good	J	VVS2	62.8	57	336	3.94	3.96	2.48
7	0.24	Very Good	I	VVS1	62.3	57	336	3.95	3.98	2.47
8	0.26	Very Good	Н	SI1	61.9	55	337	4.07	4.11	2.53
9	0.22	Fair	Е	VS2	65.1	61	337	3.87	3.78	2.49
10	0.23	Very Good	Н	VS1	59.4	61	338	4	4.05	2.39

Price: (\$326--\$18,823) Carat: (0.2--5.01) Cut: (Fair, Good, Very Good, Premium, Ideal) Color: (J (worst) to D (best)) Clarity: (I1 (worst), SI2, SI1, VS2, VS1, VVS2, VVS1, IF (best)) Size in x direction in mm (0--10.74) Size in y direction in mm (0--58.9) Size in z direction in mm (0--31.8) Depth: z / mean(x, y) = 2 \* z / (x + y) (43--79) (%)Table: width of top of diamond relative to widest point (43--95) (%)

Color (D, E, F, G, H, I, J)  $\rightarrow$  (1, 2, 3, 4, 5, 6, 7) D: colorless ~ Z: light yellow or brown

Cut Rating	Numerical value	
Premium	1	
Ideal	2	
Very Good	3	
Good	4	
Fair	5	

Clarity Rating	Numerical value
IF—Internally Flawless	1
VVS1,2-Very, Very Slightly Included 1,2	2
VS1,2—Very Slightly Included 1,2	3
SI1,2—Slightly Included 1,2	4
I1—Included 1	5





#### About the data (Description of attributes)

This classic dataset contains the prices and other attributes of almost 54,000 diamonds. There are 10 attributes included in the dataset including the target ie. price.

- carat (0.2-5.01): The carat is the diamond's physical weight measured in metric carats. One carat equals 0.20 gram and is subdivided into 100 points.
- cut (Fair, Good, Very Good, Premium, Ideal): The quality of the cut. The more precise the diamond is cut, the more captivating the diamond is to the eye thus of high grade.
- color (from J (worst) to D (best)): The colour of gem-quality diamonds occurs in many hues. In the range from colourless to light yellow or light brown. Colourless diamonds are the rarest. Other natural colours (blue, red, pink for example) are known as "fancy," and their colour grading is different than from white colorless diamonds.
- clarity (I1 (worst), SI2, SI1, VS2, VS1, VVS2, VVS1, IF (best)): Diamonds can have internal characteristics known as inclusions or external characteristics known as blemishes. Diamonds without inclusions or blemishes are rare; however, most characteristics can only be seen with magnification.
- depth (43-79): It is the total depth percentage which equals to z / mean(x, y) = 2 \* z / (x + y). The depth of the diamond is its height (in millimetres) measured from the culet (bottom tip) to the table (flat, top surface) as referred in the labelled diagram above.
- **table (43-95):** It is the width of the top of the diamond relative to widest point. It gives diamond stunning fire and brilliance by reflecting lights to all directions which when seen by an observer, seems lustrous.
- price (\$\$326 \$18826): It is the price of the diamond in US dollars. It is our very target column in the dataset.
- x (0 10.74): Length of the diamond (in mm)
- y (0 58.9): Width of the diamond (in mm)
- z (0 31.8): Depth of the diamond (in mm)

star length

lower airdle/half

facet length

girdle

thickness

 $\wedge$ 

table size

total depth

 $culet \rightarrow k$ 

crown height

pavilion depth

crown

angle

pavilion

angle

#### **Example: Moneyball**

- Kaggle (Datasets/Moneyball)
  - MLB statistics 1962-2012
  - Billy Beane and Paul DePodesta, Oakland Athletics, 2002
  - 1,232 data with 15 features
    - Player: Batting average (BA, 타율), runs batted in (RBI, 타점)
    - Win 95 games to make the playoffs, score 133 more runs than opponents
    - On-base percentage (OBP, 출루율)
    - Slugging percentage(장타율) (SLG)=(1B+2B\*2+3B\*3+HR\*4)/AB(At Bat, 타석)
    - On-base plus slugging (OPS)=OBP+SLG

Team	League	Year	RS	RA	W	OBP	SLG	BA	Playoffs	RankSeason	RankPlayoffs	G	OOBP	OSLG
ARI	NL	2012	734	688	81	0.328	0.418	0.259	0			162	0.317	0.415
ATL	NL	2012	700	600	94	0.32	0.389	0.247	1	4	5	162	0.306	0.378
BAL	AL	2012	712	705	93	0.311	0.417	0.247	1	5	4	162	0.315	0.403
BOS	AL	2012	734	806	69	0.315	0.415	0.26	0			162	0.331	0.428
СНС	NL	2012	613	759	61	0.302	0.378	0.24	0			162	0.335	0.424
CHW	AL	2012	748	676	85	0.318	0.422	0.255	0			162	0.319	0.405

### **Example: Data Collection from Indentation Testing**

Mechanistic Data Science Applications: Materials Engineering How to analyze material properties?

- Indentation testing: material testing for hardness
- Hardness: resistance to penetration of a hard indenter (related to material strength)
- Significance
  - Testing is simple, fast, relatively inexpensive, and not destructive
  - Hardness is closely related to critical mechanical properties: strength, ductility, and fatigue resistance
  - More plausible at small scales than tensile test

#### **Example: Data Collection from Indentation Testing**



### Indentation Tests: Vary by Sample Size and Shape



Data Generation and Collection - 23

## Calculating Hardness from Load-Displacement Data

- What data we collect from indentation test?
- Hardness is calculated with the maximum applied load and the indenter contact area
  - P = applied load
  - h = indentation depth
  - S = slope of unloading curve
  - *C* = curvature of load-displacement curve
  - $P_m$  = maximum load
  - $h_m$  = maximum indentation depth
  - $h_c$  = critical indentation depth (difficult to measure and must be calculated)
  - A = contact area (depends on indenter shape)
  - *H* = hardness (GPa)
  - *E* = elastic modulus (GPa)
  - $\varepsilon$  = constant (depends on indenter geometry)
  - $W_e$  = Elastic work done
  - $W_p$  = Plastic work done



• 
$$h_c = h - \frac{\varepsilon P_m}{s}$$
  
•  $A = f(h_c) = 24.56h^2$   
•  $H = \frac{P_m}{A}$   
•  $E^* = \frac{s}{2h\beta} \sqrt{\frac{\pi}{24.56}}$  ( $\beta = 1.034$  for Berkovich indenter)

#### Indentation Data from Different Sources

- Data modality: data from different sources (modes)
- Data sources
  - Experiment: Instrumented indentation
  - Imaging: Scanning Electron Microscope (SEM)
  - Sensors: Load and displacement sensing using Linear variable differential transformer (LVDT)
  - Modeling and simulation: Finite element, Atomistic simulation



#### Indentation Data from Different Sources

- Load-displacement data can be found through physical experiments and computer simulations
  - Both experiment and simulations produce the same indentation
  - Deviations in modality require data calibration



#### Indentation Data at Different Scales

- Data Fidelity: Resolution of data
  - High fidelity data is more accurate, but expensive
  - Machine learning improves the accuracy of low fidelity data, translating it to high fidelity with less cost
  - High and low fidelity are relative



Data Generation and Collection - 27

#### **Indentation Database**

Nanoindentation database summary:

Data Modality: Experimental and simulation data collected for each material

Material	Experiment	Computation			
Al-6061 alloy	7 experiment	2D FEM (Axisymmet each for conical inder	<b>ric)</b> : 100 simulations Iter half angle of		
Al-7075 alloy	7 experiment	50,60,70,80° <b>3D FEM:</b> 15 simulations for Berkovich indenter			
3D printed Ti-6Al-4V alloys (six samples)	144 experiments for each sample	Not available			

\*Load-displacement curves are available for each experiment and simulation.



Lu, L., et al., Proceedings of the National Academy of Sciences, 2020, 117(13), 7052-7062. Database source: https://github.com/lululxvi/deep-learning-for-indentation

### Working with Noisy Data and Outliers

- Noise in the data is very common
- Source of noise: human error in measurement, sensor fluctuation and so on
- Outliers are data coming from same source but vary significantly from other measurement



How do we know if this outlier to ignore or not?

- Regression model can give idea on the data trend and help identify outliers or noise from the data.
- Type of regression:
  - Least square method
  - Lasso regression
  - Ridge regression
  - Elastic net regression, etc.

#### Challenges: Data

- Data on demand
  - Do not have enough data to run an ML model
  - Produce the data by using physics-based simulations
  - Issue: extensible or adaptive sampling is critical
- Data in hand
  - Use of the historical data, stored in different places and formats created by different software versions
  - Issue: reliability, need to be converted to metadata
    - Non-standard formats without proper access (op2, d3plot, bdf, ...)
    - Non-uniform data (shell, solid, mesh, time series, text, ...)
    - Inconsistent data (1d, 2d and 3d mixed)
    - Highly dirty data (oscillations, instability, ...)
- Data in flight
  - Internet of Things (IoT) sensors: large amounts of fast data from operation
  - Issue: volume and quality of data

Vendor/Product	File Format
ANSYS	.cdb, .rst, .rth, .rfl
ABAQUS	.inp, .fil, .odb
CFD General Notation System	.cgns
CFX	.res
Ensight	.case, .encas
ESI PAM CRASH	.ERFH5 files
Fe-Safe	.fer or .csv
femfat	.dma
Fluent	.msh, .cas, .dat
LS-DYNA	.key, .d3plot
MSC.MARC	.t16, .t19
NASTRAN	.OP2, .BDF
nCode	.unv, .csv
OpenFoam	.ControlDict
OptiStruct	.OP2
Pro/Mechanica	design study files

Vendor/Product	File Format
SDRC	universal files
StarCCM+	.ccm
Tecplot	binary <mark>f</mark> iles
Neutral format	STL
Catia V4	*.model
Catia V5	*.CATPart, *.CATProduct
Catia V6	*.3D XML
CGR	*.cgr
Creo - Pro/E	*.prt, *.asm
IGES	*.igs
Inventor	*.iam, *.ipt
JT	.jt
Siemens /NX	*.prt
Solid Edge	*.par, *.asm, *.psm
SolidWorks	*.sldprt, *.sldasm
STEP	*.stp, *.step

#### From Model-centric to Data-centric AI



<u> https://www.youtube.com/watch?v=06-AZXmwHj</u>

#### Homework #1: Data Visualization (1)

- Diamond
  - Fig.1.12 price vs. (a) carat, (b) separated by cut



Figure 12 Diamond price vs carat (a) parameters combined (b) separated by clarity.

#### Homework #1: Data Visualization (2)

- Moneyball
  - Fig.3.10 (a) RS vs. BA, OBP, SLG, OPS, (b) W vs. BA, OBP, SLG, OPS

