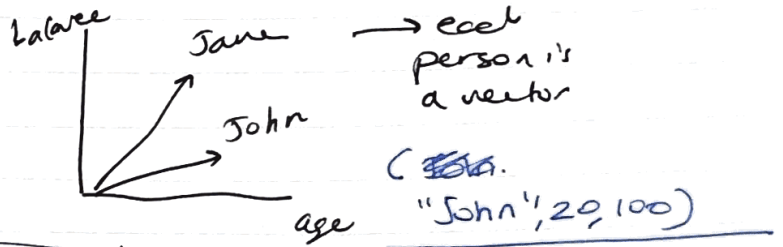- Introduction to Data Science
  - how we represent data linked to info we draw

★ | records | → ( m-dimensional points / vectors ) → just m tuples

ex. name, age, balance



balance

Jane → each person is a vector

John

( ex.
"John", 20, 100)

age

★ | Graphs | → nodes connected by edges

→ each node represented by specific col. & row

ex.

social networks



encode nodes & connection

Adjacency matrix → symmetric = redundancy

```
   (1)(2)(3)
(1)( 0  1  1 )
(2)( 1  0  1 )
(3)( 1  1  0 )
```

0 → no connection between node & itself

each node is specific col / row of matrix

※ no connection with itself

※ generally nodes are not connected to themselves



more space efficient

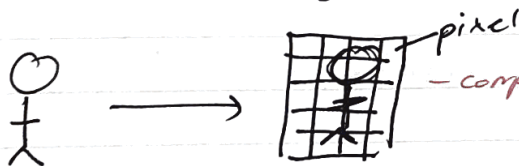Adjacency list
1: {2,3}
2: {1,3}
3: {1,2}

# list the nodes it's connected to

— how info flows is directly a function of those edges / which nodes break up the network if it's disconnected?

★ Data Representation — | Images |

- greyscale values (0, 1)
- color rgb values



pixel

— computer = matrix of pixels

interested in the context

★ Data Representation — | Text |

↳ split it up into words/sentences? → list of words

ex. list of characters: DNA sequence

→ (roots) of words that capture same concept

★ Data Representation — | Time Series |
  ↳ data @ specific intervals of time
anything of the above

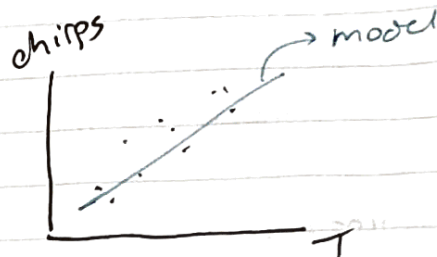ski, skiing, skied

all same concept

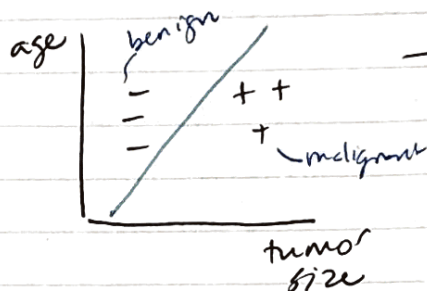· Types of learning (Supervised vs. Unsupervised)

**Supervised Learning**

| chirps | temp |
|--------|------|
| 10 | 40 |
| 5 | 37 |
| 17 | 53 |
| 55 | 103 |
| 40 | 78 |

chirps

model → T

goal:
estimate chirps from T
estimate T from chirp

*classification : relationship between two variables

| age | size | mal. |
|-----|------|------|
| 20 | 12 | 0 |
| 23 | 15 | 1 |
| 47 | 20 | 1 |

age

benign
+ +
+
malignant

tumor size

→ can classify as malignant or benign

- get a new point -
you can say w/ reasonable certainty what class it belongs to

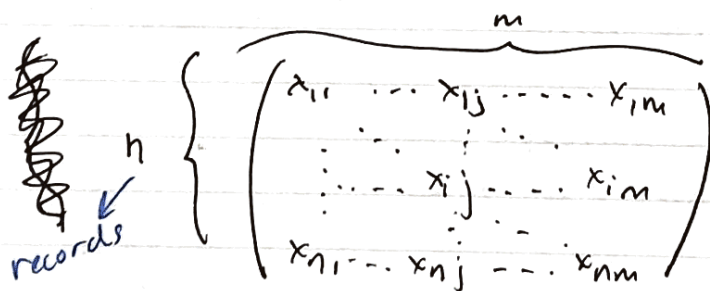**Unsupervised Learning** → goal: not to predict/classify

↳ study data structure — can you identify patterns?

Ex. clustering → try to group points together → find a higher level feature

Dataset: collection of articles, are those articles covering the same topics?

**◊ Distance & Similarity**

Data:

n records

$$\begin{pmatrix} x_{11} & \cdots & x_{1j} & \cdots & x_{1m} \\ \vdots & & \vdots & & \vdots \\ \cdots & & x_{ij} & \cdots & x_{im} \\ \vdots & & \vdots & & \vdots \\ x_{n1} & \cdots & x_{nj} & \cdots & x_{nm} \end{pmatrix}$$

m

- each column is a feature
(m distinct attributes)

feature space: all possible values for set of features in data

· Distance — to uncover structure from data,
　　 need a way to compare data points

　( dissimilarity function — takes two objects
　　　　　　　　　　 returns LARGE value for ↑
　　　　　　 dissimilarity (w/ respect to the function)

"d" is a distance function iff:
①  $d(i,j) = 0$ iff $i=j$

②  $d(i,j) = d(j,i)$　　# symmetry

# triangle inequality ↙ ↖ ③ $d(i,j) \leq d(i,k) + d(k,j)$
- if you go through
third point
distance between i,j
is necessarily smaller
over i,k → k,j

· Why a distance function?
　# it's (intuitive) → always want
　　　　　↗　　　　　 to get to a place
　The　restrictions　where it's more
　extra　　　　　 intuitive
　tend to this

Minkowski Distance — For $x, y$ points in $d$-dimensional space
　　　 d# of feature　　　　　　 d# of feature　↳ d attributes
　$x = [x_1, \ldots x_d]$　and　$y = [y_1, \ldots y_d]$

　　　　　　　　　　　　　　　　　　　　　, and raise to $p^{th}$ power
P ≥ 1　　$h_p(x,y) = \left( \sum_{i=1}^{d} |x_i - y_i|^p \right)^{1/p}$ → take the
↙　　　　　　　　　　　　　　　　　　　　　　　　　　　　 $p^{th}$ root
p is a　　　　　　　　　　　　　　　　　　　　　　　　　　 of the
parameter　　 Minkowski　　　　　　　 take　　　　　　 thing
(something to　 distance　　　　　 pairwise
tune)　　　　　　　　　　　　　 difference between
　　　　　　　　　　　　　　　　　　　 $x_i$ and $y_i$

when $p=2$ → Euclidean Distance
when $p=1$ → Manhattan Distance

Ex. $d=2$, $(p=2)$ (Euclidean distance)
　　　　　　 square root →p=2 (Pythagoren!!)
　　　　　　 of the sum



$l_p(x,y) = \left( \sum_{i=1}^{d} |x_i - y_i|^p \right)^{1/p}$

$p=1$
　　　　　　　 sum=2 → distance
　　　　　　　 ↳ sum　 must go through (Manhattan
　　　　　　 along the　　 the grid　　 Distance)
　　　　　　　 grid

Ex.  $d = 3 \rightarrow$ add another term to the sum

$\boxed{P = 2}$

B $(1,1,1)$

$\rightarrow$ length ?

$\triangleq \sqrt{3}$

$(1^2 + 1^2 + 1^2)^{\frac{1}{2}} = \sqrt{3}$

$(0,0,0)$

$\boxed{P = 1}$

B

$\rightarrow$ length ?

$(1 + 1 + 1)^{\frac{1}{2}} = 3$
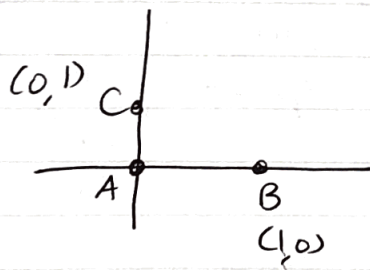
Is $L_p$ a distance function when $0 < p < 1$ ?

— no because smaller distances would have greater value

recall : important axiom is triangle inequality

↳ proof by contradiction (counter-example)

$(0,1)$ C

A    B
  C $(1,0)$

$D(B, A) = D(A, C) = 1$ ✓

$D(B, C) = 2^{\frac{1}{p}}$

But... if $p < 1$, then $\frac{1}{p} > 1$

So $D(B, C) > D(B, A) + D(A, C)$

which violates the triangle inequality

• Cosine similarity — Similarity function gives larger #'s for more similar objects

$s(x, y) = \cos(\theta)$

↳ $\theta$ is angle between $x, y$

... points further apart have very large angle $\theta$ (less similar)

↳ to get dissimilarity ...

$d(x, y) = \dfrac{1}{s(x, y)}$     or   $d(x, y) = k - s(x, y)$ for some $k \longrightarrow$ shit is ten here

↳ high similarity

$\cos(0°) = 1$   $\boxed{k = 1}$   $1 - 1 = 0$ ↳ no dissimilarity when same point

full request ↳ not many ASAP

→ when should you use cosine (dis)similarity over euclidean distance?

└ when **not** interested into the magnitude of your vectors

## Jaccard Similarity → represents documents

|   | $w_1$ | $w_2$ | $\cdots$ | $w_d$ |
|---|---|---|---|---|
| $x$ | 1 | 0 | $\cdots$ | 1 |
| $y$ | 1 | 1 | $\cdots$ | 0 |

↓                    ↘ one has it

if they have same word it's zero            └ 1

└ apply distance     ↗ manhattan distance

$$L_1(x,y) = \sum_{i=1}^{d} |x_i - y_i| \quad p = 1$$

└ count # of words that the documents differ by

→ only be 1 if $x_i \neq y_i$

BUT... consider...

|   | $w_1$ | $w_2$ | $w_3$ | $w_4$ |
|---|---|---|---|---|
| $x$ | 1 | 1 | 0 | 1 |
| $y$ | 1 | 1 | 1 | 0 |

only differ by last 2 words

|   | $w_1$ | $w_2$ |
|---|---|---|
| $x$ | 0 | 1 |
| $y$ | 1 | 0 |

completely different

— manhattan distance for BOTH = 2

— need to account also for SIMILARITY (intersection)

★ Jaccard Similarity → accounts for size of intersection

$$J_{sim}(x,y) = \frac{|x \cap y|}{|x \cup y|} \quad \text{→ ratio between intersection and the union}$$

if $x = y$,
union = intersection
$J_{sim} = 1$

if $x \neq y$
intersection = 0
$J_{sim} = 0$

} large value if similar

$$J_{Dist}(x,y) = 1 - \frac{|x \cap y|}{|x \cup y|}$$

★ try to prove this is a distance function

prove
$$d(i,j)$$
$$\leq$$
$$d(i,t) + d(t,j)$$