



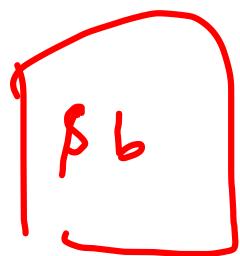
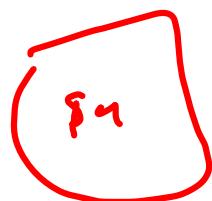
scrapping

Basics ml

Note

Travago

your



`bsr.findall('a')`



`bsr.find('div').findall('a')`

}

`attrs = { 'attr': value }`

`name`

`<a>`

Datagram

columns

values

```
<html>
  > <head>...</head>
  > <body>
    > <table>
      > <thead>
        > <tr>
          <td>Alfreds Futterkiste</td>
          <td>Maria Anders</td>
          <td>Germany</td>
        </tr>
      </thead>
      > <tbody>
        > <tr>
          <td>1 Centro comercial Moctezuma</td>
          <td>1 Francisco Chang</td>
          <td>1 Mexico</td>
        </tr>
        > <tr>
          <td>2 Centro comercial Moctezuma</td>
          <td>2 Francisco Chang</td>
          <td>2 Mexico</td> == $0
        </tr>
      </tbody>
    </table>
  </body>
</html>
```

```
<html>
  <head>...</head>
  <body>
    <table>
      <thead>
        <tr>
          <td>Alfreds Futterkiste</td>
          <td>Maria Anders</td>
          <td>Germany</td>
        </tr>
      </thead>
      <tbody>
        <tr>
          <td>1 Centro comercial Moctezuma</td>
          <td>1 Francisco Chang</td>
          <td>1 Mexico</td>
        </tr>
        <tr>
          <td>2 Centro comercial Moctezuma</td>
          <td>2 Francisco Chang</td>
          <td>2 Mexico</td>
        </tr>
      ...
    </tbody>
  </table>
</body>
</html>
```

values

1 1 1

2 2 2

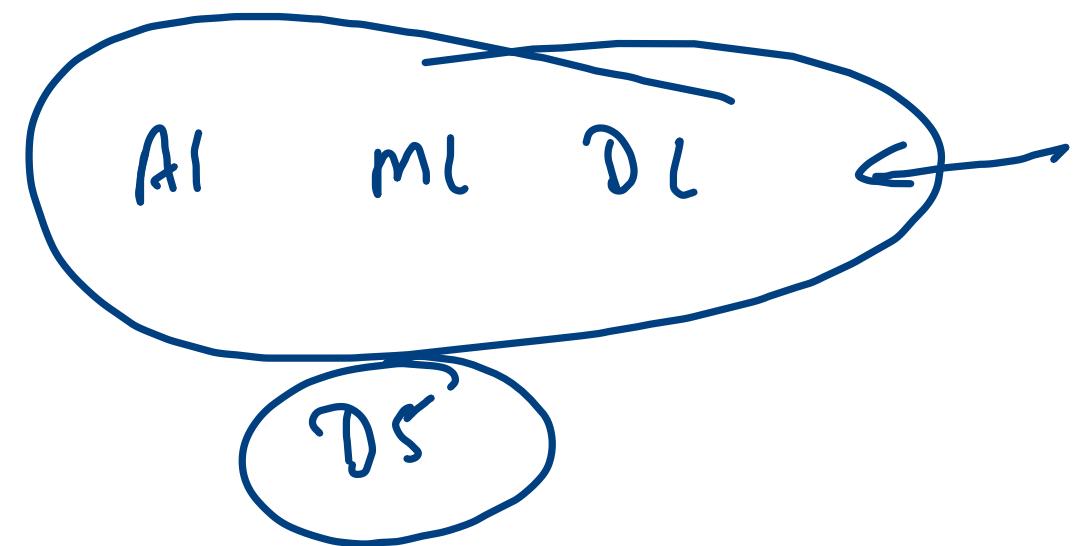
1 2

```
✓ values = bs.tbody  
values
```

```
✓ pd_values = []  
  
for tr in values.findAll('tr'):
    row = []
    for td in tr.findAll('td'):
        row.append(td.text)
    pd_values.append(row)  
  
pd_values
```

```
[[1 Centro comercial Moctezuma', '1 Francisco Chang', '1 Mexico'],
 ['2 Centro comercial Moctezuma', '2 Francisco Chang', '2 Mexico']]
```

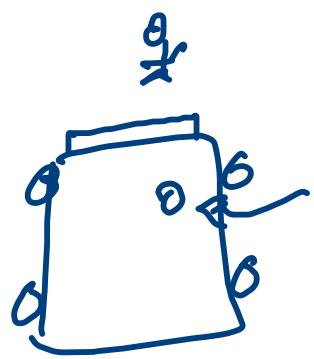
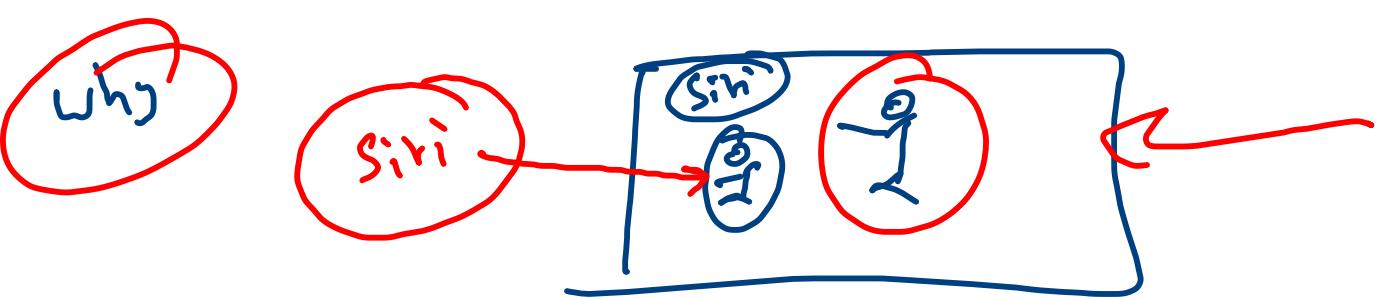
pd_values = [[1 Centro, 1 Francisco, 1 Mexico],
 [2 Centro, 2 Francisco, 2 Mexico]]



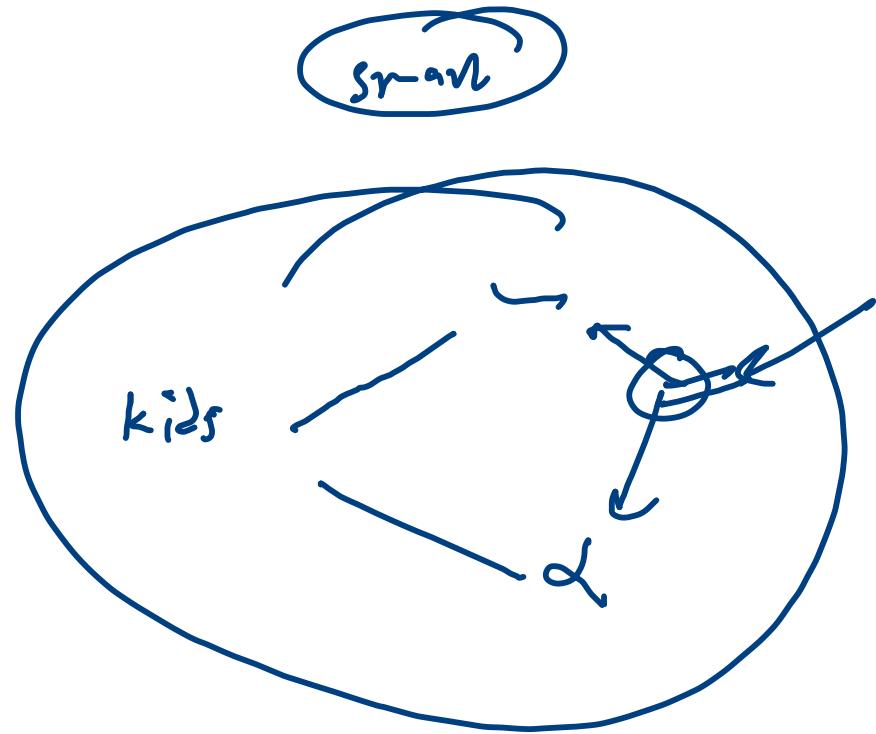
AI → Machine
Intelligence
Artificial
Intelligence

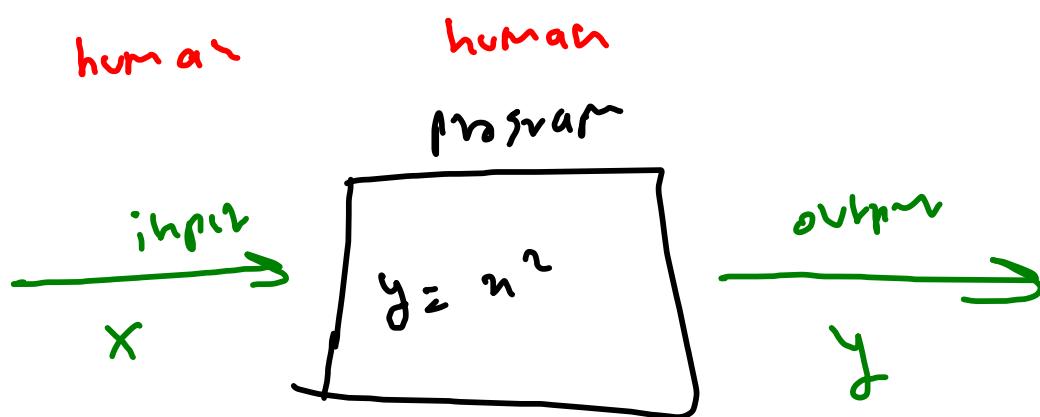
ML → ? stats

DL → neural networks



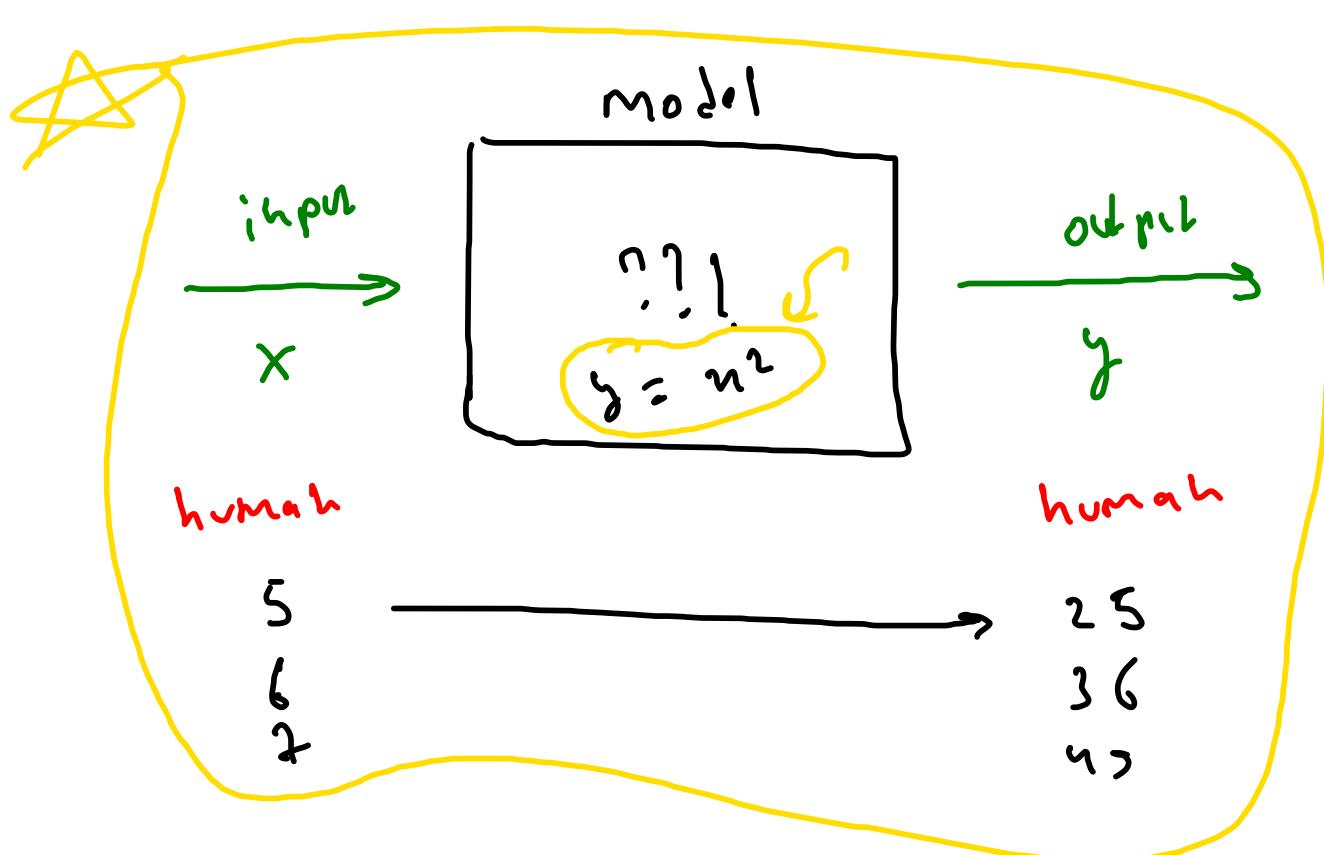
automation





garban is
garban on

DP

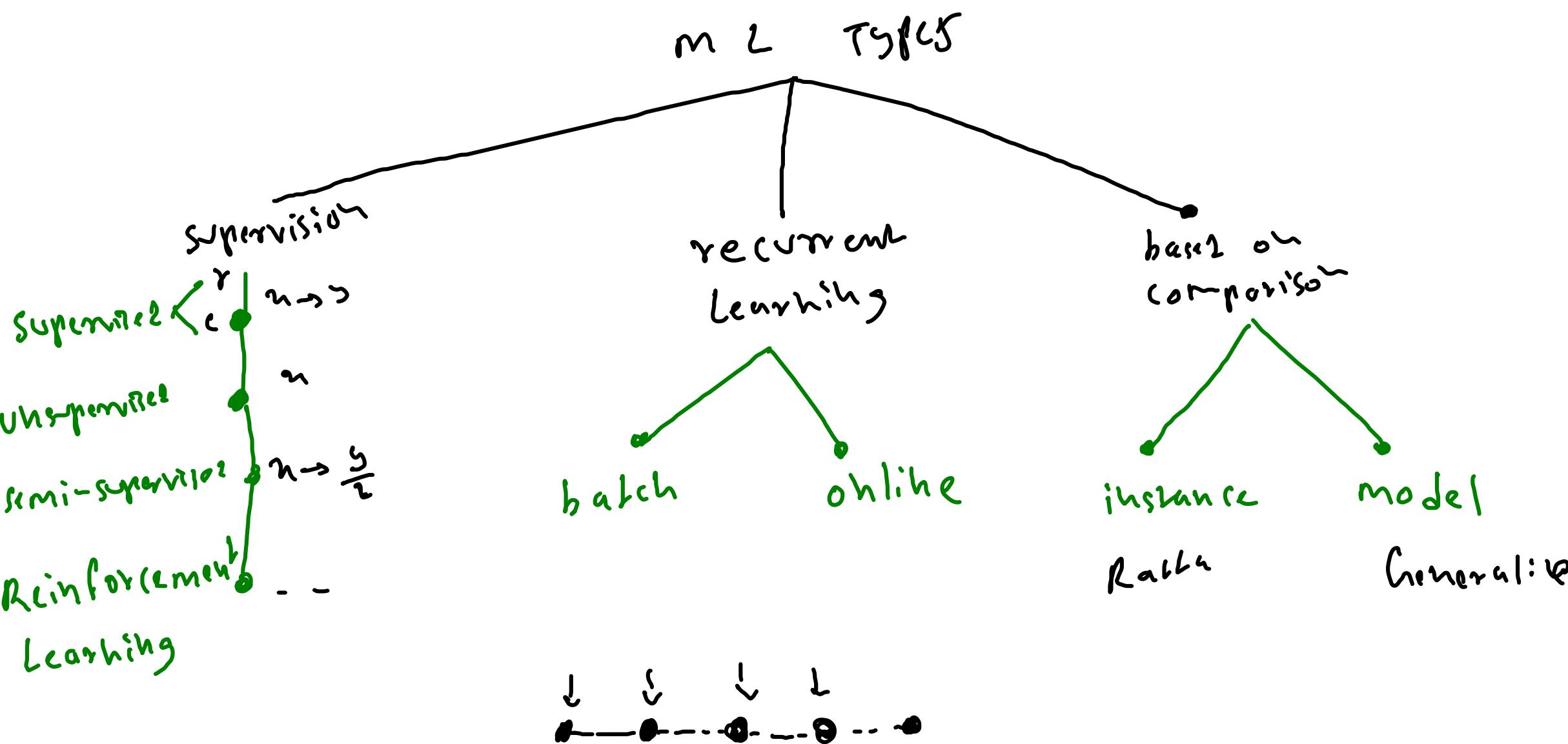


$$x^2$$

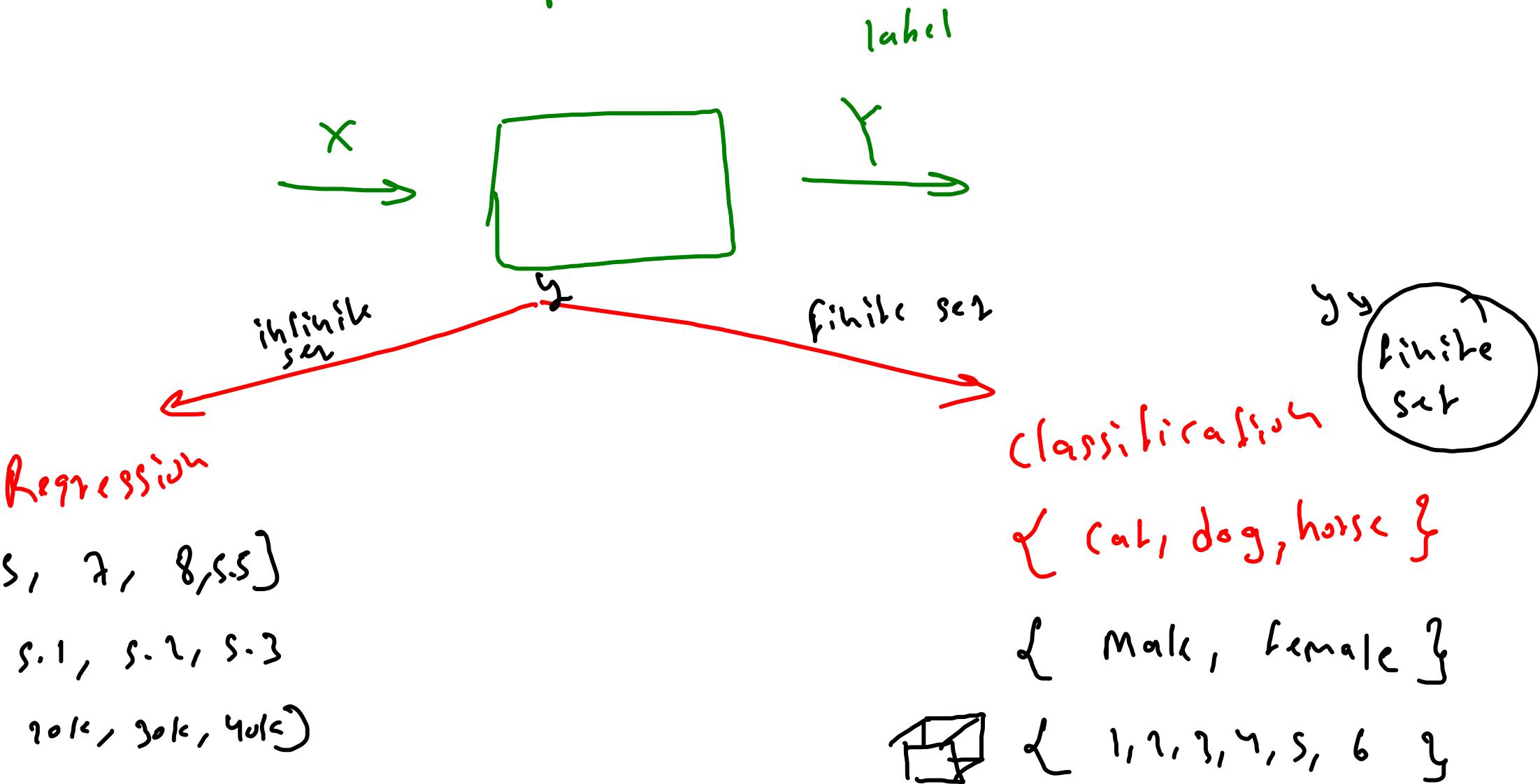
$$20+x$$

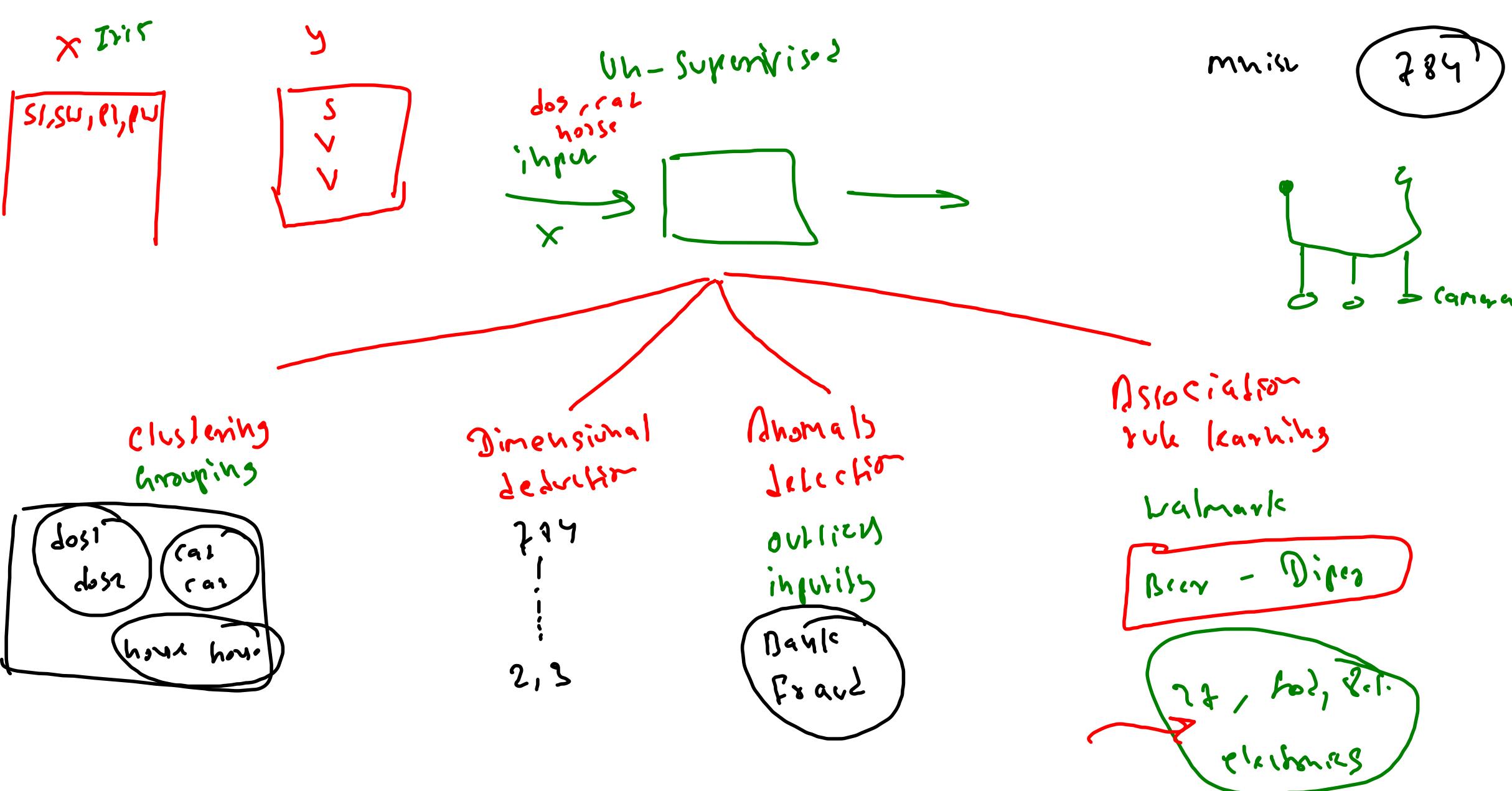
$$30-x$$

ML TYPES

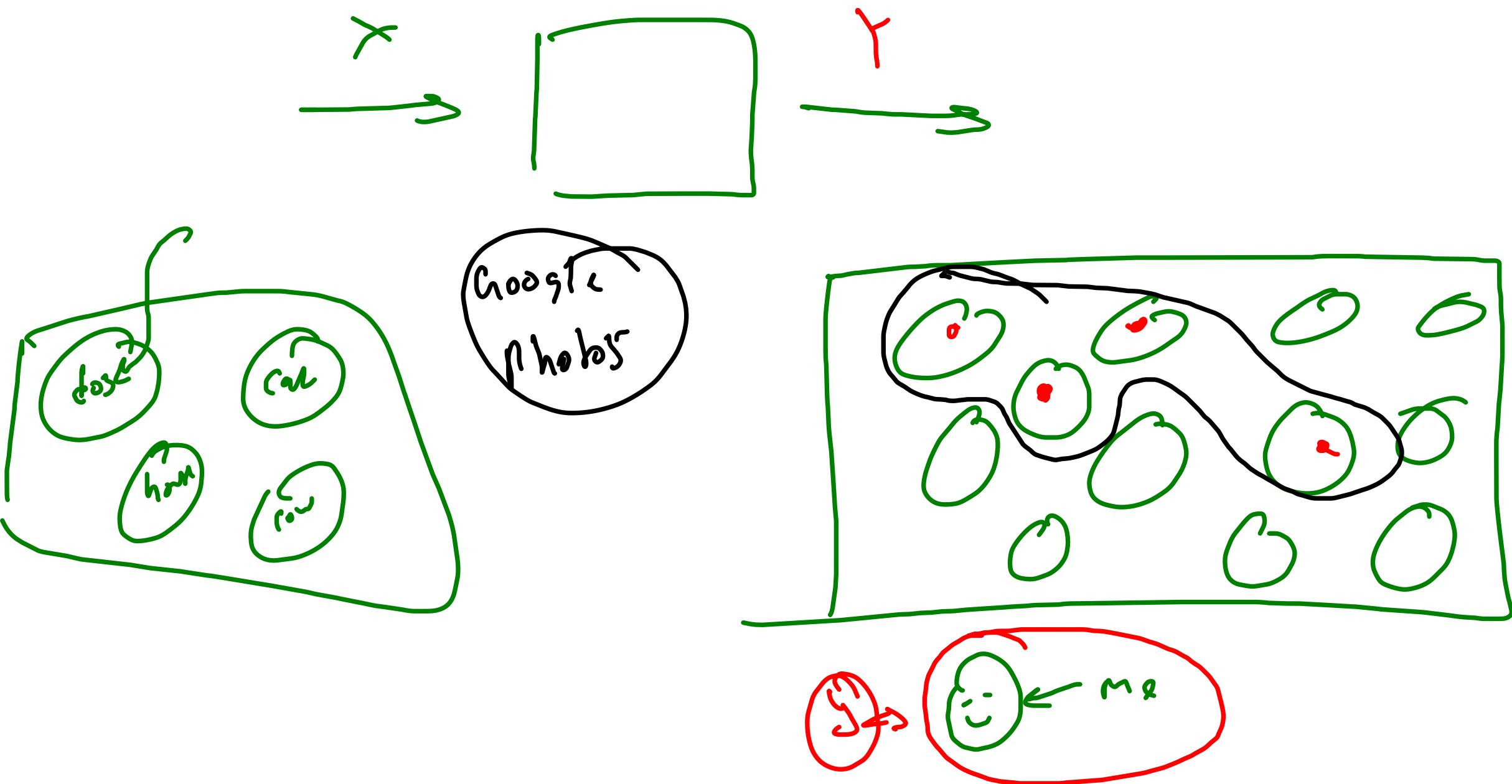


Supervised

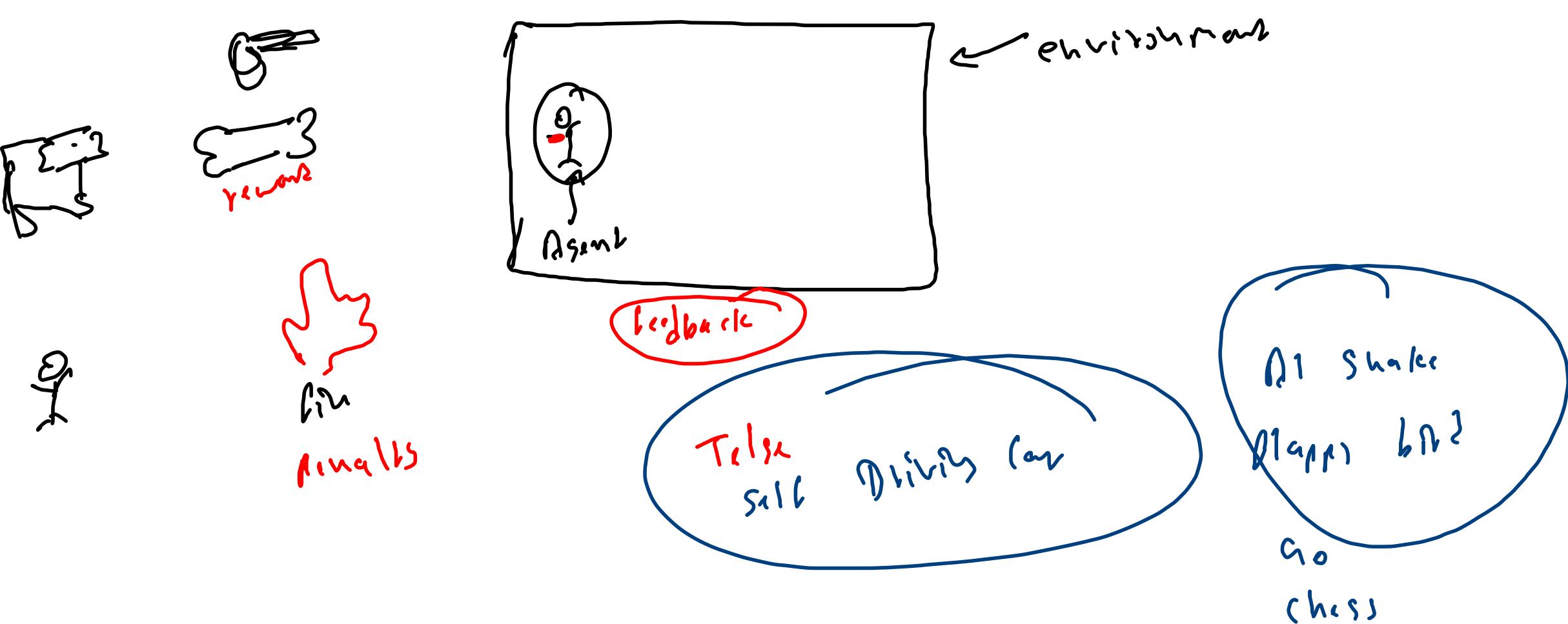




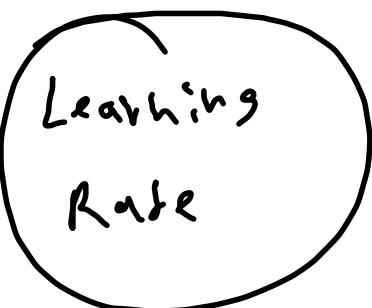
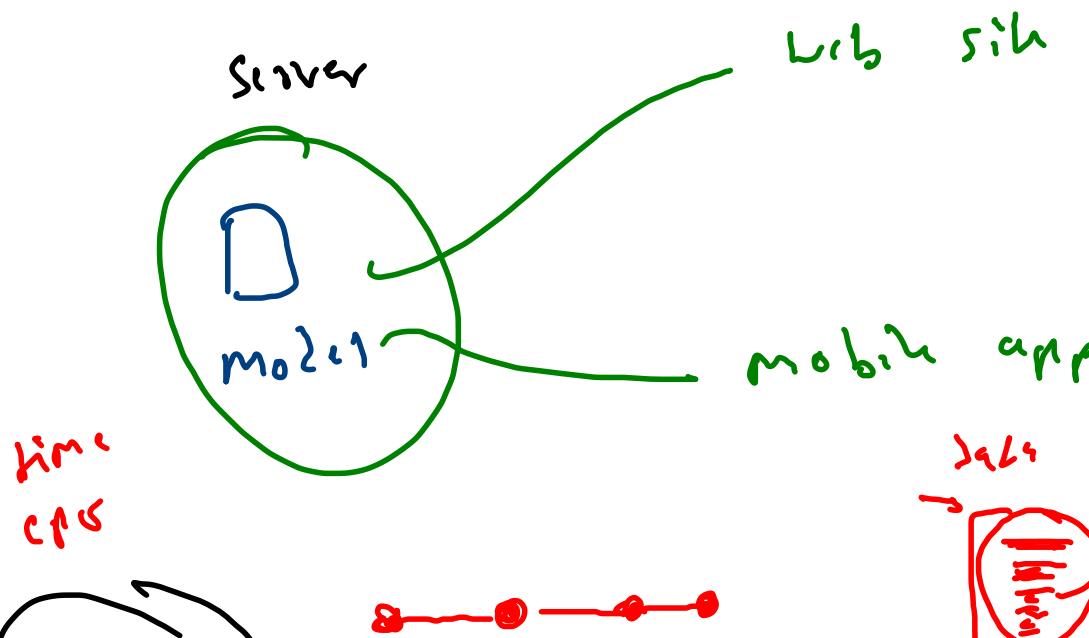
Semi-supervised



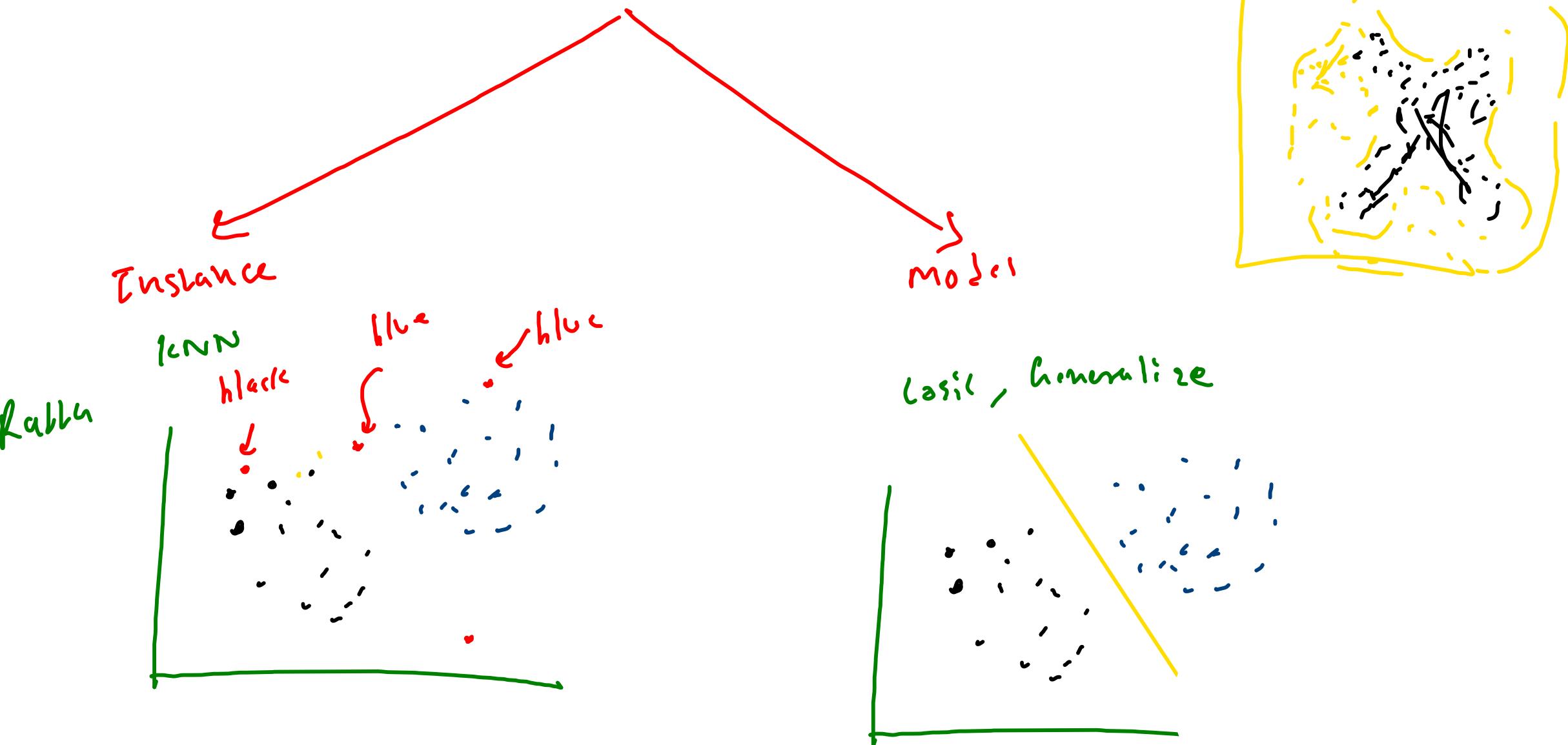
Reinforcement Learning



Recurrent Learning

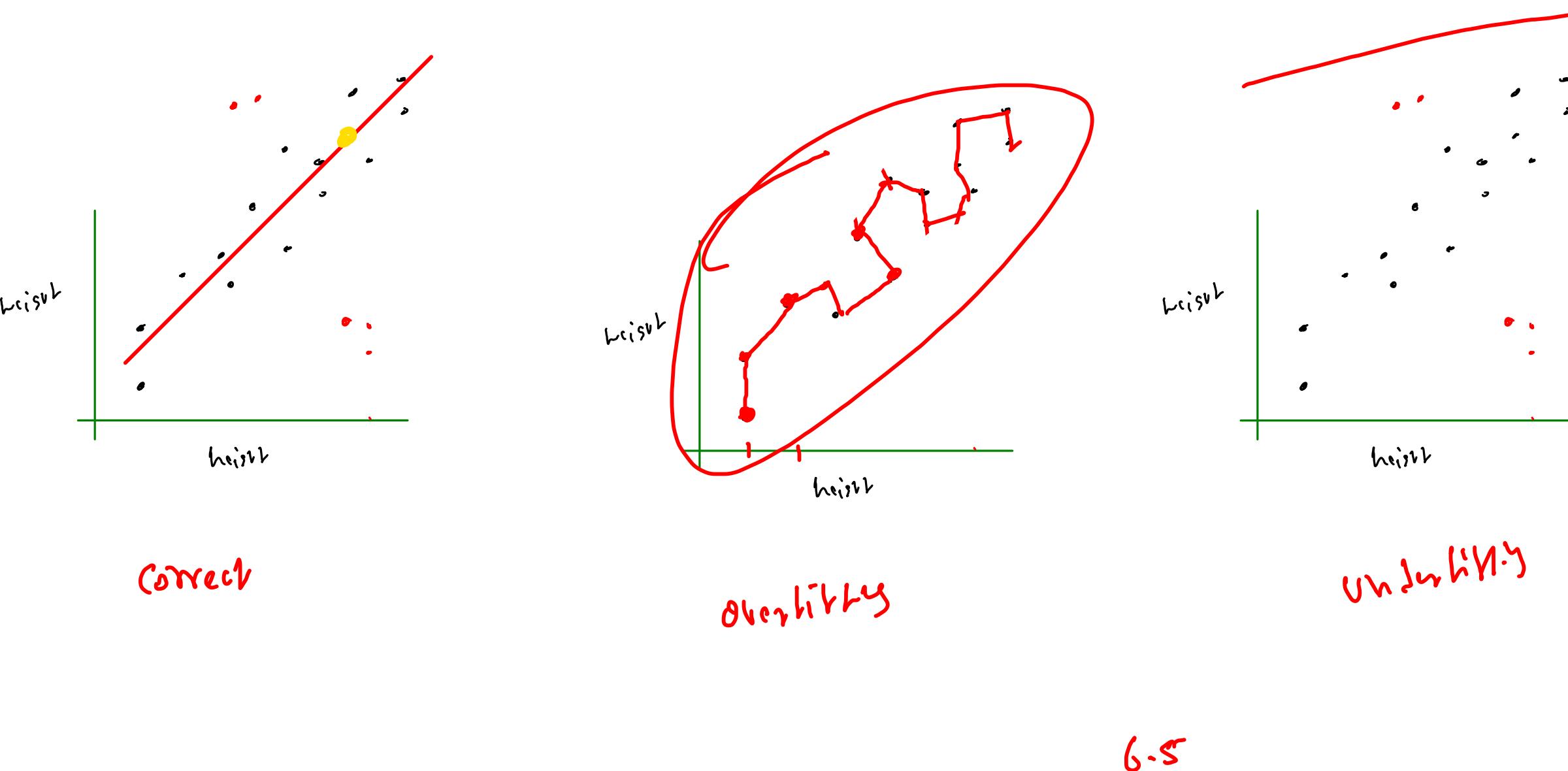


Based on Comparison



challenges

- insufficient data → [specific]
- non-representable data
- poor quality
- irrelevant features
- underfitting
- overfitting



ML DL

, frame problem

• Data gather

• preprocessing \Rightarrow

human readable

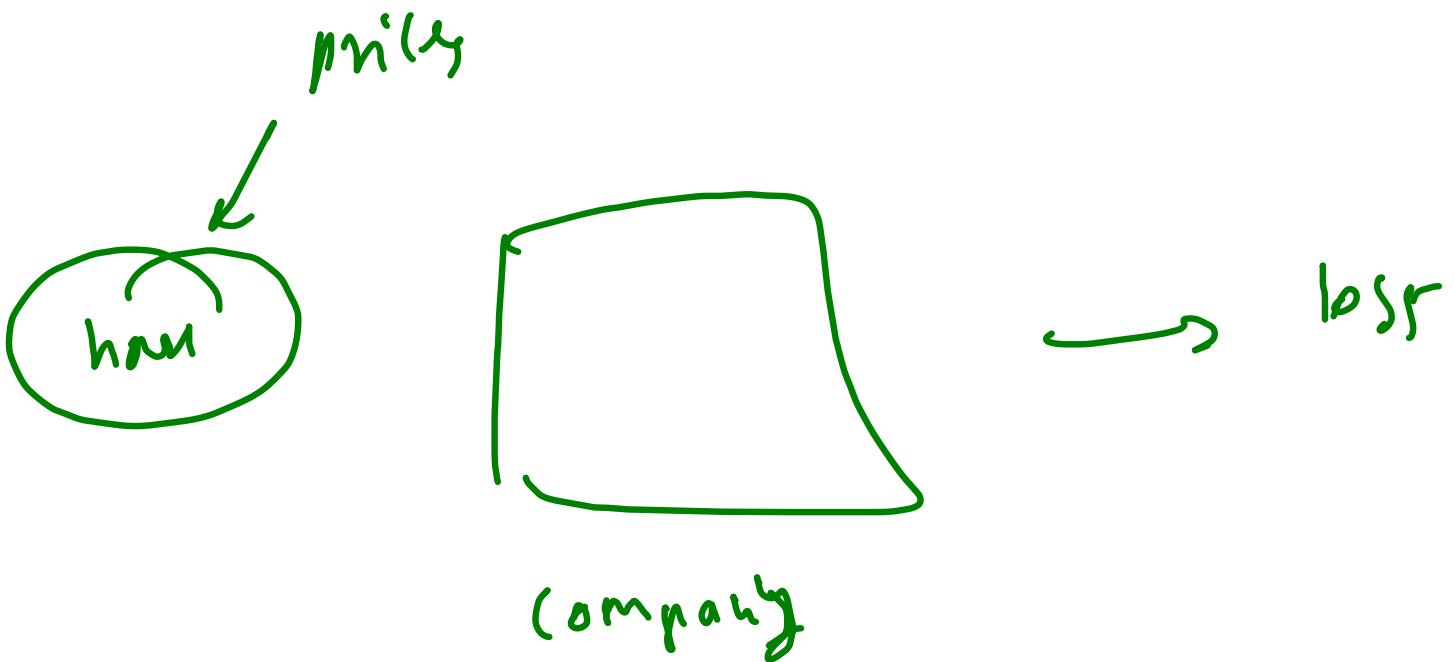
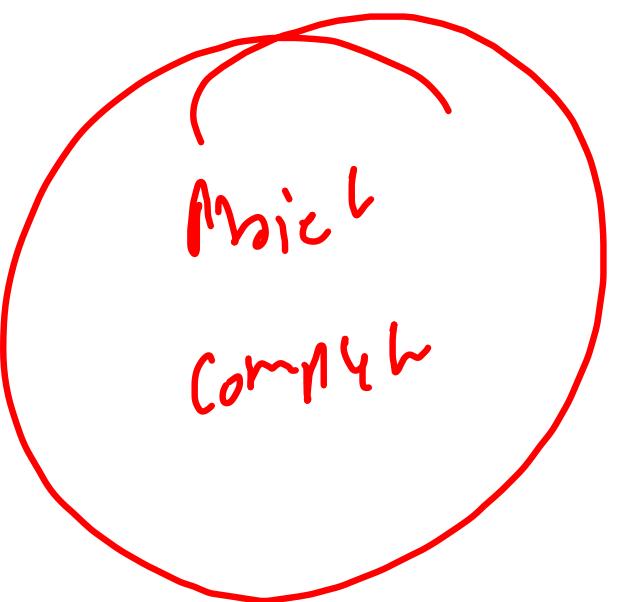
• EDA

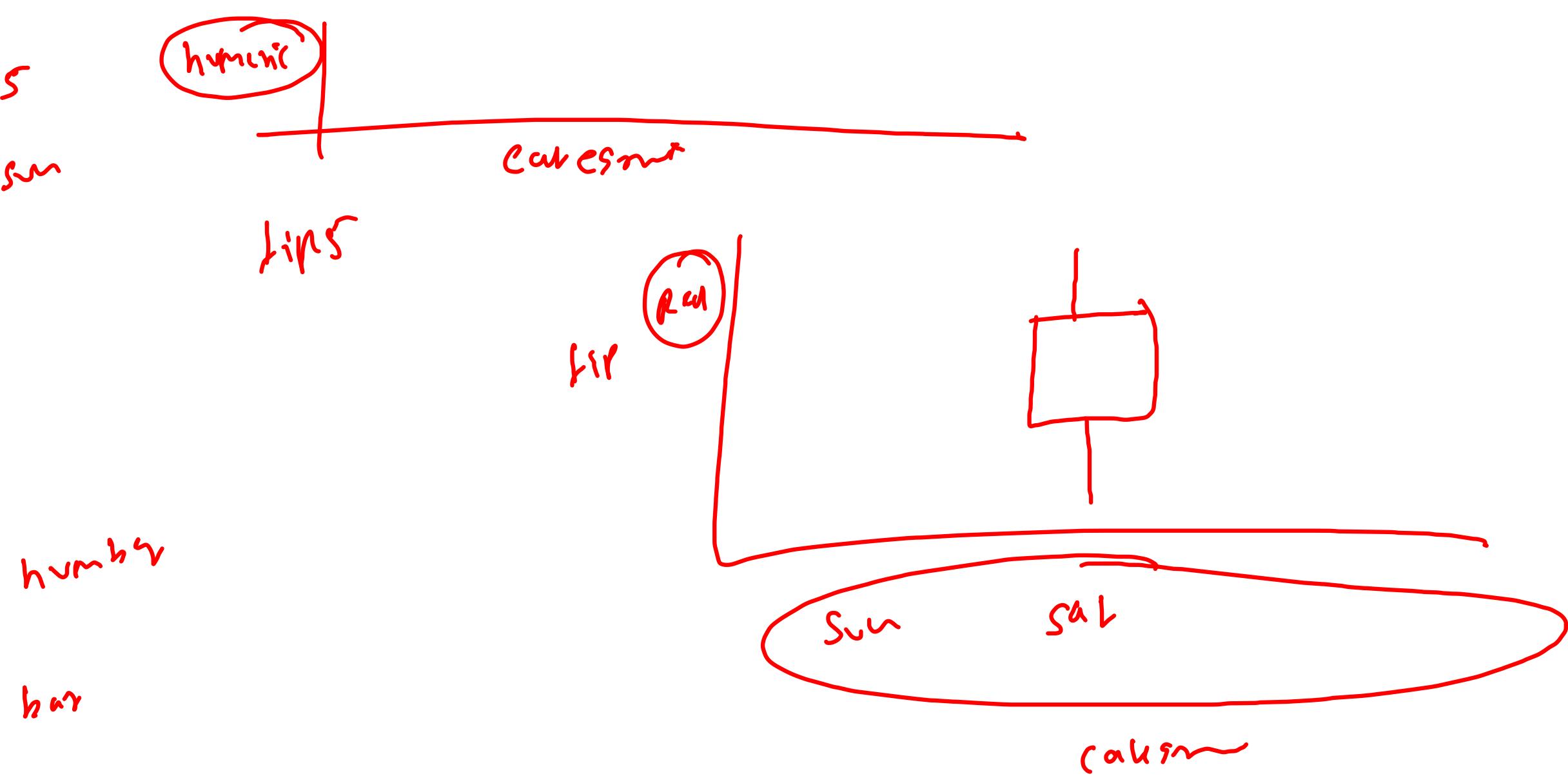
• feature engineering \rightarrow model

• evaluation

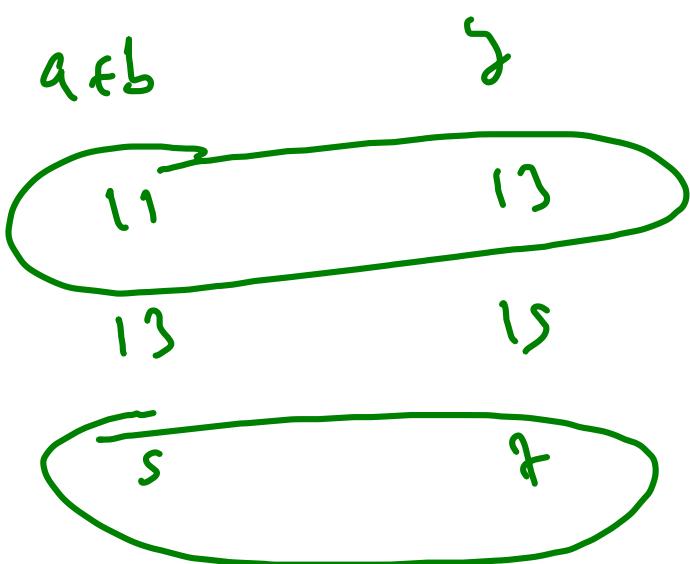
• deployment

• maintenance



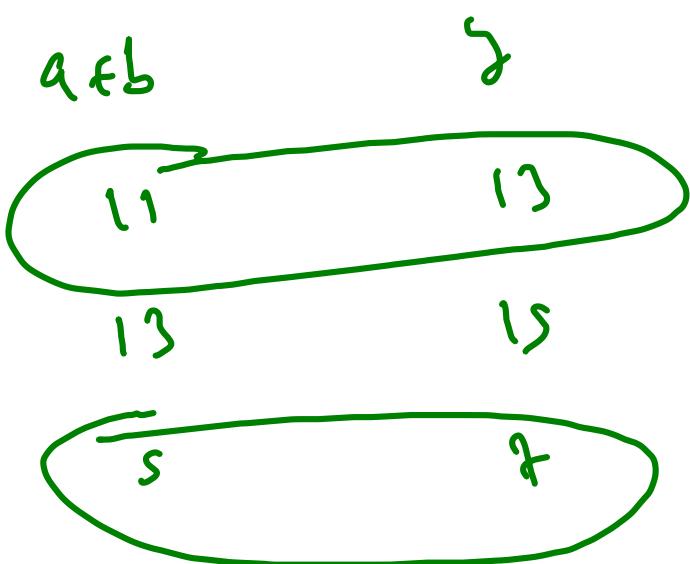


| | | | |
|----|---|-------|----|
| a | b | $a+b$ | y |
| 10 | | 11 | 13 |
| 7 | | 13 | 15 |
| -7 | | 5 | 2 |

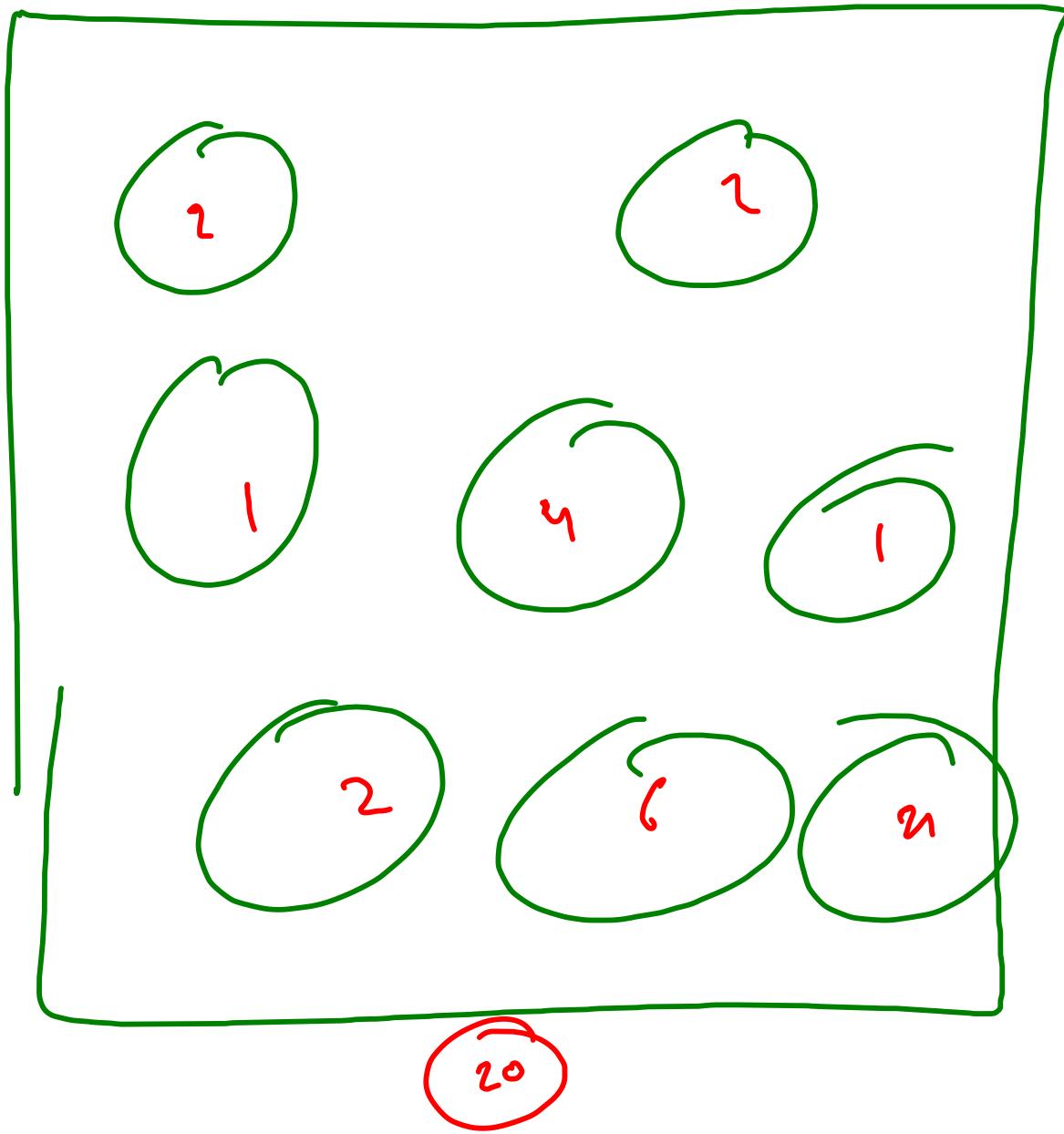
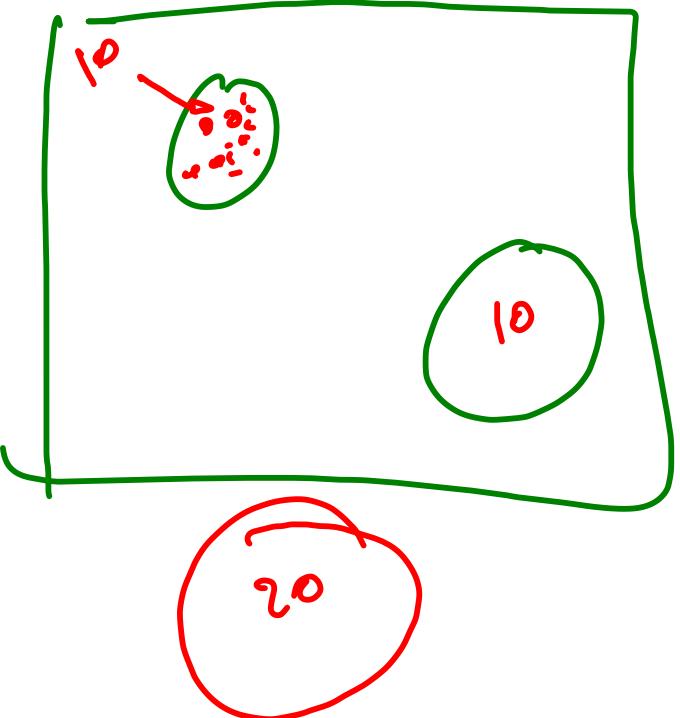


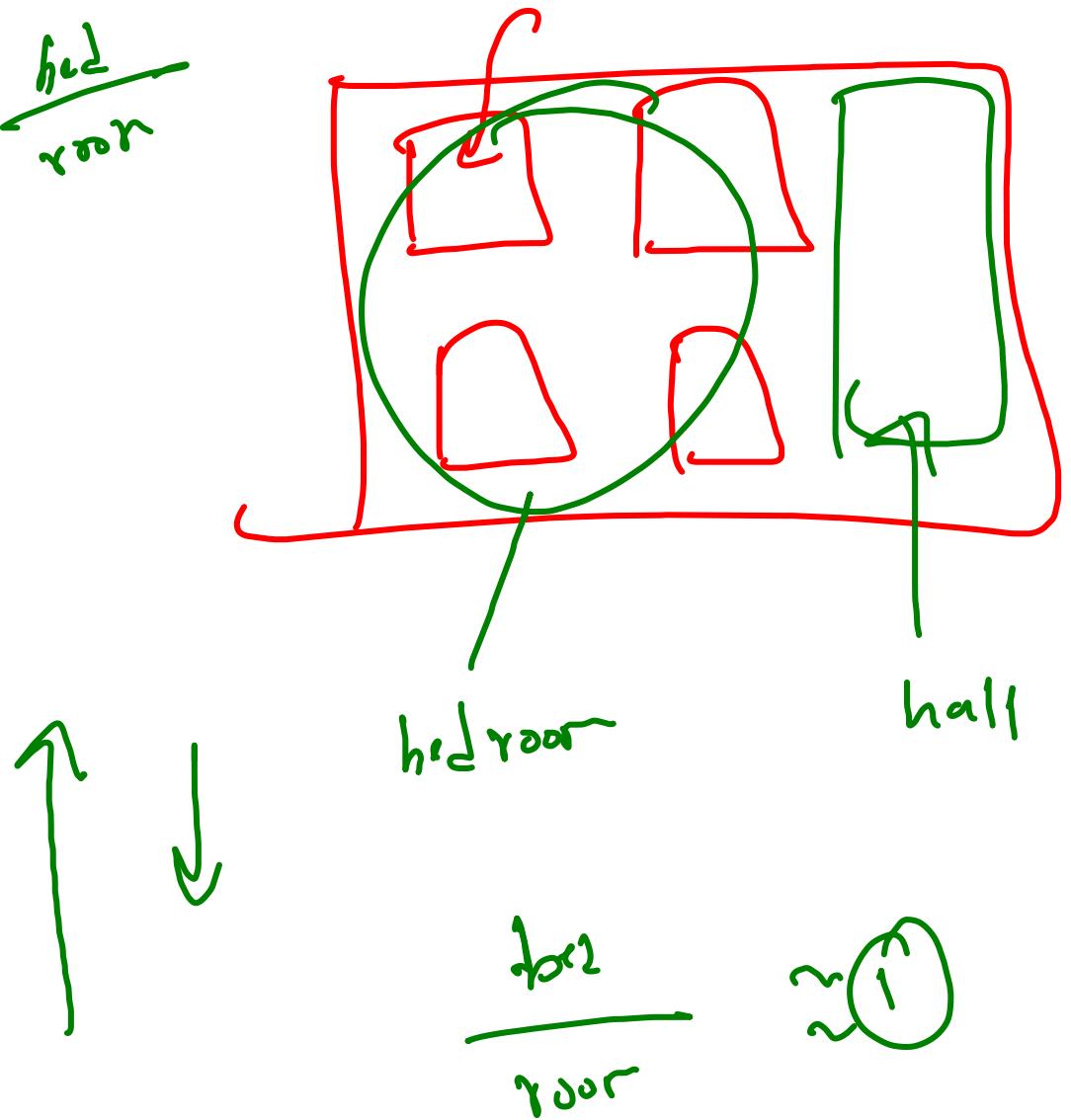
$$\begin{array}{r}
 a+b \quad 2 \\
 5 \xrightarrow{+2} 7 \\
 11 \xrightarrow{+2} 13 \\
 13 \xrightarrow{+2} 15
 \end{array}$$

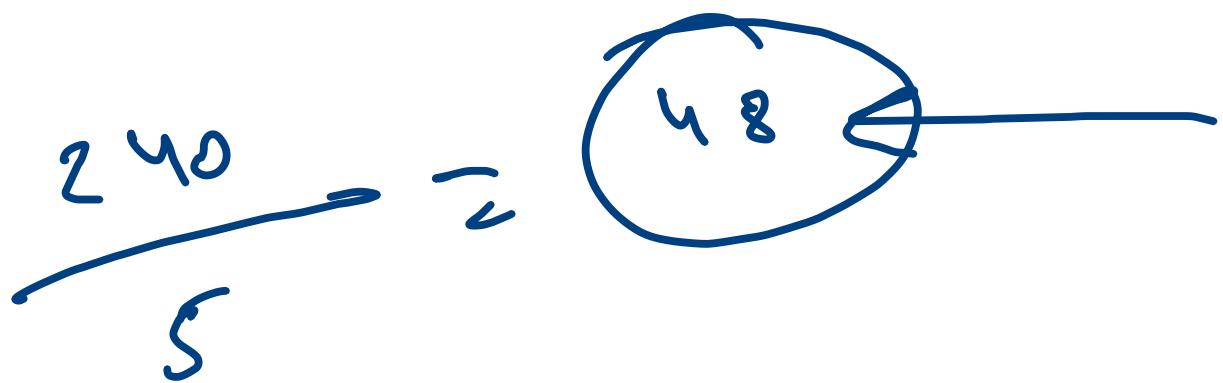
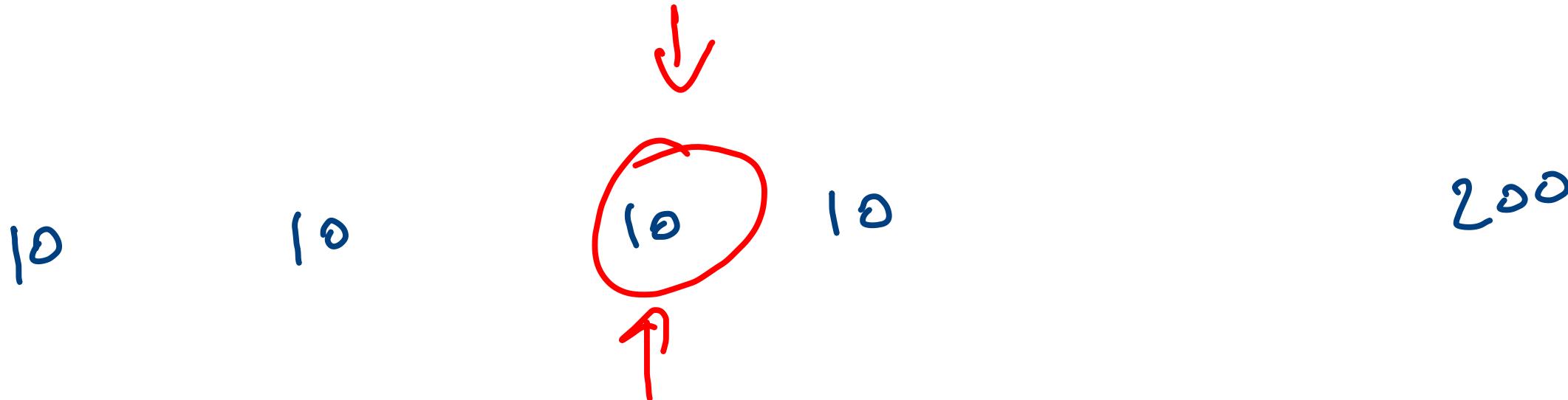
| | | | |
|---|----|-------|----|
| a | b | $a+b$ | y |
| 1 | 10 | 11 | 13 |
| 6 | 7 | 13 | 15 |
| 9 | -7 | 5 | 2 |

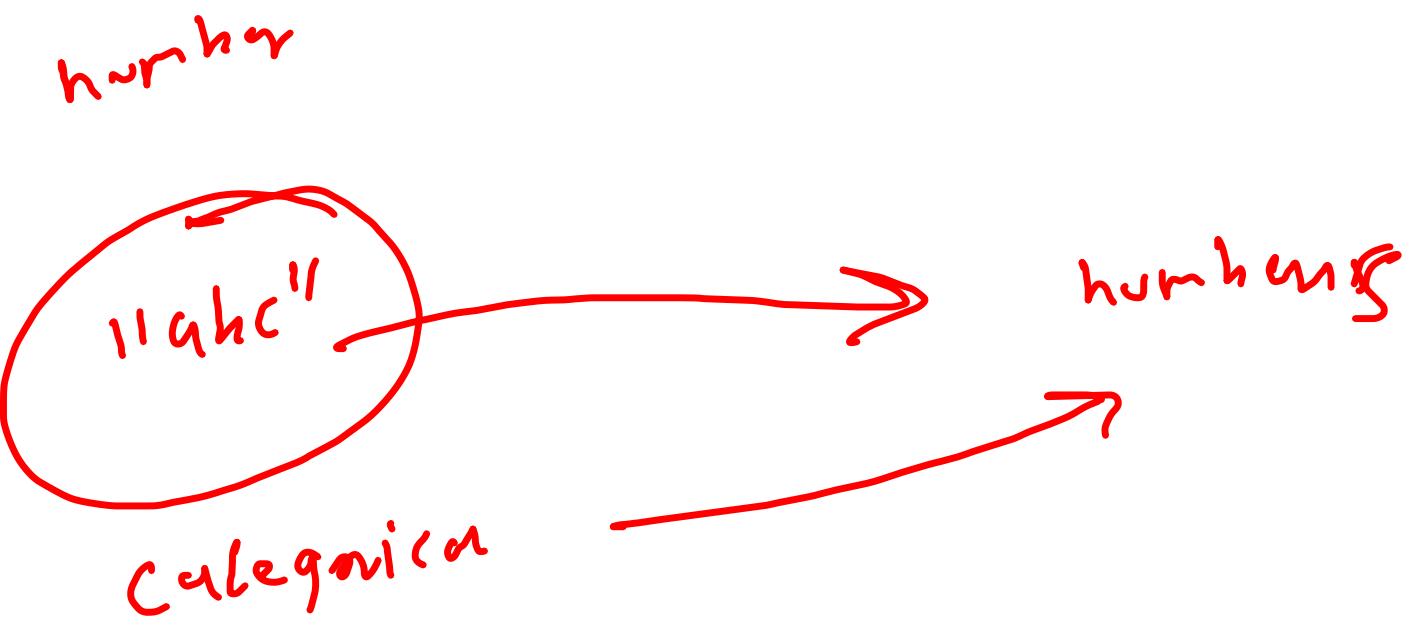


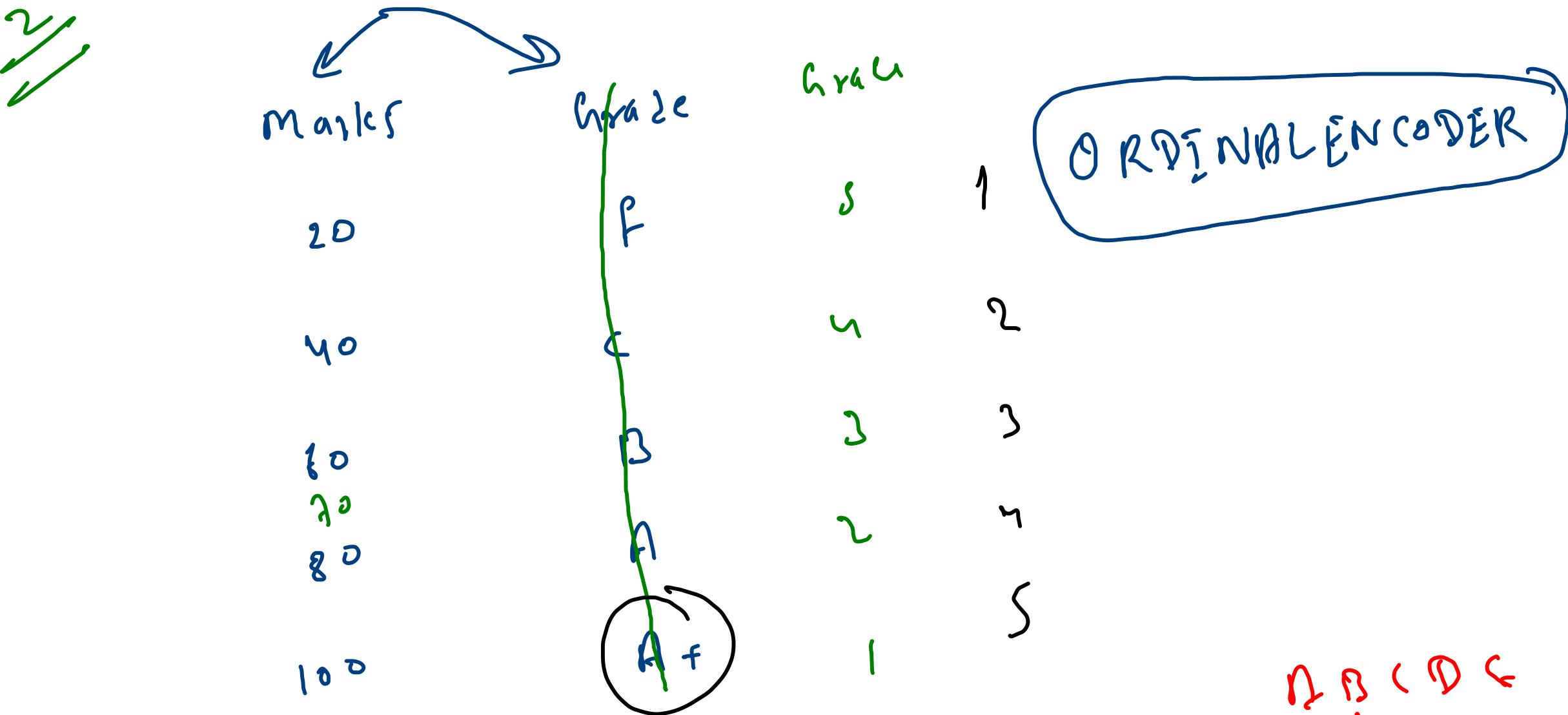
$$\begin{array}{r}
 a+b \\
 5 \xrightarrow{+2} 7 \\
 11 \xrightarrow{+2} 13 \\
 13 \xrightarrow{+2} 15
 \end{array}$$











A B C D E
 ↓
 1 2 3 und
 5 4 3 2 1

ON-NOT ENCODER
giver - receiver

Marks

20

40

60

80

100

State

noa

Mumbai

Kolkata

Delhi

Delhi

| noa | Mumbai | Kolkata | Delhi |
|-----|--------|---------|-------|
| 0 | 1 | 0 | 0 |
| 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 1 |
| 0 | 0 | 0 | 1 |

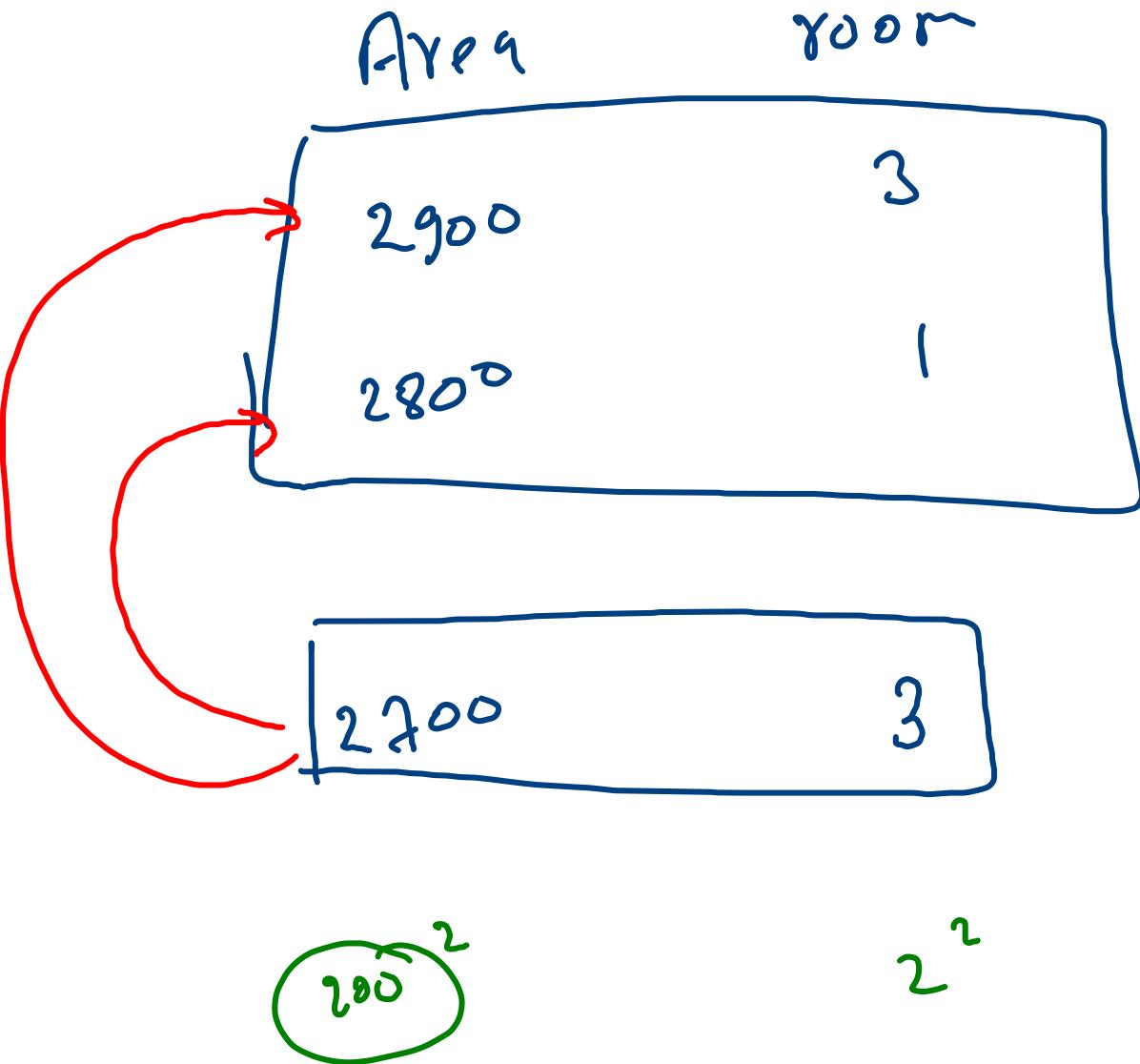
ordination [

[A, B, C], [x, y], [3 0 1 2]

]

2D

| row | col 1 | col 2 | col 3 |
|-----|-------|-------|-------|
| A | x | 0 | L |
| B | 2 | 3 | 1 |
| C | y | 2 | 2 |
| D | z | | |



$$200^2 + 0^2$$

greater

$$100^2 + 2^2$$

smaller

Standardization $\rightarrow M^T$

Normalization $\rightarrow [u, y]$

minmax

- robust scalar

