

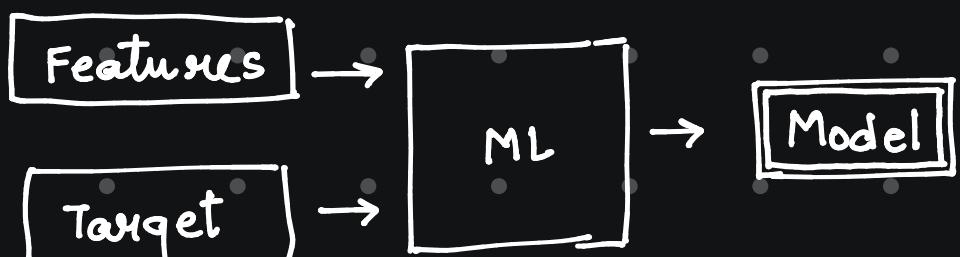
machine learning zoomcamp

1.1 Introduction to machine learning

'FEATURES'

→ age (year)
→ make (BMW...) → 'TARGET'
→ mileage

MODEL TRAINING



DATA → Expertise → PATTERNS
DATA → ML MODEL → PATTERNS

ML → process of extracting patterns from data.

PREDICTIONS

Features + Model → Predictions

1.2 ML vs Rule-Based Systems

Email System

→ SPAM detections

RULES

→ sender = 'promotions@online.com'
→ title contains = tax review and sender domain is 'online.com'

```

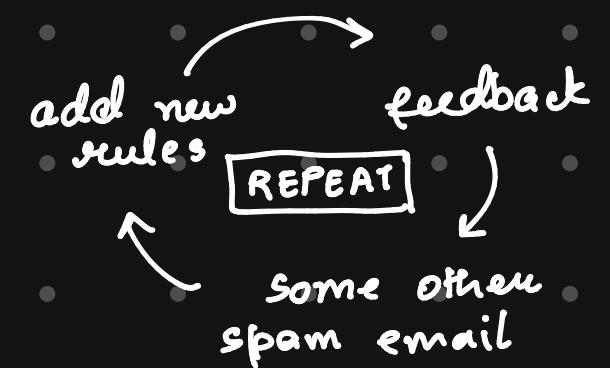
def detect_spam(email):
    if email.sender == 'promotions@online.com':
        return SPAM
    if contains(email.title, ['tax', 'review']) and
       domain(email.sender, 'online.com'):
        return SPAM
    return GOOD
  
```

PROBLEMS

→ other kind of SPAM messages.

↳ add new rules

→ if the Email Body contains a word 'deposit'
↳ if sender domain = 'test.com' then SPAM
↳ if body >= 100 words then SPAM



SOLUTION → use machine Learning

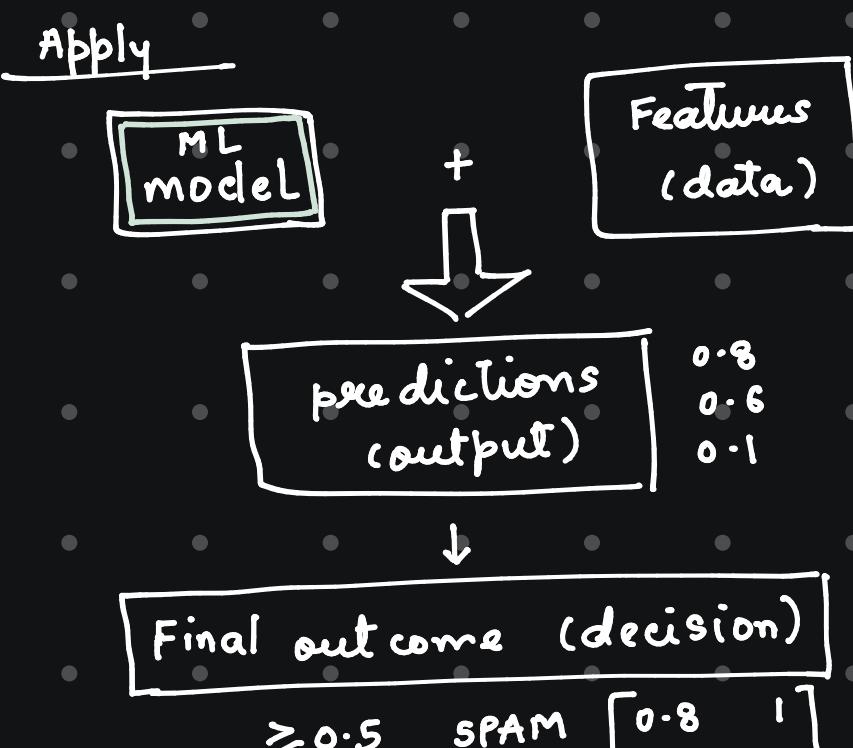
1) Get Data → Button to mark emails as 'SPAM' or 'NOT SPAM'

2) Define and calculate features → email title? body? length, particular Sender? • specific domain, description contains some words like 'deposit'! [COMING FROM THE RULES]

3) Train and use the model

features + target } train an ML model

process of training as 'fitting'



Rule Based Systems

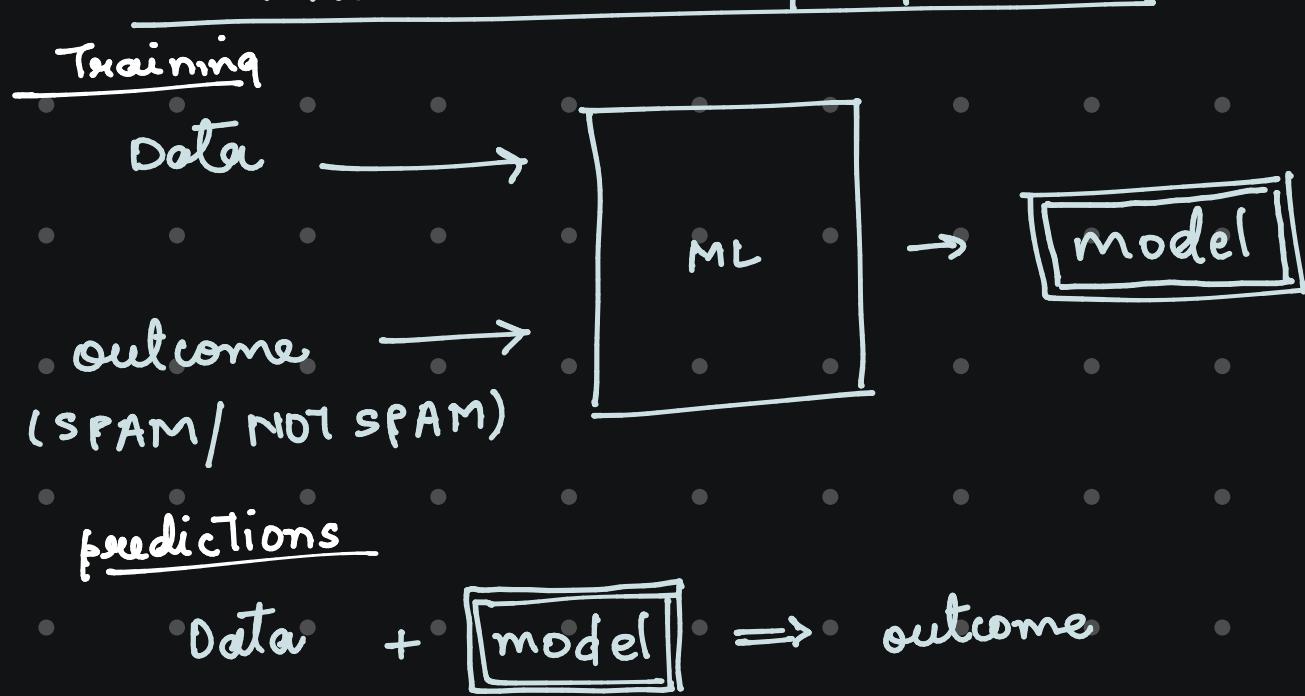
Data → software → outcome
code → software → outcome

* difficult to maintain in the long term

| Features (data) | Target (desired output) |
|--------------------|-------------------------|
| [1, 1, 0, 0, 1, 1] | 1 |
| [0, 0, 0, 1, 0, 1] | 0 |
| [1, 1, 1, 0, 1, 0] | 1 |
| [1, 0, 0, 0, 0, 1] | 1 |

based on the features

Machine learning approach



1.3 Supervised machine learning

| Features (data) observations | Target (desired output) |
|---------------------------------|-------------------------|
| [1, 1, 0, 0, 1, 1] | 1 |
| [0, 0, 0, 1, 0, 1] | 0 |
| [1, 1, 1, 0, 1, 0] | 1 |
| [1, 0, 0, 0, 0, 1] | 1 |

feature matrix X
vector y

feature matrix
 $g(X) \approx y$
model
target variable

the goal of machine learning is to be able to train to come up to this function g that takes in the feature matrix and produces output that is as close as possible to the target variable.

Types of supervised machine learning

1) Regression $g(\mathbf{X}) \rightarrow -\infty \text{ to } +\infty$

- ↳ car price, house price
[output is a number]

2) classification

- [output is a category]

↳ email spam, identify objects in picture

sub categories

- multiclass [car / cat / dog / ...]

- binary [spam / not spam]

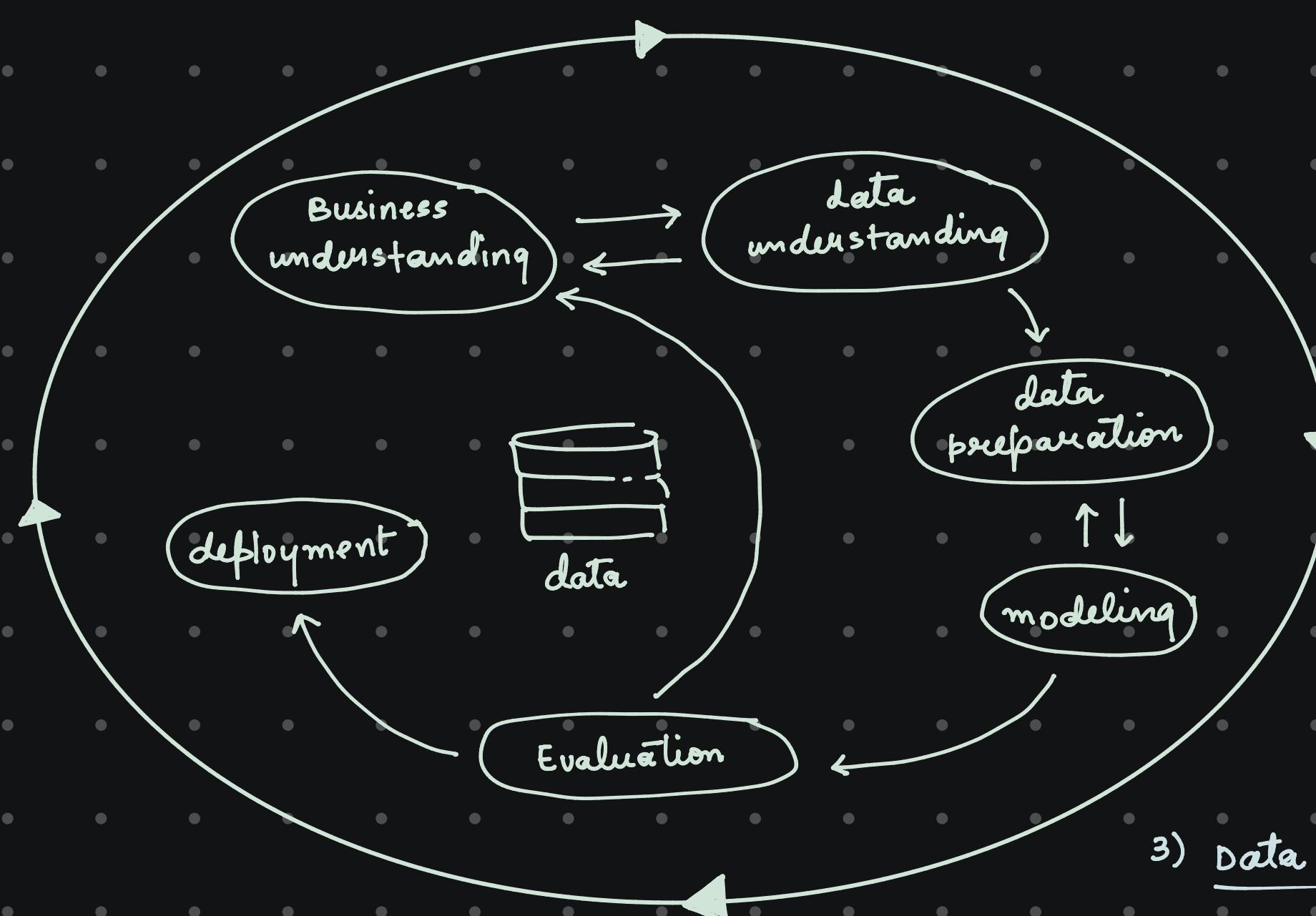
3) Ranking $g(\mathbf{X}) \rightarrow \text{probability}$

- ↳ recommending systems

1.4 CRISP - DM ML Process

CRISP - DM : cross industry standard processing data mining

- ↳ methodology for organizing ML projects.



4) Modeling

→ Train the models:

the actual machine learning happens here

- ↳ try different models (logistic regression, ...)
- ↳ select the best one (decision tree, neural network, ...)

→ sometimes, we may go back to data preparation:

- add new features
- fix data issues

5) Evaluation

→ measure how well the model solves the business problem.

Is the model good enough?

- have we reached the goal?
- Do our metrics improve?

Do a retrospective:

- was the goal achievable?
- Did we solve / measure the right thing?

After that we may decide to:

- Go back and adjust the goal
- Roll the model to more users / all users
- stop working on the project

supervised machine learning is about teaching computer / machine an algorithm by showing different examples (\mathbf{X} feature matrix) and the output predictions (y target variable).

Goal is to come up with a function g such that $g(\mathbf{X}) \approx y$.

ML Projects

- understand the problem
- collect the data
- Train the model
- use it

1) Business understanding

- identify the business problem
- understand how we can solve it

→ how would we measure the success of our project?
reduce the amount of spam by 50%
⇒ Do we actually need ML?

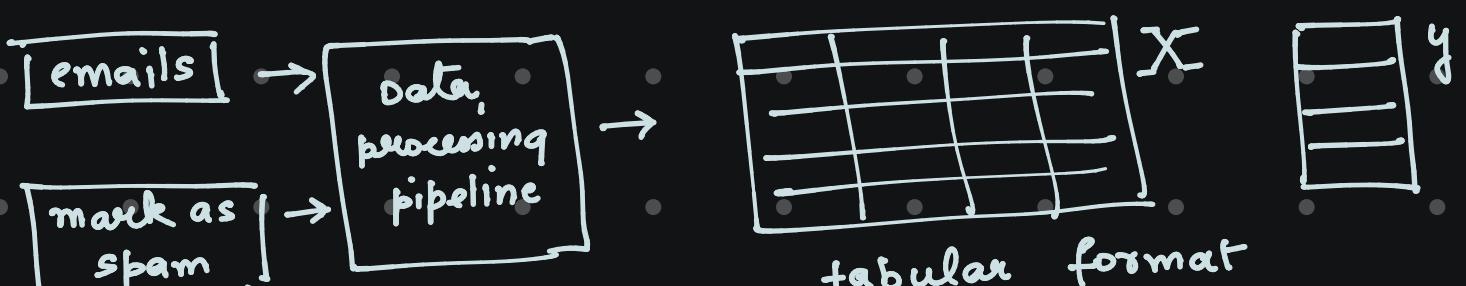
2) Data understanding

- analyze available data sources
- decide if we need to get more data?
- Identify the data sources
 - ↳ it may influence the goal
 - ↳ we may go back to the previous step and adjust it.

3) Data Preparation

- Transform the data so it can be put into a ML algorithm.
- ↳ extracting different features

→ clean the data → Build the pipeline → convert into tabular format.



6) Evaluation + Deployment

often happens together: 5% users

- online evaluation: evaluation of live users.
- it means: deploy the model, evaluate it.

Deployment

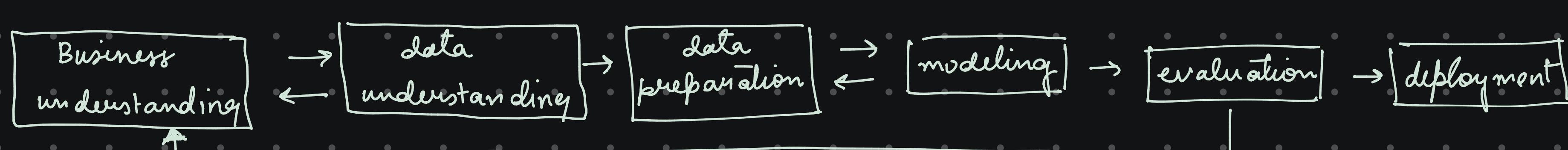
→ deploy the model to production.

- ↳ roll the model to all users.

- ↳ Proper monitoring

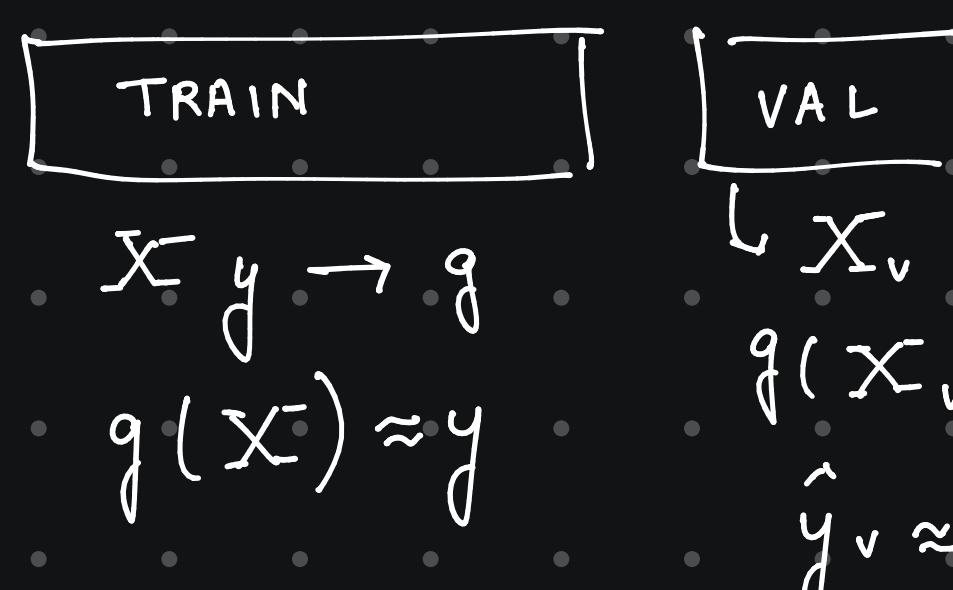
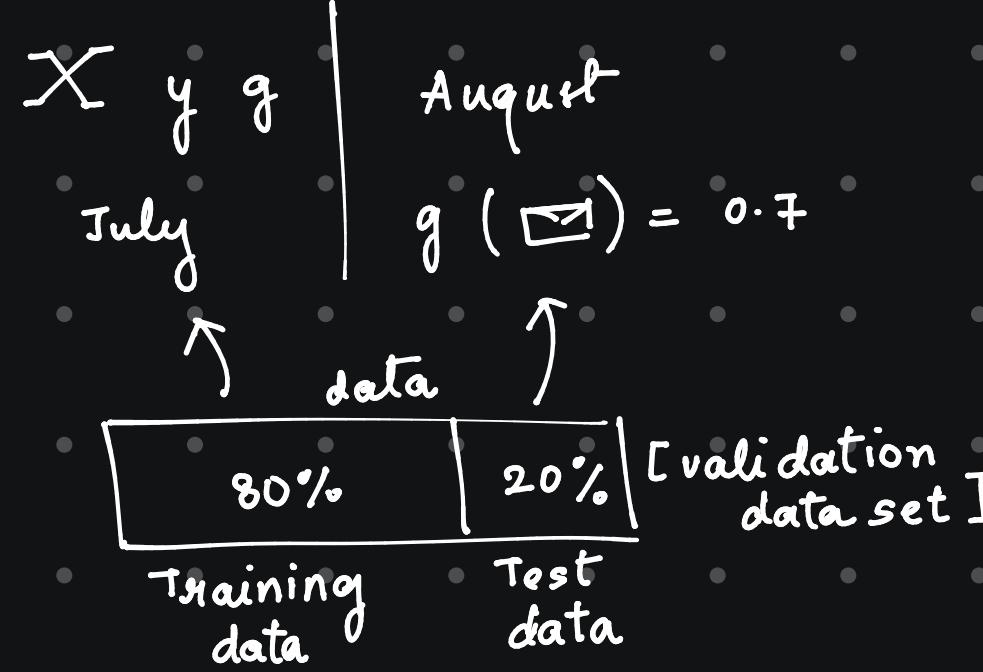
- ↳ Ensuring the quality and maintainability.

* ML Projects may require many iterations!



1.5 Model Selection Process

selecting the best model

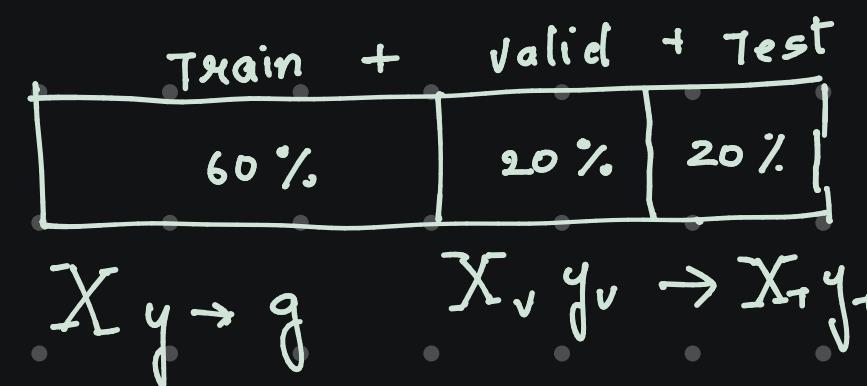


| \hat{y}_v | y_v | | |
|-------------|-------|---|---|
| 0.8 | 1 | 1 | ✓ |
| 0.7 | 1 | 0 | ✗ |
| 0.1 | 0 | 0 | ✗ |
| 0.6 | 1 | 1 | ✓ |

total = 4
correct = 2
 $\frac{2}{4} = 0.5$
50% correct model

multiple comparison problem

In statistics when we perform the same comparison many times so when we use many different models and evaluate them against the same validation data set one of the models can get particularly lucky (just by chance)



- ① split the data
- ② Train the model (60% data)
- ③ validate the model (20% data)
- ④ select the best model
- ⑤ Test the model (20% data)

1.7 Introduction to NumPy

import numpy as np

creating arrays → np.zeros(5)
array([0, 0, 0, 0, 0])

np.ones(5)

array([1, 1, 1, 1, 1])

np.full(5, 2.5)

array([2.5, 2.5, 2.5, 2.5, 2.5])

a = np.array([1, 2, 3, 4])

a = array([1, 2, 3, 4])

a[3] = 4

a[2] = 10

a = [1, 10, 3, 4]

np.arange(10) [0, 1, 2, 3, 4, 5, 6, 7, 8, 9]

np.arange(3, 10)

np.linspace(0, 1, 11)

multi-dimensional arrays

now ↓ column
np.zeros((5, 2))

n = np.array([
[1, 2, 3],
[4, 5, 6],
[7, 8, 9]]
])

n[:, 1]
array([2, 5, 8])

n[0, 1] = 2
n[2]
array([7, 8, 9])

Element-wise operations

a = np.arange(5)

array([0, 1, 2, 3, 4])

a * 2

array([0, 2, 4, 6, 8])

(10 + (a * 2)) ** 2 / 100

comparison operations

a = array([0, 1, 2, 3, 4])
a >= 2

array([False, False, True, True, True])

b = array([1, 0, 3, 4, 5])

a > b

array([False, True, False, False, False]) → returns all the elements of array

a[a > b] → elements of array a for which the condition [a > b] is True.

Randomly-generated arrays

np.random.rand(5, 2)

np.random.seed(2) → get same random numbers.

np.random.randn(5, 2) → normal distribution

np.random.randint(low=0, high=100, size=(5, 2))

a.min()

a.max()

a.sum()

a.mean()

a.std()

multiplication

vector × vector
[dot product]

$$u \cdot v = \sum_{i=1}^n u_i v_i$$

$$\begin{bmatrix} 2 \\ 4 \\ 5 \\ 6 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 0 \\ 2 \\ 1 \end{bmatrix} = 2 \cdot 1 + 4 \cdot 0 + 5 \cdot 2 + 6 \cdot 1 = 14$$

matrix × vector

$$u \rightarrow \begin{bmatrix} 2 & 4 & 5 & 6 \\ 1 & 2 & 1 & 2 \\ 3 & 1 & 2 & 1 \end{bmatrix}, v \rightarrow \begin{bmatrix} 1 \\ 0.5 \\ 2 \end{bmatrix} = \begin{bmatrix} u_1^T v \\ u_2^T v \\ u_3^T v \end{bmatrix}$$

row of matrix × vector

U · dot (V)

1.9 Introduction to Pandas

```
import numpy as np  
import pandas as pd  
  
df = pd.DataFrame(data, columns=columns)  
data can also be a list of dictionaries.  
(Keys → column name)  
df.head() → First 5 Rows  
df.head(n=2)  
  
df.Make ← column name  
df['Make']  
  
df['vehicle_Style'].str.lower()  
df['vehicle_style'].str.replace(' ', '_')  
  
df.MSRP.values  
df.to_dict(orient='records')
```

- to delete a column
del df['id']
- range index
df.index
- df.Make.index → returns a new df
- df.loc[1] → index
df.loc[[1, 2]] → df.reset_index()
df.loc[1, 2] → create an index column to store the previous indices
df.index = ['a', 'b', 'c', 'd', 'e']
- positional index
df.iloc[1]
df.iloc[[1, 2, 4]]
- df.MSRP.max()
df.MSRP.min()
df.MSRP.mean()
df.MSRP.describe()
df.describe().round(2)
- df.Make.unique()
df.unique()
- df.groupby('Transmission-type').MSRP.mean()

```
df['Engine HP'] / 100  
df['year'] >= 2015  
df[  
    df['year'] >= 2015  
]  
df[  
    (df['year'] >= 2015) &  
    (df['Make'] == 'Nissan')  
]  
  
SELECT  
    transmission_type,  
    AVG(MSRP)  
FROM  
    cars  
GROUP BY  
    transmission_type
```