

A review on machine learning concepts

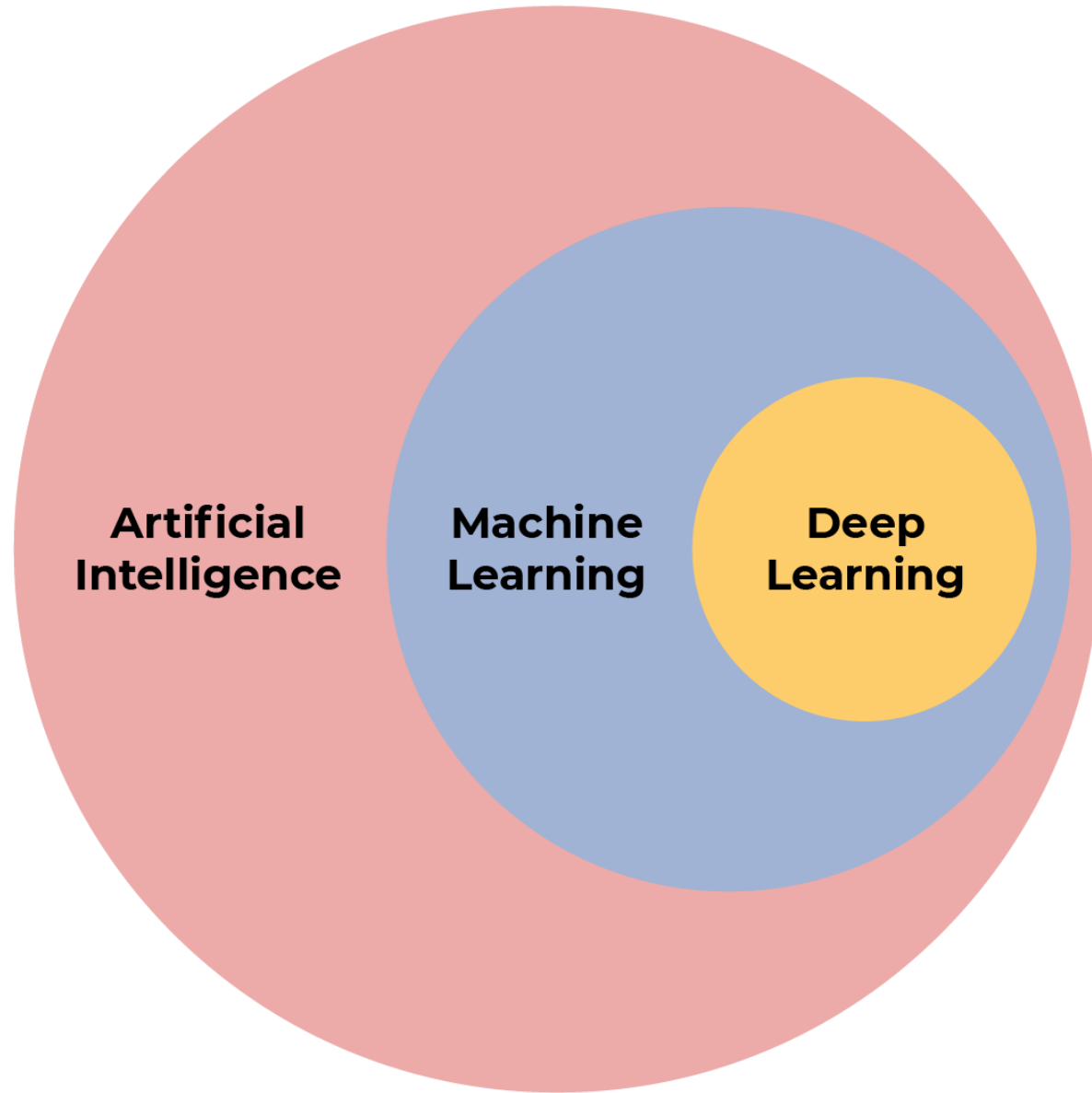
Ali Kohan



[Github.com/alikohan](https://github.com/alikohan)



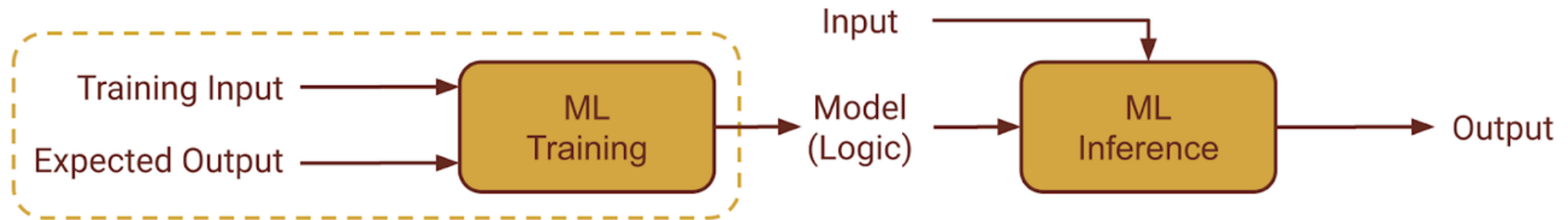
linkedin.com/in/alikohan/



Traditional Programs: Define algo/logic to compute output



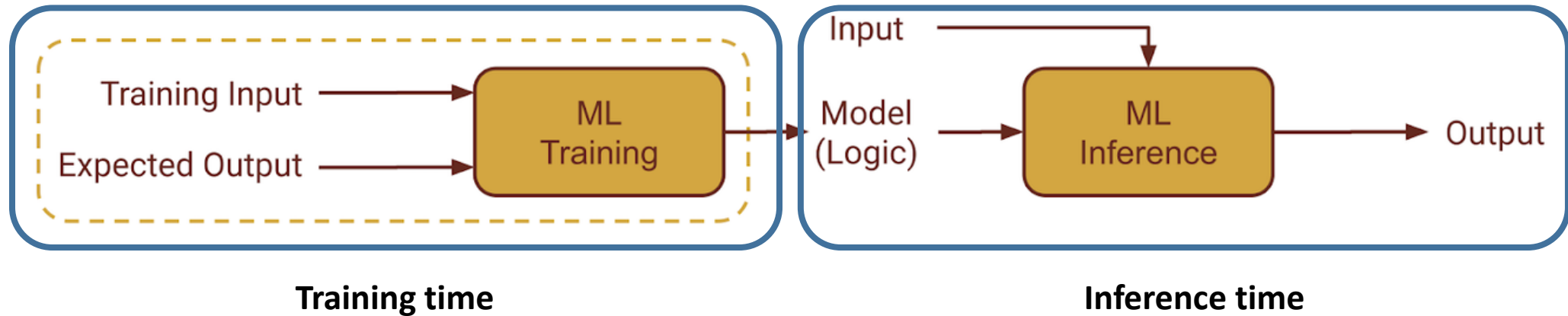
Machine Learning: Learn model/logic from data

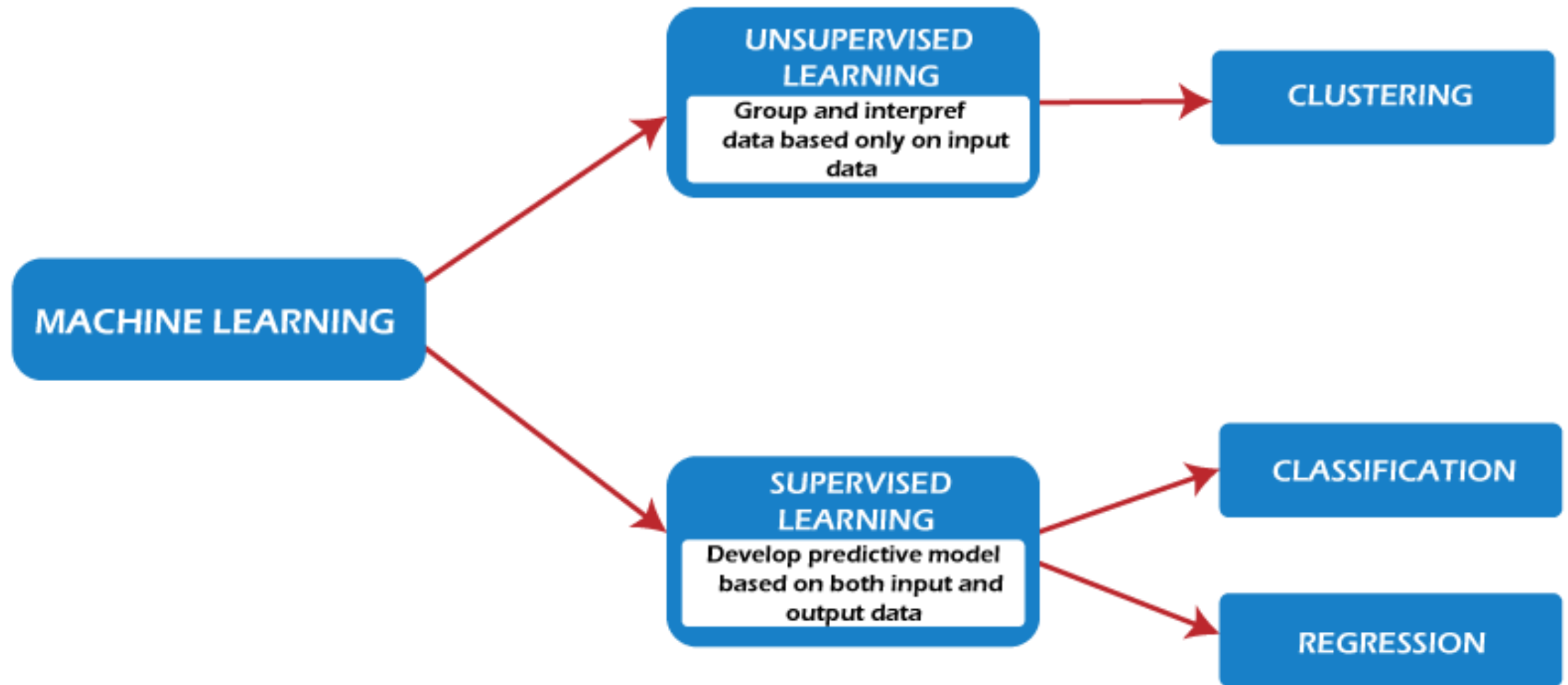


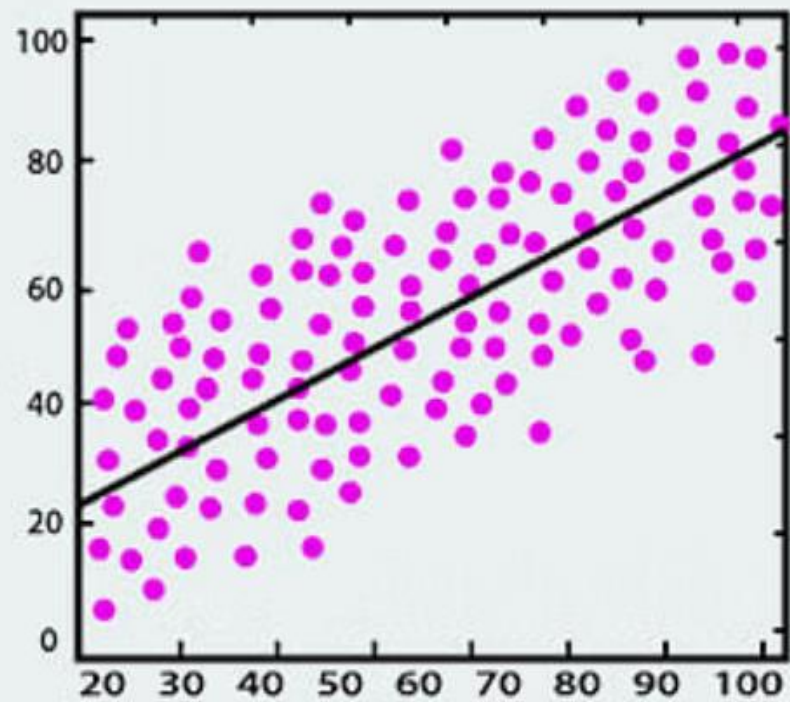
Traditional Programs: Define algo/logic to compute output



Machine Learning: Learn model/logic from data

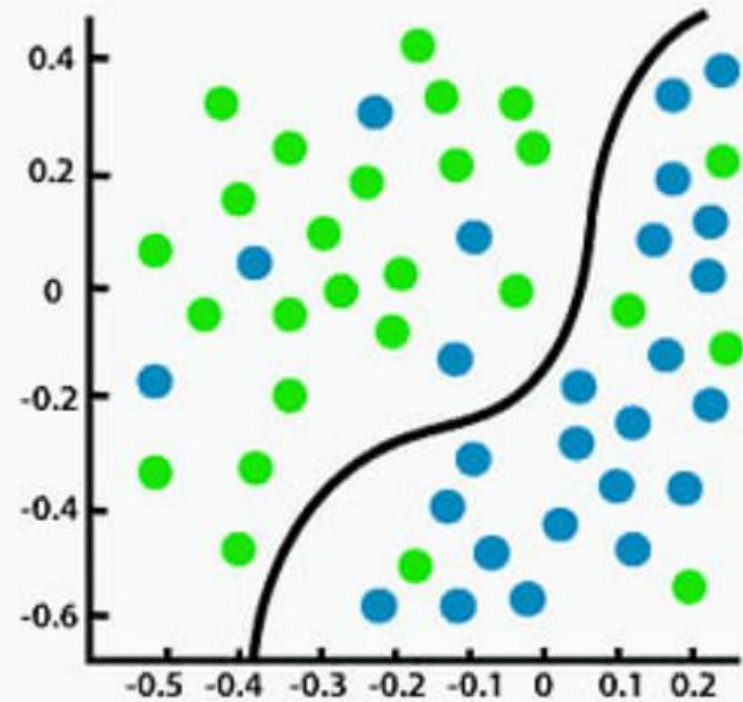






Regression

versus



Classification

Preliminary: Classification

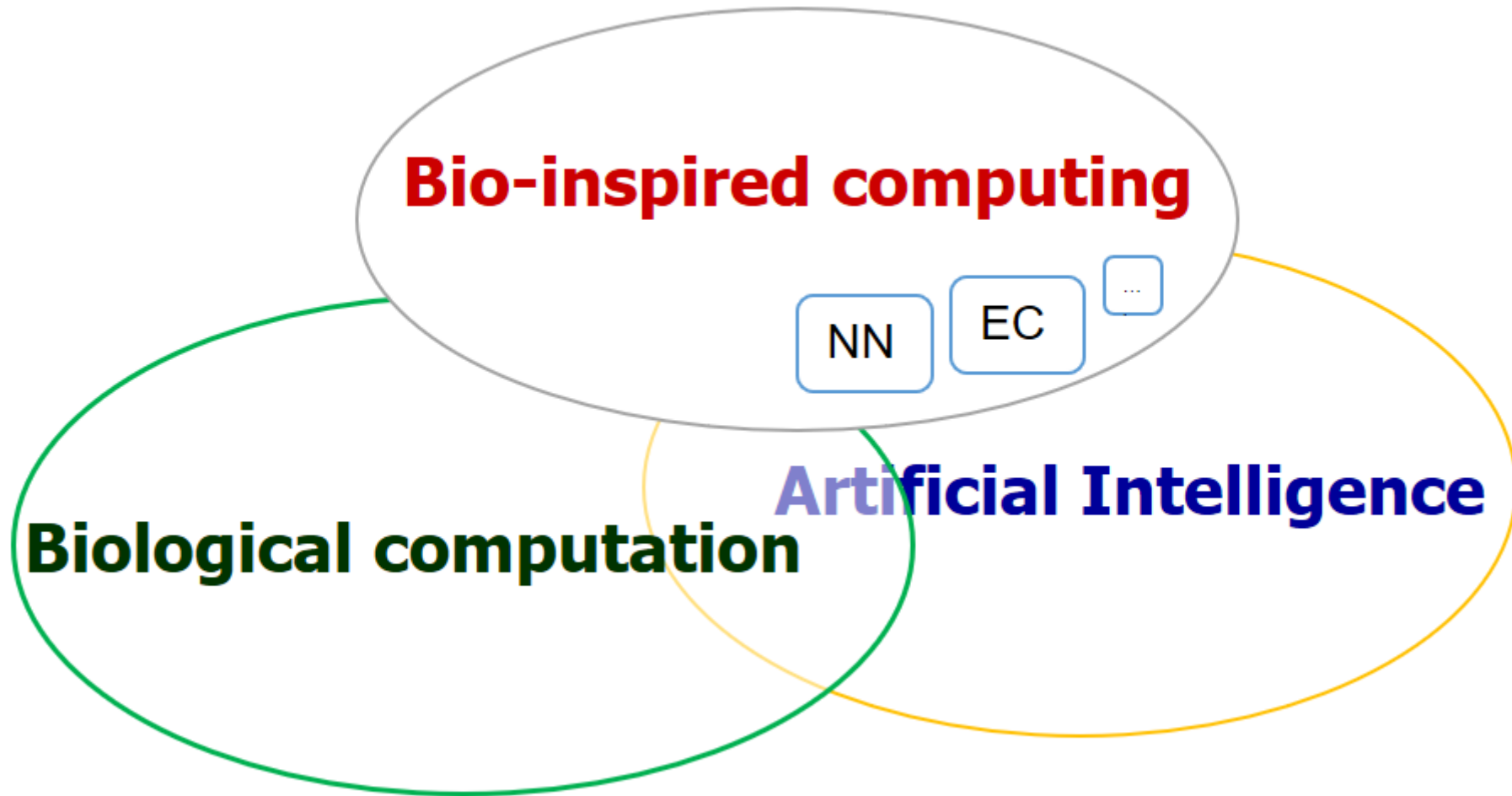
Column Name	Description
Email ID	Unique identifier for each email. Not used in classification but helps track emails.
Subject Length (words)	Number of words in the subject. Spam emails often have short, catchy subject lines.
Contains "Free" (Yes/No)	Indicates if the word "Free" is in the email. Spam frequently uses such words to grab attention.
Number of Links	Counts the number of links in the email. Spam usually contains multiple links to suspicious sites.
Has Attachment (Yes/No)	Indicates if the email has an attachment. Spam often includes attachments that may be malicious.
Class (Spam/Not Spam)	Target label that indicates whether the email is classified as spam or not spam. This is the prediction goal for the model.

Email ID	Subject Length	Contains "Free"	Number of Links	Has Attachment	Class (Spam/Not Spam)
1	5	Yes	3	Yes	Spam
2	12	No	1	No	Not Spam
3	8	Yes	0	No	Spam
4	20	No	2	Yes	Not Spam
5	6	Yes	5	Yes	Spam
...

Preliminary: Classification

Column Name	Description
Email ID	Unique identifier for each email. Not used in classification but helps track emails.
Subject Length (words)	Number of words in the subject. Spam emails often have short, catchy subject lines.
Contains "Free" (Yes/No)	Indicates if the word "Free" is in the email. Spam frequently uses such words to grab attention.
Number of Links	Counts the number of links in the email. Spam usually contains multiple links to suspicious sites.
Has Attachment (Yes/No)	Indicates if the email has an attachment. Spam often includes attachments that may be malicious.
Class (Spam/Not Spam)	Target label that indicates whether the email is classified as spam or not spam. This is the prediction goal for the model.

Email ID	Subject Length	Contains "Free"	Number of Links	Has Attachment	Class (Spam/Not Spam)
1	5	Yes	3	Yes	Spam
2	12	No	1	No	Not Spam
3	8	Yes	0	No	Spam
4	20	No	2	Yes	Not Spam
5	6	Yes	5	Yes	Spam
...

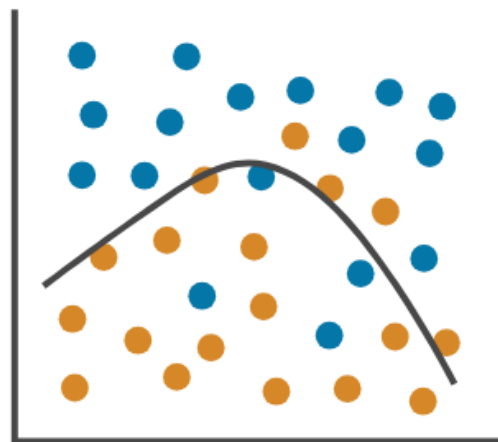


Classification

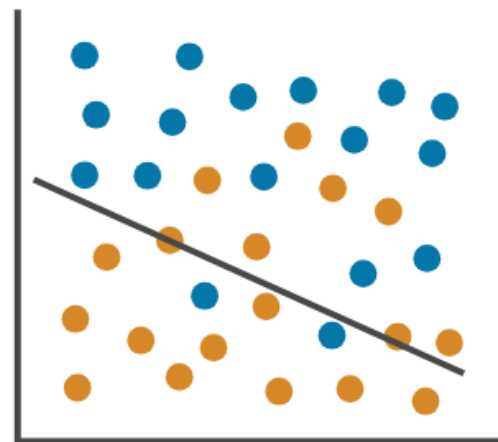
Overfitting



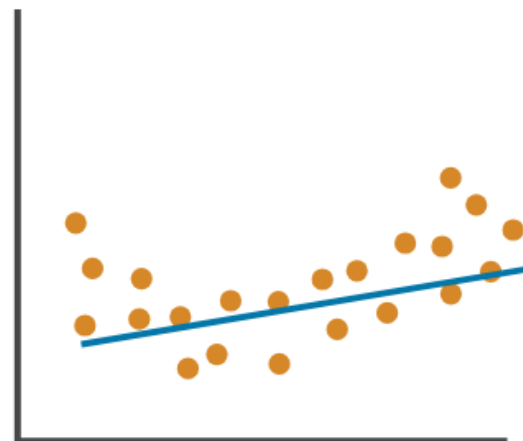
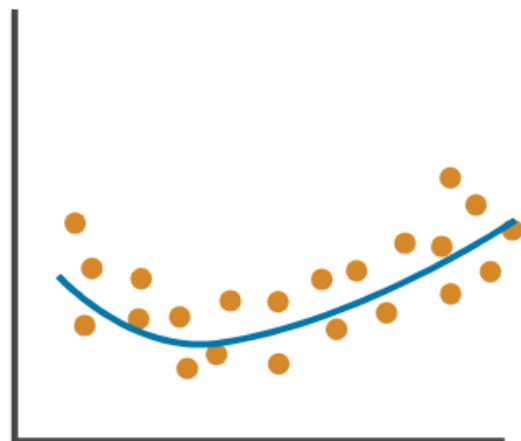
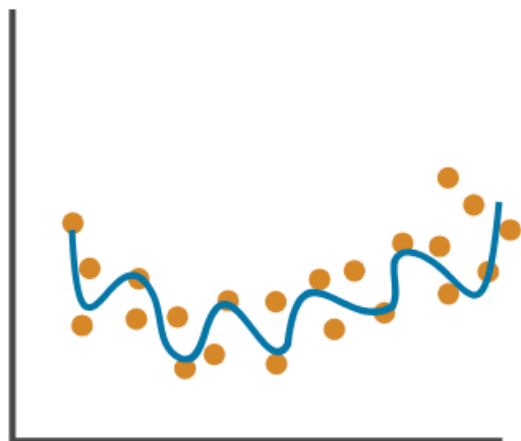
Right Fit

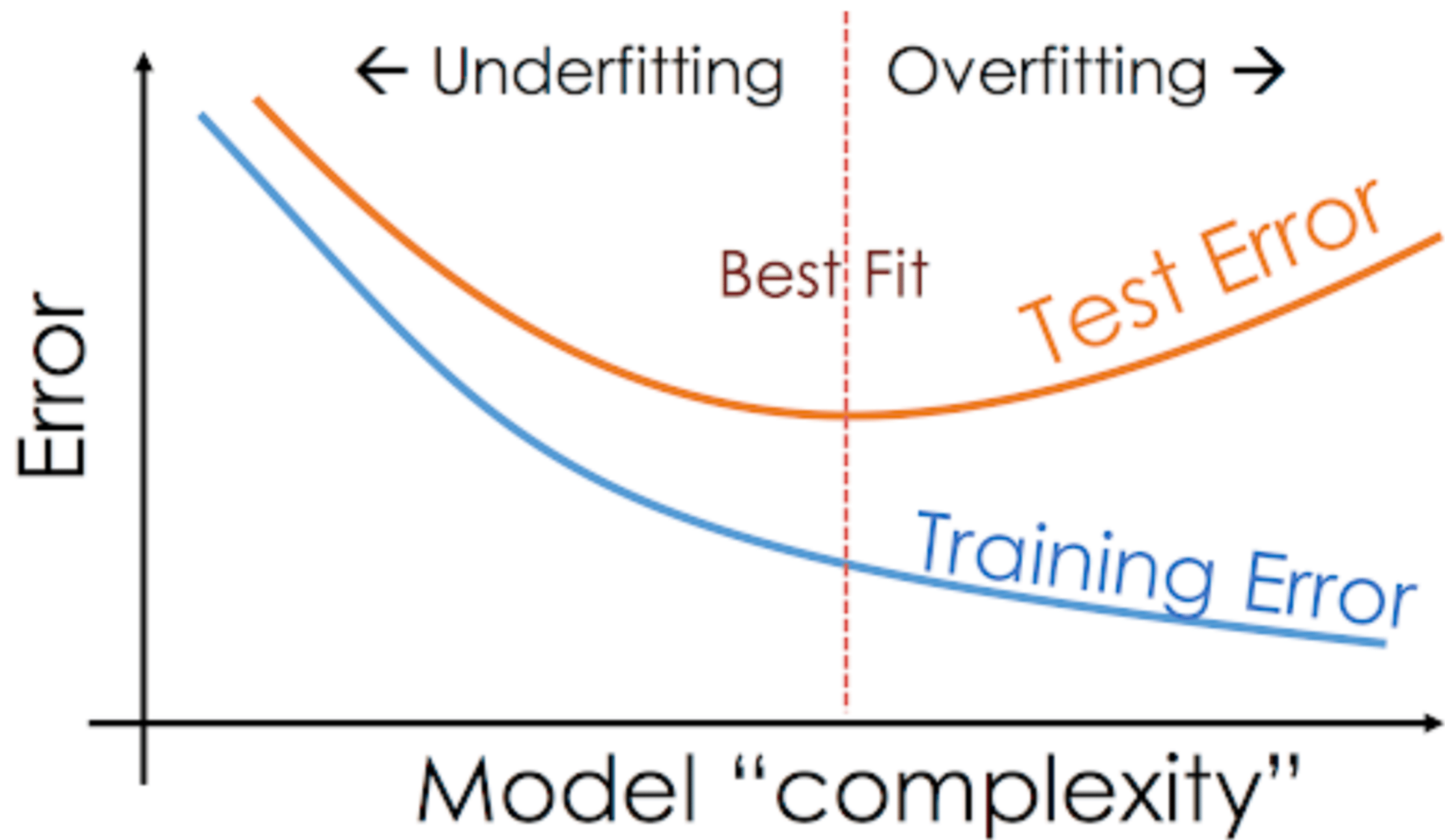


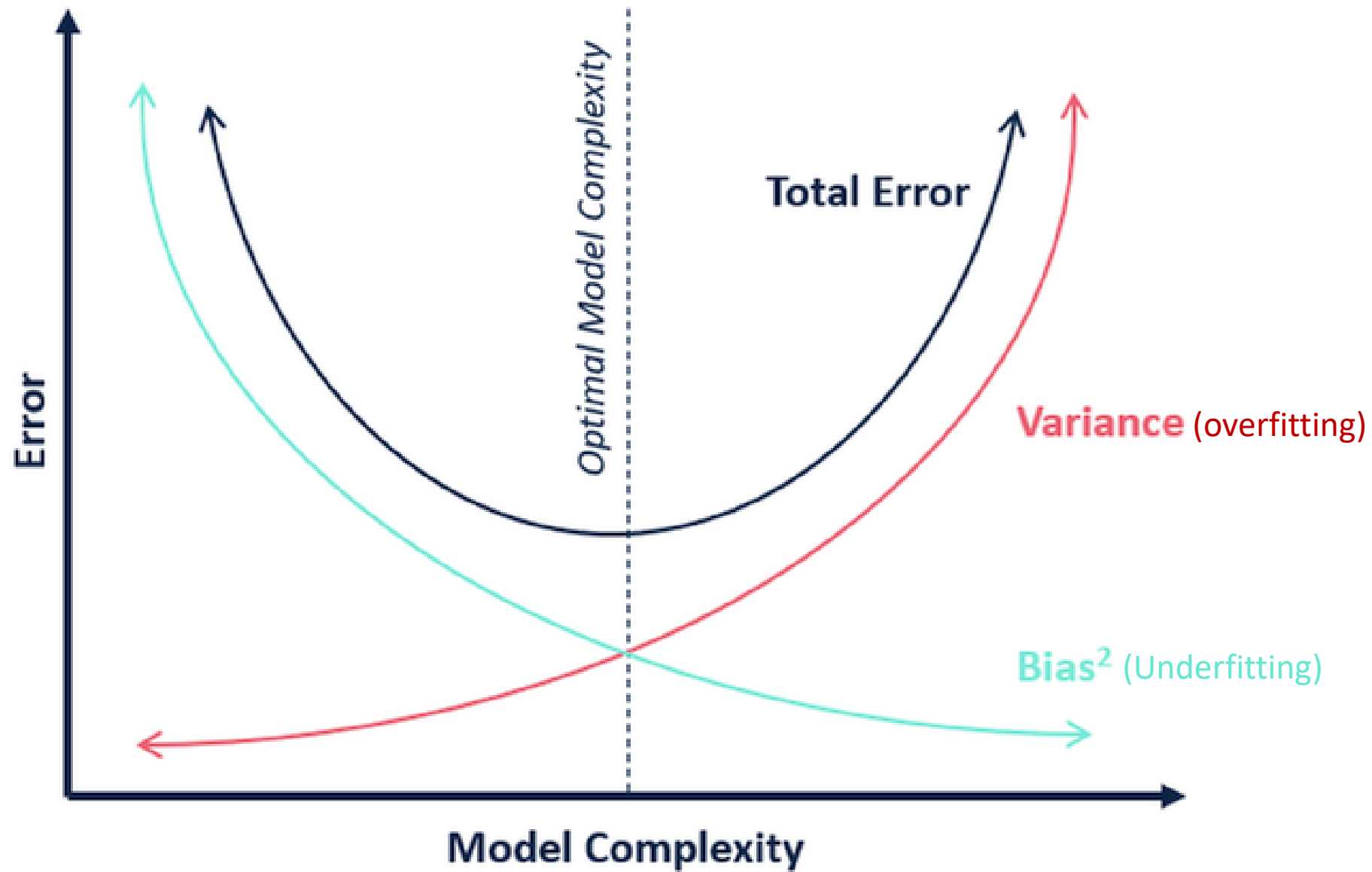
Underfitting



Regression







All Data

```
graph TD; A[All Data] --> B[Training]; A --> C[Validation]; A --> D[Test];
```



Training

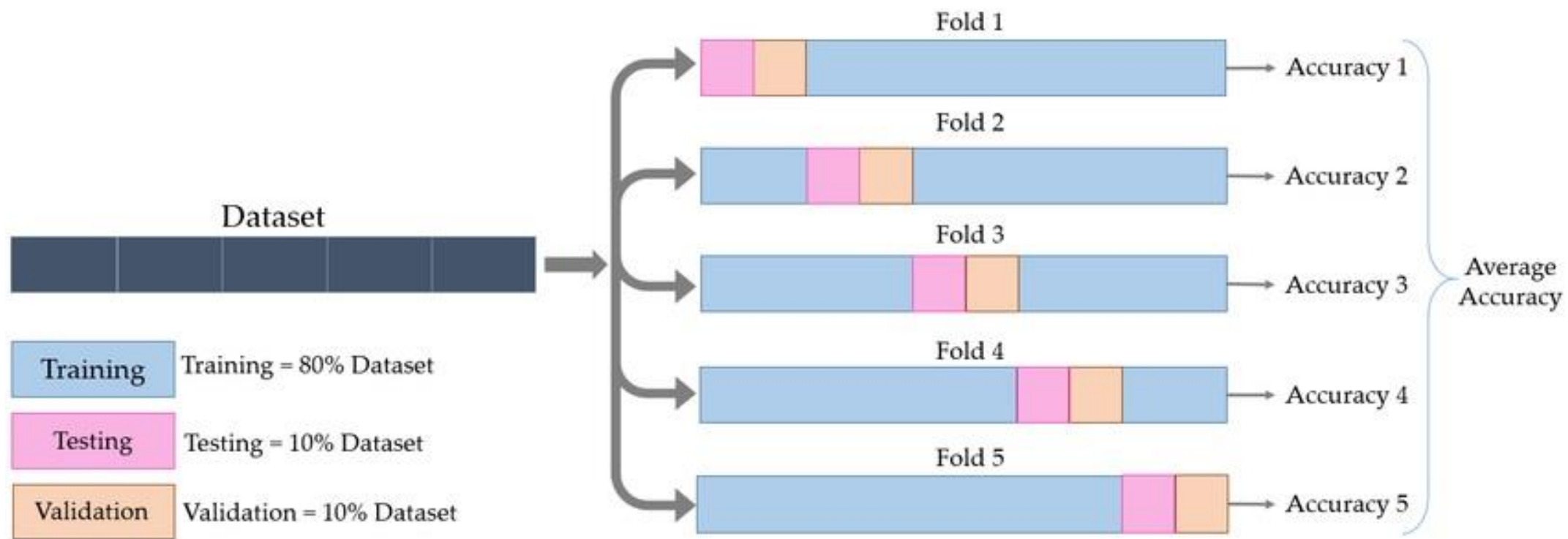
Validation

Test

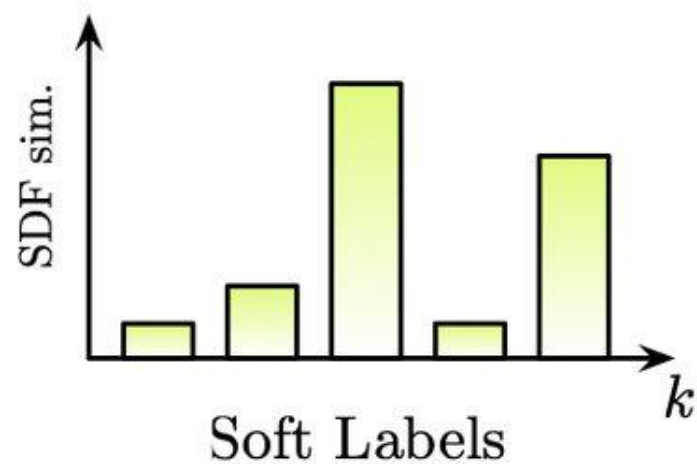
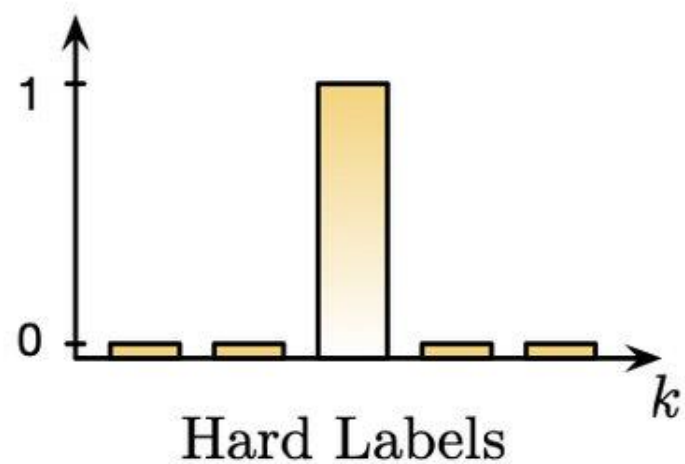
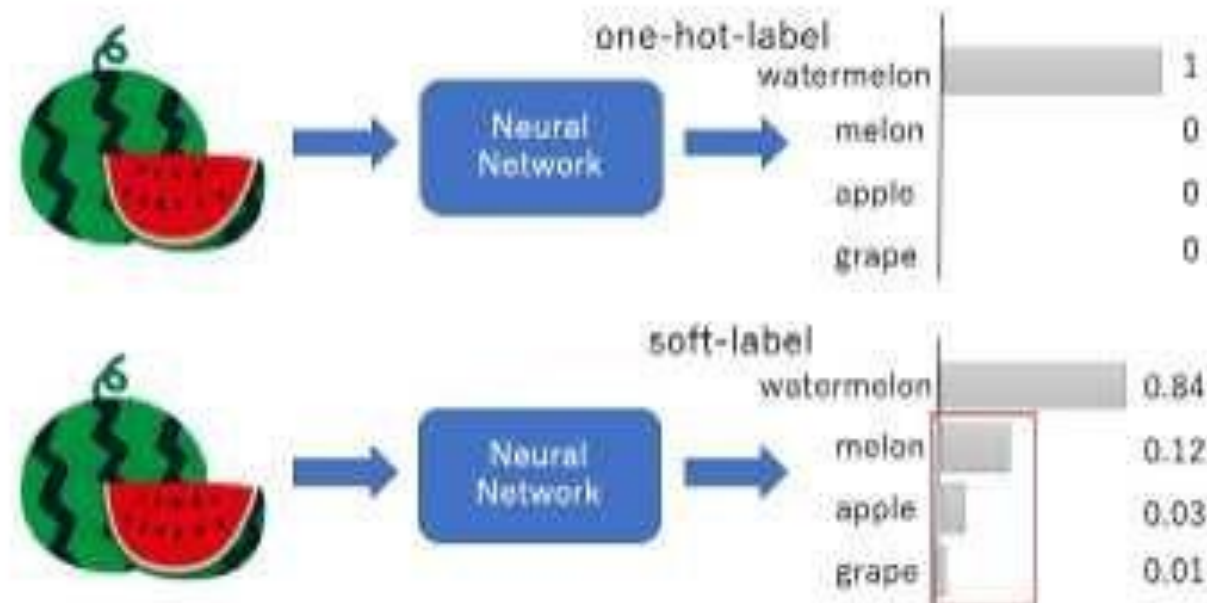
Models learn the task

Which model
is the best?

How good
is this
model truly?



Soft label vs hard label



Explore

Exploit



Higher uncertainty

Lower uncertainty

Explore



Exploit

