Scan the QR code to mark your attendance

**Attendance**

SPAI

# Learning Objectives

- ✅ Train, Test and Validate

- ✅ Cross Validation

- ✅ Understanding Bias & Variance in ML models

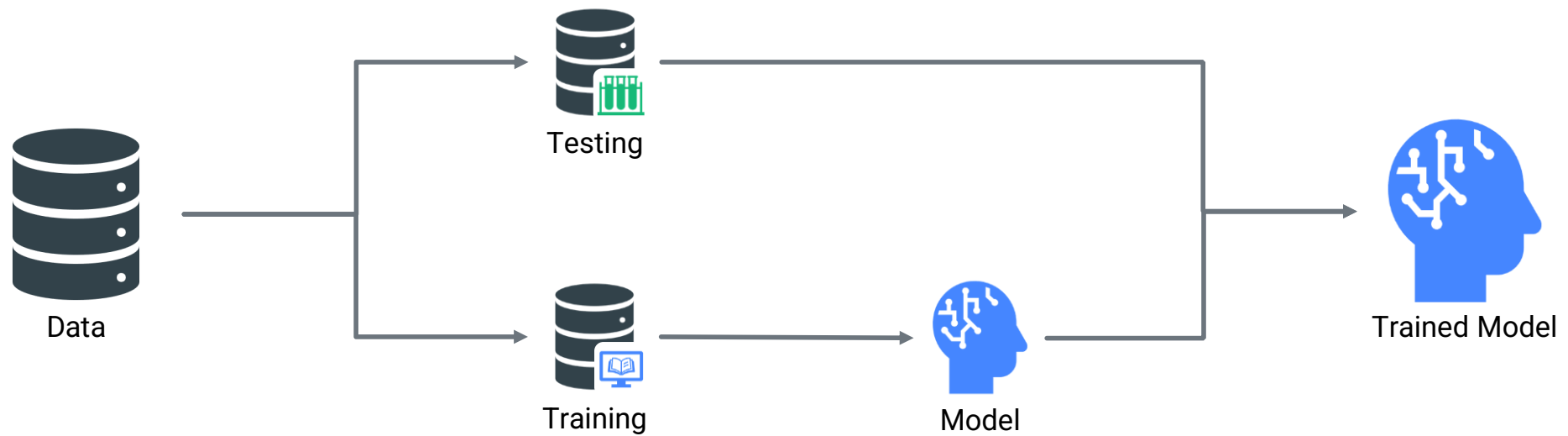- ✅ Interpreting Model Complexity using learning curve

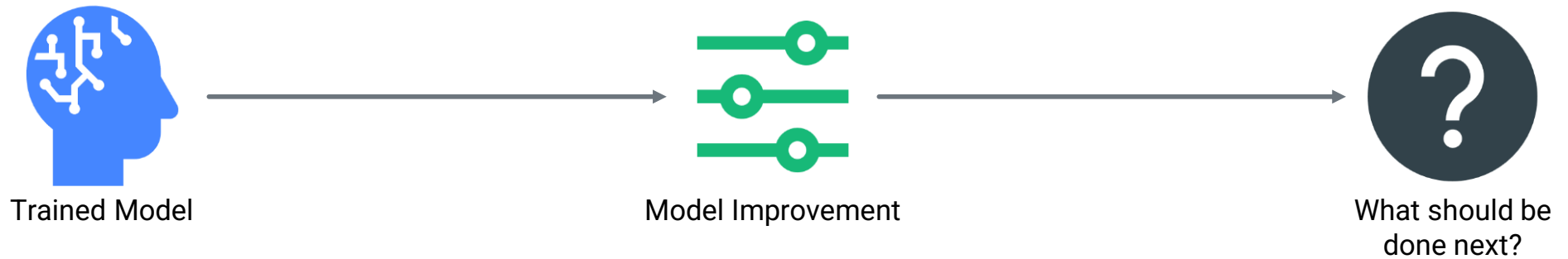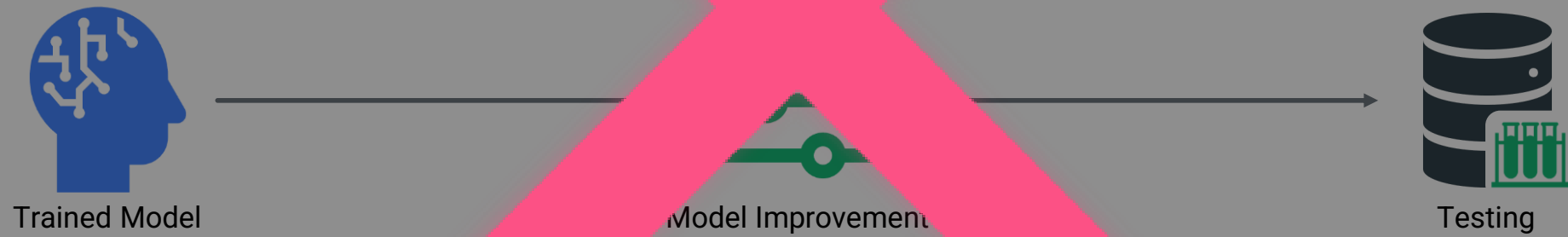Predict if customers will switch to your telco based on certain features

Scenario

SPAI

# We learnt previously...



Data → Testing, Training → Model → Trained Model

# Improving Model



Trained Model        Model Improvement        What should be done next?

SPAI

# Improving Model



Trained Model → Model Improvement → Testing

# Why can't we use the testing set to evaluate the model?

The testing set allows us to get a realistic representation of the performance of the model

Recall

# Why can't we use the testing set

**?**

- ☑ Testing set allows us to test the model in an unbiased way

- ☑ However, we are improving our model based on the results of our testing set

- ☑ This can result in biases in our model, trying to "suit" our testing set
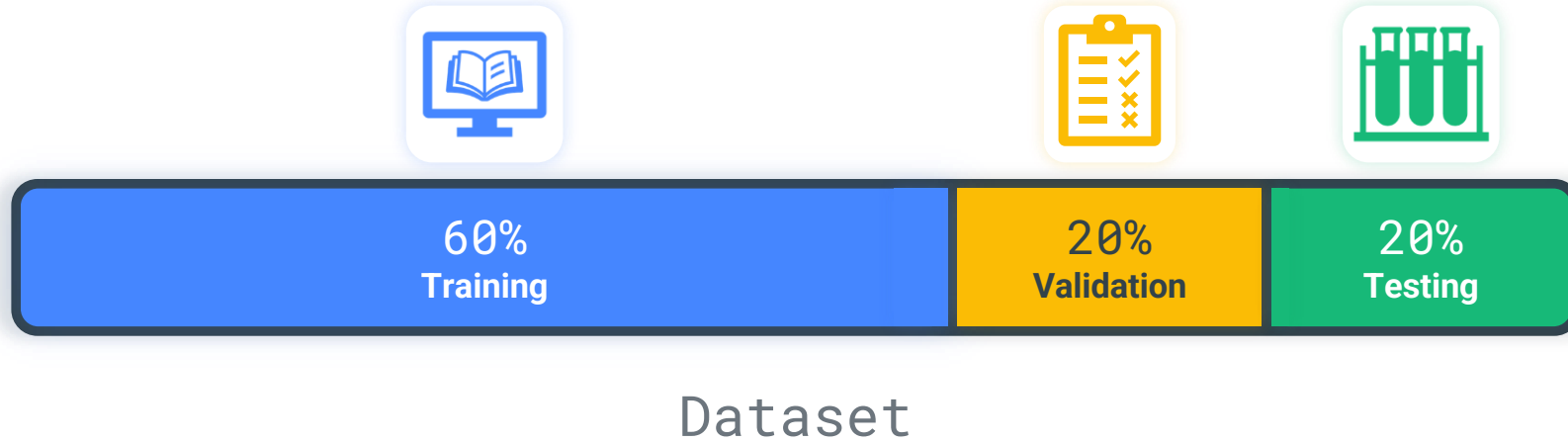
# Train, Validate and Test

"

Using the same ideology of training and testing, we now further split our data into 3 sets

**What is it?**

?

SPAI

# What is Train Test Validate?



| 60%<br>Training | 20%<br>Validation | 20%<br>Testing |
| --- | --- | --- |

Dataset

SPAI

Training Data

Usually takes up **60%** of the dataset

Dataset only used to train the model

NEVER use model score from this dataset to judge the performance of the model

SPAI

Validation Data

Usually takes up to **20%** of the dataset

Dataset only used to evaluate model performance

NEVER use this dataset as a conclusion of the model performance
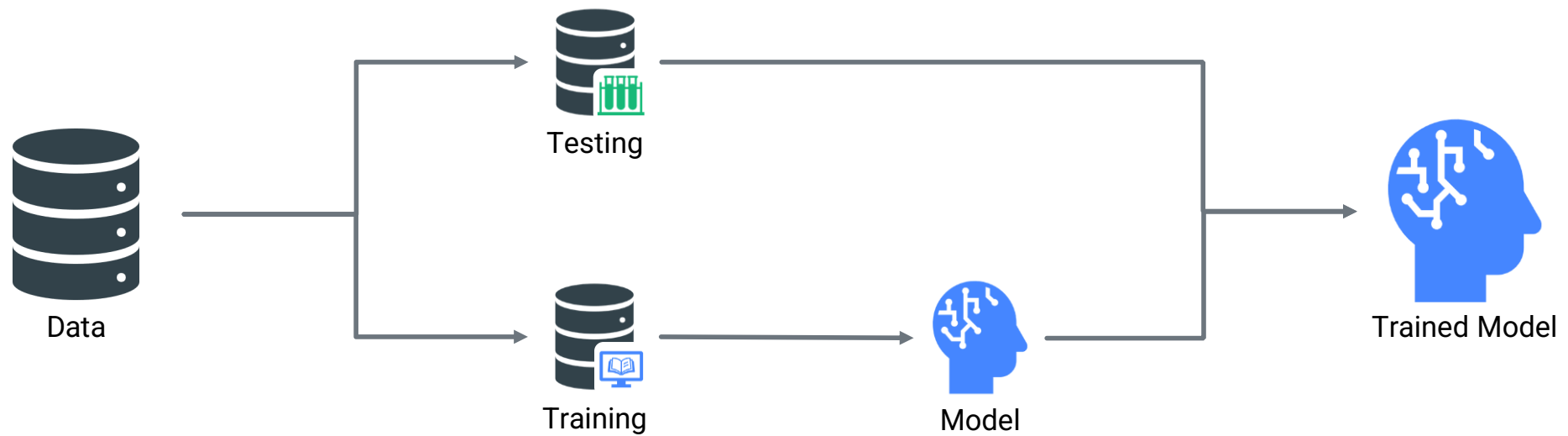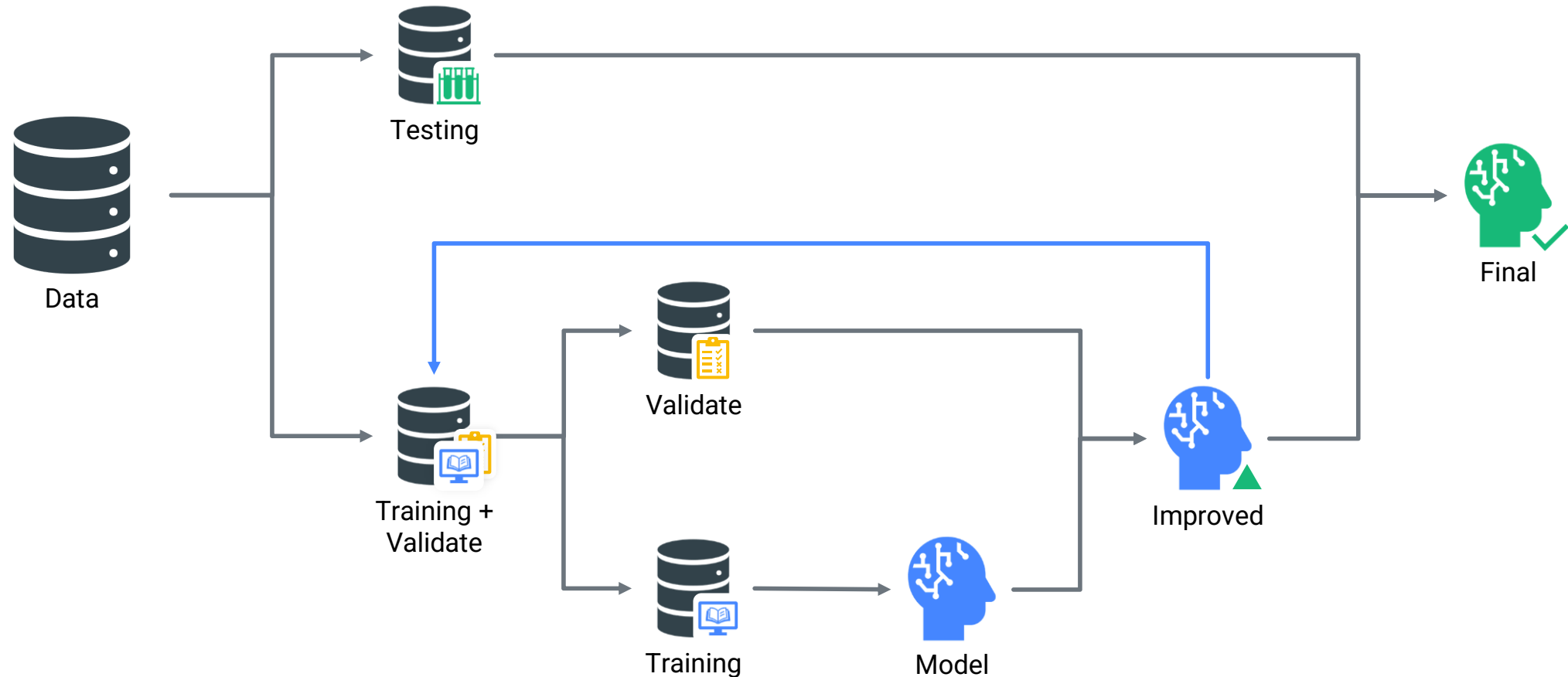
SPAI

Testing Data

Usually takes up **20%** of the dataset

Dataset only used to evaluate final model performance

NEVER make changes on the model based of the performance from this dataset

SPAI

# We learnt previously…



Data → Testing / Training → Model → Trained Model

SPAI

# Actual Diagram



Data

Testing

Training + Validate

Validate

Training

Model

Improved

Final

SPAI

Very waste of data.

Only 60% of data is used to train the model

**Problem**

SPAI

# Cross Validation (CV)

# Why Cross Validate

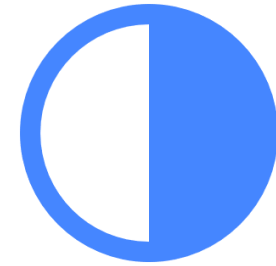- ☑ Helps us make full use of our dataset

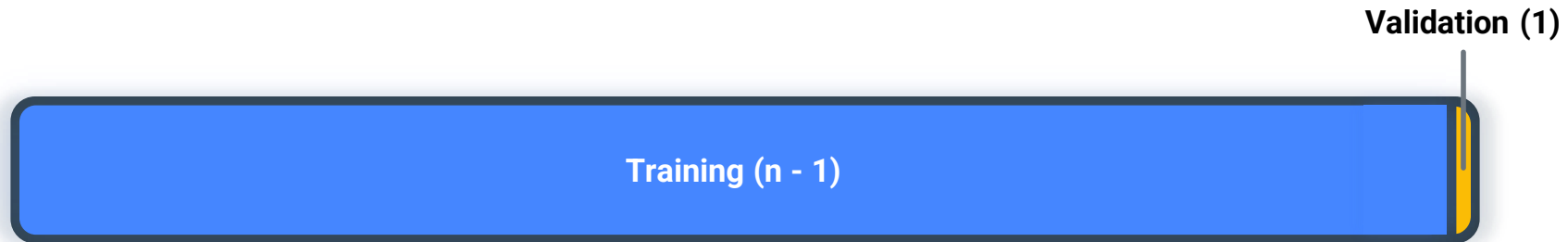- ☑ Training model is a random process

# Types of CV



LOO

K-fold

Stratified

# Leave one out CV



**Validation (1)**

**Training (n - 1)**

Training &
Validation Set

*Where n is the number of rows of data*

SPAI

# Leave one out CV



Training (n - 1) — Validation (1)

Training (n - 1) — Validation (1)

Training (n - 1) — Validation (1)

Training (n - 1) — Validation (1)

n times

SPAI

LOO

Best use of data as guarantees 100% use of it

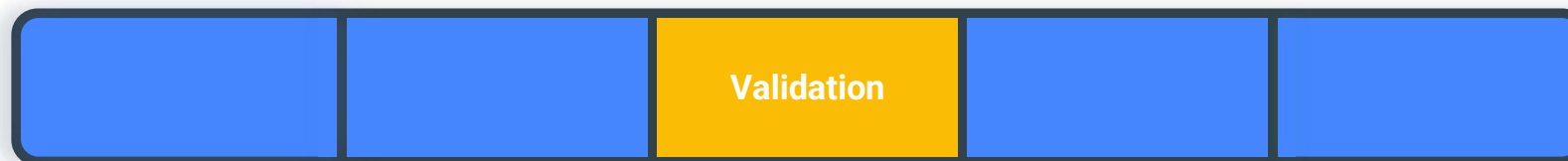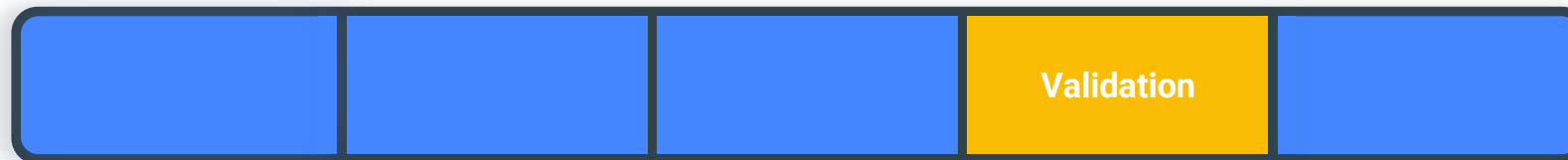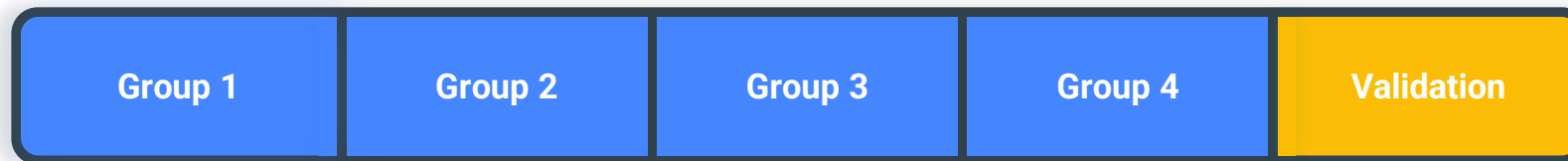Very time consuming and resource intensive

Not used often

SPAI

# K-Fold CV



Training & Validation Set

# K-Fold CV

| Group 1 | Group 2 | Group 3 | Group 4 | Group 5 |

# K-Fold CV



| Group 1 | Group 2 | Group 3 | Group 4 | Validation |

| | | | Validation | |

| | | Validation | | |

⋮ k times
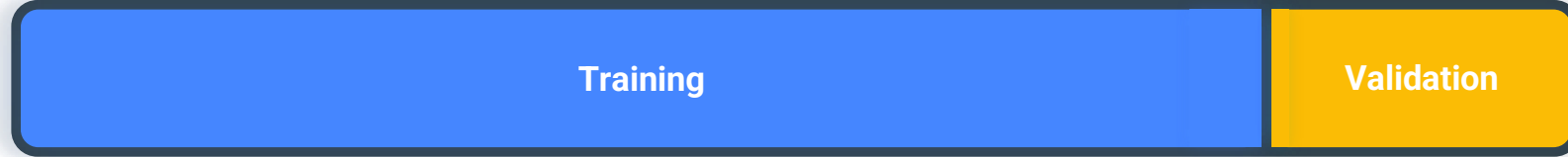
*Where k = 5*

SPAI

K-Fold

Good use of data as guarantees 100% use of it

Very efficient and fast process
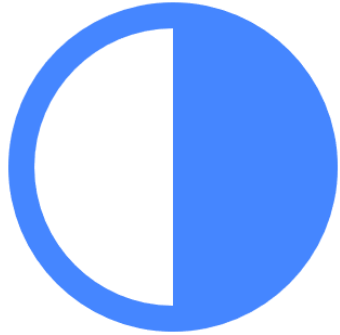
Most used method of cross validation

# Stratified CV



Training & Validation Set

# Stratified CV



*Where ratio of classes is 5:2 and there are only 2 classes*

**Stratified**

Good use of data as guarantees 100% use of it

Slightly slower than K-fold but still highly efficient

Useful for imbalanced classes

# Knowledge Check

```
scores = cross_validate(LogisticRegression(), x_train, y_train, cv=3)
```

> How many groups will the dataset be split into?

A.  1
B.  2
C.  3
D.  4

SPAI

# Knowledge Check

```
cross_validate(DecisionTreeClassifier(), x_train, y_train)
```

```
> What is the default value of cv?

A.  3
B.  5
C.  8
D.  10
```

SPAI

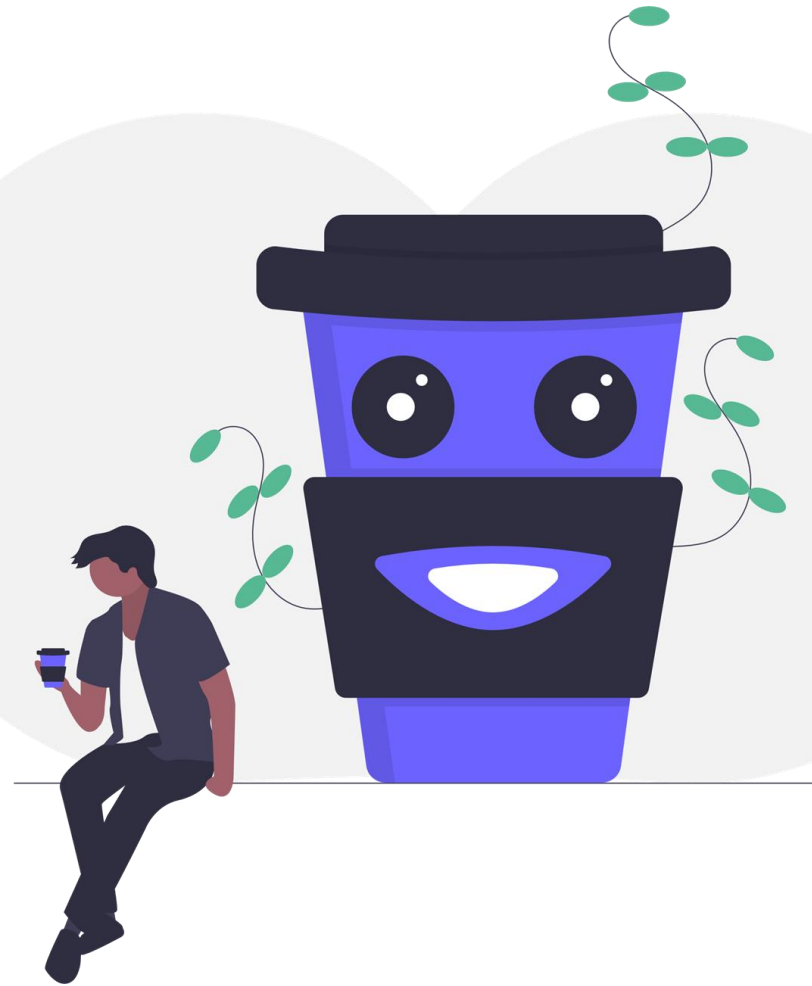# Practice Time!

# 10 Minutes

Please attempt exercise 1
We will go through the exercises later

SPAI

# Times up

We will now go through the exercises

# Break & QnA

## 10 Minutes

# Bias and Variance

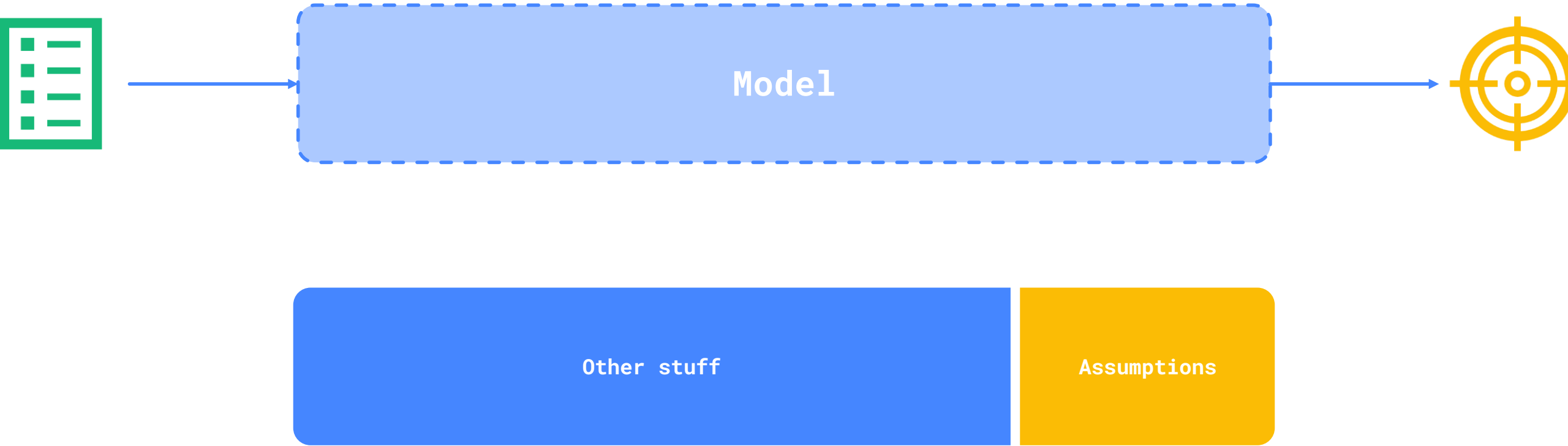Amount of assumptions made by a model to make the target function easier to learn

**What is Bias**

# What is Bias?

# What is Bias?



Model

Other stuff | Assumptions

SPAI

# Why model make assumptions?

- ☑ Makes it easier to learn and predict

- ☑ Results in faster learning speed

Bias

**Low Bias:** Less assumptions made

**High Bias:** More assumptions made

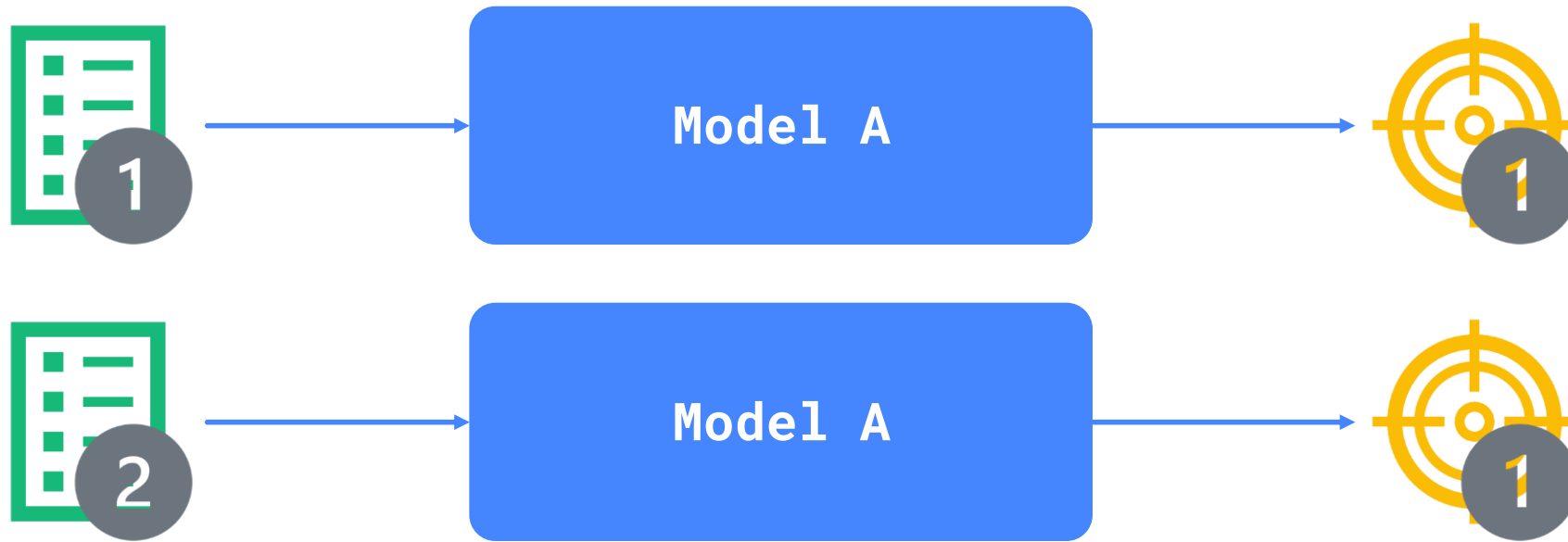Amount changes to the estimate of the target function if different training data was used
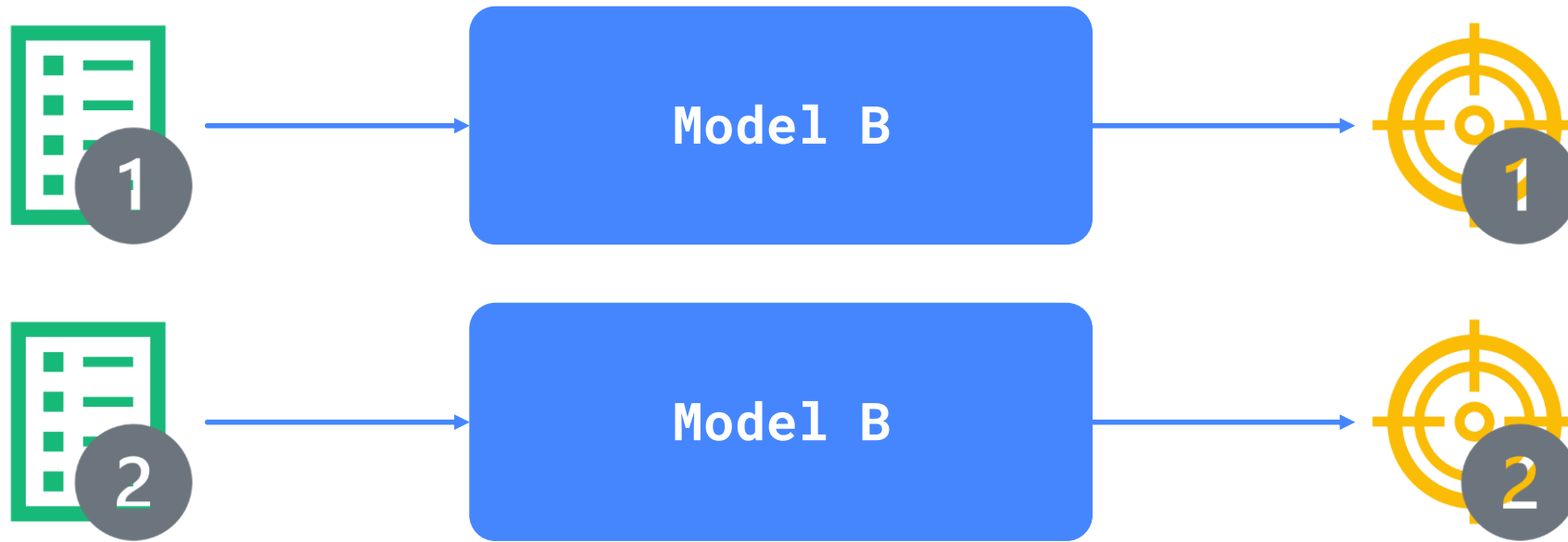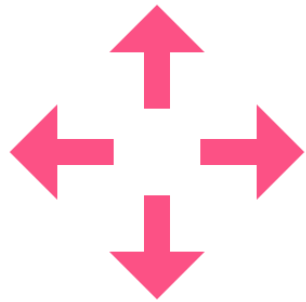
**What is Variance**

# What is Variance?



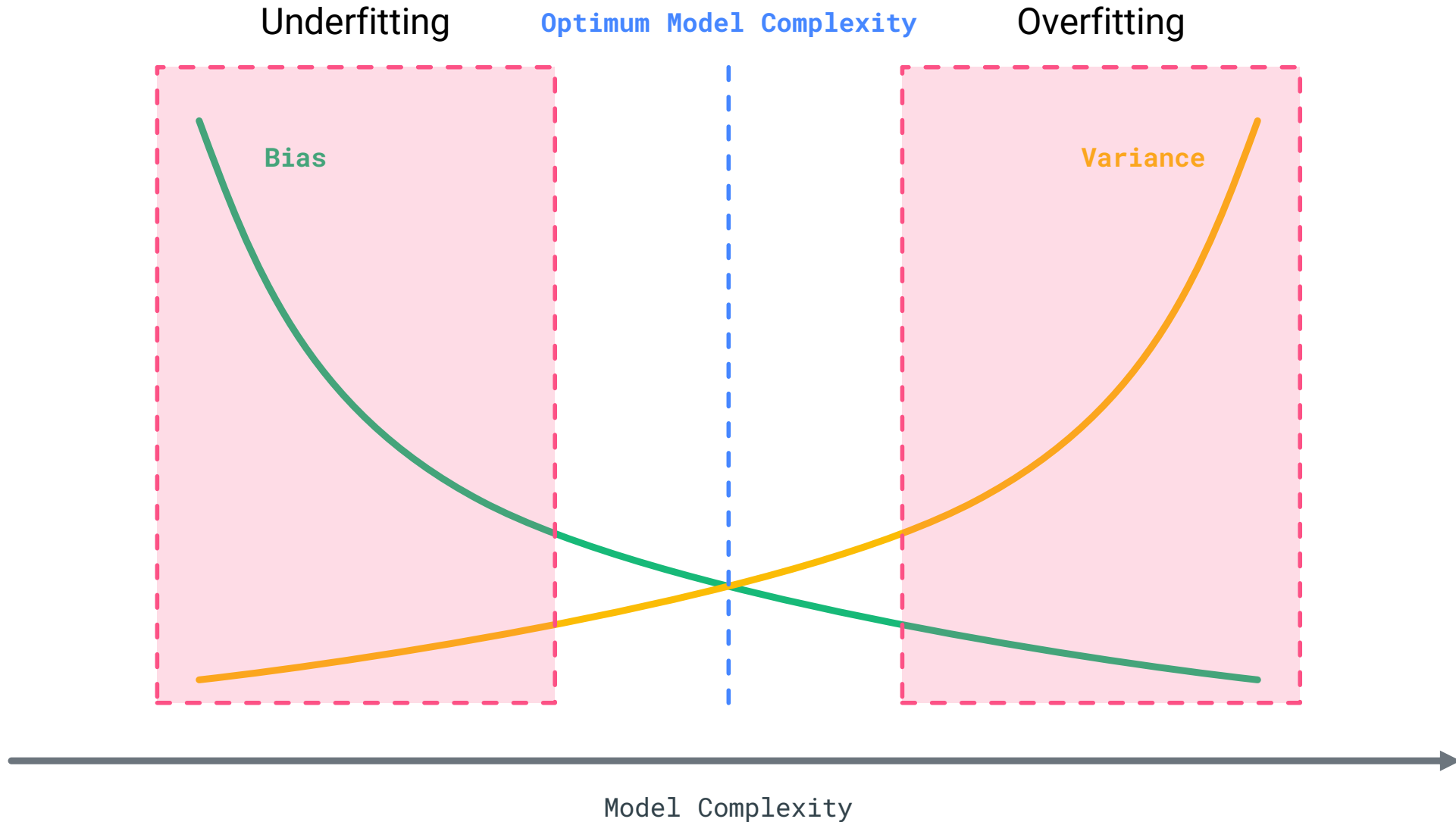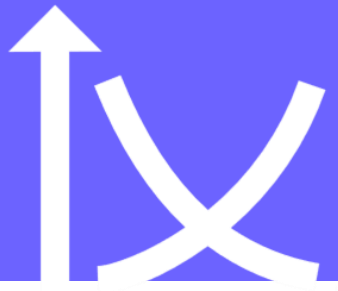Model

# What is Variance?

# What is Variance?

Variance

**Low** Variance: Small changes to training dataset results in small changes to prediction

**High** Variance: Small Changes to training dataset results in large changes to prediction
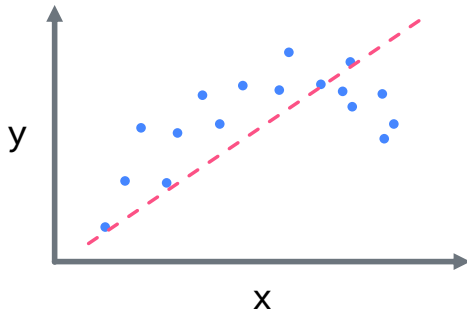
SPAI

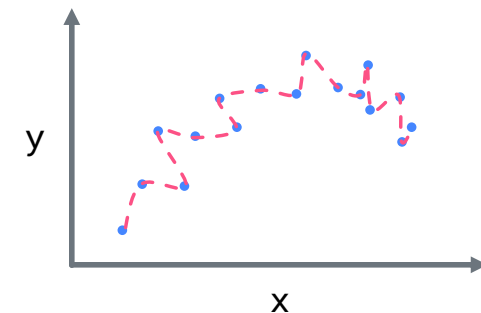# Bias & Variance together

# Under/Over fitting
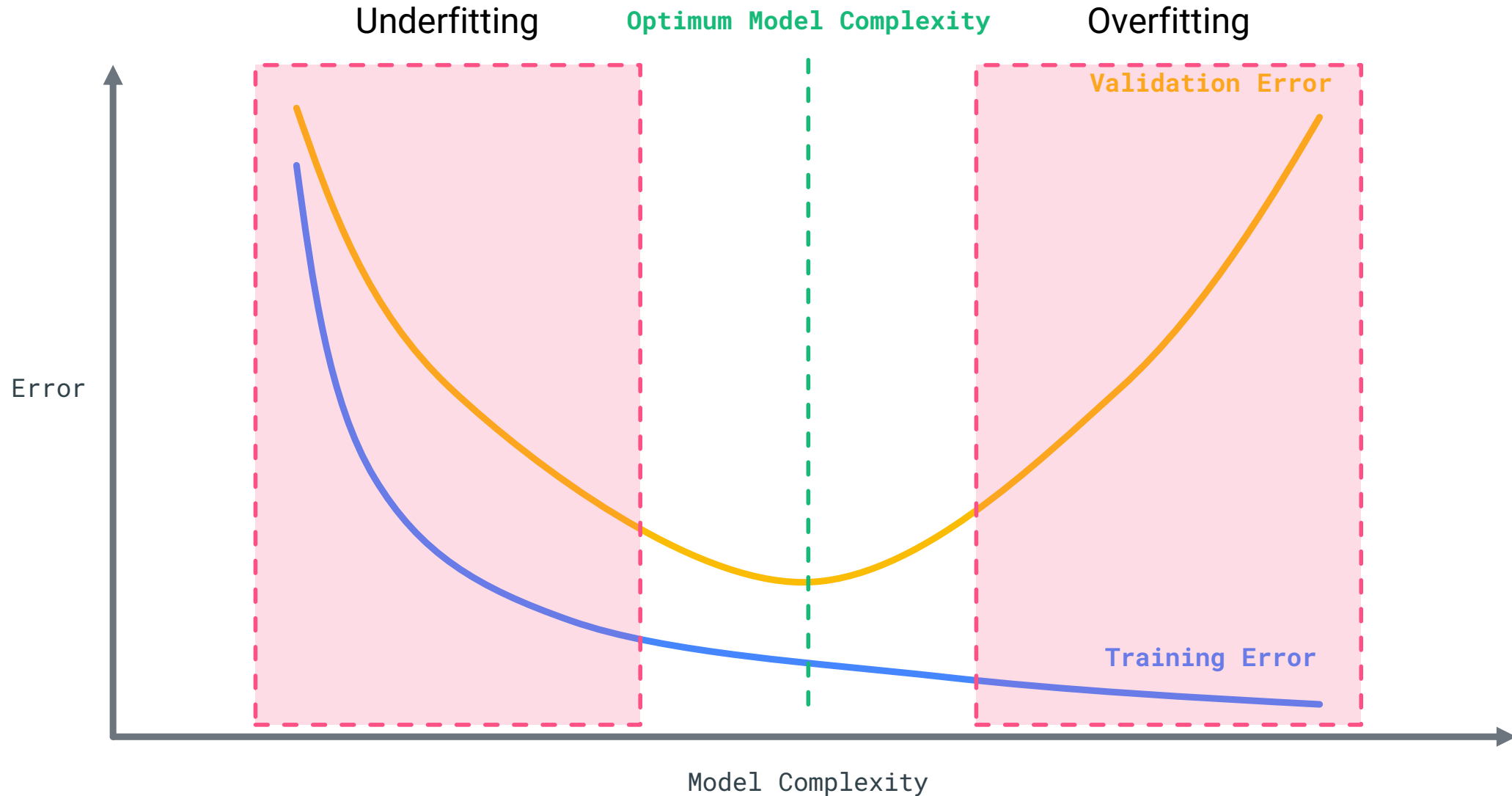
## Underfitting

- High Bias
- Low Variance



## Overfitting

- Low Bias
- High Variance

# Underfitting & Overfitting



Underfitting     Optimum Model Complexity     Overfitting

Validation Error

Error

Training Error

Model Complexity

# Knowledge Check

```
> Which of the following statements are True

A.  A model with low bias and high variance is an underfitted model

A.  When a model changes drastically with small changes on its training set, it is said to
    have high bias

A.  When a model changes drastically with small changes on its training set, it is said to
    have high variance

A.  It is best when model have high bias and high variance
```

# Knowledge Check

| | Training Accuracy | Testing Accuracy | Training F1_Score | Testing F1_Score |
|---|---|---|---|---|
| 0 | 0.998225 | 0.78967 | 0.996643 | 0.552586 |

```
> This model has…

A.  High Bias, High Variance
B.  High Bias, Low Variance
C.  Low Bias, High Variance
D.  Low Bias, Low Variance
```
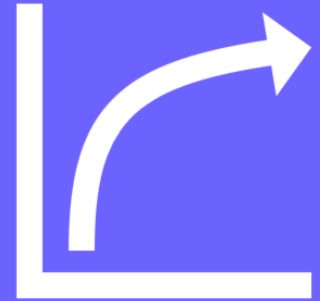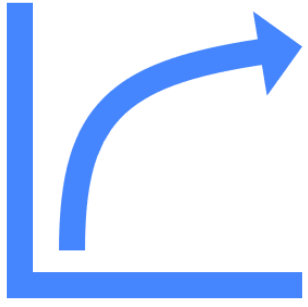
SPAI

# Knowledge Check

| | Training Accuracy | Testing Accuracy | Training F1_Score | Testing F1_Score |
|---|---|---|---|---|
| 0 | 0.503239 | 0.485269 | 0.34748 | 0.337036 |

> This model is Overfitting

A. True
B. False

**Purpose**: Allows us to understand why a model performs a certain way

# Model Learning Curves

A learning curve is the correlation between a model's score against the amount of data it is given

**What is a Learning Curve?**
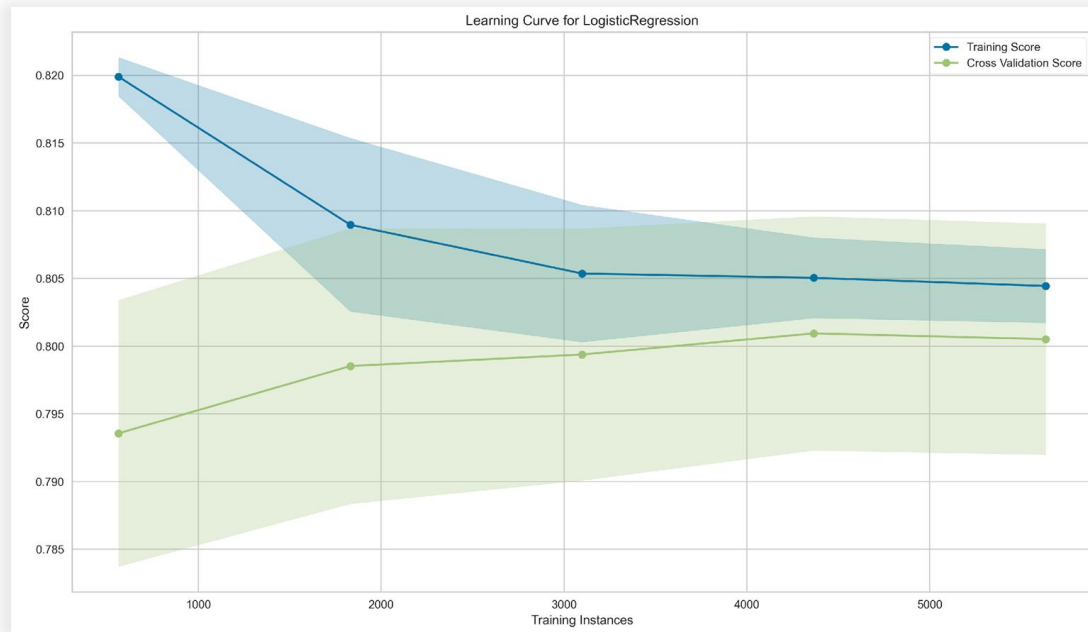
# Why plot learning curves?

☑ Allows us to see if a model is over/underfitting

☑ From there, we can make useful decisions on how we can improve our model
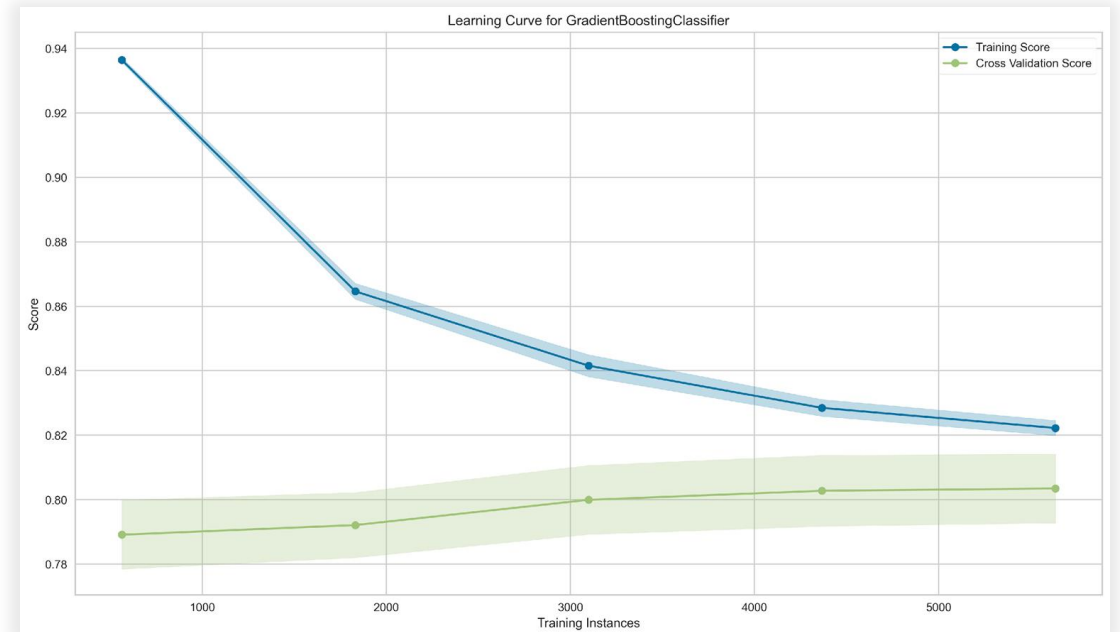
# Good Fit Characteristics



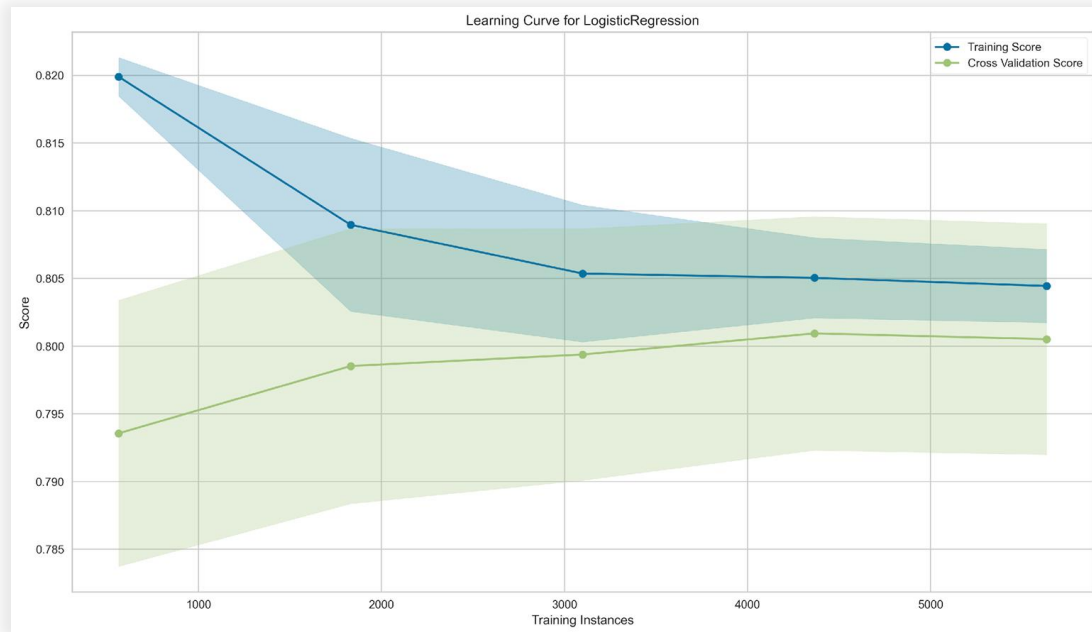Learning Curve for LogisticRegression

Lines moves towards each other
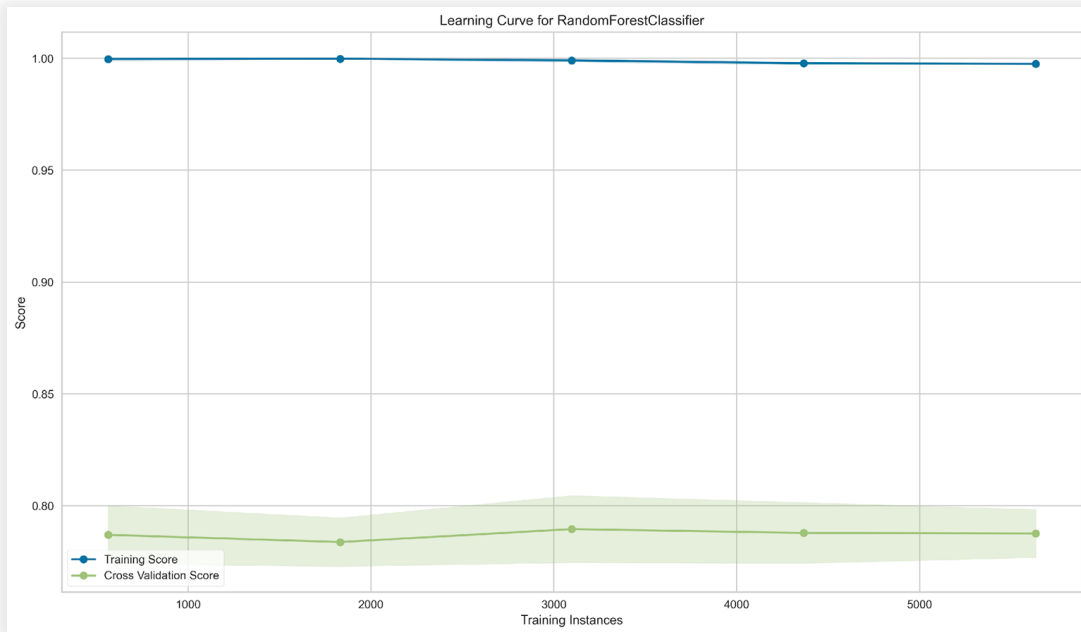
Lines maintain small space between each other

Score for BOTH are generally high

# Good Fit Characteristics

# Overfit Characteristics
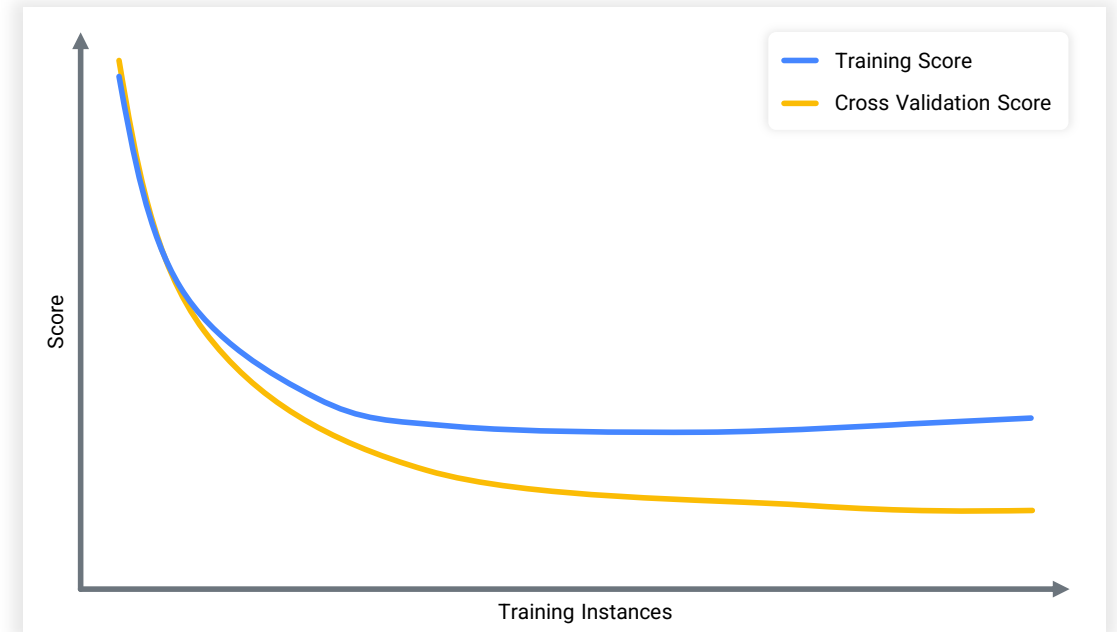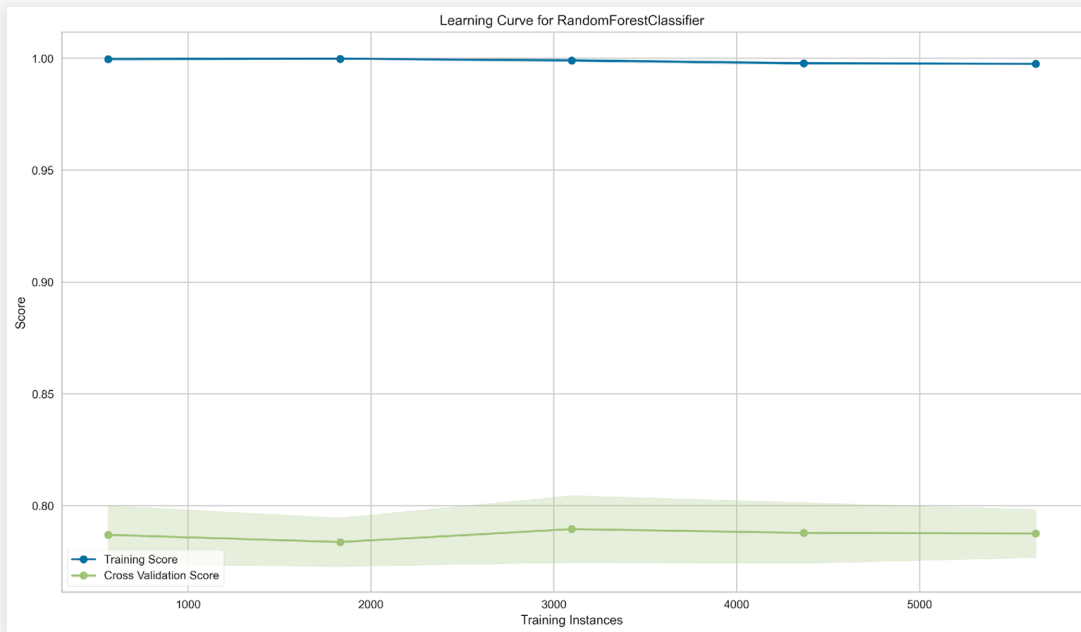


Learning Curve for RandomForestClassifier

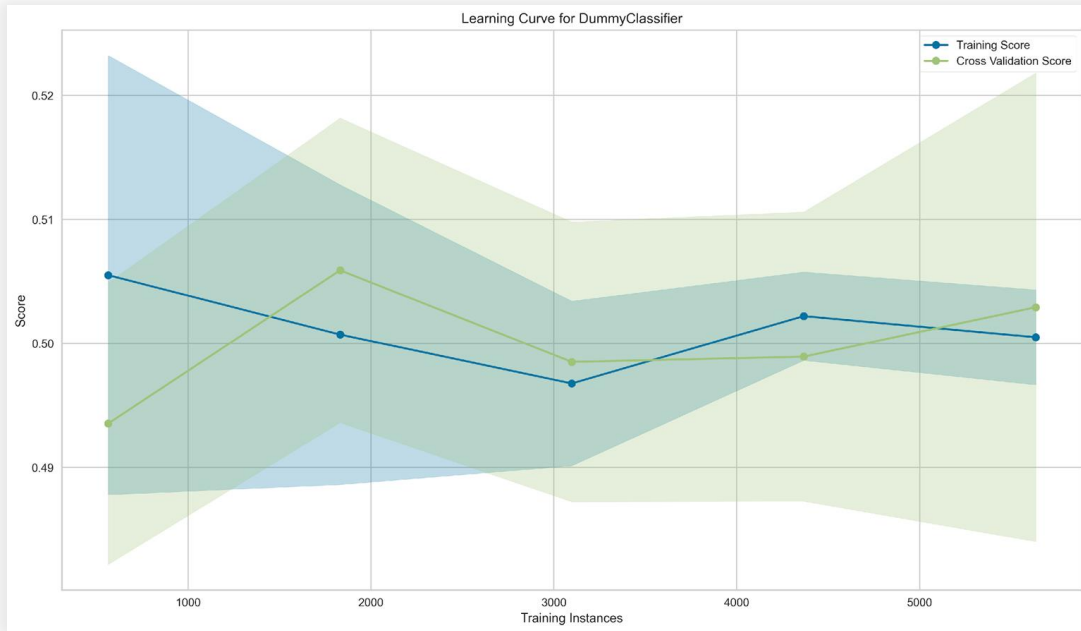Lines are very far apart from one another

Lines would sometimes cross each other

Scores for training set would be significantly higher than training set

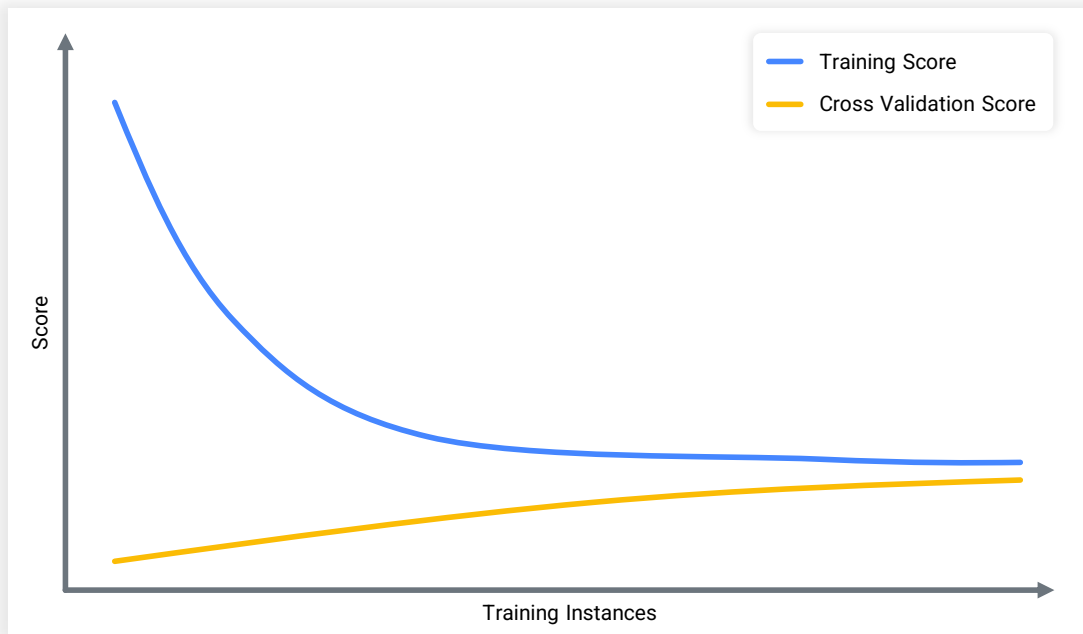# Overfit Characteristics

# Underfit Characteristics



Lines are relatively close to each other

Scores for validation training set are both low

# Knowledge Check



> This model has…

A. High Bias, High Variance
B. High Bias, Low Variance
C. Low Bias, High Variance
D. Low Bias, Low Variance

SPAI

# Practice Time!
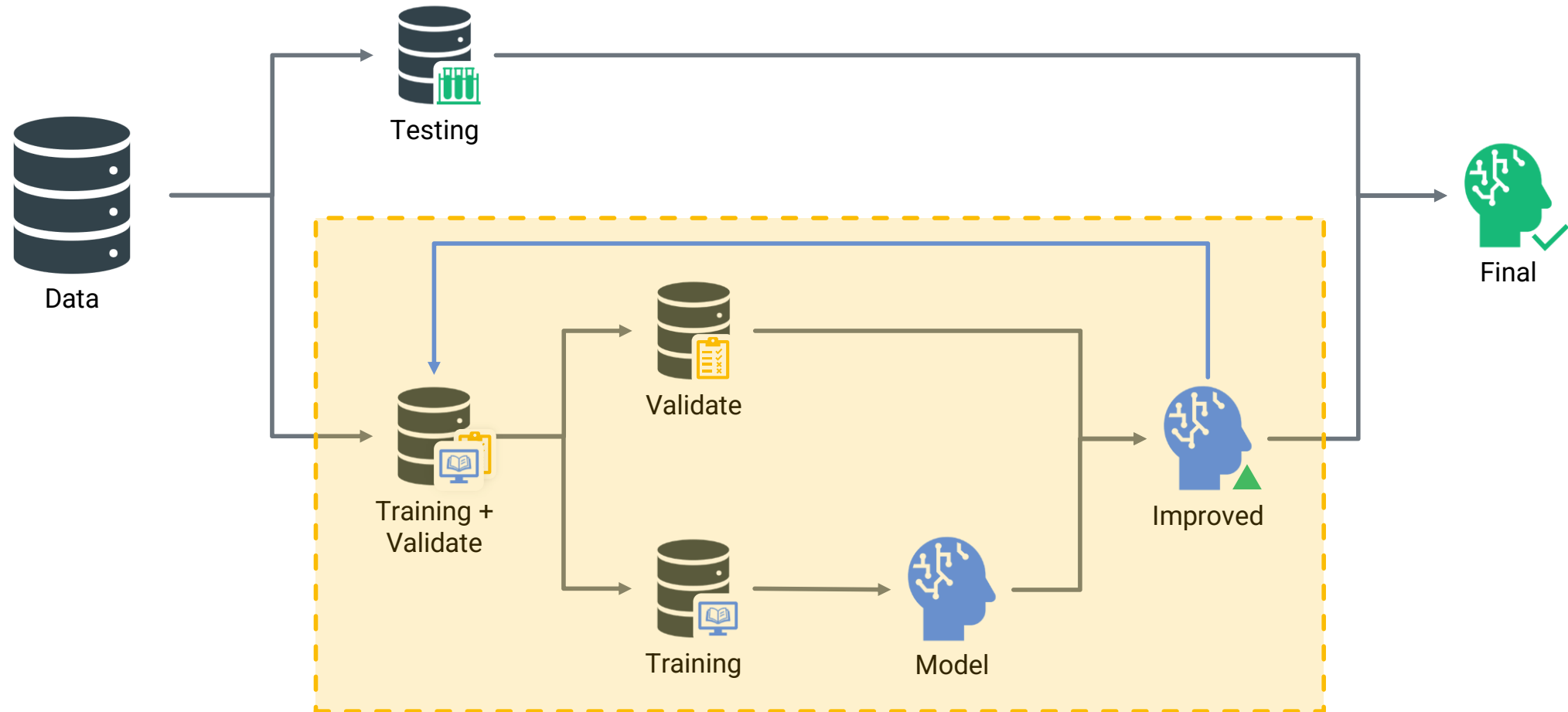
# 5 Minutes

Please attempt exercise 2
We will go through the exercise later
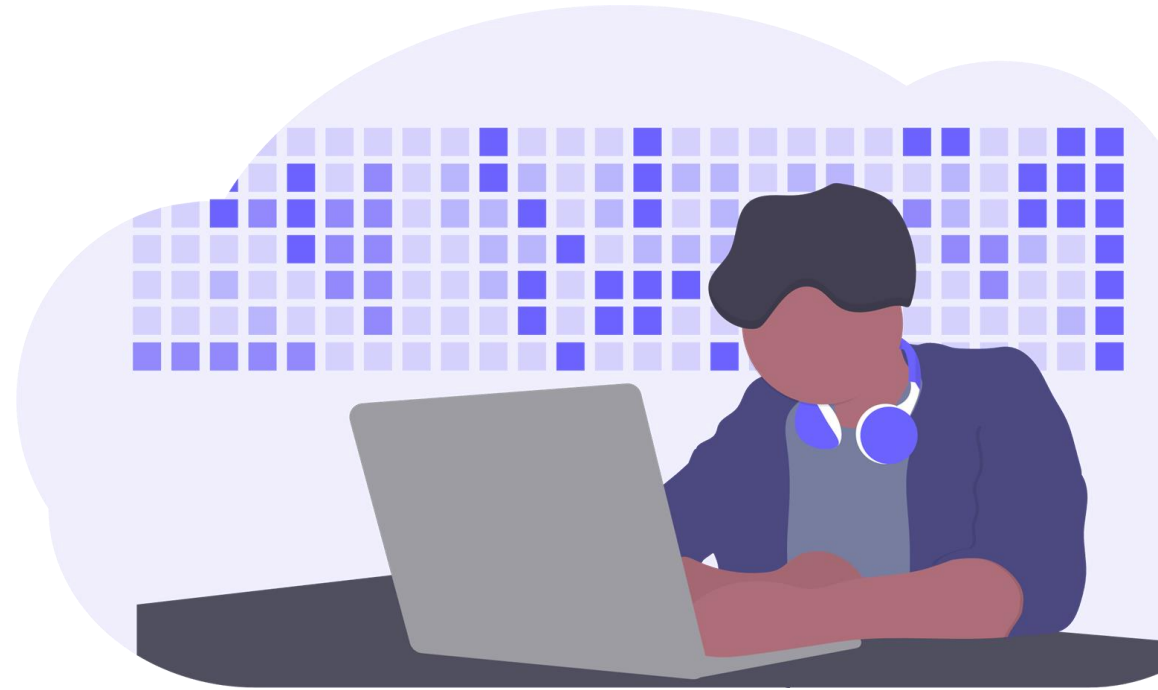
SPAI

# Times up
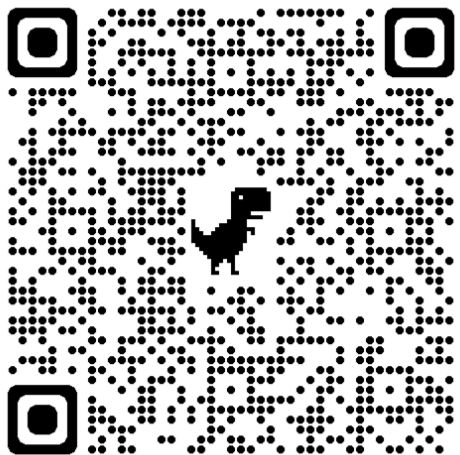
We will now go through the exercises

# Actual Diagram



Data

Testing

Training + Validate

Validate

Training

Model

Improved

Final

SPAI

# SPAI
**An AI Singapore Student Chapter**

# Thank You

# Scan the QR code to mark your attendance

**Attendance**

SPAI