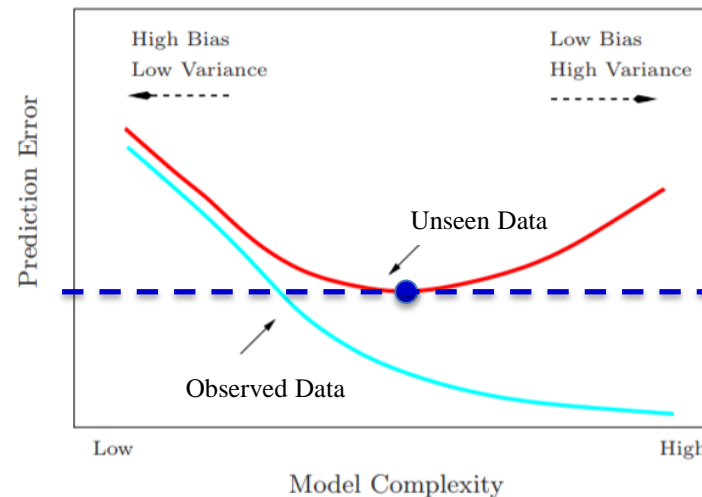# – Module 6 –
# Model Assessment and Selection II

## Outline

- Cross-Validation
  - ➢ *K*-fold Cross-Validation
  - ➢ Leave-one-out Cross-Validation
- Learning Curves

# Generalization

- The <u>central challenge</u> in <u>machine learning</u> is that the algorithm must <u>perform well</u> on <u>new</u>, <u>previously unseen inputs</u> and <u>not just</u> the <u>training data</u> on which the <u>model</u> was <u>trained</u>
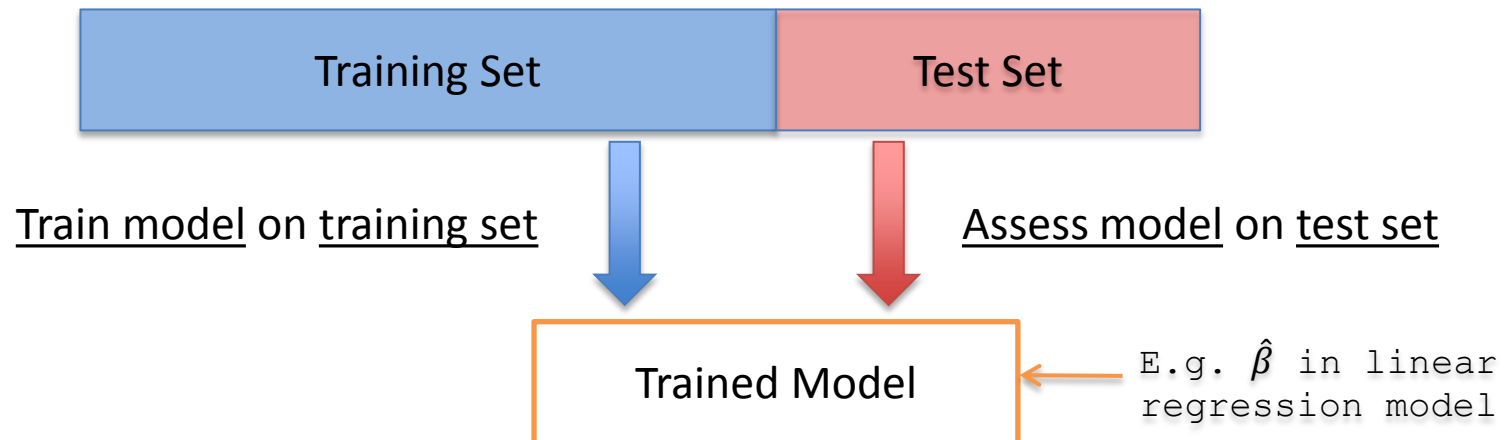


[Source: The Elements of Statistical Learning, ISBN: 978-0387848570]

> A <u>hypothesis</u> $h$ with a <u>low error rate</u> on the <u>training set</u> (observed data) <u>does not mean</u> that it will <u>generalize well</u> to <u>unseen data</u>

- In general, as the <u>model complexity</u> <u>increases</u>, the <u>bias</u> tends to <u>decrease</u> and the <u>variance</u> tends to <u>increase</u>
- The <u>goal</u> is to choose the <u>model complexity</u> to <u>trade</u> <u>bias</u> off with <u>variance</u> in such a way so as to **minimize** the **prediction error** for <u>unseen data</u>
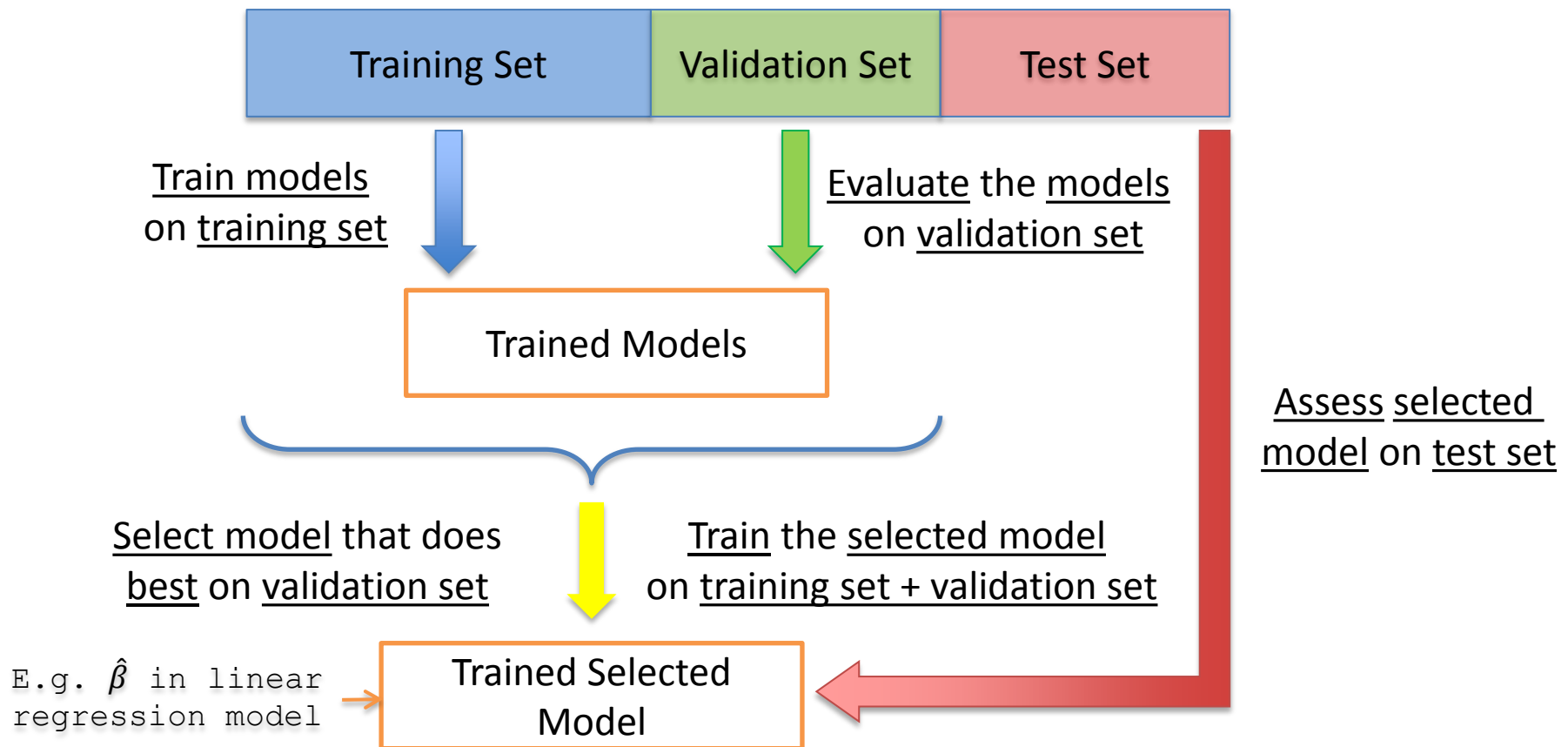
# Model Selection and Assessment (1)

- Recall in <u>machine learning</u>, there are <u>two goals</u>:
  - **Model selection**: <u>estimating</u> the <u>performance</u> of <u>different models</u> (or <u>different model complexities</u>) in order to <u>choose</u> the <u>best one</u>
  - **Model assessment**: having <u>chosen</u> a <u>final model</u>, <u>estimating its</u> <u>prediction error</u> on <u>new data</u>

- If the <u>model</u> (and its <u>complexity</u>) to <u>use</u> is already <u>known</u> ➔ only <u>model assessment</u> is needed

| Training Set | Test Set |
|:---:|:---:|

<u>Train model</u> on <u>training set</u>

<u>Assess model</u> on <u>test set</u>

Trained Model ← E.g. $\hat{\beta}$ in linear regression model

# Model Selection and Assessment (2)

- If the <u>model</u> (and its <u>complexity</u>) to <u>use</u> is <u>not known</u> ➜ <u>need</u> both <u>model selection</u> and <u>model assessment</u>

| Training Set | Validation Set | Test Set |
|---|---|---|

<u>Train models</u> on <u>training set</u>

<u>Evaluate</u> the <u>models</u> on <u>validation set</u>

Trained Models

<u>Assess selected</u> <u>model</u> on <u>test set</u>

<u>Select model</u> that does <u>best</u> on <u>validation set</u>

<u>Train</u> the <u>selected model</u> on <u>training set + validation set</u>

E.g. $\hat{\beta}$ in linear regression model

Trained Selected Model

# Evaluating and Choosing the Best Hypothesis (1)

- We want to <u>learn</u> a <u>model</u> (or <u>hypothesis</u> $h$), that <u>fits</u> the <u>future data</u> <u>best</u>

- Assumptions:
  - <u>future data</u> **stationarity**:
    - there is a **probability distribution** over <u>samples</u> (i.e., parameters such as <u>mean</u> and <u>variance</u>) that remains <u>stationary</u> (do not change) <u>over time</u>

  - <u>best fit</u>:
    - the **error rate** of a <u>hypothesis</u> is defined as the <u>proportion</u> of <u>mistakes</u> it makes (i.e., the proportion of times that $\hat{y} \neq y$ for a $(x,y)$ sample)

# Evaluating and Choosing the Best Hypothesis (2)

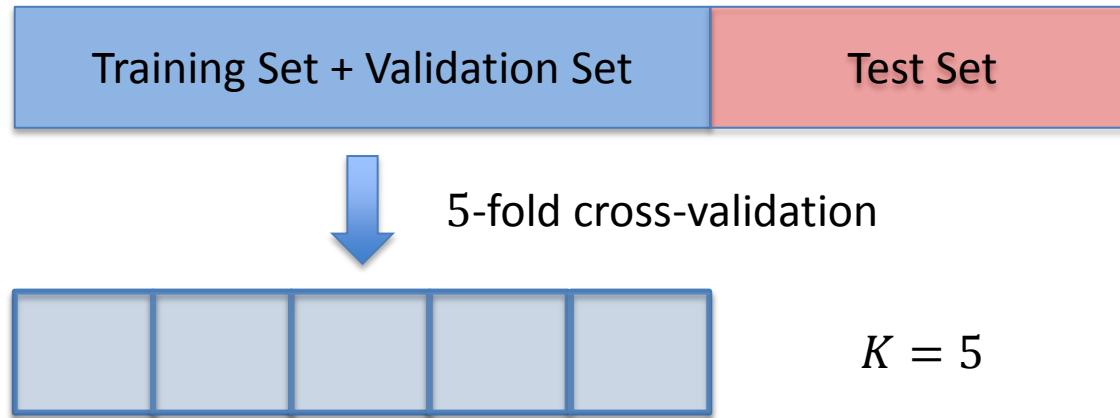- To get an <u>accurate</u> <u>evaluation</u> of a <u>hypothesis</u>, we need to <u>test</u> it on a set of <u>samples</u> that it has <u>not seen</u> yet. The <u>simplest approach</u>, sometimes called **holdout cross-validation**, is one we have already seen
  - (randomly) <u>split</u> the <u>available data</u> into
    - a <u>training set</u> from which the <u>learning algorithm</u> <u>produces</u> $h$
    - a <u>test set</u> on which the <u>accuracy</u> of $h$ is <u>evaluated</u>
      (**<span style="color:red">CAUTION</span>**: <u>NEVER</u> use <u>test set</u> to <u>train</u> your <u>model</u>)

  - However, <u>holdout cross-validation</u> has the <u>disadvantage</u> that it <u>fails</u> to <u>use</u> <u>all</u> the <u>available data</u>
    - if we <u>use</u> <u>half the data</u> for the <u>test set</u>, then we are only <u>training on</u> <u>half the data</u>
      ➔ may get a <u>poor hypothesis</u>
    - if we <u>reserve</u> only <u>10%</u> of the <u>data</u> for the <u>test set</u>
      ➔ may get a <u>poor estimate</u> of the <u>actual accuracy</u>

# Cross-Validation (1)

- Dividing a <u>dataset</u> into a <u>fixed training set</u> and a <u>fixed validation set</u> can be <u>problematic</u> when the given <u>dataset</u> is <u>small</u>
  - ➔ results in the <u>validation set</u> being <u>too small</u>
  - ➔ a <u>small</u> <u>validation set</u> is <u>unable</u> to <u>accurately</u> <u>estimate</u> the <u>prediction error</u>

- A <u>widely used</u> <u>method</u> for <u>estimating</u> <u>prediction error</u> is **cross-validation**
  - a <u>procedure</u> that <u>repeats</u> <u>training and validation</u> on <u>different</u> (randomly chosen) <u>subsets</u> of the <u>original training set</u>
  - incurs <u>increased</u> <u>computational cost</u>

    - ➔ $K$-**fold cross-validation**
    - ➔ **leave-one-out cross-validation** (LOOCV)

# Cross-Validation (2)

- $K$-fold cross-validation allows us to make more out of the data and still get an accurate estimate
  - the idea is that each sample serves double-duty: as training set and validation set

| Training Set + Validation Set | Test Set |
|---|---|

↓

5-fold cross-validation

$K = 5$

- Typical values for $K$ are $K = 5$ and $K = 10$
  - enough to give an estimate that is statistically likely to be accurate
  - but at a cost of $K$ times longer computation time
- The extreme is $K = N$, where $N$ is the number of samples in the training set, known as leave-one-out cross-validation

# Cross-Validation (3)

- $K$-fold cross-validation:



| | Training set |
| | Validation set |

- sub-divide the <u>training set</u> into $K$ <u>non-overlapping subsets</u> <u>of equal size</u> (e.g., $K = 5$)
- <u>perform $K$ trials</u> of <u>learning</u> - on <u>trial $k$,</u>
  - <u>all subsets except</u> the $k$-th subset is <u>used</u> as the **training set** to <u>train</u> the <u>model</u>
  - the $k$-<u>th subset</u> is used as the **validation set** to <u>evaluate</u> the <u>trained model</u>
- the **cross-validation estimate** of **prediction error** of the given <u>model</u> is the <u>average</u> of the <u>cross-validation estimates</u> across $K$ <u>trials</u>
  - ➔ it is <u>expected</u> that the <u>average cross-validation estimate</u> of the $K$ <u>trials</u> should be a <u>better estimate</u> than that from a <u>single trial</u>

# Cross-Validation (4)

- Specifically, for a $p$-th <u>degree polynomial regression</u>
  - <u>divide</u> the <u>training set</u> into $K$ <u>equal subsets</u>
  - for <u>each</u> $k = 1, 2, \ldots, K$
    - <u>train</u> the <u>model</u> with <u>complexity</u> $p$ on the $K - 1$ subsets (<u>training set</u>) with $k$-th <u>subset removed</u>, to obtain the <u>least squares estimate</u>, $\hat{\beta}^{-k}$
    - <u>evaluate</u> the <u>trained model</u> using $\hat{\beta}^{-k}$ on the $k$-th subset (<u>validation set</u>) to obtain **cross-validation estimate**, $Err_{CV}^k$
      (which can be <u>mean absolute error</u> (MAE), <u>root mean squared error</u> (RMS), etc.)
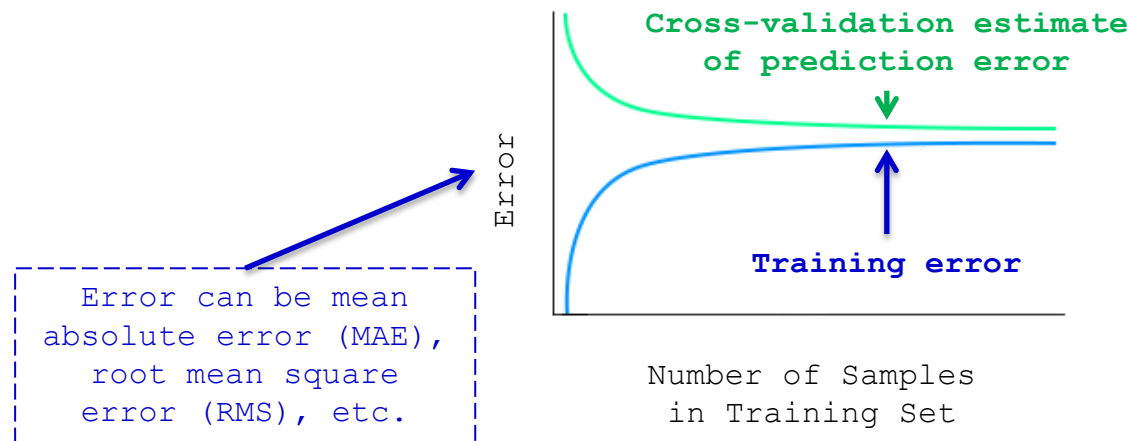  - <u>compute</u> the <u>average cross-validation estimate</u>, $Err_{CV}$, <u>across</u> $k$

$$Err_{CV} = \frac{1}{K} \sum_{k=1}^{K} Err_{CV}^k$$

# Cross-Validation (5)

- To <u>select</u> and <u>assess</u> the <u>final model complexity</u> $p$ (for $p$-th degree <u>polynomial regression</u> model)

  - <u>Final</u> <u>model selection</u>:
    - <u>repeat</u> the $K$-fold cross-validation for <u>different values</u> of $p$
    - <u>select</u> the <u>value</u> of $p$ that gives the <u>smallest</u> $Err_{CV}$

  - <u>Final</u> <u>model assessment</u>:
    - <u>train</u> <u>final model</u> with <u>selected complexity</u> $p$ on <u>all</u> of the <u>training set</u> (<u>training set + validation set</u>)
    - <u>evaluate</u> <u>trained model</u> on the <u>test set</u> to obtain the <u>final prediction error</u>

# Learning Curves (1)
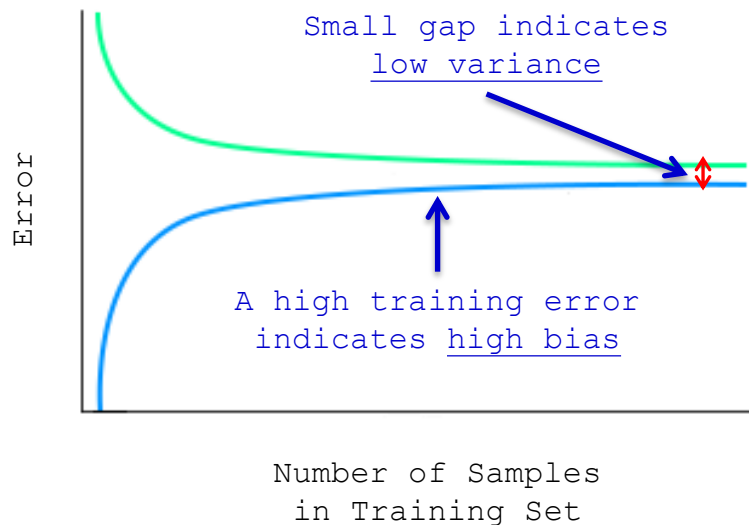
- The **learning curve** for a <u>model</u> is a <u>plot</u> showing the <u>training error</u> and the <u>cross-validation estimate of prediction error</u> as a <u>function of</u> the <u>number of samples</u> $N$ in the <u>training set</u>



**Cross-validation estimate of prediction error**

**Training error**

Error (y-axis)

Number of Samples in Training Set

Error can be mean absolute error (MAE), root mean square error (RMS), etc.
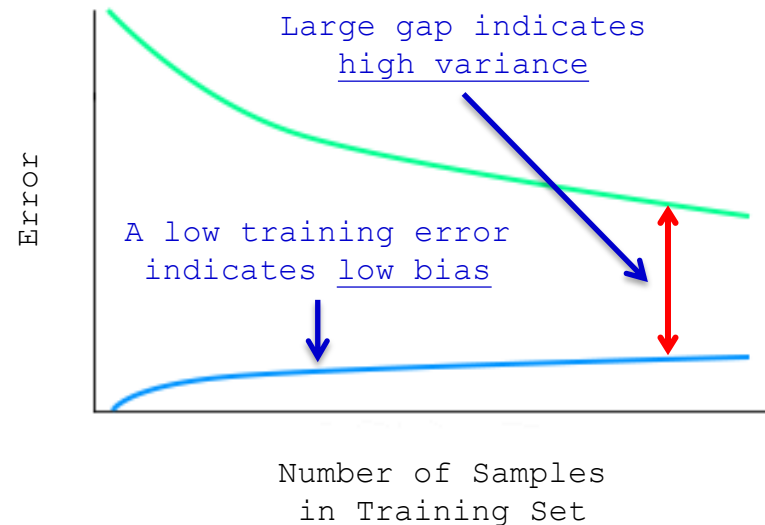
- The <u>learning curve</u> plot is a useful <u>tool</u> to <u>diagnose</u> <u>bias and variance</u> and can help <u>determine</u> whether a <u>model</u> is <u>underfitting</u> (low complexity) or <u>overfitting</u> (high complexity)
- More importantly, a <u>learning curve</u> can also provide <u>insight</u> as to whether <u>more samples</u> need to be <u>collected</u>

# Learning Curves (2)

**<u>Low Complexity Model</u> (underfit)**

Small gap indicates
<u>low variance</u>

Error

A high training error
indicates <u>high bias</u>

Number of Samples
in Training Set

**<u>High Complexity Model</u> (overfit)**

Large gap indicates
<u>high variance</u>

Error

A low training error
indicates <u>low bias</u>

Number of Samples
in Training Set

$$gap = cross\_validation\ estimate - training\ error$$

- Ideally, a <u>model</u> with <u>low bias</u> and <u>low variance</u> is desired
  - for <u>high complexity</u> models (low bias), the <u>cross-validation</u> <u>estimate</u> <u>could</u> <u>converge</u> <u>towards</u> the <u>training error</u> (low variance) if <u>more</u> <u>training samples</u> were <u>added</u> (i.e., increase $N$)