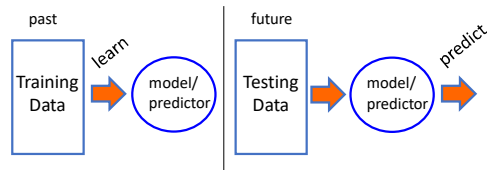


Machine Learning is...

Machine learning is about predicting the future based on the past.



7

5

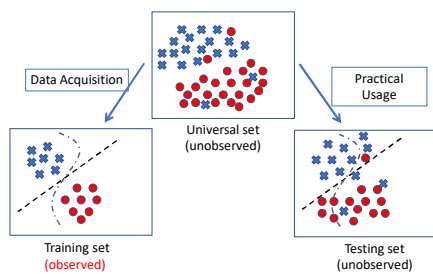
Machine Learning

- **Learning from Data.**
- The "learning" part of machine learning means that ML algorithms attempt to **optimize** along a certain dimension; i.e. they usually try to **minimize error or maximize the likelihood of their predictions** being true.
- Optimizing an **error/loss/cost function**.
- Learning a **mathematical equation** representing the relationship between the inputs and output.

6

6

Training and Testing

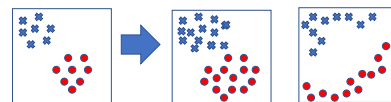


12

7

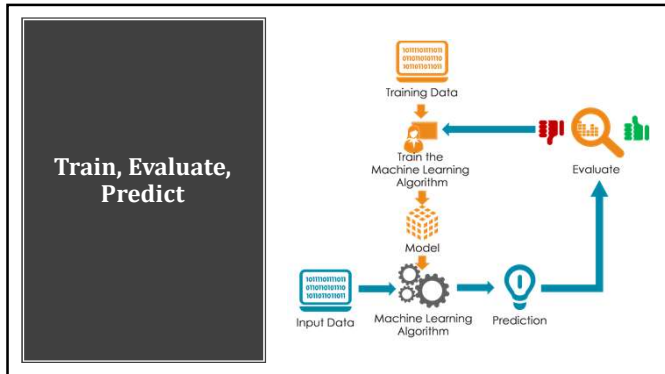
Training and Testing

- **Training** is the process of **making the system able to learn**.
- No free lunch rule:
 - Training set and Testing set come from the **same distribution**
 - Need to make **some assumptions or bias**



13

8



9

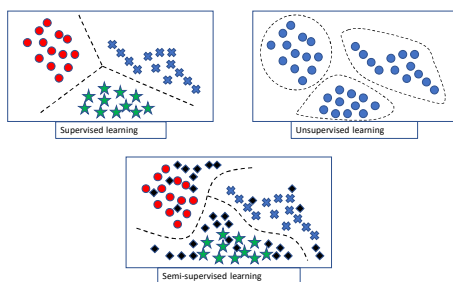
Algorithms: Types of Learning

- **Supervised Learning : Labelled Data** ($\{x_n \in R^d, y_n \in R\}_{n=1}^N$)
 - Prediction
 - Classification (discrete labels), Regression (real values)
- **Unsupervised Learning: Un-labelled Data** ($\{x_n \in R^d\}_{n=1}^N$)
 - Clustering
 - Probability distribution estimation
 - Finding association (in features)
 - Dimension reduction
- **Semi-Supervised Learning : Only a part of data is labelled**
- **Reinforcement Learning: Rewards from sequence of actions**
 - Decision making (robot, chess machine)

14

10

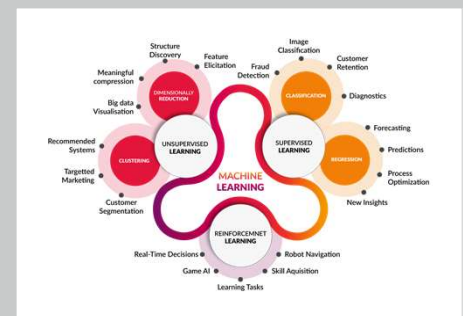
Algorithms: Types of Learning



11

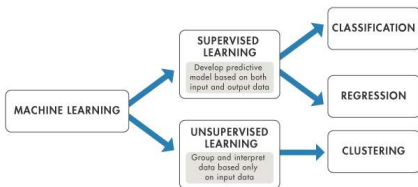
11

Types of Learning



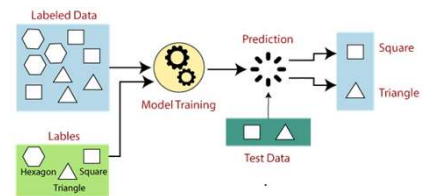
12

In this Course



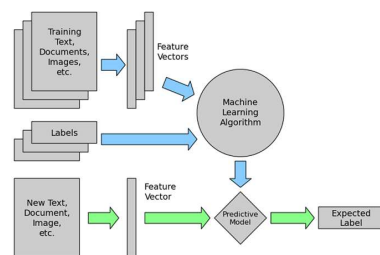
13

Supervised Learning



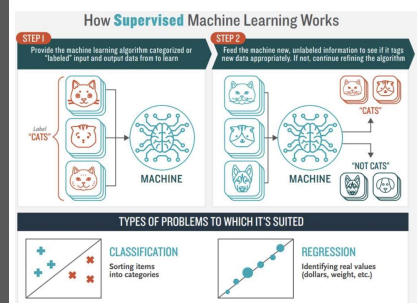
14

Supervised Learning Model

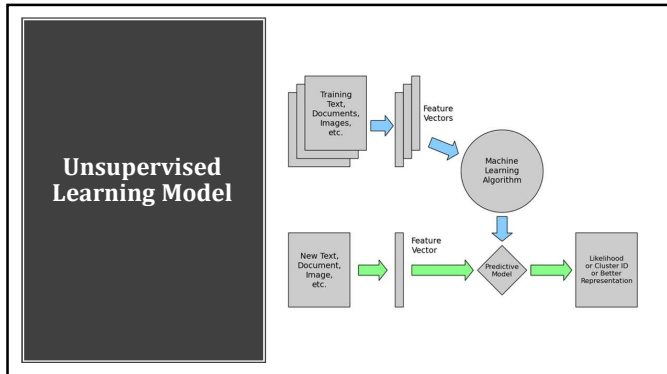


15

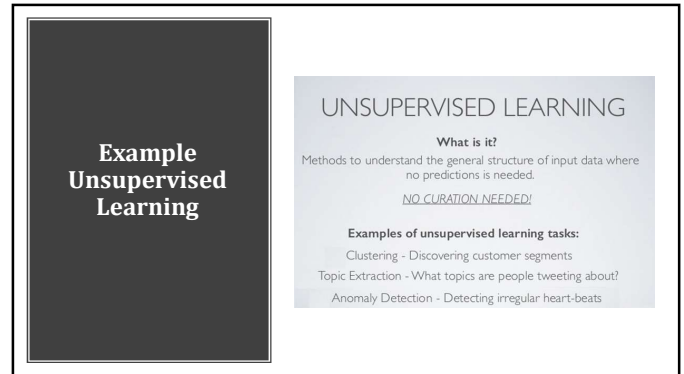
Example Supervised Learning



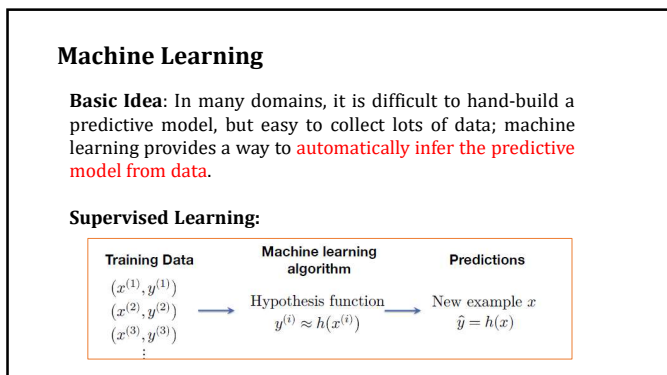
16



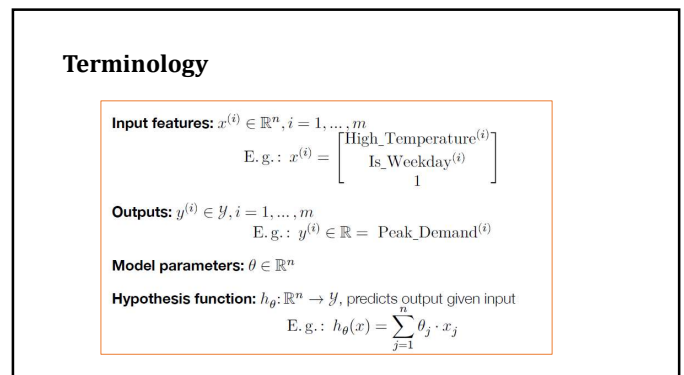
17



18



19



20

Terminology

Loss function: $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$, measures the difference between a prediction and an actual output

$$\text{E.g.: } \ell(\hat{y}, y) = (\hat{y} - y)^2$$

The machine learning optimization problem:

$$\underset{\theta}{\text{minimize}} \quad \sum_{i=1}^m \ell(h_{\theta}(x^{(i)}), y^{(i)})$$

Virtually every machine learning algorithm has this form, just specify

- What is the hypothesis function?
- What is the loss function?
- How do we solve the optimization problem?

21

Example of ML Algorithms

- **Least Squares:** {linear hypothesis, squared loss, (usually) analytical solution}
- **Linear Regression:** {linear hypothesis}
- **Support Vector Machine:** {linear or kernel hypothesis, hinge loss}
- **Neural Network:** {Composed non-linear function, (usually) gradient descent}
- **Decision Tree:** {Hierarchical axis-aligned halfplanes, greedy optimization}
- **Naïve Bayes:** {Linear hypothesis, joint probability under certain independence assumptions, analytical solution}

22

Loss vs. Error vs. Cost Function

- The loss function computes the error for a single training example, while the cost function is the average of the loss functions of the entire training set.
- If we have m training data like this $\{(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots, (\mathbf{x}_m, \mathbf{y}_m)\}$:
 - $\hat{\mathbf{y}}_i$ = output of the model for training example \mathbf{x}_i
 - \mathbf{y}_i = expected output/true value for training example \mathbf{x}_i
- The **loss function** $L(\hat{\mathbf{y}}_i, \mathbf{y}_i)$ defines the error/difference between $\hat{\mathbf{y}}_i$ and \mathbf{y}_i for the single training example \mathbf{x}_i .
- This means, loss refers to **error in model output for an individual sample**.
- If we want to find loss over **all the training examples present in a training-set**, we refer to it as the **cost function** (i.e. total or average loss over all training examples).
- This is the estimate of **total error** computer for the whole training-set.

23

Loss/Error/Cost – Objective Function

The terms cost and loss functions are synonymous some people also call it error function.

The more general scenario is to define an objective function first, which we want to optimize. This objective function could be to

1. maximize the posterior probabilities (e.g., naïve Bayes)
2. maximize a fitness function (genetic programming)
3. maximize the total reward/value function (reinforcement learning)
4. maximize information gain/minimize child node impurities (CART decision tree classification)
5. minimize a mean squared error cost (or loss) function (CART, decision tree regression, linear regression)
6. maximize log-likelihood or minimize cross-entropy loss (or cost) function (ANN), minimize hinge loss (support vector machine)

24

