

**Open Elective Course [OE]**

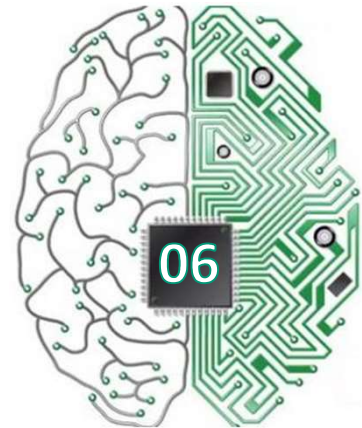
Course Code: CSO507

Winter 2023-24

**Lecture#**

# Deep Learning

## Unit-1: Machine Learning Basics [Part-II]

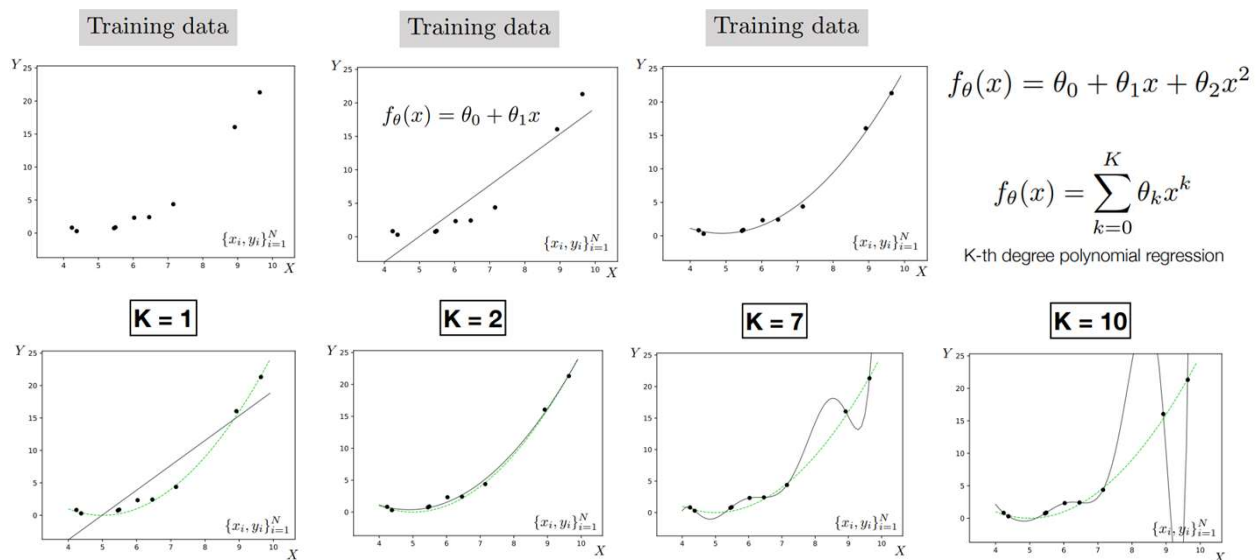
**Course Instructor:****Dr. Monidipa Das**

Assistant Professor

Department of Computer Science and Engineering

Indian Institute of Technology (Indian School of Mines) Dhanbad, Jharkhand 826004, India

## Example (revisited from previous lecture)



Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

# Fitting a model

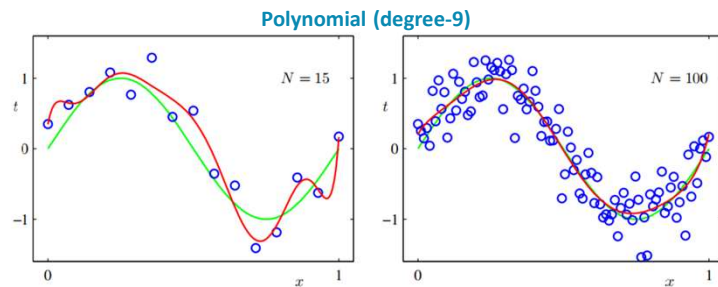


- **Underfitting?**

- Complexify model
  - Add more parameters (more features, more layers, etc.)
- Train longer

- **Overfitting?**

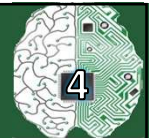
- Perform regularization
- Get more data



**Occam's razor:** Among competing hypotheses that explain known observations equally well, choose the simplest one

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

# Regularization

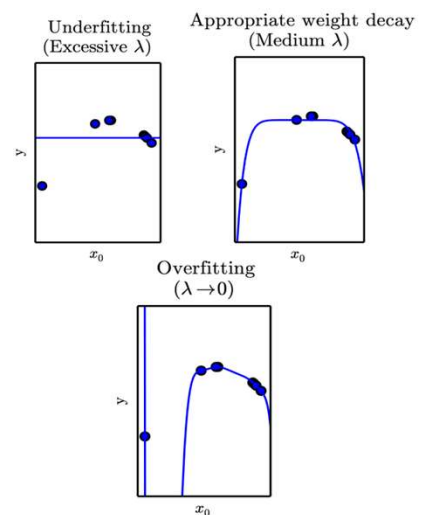


- Procedure aims to avoid the model to overfit the data
- Regularization is able to control the performance of an algorithm
  - Intended to reduce test error but not training error
- **Example: weight decay** for regression problem

$$J(\mathbf{w}) = \text{MSE}_{\text{train}} + \lambda \mathbf{w}^T \mathbf{w}$$

- $J(\cdot)$ : cost function **to be minimized** on training
  - $\lambda$ : control factor of the preference for smaller weight ( $\lambda \geq 0$ )
- **Trades-off** between fitting the training data and being small  $\mathbf{w}$

## Effect of regularization on models with degree 9



Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

# Hyperparameters vs. parameters



**Hyperparameters are higher-level properties for a model**

- Decides model's capacity (or complexity)
- Not be learned during training
- e.g., degree of regression model, weight decay

**Parameters are properties of the training data**

- Learned during training by a ML model
- e.g., weights of regression model

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

## Setting Hyperparameters with Validation Set



- **Setting hyperparameters in training step is not appropriate**
  - Hyperparameters will be set to yield overfitting
  - (e.g., higher degree of regression model,  $\lambda \rightarrow 0$ )
- **Test set will not be seen for training nor model choosing** (hyperparameter setting)
- So, we need **validation set** that the training algorithm does not observe
  - 1. Split validation set from training data
  - 2. Train a model with training data (not including validation set)
  - 3. Validate a model with validation set, update hyperparameters

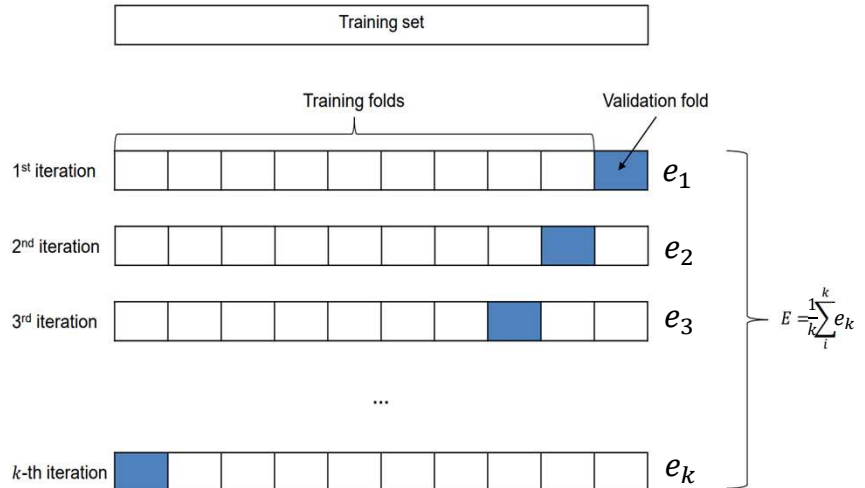
Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

# $k$ -fold Cross Validation



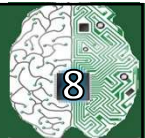
- $k$  - fold cross-validation

- Partition the training data into  $k$  non-overlapping subsets
- On trial  $i$ ,  $i$ th subset of data is used as the test set
- Rest of the data is used as the training set



Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

## $k$ -fold cross validation algorithm



Define  $\text{KFoldXV}(\mathbb{D}, A, L, k)$ :

**Require:**  $\mathbb{D}$ , the given dataset, with elements  $z^{(i)}$

**Require:**  $A$ , the learning algorithm, seen as a function that takes a dataset as input and outputs a learned function

**Require:**  $L$ , the loss function, seen as a function from a learned function  $f$  and an example  $z^{(i)} \in \mathbb{D}$  to a scalar  $\in \mathbb{R}$

**Require:**  $k$ , the number of folds

Split  $\mathbb{D}$  into  $k$  mutually exclusive subsets  $\mathbb{D}_i$ , whose union is  $\mathbb{D}$ .

**for**  $i$  from 1 to  $k$  **do**

$f_i = A(\mathbb{D} \setminus \mathbb{D}_i)$

Train  $A$  on dataset without  $\mathbb{D}_i$

**for**  $z^{(j)}$  in  $\mathbb{D}_i$  **do**

$e_j = L(f_i, z^{(j)})$

Determine errors for samples in  $\mathbb{D}_i$

**end for**

**end for**

**Return**  $e$

Return vector of errors  $e$  for samples in  $\mathbb{D}$

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

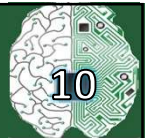
# The No-Free Lunch Theorem (NFLT)



- No free lunch theorem states:
  - Averaged over all distributions, every algorithm has same error classifying unobserved points
    - i.e., no ML algorithm universally better than any other
      - Most sophisticated algorithm has same error rate that merely predicts that every point belongs to same class
- We don't seek a universal learning algorithm because there isn't one

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

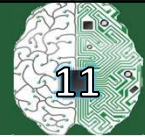
# Generalization Error



- Central challenge of ML is that **the algorithm must perform well on new, previously unseen inputs**
  - Not just those on which our model has been trained
- Ability to perform well on previously unobserved inputs is called **generalization**
- **Generalization error (also called Test Error) definition**
  - Expected value of the error on a new input
    - Expected value is computed as average over inputs taken from distribution encountered in practice
- **How can we affect performance when we observe only the training set?**

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

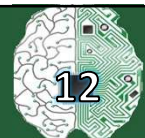
# Statistics provides tools for ML



- The field of statistics provides many tools to achieve the ML goal of solving a task not only on the training set but also to generalize
- Foundational concepts such as
  - Parameter estimation
  - Bias
  - Variance
- They characterize notions of generalization, over- and under-fitting

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

## Point Estimation and Function Estimation



- **Point Estimation:** attempt to provide the single best prediction of some quantity of interest

### Point estimator

$$\hat{\theta}_m = g(x^1, \dots, x^m)$$

$\hat{\theta}$ : point estimator for the property of a model (e.g., expectation)

$m$ : number of data elements

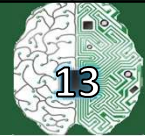
$\{x^1, \dots, x^m\}$ : independent and identically distributed (i.i.d.) data points

$g(\cdot)$ : any estimation function for the given data points

- **Function Estimation:** Point estimation that refers to estimation of relationship between input and target variables
- Most commonly studied properties of point estimators: **Bias** and **Variance**

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

# Bias and Variance of an estimator



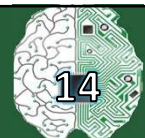
- The **bias of an estimator**  $\hat{\theta}_m = g(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)})$  for parameter  $\theta$  is defined as

$$\text{bias}(\hat{\theta}_m) = E[\hat{\theta}_m] - \theta$$

- The estimator is unbiased if  $\text{bias}(\hat{\theta}_m) = 0$
- The **variance of an estimator**  $\text{Var}(\hat{\theta})$  is how much we expect the estimator to vary as a function of the data sample  $\text{Var}(\hat{\theta}_m) = E[(\hat{\theta}_m - E[\hat{\theta}_m])^2]$
- The square root of the variance is called the standard error, denoted  $SE(\hat{\theta}_m)$
- Bias and Variance measure **two different sources of error of an estimator**

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

# Bias vs Variance Tradeoffs

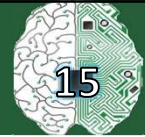


- Bias error is an error from erroneous assumptions in the learning algorithm.
  - High bias can cause an algorithm to miss the relevant relations between features and target outputs (underfitting).
- Variance is an error from sensitivity to small fluctuations in the training set.
  - High variance can cause an algorithm to model the random noise in the training data, rather than the intended outputs (overfitting).

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad



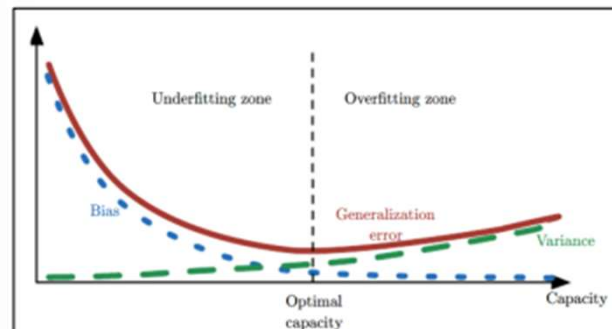
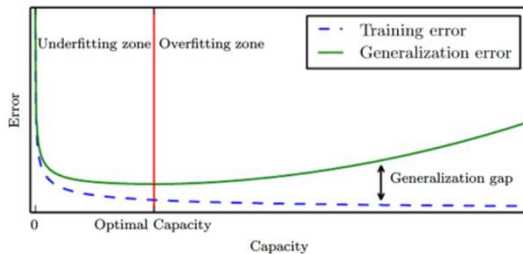
# Underfit-Overfit : Bias-Variance



- Relationship of bias-variance to capacity is similar to underfitting and overfitting relationship to capacity

## Bias-Variance to capacity

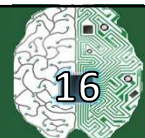
### Model complexity to capacity



- Both have a U-shaped curve of generalization Error as a function of capacity

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

# Negotiating between bias -variance trade-off

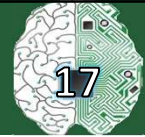


- How to choose between two algorithms, one with a large bias and another with a large variance?
  - Most common approach is to
    - use cross-validation**
    - Alternatively we can **minimize Mean Squared Error** which incorporates both bias and variance

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad



# Mean Squared Error



- Mean Squared Error of an estimate is

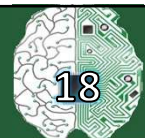
$$\begin{aligned} \text{MSE} &= E[(\hat{\theta}_m - \theta)^2] \\ &= \text{Bias}(\hat{\theta}_m)^2 + \text{Var}(\hat{\theta}_m) \end{aligned}$$

- Minimizing the MSE keeps both bias and variance in check

$$\begin{aligned} E[(\hat{\theta} - \theta)^2] &= E[(\hat{\theta} - E[\hat{\theta}] + E[\hat{\theta}] - \theta)^2] \\ &= E[\hat{\theta}^2 - 2\hat{\theta}E[\hat{\theta}] + E[\hat{\theta}]^2 - 2E[\hat{\theta}]\theta + 2E[\hat{\theta}]\theta - 2\theta E[\hat{\theta}] + \theta^2] \\ &= E[\hat{\theta}^2] - 2E[\hat{\theta}]E[\theta] + E[\hat{\theta}]^2 - 2E[\hat{\theta}]\theta + 2E[\hat{\theta}]\theta - 2\theta E[\hat{\theta}] + \theta^2 \\ &= E[(\hat{\theta} - E[\hat{\theta}])^2] + (E[\hat{\theta}] - \theta)^2 \\ &\quad + 2(E[\hat{\theta}]E[\theta] - \hat{\theta}\theta - E[\hat{\theta}]^2 + E[\hat{\theta}]\theta) \\ &= \underbrace{E[(\hat{\theta} - E[\hat{\theta}])^2]}_{\text{Var}(\hat{\theta})} + \underbrace{(E[\hat{\theta}] - \theta)^2}_{\text{bias}(\hat{\theta})^2} \end{aligned}$$

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

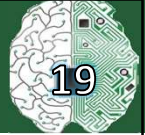
# Bayesian statistics



- Bayesian perspective**
  - Uses probability to reflect degrees of certainty of states of knowledge
  - The dataset is directly observed and so is not random
  - Parameter  $\theta$  is represented as random variable
- The prior**
  - We represent our knowledge of  $\theta$  using the prior probability distribution, notation with  $p(\theta)$ , before observing data
  - Select broad priori distribution (with high degree of uncertainty), such as finite range of volume, with a uniform distribution, or Gaussian.

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

# Mathematical description

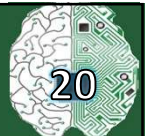


- Set of data samples  $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$
- The dataset is directly observed and so is not random
- Parameter  $\theta$  is represented as random variable
- Combine the data likelihood with the prior via Bayes' rule:

$$\underbrace{p(\theta | x^{(1)}, \dots, x^{(m)})}_{\text{Bayesian inference}} = \frac{\underbrace{p(x^{(1)}, \dots, x^{(m)} | \theta)}_{\text{Likelihood}} \underbrace{p(\theta)}_{\text{Prior}}}{p(x^{(1)}, \dots, x^{(m)})}$$

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

## Relative to Maximum Likelihood Estimate



Make prediction using a **full distribution over  $\theta$**

After observing  $m$  samples, predict distribution over the next data sample,  $x^{(m+1)}$ , is given by:

$$p(x^{(m+1)} | x^{(1)}, \dots, x^{(m)}) = \int p(x^{(m+1)} | \theta) p(\theta | x^{(1)}, \dots, x^{(m)}) d\theta$$

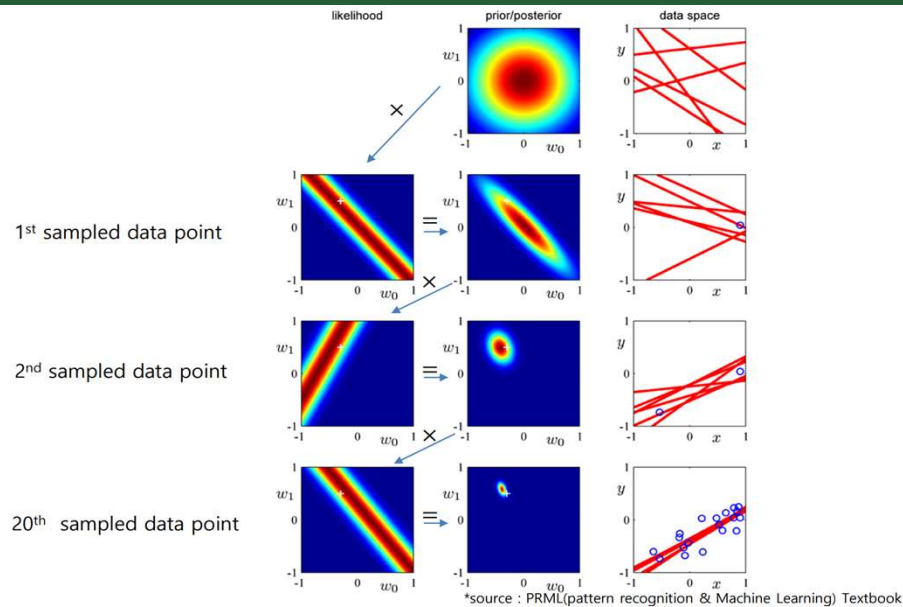
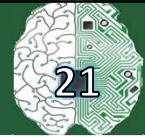
Prior distribution has influence by **shifting probability toward the parameter space**

Bayesian method typically **generalize much better**

But high computational cost

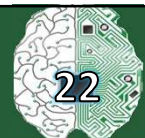
Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

# Example: Bayesian statistics



Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

## Maximum A Posterior (MAP) Estimation



Chose the point of maximal posterior probability

$$\theta_{MAP} = \arg \max_{\theta} \underbrace{p(\theta|x)}_{\text{posterior}} = \arg \max_{\theta} \underbrace{\log p(x|\theta)}_{\text{likelihood}} + \underbrace{\log p(\theta)}_{\text{prior}}$$

Similar with weight decay term

Has the advantage of leveraging information that is brought by the prior

Additional information helps the variance of MAP estimation

But it increase bias

Regularized estimation strategies can be interpreted as making the MAP approximation

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

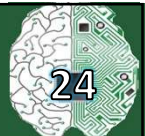
## Exercise (1)



1. How can we prevent underfitting?
  - a) Increase the number of data samples
  - b) Increase the number of features
  - c) Decrease the number of features
  - d) Decrease the number of data samples
  
2. If your ML model (say Neural Network) seems to have high variance, which of the following would be promising to try?
  - a) Increasing the number of parameters/features
  - b) Get more training data
  - c) Get more test data
  - d) Add regularization

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

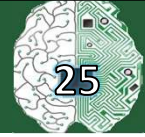
## Exercise (2)



3. What happens when you increase the regularization (hyper)parameter  $\lambda$ ?
  - a) The weights are pushed towards becoming smaller (closer to 0)
  - b) The weights are pushed towards becoming larger (further from 0)
  - c) The total number of weights increased
  - d) None of the above, because there would be no impact on weights
  
4. You are working on building a classifier model to distinguish apples, bananas, and oranges. Suppose your classifier obtains a training set error of 0.3%, and validation set error of 8%. Which of the following would you try?
  - a) Increase regularization (hyper)parameter lambda
  - b) Decrease regularization (hyper)parameter lambda
  - c) Increase the number of training data samples
  - d) Increase the number of parameters/features

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

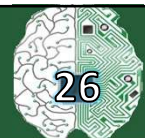
# Recipe for Machine Learning



- All Machine Learning (ML) is an instance of a recipe:
  - Specification and representation of a dataset
  - Math model of the ML architecture
  - Evaluation (training/generalization) with a cost function
  - Optimization procedure
- Example of building an **ML model for linear regression** is shown next

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

## Example: Linear Regression



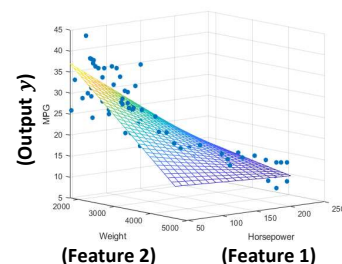
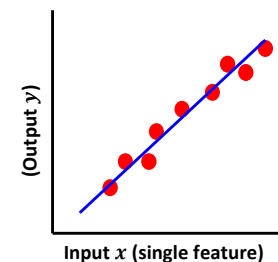
- Linear regression is like fitting a line or (hyper)plane to a set of points
- **Given:** Training data with  $N$  input-output pairs  $\{(x_n, y_n)\}_{n=1}^N$ ,  $x_n \in \mathbb{R}^D$ ,  $y_n \in \mathbb{R}$   
**Goal:** Learn a model to predict the output for new test inputs

- **Model:** A function that approximates the I/O relationship to be a linear model

$$y_n \approx f(x_n) = \mathbf{w}^T \mathbf{x}_n \quad (n = 1, 2, \dots, N)$$

- **Cost Function:** the total error or “loss” of this model over the training data:

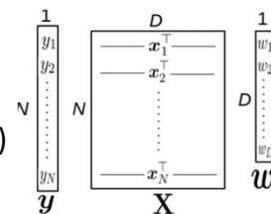
$$L(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \ell(y_n, \mathbf{w}^T \mathbf{x}_n)$$



### Linear Regression with Squared Loss

$$L(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N (y_n - \mathbf{w}^T \mathbf{x}_n)^2$$

- **Optimization:** find the  $\mathbf{w}$  that optimizes (minimizes) the above squared loss



Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

# Cost Function



- Typically the cost function can be written as an average over a training set

$$J(\theta) = E_{(\mathbf{x}, y) \sim \hat{p}_{\text{data}}} (L(f(\mathbf{x}; \theta), y))$$

– Where

- $L$  is the per-example loss function
- $f(\mathbf{x}; \theta)$  is the predicted output when the input is  $\mathbf{x}$
- $\hat{p}_{\text{data}}$  is the empirical distribution

– In supervised learning  $y$  is target output

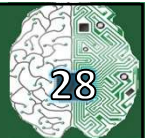
- Trivial to extend to cases:

- Where parameters  $\theta$  and input  $\mathbf{x}$  are arguments or
- Exclude output  $y$  as argument

– For regularization or unsupervised learning

- Typical examples of cost functions will be discussed along with the ML tasks in subsequent lectures

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad



## Questions?

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad