# Linear Regression and Gradient Descent

**Prof. Mingkui Tan**

SCUT Machine Intelligence Laboratory (SMIL)

# Contents

# Contents

**What is Machine Learning?**

Machine Learning composes of three parts:

- Data

- Model

- Loss Function
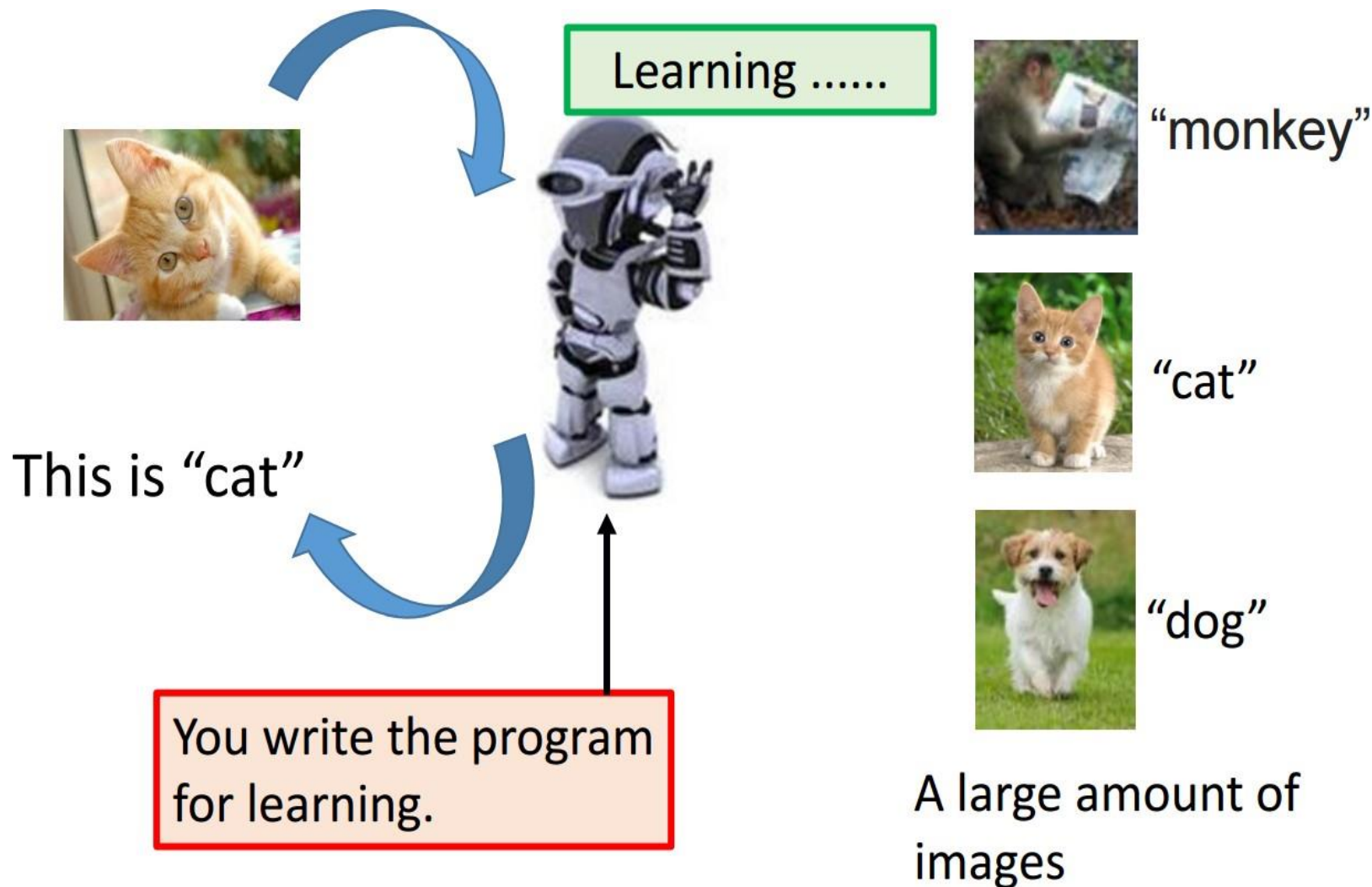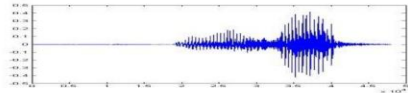
Speech Recognition

# Introduction to Machine Learning



Learning ......

This is "cat"

You write the program for learning.

"monkey"

"cat"

"dog"

A large amount of images

Image Recognition

## Machine Learning ≈ Looking for a Function

■ Speech Recognition

$$f\left( \vcenter{\hbox{[waveform]}} \right) = \text{"How are you"}$$

■ <span style="color:red">Image Recognition</span>

$$f\left( \vcenter{\hbox{[cat image]}} \right) = \text{"Cat"}$$

■ Playing Go

$$f\left( \vcenter{\hbox{[Go board]}} \right) = \text{"5-5"} \quad \text{(next move)}$$

■ Dialogue System

$$f\left( \underset{\text{(what the user said)}}{\text{"Hi"}} \right) = \underset{\text{(system response)}}{\text{"Hello"}}$$

# Introduction to Machine Learning

A set of function | Model $f_1, f_2 \cdots$

$f_1(\quad) = $ "cat"

$f_2(\quad) = $ "money"

$f_1(\quad) = $ "dog"

$f_2(\quad) = $ "snake"

Image Recognition

# Three Main Elements of Machine Learning

**Data** → Different application have different data
Such as face detection, financial application

**Model** → Define the model according to specific problem
Such as recommendation system

**Model evaluation** → Use loss function ( Hinge loss )
(Logistic loss)
(Softmax loss)

**Machine Learning is so simple…**

| Step 1: Define a set of functions | → | Step 2: Goodness of function | → | Step 3: Pick the best function |

Just like putting an elephant into the fridge…

Data:

$$\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$$

**x** is the input, which is usually presented as a <span style="color:red">column vector</span>
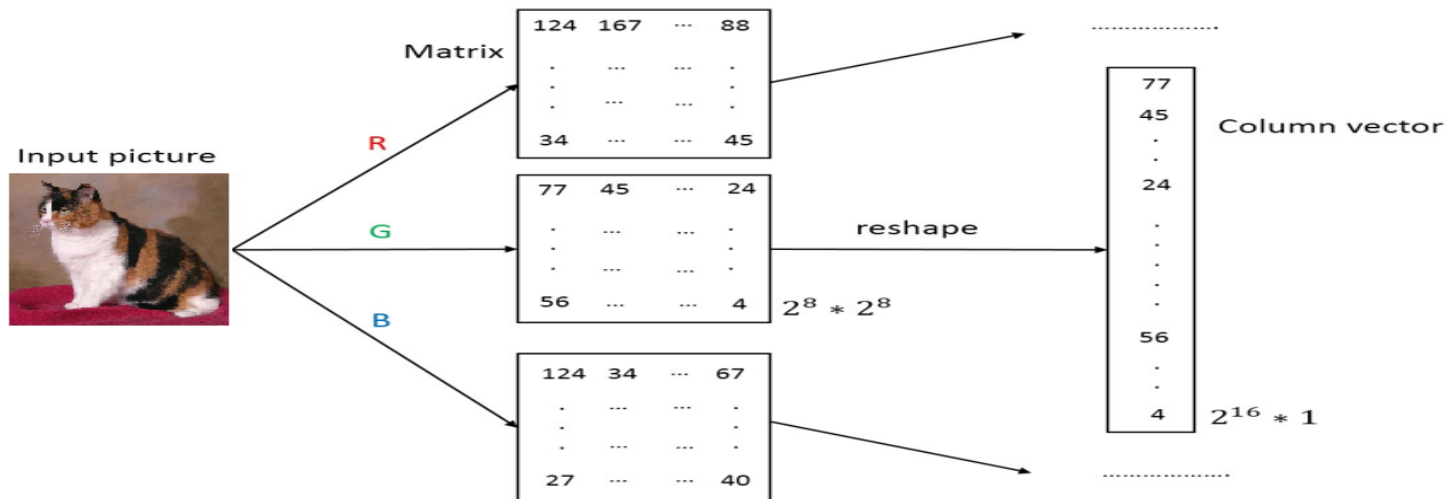
$y$ is the output, for example, a person's name

$n$ is the number of samples

For example, **x** can be a picture stored as a matrix:

# Introduction to Machine Learning

■ Use a function to predict $y$:

$$\hat{y} = f(x)$$

■ However, the prediction may be inconsistent with the ground-truth

■ Calculate the difference by loss function:

$$\mathcal{L}_{\mathcal{D}}(\mathbf{W}) = \sum_{i=1}^{n} l(\hat{y}_i, y_i)$$

where $\mathcal{D}$ refers to data and $\mathbf{W}$ refers to parameter

# Regression

## Loss:

■ Absolute value loss:

$$l(\hat{y}_i, y_i) = |\hat{y}_i - y_i|$$

■ Least squares loss:

$$l(\hat{y}_i, y_i) = \frac{1}{2}(\hat{y}_i - y_i)^2$$

## Total loss function:

$$\mathcal{L}_{\mathcal{D}}(\mathbf{W}) = \sum_{i=1}^{n} l(\hat{y}_i, y_i)$$
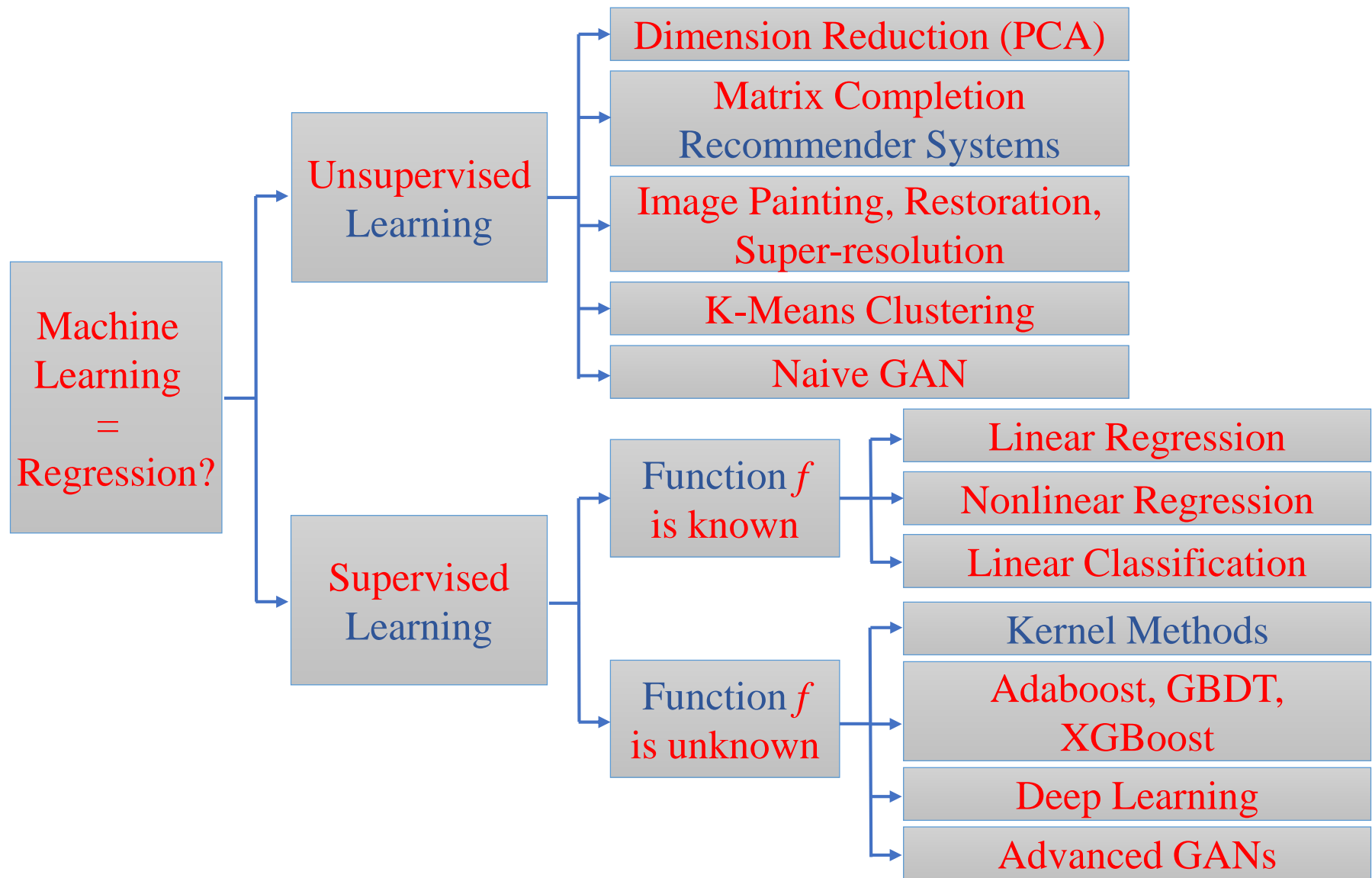
# Regression

■The smaller value of $\mathcal{L}_{\mathcal{D}}$ is better, and loss function $\mathcal{L}_{\mathcal{D}}$ plays a major role in machine learning

## Target:

■ Find the best $f$ by solving the following optimization problem:

$$f^* = \underset{f}{\operatorname{argmin}} \sum_{i=1}^{n} l(f(x), y_i)$$

# Machine Learning



Machine Learning = Regression?

**Unsupervised Learning**
- Dimension Reduction (PCA)
- Matrix Completion / Recommender Systems
- Image Painting, Restoration, Super-resolution
- K-Means Clustering
- Naive GAN

**Supervised Learning**
- Function $f$ is known
  - Linear Regression
  - Nonlinear Regression
  - Linear Classification
- Function $f$ is unknown
  - Kernel Methods
  - Adaboost, GBDT, XGBoost
  - Deep Learning
  - Advanced GANs

■**Supervised learning is the machine learning task of inferring a function from <span style="color:red">labeled training data</span>**
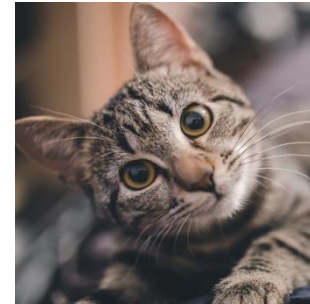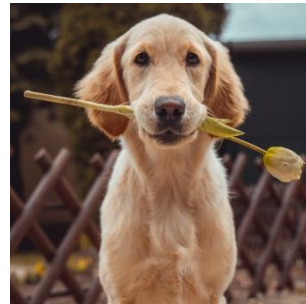
Labeled data



cat                  dog

Unlabeled data

# Dataset for Supervised Learning

## Libsvm dataset

- It contains many classification, regression, multi-label and string data sets
  https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/
- You can use LIBSVM, a package, with these sets
  http://www.csie.ntu.edu.tw/~cjlin/libsvm
- You can also use LIBLINEAR, a linear classifier, with the sets

  https://www.csie.ntu.edu.tw/~cjlin/liblinear/#document
- Other tutorials you can read are as follows:

Tools:  https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/

Guide:  https://www.csie.ntu.edu.tw/~cjlin/papers/guide/guide.pdf

# Introduction to the Format of LIBSVM

## Two properties of data:

- The number of features is large
- Each instance is sparse for most feature values are zero

## Sparse format:

<label1> <index1>:<value1> <index2>:<value2> …
<label2> <index1>:<value1> <index2>:<value2> …

- An example for classification:

  +1 1:2 4:5 \n-1 2:4 \n

  translate to: The points $(2,0,0,5)$ and $(0,4,0,0)$ are assigned to class +1 and class -1 respectively

# Contents

# Linear Regression

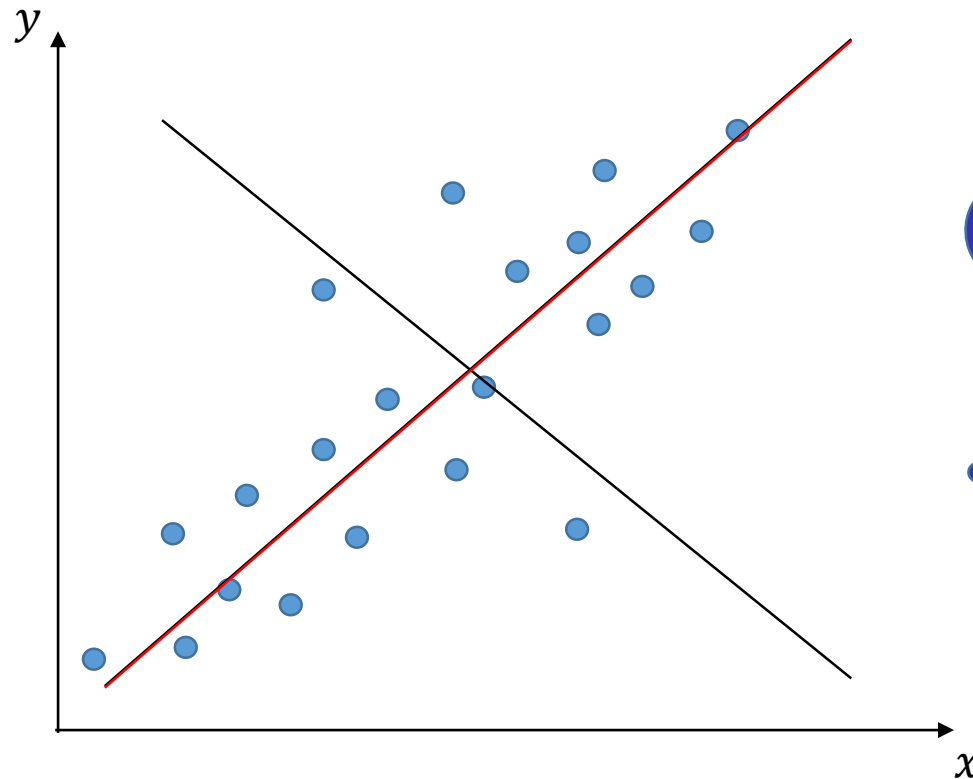Simple linear regression describes the linear relationship between a variable $x$ and a response variable $y$

Simple linear 1-D regression

Which one is better?

# Problem Setup for Regression

- **Inputs**

  Input space: $\mathcal{X} = \{\mathbf{x}_i\}_{i=1}^N, \mathbf{x}_i \in \mathbb{R}^m$

  $N$ is the number of data samples

  $\mathbf{x}_i$ includes $m$ features

- **Outputs**

  Output space: $\mathcal{Y} = \{y_i\}_{i=1}^N, y_i \in \mathbb{R}$

- **Goal**

  Learn a hypothesis / model $f: \mathcal{X} \rightarrow \mathcal{Y}$
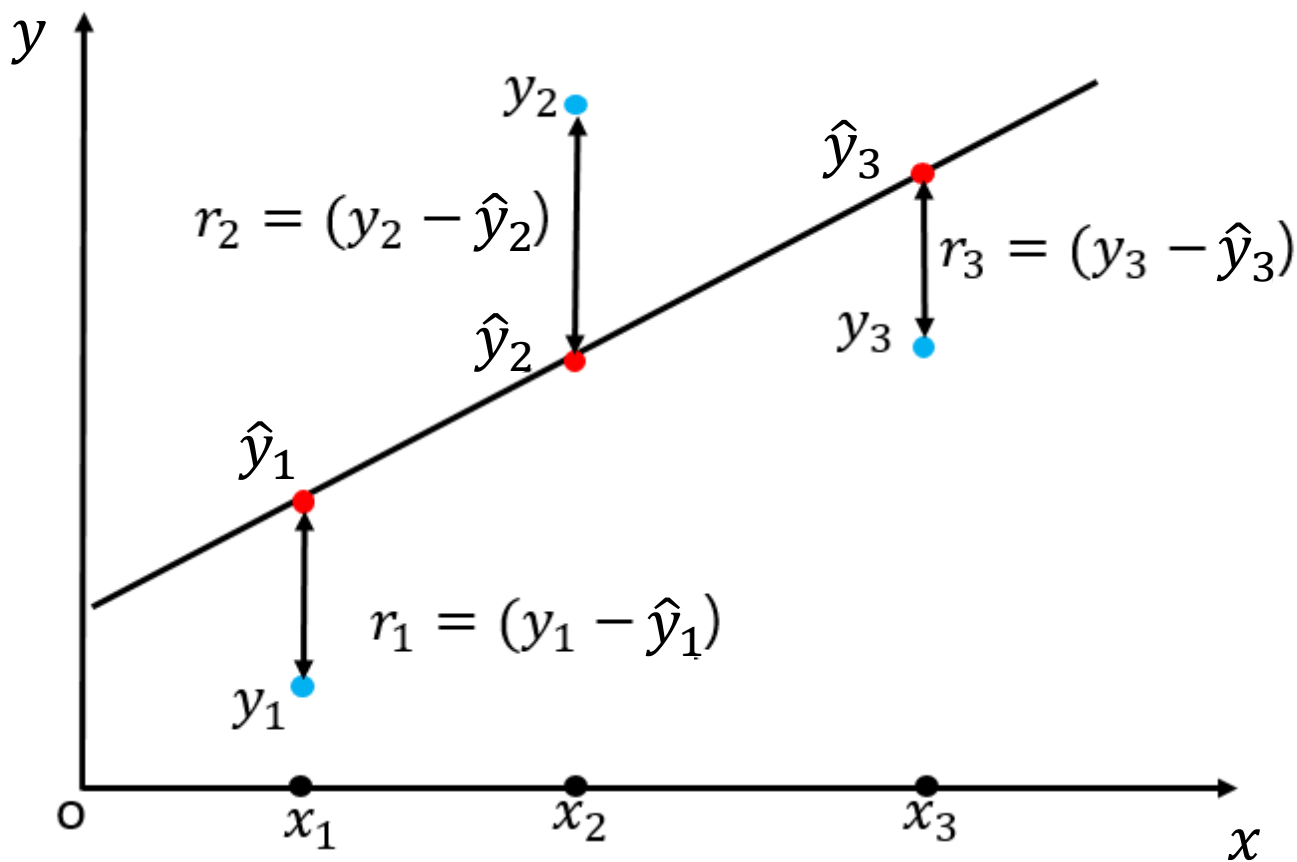
# Linear Regression

Learn $f(\mathbf{x}; \mathbf{w}, b)$ with

- Parameters: $\mathbf{w} \in \mathbb{R}^m, b \in \mathbb{R}$

- Input: $\mathbf{x}$ where $x_i \in \mathbb{R}$ , features for $i \in \{1, \cdots, m\}$

- Model Function:

$$f(\mathbf{x}; \mathbf{w}, b) = w_1 x_1 + \cdots + w_m x_m + b$$

$$= \sum_{i=1}^{m} w_i x_i + b$$

$$= \mathbf{w}^{\mathrm{T}} \mathbf{x} + b$$

■ **What makes a good model?**

# Performance Measure for Regression

■ Least squared loss

$$\mathcal{L}_{\mathcal{D}}(\mathbf{w}, b) = \frac{1}{2}\sum_{i=1}^{n}(y_i - f(\mathbf{x}_i; \mathbf{w}, b))^2$$

$$= \frac{1}{2}\sum_{i=1}^{n}(y_i - \hat{y}_i)^2$$

Training: find minimizer of least squared loss

$$\mathbf{w}^*, b^* = \operatorname*{argmin}_{\mathbf{w}, b} \mathcal{L}_{\mathcal{D}}(\mathbf{w}, b)$$

# Contents

# Matrix Presentation for Loss Function

In order to simplify our proof, we introduce augmented matrix and augmented vector and still represent them by **w** and **X**.

$$\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_i, \ldots, \mathbf{x}_n)^{\mathrm{T}}$$

i.e.

$$\mathbf{x}_i = (1, x_{i1}, x_{i2}, \ldots, x_{in})$$

$$\mathbf{w} = (b, w_1, w_2, \ldots, w_n)^{\mathrm{T}}$$

**Loss function:**

$$\mathcal{L}_{\mathcal{D}}(\mathbf{w}) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2$$

where $\mathbf{X} = \begin{pmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{pmatrix} = \begin{pmatrix} 1 & x_{11} & \cdots & x_{1n} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{nn} \end{pmatrix}, \mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$

# Matrix Presentation for Loss Function

- **Proof:**

$$\mathcal{L}_D(\mathbf{w}) = \frac{1}{2}\sum_{i=1}^{n}(y_i - \mathbf{x}_i\mathbf{w})^2$$

$$= \frac{1}{2}\begin{bmatrix} y_1 - \mathbf{x}_1^{\mathrm{T}}\mathbf{w} \\ \vdots \\ y_n - \mathbf{x}_n^{\mathrm{T}}\mathbf{w} \end{bmatrix}^{\mathrm{T}}\begin{bmatrix} y_1 - \mathbf{x}_1^{\mathrm{T}}\mathbf{w} \\ \vdots \\ y_n - \mathbf{x}_n^{\mathrm{T}}\mathbf{w} \end{bmatrix}$$

$$= \frac{1}{2}\left(\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} - \begin{bmatrix} \mathbf{x}_1^{\mathrm{T}} \\ \vdots \\ \mathbf{x}_n^{\mathrm{T}} \end{bmatrix}\mathbf{w}\right)^{\mathrm{T}}\left(\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} - \begin{bmatrix} \mathbf{x}_1^{\mathrm{T}} \\ \vdots \\ \mathbf{x}_n^{\mathrm{T}} \end{bmatrix}\mathbf{w}\right)$$

$$= \frac{1}{2}(\mathbf{y} - \mathbf{Xw})^{\mathrm{T}}(\mathbf{y} - \mathbf{Xw})$$

$$= \frac{1}{2}\|\mathbf{y} - \mathbf{Xw}\|_2^2$$

# Analytical Solution

How to address the linear regression question?

■ Closed-form solution to linear regression:

$$\mathcal{L}_D(\mathbf{w}) = \frac{1}{2}(\mathbf{y} - \mathbf{Xw})^{\mathrm{T}}(\mathbf{y} - \mathbf{Xw}) \text{ , Let } \mathbf{a} = \mathbf{y} - \mathbf{Xw},$$

$$\frac{\partial \mathcal{L}_{\mathcal{D}}(\mathbf{w})}{\partial \mathbf{w}} = \frac{\partial \mathbf{a}}{\partial \mathbf{w}} \frac{\partial (\frac{1}{2} \mathbf{a}^T \mathbf{a})}{\partial \mathbf{a}}$$

$$= \frac{1}{2} \frac{\partial \mathbf{a}}{\partial \mathbf{w}} (2\mathbf{a})$$

$$= \frac{\partial (\mathbf{y} - \mathbf{Xw})}{\partial \mathbf{w}} (\mathbf{y} - \mathbf{Xw})$$

$$= -\mathbf{X}^{\mathrm{T}}(\mathbf{y} - \mathbf{Xw})$$

Since $\mathcal{L}_{\mathcal{D}}(\mathbf{w})$ is a convex function, $\frac{\partial \mathcal{L}_{\mathcal{D}}(\mathbf{w})}{\partial \mathbf{w}} = 0$ derive $\mathbf{w}^*$

# Analytical Solution

- Assuming $\left|\mathbf{X}^{\mathrm{T}}\mathbf{X}\right| \neq 0$

- Let
$$\frac{\partial \mathcal{L}_{\mathcal{D}}(\mathbf{w})}{\partial \mathbf{w}} = -\mathbf{X}^{\mathrm{T}}\mathbf{y} + \mathbf{X}^{\mathrm{T}}\mathbf{X}\mathbf{w} = 0$$

$$\implies \mathbf{X}^{\mathrm{T}}\mathbf{X}\mathbf{w} = \mathbf{X}^{\mathrm{T}}\mathbf{y}$$

$$\implies \mathbf{w} = (\mathbf{X}^{\mathrm{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathrm{T}}\mathbf{y}$$

- Solve the optimal parameter $\mathbf{w}^{*}$

$$\mathbf{w}^{*} = \operatorname*{argmin}_{\boldsymbol{w}} \mathcal{L}_{\mathcal{D}}(\mathbf{w}) = (\mathbf{X}^{\mathrm{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathrm{T}}\mathbf{y}$$

There are two challenges left to address about the analytical solution $\mathbf{w}^* = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$ :

■ Many matrices are not invertible

Necessary and Sufficient Condition:

If $\mathbf{X}$ is a matrix of $m$ rows and $n$ columns $(n \leq m)$,

$$\left|\mathbf{X}^T\mathbf{X}\right| \neq 0 \iff rank(\mathbf{X}) = n$$

■ The inverse of a large matrix needs huge memory, which takes $O(m^3)$ to compute.

# Contents

# Gradient Descent

■ Get the best **w** by minimizing a loss function $\mathcal{L}_{\mathcal{D}}(\mathbf{w})$

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \mathcal{L}_{\mathcal{D}}(\mathbf{w})$$

# Descent Direction

- We use $\mathbf{d} = -\dfrac{\partial \mathcal{L}_{\mathcal{D}}(\mathbf{w})}{\partial \mathbf{w}}$ as the direction of optimization
- Gradient (vector of partial derivatives)

$$\frac{\partial \mathcal{L}_{\mathcal{D}}(\mathbf{w})}{\partial \mathbf{w}} = \begin{bmatrix} \dfrac{\partial \mathcal{L}_{\mathcal{D}}(\mathbf{w})}{\partial w_1} \\ \dfrac{\partial \mathcal{L}_{D}(\mathbf{w})}{\partial w_2} \\ \vdots \\ \dfrac{\partial \mathcal{L}_{D}(\mathbf{w})}{\partial w_m} \end{bmatrix}$$

(We always write a vector into column form)

- Why $\mathcal{L}_{\mathcal{D}}(\mathbf{w}') = \mathcal{L}_{D}(\mathbf{w} + \eta\mathbf{d}) \leq \mathcal{L}_{\mathcal{D}}(\mathbf{w}), \ \eta \to 0^+$ ?

# Descent Direction

By Taylor expansion, when $\eta \to 0^+$:

$$\mathcal{L}_{\mathcal{D}}(\mathbf{w} + \eta\mathbf{d}) = \mathcal{L}_{\mathcal{D}}(\mathbf{w}) + \left(\frac{\partial \mathcal{L}_{\mathcal{D}}(\mathbf{w})}{\partial \mathbf{w}}\right)^{\mathrm{T}} \eta\mathbf{d} + o(\eta\mathbf{d})$$

$$= \mathcal{L}_{\mathcal{D}}(\mathbf{w}) + \eta' \left(\frac{\partial \mathcal{L}_{\mathcal{D}}(\mathbf{w})}{\partial \mathbf{w}}\right)^{\mathrm{T}} \mathbf{d}$$

Note that $\eta' > 0$ and

$$\eta' \left(\frac{\partial \mathcal{L}_{\mathcal{D}}(\mathbf{w})}{\partial \mathbf{w}}\right)^{\mathrm{T}} \mathbf{d} = -\eta' \mathbf{d}^{\mathrm{T}} \mathbf{d} \leq 0$$
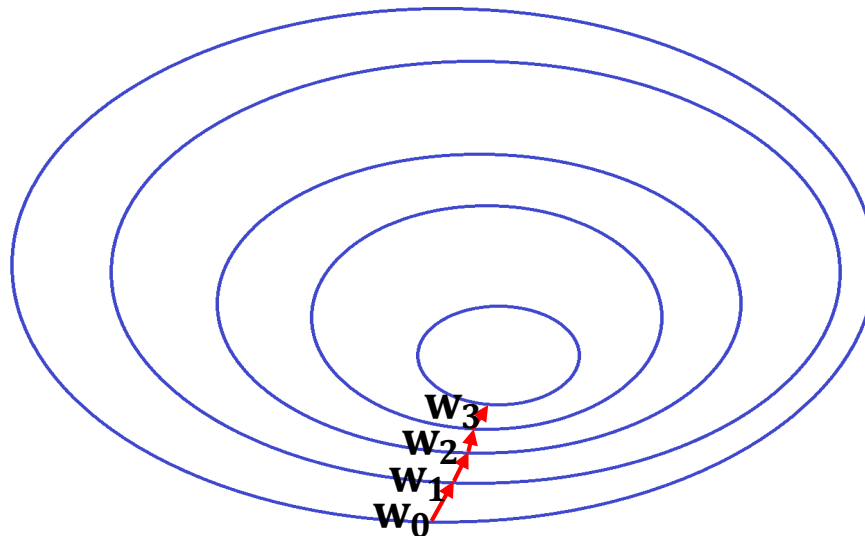
We have:

$$\mathcal{L}_{\mathcal{D}}(\mathbf{w}') = \mathcal{L}_{\mathcal{D}}(\mathbf{w} + \eta\mathbf{d}) \leq \mathcal{L}_{\mathcal{D}}(\mathbf{w})$$

# Gradient Descent: Update Parameters

Minimize loss by repeated gradient steps (when no closed form):

- Compute gradient of loss with respect to parameters $\frac{\partial \mathcal{L}_{\mathcal{D}}(\mathbf{w})}{\partial \mathbf{w}}$

- Update parameters with learning rate $\eta$

$$\mathbf{w}' = \mathbf{w} - \eta \frac{\partial \mathcal{L}_{\mathcal{D}}(\mathbf{w})}{\partial \mathbf{w}}$$

# Appropriate Value of Learning Rate

**Learning rate $\eta$** has a large impact on convergence

- Too large $\eta \Rightarrow$ oscillate and may even diverge

- Too small $\eta \Rightarrow$ too slow to converge

Adaptive learning rate (For example) :

- Set larger learning rate at the beginning

- Use relatively smaller learning rate in the later epochs

- Decrease the learning rate:

$$\eta^{t+1} = \frac{\eta^t}{t + 1}$$

# Thank You