

# Machine Learning

## 1. Introduction

**Definition:** If a machine's performance, denoted by  $P$ , on the task  $T$  improves with experience  $E$ , then the machine is said to learn from  $E$  with respect to  $T$  and  $P$ .

**Taxonomy of ML:**

data type	learning type
1. Supervised Learning	1. Offline Learning
2. Unsupervised Learning	2. Online Learning
3. Reinforcement Learning	
4. Semi-supervised Learning	

**Supervised Learning:** It could help us with two tasks in general. The first is regression and the second is classification. In regression, the output is a continuous value, while in classification, the output is a discrete value.

**Unsupervised Learning:** It is used to find the hidden patterns in the data. It is used in clustering and association.

**Reinforcement Learning:** The agent learns from the feedback it receives from the environment after each action it takes.

## 2. Mathematical Foundation

### 2.1. Derivative

- For a function  $f : \mathbb{R}^d \mapsto \mathbb{R}$ , the derivative of  $f$  with respect to  $\mathbf{x}$  is:

$$\frac{df}{d\mathbf{x}} = \nabla_{\mathbf{x}} f(\mathbf{x}) = \left[ \frac{df}{dx_1} \quad \frac{df}{dx_2} \quad \cdots \quad \frac{df}{dx_d} \right]^T$$

- For a fixed vector  $\mathbf{b} \in \mathbb{R}^d$ , consider  $f(\mathbf{x}) = \mathbf{x}^T \mathbf{b}$  (dot product), then the derivative of  $f(\mathbf{x})$  with respect to  $\mathbf{x}$  is:

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = \mathbf{b}$$

- For a matrix  $\mathbf{A} \in \mathbb{R}^{d \times d}$ , consider  $f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x}$  (quadratic form), then the derivative of  $f(\mathbf{x})$  with respect to  $\mathbf{x}$  is:

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = (\mathbf{A} + \mathbf{A}^T) \mathbf{x}$$

In most cases, the matrix  $\mathbf{A}$  is symmetric, so the derivative of  $f(\mathbf{x})$  with respect to  $\mathbf{x}$  is:

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = 2\mathbf{A}\mathbf{x}$$

### 2.2. Maximum Likelihood Estimation

For an example of multiple cases from the same probability distribution, parameterized by  $\hat{\theta}$ , there exists a  $\theta$  that maximized the likelihood function of the joint distribution of all the cases, written:

$$\begin{aligned}\hat{\theta}_{\text{ML}} &= \operatorname{argmax}_{\theta \in (0,1)} p_{X_1, \dots, X_m}(X_1, X_2, \dots, X_m; \theta) \\ &= \operatorname{argmax}_{\theta \in (0,1)} \log p_{X_1, \dots, X_m}(X_1, X_2, \dots, X_m; \theta)\end{aligned}$$

Property:

- Consistent: As the number of samples increases, the estimated parameter converges to the true parameter.
- Asymptotically Normal: The estimated parameter is normally distributed as the number of samples increases, and the variance of the distribution decreases as the number of samples increases.

### 2.3. Solution to Linear System

- Rouché-capelle theorem: For a linear system  $\mathbf{X}\mathbf{w} = \mathbf{y}$  where  $\mathbf{X} \in \mathbb{R}^d$  and the augmented matrix  $\tilde{\mathbf{X}} = [\mathbf{X} \ \mathbf{y}]$ 
  1. The system has a unique solution, iff.  $\operatorname{rank}(\mathbf{X}) = \operatorname{rank}(\tilde{\mathbf{X}}) = d$ .
  2. The system has no solution, iff.  $\operatorname{rank}(\mathbf{X}) < \operatorname{rank}(\tilde{\mathbf{X}})$ .
  3. The system has infinitely many solutions, iff.  $\operatorname{rank}(\mathbf{X}) = \operatorname{rank}(\tilde{\mathbf{X}}) < d$ .

## 3. Linear Regression(With offset)

We define:

- $m$ : number of training examples
- $d$ : number of features
- $\mathbf{X} \in \mathbb{R}^{m \times (d+1)}$ : input matrix, put all the data vector into a matrix
- $\mathbf{y} \in \mathbb{R}^m$ : target vector, each entry corresponds to the output of one input vector
- $\mathbf{w} \in \mathbb{R}^{d+1}$ : weight vector
- $b \in \mathbb{R}$ : offset number

### 3.1. Task

Train the model(defined by the parameter, here is a vector) to predict the output vector  $\mathbf{y}$  given a new input vector  $\mathbf{x}$ .

### 3.2. Procedure

Here we just use the offset version

1. We simply want to find the  $\mathbf{w}$  that satisfied:  $\mathbf{y} = \mathbf{X}\mathbf{w}$  In this step, we may decide whether the linear system has a solution or not. If not, we may use the least square solution.
2. Define notation  $f_{\mathbf{w},b}(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i + b$ , and define the loss for each example

$$e_i = f_{\mathbf{w},b}(\mathbf{x}_i) - y_i$$

Then the loss function is  $L(\mathbf{w}, b) = \frac{1}{m} \sum_{i=1}^m e_i^2 = \frac{1}{m} \sum_{i=1}^m (f_{\mathbf{w},b}(\mathbf{x}_i) - y_i)^2$

3. The least square solution that minimize the loss function is  $\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ , (without offset, just use the original input matrix  $\mathbf{X}$ ). Be careful whether the matrix  $\mathbf{X}^T \mathbf{X}$  is invertible or not, i.e.,  $\mathbf{X}$  is full column rank.
4. Then we could predict new input  $\mathbf{y}_{\text{new}} = \mathbf{x}_{\text{new}}^T \mathbf{w}$ .

#### Obtaining the LSE Solution

We have loss function:

$$L(\mathbf{w}, b) = \frac{1}{m} \sum_{i=1}^m e_i^2 = \frac{1}{m} \sum_{i=1}^m (f_{\mathbf{w}, b}(x_i) - y_i)^2 = (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y})$$

Then we open the bracket and get:

$$L(\mathbf{w}, b) = \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{w}^T \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \mathbf{w} + \mathbf{y}^T \mathbf{y}$$

The middle two terms are symmetric and since they are scalar, they actually equal to each other, so we have:

$$L(\mathbf{w}, b) = \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{y}^T \mathbf{y}$$

Then we take the derivative with respect to  $\mathbf{w}$  and set it to zero:

$$2\mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{X}^T \mathbf{y} = 0$$

Then we get:

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

If  $\mathbf{X}^T \mathbf{X}$  is invertible, we have:

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

## 4. Linear Regression with Multiple Outputs

We define:

- $m$ : number of training examples
- $d$ : number of features
- $h$ : number of outputs features
- $\mathbf{X} \in \mathbb{R}^{m \times (d+1)}$ : input matrix,  $\mathbf{X} = [\mathbf{1} \ \mathbf{x}']$
- $\mathbf{Y} \in \mathbb{R}^{m \times h}$ : output matrix
- $\mathbf{W} \in \mathbb{R}^{(d+1) \times h}$ : design matrix,  $\mathbf{W} = \begin{bmatrix} \mathbf{b}^T \\ \mathbf{w}' \end{bmatrix}$
- $\mathbf{b} \in \mathbb{R}^h$ : offset vector, we ignore it and put it into the design matrix

### 4.1. Task

Train the model (defined by the parameter, here is a matrix) to predict the output vector  $\mathbf{y}$  given a new input vector  $\mathbf{x}$ .

### 4.2. Procedure

The procedure is similar to the single output case. For the first step, we could divide the output matrix into  $h$  single output vectors, try to solve the linear system for each output vector. If the matrix  $\mathbf{X}$  is full column rank, we may get the least square solution.

1. Define loss function:  $\text{Loss}(\mathbf{W}) = \text{Loss}(\mathbf{W}, \mathbf{b}) = \sum_{k=1}^h (\mathbf{X}\mathbf{w}_k - \mathbf{y}^k)^T (\mathbf{X}\mathbf{w}_k - \mathbf{y}^k)$
2. The least square solution is  $\mathbf{W} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$ , where  $\mathbf{X} = [\mathbf{1} \ \mathbf{x}']$ .

### 4.3. MLE and Linear Regression

We have  $y_i = \mathbf{w}^T \mathbf{x}_i + e_i$  ( $e_i$  is in the vector  $\mathbf{w}$ ). The error term  $e_i$  is of gaussian distribution and therefore  $y_i | \mathbf{x}_i; \mathbf{w}, \sigma^2 \sim N(\mathbf{w}^T \mathbf{x}_i, \sigma^2)$ , where  $\sigma^2$  is the variance of the error term. We could write the likelihood function as:

$$L(\mathbf{w}, \sigma^2 | \{y_i, \mathbf{x}_i\}) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - \mathbf{w}^T \mathbf{x}_i)^2}{2\sigma^2}\right)$$

Solve the MLE problem

1.  $\log(L(\mathbf{w}, \sigma^2) | \{y_i, \mathbf{x}_i\}) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum (y_i - \mathbf{w}^T \mathbf{x}_i)^2$
2. use the derivative to find the maximum likelihood estimator:

$$\frac{\sigma}{\sigma \mathbf{w}} \log(L(\mathbf{w}, \sigma^2) | \{y_i, \mathbf{x}_i\}) = \frac{1}{\sigma^2} \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i) \mathbf{x}_i = 0$$

3. We then have

$$\sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i) \mathbf{x}_i = 0 \Rightarrow \mathbf{X}^T (\mathbf{y} - \mathbf{w}^T \mathbf{x}) = 0$$

4. we could get the least square solution:  $\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ . In fact, the least square solution is the MLE solution. (The deduction process is ignored here)

## 5. Linear Classification

Main idea: To treat binary classification as linear regression in which the output  $y_i$  is binary.  $y_i \in \{-1, +1\}$ .

- Learning and training part: Similar procedure, obtain the weight vector  $\mathbf{w}$ .
- Prediction part:  $y_{\text{new}} = \text{sign}(\mathbf{x}_{\text{new}}^T \mathbf{w}) = \text{sign}\left(\begin{bmatrix} 1 \\ \mathbf{x}'_{\text{new}} \end{bmatrix}^T \mathbf{w}\right) \in \{-1, +1\}$

The sign function is defined as:

$$\text{sign}(x) = \begin{cases} +1 & \text{if } x \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

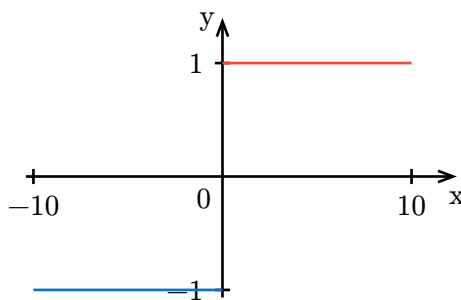


Figure 1: Sign function.

### 5.1. Python Demo

```
1 import numpy as np
2 from numpy.linalg import inv
3
```



```

4  X = np.array([[1, -7], [1, -5], [1, 1], [1, 5]])
5  y = np.array([[ -1], [ -1], [ 1], [ 1]])
6  #linear regression for classification
7  w = inv(X.T @ (X)) @ (X.T @ (y))
8  print(w, "\n")
9
10 #predict
11 X_new = np.array([[1, 2]])
12 y_predict_new = np.sign(X_new @ w)
13 print(y_predict_new)
14 # expected output: [[-1.]]

```

## 5.2. Multi-class Classification

Idea: one-hot encoding, for classes  $\{1, 2, \dots, C\}$ , where  $C > 2$  is the number of classes. The corresponding label vector is

$$\begin{aligned}
 \mathbf{y}_{c1} &= [1 \ 0 \ 0 \ \dots \ 0] \\
 \mathbf{y}_{c2} &= [0 \ 1 \ 0 \ \dots \ 0] \\
 &\vdots \\
 \mathbf{y}_{cC} &= [0 \ 0 \ 0 \ \dots \ 1]
 \end{aligned}$$

We store the class vectors of datasets into a label matrix  $\mathbf{Y} \in \mathbb{R}^{m \times C}$ , where  $m$  is the number of training examples. It is a binary matrix. Essentially we are doing  $C$  separate linear classification problems with class  $k$  and other classes as a single class.

- Learning and training part: Similar procedure, obtain the weight matrix  $\mathbf{W} \in \mathbb{R}^{(d+1) \times C}$ .
- Prediction part:

$$\mathbf{y}_{\text{new}} = \underset{k \in \{1, 2, \dots, C\}}{\operatorname{argmax}} \left( \mathbf{x}_{\text{new}}^T \mathbf{W}[:, k] \right)$$

in which the  $\mathbf{W}[:, k] \in \mathbb{R}^{d+1}$  is the  $k$ -th column of the weight matrix.

### 5.2.1. Python Demo

```

1  import numpy as np
2  from numpy.linalg import inv
3  from sklearn.preprocessing import OneHotEncoder
4  # manually encode the label matrix
5  X = np.array([[1, 1, 1], [1, -1, 1], [1, 1, 3], [1, 1, 0]])
6  Y_class = np.array([[1], [2], [1], [3]])
7  Y = np.array([[1, 0, 0], [0, 1, 0], [1, 0, 0], [0, 0, 1]])
8

```

```

9  # one-hot encoding function
10 Y_onehot = OneHotEncoder().fit_transform(Y_class).toarray()
11
12 # learning
13 W = inv(X.T @ X) @ X.T @ Y # could use Y_onehot instead of Y
14
15 # predicting
16 X_new = np.array([[1,0,-1]])
17 Y_predict = X_new @ W
18 print(np.argmax(Y_predict)+1)

```

## 6. Polynomial Regression

### 6.1. Motivation

- The data may come from a polynomial function, but the linear model may not be able to fit the data well.
- For classification, the  $\oplus$  dataset is not linearly separable, and polynomial function could separate the data.

### 6.2. Idea

- We need some knowledge about polynomial functions.
  1. Let  $f_w(x) = w_0 + w_1x_1 + \dots + w_dx_d$ , then it has  $d$  variable and 1 degree. Let  $f_w(x) = w_0 + w_1x + w_2x_2 + w_3x_1x_2 + w_4x_1^2 + w_5x_2^2$ , then it has 2 variables and 2 degree. each term is called monomial.
  2. For a polynomial function with  $d$  degree and  $n$  variables, the number of monomials is  $\binom{d+n}{n}$

#### Training

- For  $m$  data,  $n$  order and  $d$  variables, the polynomial matrix is

$$P = \begin{bmatrix} p_1^T \\ p_2^T \\ \vdots \\ p_m^T \end{bmatrix} \in \mathbb{R}^{m \times \binom{d+n}{n}}$$

and  $p_i = [1 \ m_1 \ m_2 \ \dots \ m_k]^T$ ,  $m_i$  is monomial and  $k = \binom{d+n}{n}$

- weight vector  $w = P^T(P P^T)^{-1}y$  (Dual form of the least square solution before)

#### Predicting

- $y_{\text{new}} = p_{\text{new}}^T w$

### 6.3. Polynomial Classification

- Similar to the linear Classification
- For simple output,  $y_{\text{new}} = \text{sign}(p_{\text{new}}^T w)$
- For multiple outputs,  $y_{\text{new}} = \underset{k \in \{1,2,\dots,c\}}{\text{argmax}} (p_{\text{new}}^T W[:, k])$

## 7. Ridge Regression

### 7.1. Motivation

In real life, the features  $d$  is large and the number of examples  $m$  is small, so the matrix  $\mathbf{X}^T \mathbf{X} \in \mathbb{R}^{d \times d}$  is hard to invert. So we need to stabilize and robustify the solution.

### 7.2. Idea

Recall the loss function of linear regression and introduce the regularization term (ridge regression version):

$$\begin{aligned} L(\mathbf{w}, b) &= \frac{1}{m} \sum_{i=1}^m (f_{\mathbf{w},b}(\mathbf{x}_i) - y_i)^2 = (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) \\ J(\mathbf{w}, b) &= L(\mathbf{w}, b) + \lambda \|\mathbf{w}\|_2^2 \\ &= \frac{1}{m} \sum_{i=1}^m (f_{\mathbf{w},b}(\mathbf{x}_i) - y_i)^2 + \lambda \|\mathbf{w}\|_2^2 \\ &= (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda \mathbf{w}^T \mathbf{w} \end{aligned}$$

where  $\lambda$  is the regularization parameter. We need to minimize the loss function, that is to find:

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} ((\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda \mathbf{w}^T \mathbf{w})$$

and the result is

$$\mathbf{w}^* = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_{d+1})^{-1} \mathbf{X}^T \mathbf{y}$$

The term  $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}$  is always invertible because it is positive definite. The predicting part is the same as the linear regression.

To prove  $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}$  is invertible, we need to prove it is positive definite. By definition of positive definite, we need to prove:  $\forall \mathbf{x}, \mathbf{x}^T \mathbf{A} \mathbf{x} \leftrightarrow \mathbf{A}$  is positive definite. We have

$$\begin{aligned} \mathbf{x}^T (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}) \mathbf{x} &= \mathbf{x}^T \mathbf{X}^T \mathbf{X} \mathbf{x} + \mathbf{x}^T \lambda \mathbf{I} \mathbf{x} \\ &= (\mathbf{X} \mathbf{x})^T (\mathbf{X} \mathbf{x}) + \mathbf{x}^T \mathbf{x} \end{aligned}$$

Since every term is positive, then the result is larger or equal to zero.  $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}$  is positive definite and therefore always invertible.

However, this term is in  $\mathbb{R}^{(d+1) \times (d+1)}$  so it is hard to find the inverse of it. So we could use the dual form of the solution which needs to find the inverse of a matrix in  $\mathbb{R}^{m \times m}$ .

$$\mathbf{w} = \mathbf{X}^T (\mathbf{X} \mathbf{X}^T + \lambda \mathbf{I}_m)^{-1} \mathbf{y}$$

The proof of the dual form need to use Woodbury formula:

$$(\mathbf{I} + \mathbf{U} \mathbf{V})^{-1} = \mathbf{I} - \mathbf{U} (\mathbf{I} + \mathbf{U} \mathbf{V})^{-1} \mathbf{V}$$

**Primal-Dual Equivalence Proof:**

Starting with dual form solution:

$$\mathbf{w}^* = \mathbf{X}^T (\mathbf{X} \mathbf{X}^T + \lambda \mathbf{I})^{-1} \mathbf{y}$$

Develop equivalence through Woodbury identity  $(\mathbf{I} + \mathbf{U} \mathbf{V})^{-1} = \mathbf{I} - \mathbf{U}(\mathbf{I} + \mathbf{U} \mathbf{V})^{-1} \mathbf{V}$ :

$$\begin{aligned} & \mathbf{X}^T (\mathbf{X} \mathbf{X}^T + \lambda \mathbf{I})^{-1} \mathbf{y} \\ &= \lambda^{-1} \mathbf{X}^T (\mathbf{I} + \lambda^{-1} \mathbf{X} \mathbf{X}^T)^{-1} \mathbf{y} \\ &= \lambda^{-1} \mathbf{X}^T \left[ \mathbf{I} - \lambda^{-1} \mathbf{X} (\mathbf{I} + \lambda^{-1} \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \right] \mathbf{y} \\ &= \lambda^{-1} \left[ \mathbf{X}^T \mathbf{y} - \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \right] \\ &= \lambda^{-1} \left[ \mathbf{I} - \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \right] \mathbf{X}^T \mathbf{y} \\ &= \lambda^{-1} \left[ \mathbf{I} - (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}) (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} + \lambda \mathbf{I} (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \right] \mathbf{X}^T \mathbf{y} \\ &= (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \end{aligned}$$

Here  $\mathbf{U} \equiv \lambda^{-1} \mathbf{X}$ ,  $\mathbf{V} \equiv \mathbf{X}^T$  Final form matches primal solution. QED.

### 7.3. Python Demo

Refer to [Python-Demo-4-Algorithm](#)

## 8. Gradient Descent

### 8.1. Motivation

In linear regression and the other few models mentioned above, the optimal solution could be simply found by solving the equation. However, in many cases, to minimize the loss function with respect to parameter  $\mathbf{w}$  is hard. So we want an algorithm that could iteratively find the optimal solution, that is the gradient descent algorithm.

### 8.2. Idea

- Task: Minimize the loss function  $C(\mathbf{w})$  in which  $\mathbf{w} = [w_1 \ \dots \ w_d]^T$
- Gradient of  $\mathbf{w}$ :  $\nabla_{\mathbf{w}} C(\mathbf{w}) = \left[ \frac{\partial C}{\partial w_1} \ \frac{\partial C}{\partial w_2} \ \dots \ \frac{\partial C}{\partial w_d} \right]^T$ . It's a function or say a vector of  $\mathbf{w}$ . And the direction of the gradient is the direction of the fastest **increase** of the function, while the opposite direction is the direction of the fastest **decrease** of the function.
- Algorithm
  1. Initial  $\mathbf{w}$  with learning rate  $\eta > 0$
  2.  $\mathbf{w}_{k+1} = \mathbf{w}_k - \eta \nabla_{\mathbf{w}} C(\mathbf{w}_k)$
  3. Repeat the above step until the convergence condition is satisfied.
- Convergence Criterias:
  1. The absolute or percentage change of the loss function is smaller than a threshold.
  2. The absolute or percentage change of parameter  $\mathbf{w}$  is smaller than a threshold.
  3. Set a maximum number of iterations.



- Notice that according to multivariate calculus, if  $\eta$  is not too large,  $C(\mathbf{w}_{k+1}) < C(\mathbf{w}_k)$ . But GD could only find the local minimum.

### 8.3. Variation of GD

#### 8.3.1. Change the learning rate

1. Decreasing learning rate

$$\eta = \frac{\eta_0}{1+k}$$

where  $k$  is the number of iterations or other form of changes. It could help the algorithm to converge faster at the beginning and avoid oscillation at the end.

2. Adaptive learning rate

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \frac{\eta}{\sqrt{G_k} + \varepsilon} \nabla_{\mathbf{w}} C(\mathbf{w}_k)$$

where  $G_k = \sum_{i=0}^k \|\nabla_{\mathbf{w}} C(\mathbf{w}_i)\|_2^2$  and  $\varepsilon$  is a small positive number to avoid zero denominator. This method gives larger update to smaller gradient and vice versa, adjusting the learning rate according to the history information. But the learning rate may shrink too fast.

#### 8.3.2. Different gradient

1. Momentum-based GD

$$\begin{aligned} \mathbf{v}_k &= \beta \mathbf{v}_{k-1} + (1 - \beta) \nabla_{\mathbf{w}} C(\mathbf{w}_{k-1}) \\ \mathbf{w}_{k+1} &= \mathbf{w}_k - \eta \mathbf{v}_k \end{aligned}$$

where  $\beta \in (0, 1)$  is the momentum parameter and  $\mathbf{v}_0 = \nabla_{\mathbf{w}} C(\mathbf{w}_0)$ . It converges fast but may overshoot the optimal point.

2. Nesterov Accelerated Gradient (NAG)

$$\begin{aligned} \mathbf{v}_k &= \beta \mathbf{v}_{k-1} + \eta \nabla_{\mathbf{w}} C(\mathbf{w} + \beta \mathbf{w}_{k-1}) \\ \mathbf{w}_{k+1} &= \mathbf{w}_k - \mathbf{v}_k \end{aligned}$$

where  $\mathbf{v}_0 = 0$ . It works by anticipating the next direction of the optimizer. It is fast but complex.

#### 8.3.3. Design of loss function

The idea is only calculate the loss of some sample from the dataset, reduce the computation cost.

1. Batch GD: Use all the data to calculate the gradient.
2. Stochastic GD: Use one randomly chosen data to calculate the gradient.
3. Mini-batch GD: Use a small batch of randomly chosen data to calculate the gradient.

## 9. Logistic Regression

### 9.1. Motivation

There are some possible issues for classification problems:

1. Noises: Lead to unseparable data.
2. Mediocre generalization: Only find barely boundary.

3. Overfitting: The model is too complex and fits the noise in the data.

In that case, we want a model that output our confidence of the prediction and that's why we need logistic regression.

## 9.2. Idea

- Logistic function:

$$g(z) = \frac{1}{1 + e^{-z}}$$

- Logistic Regression: In binary classification task, we set our prediction

$$\Pr(y = 1|\mathbf{x}, \boldsymbol{\theta}, \theta_0) = g(\langle \boldsymbol{\theta}, \mathbf{x} \rangle + \theta_0) = \frac{1}{1 + e^{-(\langle \boldsymbol{\theta}, \mathbf{x} \rangle + \theta_0)}}$$

This is known as class conditional probability.

How could we derive the form of the class conditional probability? We first define the log-odds function of both classes as a affine function of the input vector  $\mathbf{x}$ :

$$\log \frac{\Pr(y = 1|\mathbf{x}, \boldsymbol{\theta}, \theta_0)}{\Pr(y = -1|\mathbf{x}, \boldsymbol{\theta}, \theta_0)} = \boldsymbol{\theta} \mathbf{x} + \theta_0$$

Set  $\Pr(y = -1|\mathbf{x}, \boldsymbol{\theta}, \theta_0) = 1 - \Pr(y = 1|\mathbf{x}, \boldsymbol{\theta}, \theta_0)$ , then:

$$\log \frac{\Pr(y = 1|\mathbf{x}, \boldsymbol{\theta}, \theta_0)}{1 - \Pr(y = 1|\mathbf{x}, \boldsymbol{\theta}, \theta_0)}$$

Let  $\Pr(y = 1|\mathbf{x}, \boldsymbol{\theta}, \theta_0) = a$ ,  $\langle \boldsymbol{\theta}, \mathbf{x} \rangle + \theta_0 = b$ , then it is equivalent to:

$$\log \frac{a}{1-a} = b \Rightarrow \frac{a}{1-a} = e^b \Rightarrow a = (1-a)e^b \Rightarrow a = \frac{e^b}{1+e^b} \Rightarrow a = \frac{1}{1+e^{-b}}$$

$$\text{So } \Pr(y = 1|\mathbf{x}, \boldsymbol{\theta}, \theta_0) = \frac{1}{1+e^{-(\langle \boldsymbol{\theta}, \mathbf{x} \rangle + \theta_0)}}$$

- $\langle \boldsymbol{\theta}, \mathbf{x} \rangle + \theta_0 = 0$  is the decision boundary. If  $\langle \boldsymbol{\theta}, \mathbf{x} \rangle + \theta_0 > 0$ , then  $\Pr(y = 1|\mathbf{x}, \boldsymbol{\theta}, \theta_0) > 0.5$  and vice versa.

How to derive this: If a data lies on  $\langle \boldsymbol{\theta}, \mathbf{x} \rangle + \theta_0 = 0$ , then we are not sure about the class of the data. So we simply put  $\Pr(y = 1|\mathbf{x}, \boldsymbol{\theta}, \theta_0) = \Pr(y = -1|\mathbf{x}, \boldsymbol{\theta}, \theta_0) = \frac{1}{2}$  here.

**For how do we measure the performance, we could use likelihood function:**

- For a single input  $(\mathbf{x}_t, y_t)$ , the likelihood function is

$$L(\boldsymbol{\theta}, \theta_0 | (\mathbf{x}_t, y_t)) := \Pr(y_t | \mathbf{x}_t, \boldsymbol{\theta}, \theta_0) = g(y_t(\langle \mathbf{x}_t, \boldsymbol{\theta} \rangle + \theta_0))$$

- For the whole dataset, the likelihood function is

$$L(\boldsymbol{\theta}, \theta_0 | \mathbf{X}, \mathbf{y}) = \prod_{t=1}^m g(y_t(\langle \mathbf{x}_t, \boldsymbol{\theta} \rangle + \theta_0))$$

Through maximizing this likelihood function, we could get the optimal parameter  $\boldsymbol{\theta}$  and  $\theta_0$ . Our strategy is stochastic gradient descent.

**Procedure:**

1. Goal

$$\begin{aligned} & \max \prod_{t=1}^m g(y_t(\langle \mathbf{x}_t, \boldsymbol{\theta} \rangle + \theta_0)) \\ &= \max \sum_{t=1}^m \log g(y_t(\langle \mathbf{x}_t, \boldsymbol{\theta} \rangle + \theta_0)) \end{aligned}$$

2. We transform it into the minimization problem for better generality: [Minimization and Maximization in ML](#):

$$\begin{aligned} & \min - \sum_{t=1}^m \log g(y_t(\langle \mathbf{x}_t, \boldsymbol{\theta} \rangle + \theta_0)) \\ &= \min \sum_{t=1}^m \log(1 + e^{-y_t(\langle \mathbf{x}_t, \boldsymbol{\theta} \rangle + \theta_0)}) \end{aligned}$$

3. Take the gradients w.r.t.  $\theta_0, \boldsymbol{\theta}$  of the function:

$$\begin{aligned} \frac{d}{d\theta_0} \log(1 + e^{-y_t(\langle \mathbf{x}_t, \boldsymbol{\theta} \rangle + \theta_0)}) &= -y_t \frac{e^{-y_t(\langle \mathbf{x}_t, \boldsymbol{\theta} \rangle + \theta_0)}}{1 + e^{-y_t(\langle \mathbf{x}_t, \boldsymbol{\theta} \rangle + \theta_0)}} = -y_t [1 - \Pr(y_t | \mathbf{x}_t, \boldsymbol{\theta}, \theta_0)] \\ \frac{d}{d\boldsymbol{\theta}} \log(1 + e^{-y_t(\langle \mathbf{x}_t, \boldsymbol{\theta} \rangle + \theta_0)}) &= -y_t [1 - \Pr(y_t | \mathbf{x}_t, \boldsymbol{\theta}, \theta_0)] \mathbf{x}_t \end{aligned}$$

4. Set the gradient to zero:  $\sum_{t=1}^m (-y_t [1 - \Pr(y_t | \mathbf{x}_t, \boldsymbol{\theta}, \theta_0)]) = 0$  and  $\sum_{t=1}^m (-y_t [1 - \Pr(y_t | \mathbf{x}_t, \boldsymbol{\theta}, \theta_0)]) \mathbf{x}_t = \mathbf{0}$

Note:

1. SGD has no significant change when the gradient of the full objective is zero.
2. We could add the regularization term into the function we concern:

$$\min \sum_{t=1}^m \log(1 + e^{-y_t(\langle \mathbf{x}_t, \boldsymbol{\theta} \rangle + \theta_0)}) + \frac{\lambda}{2} \|\boldsymbol{\theta}\|^2$$

### 9.3. Multi-Class Logistic Regression

For multi-class classification, we just replace

$$\Pr(y = 1 | \boldsymbol{\theta}, \theta_0, \mathbf{x}) = \frac{1}{1 + e^{-(\langle \mathbf{x}, \boldsymbol{\theta} \rangle + \theta_0)}}$$

with softmax function:

$$\Pr(y = k | \boldsymbol{\theta}, \theta_0, \mathbf{x}) = \frac{e^{\langle \mathbf{x}, \boldsymbol{\theta}_k \rangle + \theta_{0,k}}}{\sum_{j=1}^C e^{\langle \mathbf{x}, \boldsymbol{\theta}_j \rangle + \theta_{0,j}}}$$

where  $C$  is the number of classes.

**Claim:**

Two-class logistic regression is a special case of multi-class logistic regression.

**Proof:**

Without loss of generality, we assume  $\theta_1 = \mathbf{0}$ ,  $\theta_{0,1} = 0$ , using the equation defined for multi-class logistic regression, we have:

$$\begin{aligned}\Pr(y = 1 | \mathbf{x}, \theta_1, \theta_{0,1}) &= \frac{e^{\langle \mathbf{x}, \theta_1 \rangle + \theta_{0,1}}}{e^{\langle \mathbf{x}, \theta_1 \rangle + \theta_{0,1}} + e^{\langle \mathbf{x}, \theta_2 \rangle + \theta_{0,2}}} \\ \Pr(y = 2 | \mathbf{x}, \theta_2, \theta_{0,2}) &= \frac{e^{\langle \mathbf{x}, \theta_2 \rangle + \theta_{0,2}}}{e^{\langle \mathbf{x}, \theta_1 \rangle + \theta_{0,1}} + e^{\langle \mathbf{x}, \theta_2 \rangle + \theta_{0,2}}} \\ &= \frac{e^{\langle \mathbf{x}, \theta_2 \rangle + \theta_{0,2}}}{1 + e^{\langle \mathbf{x}, \theta_2 \rangle + \theta_{0,2}}} \\ &= \frac{1}{1 + e^{-(\langle \mathbf{x}, \theta_2 \rangle + \theta_{0,2})}}\end{aligned}$$

The equation below appears in the same form of two-class logistic regression. QED.

## 9.4. Python Demo

Refer to [Python-Demo-4-Algorithm](#)

## 10. Support Vector Machine

### 10.1. Get familiar with Logistic Regression

Logistic regression is a binary classification model.

$$f_{\mathbf{w}}(\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})}$$

The loss function in logistic regression is

$$L(\mathbf{w}) = -\sum_{i=1}^m y_i \log(f_{\mathbf{w}}(\mathbf{x}_i)) + (1 - y_i) \log(1 - f_{\mathbf{w}}(\mathbf{x}_i)) + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$

Here our class is  $\{0, 1\}$ . The loss function shows that if the label is 1, then the loss is  $-\log(f_{\mathbf{w}}(\mathbf{x}_i))$ , meaning that  $f_{\mathbf{w}}(\mathbf{x}_i)$  or say  $-\mathbf{w}^T \mathbf{x}$  should be as large as possible and vice versa. Note that the  $\lambda$  could be understood as how much we want to penalize the large  $\mathbf{w}$ , if we care more about the loss of each example, we just set  $\lambda$  small.

### 10.2. Task

Given a dataset  $D = \{(\mathbf{x}_t, y_t) \in \mathbb{R}^d \times \{-1, +1\} : t = 1, 2, \dots, n\}$ , our goal is to learn from the data and predict the class of a new input  $\mathbf{x}_{\text{new}}$ .

### 10.3. Idea

#### 10.3.1. Preliminary

- Margin of classifier: We have sign function

$$\text{sign}(\boldsymbol{\theta}^T \mathbf{x}) = \begin{cases} +1 & \text{if } \boldsymbol{\theta}^T \mathbf{x} \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

For a correct classification, we have  $\text{sign}(\boldsymbol{\theta}^T \mathbf{x}) = y_t$ . So we have  $y_t \boldsymbol{\theta}^T \mathbf{x} > 0$ . We define the margin of the classifier as:

$$y_t \boldsymbol{\theta}^T \mathbf{x}$$

If it's positive, then the sample is correctly classified and vice versa. What's more, the larger the margin is, the better the classifier is (we could infer this from the graph of sign function).

- Linearly separable: A dataset  $\mathbf{D} = \{(\mathbf{x}_t, y_t) \in \mathbb{R}^d \times \{-1, +1\} : t = 1, 2, \dots, n\}$  is linearly separable if

$$\exists \boldsymbol{\theta}, \forall t, y_t \boldsymbol{\theta}^T \mathbf{x} > 0$$

- Minimum margin:

$$\gamma = \min_t y_t \boldsymbol{\theta}^T \mathbf{x} > 0$$

- $\gamma$ -linearly separable: A dataset  $\mathbf{D}$  is  $\gamma$ -linearly separable if for some  $\gamma$ ,

$$\exists \boldsymbol{\theta}, \forall t, y_t \boldsymbol{\theta}^T \mathbf{x} \geq \gamma$$

- Geometric margin: the smallest distance over all samples to the decision boundary

$$\gamma_{\text{geom}} = \frac{\gamma}{\|\boldsymbol{\theta}\|} = \frac{\min_t y_t \boldsymbol{\theta}^T \mathbf{x}}{\|\boldsymbol{\theta}\|}$$

The inverse of the margin  $\gamma_{\text{geom}}^{-1}$  could represent the difficulty of the classification problem.

### 10.3.2. Optimization

For a  $\gamma$ -linearly separable dataset, we have the following optimization problem:

$$\max_{\boldsymbol{\theta}' \in \mathbb{R}^d} \frac{\gamma}{\|\boldsymbol{\theta}'\|} \text{ s.t. } \forall t, y_t \mathbf{x}^T \boldsymbol{\theta}' \geq \gamma \Leftrightarrow \min_{\boldsymbol{\theta}' \in \mathbb{R}^d} \frac{\|\boldsymbol{\theta}'\|}{\gamma} \text{ s.t. } \forall t, y_t \mathbf{x}^T \boldsymbol{\theta}' \geq \gamma$$

We could transform the term we optimize to be  $\left\| \frac{\boldsymbol{\theta}'}{\gamma} \right\|$ , which is the inverse of the geometric margin.

$$\min_{\boldsymbol{\theta}' \in \mathbb{R}^d} \left\| \frac{\boldsymbol{\theta}'}{\gamma} \right\| \text{ s.t. } \forall t, y_t \mathbf{x}^T \frac{\boldsymbol{\theta}'}{\gamma} \geq 1$$

Since scaling the  $\boldsymbol{\theta}'$  by constant does not change the decision boundary, the goal here simply put:

$$\min_{\boldsymbol{\theta}' \in \mathbb{R}^d} \frac{1}{2} \|\boldsymbol{\theta}'\|^2 \text{ s.t. } \forall t, y_t \mathbf{x}^T \boldsymbol{\theta}' \geq 1$$

It is so called the primal SVM problem.

### 10.3.3. Proof of Uniqueness

**Claim:**

The primal SVM problem has a unique solution.

**Proof:**

1. Suppose we have two distinct solutions  $\theta_1, \theta_2$
2. Then  $\|\theta_1\| = \|\theta_2\|$
3. Let  $\bar{\theta} = \frac{\theta_1 + \theta_2}{2}$
4. By the linearity of inner product, we have

$$y_t x^T \bar{\theta} = y_t x^T \left( \frac{\theta_1 + \theta_2}{2} \right) = \frac{1}{2} y_t x^T \theta_1 + \frac{1}{2} y_t x^T \theta_2 \geq 1$$

5. By the triangle inequality, we have

$$\|\bar{\theta}\| = \left\| \frac{1}{2}(\theta_1 + \theta_2) \right\| \leq \frac{1}{2}\|\theta_1\| + \frac{1}{2}\|\theta_2\| = \|\theta_1\| = \|\theta_2\|$$

If  $<$  hold, that's a contradiction to the assumption that  $\theta_1, \theta_2$  are solutions. If  $=$  hold, then  $\theta_1, \theta_2$  are the same solution.

6. Then the primal SVM problem has a unique solution.

### 10.3.4. SVM with offset

Sometimes the datas do not “centered around zero”, so we could add a constant term to the decision boundary:

$$\theta^T x + \theta_0 = 0$$

Then the optimization goal becomes:

$$\min_{(\theta, \theta_0) \in \mathbb{R}^d \times \mathbb{R}} \frac{1}{2} \|\theta\|^2 \quad s.t. \quad \forall t, y_t (\theta^T x + \theta_0) \geq 1$$

Remarks:  $\theta_0$  only appears in the constraint. In that case,  $\theta_0$  is not in the vector that we optimize and it has no constraint to the “place” of the decision boundary, compared to  $\tilde{x} = [x, 1]$

### 10.3.5. Evaluation

The correctness of the classification model could be:

$$\frac{\text{correct predictions}}{\text{all predictions}} = \frac{1}{n} \sum_1^n I\{y_i, f(x, (\theta, \theta_0))\}$$

$I$  returns 1 if the parameters are the same and 0 otherwise.

#### 10.3.5.1. Leave-one-out Cross Validation

To access the robustness of the model:

$$\text{LOOCV} = \frac{1}{n} \sum_1^n I\{y_t, f(x, (\theta^{-t}, \theta_0^{-t}))\}$$

The parameters  $\theta^{-t}, \theta_0^{-t}$  are obtained training the model with the dataset that removes the  $t$ -th sample. We have the proposition:

**Proposition:**

$$\text{LOOCV} \leq \frac{N}{n}$$

where  $N$  is the number of support vectors and  $n$  is the size of dataset.

**Proof:**

1. for all datapoint that is not a support vector, we have  $y_t f(\mathbf{x}, (\theta^{-t}, \theta_0^{-t})) = 0$ , so the term is not counted.
2. for all datapoint that is a support vector, we have  $y_t f(\mathbf{x}, (\theta^{-t}, \theta_0^{-t})) \leq 1$ , so the term is counted.
3. Then  $\text{LOOCV} \leq \frac{N}{n}$  by sum over all datapoints.

If the support vectors are few, then the loss is small and the model generalizes well.

### 10.3.5.2. Allowing misclassification

The primal form of SVM with offset and slack:

$$\min_{\theta, \theta_0, \varepsilon \in \mathbb{R}^d \times \mathbb{R} \times \mathbb{R}_+^n} \frac{1}{2} \|\theta\|^2 + C \sum_{i=1}^n \varepsilon_i \quad \text{s.t.} \quad \forall t, y_t (\mathbf{x}^T \theta + \theta_0) \geq 1 - \varepsilon_t$$

We could set the slack variable to be large if we want to allow more misclassification for specific datapoint and here will still be penalty since we have sum over the slack variable. For the constant,

$$C \rightarrow \infty \Leftrightarrow \text{No misclassification}, C \rightarrow 0 \Leftrightarrow \text{Allow misclassification}$$

## 11. Feature Engineering

Machine learning algorithms prefer numerical data and fixed length feature vectors. Before, in computer vision, people need to manually extract the features from the image then put the features to linear model like SVM or other models for training. In that case, the feature extraction is the most important part since in the model training part, all we could do is to tune the parameters. Later in deep learning, the feature extraction part is done by neural network, which is more efficient and flexible.

### 11.1. Tabular Data

- Int/float: Could be used directly.
- Categorical: One-hot encoding. When there is long tail, we could put the rare category into a new category-unknown.
- Date/time: A feature list such as: [year, month, day, day\_of\_year, day\_of\_month, day\_of\_week].
- Feature Combination: Put the cartesian product of the features into the feature vector, we could obtain the relevance between the features.

### 11.2. Text Data

- Represent as token features

1. Bag of words: Count the number of times each word appears in the document.
  2. Word-Embedding: Vectorizing the words such that similar words's distance is small (could be represented by cosine similarity).
- Pre-trained LLM: Use the pre-trained LLM to extract the features from the text.

### 11.3. Image Data

- Manually: SIFT, color feature, edge feature, etc.
- Pre-trained deep neural network: Put the image into the pre-trained deep neural network and extract the features from the second last layer (last layer is the classification layer).

## 12. Training Set & Model Complexity

This section is the summary of lecture 2 of Li Hongyi's course: [What to do if my network fails to train](#)

### 12.1. Context

For a classification problem, we could define the prediction function as

$$f_{h(x)} = \begin{cases} 1 & \text{if } e(x) \geq h \\ 0 & \text{otherwise} \end{cases}$$

Here, the  $e(x)$  is a function that return specific features of the input  $x$ . The  $h$  is a threshold that we set. So here the parameter of the model is  $h$ .

- Model Complexity: We could measure the complexity of the model by the size of the parameter set, that is  $|\mathcal{H}|$ .
- Loss function (Given dataset  $D$  and parameter  $h$ ):

$$L(h) = \sum_{i=1}^m I\{y_i, f_{h(x_i)}\}$$

So then we could define the global optimal parameter  $h^{\text{all}} = \operatorname{argmin} L(h, D^{\text{all}})$  and the local optimal parameter  $h^{\text{train}} = \operatorname{argmin} L(h, D^{\text{train}})$ .

### 12.2. Goal

**Goal:** We want the difference between  $L(h^{\text{train}}, D^{\text{all}}), L(h^{\text{all}}, D^{\text{all}})$  to be small, which means the model generalizes well. We could put a small constant to bound it:

$$L(h^{\text{train}}, D^{\text{all}}) - L(h^{\text{all}}, D^{\text{all}}) \leq \delta$$

**Claim:**  $\forall h \in \mathcal{H}, |L(h, D^{\text{all}}) - L(h, D^{\text{train}})| \leq \frac{\delta}{2} \Rightarrow L(h^{\text{train}}, D^{\text{all}}) - L(h^{\text{all}}, D^{\text{all}}) \leq \delta$

In the claim above, we call a training dataset  $D^{\text{set}}$  a good dataset if it satisfies the condition

$$\forall h \in \mathcal{H}, |L(h, D^{\text{all}}) - L(h, D^{\text{train}})| \leq \frac{\delta}{2}$$

We want to know the probability of we using a bad dataset.

1. For a bad dataset,  $\exists h \in \mathcal{H}, L(h^{\text{train}}, D^{\text{all}}) - L(h^{\text{all}}, D^{\text{all}}) > \frac{\delta}{2} = \varepsilon$
2. So for different  $h$ , there maybe some datasets that is bad because the difference between  $L(h^{\text{train}}, D^{\text{all}}), L(h^{\text{all}}, D^{\text{all}})$  is larger than  $\varepsilon$ .



3. Therefore,

$$P(\text{Bad Training Set}) = \bigcup P(\text{Bad Training Set}|h) \leq \sum_{h \in \mathcal{H}} P(\text{Bad Training Set}|h)$$

4. We introduce Hoeffding's inequality:

$$P(\text{Bad Training Set} \mid h) \leq 2 \exp(-2N\varepsilon^2)$$

where  $N$  is the size of dataset. Then we have

$$P(\text{Bad Training Set}) \leq 2|\mathcal{H}| \exp(-2N\varepsilon^2) \quad (1)$$

By the inequality (1), we could deduce, firstly, what is the probability of a bad dataset? Secondly, what size of dataset is required to achieve a certain probability of good dataset?

### 12.3. Trade-off

From (1), we could obtain the size of “credible” dataset  $N \geq \frac{\log\left(\frac{2|\mathcal{H}|}{\delta}\right)}{2\varepsilon^2}$ , that means, if we want a small  $N$ , then we need small  $|\mathcal{H}|$  or large  $\delta$ .

- Small  $|\mathcal{H}|$ : The model is simple, so the difference between  $L(h^{\text{train}}, D^{\text{all}})$ ,  $L(h^{\text{all}}, D^{\text{all}})$  is small. However, the model could be bad because the model could not fit the reality well. Though the difference is small, it is of no use.
- Large  $\delta$ : Large  $\delta$  means we are willing to accept larger difference between our perception and reality. However, if the difference is large, the perception (model) is bad.

The solution for the trade-off is deep learning.

## 13. Model Selection

### 13.1. Generalization of Supervised Learning Algorithms

A supervised learning algorithm contains the following key elements:

- Unknown target function  $f : \mathcal{X} \mapsto \mathcal{Y}$
- Training examples  $\mathcal{D} = \{(x_i, y_i); i \in \{1, 2, \dots, m\}\}$
- Hypothesis set  $\mathcal{H}$
- Learning algorithm  $\mathcal{A}$
- Final hypothesis  $\mathcal{H}$

Usually, we do not know the performance of the model on new data and we need to evaluate it. For infinite data space, we just evaluate the model on new data. However, in practice, we only have a finite dataset  $\mathcal{D}$ . So we need to split the dataset into training set and test set.

- Data Partition:
  1. Training set, used to fit the model
  2. Validation set, used to fit the hyperparameters
  3. Test set, used to evaluate the model

### 13.2. Error Metric & Regularizer

#### 13.2.1. Error Metric

Generally speaking, an error metric is a function that add penalty to the model if it classify (predict) the data wrong, denoted by:

$$E : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$$

Here we have  $y$  as the ground truth and  $y'$  as the prediction.

- Regression
  - Square error:

$$e(y, y') = (y - y')^2$$

- Absolute error:

$$e(y, y') = |y - y'|$$

- Classification
  - Zero-one error:

$$e(y, y') = \mathbf{I}\{y \neq y'\}$$

- Weighted Misclassification error: Used when we want to give different penalty to different types of errors, such as false positive is more serious than false negative

$$e_{\beta(y, y')} = (1 - \beta)\mathbf{I}\{y = 1, y' = 0\} + \beta\mathbf{I}\{y = 0, y' = 1\}, \beta > .5$$

- Balanced Error Rate:

$$e(y, y') = \frac{1}{2} \left( \frac{\mathbf{I}\{y = 1, y' = 0\}}{n\{y = 1\}} + \frac{\mathbf{I}\{y = 0, y' = 1\}}{n\{y = 0\}} \right)$$

- Confusion Matrix: A table that shows the number of true positives, false positives, true negatives, and false negatives.
  1. True Positive:  $\mathbf{I}\{y = 1, y' = 1\}$ , False Positive:  $\mathbf{I}\{y = 0, y' = 1\}$ , True Negative:  $\mathbf{I}\{y = 0, y' = 0\}$ , False Negative:  $\mathbf{I}\{y = 1, y' = 0\}$
  2. Precision:  $\frac{TP}{TP+FP}$
  3. Recall:  $\frac{TP}{TP+FN}$
  4. F1-Score:  $\frac{2}{\text{Precision}^{-1} + \text{Recall}^{-1}}$ , useful when exists the trade off between precision and recall.
  5. Accuracy:  $TP + \frac{TN}{TP+TN+FP+FN}$
  6. AUC, ROC: ROC is a plot of the true positive rate against the false positive rate. AUC is the area under the ROC curve. When AUC is 1, the model is perfect. When AUC is 0.5, the model is random. When AUC is larger than 0.5, the model is better than random.

This metric is not sensitive to **imbalanced dataset**, for example, some diseases are not common so the number of patients is small and comprise a very small part of the dataset, so if the model predict negative for all datapoint, it could perform well when using traditional accuracy metrics, which is not the case. Besides, AUC could reflect the overall performance of the model.

### 13.2.2. Regularizer

In ridge regression, we use  $l_2$  norm regularization to penalize the large  $w$ . From this example, we could get the insight: Objective function = Error + Regularizer.

- Why regularization?

1. To avoid overfitting
  2. Impose prior structural knowledge
  3. Reduce the variance of the model
  4. Improve interpretability: simpler model, automatic feature selection
- Why not?
    1. Gauss-Markov Theorem: Least square estimation is best linear unbiased estimator (BLUE) in some cases.
    2. Introduce regularization would also introduce bias.
  - expected prediction error =  $\text{bias}^2 + \text{variance} + \text{noise}$
  - No free lunch: If the bias is low, we fit the data well but the variance is high, which means the model is overfitting. If the bias is high, we fit the data poorly but the variance is low, which means the model is underfitting.
  - $l_1$  regularization (used in Lasso Regression):

$$L(h) = \text{Mean Square Loss} + \lambda \sum_{j=1}^n |w_j|$$

It penalize less than  $l_2$  regularization. It is used when we want to perform feature selection. However, absolute function is not differentiable, so there is no closed form solution.

### 13.3. Generalizing & Overfitting

	Low test set error	High test set error
Low training set error	Generalize	Overfit
High training set error	Error	Underfit

- Fixing Underfit: Increase the model complexity (more features, more complex model), more data.
- Fixing Overfit: Reduce the model complexity (reduce features, less complex model), more data.

Generally, it is not possible to overfit and underfit at the same time.

### 13.4. Data Partition: Train/Validation/Test

#### 13.4.1. Validation Set

##### Procedure

---

1. Split the dataset into training set and validation set:  $\mathcal{D} = \mathcal{D}_{\text{train}} + \mathcal{D}_{\text{val}}$
2. Pick  $m$  different models  $\phi : \phi_1, \phi_2, \phi_3$
3. Train the model on the training set and get hypothesis for each  $\phi_i$ :

$$\mathcal{H} = \{h_i = \text{argmin } L(h, \mathcal{D}_{\text{train}}) \mid i \in [1, m]\}$$

Compute error of each hypothesis on validation set  $\mathcal{D}_{\text{val}}$  and choose the lowest:

$$\hat{h}^* = \operatorname{argmin}_{h \in \mathcal{H}} L(h, D_{\text{val}})$$

4. Retrain the model on the whole data set to obtain useful information)

### 13.4.2. Cross-Validation

#### Procedure

1. Split the dataset into  $k$  folds:  $\mathcal{D} = \mathcal{D}_1 + \mathcal{D}_2 + \dots + \mathcal{D}_k$ , we now have  $k$  pairs of training set and validation set:  $\{(\mathcal{D}_1, \mathcal{D}'_1), (\mathcal{D}_2, \mathcal{D}'_2), \dots, (\mathcal{D}_k, \mathcal{D}'_k)\}$
2. Pick  $m$  different models  $\phi : \phi_1, \phi_2 \dots$
3. For each  $\phi_i, k$ 
  1. Train the model on the training set and get hypothesis  $h_{\phi_i, \mathcal{D}_k}$
  2. Compute error of the hypothesis on validation set  $\mathcal{D}'_k$
4. Choose the best model  $\phi^* = \operatorname{argmin}_{\phi \in \Phi} \frac{1}{k} \sum_{j=1}^k E_{\mathcal{D}'_j}(h_{\phi_i, \mathcal{D}_j})$
5. Retrain the model on the whole data set to obtain useful information

The technique for folds is called **k-fold cross-validation**. Besides, we could also use **leave-one-out cross-validation** to split the dataset.

### 13.4.3. LI Hungyi: Why validation set $\rightarrow$ still overfitting?

The context here: We have 3 different models  $\mathcal{H}_1, \mathcal{H}_2, \mathcal{H}_3$  and obtain their optimal parameters  $h_1^*, h_2^*, h_3^*$  on the training set. Then we put the optimal parameters  $h_1^*, h_2^*, h_3^*$  to the validation set to evaluate the performance of the model and get the best model

$$h^* = \operatorname{argmin}_{h \in \{h_1^*, h_2^*, h_3^*\}} L(h, D_{\text{val}})$$

Then we put the best model  $h^*$  to the test set to evaluate the performance of the model. However, sometimes the model would still overfit the data and perform badly on the test set.

The process of model selection is actually another training process. So there is probability that the validation set is not a good dataset. Recall that:

$$P(D_{\text{val}} \text{ is bad}) \leq 2|\mathcal{H}| \exp(-2N\epsilon^2)$$

where  $N$  is the size of dataset. If the number of parameters is too large, then the upper bound is large. So the probability of the validation set is bad could large, leading to a bad generalized model.

## 14. Ensemble Learning

### 14.1. Context

For a decision tree model, usually the prediction performance is bad and the model does not generalize well. But a decision tree model is simple and easy to interpret. So we could combine multiple decision tree models to get a better prediction performance.

## 14.2. Ensemble Learning

- The base learner used for ensemble should be:
  1. Comparatively accurate
  2. Diverse
- Error rate of a majority vote ensemble: Assume we have  $N$  independent base learners, each of them has error rate  $\varepsilon$ . Then the error rate of a majority vote ensemble is:

$$P(\text{Error}) = \sum_{k=0}^{\lfloor \frac{N}{2} \rfloor} \binom{N}{k} (1 - \varepsilon)^k \varepsilon^{N-k} \leq \exp\left(-\frac{N(1 - 2\varepsilon)^2}{2}\right)$$

When  $\varepsilon < 0.5$ , the upper bound of the error rate would converge to 0 as  $N$  (number of base learners) increases.

### Upper Bound for Error Rate

First, we introduce the Hoeffding's inequality. Let  $\{X_i | 1 \leq i \leq n, \forall i, a \leq X_i \leq b\}$  be a set of bounded random variables, and  $-\infty < a \leq b < \infty$ . Then we have:

$$P\left(\frac{1}{n} \sum X_i - \frac{1}{n} \sum E(X_i) \leq -\delta\right) \leq \exp\left(-\frac{2N\delta^2}{(b-a)^2}\right)$$

Now assume the set of random variables is our base learner and if they give correct prediction then the value is 1.  $X_i \in [0, 1] \Rightarrow a = 0, b = 1$  and  $E[X_i] = 1 - \varepsilon$

$$\begin{aligned} P\left(\frac{1}{n} \sum X_i < \frac{1}{2}\right) &= P\left(\frac{1}{n} \sum X_i - (1 - \varepsilon) < \frac{1}{2} - (1 - \varepsilon)\right) \\ &\Rightarrow \delta = \varepsilon - \frac{1}{2}, (b-a)^2 = 1 \end{aligned}$$

Therefore,

$$P(\text{Error}) = P\left(\frac{1}{n} \sum X_i < \frac{1}{2}\right) \leq \exp\left(-\frac{2N(\varepsilon - \frac{1}{2})^2}{1}\right) = \exp\left(-\frac{N(1 - 2\varepsilon)^2}{2}\right)$$

## 14.3. Bagging

**Goal:** Obtain a set of diverse and high-performance base learners.

**Methods:** Different models, different hyperparameters, different training Algorithms

### 14.3.1. How to bagging

1. Sampling: The method is bootstrap sampling, which means we sample with replacement from the dataset. Each sample are of same size with slight variation.
2. Training: Train each base learner with the sampled dataset.
3. Aggregation: Majority vote for classification problem, average for regression problem.

Each bootstrap sampling process could only cover roughly 63.2% of the dataset. So we could use the remaining samples to evaluate the performance of the model. This is called out-of-bag (OOB) error. OOB error could help for hyperparameter tuning, early stopping and cross-validation.

Assume we want to learn the function  $g$ , and trained result is  $f'$ , the expected model is  $E[f']$ . The error is:

$$\begin{aligned} E[(g - f')^2] &= E[(g - E[f'] + E[f'] - f')^2] \\ &= (g - f')^2 + E[(f' - E[f'])^2] \end{aligned}$$

The first term is the bias of the model and the second term is the variance of the model. The following are several conditions we have:

1. The expectation of bagged tree is equal to the expectation of the individual tree.

$$E[f'] = E\left[\frac{1}{B} \sum_{b=1}^B f'_b\right] = \frac{1}{B} \sum_{b=1}^B E[f'_b] = E[f'_b]$$

2. Bagged tree has the same bias as the individual tree.
3. Each tree is identically distributed but might not be independent identically distributed.

Therefore, we have the following:

**Claim:**

Correlated learners could not effectively reduce the variance.

**Proof:**

If we have a size  $B$  bagged model, the change of variance has 2 cases. For i.i.d random variables  $X_1, X_2, \dots, X_n$ , we have:

$$\text{Var}(\bar{X}) = \text{Var}\left(\frac{1}{B} \sum_i X_i\right) = \frac{1}{B^2} \text{Var}\left(\sum_i X_i\right) = \frac{1}{B^2} B \text{Var}(X_i) = \frac{1}{B} \text{Var}(X_i) = \frac{\sigma^2}{B}$$

For i.d random variables  $X_1, X_2, \dots, X_n$ , let say the pairwise correlation between the base learners is  $\rho$ . The variance is:

$$\text{Var}(\bar{X}) = \frac{1}{B^2} \text{Var}\left(\sum_i X_i\right) = \frac{\sum_i \text{Var}(X_i) + \sum_{i \neq j} \text{Cov}(X_i, X_j)}{B^2} = \left(\rho + \frac{1 - \rho}{B}\right) \sigma^2$$

Therefore, when  $\rho \rightarrow 0, \sigma^2 \rightarrow \frac{\sigma^2}{B}$ . When  $\rho \rightarrow 1, \sigma^2 \rightarrow \sigma^2$ . When the correlation between the base learners is large, the variance does not decrease much.

## 14.4. Random Forest

## 14.5. Boosting

## 14.6. Feature Selection

# 15. Multi Layer Perceptron

## 15.1. Perceptron

- Output

It's similar to SVM. Given an input  $\mathbf{x}$ , weight  $\mathbf{w}$  and bias  $b$ , the model output:

$$o = \sigma(\mathbf{w}^T \mathbf{x} + b), \text{ where } \sigma(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$$

We could also output the probability using softmax function.

- Training We have algorithm to train the model:

1	Algo
2	initialize $w = 0, b = 0$
3	repeat
4	if $y_i (< w, x_i > + b) \leq 0$
5	$w \leftarrow w + y_i x_i, b \leftarrow b + y_i$
6	end if
7	until all classified correctly

The loss function is  $\ell(\mathbf{x}, y, \mathbf{w}) = \max(0, -y_i \langle \mathbf{x}, \mathbf{w} \rangle)$  and the operation in loop is actually a gradient descent using 1 data each time.

#### Convergence Theorem

1. The data lies in an area of radius  $R$
2.  $\exists \rho, s.t. y(\mathbf{w}^T \mathbf{x} + b) \geq \rho$ , where  $\|\mathbf{w}\|^2 + b^2 \leq 1$
3. The model converges after  $\frac{R^2+1}{\rho^2}$  iterations.

The 1 layer perceptron is a linear classifier and could not be used to solve non-linear problems, such as XOR problem.

## 15.2. Multi Layer Perceptron

First we consider the situation of only 1 hidden layer.

1. Input layer:  $\mathbf{x} \in \mathbb{R}^m$
2. Hidden layer:  $\mathbf{W}_1 \in \mathbb{R}^{m \times n}, b_1 \in \mathbb{R}^n$
3. Output layer:  $\mathbf{w}_2 \in \mathbb{R}^{n \times 1}, b_2 \in \mathbb{R}$

Procedure:

1.  $\mathbf{h} = \sigma(\mathbf{W}_1^T \mathbf{x} + b_1)$
2.  $o = \mathbf{w}_2^T \mathbf{h} + b_2$

$\sigma$  is the activation function. We could use sigmoid function, tanh function, relu function, etc.

1. Sigmoid function:  $\sigma(x) = \frac{1}{1+e^{-x}}$
2. Tanh function:  $\tanh(x) = \frac{1-e^{-2x}}{1+e^{-2x}}$
3. ReLU function:  $\text{relu}(x) = \max(0, x)$

Secondly, to do task of multi-class classification, we could use softmax function.

1. Input layer:  $\mathbf{x} \in \mathbb{R}^m$

2. Hidden layer:  $\mathbf{W}_1 \in \mathbb{R}^{m \times n}$ ,  $b_1 \in \mathbb{R}^n$  (we could have many layers here, the number of layers and the size of each layer are hyperparameters)
3. Output layer:  $\mathbf{W}_2 \in \mathbb{R}^{n \times k}$ ,  $b_2 \in \mathbb{R}^k$

Procedure:

1.  $\mathbf{h} = \sigma(\mathbf{W}_1^T \mathbf{x} + b_1)$
2.  $\mathbf{O} = \mathbf{W}_2^T \mathbf{h} + b_2$
3.  $\mathbf{y} = \text{softmax}(\mathbf{O})$

When designing a MLP, the size of layers are important. The size of layers could be determined by the following rules:

1. The size of first hidden layer could be larger to track more complex features.
2. The size of following hidden layers could be reduced step by step to get the same size of expected output.