

# Week3 - Linear Regression Models

- Polynomial is a specific example of a broad class of function call Linear Regression Models
- Having the property of being linear functions of the adjustable parameters
- The simplest form of the linear regression models which is also linear to the input variable:  
 $y(x, \mathbf{w}) = w_0 + w_1x + w_2x + \dots + w_{M-1}x$  (the simplest basis function is  $\phi_i(x) = x$ .)

## 1. Basis Function

- Usually taking linear combination of a fixed set of **nonlinear functions of the input variables**, known as **basis functions**.
- **Aim**: predict the value of continuous target variable  $t$
- **Input**: d-dimensional vector  $\mathbf{x}$  of input variables
- Generalised form:  
$$y(x, \mathbf{w}) = w_0 + w_1\phi_1(x) + w_2\phi_2(x) + \dots + w_{M-1}\phi_{M-1}(x) = w_0 + \sum_{j=1}^{M-1} w_j\phi_j(x) = \sum_{j=0}^{M-1} w_j\phi_j(x)$$
  - If  $\phi_i(x)$  is linear of  $x$ , it is called **linear basis function**
  - Otherwise, it is called **non-linear basis function**
- Basis function provide a way of making a linear regression model nonlinear
- Dummy basis function:  $\phi_0(x) = 1$

### Popular Basis functions

- Gaussian:  $\phi_j(x) := \exp\left(-\frac{(x - \mu_j)^2}{2s^2}\right)$
- Sigmoidal:  $\phi_j(x) := \sigma\left(\frac{x - \mu_j}{s}\right)$ 
  - Logistic sigmoid function:  $\sigma(a) := \frac{1}{1 + \exp(-a)}$
  - Hyperbolic tangent function:  $\tanh(a) = 2\sigma(2a) - 1 = \frac{1 - e^{-2a}}{1 + e^{-2a}}$

## 2. Loss function

We fit a regression model to the training data by minimising an error function that measures the misfit between the predictions  $y(\mathbf{x}, \mathbf{w})$  for each data training point  $x_n$  and the target value  $t_n$ .

- Calculate the misfit: SSE - the sum of the squares of the errors

$$E(\mathbf{w}) := \frac{1}{2} \sum_{n=1}^N [y(x_n, \mathbf{w}) - t_n]^2$$

It shifts from minimising the error function to finding  $\mathbf{w}$  that satisfies the function:

$$\mathbf{w} := \arg \min_{\mathbf{w}} E(\mathbf{w})$$

- Components:
  - **Noise**: Data collected by sampling from ensemble have inevitably included noise:

$\epsilon = N(0, \sigma^2)$ . Later we will find it represents the irreducible minimum value of the loss function.

- Model Function:  $y(\mathbf{x}, \mathbf{w})$
- Observed target value:  $\mathbf{t} := \{t_n\}$  where  $n \in [1, N]$
- Assume the training data points are drawn independently (**i.i.d**) from the underlying ensemble:  $\mathbf{x} := \{x_n\}$ , where  $n \in [1, N]$

Then, The likelihood:

$$p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \sigma^2) = \prod_{n=1}^N N(t_n|y(x_n, \mathbf{w}), \sigma^2)$$

The log-likelihood function:

$$\ell(\mathbf{w}) = \log(p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \sigma^2)) = \sum_{n=1}^N \log N(t_n|y(x_n, \mathbf{w}), \sigma^2)$$

Let us use precision to replace variance:  $\beta^{-1} = \sigma^2$

$$\ell(\mathbf{w}) = \frac{N}{2} \ln \beta - \frac{N}{2} \ln 2\pi - \beta E(\mathbf{w})$$

Note that only the last term relevant to the function  $\mathbf{w}$ .

Then differentiates the equation with respect to  $\mathbf{w}$ , we can have the gradient of the log likelihood:

$$\nabla \ell(\mathbf{w}) = \beta \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(x_n)\} \phi(x_n)^T$$

Let the gradient equals to zero, we can find the best  $\mathbf{w}$  that maximise the log-likelihood function:

$$\nabla \ell(\mathbf{w}) = \begin{bmatrix} \frac{\partial \ell(\mathbf{w})}{\partial w_0} \\ \frac{\partial \ell(\mathbf{w})}{\partial w_1} \\ \vdots \\ \frac{\partial \ell(\mathbf{w})}{\partial w_{M-1}} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \Rightarrow \mathbf{w} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

We can find the parameter by solving the matrix, elegant but expensive

### 3. Iterative Optimization Algorithms

- Step by step:
  - Initialise the parameters at random
  - Iteratively moves towards the negative direction of the objective function gradient to find the optimal parameters

Gradient: 梯度是一个矢量，其的方向取决于其方向导数中最大值的方向，梯度的值是方向导数的最大值（最陡）。

需要注意的是，梯度指的是增长最快的方向。因此对于梯度下降，我们需要朝某点梯度的反方向前进才能得到其驻点。

- **Gradient Descent Algorithms (with Batch technique: BGD):**
  - Pseudocode:

```

# Gradient Descent
Initialise w and t = 1, t_max = N
Initialise eta and epsilon #learning rate & threshold
While (diff > epsilon and t <= t_max)
do{
    w_new := w - eta * gradient(w,X)
    diff := gradient(w_new, X) - gradient(w,X)
    if diff < 0 then break
    else eta = eta * 0.5 #control learning rate
    w := w_new # update
    t = t + 1 #increment
}

```

- Error function:  $E(\mathbf{w}) = \sum_{n=1}^n E_n(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N [t_n - \mathbf{w} \cdot \phi(x_n)]^2$
- Gradient:  $\nabla E(\mathbf{w}) = - \sum_{n=1}^N [t_n - \mathbf{w} \cdot \phi(x_n)] \cdot \phi(x_n)$
- An algorithm that minimises functions. It starts with an initial set of parameter value  $\mathbf{w}^{(0)}$ , and iteratively moves towards optimal parameters which achieve the minimum value of the objective function.
- Process the **entire training set** at each iteration.
- Computationally **costly** for large data sets
- Not guarantee to converge: sensitive to starting points, may reach local minima instead of global minima
- **Stochastic Gradient Descent Algorithms (SGD):**
  - Pseudocode

```

# Stochastic Gradient Descent
Initialise w and t = 1, t_max = N
Initialise eta and epsilon #learning rate & threshold
While (diff > epsilon and t <= t_max)
do{
    for each training point x_i
    w_new := w - eta * gradient(w,X)
    diff := gradient(w_new, X) - gradient(w,X)
    if diff < 0 then break
    else eta = eta * 0.5 # control learning rate
    w := w_new # update
    t = t + 1 #increment
}

```

- Error function:  $E(\mathbf{w}) = \sum_{n=1}^n E_n(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N [t_n - \mathbf{w} \cdot \phi(x_n)]^2$
- Gradient:  $\nabla E_n(\mathbf{w}) = -[t_n - \mathbf{w}^T \cdot \phi(x_n)] \cdot \phi(x_n)$
- Process **one data points** at each iteration  $E(\mathbf{w}) = \sum_{n=1}^n E_n(\mathbf{w})$
- **Incrementally update** model parameters after each such presentation
  - known as **least-mean-squares (LMS)**:  $\mathbf{w}^t = \mathbf{w}^{t-1} - \eta \cdot \nabla E_n(\mathbf{w}^{t-1})$
  - The value of  $\eta$  needs to be chosen with care to **ensure that the algorithm**

**converges.**

- Not necessary to process the entire data set
- **Faster than BGD** to get  $\mathbf{w}$  close to the minimum but not **never converge** to the minimum (keep **oscillating around** the minimum of the objective function)
- Sensitive to **feature scaling**, need to standardise training data
- **Problem:**
  - Need to compute the gradient of the error function  $\nabla E(\mathbf{w})$ .
  - Need to consider the **starting point**  $\mathbf{w}^{(0)}$ , the **learning rate**  $\eta$ , the **threshold**  $\epsilon$ , the **maximum iteration number**  $t_{max}$ 
    - Starting point  $\mathbf{w}^{(0)}$ : reach local minima
    - Learning rate (step-size)  $\eta$ :
      - Too large:
      - Too small: too slow to convergence
    - Threshold  $\epsilon$ : criteria to stop the update process
    - Maximum iteration number  $t_{max}$ : criteria to stop the update process

## 4. Regularisation

- The total error function with regularization term to control over-fitting:  $E_D(\mathbf{w}) + \lambda E_W(\mathbf{w})$
- The general regularization term:  $E_W(\mathbf{w}) = \frac{1}{2} \sum_{j=0}^{M-1} |w_j|^q$ 
  - The simplest regularization term given by the sum-of-squares of the weight vector elements(parameters):  $E_W(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} = \frac{1}{2} \sum_{j=0}^{M-1} |w_j|^2$
  - Different  $q$  gives different regularization function. **The choice of  $q$  depend on problems.**

### L2 Ridge Regularisation: $q = 2$

- Known as **weight decay**: encourage weight values to **decay towards zero**, unless supported by the data.
- Advantage: the error function remains a quadratic function of  $\mathbf{w}$ , the exact minimizer can be found in closed form:  $\mathbf{w} = (\lambda \mathbf{I} + \Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$

### L1 Lasso Regularisation: $q = 1$

- $E_W(\mathbf{w}) = \frac{1}{2} \sum_{j=0}^{M-1} |w_j|$
- If  $\lambda$  is sufficiently large, some of the coefficient  $w_j$  are driven to zero, leading to a **sparse model in which the corresponding basis functions play no role**
- As  $\lambda$  is increased, an increasing number of parameters are driven to zero.
- Can be use for **feature selection**

## Key Difference: L2 v.s L1

- **L2 Ridge:** Mainly for **prevent overfitting, not useful for large number of input variables.**
  - includes all of the input variables in the model
  - Parameter shrinkage (limit the magnitude)
  - Reduce model complexity
- **L1 Lasso:** Mainly for **sparse solution, useful for large number of input variables.**
  - Parameter shrinkage
  - Some coefficients are driven to zero (**feature selection:** some input variables are excluded)

Note: Limiting the number of basis functions in order to avoid over-fitting has the side effect of limiting the flexibility of the model to capture interesting and important trends in the data.