

Regression

A training set of P input/output observation pairs:

$$\{(\bar{x}_1, y_1), (\bar{x}_2, y_2), \dots, (\bar{x}_P, y_P)\}$$

where

$$\bar{x}_p = \begin{bmatrix} x_{1,p} \\ x_{2,p} \\ \vdots \\ x_{N,p} \end{bmatrix}$$

Then we have a bias and N associated weights,

$$b + \bar{x}_p \bar{w} \approx y_p, \quad p=1, \dots, P.$$

The Least Squares cost function for linear regression:

For a given set of parameters (b, \bar{w}) this cost function computes the total squared error between the associated hyperplane and the data.

$$g(b, \bar{w}) = \sum_{p=1}^P (b + \bar{x}_p^T \bar{w} - y_p)^2$$

$$\therefore \underset{b, \bar{w}}{\text{minimize}} \sum_{p=1}^P (b + \bar{x}_p^T \bar{w} - y_p)^2$$

Minimization of the Least Squares cost function.

Assume

$$\tilde{x}_p = \begin{bmatrix} 1 \\ \bar{x}_p \end{bmatrix} \quad \tilde{w} = \begin{bmatrix} b \\ \bar{w} \end{bmatrix}$$

$$\therefore g(\tilde{w}) = \sum_{p=1}^P (\tilde{x}_p^T \tilde{w} - y_p)^2$$

$$\nabla g(\tilde{w}) = 2 \sum_{p=1}^P \tilde{x}_p (\tilde{x}_p^T \tilde{w} - y_p) = 2 \left(\sum_{p=1}^P \tilde{x}_p \tilde{x}_p^T \right) \tilde{w} - 2 \sum_{p=1}^P \tilde{x}_p y_p$$

$$\left(\sum_{p=1}^P \tilde{x}_p \tilde{x}_p^T \right) \tilde{w} = \sum_{p=1}^P \tilde{x}_p y_p$$

$$\Rightarrow \tilde{w}^* = \left(\sum_{p=1}^P \tilde{x}_p \tilde{x}_p^T \right)^{-1} \sum_{p=1}^P \tilde{x}_p y_p$$

Mean squared error (MSE): $MSE = \frac{1}{P} \sum_{p=1}^P (b^* + \tilde{x}_p^T \tilde{w}^* - y_p)^2 \rightarrow$ efficacy of a learned model.

more generally, $b + f(x_p)w \approx y_p, p=1, \dots, P$

$$\underset{b, w}{\text{minimize}} \sum_{p=1}^P (b + f_p w - y_p)^2$$

$$\text{Assume } \tilde{f}_p = \begin{bmatrix} 1 \\ f_p \end{bmatrix}, \tilde{w} = \begin{bmatrix} b \\ w \end{bmatrix}$$

$$g(\tilde{w}) = \sum_{p=1}^P (\tilde{f}_p^T \tilde{w} - y_p)^2$$

$$\therefore \tilde{w}^* = \left(\sum_{p=1}^P \tilde{f}_p \tilde{f}_p^T \right)^{-1} \sum_{p=1}^P \tilde{f}_p y_p$$

"tanh" function (step function)

$$\tanh(b + \tilde{x}^T \tilde{w}) = 2\sigma(b + \tilde{x}^T \tilde{w}) - 1$$

$$\tanh(b + \tilde{x}^T \tilde{w}) \approx \begin{cases} +1 & \text{if } b + \tilde{x}^T \tilde{w} > 0 \\ -1 & \text{if } b + \tilde{x}^T \tilde{w} \leq 0 \end{cases}$$

$$\therefore \tanh(y_p(b + \tilde{x}^T \tilde{w})) \approx 1$$

$$\Rightarrow 1 + e^{-y_p(b + \tilde{x}^T \tilde{w})} \approx 1$$

$$\Rightarrow \log(1 + e^{-y_p(b + \tilde{x}^T \tilde{w})}) \approx 0$$

$$\therefore \underset{b, \tilde{w}}{\text{minimize}} \sum_{p=1}^P \log(1 + e^{-y_p(b + \tilde{x}_p^T \tilde{w})})$$

Chapter 4 Automatic feature design for regression.

4.1 Automatic feature design for the ideal regression scenario

4.1.1 continuous function approximation

A basis is a set of basic feature transformation $\{f_m(\bar{x})\}_{m=1}^{\infty}$ such that all points \bar{x} in the unit hypercube we can express $y(\bar{x})$ perfectly as

$$\sum_{m=0}^{\infty} f_m(\bar{x}) w_m = y(\bar{x})$$

for large enough M we can approximate y over its input domain as:

$$\sum_{m=0}^M f_m(\bar{x}) w_m \approx y(\bar{x})$$

This approximation can be made as finely as desired by increasing the number of basis features M , i.e., by increasing M we can design automatically perfect features to represent $y(\bar{x})$.

4.1.2 Common base for continuous function approximation.

Polynomial basis: ($N=1$)

$$\begin{cases} f_0(x) = 1 & m=0 \\ f_m(x) = x^m & \text{for all } m \geq 1 \end{cases}$$

Fourier basis / sinusoidal basis ($N=1$)

$$\begin{cases} f_0(x) = 1 & m=0 \\ f_{2m-1}(x) = \cos(2\pi m x) & \text{for all } m \geq 1 \\ f_{2m}(x) = \sin(2\pi m x) & \text{for all } m \geq 1 \end{cases}$$

multilayer Neural network (N dimensional) basis
Activation function:

① The hyperbolic tangent function:

$$f_m(\bar{x}) = \tanh(c_m + \bar{x}^T \bar{v}_m) \text{ for all } m \geq 1$$

② The Max or hinge function:

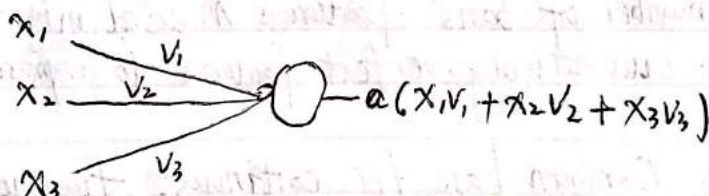
$$f_m(\bar{x}) = \max(0, c_m + \bar{x}^T \bar{v}_m) \text{ for all } m \geq 1$$

Graphical representation of a neural network

It is common to represent the weighted sum of M neural network basis features.

$$r = b + \sum_{m=1}^M f_m(\bar{x}) w_m$$

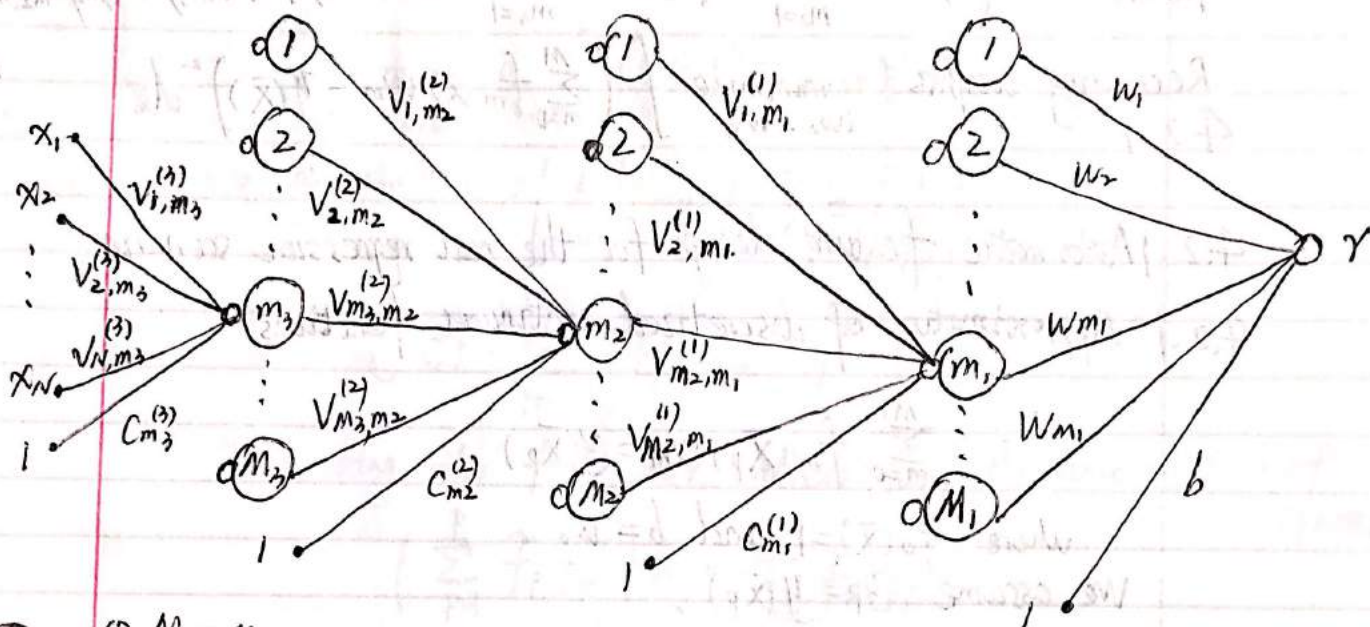
A simple example.



where $a(\cdot)$ is any activation function.

Time $\left\{ \begin{array}{l} \text{weighted edges} \\ \text{summation unit} \\ \text{activation unit} \end{array} \right.$

Three hidden layer neural network.



- ① $M = m_1$
- ② input \bar{x} is shown entry-wise on the left and output r on the right.
- ③ In between the three hidden layers, from right to left each consist of m_1, m_2 and m_3 hidden units respectively.
- ④ A basis with $L > 2$ or 3 hidden-layers is usually called a deep network. A single basis element of a feed forward neural network with two hidden layers. Doing this with the tanh basis:

$$f_m(\bar{x}) = \tanh(C_m^{(1)} + \sum_{m_2=1}^{m_2} \tanh(C_{m_2}^{(2)} + \bar{x}^T \bar{V}_{m_2}^{(2)}) V_{m_2, m}^{(1)})$$

or hinge function:

$$f_m(\bar{x}) = \max(0, C_m^{(1)} + \sum_{m_2=1}^{m_2} \max(0, C_{m_2}^{(2)} + \bar{x}^T \bar{V}_{m_2}^{(2)}) V_{m_2, m}^{(1)})$$

A neural network with three hidden layers hinge basis function:

$$f_m(\bar{x}) = \max\left(0, c_m^{(1)} + \sum_{m_2=1}^{m_2} \max\left(0, c_{m_2}^{(2)} + \sum_{m_3=1}^{m_3} \max\left(0, c_{m_3}^{(3)} + \bar{x}^T V_{m_3}^{(3)}\right) V_{m_3, m_2}^{(2)}\right) V_{m_2, m}^{(1)}\right)$$

Recovering weights: minimize $\int \left(\sum_{m=0}^M f_m(\bar{x}) w_m - y(\bar{x}) \right)^2 d\bar{x}$
 w_0, \dots, w_M

4.2. Automatic feature design for the real regression scenario.

4.2.1 Approximation of discretized continuous functions

$$\sum_{m=0}^M f_m(\bar{x}_p) w_m = y(\bar{x}_p)$$

where $f_0(\bar{x}) = 1$ and $b = w_0$.

We assume $y_p = y(\bar{x}_p)$,

feature vector $\bar{f}_p = [f_1(\bar{x}_p) \ f_2(\bar{x}_p) \ \dots \ f_M(\bar{x}_p)]^T$

weight vector $\bar{w} = [w_1 \ w_2 \ \dots \ w_M]^T$

$$\therefore b + \bar{f}_p^T \bar{w} \approx y_p$$

$$\therefore \text{minimize}_{b, \bar{w}} \sum_{p=1}^P (b + \bar{f}_p^T \bar{w} - y_p)^2$$

4.2.2 The real regression scenario

The general instance of regression is a function approximation problem based on noisy samples of the underlying function.

Exp. Regression with fixed basis of features.

Employing a degree D polynomial or Fourier basis for a scalar input.

$$\bar{f}_p = [x_p \ x_p^2 \ \dots \ x_p^D]^T \quad \text{or} \quad \bar{f}_p = [\cos(2\pi x_p) \ \sin(2\pi x_p)]^T$$

$$= [\cos(2\pi DX_p) \quad \sin(2\pi DX_p)]^T$$

$$\therefore \underset{b, \tilde{w}}{\text{minimize}} \sum_{p=1}^P (b + \tilde{f}_p^T \tilde{w} - y_p)^2$$

We assume $\tilde{w} = \begin{bmatrix} b \\ \tilde{w} \end{bmatrix}$ and $\tilde{f}_p = \begin{bmatrix} 1 \\ \tilde{f}_p \end{bmatrix}$

$$\therefore g(\tilde{w}) = \sum_{p=1}^P (\tilde{f}_p^T \tilde{w} - y_p)^2$$

checking the first order condition ^{then} gives the linear system of equations

$$\left(\sum_{p=1}^P \tilde{f}_p \tilde{f}_p^T \right) \tilde{w} = \sum_{p=1}^P \tilde{f}_p y_p$$

that when solved recovers an optimal set of parameters \tilde{w} .

Exp. Regression with a basis of single hidden layer neural network features.

$$\text{feature vector: } \tilde{f}_p = [a(c_1 + \bar{x}_p^T \bar{v}_1) \quad a(c_2 + \bar{x}_p^T \bar{v}_2) \quad \dots \quad a(c_m + \bar{x}_p^T \bar{v}_m)]^T$$

$$\underset{b, \tilde{w}}{\text{minimize}} \sum_{p=1}^P (b + \tilde{f}_p^T \tilde{w} - y_p)^2$$

$$\nabla g = \left[\frac{\partial}{\partial b} g \quad \frac{\partial}{\partial w_1} g \quad \dots \quad \frac{\partial}{\partial w_m} g \quad \frac{\partial}{\partial c_1} g \quad \dots \quad \frac{\partial}{\partial c_m} g \quad \nabla_{\bar{v}_1}^T g \quad \dots \quad \nabla_{\bar{v}_m}^T g \right]^T$$

$$\text{We have a vector of length } Q = m \cdot (N+1) + m+1 \\ = m \cdot (N+2) + 1$$

4.3 Cross-validation for regression

An effective framework for choosing the proper value for m automatically and intelligently so as to prevent the problem of underfitting/overfitting

4.3.1 Diagnosing the problem of overfitting/underfitting.

Criterion: the number m of basis features used should be such that the corresponding model fits well to both the current dataset as well as to new data we will receive in the future.

4.3.2 Hold out cross-validation.

In general the larger/smaller the original dataset, the larger/smaller the portion of original data that should be assigned to the testing set (between $1/10$ to $1/3$ typically)



original data



random splitting

training



testing

k non-overlapping sets (here $k=3$)

4.3.3 Hold out calculation.

$$\Omega_{\text{train}} = \{p \mid (\bar{x}_p, y_p) \text{ belongs to the training set}\}$$

$$\Omega_{\text{test}} = \{p \mid (\bar{x}_p, y_p) \text{ belongs to the testing set}\}$$

We then choose a basis type (e.g., polynomial, Fourier, neural network) and choose a range for the number of basis features over which we search for an ideal value for m .

Procedure: 1) Form the corresponding feature vector $F_p = [f, (\bar{x}_p)]$

$$f_2(\bar{x}_p) \dots f_m(\bar{x}_p)]^T$$

2) fit a corresponding model to the training set by LS.

$$\text{minimize}_{b, \bar{w}, \Theta} \sum_{p \in \Omega_{\text{train}}} (b + \bar{f}_p^T \bar{w} - y_p)^2$$

Denoting a solution to the problem above as $(b_m^*, \bar{w}_m^*, \Theta_m^*)$

$$\text{Training error} = \frac{1}{|\Omega_{\text{train}}|} \sum_{p \in \Omega_{\text{train}}} (b_m^* + \bar{f}_p^T \bar{w}_m^* - y_p)^2$$

$$\text{Testing error} = \frac{1}{|\Omega_{\text{test}}|} \sum_{p \in \Omega_{\text{test}}} (b_m^* + \bar{f}_p^T \bar{w}_m^* - y_p)^2$$

We choose the one that provides the lowest testing error denoted by m^* .

3) Form the feature vector $\bar{f}_p = [f_1(\bar{x}_p) f_2(\bar{x}_p) \dots f_m(\bar{x}_p)]^T$ for all the points in the entire dataset and solve the LS problem over the entire dataset to form the final model

4.3.4 k-fold cross-validation

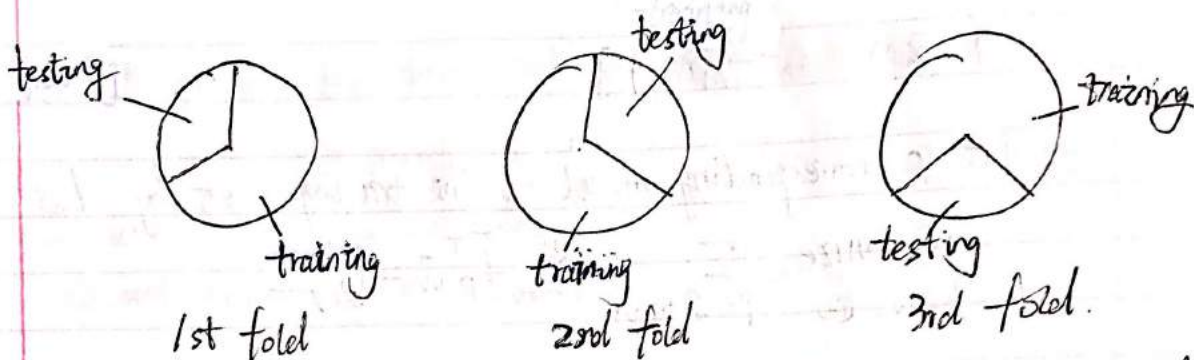
Flaw of Hold out method. Having been chosen at random, the points assigned to the training set may not adequately describe the original data.



original data



random splitting



By averaging the testing errors and choosing the model with minimum associated average test error, and determine an excellent model for the phenomenon.

Supplement:

Comparison:
Scalar x

Vector $\bar{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix}$

Approximation $y(x) \approx \sum_{m=0}^M f_m(x) w_m$

$y(\bar{x}) \approx \sum_{m=0}^M f_m(\bar{x}) w_m$

Polynomial $f_m(x) = x^m = f_m(\bar{x}) = x_1^{m_1} x_2^{m_2} \dots x_N^{m_N}$

Single hidden layer $f_m(x) = a(c_m + x v_m)$ $f_m(\bar{x}) = a(c_m + \bar{x}^T \bar{v}_m)$

For a general N -dimensional input \bar{x} ,

1) Polynomial: $f_m(\bar{x}) = x_1^{m_1} x_2^{m_2} \dots x_N^{m_N}$

where $0 \leq m_1 + m_2 + \dots + m_N \leq D$,

That means polynomial basis-feature includes all monomials of form $f_m(\bar{x})$

Exp. $N=2, D=3$

$$\begin{array}{ccc} \frac{x_1}{D_1} & \frac{x_1 x_2}{D_2} & \frac{x_1^2 x_2}{D_3} \\ \frac{x_2}{D_1} & \frac{x_1^2}{D_2} & \frac{x_1 x_2^2}{D_3} \\ \frac{x_2^2}{D_1} & \frac{x_2^3}{D_2} & \frac{x_2^3}{D_3} \end{array}$$

$$M = 9 + 1 = 10$$

$$M = \binom{N+D}{D} = \frac{(N+D)!}{D! (N+D-D)!} = \frac{(N+D)!}{D! N!}$$

2) Fourier:

$$f_m(\bar{x}) = e^{2\pi i m_1 x_1} e^{2\pi i m_2 x_2} \dots e^{2\pi i m_N x_N}$$

where $-D \leq m_1, m_2, \dots, m_N \leq D$.

$$M = (2D+1)^N$$

$$\therefore \cos(\alpha) = \frac{1}{2}(e^{j\alpha} + e^{-j\alpha}) \text{ and } \sin(\alpha) = \frac{1}{2j}(e^{j\alpha} - e^{-j\alpha})$$

We can write the partial Fourier expansion.

$$w_0 + \sum_{m=1}^M (\cos(2\pi m x) w_{2m-1} + \sin(2\pi m x) w_{2m})$$

equivalently as

$$\sum_{m=-M}^M e^{2\pi i m x} w_m$$

where $w'_m = \begin{cases} \frac{1}{2}(w_{2m-1} - j w_{2m}) & \text{if } m > 0 \\ w_0 & \text{if } m = 0 \\ \frac{1}{2}(w_{-2m-1} - j w_{2m}) & \text{if } m < 0 \end{cases}$