

UNIT - II

Regression & Applications

Regression :

- Main task of analytics is Prediction
- + induction model is used to predict and assigns labels to a new, unlabeled, object
- Prediction is used to
 - reduce cost
 - Increase profits
 - Improve Product & service quality
 - Improve customer satisfaction
 - Reduce environmental damage.

→ Test data

- ↳ used to test the performance of the induced model.

→ deduction

- Predict Correct label.

Regression task:

A predictive task whose aim is to assign a quantitative value to a new, unlabeled object, given the values of its predictive attributes.

Regression methods used in many different domains

① stock market

② Transport

③ Higher Education

④ Survival analysis

⑤ Macro economics

Types:

Linear Regression

Ridge

Lasso

Principal Components Regression

Partial Least squares

Linear Regression

- oldest & simplest regression alg.

- Induce good regression models, which are easily interpretable

$$\text{height} = 128.017 + 0.611 \times \text{weight}$$

- instance 'x' associated with only one attribute
- 'y' is associated with the height
- 2 parameters $\hat{\beta}_0 \downarrow \hat{\beta}_1$
 - \downarrow associated with x (weight)
 - Simple linear reg
 - Multiple ...
- 2 types

$$\underset{\hat{\beta}_0, \hat{\beta}_1}{\operatorname{argmin}} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \underset{\hat{\beta}_0, \hat{\beta}_1}{\operatorname{argmin}} \sum_{i=1}^n (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_{i1}))^2$$

- Multivariate linear regression - the linear model generalized for any no. of predictive attributes

$$\hat{y} = \beta_0 + \sum_{j=1}^p \beta_j x_{ij}$$

β_j slope of linear model

value of \hat{y} , when $x_j = 0$

x_j - j th attribute of some object x represented as tuple $x = (x_1, \dots, x_j, \dots, x_p)$.

$\beta_0, \beta_1, \dots, \beta_p$ - estimated using an appropriate optimization method to minimize objective function

$$\underset{\beta_0 - \hat{\beta}_p}{\operatorname{argmin}} \sum_{i=1}^n (y_i - \hat{y}_i)^2 =$$

$$\underset{\hat{\beta}_0 - \hat{\beta}_p}{\operatorname{argmin}} \sum_{i=1}^n \left[y_i - \left(\hat{\beta}_0 + \sum_{j=1}^p \hat{\beta}_j x_{ij} \right) \right]^2$$

Lasso Regression:

- Least absolute shrinkage & selection operator regression alg.
- deal efficiently with high dimensional data sets.
- Does not predict just performs attribute selection.
- Complexity measured by predictive attributes.
- Usually produces sparse solutions.
- sparse means - large no. of predictive attributes have zero weight.
- Performs shrinkage.

$$\underset{\hat{\beta}_0 - \hat{\beta}_1}{\operatorname{argmin}} \left[\sum_{i=1}^n \text{error}(y_i, \hat{y}_i) + \lambda * \sum_{j=1}^p |\hat{\beta}_j| \right]$$

Principal Components Regression

- PCR creates linear combinations of predictive attributes.
- 1st principal Component: the first linear combination, is the most variance of all possible linear combinations.
- just evaluate attribute without considering target attributes.
- Used to predictive attributes in the formulation of the multivariate linear regression problem.

Partial Least Squares Regression:

- PLs starts by evaluating the correlation of each predictive attribute with the target attribute.

- 1st component is linear combination of the predictive attributes.
- gives similar result as PCR.