

Data Analysis 3

Regression

Jan Platoš

December 9, 2018

Department of Computer Science
Faculty of Electrical Engineering and Computer Science
VŠB - Technical University of Ostrava



Regression

Regression

- In several cases the class label is numerical.
- The goal is to minimize the squared error of prediction.
- The predicted class label is also referred to as the response variable, dependent variable or regressand.
- The feature variables are referred to as explanatory variables, input variables, predictor variables, independent variables or regressors.
- The prediction process is referred to as regression modeling.

Regression - Linear Regression

- Let D be an $n \times d$ data matrix.
- The feature vector \bar{X}_i , i -th row of D , is the d -dimensional input vector.
- The corresponding response variable is y_i .
- In linear regression, the dependence of each response variable y_i on the \bar{X}_i is modeled as a linear relationship:

$$y_i \approx \bar{W} \cdot \bar{X}_i \quad \forall i \in \{1, \dots, n\}$$

- $\bar{W} = (w_1, \dots, w_d)$ is a d -dimensional vector of coefficients that needs to be learned from the training data to minimize the unexplained error

$$E = \sum_{i=1}^n (\bar{W} \cdot \bar{X}_i - y_i)^2$$

Regression - Linear Regression

- The bias b may be modeled:
 - as a part of the \bar{W} and artificial dimension in training data that is set to 1.
 - removed due to mean-centered the data matrix and response variables.
- The data are normalized/standardized to ensure similar scaling and weighting for all attributes.

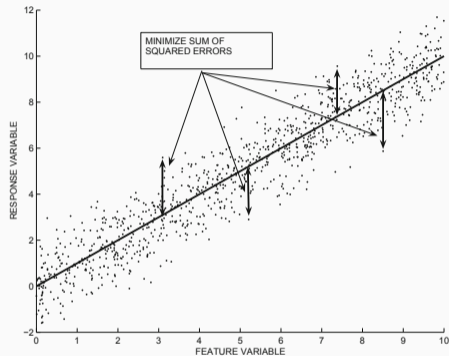


Figure 1: Linear regression

Regression - Linear Regression

- The objective function O , squared error of prediction, the have to be minimized by determination of W is defined as (where $\bar{y} = (y_1, \dots, y_n)$):

$$O = \sum_{i=1}^n (\bar{W} \cdot \bar{X}_i - y_i)^2 = \left\| D\bar{W}^T - \bar{y} \right\|^2$$

- The gradient of O with respect to \bar{W} is a vector $2D^T(D\bar{W}^T) = D^T\bar{y}$.
- Setting the gradient equal to o we get:

$$D^T D\bar{W}^T = D^T\bar{y}$$

- When $D^T D$ is invertible then $\bar{W}_T = (D^T D)^{-1} D^T\bar{y}$.
- otherwise we may use pseudo-inverse $D^+ = (D^T D)^{-1} D^T$ and then $\bar{W}_T = D^+\bar{y}$.

Regression - Linear Regression Regularization

- The objective function O minimizes the SSE:

$$O = \left\| D\bar{W}^T - \bar{y} \right\|^2$$

- *Ridge regression* reduce the size of the coefficient and minimizes chaotic behavior.

$$O = \left\| D\bar{W}^T - \bar{y} + \lambda \|\bar{W}\|^2 \right\|^2$$

- *Lasso regression* eliminates small weight (produces sparse model).

$$O = \left\| D\bar{W}^T - \bar{y} + \lambda \sum_{i=1}^d |w_i| \right\|^2$$

- Mixing model between *Lasso* and *Ridge* is called *ElasticNet*.

- Intuitively, we expect that a constant change in a feature variable leads to the constant change in the response variable.
- This is not true in many cases, e.g. the height of the person is not linearly dependent on the age of a person.
- Moreover, such features will never be negative.
- The generalized linear models (GLM) solves this problems.
- Each responsible variable y_i is modeled as an outcome of a probability distribution with mean $f(\bar{W} \cdot \bar{X}_i)$.
- The function $f(\cdot)$ is referred to as the *mean function* and its inverse as *link function*.
- The selection of the mean/link function and corresponding probability distribution should maximize effectiveness and interpretability of the model.

Regression - Generalized Linear Models

- Intuitively, we expect that a constant change in a feature variable leads to the constant change in the response variable.
- This is not true in many cases, e.g. the height of the person is not linearly dependent on the age of a person.
- Moreover, such features will never be negative.
- The generalized linear models (GLM) solves this problems.
- Each responsible variable y_i is modeled as an outcome of a probability distribution with mean $f(\bar{W} \cdot \bar{X}_i)$.
- The function $f(\cdot)$ is referred to as the *mean function* and its inverse as *link function*.
- The selection of the mean/link function and corresponding probability distribution should maximize effectiveness and interpretability of the model.
- the response variable is modeled using probability, the \bar{W} is determined using maximum likelihood approach.

Regression - Nonlinear and polynomial regression

- Linear regression cannot capture nonlinear relationships.
- Linear approach may be applied on the derived features.
- The derivation means an application of non-linear functions on the each input points.
- The new set of points may have different number of dimensions.

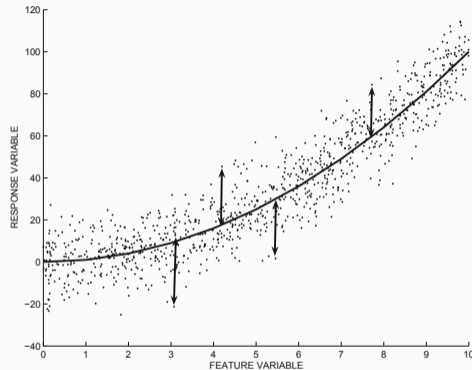


Figure 2: Linear regression

Regression - Nonlinear and polynomial regression

- The new set of m features denoted as $h_1(\bar{X}_i) \dots h_m(\bar{X}_i)$ for the i -th data point.
- The $h(\cdot)$ represents nonlinear transformation from the d -dimensional input feature space into 1-dimensional space.
- The size of the new dataset D_h is $n \times m$.
- The linear relationship is then defined as :

$$y = \sum_{i=1}^m w_i h_i(\bar{X})$$

- The *polynomial regression* expands the number of features by factor r

$$\bar{X} = (x_1, \dots, x_d) \Rightarrow \bar{X}^h = (x_1, x_1^2, x_1^3, \dots, x_1^r, x_2, \dots, x_d^r)$$

- The Kernel trick is allows by the reformulation of the regression problem with dot-products.

Regression - Regression Trees

- In reality, local linear regression may be quite effective even when the relationships is nonlinear.
- This is used in Regression Trees.
- Each test instance is classified with its locally optimized linear regression by determining its appropriate partition.
- The partition is determined using split criteria in the internal nodes, i.e. the same as the Decision trees.
- The general strategy of tree construction is the same as for Decision Trees.
- The splits are univariate (single variable/axis parallel).
- The changes are done in splitting criterion determination and in the pruning.
- The number of points used for training need to be high to avoid over-fitting

Splitting criterion

- Due to numeric nature of the class variable, error-based measure have to be used instead of entropy or Gini index.
- The regression modeling is applied on each child resulting from potential split.
- The aggregated squared error of prediction of all training points is computed.
- The split point with the minimum aggregated error is selected.
- The complete regression modeling is computationally very expensive.
- An average variance of the numeric class variable may be used instead.
- The linear regression models are constructed at the leaf nodes after the tree is created.
- This results in larger trees but it its computational expensiveness is much lower.

Pruning criterion

- A portion of the training data is not used during construction phase.
- This set is used for evaluation of the squared error of the prediction.
- Leaf nodes are iteratively removed if the accuracy not decreases.

Regression - Assessing Model Effectiveness

- The effectiveness of the linear regression models can be evaluated with a measure known as R^2 -statistics or *coefficient of determination*.
- The standard Sum of Squared Error is defined for a model $g(\bar{X})$ as:

$$SSE = \sum_{i=1}^n (y_i - g(\bar{X}_i))^2$$

- The Squared Error of the response variable about its mean is defined as:

$$SST = \sum_{i=1}^n \left(y_i - \sum_{j=1}^n \frac{y_j}{n} \right)^2$$

- The R^2 -statistics is then defined as:

$$R^2 = 1 - \frac{SSE}{SST}$$

- The value is always between 0 and 1 and higher are more desirable.
- For high dimension data, adjusted version is more accurate:

$$R^2 = 1 - \frac{(n - d)SSE}{(n - 1)SST}$$

- The R^2 -statistics is not applicable on the nonlinear models.
- The nonlinear regression may be evaluated using pure SSE.

Questions?