

Machine Learning (DLMD SML01)

Nghia Duong-Trung¹

¹ German Research Center for Artificial Intelligence (DFKI GmbH)

Internal Use @ IU International University of Applied Science, Berlin Campus

SPONSORED BY THE



Federal Ministry
of Education
and Research

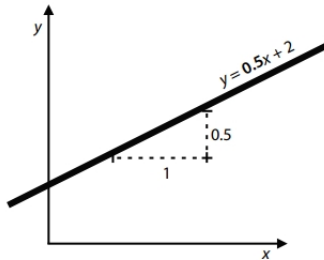
Section 3: Regression

- ▶ Textbook: check it in MyCampus
- ▶ Basic concepts
- ▶ Linear & Nonlinear & Logistic & Quantile regression
- ▶ Regularization
- ▶ Transfer task

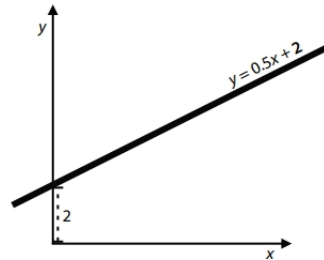
Basic concepts

The equation of a line

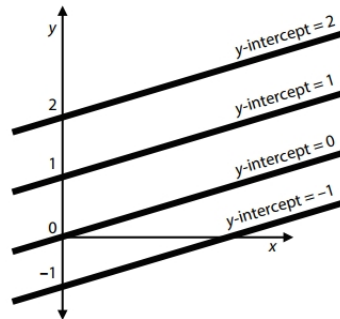
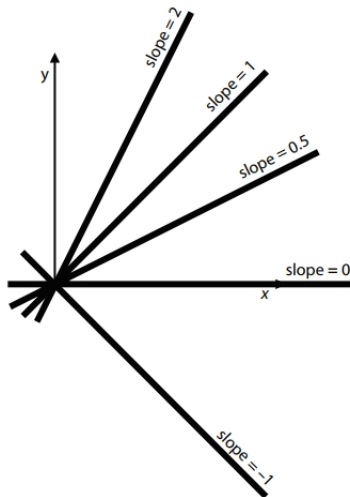
- ▶ Two components:
 - ▶ The slope
 - ▶ The y-intercept
- ▶ $y = 0.5x + 2$

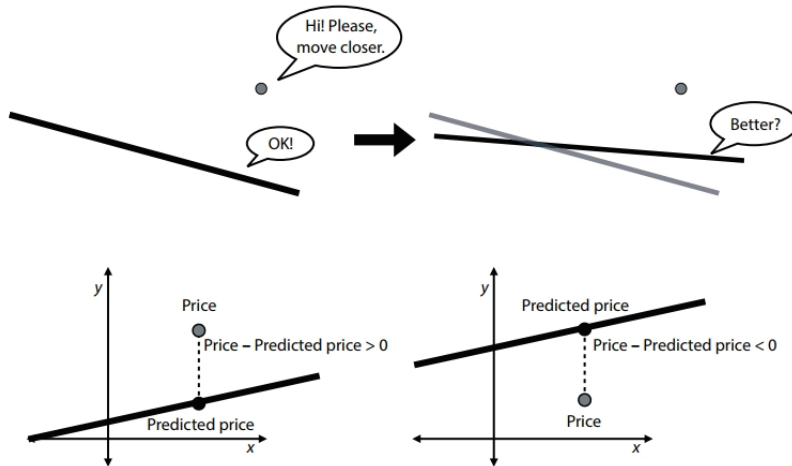


Slope = 0.5



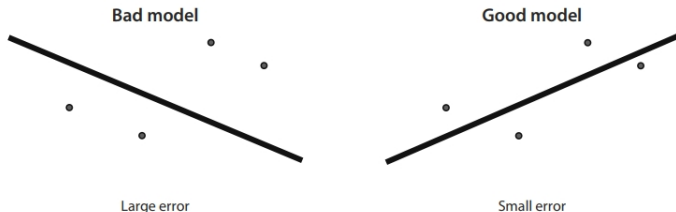
y-intercept = 2



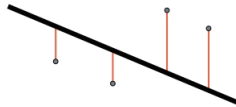


The error function

- ▶ The error function is a metric that tells how our model is doing
- ▶ Error function is also sometimes called loss function or cost function in the literature
- ▶ how do we define a good error function for linear regression models?
 - ▶ Absolute error is the sum of vertical distances from the line to the points in the dataset, and the
 - ▶ Square error is the sum of the squares of these distances

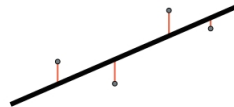


Large absolute error



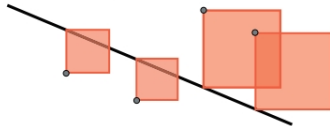
$$\text{Error} = \text{---} + \text{---} + \text{---} + \text{---}$$

Small absolute error



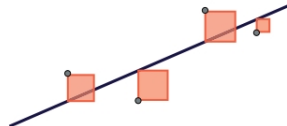
$$\text{Error} = \text{---} + \text{---} + \text{---} + \text{---}$$

Large square error



$$\text{Error} = \text{---} + \text{---} + \text{---} + \text{---}$$

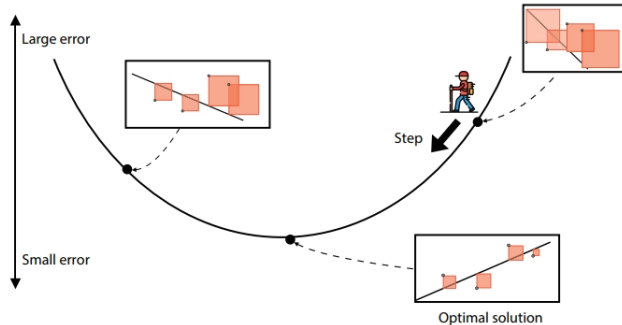
Small square error



$$\text{Error} = \text{---} + \text{---} + \text{---} + \text{---}$$

Gradient descent

- ▶ How to decrease an error function by slowly descending from a mountain
- ▶ This process uses derivatives to minimize the error function

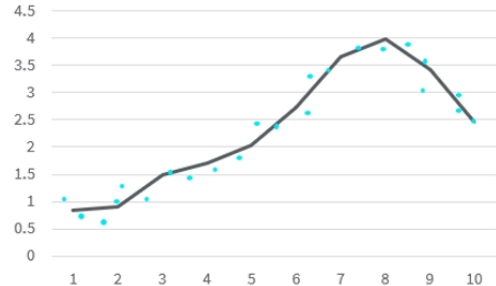
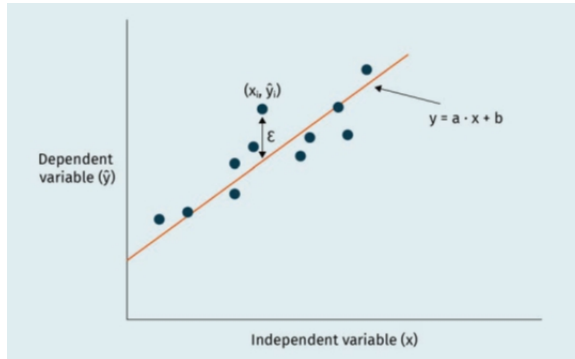


Linear & Nonlinear & Logistic & Quantile regression

Linear & Nonlinear regression

- ▶ Regression: This statistical method attempts to determine the strength and character of the relationship between dependent variable(s) and independent variable(s)
- ▶ In linear regression, the relationship is a straight-line equation:
$$y = a_1x_1 + a_2x_2 + \dots + a_nx_n + b$$
 - ▶ where $\{x_1, x_2, \dots, x_n\}$ are the independent variables of the dataset, and $\{a_1, a_2, \dots, a_n\}$ are the coefficients of these variables, implying the weight that each independent variable shares in the resultant target variable y . The b value is the constant or bias term
- ▶ in nonlinear regression, the relationship is a nonlinear equation (e.g., polynomial, exponential): $y = a_1x_1^2 + a_2x_2^3 + \dots + a_nx_n^{n+1} + b$

Linear & Nonlinear regression

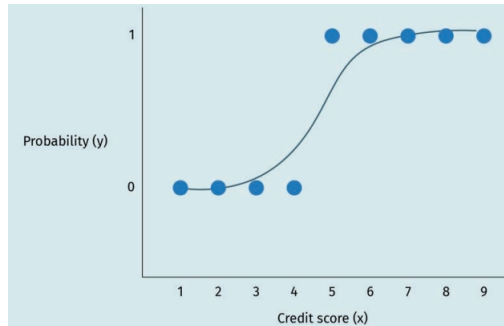


Simple linear regression

- ▶ $y = ax + b$
- ▶ $\hat{y} = y + \epsilon$
- ▶ For each data point, the difference between the predicted value y and the actual observation \hat{y} is the residue ϵ
- ▶ The evaluation metric for the simple regression model is the sum of the squared residues (i.e., sum of the squared errors E), which we aim to minimize
 - ▶ $E = \sum(\hat{y} - y)^2 = \sum(\hat{y} - (ax + b))^2$

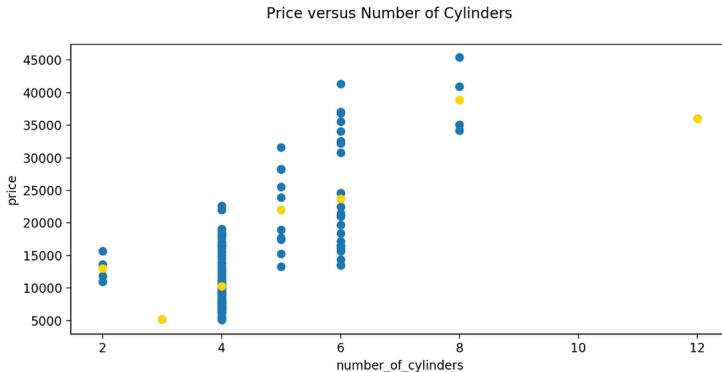
Logistic regression

- ▶ Logistic regression is considered an extension of the linear regression analysis, and its corresponding model can be used to classify input data records into a set of given categories or discrete values that form the dependent variable
- ▶ This statistical model uses a logistic function to model a binary dependent variable



Quantile regression

- In a regression model, one is normally interested in estimating the conditional mean of the response variable

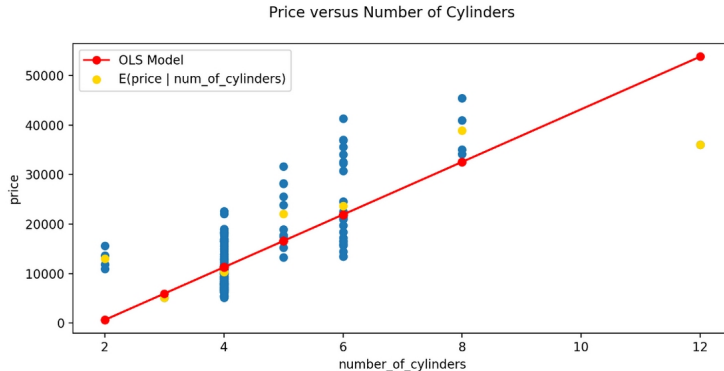


The gold dots represent the observed mean price conditioned upon the number of cylinders

Quantile regression

- If we wanted to estimate this conditional mean price using a regression model, we could employ the following linear model:

$$E(\text{price} | \text{num-of-cylinders}_i) = b + a_i * \text{num-of-cylinders}_i$$

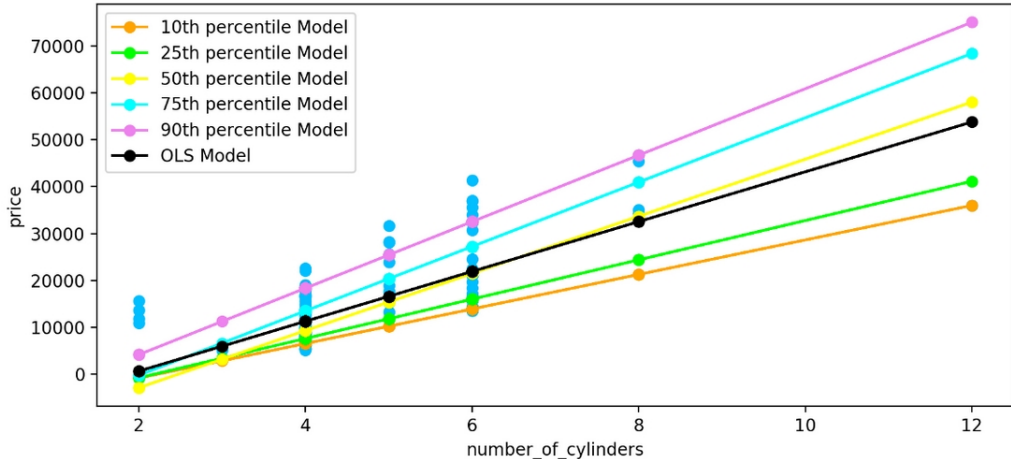


- ▶ In some data sets, the mean is not a good exemplar of the data.
 - ▶ The data are highly skewed to the left or to the right. Such is often the case in insurance claims or healthcare costs data where most claims are small valued but the data has a long tail of claims of increasing value
 - ▶ The data is multi-modal i.e. it has more than one highly frequently occurring value
 - ▶ The data contains influential outliers
- ▶ The mean does not adequately represent the nature of the data. One may be better served by estimating the conditional median

- ▶ A model for estimating the median price of automobiles
 $\text{Median}(\text{price} | \text{num-of-cylinders}_i) = b + a_i * \text{num-of-cylinders}_i$, where
 $\text{Probability}(\text{price} \leq b + a_i * \text{num-of-cylinders}_i) = 0.5$
- ▶ We may even go a step further. In a data set, the median is the 0.5 quantile (or 50th percentile) point meaning that 50% of the data points are less than the value of the median
- ▶ Similarly, there are other quantile points that can be defined
 - ▶ The 0.1 quantile point (10th percentile) is the value such that only 10% of the data set is smaller than this value

Quantile regression

Price versus Number of Cylinders



Regularization

- ▶ Regression analysis tries to develop the best fit between the independent variables and the dependent variable such that the loss function is minimized
- ▶ the developed model can **overfit** the training dataset and lose the value of generalization when applied to a different testing dataset
- ▶ To avoid the issues of overfitting and outliers and to have a more robust model, we penalize the loss function by adding a penalty term to the regression model

- ▶ Ridge Regression or L_2 L2regularization

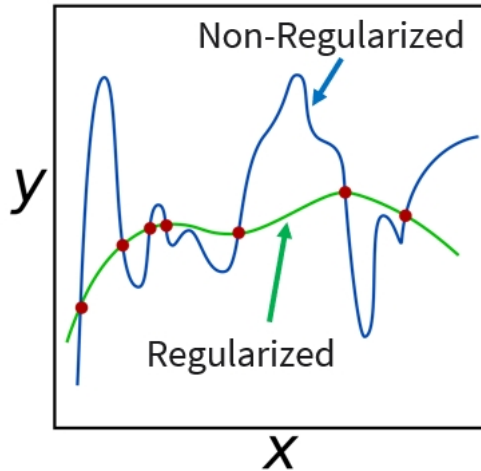
$$E = \sum (\hat{y} - y)^2 + \lambda \sum w^2 \quad (1)$$

- ▶ Lasso Regression or L_1 L2regularization

$$E = \sum (\hat{y} - y)^2 + \lambda \sum |w| \quad (2)$$

- ▶ Selecting the optimum value of the regularization parameter λ is a trade-off

Regularization



- ▶ `https://www.anaconda.com/products/distribution`
- ▶ `https://scikit-learn.org/stable/`
- ▶ `https://www.tutorialspoint.com/scikit_learn/index.htm`

Questions?