

Models and Patterns

Model

- High Level global description of a data set
- It takes a large sample perspective
 - Summarizing data in convenient, concise way
- Basic Models
 - Linear regression models
 - Mixture models
 - Markov models

Pattern

- Local Feature of the Data that holds for few records/variables
 - E.g., Mode or gap in pdf, Inflexion point in regression curve
- Departure from run of data
- Identify members with unusual properties
- Outliers in a database

Model Structures for Prediction

- Predict response variable from given values of others
- Response variable Y given p predictor variables X_1, \dots, X_p
- When Y is quantitative the task is known as regression
- When Y is categorical, it is known as *classification learning* or *supervised classification*

Regression with Linear Structure

- Response variable is a linear function of predictor variables

Model Constructed from data X
Hyperplane in p -dimensions

$$\hat{Y} = a_0 + \sum_{j=1}^p a_j X_j$$

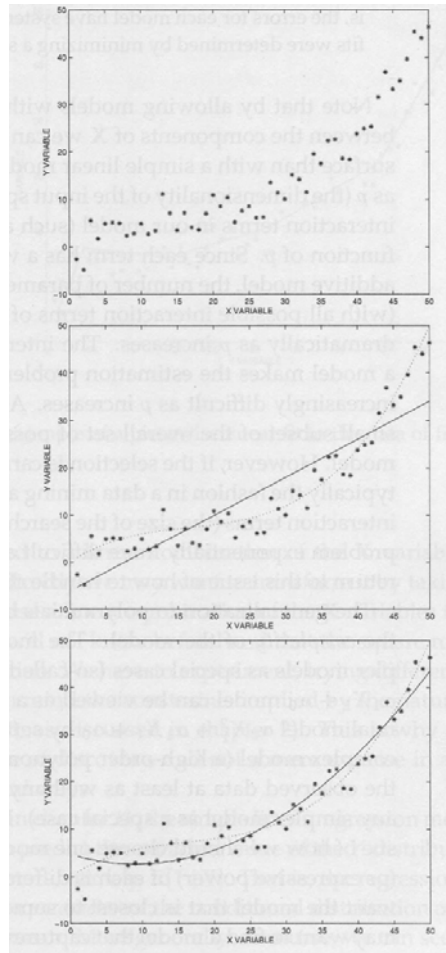
- Estimation of parameters a is straightforward
- Generalizing beyond linear functions

$$\hat{Y} = a_0 + \sum_{j=1}^p a_j f_j(X_j)$$

- Although nonlinear in variables, still linear in parameters

Regression Example

Model parameters estimated by minimizing Sum of Squared errors



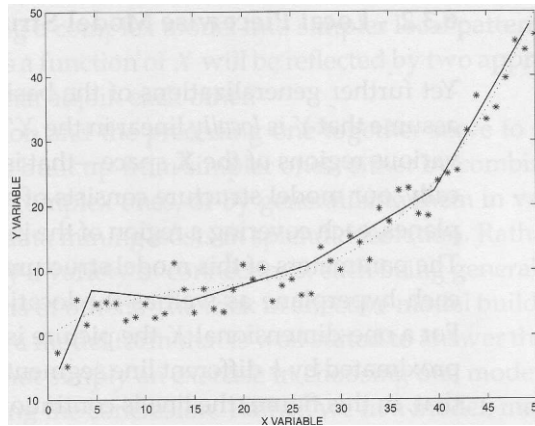
Fifty Data Points simulated from 3rd order polynomial equation
 $y = 0.001x^3 - 0.05x^2 + x + e$
 e is additive Gaussian noise with std dev 3 in range[1,50]

Fit of the model $aX+b$

Fit of the model aX^2+bX+c

Local Piecewise Model Structures for Regression

- Another generalization of basic linear model
- Assume Y is locally linear in the X s
- Curve is approximated by k linear segments



Linear Fit with $k = 5$

- If discontinuities are undesirable-- enforce continuity of various orders at end of segments
- Splines (each segment is a low degree quadratic or polynomial) ⁶

Nonparametric Local Models

Kernel Regression With Triangular Kernels

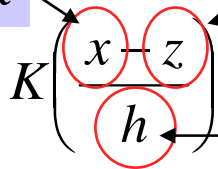
Retain data points.

Leave estimation of predicted value of Y until prediction is actually required

Weight data objects based on how similar they are to new object

Weight function that decays slowly with decreasing similarity will lead to a smooth estimate

Data set point

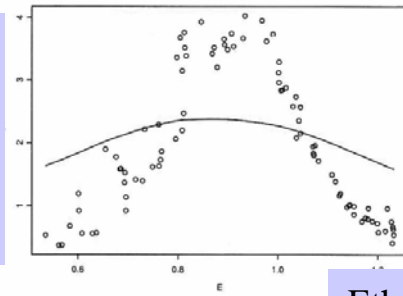


New point

Bandwidth, larger value leads to smoother estimate

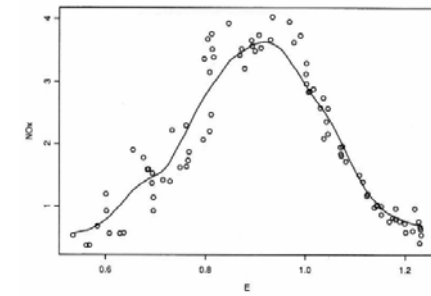
Related to nearest-neighbor methods

Nitrous Oxide
In emission

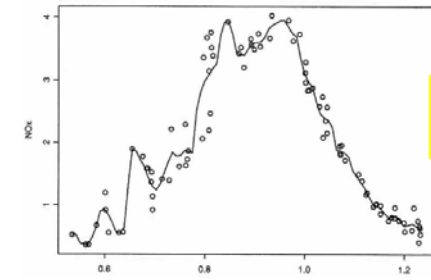


$h = 0.5$

Ethanol



$h = 0.1$



$h = 0.02$

Stochastic Components of Model Structures

- For any given vector of predictor variables more than one value of Y can be observed
- A distribution of values of y at each value of X
- Variables of X are insufficient
- It is a random component of the variation
- Regression model can be extended to include a stochastic component

$$y = g(x; \theta) + e$$

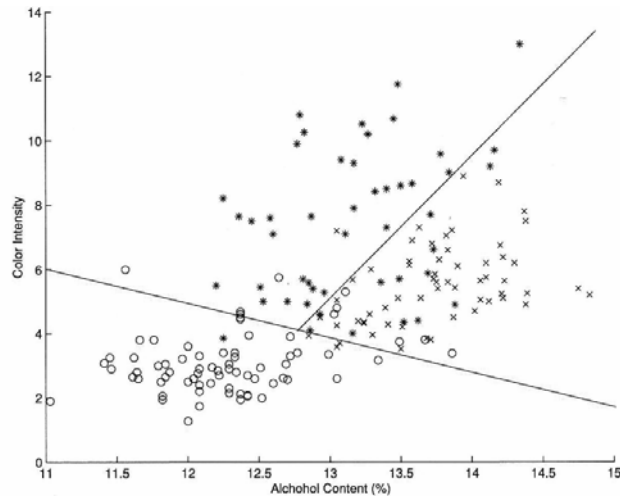
Parameters of
Model structure

Random variable with constant variance σ^2 and zero-mean

Predictive Models for Classification

Y is a categorical variable, taking a few possible values

Linear Decision Boundaries



Combine multiple
Simple models
Piecewise Linear Decision Boundaries

