

Introduction to Statistical Learning



Chapter 1: Introduction

stepwise, ridge, PCA, LASSO.
 ↑
 Nonlinear.

doubt: Where to Classify KNN?
 In chart.

doubt: KNN when to Classify?

linear ← Regression

logistic regression

Classification

linear discriminant analysis

Statistical learning

(tools to understand)

→ set of approaches for estimating f , minimize reducible error

Supervised

$X \rightarrow f(\cdot) \rightarrow Y$

S. b.

\hat{y} & Y are "close"

Unsupervised

Find relationships, patterns & structures in data.

Clustering

KMeans

hierarchical

PCA

natural output is available

Supervised Learning

Unsupervised Learning

natural output not available

X \ Y	Y	
	continuous/quantitative	categorical/qualitative
continuous	Regression (Wage)	Classification
categorical	?	?

doubt

clustering problem, (types of customers)

No single approach will perform well in all possible applications
 → model, intuition, assumption, trade-off

Notation:

n := sample size.

p := no. of features.

x_{ij} := value of j th variable for i th observation.

$\Rightarrow i = 1(1)n \quad j = 1(1)p$

choosing best method:

- ① cross validation
- ② bootstraps

$$X_{n \times p} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix}$$

$$x_i = \begin{bmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{ip} \end{bmatrix} \quad (\text{column vector})$$

data:

$\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$