

EXPLORATORY DATA ANALYSIS I

Roberta De Vito



BROWN
Public Health

MATRIX FORM OF THE DATA

MATRIX FORM OF THE DATA

- $(y_i, \underline{x}_i), \quad i = 1, \dots, n$

MATRIX FORM OF THE DATA

- $(y_i, \underline{x}_i), \quad i = 1, \dots, n$
- $\underline{x}_i = (x_{i1}, \dots, x_{ip})^\top$

MATRIX FORM OF THE DATA

- $(y_i, \underline{x}_i), \quad i = 1, \dots, n$
- $\underline{x}_i = (x_{i1}, \dots, x_{ip})^\top$
- the matrix form for the data

$$\begin{bmatrix} y_1 \\ \vdots \\ y_i \\ \vdots \\ y_n \end{bmatrix} \quad \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ \vdots & \vdots & \vdots & \vdots \\ x_{i1} & x_{i2} & \dots & x_{ip} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix}$$

MATRIX FORM OF THE DATA

Examples?

MATRIX FORM OF THE DATA

- teens example

$$\begin{array}{ccc} y_1 & \dots & y_n \\ teendep_1 & \dots & teendep_n \end{array}$$

- variables

$$X = \begin{pmatrix} yes & yes & \dots & none \\ \vdots & \vdots & \vdots & \vdots \\ no & yes & \dots & 3/week \\ \vdots & \vdots & \ddots & \vdots \\ yes & no & \dots & 6/week \end{pmatrix}$$

- $i = 1, \dots, n$

MATRIX FORM OF THE DATA

- teens example

$$\begin{array}{ccc} y_1 & \dots & y_n \\ teendep_1 & \dots & teendep_n \end{array}$$

- variables

$$X = \begin{pmatrix} \text{jail} & \text{act} & \dots & \text{smok} \\ \text{yes} & \text{yes} & \dots & \text{none} \\ \vdots & \vdots & \vdots & \vdots \\ \text{no} & \text{yes} & \dots & \text{3/week} \\ \vdots & \vdots & \ddots & \vdots \\ \text{yes} & \text{no} & \dots & \text{6/week} \end{pmatrix}$$

- $i = 1, \dots, n$

MATRIX FORM OF THE DATA

- There is no Y
- variables

$$X = \begin{pmatrix} yes & yes & \dots & none \\ \vdots & \vdots & \vdots & \vdots \\ no & yes & \dots & 3/week \\ \vdots & \vdots & \ddots & \vdots \\ yes & no & \dots & 6/week \end{pmatrix}$$

The variables can be continuous or categorical

SUPERVISED VS UNSUPERVISED

Supervised

- distinguish X and Y
- no matter if the relationship is linear or not linear
- example: depression

$$\begin{bmatrix} y_1 \\ \vdots \\ y_i \\ \vdots \\ y_n \end{bmatrix} \quad \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ \vdots & \vdots & \vdots & \vdots \\ x_{i1} & x_{i2} & \dots & x_{ip} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix}$$

SUPERVISED VS UNSUPERVISED

Unsupervised

- all the variables have the same role (No outcome)
- Do they correlate?

$$r = \frac{\sum (x_1 - \mu_1)(x_2 - \mu_2)}{\sqrt{\sum (x_1 - \mu_1)^2 \sum (x_2 - \mu_2)^2}}$$

- example: diet with nutrients

$$\begin{bmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ \vdots & \vdots & \vdots & \vdots \\ x_{i1} & x_{i2} & \dots & x_{ip} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix}$$

SUPERVISED VS UNSUPERVISED

	Unsupervised	Supervised
Continuous	<div>Clustering</div> <div>Dimension Reduction</div> <ul style="list-style-type: none">• SVD• PCA• K-mean	<div>Regression</div> <ul style="list-style-type: none">• Linear• Polynomial <div>Decision tree</div> <div>Random Forest</div>
Categorical	<div>Association Analysis</div> <div>Hidden Markov Model</div>	<div>Logistic Regression</div> <div>SUM</div>

EXPECTED VALUE

EXPECTED VALUE

- Discrete : $\sum_{r=1}^k x_r f_x(x_i; \theta) = \sum x_r p_r$

EXPECTED VALUE

- Discrete : $\sum_{r=1}^k x_r f_x(x_i; \theta) = \sum x_r p_r$
- Continuous $\int x f_x(x; \theta)$

EXPECTED VALUE

- Discrete : $\sum_{r=1}^k x_r f_x(x_i; \theta) = \sum x_r p_r$
- Continuous $\int x f_x(x; \theta)$
- Example: the toss

$$E[X] = 1 * \frac{1}{6} + 2 * \frac{1}{6} + 3 * \frac{1}{6} + 4 * \frac{1}{6} + 5 * \frac{1}{6} + 6 * \frac{1}{6} = 3.5$$

- Example with freezer (-18) and oven (100)

$$E[X] = \frac{100 - 18}{2} = 41$$

THE VARIANCE

THE VARIANCE

- $Var[X] = \sigma^2 = E[(X - \mu)^2]$

THE VARIANCE

- $Var[X] = \sigma^2 = E[(X - \mu)^2]$
- $Var[X] = E[X^2] - E[X]^2$

THE VARIANCE

- $Var[X] = \sigma^2 = E[(X - \mu)^2]$
- $Var[X] = E[X^2] - E[X]^2$
- Example with freezer (-18) and oven (100)

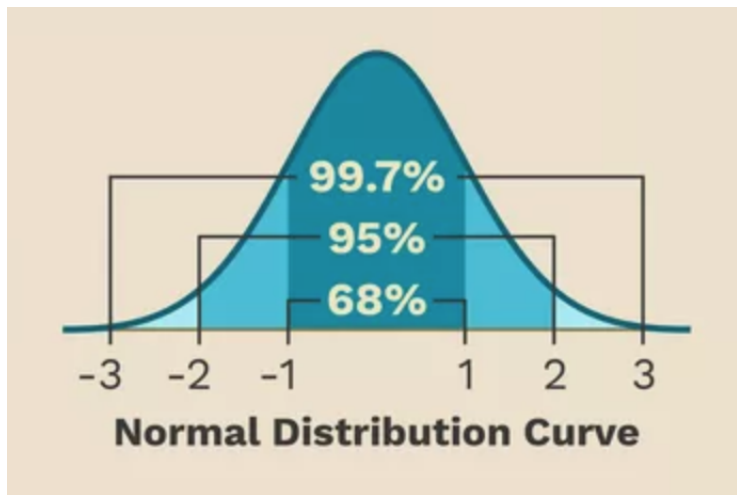
$$V[X] = 6962$$

Var function in R

STANDARD DEVIATION

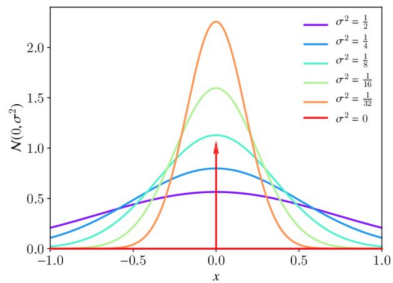
STANDARD DEVIATION

$$s_x = \sigma = \frac{\sum_{i=1}^n (x_i - \mu)^2}{n - 1}$$



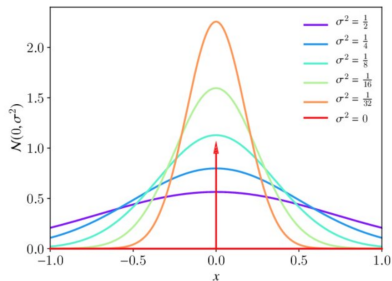
THE NORMAL DISTRIBUTION

$$X \sim N(\mu, \sigma^2)$$



THE NORMAL DISTRIBUTION

$$X \sim N(\mu, \sigma^2)$$



Central limit theorem: the average of many samples converges to a normal distribution as the number of samples increases

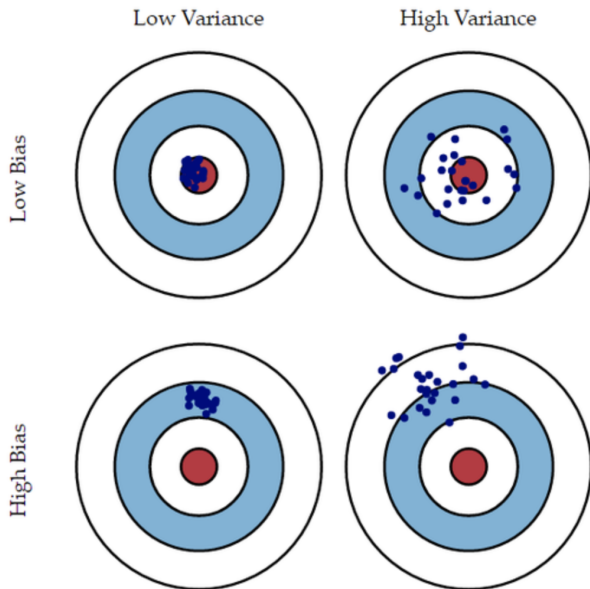
1. PREDICTION

For any estimate $\hat{f}(x)$ of $f(x)$, we have

$$Y - \hat{Y} = f(x) + \epsilon - \hat{f}(x) = \underbrace{(f(x) - \hat{f}(x))}_{\text{reducible}} + \underbrace{\epsilon}_{\text{irreducible}}$$

$$E[(Y - \hat{Y})^2 | X = x] = \underbrace{[f(x) - \hat{f}(x)]^2}_{\text{reducible}} + \underbrace{\text{Var}(\epsilon)}_{\text{irreducible}}$$

1. PREDICTION



2. INFERENCE

- What is inference?

2. INFERENCE

- What is inference?
- from the population to the sample

2. INFERENCE

- What is inference?
- from the population to the sample
- AIM: to arrive at the true
- $X_1, \dots, X_n \rightarrow x_1, \dots, x_n$
- estimate θ giving a function f
- example: the average of the students' age at Brown University

PROBLEM: we will never discover the true θ , we can discover the law f

INFERENCE PROBLEM

1. sample
2. function

Example: Bernoulli

$$\{x_1, x_2, x_3, \dots, x_n\} = \{1, 1, 1, \dots, 0\}, \quad 4C, 6H$$

the law

$$f(x_1, \dots, x_n) = \theta^{\sum x_i} (1 - \theta)^{n - \sum x_i} \quad 0 \leq \theta \leq 1$$

INFERENCEAL PROBLEM

1. Exact estimation

$$\frac{\sum x_i}{n} = \frac{4}{10}$$

2. confidence intervals

$$(0.4 - \delta_1; 0.4 + \delta_1)$$

3. Hypothesis testing

$$\begin{aligned} H_0 : \quad \theta &= \frac{1}{2} \\ H_1 : \quad \theta &= \frac{3}{4} \end{aligned}$$

CONSIDERATIONS

- When interested in prediction, want flexible model; interpretability less important
- When interested in inference, want interpretable (often less flexible) model

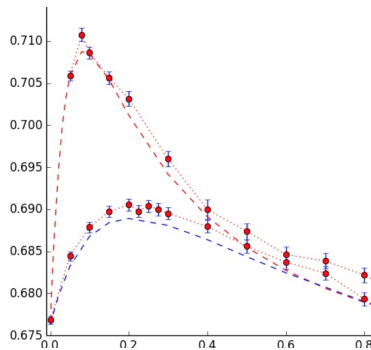
PARAMETRIC VS NONPARAMETRIC

Parametric

- assumption on f , then fit the model
- Pros: easy to estimate, simple
- Cons: far from the true f , our estimates can be poor, more parameters, overfitting

PARAMETRIC VS NONPARAMETRIC

Nonparametric



- no assumption on f
- problem of overfitting