# CMSC 471
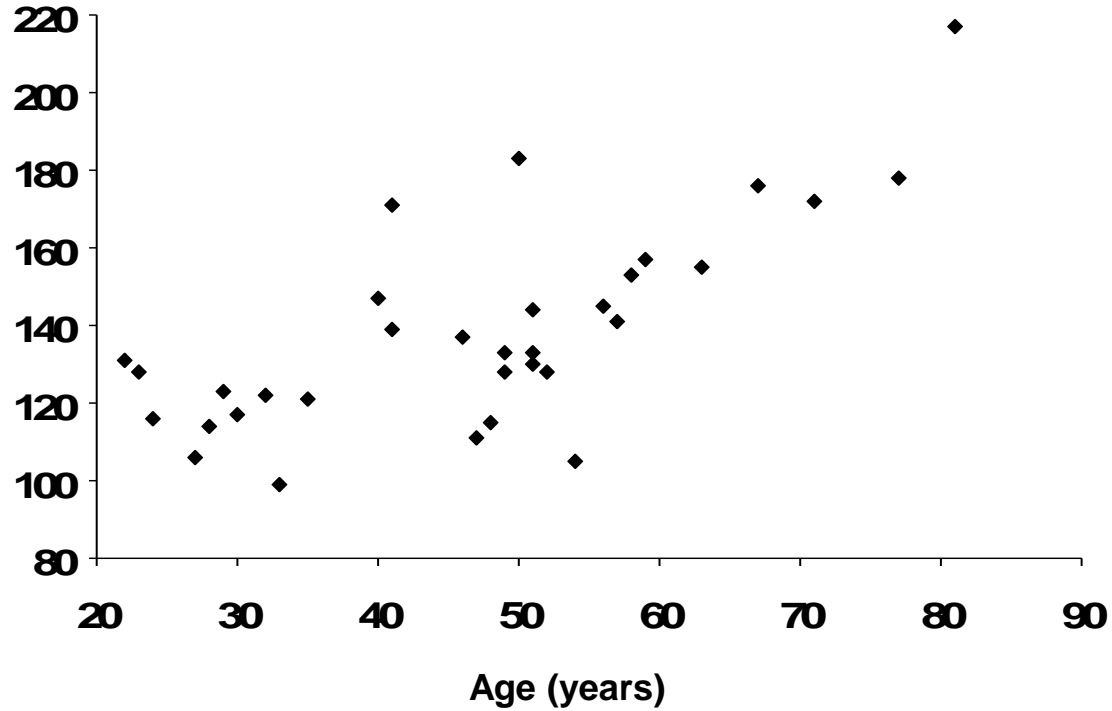
ML: Regression

# Simple linear regression

**Table 1    Age and systolic blood pressure (SBP) among 33 adult women**
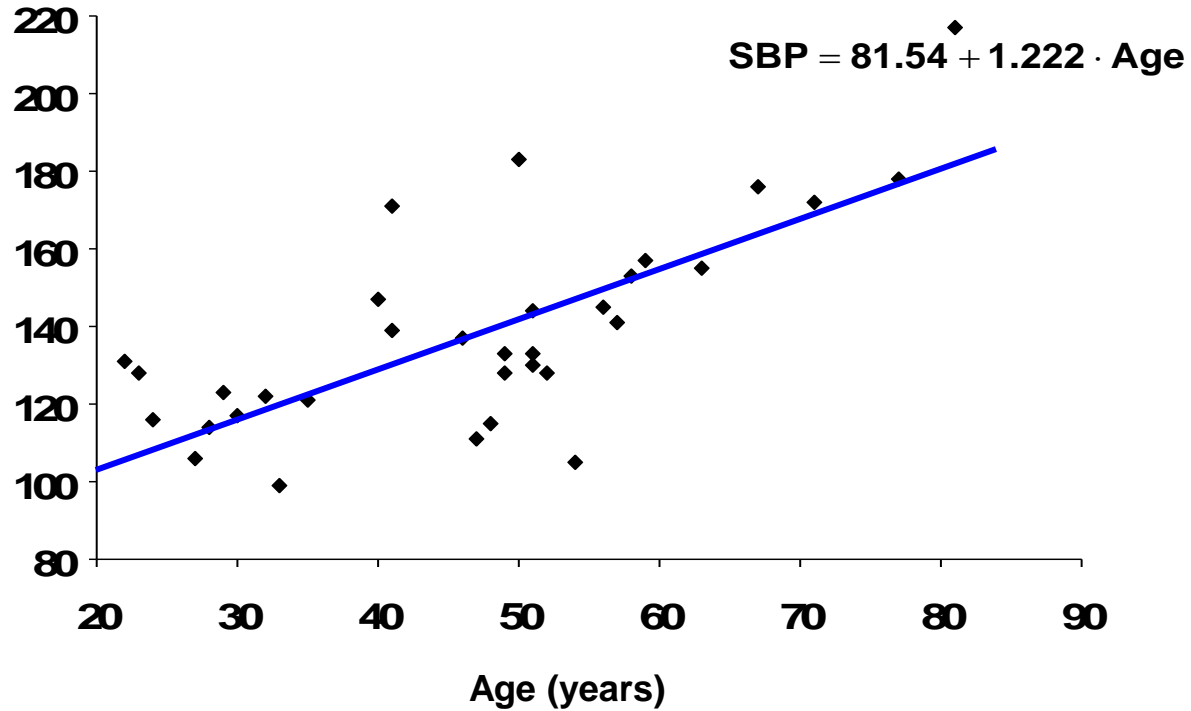
| Age | SBP | Age | SBP | Age | SBP |
|-----|-----|-----|-----|-----|-----|
| 22 | 131 | 41 | 139 | 52 | 128 |
| 23 | 128 | 41 | 171 | 54 | 105 |
| 24 | 116 | 46 | 137 | 56 | 145 |
| 27 | 106 | 47 | 111 | 57 | 141 |
| 28 | 114 | 48 | 115 | 58 | 153 |
| 29 | 123 | 49 | 133 | 59 | 157 |
| 30 | 117 | 49 | 128 | 63 | 155 |
| 32 | 122 | 50 | 183 | 67 | 176 |
| 33 | 99 | 51 | 130 | 71 | 172 |
| 35 | 121 | 51 | 133 | 77 | 178 |
| 40 | 147 | 51 | 144 | 81 | 217 |

adapted from Colton T. Statistics in Medicine. Boston: Little Brown, 1974

SBP (mm Hg)

$$SBP = 81.54 + 1.222 \cdot Age$$
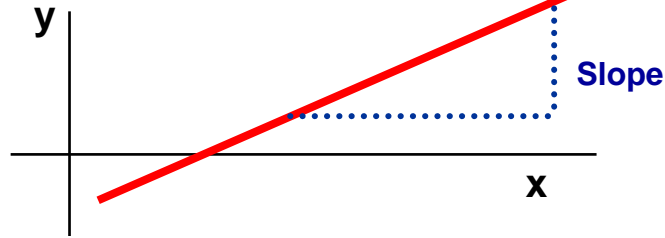
adapted from Colton T. Statistics in Medicine. Boston: Little Brown, 1974

# Simple linear regression

- Relation between 2 continuous variables (SBP and age)



$$y = \alpha + \beta_1 x_1$$

- Regression coefficient $\beta_1$
  - Measures association between y and x
  - Amount by which y changes on average when x changes by one unit
  - Least squares method

# Fitting equation to the data

- Finding the best curve = learning the best parameters
- Which is the best parameter?
  - One with minimum error
- Calculating error: Mean squared error

$$\mathrm{MSE} = \frac{1}{n} \sum_{i=1}^{n} \left( Y_i - \hat{Y}_i \right)^2.$$

actual

predicted

# Multiple linear regression

- Relation between a continuous variable and a set of i continuous variables

$$\mathbf{y} = \boldsymbol{\alpha} + \boldsymbol{\beta_1}\mathbf{x_1} + \boldsymbol{\beta_2}\mathbf{x_2} + \boldsymbol{...} + \boldsymbol{\beta_i}\mathbf{x_i}$$

- Partial regression coefficients $\beta_i$
  - Amount by which y changes on average when $x_i$ changes by one unit and all the other $x_i$s remain constant
  - Measures association between $x_i$ and y adjusted for all other $x_i$

- Example
  - SBP *versus* age, weight, height, etc

# Multiple linear regression

$$\mathbf{y} \qquad = \qquad \boldsymbol{\alpha} + \boldsymbol{\beta_1 x_1} + \boldsymbol{\beta_2 x_2} + \boldsymbol{...} + \boldsymbol{\beta_i x_i}$$

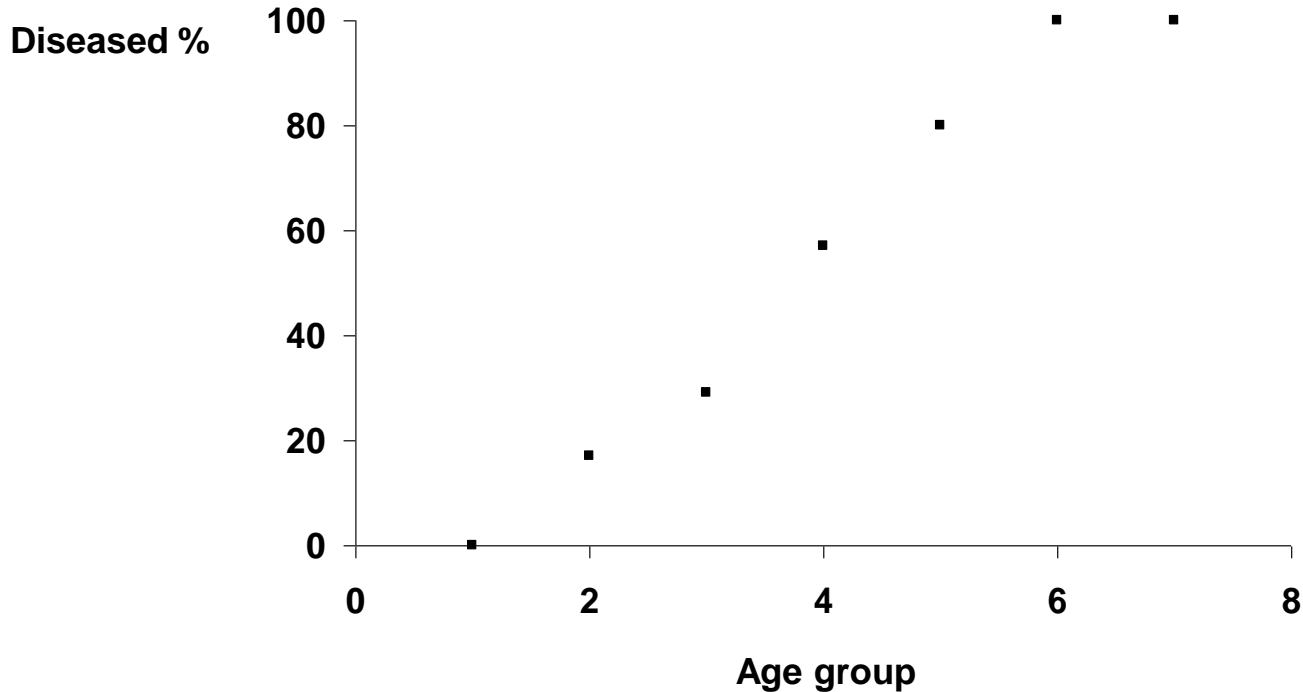| | |
|---|---|
| Predicted | Predictor variables |
| Response variable | Explanatory variables |
| Outcome variable | Covariables |
| Dependent | Independent variables |

# Logistic regression

Table 2    Prevalence (%) of signs of CD according to age group

| Age group | # in group | Diseased # | Diseased % |
|:---:|:---:|:---:|:---:|
| 20 - 29 | 5 | 0 | 0 |
| 30 - 39 | 6 | 1 | 17 |
| 40 - 49 | 7 | 2 | 29 |
| 50 - 59 | 7 | 4 | 57 |
| 60 - 69 | 5 | 4 | 80 |
| 70 - 79 | 2 | 2 | 100 |
| 80 - 89 | 1 | 1 | 100 |

# Dot-plot: Data from Table 2

# Logistic function (1)

**Probability of disease**

# Logistic function (1)

**Probability of disease**



$$P(y|x) = \frac{e^{\alpha + \beta x}}{1 + e^{\alpha + \beta x}}$$

$x$

# Fitting equation to the data

- Can not use MSE
- Estimate [Maximum likelihood](#)
- Likelihood function
  - Estimates parameters $\alpha$ and $\beta$
  - Practically easier to work with log-likelihood

$$L(\mathrm{B}) = \ln\left[l(\mathrm{B})\right] = \sum_{i=1}^{n}\left\{y_i \ln\left[\pi(x_i)\right] + (1 - y_i)\ln\left[1 - \pi(x_i)\right]\right\}$$