

# **ACIT4830 – Special Robotics and Control Subject**

## **Topic5 – Regression methods**

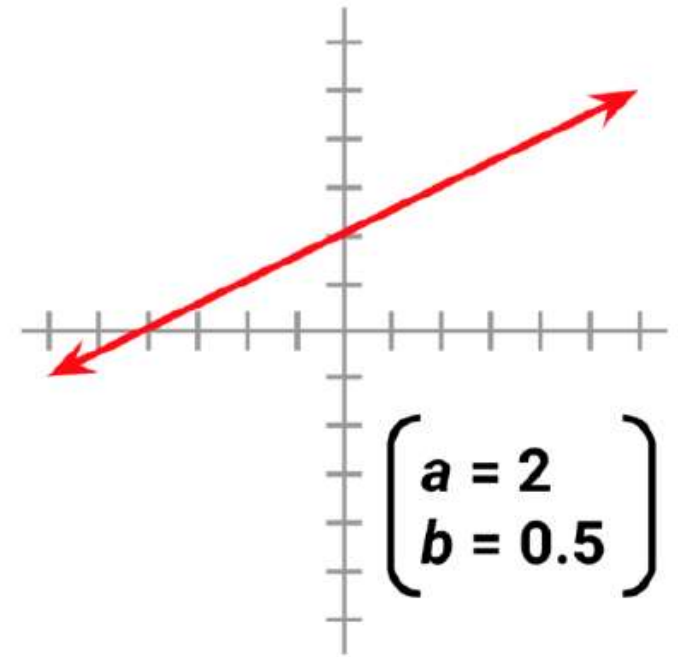
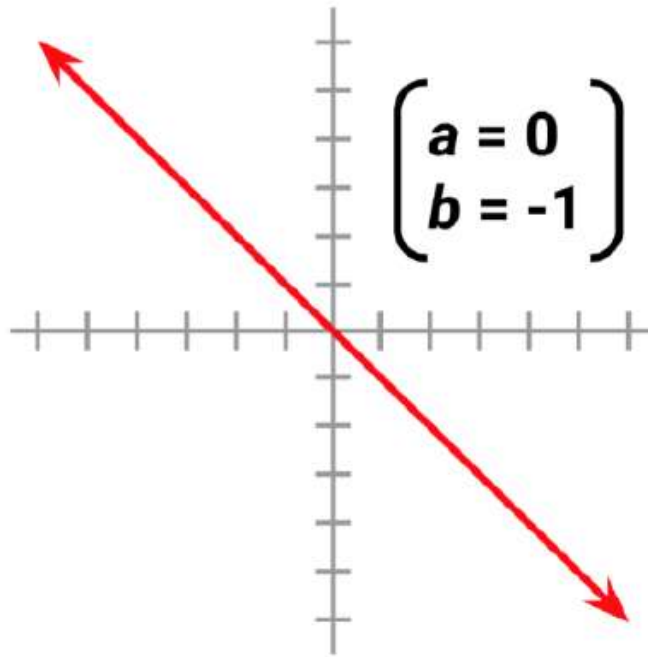
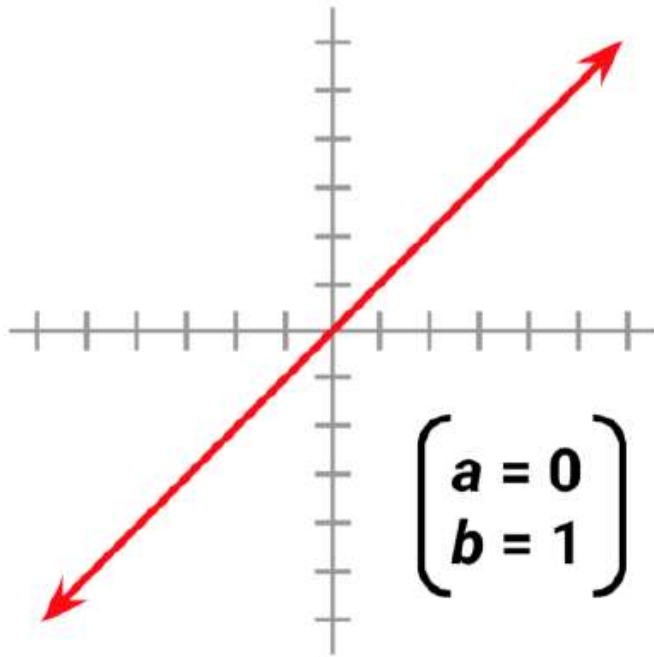
**Evi Zouganeli**  
**OsloMet – Oslo Metropolitan University**  
**([evizou@oslomet.no](mailto:evizou@oslomet.no))**

# Content

(Chapter 6 – Regression methods)

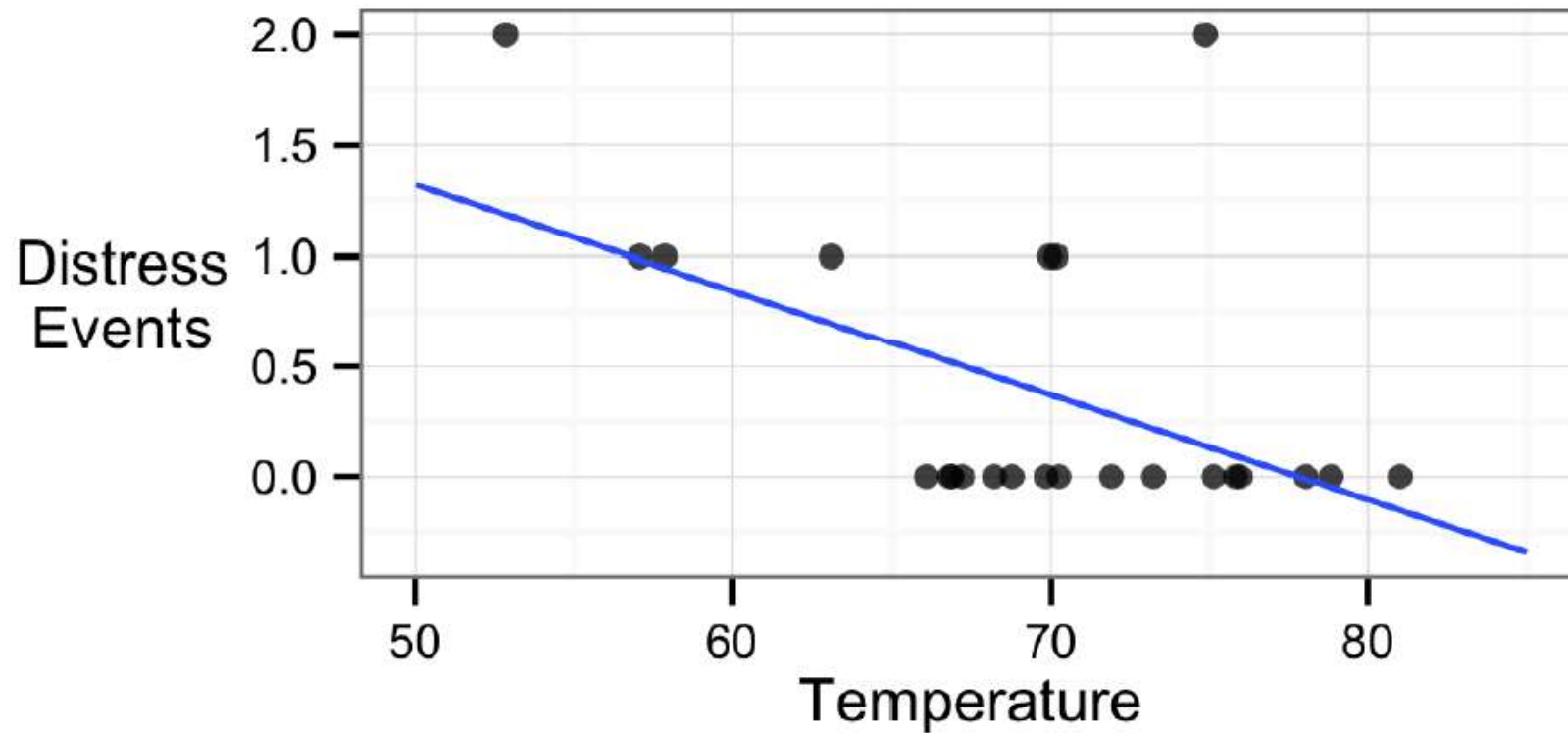
- Linear and Multiple regression
- Non-linear regression
- Regression trees
- Model trees

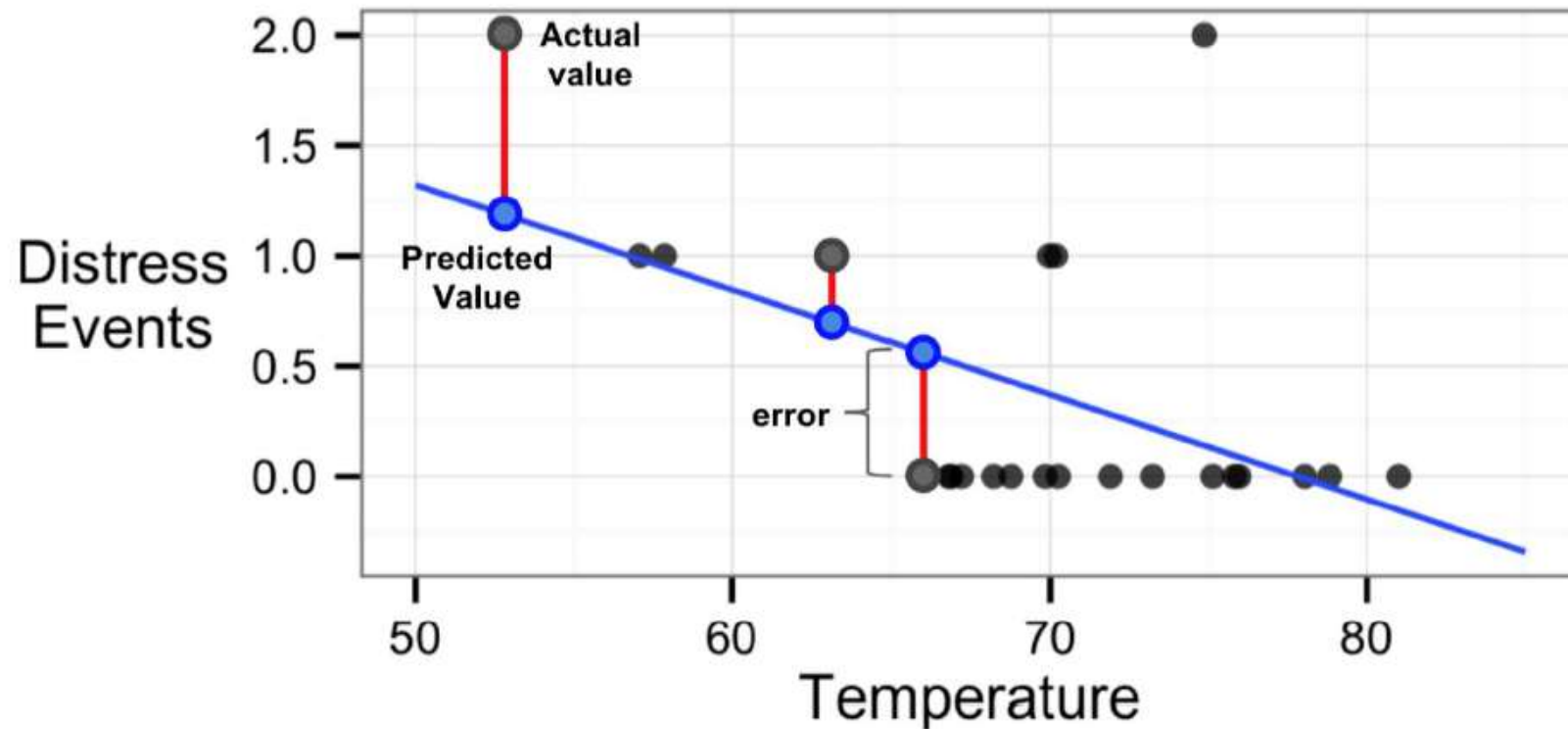
# Linear Regression



$$y = \alpha + \beta x$$

# Linear Regression





Minimize error:  
sum least squares

$$\sum (y_i - \hat{y}_i)^2 = \sum e_i^2$$

# Linear Regression

$$y = \alpha + \beta x$$

$$a = \bar{y} - b\bar{x}$$

$$b = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$b = \frac{\text{Cov}(x, y)}{\text{Var}(x)}$$

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

variance

$$\text{Cov}(X, Y) = \frac{\sum (X_i - \bar{X})(Y_j - \bar{Y})}{n}$$

covariance

Pearson's correlation  
(-1 to +1)

$$\rho_{x,y} = \text{Corr}(x, y) = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y}$$

# Multiple Regression

Strengths	Weaknesses
<ul style="list-style-type: none"><li>• By far the most common approach for modeling numeric data</li><li>• Can be adapted to model almost any data</li><li>• Provides estimates of the strength and size of the relationships among features and the outcome</li></ul>	<ul style="list-style-type: none"><li>• Makes strong assumptions about the data</li><li>• The model's form must be specified by the user in advance</li><li>• Does not do well with missing data</li><li>• Only works with numeric features, so categorical data require extra processing</li><li>• Requires some knowledge of statistics to understand the model</li></ul>

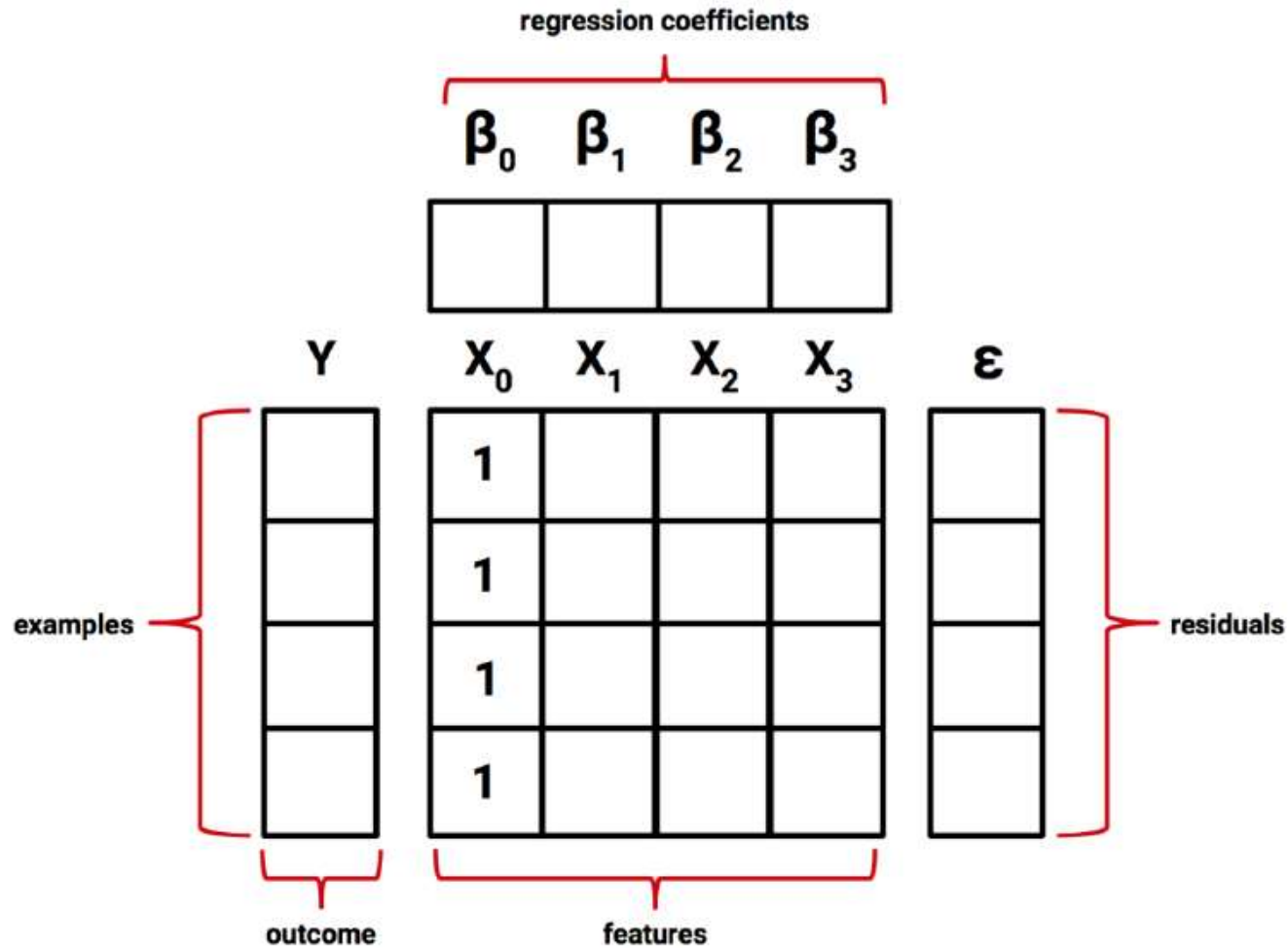


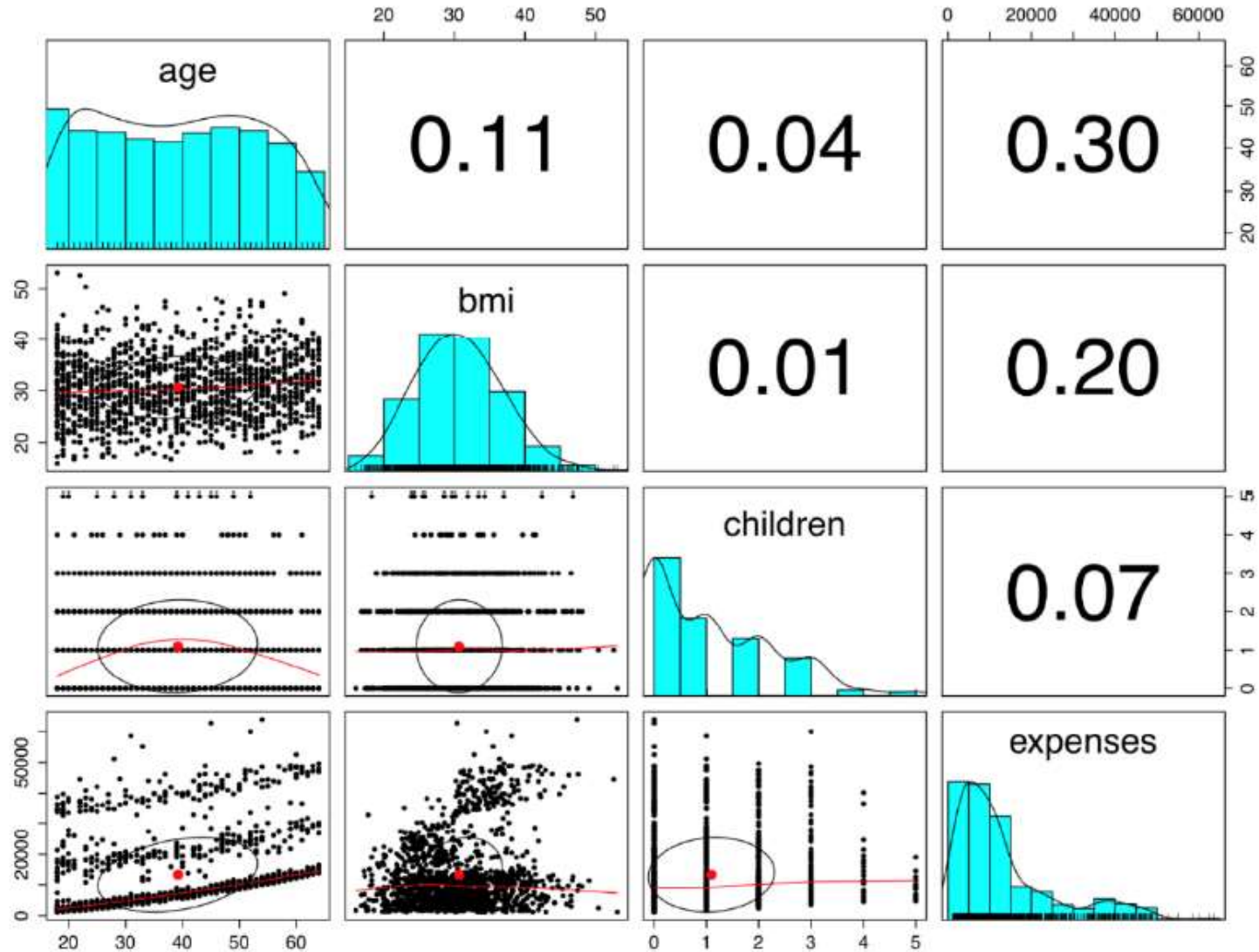
# Multiple Regression

$$y = \alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_i x_i + \varepsilon$$

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_i x_i + \varepsilon$$

# Multiple Regression





# Non-linear Regression

$$y = \alpha + \beta_1 x + \beta_2 x^2$$

# Model adjustment

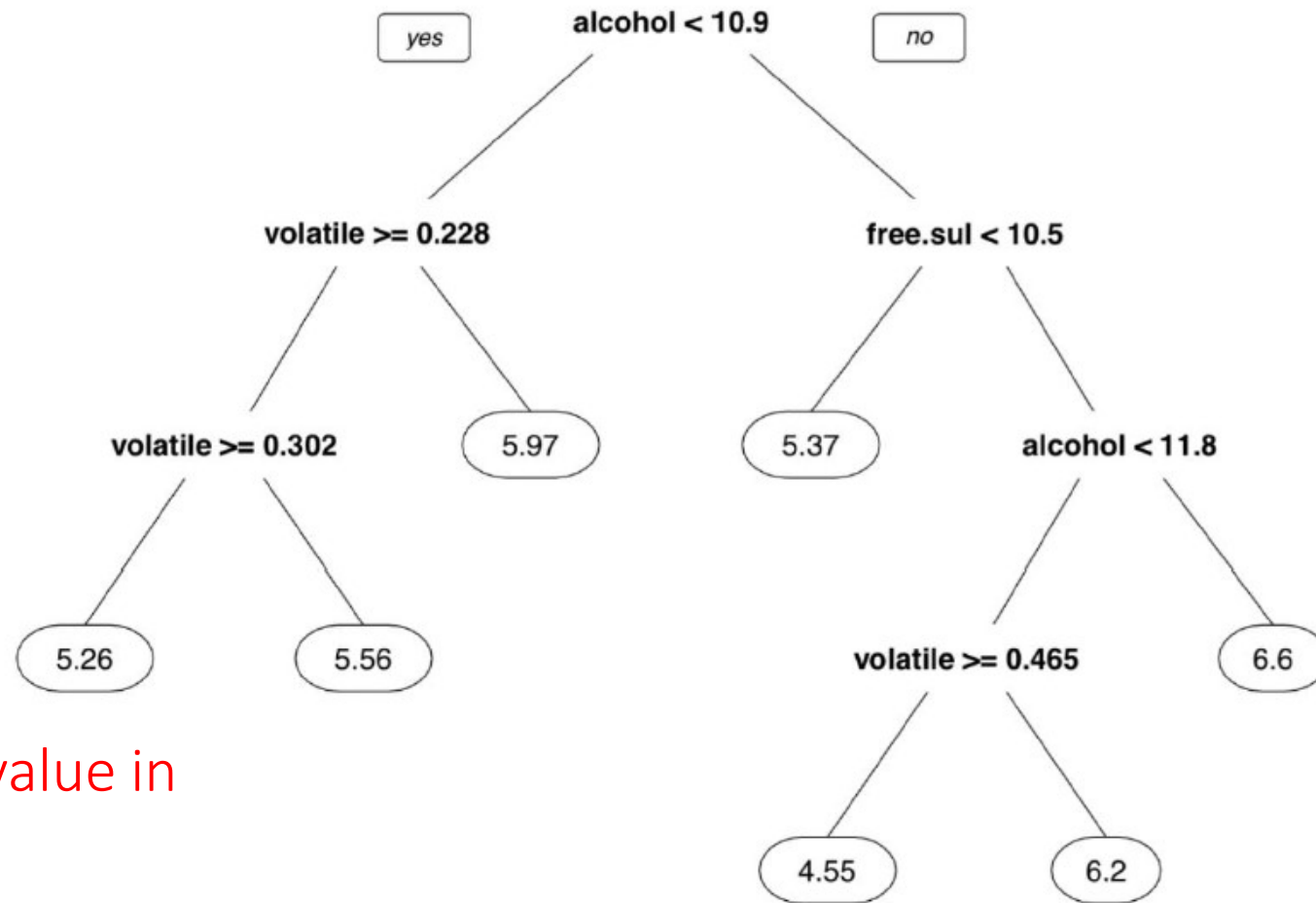
– medical expences example in hands-on

- Added a non-linear term for age
- Created an indicator for obesity
- Specified an interaction between obesity and smoking

# Adding Regression to Trees

Strengths	Weaknesses
<ul style="list-style-type: none"><li>• Combines the strengths of decision trees with the ability to model numeric data</li><li>• Does automatic feature selection, which allows the approach to be used with a very large number of features</li><li>• Does not require the user to specify the model in advance</li><li>• May fit some types of data much better than linear regression</li><li>• Does not require knowledge of statistics to interpret the model</li></ul>	<ul style="list-style-type: none"><li>• Not as commonly-used as linear regression</li><li>• Requires a large amount of training data</li><li>• Difficult to determine the overall net effect of individual features on the outcome</li><li>• May be more difficult to interpret than a regression model</li></ul>

# Regression Tree



mean value in  
subset

# Regression Trees

original data	1	1	1	2	2	3	4	5	5	6	6	7	7	7	7
split on feature A	1	1	1	2	2	3	4	5	5	6	6	7	7	7	7
split on feature B	1	1	1	2	2	3	4	5	5	6	6	7	7	7	7
	$T_1$							$T_2$							

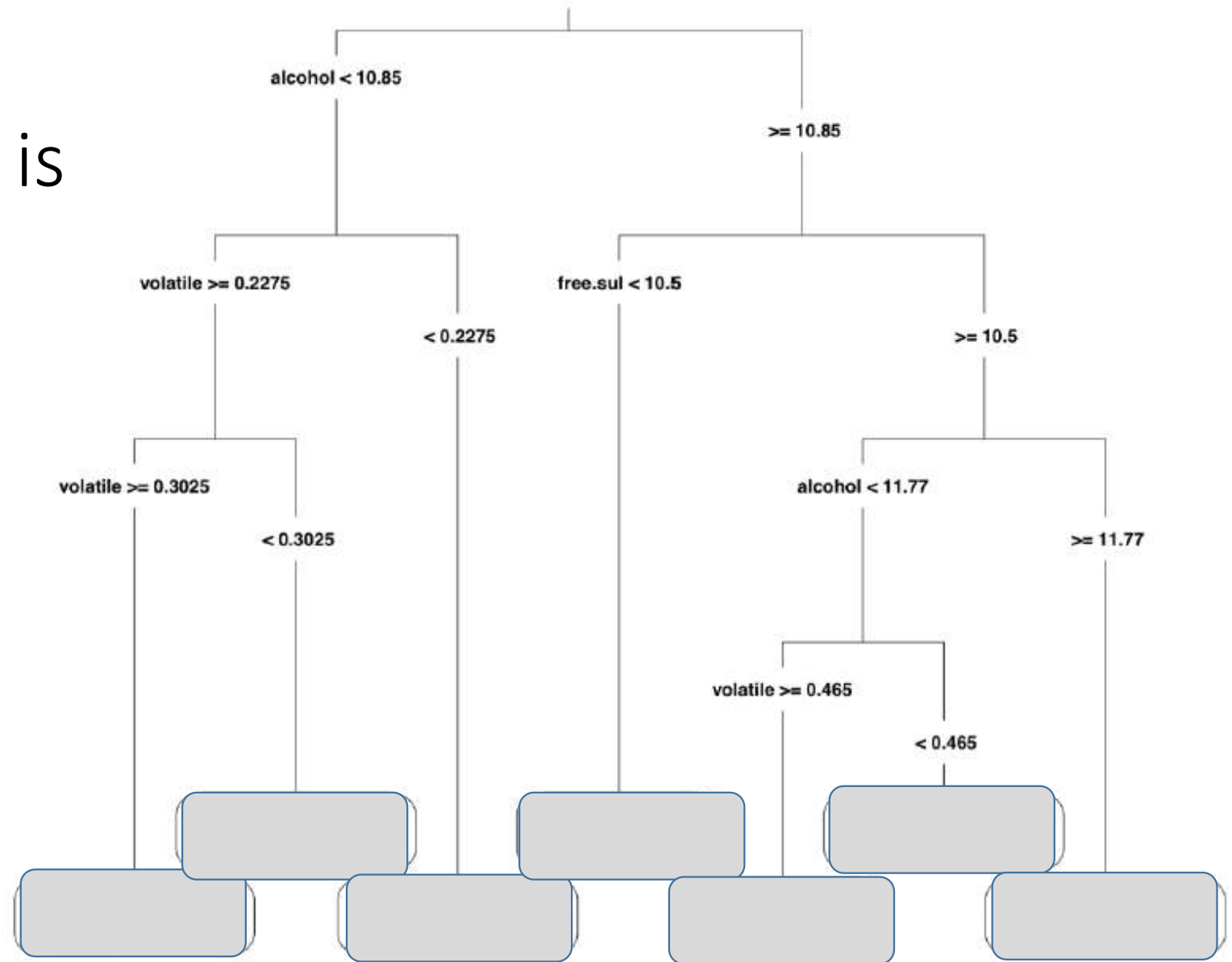
Standard Deviation Reduction

$$SDR = sd(T) - \sum_i \frac{|T_i|}{|T|} \times sd(T_i)$$

Choose the feature and split that **maximizes SDR**



Model Tree:  
every branch is  
a separate  
regression  
model



# Performance evaluation

Mean Absolute Error

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |e_i|$$