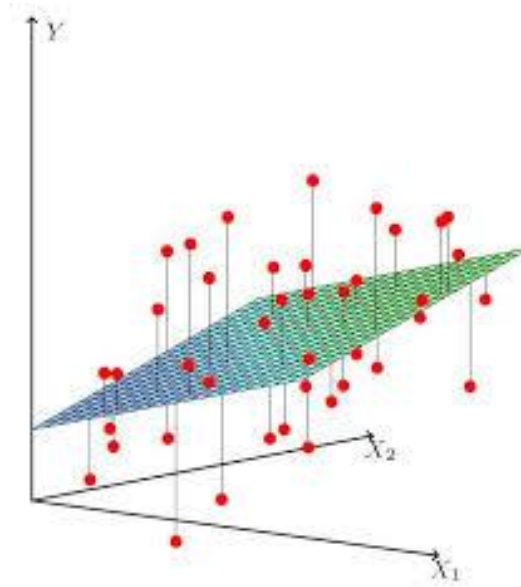




Part 11 – Linear Regression in Machine Learning

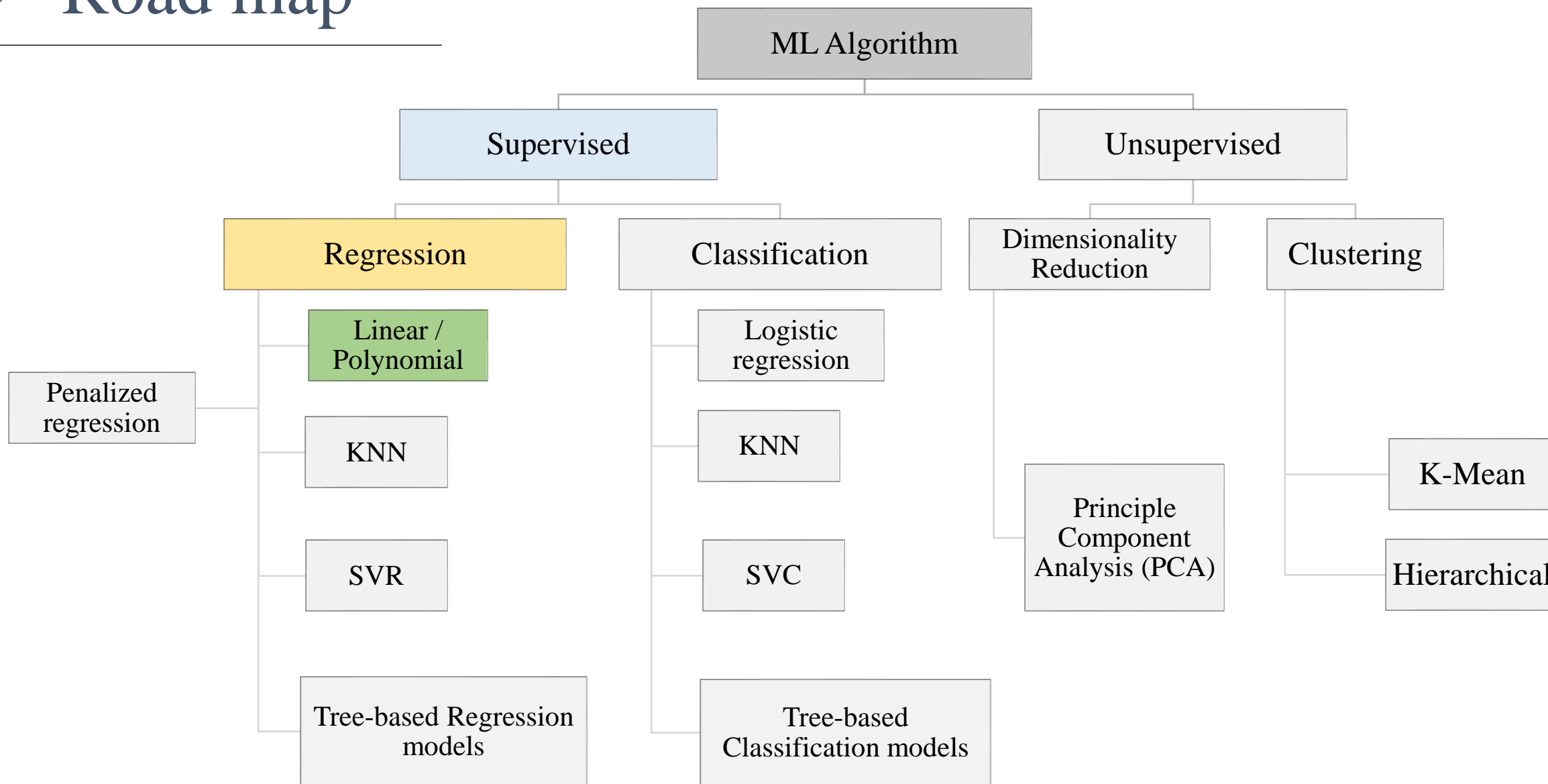


Prof. Pedram Jahangiry



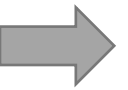


Road map



Part I

Linear regression: Machine Learning approach



Linear Models

The linear model is an important example of a **parametric** model

- We have a collection of labeled examples $\{(X_i, y_i)\}_{i=1}^N$, where
 - N is the size of the collection
 - X_i is the **D-dimensional** feature vector
 - y_i is a real-valued target

$$f_{w,b}(X) = \mathbf{W}X + \mathbf{b}$$

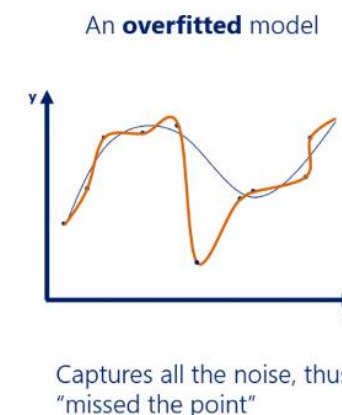
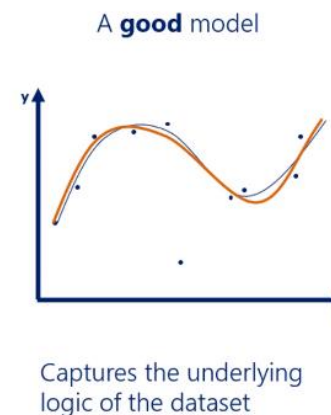
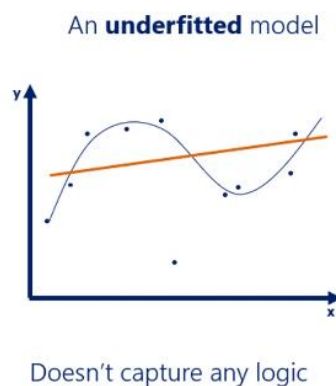
The model $f_{w,b}$ is a linear combination of features and parameterized by **W** and **b**

- **W** is a D-dimensional vector of parameters
- **b** is a real number

→ Linear Models (cont'd)

$$f_{w,b}(X) = \textcolor{red}{W}X + \textcolor{green}{b}$$

- The model is specified with $D+1$ parameter.
- We estimate the parameters (W^*, b^*) by fitting the model to training data.
- Although it is almost never correct, a linear model often serves as a **simple** and **interpretable** approximation the unknown true $f(X)$.
- It may seem overly simplistic, but linear regression is extremely useful both conceptually and practically.
- Linear regression models rarely overfit.



→ The optimization problem

The optimization problem is defined as:

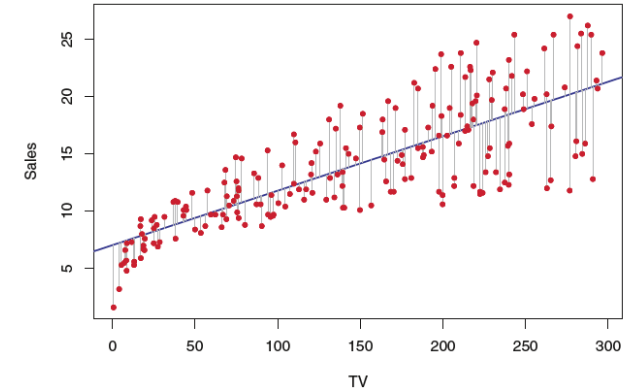
$$\text{Min}_{\mathbf{w}, \mathbf{b}} \text{MSE} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2 = \frac{1}{N} \sum_{i=1}^N \left(y_i - f_{\mathbf{w}, \mathbf{b}}(X_i) \right)^2$$

$\left(y_i - f_{\mathbf{w}, \mathbf{b}}(X_i) \right)^2$ is called the **objective function**, or the **loss function**! Or the **squared error loss**.

Why quadratic loss function? Why not using absolute value or cube?

1. More convenient (well-behaved derivative).
2. There exists a closed form solution.

The solution to this optimization problem is \mathbf{w}^* and \mathbf{b}^* . Now we can make predictions!



Linear Regression Evaluation Metrics

$$R^2 = 1 - \frac{\sum (y_i - \hat{y})^2}{\sum (y_i - \bar{y})^2} = 1 - \frac{SS_{residuals}}{SS_{total}}$$

$$MAE = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}|$$

$$MSE = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y})^2$$

$$RMSE = \sqrt{MSE} = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y})^2}$$

Adjusted

versions

$$\text{Adjusted } R^2 = 1 - (1 - R^2) * \frac{n - 1}{n - k - 1}$$

$$AIC = \frac{2K}{N} - \frac{2 \ln(\hat{L})}{N}$$

$$BIC = \ln(N) K - 2 \ln(\hat{L})$$

- **AIC**: Akaike information criterion
- **BIC**: Bayesian information criterion
- K : number of estimated parameters
- \hat{L} : Maximum value of the likelihood function

Part II

Linear regression: **Econometrics** approach
(this part is optional)

GAUSS MARKOV ASSUMPTIONS FOR REGRESSION

THE GAUSS-MARKOV ASSUMPTIONS

The following is a summary of the five Gauss-Markov assumptions that we used in this chapter. Remember, the first four were used to establish unbiasedness of OLS, whereas the fifth was added to derive the usual variance formulas and to conclude that OLS is best linear unbiased.

Assumption MLR.1 (Linear in Parameters)

The model in the population can be written as

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + u,$$

where $\beta_0, \beta_1, \dots, \beta_k$ are the unknown parameters (constants) of interest and u is an unobserved random error or disturbance term.

Assumption MLR.2 (Random Sampling)

We have a random sample of n observations, $\{(x_{i1}, x_{i2}, \dots, x_{ik}, y_i) : i = 1, 2, \dots, n\}$, following the population model in Assumption MLR.1.

Assumption MLR.3 (No Perfect Collinearity)

In the sample (and therefore in the population), none of the independent variables is constant, and there are no *exact linear* relationships among the independent variables.

Assumption MLR.4 (Zero Conditional Mean)

The error u has an expected value of zero given any values of the independent variables. In other words,

$$E(u|x_1, x_2, \dots, x_k) = 0.$$

Assumption MLR.5 (Homoskedasticity)

The error u has the same variance given any value of the explanatory variables. In other words,

$$\text{Var}(u|x_1, \dots, x_k) = \sigma^2.$$

Standard assumptions for the multiple regression model

Assumption MLR.1

Linear in Parameters

The model in the population can be written as

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + u, \quad [3.31]$$

where $\beta_0, \beta_1, \dots, \beta_k$ are the unknown parameters (constants) of interest and u is an unobserved random error or disturbance term.

Assumption MLR.2

Random Sampling

We have a random sample of n observations, $\{(x_{i1}, x_{i2}, \dots, x_{ik}, y_i) : i = 1, 2, \dots, n\}$, following the population model in Assumption MLR.1.

Standard assumptions for the multiple regression model

Assumption MLR.3

No Perfect Collinearity

In the sample (and therefore in the population), none of the independent variables is constant, and there are *no exact linear relationships* among the independent variables.

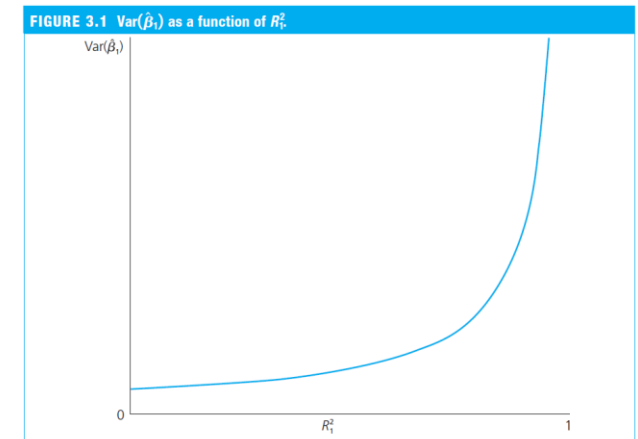
1. The assumption only rules out **perfect collinearity/correlation** between explanatory variables; **imperfect correlation is allowed**
2. If an explanatory variable is a perfect linear combination of other explanatory variables it is superfluous and may be eliminated
3. MLR.3 fails if $n < k + 1$. Intuitively, this makes sense: to estimate $k + 1$ parameters, we need at least $k + 1$ observations.

Detecting multicollinearity

Multicollinearity may be detected through **Variance Inflation Factors**:

$$VIF_j = 1/(1 - R_j^2)$$

$$\text{Var}(\hat{\beta}_j) = \frac{\sigma^2}{\text{SST}_j(1 - R_j^2)} \longrightarrow \text{Var}(\hat{\beta}_j) = \frac{\sigma^2}{\text{SST}_j} \cdot VIF_j$$



As an arbitrary rule of thumb, the variance inflation factor should not be larger than 10

Standard assumptions for the multiple regression model (cont.)

Assumption MLR.4

Zero Conditional Mean

The error u has an expected value of zero given any values of the independent variables. In other words,

$$E(u|x_1, x_2, \dots, x_k) = 0. \quad [3.36]$$

- The value of the **explanatory variables** must contain no information about the mean of the unobserved factors
- In a multiple regression model, the **zero conditional mean assumption** is much **more likely to hold** because fewer things end up in the error.

Theorem 3.1 (Unbiasedness of OLS)

THEOREM 3.1

UNBIASEDNESS OF OLS

Under Assumptions MLR.1 through MLR.4,

$$E(\hat{\beta}_j) = \beta_j, j = 0, 1, \dots, k, \quad [3.37]$$

for any values of the population parameter β_j . In other words, the OLS estimators are unbiased estimators of the population parameters.



Unbiasedness is an **average property in repeated samples**;
In a given sample, the estimates may still be far away from the true values!

Standard assumptions for the multiple regression model (cont.)

Assumption MLR.5

Homoskedasticity

The error u has the same variance given any value of the explanatory variables. In other words,
 $\text{Var}(u|x_1, \dots, x_k) = \sigma^2$.

- The value of the explanatory variables must contain no information about the **variance** of the unobserved factors
- Example: Wage equation

$$\text{Var}(u_i | \text{educ}_i, \text{exper}_i, \text{tenure}_i) = \sigma^2$$

← This assumption may also be hard to justify in many cases

THEOREM 3.2

SAMPLING VARIANCES OF THE OLS SLOPE ESTIMATORS

Under Assumptions MLR.1 through MLR.5, conditional on the sample values of the independent variables,

$$\text{Var}(\hat{\beta}_j) = \frac{\sigma^2}{\text{SST}_j(1 - R_j^2)} \quad [3.51]$$

for $j = 1, 2, \dots, k$, where $\text{SST}_j = \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2$ is the total sample variation in x_j , and R_j^2 is the R -squared from regressing x_j on all other independent variables (and including an intercept).

The sampling variability of the estimated regression coefficients depends on 4 things:

1. Variability of the unobserved factors (σ^2)
2. Variation in the explanatory variable $\text{var}(X_j)$ or SST_j
3. Number of observations n
4. Linear relationships among the independent variables (R^2)

Testing for Heteroskedasticity

There are many tests for heteroskedasticity; two popular:

☐ Breusch-Pagan test

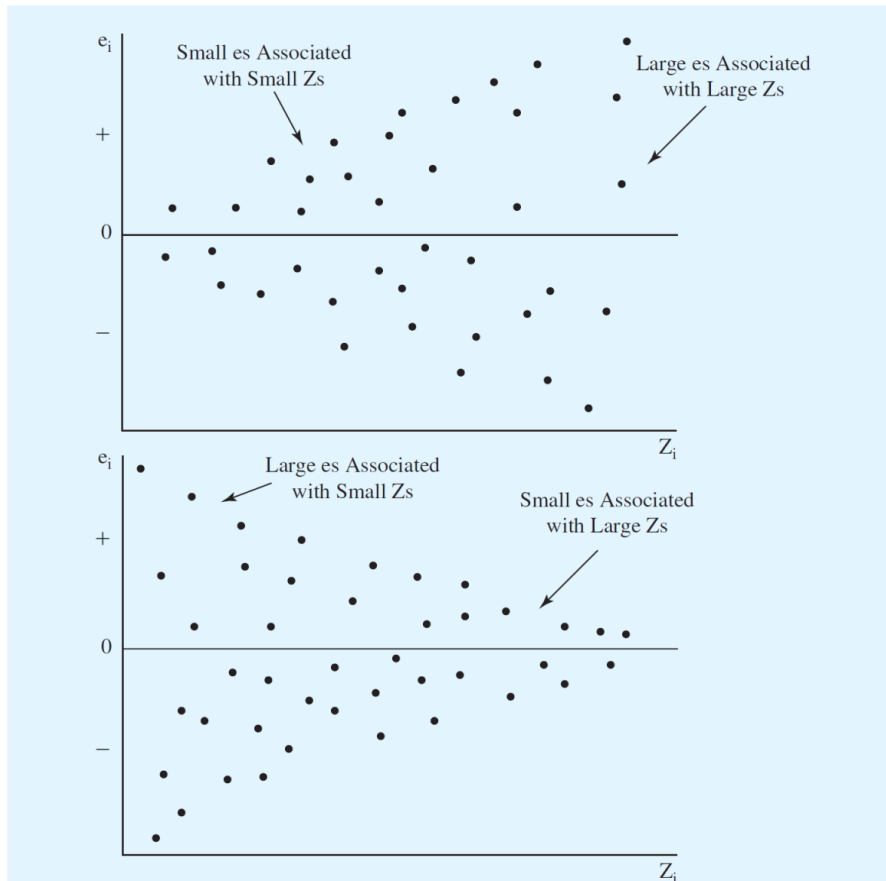
☐ White test

Before testing for heteroskedasticity, start with asking:

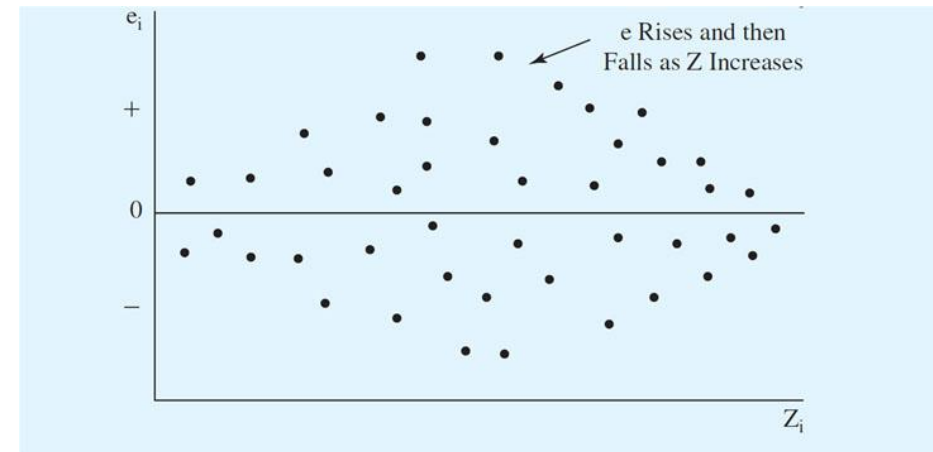
1. Are there any obvious specification errors?
2. Are there any early warning signs of heteroskedasticity?
3. Does a graph of the residuals show any evidence of heteroskedasticity?

Testing for Heteroskedasticity (cont'd)

Eyeballing Residuals for Possible Heteroskedasticity



If you plot the residuals of an equation with respect to a potential explanatory variable Z , a **pattern** in the residuals is an indication of possible **heteroskedasticity**.



The Breusch-Pagan Test for Heteroskedasticity:

Steps:

1. Estimate the model $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + u$ by OLS, as usual. Obtain the squared OLS residuals \hat{u}
2. Run the regression in $\hat{u}^2 = \delta_0 + \delta_1 x_1 + \delta_2 x_2 + \dots + \delta_k x_k + \text{error}$ Keep the R-squared from this regression $R_{\hat{u}^2}^2$
3. Form either the **F statistic** or the **LM statistic** and compute the p -value. If the p -value is sufficiently small, that is, below the chosen significance level, then we reject the **null hypothesis of homoskedasticity**.

$H_0 : Var(u|x_1, x_2, \dots, x_k) = Var(u|x) = \sigma^2$ \longrightarrow $H_0 : \delta_1 = \delta_2 = \dots = \delta_k = 0$ Regress squared residuals on all explanatory variables and test whether this regression has explanatory power.

$$F = \frac{R_{\hat{u}^2}^2 / k}{1 - R_{\hat{u}^2}^2 / (n - k - 1)}$$

$$LM = n \cdot R_{\hat{u}^2}^2 \sim \chi_k^2$$

A large **F statistic** or a large **Lagrange multiplier** statistic, (LM) lead to rejection of the null hypothesis.

The White Test for Heteroskedasticity

Steps:

1. Estimate the model $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + u$ by OLS, as usual. Obtain the squared OLS residuals \hat{u}
2. Run the regression in $\hat{u}^2 = \delta_0 + \delta_1 x_1 + \delta_2 x_2 + \delta_3 x_3 + \delta_4 x_1^2 + \delta_5 x_2^2 + \delta_6 x_3^2 + \delta_7 x_1 x_2 + \delta_8 x_1 x_3 + \delta_9 x_2 x_3 + \text{error}.$ Keep the R-squared from this regression $R_{\hat{u}^2}^2$
3. Form either the **F statistic** or the **LM statistic** and compute the p -value. If the p -value is sufficiently small, that is, below the chosen significance level, then we reject the **null hypothesis of homoskedasticity**.

$H_0 : Var(u|x_1, x_2, \dots, x_k) = Var(u|x) = \sigma^2$ \longrightarrow $H_0 : \delta_1 = \delta_2 = \dots = \delta_9 = 0$ Regress squared residuals on all explanatory variables, **their squares**, and **interactions** (here: example for $k=3$)

$$F = \frac{R_{\hat{u}^2}^2 / k}{1 - R_{\hat{u}^2}^2 / (n - k - 1)}$$

$$LM = n \cdot R_{\hat{u}^2}^2 \sim \chi_k^2$$

A large **F statistic** or a large **Lagrange multiplier** statistic, (LM) lead to rejection of the null hypothesis.

Remedies for Heteroskedasticity

- ❑ If heteroskedasticity is found, the first thing to do is examine the equation carefully for specification errors.
- ❑ If there are no obvious specification errors, the heteroskedasticity is probably pure in nature and one of the following remedies should be considered.
 1. Redefining the Variables
 2. Heteroskedasticity-Corrected Standard Errors
 3. Weighted Least Square Estimation!

