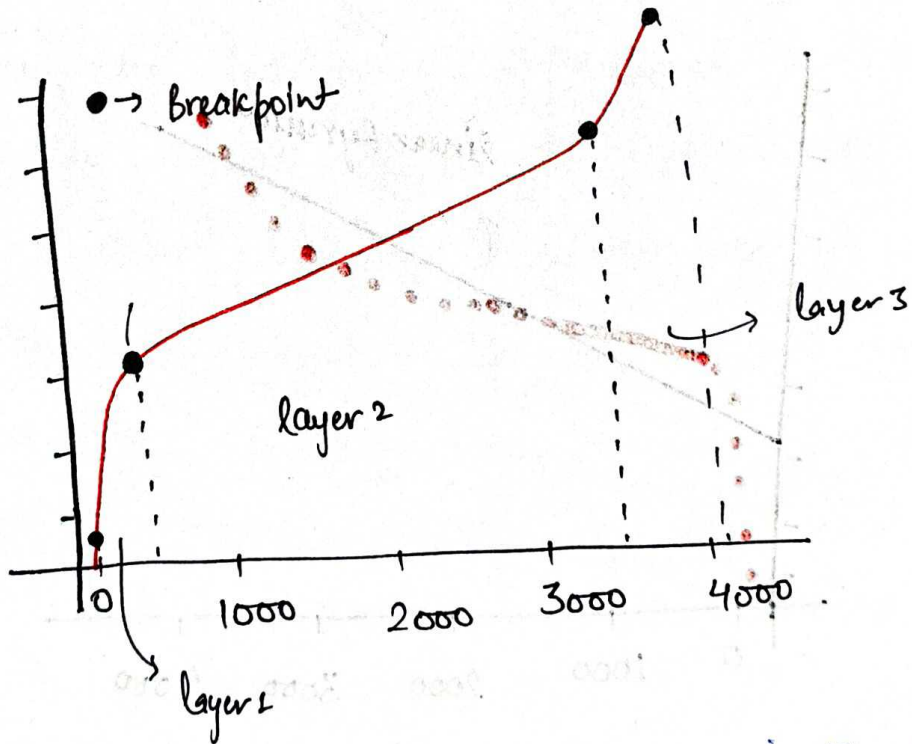
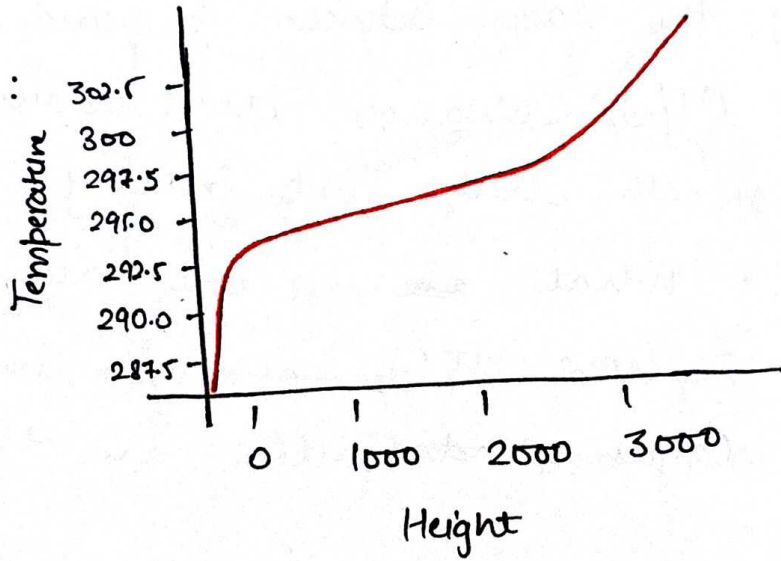


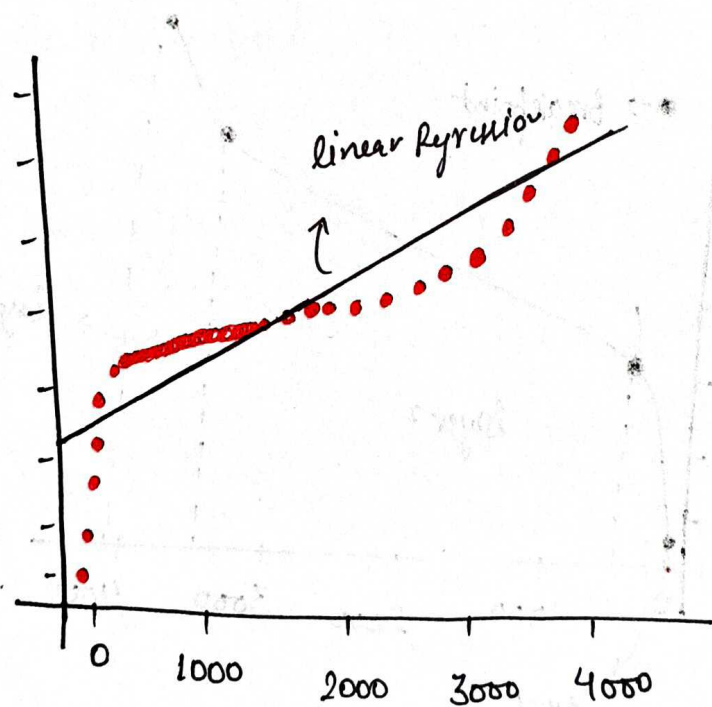
Piece wise Regression

Problem:



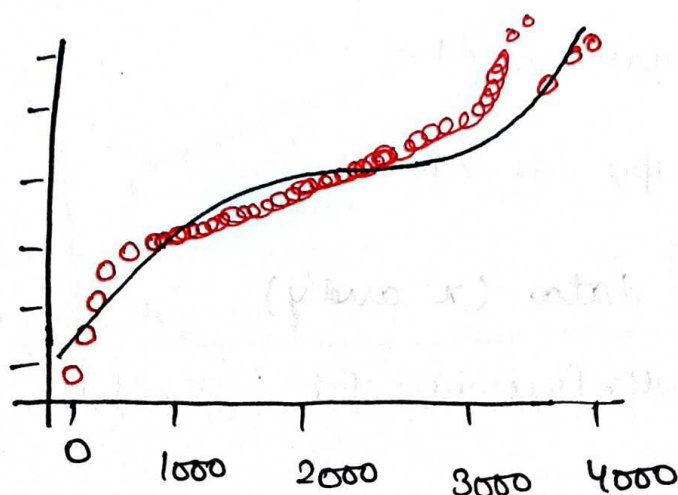
In fact, we can visually divide it into 3 different layer (or piece). Above we can see the presentation of these layers and its approximate height (x-vals). The break points are the points where the data

completely changes its behaviour. For instance, break point 1 (let's denote it b_1) is our start point; the layer between b_1 and b_2 has a similar $\Delta T / \Delta z$ changes its behaviour to a more smooth change not increasing exponentially like before. What ~~we~~ we are trying to understand is how $\Delta T / \Delta z$ changes inside of each these layers and where are the ~~best~~ break points.



if we plot our data and try to fit a regression on it we fail miserably. However, it doesn't mean that a linear regression isn't the best choice for our problem,

Polynomial Regression



The problem of polynomial regression is that you lose the interpretability (human can ~~under~~ understand the cause of a decision.) of the model when adding the polynomial terms (quadratic, cubic, etc). So if you don't care so much about interpretability, you can stop reading here. However if you need interpretability to deeply understand the problem, piecewise linear regression is your buddy.

Piecewise Linear Regression

$$y(x) = \begin{cases} \eta_1 + \beta_1 (x - b_1), & b_1 < x \leq b_2 \\ \eta_2 + \beta_2 (x - b_2), & b_2 < x \leq b_3 \\ \dots \\ \eta_n + \beta_n (x - b_{n-1}), & b_{n-1} < x \leq b_n \end{cases}$$

Code:

```
import pwlfit
```

```
import pandas as pd
```

```
import numpy as np
```

```
# fit your data (x and y)
```

```
myPWLfit = pwlfit.Piecewiselinefit(x, y) → df[column]
```

```
# fit the data for n line segments
```

```
z = myPWLfit.fit(3)
```

```
# Calculate slopes
```

```
slopes = myPWLfit.cal_slopes()
```

```
# predict for the determined points
```

```
xHat = df[column]
```

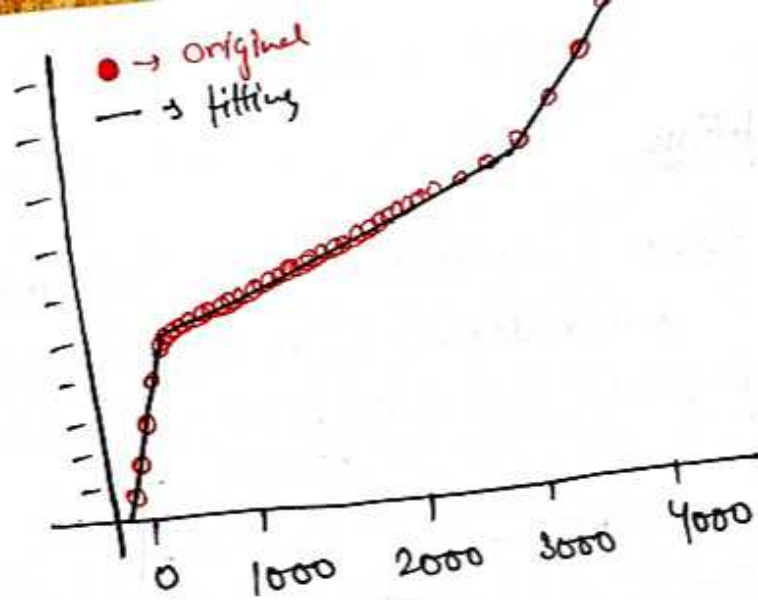
```
yHat = myPWLfit.predict(xHat)
```

```
# Calculate statistic
```

```
p = myPWLfit.p_value (method = 'non-linear', step-size = 1e-4)
```

```
se = myPWLfit.se
```

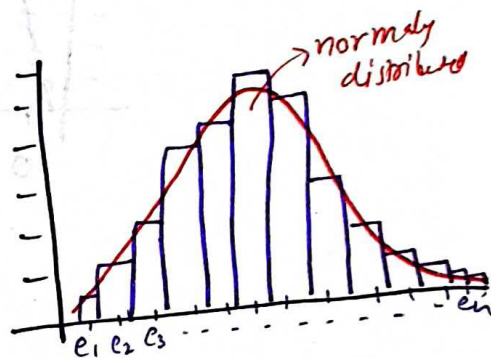
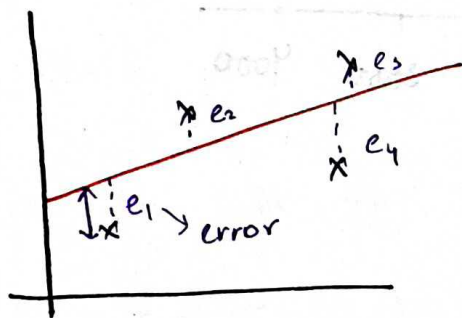
↓
Standard
error



2. Normality Of Residual

The Assumption

The error term (residual) are assumed to follow a normal distribution with a mean of zero and a constant variance.



What happens when this assumption is violated?

1. Inaccurate hypothesis tests: The t-tests and F-tests used to assess the significance of the regression coefficients and the overall model rely on the normality assumption. If the residuals are not normally distributed, these tests may produce inaccurate results, leading to incorrect inferences about the significance of the independent variables.
2. Invalid Confidence Intervals: The confidence interval for the regression coefficients ~~are~~

and the overall are based on the assumption of normally distributed residuals. If the normality assumption is violated, the confidence intervals may not be accurate, affecting the interpretation of the effect sizes and the precision of the estimates.

3. Model performance: The violation of the normality assumption may indicate that the chosen model is not the best fit for the data, potentially leading to reduced predictive accuracy.

How to check this assumption?

1. Histogram of Residuals: Plot a histogram of a residual to visually assess their visually across their distribution. If the histogram resembles a bell shaped curve. It suggested that the residual are normally distributed.

2. Q-Q plots

3. Statistical Test: Statistical tests like Omnibus test, Jarque - Bera test or even Shapiro Wilk test can test this assumption.

What to do when the assumption fails?

1. Model selection technique: Employ model selection techniques like cross-validation, AIC, or BIC to choose the best model among different candidate models that can handle non-normal residuals.
2. Robust regression: Use robust regression techniques that are less sensitive to the distribution of the residuals, such as M-estimation, least Median of square (LMS), or least Trimmed Square (LTS). (Transformation may also help)
3. Non-parametric or semi-parametric methods.
4. Use bootstrapping

Remember that the normality of Residual assumption is not always critical for Linear Regression, especially, when the sample size is large, due to Central Limit Theorem (more than 30)

→ automatic Normal
if $n > 30$

no. of observations

Omnibus Test


The omnibus test is a statistical test used to check if the residuals from a linear regression model follow normal distribution. The test is based on the skewness and kurtosis of the residuals. Here is a step by step guide on how to conduct the omnibus test.

1. Decide the Null and Alternate Hypothesis: The Null Hypothesis states that the residuals are normally distributed and the Alternate Hypothesis says that the residuals are not normally distributed.
2. Fit the linear Regression model: Fit the linear regression model to your data to obtain the predicted values.
3. Calculate the Residuals: Compute the residuals (error terms) by subtracting the predicted values from the observed values of the dependent variable.
4. Calculate the skewness: Calculate the skewness of the residuals. Skewness measures the asymmetry of the distribution. For a normal distribution skewness is expected to be close to zero.

5. Calculate the kurtosis: Calculate the kurtosis of the residuals, kurtosis the 'tailedness' of the distribution. For a normal distribution, kurtosis is expected to be close to zero (in excess kurtosis terms).

6. Calculate the Omnibus test statistic: Compute the Omnibus test statistic (K^2) using the skewness and kurtosis values. The formula for the Omnibus test statistic is:

chi-square $[df=2]$


$$K^2 = n \left[\frac{(\text{Skewness})^2}{6} + \frac{(\text{kurtosis})^2}{24} \right]$$

n = number of observation

7. Determine the p-value: The Omnibus test statistic follows a chi-square distribution with 2 degree of freedom. Use this distribution to calculate the p-value corresponding to the test statistic.

8. Compare the p-value to the significance level:

Compare the p-value obtained in step 6 to your chosen significance level (e.g., 0.05). If the p-value is greater than the significance

level, you can accept the null hypothesis that the residuals are normally distributed. If the p-value is similar than the significance level, you reject the null hypothesis, suggesting that the residual may not follow a normal distribution.

Shapiro - Wilk Test

The Shapiro-Wilk test is a hypothesis test that is applied to a sample with a null hypothesis that the sample has been generated from a normal distribution. If the p-value is low, we can reject such a null hypothesis and say that the sample has not been generated from a normal distribution.

But it has one flaw: It doesn't work well with large datasets. The maximum allowed size for a data set depends on the implementation (5000).

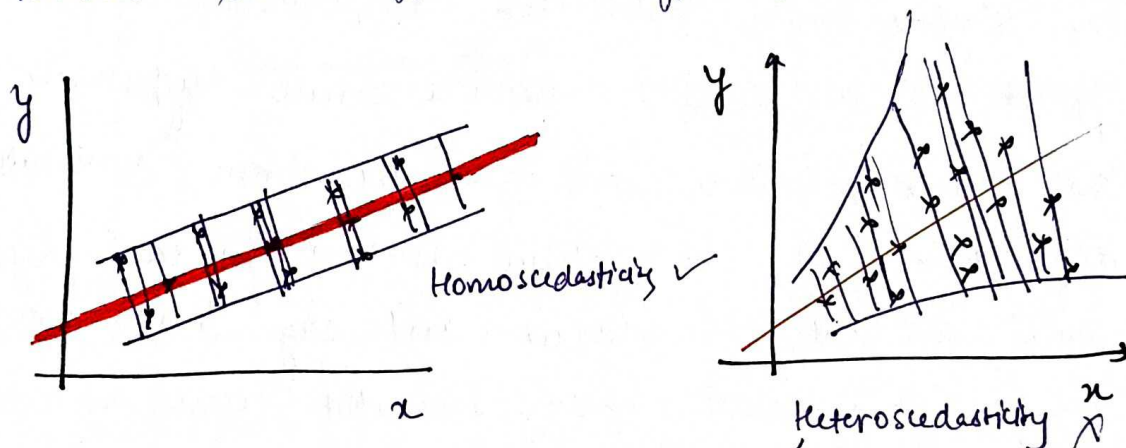
(i) p-value

(ii) statistic is nearly ~~to~~ zero for ~~too~~ normality.

3. Homoscedasticity

The Assumption

The spread the error term (residuals) should be constant across all level of the independent variable. If the spread of the residuals changes systematically, it leads to Heteroscedasticity, which can affect the efficiency of the estimate.



→ 1. Problem: Coefficient of Standard Error (SE) is not reliable. If Standard Error is not reliable then t -statistic is reliable and p -value also reliable.

2. Problem: Using not reliable standard error to find Confidence Interval.
not reliable

What happens when this assumption is violated?

1. Inefficient estimates: While the parameter estimate (coefficient) are still unbiased, they are no longer the best linear unbiased estimator (BLUE) under heteroscedasticity. The inefficiency of the estimates implies that the standard errors are then they should be, which may reduce the statistical power of hypothesis.
2. Inaccurate hypothesis tests: The t-tests and F-tests used to assess the significance of the regression coefficient and the overall model rely on the assumption of homoscedasticity, these tests may produce misleading results, leading to incorrect inference about the significance of the independent variables.
3. Invalid Confidence Intervals: The confidence intervals for the regression coefficients are based on the assumption of homoscedastic residuals. If the homoscedasticity assumption is violated, the confidence intervals may not be accurate, affecting the interpretation of the effect size and precision of the estimates.

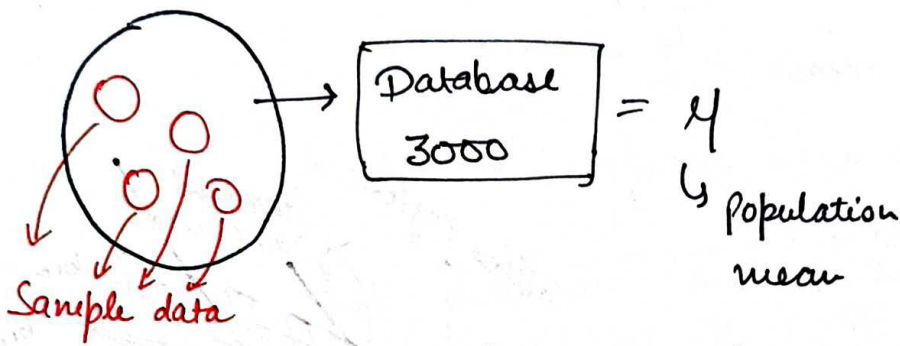
How to check this assumption?

1. Residual Plot:
2. Breusch-Pagan test: This is a formal statistical test for heteroscedasticity. The null hypothesis is that the error variance

What to do when the Assumption fails?

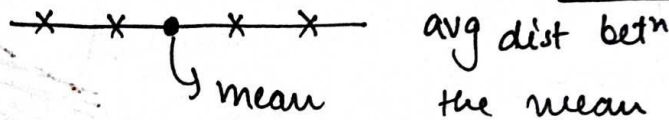
1. Transformation
2. Weighted least square (WLS): Use a weighted least squares approach, which assigns different weights to the observations based on the magnitude of their residuals. This method can help account for heteroscedasticity by giving more importance to observations with smaller residuals and less importance to those with larger residuals.
3. Robust standard errors: Calculate robust (heteroscedasticity-consistent) standard errors for the regression coefficients. These standard errors are more reliable under heteroscedasticity and can be used to perform more accurate hypothesis tests and construct valid confidence intervals.

Standard Error



$$\mu = \frac{\sum x_i}{n} \quad \text{std Deviation}$$

$$\sigma^2 = \frac{\sum (x_i - \bar{x})^2}{n}$$



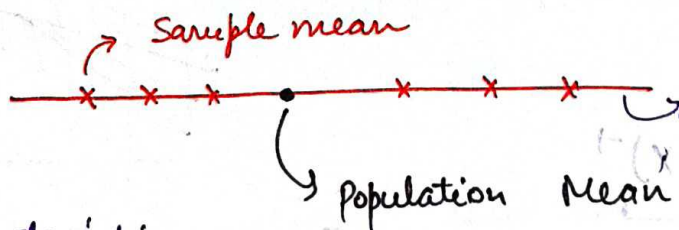
* If standard deviation \downarrow (small) means ^{that} all points are close (***x)

* If standard deviation \uparrow (large) means ^{not} all points are not close (+ x - x)

sample data 1 $\boxed{100} = \bar{x}_1$

sample data 2 $\boxed{100} = \bar{x}_2$

sample data 3 $\boxed{100} = \bar{x}_3$



Find standard deviation of this line is called standard error.

* Standard deviation of the Sample means is standard error

Central Limit Theorem $\boxed{SE = \frac{\sigma}{\sqrt{n}}}$

pop std deviation \rightarrow

no. of observation \rightarrow

\hookrightarrow when extracting sample from Population at many times and find off mean of each sample individually and That distribution is called sample distribution because we are observing

Sample distribution more than 30 and also called
 Normal Distribution $(\mu, \sigma/\sqrt{n})$
 μ pop mean σ/\sqrt{n} std error

Example

* Sample dataset (100)

exp	salary
-	-
-	-
-	-
-	-

$SE(b_2)$

* Another sample dataset (100)

exp	salary
-	-
-	-
-	-

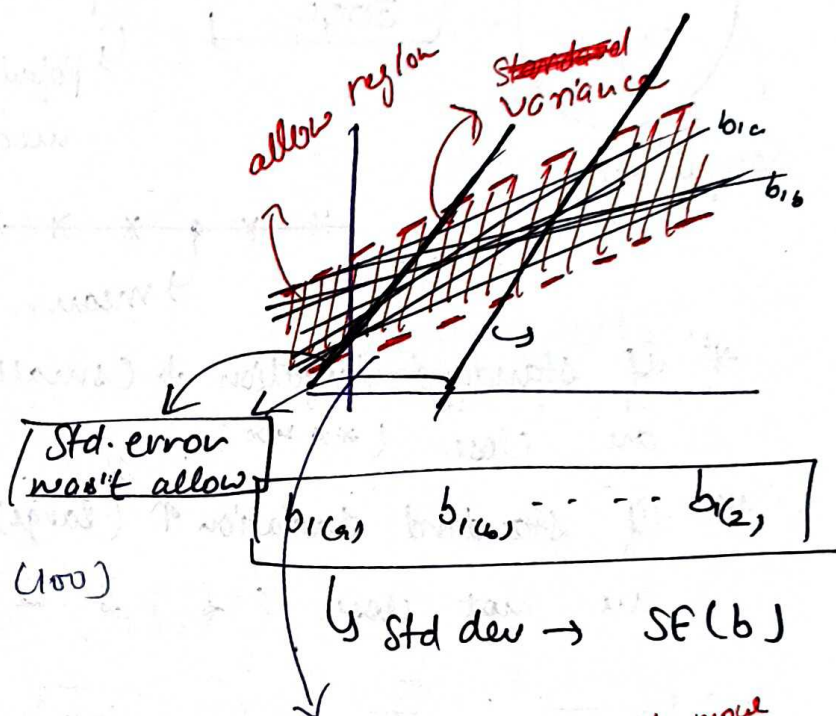
$SE(b_2)$

$$\text{Var}(\beta) = \sigma^2 \underbrace{(X^T X)^{-1}}_{3 \times 3}$$

$$SE = \sqrt{a \rightarrow \text{var}(\beta_0)}$$

$$SE = \sqrt{b \rightarrow \text{var}(\beta_1)}$$

$$SE = \sqrt{c \rightarrow \text{var}(\beta_2)}$$



$$\sigma^2 \begin{bmatrix} - & - & - \\ - & - & - \\ - & - & - \end{bmatrix} \Rightarrow \begin{bmatrix} a & - & - \\ - & b & - \\ - & - & c \end{bmatrix}$$

$$SE(\beta) = \sqrt{\text{diagonal}(\sigma^2 (X^T X)^{-1})}$$