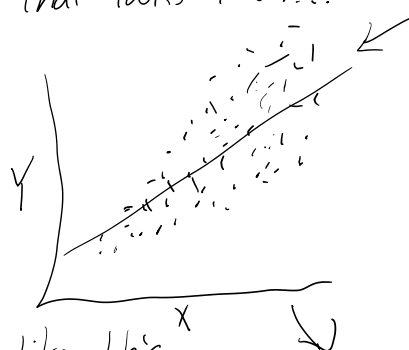## Diagnostics

When doing regression we assume follows a <u>linear trend</u> and that <u>residuals have no pattern</u> and are normally distributed. $\rightarrow (y_i - \hat{y}_i) = resid_i$

residual $\hat{\varepsilon}_i^\bullet$

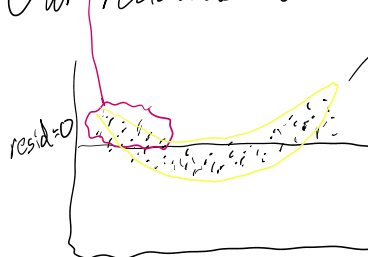One of the best ways to confirm assumptions is to plot residuals against fitted value.
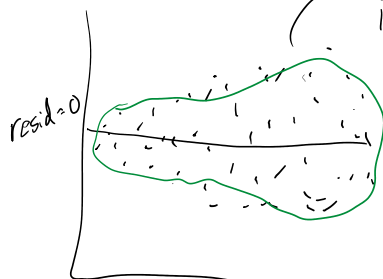What if we have data that looks like....
(non-linear)

non-homoscedastic
or
heteroscedastic

Scedasticity refers to randomness

Y

X

Y

X

Our residuals will look like this....

This banana shape is an indicator of non-linearity

resid=0

This pear shape is an indicator non-constant variance.

resid=0

We want <u>no pattern</u>

gaps are okay but maybe indicate different groups

Diagnosing normality of residuals can be done with a simple histogram or more methodically with a quartile-quartile plot (qq plot)

A qqplot will find the quartiles of the data and plot them against theoretical quartiles.
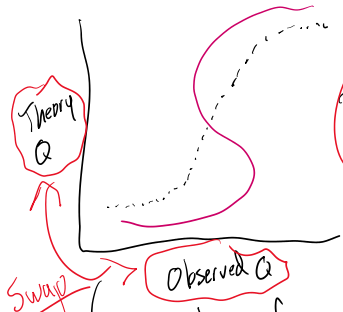
Bad QQplot          Good QQplot          Bad QQplot
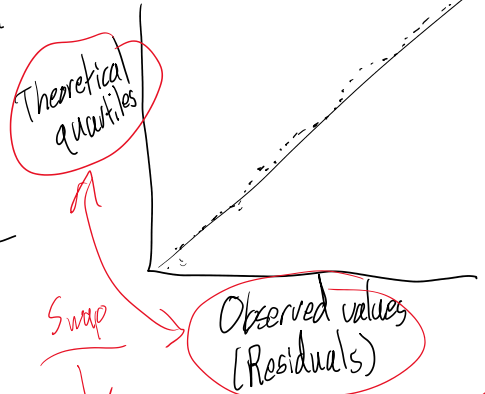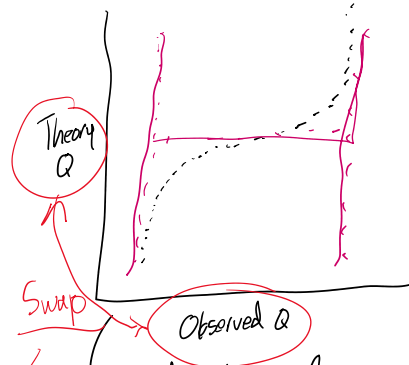
theoretical quartiles.

**Bad QQplot**

Theory Q

Swap

Observed Q

Indicative of short tails

Correction from lecture: observed quantiles on y-axis, theoretical on X-axis

**Good QQplot**

Theoretical quantiles

Swap

Observed values (Residuals)

observed quantiles on y axis, theoretical on X-axis

**Bad QQplot**

Theory Q

Swap

Observed Q

Indicative of heavy tails

**Correct**

observed Quantiles

Theoretical Quantiles

We can sometimes resolve these issues by transforming the data and/or adding more terms to the model

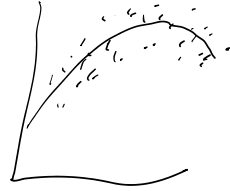rather than $y_i = b_0 + b_1 X + \varepsilon_i$ ← resid...

We can add a squared term for X.

we might consider $y_i = b_0 + b_1 X + b_2 X^2$ →

But we can all sorts of transformations on X or Y.

$$\log(y_i) = b_0 + b_1 X + b_2 X^2$$

$$\sqrt{y_i} = b_0 + b_1 X + b_2 \log(X) + b_3 \sqrt{X}$$

We want to find the $\hat{b_0}, \hat{b_1}, \hat{b_2}, \ldots$ that best describe the data with linear relations between outcome and predictors.

→ I.E: The relationship between $\sqrt{y_i}$ and $\log(X)$ is still linear.

→ IE: The relationship between $y_i$ and logits is still linear.

One transformation that is common but has "little impact" on the model

We can take $X$ and $Y$ and "normalize" or "standardize" them.

$$Z_{Y_i} = \frac{Y_i - \overline{Y}}{S_Y} \qquad Z_{X_i} = \frac{X_i - \overline{X}}{S_X}$$

↳ mean 0, sd 1       ↳ mean 0, sd 1.

Benefit: Slope interpretation becomes: "For a 1 standard deviation increase in $X$, we expect a $b_1$ standard deviation increase in $Y$."

Benefit: $b_0 = 0$



fit line
standardize
o       Original Data
fit line