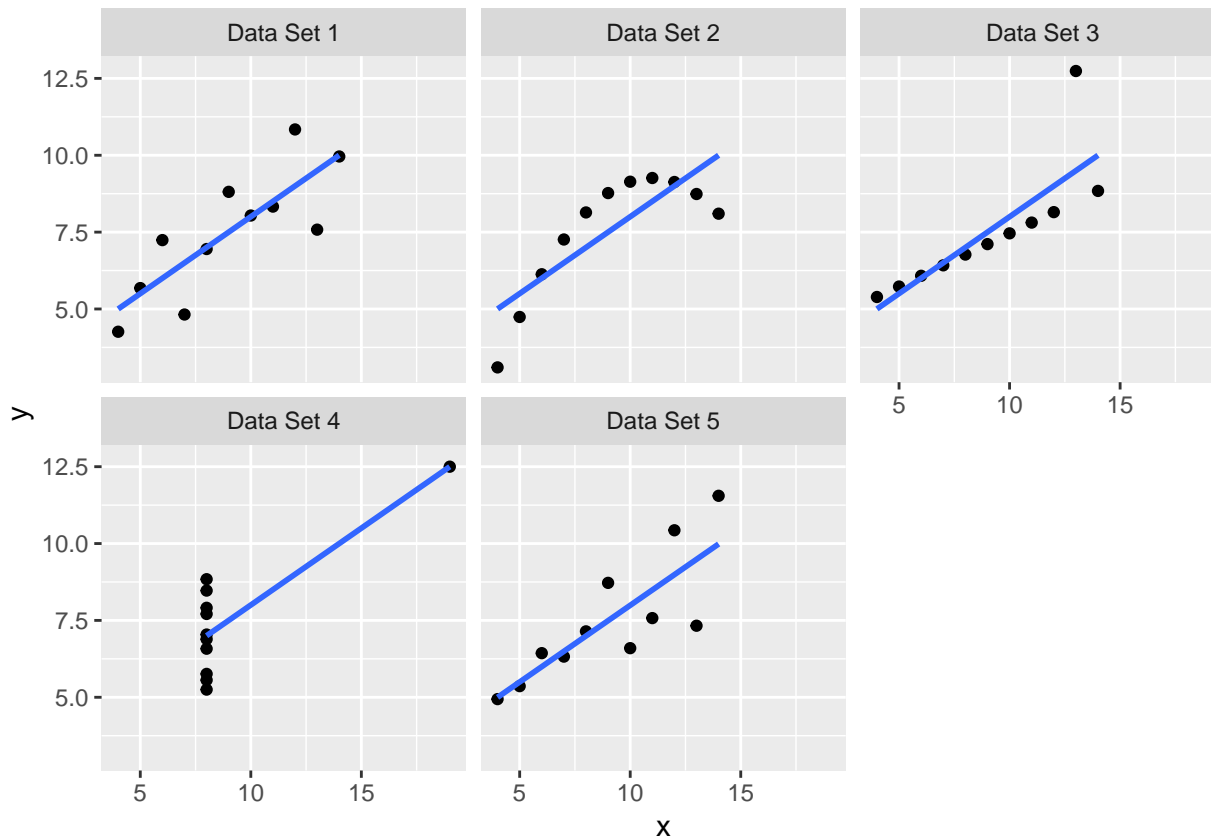# Linear Regression: Conditions for Inference, Residual Diagnostics
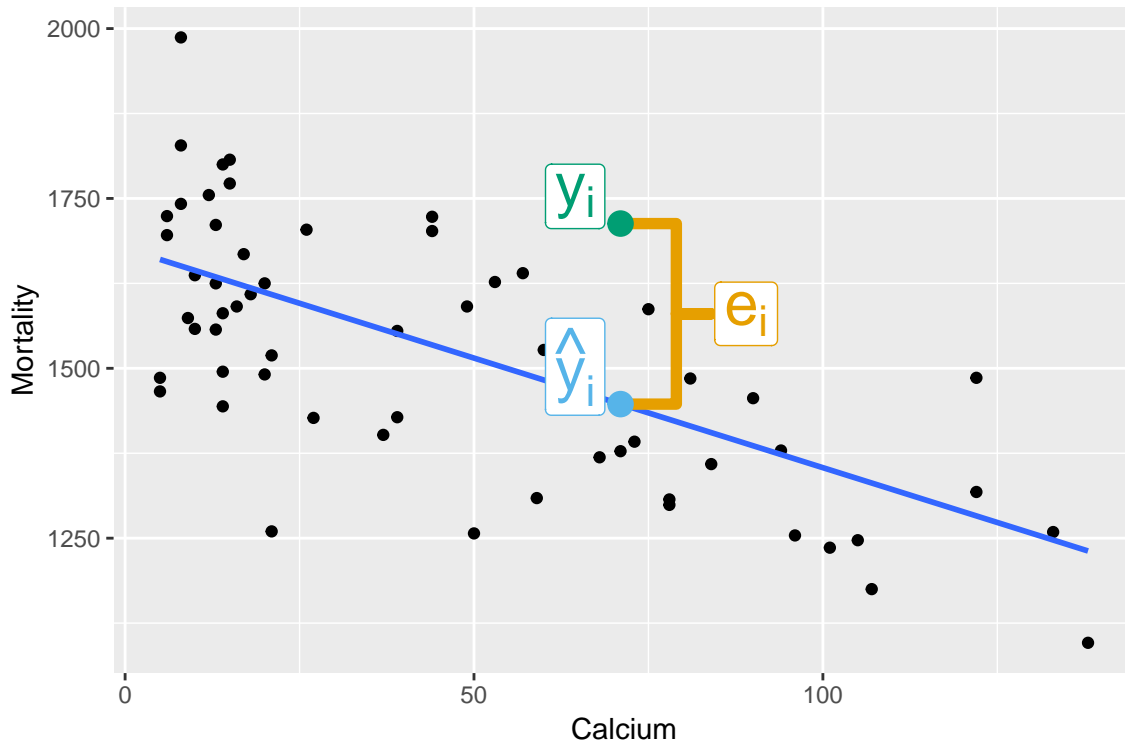
Anscombe Lab Wrap-Up

## All 5 Have Essentially the Same Estimated Intercept Slope, Intercept, $R^2$, and Residual Standard Deviation!



- Briefly, **conditions for linear regression** (see last page for more detailed summary):
  - Sample **representative** of population
  - No **outliers** (points that don't fit the trend)
  - **Linear** relationship
  - **Independent** observations
  - **Normally** distributed residuals
  - **Equal variability** of residuals
- **Use plots** to help diagnose the appropriateness of a linear model:
  - Scatter plot of explanatory (x axis) vs. response (y axis)
  - Scatter plot of predicted (x axis) vs. residual (y axis)
  - Histogram or density plot of residuals (x axis)
- Checks of whether the sample is representative and whether the observations are independent come from thinking about data collection process, not plots.
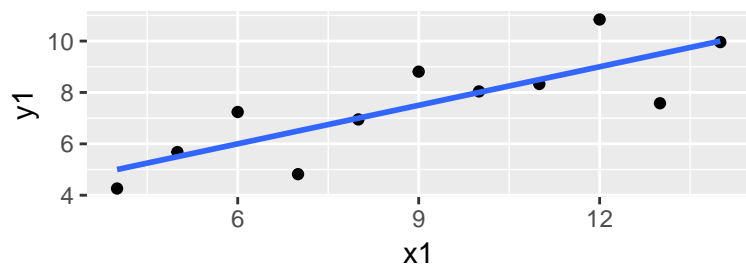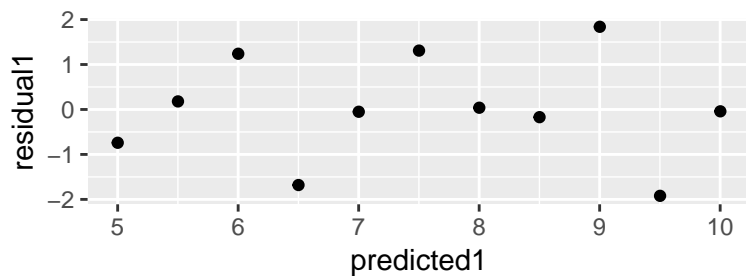
# A Reminder about Residuals



- Residuals give the vertical distance between a data point and the line of best fit
- Positive if point above line, negative otherwise
- Residual = *Observed* - *Predicted*
- $e_i = y_i - \widehat{y}_i$ ($e$ stands for error)

# Anscombe Quintet: Data Set 1 (The Way Life Should Be)

```r
ggplot(data = anscombe, mapping = aes(x = x1, y = y1)) +
  geom_point() +
  geom_smooth(method = "lm", se = FALSE)
```
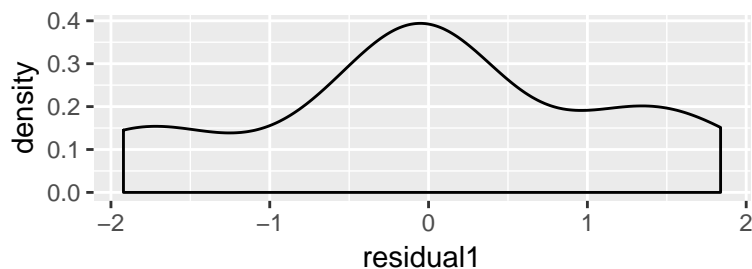


```r
linear_fit1 <- lm(y1 ~ x1, data = anscombe)
anscombe <- anscombe %>% mutate(
  predicted1 = predict(linear_fit1),
  residual1 = residuals(linear_fit1)
)
ggplot(data = anscombe, mapping = aes(x = predicted1, y = residual1)) +
  geom_point()
```
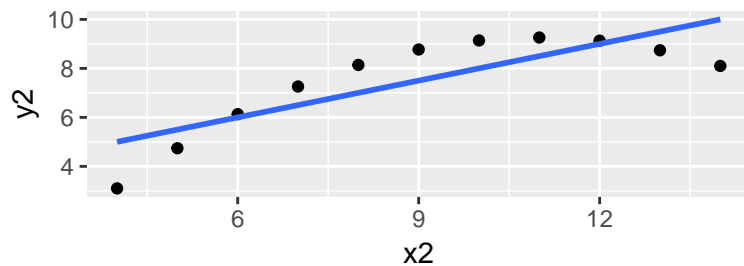


```r
ggplot(data = anscombe, mapping = aes(x = residual1)) +
  geom_density()
```



- **Outliers**?

- **Linear** relationship?

- **Normally** distributed residuals?

- **Equal variability** of residuals?

## Anscombe Quintet: Data Set 2

```r
ggplot(data = anscombe, mapping = aes(x = x2, y = y2)) +
  geom_point() +
  geom_smooth(method = "lm", se = FALSE)
```



```r
linear_fit2 <- lm(y2 ~ x2, data = anscombe)
anscombe <- anscombe %>% mutate(
  predicted2 = predict(linear_fit2),
  residual2 = residuals(linear_fit2)
)
ggplot(data = anscombe, mapping = aes(x = predicted2, y = residual2)) +
  geom_point()
```



```r
ggplot(data = anscombe, mapping = aes(x = residual2)) +
  geom_density()
```
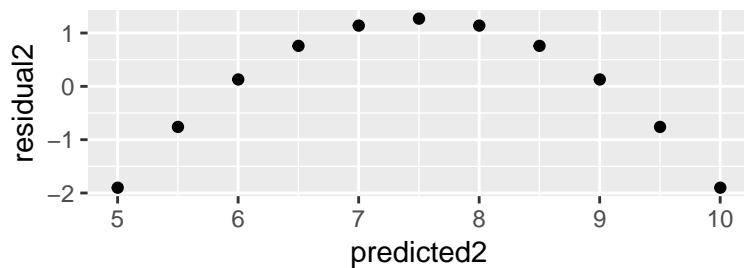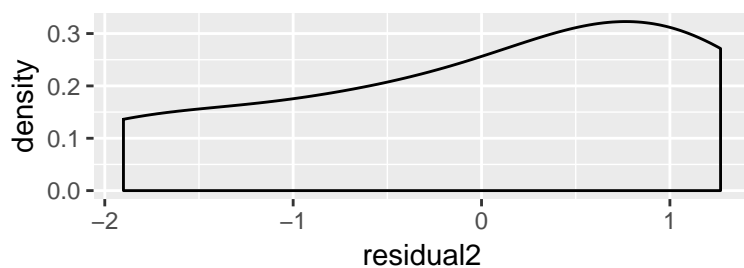


- **Outliers**?

- **Linear** relationship?

- **Normally** distributed residuals?

- **Equal variability** of residuals?

# Anscombe Quintet: Data Set 3

```r
ggplot(data = anscombe, mapping = aes(x = x3, y = y3)) +
  geom_point() +
  geom_smooth(method = "lm", se = FALSE)
```



```r
linear_fit3 <- lm(y3 ~ x3, data = anscombe)
anscombe <- anscombe %>% mutate(
  predicted3 = predict(linear_fit3),
  residual3 = residuals(linear_fit3)
)
ggplot(data = anscombe, mapping = aes(x = predicted3, y = residual3)) +
  geom_point()
```
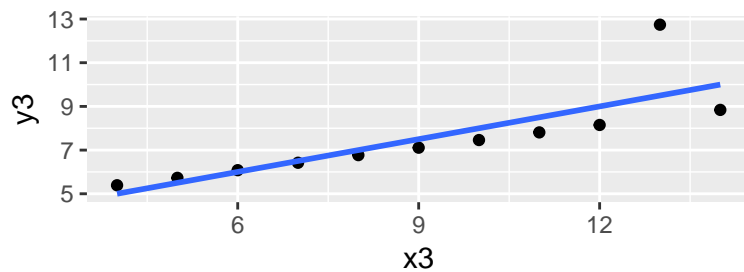


```r
ggplot(data = anscombe, mapping = aes(x = residual3)) +
  geom_density()
```
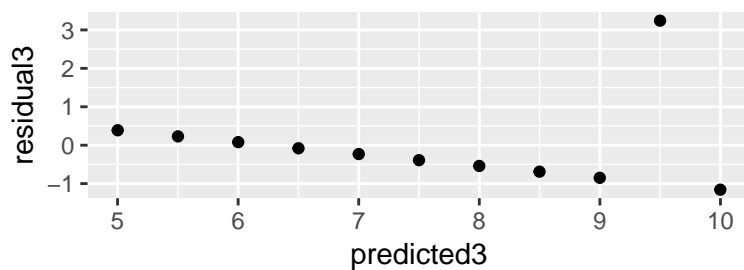


- **Outliers**?


- **Linear** relationship?


- **Normally** distributed residuals?


- **Equal variability** of residuals?

# Anscombe Quintet: Data Set 4

```r
ggplot(data = anscombe, mapping = aes(x = x4, y = y4)) +
  geom_point() +
  geom_smooth(method = "lm", se = FALSE)
```
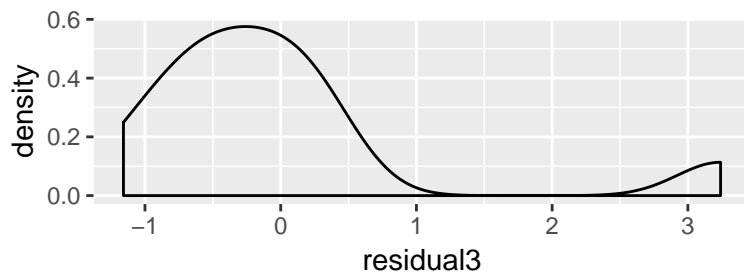


```r
linear_fit4 <- lm(y4 ~ x4, data = anscombe)
anscombe <- anscombe %>% mutate(
  predicted4 = predict(linear_fit4),
  residual4 = residuals(linear_fit4)
)
ggplot(data = anscombe, mapping = aes(x = predicted4, y = residual4)) +
  geom_point()
```
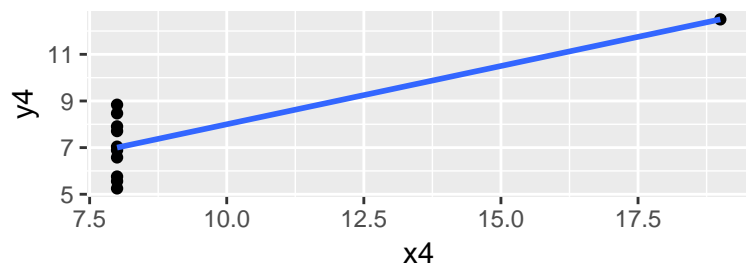


```r
ggplot(data = anscombe, mapping = aes(x = residual4)) +
  geom_density()
```
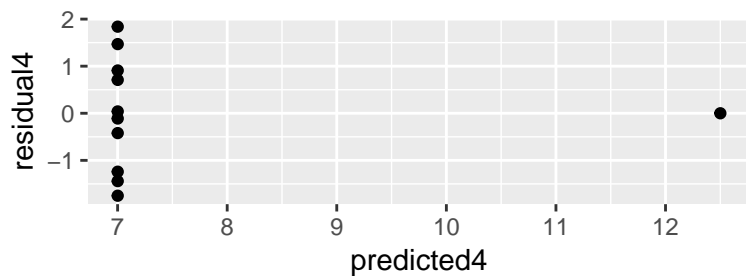


- **Outliers**?

- **Linear** relationship?

- **Normally** distributed residuals?

- **Equal variability** of residuals?

# Anscombe Quintet: Data Set 5

```r
ggplot(data = anscombe, mapping = aes(x = x5, y = y5)) +
  geom_point() +
  geom_smooth(method = "lm", se = FALSE)
```
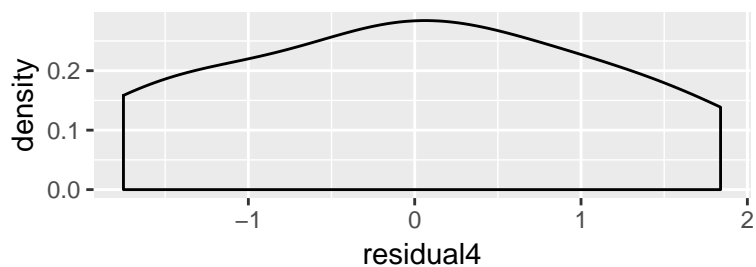


```r
linear_fit5 <- lm(y5 ~ x5, data = anscombe)
anscombe <- anscombe %>% mutate(
  predicted5 = predict(linear_fit5),
  residual5 = residuals(linear_fit5)
)
ggplot(data = anscombe, mapping = aes(x = predicted5, y = residual5)) +
  geom_point()
```
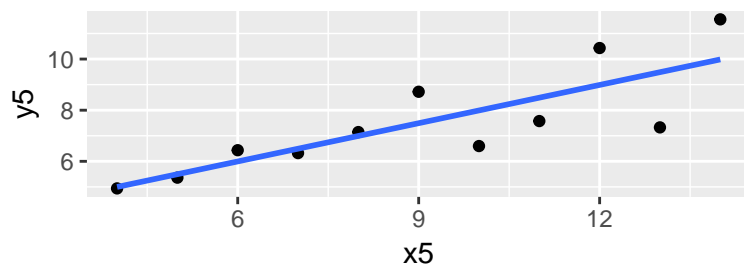


```r
ggplot(data = anscombe, mapping = aes(x = residual5)) +
  geom_density()
```
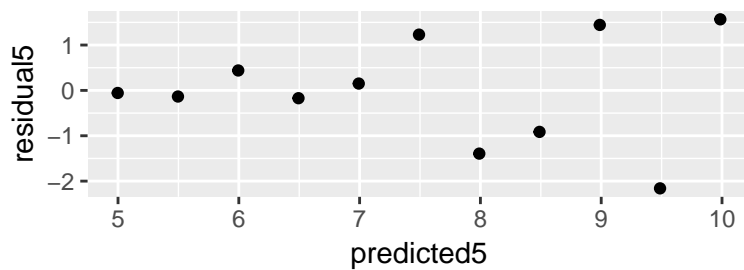


- **Outliers**?


- **Linear** relationship?


- **Normally** distributed residuals?


- **Equal variability** of residuals?

# Regression Conditions

Think of a helpful leprechaun named **R**obert **O'Line**:



| Condition | How Important? | How to Check? |
|---|---|---|
| **R**epresentative sample | Critical | Think about data collection (randomization?) |
| No **O**utliers | Very Important | Plots: <ul><li>Scatter Plot of explanatory variable vs response variable (no points stand out)</li><li>Scatter plot of predicted value vs residuals (no points stand out)</li><li>histogram or density plot of residuals (no outliers)</li></ul> |
| **L**inear relationship | Very Important | Plots: <ul><li>Scatter Plot of explanatory variable vs response variable (pattern is linear)</li><li>Scatter plot of predicted value vs residuals (no curved patterns)</li></ul> |
| **I**ndependent observations | Very Important | <ul><li>Think about data collection (randomization?)</li><li>Situations where observations are **not** independent:<ul><li>Observations collected over time (e.g., monthly unemployment measurements over time)</li><li>Observations collected in space (e.g., number of pitcher plants in each square meter of a bog)</li><li>Multiple observations on the same person (e.g., baseline and follow-up measurements of health in a clinical trial)</li></ul></li></ul> |
| **N**ormal distribution for residuals | Somewhat Important | Plots: <ul><li>histogram or density plot of residuals (approximately symmetric, no outliers)</li></ul> |
| **E**qual variability of residuals about the line as the explanatory variable changes. | Somewhat Important | Plots: <ul><li>Scatter Plot of explanatory variable vs response variable (same amount of vertical spread around line for all values of $x$)</li><li>Scatter plot of predicted value vs residuals (same amount of vertical spread for all values of $x$)</li></ul> |