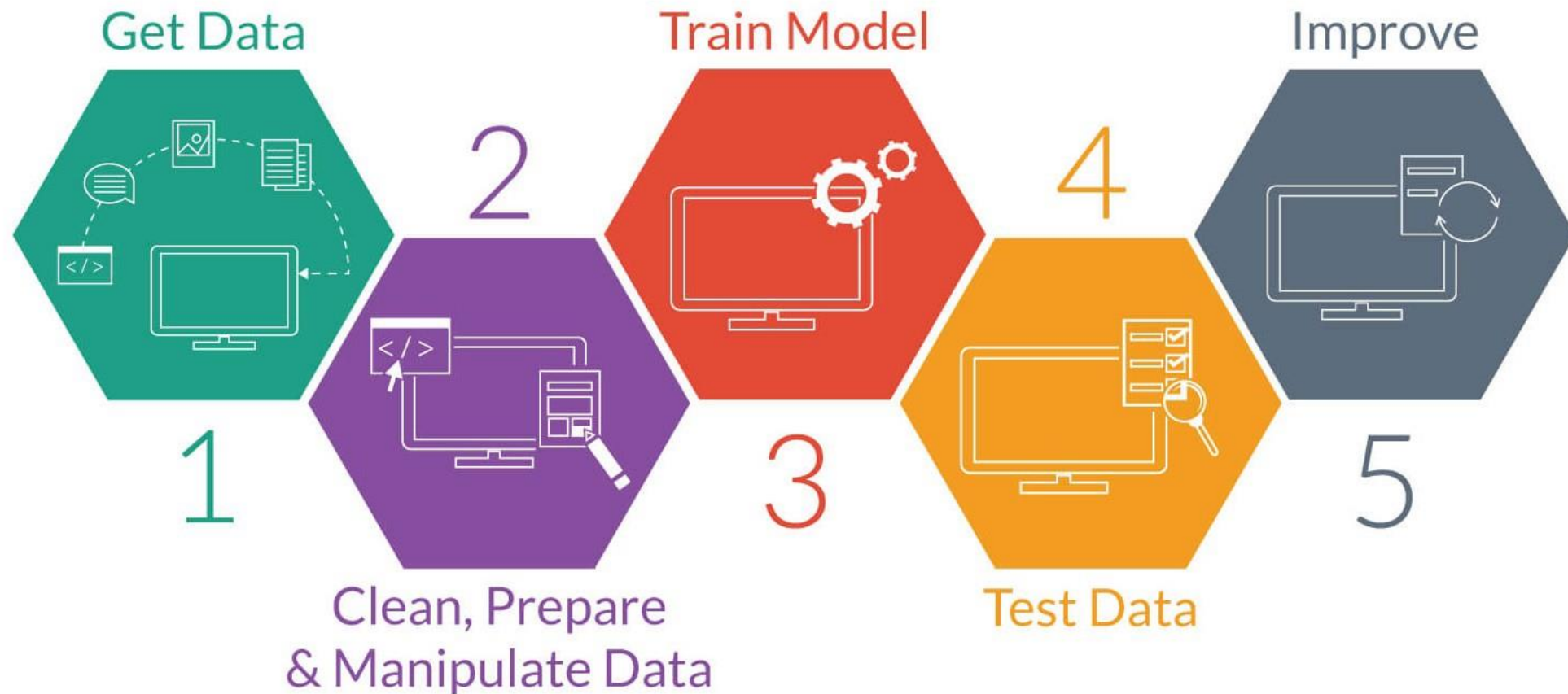




Feature Selection In Machine Learning

dataaspirant.com

Building a Machine Learning model

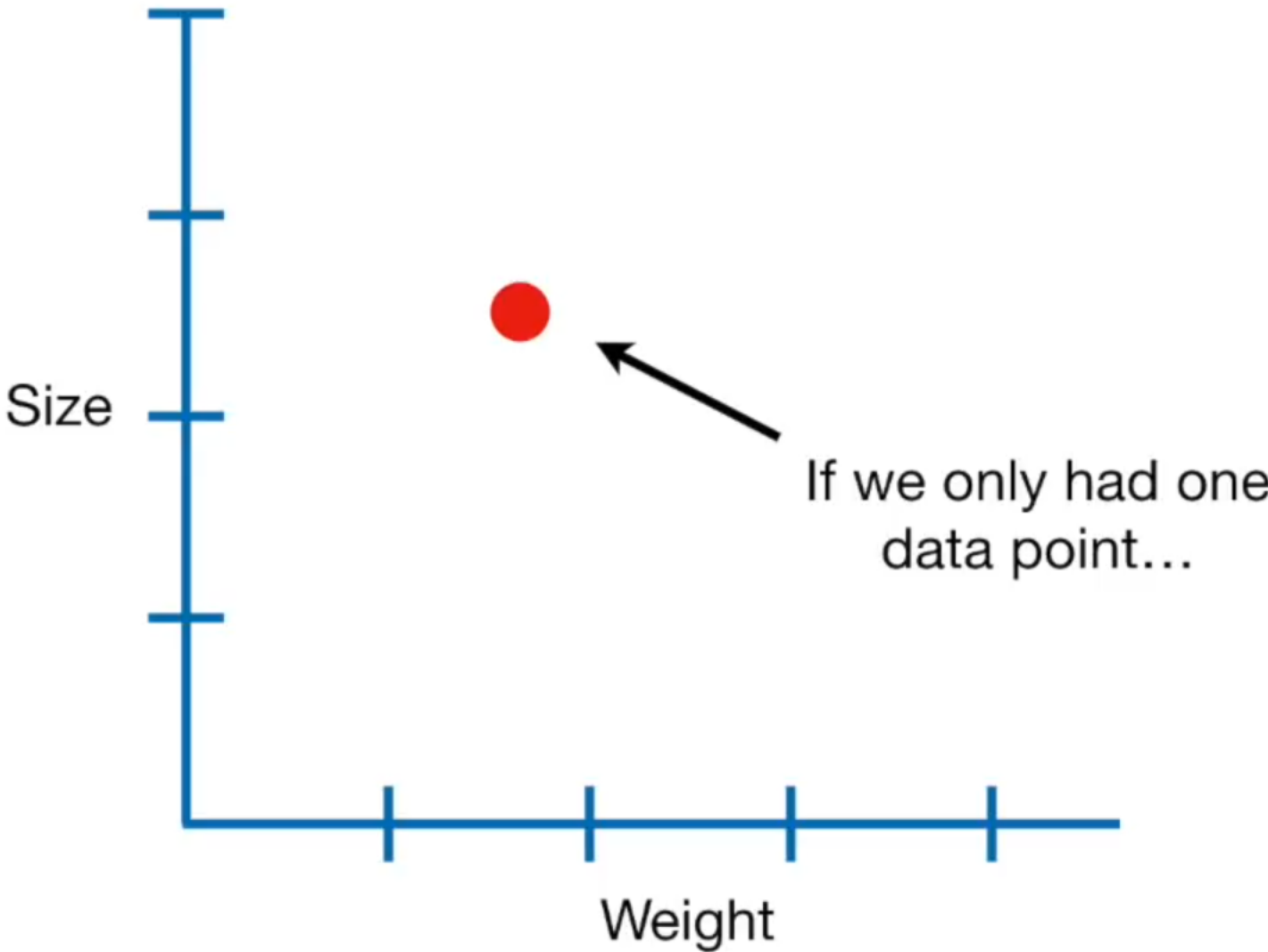


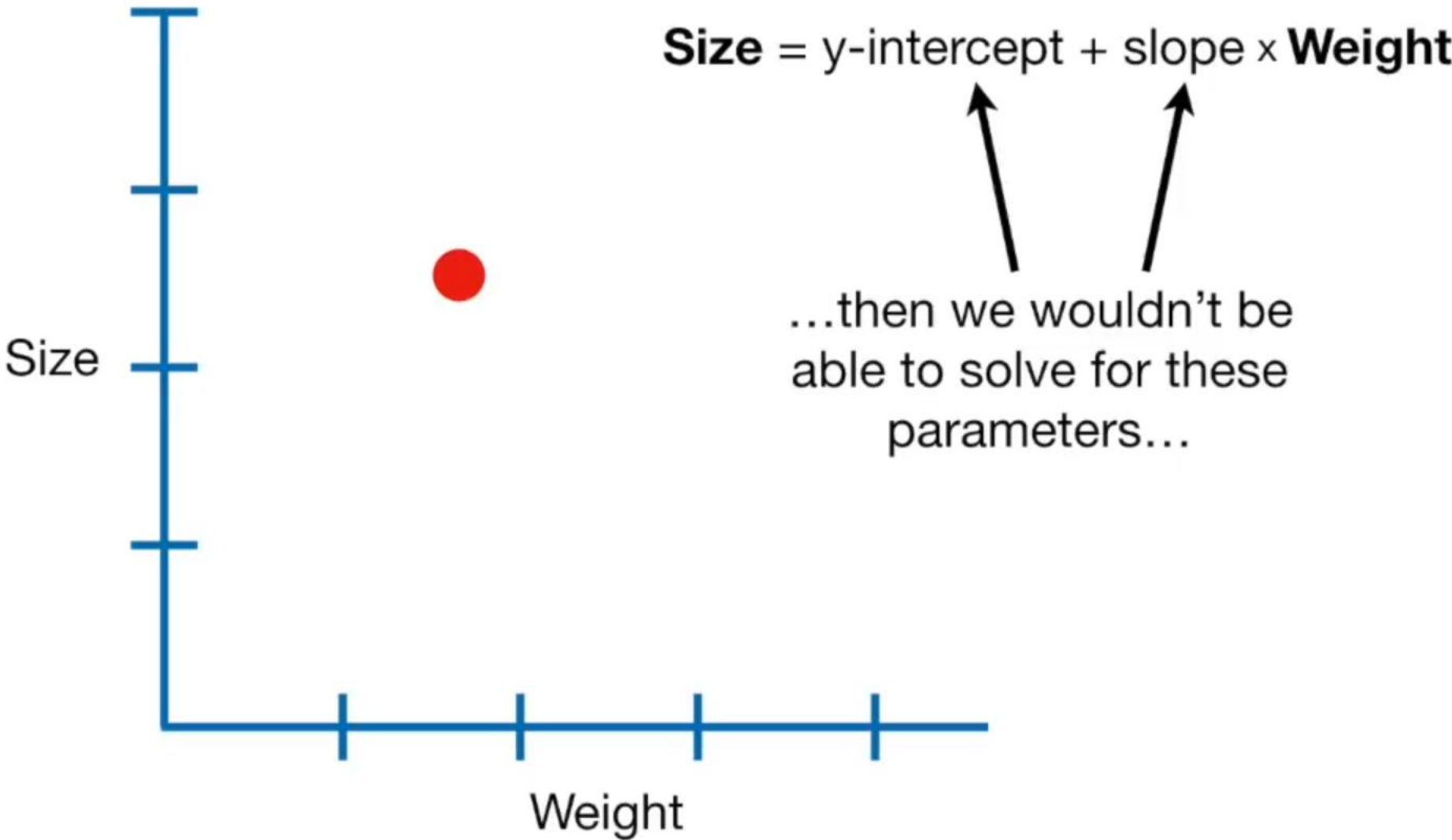
Data Preparation

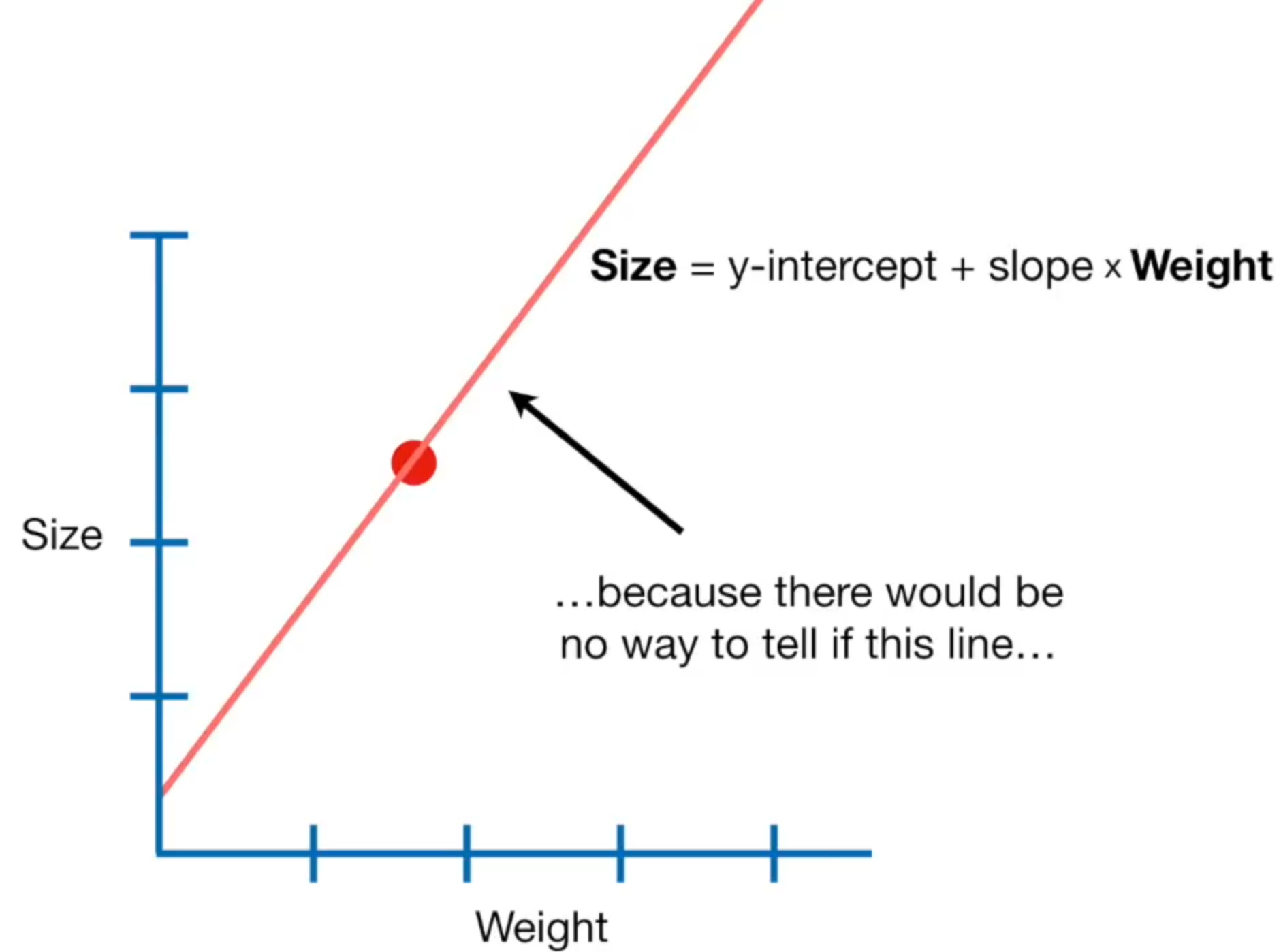
- Combine data from various sources
- Check data schema:
 - Usually, the source of the data has schema describing the data
- Clean your Data: identify and handle errors in your data.
 - Handle Missing data
 - Handle outliers
- Data Transformation: change the scale of some/all variables. Why?? We will discuss it later.
- Feature Selection: select these features that are most relevant to your task.
- Feature Engineering: combine features, derive new variables, dimensionality reduction, etc.

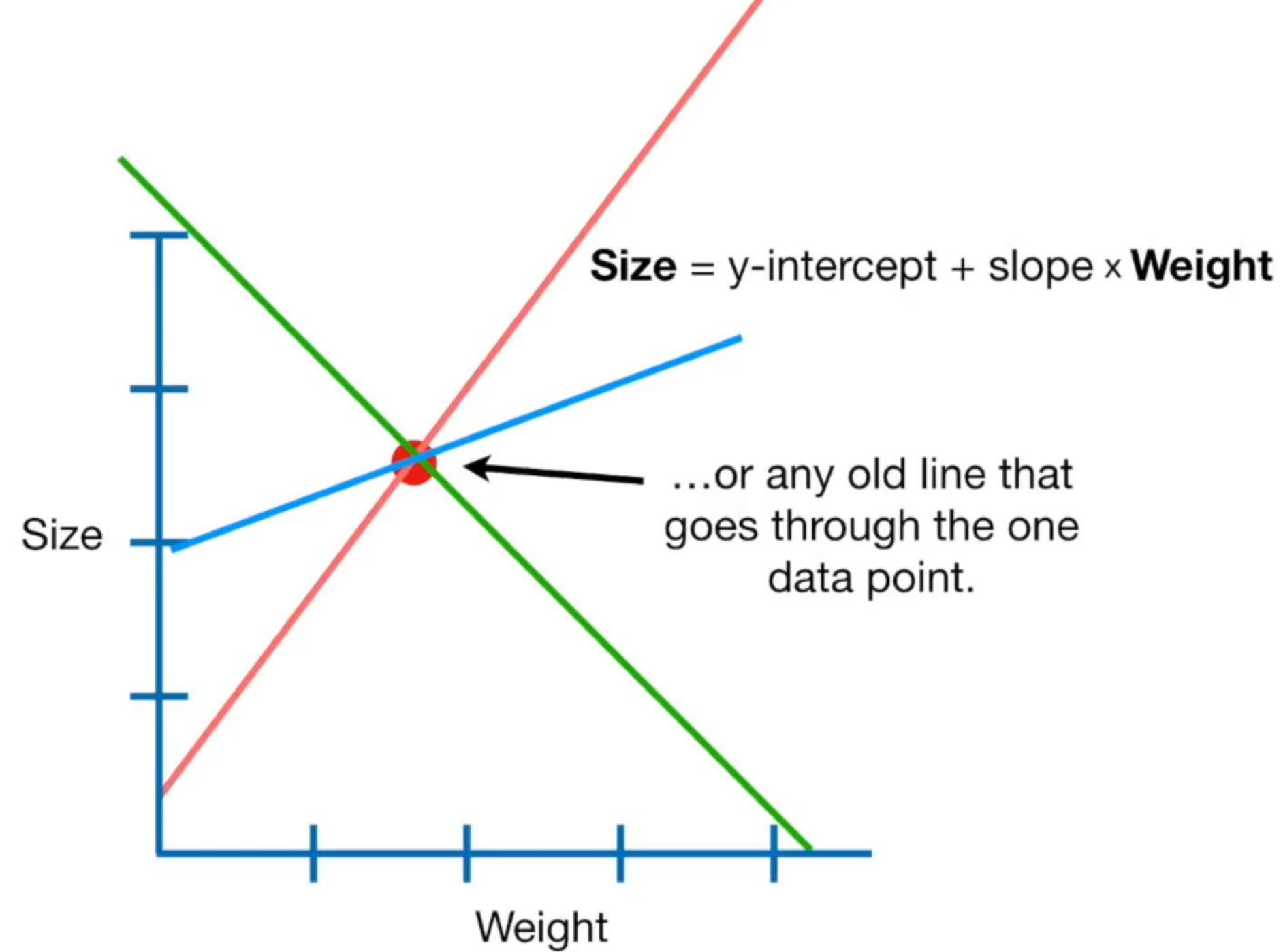
Exploratory Data Analysis

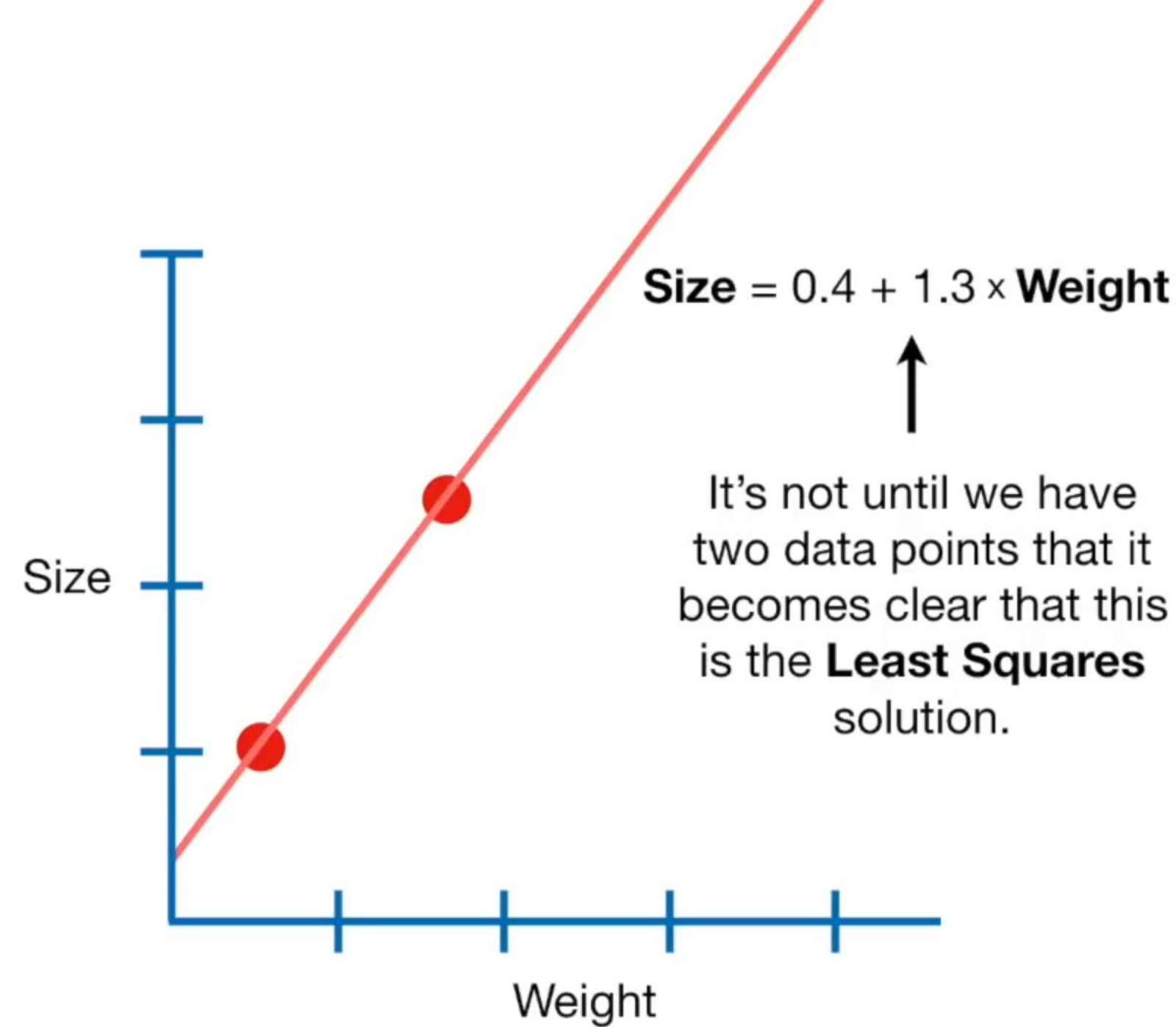
- We can't prepare the data without first exploring the data.
- Quantify missing data
- Identify numerical and categorical variables
- Determine unique values (cardinality) in categorical features
- Check rare/ dominant categories in categorical features
- Highlight outliers
- Identify linear relationships
- Identify a normal distribution
- Check histograms





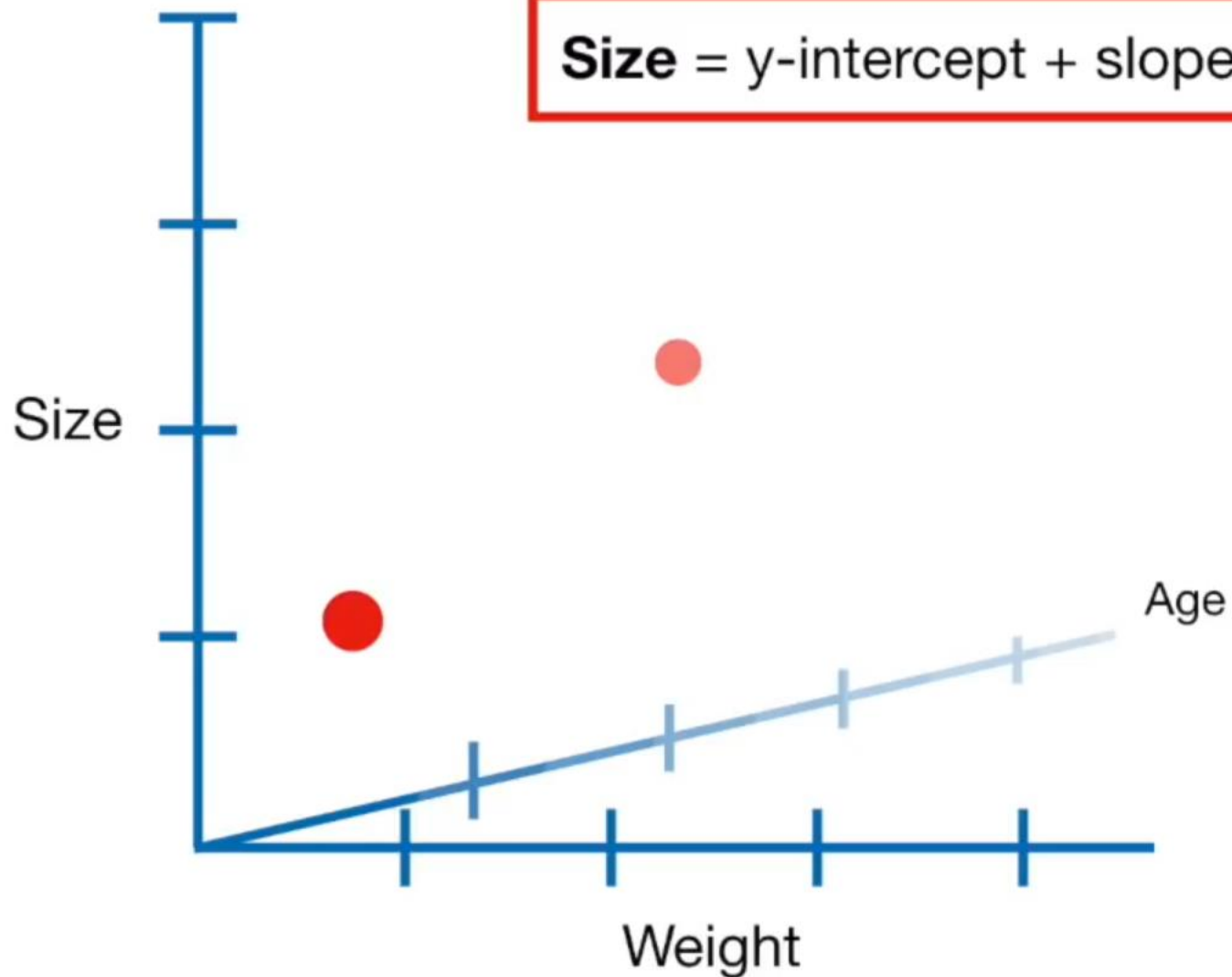






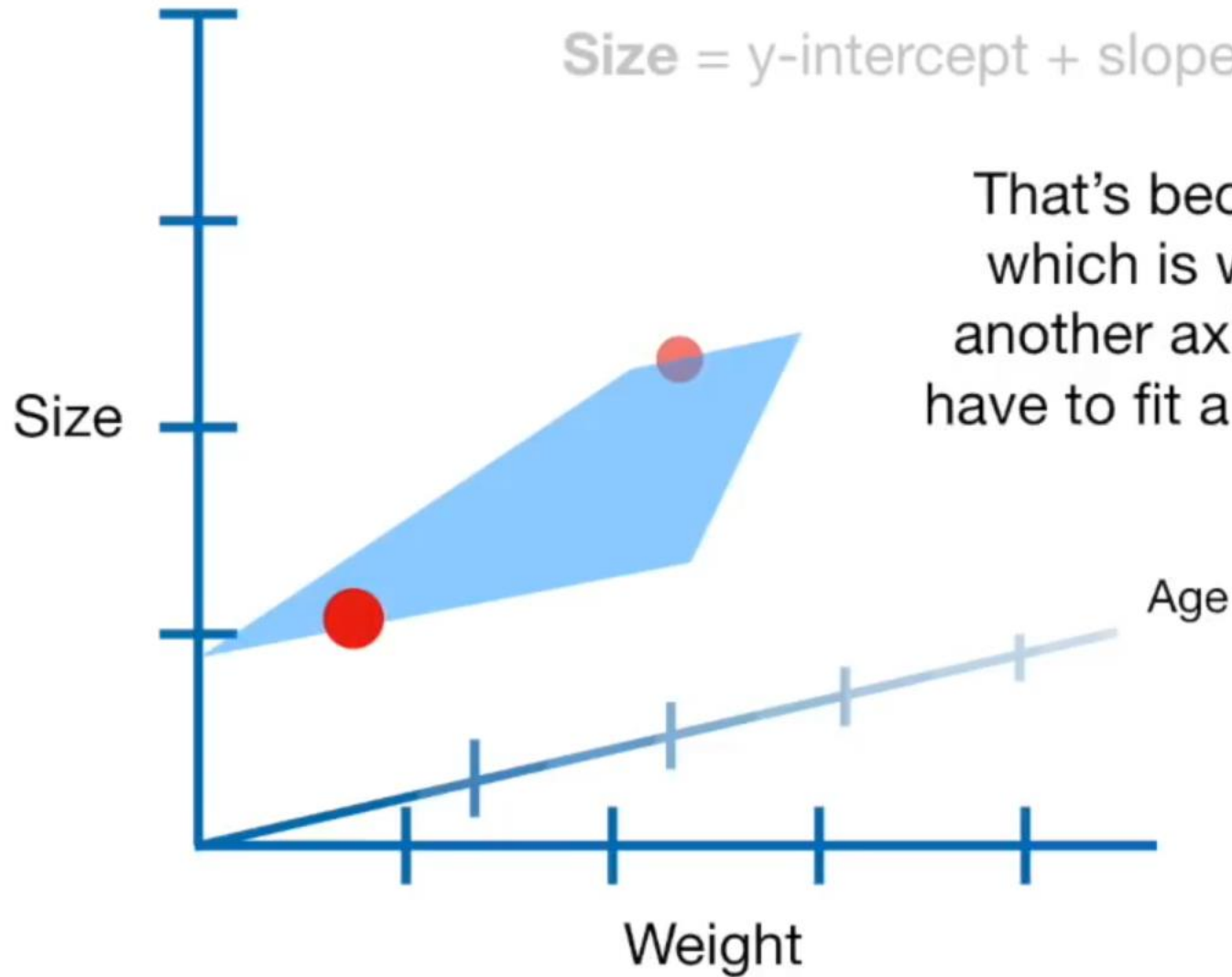
Now let's look at an equation that has three parameters to estimate.

$$\text{Size} = \text{y-intercept} + \text{slope1} \times \text{Weight} + \text{slope2} \times \text{Age}$$



$$\text{Size} = \text{y-intercept} + \text{slope1} \times \text{Weight} + \text{slope2} \times \text{Age}$$

That's because in three dimensions, which is what we get when we add another axis to our graph for **Age**, we have to fit a plane to the data instead of just a line...



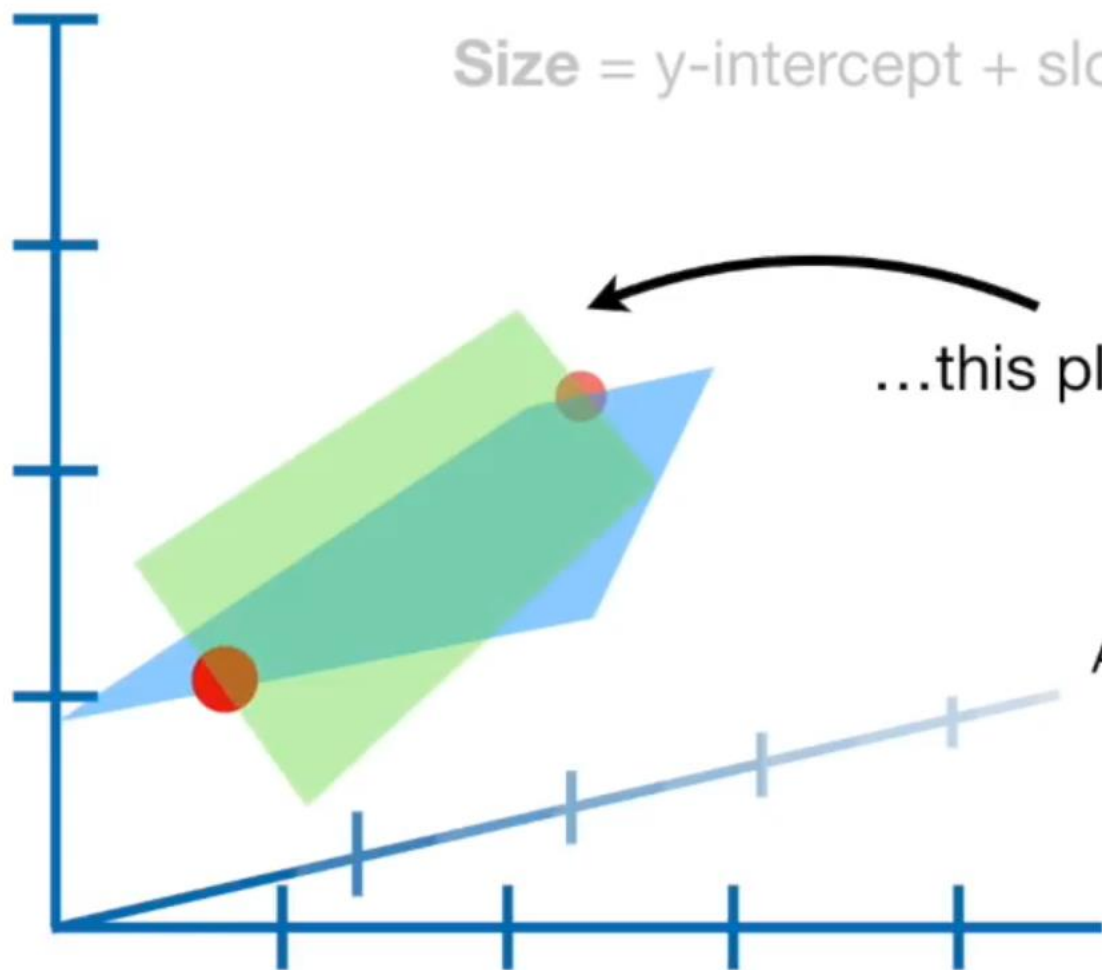
$$\text{Size} = \text{y-intercept} + \text{slope1} \times \text{Weight} + \text{slope2} \times \text{Age}$$

...this plane...

Size

Age

Weight



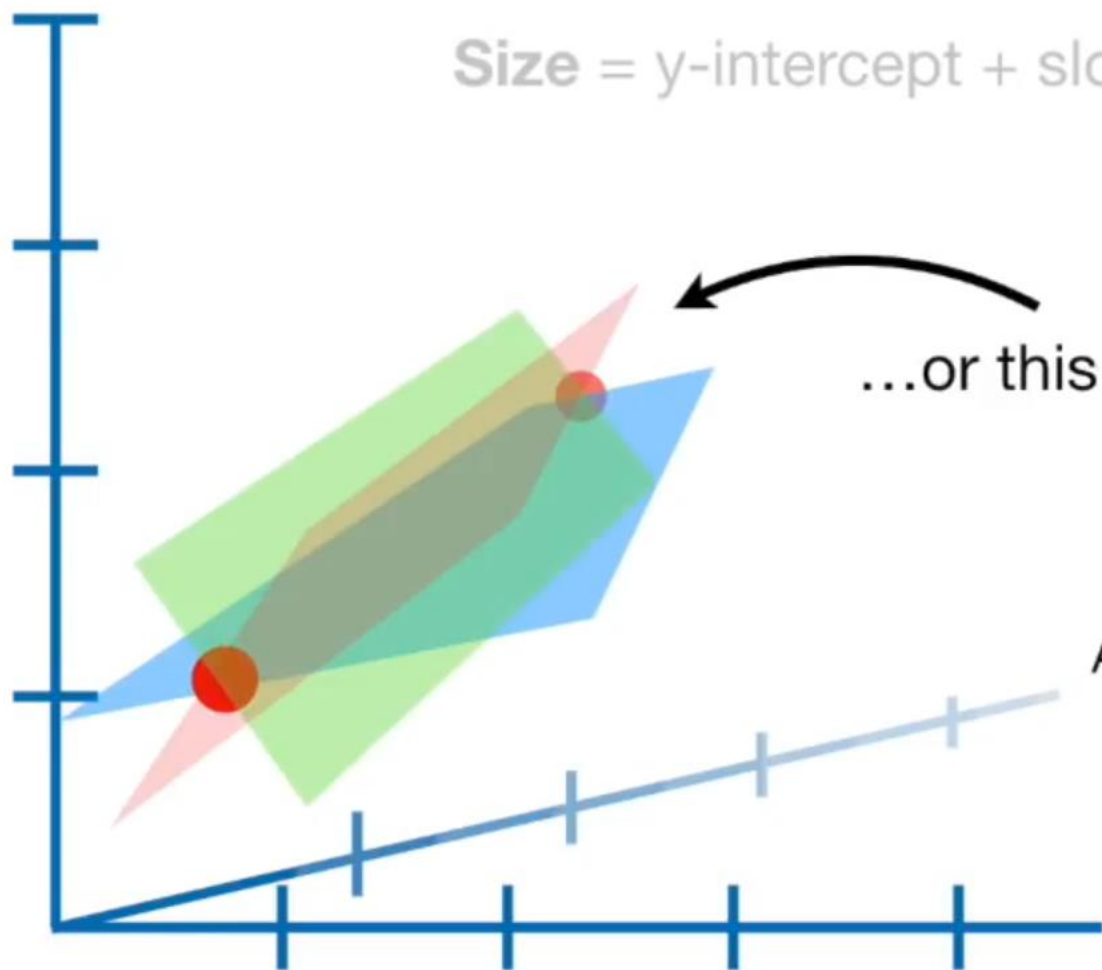
$$\text{Size} = \text{y-intercept} + \text{slope1} \times \text{Weight} + \text{slope2} \times \text{Age}$$

...or this plane.

Size

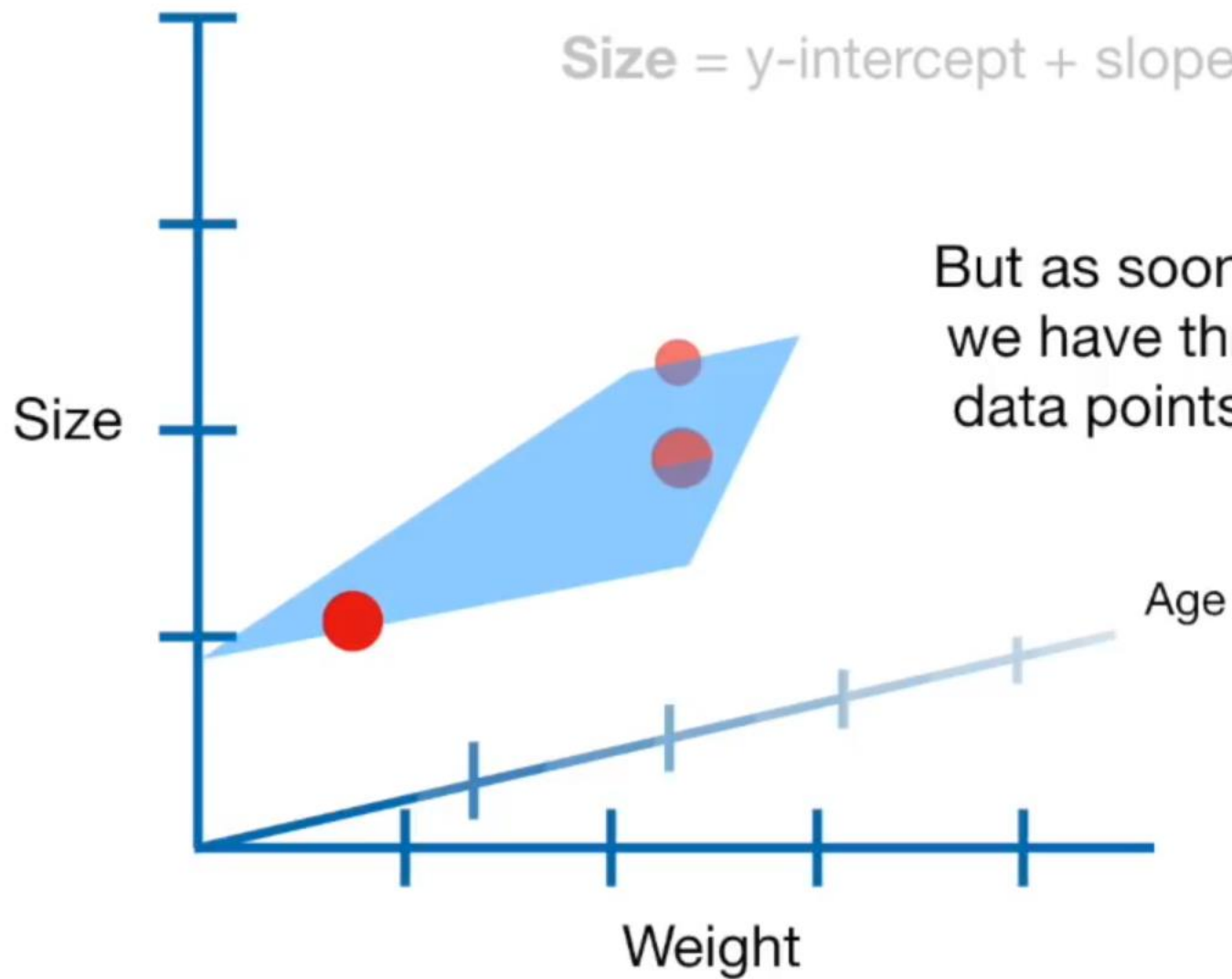
Age

Weight



$$\text{Size} = \text{y-intercept} + \text{slope1} \times \text{Weight} + \text{slope2} \times \text{Age}$$

But as soon as
we have three
data points...



$$\mathbf{Size} = \text{y-intercept} + \text{slope1} \times \mathbf{Weight} + \text{slope2} \times \mathbf{Age} + \text{slope3} \times \mathbf{fur\ color}$$

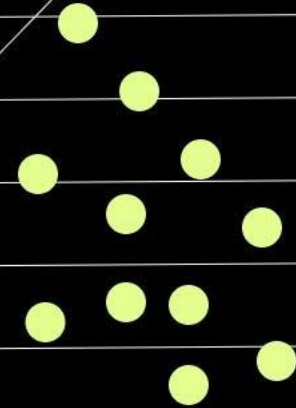
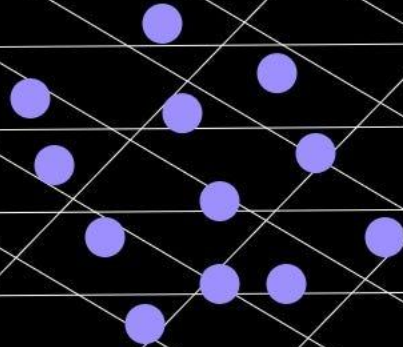
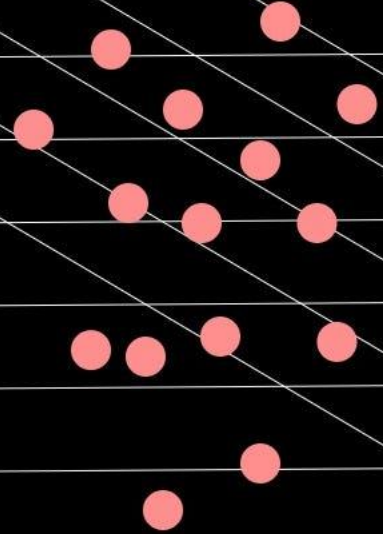


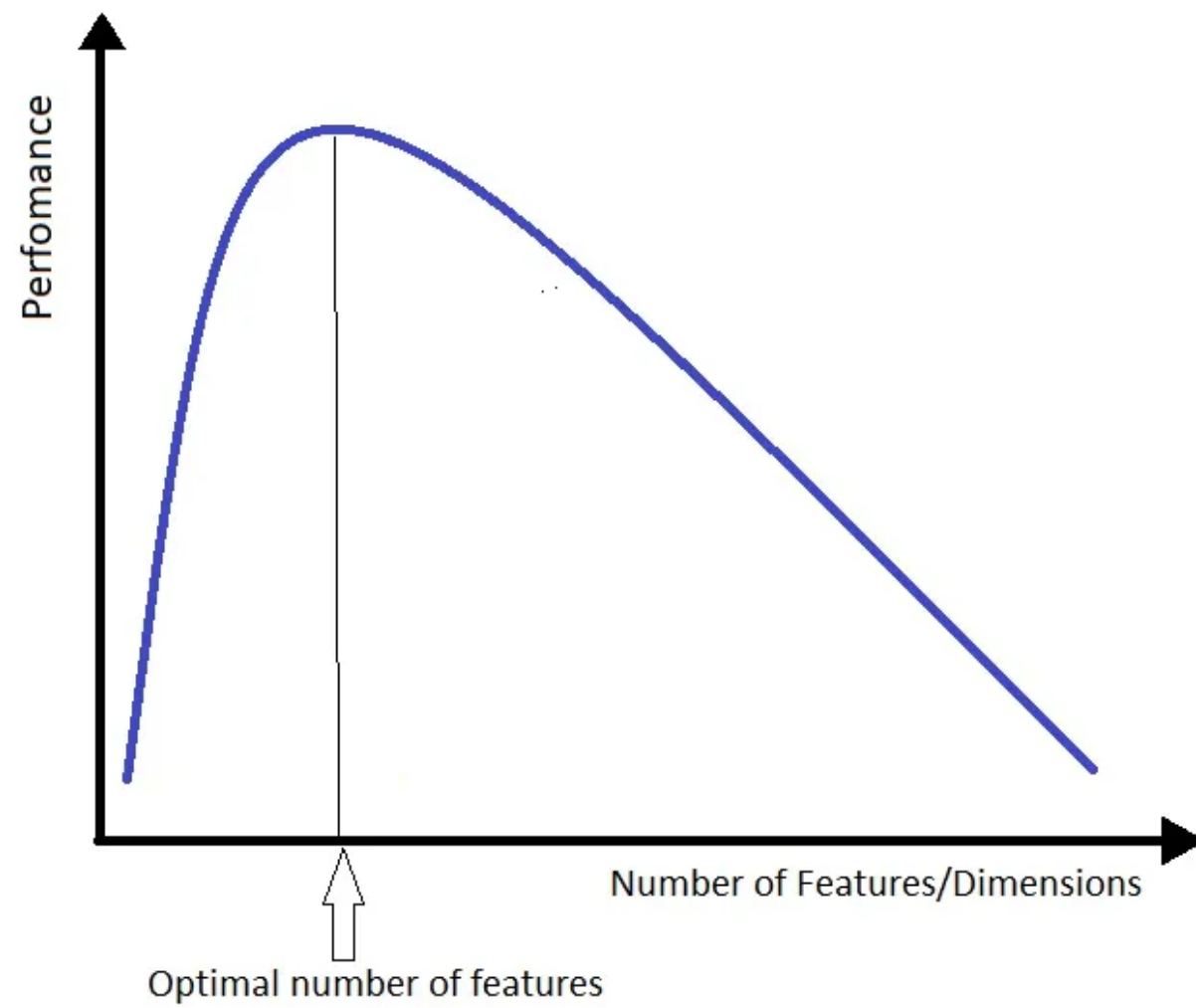
If we have an equation with four
parameters...

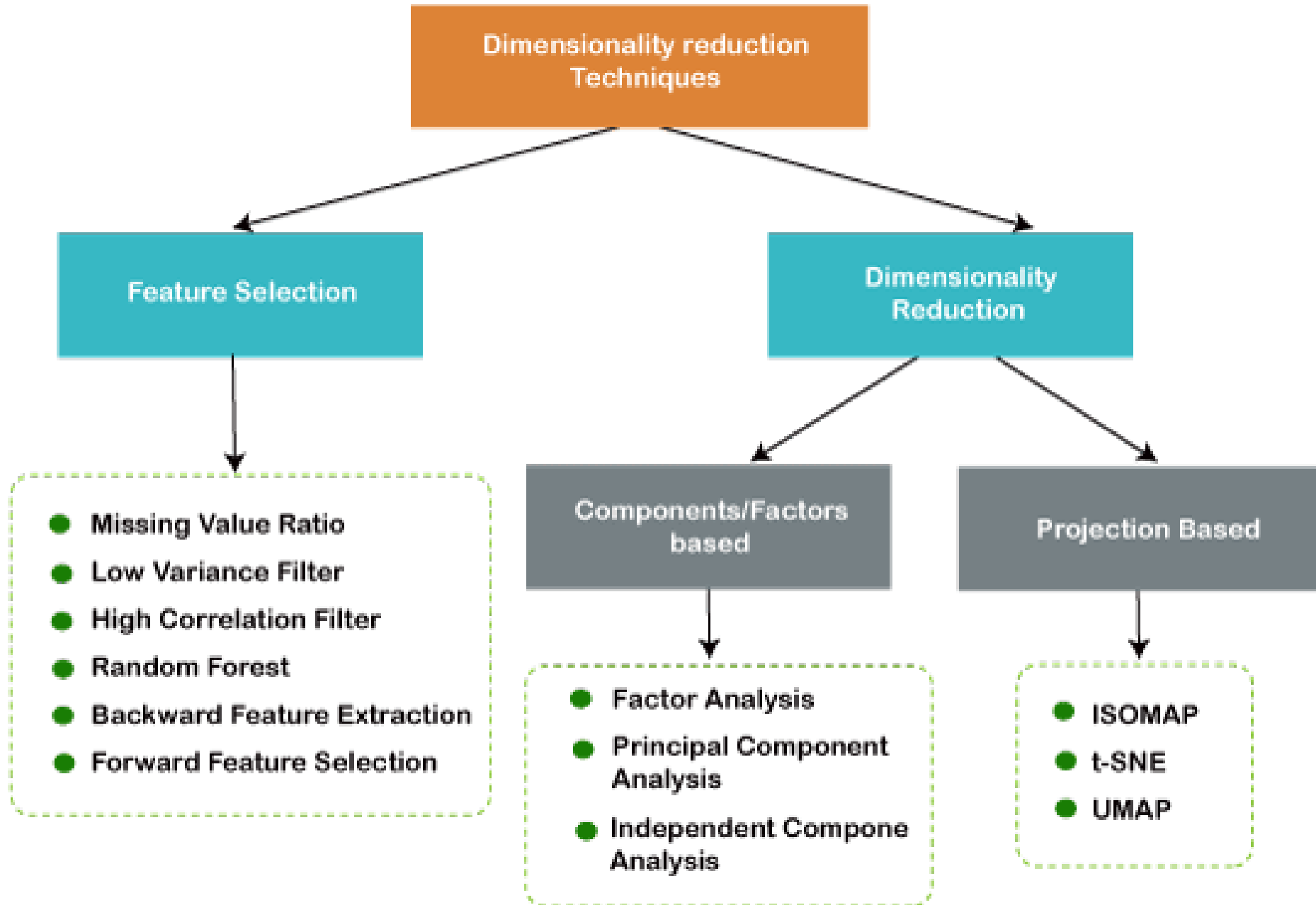
...then **Least Squares** needs at least four
data points to estimate all four
parameters.

curse of

dimensionality

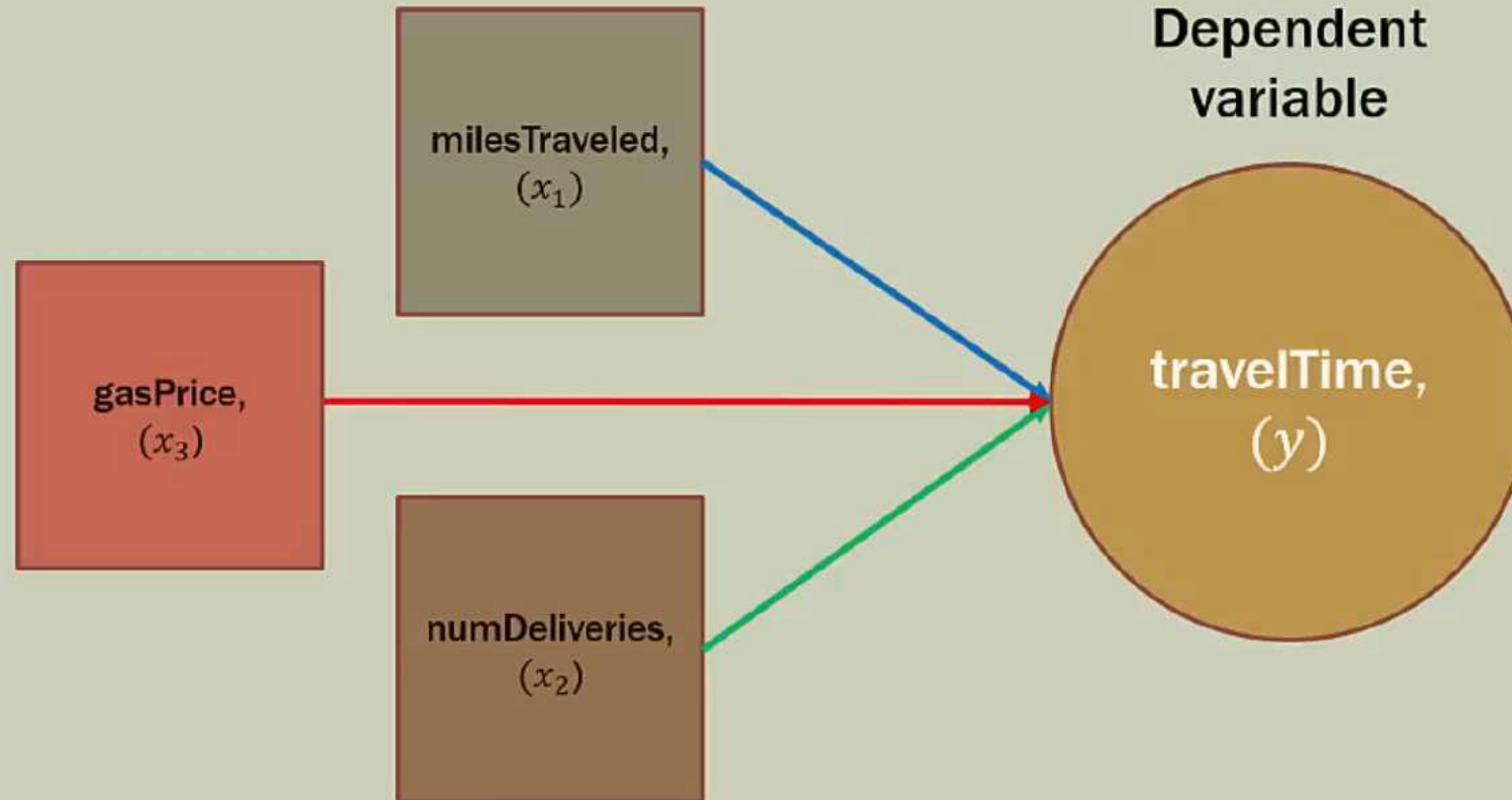






SKETCHING OUT RELATIONSHIPS

Independent variables



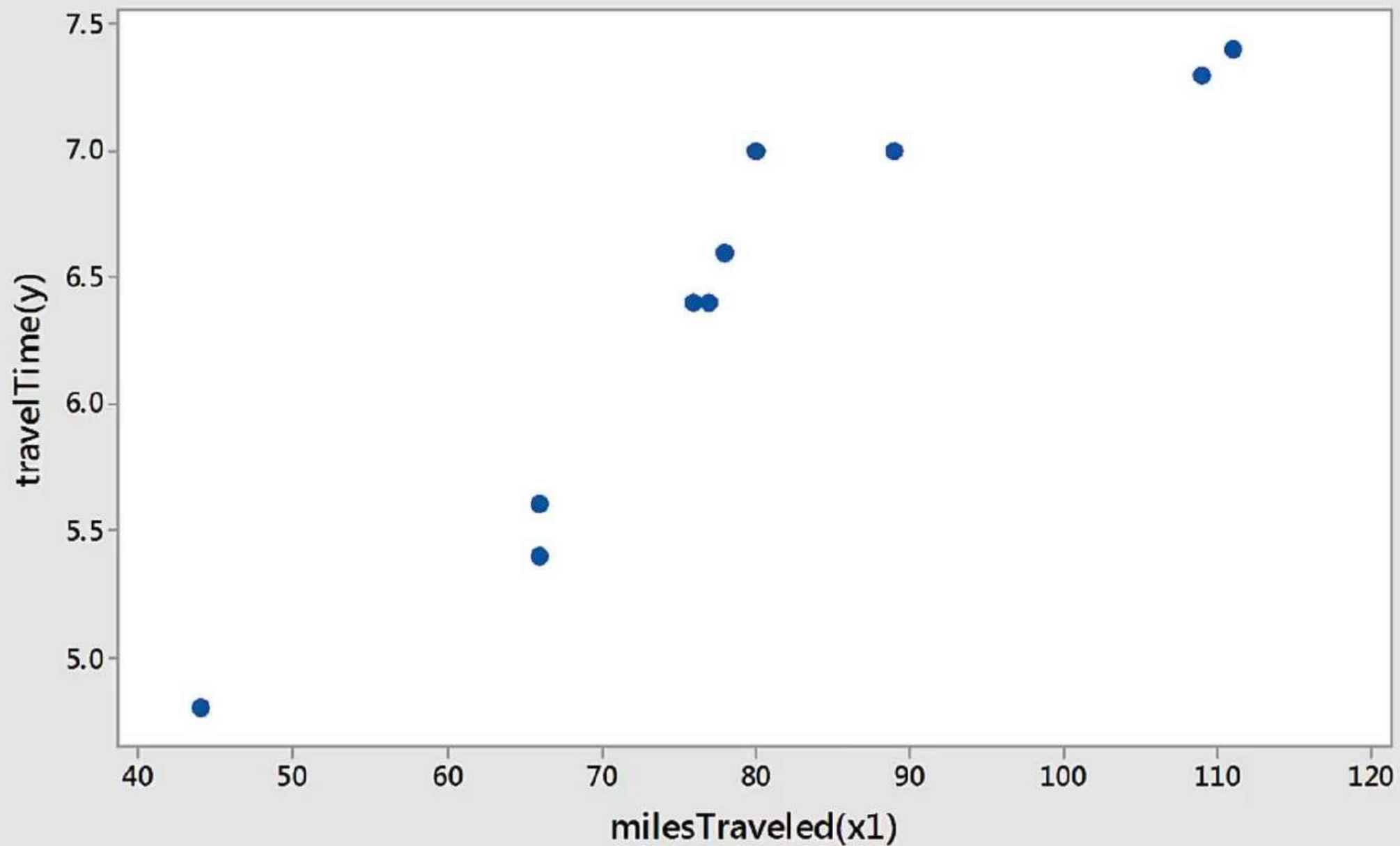
Multiple regression
many-to-one

IV TO DV SCATTERPLOTS

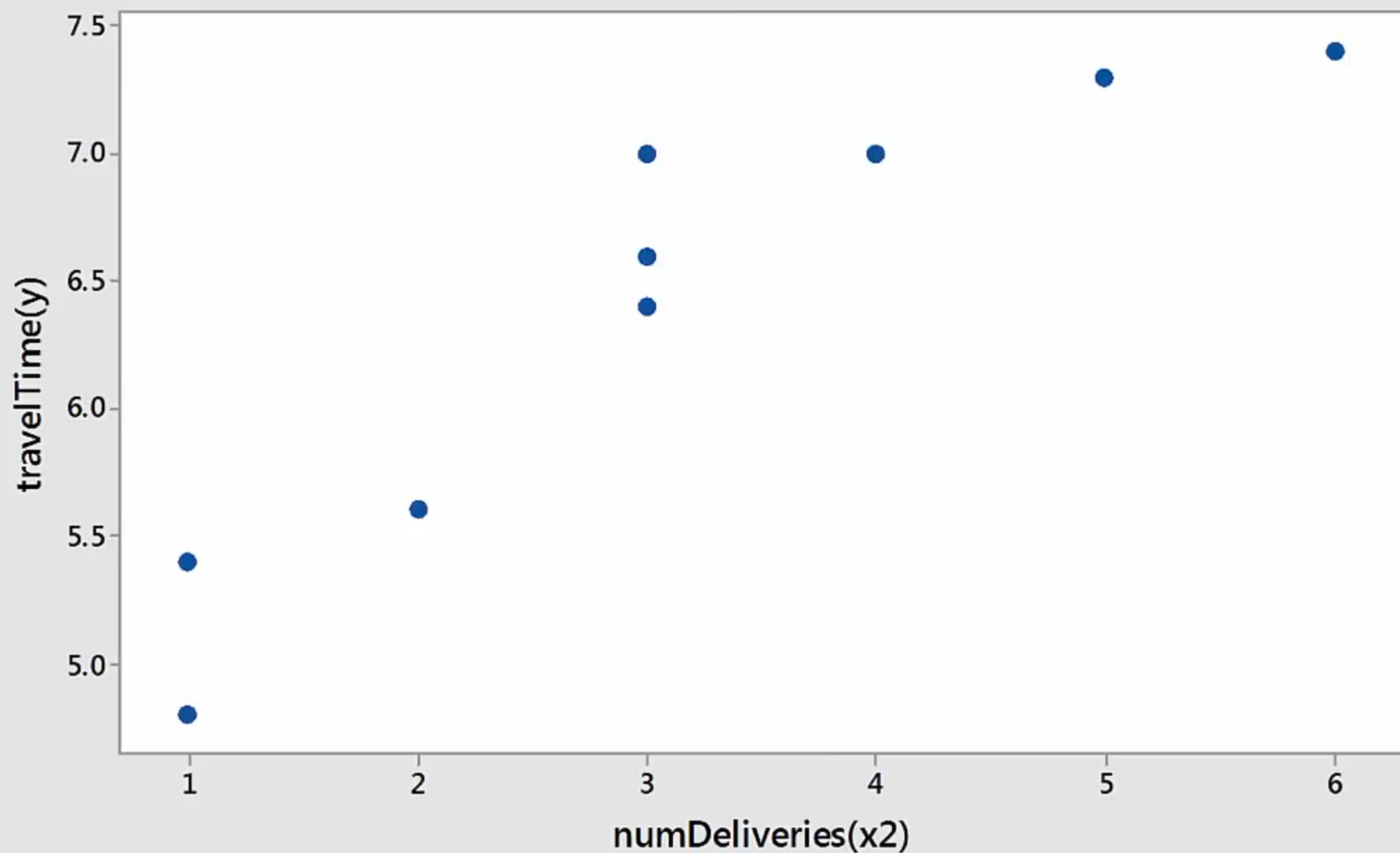
 [Subscribe](#)

Relevancy
Check

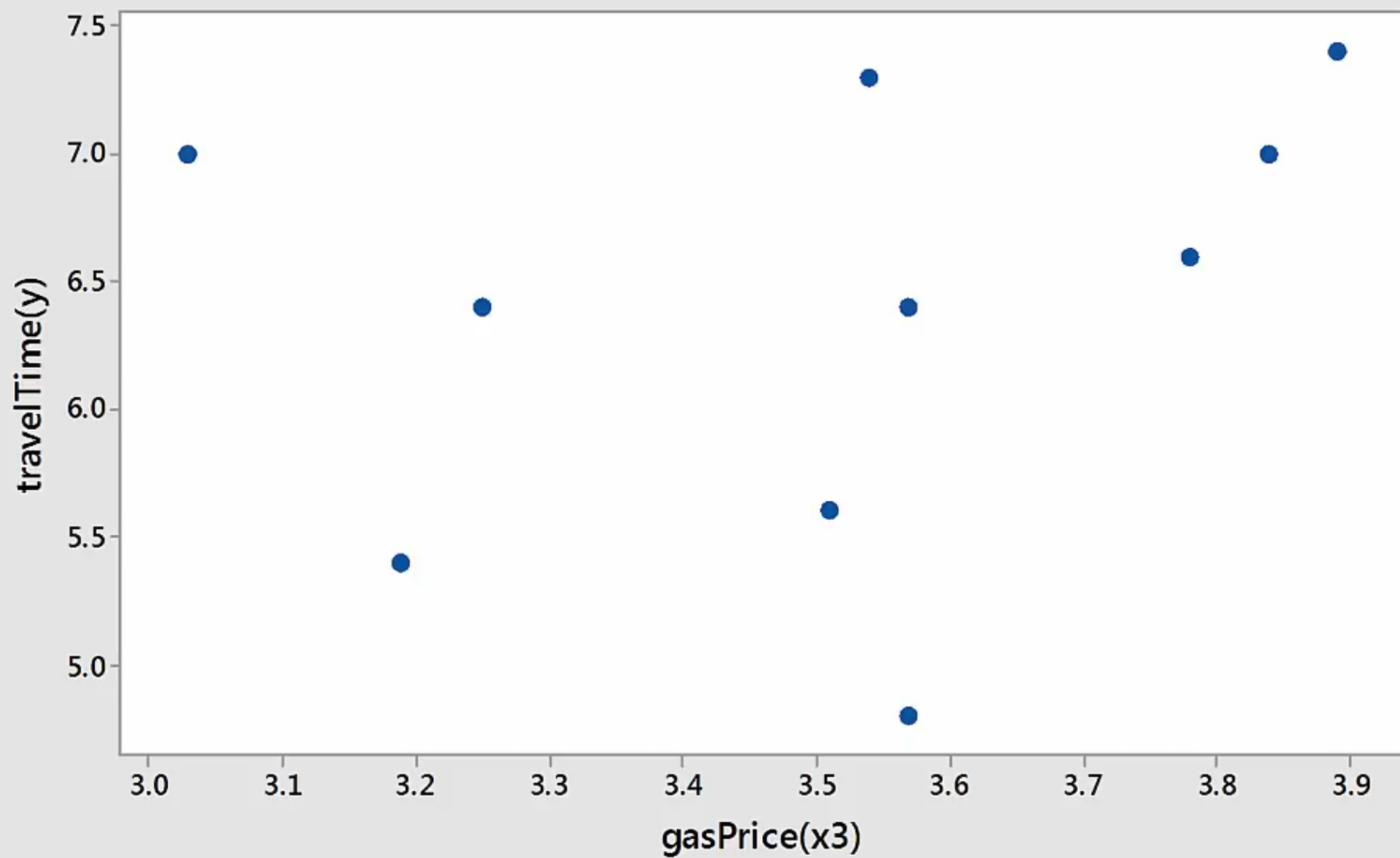
Scatterplot of travelTime(y) vs milesTraveled(x1)



Scatterplot of travelTime(y) vs numDeliveries(x2)

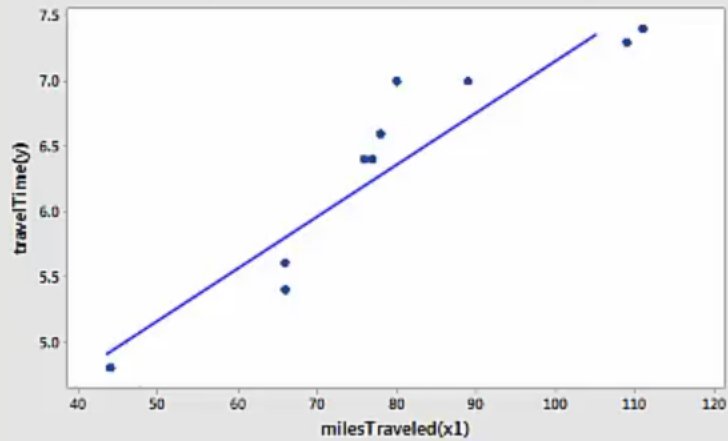


Scatterplot of travelTime(y) vs gasPrice(x3)

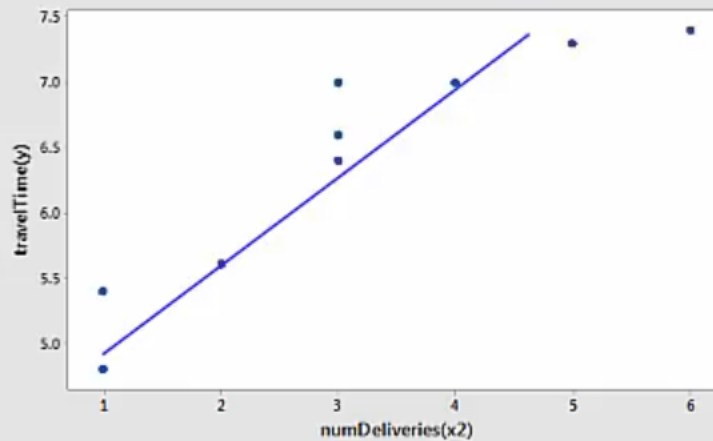


DV VS IV SCATTERPLOTS

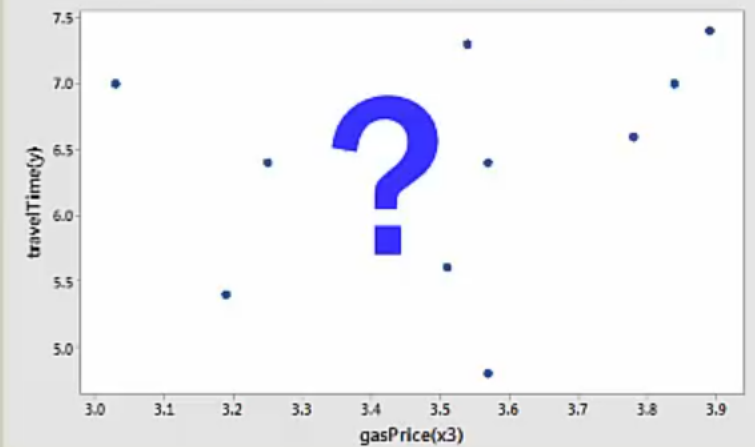
Scatterplot of travelTime(y) vs milesTraveled(x1)



Scatterplot of travelTime(y) vs numDeliveries(x2)

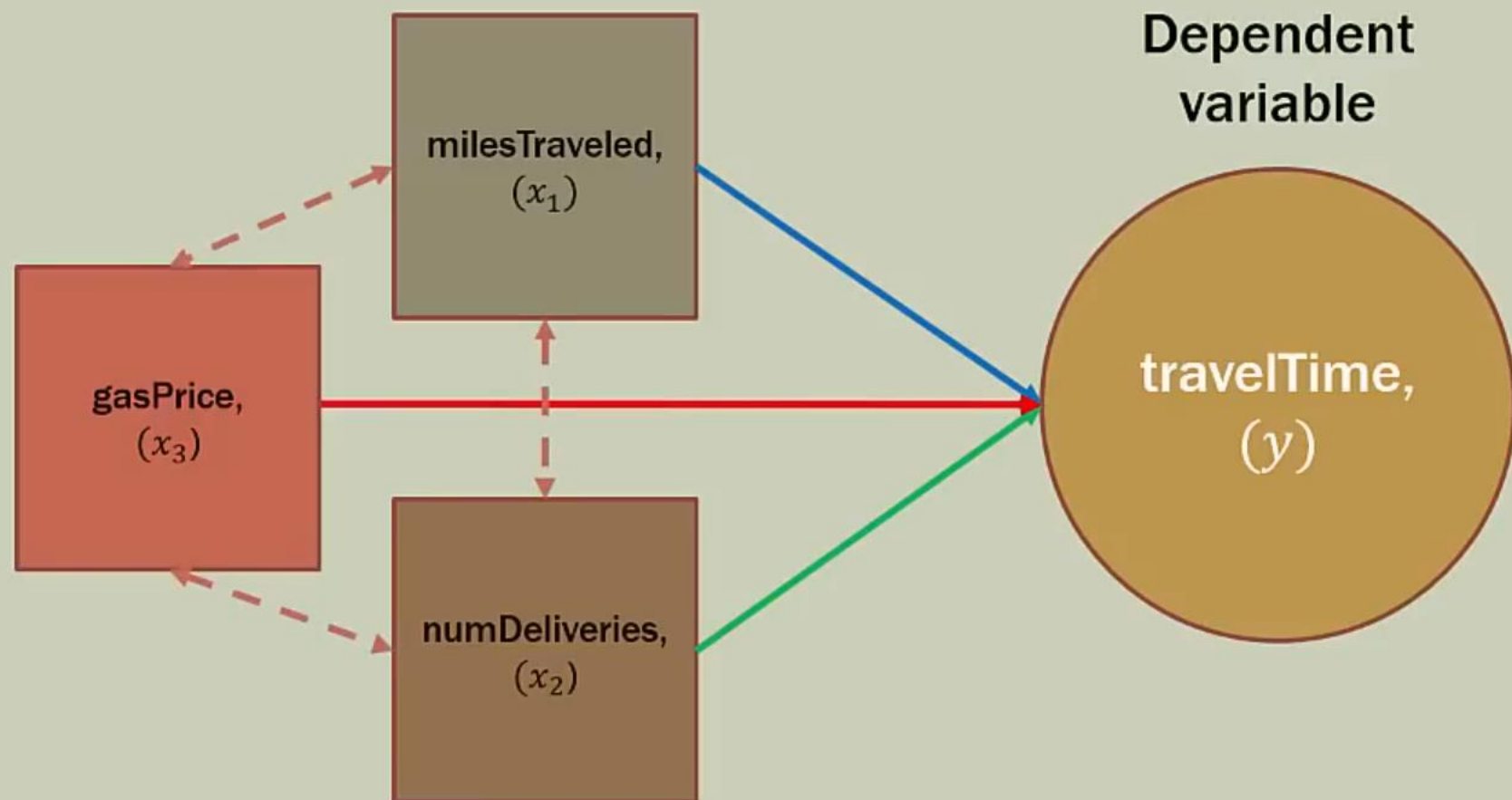


Scatterplot of travelTime(y) vs gasPrice(x3)



SKETCHING OUT RELATIONSHIPS

Independent variables



Dependent variable

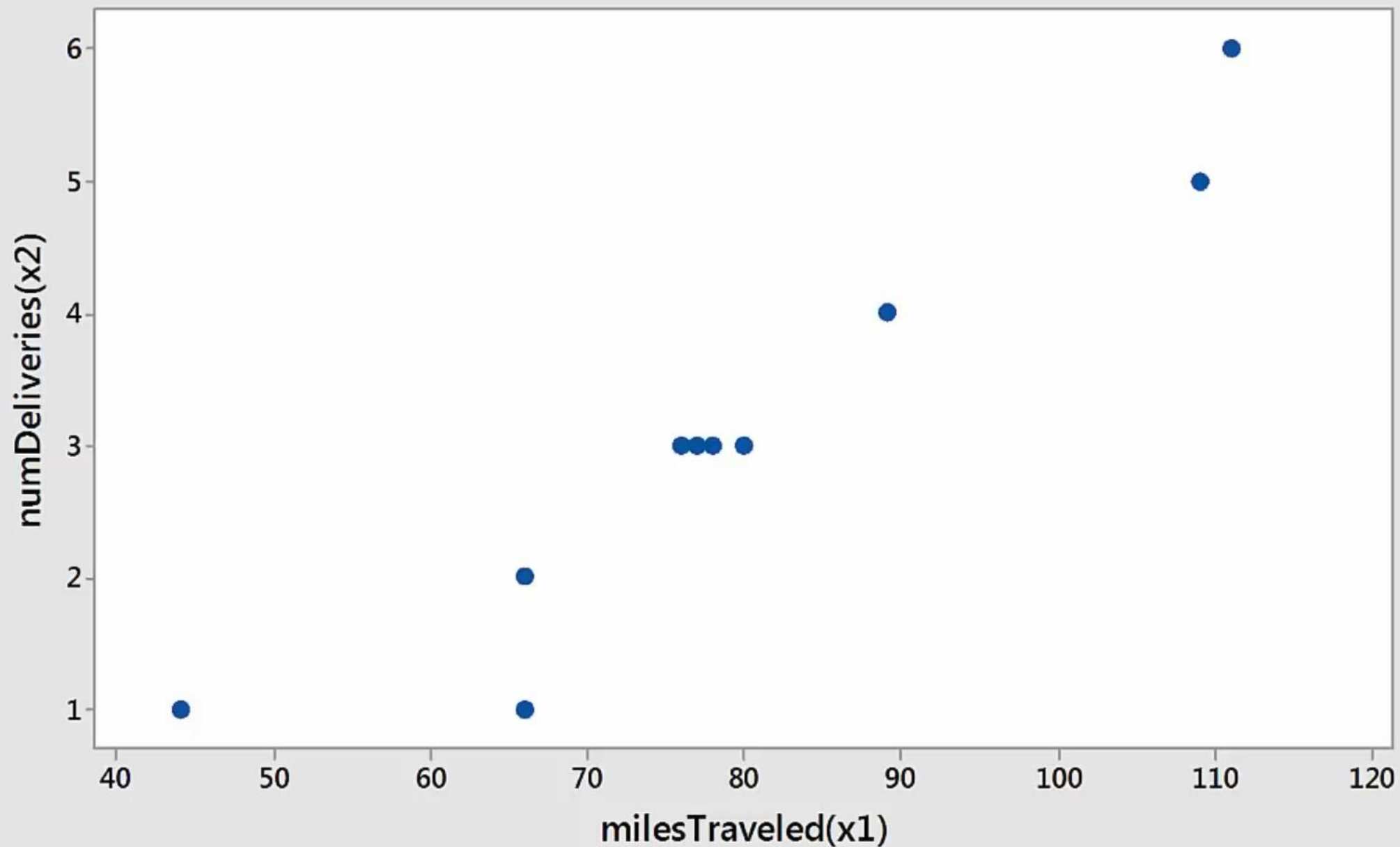
Multiple regression
many-to-one

IV TO IV SCATTERPLOTS

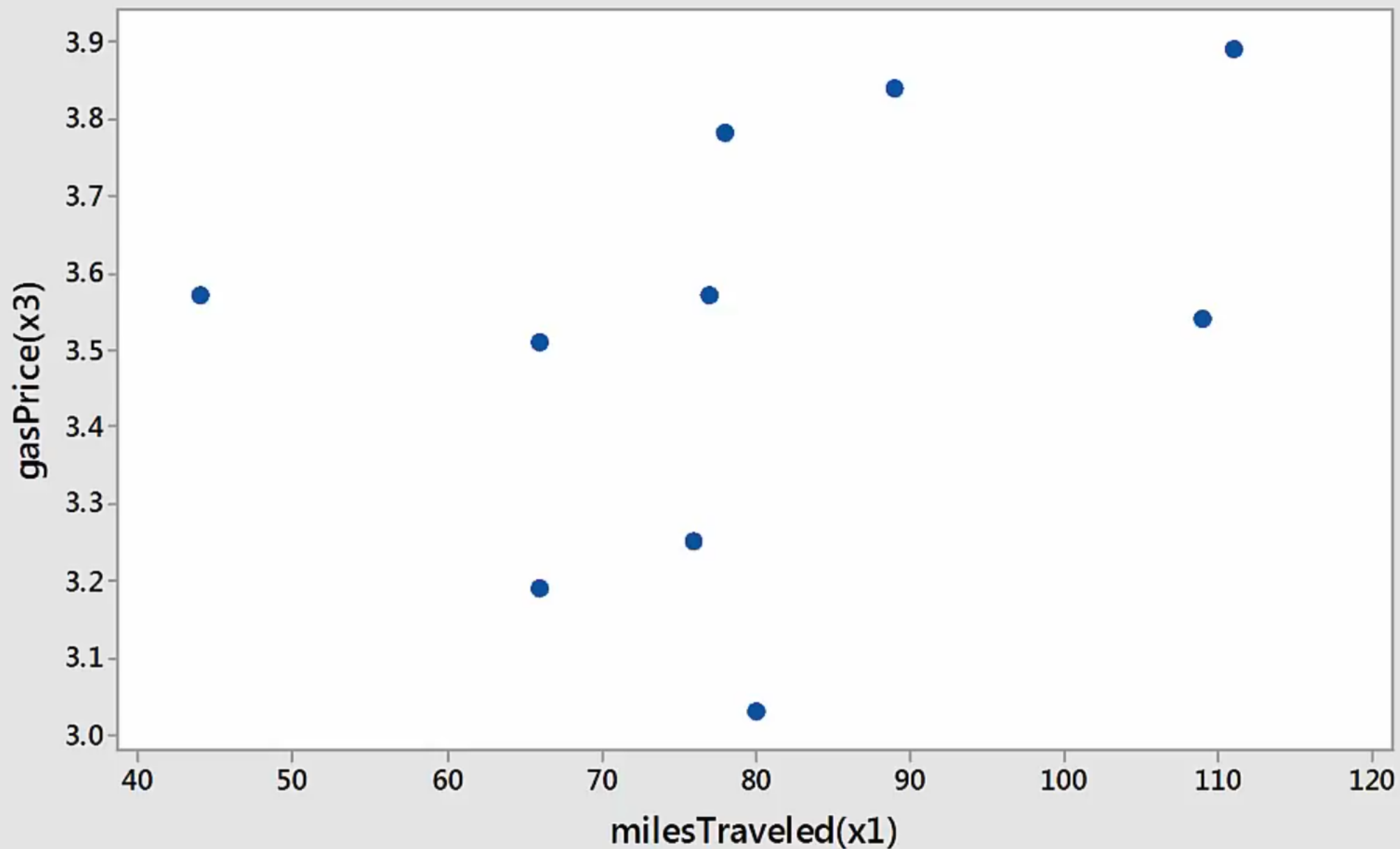
 [Subscribe](#)

Multicollinearity
Check

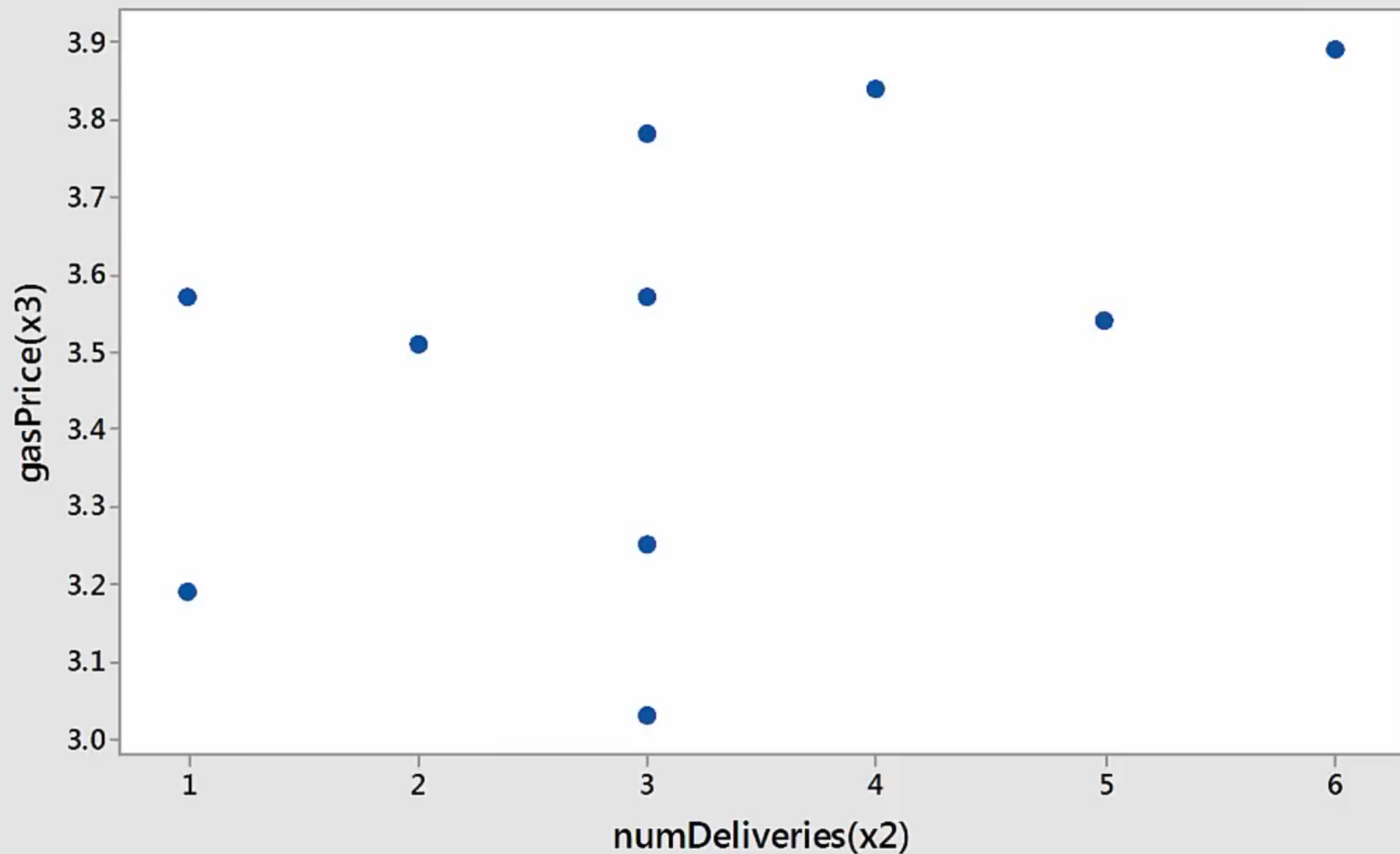
Scatterplot of numDeliveries(x2) vs milesTraveled(x1)



Scatterplot of gasPrice(x3) vs milesTraveled(x1)



Scatterplot of gasPrice(x3) vs numDeliveries(x2)



IV SCATTERPLOTS (MULTICOLLINEARITY)

