# PSTAT 126: Regression Analysis
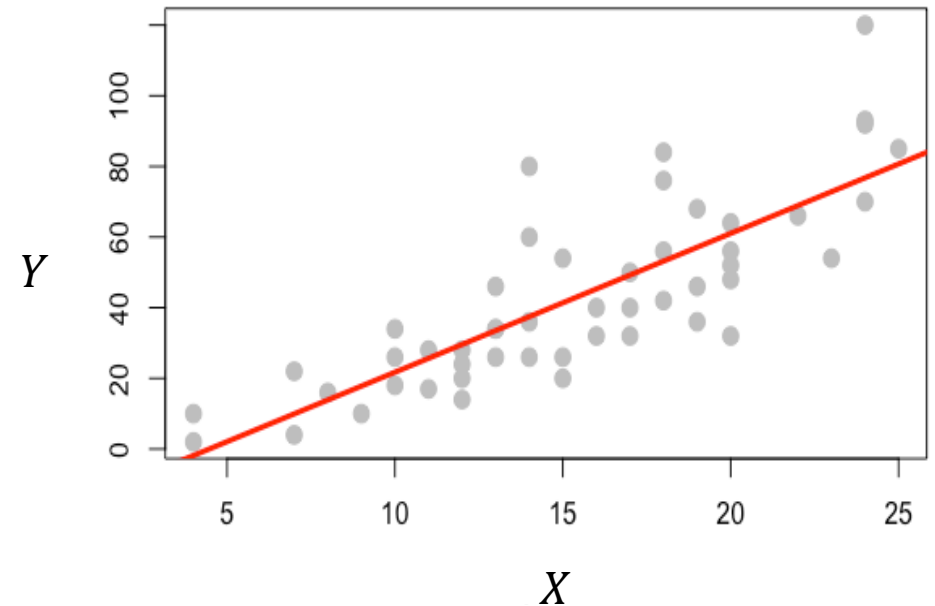# Department of Statistics and Applied Probability
# University of California, Santa Barbara

# Regression and its Applications

There are many areas of human endeavor in which we would like to learn and model, from relevant but noisy data, an unknown functional relationship between a variable $X$ (or variables) and a variable $Y$, the values of which we think of as dependent, in some sense, on those of $X$. The ability to do this has key applications in such areas as, among others:

- Science & Medicine

- Technology & Industry

- Economics & Finance

- Sociology & Behavioral Sciences

- Public Policy

The study of how best to do this, including which mathematical and statistical methods and algorithms to use, is the subject of **Regression**.
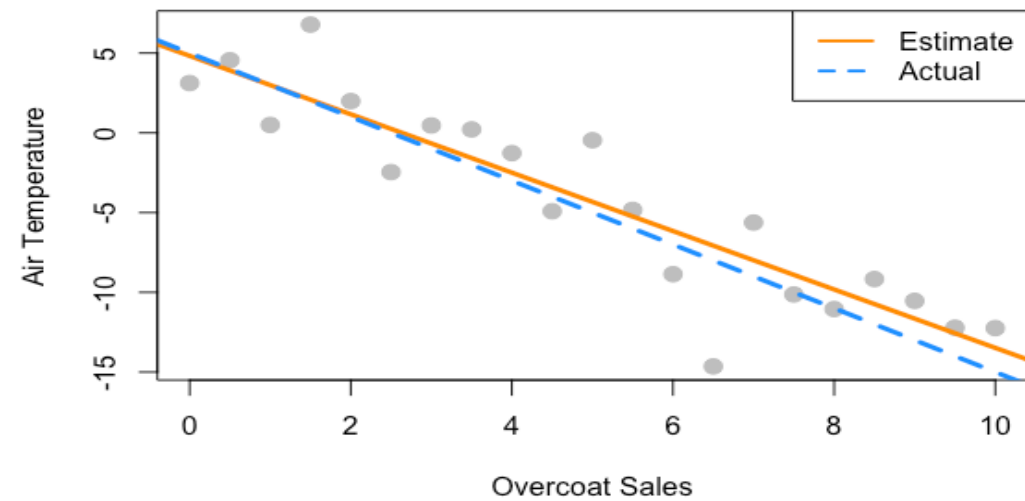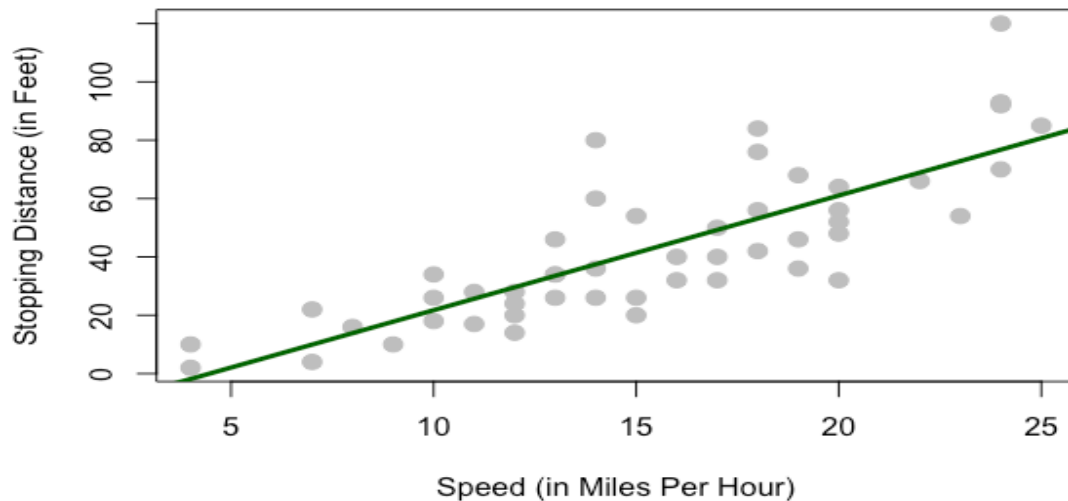
# Uses of Regression: Explanation and Insight

**Explanation and insight**:

Modeling the relationship between an input/inputs and an outcome, given observed, sampled data, in order to gain deeper understanding into that relationship.

What is the functional relationship between the stopping distance of a car (that is, the safe stopping distance, without the driver's loss of control) and the car's speed?



The graphic shows an example of linear regression – regression for which the functional relationship between X and Y is, or is presumed to be, linear in an appropriate sense.

# Uses of Regression: Prediction

**Prediction**:
Given a new input value, not previously sampled, estimate the corresponding outcome/output value using the trained regression model.

Given one's high school and/or college GPA, can SAT and/or GRE scores be predicted?

# History of Regression

- The mathematicians Legendre (1805) and Gauss (1809) were the first known to have used the technique of statistical regression (that is, the method of least squares) as such, in order to find the best linear fit to a finite set of data points.

- They applied the method to analyze and predict planetary motion.

- Using the normal (or Gaussian) distribution to describe the behavior of errors, Gauss also developed a formula for this distribution, which plays such an important role in modeling errors in (linear) regression.

- Techniques for Linear Regression can rightly be viewed as Artificial Intelligence/Machine Learning methods and indeed as, historically speaking, perhaps the original versions of the types of Machine Learning algorithms so widely used today.
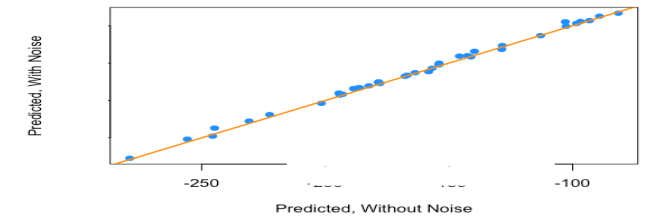


*NASA / Johnson Space Center*

# Goals of Regression

Let P be any population. This population could be virtually any set of objects of interest, including people, cities, companies, biological cells, or stars in the night sky, for example. For any given population, we may be interested in the relationship between two variables of interest, a so-called predictor variable $X$ -- also called the explanatory or independent variable -- and a response variable Y (also known as the dependent or target variable). For example, $X$ and $Y$ could be the respective

- Height and weight of people in P

- Distance from Earth of a set of stars and their corresponding brightness

- Education level and average income in the population of a given city.

In order to understand and explain the interaction between $X$ and $Y$, which we think of as random variables, we would like to find an approximate functional relationship $f(X) \approx Y$ between them. Note that, for us, the function $f$ we will attempt to learn will be assumed deterministic (non-random), and we will have

$$Y = f(X) + \epsilon,$$

with the noise term $\epsilon$, also a random variable and such that the conditional expectation $E[\epsilon|X = x] = 0,$ representing random error or variation in the model.
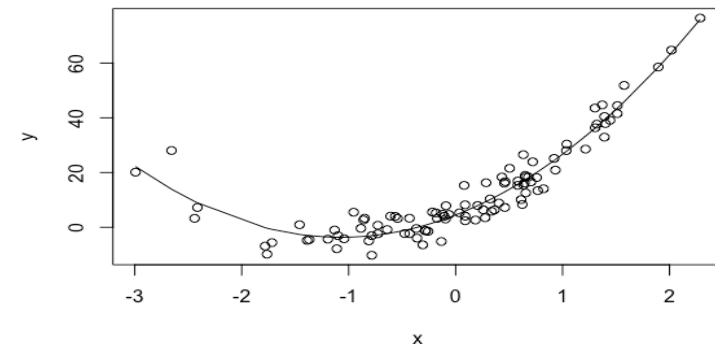
- We are essentially always interested in determining $f$ in the context of regression models, but interest in determining more about the random error $\epsilon$ may depend on context.

- Indeed, for the purpose of explanation and insight concerning the relationship between $X$ and $Y$, more information about the nature of $\epsilon$, including its variance, may be of significant interest, whereas, when applying the model expressly for prediction, additional information about $\epsilon$ may be of less value.

# Regression and the Mean Function

To determine a functional relationship between predictor $X$ and response $Y$, our goal is to learn the conditional expectation function $E[Y|X]$ – or, at least, a reasonably close approximation of it. We call $E[Y|X]$ the **regression** or **mean function**.

**Why is the mean function $E[Y|X]$ so important here?**
It clearly gives you the mean value of $Y$ given $X=x$. But we can go further than this. We want to find a a function minimizing the difference between $f(X)$ and $Y$, on average. So, this would suggest looking at the absolute value of the difference $f(X) - Y$, i.e., $|f(X) - Y|$, and then considering the mean or expectation $E[|f(X) - Y|]$. However, in part because the absolute value function is not smooth as it is not differentiable at 0 (spaces of functions defined by the square having other nice mathematical properties as well), it is more convenient to consider $E[(f(X) - Y)^2]$.



$E[Y|X]$ is the function that minimizes this squared error among all candidate functions $f$.
In fact it can be shown that

$$E[(f(X) - Y)^2] = E[(f(X) - E[Y|X])^2] + E[(Y - E[Y|X])^2], \qquad (1)$$

for any candidate function $f$, where $E[(Y - E[Y|X])^2]$ depends on $X$ and $Y$ but not $f$. Equation (1) holds whether $X$ is a scalar or vector-valued variable. Equation (1) says that that, for any function $f$, the expectation of the square of the difference between $f(X)$ and $Y$ is equal to the expectation of the square of the difference between $f$ and the mean function (plus a nonnegative constant, as shown in (1)).

- Since $(f(X) - Y)^2 \geq 0$ for any function $f$ we can minimize the magnitude of the error of approximating $f$ by Y on the left-hand side of (1) by in fact taking $f(X) = E[Y|X]$.
- This means that the function of $X=x$ that approximates the behavior of the response $Y$ with the smallest error on average is in fact the mean $E[Y|X]$ function itself.
- So, it is the mean function which gives us the "best" representation of the functional relationship between $X$ and $Y$ in the sense described.

Hence, it is the mean function $E[Y|X]$ that we would like to use regression methods and algorithms to determine or at least closely approximate in order to identify and understand any functional relationship between $X$ and $Y$.

# Linear Regression

Our goal in this course is to study specifically **Linear Regression**, which is regression for which $E[Y|X]$ is or may be presumed to be closely approximated by a linear function (i.e., more technically, a function selected from a finite-dimensional, linear space of candidate functions).

The linear case is of great interest because

- from the point-of-view of mathematical structure, it is relatively simple (shades of Occam's razor)

- it robustly describes many situations arising in applications

- it is the model base case for investigations into nonlinear regression (indeed, somewhat paradoxically, the linear regression model itself encompasses many seemingly "nonlinear" cases as well, as we shall see).

So, for the first part of the course we will be considering models of the relatively simple form

$$E[Y|X = x] = \beta_0 + \beta_1 x, \qquad\qquad (2)$$

where $x$ is a fixed, scalar value (real number), and $Y$ is a scalar-valued continuous random variable. The numbers $\beta_0, \beta_1$ are parameters which, as we shall see, it is the goal of canonical regression algorithms to compute. When the regression function can be represented as in (2) it is called **Simple Linear Regression** (see the next slide) because only one predictor variable is involved and the predictor appears within a linear term only. Later on, we will augment this framework by adding additional predictor variable terms on the right in (2). This is called **Multiple Linear Regression.** Note that any representation of the function $E[Y|X]$ in the form as on the RHS of (2) will be unique for either simple -- or multiple – regression (at least for the kinds of typical continuous probability distributions we are interested in in this course).

# Simple Linear Regression (SLR) Model

But what are the methods of regression that enable us to determine the parameters $\beta_0$ and $\beta_1$ (or close approximations of these parameters), given that in general we have no ready or direct access to the actual values of the function $E[Y|X]$?

The answer of course involves sampling. For this, let $x_1, x_2, \ldots, x_N$ be $N$ given fixed, real numbers. We could think of these numbers as sampled from the predictor $X$, but, in keeping with what seems to be fairly standard expository practice in textbooks on basic regression, we usually downplay or suppress the explicit role of the underlying variable $X$. Now, given these $N$ values $x_n, n = 1, \ldots, N$, write

$$Y_n = \beta_0 + \beta_1 x_n + \epsilon_n, n = 1, \ldots, N, \qquad (3)$$

$$\epsilon_n \sim N(0, \sigma^2), n = 1, \ldots, N. \qquad (4)$$

Here, the $\epsilon_n$ are $N$ independent, real-valued, normally-distributed random variables (i.i.d.), with $N(0, \sigma^2)$ being the normal (Gaussian) distribution with mean 0 and variance $\sigma^2$. The $\epsilon_n$ represent random variation or noise in the model, and we shall have more to say later about our assumptions concerning the $\epsilon_n$. We call (3)-(4) are our **Simple Linear Regression (SLR) Model**. The goal of regression is it to identify the scalar parameters $\beta_0$ and $\beta_1$ and also, often, $\sigma$ as well, or, more commonly, close approximations of these three parameters. The SLR model above in (3)-(4) is the formal model we will now generally work with until we get to Multiple Linear Regression.

In (3)-(4) we assume, as already noted, that each $x_n$ is a known constant (say the outcome of an experiment after the $n$th trial). So, for each $n$, we actually can write

$$E[Y_n] = E[Y_n|X = x_n] = \beta_0 + \beta_1 x_n. \qquad (5)$$

# Simple Linear Regression Model (cont'd)

Our SLR model: Given $N$ values $x_n, n = 1, \ldots, N$, write

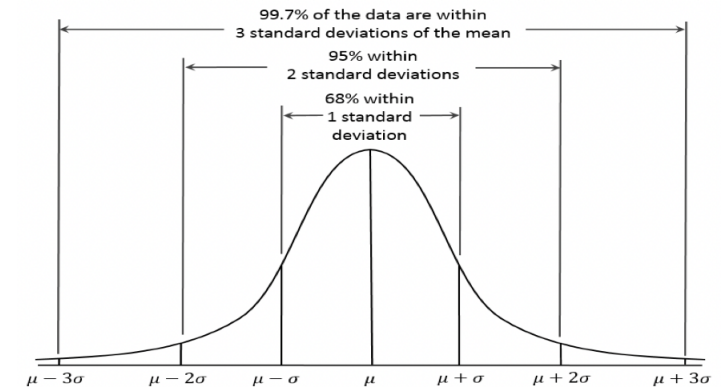$$Y_n = \beta_0 + \beta_1 x_n + \epsilon_n, n = 1, \ldots, N, \qquad (6)$$

$$\epsilon_n \sim N(0, \sigma^2), n = 1, \ldots, N, \qquad (7)$$

the $\epsilon_n$ being $N$ independent, normally-distributed random variables (i.i.d.), with $N(0, \sigma^2)$ being the normal distribution with mean 0 and variance $\sigma^2$. So independence of the $\epsilon_n$ for us means *mutual independence* so that the corresponding joint and respective individual probability density functions satisfy

$$f_{\epsilon_1, \ldots, \epsilon_N}(z_1, \ldots, z_N) = f_{\epsilon_1}(z_1) \ldots f_{\epsilon_N}(z_N), \qquad (8)$$

where

$$f_{\epsilon_n}(z) = N(0, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{z^2}{2\sigma^2}\right), \text{ for each } n. \qquad (9)$$



Normal Distribution

The $Y_n$ satisfy similar conditions but with different means. Note that the error $\epsilon_n$ is distributed symmetrically about $E[Y_n|X = x_n] = \beta_0 + \beta_1 x_n$. We also note that the i.i.d. assumption is, while a common assumption, a strong assumption and its full strength is not always necessary in the context of regression analysis as we study in this course.

# First Steps with R

At this point, let's see how the R language can be applied in the context of an actual data set to generate a simple linear regression model.
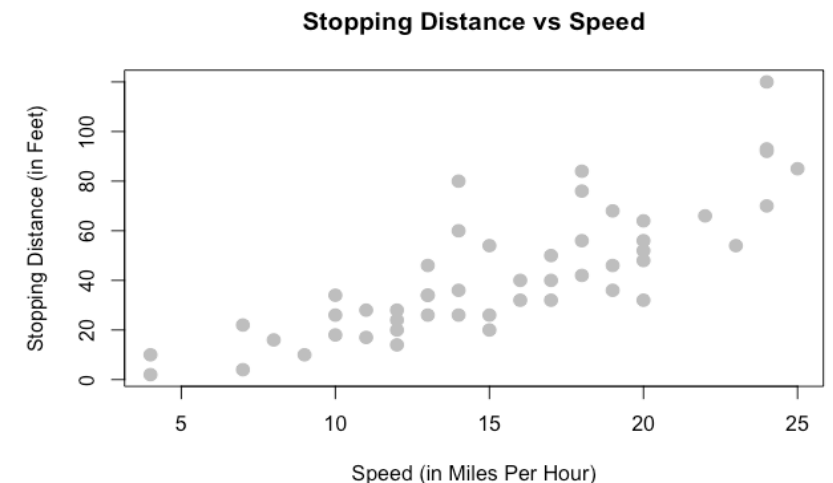
R is a language and environment for statistical computing and graphics, an integrated suite of software facilities for data manipulation. R is similar to the S language and environment which was developed at Bell Laboratories (formerly AT&T, now Lucent Technologies) by John Chambers and colleagues. R can be considered as a different implementation of S.

R provides a wide variety of statistical (linear and nonlinear modelling, classical statistical tests, time-series analysis, classification, clustering, etc.) and graphical techniques, and is highly extensible. The S language is often the vehicle of choice for research in statistical methodology, and R provides an Open Source route to participation in that activity.

We use the "cars" data set, which is "built-in" to R. It contains data gathered during the 1920s about the speed of cars and the resulting distance it takes for the car to safely come to a stop, without loss of vehicle control.
Thinking of Speed as our predictor variable X and Stopping Distance as our Response Y, we can plot the stopping distance against the speed using the R code below.

```
plot(dist ~ speed, data = cars,
    xlab = "Speed (in Miles Per Hour)",
    ylab = "Stopping Distance (in Feet)",
    main = "Stopping Distance vs Speed",
    pch  = 20,
    cex  = 2,
    col  = "grey")
```
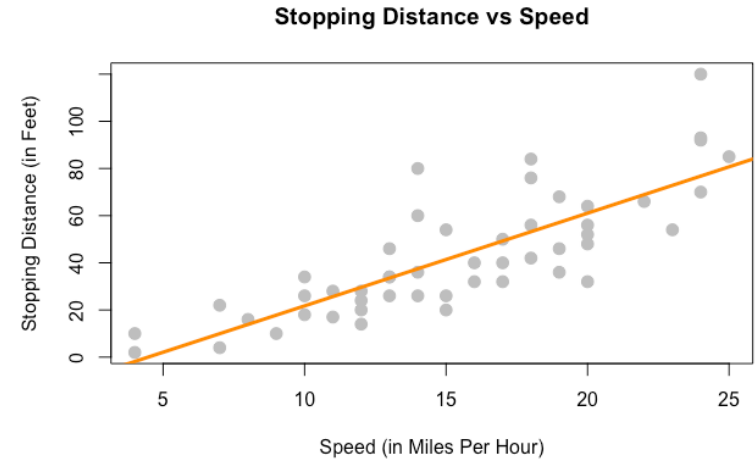


Stopping Distance vs Speed

# First Steps with R (cont'd)

stop_dist_model = **lm**(dist ~ speed, data = cars)
  stop_dist_model
## Call:
## lm(formula = dist ~ speed, data = cars)
## Coefficients:
## (Intercept)    speed
## -17.579        3.932

**Stopping Distance vs Speed**



In order to compute the regression function (regression line) for the cars example we use the lm( ) function in R. The initials stand for "linear model", and it will be perhaps our most commonly used R function in this course. We will concern ourselves with how estimates of the model parameters are computed in forthcoming slides, but for now note that R gives

$$\beta_0 = Intercept \approx -17.579$$
$$\beta_1 = Slope \approx 3.932$$

```
plot(dist ~ speed, data = cars,
    xlab = "Speed (in Miles Per Hour)",
    ylab = "Stopping Distance (in Feet)",
    main = "Stopping Distance vs Speed",
    pch  = 20,
    cex  = 2,
    col  = "grey")
abline(stop_dist_model, lwd = 3, col = "darkorange")
```