# L22: Linear Regression

## Part 1: $y = ax$ and $y = ax + b$

Lucas A. J. Bastien

E7 Spring 2017, University of California at Berkeley

March 13, 2017

Version: release

## Announcements

**Lab 08 is due on March 17 at 12 pm (noon)**

**Today:**

- Linear regression
    - Introduction
    - Specific cases:
        - $y = ax$
        - $y = ax + b$
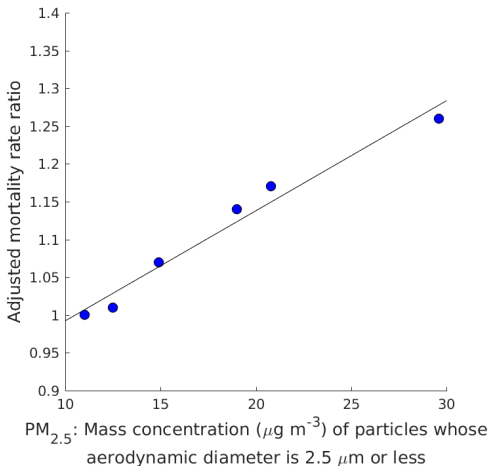        - "Transform" $y = bx^a$ into a linear form

**Wednesday:**

- Linear regression, continued
    - General approach

**Wednesday:**

- Linear regression: discussion and applications
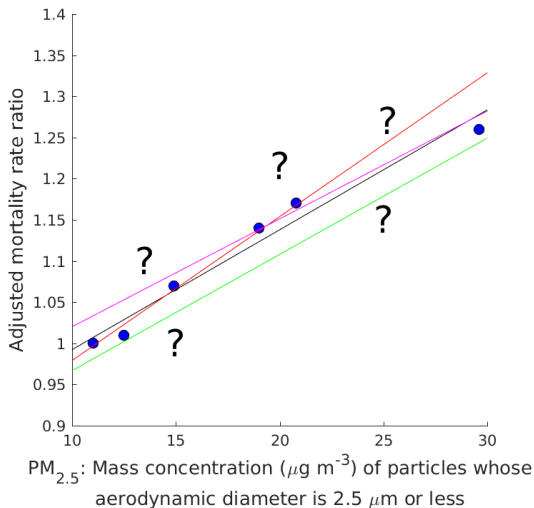
# Introduction to linear regression

**Harvard six cities study: It appears that there is a correlation between exposure to particulate matter pollution and premature mortality.** We can fit a straight line to the data reasonably well



PM$_{2.5}$: Mass concentration ($\mu$g m$^{-3}$) of particles whose aerodynamic diameter is 2.5 $\mu$m or less

# Introduction to linear regression

**How to choose the line that "fits best"?**



Adjusted mortality rate ratio vs. $PM_{2.5}$: Mass concentration ($\mu$g m$^{-3}$) of particles whose aerodynamic diameter is 2.5 $\mu$m or less

## Linear regression: goals and motivation

Given two sets of data (x- and y-values), **can we model one set as a function of the other**, using a "simple function" (*e.g.*, straight line, polynomial)?

For a given function shape (*e.g.*, straight line, polynomial), **how to obtain the "best fit"?** For example:

▶ What is the best straight line that we can fit to the data?

▶ What is the best polynomial that we can fit to the data?

In the context of linear regression, **what do "simple function" and "best fit"** mean?

**Linear regression does not necessarily match all points exactly** (it does when it can, though)

How can we **measure how good or bad a fit** is?

# What type of "simple function" can be fitted?

**Answer:** any function of the form:

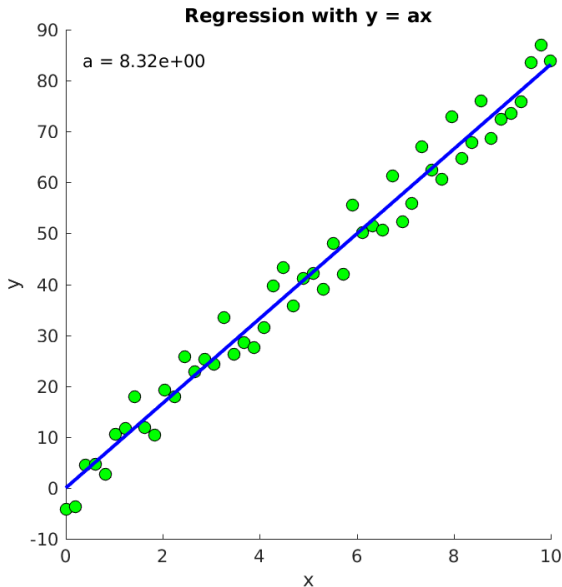$$y = a_1 f_1(x) + a_2 f_2(x) + \ldots a_n f_n(x)$$

where:

- The $a_i$'s are the coefficients to be determined using linear regression
- The $f_i$'s are real-valued functions
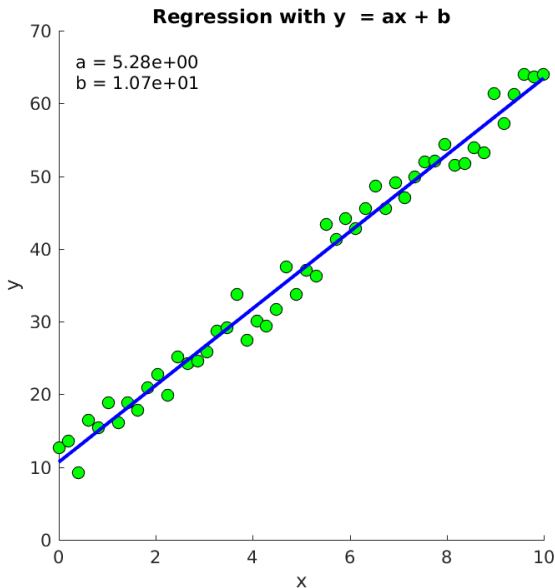
(*i.e.* any **linear** combination of known functions)

## Examples:

- $y = ax$                                                      (coefficient: $a$)
- $y = ax + b$                                       (coefficients: $a$ and $b$)
- $y = a_0 + a_1 x + a_2 x^2 + a_3 x^3$    (coefficients: $a_0$, $a_1$, $a_2$, and $a_3$)
- $y = a\cos(x) + b\sin(x) + c$         (coefficients: $a$, $b$, and $c$)
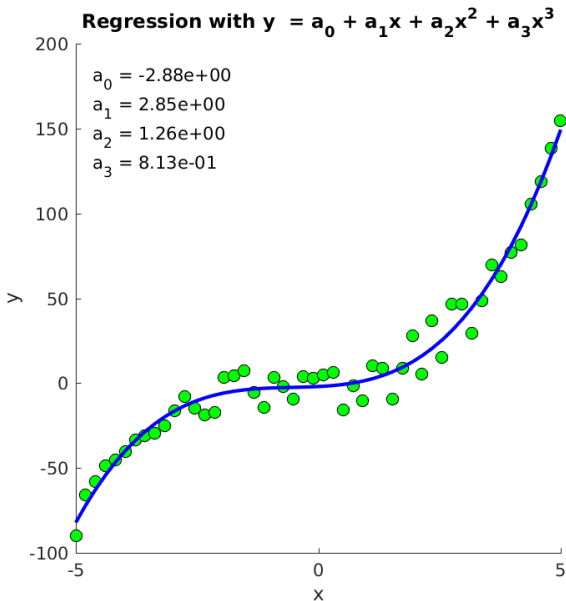
# Example of linear regression: $y = ax$



**Regression with y = ax**

$a = 8.32e+00$

# Example of linear regression: $y = ax + b$



**Regression with y = ax + b**

a = 5.28e+00
b = 1.07e+01

# Example of linear regression: $y = a_0 + a_1x + a_2x^2 + a_3x^3$



Regression with $y = a_0 + a_1x + a_2x^2 + a_3x^3$

$a_0$ = -2.88e+00
$a_1$ = 2.85e+00
$a_2$ = 1.26e+00
$a_3$ = 8.13e-01

Regression with y = a×cos(x) + b×sin(x) + c

a = 8.37e+00
b = 2.57e+00
c = -4.17e+00

## Today and Wednesday

**Today:**

- ▶ We fit lines of equation $y = ax$ and $y = ax + b$
- ▶ We derive the formulas for $a$ and $b$ by hand
- ▶ We calculate $a$ and $b$ "manually"

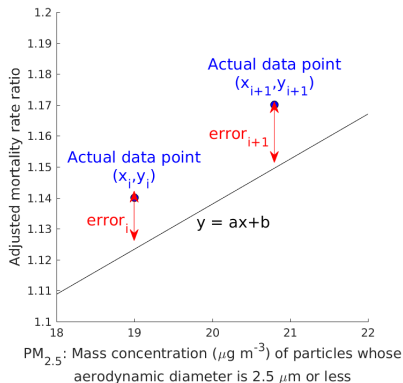**Next lecture (L23, Wednesday March 15):**

- ▶ We fit any line of equation $y = a_1 f_1(x) + a_2 f_2(x) + \ldots a_n f_n(x)$
- ▶ We let Matlab do most of the work

# Definition of "best fit" for linear regression

Consider that we have:

- A set of $x$- and $y$-data ($n$ data points)
- A set of predicted $y$ values

**How good is the fit?**



PM$_{2.5}$: Mass concentration ($\mu$g m$^{-3}$) of particles whose aerodynamic diameter is 2.5 $\mu$m or less

**Total squared error $E_2$:**
("square error" for short)

$$E_2 = \sum_{i=1}^{n} (\text{error}_i)^2$$

$$= \sum_{i=1}^{n} (y_{i,\text{predicted}} - y_i)^2$$

In this example:

$$E_2 = \sum_{i=1}^{n} (ax_i + b - y_i)^2$$

**Total squared error $E_2$:**
("square error" for short)

$$E_2 = \sum_{i=1}^{n} (\text{error}_i)^2$$

$$= \sum_{i=1}^{n} (y_{i,\text{predicted}} - y_i)^2$$

**Best fit:**

**(i.e. the best choice of parameters)**

**The fit that minimizes the total squared error**

# Fitting a line of equation $y = ax$: manual derivation

**Objective: given a set of $x$- and $y$-data, what is the value of $a$ such that the line of equation $y = ax$ is the best possible fit to the data?**

The total square error $E_2$ is a function of $a$:

$$E_2(a) = \sum_{i=1}^{n}(y_{i,\text{predicted}} - y_i)^2$$

$$= \sum_{i=1}^{n}(ax_i - y_i)^2$$

If $E_2$ is minimum for $a = a_{\min}$, then $E_2'(a_{\min}) = 0$

$$E_2'(a) = \sum_{i=1}^{n}2x_i(ax_i - y_i)$$

$$= \sum_{i=1}^{n}2ax_i^2 - 2x_iy_i$$

# Fitting a line of equation $y = ax$: manual derivation

$$E_2'(a) = \sum_{i=1}^{n} 2ax_i^2 - 2x_iy_i$$

Assuming that at least one of the $x_i$ is non-zero *i.e.* :

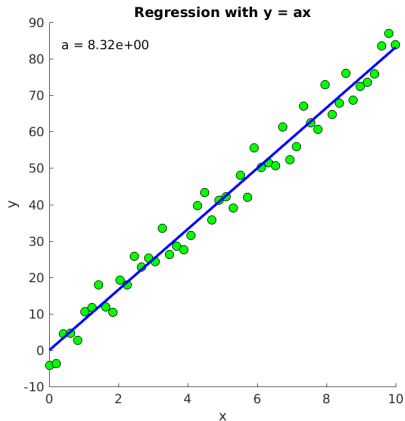$$\sum_{i=1}^{n} x_i^2 \neq 0$$

we have:

$$E_2'(a) = 0 \iff a = \frac{\sum_{i=1}^{n} x_iy_i}{\sum_{i=1}^{n} x_i^2}$$

**The value of $a$ that yields the best fit for the line $y = ax$ is:**

$$a = \frac{\sum_{i=1}^{n} x_iy_i}{\sum_{i=1}^{n} x_i^2}$$

# Fitting a line of equation $y = ax$: example



$$a = \frac{\sum_{i=1}^{n} x_i y_i}{\sum_{i=1}^{n} x_i^2}$$

```
% Determine and plot the linear regression line
a = sum(x.*y) / sum(x.*x);
plot(x, a*x, 'b', 'LineWidth', 2)
```

## Fitting a line of equation $y = ax + b$: manual derivation

**Objective: given a set of $x$- and $y$-data, what are the values of $a$ and $b$ such that the line of equation $y = ax + b$ is the best possible fit to the data?**

The total square error $E_2$ is a function of $a$ and $b$:

$$E_2(a, b) = \sum_{i=1}^{n} (y_{i,\text{predicted}} - y_i)^2$$
$$= \sum_{i=1}^{n} (ax_i + b - y_i)^2$$

If $E_2$ is minimum for $a = a_{\min}$ and $b = b_{\min}$, then

$$\left. \frac{\partial E_2}{\partial a} \right|_{(a_{\min}, b_{\min})} = 0 \quad \text{and} \quad \left. \frac{\partial E_2}{\partial b} \right|_{(a_{\min}, b_{\min})} = 0$$

# Fitting a line of equation $y = ax + b$: manual derivation

$$\frac{\partial E_2}{\partial a} = \sum_{i=1}^{n} 2x_i(ax_i + b - y_i)$$

$$\frac{\partial E_2}{\partial b} = \sum_{i=1}^{n} 2(ax_i + b - y_i)$$

Solve the following system of two equations and two unknowns ($a$ and $b$):

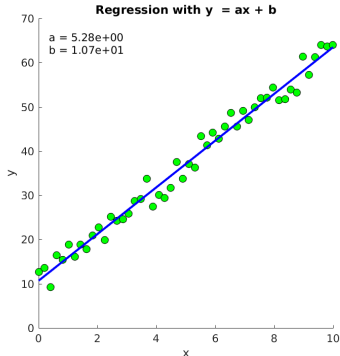$$\frac{\partial E_2}{\partial a} = 0 \quad \text{and} \quad \frac{\partial E_2}{\partial b} = 0$$

**to obtain the values of $a$ and $b$ that yield the best fit for the line $y = ax + b$:**

$$a = \frac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^{n}(x_i - \bar{x})^2} \quad \text{and} \quad b = \bar{y} - a\bar{x}$$

where $\bar{x}$ and $\bar{y}$ are the mean values of the $x$- and $y$-data, respectively

# Fitting a line of equation $y = ax + b$: example



$$a = \frac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^{n}(x_i - \bar{x})^2}$$

$$b = \bar{y} - a\bar{x}$$

```
% Determine and plot the linear regression line
x_mean = mean(x);
y_mean = mean(y);
a = sum((x-x_mean).*(y-y_mean)) / sum((x-x_mean).^2);
b = y_mean - a*x_mean;
plot(x, a*x+b, 'b', 'LineWidth', 2)
```

# Coefficient of determination

The coefficient of determination is often written $r^2$ or $R^2$ and pronounced "R squared". **It is another metric (in addition to the total square error $E_2$) that measures the goodness of fit**

$$r^2 = 1 - \text{fraction of variance unexplained by the regression model}$$
$$= 1 - \frac{\text{variance in } y \text{ not explained by the model}}{\text{variance in } y \text{ in the data}}$$
$$= 1 - \frac{\sum_{i=1}^{n}(y_i - y_{\text{predicted,i}})^2}{\sum_{i=1}^{n}(y_i - \bar{y})^2}$$

**Property:** $r^2 \leqslant 1$

**The closer $r^2$ is to 1, the better the fit**

$\bar{y}$ is the mean of the $y$-data

# Making problems linear: example with the power law

**Sometimes one wants to fit a line whose equation is not a linear combination of known functions**

**In some cases, it is possible to make the problem linear through mathematical manipulations**

**For example:** fit the line of equation $y = bx^a$ (power law)

$$y = bx^a$$
$$\ln(y) = \ln(b) + a\ln(x)$$
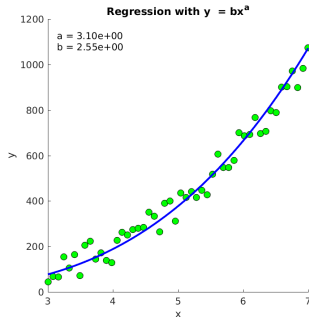$$Y = c + aX \qquad \text{with} \quad X = \ln(x)$$
$$Y = \ln(y)$$
$$c = \ln(b)$$

1. Calculate $X$ and $Y$ from $x$ and $y$
2. Use linear regression on the $X$- and $Y$- data to determine the coefficients $a$ and $c$
3. Calculate $b = \exp(c)$

# Making problems linear: example with the power law



$$a = \frac{\sum_{i=1}^{n}(X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^{n}(X_i - \bar{X})^2}$$

$$c = \bar{Y} - a\bar{X}$$

$$b = \exp(c)$$

```
% Determine and plot the linear regression line
logx = log(x);
logy = log(y);
logx_mean = mean(logx);
logy_mean = mean(logy);
a = sum((logx-logx_mean).*(logy-logy_mean)) / ...
    sum((logx-logx_mean).^2);
c = logy_mean - a*logx_mean;
b = exp(c);
plot(x, b*x.^a, 'b', 'LineWidth', 2)
```