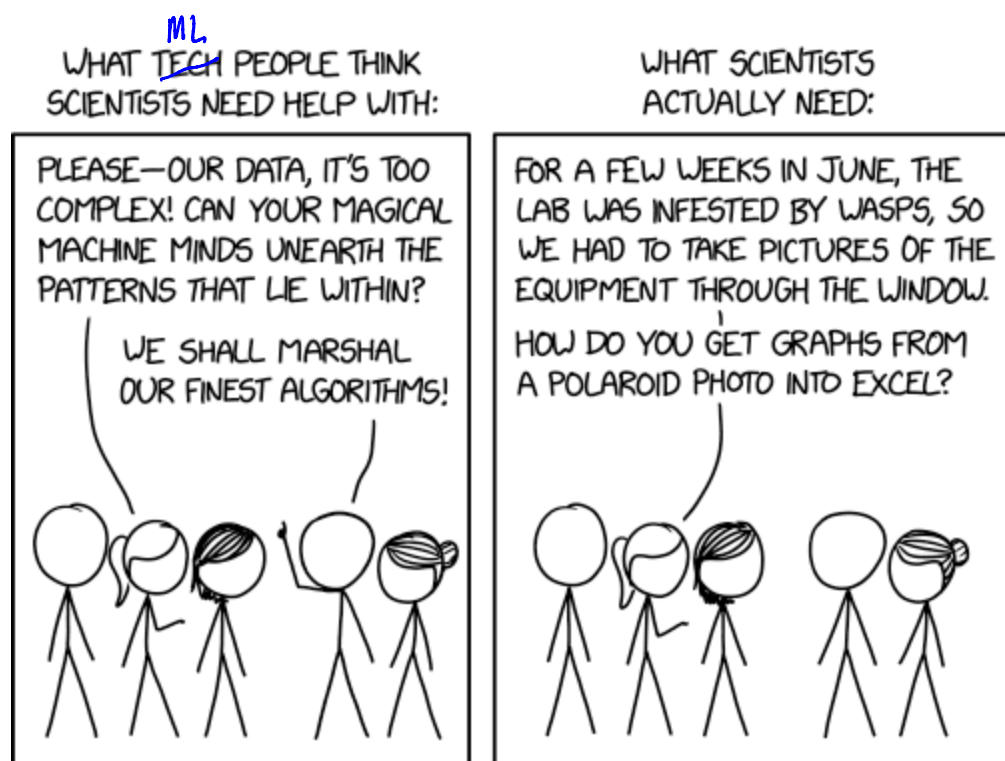


# Chapter 1: Introduction

*Statistical learning* refers to a vast set of tools for understanding data.



<https://xkcd.com/2341/>

**Alternative text:** I vaguely and irrationally resent how useful WebPlotDigitizer is.

These tools can broadly be thought of as

Supervised  
↓  
predicting or estimating  
an output based on  
one or more inputs.

or

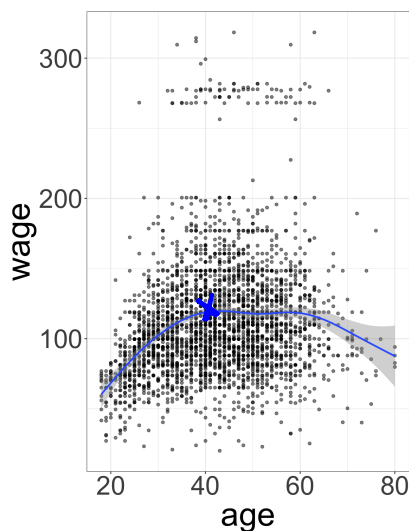
Unsupervised  
↓  
inputs w/ no supervising outputs  
can still learn about the structure of data.

Examples:

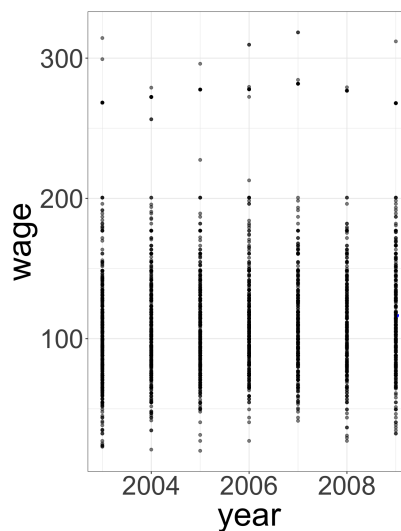
**Wage data**

year	age	maritl	race	edu- cation	region	job- class	health	health_ins	logwage	wage
2006	18	1. Never Mar- ried	1. White	1. < HS Grad	2. Mid- dle At- lantic	1. Indus- trial	1. <=Good	2. No	4.318063	75.04315
2004	24	1. Never Mar- ried	1. White	4. Col- lege Grad	2. Mid- dle At- lantic	2. Infor- ma- tion	2. >=Very Good	2. No	4.255273	70.47602
2003	45	2. Mar- ried	1. White	3. Some Col- lege	2. Mid- dle At- lantic	1. Indus- trial	1. <=Good	1. Yes	4.875061	130.98218

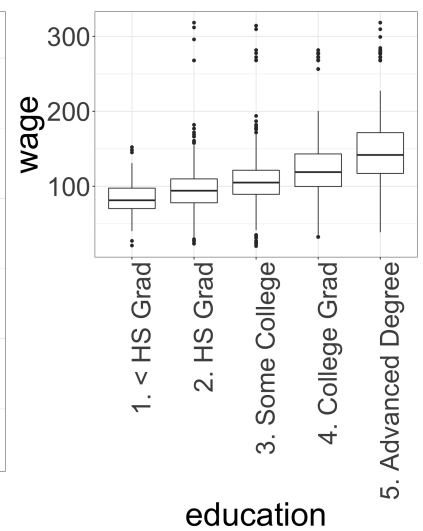
Factors related to wages for a group of males from the Atlantic region of the United States. We might be interested in the association between an employee's age, education, and the calendar year on his wage. *relationship*



*wage looks to increase w/age but then decreases after age 60.*



*slow but slight increase in wage over time. lot of variability.*



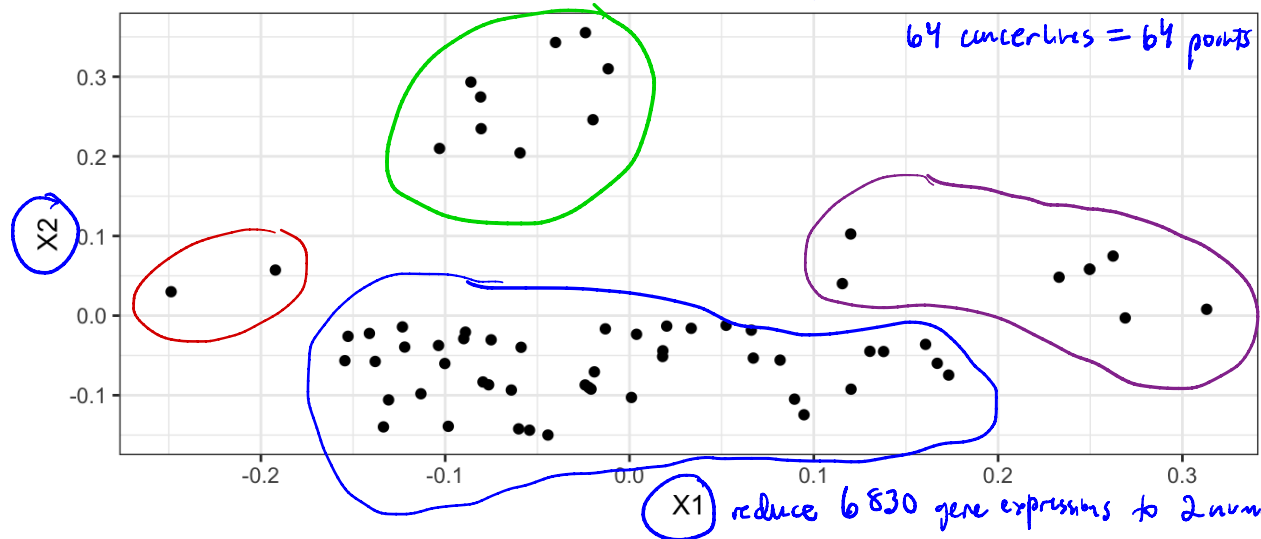
*wage typically higher for individuals w/ greater education levels.*

*could use 1 factor to predict wage, but lots of variability. would be better (more accurate) to combine age, education, & year and also account for non-linear relationship w/ age and wage.*

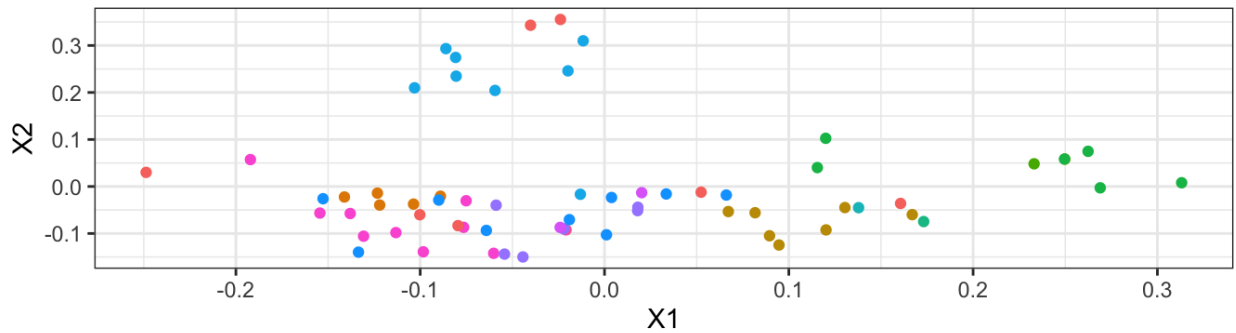
## Gene Expression Data

Consider the NCI60 data, which consists of 6,830 gene expression measurements for 64 cancer lines. We are interested in determining whether there are groups among the cell lines based on their gene expression measurements.

*we have no known output (cancer type), instead we want to look for structure in data*



*reduce 6830 gene expressions to 2 numbers "principal components" to describe the structure ("dimension reduction")*

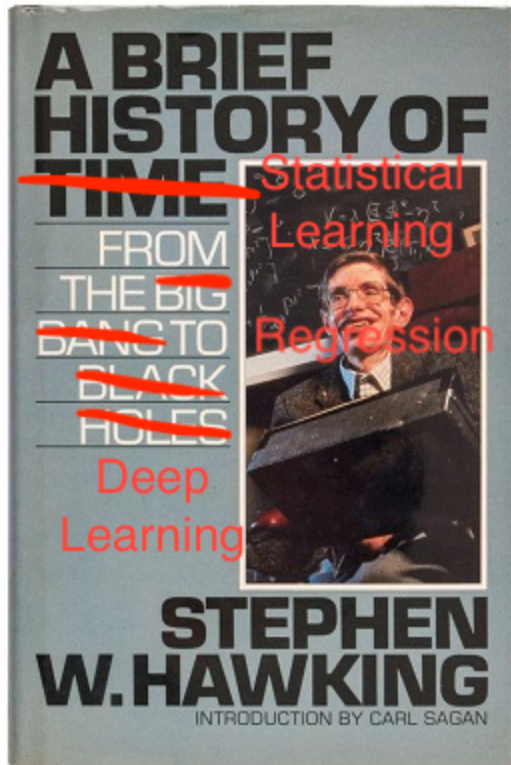


*"true" cancer types*

- |          |               |               |            |           |
|----------|---------------|---------------|------------|-----------|
| ● BREAST | ● K562A-repro | ● MCF7A-repro | ● NSCLC    | ● RENAL   |
| ● CNS    | ● K562B-repro | ● MCF7D-repro | ● OVARIAN  | ● UNKNOWN |
| ● COLON  | ● LEUKEMIA    | ● MELANOMA    | ● PROSTATE |           |

*cell lines w/ same cancer type are close in 2D representation and our clustering (top) was able to find some of these types*

# 1 A Brief History



Although the term “statistical machine learning” is fairly new, many of the concepts are not. Here are some highlights:

early 19th century - Legendre and Gauss publish method of least squares  $\Rightarrow$  linear regression

1936 - Linear discriminant analysis

1940s - Logistic regression

1960s - Bayesian Methods (1980s popularized)

1970s - generalized linear regression (includes linear + logistic)

$\rightarrow$

1980s - Breiman + Friedman introduced - classification + regression trees (random forest) + cross-validation

1990s - ML Boom! Shift to data-driven approach

Support vector machines  
recurrent neural nets

2000s - kernel methods, unsupervised learning becomes more popular

2010s - “deep learning”

non-linear methods too computationally complex at this point

more data

more computational complexity.

## 2 Notation and Simple Matrix Algebra

I'll try to keep things consistent notationally throughout this course. Please call me out if I don't!

$n$  - number of distinct data points or observations in our sample.

$p$  - # of variables available for making predictions.

e.g. Wage data has  $p=12$  variables +  $n=3,000$  people.

$x_{ij}$  - value of the  $j$ th variable for  $i$ th individual.

$$i = 1, \dots, n$$

$$j = 1, \dots, p$$

$X$  -  $n \times p$  matrix whose  $(i,j)$ th element is  $x_{ij}$

$$X = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{pmatrix}$$

$$\underline{x}_i = \underline{x}_i = i^{\text{th}} \text{ row of } X \text{ (vector of length } p) = \begin{pmatrix} x_{i1} \\ \vdots \\ x_{ip} \end{pmatrix}$$

$$\underline{x}_i^T = \underline{x}_i' = (x_{i1} \dots x_{ip}) \text{ "transpose"}$$

$y$  - variable on which we wish to make a prediction

$$y_i = i^{\text{th}} \text{ observation of } y$$

$a, A, A$  - scalar, matrix, random variables

$a \in \mathbb{R}$   $\leftarrow$  indicates dimension.

$$\underline{A} \in \mathbb{R}^{r \times s} = r \times s \text{ matrix}$$

$$\underline{y} \in \mathbb{R}^n$$

Matrix multiplication

Let  $A \in \mathbb{R}^{r \times d}$  and  $B \in \mathbb{R}^{d \times s}$  then the product of  $A$  and  $B$  is " $AB$ "  $\rightarrow$  multiply rows of  $A$  by columns of  $B$  (elementwise)

$$(AB)_{ij} = \sum_{k=1}^d a_{ik} b_{kj}$$

e.g.  $A = \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix}$ ,  $B = \begin{pmatrix} 5 & 6 \\ 7 & 8 \end{pmatrix}$

$$AB = \begin{pmatrix} 1 \times 5 + 2 \times 7 & 1 \times 6 + 2 \times 8 \\ 3 \times 5 + 4 \times 7 & 3 \times 6 + 4 \times 8 \end{pmatrix} = \begin{pmatrix} 19 & 22 \\ 43 & 50 \end{pmatrix} \leftarrow \text{result is } r \times s \text{ matrix}$$