

Predictive Analytics Lecture 2

Adam Kapelner

Stat 422/722

at The Wharton School of the University of Pennsylvania

January 17 & 18, 2017

No Inference Possible as of Now

We haven't spoken about t or F tests. Why is that?

In order to have inference, we need to make explicit random variable model assumptions

$$Y \sim g(\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p, \sigma^2, \dots)$$

must be assumed to be something like

$$Y \sim \mathcal{N}(\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p, \sigma^2)$$

Is this a reasonable thing to do?

Back to Modeling

We said before that our model for Y was

$$Y = f(x_1, \dots, x_p) + \mathcal{E}$$

assuming we can know the model, there still is \mathcal{E} . Where does it come from? According to determinism a la Laplace, if one knew all the causal information, there would be no error

$$y = t(z_1, z_2, \dots)$$

i.e t is the deterministic true mathematical model.

Universal determinism and Laplace's demon

Laplace writes:

We ought then to regard the present state of the universe as the effect of its anterior state and as the cause of the one which is to follow. Given for one instant an intelligence which could comprehend all the forces by which nature is animated and the respective situation of the beings who compose it – an intelligence sufficiently vast to submit these data to analysis – it would embrace in the same formula the movements of the greatest bodies of the universe and those of the lightest atom; for it, nothing would be uncertain and the future, as the past would be present to its eyes.

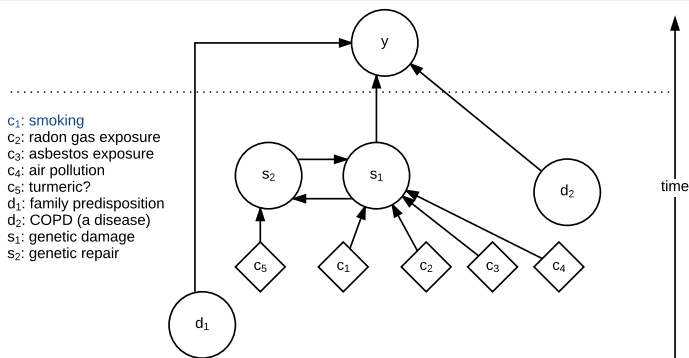
(1814: 4)

The vast intelligence here described has come to be known as Laplace's demon. The idea is obviously founded on that of a human scientist (perhaps Laplace himself) using Newtonian mechanics to calculate the future paths of planets and comets. Extrapolating from this success, it was natural to suppose that a sufficiently vast intelligence could calculate the entire future course of the universe. Laplace himself relates his vast intelligence to human successes in astronomy. As he says:

The human mind offers, in the perfection which it has been able to give to astronomy, a feeble idea of this intelligence. Its discoveries in mechanics and geometry, added to that of universal gravity, have enabled it to comprehend in the same analytical expressions the past and future states of the system of the world.

(Laplace 1814: 4)

Example Lung Cancer Causal Model



Arrows represent causal directions and diamond boxes represent “manipulable” variables (more on this soon). What functions for the response would be deterministic?

$$y = t(d_1, d_2, s_1), y = t(d_1, d_2, s_2, c_1, c_2, c_3, c_4), y = t(d_1, d_2, c_1, c_2, c_3, c_4, c_5)$$

The Root Cause of Randomness

But let's say we only have information about c_1 (a contributory cause, one among many co-occurrent causes). Since we don't have all the inputs (nor the information of the states of the co-occurrent causes), we cannot be sure of y . Hence we'll employ a statistical model,

$$Y \sim \text{Bernoulli}(f(c_1))$$

where we saw before that $f(c_1 = 1) = 16\%$ and $f(c_1 = 0) = 0.4\%$ (probabilistic causation). Thus, the response is stochastic only because we lack information. For regression,

$$Y = f(x_1, \dots, x_p) + \underbrace{t(z_1, z_2, \dots) - f(x_1, \dots, x_p)}_{\varepsilon}$$

(i.e. the "noise" is due to ignorance)

Note... some believe that there is still intrinsic randomness in the universe even with all relevant information known. But we are punting on the actual philosophy...

t is Difficult to Model

In order to get t , you'll need to know all these functions explicitly:

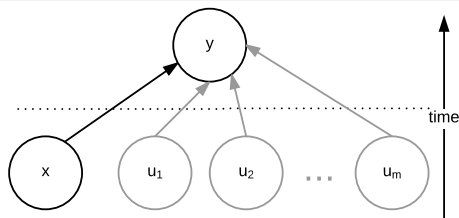
$$y = t_y(d_1, d_2, s_1)$$

$$s_1 = t_{s_1}(c_1, c_2, c_3, c_4, s_2)$$

$$s_2 = f_{s_2}(c_5, s_1)$$

which means that even if you know all the values of variables, you may not be able to properly model the response since you do not know the functional forms t_y , t_{s_1} and t_{s_2} .

A “Nice” Type of Ignorance



In the situation where the true model is

$$y = g(x) + h_1(u_1) + h_2(u_2) + \dots + h_m(u_m)$$

and x is observed but u_1, \dots, u_m are the “unknowns”.

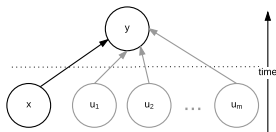
$$h_1(u_1) + h_2(u_2) + \dots + h_m(u_m) \xrightarrow{\mathcal{D}} \mathcal{N} \left(\sum_{k=1}^m \mu_k, \sum_{k=1}^m \sigma_k^2 \right)$$

as the number of unseen variables increase (central limit theorem).

The Normal Homoskedastic Error Model

Let $c = \sum_{k=1}^m \mu_k$ and $\sigma^2 = \sum_{k=1}^m \sigma_k^2$, then

$$y = \underbrace{g(x) + c}_{f(x)} + \mathcal{E}, \quad \text{s.t.} \quad \mathcal{E} = \sum_{k=1}^m h_k(u_k) - c \sim \mathcal{N}(0, \sigma^2)$$



Also, since x does not affect the other variables in any way, it cannot have an influence on their spread, hence σ^2 is not a function of x . Thus the error spread is the same everywhere across the range of x (homoskedasticity).

Parametric Worldview

We are back to the fundamental statistical problem, $Y = f(x) + \mathcal{E}$ where now we are more “okay” with the noise being normal and homoskedastic for all x .

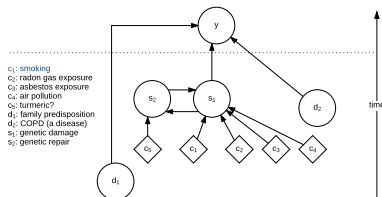
We now invoke the parametric worldview. Within that parametric worldview, we will buy into the linear model. Thus,

$$Y \sim \mathcal{N}(\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p, \sigma^2)$$

But there is one more assumption...

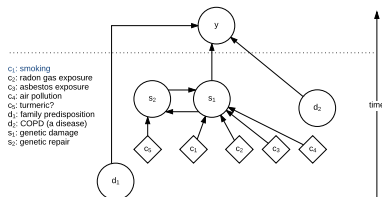
Independence

We now assume that each response is independent of every other response. Second person:



No effect of first person's y_1 (nor any of the unobserved variables which generate the \mathcal{E}_1) on the second person's y (or \mathcal{E}_2).

First person:



If there are, we need to observe them and rotate them into our estimate of $f(x)$. Examples for this cigarette case?

The Classic OLS Assumptions

Preassuming

- linearity (the parametric assumption)

we then further assume

- independence (most important)
- homoskedasticity (less important)
- normality of \mathcal{E} (least important if n is large)

in order to get inference.

A Different Means of Estimation

Last time, we were working on creating a fit \hat{f} that means we need estimates of all the parameters:

$$\hat{f}(x_1, x_2, \dots, x_p) = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_p x_p$$

where the unknown parameters were $\beta_0, \beta_1, \dots, \beta_p$. Our strategy last time was to minimize SSE via a calculus to obtain $\{\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p\}$.

Why was this arbitrary?

Given the three new assumptions (that it's not too much of a stretch to buy into usually), we now have a completely specified joint probability distribution for our observed data,

$$\mathbb{P}(Y_1 = y_1, Y_2 = y_2, \dots, Y_n = y_n \mid \mathbf{X}_1 = \mathbf{x}_1, \mathbf{X}_2 = \mathbf{x}_2, \dots, \mathbf{X}_n = \mathbf{x}_n)$$

where $\mathbf{x}_i := [x_{i1}, x_{i2}, \dots, x_{ip}]$ i.e. the vector of all known measurements / covariates.

What's a probability? What's a likelihood?

In general, a parametric density function / mass function of a r.v. looks like the following:

$$\mathbb{P}(x; \theta) = \dots$$

where θ are the tuning knobs on the model. We ask the question “what’s the probability of this realization x (the data) assuming the density was parameterized at θ ”? Now we ask the inverse question:

$$\mathcal{L}(\theta; x) = \dots$$

that is “what’s the likelihood of these parameters assuming we saw x (the data) come out the way it did”? The $\mathcal{L}()$ denotes the **likelihood function**. Of course, probability and likelihood are exactly the same numerically,

$$\mathbb{P}(x; \theta) = \mathcal{L}(\theta; x) = \dots$$

but conceptually they couldn’t be further apart!

Maximum Likelihood Estimation (MLE)

Why not just ask the very common-sense question, what θ maximizes the probability of seeing what we observe? That would be a good guess as to what θ is.

$$\hat{\theta} := \arg \max_{\theta \in \Theta} \{ \mathcal{L}(\theta; x) \}$$

where Θ represents the space the parameter lives in. In our situation, Θ represents all real numbers in p dimensions. Let's do this in our example. The first step:

$$\begin{aligned} & \mathbb{P}(Y_1 = y_1, Y_2 = y_2, \dots, Y_n = y_n \mid \mathbf{X}_1 = \mathbf{x}_1, \mathbf{X}_2 = \mathbf{x}_2, \dots, \mathbf{X}_n = \mathbf{x}_n) \\ &= \prod_{i=1}^n \mathbb{P}(Y_i = y_i \mid \mathbf{X}_i = \mathbf{x}_i) \end{aligned}$$

How so? Each observation is independent of every other. Recall $\mathbb{P}(ABC) = \mathbb{P}(A)\mathbb{P}(B)\mathbb{P}(C)$ if A , B and C are independent.

MLE of the Linear Model Parameters

We can continue,

$$\begin{aligned} &= \prod_{i=1}^n \mathbb{P}(Y_i = y_i \mid \mathbf{X}_1 = \mathbf{x}_i) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} (y_i - \mathbb{E}[Y_i \mid \mathbf{x}_i])^2\right) \end{aligned}$$

How? Normality and homoskedasticity of \mathcal{E} .

$$= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} (y_i - (\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}))^2\right)$$

How? Linearity of $\mathbb{E}[Y_i \mid \mathbf{x}_i]$. Now we wish to maximize the above over all possible $\beta_0, \beta_1, \dots, \beta_p$.

MLE of the Linear Model Parameters

Then, by some precalc tricks,

$$\begin{aligned} &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} \mathcal{E}_i^2\right) \\ &= \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^n \exp\left(\sum_{i=1}^n -\frac{1}{2\sigma^2} \mathcal{E}_i^2\right) \\ &= \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^n \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n \mathcal{E}_i^2\right) \end{aligned}$$

Pick $\{\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p, \hat{\sigma}^2\}$ such that the above is minimized. The solutions are called the “maximum likelihood estimates (MLE’s)”.

Amazing coincidence

Using calculus, the solution to $\{\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p\}$ is equivalent to minimizing SSE... What a coincidence!!

Note: $\hat{\sigma}^2 = \frac{1}{n}SSE$

The Likelihood Ratio (LR)

Imagine two models: (a) the “full” model where $\theta \in \Theta$ and (b) a reduced model where $\theta \in \Theta_R \subset \Theta$. The reduced space has q less degrees of freedom for θ to operate. Consider the ratio of the likelihoods

$$LR := \max_{\theta \in \Theta} \mathcal{L}(\theta; x) / \max_{\theta \in \Theta_R} \mathcal{L}(\theta; x)$$

representing how much more probable the full model is over the restricted model. But is this increase in probability **statistically significant**? It turns out as n gets large and under pretty forgiving conditions,

$$Q := 2 \ln(LR) \xrightarrow{\mathcal{D}} \chi_q^2$$

Testing the Simple Reduced Model

Let's test our "naive model" from Lecture 1 (always predicting $\hat{y} = \bar{y}$) versus having a model having many predictors in a linear model.

$$\begin{aligned}
 LR &= \frac{\max_{\beta_0, \beta_1, \dots, \beta_p, \sigma^2} \mathcal{L}(\beta_0, \beta_1, \dots, \beta_p; y_1, \dots, y_n, \mathbf{x}_1, \dots, \mathbf{x}_n)}{\max_{\beta_0, \sigma^2} \mathcal{L}(\beta_0, \beta_1 = 0, \dots, \beta_p = 0; y_1, \dots, y_n, \mathbf{x}_1, \dots, \mathbf{x}_n)} \\
 &= \frac{\left(\frac{1}{\sqrt{2\pi\hat{\sigma}^2}}\right)^n \exp\left(-\frac{1}{2\hat{\sigma}^2} SSE\right)}{\left(\frac{1}{\sqrt{2\pi\hat{\sigma}_0^2}}\right)^n \exp\left(-\frac{1}{2\hat{\sigma}_0^2} SSE_0\right)} \\
 &= \left(\frac{SSE_0}{SSE}\right)^{n/2} \frac{\exp\left(-\frac{n}{2SSE} SSE\right)}{\exp\left(-\frac{n}{2SSE_0} SSE_0\right)}
 \end{aligned}$$

Testing the Simple Reduced Model

Now we build the Q statistic:

$$Q = 2 \ln \left(\left(\frac{SSE_0}{SSE} \right)^{n/2} \right) = n \ln \left(\frac{SSE_0}{SSE} \right) \xrightarrow{\mathcal{D}} \chi_p^2$$

This can be used to test

$$H_0 : \beta_1 = 0, \beta_2 = 0, \dots, \beta_p = 0$$

$$H_a : \text{at least one is non-zero}$$

There is another test for this you've learned about?

Omnibus F-test

$$F = \frac{\frac{SSE_0 - SSE}{p}}{\frac{SSE}{n-p}} = \frac{SSE_0 - SSE}{SSE} \frac{n-p}{p} = \left(\frac{SSE_0}{SSE} - 1 \right) \frac{n-p}{p} \sim F_{p, n-p}$$

Both tests use the same test statistic, namely SSE_0/SSE (up to constants and a monotonic transformation). It is a harder proof to demonstrate they have the same power for the same n and α (but they do).

Some points

- The likelihood ratio test / F test can also test any subset of the predictors (even one).
- Thus, we now have inference for every predictor or subset of predictors i.e.
 - Hypothesis testing
 - Confidence intervals

What does inference buy you?

Previously,

$$Y \sim g(\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p, \sigma^2, \dots)$$

Do not assume OLS assumptions. We picked L2 loss and minimized to get $\{\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p\}$. What do these numbers means?

$$Y \stackrel{ind}{\sim} \mathcal{N}(\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p, \sigma^2)$$

Assume OLS assumptions. Using MLE, we wind up minimizing L2 loss and get the same $\{\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p\}$. What do these numbers means? Same thing, except now ... we can “test” each value and provide confidence intervals for each value. You know their stability.

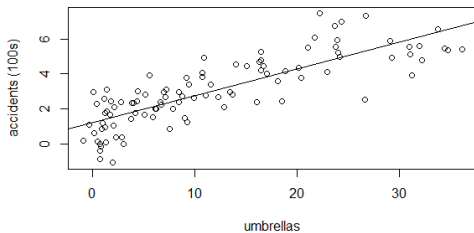
What you want to say about $\hat{\beta}_j$

[Interpret stolen bases in baseball dataset in JMP].

A change in x_j of +1 causes / induces a β_j difference in its mean response y .

Umbrella Sales and Car Accidents

Consider a simple example. x : umbrella sales and y : car accidents.
What would the relationship look like?



Does 100 more umbrellas sold *cause* 15.3 more car accidents (on average)? No... only an association (assessed by a linear correlation).

Correlation Does Not Imply Causation

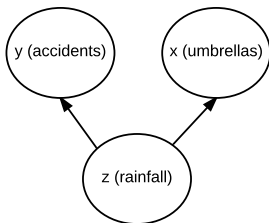
What can correlation mean?

- It is a coincidence. How can this be?
- They are consequence from of a common cause (the **lurking** or **counfounding** variable). How can this be?
- There is causation
 - x causes y (possibly with intermediates)
 - y causes x (possibly with intermediates)
 - x and y cause each other (cyclic)

(recall time-boundedness property)

Controlling for the Confounder

The confounding variable is likely $z = \text{rainfall}$.



The illustration shows that if you change x obviously y doesn't change whatsoever (causes always precede their dependent effects an assumption known as temporal boundedness)

[Show regression in R]

A Proper Interpretation of $\hat{\beta}_j$

$\hat{\beta}_j$ estimates β_j . Imagine n is large and the confidence interval is really small. So basically, $\hat{\beta}_j = \beta_j \neq 0$. Interpretation?

Another object naturally observed with exactly the same features except that x_j is increased by 1 unit will have a β_j difference in its mean response y .

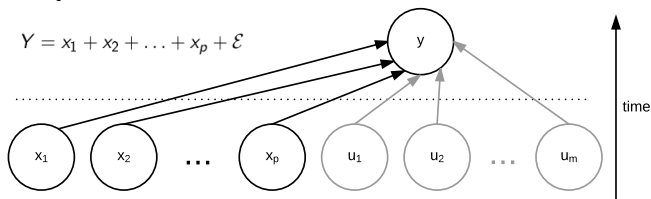
Another one: $\hat{\beta}_j$ estimates β_j . Imagine n is not so large and the confidence interval is not small but we are still convinced $\beta_j \neq 0$. Interpretation?

Another object naturally observed with exactly the same features except that x_j is increased by 1 unit will have a $\hat{\beta}_j \pm \text{SE} [\hat{\beta}_j]$ difference in its mean response y . (Not much difference except accounting for parameter estimation error).

When can you say “causes”?

When can the interpretation be as follows? ~~Another object naturally observed with exactly the same features except for a change~~ **If this object in front of us has its x_j changed by +1, it will have** **cause** a $\hat{\beta}_j \pm \text{SE} [\hat{\beta}_j]$ difference in its mean response y .

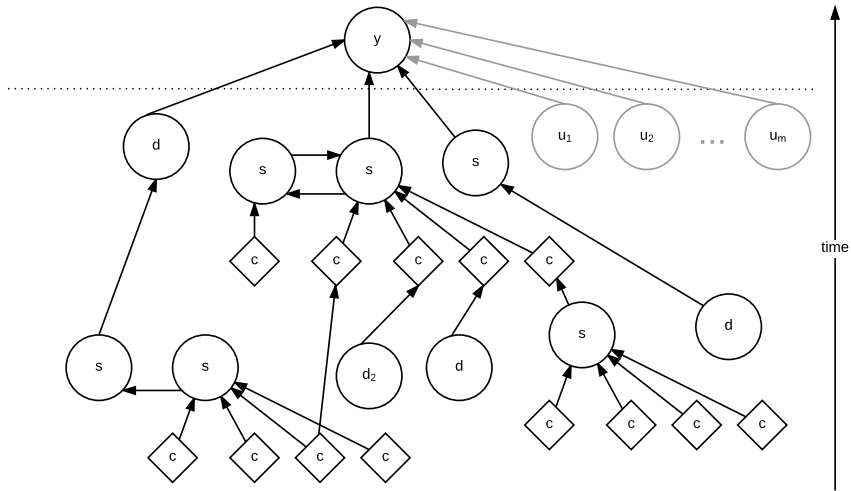
- 1 If we can just assume the model looks as follows:



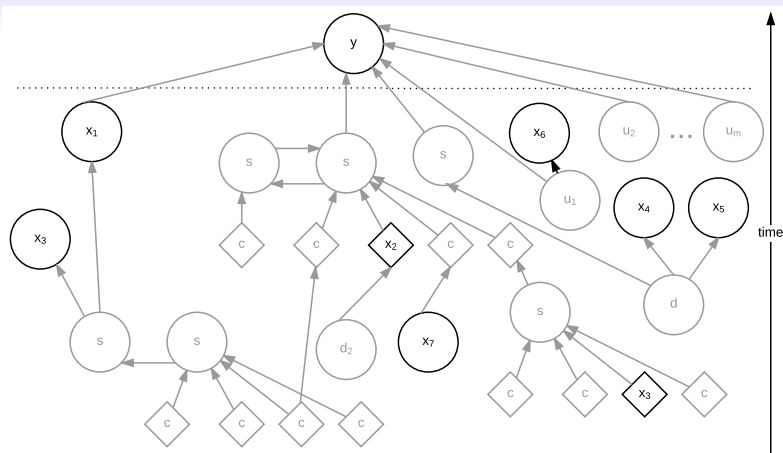
(for all p features ... how can the illustration be updated for one variable?)

- 2 If we've run a randomized experiment manipulating x_j among the objects AND assuming an linear additive effect of x_j on y .

Consider a Realistic Model



Consider Realistic Predictors



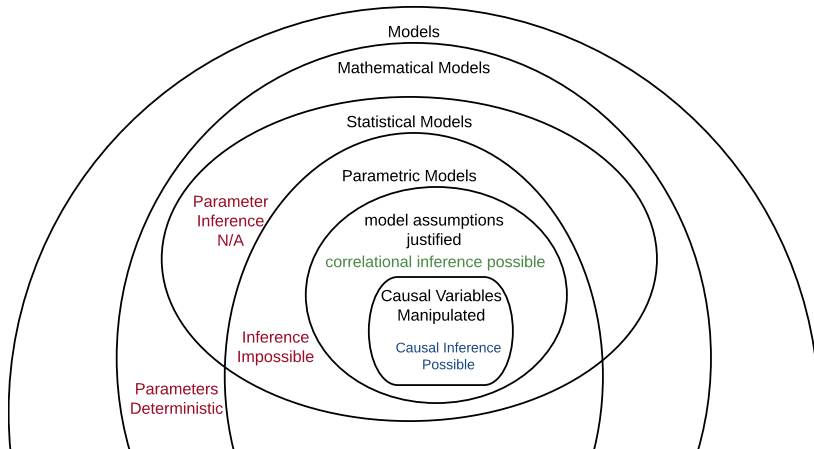
Grey variables and known to be dependent but the values are unknown and the u_k 's are the “unknown unknowns”.

Consider Realistic Predictors

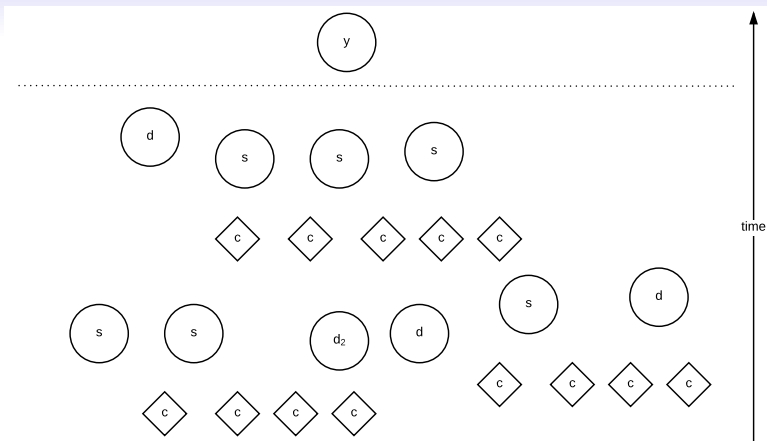
Observations from the previous illustration

- Maybe some of the predictors x_1, \dots, x_p are causal, but most are likely not.
- Of the ones that are not causal due to a confounder, you may have an idea of the lurking variables but it is unlikely you can measure them. Think college GPA vs SAT with confounder true IQ / ability.
- If some variables are causal, it is unlikely they have an additive causal effect; their effect is likely moderated by many other interacting variables possibly in non-linear ways.
- A linear model for y on x_1, \dots, x_p is likely far from the truth (not related to our discussion on causality).

Inference and Causality



Sidebar: Theories are Hard...

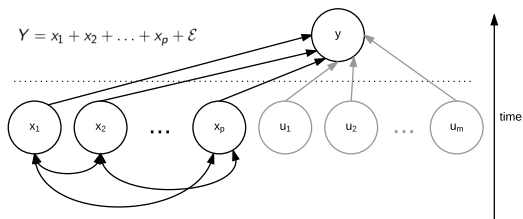


Maybe we know the predictors and the dependence, but don't know the causal dependencies. How many theories are possible? 23 variables + the response ... And that's not even counting the unknown unknowns...

More on OLS Coefficient Interpretation

The linear regression coefficient interpretation again: another object **naturally observed** with exactly the same features except that x_j is increased by 1 unit will have a $\hat{\beta}_j \pm \text{SE}[\hat{\beta}_j]$ difference in its mean response y .

What do we mean by naturally observed? This other object is realized from the same joint distribution as all other observations. This means that whatever *multicollinearity* / *covariance structure* exists between the predictors, $\{\text{Cov}[X_j, X_k]\}$, will give rise to the predictor values in the other object.



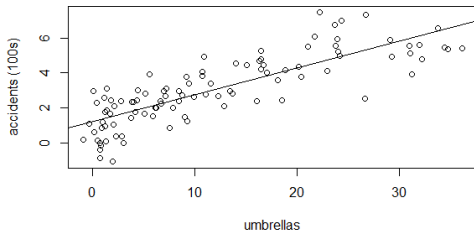
The Hidden “Fifth” OLS Assumption

So this language “... exactly the same features except that x_j is increased by ...” is kind of absurd in the context of a strong covariance structure as ... i.e. it will be very rare to observe an observation with x_j different without any other predictor values different. Example from baseball dataset?

There is room to argue that to have these interpretations be at all realistic, we must assume there is not a strong multicollinearity structure between x_j and the other predictors.

But Real Correlations Still Rock

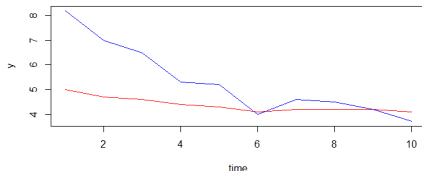
We've been beating up on correlations and their interpretations e.g. the following:



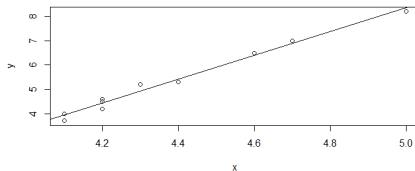
But even though higher umbrella sales do not “cause” accidents, can they still predict them? Yes, R^2 is totally agnostic to (a) if your model is true and (b) if your variables are causal or not. Predictors *truly* correlated (a causal link exists) to the response contain information about the value of the response and it doesn't matter through what channel it provides that information.

Fake / Spurious Correlations

x is margarine consumption per capita in America measured yearly for 10 years from 2000-2009, y is the divorce rate in Maine per 1000 people measured yearly for 10 years from 2000-2009



Are they linearly predictive of one another?



$R^2 \approx 99\%$ and F test has a $p_{val} \approx 1 \times 10^{-8}$. [R demo] Be careful about featurization... try to at least have some inkling of an idea for a causal dependency for the response on the predictors...

Testing Multiple Predictors at Once

[JMP Baseball data]

Dataframe Design

We spoke a lot about featurization i.e. selecting the columns in the dataframe (these are the predictors to measure). Once we did this, we can then go out and sample observations and then measure each for their predictor values.

But we didn't speak at all about selecting the observations themselves!

Modeling Categorical Responses

Previously the response y was continuous and via the OLS assumptions we obtained the statistical model,

$$Y \stackrel{ind}{\sim} \mathcal{N}(\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p, \sigma^2)$$

If the response y is categorical, can we still use this? No... the only elements in the support of the r.v. Y are the levels only.

First, assume Y is binary i.e. zero or one. We spoke about this statistical model before,

$$Y \sim \text{Bernoulli}(f(x_1, \dots, x_p))$$

since $\mathbb{E}[Y \mid x_1, \dots, x_p] = f(x_1, \dots, x_p)$, then f is still the conditional expectation function like before except now it varies only within $[0, 1]$ and it is the same as $\mathbb{P}(Y = 1 \mid x_1, \dots, x_p)$.

