# MS&E228: Lecture 2 Causality via Experiments

Vasilis Syrgkanis
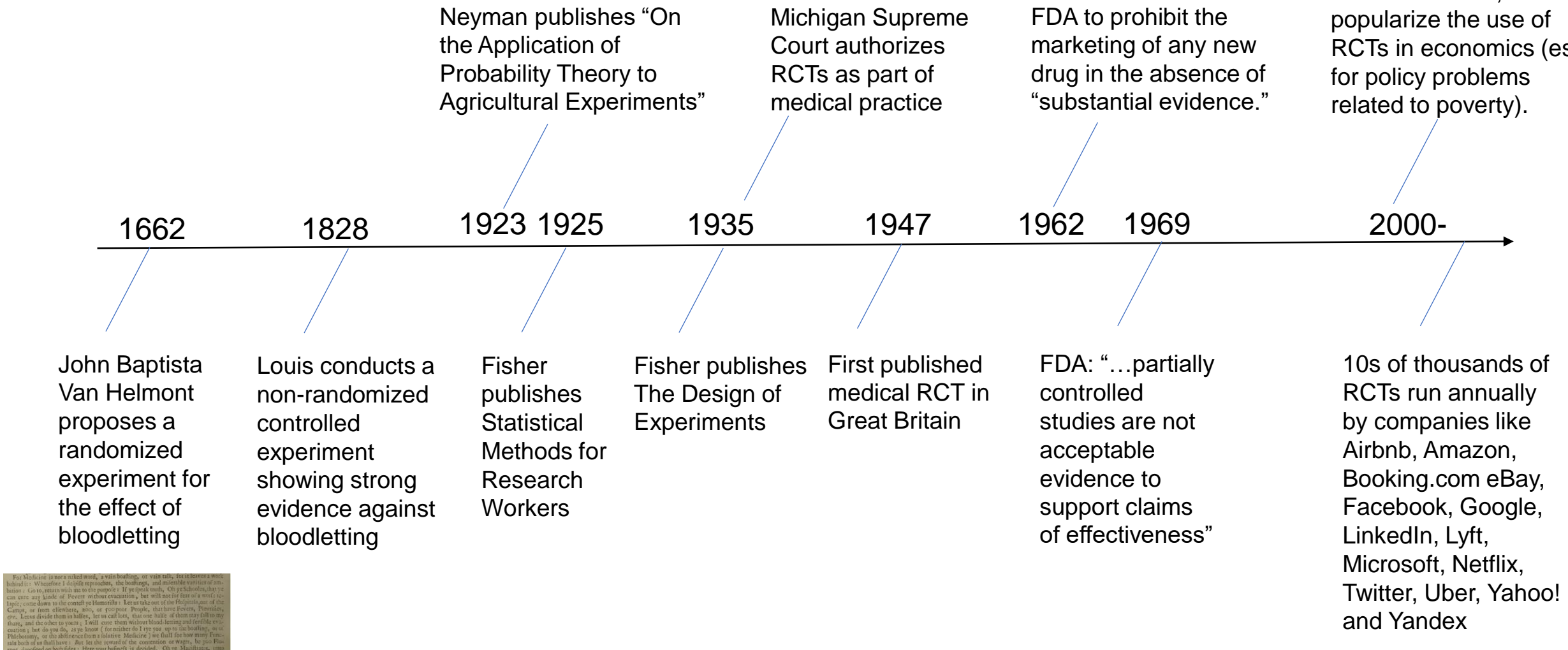
MS&E, Stanford

# Randomized Control Trial (RCT)

- $n$ units are drawn from a population

- Each unit is assigned at random with some probability $p$ to one of two groups {control, treatment}

- Each unit in the control group receives treatment $D = 0$

- Each unit in the treatment group receives treatment $D = 1$

- At the end of the experiment an outcome $y$ is measured for each unit

# A Brief History of Experimentation

Abhijit Banerjee, Esther Duflo, and Michael Kremer, popularize the use of RCTs in economics (esp. for policy problems related to poverty).

Neyman publishes "On the Application of Probability Theory to Agricultural Experiments"

Michigan Supreme Court authorizes RCTs as part of medical practice

FDA to prohibit the marketing of any new drug in the absence of "substantial evidence."

| 1662 | 1828 | 1923 | 1925 | 1935 | 1947 | 1962 | 1969 | 2000- |

John Baptista Van Helmont proposes a randomized experiment for the effect of bloodletting

Louis conducts a non-randomized controlled experiment showing strong evidence against bloodletting

Fisher publishes Statistical Methods for Research Workers

Fisher publishes The Design of Experiments

First published medical RCT in Great Britain

FDA: "…partially controlled studies are not acceptable evidence to support claims of effectiveness"

10s of thousands of RCTs run annually by companies like Airbnb, Amazon, Booking.com eBay, Facebook, Google, LinkedIn, Lyft, Microsoft, Netflix, Twitter, Uber, Yahoo! and Yandex

Credit: Ron Kohavi, History of Controlled Experiments

RCTs are the gold standard for measuring the "causal effect" of a "treatment" on an "outcome"

# Goal #1: Mathematical Definition of Causality

- We want to formally (mathematically) define the causal effect of a binary treatment on a scalar outcome of interest

**Examples**
- Effect of seed $A$ vs $B$ on crop yield
- Effect of completion of a job training program on observed wage
- Effect of drug vs placebo on overall survival
- Effect of an ad impression on conversion (purchase)
- Effect of eating avocados on longevity

# Causality via Potential Outcomes

- Nature generates two latent (unobserved) random outcomes
$$Y^{(1)}, Y^{(0)}$$

- $Y^{(d)}$: potential outcome that would have been observed if unit received treatment $d \in \{0,1\}$

**Example**

- $Y^{(0)}$: wage if you don't participate in a training program
- $Y^{(1)}$: wage if you participate in a training program

# Causality via Potential Outcomes

- $Y^{(0)}, Y^{(1)}$ are called "counterfactuals" as they can never be simultaneously observed

- Fundamental problem of causal inference

**Example**

- We don't have two replicas of each unit

- We cannot observe your wage with the training program and without

# Causality via Potential Outcomes

- The individual treatment effect is: $\text{ITE} := Y^{(1)} - Y^{(0)}$
- Un-attainable due to fundamental problem of causal inference
- Average Treatment Effect (ATE)

$$\delta := E\left[Y^{(1)} - Y^{(0)}\right]$$

**Example**

- $ITE :=$ difference in wage of your potential outcome with the training program and without
- $ATE :=$ average difference in the population

# Treatment Assignment and Observed Outcome

- Each unit receives treatment $D \in \{0,1\}$, we observe
$$Y \equiv Y^{(D)}$$

- Given data $(D, Y)$ what quantities can we "identify" $\equiv$ "measure if we had access to infinite data"

- Can we measure the ATE?

# The Effect of Eating Avocados on Health



high

Life expectancy

low

Avocado
Consumption

- Can identify average outcome within each treatment group

$$E[Y|D=d] = E\big[Y^{(D)}\big|D=d\big] = E\big[Y^{(d)}\big|D=d\big]$$

- When is this representative of average un-conditional counterfactual outcome
$$E\big[Y^{(d)}\big]$$

- In general
$$E\big[Y^{(d)}\big|D=d\big] \neq E\big[Y^{(d)}\big]$$

- Known as "selection bias"

# The Effect of Eating Avocados on Health

high

Life expectancy

low

$E[Y^{(0)}]$

$E[Y^{(0)}|D = 0]$

Avocado
Consumption

- Can estimate average outcome within each treatment group

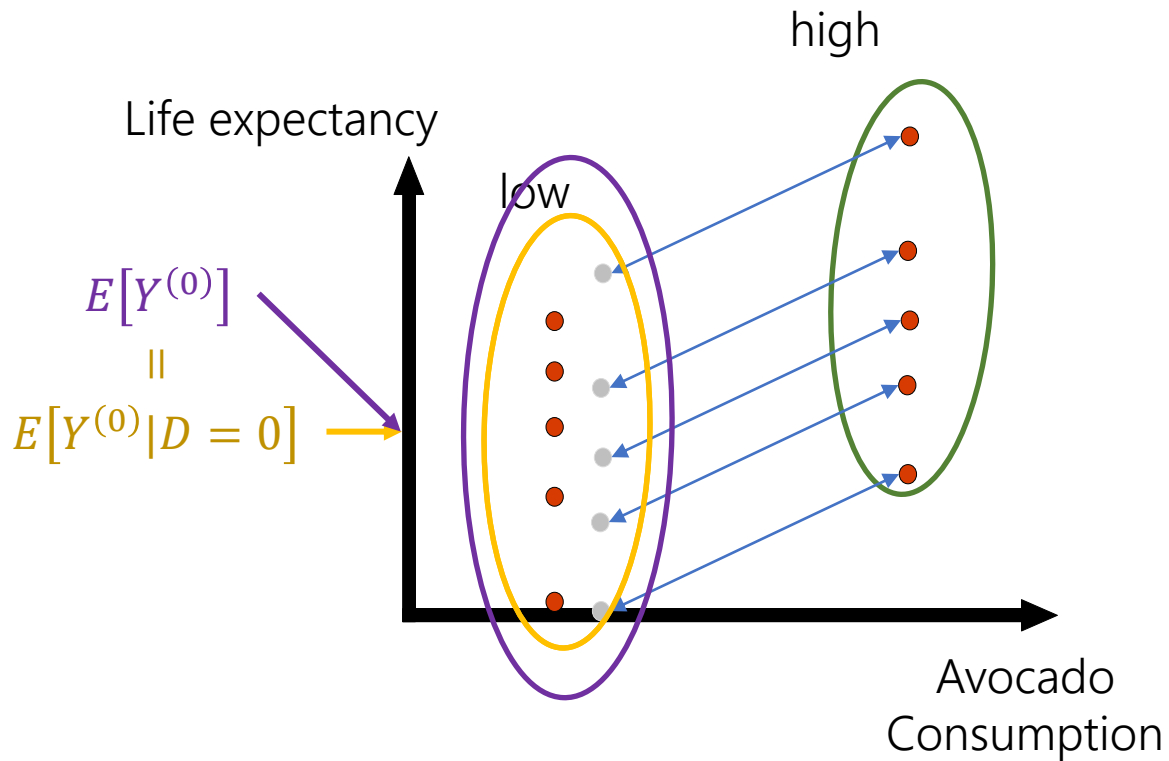$$E[Y|D = d] = E[Y^{(D)}|D = d] = E[Y^{(d)}|D = d]$$

- When is this representative of average un-conditional counterfactual outcome
$$E[Y^{(d)}]$$

- In general
$$E[Y^{(d)}|D = d] \neq E[Y^{(d)}]$$

- Known as "selection bias"

# The Effect of Eating Avocados on Health



- Can estimate average outcome within each treatment group

$$E[Y|D=d] = E[Y^{(D)}|D=d] = E[Y^{(d)}|D=d]$$

- When is this representative of average un-conditional counterfactual outcome
$$E[Y^{(d)}]$$

- If treatment was fully randomly assigned (e.g. a randomized experiment)
$$Y^{(d)} \perp D$$

- Then no selection bias:
$$E[Y^{(d)}|D=d] = E[Y^{(d)}]$$

# Identification of ATE under Random Assignment

Suppose that treatment is randomly assigned (i.e. RCT)

$$Y^{(d)} \perp D$$

and $0 < \Pr(D = 1) < 1$.

Then average **observed** outcome in treatment group $d \in \{0,1\}$ recovers average **potential** outcome for treatment $d$

$$E[Y|D = d] = E\left[Y^{(D)}|D = d\right] = E\left[Y^{(d)}|D = d\right] = E\left[Y^{(d)}\right]$$

Hence, average **predictive** effect recovers the average **treatment** effect

$$\pi := E[Y|D = 1] - E[Y|D = 0]$$
$$= E\left[Y^{(1)}\right] - E\left[Y^{(0)}\right] =: \delta$$

# Statistical Inference

# What can we do with finite samples

- Assume we have access to $n$ samples $(Y_i, D_i)$

- Drawn i.i.d. from the distribution of the random variables $(Y, D)$

- Want estimate $\hat{\delta}$ of $\delta = E[Y|D = 1] - E[Y|D = 0]$ (identification)

- Denote with $E_n[\cdot]$ the empirical average $E_n[X] = \frac{1}{n}\sum_{i=1}^{n} X_i$

- Estimate $\hat{\delta}$ by using empirical averages for each sub-population

$$\hat{\theta}_d := \frac{E_n[Y\,1\{D = d\}]}{E_n[1\{D = d\}]} = \frac{1}{\#\{i : D_i = d\}}\sum_{i:D_i=d} Y_i$$

$$\hat{\delta} = \hat{\theta}_1 - \hat{\theta}_0$$

# How accurate is the estimate?

Under mild regularity conditions

$$\sqrt{n}\{\hat{\theta}_d - \theta_d\}_{d \in \{0,1\}} \overset{a}{\sim} N(0, V)$$

where

$$V = \begin{pmatrix} \dfrac{\text{Var}(Y|D=0)}{P(D=0)} & 0 \\ 0 & \dfrac{\text{Var}(Y|D=1)}{P(D=1)} \end{pmatrix}$$

Hence

$$\sqrt{n}(\hat{\delta} - \delta) \overset{a}{\sim} N(0, V_{11} + V_{22})$$

$X \overset{a}{\sim} Y$ means that as $n \to \infty$:

$$\sup_{R \in \mathcal{R}} |P(X \in R) - P(Y \in R)| \approx 0$$

where $\mathcal{R}$ set of all hyper-rectangles

# Proof Sketch

Trivially we can write $\theta_d = E[Y^{(d)}] \frac{E_n[1(D=d)]}{E_n[1(D=d)]}$

$$\hat{\theta}_d - \theta_d = \frac{E_n[Y^{(d)} 1(D=d)]}{E_n[1(D=d)]} - \theta_d = \frac{E_n[(Y^{(d)} - E[Y^{(d)}]) 1(D=d)]}{E_n[1(D=d)]}$$

By Law of Large Numbers (LLN)

$$\hat{\theta}_d - \theta_d \approx \frac{E_n[(Y^{(d)} - E[Y^{(d)}]) 1(D=d)]}{P(D=d)}$$

Difference is average of the i.i.d. mean zero r.v.s: $\dfrac{(Y_i^{(d)} - E[Y^{(d)}]) 1(D_i=d)}{P(D=d)}$

With zero covariance and variance: $\dfrac{E[(Y_i^{(d)} - E[Y^{(d)}])^2 1(D_i=d)]}{P(D=d)^2} = \dfrac{Var(Y|D=d)}{P(D=d)}$

Statement follows by Central Limit Theorem (CLT)

# Variance estimate

- Same statement also holds with consistent estimate of variance

$$\widehat{V} = \begin{pmatrix} \dfrac{\text{Var}_n(Y|D=0)}{E_n[1-D]} & 0 \\ 0 & \dfrac{\text{Var}_n(Y|D=1)}{E_n[D]} \end{pmatrix}$$

# Confidence Interval

- $X \overset{a}{\sim} Y \equiv \sup\limits_{[\ell,u]} |P(X \in [\ell,u]) - P(Y \in [\ell,u])| \approx 0$
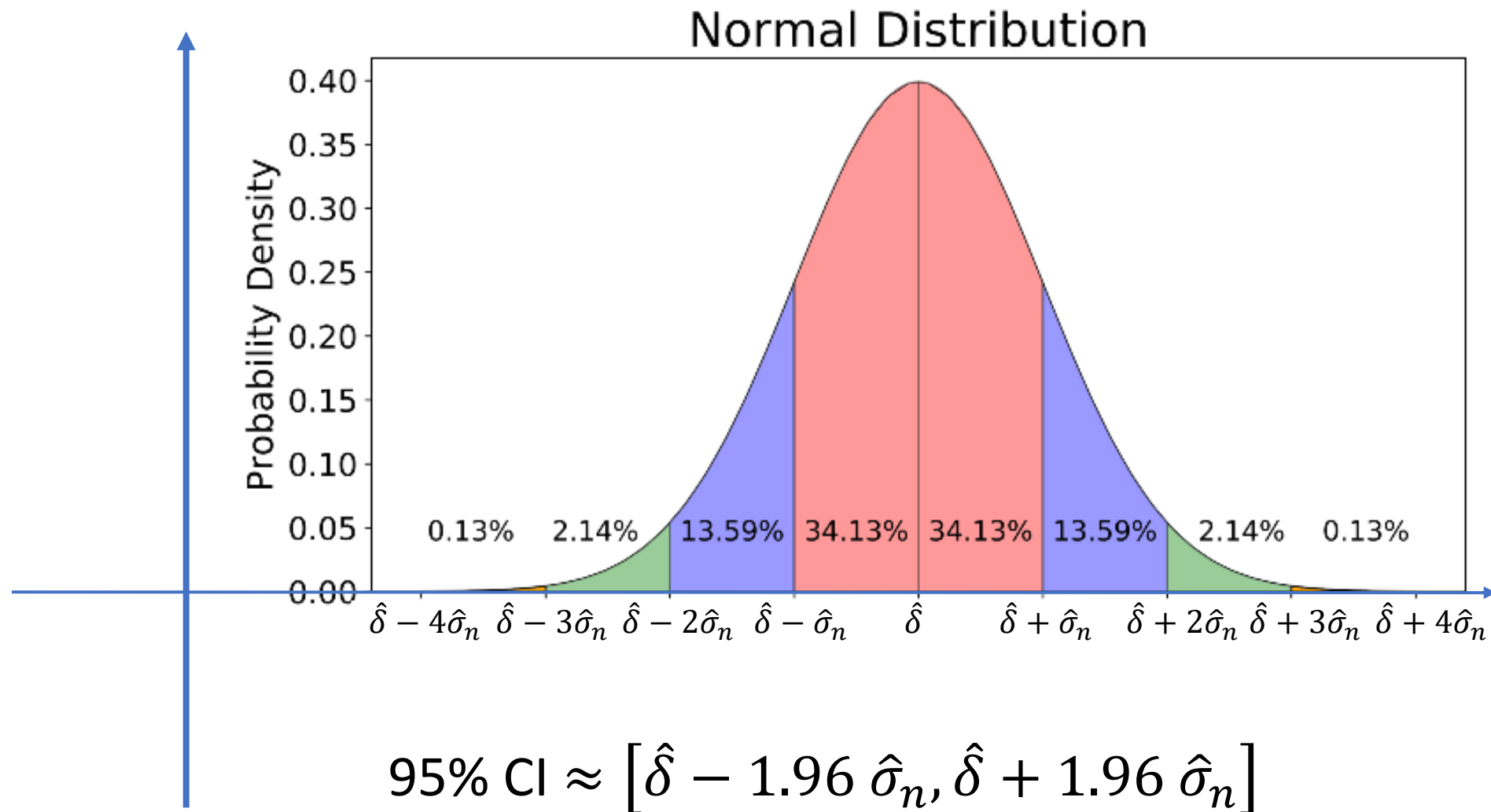
- If we consider $[\ell,u]$ the $\left(\frac{\alpha}{2}, 1 - \frac{\alpha}{2}\right)$ quantile of $N(0,\hat{V})$ then

$$P\left(\sqrt{n}(\hat{\delta} - \delta) \in [\ell,u]\right) \approx 1 - \alpha$$

**Standard error**

- Equivalently, let $z_\alpha$ the $\alpha$ quantile of $N(0,1)$ and $\boxed{\hat{\sigma}_n = \sqrt{\hat{V}/n}}$ then:

$$P\left(\delta \in \left[\hat{\delta} - z_{1-\frac{\alpha}{2}}\hat{\sigma}_n, \hat{\delta} + z_{1-\frac{\alpha}{2}}\hat{\sigma}_n\right]\right) \approx 1 - \alpha$$

# Confidence Interval



$$95\% \text{ CI} \approx \left[\hat{\delta} - 1.96\,\hat{\sigma}_n, \hat{\delta} + 1.96\,\hat{\sigma}_n\right]$$

Let's try it out!

# Relative effect

- Many times (e.g. vaccine trials) we are interested in relative effect

$$RE = \frac{E\left[Y^{(1)} - Y^{(0)}\right]}{E\left[Y^{(0)}\right]} = \frac{\theta_1 - \theta_0}{\theta_0}$$

- We can construct a plug-in estimate

$$\widehat{RE} = \frac{\hat{\theta}_1 - \hat{\theta}_0}{\hat{\theta}_0} = \frac{\hat{\theta}_1}{\hat{\theta}_0} - 1$$

**Delta method:** for any function $f$, with $G = \nabla f(\theta)$

$$\sqrt{n}\left(f(\hat{\theta}) - f(\theta)\right)$$
$$\approx G\sqrt{n}(\hat{\theta} - \theta) \overset{a}{\sim} N(0, G'VG)$$

- By delta method with $G = \left(-\theta_1/\theta_0^2, 1/\theta_0\right)$

$$\sqrt{n}\left(\widehat{RE} - RE\right) \overset{a}{\sim} N\left(0, \frac{\theta_1^2 V_{11}}{\theta_0^4} + \frac{V_{22}}{\theta_0^2}\right)$$

# Example: Pfizer Vaccine

| Efficacy Endpoint Subgroup | BNT162b2 Nᵃ=19965 Cases n1ᵇ Surveillance Timeᶜ (n2ᵈ) | Placebo Nᵃ=20172 Cases n1ᵇ Surveillance Timeᶜ (n2ᵈ) | Vaccine Efficacy % (95% CI)ᵉ |
|---|---|---|---|
| Overall | 9 | 169 | 94.6 (89.6, 97.6) |
| | 2.332 (18559) | 2.345 (18708) | |
| Age group (years) | | | |
| 16 to 17 | 0 | 1 | 100.0 (-3969.9, 100.0) |
| | 0.003 (58) | 0.003 (61) | |
| 18 to 64 | 8 | 149 | 94.6 (89.1, 97.7) |
| | 1.799 (14443) | 1.811 (14566) | |
| 65 to 74 | 1 | 14 | 92.9 (53.2, 99.8) |
| | 0.424 (3239) | 0.423 (3255) | |
| ≥75 | 0 | 5 | 100.0 (-12.1, 100.0) |
| | 0.106 (805) | 0.109 (812) | |

# Approximate Confidence Interval

- Outcomes are binary $y \in \{0,1\}$
- Distribution of outcome for each $d$ is Bernoulli with success $p_d$
- For each subpopulation, mean outcome $\hat{p}_d = \frac{\text{Cases}_d}{N_d}$ is estimate of $p_d$

$$\text{Vaccine Efficacy } (VE) = -\widehat{RE} = \frac{\hat{p}_0 - \hat{p}_1}{\hat{p}_0}$$

- An estimate of the variance of $y$ is $\hat{p}_d(1 - \hat{p}_d)$
- 95% confidence interval can be derived by delta method
- Alternative: approximate bootstrap; simulate $\hat{p}_d + N\left(0, \hat{V}_d/n\right)$

Let's try it out!

# Pre-Treatment Covariates

Precision and Heterogeneity

# Pre-Treatment Covariates

- Assume we have covariates $W$ that correspond to variables determined prior to treatment assignment (e.g. age, income groups)

- How can we use them?

- Heterogeneity: how does the effect vary with these covariates

- Formalized by the Conditional Average Treatment Effect (CATE)

$$\delta(W) := E\big[Y^{(1)} - Y^{(0)}\big|W\big]$$

# Identification of CATE under Random Assignment

Suppose that treatment is randomly assigned (i.e. RCT)

$$\left(Y^{(d)}, W\right) \perp D$$

and $0 < \Pr(D = 1) < 1$.

Then conditional **observed** outcome in treatment group $d \in \{0,1\}$ recovers conditional **potential** outcome for treatment $d$

$$E[Y|D = d, W] = E\left[Y^{(d)}\middle|D = d, W\right] = E\left[Y^{(d)}\middle|W\right]$$

Hence, conditional **predictive** effect recovers the CATE

$$\pi(W) := E[Y|D = 1, W] - E[Y|D = 0, W]$$

$$= E\left[Y^{(1)}\middle|W\right] - E\left[Y^{(0)}\middle|W\right] =: \delta(W)$$

# If we only care about ATE are co-variates useful?

# Co-variates for Sanity Check

- Since treatment is supposed to be independent of co-variates $W$
$$W|D = 1 \sim W|D = 0$$
- For instance, $E[W|D = 1] = E[W|D = 0] = E[W]$
- $D$ does not predict any covariate
- Equivalently
$$D|W \sim D$$
- $D$ is not predictable by any covariate
- We can test all these conditions on the samples to uncover violations of random assignment
- These are typically referred to as co-variate balance tests

# Co-variates for Precision

- Even if we are only interested on ATE covariates can be valuable for precision

- Suppose variance of $y$ is large but can be explained largely by $W$

- Then we can use $W$ to remove all the explained variation from $y$

- Then perform our ATE analysis on the remnant variation

- This is oftentimes performed in practice via ordinary linear regression of $y$ on the vector $(1, D, W)$ (after centering $W$, i.e. $E[W] = 0$)

# Is this consistent?

- Suppose that the conditional expectation function (CEF) of the outcome is indeed linear, with $(D, 1, W)$

$$E[Y \mid D, W] = D\alpha + \alpha_0 + W'\beta$$

- Then note that

$$E[Y(0)] = E\big[E[Y|D = 0, W]\big] = \alpha_0$$
$$E[Y(1)] = E\big[E[Y|D = 1, W]\big] = \alpha + \alpha_0$$

- Baseline outcome is coefficient associated with the intercept 1

- Average effect is coefficient associated with treatment $D$

- Next lecture: this does not require the linear CEF assumption

# Analysis of Variance (ANOVA)

Sneak peek into some material from next lecture

# Variance of Estimate

- The OLS theory that we will cover in the next section yields

$$\sqrt{n}(\hat{\alpha} - \alpha) \overset{a}{\sim} N(0, V_\alpha) \qquad\qquad V_\alpha = \frac{E[\epsilon^2 \widetilde{D}^2]}{E[\widetilde{D}^2]^2}$$

- $\epsilon$ is residual outcome:

$$\epsilon = y - D\alpha - \alpha_0 - W'\beta \qquad E[\epsilon|D, W] = 0 \text{ (by Linear CEF)}$$

- $\widetilde{D}$ is residual treatment (removing whatever is linearly predictable from $(1, W)$)

$$\widetilde{D} = D - E[D]$$

# Variance without adjustment

- OLS theory that we will cover in the next section yields

$$V_\alpha = \frac{E[\epsilon^2 \widetilde{D}^2]}{E[\widetilde{D}^2]^2}$$

- $\epsilon$ is residual outcome: $\epsilon = y - D\alpha - \alpha_0 - W'\beta$ with $E[\epsilon|D,W] = 0$

- Two means estimate is equivalent to OLS without $W$. OLS theory gives variance $\bar{V}_\alpha$ of same form but with residual:

$$\bar{\epsilon} = y - D\alpha - \alpha_0 = W'\beta + \epsilon$$

$$\begin{aligned}
E[\bar{\epsilon}^2 \widetilde{D}^2] &= E[(W'\beta + \epsilon)^2 \widetilde{D}^2] \\
&= E[(W'\beta)^2 \widetilde{D}^2] + E[\epsilon^2 \widetilde{D}^2] + 2E[\beta'W\epsilon\widetilde{D}^2] \\
&= E[(W'\beta)^2 \widetilde{D}^2] + E[\epsilon^2 \widetilde{D}^2] + 2E[\beta'WE[\epsilon \mid D, X]\widetilde{D}^2] \\
&= E[(W'\beta)^2 \widetilde{D}^2] + E[\epsilon^2 \widetilde{D}^2]
\end{aligned}$$

$$\bar{V}_\alpha \geq V_\alpha$$

Variance of OLS estimate with extra co-variates (adjusted) is weakly smaller than two-means estimate (un-adjusted)

# Heteroskedasticity Robust Variance

- Variance formula

$$V_\alpha = \frac{E\left[\epsilon^2 \widetilde{D}^2\right]}{E\left[\widetilde{D}^2\right]^2}$$

- is valid even when the linear CEF assumption is violated

- Inference is asymptotically valid!

- Important to note that this formula is known as the "heteroskedasticity robust variance formula" (HC0)

- Many software packages make the simplification that the residual $\epsilon$ is independent of $D, W$, leading to $V_\alpha = E\left[\epsilon^2\right]/E\left[\widetilde{D}^2\right]$. This is incorrect in most cases!

# Precision Beyond Linear CEF

- The precision statement invoked the property that the residual of the OLS regression of $y$ on $D, X$ is mean zero conditional on $D, X$

- If linear CEF is violated, then all we know is the orthogonality property

$$E\left[\epsilon \begin{pmatrix} D \\ 1 \\ W \end{pmatrix}\right] = 0 \left[\text{FOC of: } \min_{\alpha, \alpha_0, \beta} E\left[(y - D\alpha - \alpha_0 - W'\beta)^2\right]\right]$$

- This is not sufficient to argue that the cross-term vanishes

$$E\left[\beta' W \epsilon \widetilde{D}^2\right] = E\left[\beta' E[W\epsilon \mid D]\widetilde{D}^2\right]$$

- Note that we only need that:

$$E[W\epsilon \mid D] = 0$$

Let's try it out!

# OLS with Interactive Terms

- It is advisable that instead of running OLS of $y$ on $D, 1, W$ we also include interaction terms, i.e. $y$ on $D, 1, W, DW$ [Lin'13]
- In the absence of any model assumptions, the coefficient of $D$ and of the intercept, recover the ATE and the mean baseline outcome
- These interactive terms enforce the residual $\epsilon$ of OLS to satisfy the stronger orthogonality property with $X = (1, W)$

$$E\left[\epsilon \begin{pmatrix} X \\ DX \end{pmatrix}\right] = 0$$

- $E[\epsilon DX] = 0 \Rightarrow E[\epsilon X \mid D = 1] = 0 \Rightarrow E[\epsilon X \mid D = 0] = 0$
- Interaction term in $\bar{V}_\alpha$ is zero and we get that $\bar{V}_\alpha \geq V_\alpha$ without assumptions
- OLS with interactive terms always has weakly smaller variance than two means estimate!

Let's try it out!

# Example: Heterogeneous Effects

- Suppose that:

$$E[\,y \mid D, X\,] = D\underset{\text{ATE}}{\boxed{\alpha}} + \alpha_0 + D\underset{\text{effect modifier}}{\boxed{W'\gamma}} + W'\beta$$

- What OLS is estimating is the solution to:

$$E\left[\left(y - D\tilde{\alpha} - \tilde{\alpha}_0 - W'\tilde{\beta}\right)\begin{pmatrix}D\\X\end{pmatrix}\right] = 0$$

$$E\left[\left(D\alpha + \alpha_0 + DW'\gamma + W'\beta - D\tilde{\alpha} - \tilde{\alpha}_0 - W'\tilde{\beta}\right)\begin{pmatrix}D\\X\end{pmatrix}\right] = 0$$

- Since $E[W] = 0$ and $W \perp D$: $\alpha = \tilde{\alpha}, \alpha_0 = \tilde{\alpha}_0$
- But $\tilde{\beta} = \beta + E[D]\gamma$
- Residual of OLS is $\epsilon = \widetilde{D}W'\gamma + \nu$ with $E[\nu|D,W] = 0$
- Then interaction term is:

$$\mathrm{E}\left[\beta'W\epsilon\widetilde{D}^2\right] = E\left[\widetilde{D}^3\right]\beta'E[WW']\gamma \neq 0$$
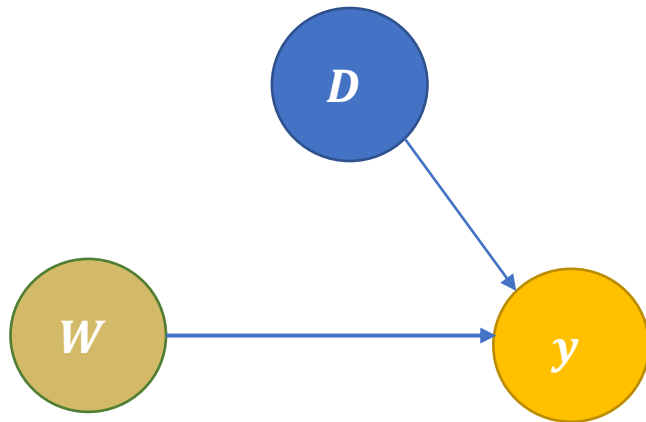
Even if you only care about ATE, if you have $p$ covariates and $p \ll n$ run OLS with interactive terms!

Guaranteed improved precision, plus can uncover potential dimensions of heterogeneity
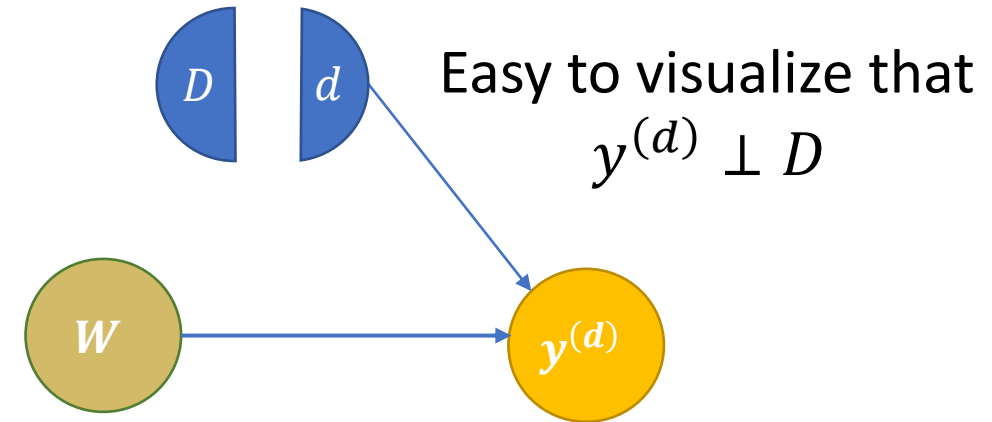
# RCTs and Causal Diagrams

# RCTs and Causal Diagrams

- Causal diagrams can help visualize how our assumptions imply the identification of a causal effect

- First instances in work of Sewall and Philip Wright'28

- Pioneered and fully developed by Pearl and Robins [80s-90s]



Causal Diagram (Graph)

Easy to visualize that
$$y^{(d)} \perp D$$

Single World Intervention Graph (SWIG)

# Limitations of RCTs

# Externalities, Stability and Equilibrium Effects

- Stable Unit Treatment Value Assumption (SUTVA)

- Implicit in our notation the potential outcome of a unit depends only on its own treatment

- This can be violated due to what is known as spillover effects, externalities or general equilibrium effects

- In vaccine example: if large fraction of population is treated, then the effect of vaccinating an additional unity changes, due to herd immunity

- In labor: if we make an intervention that incentivizes a large fraction of population to attend college, the college-wage premium will likely decrease

# Ethical, Practical and Generalizability Concerns

- **Ethical Concerns:** Many potential trial would correctly be judged unethical by a human subject trial review board; Key principles [78 Belmont report]: (i) respect for persons, (ii) beneficence, (iii) justice

- **Practical Concerns:** Practically infeasible due to high cost (e.g. expensive treatment, expensive data collection, low signal regime, long-term outcomes)

- **Generalizability:** Local population used for RCT might not generalize to broader population