

2. Randomization Inference

ISS5096 || ECI

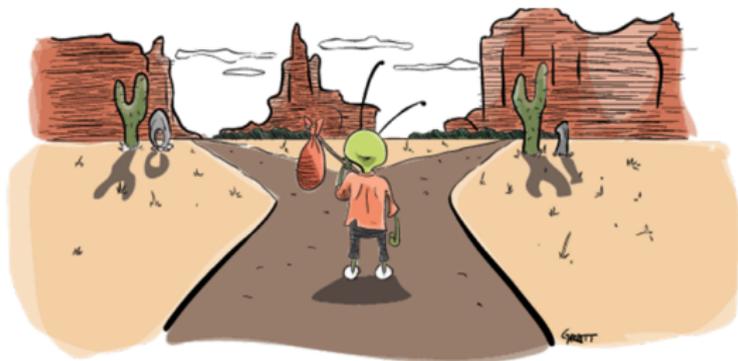
Jaewon (“Jay-one”) Yoo

National Tsing Hua University

Outline

1. Randomized Experiments
2. Randomization Inference
3. Application of Randomization Inference
4. In-Class Exercise

Where are we? Where are we going?



Source: *Chapter 1 of Mastering Metrics (Textbook 2)* by J. Angrist & J. Pischke

- Last time: defining causal effects as **contrasts of counterfactuals**.
- What can we learn about these contrasts from randomized experiments?
 - Message: randomization allows for inference under practically no assumptions.
- Useful to have notations for vector of all r.v.s:
 - Treatment: $\mathbf{D} = \{D_1, D_2, \dots, D_n\}$
 - Potential outcomes: $\mathbf{Y}(1) = \{Y_1(1), \dots, Y_n(1)\}$
 - Covariates: $\mathbf{X} = \{X_1, \dots, X_n\}$

1/ Randomized Experiments

Randomized Experiments

- **Experiment:** method/design where the researcher controls the treatment assignment.
 - p_i (i.e., the probability of treatment assignment) is controlled by and known to the researcher in an experiment.
- **Randomized experiment** is an experiment with two properties:
 1. **Positivity:** assignment is probabilistic (i.e., $0 < p_i < 1$).
 - No deterministic assignment.
 2. **Unconfoundedness:** $\mathbb{P}[D_i = 1 | Y(1), Y(0)] = \mathbb{P}[D_i = 1]$
 - Treatment assignment does not depend on any potential outcomes.
 - If patients were assigned to treatment group based on how researchers anticipate the patients will respond to the medication?

Context: Effect of AI Assistant on Learning Outcomes

- Q: Does AI assistant help students achieving higher academic performance?
 - Difficult with observational studies: usage of AI assistant correlated with lots of stuff!
- Randomized controlled trial can be helpful.
- Setup:
 - Units: students i
 - Treatment: assignment to a tutor with AI assistant ($D_i = 1$) or not ($D_i = 0$)
 - Outcome: student passes a standardized test ($Y_i = 1$) or not ($Y_i = 0$)
- If AI assistant \rightsquigarrow performance, we should see a difference between the treatment and control groups.

Why Randomize?

- Randomization makes treated and control groups **comparable!**
 - If both groups are random samples from all units in the study.
 - \rightsquigarrow **Balanced** on all variables: roughly, $N_{men}^T \approx N_{men}^C$, etc.
 - True for all observed & **unobserved pretreatment** variables.
 - Most importantly: potential outcomes are comparable by unconfoundedness:

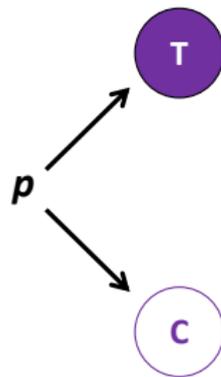
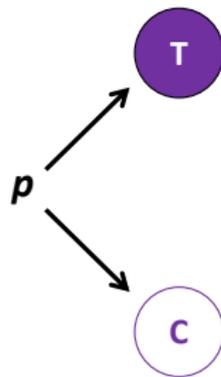
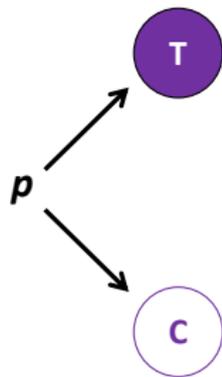
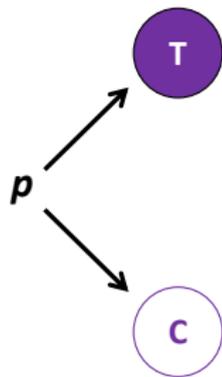
$$\mathbb{P}(Y_i(1) = 1 | D_i = 1) = \mathbb{P}(Y_i(1) = 1) = \mathbb{P}(Y_i(1) = 1 | D_i = 0)$$

- Caveat: groups are not comparable on **post-treatment** variables.
 - Recall consistency: $Y_i = D_i \cdot Y_i(1) + (1 - D_i) \cdot Y_i(0)$
 - $Y_i(1) \perp\!\!\!\perp D_i$ but not $Y_i \perp\!\!\!\perp D_i$
- Really talking about **ideal** randomized experiment:
 - Full compliance, no missing data.
 - Important to admit limitations: *external validity (i.e., generalizability), sample selection, Hawthorne effect*, etc.

Types of Experiments

- Experiments can be classified by their **assignment mechanism**.
 - What (random) function do we use to assign treatment?
- **Bernoulli (coin flips) experiment:**
 - Each unit is assigned $D_i = 1$ with prob. p_i independently.
 - Downside: “bad” randomization is possible (i.e., all assigned into treatment or control group)
- **Completely randomized experiment:**
 - Randomly sample n_1 units from the population to be treated.
 - For any given i , $p_i = \mathbb{P}(D_i = 1) = \frac{n_1}{n}$

Bernoulli Assignment



Completely Randomized Design



- Begin with $N = 6$, and say, we want to have $N_t = 3$
- Randomly pick 3 from $\{1, 2, 3, 4, 5, 6\}$: 2, 4, 5
- Fixed number of treated units induces dependence between D_i and D_j
 - Knowing unit 2 is treated \rightsquigarrow unit 3 is less likely to be treated.
 - Makes variance calculations tricky (we will come back to this)
- We can also randomize within groups (**block/stratified randomization**).
 - When we have 2 units in each block, this is referred to as a **pair-matched design**.

Example Data from AI RCT

Student	AI Assistant	Pass the Exam?	$Y_i(0)$	$Y_i(1)$
	D_i	Y_i		
1	1	0	?	0
2	1	0	?	0
3	0	1	1	?
4	1	0	?	0
5	1	1	?	1
6	0	1	1	?
7	0	0	0	?
8	1	1	?	1
9	0	1	1	?
10	0	0	0	?

- Students passed the exam 2/5 times with AI assistant vs. 3/5 times w/o AI assistant!
- Very small sample size \rightsquigarrow can we learn anything from this data?

2/ Randomization Inference

What is Randomization Inference?

- **Randomization inference (RI):** Inference based on different possible randomization of treatment.
 - R. Fisher: randomization is the “reasoned basis for inference.”
 - We can generate exact p-values for tests of a “sharp” null hypothesis.
 - Also called: **design-based inference**.
- Allows us to make **exact, distribution-free** inferences.
 - No reliance on normality, etc.
 - No reliance on large-sample approximations.
 - \rightsquigarrow non-parametric, but less flexible.

Brief Review of Hypothesis Testing

- RI is a type of hypothesis testing in a randomized experiment, so it is helpful to review!
 1. *Choose a null hypothesis:*
 - $H_0 : \beta_1 = 0$
 - Claim \Rightarrow no average treatment effect.
 - This is a claim we would like to reject.
 2. *Choose a test statistic:*
 - e.g., $Z_j = (X_j - \bar{X}) / (s / \sqrt{n})$
 3. *Determine the distribution of the test statistic under the null.*
 - Statistical thought experiment: if we knew the truth, what data should we expect?
 4. *Calculate the probability of the test statistics under the null.*
 - What is this called? **p-value**

Sharp Null Hypothesis of No Effect

- **Sharp null hypothesis:**

$$H_0 : \tau_i = Y_i(1) - Y_i(0) = 0 \quad \text{for all } i \quad (1)$$

- Specified for each unit i & assumes zero effect for every unit (hence, *sharp* vs. *weak* hypothesis, which assumes ATE = 0.)
- **What if treatment affected no one at all?**
- Implies no **average** treatment effect, but no ATE \Rightarrow sharp null.
 - Take a simple example with only two units: $\tau_1 = 1$ and $\tau_2 = -1$.
 - Here, $\tau = 0$, but the sharp null is violated.
- If the sharp null is true, than we know all the potential outcomes:

$$Y_i(1) = Y_i(0) = Y_i \quad (2)$$

Life Under the Sharp Null

We can use the sharp null (i.e., $Y_i(1) = Y_i(0) = Y_i$) to fill in the missing potential outcomes:

Student	AI Assistant D_i	Pass the Exam? Y_i	$Y_i(0)$	$Y_i(1)$
1	1	0	?	0
2	1	0	?	0
3	0	1	1	?
4	1	0	?	0
5	1	1	?	1
6	0	1	1	?
7	0	0	0	?
8	1	1	?	1
9	0	1	1	?
10	0	0	0	?

Life Under the Sharp Null

We can use the sharp null (i.e., $Y_i(1) = Y_i(0) = Y_i$) to fill in the missing potential outcomes:

Student	AI Assistant D_i	Pass the Exam? Y_i	$Y_i(0)$	$Y_i(1)$
1	1	0	0	0
2	1	0	0	0
3	0	1	1	1
4	1	0	0	0
5	1	1	1	1
6	0	1	1	1
7	0	0	0	0
8	1	1	1	1
9	0	1	1	1
10	0	0	0	0

Test Statistic

Test statistic

A test statistic is a known, scalar quantity calculated from the treatment assignments, observed outcomes, and possibly covariates: $T(\mathbf{D}, \mathbf{Y}, \mathbf{X})$.

- Test statistics measure how unusual the data is under the null.
- Typically measures the relationship between two variables (in causal inference).
- We want a test statistic with high **statistical power**:
 - Has large values when the null is likely false.
 - These large values are unlikely when the null is true.
- These will help us perform a test of the sharp null.

Null/Randomization Distribution

- What is the distribution of the test statistic under the sharp null?
 - If there was no effect (i.e., sharp null), what test statistics would we expect over different randomizations?
- **Key insight of RI:** Sharp null \rightsquigarrow treatment assignment doesn't matter.
 - Shuffling treatment vector won't change the outcomes!
 - $Y_i(1) = Y_i(0) = Y_i$
- **Randomization distribution:** distribution of T under the sharp null.

Calculate P-value

- How often would we get a test statistic this big or bigger if the sharp null holds?
- **Exact p-values:**

$$Pr(T \geq T^{obs}) = \frac{1}{K} \sum_{k=1}^K \mathbb{I}(T(\mathbf{d}, \mathbf{Y}, \mathbf{X}) \geq T^{obs}) \quad (3)$$

- How often T under different randomization larger than the T^{obs} divided by total number of randomizations (K)?
- Can be compared to a chosen threshold, α , to determine whether to reject the sharp null.

Randomization Inference Step-by-Step

1. Choose a sharp null hypothesis and a test statistic.
2. Calculate observed test statistic: $T^{obs} = T(\mathbf{D}, \mathbf{Y}, \mathbf{X})$.
3. Randomly select different treatment assignment vector $\tilde{\mathbf{D}}_1$.
 - e.g., if $\mathbf{D} = \{1, 1, 1, 0, 0, 0\}$, then $\tilde{\mathbf{D}}_1$ could be $\{1, 1, 0, 1, 0, 0\}$.
4. Calculate $\tilde{T}_1 = T(\tilde{\mathbf{D}}_1, \mathbf{Y}, \mathbf{X})$.
5. Repeat **steps 3** and **4** to get $\tilde{T} = \{\tilde{T}_1, \dots, \tilde{T}_K\}$ (randomization dist.)
6. Calculate the p-value: $p = \frac{1}{K} \sum_{k=1}^K \mathbb{I}(\tilde{T}_k \geq T^{obs})$.

Difference in Means as a Test Statistic

- Many different types of test statistics with different strengths/benefits.
- Natural (if not optimal): absolute difference-in-means estimator:

$$T_{\text{diff}} = \left| \frac{1}{n_1} \sum_{i=1}^N D_i Y_i - \frac{1}{n_0} \sum_{i=1}^N (1 - D_i) Y_i \right| \quad (4)$$

- Larger values of T_{diff} are evidence against the sharp null.
- A good estimator for constant treatment effects (across all units) with relatively few outliers in the potential outcomes.

3/ Application of Randomization Inference

Example: Encouraging Donation to NTHU

- Suppose we are targeting and encouraging 6 people to make a donation to National Tsing Hua University.
- As an encouragement, we send 3 of them an email with inspirational stories of learning from our alumni.
- Afterwards, we observe them giving between \$0 and \$5.
- Simple example to show the steps of RI in a concrete case.

Randomization Distribution

Unit	Email D_i	Donation in $US\$$ Y_i	$Y_i(0)$	$Y_i(1)$
Brian	1	3	(3)	3
Desi	1	5	(5)	5
Medjine	1	0	(0)	0
Natasha	0	4	4	(4)
Fifi	0	0	0	(0)
Matthew	0	1	1	(1)

$$T_{\text{diff}} = |8/3 - 5/3| = 1$$

with the observed treatment assignment, $D = \{1, 1, 1, 0, 0, 0\}$.

Randomization Distribution

Unit	Email \tilde{D}_i	Donation in US\$ Y_i	$Y_i(0)$	$Y_i(1)$
Brian	1	3	(3)	3
Desi	1	5	(5)	5
Medjine	0	0	0	(0)
Natasha	1	4	(4)	4
Fifi	0	0	0	(0)
Matthew	0	1	1	(1)

$$\tilde{T}_{\text{diff}} = |12/3 - 1/3| = 3.67 \text{ if } \tilde{D} = \{1, 1, 0, 1, 0, 0\}$$

$$\tilde{T}_{\text{diff}} = |8/3 - 5/3| = 1 \text{ if } \tilde{D} = \{1, 1, 1, 0, 0, 0\}$$

$$\tilde{T}_{\text{diff}} = |9/3 - 4/3| = 1.67 \text{ if } \tilde{D} = \{0, 1, 1, 1, 0, 0\}$$

Randomization Distribution

D_1	D_2	D_3	D_4	D_5	D_6	Diff. in Means
1	1	1	0	0	0	1.00
1	1	0	1	0	0	3.67
1	1	0	0	1	0	1.00
1	1	0	0	0	1	1.67
1	0	1	1	0	0	0.33
1	0	1	0	1	0	2.33
1	0	1	0	0	1	1.67
1	0	0	1	1	0	0.33
1	0	0	1	0	1	1.00
1	0	0	0	1	1	1.67
0	1	1	1	0	0	1.67
0	1	1	0	1	0	1.00
0	1	1	0	0	1	0.33
0	1	0	1	1	0	1.67
0	1	0	1	0	1	2.33
0	1	0	0	1	1	0.33
0	0	1	1	1	0	1.67
0	0	1	1	0	1	1.00
0	0	1	0	1	1	3.67
0	0	0	1	1	1	3.67

Application in R

```
1 > library(ri) # loading randomization inference package, ri.
2 > y <- c(3, 5, 0, 4, 0, 1)
3 > D <- c(1, 1, 1, 0, 0, 0)
4
5 # Diff. in means as a test stat.
6 > T_obs <- abs(mean(y[D == 1]) - mean(y[D == 0]))
7
8 # genperms() to generate all possible treatment assignments, tilde D.
9 > D_bold <- ri::genperms(D) # 20 diff. ways
10 > D_bold[, 1:10]
```

```
## [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9] [,10]
##    1    1    1    1    1    1    1    1    1    1
##    1    1    1    1    0    0    0    0    0    0
##    1    0    0    0    1    1    1    0    0    0
##    0    1    0    0    1    0    0    1    1    0
##    0    0    1    0    0    1    0    1    0    1
##    0    0    0    1    0    0    1    0    1    1
```

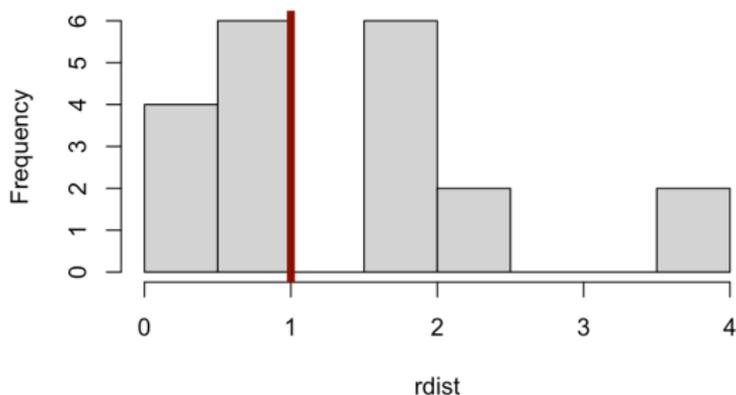
Application in R

```
1 # Calculating the test statistics to obtain the rand. dist.
2 > rdist <- rep(NA, times = ncol(D_bold))
3 > for (i in seq_len(ncol(D_bold))) {
4     D_tilde <- D_bold[,i]
5     rdist[i] <- abs(mean(y[D_tilde == 1]) - mean(y[D_tilde == 0]))
6 }
7 > round(rdist, digits = 3) # print randomization distribution

## [1] 1.000 3.667 1.000 1.667 0.333 2.333 1.667 0.333 1.000 1.667
## [11] 1.667 1.000 0.333 1.667 2.333 0.333 1.667 1.000 3.667 1.000
```

Example: Computation in R Cont.

```
1 > hist(x = rdist) # visualize the randomization distribution
2 > abline(v = 1, col = "darkred", lwd = 5) # Red line = observed T
```



```
1 > mean(rdist >= T_obs) # Computing p-value
```

```
## [1] 0.8
```

- α levels? Not enough evidence to reject the sharp null of no effect! (or not enough evidence to say encouragement helped)

Computational Burden in RI

Computing the exact randomization distribution is not always possible:

- $n = 6$ and $n_1 = 3 \rightsquigarrow 20$ assignment vectors (6C_3).
- $n = 10$ and $n_1 = 5 \rightsquigarrow 252$ vectors (${}^{10}C_5$).
- $n = 100$ and $n_1 = 50 \rightsquigarrow 1.009 \times 10^{29}$ vectors (${}^{100}C_{50}$).
- Workaround: sampling!
 - Take K samples from the treatment assignment.
 - Compute the randomization distribution in the K samples.
 - Tests are no longer exact, but bias is under your control!
(we have control over K)

Other Test Statistics

- The difference in means is great for when effects are:
 - constant (i.e., $\tau = \tau_i$ for all i) and additive (e.g., $Y(1) - Y(0)$ and not $Y(1)/Y(0)$)
 - few outliers in the data
- Outliers \rightsquigarrow more variation in the randomization distribution.
- What about alternative test statistics?

Transformations

- What if there was a constant multiplicative effect: $Y_i(1)/Y_i(0) = C$?
- T_{diff} will have low power in this case as it does not capture the nature of the causal effect.
- \rightsquigarrow Transform the observed outcome using the natural logarithm:

$$T_{\log} = \left| \frac{1}{n_1} \sum_{i=1}^N D_i \log(Y_i) - \frac{1}{n_0} \sum_{i=1}^N (1 - D_i) \log(Y_i) \right| \quad (5)$$

- Log-transforming the outcomes linearizes the multiplicative nature of the causal effect.
- Recall $\log\left(\frac{a}{b}\right) = \log(a) - \log(b)$.
- Also useful for skewed distribution of the outcomes.

p.s., See [Chen & Roth \(2024, QJE\)](#): log-like transformations, e.g., $\log(1 + Y)$, with zero-valued outcomes \rightsquigarrow ATE is arbitrarily sensitive to units of Y .

Difference in Median/Quantile

- To protect against outliers: quantiles.
- Difference in medians:

$$T_{\text{med}} = |\text{med}(\mathbf{Y}_t) - \text{med}(\mathbf{Y}_c)| \quad (6)$$

- where $\mathbf{Y}_t = Y_i$ for $i : D_i = 1$ and $\mathbf{Y}_c = Y_i$ for $i : D_i = 0$
- Recall that the median is the 0.5 quantile.
- Could use other quantiles (e.g., the 0.25 quantile or the 0.75 quantile).

Two-Sided or One-Sided?

- So far, we have defined all test statistics as **absolute values**.
- \rightsquigarrow testing against a **two-sided** alternative hypothesis:

$$H_0 : \tau_j = 0 \quad \forall i \quad H_1 : \tau_j \neq 0 \text{ for some } i$$

- What about a **one-sided** alternative?

$$H_0 : \tau_j = 0 \quad \forall i \quad H_1 : \tau_j > 0 \text{ for some } i$$

- For these, use a test statistic that is bigger under the alternative:

$$T_{\text{diff}}^* = \bar{Y}_t - \bar{Y}_c$$

- No absolute value \rightsquigarrow only large *positive* values count as evidence.

Other Sharp Nulls

- Sharp null of no effect is not the only sharp null.
- Sharp null in general is one of a **constant additive effect**: $H_0 : \tau_i = 0.2$.
 - Implies that $Y_i(1) = Y_i(0) + 0.2$.
 - Can still calculate all the potential outcomes!
- More generally, we could have $H_0 : \tau_i = \tau_0$ for a fixed τ_0 .

Testing Non-Zero Sharp Nulls

- Suppose that we had: $H_0 : \tau_i = Y_i(1) - Y_i(0) = 1$

Unit	Email D_i	Donation Y_i	$Y_i(0)$	$Y_i(1)$	Adjusted $Y_i - D_i\tau_0$
Brian	1	3	(2)	3	2
Desi	1	5	(4)	5	4
Medjine	1	0	(-1)	0	-1
Natasha	0	4	4	(5)	4
Fifi	0	0	0	(1)	0
Matthew	0	1	1	(2)	1

- Assignments will now affect Y_i .
- Solution: use **adjusted outcomes**, $Y_i^* = Y_i - D_i\tau_0$.
- Now, just test sharp null of no effect for Y_i^* .
 - $Y_i^*(1) = Y_i(1) - 1 \times 1 = Y_i(0)$
 - $Y_i^*(0) = Y_i(0) - 0 \times 1 = Y_i(0)$
 - $\tau_i^* = Y_i^*(1) - Y_i^*(0) = 0$

Point Estimates via RI

- Is it possible to get point estimates?
- Not really the point of RI, but still possible:
 1. Create a grid of possible sharp null hypotheses.
 2. Calculate p-values for each sharp null.
 3. Pick the value that is “least surprising” under the null.
- Usually this means selecting the τ_0 with the **highest p-value**.

Using Regression in RI

- We can also use covariates to improve power.
- For instance, run an OLS regression:

$$(\hat{\beta}_0, \hat{\beta}_D, \hat{\beta}_X) = \arg \min_{\beta_0, \beta_D, \beta_X} \sum_{i=1}^n (Y_i - \beta_0 - \beta_D \cdot D_i - \beta_X \cdot X_i)^2$$

and use $T_{\text{ols}} = \hat{\beta}_D$ as our test statistic.

- RI is justified **even if the model is wrong!**
 - OLS is just another way to generate a scalar test statistic.
 - Inference comes from randomization, not from distributional assumptions.
- If the model is predictive of POs, then T_{ols} will have **higher power**.
 - Intuition: controlling for X_i reduces residual variance \rightsquigarrow easier to detect smaller effects.

Next Up

- Inference of the Average Treatment Effect (ATE)
- J. Neyman's way of thinking about the question of causal effect.

4/ In-Class Exercise

In-Class Exercise: Randomization Inference

A small online retailer A/B tested a new product page on 10 customers (5 treated, 5 control; Y = number of items purchased). **This R script** implements RI with $T_{\text{diff}} = |\bar{Y}_t - \bar{Y}_c|$ (two-sided p-value = 0.95).

1. **Different test statistic.** Replace diff-in-means with **diff-in-medians**:
 $T_{\text{med}} = |\text{med}(Y_t) - \text{med}(Y_c)|$. Compare the p-value to the baseline. Look at the data, why do they differ?
2. **One-sided test.** Using your diff-in-medians code from Task 1, remove the absolute value from **both** T^{obs} and the randomization distribution. Compare the p-value to Task 1. Why did it change?
3. **Changing an outcome.** Pick one person and change their Y value. Re-run RI. Which of the following changed, and why?
 - (i) The set of possible treatment assignments
 - (ii) The randomization distribution of T_{diff}
 - (iii) The observed test statistic T^{obs}

Have a great weekend! :)

Contact Information:

jaewon.yoo@iss.nthu.edu.tw

<https://j1yoo.github.io/>

