Conditioning on and controlling for variates via cumulative differences

Measuring calibration, reliability, biases, and other treatment effects

Mark Tygert

Meta

mark@tygert.com

Outline

- 1 Calibration and reliability
 - Probabilistic predictions
 - Cumulative differences
 - Binned responses
- 2 Fairness and matching scores
 - Deviation of a subpopulation from the full pop.
 - Deviation between two disjoint subpopulations
 - Controlling for multiple covariates

3 Extensions

- Weighted sampling
- Avoiding randomization
- Interactive plots

イロト イボト イヨト イヨト

Probabilistic predictions Cumulative differences Binned responses

Calibration and reliability

Calibration and reliability

Mark Tygert Conditioning & controlling via cumulative differences 3/70

イロト イボト イヨト

DQC2

Calibration and reliability

Fairness and matching scores Extensions Probabilistic predictions Cumulative differences Binned responses

Probabilistic predictions

Probabilistic predictions

Mark Tygert Conditioning & controlling via cumulative differences 4/70

イロト イヨト イヨト

Probabilistic predictions Cumulative differences Binned responses

Predictions with a given probability

- 1 Consider an S = 30% chance of snow.
- ② The reality will be either R = 1 (it actually did snow) or R = 0 (it did not snow).
- **3** A random variable R which follows the Bernoulli distribution with expected value S takes R = 1 with probability S and R = 0 with probability 1 - S.
- Given *n* independent random results, also known as "responses," *R*₁, *R*₂, ..., *R_n*, the prediction is perfectly calibrated if *S* = 30% of *R*₁, *R*₂, ..., *R_n* are equal to 1, and the rest are equal to 0.

イロト イポト イヨト イヨト 三日

Probabilistic predictions Cumulative differences Binned responses

Predictions with different probabilities

- **1** Probabilities of success S_1, S_2, \ldots, S_n ; the probabilities are the expected values of Bernoulli distributions.
- 2 Responses R₁, R₂, ..., R_n; the responses are independent Bernoulli random variables (R_k and S_k come as a pair).
- ③ Probabilities of success are also known as "scores" and will be viewed as deterministic during this tutorial.
- ④ Responses are also known as "results" or "outcomes" and viewed as random and probabilistically independent.
- **5** Goal: gauge deviation of the responses from the scores.

イロト イポト イヨト イヨト 三日

Probabilistic predictions Cumulative differences Binned responses

Without loss of generality

- **1** Re-order the indices such that $S_1 \leq S_2 \leq \cdots \leq S_n$, preserving the pairing of R_k with S_k for every $k = 1, 2, \ldots, n$.
- 2 Randomly perturb the scores such that $S_1 < S_2 < \cdots < S_n$ if the inequalities were not already strict.
- ③ There exists an alternative to this random perturbation that is more complicated but that avoids randomization altogether. Anyone interested can consult the later section, "Extensions," specifically its subsection, "Avoiding randomization."

イロト イポト イヨト イヨト 三日

Calibration and reliability

Fairness and matching scores Extensions Probabilistic predictions Cumulative differences Binned responses

Cumulative differences

Cumulative differences

Mark Tygert Conditioning & controlling via cumulative differences 8/70

イロト イボト イヨト イヨト

DQC2

Э

Probabilistic predictions Cumulative differences Binned responses

Cumulative aggregation

The cumulative difference of responses from scores is

$$C_{k} = \frac{1}{n} \sum_{j=1}^{k} (R_{j} - S_{j})$$
(1)

for
$$k = 1, 2, ..., n$$
.

The expected slope of a graph of C_k versus k/n from j = k - 1 to j = k is

$$\frac{\mathbb{E}[C_k - C_{k-1}]}{k/n - (k-1)/n} = \mathbb{E}[R_k] - S_k$$
(2)

for k = 1, 2, ..., n; this $\mathbb{E}[R_k] - S_k$ is the expected difference from perfect calibration.

Mark Tygert Conditioning & controlling via cumulative differences 9/70

イロト イヨト イヨト

Probabilistic predictions Cumulative differences Binned responses

Miscalibration is slope

- The slope of a secant line connecting two points on the graph of C_k versus k/n becomes the average miscalibration over the long range of k between those points. This follows from combining the law of large numbers with the last formula on the previous slide.
- Slope is easy to perceive with quantitative precision even when the constant offsets of the secant lines are irrelevant.
- 3 Miscalibration at an index k of interest determines the slope, unpolluted by accumulation from indices other than that of interest (the other indices affect only the constant offset).

イロト イ伊ト イヨト イヨト

Probabilistic predictions Cumulative differences Binned responses

Empirical plot with $n = 2^{15} = 32,768$



Mark Tygert Conditioning & controlling via cumulative differences 11 / 70

イロト イヨト イヨト

nan

3

Calibration and reliability

Fairness and matching scores Extensions Probabilistic predictions Cumulative differences Binned responses

Plot of the ground truth



Mark Tygert Conditioning & controlling via cumulative differences 12 / 70

イロト イヨト イヨト

990

Э

Probabilistic predictions Cumulative differences Binned responses

Kolmogorov-Smirnov and Kuiper metrics

- Good calibration corresponds to a flat, horizontal graph.
- 2 Deviation from flat and horizontal measures miscalibration.
- Two scalar statistics which summarize the deviations from 0 are the max. absolute deviation and the range of deviations.
- The max. absolute deviation is Kolmogorov's and Smirnov's

$$C_{\text{MAD}} = \max_{1 \le k \le n} |C_k|.$$
(3)

5 The range of deviations is Kuiper's

$$C_{\text{range}} = \max_{0 \le k \le n} C_k - \min_{0 \le k \le n} C_k, \tag{4}$$

where $C_0 = 0$.

Mark Tygert Conditioning & controlling via cumulative differences 13 / 70

イロト イ伊ト イヨト イヨト

Probabilistic predictions Cumulative differences Binned responses

Another expression for the Kuiper metric

The absolute value of the total miscalibration $\sum_{j \in I} (R_j - S_j)/n$ over any interval I of indices is less than or equal to C_{range} ; indeed, the Kuiper metric C_{range} is equal to the maximum of the absolute value of the total miscalibration over any interval of indices:

$$C_{\text{range}} = \max_{I} \left| \frac{1}{n} \sum_{j \in I} (R_j - S_j) \right|, \qquad (5)$$

where the maximum is taken over every interval I of indices; summing over the index $j \in I$ means summing over the indices $j = \min I$, $1 + \min I$, $2 + \min I$, $3 + \min I$, ..., $\max I$.

Mark Tygert Conditioning & controlling via cumulative differences 14 / 70

イロト イヨト イヨト

Probabilistic predictions Cumulative differences Binned responses

Significance testing

- The P-value for a given assertion asserting a null hypothesis is 1 the probability that a specified measure of deviation from the assertion is greater than or equal to that measure for the actual observations. The metric could be C_{range} , for example.
- 2 The null hypothesis relevant here is of perfect calibration, namely that every response R_k comes from the Bernoulli distribution whose expected value is the associated score S_k .
- A small P-value indicates that the observations are inconsistent with the assertion (the null hypothesis).
- ④ If P is the P-value, then 1 P is the attained confidence level for believing the assertion (the null hypothesis). This value 1 - P is not exactly the probability of the null hypothesis being true, but is a monotonically increasing function of such a probability, for most formalizations of such a probability.

DQC2

Э

Probabilistic predictions Cumulative differences Binned responses

P-values

Under the null hypothesis of perfect calibration E[R_k] = S_k for k = 1, 2, ..., n, C_{MAD}/σ converges in distribution to the max. absolute value of the standard Brownian motion over [0, 1], and C_{range}/σ converges in distribution to the Brownian motion's range, where σ² is the total expected variance

$$\sigma^2 = \sum_{k=1}^n \frac{S_k(1-S_k)}{n^2},$$
(6)

assuming that $\max_{1 \le k \le n} S_k(1 - S_k) / \sum_{j=1}^n S_j(1 - S_j)$ converges to 0 as *n* becomes arbitrarily large.

⁽²⁾ $C_{\rm MAD}$ is $0.01243/\sigma = 5.512$, $C_{\rm range}$ is $0.02291/\sigma = 10.16$, and the corresponding asymptotic P-values are 7.1E–08 and less than the machine precision, for the earlier example.

Mark Tygert Conditioning & controlling via cumulative differences 16 / 70

(日)

DQ C

Probabilistic predictions Cumulative differences Binned responses

Why Brownian motion?

- Under the null hypothesis of perfect calibration, the graph of the line segments connecting (S₀, C₀), (S₁, C₁), ..., (S_n, C_n) are the definition of a driftless random walk (driftless means that the expected value of C_k C_{k-1} = (R_k S_k)/n is 0 for all k = 1, 2, ..., n, which follows from the null hypothesis that the expected value of R_k is S_k). The graph of the line segments converges in distribution to the graph of Brownian motion, under the condition on the previous slide (as n → ∞).
- 2 The variance of a Bernoulli distribution whose expected value is S_k is $S_k(1 - S_k)$. The variance of a sum of independent random variables is the sum of the variances of the random variables. Therefore, the variance of $C_n = \sum_{k=1}^n (R_k - S_k)/n$ is σ^2 defined on the previous slide. Dividing by σ standardizes the Brownian motion to be over the unit interval [0, 1].

Mark Tygert Conditioning & controlling via cumulative differences 17 / 70

- コット (中) (日) (日) (日)

SQR

Probabilistic predictions Cumulative differences Binned responses

Error bars

- The standard deviation of the standard Brownian motion over the unit interval [0, 1] ranges from 0 at 0 to 1 at 1.
- ⁽²⁾ The cumulative graphs plot the cumulative differences without normalizing by σ ; standardizing the Brownian motion associated with the cumulative graphs would require normalization by σ .
- **3** The cumulative plots (including the two from earlier slides) display a triangle at the origin whose tip-to-tip height is 4σ , corresponding to a roughly 95% confidence interval.
- ④ The triangle at the origin indicates the length scale of statistically significant deviations from zero, effectively showing the effect of normalization by σ.

Mark Tygert Conditioning & controlling via cumulative differences 18 / 70

イロト イボト イヨト

SQR

Probabilistic predictions Cumulative differences Binned responses

Binned responses

Binned responses

Mark Tygert Conditioning & controlling via cumulative differences 19 / 70

《日》《圖》《臣》《臣》

990

Э

Probabilistic predictions Cumulative differences Binned responses

Binning (instead of the newer cumulative approach)

 $S_1^1 < S_1^2 < \cdots < S_1^{n_1}$ will be the least n_1 of the scores, $S_2^1 < S_2^2 < \cdots < S_2^{n_2}$ will be the next n_2 of the scores, ..., $S_m^1 < S_m^2 < \cdots < S_m^{n_m}$ will be the greatest n_m of the scores. Maintain $n = \sum_{i=1}^{m} n_i$ and calculate the average score

$$\tilde{S}_j = \frac{1}{n_j} \sum_{k=1}^{n_j} S_j^k \tag{7}$$

and the average response

$$\tilde{R}_j = \frac{1}{n_j} \sum_{k=1}^{n_j} R_j^k \tag{8}$$

for $i = 1, 2, \ldots, m$.

Mark Tygert

イロト イヨト イヨト Conditioning & controlling via cumulative differences 20 / 70

Probabilistic predictions Cumulative differences Binned responses

The "reliability diagram" is the graph of the average response \hat{R}_j versus the average score \tilde{S}_j , together with a diagonal line from the origin (0,0) to the point (1,1); the diagonal line corresponds to perfect calibration, for which the response is equal to the score. The reliability diagram is the traditional graphical method for assessing calibration.

イロト イポト イヨト イヨト

Probabilistic predictions Cumulative differences Binned responses

Same range of probabilities for each bin, 4-32 bins



Mark Tygert

イロト イボト イヨト Conditioning & controlling via cumulative differences 22 / 70

nac

Probabilistic predictions Cumulative differences Binned responses

Same range of probabilities for each bin, 64-256 bins



Mark Tygert

イロト イヨト イヨト Conditioning & controlling via cumulative differences 23 / 70

nac

Probabilistic predictions Cumulative differences Binned responses

Same number of responses for each bin, 4-32 bins



Mark Tygert

<ロト < 団ト < 巨ト < 巨ト -Conditioning & controlling via cumulative differences 24 / 70

nac

Э

Probabilistic predictions Cumulative differences Binned responses

Same number of responses for each bin, 64-256 bins



Mark Tygert

イロト イボト イヨト Conditioning & controlling via cumulative differences 25 / 70

nac

Probabilistic predictions Cumulative differences Binned responses

Empirical calibration errors

The empirical calibration errors (ECEs) are the Riemann sums

$$ECE^{1} = \sum_{j=1}^{m} (S_{j+1}^{1} - S_{j}^{1}) \left| \tilde{R}_{j} - \tilde{S}_{j} \right| = \sum_{j=1}^{m} (S_{j+1}^{1} - S_{j}^{1}) \left| \sum_{k=1}^{n_{j}} \frac{R_{j}^{k} - S_{j}^{k}}{n_{j}} \right|$$
(9)

and

$$ECE^{2} = \sum_{j=1}^{m} (S_{j+1}^{1} - S_{j}^{1}) \left| \tilde{R}_{j} - \tilde{S}_{j} \right|^{2} = \sum_{j=1}^{m} (S_{j+1}^{1} - S_{j}^{1}) \left| \sum_{k=1}^{n_{j}} \frac{R_{j}^{k} - S_{j}^{k}}{n_{j}} \right|^{2}$$
(10)

where $S_{m+1}^1 = 1$, and the bin width $(S_{j+1}^1 - S_j^1)$ can be replaced with 1/m when the average scores are roughly equispaced on [0, 1]. The ECEs are the traditional metrics for assessing calibration.

Mark Tygert Conditioning & controlling via cumulative differences 26 / 70

イロト イポト イヨト

Probabilistic predictions Cumulative differences Binned responses

Integrated calibration index (ICI) = ECE

While "ECE" is the accepted standard terminology, biostatisticians have recently re-introduced the ECE, referring to the ECE as the "integrated calibration index" (ICI). There is now a large literature that uses "ICI" rather than "ECE."

Beware, too, that the abbreviation, "ECE," can refer to "empirical," "estimated," "expected," or "experimental" calibration errors; all these possibilities for the first letter in "ECE" pertain to the same mathematical formulas from the previous slide.

Mark Tygert Conditioning & controlling via cumulative differences 27 / 70

イロト イポト イヨト イヨト 三日

Probabilistic predictions Cumulative differences Binned responses

ECEs vary wildly with the choice of bins

 $n = 2^{15}$; earlier, the scores were square rooted from equispaced.

Mark Tygert

イロト イ伊ト イヨト イヨト Conditioning & controlling via cumulative differences 28 / 70

nac

Э

Probabilistic predictions Cumulative differences Binned responses

Comparison of ECEs to cumulative summary statistics

- If the number of observations per bin stays bounded on a set interval of scores as the max. width of a bin converges to 0, then the ECEs hit a noise floor, which prevents the ECEs from distinguishing a fixed imperfectly calibrated distribution from perfect calibration, even in the limit of infinite sample size n.
- This highlights a trade-off inherent to binning (or kernel density estimation) decreasing the width of bins or kernels resolves finer variations as a function of score at the expense of averaging away less noise, while increasing the width of bins sacrifices resolving power but boosts statistical confidence.
- 3 The cumulative statistics C_{range} and C_{MAD} of Kuiper and of Kolmogorov and Smirnov have no such explicit trade-off.

Mark Tygert Conditioning & controlling via cumulative differences 29 / 70

イロト イ伊ト イヨト イヨト

Deviation of a subpopulation from the full pop. Deviation between two disjoint subpopulations Controlling for multiple covariates

Fairness and matching scores

Fairness and matching scores

Mark Tygert Conditioning & controlling via cumulative differences 30 / 70

イロト イヨト イヨト

Deviation of a subpopulation from the full pop. Deviation between two disjoint subpopulations Controlling for multiple covariates

Deviation of a subpopulation from the full pop.

Deviation of a subpopulation from the full pop.

Mark Tygert Conditioning & controlling via cumulative differences 31 / 70

イロト イポト イヨト イヨト

Deviation of a subpopulation from the full pop. Deviation between two disjoint subpopulations Controlling for multiple covariates

Fairness for a subpopulation

- Goal: gauge deviation of the responses of the subpopulation from those of the full population.
- 2 Compare only individuals who are comparable.
- 3 Use the scores to match up individuals being compared.
- ④ Deviation of the subpopulation from the full population may indicate inequities in or other effects on outcomes or treatment of members of the subpopulation.
- Somparing individuals conditional on their scores ensures that the score variate is not the cause of differences detected.
- Examples of scores in medicine and the social sciences include age, income, or predicted probability (as with calibration).

イロト イヨト イヨト

Deviation of a subpopulation from the full pop. Deviation between two disjoint subpopulations Controlling for multiple covariates

Respecting the subpopulation's sampling

- In addition to the scores S₁ < S₂ < ··· < S_p and corresponding responses R₁, R₂, ..., R_p for a full population of p individuals, consider a subpopulation specified by n indices i₁ < i₂ < ··· < i_n.
- 2 Define \tilde{R}_{i_k} to be the average of the full population's responses whose corresponding scores are closer to S_{i_k} than to any other of the subpopulation's scores S_{i_1} , S_{i_2} , ..., S_{i_n} .
- ③ Calculate the cumulative differences

Mark Tygert

$$C_j = \frac{1}{n} \sum_{k=1}^{j} \left(R_{i_k} - \tilde{R}_{i_k} \right) \tag{11}$$

for j = 1, 2, ..., n.

Conditioning & controlling via cumulative differences 33 / 70

イロト イポト イヨト イヨト 三日

Deviation of a subpopulation from the full pop. Deviation between two disjoint subpopulations Controlling for multiple covariates

Slope is deviation of the subpop. from the full pop.

The expected slope of a graph of C_k versus k/n from j = k - 1 to j = k is

$$\frac{\mathbb{E}[C_k - C_{k-1}]}{k/n - (k-1)/n} = \mathbb{E}[R_{i_k}] - \mathbb{E}[\tilde{R}_{i_k}]$$
(12)

for k = 1, 2, ..., n; this $\mathbb{E}[R_{i_k}] - \mathbb{E}[\tilde{R}_{i_k}]$ is the expected deviation of the subpopulation from the full population.

The slope of a secant line connecting two distant points on the graph is the average deviation of the subpopulation from the full population over the subpopulation's indices between those points.

This analyzes the difference of the subpopulation's responses from the full population's while "controlling for" (also known as "conditioning on") scores, matching individuals with similar scores.

Mark Tygert Conditioning & controlling via cumulative differences 34 / 70

イロト イポト イヨト イヨト 三日

Deviation of a subpopulation from the full pop. Deviation between two disjoint subpopulations Controlling for multiple covariates

Wild boars vs. full ImageNet

The next two slides display reliability diagrams and cumulative graphs for the difference in accuracy of classification for a subpopulation from the full population in popular training data from computer vision known as ImageNet. The black points and lines correspond to the images of wild boars in ImageNet, while the gray points and lines correspond to the full data set. The score is the predicted probability corresponding to the class predicted to be most likely; the score is a measure of confidence in the prediction. The reliability diagrams vary significantly as the number of bins varies, while the cumulative graphs are reasonably easy to interpret; in particular, the drop-off for the highest scores is easy to quantify in the cumulative graphs. The reliability diagrams might even look inconsistent with each other without looking at the cumulative graphs (the latter resolve the inconsistencies at a glance).

Mark Tygert Conditioning & controlling via cumulative differences 35 / 70

イロト イヨト イヨト

Deviation of a subpopulation from the full pop. Deviation between two disjoint subpopulations Controlling for multiple covariates

Wild boars vs. full ImageNet (same width for each bin)

Mark Tygert

イロト イヨト イヨト Conditioning & controlling via cumulative differences 36 / 70

nac

Deviation of a subpopulation from the full pop. Deviation between two disjoint subpopulations Controlling for multiple covariates

Wild boars vs. full ImageNet (adaptive widths of bins)

Mark Tygert

イロト イ伊ト イヨト イヨト Conditioning & controlling via cumulative differences 37 / 70

nac

3

Deviation of a subpopulation from the full pop.

Deviation between two disjoint subpopulations Controlling for multiple covariates

Wild boars' empirical calibration errors

What is the right choice of bins? Which is the right value for m?

Mark Tygert

イロト イポト イヨト イヨト Conditioning & controlling via cumulative differences 38 / 70

nac

Э

Deviation of a subpopulation from the full pop. **Deviation between two disjoint subpopulations** Controlling for multiple covariates

Deviation between two disjoint subpopulations

Deviation between two disjoint subpopulations

Mark Tygert Conditioning & controlling via cumulative differences 39 / 70

・ロト ・ 同ト ・ ヨト ・ ヨト

Caveats

Deviation of a subpopulation from the full pop. **Deviation between two disjoint subpopulations** Controlling for multiple covariates

- Comparing a subpopulation to the full population while controlling for scores is always legitimate — every member of the subpopulation matches up with at least one member of the full population, namely, that very same member (that is the very definition of "sub-population"!).
- ② Comparing different subpopulations directly might be absurd. Indeed, the scores for one subpopulation might even all be less than all the scores for another subpopulation; there is no way to match up members according to score in such a case.
- ③ Nevertheless, it can make sense to compare subpopulations directly (rather than via the individual subpopulations' deviations from the full population). An example is comparing the subpopulation of individuals who received a certain medical treatment to the untreated subpopulation.

Mark Tygert Conditioning & controlling via cumulative differences 40 / 70

・ロト ・ 一下 ・ ト ・ ト ・

Э

Deviation of a subpopulation from the full pop. **Deviation between two disjoint subpopulations** Controlling for multiple covariates

Finest-possible binning

- - The crosses ("x") indicate the scores for subpopulation 0 while the circles ("o") indicate the scores for subpopulation 1.
 - The averages of the scores for subpopulation k corresponding to the indicated blocks of observed scores are S^k₀, S^k₁, ..., S^k₉, for k = 0, 1.
 - 3 The averages of the responses for subpop. k corresponding to the indicated blocks of observed scores are R₀^k, R₁^k, ..., R₉^k, for k = 0, 1.

イロト イポト イヨト

Deviation of a subpopulation from the full pop. **Deviation between two disjoint subpopulations** Controlling for multiple covariates

Average forward and backward differences

Define D_{2k} to be the average of the even indexed forward and backward differences and D_{2k+1} to be the average of the odd indexed forward and backward differences:

Conditioning & controlling via cumulative differences 42/70

イロト イ伊ト イヨト イヨト

Deviation of a subpopulation from the full pop. **Deviation between two disjoint subpopulations** Controlling for multiple covariates

Accumulation of averaged differences

Define the cumulative differences

$$C_j = \frac{1}{n} \sum_{k=0}^{j-1} D_k \tag{13}$$

for j = 1, 2, ..., n.

The graph of C_k versus k/n and the metrics C_{MAD} and C_{range} — defined via the same formulas (3) and (4) as on slide 13 above — admit interpretations analogous to those discussed on earlier slides for the related settings treated there.

Mark Tygert Conditioning & controlling via cumulative differences 43 / 70

イロト イヨト イヨト

Deviation of a subpopulation from the full pop. Deviation between two disjoint subpopulations Controlling for multiple covariates

Controlling for multiple covariates

Controlling for multiple covariates

Mark Tygert Conditioning & controlling via cumulative differences 44 / 70

イロト イポト イヨト イヨト

Deviation of a subpopulation from the full pop. Deviation between two disjoint subpopulations Controlling for multiple covariates

Multiple covariates require multiple dimensions

- "Covariates" is yet another name for the "input variables,"
 "independent variables," or "control variables" in a regression.
- ② Earlier slides considered only a single real-valued covariate; the scores specified the values of the covariate.
- When there are multiple real-valued covariates, their values fill a multidimensional vector space; denoting by d the number of scalar covariates, the vector space of all their possible values will be d-dimensional.
- There is a natural notion of locality in d dimensions that looks over all possible length scales, hierarchically organizing d-dimensional space into the canonical "dyadic tree" defined on the following slide.

イロト イヨト イヨト

SQR

Deviation of a subpopulation from the full pop. Deviation between two disjoint subpopulations Controlling for multiple covariates

Dyadic trees

- The dyadic tree in a d-dimensional vector space is the partition of space obtained by recursively subdividing equally along every coordinate axis during every split.
- In one dimension, the canonical dyadic tree is the balanced binary tree. In two dimensions, the canonical dyadic tree is the canonical quad tree. In three dimensions, the canonical dyadic tree is the canonical oct tree.

four levels of the canonical quad tree:

< □ ▷ < @ ▷ < \bar{B} ▷ < \bar{B} ▷ < \bar{B} ▷ \bar{B} ▷ \bar{B} \bar{B} \bar{B} \bar{D} \bar{D}

Deviation of a subpopulation from the full pop. Deviation between two disjoint subpopulations Controlling for multiple covariates

Hilbert space-filling curves

The Hilbert curve puts a total order on the unit hypercube $[0, 1]^d$ via the depth-first traversal of the canonical 2^d -ary (dyadic) tree. With d = 2, an approximation with 255 line segments is

Mark Tygert Conditioning & controlling via cumulative differences 47 / 70

< ロ > < 同 > < 回 > < 回 >

Deviation of a subpopulation from the full pop. Deviation between two disjoint subpopulations Controlling for multiple covariates

Reduction to one dimension

A space-filling curve h (e.g., Hilbert's) maps continuously from the unit interval [0, 1] onto (surjectively) the unit hypercube $[0, 1]^d$. This ensures that, given a function f on the unit hypercube $[0, 1]^d$, local averages of the composition $f \circ h$ are also local averages of f; as usual, the definition of the composition is $(f \circ h)(t) = f(h(t))$ for all t in the unit interval [0, 1].

Given points in the unit hypercube $[0,1]^d$, use the Hilbert curve to map each of the points to a score in the unit interval [0,1], then use the methods of the earlier slides with these scalar scores.

イロト イヨト イヨト

Deviation of a subpopulation from the full pop. Deviation between two disjoint subpopulations Controlling for multiple covariates

Hilbert-curve score for folding or normal as brightness

Mark Tygert Conditioning & controlling via cumulative differences 49 / 70

Deviation of a subpopulation from the full pop. Deviation between two disjoint subpopulations Controlling for multiple covariates

Cumulative diffs. between folding and normal mailings

Mark Tygert Conditioning & controlling via cumulative differences 50 / 70

イロト イヨト イヨト

Deviation of a subpopulation from the full pop. Deviation between two disjoint subpopulations Controlling for multiple covariates

Description of the data (skip these details)

- (1) In 1994, a national veterans organization mailed folding cards to $n_0 = 1,236$ prospective donors and mailed normal cards to $n_1 = 15,866$.
- The responses (indicating whether the recipients responded to the solicitation or not), together with covariates such as age and average household income in the associated Census block, formed the core of the data for the 1998 KDD Cup.
- 3 The combined sample size derived from the diagram on slide 41 is n = 2,279 (rather than the n = 20 on slide 41).

イロト イポト イヨト イヨト

Weighted sampling Avoiding randomization Interactive plots

Extensions

Extensions

Mark Tygert Conditioning & controlling via cumulative differences 52 / 70

《日》《圖》《臣》《臣》

990

3

Weighted sampling Avoiding randomization Interactive plots

Weighted sampling

Weighted sampling

Mark Tygert Conditioning & controlling via cumulative differences 53 / 70

イロト イポト イヨト イヨト

DQC2

1

Weighted sampling Avoiding randomization Interactive plots

Graphs for weighted samples (skip if low on time)

All methods generalize to weighted sampling; in weighted samples each pair of score S_k and response R_k comes with a positive real number W_k . The ordinates (vertical coordinates) become

$$C_{j} = \frac{\sum_{k=1}^{j} (R_{i_{k}} - \tilde{R}_{i_{k}}) W_{i_{k}}}{\sum_{k=1}^{n} W_{i_{k}}}$$
(14)

(with $\tilde{R}_{i_k} := S_{i_k}$ and $i_k := k$ when assessing calibration) and the abscissae (horizontal coordinates) become

$$A_{j} = \frac{\sum_{k=1}^{j} W_{i_{k}}}{\sum_{k=1}^{n} W_{i_{k}}}$$
(15)

for j = 1, 2, ..., n.

Mark Tygert Conditioning & controlling via cumulative differences 54 / 70

イロト イポト イヨト イヨト 三日

Weighted sampling Avoiding randomization Interactive plots

P-values for weighted samples (skip if low on time)

- **(1)** The interpretation of slope in the graph of C_k versus A_k is the same as in the case of unweighted sampling.
- ² P-values for the maximum absolute value of C_1, C_2, \ldots, C_n and for the range of C_0, C_1, \ldots, C_n arise similarly, courtesy of convergence in distribution to the maximum absolute value or range of the standard Brownian motion over [0, 1], as before.
- **3** For example, formula (6) on slide 16 above, giving the standard deviation under the null hypothesis of perfect calibration (so that $\tilde{R}_{i_k} := S_{i_k}$ and $i_k := k$), becomes

$$\sigma = \frac{\sqrt{\sum_{k=1}^{n} S_k (1 - S_k) (W_k)^2}}{\sum_{k=1}^{n} W_k}.$$
 (16)

Mark Tygert Conditioning & controlling via cumulative differences 55 / 70

イロト イヨト イヨト

Weighted sampling Avoiding randomization Interactive plots

Avoiding randomization

Avoiding randomization

Mark Tygert Conditioning & controlling via cumulative differences 56 / 70

・ロト ・ 同ト ・ ヨト ・ ヨト

nac

Weighted sampling Avoiding randomization Interactive plots

Slight random perturbations (skip if low on time)

- (1) If any of the inequalities $S_1 \leq S_2 \leq \cdots \leq S_n$ were not strict, one possibility would be to perturb the scores at random such that the scores become all distinct from each other, so that $S_1 < S_2 < \cdots < S_n$. Slide 7 suggested this approach.
- 2 To describe another possibility, we can keep the notation $S_1 < S_2 < \cdots < S_n$ while denoting the responses and weights for n_k repetitions of score S_k by $R_k^{(1)}$, $R_k^{(2)}$, \ldots , $R_k^{(n_k)}$ and $W_k^{(1)}$, $W_k^{(2)}$, \ldots , $W_k^{(n_k)}$ (so then the total number of observations is $\sum_{k=1}^n n_k$).

Mark Tygert Conditioning & controlling via cumulative differences 57 / 70

イロト (雪) (ヨ) (ヨ) - ヨ

Weighted sampling Avoiding randomization Interactive plots

A new data set (skip if low on time)

We construct a new data set, with the response being the weighted average

$$R_{k} = \frac{\sum_{j=1}^{n_{k}} R_{k}^{(j)} W_{k}^{(j)}}{\sum_{j=1}^{n_{k}} W_{k}^{(j)}}$$
(17)

and the weight being the sum

$$W_k = \sum_{j=1}^{n_k} W_k^{(j)}$$
(18)

for k = 1, 2, ..., n.

Mark Tygert Conditioning & controlling via cumulative differences 58 / 70

イロト イヨト イヨト

Weighted sampling Avoiding randomization Interactive plots

Cumulative differences stay the same (skip if low on time)

I The cumulative differences for the new data set are

$$B_{\ell} = \frac{\sum_{k=1}^{\ell} (R_k - r(S_k)) W_k}{\sum_{k=1}^{n} W_k}$$
(19)

for $\ell = 1, 2, ..., n$, with $r(S_k) = S_k$ for gauging calibration.

② The cumulative differences for the original data set with the scores perturbed infinitesimally at random are

$$C_{\ell} = \frac{\sum_{k=1}^{\ell} \sum_{j=1}^{n_k} \left(R_k^{(j)} - r(S_k) \right) W_k^{(j)}}{\sum_{k=1}^{n} \sum_{j=1}^{n_k} W_k^{(j)}}$$
(20)

for $\ell = 1, 2, ..., n$, where the responses $R_k^{(1)}, R_k^{(2)}, ..., R_k^{(n_k)}$ & weights paired with them are randomly permuted for each k. 3 So $B_\ell = C_\ell$ for all $\ell = 1, 2, ..., n$, due to the previous slide.

Mark Tygert Conditioning & controlling via cumulative differences 59 / 70

Weighted sampling Avoiding randomization Interactive plots

Ties in scores can be treated as weighted samples

The aggregated abscissae (horizontal coordinates) are

$$A_{\ell} = \frac{\sum_{k=1}^{\ell} W_{k}}{\sum_{k=1}^{n} W_{k}} = \frac{\sum_{k=1}^{\ell} \sum_{j=1}^{n_{k}} W_{k}^{(j)}}{\sum_{k=1}^{n} \sum_{j=1}^{n_{k}} W_{k}^{(j)}}$$
(21)

for $\ell = 1, 2, ..., n$.

Thus, the horizontal coordinates in the cumulative graphs for the original and new data sets are the same, as are the vertical coords. (as shown on the previous slide), at least when the graph consists of straight line segments connecting $(A_{k-1}, B_{k-1}) = (A_{k-1}, C_{k-1})$ to $(A_k, B_k) = (A_k, C_k)$ for k = 1, 2, ..., n, with $A_0 = B_0 = C_0 = 0$. The original data set's cumulative graph which interpolates linearly from every perturbed score to the next greatest perturbed score will be almost the same, aside from having slightly larger random excursions between scores that were the same prior to perturbation.

Mark Tygert Conditioning & controlling via cumulative differences 60 / 70

・ロト ・ 同ト ・ ヨト ・ ヨト

JOC P

Weighted sampling Avoiding randomization Interactive plots

Advantages and disadvantages (skip if low on time)

	randomized	non-randomized
pros	displays every single member of the original data set	horizontal axis has no randomization
cons	horizontal axis is randomized to some extent (the ordering for repeated scores must be randomized, obviously)	displays only averaged responses for repeated (degenerate) scores

Mark Tygert Conditioning & controlling via cumulative differences 61 / 70

イロト イヨト イヨト

DQC2

1

Weighted sampling Avoiding randomization Interactive plots

Interactive plots

Interactive plots

Mark Tygert Conditioning & controlling via cumulative differences 62 / 70

・ロト ・ 同ト ・ ヨト ・ ヨト

DQC2

Э

Weighted sampling Avoiding randomization Interactive plots

Interactive traversal of the Hilbert curve

Census Bureau's weighted sampling of counties in California: the county is the subpopulation of the full state pop.; responses are given by "Variable". (Click to play the movie.)

Mark Tygert Conditioning & controlling via cumulative differences 63 / 70

イロト イヨト イヨト

Weighted sampling Avoiding randomization Interactive plots

Concluding thoughts

Concluding thoughts

Mark Tygert Conditioning & controlling via cumulative differences 64 / 70

イロト イポト イヨト イヨト

DQC2

1

Weighted sampling Avoiding randomization Interactive plots

Ideal limit of the empirical Kuiper metric

Given 3 real-valued random variables, X, Y, Z, the Kuiper metric is the following summary stat of differences between the regression $\mathbb{E}[Y|X]$ of Y on X and the regression $\mathbb{E}[Z|X]$ of Z on X; this stat compares Y and Z while controlling for X (say X is age or income):

$$\max_{-\infty \le a \le b \le \infty} \left| \mathbb{E}_{a \le X \le b} \left[\mathbb{E}[Y|X] - \mathbb{E}[Z|X] \right] \right|,$$
(22)

where
$$\mathbb{E}_{a \leq X \leq b} \left[\mathbb{E}[Y|X] - \mathbb{E}[Z|X] \right]$$

= $\mathbb{E} \left[\left(\mathbb{E}[Y|X] - \mathbb{E}[Z|X] \right) \cdot \mathbb{1} \{ a \leq X \leq b \} \right]$, (23)

$$\mathbb{1}\{a \le X \le b\} = 1 \text{ when } a \le X \le b, \text{ and}$$
(24)

$$\mathbb{1}\{a \le X \le b\} = 0 \text{ when } X < a \text{ or } X > b. \tag{25}$$

Calibration is the special case with Z = X (making $\mathbb{E}[Z|X] = X$).

Mark Tygert Conditioning & controlling via cumulative differences 65 / 70

イロト イ伊ト イヨト イヨト

Weighted sampling Avoiding randomization Interactive plots

Matched-pair analysis

The simplest, most common paired samples consist of observations from two subpopulations, with each observed response from one subpopulation corresponding to an observed response from the other subpopulation at the same value(s) of the covariate(s). Such a pair of observed responses (one from each subpopulation) at the same value(s) of the covariate(s) is known as a "matched pair," with the matching based on the value(s) of the covariate(s). Generalization of all of this tutorial's cumulative methodologies to the analysis of paired samples should be obvious; full details are in the paper available at https://arxiv.org/abs/2305.11323

Mark Tygert Conditioning & controlling via cumulative differences 66 / 70

イロト イ伊ト イヨト イヨト

Weighted sampling Avoiding randomization Interactive plots

Open-access articles and open-source software

Open-access references are available at the top of http://tygert.com/research.html

Open-source software is available at the bottom of http://tygert.com/software.html

Mark Tygert Conditioning & controlling via cumulative differences 67 / 70

イロト イ伊ト イヨト イヨト

Weighted sampling Avoiding randomization Interactive plots

Randomized controlled trials and A/B tests

The cumulative approach is relevant for analyzing the treatment effect in any randomized controlled trial. Conditioning on covariates traditionally can lead to issues with Simpson's Paradox, due to the need to choose the sizes of the bins for binning ("binning" is also known as "bucketing," "segmenting," or "stratifying"). Cumulative stats avoid Simpson's Paradox entirely.

A/B tests are a form of randomized controlled trial that is especially popular in the data-driven industries.

イロト イポト イヨト イヨト 三日

Weighted sampling Avoiding randomization Interactive plots

Observational studies

Observational studies can also benefit from cumulative statistics. "Adjusting" for covariates is common in such studies. Adjustment often involves conditioning on (controlling for) covariates. Conditioning on confounding covariates facilitates causal or counterfactual analysis.

イロト イ伊ト イヨト イヨト

Weighted sampling Avoiding randomization Interactive plots

Conclusion

- Kolmogorov, Kuiper, Smirnov, and Wiener had good reason to introduce cumulative statistics, and Hilbert and Peano had good reason to introduce their space-filling curves.
- Classical non-cumulative statistics have an explicit trade-off between statistical confidence and resolving power (resolution of variations as a function of score). The classical metrics and reliability diagrams are typically misleading (or manipulated) due to the strong dependence on the arbitrary choice of bins.
- 3 Meta's Fairness Flow promotes the cumulative metrics as the preferred method for detecting differences while controlling for variates, thanks to the hard work of the Responsible AI team.
- 4 There is plenty more to come in this space stay tuned!

イロト イボト イヨト イヨト