

Generative Trees: Adversarial and Copycat

Richard Nock

Mathieu Guillaume-Bert

Why this work

- Tabular data = prevalent+ for value potential
Chui *et al.* 2018
- Generative SOTA for tabular data = DL based **unlike** SOTA for classification (tree based)
- Dissatisfaction, new directions needed for generative approaches
Camino *et al.* 2020

Loss functions

Discriminative (for class probability estimation)
 Objective: learn posterior η to estimate $P[Y=1|X]$
 Ideal: $\eta^* = \pi \cdot (dP/dM)$
 Bayes posterior $M = \pi \cdot P + (1-\pi) \cdot N$
 Loss: $\ell: \{-1, 1\} \times [0, 1] \rightarrow \bar{\mathbb{R}}$
 true class posterior estimate

$$\underline{L}(p) = \inf_u \mathbb{E}_{Y \sim B(p)} [\ell(Y, u)] \text{ Bayes risk}$$

Calibrated posterior:
 $\tilde{\eta} = \pi \cdot \frac{dP_{\tilde{\eta}}}{dM_{\tilde{\eta}}}$
 agreement on the level sets of the posterior
 ↪ Best calibrated posterior: Bayes' $\tilde{\eta}^*$ worst: prior π
 ↪ Important: any **decision tree** with local posterior predictions at the leaves is calibrated

Statistical information of calibrated $\tilde{\eta}$
 $\Delta\underline{L}(\tilde{\eta}, M_{\tilde{\eta}}) = L(\pi) - \mathbb{E}_{X \sim M_{\tilde{\eta}}} [\underline{L}(\tilde{\eta}(X))]$
 CART, C4.5, etc. (splitting criterion)

The better $\tilde{\eta}$, the higher $\Delta\underline{L}(\tilde{\eta}, M_{\tilde{\eta}})$

Starting point

- Map key frameworks for tabular data classification onto generative world
 - ↳ **Loss functions:** properness Savage, 1971
 - ↳ **Models:** tree-based Breiman *et al.*, 1984
 - ↳ **Algorithms:** boosting Kearns & Mansour, 1996

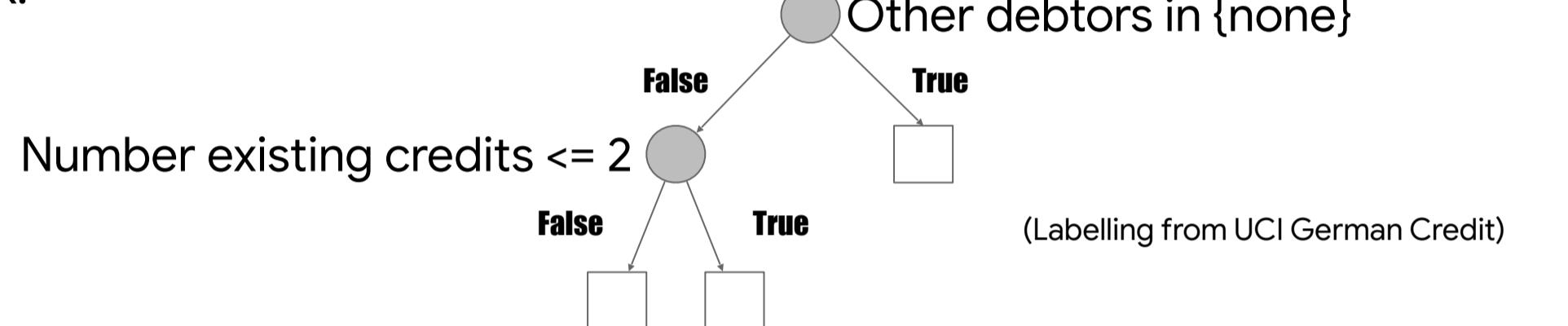
What we get (summary)

- GAN-game losses from discriminator, tight if discriminator calibrated, chi square “universal” to train generator
- New **tree-based generative models:** XAI+, density in $O(\text{depth})$, likelihood w/ missing features in $O(\text{size})$
- Boosting algorithm for adversarial training + geometric conv. of chi square under weak generative assumption**
- New copycat training** (unknown for DL), “boosting for free” convergence if boosted discriminator (C4.5, etc.)

Models

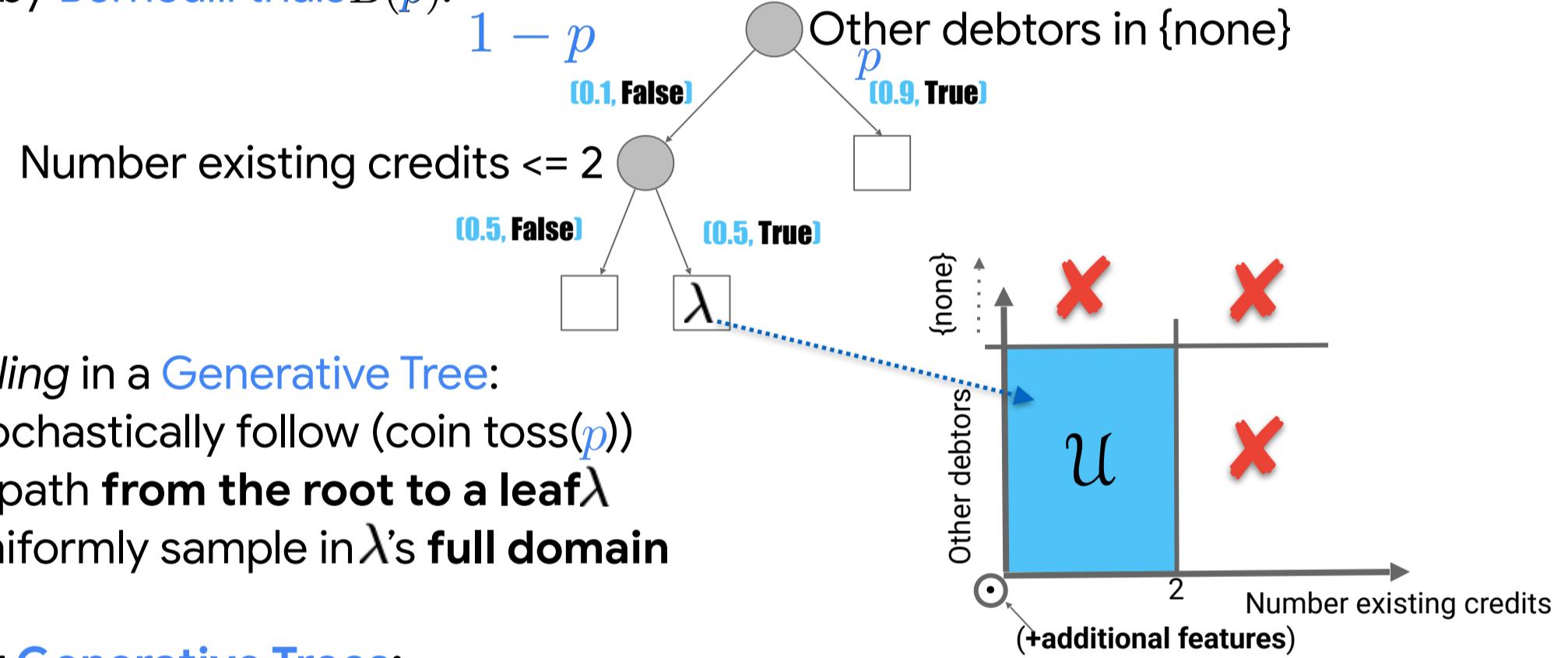
Models

Tree: a tree is a binary directed tree whose internal nodes are labeled with a test on an observation variable and outgoing arcs are labeled with truth values. Leaves are blank.



Decision Tree: A *decision tree h* is a tree in which leaves are labeled by values in $[0,1]$.

Generative Tree: A *generative tree (GT) G* is a tree in which outgoing arcs are labeled by Bernoulli trials $B(p)$.



Sampling in a **Generative Tree**:

- (1) Stochastically follow (coin toss \mathcal{U}) a path from the root to a leaf,
- (2) Uniformly sample in λ 's full domain

Pros for Generative Trees:

- Trainable from data w/ missing values (see algorithms)
- XAI+ as “easy” to interpret as a **decision tree** (see experiments)
- Density computable in $O(\text{depth})$ at any point,
- Likelihood of missing values given partial observation in $O(\text{size})$, $P[A|B]$ equivalently easy
- Many metrics involving a GT computable exactly (no sampling required!): chi square, etc.

Cons for Generative Trees:

- Axis // partition of the support (restrictive)
- Support / domain closed (workarounds exist, e.g. Box-Muller transform for Gaussians)

Algorithms

Adversarial training

Algorithm TD-GEN(G, h)

Input: current generator G , current discriminator h ;
Output: G with a new split;

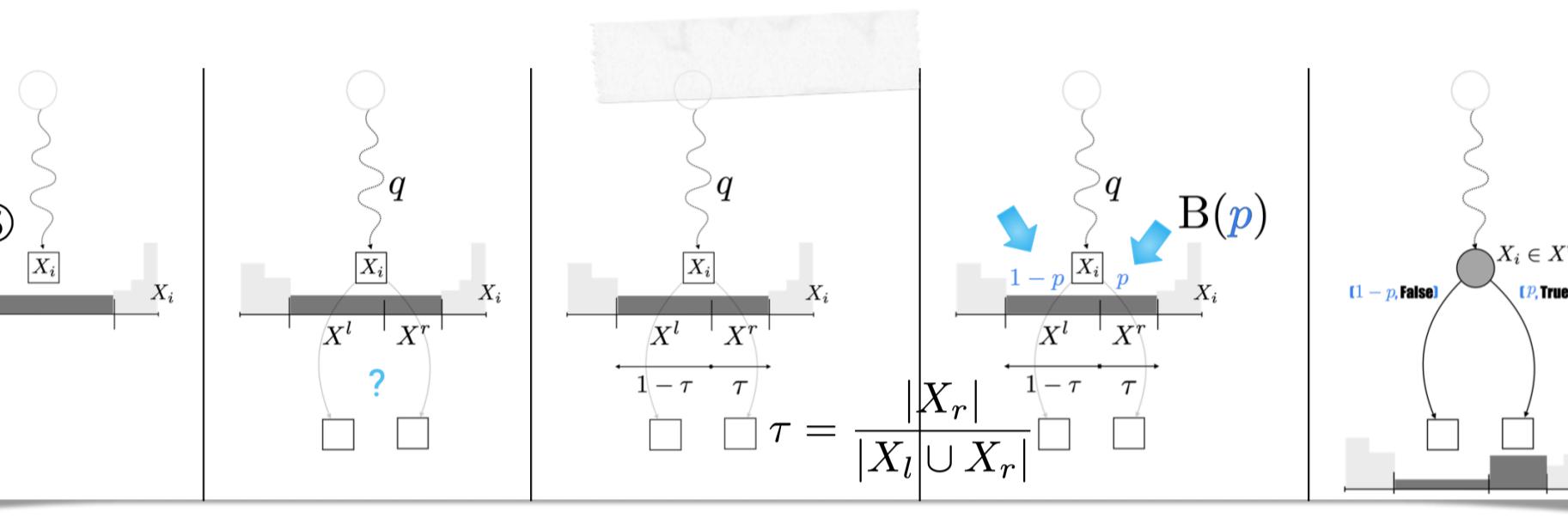
Step 1 : pick $\mathbb{S} \in \Lambda(G)$, $i \in [d]$; // leaf and variable for the current split

Step 2 : choose (X^l, X^r) and compute τ ; // split choice

Step 3 : compute p as parameters depend on G, h and τ

$$p \leftarrow \text{CLAMP} \left(\frac{\mu_{LL} - \mu_{LR}}{\mu_{LL} + \mu_{RR} - 2\mu_{LR}} \right); \quad (\text{see paper})$$

Step 4 : replace \mathbb{S} by a split as designed in Steps 1,2 w/ Bernoulli probability p as in (16);



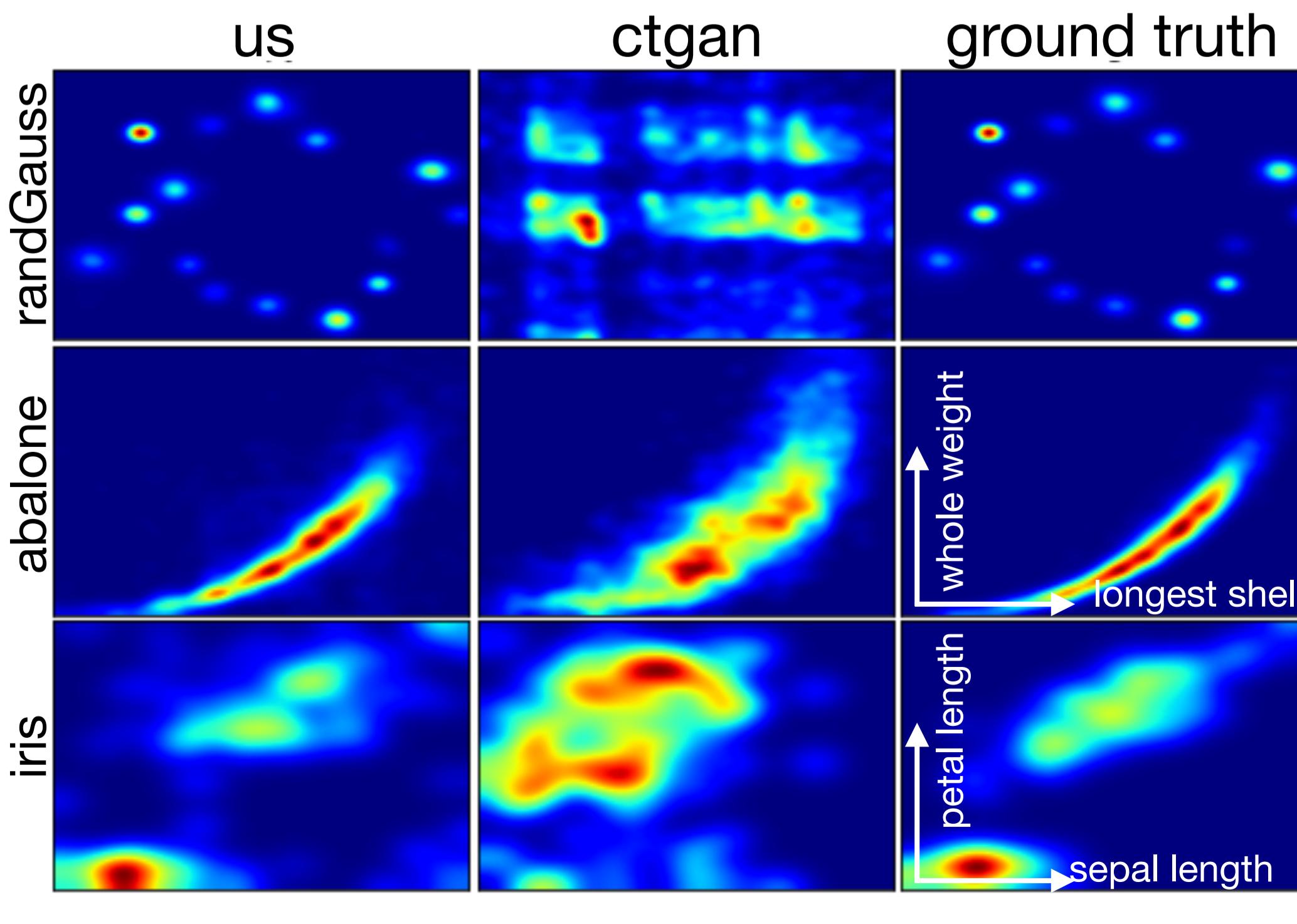
Convergence if p in $(0,1)$

↳ If $\chi^2(N_{\tilde{\eta}}||P_{\tilde{\eta}}) \geq \delta$ and $|\tau - p| \geq \varepsilon$

$$\text{new: } \chi^2(N_{\tilde{\eta}}' || P_{\tilde{\eta}}) \leq \frac{1}{1 + \delta\varepsilon^2} \cdot \chi^2(N_{\tilde{\eta}} || P_{\tilde{\eta}})$$

old: $\chi^2(N_{\tilde{\eta}}' || P_{\tilde{\eta}}) \leq \frac{1}{1 + \delta q^2 (\tau + (1 - 2\tau)p)^2} \cdot \chi^2(N_{\tilde{\eta}} || P_{\tilde{\eta}})$

Experiments



Missing data imputation

Setup:

- Remove $q\%$ values in dataset (*Missing Completely At Random*), 5-fold CV
- Impute missing values w/ algorithm, compare with domain using OT (W_2^2) metric

Some results (more in paper)

	NORM	CART	RF	Us vs Mice	CART RF	Us vs Mice	CART RF	Us vs Mice	CART	RF
($q=1\%$)% circGauss	$U(0.003)$	U	$U(0.07)$	5% $U(0.04)$	10% $U(0.05)$	10% $U(0.05)$	10% $U(0.05)$	10% $U(0.05)$	10% $M(0.02)$	10% $M(0.02)$
10% circGauss	$U(0.03)$	U(0.02)	$U(0.007)$	$U(0.04)$	$U(0.01)$	$U(0.05)$	$U(0.05)$	$U(0.05)$	$U(0.04)$	$U(0.04)$
20% circGauss	$U(0.0005)$	$U(0.001)$		$U(0.01)$	$U(0.001)$	$U(0.05)$	$U(0.05)$	$U(0.05)$	$U(0.01)$	$U(0.01)$
50% circGauss	$U(0.04)$	U		$U(0.01)$	$U(0.001)$	$U(0.05)$	$U(0.05)$	$U(0.05)$	$U(0.01)$	$U(0.01)$

Observations:

- CTGs can beat mice if small dim (unexpected) & always with 1 tree (vs many for mice)
 - training time: CTGs << mice (especially w/ RFs)
 - CTGs provide generator (mice: no byproduct)
- us = copycat training a GT G for max 10K splits
 SOTA: mice (van Buuren, 2018)

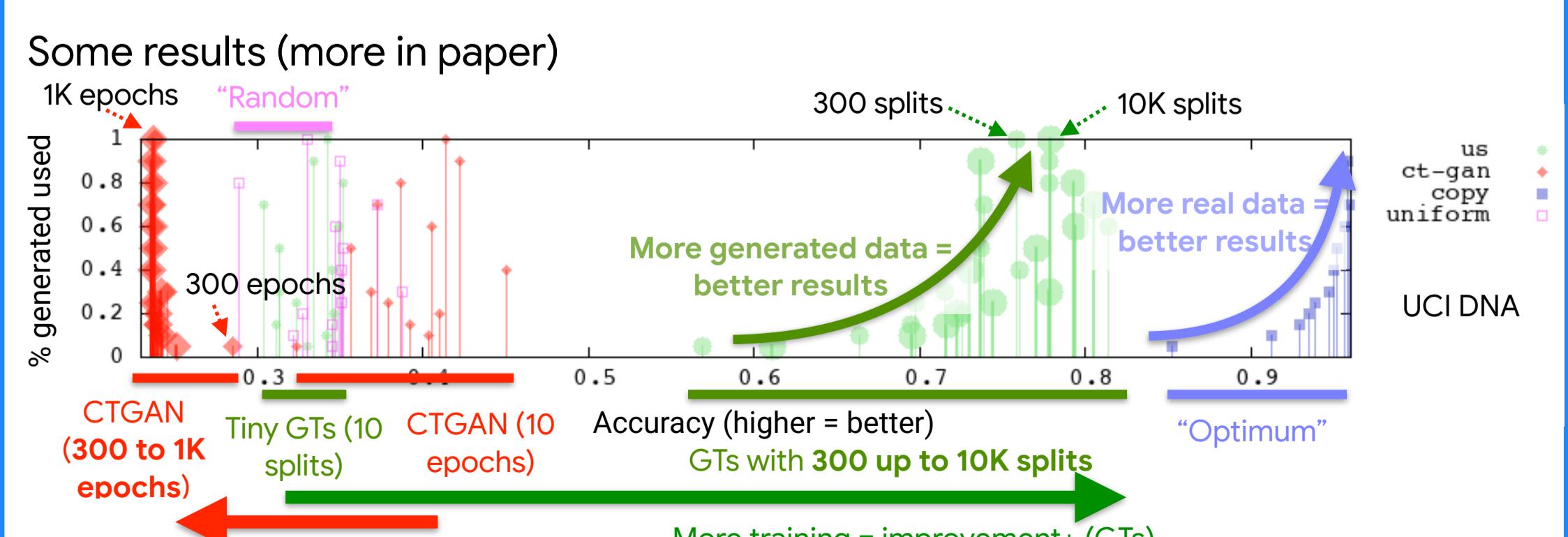
- After initial guess, round-robin updates one column's missing values given the others using regression/classification **method**, iterates 5 times
- We used **method** in {NORM, CART, Random Forests = RF with 100 trees each}

Generated data augmentation

Setup:

- Take a **supervised** domain (e.g. UCI iris), 5-fold CV
- Train generator, generate additional q in {5, 10, ...100}% (of training) examples, train supervised model (GBDT or RFs) from all data, get accuracy/RMSE on test
- 2 baselines: **random** = +uniform data (“worst”), **copy** = +real data (“optimum”)

Some results (more in paper)



... and more

Additional experiments in paper:

- Training from just generated data (vs CTGAN)
- Distinguishing real from fake data (vs CTGAN)

XAI by the example: crop of a GT learned on Stanford Open Policing (Hartford)

Code snippet showing generated data for Stanford Open Policing (Hartford):

```

  ...
  | search vehicle (CONTINUOUS) in [-76.8665, -72.7167]; [-1, -1] | ]--#2
  |-0.0366, [ search vehicle (NOMINAL) in (TRUCK); [-1, -1] ] ]--#100
  |-0.1212, [ lat (CONTINUOUS) in (40.7067, 41.3291); [-1, -1] ] ]--#3010 (sampling)
  |-0.8788, [ lat (CONTINUOUS) in (41.7329, 42.3426); [-1, -1] ] ]--#3011 (sampling)
  |-0.8276, [ lat (CONTINUOUS) in (41.7329, 41.8060); [-1, -1] ] ]--#3338 (sampling)
  |-0.8333, [ warning issued (NOMINAL) = (FALSE); [-1, -1] ] ]--#3756 (sampling)
  |-0.9500, [ raw_subject_race_code (NOMINAL) = (W, B); [-1, -1] ] ]--#1848 (sampling)
  |-0.7368, [ reason_for_stop (NOMINAL) in (STOP SIGN, DEFECTIVE LIGHTS, CELLPHONE, SUSPENDED LICENSE, OTHER); [-1, -1] ] ]--#788 (sampling)
  |-0.5000, [ district (NOMINAL) in (BRYNNSKROG, NORTHEMADOWNS); [-1, -1] ] ]--#3011 (sampling)
  |-0.4286, [ subject_age (INTEGER) in (14, 15, ..., 29); [-1, -1] ] ]--#9119 (sampling)
  |-0.5714, [ subject_age (INTEGER) in (30, 31, ..., 94); [-1, -1] ] ]--#1119 (sampling)
  |-0.0553, [ raw_subject_race_code (NOMINAL) in (ASIAN, HISPANIC, PAKISTANI, FROGGOLLOW, BEHINDTHEROADS, SOUTHGREN, OTHER); [-1, -1] ] ]--#14185 (sampling)
  |-0.1667, [ raw_subj_race_code (NOMINAL) in (TRUE); [-1, -1] ] ]--#3757 (sampling)
  |-0.9634, [ search vehicle (NOMINAL) in (FALSE); [-1, -1] ] ]--#1724 (sampling)
  |-0.2625, [ raw_search_authorization_code (NOMINAL) in (0, 1); [-1, -1] ] ]--#3757 (sampling)
  |-0.0870, [ lat (CONTINUOUS) in (40.7067, 41.6730); [-1, -1] ] ]--#3757 (sampling)
  |-0.9130, [ lat (CONTINUOUS) in (41.6730, 42.3426); [-1, -1] ] ]--#3757 (sampling)
  ...
  
```

↓ disparities for young (>30 yrs) vs not-young (>=30) on “car/driver search” w/o warning issued in specific area (Each young age > 3x more likely than each not-young)

Conclusion and future work:

- Contribution on losses applies beyond our framework to any calibrated disc.
- Interesting avenues on alleviating the Cons for **Generative Trees**
- Experiments display edge vs CTGAN + potential use for missing data imputation