

Generative Trees: Adversarial and Copycat

Richard Nock

Mathieu Guillame-Bert



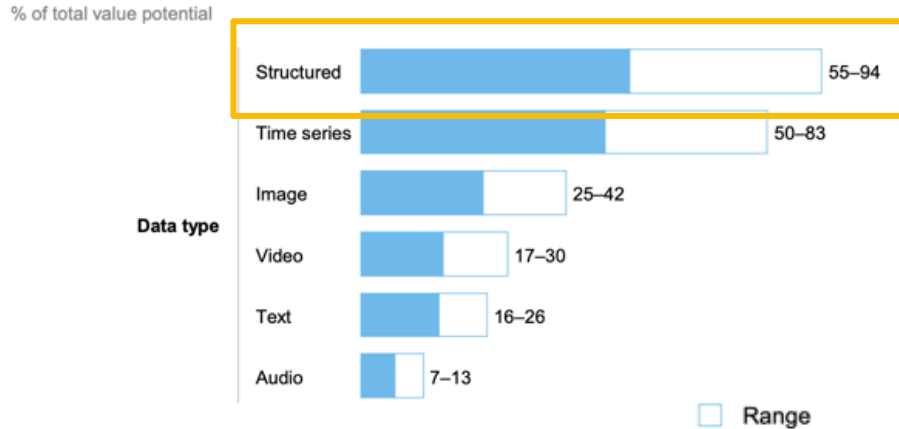
Google Research

{richardnock,gbm}@google.com



Why this work ?

Tabular data: important but scarce generative



SOURCE: McKinsey Global Institute analysis

- Modern generative techniques = Neural Networks (NN) / Deep Learning (DL) based
- On supervised learning side, the best techniques are (still) **not** DL-based but tree-based ; competing requires sophisticated+ DL techs
- “Lack of novelty” in state of the art (SOTA) modern generative approaches for tabular data
- Unconvincing results for DL + tabular pipelines

Camino *et al.*, ICBINB@NeurIPS'20

Google Research

Tabular data, *Supervised*

Losses: **proper**

Savage, JASA'71

Models: **tree-based**

Breiman *et al.* '84

Algorithms: **boosting**

Kearns & Mansour, STOC'96

Tabular data, *Supervised*

Losses: **proper**

Savage, JASA'71

Models: **tree-based**

Breiman *et al.* '84

Algorithms: **boosting**

Kearns & Mansour, STOC'96

This paper, *generative*, tabular data

Background: GAN game

Losses: designed from ***discriminator*** & in the **proper** framework

Models: ***tree-based***

Algorithms: **boosting**

↳ adversarial

↳ ***copycat***

Loss functions

GAN framework in a

Measure-based loss, crafted from *generator*

$$\mathbb{I}_f(\underset{\substack{\text{(real)} \\ \downarrow}}{\mathbf{P}}, \underset{\substack{\uparrow \\ \text{(fake)}}}{\mathbf{N}}) \doteq \int f\left(\frac{d\mathbf{P}}{d\mathbf{N}}\right) d\mathbf{N} \quad f\text{-divergence}$$

“Information of Binary Task”

GAN framework in a

Measure-based loss, crafted from *generator*

$$\mathbb{I}_f(\underset{\substack{\text{(real)} \\ \downarrow}}{\mathbf{P}}, \underset{\substack{\text{(fake)} \\ \uparrow}}{\mathbf{N}}) \doteq \int f\left(\frac{d\mathbf{P}}{d\mathbf{N}}\right) d\mathbf{N} \quad f\text{-divergence}$$

“Information of Binary Task”

↳ variational formulation

$$\mathbb{I}_f(\mathbf{P}, \mathbf{N}) \geq \sup_{\tilde{h}} \left\{ \overset{\text{H}}{\mathbb{E}_{\mathbf{P}}[\tilde{h}(\mathbf{X})]} + \overset{\text{G}}{\mathbb{E}_{\mathbf{N}}[f^* \circ \tilde{h}(\mathbf{X})]} \right\}$$

Not an equality in general

GAN framework in a

Measure-based loss, crafted from *generator*

$$\mathbb{I}_f(\underset{\substack{\text{(real)} \\ \downarrow}}{\mathbf{P}}, \underset{\substack{\text{(fake)} \\ \uparrow}}{\mathbf{N}}) \doteq \int f\left(\frac{d\mathbf{P}}{d\mathbf{N}}\right) d\mathbf{N} \quad f\text{-divergence}$$

“Information of Binary Task”

↳ variational formulation

$$\mathbb{I}_f(\mathbf{P}, \mathbf{N}) \geq \sup_{\tilde{h}} \left\{ \overset{\text{H}}{\mathbb{E}_{\mathbf{P}}[\tilde{h}(\mathbf{X})]} + \overset{\text{G}}{\mathbb{E}_{\mathbf{N}}[f^* \circ \tilde{h}(\mathbf{X})]} \right\}$$

Not an equality in general

↳ *discriminator* hidden in \tilde{h} , seeks to increase the IBT by discriminating **real** vs **fake** (w/ **H** + **G**)

GAN framework in a

Measure-based loss, crafted from *generator*

$$\mathbb{I}_f(\underset{\substack{\text{(real)} \\ \downarrow}}{\mathbf{P}}, \underset{\substack{\uparrow \\ \text{(fake)}}}{\mathbf{N}}) \doteq \int f\left(\frac{d\mathbf{P}}{d\mathbf{N}}\right) d\mathbf{N} \quad f\text{-divergence}$$

“Information of Binary Task”

↳ variational formulation

$$\mathbb{I}_f(\mathbf{P}, \mathbf{N}) \geq \sup_{\tilde{h}} \left\{ \overset{\mathbf{H}}{\mathbb{E}_{\mathbf{P}}[\tilde{h}(\mathbf{X})]} + \overset{\mathbf{G}}{\mathbb{E}_{\mathbf{N}}[f^* \circ \tilde{h}(\mathbf{X})]} \right\}$$

Not an equality in general

↳ *discriminator* hidden in \tilde{h} , seeks to increase the IBT by discriminating **real** vs **fake** (w/ **H** + **G**)

↳ *generator* = \mathbf{N} , seeks to decrease the IBT by generating **fake** data that looks like **real** (w/ **G**)

GAN framework in a

Measure-based loss, crafted from *generator*

$$\mathbb{I}_f(\underset{\text{(fake)}}{\mathbb{P}}, \underset{\text{(real)}}{\mathbb{N}}) \doteq \int f\left(\frac{d\mathbb{P}}{d\mathbb{N}}\right) d\mathbb{N} \quad f\text{-divergence}$$

“Information of Binary Task”

↳ variational formulation

$$\mathbb{I}_f(\mathbb{P}, \mathbb{N}) \geq \sup_{\tilde{h}} \left\{ \mathbb{E}_{\mathbb{P}}[\tilde{h}(X)] - \mathbb{E}_{\mathbb{N}}[f^* \circ \tilde{h}(X)] \right\}$$

↑ \tilde{h}

Not an equality in general

↳ *discriminator* hidden in \tilde{h} , seeks to increase the IBT by discriminating *real* vs *fake* (w/ **H** + **G**)

↳ *generator* = \mathbb{N} , seeks to decrease the IBT by generating *fake* data that looks like *real* (w/ **G**)

Our framework

Properness: $\mathbb{P}[Y = 1]$ $\mathbb{P}[X|Y = 1]$ $\mathbb{P}[X|Y = -1]$

↳ *Binary task* $\mathbf{B} \doteq (\pi, \mathbb{P}, \mathbb{N})$

↳ *Mixture* $\mathbf{M} \doteq \pi \cdot \mathbb{P} + (1 - \pi) \cdot \mathbb{N}$

GAN framework in a

Measure-based loss, crafted from *generator*

$$\mathbb{I}_f(\underset{\text{(fake)}}{\mathbb{P}}, \underset{\text{(real)}}{\mathbb{N}}) \doteq \int f\left(\frac{d\mathbb{P}}{d\mathbb{N}}\right) d\mathbb{N} \quad \textit{f-divergence}$$

“Information of Binary Task”

↳ variational formulation

$$\mathbb{I}_f(\mathbb{P}, \mathbb{N}) \geq \sup_{\tilde{h}} \left\{ \mathbb{E}_{\mathbb{P}}[\tilde{h}(X)] - \mathbb{E}_{\mathbb{N}}[f^* \circ \tilde{h}(X)] \right\}$$

↑ \tilde{h}

Not an equality in general

↳ *discriminator* hidden in \tilde{h} , seeks to increase the IBT by discriminating *real* vs *fake* (w/ **H** + **G**)

↳ *generator* = \mathbb{N} , seeks to decrease the IBT by generating *fake* data that looks like *real* (w/ **G**)

Our framework

Properness: $\mathbb{P}[Y=1]$ $\mathbb{P}[X|Y=1]$ $\mathbb{P}[X|Y=-1]$

↳ *Binary task* $\mathbf{B} \doteq (\pi, \mathbb{P}, \mathbb{N})$

↳ *Mixture* $\mathbf{M} \doteq \pi \cdot \mathbb{P} + (1 - \pi) \cdot \mathbb{N}$

↳ *discriminator* learns a *posterior* $\eta : \mathcal{X} \rightarrow [0, 1]$
↑ $\hat{\mathbb{P}}[Y=1|X]$

↳ “*ideal*” posterior computes $\mathbb{P}[Y=1|X]$

$$\eta^* = \pi \cdot \frac{d\mathbb{P}}{d\mathbf{M}} \quad \textbf{Bayes posterior}$$

Google Research

GAN framework in a

Measure-based loss, crafted from *generator*

$$\mathbb{I}_f(\underset{\text{(fake)}}{\mathbb{P}}, \underset{\text{(real)}}{\mathbb{N}}) \doteq \int f\left(\frac{d\mathbb{P}}{d\mathbb{N}}\right) d\mathbb{N} \quad \text{f-divergence}$$

“Information of Binary Task”

↳ variational formulation

$$\mathbb{I}_f(\mathbb{P}, \mathbb{N}) \geq \sup_{\tilde{h}} \left\{ \mathbb{E}_{\mathbb{P}}[\tilde{h}(X)] - \mathbb{E}_{\mathbb{N}}[f^* \circ \tilde{h}(X)] \right\}$$

Not an equality in general

↳ *discriminator* hidden in \tilde{h} , seeks to increase the IBT by discriminating *real* vs *fake* (w/ **H** + **G**)

↳ *generator* = \mathbb{N} , seeks to decrease the IBT by generating *fake* data that looks like *real* (w/ **G**)

Our framework

Properness: $\mathbb{P}[Y=1]$ $\mathbb{P}[X|Y=1]$ $\mathbb{P}[X|Y=-1]$

↳ *Binary task* $\mathbf{B} \doteq (\pi, \mathbb{P}, \mathbb{N})$

↳ *Mixture* $\mathbf{M} \doteq \pi \cdot \mathbb{P} + (1 - \pi) \cdot \mathbb{N}$

↳ discriminator learns a *posterior* $\eta : \mathcal{X} \rightarrow [0, 1]$
 $\uparrow \hat{\mathbb{P}}[Y=1|X]$

↳ “*ideal*” posterior computes $\mathbb{P}[Y=1|X]$

$$\eta^* = \pi \cdot \frac{d\mathbb{P}}{d\mathbf{M}} \quad \text{Bayes posterior}$$

↳ a **loss** can be decomposed in two *partial losses*

$$\ell(y, u) \doteq \mathbb{I}[y=1] \cdot \ell_1(u) + \mathbb{I}[y=-1] \cdot \ell_{-1}(u)$$

↑ estimated posterior in [0,1]
 ↑ true label / class in {-1,1}

GAN framework in a

Measure-based loss, crafted from *generator*

$$\mathbb{I}_f(\underset{\text{(fake)}}{\mathbb{P}}, \underset{\text{(real)}}{\mathbb{N}}) \doteq \int f\left(\frac{d\mathbb{P}}{d\mathbb{N}}\right) d\mathbb{N} \quad \text{\textit{f-divergence}}$$

“Information of Binary Task”

↳ variational formulation

$$\mathbb{I}_f(\mathbb{P}, \mathbb{N}) \geq \sup_{\tilde{h}} \left\{ \mathbb{E}_{\mathbb{P}}[\tilde{h}(X)] - \mathbb{E}_{\mathbb{N}}[f^* \circ \tilde{h}(X)] \right\}$$

Not an equality in general

↳ *discriminator* hidden in \tilde{h} , seeks to increase the IBT by discriminating *real* vs *fake* (w/ **H** + **G**)

↳ *generator* = \mathbb{N} , seeks to decrease the IBT by generating *fake* data that looks like *real* (w/ **G**)

Our framework

Properness: $\mathbb{P}[Y=1]$ $\mathbb{P}[X|Y=1]$ $\mathbb{P}[X|Y=-1]$

↳ *Binary task* $\mathbf{B} \doteq (\pi, \mathbb{P}, \mathbb{N})$

↳ *Mixture* $\mathbf{M} \doteq \pi \cdot \mathbb{P} + (1 - \pi) \cdot \mathbb{N}$

↳ *discriminator* learns a *posterior* $\eta : \mathcal{X} \rightarrow [0, 1]$
 $\uparrow \hat{\mathbb{P}}[Y=1|X]$

↳ “*ideal*” posterior computes $\mathbb{P}[Y=1|X]$

$$\eta^* = \pi \cdot \frac{d\mathbb{P}}{d\mathbf{M}} \quad \text{\textbf{Bayes posterior}}$$

↳ a **loss** can be decomposed in two *partial losses*

$$\ell(y, u) \doteq \mathbb{I}[y=1] \cdot \ell_1(u) + \mathbb{I}[y=-1] \cdot \ell_{-1}(u)$$

↑ estimated posterior in [0,1]

↑ true label / class in {-1,1}

↳ a loss is *symmetric* iff $\ell_1(u) = \ell_{-1}(1 - u)$

GAN framework in a

Measure-based loss, crafted from *generator*

$$\mathbb{I}_f(\underset{\text{(fake)}}{\mathbb{P}}, \underset{\text{(real)}}{\mathbb{N}}) \doteq \int f\left(\frac{d\mathbb{P}}{d\mathbb{N}}\right) d\mathbb{N} \quad \text{f-divergence}$$

“Information of Binary Task”

↳ variational formulation

$$\mathbb{I}_f(\mathbb{P}, \mathbb{N}) \geq \sup_{\tilde{h}} \left\{ \mathbb{E}_{\mathbb{P}}[\tilde{h}(X)] - \mathbb{E}_{\mathbb{N}}[f^* \circ \tilde{h}(X)] \right\}$$

↑ \tilde{h}

Not an equality in general

↳ *discriminator* hidden in \tilde{h} , seeks to increase the IBT by discriminating *real* vs *fake* (w/ **H** + **G**)

↳ *generator* = \mathbb{N} , seeks to decrease the IBT by generating *fake* data that looks like *real* (w/ **G**)

Our framework

Properness: $\mathbb{P}[Y=1]$ $\mathbb{P}[X|Y=1]$ $\mathbb{P}[X|Y=-1]$

↳ *Binary task* $\mathbb{B} \doteq (\pi, \mathbb{P}, \mathbb{N})$

↳ *Mixture* $\mathbb{M} \doteq \pi \cdot \mathbb{P} + (1 - \pi) \cdot \mathbb{N}$

↳ discriminator learns a *posterior* $\eta : \mathcal{X} \rightarrow [0, 1]$
 $\uparrow \hat{\mathbb{P}}[Y=1|X]$

↳ “*ideal*” posterior computes $\mathbb{P}[Y=1|X]$

$$\eta^* = \pi \cdot \frac{d\mathbb{P}}{d\mathbb{M}} \quad \text{Bayes posterior}$$

↳ a **loss** can be decomposed in two *partial losses*

$$\ell(y, u) \doteq \mathbb{I}[y=1] \cdot \ell_1(u) + \mathbb{I}[y=-1] \cdot \ell_{-1}(u)$$

↑ estimated posterior in [0,1]
 ↑ true label / class in {-1,1}

↳ a loss is *symmetric* iff $\ell_1(u) = \ell_{-1}(1 - u)$

↳ a loss is **strictly proper** iff Bayes posterior solely realises the **inf** of

$$\underline{L}(p) \doteq \inf_u \mathbb{E}_{Y \sim B(p)}[\ell(Y, u)] \quad \text{Google Research}$$

↑ **Bayes risk** (concave)

GAN framework in a

Measure-based loss, crafted from *generator*

$$\mathbb{I}_f(\underset{\text{(fake)}}{\mathbb{P}}, \underset{\text{(real)}}{\mathbb{N}}) \doteq \int f\left(\frac{d\mathbb{P}}{d\mathbb{N}}\right) d\mathbb{N} \quad \text{f-divergence}$$

“Information of Binary Task”

↳ variational formulation

$$\mathbb{I}_f(\mathbb{P}, \mathbb{N}) \geq \sup_{\tilde{h}} \left\{ \mathbb{E}_{\mathbb{P}}[\tilde{h}(X)] - \mathbb{E}_{\mathbb{N}}[f^* \circ \tilde{h}(X)] \right\}$$

↑ \tilde{h}

Not an equality in general

↳ *discriminator* hidden in \tilde{h} , seeks to increase the IBT by discriminating *real* vs *fake* (w/ **H** + **G**)

↳ *generator* = \mathbb{N} , seeks to decrease the IBT by generating *fake* data that looks like *real* (w/ **G**)

Our framework

Properness: $\mathbb{P}[Y=1]$ $\mathbb{P}[X|Y=1]$ $\mathbb{P}[X|Y=-1]$

↳ Binary task $\mathbf{B} \doteq (\pi, \mathbb{P}, \mathbb{N})$

↳ Mixture $\mathbf{M} \doteq \pi \cdot \mathbb{P} + (1 - \pi) \cdot \mathbb{N}$

↳ discriminator learns a *posterior* $\eta : \mathcal{X} \rightarrow [0, 1]$
 $\uparrow \hat{\mathbb{P}}[Y=1|X]$

↳ “*ideal*” posterior computes $\mathbb{P}[Y=1|X]$

$$\eta^* = \pi \cdot \frac{d\mathbb{P}}{d\mathbf{M}} \quad \text{Bayes posterior}$$

↳ a **loss** can be decomposed in two *partial losses*

$$\ell(y, u) \doteq \mathbb{I}[y=1] \cdot \ell_1(u) + \mathbb{I}[y=-1] \cdot \ell_{-1}(u)$$

↑ estimated posterior in [0,1]
 ↑ true label / class in {-1,1}

↳ a loss is *symmetric* iff $\ell_1(u) = \ell_{-1}(1 - u)$

↳ a loss is **strictly proper** iff Bayes posterior solely realises the **inf** of

$$\underline{L}(p) \doteq \inf_u \mathbb{E}_{Y \sim \mathbf{B}(p)}[\ell(Y, u)]$$

↑ **Bayes risk** (concave)

log-, square-, Matusita, etc.

GAN framework in a

Measure-based loss, crafted from *generator*

$$\mathbb{I}_f(\underset{\substack{\text{(real)} \\ \downarrow}}{\mathbb{P}}, \underset{\substack{\text{(fake)} \\ \uparrow}}{\mathbb{N}}) \doteq \int f\left(\frac{d\mathbb{P}}{d\mathbb{N}}\right) d\mathbb{N} \quad f\text{-divergence}$$

“Information of Binary Task”

↳ variational formulation

$$\mathbb{I}_f(\mathbb{P}, \mathbb{N}) \geq \sup_{\tilde{h}} \left\{ \boxed{\mathbb{E}_{\mathbb{P}}[\tilde{h}(X)]}^{\text{H}} + \boxed{\mathbb{E}_{\mathbb{N}}[f^* \circ \tilde{h}(X)]}^{\text{G}} \right\}$$

Not an equality in general

↳ *discriminator* hidden in \tilde{h} , seeks to increase the IBT by discriminating *real* vs *fake* (w/ **H** + **G**)

↳ *generator* = \mathbb{N} , seeks to decrease the IBT by generating *fake* data that looks like *real* (w/ **G**)

Our framework

Partial losses ℓ_1, ℓ_{-1} , Bayes posterior η^* & risk \underline{L}



GAN framework in a

Measure-based loss, crafted from *generator*

$$\mathbb{I}_f(\underset{\text{(fake)}}{\mathbb{P}}, \underset{\text{(real)}}{\mathbb{N}}) \doteq \int f\left(\frac{d\mathbb{P}}{d\mathbb{N}}\right) d\mathbb{N} \quad \textit{f-divergence}$$

“Information of Binary Task”

↳ variational formulation

$$\mathbb{I}_f(\mathbb{P}, \mathbb{N}) \geq \sup_{\tilde{h}} \left\{ \mathbb{E}_{\mathbb{P}}[\tilde{h}(X)] - \mathbb{E}_{\mathbb{N}}[f^* \circ \tilde{h}(X)] \right\}$$

Not an equality in general

↳ *discriminator* hidden in \tilde{h} , seeks to increase the IBT by discriminating *real* vs *fake* (w/ **H** + **G**)

↳ *generator* = \mathbb{N} , seeks to decrease the IBT by generating *fake* data that looks like *real* (w/ **G**)

Our framework

Partial losses l_1, l_{-1} , Bayes posterior η^* & risk \underline{L}

↳ posterior $\tilde{\eta}$ is said **calibrated** iff satisfies

$$\tilde{\eta} = \pi \cdot \frac{d\mathbb{P}_{\tilde{\eta}}}{d\mathbb{M}_{\tilde{\eta}}}$$

σ -algebra coarsened to the level sets of $\tilde{\eta}$

↳ ex: the prior π , Bayes posterior η^* are calibrated

↳ **any decision tree** (w/ empirical posterior prediction at the leaves) is **calibrated**

GAN framework in a

Measure-based loss, crafted from *generator*

$$\mathbb{I}_f(\underset{\text{(fake)}}{P}, \underset{\text{(real)}}{N}) \doteq \int f\left(\frac{dP}{dN}\right) dN \quad \text{\textit{f-divergence}}$$

“Information of Binary Task”

↳ variational formulation

$$\mathbb{I}_f(P, N) \geq \sup_{\tilde{h}} \left\{ \mathbb{E}_P[\tilde{h}(X)] - \mathbb{E}_N[f^* \circ \tilde{h}(X)] \right\}$$

Not an equality in general

↳ *discriminator* hidden in \tilde{h} , seeks to increase the IBT by discriminating **real** vs **fake** (w/ **H** + **G**)

↳ *generator* = N , seeks to decrease the IBT by generating **fake** data that looks like **real** (w/ **G**)

Our framework

Partial losses l_1, l_{-1} , Bayes posterior η^* & risk \underline{L}

↳ posterior $\tilde{\eta}$ is said **calibrated** iff satisfies

$$\tilde{\eta} = \pi \cdot \frac{dP_{\tilde{\eta}}}{dM_{\tilde{\eta}}}$$

σ -algebra coarsened to the level sets of $\tilde{\eta}$

↳ ex: the prior π , Bayes posterior η^* are calibrated

↳ **any decision tree** (w/ empirical posterior prediction at the leaves) is **calibrated**

↳ for any calibrated $\tilde{\eta}$, its *statistical information* is

$$\Delta \underline{L}(\tilde{\eta}, M_{\tilde{\eta}}) = \underline{L}(\pi) - \mathbb{E}_{X \sim M_{\tilde{\eta}}} [\underline{L}(\tilde{\eta}(X))]$$

CART, C4.5, etc.
(splitting criterion)

Google Research

GAN framework in a

Measure-based loss, crafted from *generator*

$$\mathbb{I}_f(\underset{\substack{\text{(real)} \\ \downarrow}}{\mathbb{P}}, \underset{\substack{\uparrow \\ \text{(fake)}}}{\mathbb{N}}) \doteq \int f\left(\frac{d\mathbb{P}}{d\mathbb{N}}\right) d\mathbb{N} \quad f\text{-divergence}$$

“Information of Binary Task”

↳ variational formulation

$$\mathbb{I}_f(\mathbb{P}, \mathbb{N}) \geq \sup_{\tilde{h}} \left\{ \mathbb{E}_{\mathbb{P}}[\tilde{h}(X)] - \mathbb{E}_{\mathbb{N}}[f^* \circ \tilde{h}(X)] \right\}$$

↑ \tilde{h}

Not an equality in general

↳ *discriminator* hidden in \tilde{h} , seeks to increase the IBT by discriminating **real** vs **fake** (w/ **H** + **G**)

↳ *generator* = \mathbb{N} , seeks to decrease the IBT by generating **fake** data that looks like **real** (w/ **G**)

Our framework

Partial losses l_1, l_{-1} , Bayes posterior η^* & risk \underline{L}
For any calibrated $\tilde{\eta}$, its *statistical information*:

$$\Delta \underline{L}(\tilde{\eta}, M_{\tilde{\eta}}) = \underline{L}(\pi) - \mathbb{E}_{X \sim M_{\tilde{\eta}}}[\underline{L}(\tilde{\eta}(X))]$$



GAN framework in a

Measure-based loss, crafted from *generator*

$$\mathbb{I}_f(\underset{\text{(fake)}}{\mathbf{P}}, \underset{\text{(real)}}{\mathbf{N}}) \doteq \int f\left(\frac{d\mathbf{P}}{d\mathbf{N}}\right) d\mathbf{N} \quad f\text{-divergence}$$

“Information of Binary Task”

↳ variational formulation

$$\mathbb{I}_f(\mathbf{P}, \mathbf{N}) \geq \sup_{\tilde{h}} \left\{ \mathbb{E}_{\mathbf{P}}[\tilde{h}(\mathbf{X})] - \mathbb{E}_{\mathbf{N}}[f^* \circ \tilde{h}(\mathbf{X})] \right\}$$

Not an equality in general

↳ *discriminator* hidden in \tilde{h} , seeks to increase the IBT by discriminating **real** vs **fake** (w/ **H** + **G**)

↳ *generator* = \mathbf{N} , seeks to decrease the IBT by generating **fake** data that looks like **real** (w/ **G**)

Our framework

Partial losses ℓ_1, ℓ_{-1} , Bayes posterior η^* & risk \underline{L}

For any calibrated $\tilde{\eta}$, its *statistical information*:

$$\Delta \underline{L}(\tilde{\eta}, \mathbf{M}_{\tilde{\eta}}) = \underline{L}(\pi) - \mathbb{E}_{\mathbf{X} \sim \mathbf{M}_{\tilde{\eta}}}[\underline{L}(\tilde{\eta}(\mathbf{X}))]$$

Theorem: for any calibrated $\tilde{\eta}$ and any *strictly proper symmetric and differentiable* loss ℓ

$$\mathbb{I}_{f^\pi}(\mathbf{P}_{\tilde{\eta}}, \mathbf{N}_{\tilde{\eta}}) = \underset{\substack{\uparrow \\ \text{Not shown for readability}}}{\mathbf{H}} + \underset{\uparrow}{\mathbf{G}} = \Delta \underline{L}(\tilde{\eta}, \mathbf{M}_{\tilde{\eta}})$$

+ if density ratio fct $\ell_{-1}^{\text{DR}}(\rho) \doteq \ell_{-1}\left(\frac{1}{1+\rho}\right)$ cvx, then

$$\mathbf{G} \leq \underline{L}(\pi) - (1 - \pi) \cdot \ell_{-1}\left(\frac{\pi}{1 + (1 - \pi) \cdot \chi^2(\mathbf{N}_{\tilde{\eta}} \parallel \mathbf{P}_{\tilde{\eta}})}\right)$$

chi square

Google Research

GAN framework in a

Measure-based loss, crafted from *generator*

$$\mathbb{I}_f(\underset{\text{(fake)}}{\mathbf{P}}, \underset{\text{(real)}}{\mathbf{N}}) \doteq \int f\left(\frac{d\mathbf{P}}{d\mathbf{N}}\right) d\mathbf{N} \quad f\text{-divergence}$$

"Information of Binary Task"

↳ variational formulation

$$\mathbb{I}_f(\mathbf{P}, \mathbf{N}) \geq \sup_{\tilde{h}} \left\{ \mathbb{E}_{\mathbf{P}}[\tilde{h}(\mathbf{X})] - \mathbb{E}_{\mathbf{N}}[f^* \circ \tilde{h}(\mathbf{X})] \right\}$$

↑ \tilde{h}

Not an equality in general

↳ *discriminator* hidden in \tilde{h} , seeks to increase the IBT by discriminating **real** vs **fake** (w/ **H** + **G**)

↳ *generator* = \mathbf{N} , seeks to decrease the IBT by generating **fake** data that looks like **real** (w/ **G**)

Our framework

Partial losses ℓ_1, ℓ_{-1} , Bayes posterior η^* & risk \underline{L}

For any calibrated $\tilde{\eta}$, its *statistical information*:

$$\Delta \underline{L}(\tilde{\eta}, \mathbf{M}_{\tilde{\eta}}) = \underline{L}(\pi) - \mathbb{E}_{\mathbf{X} \sim \mathbf{M}_{\tilde{\eta}}}[\underline{L}(\tilde{\eta}(\mathbf{X}))]$$

Theorem: for any proper symmetric

True for all tested losses; proof of partial ppty in general case

$$\mathbb{I}_{f^\pi}(\mathbf{P}_{\tilde{\eta}}, \mathbf{N}_{\tilde{\eta}})$$

+ if density ratio fct $\ell_{-1}^{\text{DR}}(\rho) \doteq \ell_{-1}\left(\frac{1}{1+\rho}\right)$ cvx, then

$$\mathbf{G} \leq \underline{L}(\pi) - (1 - \pi) \cdot \ell_{-1}\left(\frac{\pi}{1 + (1 - \pi) \cdot \chi^2(\mathbf{N}_{\tilde{\eta}} \parallel \mathbf{P}_{\tilde{\eta}})}\right)$$

chi square



GAN framework in a

Measure-based loss, crafted from *generator*

$$\mathbb{I}_f(\underset{\text{(fake)}}{\mathbb{P}}, \underset{\text{(real)}}{\mathbb{N}}) \doteq \int f\left(\frac{d\mathbb{P}}{d\mathbb{N}}\right) d\mathbb{N} \quad f\text{-divergence}$$

"Information of Binary Task"

↳ variatio

Summary

Not an equality in general

↳ *discriminator* hidden in \tilde{h} , seeks to increase the IBT by discriminating **real** vs **fake** (w/ **H** + **G**)

↳ *generator* = \mathbb{N} , seeks to decrease the IBT by generating **fake** data that looks like **real** (w/ **G**)

Our framework

Partial losses ℓ_1, ℓ_{-1} , Bayes posterior η^* & risk \underline{L}

For any calibrated $\tilde{\eta}$, its *statistical information*:

$$\Delta \underline{L}(\tilde{\eta}, M_{\tilde{\eta}}) = \underline{L}(\pi) - \mathbb{E}_{X \sim M_{\tilde{\eta}}} [\underline{L}(\tilde{\eta}(X))]$$

strictly

s ℓ

$$\underline{L}(\tilde{\eta}, M_{\tilde{\eta}})$$

Not shown for readability

+ if density ratio fct $\ell_{-1}^{DR}(\rho) \doteq \ell_{-1}\left(\frac{1}{1+\rho}\right)$ cvx, then

$$\mathbf{G} \leq \underline{L}(\pi) - (1 - \pi) \cdot \ell_{-1}\left(\frac{\pi}{1 + (1 - \pi) \cdot \chi^2(\mathbb{N}_{\tilde{\eta}} || \mathbb{P}_{\tilde{\eta}})}\right)$$

chi square



GAN framework in a

Measure-based loss, crafted from *generator*

(real) \downarrow
 $\mathbb{I}_f(\mathbb{P}, \mathbb{N}) \doteq \int f\left(\frac{d\mathbb{P}}{d\mathbb{N}}\right) d\mathbb{N}$ *f*-divergence
 (fake) \uparrow "Information of Binary Task"

↳ variational formulation

$$\mathbb{I}_f(\mathbb{P}, \mathbb{N}) \geq \sup_{\tilde{h}} \left\{ \mathbb{E}_{\mathbb{P}}[\tilde{h}(X)] - \mathbb{E}_{\mathbb{N}}[f^* \circ \tilde{h}(X)] \right\}$$

H + G

Not an equality in general

↳ *discriminator* hidden in \tilde{h} , seeks to increase

Gets the discriminator's loss from the generator's

fake (w/ H + G)
 use the IBT by
 real (w/ G)
 amson, JMLR'11

Our framework

Partial losses ℓ_1, ℓ_{-1} , Bayes posterior η^* & risk \underline{L}

For any calibrated $\tilde{\eta}$, its *statistical information*:

$$\Delta \underline{L}(\tilde{\eta}, M_{\tilde{\eta}}) = \underline{L}(\pi) - \mathbb{E}_{X \sim M_{\tilde{\eta}}}[\underline{L}(\tilde{\eta}(X))]$$

Theorem: for any calibrated $\tilde{\eta}$ and any *strictly proper symmetric and differentiable* loss ℓ

$$\mathbb{I}_{f^\pi}(\mathbb{P}_{\tilde{\eta}}, \mathbb{N}_{\tilde{\eta}}) = \mathbf{H} + \mathbf{G} = \Delta \underline{L}(\tilde{\eta}, M_{\tilde{\eta}})$$

H + G = $\Delta \underline{L}(\tilde{\eta}, M_{\tilde{\eta}})$
Not shown for readability

+ if density ratio fct $\ell_{-1}^{DR}(\rho) \doteq \ell_{-1}\left(\frac{1}{1+\rho}\right)$ cvx then

Gets the generator's loss from the discriminator's

$$\chi^2(\mathbb{N}_{\tilde{\eta}} || \mathbb{P}_{\tilde{\eta}})$$

chi square
 Research

GAN framework in a

Measure-based loss, crafted from *generator*

$$\mathbb{I}_f(\underset{\text{(fake)}}{P}, \underset{\text{(real)}}{N}) \doteq \int f\left(\frac{dP}{dN}\right) dN \quad \text{f-divergence}$$

"Information of Binary Task"

↳ variational formulation

$$\mathbb{I}_f(P, N) \geq \sup_{\tilde{h}} \left\{ \mathbb{E}_P[\tilde{h}(X)] - \mathbb{E}_N[f^* \circ \tilde{h}(X)] \right\}$$

Not an equality, in general

↳ discriminator hidden in \tilde{h} , seeks to increase

Loose approximation (inequality) via variational formulation

fake (w/ **H** + **G**)
 use the IBT by
 real (w/ **G**)
 amson, JMLR'11

Our framework

Partial losses ℓ_1, ℓ_{-1} , Bayes posterior η^* & risk \underline{L}

For any calibrated $\tilde{\eta}$, its *statistical information*:

$$\Delta \underline{L}(\tilde{\eta}, M_{\tilde{\eta}}) = \underline{L}(\pi) - \mathbb{E}_{X \sim M_{\tilde{\eta}}}[\underline{L}(\tilde{\eta}(X))]$$

Theorem: for any calibrated $\tilde{\eta}$ and any strictly proper symmetric and differentiable loss ℓ

$$\mathbb{I}_{f^\pi}(P_{\tilde{\eta}}, N_{\tilde{\eta}}) = \mathbf{H} + \mathbf{G} = \Delta \underline{L}(\tilde{\eta}, M_{\tilde{\eta}})$$

+ if density ratio fct $\ell_{-1}^{DR}(\rho) \doteq \ell_{-1}\left(\frac{1}{1+\rho}\right)$ cvx, then

Tight characterisation (all equalities), and one loss "to train against them all": the chi square

$\chi^2(N_{\tilde{\eta}} || P_{\tilde{\eta}})$ chi square

Google Research

: Adversarial and Copycat

GAN framework in a

Measure-based loss, crafted from *generator*

$$\mathbb{I}_f(\underset{\text{(fake)}}{\mathbf{P}}, \underset{\text{(real)}}{\mathbf{N}}) \doteq \int f\left(\frac{d\mathbf{P}}{d\mathbf{N}}\right) d\mathbf{N} \quad f\text{-divergence}$$

"Information of Binary Task"

↳ variational formulation

$$\mathbb{I}_f(\mathbf{P}, \mathbf{N}) \geq \sup_{\tilde{h}} \left\{ \mathbb{E}_{\mathbf{P}}[\tilde{h}(\mathbf{X})] - \mathbb{E}_{\mathbf{N}}[f^* \circ \tilde{h}(\mathbf{X})] \right\}$$

Not an equality in general

↳ *discriminator* hidden in \tilde{h} , seeks to increase

"No" assumption necessary

ake (w/ **H** + **G**)
 se the IBT by
 e real (w/ **G**)
 amson, JMLR'11

Our framework

Partial losses ℓ_1, ℓ_{-1} , Bayes posterior η^* & risk \underline{L}

For any calibrated $\tilde{\eta}$, its *statistical information*:

$$\Delta \underline{L}(\tilde{\eta}, \mathbf{M}_{\tilde{\eta}}) = \underline{L}(\pi) - \mathbb{E}_{\mathbf{X} \sim \mathbf{M}_{\tilde{\eta}}}[\underline{L}(\tilde{\eta}(\mathbf{X}))]$$

Theorem: for any calibrated $\tilde{\eta}$ and any *strictly proper symmetric and differentiable* loss ℓ

$$\mathbb{I}_{f^\pi}(\mathbf{P}_{\tilde{\eta}}, \mathbf{N}_{\tilde{\eta}}) = \underset{\text{Not shown for readability}}{\mathbf{H}} + \mathbf{G} = \Delta \underline{L}(\tilde{\eta}, \mathbf{M}_{\tilde{\eta}})$$

+ if density ratio fct $\ell_{-1}^{\text{DR}}(\rho) \doteq \ell_{-1}\left(\frac{1}{1+\rho}\right)$ cvx, then

Discriminator calibrated

$$\pi \cdot \chi^2(\mathbf{N}_{\tilde{\eta}} || \mathbf{P}_{\tilde{\eta}})$$

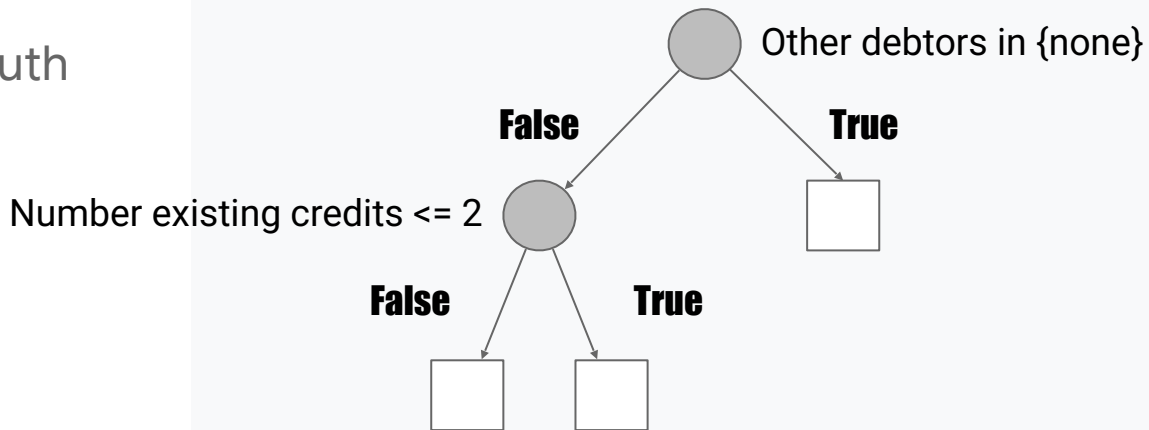
chi square
 ogle Research

: Adversarial and Copycat

Models

Tree

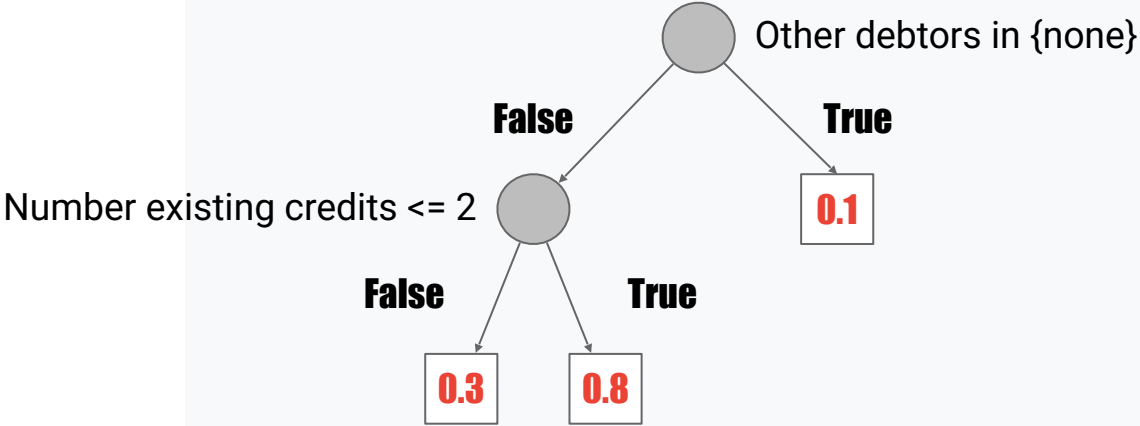
A **tree** is a binary directed tree whose internal nodes are labeled with a test on an observation variable and outgoing arcs are labeled with truth values. Leaves are blank.



(Labelling from UCI German Credit)

Decision Tree (DT)

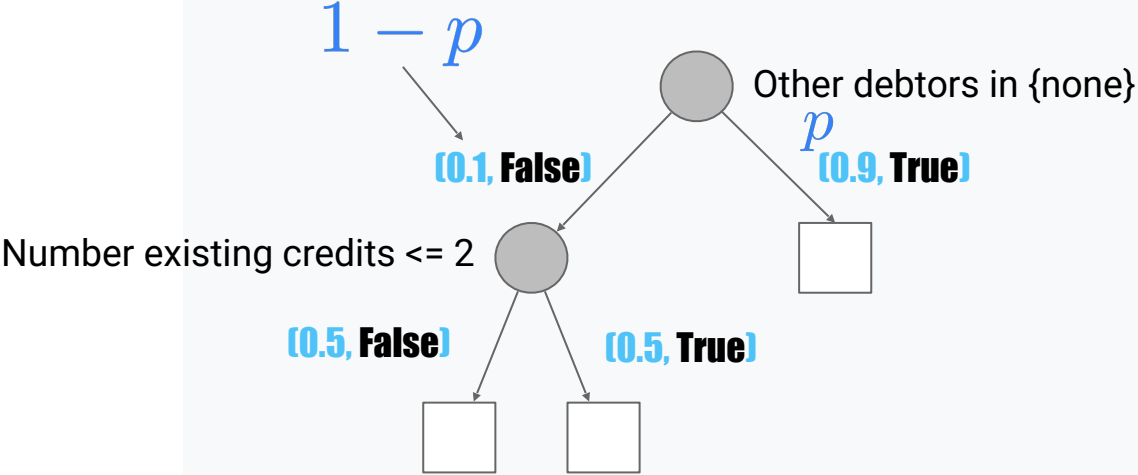
A *decision tree* h is a tree in which leaves are labeled by values in $[0,1]$



(Labelling from UCI German Credit)

Generative Tree (GT)

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

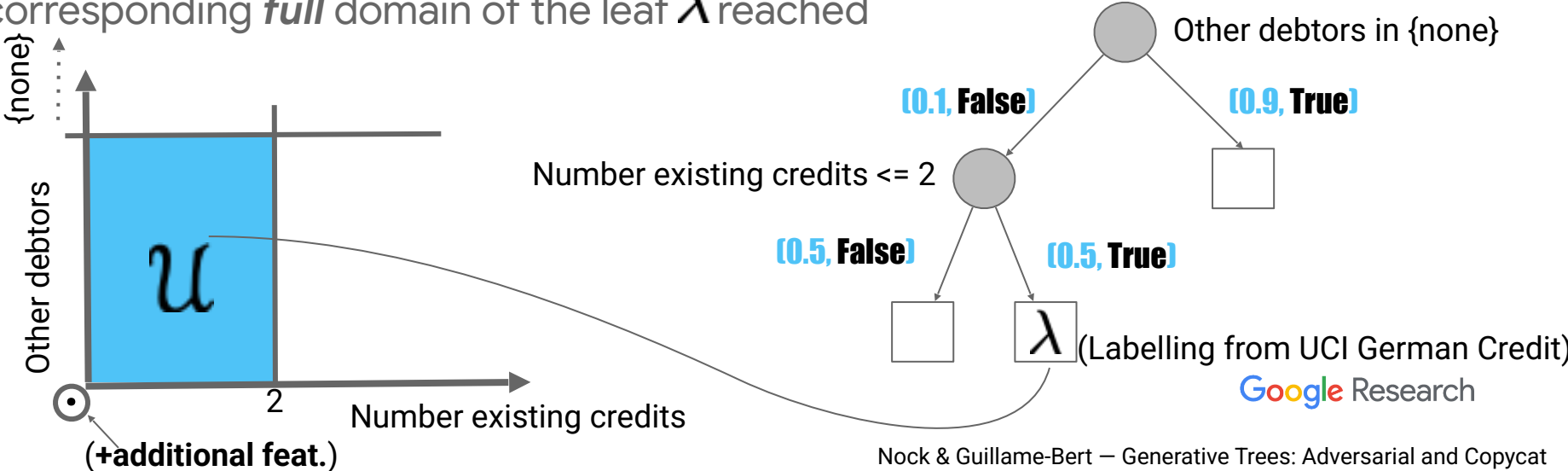


(Labelling from UCI German Credit)

Key routines

For a **decision tree h** : for a given observation $x \in \mathcal{X}$, return the leaf $\lambda(x)$ whose path in the tree is satisfied by x

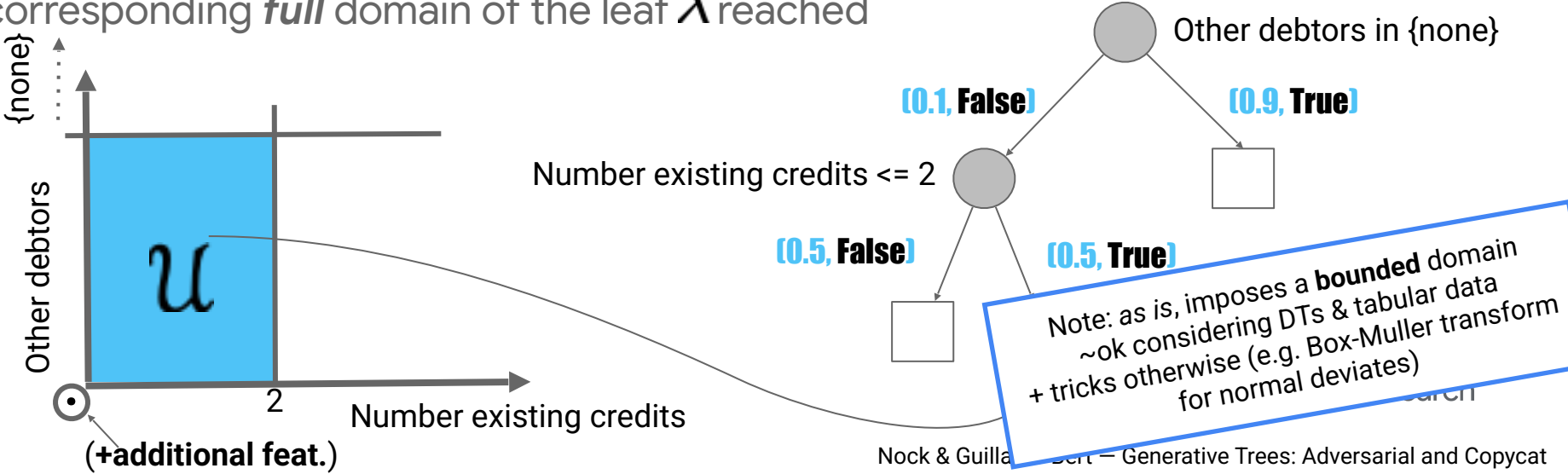
For a **generative tree G** : sample a path (wrt “Bernoullis”) and sample uniformly in the corresponding **full** domain of the leaf λ reached



Key routines

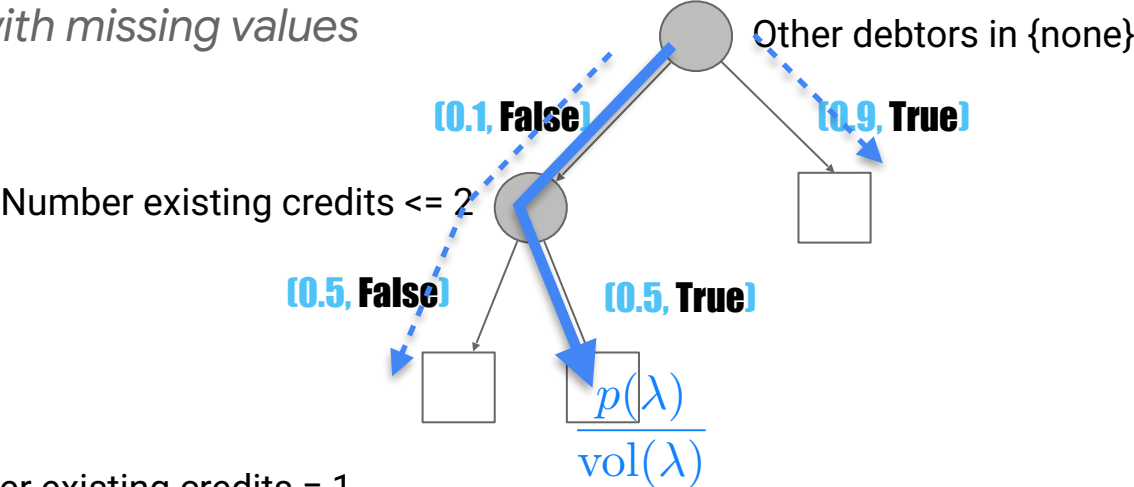
For a **decision tree** h : for a given observation $x \in \mathcal{X}$, return the leaf $\lambda(x)$ whose path in the tree is satisfied by x

For a **generative tree** G : sample a path (wrt “Bernoullis”) and sample uniformly in the corresponding **full** domain of the leaf λ reached



Additional conveniences of generative trees

- ↳ For any observation, *local density* computable in $O(\text{depth}(\mathbf{G}))$
- ↳ If missing values, *likelihood* | observed values & generator \mathbf{G} available in $O(\text{size}(\mathbf{G}))$
- ↳ XAI / fairness: “as easy” to interpret as a **decision tree**
- ↳ Easily trainable from data *with missing values*



Solid blue arrow: Other debtors = guarantor, Number existing credits = 1, ...
Dashed blue arrow: Other debtors = ?, Number existing credits = 3, ...

Algorithms

Adversarial

- ↳ GAN-style (for “vs” training)
- ↳ simple (leaf→feature→split→ p in $B(p)$ →repeat)

Boosting compliance in generative framework:

- ↳ a weak *generative* assumption = non total independence between data generation (**G**) and classification (**h**)
- ↳ most “expensive” computational bit = the computation of Bernoulli p 's
- ↳ geometric convergence of the chi square

$$\chi^2 (N_{\tilde{\eta}}^{\text{new}} || P_{\tilde{\eta}}) \leq \frac{1}{1+Q} \cdot \chi^2 (N_{\tilde{\eta}}^{\text{old}} || P_{\tilde{\eta}})$$

↑
(details in paper)

Adversarial

- ↳ GAN-style (for “vs” training)
- ↳ simple (leaf→feature→split→ p in $B(p)$ →repeat)

Boosting compliance in generative framework:

- ↳ a weak *generative* assumption = non total independence between data generation (G) and classification (h)
- ↳ most “expensive” computational bit = the computation of Bernoulli p 's
- ↳ geometric convergence of the chi square

$$\chi^2(N_{\tilde{\eta}}^{\text{new}} || P_{\tilde{\eta}}) \leq \frac{1}{1+Q} \cdot \chi^2(N_{\tilde{\eta}}^{\text{old}} || P_{\tilde{\eta}})$$

Copycat

- ↳ Powerful boosting DT induction algorithms for discriminator h . Can we rely on them to train G ?

Adversarial

- ↳ GAN-style (for “vs” training)
- ↳ simple (leaf→feature→split→ p in $B(p)$ →repeat)

Boosting compliance in generative framework:

- ↳ a weak *generative* assumption = non total independence between data generation (G) and classification (h)
- ↳ most “expensive” computational bit = the computation of Bernoulli p 's
- ↳ geometric convergence of the chi square

$$\chi^2(N_{\tilde{\eta}}^{\text{new}} || P_{\tilde{\eta}}) \leq \frac{1}{1+Q} \cdot \chi^2(N_{\tilde{\eta}}^{\text{old}} || P_{\tilde{\eta}})$$

Copycat

- ↳ Powerful boosting DT induction algorithms for discriminator h . Can we rely on them to train G ?

Train G at “0” additional cost & with guarantees

- ↳ GT G and DT h share a **tree** (graph)
- ↳ G copies h 's tree at induction time & completes it (p) for *hardest* current generator
- ↳ G = balanced distribution of the weak learning assumption in Kearns & Mansour, STOC'96.
- ↳ trivial computations for G + geometric convergence in density ratio loss *for free* from boosting

Google Research

(details in paper)

Experiments

Summary

Experiments carried out with **Copycat training** (fast, simple, little hyperparameter tuning required, ...), using Kearns and Mansour's optimal **top-down algorithm**;

1 classical toy generative problem + 4 more experiments against SOTA

↳ **Toy**: 2D heatmaps of densities (vs CTGAN)

↳ **Missing data imputation**: predict missing values in a dataset (vs MICE)

↳ **Gen-discrim**: discriminate between fake and real examples (vs CTGAN)

↳ **Train-gen** (*supervised data*): train model over fake data, test over real (vs CTGAN)

↳ **Gen-aug** (*supervised data*): augment real with generated + Train-gen (vs CTGAN)

(details in paper)

Google Research

CTGAN: Xu, Skoularidou, Cuesta-Infante & Veeramachaneni, NeurIPS'19

MICE: van Buuren, "Flexible imputation of missing data", Chapman & Hall / CRC, 2018

XAI in a

```
[#1:root]
-[0.0489, [ lng (CONTINUOUS) in [-76.8665, -72.7167]; |-1|-1| ] ]--[#2]
|-[0.0366, [ search_vehicle (NOMINAL) in {TRUE}; |-1|-1| ] ]--[#100]
| |-[0.1212, [ lat (CONTINUOUS) in [40.7067, 41.7329]; |-1|-1| ] ]--[#3010 (sampling)]
| | \-[0.8788, [ lat (CONTINUOUS) in [41.7329, 42.3426]; |-1|-1| ] ]--[#3011]
| | |-[0.8276, [ lat (CONTINUOUS) in [41.7329, 41.8060]; |-1|-1| ] ]--[#3338]
| | | |-[0.8333, [ warning_issued (NOMINAL) in {FALSE}; |-1|-1| ] ]--[#3756]
| | | | |-[0.9500, [ raw_subject_race_code (NOMINAL) in {W, B}; |-1|-1| ] ]--[#4184]
| | | | | |-[0.2632, [ reason_for_stop (NOMINAL) in {TrafficControlSignal, Other}; |-1|-1| ] ]--[#4380 (sampling)]
| | | | | | \-[0.7368, [ reason_for_stop (NOMINAL) in {StopSign, DefectiveLights, CellPhone, SuspendedLicense, Registration,
| | | | | | |-[0.5000, [ district (NOMINAL) in {BARRYSQUARE, NORTHMEADOWS}; |-1|-1| ] ]--[#5788]
| | | | | | | |-[0.4286, [ subject_age (INTEGER) in {14, 15, ..., 29}; |-1|-1| ] ]--[#9118 (sampling)]
| | | | | | | | \-[0.5714, [ subject_age (INTEGER) in {30, 31, ..., 94}; |-1|-1| ] ]--[#9119 (sampling)]
| | | | | | | | \-[0.5000, [ district (NOMINAL) in {SOUTHWEST, ASYLUMHILL, PARKVILLE, FROGHOLLOW, BEHINDTHEROCKS, SOUTHGREEN,
| | | | | | | | \-[0.0500, [ raw_subject_race_code (NOMINAL) in {A, I}; |-1|-1| ] ]--[#4185 (sampling)]
| | | | | | | | \-[0.1667, [ warning_issued (NOMINAL) in {TRUE}; |-1|-1| ] ]--[#3757 (sampling)]
| | | | | | | | \-[0.1724, [ lat (CONTINUOUS) in [41.8060, 42.3426]; |-1|-1| ] ]--[#3339 (sampling)]
| | | | | | | | \-[0.9634, [ search_vehicle (NOMINAL) in {FALSE}; |-1|-1| ] ]--[#101]
| | | | | | | | |-[0.0265, [ raw_search_authorization_code (NOMINAL) in {C, I}; |-1|-1| ] ]--[#3339 (sampling)]
| | | | | | | | | |-[0.0870, [ lat (CONTINUOUS) in [40.7067, 41.6730]; |-1|-1| ] ]--[#3339 (sampling)]
| | | | | | | | | | \-[0.9130, [ lat (CONTINUOUS) in [41.6730, 42.3426]; |-1|-1| ] ]--[#3339 (sampling)]
```

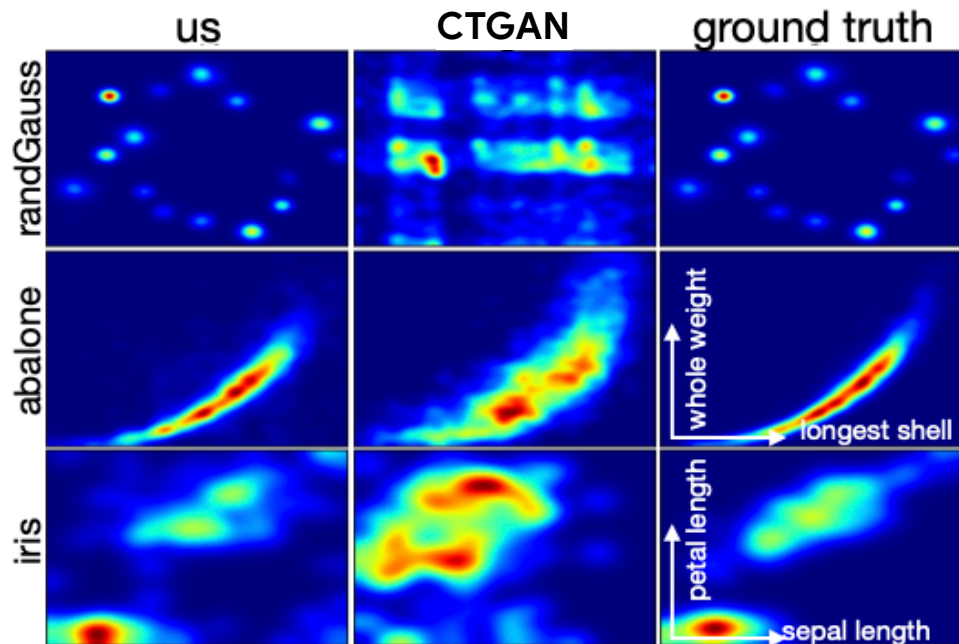
↑↑ disparities in density **young vs not-young** on “car/driver based search” in **specific area**

Example of **generative tree** learned on Stanford Open Policing / Hartford (more examples in paper)

Toy 2D heat maps

↳ **Setup:** generate data, compare with ground truth (10 000 nodes GT, 1K epoch CTGANs)

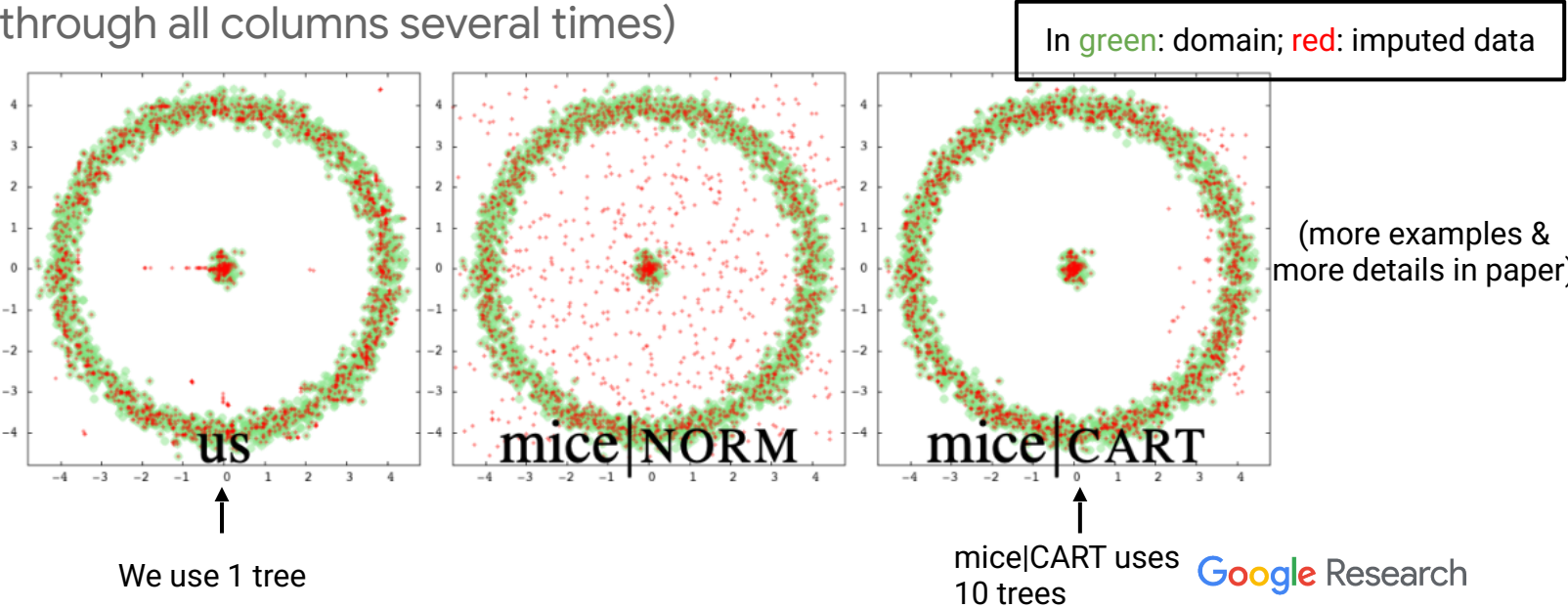
↳ Some results:



Google Research

Missing Data Imputation

↳ **Summary:** synthetic data, remove $q\%$ features (Missing Completely At Random), impute w/ GT vs SOTA = mice (CART: use decision trees to predict missing in one column given the others, cycle through all columns several times)

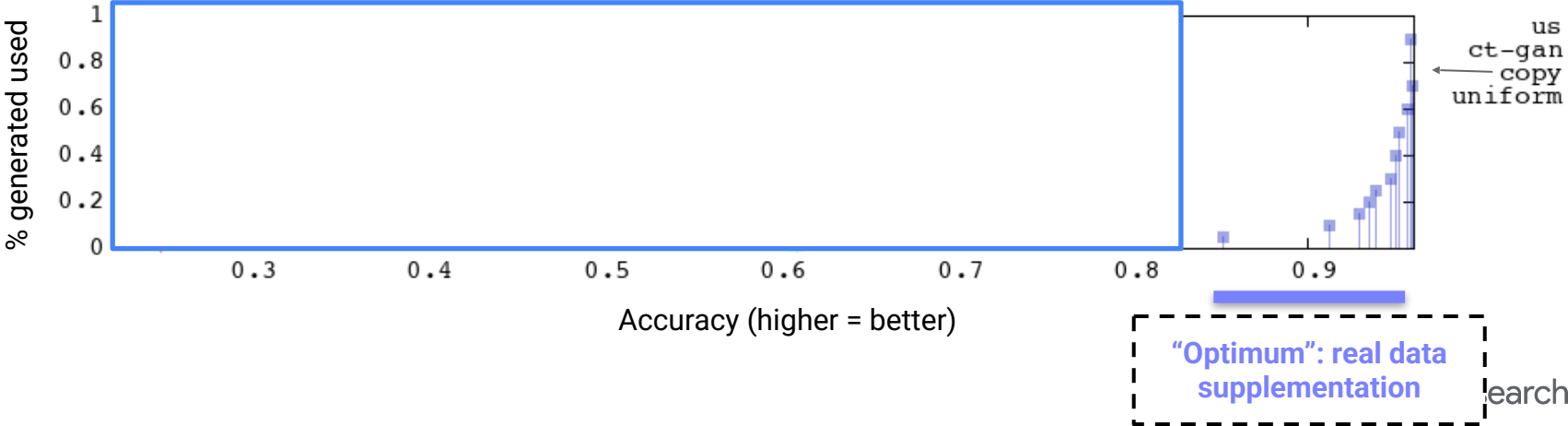


Gen-aug

↳ **Summary:** use part of real data to train generator, supplement remaining training data with varying % of generated data, train supervised classifier for the task, evaluate accuracy on test data

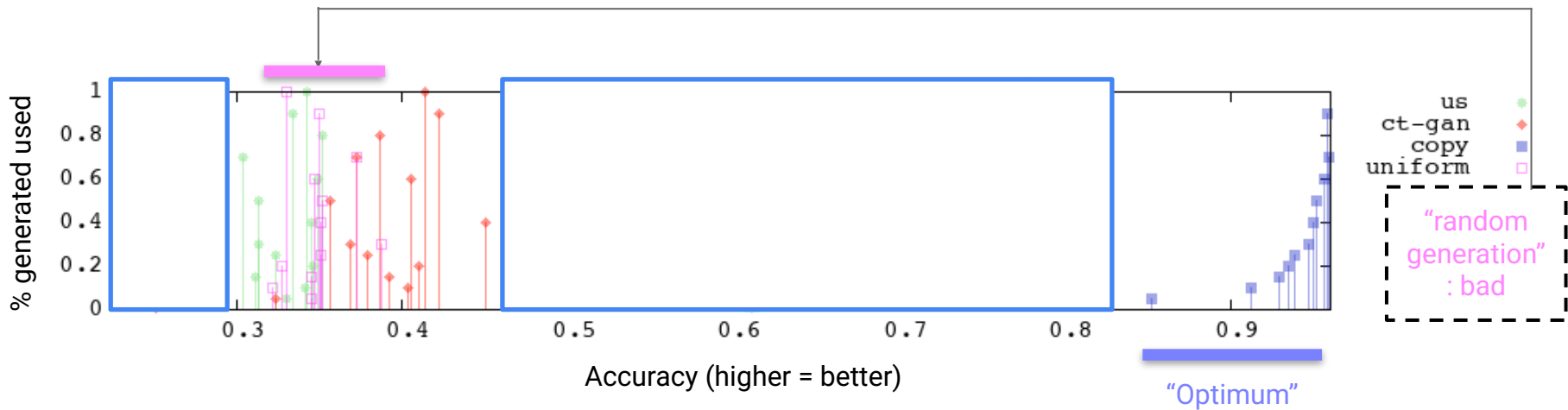
Gen-aug

↳ **Summary:** use part of real data to train generator, supplement remaining training data with varying % of generated data, train supervised classifier for the task, evaluate accuracy on test data — example of UCI DNA, 181 binary features



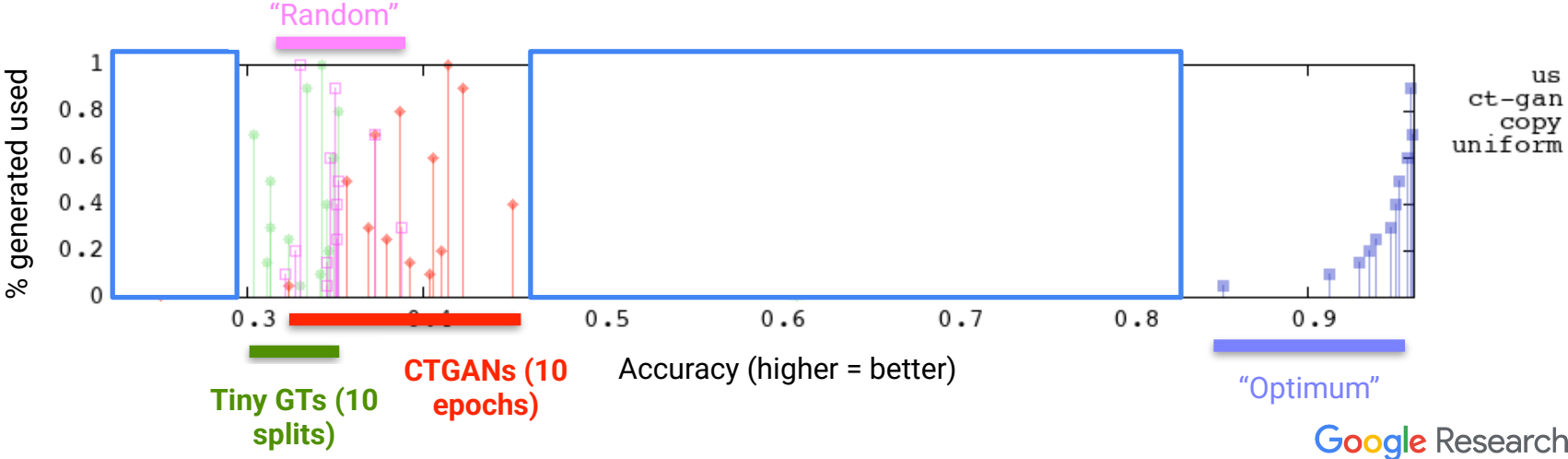
Gen-aug

↳ **Summary:** use part of real data to train generator, supplement remaining training data with varying % of generated data, train supervised classifier for the task, evaluate accuracy on test data — example of UCI DNA, 181 binary features



Gen-aug

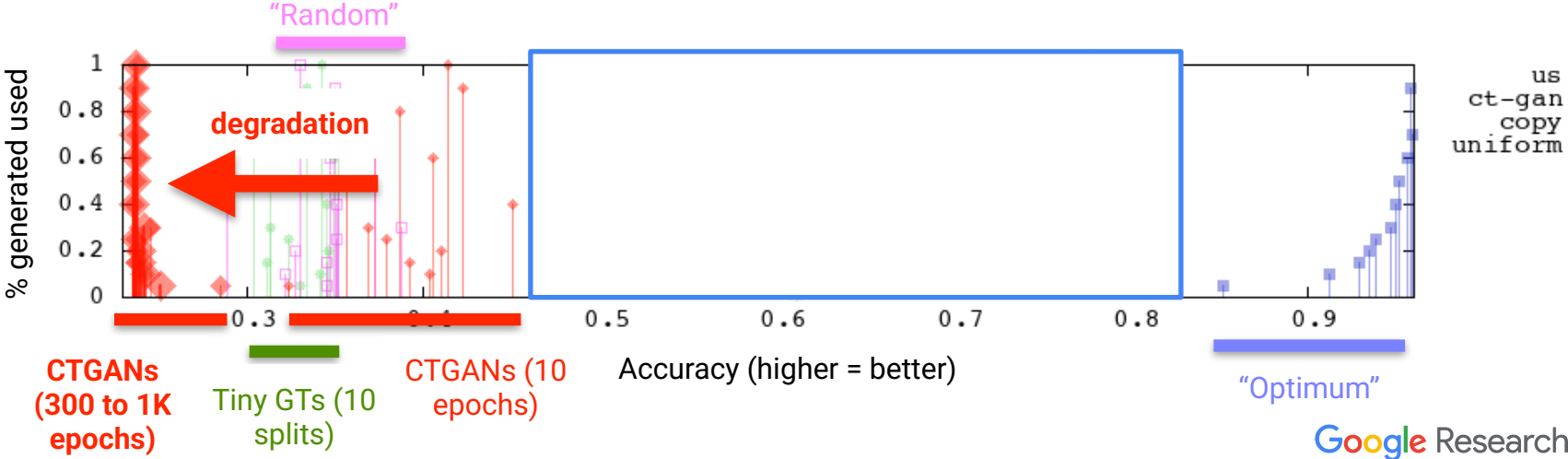
↳ **Summary:** use part of real data to train generator, supplement remaining training data with varying % of generated data, train supervised classifier for the task, evaluate accuracy on test data — example of UCI DNA, 181 binary features



Google Research

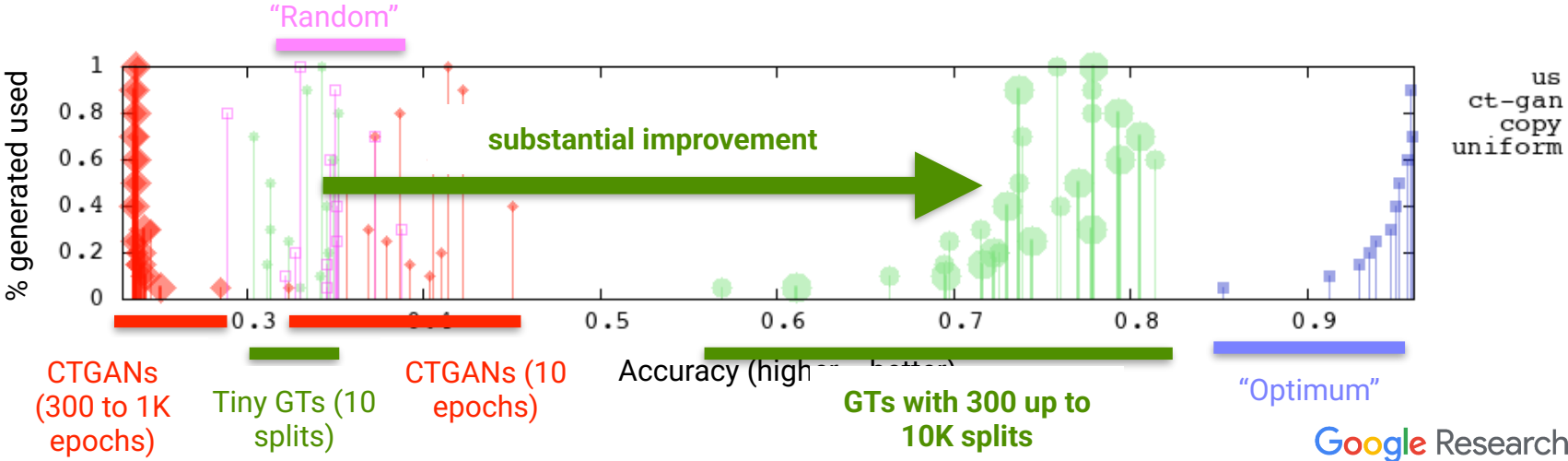
Gen-aug

↳ **Summary:** use part of real data to train generator, supplement remaining training data with varying % of generated data, train supervised classifier for the task, evaluate accuracy on test data — example of UCI DNA, 181 binary features



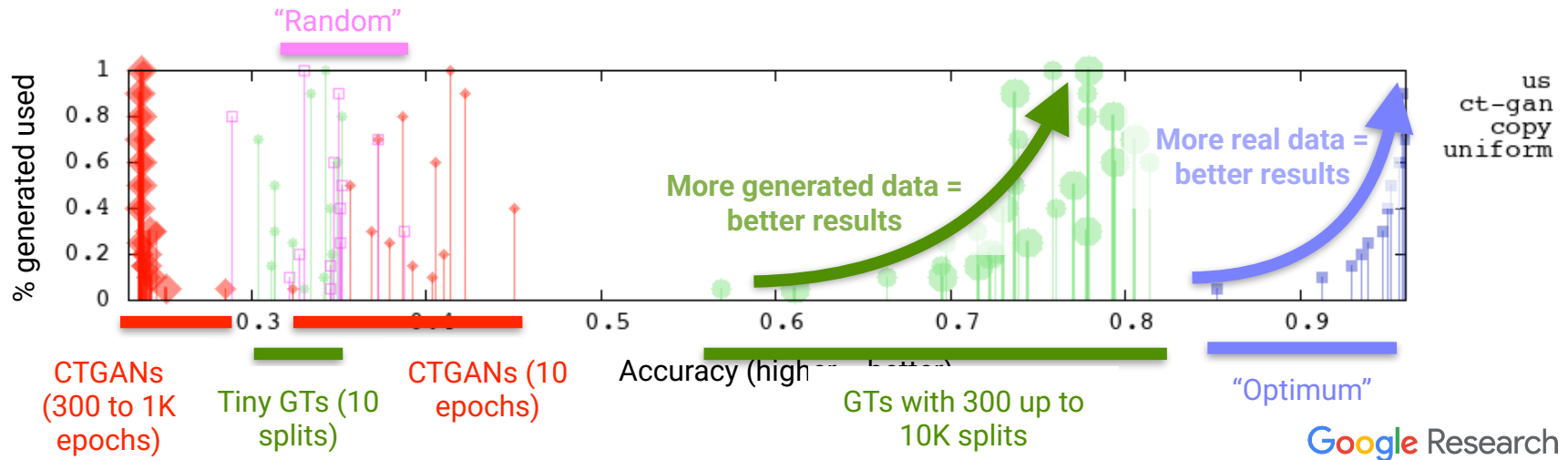
Gen-aug

↳ **Summary:** use part of real data to train generator, supplement remaining training data with varying % of generated data, train supervised classifier for the task, evaluate accuracy on test data — example of UCI DNA, 181 binary features



Gen-aug

↳ **Summary:** use part of real data to train generator, supplement remaining training data with varying % of generated data, train supervised classifier for the task, evaluate accuracy on test data — example of UCI DNA, 181 binary features



↳ See paper for more results

Conclusion / future work

Our contributions

- ↳ new *tight* formulation of the GAN losses *from the supervised side* (properness) if discriminator calibrated, gives the *chi square* as a “default” generator training loss
- ↳ new *generative models* & adversarial training w/ *boosting compliant convergence*
- ↳ new *cheap* training for *generative models* (*copycat*) + “boosting for free” convergence

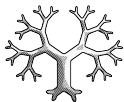
Future work includes

- ↳ XAI / fairness: constrained induction of generative models
- ↳ privacy
- ↳ lots of formal questions (generalisation, pruning generators, ensembles of GTs, etc.)

Thank you !



{richardnock,gbm}@google.com



See also: **Yggdrasil** <https://github.com/google/yggdrasil-decision-forests>

Decision Forests

Google Research