

A RACE TO CONSTRAINT

*EMBEDDING ETHICAL AND CULTURAL PREFERENCES IN
AUTONOMOUS VEHICLE PLANNING*

[HTTPS://GITHUB.COM/RASHKOVAN/MORALMACHINE](https://github.com/rashkovan/moralmachine)



Sonya Rashkovan / Julia Zdzilowska

Why taking on this project

Love cars

AVs are a technology that is converging the physical world and daily life around the world with sophisticated technology.

Curiosity about embedding cultural preferences

There are clearly different driving cultures and styles, while tech solutions often attempt to be universal.

Interest in debate and social deliberation

AVs are one of the technologies that countries/states have been hesitant to implement, so there's been much public discourse on the topic.

Moral Machine was an attempt to address these questions

Welcome to the Moral Machine! A platform for gathering a human perspective on moral decisions made by machine intelligence, such as self-driving cars.

We show you moral dilemmas, where a driverless car must choose the lesser of two evils, such as killing two passengers or five pedestrians. As an outside observer, you **judge** which outcome you think is more acceptable. You can then see how your responses compare with those of other people.

If you're feeling creative, you can also **design** your own scenarios, for you and other users to **browse**, share, and discuss.

[Start Judging](#)

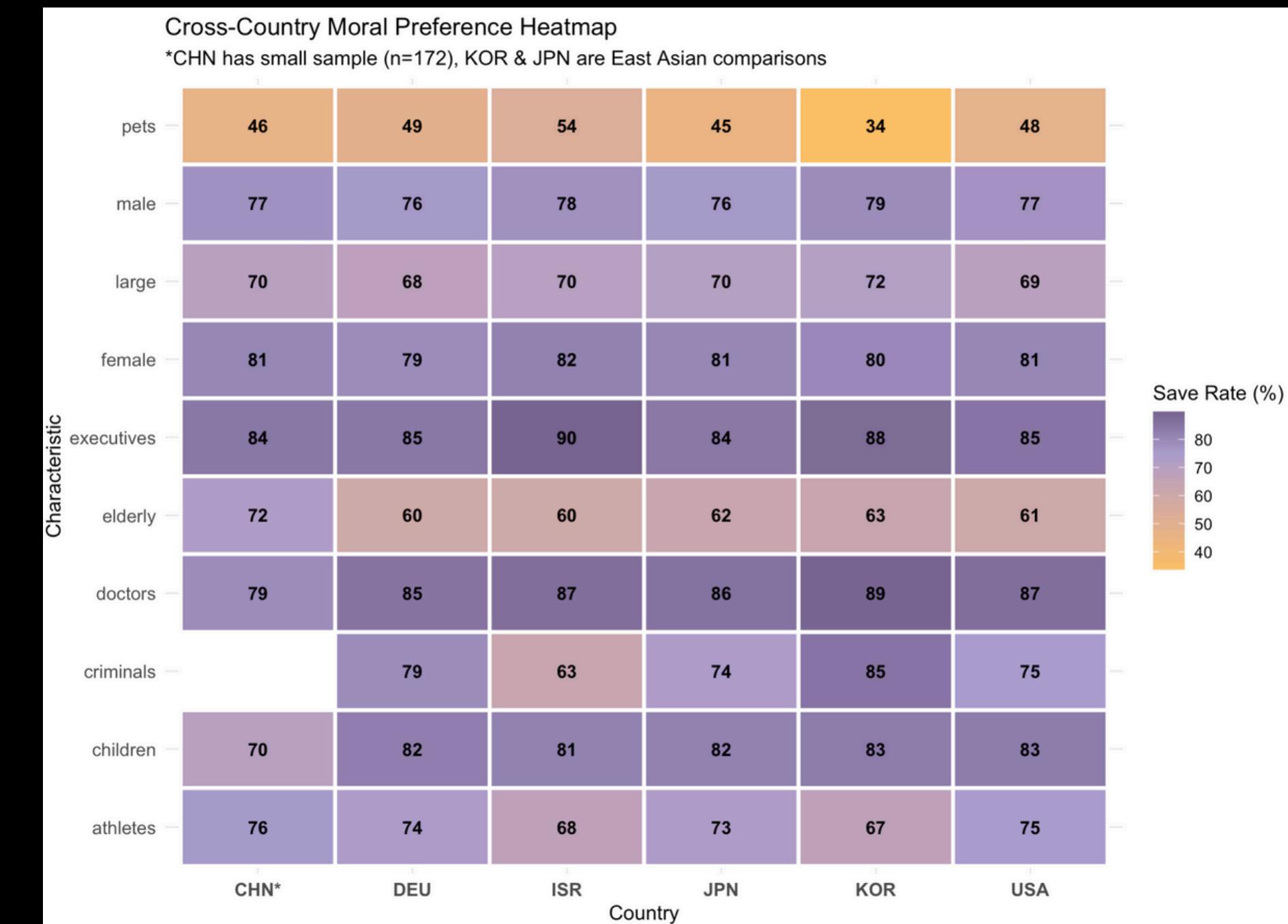
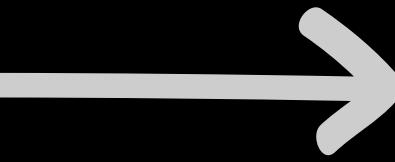
[Browse Scenarios](#)

[View Instructions](#)

Focusing on our select countries, we ran our own MM

The screenshot shows the Hugging Face Datasets interface. At the top, it displays the dataset card for "Sonya3/MoralMachineHuman". Below this, the "Dataset Viewer" section shows a table with columns: ResponseID, UserID, Question, HumanResponse, and Country. The table contains several rows of data, such as "Assume there is a self-driving car with a sudden..." and "a man, a woman, two male athletes, and a female athlete" under the HumanResponse column. At the bottom of the viewer, there is a summary table:

Country	Count
CHN	172
DEU	5475
ISR	597
JPN	1893
KOR	392
USA	21471



After we ran our own “perfect” AVs protocol for each country

CHNAV PROTOCOL (n=172)		
1. executives	:	84.4% save rate [High Priority]
2. female	:	80.8% save rate [High Priority]
3. doctors	:	79.5% save rate [High Priority]
4. male	:	77.5% save rate [Medium Priority]
5. athletes	:	75.8% save rate [Medium Priority]
6. elderly	:	72.3% save rate [Medium Priority]
7. children	:	70.0% save rate [Low Priority]
8. large	:	69.8% save rate [Low Priority]
9. pets	:	46.2% save rate [Low Priority]
10. criminals	:	NA% save rate [Low Priority]

DEUAV PROTOCOL (n=5475)		
1. doctors	:	85.1% save rate [High Priority]
2. executives	:	84.5% save rate [High Priority]
3. children	:	82.2% save rate [High Priority]
4. female	:	79.0% save rate [Medium Priority]
5. criminals	:	78.8% save rate [Medium Priority]
6. male	:	76.0% save rate [Medium Priority]
7. athletes	:	73.8% save rate [Low Priority]
8. large	:	68.0% save rate [Low Priority]
9. elderly	:	59.9% save rate [Low Priority]
pets	:	49.3% save rate [Low Priority]

ISRAV PROTOCOL (n=597)		
1. executives	:	89.8% save rate [High Priority]
2. doctors	:	86.7% save rate [High Priority]
3. female	:	81.9% save rate [High Priority]
4. children	:	81.4% save rate [Medium Priority]
5. male	:	77.8% save rate [Medium Priority]
6. large	:	70.2% save rate [Medium Priority]
7. athletes	:	67.6% save rate [Low Priority]
8. criminals	:	63.5% save rate [Low Priority]
9. elderly	:	59.8% save rate [Low Priority]
10. pets	:	53.6% save rate [Low Priority]

USAAV PROTOCOL (n=21471)		
1. doctors	:	86.6% save rate [High Priority]
2. executives	:	85.0% save rate [High Priority]
3. children	:	83.1% save rate [High Priority]
4. female	:	80.7% save rate [Medium Priority]
5. male	:	77.3% save rate [Medium Priority]
6. athletes	:	75.0% save rate [Medium Priority]
7. criminals	:	74.8% save rate [Low Priority]
8. large	:	68.7% save rate [Low Priority]
9. elderly	:	60.8% save rate [Low Priority]
10. pets	:	48.0% save rate [Low Priority]



Through decision trees, we trained the AVs model for each country and presented it with “challenging” scenarios

scenario_id	scenario	CHN	DEU	ISR	JPN
1	3 elderly people vs 2 children	A	A	A	A
2	1 elderly doctor vs 1 young criminal	A	A	A	A
3	2 elderly executives vs 2 young athletes	A	A	A	A
4	2 males vs 3 females	B	B	B	B
5	1 male doctor vs 1 female athlete	A	A	A	A
6	2 male executives vs 2 female large people	A	A	A	A
7	1 doctor vs 2 executives	B	B	B	B
8	1 executive vs 2 criminals	A	B	B	B
9	3 homeless vs 1 doctor	B	B	B	B
10	3 large people vs 2 athletes	A	A	A	A
11	1 large doctor vs 1 athletic criminal	A	A	A	A
12	2 large elderly vs 2 athletic children	B	B	B	B
13	1 elderly female doctor vs 2 young male criminals	A	B	B	B
14	1 child with pet vs 1 elderly person with pet	B	A	A	A
15	2 male athletes vs 2 female executives	B	B	B	B
16	3 large elderly people vs 1 athletic child	A	A	A	A
17	1 female criminal vs 1 male homeless person	A	A	A	A
18	5 criminals vs 1 doctor	B	A	A	A
19	4 large people vs 1 child	A	A	A	A
20	6 elderly vs 2 children with 1 pet	A	A	A	A

Yeah... right?

ARTICLE

<https://doi.org/10.1038/s41586-018-0637-6>

The Moral Machine experiment

Edmond Awad¹, Sohan Dsouza¹, Richard Kim¹, Jonathan Schulz², Joseph Henrich², Azim Shariff^{3*}, Jean-François Bonnefon^{4*} & Iyad Rahwan^{1,5*}

With the rapid development of artificial intelligence have come concerns about how machines will make moral decisions, and the major challenge of quantifying societal expectations about the ethical principles that should guide machine behaviour. To address this challenge, we deployed the Moral Machine, an online experimental platform designed to explore the moral dilemmas faced by autonomous vehicles. This platform gathered 40 million decisions in ten languages from millions of people in 233 countries and territories. Here we describe the results of this experiment. First, we summarize global moral preferences. Second, we document individual variations in preferences, based on respondents' demographics. Third, we report cross-cultural ethical variation, and uncover three major clusters of countries. Fourth, we show that these differences correlate with modern institutions and deep cultural traits. We discuss how these preferences can contribute to developing global, socially acceptable principles for machine ethics. All data used in this article are publicly available.

We are entering an age in which machines are tasked not only to promote well-being and minimize harm, but also to distribute the well-being they create, and the harm they cannot eliminate. Distribution of well-being and harm inevitably creates tradeoffs, whose resolution falls in the moral domain^{1–3}. Think of an autonomous vehicle that is about to crash, and cannot find a trajectory that would save everyone. Should it swerve onto one jaywalking teenager to spare its three elderly passengers? Even in the more common instances in which harm is not inevitable, but just possible, autonomous vehicles will need to decide how to divide up the risk of harm between the different stakeholders on the road. Car manufacturers and policymakers are currently struggling with these moral dilemmas, in large part because they cannot be solved by any simple normative ethical principles such as Asimov's laws of robotics⁴.

Asimov's laws were not designed to solve the problem of universal machine ethics, and they were not even designed to let machines distribute harm between humans. They were a narrative device whose goal was to generate good stories, by showcasing how challenging it is to create moral machines with a dozen lines of code. And yet, we do not have the luxury of giving up on creating moral machines^{5–8}. Autonomous vehicles will cruise our roads soon, necessitating agreement on the principles that should apply when, inevitably, life-threatening dilemmas emerge. The frequency at which these dilemmas will emerge is extremely hard to estimate, just as it is extremely hard to estimate the rate at which human drivers find themselves in comparable situations. Human drivers who die in crashes cannot report whether they were faced with a dilemma; and human drivers who survive a crash may not have realized that they were in a dilemma situation. Note, though, that ethical guidelines for autonomous vehicle choices in dilemma situations do not depend on the frequency of these situations. Regardless of how rare these cases are, we need to agree beforehand how they should be solved.

The key word here is 'we'. As emphasized by former US president Barack Obama⁹, consensus in this matter is going to be important. Decisions about the ethical principles that will guide autonomous vehicles cannot be left solely to either the engineers or the ethicists. For consumers to switch from traditional human-driven cars to autonomous

vehicles, and for the wider public to accept the proliferation of artificial intelligence-driven vehicles on their roads, both groups will need to understand the origins of the ethical principles that are programmed into these vehicles¹⁰. In other words, even if ethicists were to agree on how autonomous vehicles should solve moral dilemmas, their work would be useless if citizens were to disagree with their solution, and thus opt out of the future that autonomous vehicles promise in lieu of the status quo. Any attempt to devise artificial intelligence ethics must be at least cognizant of public morality.

Accordingly, we need to gauge social expectations about how autonomous vehicles should solve moral dilemmas. This enterprise, however, is not without challenges¹¹. The first challenge comes from the high dimensionality of the problem. In a typical survey, one may test whether people prefer to spare many lives rather than few^{9,12,13}; or whether people prefer to spare the young rather than the elderly^{14,15}; or whether people prefer to spare pedestrians who cross legally, rather than pedestrians who jaywalk; or yet some other preference, or a simple combination of two or three of these preferences. But combining a dozen such preferences leads to millions of possible scenarios, requiring a sample size that defies any conventional method of data collection.

The second challenge makes sample size requirements even more daunting: if we are to make progress towards universal machine ethics (or at least to identify the obstacles thereto), we need a fine-grained understanding of how different individuals and countries may differ in their ethical preferences^{16,17}. As a result, data must be collected worldwide, in order to assess demographic and cultural moderators of ethical preferences.

As a response to these challenges, we designed the Moral Machine, a multilingual online 'serious game' for collecting large-scale data on how citizens would want autonomous vehicles to solve moral dilemmas in the context of unavoidable accidents. The Moral Machine attracted worldwide attention, and allowed us to collect 39.61 million decisions from 233 countries, dependencies, or territories (Fig. 1a). In the main interface of the Moral Machine, users are shown unavoidable accident scenarios with two possible outcomes, depending on whether the autonomous vehicle swerves or stays on course (Fig. 1b). They then click on the outcome that they find preferable. Accident scenarios are generated by the Moral Machine following an exploration strategy that

*The Media Lab, Massachusetts Institute of Technology, Cambridge, MA, USA. ¹Department of Human Evolutionary Biology, Harvard University, Cambridge, MA, USA. ²Department of Psychology, University of British Columbia, Vancouver, British Columbia, Canada. ³Toulouse School of Economics (TSE-R), CNRS, Université Toulouse Capitole, Toulouse, France. ⁴Institute for Data, Systems & Society, Massachusetts Institute of Technology, Cambridge, MA, USA. *e-mail: shariff@psych.ubc.ca; jean-francois.bonnefon@tse-fr.eu; irahwan@mit.edu

New Technology > Infrastructure & Transportation

Out of Two Million People, Most Prefer That a Self-Driving Car Kill the Elderly

Questions towards the MM data set

Data provenance

"The sample is self-selected... policymakers should not embrace this data as the final word on societal preferences"
(Etienne 2021).

Also no way of checking the country provenance

High-stake issues in low-stake environment

- Binary choices with "no information, no deliberation, and no consequences"
(Etienne 2021)
- captures stereotypes about categories of people, not reasoned principles (Luft 2020).

Hypotheticals aren't real policy-making

AV ethics "is precisely the wrong question" when framed as an individual's split-second choice (Jaques 2025)



Computational ethics fail us

(Awad 2022)



Descriptive click-data as normative moral truths

The fallacy of measuring attitudes as building moral rules (Etienne 2021).

Normalizes discriminatory logics under the guise of “optimization”

By forcing choices, MM operationalizes harmful social stereotypes (Luft 2020).

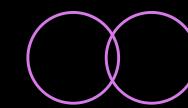
Narrows the entire ethics of AVs to a single artificial dilemma, ignoring real risks

AV ethics is about infrastructures, institutions, and power, not trolley hypotheticals (Pink 2022).

This is just not how conversations happen)

Law and social sanctions express and renegotiate shared moral norms as social differentiation increases (Durkheim 1893).

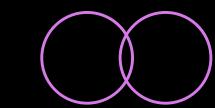
AV Stack Overview



01: Perception

Detects and classifies objects (pedestrians, vehicles, signs); Estimates positions, velocities, and uncertainties.

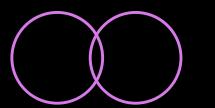
Example: Camera + LiDAR identify a pedestrian partially occluded by a parked van



02: Prediction

Forecasts future trajectories of road users; models uncertainty and multi-modal behaviors.

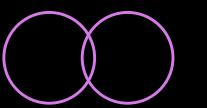
Example: Predicting whether a cyclist will continue straight, merge, or yield



03: Planning

Chooses the safest and most appropriate maneuver; integrates cultural/ethical preferences.

Example: Deciding whether to brake aggressively or yield gradually depending on local norms

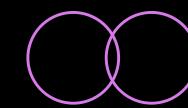


04: Control

Executes steering, throttle, braking to follow the chosen path; ensures smooth, stable actuation.

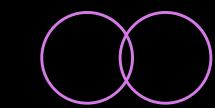
Example: Applying precise braking force to avoid wheel lock or passenger discomfort

AV Stack Overview



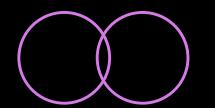
01: Perception

Detects and classifies objects (pedestrians, vehicles, signs); Estimates positions, velocities, and uncertainties.
Example: Camera + LiDAR identify a pedestrian partially occluded by a parked van



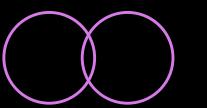
02: Prediction

Forecasts future trajectories of road users; models uncertainty and multi-modal behaviors.
Example: Predicting whether a cyclist will continue straight, merge, or yield



03: Planning

Chooses the safest and most appropriate maneuver; integrates cultural/ethical preferences.
Example: Deciding whether to brake aggressively or yield gradually depending on local norms



04: Control

Executes steering, throttle, braking to follow the chosen path; ensures smooth, stable actuation.
Example: Applying precise braking force to avoid wheel lock or passenger discomfort

Ethics mainly enter in planning/decision-making via:

- 1) **Reward function design** (soft preferences, tuned weights)
- 2) **Ethical/safety constraint module** (hard rules, veto layer)

Soft Ethics - Reward Function Structure

$$R(s, a) = - \left(C_{\text{human}} N_{\text{human harm}} + C_{\text{animal}} N_{\text{animal harm}} \right. \\ \left. + C_{\text{property}} D_{\text{property}} + C_{\text{law}} I_{\text{law viol}} + C_{\text{comfort}} A_{\text{jerk}} \right),$$

Sources of Cultural & Legal Variation:

1. Safety vs. legality vs. efficiency weighting: risk-averse vs. assertive driving cultures
2. Per-category treatment of humans: some regions allow distinctions (young/old, many/few); some (e.g., Germany) require uniform treatment
3. Hard constraints vs. soft penalties: non-discrimination, no "one for many" sacrifice → hard rule

Learning Cultural Preferences

- Handcrafted weights are brittle, subject to designer bias
- Cultural norms differ across regions and evolve over time
- **IRL:** infer reward from demonstrations
- **Preference learning:** infer utilities from pairwise judgments, e.g.

$U_\theta(\cdot)$ such that $U_\theta(x) > U_\theta(y)$

- Cultural heads on global moral backbone
- Combine with legal guardrails

Algorithm 1 Inverse Reinforcement Learning for Culture-Specific Reward Inference

Require: Demonstrations D_c (state-action trajectories) from culture c

- 0: Initialize reward parameters θ randomly
- 0: **repeat**
- 0: Compute optimal policy π_θ for reward $R(s, a; \theta)$
- 0: Compute feature expectations $\mu(\pi_\theta)$ and $\mu(D_c)$
- 0: Update parameters: $\theta \leftarrow \theta + \alpha(\mu(D_c) - \mu(\pi_\theta))$
- 0: **until** convergence

Ensure: θ_c (culture-specific reward parameters) =0

Case Studies



Legal Principles ('17 Ethics Commission, Autonomous Driving Act)

- Human dignity is inviolable → no discrimination by personal characteristics
- No numerical trade-offs → can't sacrifice one to save many
- Risk minimization only → braking, deceleration, evasive maneuvers are allowed; choosing whom to hit is not

Hard constraints + infinite penalties for human harm

Verification via attribute-swapping tests



Regulatory sandbox (C4IR)

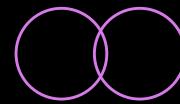
Process-oriented, not rule-prescriptive

Iterative policy adjustment (AV ethics shaped through observed behavior rather than fixed in advance)

Focus on aligning AV behavior with local traffic norms and public expectations

Co-development with industry and public stakeholders (Mobileye, GM, Cruise, Waymo)

Real-Time Decision-Making Challenges & Uncertainties

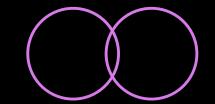


1. Computational Limits

Planner runs at 50–100 ms per cycle;
Rich ethical optimization too slow; (Monte Carlo: 10–50 ms per rollout)
Runtime uses precomputed policies + 5–10 sampled trajectories;

$$\mathbf{a}^* = \arg \max_{\mathbf{a} \in \mathcal{A}_{\text{sample}}} V_\theta(\mathbf{s}_t, \mathbf{a})$$

Region-specific policies loaded offline
(Germany vs. Israel)

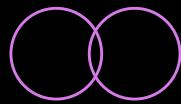


2. Layered Architecture

Nominal planner: 10–20 Hz, evaluates candidate trajectories
Safety layer: 100 Hz, monitors Time-to-Collision (TTC)

$$\text{TTC}_i = \frac{\text{distance}_i(\mathbf{s}_t)}{\max(0, \text{relative_velocity}_i(\mathbf{s}_t))}.$$

If $\text{TTC} < 0.5\text{--}1.0$ s → emergency braking override
Hard constraints must hold across both layers



3. Perception Uncertainty

Errors: misclassification, localization noise, occlusions, overconfidence
Planner uses probabilistic outputs, not hard labels
RS-NNs → belief masses over classes

```
TRIGGER_FALLBACK ← false
for each detected agent  $i$  do
    if  $P(\text{human}_i | \mathbf{z}_i) > p_{\text{threshold}}$  and  $\text{TTC}_i < t_{\text{threshold}}$  then
        TRIGGER_FALLBACK ← true
        break
    end if
end for=0
```

Fallback triggered if $P(\text{human}) > p_{\text{threshold}}$ AND $\text{TTC} < t_{\text{threshold}}$

Thresholds reflect cultural risk tolerance

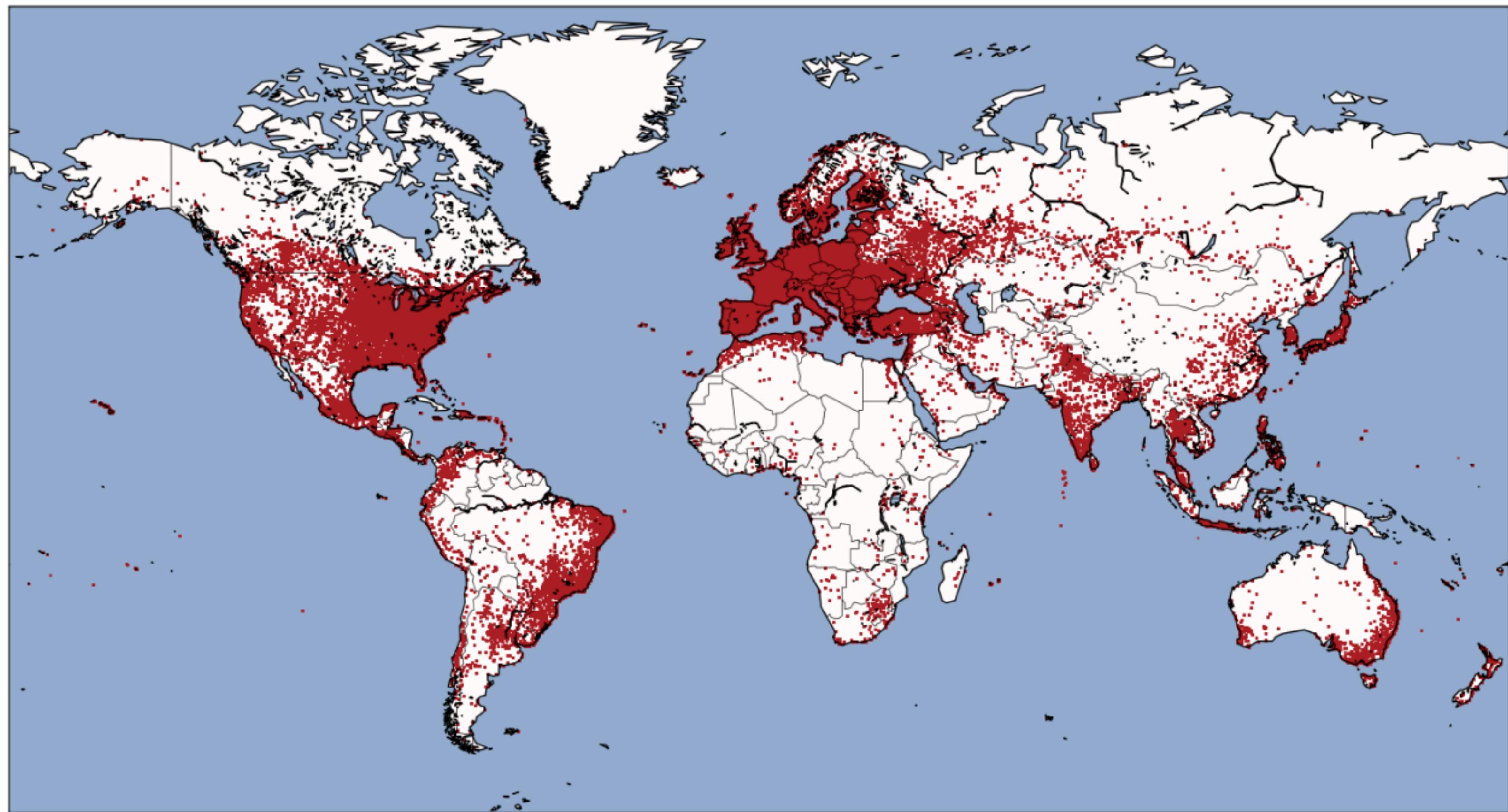
A RACE TO CONSTRAINT

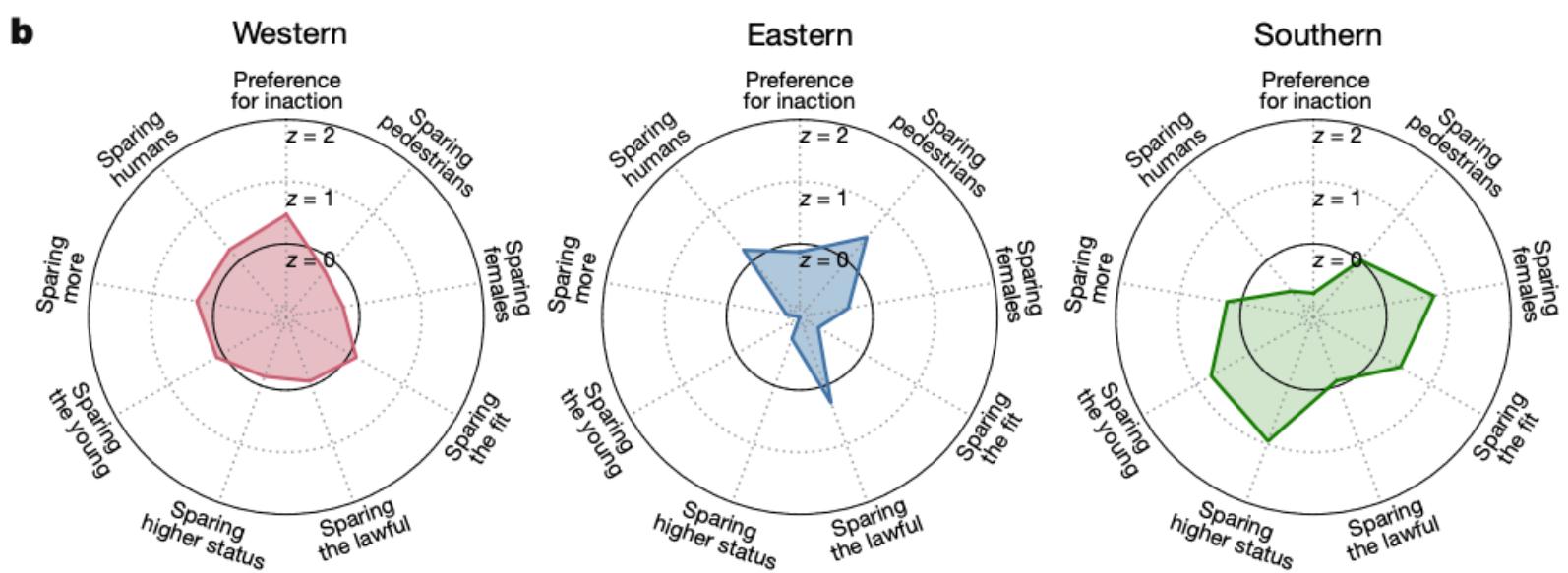
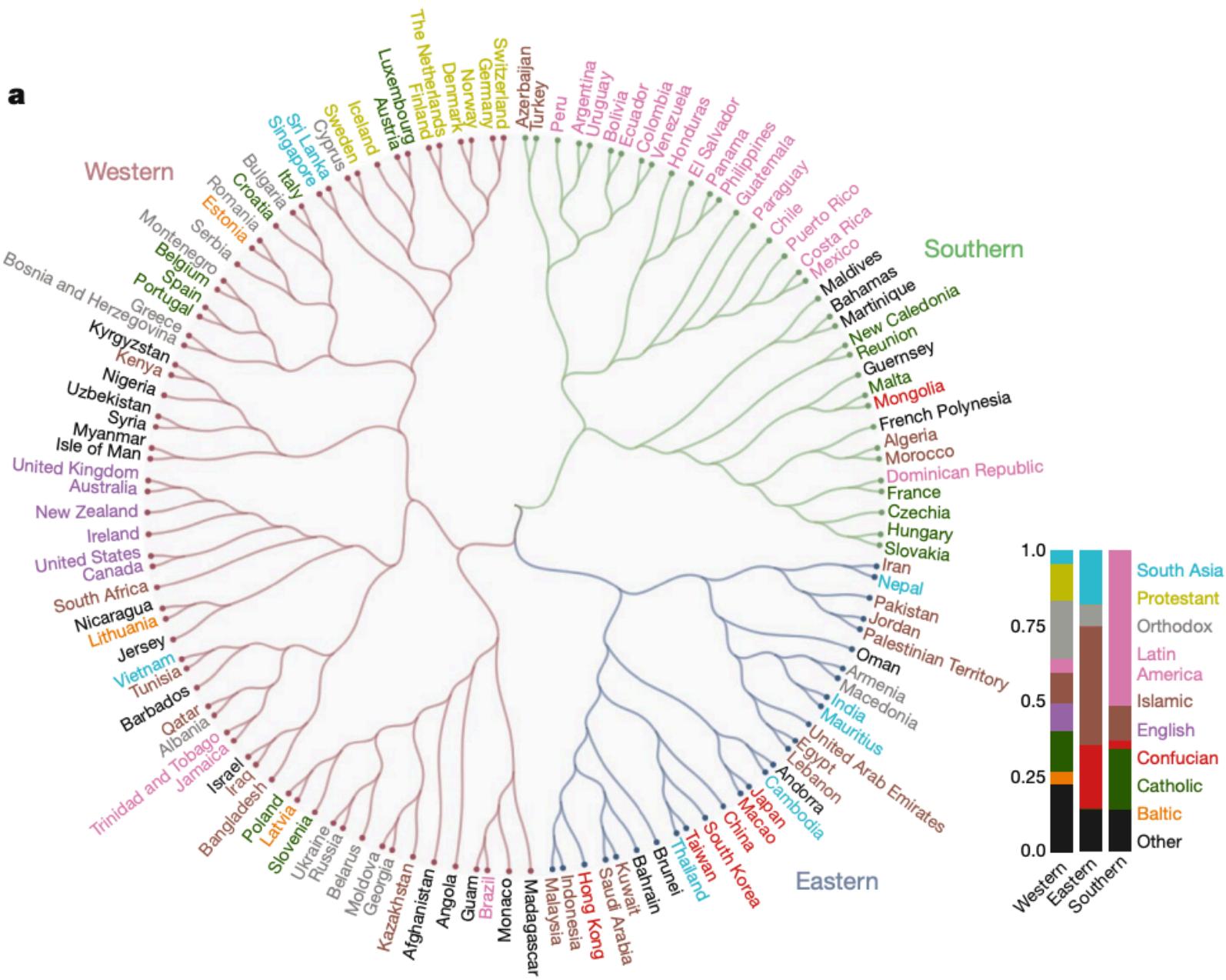
*EMBEDDING ETHICAL AND CULTURAL PREFERENCES IN
AUTONOMOUS VEHICLE PLANNING*



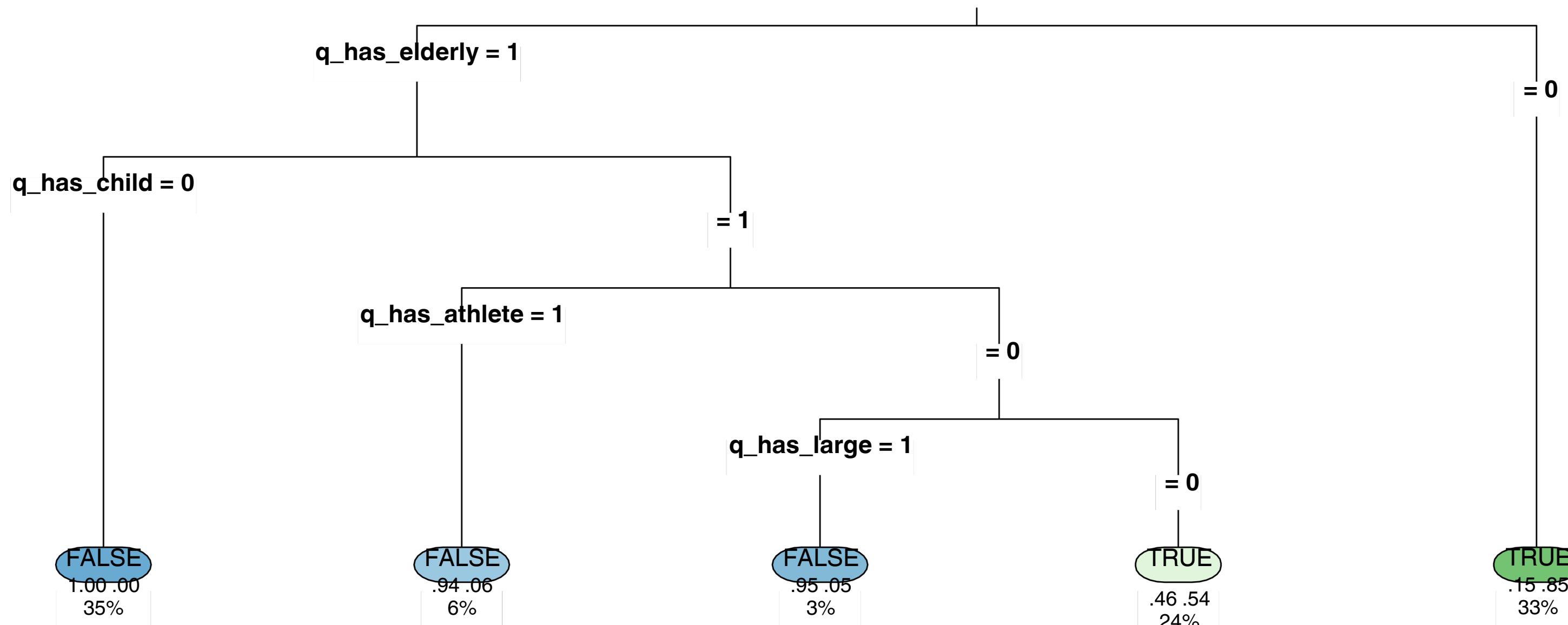
Sonya Rashkovan / Julia Zdzilowska

APPENDIX

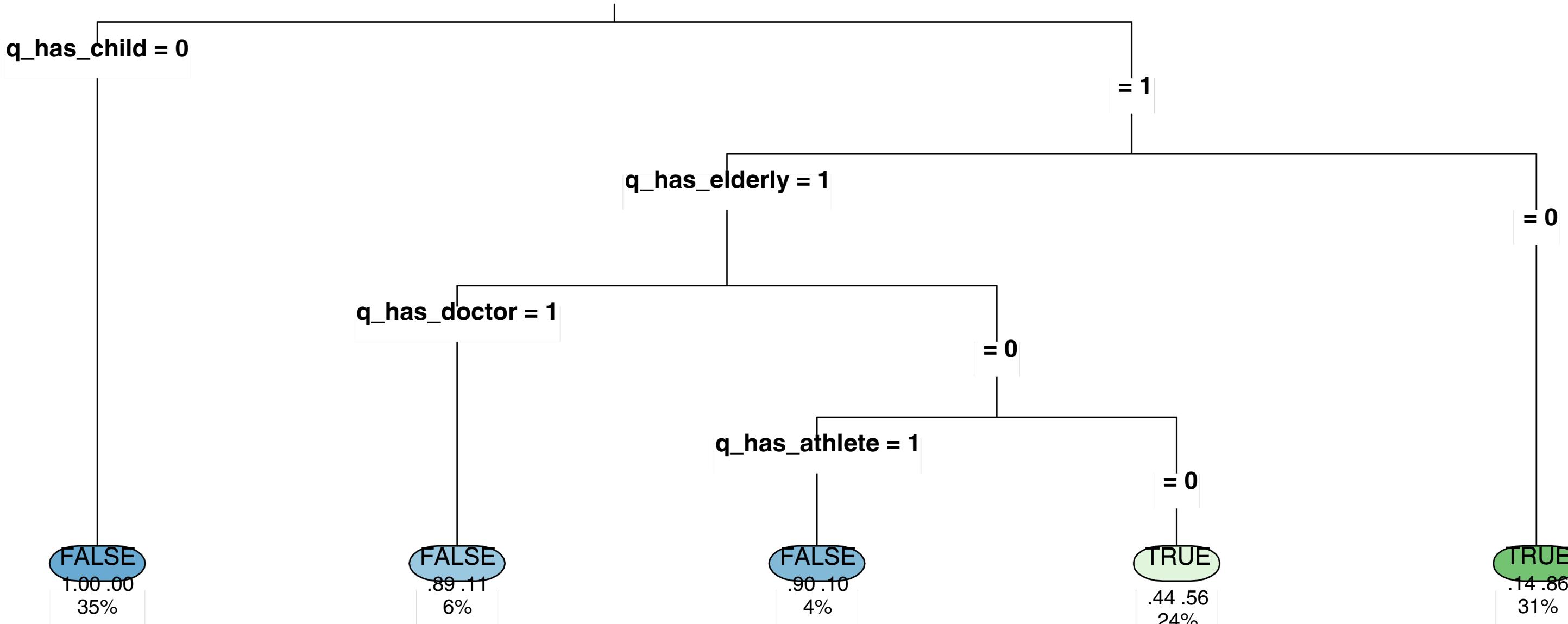
a



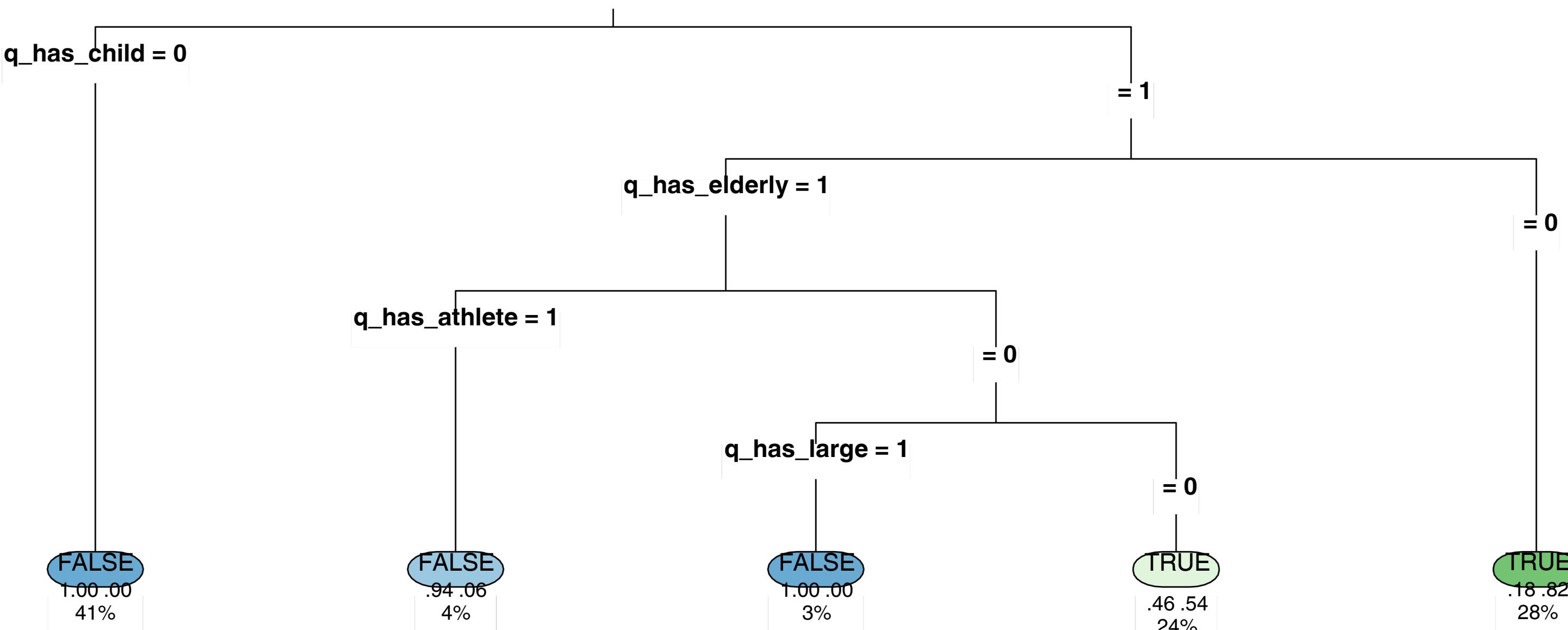
DEU AV Decision Tree: Choose Youth vs Elderly



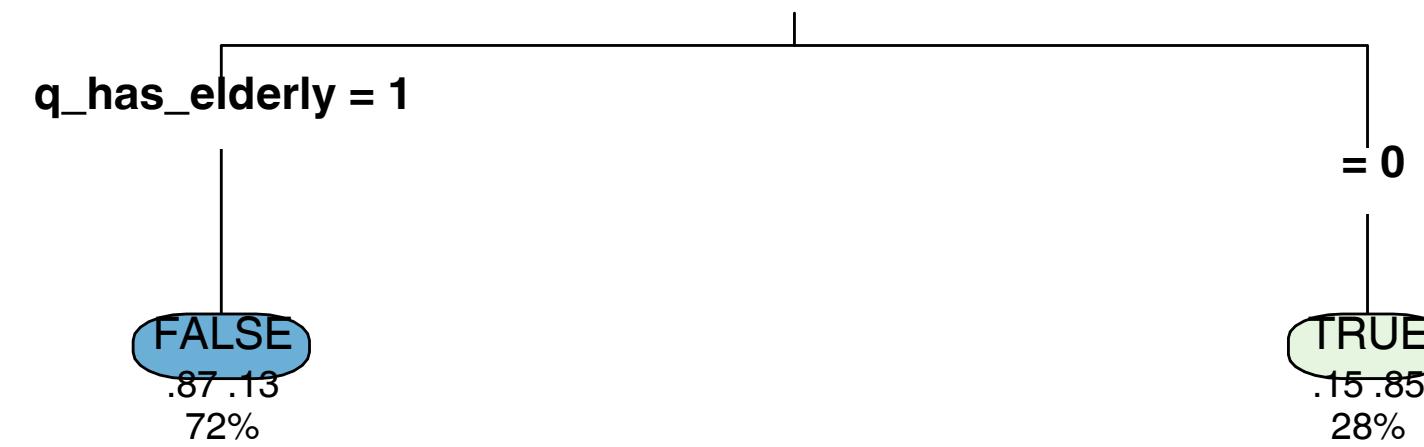
USA AV Decision Tree: Choose Youth vs Elderly



ISR AV Decision Tree: Choose Youth vs Elderly

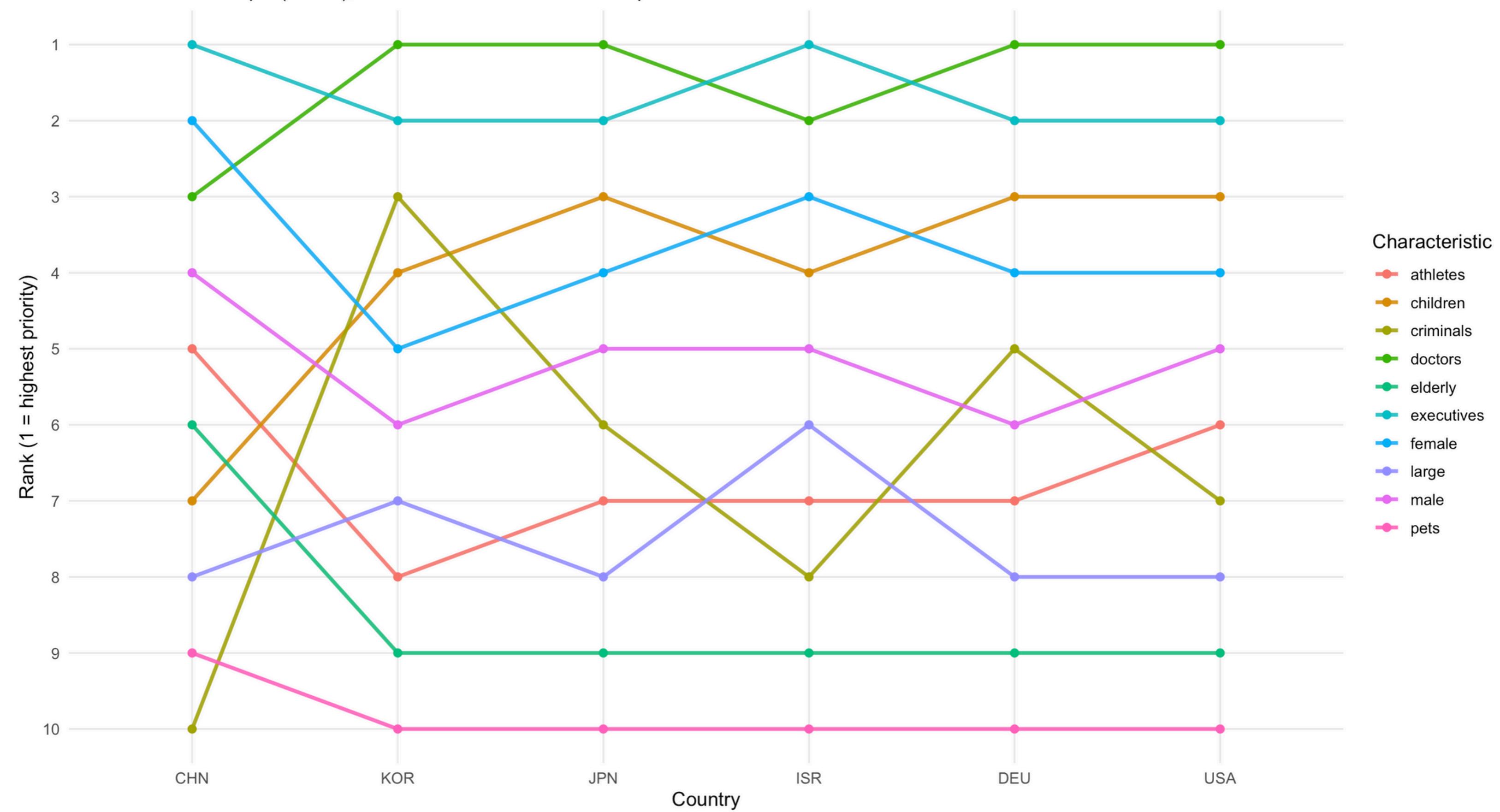


CHN AV Decision Tree: Choose Youth vs Elderly



Cross-Country Ranking of AV Priorities

*CHN has small sample (n=172); KOR & JPN are East Asian comparisons



Save Rates by Characteristic and Country

*CHN has small sample (n=172); KOR & JPN are East Asian comparisons

