

# How to ask without speech? On quantifying zero-evidence speech

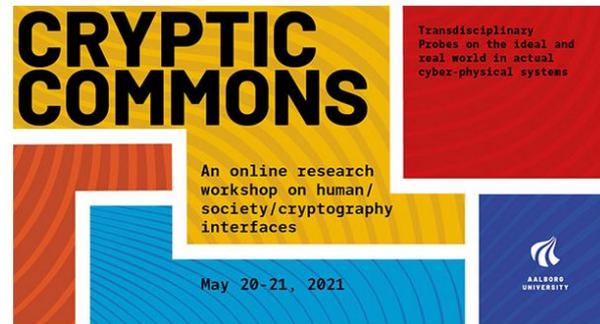
## *“The Privacy ZEBRA”*



Andreas Nautsch



2021-05-21



# Context: VoicePrivacy 2020 Challenge

- Task: audio pseudonymisation  $\Rightarrow$  modify raw audio
  - Voice biometrics should fail  
*“Same person or different person?”*
  - Speech recognition should work  
*“What was said?”*
- Metric: Zero-Evidence Biometric Recognition Assessment



THE UNIVERSITY  
of EDINBURGH



**Natalia  
Tomashenko**



**Brij M.L.  
Srivastava**



**Xin Wang**



**Emmanuel  
Vincent**



**Andreas  
Nautsch**



**Junichi  
Yamagishi**



**Nicholas  
Evans**



**Jose  
Patino**



**Jean-François  
Bonastre**



**Paul-Gauthier  
Noé**



**Massimiliano  
Todisco**

# Intuition: benefit to decision making?

- Motivation in forensic sciences
  - What is the benefit of *evidence reporting* to *decision making*?
  - How to validate?
- Empirical cross-entropy (ECE)
  - *Less uncertainty with evidence* than without?
- Strength-of-evidence: likelihood ratios
  - Which *decision* is *more supported*?

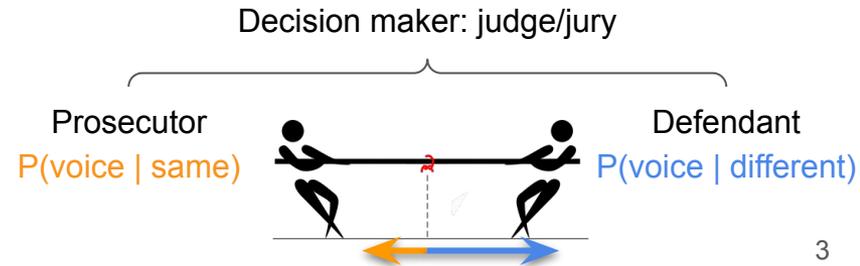
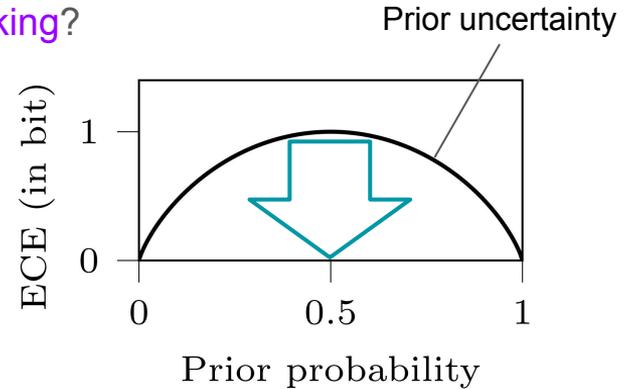


Figure based on wikimedia.org

# Textbook: empirical cross-entropy (ECE)

- The ECE step-by-step

- Prior entropy in making yes/no decision
- Posterior entropy based on scores/evidence
- !! Issue: no theoretical foundation for reference likelihoods
- Remedy: cross-entropy, law of large numbers & priors are external to the classifier

$\Theta = \{ A: \text{“same person”}, B: \text{“different person”} \}$

P: reference probability space  $H_P(\Theta) = - \sum_{\theta \in \Theta} P(\theta) \log_2 P(\theta)$

$H_P(\Theta | S) = - \sum_{\theta \in \Theta} P(\theta) \int_s P(s | \theta) \log_2 P(\theta | s) ds$   
 Set of scores

$H_{P||\tilde{P}}(\Theta | S) = - \sum_{\theta \in \Theta} P(\theta) \int_s P(s | \theta) \log_2 \tilde{P}(\theta | s) ds$

reference and classifier prior values  $P(\theta), \tilde{P}(\theta)$   
 $\pi = P(A) = \tilde{P}(A)$  and  $1 - \pi = P(B) = \tilde{P}(B)$

$P(s | \theta) \approx |S_\theta|^{-1}$

Strength-of-evidence

“the classifier”

score:  
 $P(\text{voice} | \text{same}) / P(\text{voice} | \text{different})$

$$ECE(\Theta | S) := \frac{\pi}{|S_A|} \sum_{a \in S_A} \log_2 \left( 1 + \frac{1 - \pi}{a \pi} \right) + \frac{1 - \pi}{|S_B|} \sum_{b \in S_B} \log_2 \left( 1 + \frac{b \pi}{1 - \pi} \right)$$

# Disclosure: worst-case?

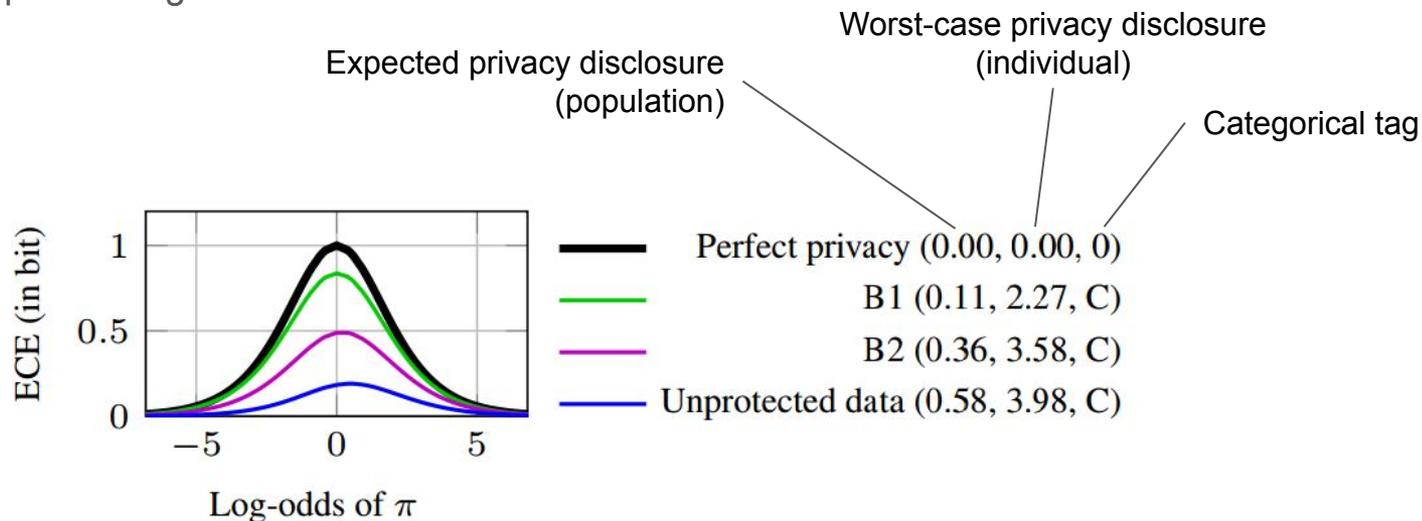
- Motivation: privacy for the individual; not for the average only
- Analogue from forensic sciences to privacy preservation
  - Prosecutor & defendant in a tug of war  $\Rightarrow$  i.e. *strength-of-evidence*
  - Decision maker: the adversary  $\Rightarrow$  i.e. what is the worst case?
- Categorical tags

Tag	Category	Posterior odds ratio (flat prior)
0	$l = 1 = 10^0$	50 : 50 (flat posterior)
A	$10^0 < l < 10^1$	more disclosure than 50 : 50
B	$10^1 \leq l < 10^2$	one wrong in 10 to 100
C	$10^2 \leq l < 10^4$	one wrong in 100 to 10 000
D	$10^4 \leq l < 10^5$	one wrong in 10 000 to 100 000
E	$10^5 \leq l < 10^6$	one wrong in 100 000 to 1 000 000
F	$10^6 \leq l$	one wrong in at least 1 000 000

Categorical scale of privacy disclosure  
(adapted from forensic sciences)

# ZEBRA framework, an example

- VoicePrivacy 2020 Challenge — audio pseudonymisation
  - Task: speech recognition should work — voice biometrics should fail
  - Unprotected data: state-of-the-art voice biometrics
  - B1: DNN baseline
  - B2: signal processing baseline



# The Privacy ZEBRA

