

CS208: Applied Privacy for Data Science

Membership Attacks

James Honaker & Salil Vadhan

School of Engineering & Applied Sciences
Harvard University

February 11, 2019



CRCS Center for Research on
Computation and Society

at Harvard John A. Paulson School of Engineering and Applied Sciences

Hypothesis Tests

A null hypothesis is a conjectured model of the world with observable implications.

Often it is a simplified model, for which there is some informational value if it can be refuted.

Null Distributions

If t is a function of the data, it has a sampling distribution. The distribution that t would obtain if the null hypothesis were true is called the *null distribution*.

- If we use the value of t to draw an inference about the null hypothesis, we call t a **test statistic**.
- We observe t^* in some observed dataset \mathbf{X}^* and reason whether it could have been a draw from the null distribution.
- If t^* is unlikely to have come from the null distribution, we **reject the null** hypothesis.
- If t^* could have been obtained from the null distribution, we **fail to reject the null**.
- Failing to reject the null, does not prove the null to be true.

Inferential Errors

Reasoning from known data to about an unknown hypothesis is called inference. Inferential errors are commonly labelled by type:

	Null True	Null False
Fail to Reject Null	Specificity	
Reject Null	Specificity	

Inferential Errors

Reasoning from known data to about an unknown hypothesis is called inference. Inferential errors are commonly labelled by type:

	Null True	Null False
Fail to Reject Null	Correct Specificity	
Reject Null		Correct Specificity

Inferential Errors

Reasoning from known data to about an unknown hypothesis is called inference. Inferential errors are commonly labelled by type:

	Null True	Null False
Fail to Reject Null	Correct Specificity	Error (Type II)
Reject Null	Error (Type I) Specificity	Correct

Inferential Errors

Reasoning from known data to about an unknown hypothesis is called inference. Inferential errors are commonly labelled by type:

	Null True	Null False
Fail to Reject Null	Correct Specificity	Error (Type II) Sensitivity
Reject Null	Error (Type I) Specificity	Correct Sensitivity

Inferential Errors

Reasoning from known data to about an unknown hypothesis is called inference. Inferential errors are commonly labelled by type:

	Null True	Null False
Fail to Reject Null	Correct Specificity $1 - \alpha$	Error (Type II) Sensitivity
Reject Null	Error (Type I) Specificity α	Correct Sensitivity

We parameterize our hypothesis test by choice of α which results in a **critical value**, c , which divides the null distribution into the rejection regions.

Example

H_0 : K -dimensional random variables \mathbf{x} and \mathbf{z} are both drawn from a standard Normal distribution with the same mean, $\mathcal{N}(\vec{\mu}, 1)$.

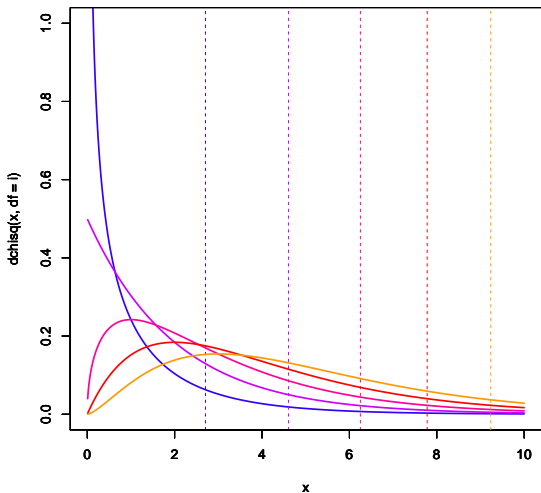
Then one test statistic is:

$$t(\mathbf{x}, \mathbf{z}) = \|\mathbf{x} - \mathbf{z}\|_2 = \sqrt{\sum_{i=1}^K (x_i - z_i)^2}$$

Which has null distribution $\chi^2(K)$.

Example

$\chi^2(K)$ Distribution with critical values for $\alpha = 0.1$:



Homer , Szelinger, Redman, Duggan, Tembe,
Muehling, Pearson, Stephan, Nelson, & Craig (2008)

Resolving individuals contributing trace amounts of DNA to highly
complex mixtures using high-density SNP genotyping microarrays.

Homer , Szelinger, Redman, Duggan, Tembe,
Muehling, Pearson, Stephan, Nelson, & Craig (2008)

Resolving individuals contributing trace amounts of DNA to highly
complex mixtures using high-density SNP genotyping microarrays.

Author Contributions

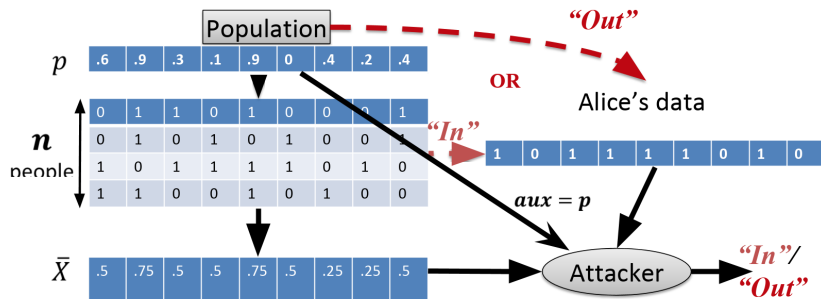
Conceived and designed the experiments: SFN DWC. Performed the
experiments: SS MR JM. Analyzed the data: NH WT DWC. Contributed
reagents/materials/analysis tools: DD JVP DS SFN DWC. Wrote the
paper: NH DWC.

Homer *et al.* (2008)

Resolving individuals contributing trace amounts of DNA to highly complex mixtures using high-density SNP genotyping microarrays.

- Membership attack on individual's inclusion in sample/dataset with published summary statistics (means).
- Membership can violate privacy if membership betrays an implicit variable.
- *"Their [Braun et al.] work showed high specificity for the test statistic of Homer et al., but with possibility of low sensitivity."* Bruce Weir's Viewpoint: Individual Genotyping in Forensics and GWAS Contexts

Dwork, Smith, Steinke, Ullman, Vadhan (2015)



$$A(y, a, p) = \begin{cases} \text{IN} & \text{if } \langle y, a \rangle - \langle p, a \rangle > T \\ \text{OUT} & \text{if } \langle y, a \rangle - \langle p, a \rangle \leq T \end{cases}$$

$$T = T_{p,a} = O\left(\sqrt{d \log(1/\delta)}\right)$$

$$p, a \in [-1, 1]; y \in \{-1, 1\}$$