

Tea: A High-level Language and Runtime System for Automating Statistical Analysis

Eunice Jun

with Maureen Daum, Jared Roesch, Sarah Chasins, Emery Berger, Rene Just, and Katharina Reinecke

Posted by u/WestEnd89 1 hour ago

Question [Q] What percentage of scores fall below below one standard deviation above the mean?

Hi all,

I'm having a bit of trouble working out answers to percentage questions about normal distributions.

Posted by u/ice_shadow 6 hours ago

Question [Q] Is classification in ML the opposite of ANOVA in classical stats?

So ANOVA and Mixed Models tell you whether a certain factor had a significant effect on the response and whether levels of a factor had a significantly different effect on the response.

From what I understand, things like Logistic Regression, discriminant analysis, kNN, SVM etc seem to use the response to try to predict the classes the data points belong to.

So are these approaches basically opposites of each other?

If the ANOVA contrasts is significant, would one of the classification approaches also be expected to perform well?

And if a classification approach has high accuracy, sensitivity, specificity then can you make the conclusion that there is a factor or well differentiated by the response?

Posted by u/WestEnd89 1 hour ago

Question [Q] What percentage of scores fall below below one standard deviation above the mean?

Hi all,

I'm having a bit of trouble working out answers to percentage questions about normal distributions.

Posted by u/ice_shadow 6 hours ago

Question [Q] Is classification in ML the opposite of ANOVA in classical stats?

So ANOVA and Mixed Models tell you whether a certain factor had a significant effect on the response and whether levels of a factor had a significantly different effect on the response.

From what I understand, things like Logistic Regression, discriminant analysis, k etc seem to use the response to try to predict the classes the data points belong

Posted by u/IzzyBee1 14 hours ago

Question [Q] What statistical test should I use?

Hi all, I'm looking for guidance from someone who knows what I'm sure don't! I have a semester of stats 101 under my belt (that I did in college) and that's more or less the extent of my knowledge.

Project background: my workplace has a moth problem. We have traps in the building and I check them once a month. Each trap location is in a different part of the building. I also have a combined data set with monthly moth catch of all locations for ~1.5 years' worth of data, but since I check monthly, it's a relatively small data set. Also, since I'm recording moth catch, the data is relatively skewed (I know I'm wrong) because many traps have caught 0 moths dur

Posted by u/banannah09 1 day ago

Question [Q] Can I do anything with this data?

Hello everyone! I've been reviewing some data for parents and children who received therapy. The way their mental health is measured is with 2 tests, so both parents and children should complete both of these tests before and after the treatment.

However... Even though 20 children received therapy, there are few cases where there is both pre and post treatment data (between 5-8 for both tests for both parents and children). I had many ideas for how I could analyse this data before, but now I'm not sure I can do anything with this aside from a few graphs (which I've already done)?

Posted by u/eddyks 19 hours ago

Question [Q] What is the best way to analyze my dataset?

Hi there,

Statistics is not my strongest point, so I was wondering if some of you could help me out a little bit.

So far I have almost 200 respondents who filled in my survey. They were given 6 sets, each having 3 statements, prior to being asked where they would buy a certain product (e.g. offline or online), which my moderator being the price of the product (high priced vs. low priced). Each set measures certain characteristics (e.g. price-conscious, time-conscious etc.). Now I want to test my hypotheses that price-conscious consumers buy high priced products rather than online.

What would be the best way to do this in SPSS?

Posted by u/WestEnd89 1 hour ago

Question [Q] What percentage of scores fall below below one standard deviation above the mean?

Hi all,

I'm having a bit of trouble with normal distributions.

Posted by u/ice_shadow 6

Question [Q] Is class ANOVA stats?

So ANOVA and Mixed Effects Models are the response and the response.

From what I understand, ANOVA and Mixed Effects Models etc seem to use the response variable.

Posted by u/IzzyBee1 14 hours ago

Question [Q] What statistical test should I use?

Hi all, I'm looking for guidance on what statistical test to use (I'm sure don't! I have a semester project (for a college) and that's more or less the whole project).

Project background: my building and I check the data every day. I also have a combined dataset of ~1.5 years' worth of data. Also, since I'm recording the data, I'm not sure if I'm wrong) because

Posted by u/banannah09 1 day ago

Question [Q] Can I do anything with this data?

When you have data from both parents and children, is there a statistical test that can be used to determine if the data is significantly different?

Does anyone know of any statistical tests that can be used to determine if the data is significantly different? (I've already done a t-test but now I'm not sure if that's the right test to use.)

1 set?

if you could help me out

They were given 6 sets, and they had to buy a certain product. The product (high priced rice-conscious, time-conscious consumers buy

Statistics is

hard

Which statistical test should I use?

Does an optimization make my program run faster?

H₁: Optimized code runs faster

H₀: Difference between run times due to chance

Which statistical test should I use?

Does an optimization make my program run faster?

Pearson's r	Welch's	Fisher's Exact
Pointbiserial	F-test	Linear regression
Kendall's T	Repeated measures	Logistic regression
Spearman's p	one-way ANOVA	MANOVA
Student's t-test	Factorial ANOVA	ANCOVA
Paired t-test	Two-way ANOVA	MANCOVA
Mann-Whitney U	Kruskal Wallis	McNemar
Wilcoxon signed rank	Friedman	Chi Square

H₁: Optimized code runs faster

H₀: Difference between run times due to chance

Which statistical test should I use?

Does an optimization make my program run faster?

Pearson's r

Welch's

Fisher's Exact

It depends!

Wilcoxon signed rank

Friedman

Chi Square

H_1 : Optimized code runs faster

H_0 : Difference between run times due to chance

Which statistical test should I use?

How do financial incentives affect users' performance?

H₁: Higher financial incentives, better user performance

H₀: Difference in performance due to chance

Which statistical test should I use?

How do financial incentives affect users' performance?

Pearson's r	Welch's	Fisher's Exact
Pointbiserial	F-test	Linear regression
Kendall's T	Repeated measures	Logistic regression
Spearman's p	one-way ANOVA	MANOVA
Student's t-test	Factorial ANOVA	ANCOVA
Paired t-test	Two-way ANOVA	MANCOVA
Mann-Whitney U	Kruskal Wallis	McNemar
Wilcoxon signed rank	Friedman	Chi Square

H₁: Higher financial incentives, better user performance

H₀: Difference in performance due to chance

Which statistical test should I use?

How do financial incentives affect users' performance?

Pearson's r

Welch's

Fisher's Exact

It depends!

Wilcoxon signed rank

Friedman

Chi Square

H_1 : Higher financial incentives, better user performance

H_0 : Difference in performance due to chance

Which statistical test should I use?

Does tea taste better with milk-then-tea or tea-then-milk?

H₁: Tea first tastes better

H₀: Difference in taste due to chance

Which statistical test should I use?

Does tea taste better with milk-then-tea or tea-then-milk?

Pearson's r	Welch's	Fisher's Exact
Pointbiserial	F-test	Linear regression
Kendall's T	Repeated measures	Logistic regression
Spearman's p	one-way ANOVA	MANOVA
Student's t-test	Factorial ANOVA	ANCOVA
Paired t-test	Two-way ANOVA	MANCOVA
Mann-Whitney U	Kruskal Wallis	McNemar
Wilcoxon signed rank	Friedman	Chi Square

H₁: Tea first tastes better

H₀: Difference in taste due to chance

Which statistical test should I use?

Does tea taste better with milk-then-tea or tea-then-milk?

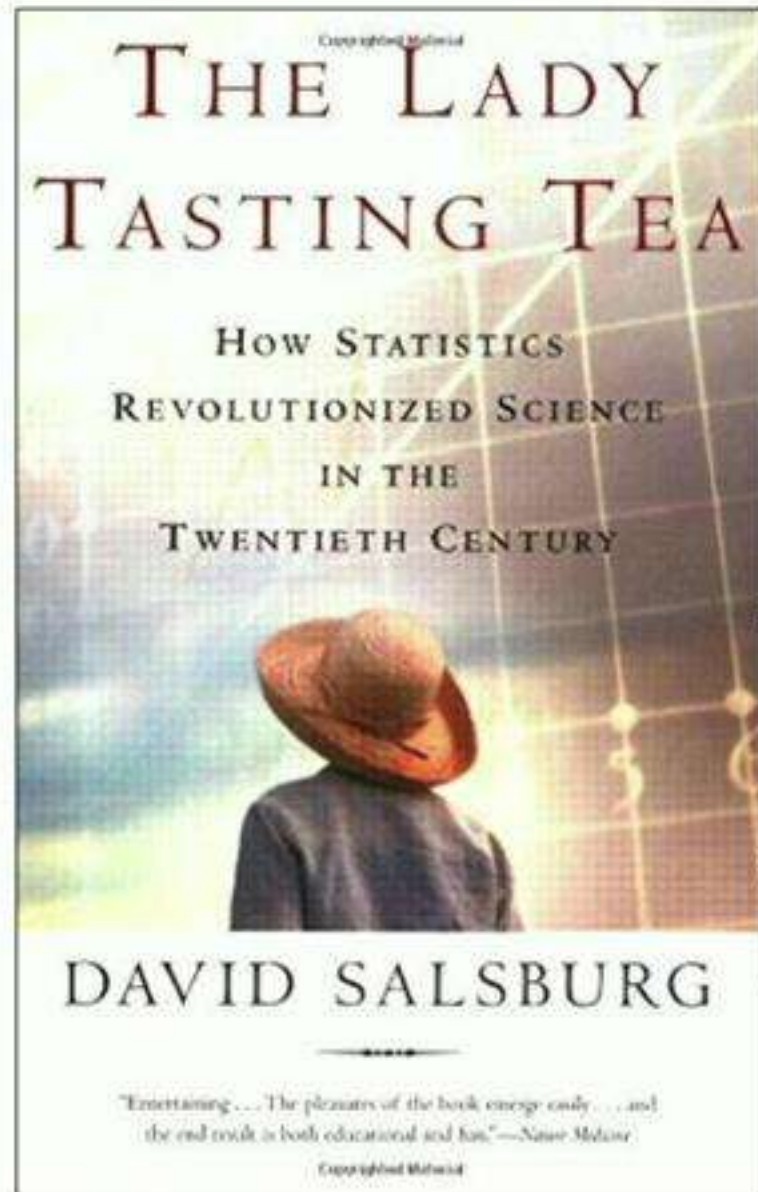
Pearson's r	Welch's	Fisher's Exact
Pointbiserial	F-test	Linear regression
Kendall's T	Repeated measures	Logistic regression
Spearman's p	one-way ANOVA	MANOVA
Student's t-test	Factorial ANOVA	ANCOVA
Paired t-test	Two-way ANOVA	MANCOVA
Mann-Whitney U	Kruskal Wallis	McNemar
Wilcoxon signed rank	Friedman	Chi Square

H_1 : Tea first tastes better

H_0 : Difference in taste due to chance

Which statistical test should I use?

Does tea taste better with milk-then-tea or tea-then-milk?



Fisher's Exact

Linear regression

Logistic regression

MANOVA

ANCOVA

MANCOVA

McNemar

Chi Square

H_1 : Tea first tastes better

H_0 : Difference in taste due to chance

tea

tea

EASY {

Does this optimization make my program execute faster?

How do financial incentives affect users' performance on a task?

Does tea taste better with milk poured first then tea or tea first then milk?

tea

EASY { Does this optimization make my program execute faster?
How do financial incentives affect users' performance on a task?
Does tea taste better with milk poured first then tea or tea first then milk?

HARD {

Pearson's r	Welch's	Fisher's Exact
Pointbiserial	F-test	Linear regression
Kendall's T	Repeated measures	Logistic regression
Spearman's p	one-way ANOVA	MANOVA
Student's t-test	Factorial ANOVA	ANCOVA
Paired t-test	Two-way ANOVA	MANCOVA
Mann-Whitney U	Kruskal Wallis	McNemar
Wilcoxon signed rank	Friedman	Chi Square

EASY



Does this optimization make my program execute faster?
How do financial incentives affect users' performance on a task?
Does tea taste better with milk poured first then tea or tea first then milk?

HARD



Pearson's r	Welch's	Fisher's Exact
Pointbiserial	F-test	Linear regression
Kendall's T	Repeated measures	Logistic regression
Spearman's p	one-way ANOVA	MANOVA
Student's t-test	Factorial ANOVA	ANCOVA
Paired t-test	Two-way ANOVA	MANCOVA
Mann-Whitney U	Kruskal Wallis	McNemar
Wilcoxon signed rank	Friedman	Chi Square

EASY	Does this optimization make my program execute faster?		
	How do financial incentives affect users' performance on a task?		
	Does tea taste better with milk poured first then tea or tea first then milk?		
HARD	Pearson's r	Welch's	Fisher's Exact
	Pointbiserial	F-test	Linear regression
	Kendall's T	Repeated measures	Logistic regression
	Spearman's p	one-way ANOVA	MANOVA
	Student's t-test	Factorial ANOVA	ANCOVA
	Paired t-test	Two-way ANOVA	MANCOVA
	Mann-Whitney U	Kruskal Wallis	McNemar
	Wilcoxon signed rank	Friedman	Chi Square

```
t.test(x y = NULL
       alternative = c("two.sided" "less" "greater")
       mu = 0 paired = FALSE var.equal = FALSE
       conf.level = 0.95 ...)
```

Tea eliminates this problem entirely.

EASY	Does this optimization make my program execute faster?		
	How do financial incentives affect users' performance on a task?		
	Does tea taste better with milk poured first then tea or tea first then milk?		
HARD	Pearson's r	Welch's	Fisher's Exact
	Pointbiserial	F-test	Linear regression
	Kendall's T	Repeated measures	Logistic regression
	Spearman's p	one-way ANOVA	MANOVA
	Student's t-test	Factorial ANOVA	ANCOVA
	Paired t-test	Two-way ANOVA	MANCOVA
	Mann-Whitney U	Kruskal Wallis	McNemar
	Wilcoxon signed rank	Friedman	Chi Square

```
t.test(x y = NULL
       alternative = c("two.sided" "less" "greater")
       mu = 0 paired = FALSE var.equal = FALSE
       conf.level = 0.95 ...)
```

Difference between Student's t-test and Paired t-test

Tea eliminates this problem entirely.

EASY	Does this optimization make my program execute faster?		
	How do financial incentives affect users' performance on a task?		
	Does tea taste better with milk poured first then tea or tea first then milk?		
HARD	Pearson's r	Welch's	Fisher's Exact
	Pointbiserial	F-test	Linear regression
	Kendall's T	Repeated measures	Logistic regression
	Spearman's p	one-way ANOVA	MANOVA
	Student's t-test	Factorial ANOVA	ANCOVA
	Paired t-test	Two-way ANOVA	MANCOVA
	Mann-Whitney U	Kruskal Wallis	McNemar
	Wilcoxon signed rank	Friedman	Chi Square

```
t.test(x y = NULL
      alternative = c("two.sided" "less" "greater")
      mu = 0 paired = FALSE var.equal = FALSE
      conf.level = 0.95 ...)
```

Difference between Student's t-test and Paired t-test

Each participant contributes **exactly one** data point

Tea eliminates this problem entirely.

EASY	Does this optimization make my program execute faster?		
	How do financial incentives affect users' performance on a task?		
	Does tea taste better with milk poured first then tea or tea first then milk?		
HARD	Pearson's r	Welch's	Fisher's Exact
	Pointbiserial	F-test	Linear regression
	Kendall's T	Repeated measures	Logistic regression
	Spearman's p	one-way ANOVA	MANOVA
	Student's t-test	Factorial ANOVA	ANCOVA
	Paired t-test	Two-way ANOVA	MANCOVA
	Mann-Whitney U	Kruskal Wallis	McNemar
	Wilcoxon signed rank	Friedman	Chi Square

```
t.test(x y = NULL
      alternative = c("two.sided" "less" "greater")
      mu = 0 paired = FALSE var.equal = FALSE
      conf.level = 0.95 ...)
```

Difference between Student's t-test and Paired t-test

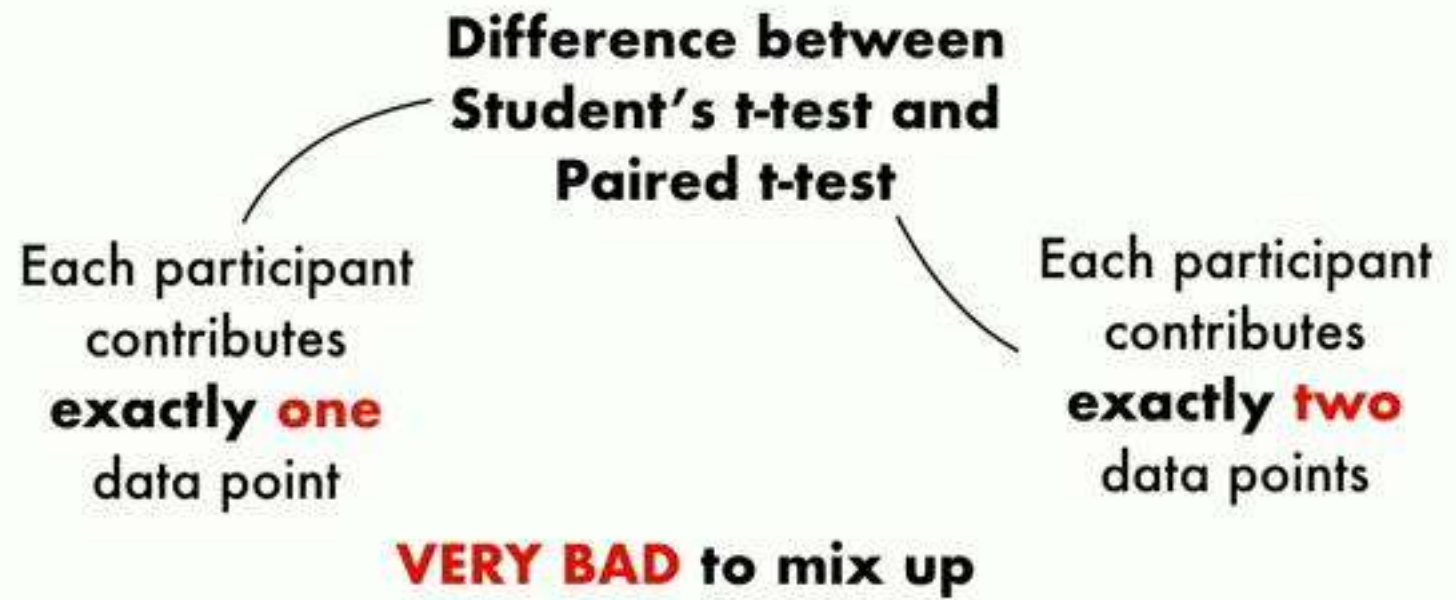
Each participant contributes **exactly one** data point

Each participant contributes **exactly two** data points

Tea eliminates this problem entirely.

EASY	Does this optimization make my program execute faster?		
	How do financial incentives affect users' performance on a task?		
	Does tea taste better with milk poured first then tea or tea first then milk?		
HARD	Pearson's r	Welch's	Fisher's Exact
	Pointbiserial	F-test	Linear regression
	Kendall's T	Repeated measures	Logistic regression
	Spearman's p	one-way ANOVA	MANOVA
	Student's t-test	Factorial ANOVA	ANCOVA
	Paired t-test	Two-way ANOVA	MANCOVA
	Mann-Whitney U	Kruskal Wallis	McNemar
	Wilcoxon signed rank	Friedman	Chi Square

```
t.test(x y = NULL
       alternative = c("two.sided" "less" "greater")
       mu = 0 paired = FALSE var.equal = FALSE
       conf.level = 0.95 ...)
```



Tea eliminates this problem entirely.

EASY	Does this optimization make my program execute faster?		
	How do financial incentives affect users' performance on a task?		
	Does tea taste better with milk poured first then tea or tea first then milk?		
HARD	Pearson's r	Welch's	Fisher's Exact
	Pointbiserial	F-test	Linear regression
	Kendall's T	Repeated measures	Logistic regression
	Spearman's p	one-way ANOVA	MANOVA
	Student's t-test	Factorial ANOVA	ANCOVA
	Paired t-test	Two-way ANOVA	MANCOVA
	Mann-Whitney U	Kruskal Wallis	McNemar
	Wilcoxon signed rank	Friedman	Chi Square

```
t.test(x y = NULL
      alternative = c("two.sided" "less" "greater")
      mu = 0 paired = FALSE var.equal = FALSE
      conf.level = 0.95 ...)
```

Difference between Student's t-test and Paired t-test

Each participant contributes **exactly one** data point

Each participant contributes **exactly two** data points

VERY BAD to mix up
Violate study design

Tea eliminates this problem entirely.

EASY	Does this optimization make my program execute faster? How do financial incentives affect users' performance on a task? Does tea taste better with milk poured first then tea or tea first then milk?		
HARD	Pearson's r	Welch's	Fisher's Exact
	Pointbiserial	F-test	Linear regression
	Kendall's T	Repeated measures	Logistic regression
	Spearman's p	one-way ANOVA	MANOVA
	Student's t-test	Factorial ANOVA	ANCOVA
	Paired t-test	Two-way ANOVA	MANCOVA
	Mann-Whitney U	Kruskal Wallis	McNemar
Wilcoxon signed rank	Friedman	Chi Square	

```
t.test(x y = NULL
      alternative = c("two.sided" "less" "greater")
      mu = 0 paired = FALSE var.equal = FALSE
      conf.level = 0.95 ...)
```

Difference between Student's t-test and Paired t-test

Each participant contributes **exactly one** data point

Each participant contributes **exactly two** data points

VERY BAD to mix up
Violate study design

 **Returns wrong statistical results!**

Tea eliminates this problem entirely.

Stats is better with Tea



Tea is correct by construction.

Tea is high-level.

Tea infers tests.

Tea *improves upon expert choices, prevents common mistakes.*

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3

Run Code

```
In [ ]: import tea
```

```
In [ ]: # Load data
tea.data("../datasets/UScrime.csv")
```

```
In [ ]: # Declare and annotate the variables of interest
variables = [
    {
        'name' : 'So',
        'data type' : 'nominal',
        'categories' : ['0', '1']
    },
    {
        'name' : 'Prob',
        'data type' : 'ratio',
        'range' : [0,1]
    }
]

tea.define_variables(variables)
```

```
In [ ]: assumptions = {
    'groups normally distributed': [['So', 'Prob']],
    'Type I (False Positive) Error Rate': 0.05,
}
```

```
import tea
tea.data('UScrime.csv')
variables = [
  {
    'name' : 'So',
    'data type' : 'nominal',
    'categories' : ['0', '1']
  },
  {
    'name' : 'Prob',
    'data type' : 'ratio',
    'range' : [0,1]
  }
]
tea.define_variables(variables)
```

data

variables

```
study_design = {
  'study type': 'observational study',
  'contributor variables': 'So',
  'outcome variables': 'Prob',
}
tea.define_study_design(study_design)
```

study design

```
assumptions = {
  'groups normally distributed': [['So', 'Prob']],
  'Type I (False Positive) Error Rate': 0.05
}
tea.assume(assumptions)
```

assertions

```
hypothesis = 'So:1 > 0'
tea.hypothesize(['So', 'Prob'], hypothesis)
```

hypothesis

**** NO STATISTICAL TEST ****

```
variables = [  
    {  
        'name' : 'So',  
        'data type' : 'nominal',  
        'categories' : ['0', '1']  
    },  
    {  
        'name' : 'Prob',  
        'data type' : 'ratio',  
        'range' : [0,1]  
    }  
]  
tea.define_variables(variables)
```

```
variables = [  
    {  
        → 'name' : 'So',  
          'data type' : 'nominal',  
          'categories' : ['0', '1']  
    },  
    {  
        → 'name' : 'Prob',  
          'data type' : 'ratio',  
          'range' : [0,1]  
    }  
]  
tea.define_variables(variables)
```

```
variables = [  
    {  
        'name' : 'So',  
        → 'data type' : 'nominal',  
        'categories' : ['0', '1']  
    },  
    {  
        'name' : 'Prob',  
        → 'data type' : 'ratio',  
        'range' : [0,1]  
    }  
]  
tea.define_variables(variables)
```

```
variables = [  
    {  
        Nominal      'name' : 'So',  
        Ordinal    → 'data type' : 'nominal',  
                    'categories' : ['0', '1']  
    },  
    {  
        Ratio      'name' : 'Prob',  
        → 'data type' : 'ratio',  
          'range' : [0,1]  
    }  
]  
tea.define_variables(variables)
```



```
variables = [  
  {  
    Nominal 🍦 'name' : 'So',  
    Ordinal 🎒 → 'data type' : 'nominal',  
    'categories' : ['0', '1']  
  },  
  {  
    Interval 🎯,  
    Ratio 🕒 → 'name' : 'Prob',  
    'data type' : 'ratio',  
    'range' : [0,1]  
  }  
]  
tea.define_variables(variables)
```

```
variables = [  
    {  
        'name' : 'So',  
        'data type' : 'nominal',  
        → 'categories' : ['0', '1']  
    },  
    {  
        'name' : 'Prob',  
        'data type' : 'ratio',  
        → 'range' : [0,1]  
    }  
]  
tea.define_variables(variables)
```

```
study_design = {  
    → 'study type': 'observational study',  
      'contributor variables': 'So',  
      'outcome variables': 'Prob',  
    }  
tea.define_study_design(study_design)
```

```
study_design = {  
    'study type': 'observational study',  
    → 'contributor variables': 'So',  
    'outcome variables': 'Prob',  
}  
tea.define_study_design(study_design)
```

```
study_design = {  
    'study type': 'observational study',  
    'contributor variables': 'So',  
    → 'outcome variables': 'Prob',  
}  
tea.define_study_design(study_design)
```

```
assumptions = {  
    'groups normally distributed': [['So', 'Prob']],  
    'Type I (False Positive) Error Rate': 0.05  
}  
tea.assume(assumptions)
```

```
hypothesis = 'So:1 > 0'  
tea.hypothesize(['So', 'Prob'], hypothesis)
```

```
hypothesis = 'So:1 > 0'  
tea.hypothesize(['So', 'Prob'], hypothesis)
```

Nominal, Ordinal:

Chocolate > Mint

Grade 1 < Grade 2

Ordinal, Ratio, Interval:

Grade ~ Temperature

Time of day ~ - Temperature


```
import tea
tea.data('UScrime.csv')
variables = [
    {
        'name' : 'So',
        'data type' : 'nominal',
        'categories' : ['0', '1']
    },
    {
        'name' : 'Prob',
        'data type' : 'ratio',
        'range' : [0,1]
    }
]
tea.define_variables(variables)
```

```
study_design = {
    'study type': 'observational study',
    'contributor variables': 'So',
    'outcome variables': 'Prob',
}
tea.define_study_design(study_design)
```

```
assumptions = {
    'groups normally distributed': [['So', 'Prob']],
    'Type I (False Positive) Error Rate': 0.05
}
tea.assume(assumptions)
```

```
hypothesis = 'So:1 > 0'
tea.hypothesize(['So', 'Prob'], hypothesis)
```

data

variables

study design

assertions

hypothesis

Statistical test selection as constraint satisfaction

```
import tea
tea.data('UScrime.csv')
variables = [
  {
    'name' : 'So',
    'data type' : 'nominal',
    'categories' : ['0', '1']
  },
  {
    'name' : 'Prob',
    'data type' : 'ratio',
    'range' : [0,1]
  }
]
tea.define_variables(variables)

study_design = {
  'study type': 'observational study',
  'contributor variables': 'So',
  'outcome variables': 'Prob',
}
tea.define_study_design(study_design)

assumptions = {
  'groups normally distributed': [['So', 'Prob']],
  'Type I (False Positive) Error Rate': 0.05
}
tea.assume(assumptions)

hypothesis = 'So:1 > 0'
tea.hypothesize(['So', 'Prob'], hypothesis)
```

✓ completeness

✓ syntax

✓ well-formed hypotheses

Statistical test selection as constraint satisfaction

```
import tea
tea.data('UScrime.csv')
variables = [
  {
    'name' : 'So',
    'data type' : 'nominal',
    'categories' : ['0', '1']
  },
  {
    'name' : 'Prob',
    'data type' : 'ratio',
    'range' : [0,1]
  }
]
tea.define_variables(variables)

study_design = {
  'study type': 'observational study',
  'contributor variables': 'So',
  'outcome variables': 'Prob',
}
tea.define_study_design(study_design)

assumptions = {
  'groups normally distributed': [['So', 'Prob']],
  'Type I (False Positive) Error Rate': 0.05
}
tea.assume(assumptions)

hypothesis = 'So:1 > 0'
tea.hypothesize(['So', 'Prob'], hypothesis)
```



logical constraints

$continuous(x) \wedge \neg categorical(x)$

$normal(x) \rightarrow \neg categorical(x)$

...

- ✓ completeness
- ✓ syntax
- ✓ well-formed hypotheses

Statistical test selection as constraint satisfaction

```
import tea
tea.data('UScrime.csv')
variables = [
  {
    'name' : 'So',
    'data type' : 'nominal',
    'categories' : ['0', '1']
  },
  {
    'name' : 'Prob',
    'data type' : 'ratio',
    'range' : [0,1]
  }
]
tea.define_variables(variables)
study_design = {
  'study type': 'observational study',
  'contributor variables': 'So',
  'outcome variables': 'Prob',
}
tea.define_study_design(study_design)
assumptions = {
  'groups normally distributed': [['So', 'Prob']],
  'Type I (False Positive) Error Rate': 0.05
}
tea.assume(assumptions)
hypothesis = 'So:1 > 0'
tea.hypothesize(['So', 'Prob'], hypothesis)
```



logical constraints

$continuous(x) \wedge \neg categorical(x)$

$normal(x) \rightarrow \neg categorical(x)$

...



MaxSat

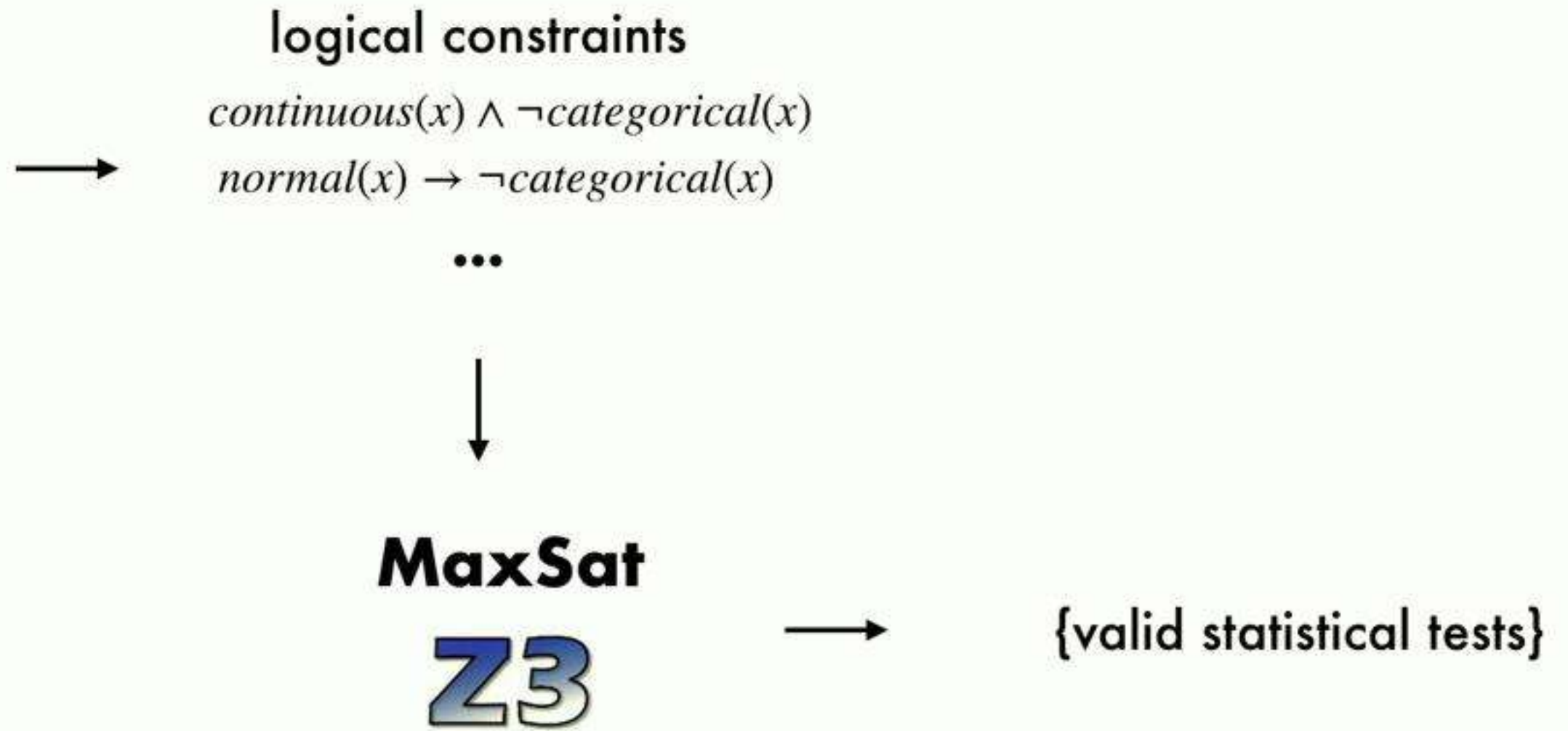
Z3

- ✓ completeness
- ✓ syntax
- ✓ well-formed hypotheses

Statistical test selection as constraint satisfaction

```
import tea
tea.data('UScrime.csv')
variables = [
  {
    'name' : 'So',
    'data type' : 'nominal',
    'categories' : ['0', '1']
  },
  {
    'name' : 'Prob',
    'data type' : 'ratio',
    'range' : [0,1]
  }
]
tea.define_variables(variables)
study_design = {
  'study type': 'observational study',
  'contributor variables': 'So',
  'outcome variables': 'Prob',
}
tea.define_study_design(study_design)
assumptions = {
  'groups normally distributed': [['So', 'Prob']],
  'Type I (False Positive) Error Rate': 0.05
}
tea.assume(assumptions)
hypothesis = 'So:1 > 0'
tea.hypothesize(['So', 'Prob'], hypothesis)
```

- ✓ completeness
- ✓ syntax
- ✓ well-formed hypotheses



Statistical test selection as constraint satisfaction

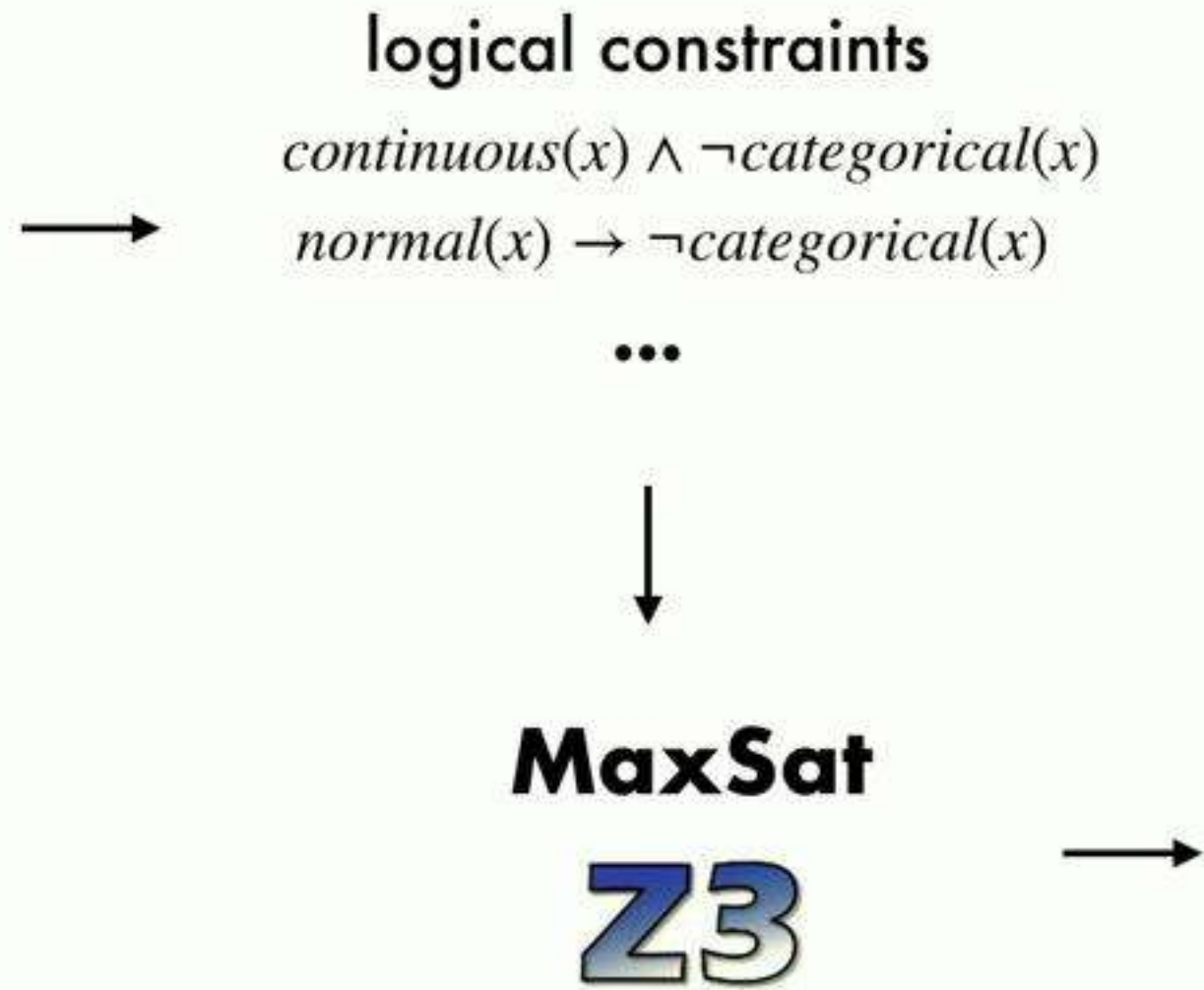
```
import tea
tea.data('UScrime.csv')
variables = [
  {
    'name' : 'So',
    'data type' : 'nominal',
    'categories' : ['0', '1']
  },
  {
    'name' : 'Prob',
    'data type' : 'ratio',
    'range' : [0,1]
  }
]
tea.define_variables(variables)

study_design = {
  'study type': 'observational study',
  'contributor variables': 'So',
  'outcome variables': 'Prob',
}
tea.define_study_design(study_design)

assumptions = {
  'groups normally distributed': [['So', 'Prob']],
  'Type I (False Positive) Error Rate': 0.05
}
tea.assume(assumptions)

hypothesis = 'So:1 > 0'
tea.hypothesize(['So', 'Prob'], hypothesis)
```

- ✓ completeness
- ✓ syntax
- ✓ well-formed hypotheses



```
Test: students_1
***Test assumptions:
Exactly two variables involved in analysis: So, Prob
Exactly one explanatory variable: So
Exactly one explained variable: Prob
Independent (not paired) observations: So
Variable is categorical: So
Variable has two categories: So
Continuous (not categorical) data: Prob
Equal variance: So, Prob
Groups are normally distributed: So, Prob

***Test results:
name = Student's T Test
test_statistic = 4.202130736875173
p_value = 0.00012364897286532775
adjusted_p_value = 6.182448633266387e-05
alpha = 0.05
df = 45
Effect size:
Cohen's d = 1.2426167296374897
A12 = 0.8366935483870968
Null hypothesis = There is no difference in means between 0 and 1 on Prob.
Interpretation = t(45) = 4.202130736875173, 6.182448633266387e-05. Reject the null hypothesis at alpha = 0.05. The mean of Prob for So = 1 is significantly greater than the mean for So = 0. The effect size is ('Cohen's d': 1.2426167296374897, 'A12': 0.8366935483870968). The effect size is the magnitude of the difference, which gives a holistic view of the results [1].
[1] Sullivan, G. M., & Feinn, R. (2012). Using effect size—or why the P value is not enough. Journal of Graduate Medical Education, 4(3), 279-282.
```

↑

{valid statistical tests}

Statistical test selection as constraint satisfaction

```

import tea
tea.data('UScrime.csv')

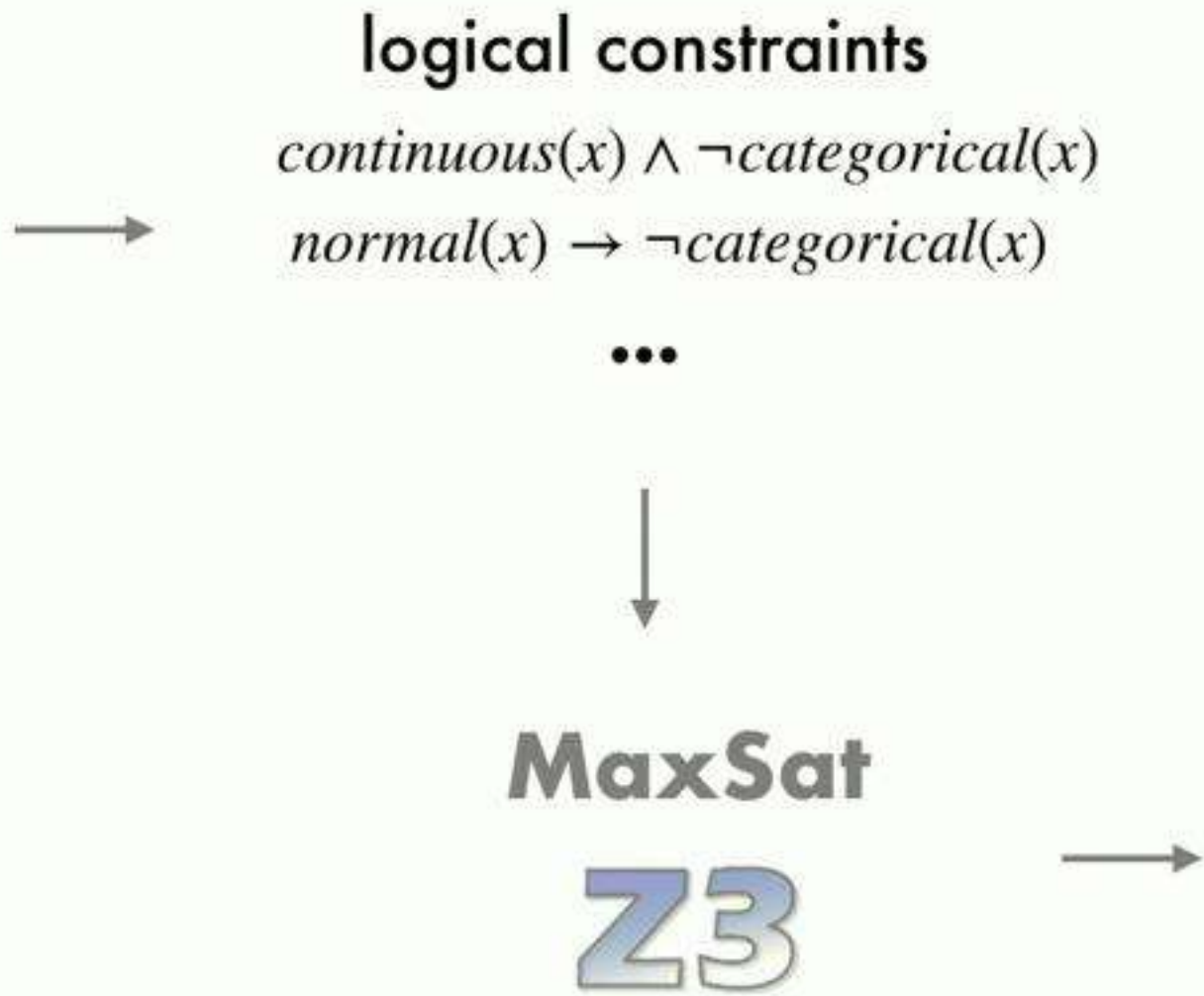
variables = [
    {
        'name' : 'So',
        'data type' : 'nominal',
        'categories' : ['0', '1']
    },
    {
        'name' : 'Prob',
        'data type' : 'ratio',
        'range' : [0,1]
    }
]
tea.define_variables(variables)

study_design = {
    'study type': 'observational study',
    'contributor variables': 'So',
    'outcome variables': 'Prob',
}
tea.define_study_design(study_design)

assumptions = {
    'groups normally distributed': [['So', 'Prob']],
    'Type I (False Positive) Error Rate': 0.05
}
tea.assume(assumptions)

hypothesis = 'So:1 > 0'
tea.hypothesize(['So', 'Prob'], hypothesis)
    
```

- ✓ completeness
- ✓ syntax
- ✓ well-formed hypotheses



```

Test: students.t
---Test assumptions:
Exactly two variables involved in analysis: So, Prob
Exactly one explanatory variable: So
Exactly one explained variable: Prob
Independent (not cased) observations: So
Variable is categorical: So
Variable has two categories: So
Continuous (not categorical) data: Prob
Equal variance: So, Prob
Groups are normally distributed: So, Prob

---Test results:
name = Student's T-Test
test_statistic = 4.202130736670173
p_value = 0.0001236489726652779
adjusted_p_value = 6.18244863206037e-05
alpha = 0.05
df = 45
Effect size:
Cohen's d = 1.2420167266374897
A12 = 0.836603548370968
Null hypothesis = There is no difference in means between 0 and 1 on Prob.
Interpretation = t(45) = 4.202130736670173, 6.18244863206037e-05. Reject the null hypothesis at alpha = 0.05. The mean of Prob for So = 1 is significantly greater than the mean for So = 0. The effect size is (Cohen's d): 1.2420167266374897, A12: 0.836603548370968. The effect size is the magnitude of the difference, which gives a holistic view of the results [1].
[1] Sullivan, G. M., & Feinn, R. (2012). Using effect size--or why the P value is not enough. Journal of Graduate Medical Education, 40(1), 278-282.
    
```

{valid statistical tests}

How do we logically represent statistical knowledge?

Statistical test applies

iff

all preconditions apply

Student's t-test

↔

bivariate
one_x_variable
one_y_variable
independent_obs
categorical
two_categories
continuous
equal_variance
groups_normal

Statistical test applies

iff

all preconditions apply

Student's t-test

↔

test properties

bivariate

one_x_variable

one_y_variable

independent_obs

variable properties

categorical

two_categories

continuous

equal_variance

groups_normal

Statistical test applies

iff

all preconditions apply

Student's t-test

↔

test properties

bivariate(xy)

one_x_variable(xy)

one_y_variable(xy)

independent_obs(xy)

variable properties

categorical(x)

two_categories(x)

continuous(y)

equal_variance(xy)

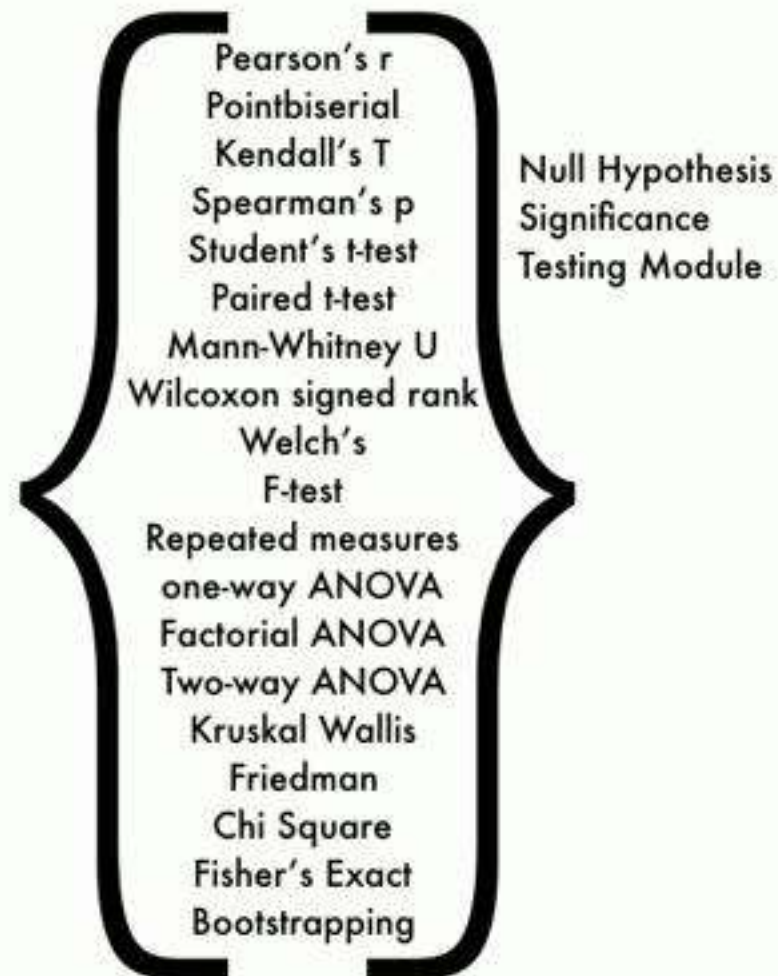
groups_normal(xy)

Statistical test applies

iff

all preconditions apply

Student's t-test



test properties

```
bivariate(xy) ^  
one_x_variable(xy) ^  
one_y_variable(xy) ^  
independent_obs(xy) ^
```

variable properties

```
categorical(x) ^  
two_categories(x) ^  
continuous(y) ^  
equal_variance(xy) ^  
groups_normal(xy)
```

Statistical test selection as constraint satisfaction

```

import tea
tea.data('UScrime.csv')

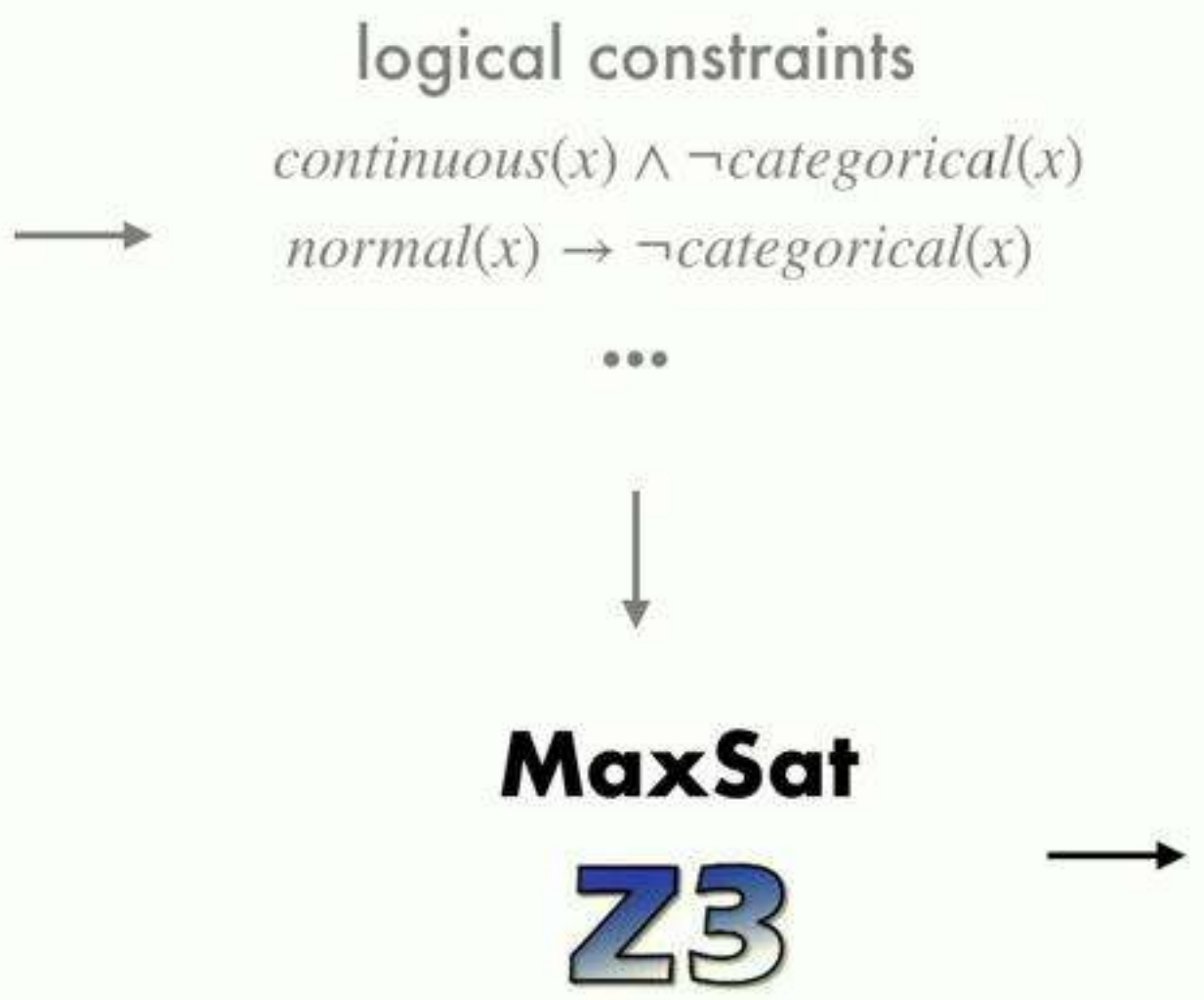
variables = [
    {
        'name' : 'So',
        'data type' : 'nominal',
        'categories' : ['0', '1']
    },
    {
        'name' : 'Prob',
        'data type' : 'ratio',
        'range' : [0,1]
    }
]
tea.define_variables(variables)

study_design = {
    'study type': 'observational study',
    'contributor variables': 'So',
    'outcome variables': 'Prob',
}
tea.define_study_design(study_design)

assumptions = {
    'groups normally distributed': [['So', 'Prob']],
    'Type I (False Positive) Error Rate': 0.05
}
tea.assume(assumptions)

hypothesis = 'So:1 > 0'
tea.hypothesize(['So', 'Prob'], hypothesis)
    
```

- ✓ completeness
- ✓ syntax
- ✓ well-formed hypotheses



```

Test: students.t
---Test assumptions:
Exactly two variables involved in analysis: So, Prob
Exactly one explanatory variable: So
Exactly one explained variable: Prob
Independent (not paired) observations: So
Variable is categorical: So
Variable has two categories: So
Continuous (not categorical) data: Prob
Equal variance: So, Prob
Groups are normally distributed: So, Prob

---Test results:
name = Student's T-Test
test_statistic = 4.202130736875173
p_value = 0.0001206489726652779
adjusted_p_value = 6.182448632060387e-05
alpha = 0.05
df = 45
Effect size:
Cohen's d = 1.2426167266374897
A12 = 0.8366035483870968
Null hypothesis = There is no difference in means between 0 and 1 on Prob.
Interpretation = t(45) = 4.202130736875173, 6.182448632060387e-05. Reject the null hypothesis at alpha = 0.05. The mean of Prob for So = 1 is significantly greater than the mean for So = 0. The effect size is (Cohen's d): 1.2426167266374897, A12: 0.8366035483870968. The effect size is the magnitude of the difference, which gives a holistic view of the results [1].
[1] Sullivan, G. M., & Feinn, R. (2012). Using effect size—or why the P value is not enough. Journal of Graduate Medical Education, 40(1), 278-282.
    
```

{valid statistical tests}

How do we formulate a MaxSat problem?

Z3

Z3

Satisfiability of logical formulas

Z3

Satisfiability of logical formulas

boolean, real number, integer, uninterpreted functions

Z3

Test clauses

```
Students_t_test ^  
bivariate(xy) ^  
one_x_variable(xy) ^  
one_y_variable(xy) ^  
independent_obs(xy) ^  
categorical(x) ^  
two_categories(x) ^  
continuous(y) ^  
equal_variance(xy) ^  
groups_normal(xy)
```

...

Other data/statistical analysis constraints

$(\text{continuous}(x) \vee \text{categorical}(x)) \wedge \neg(\text{continuous}(x) \wedge \text{categorical}(x))$

$\text{normal}(x) \rightarrow \neg\text{categorical}(x)$

$\text{continuous}(x) \vee \text{ordinal}(x) \rightarrow \text{continuous}(x)$

...

Test clauses

$\text{Students_t_test} \wedge$
 $\text{bivariate}(xy) \wedge$
 $\text{one_x_variable}(xy) \wedge$
 $\text{one_y_variable}(xy) \wedge$
 $\text{independent_obs}(xy) \wedge$
 $\text{categorical}(x) \wedge$
 $\text{two_categories}(x) \wedge$
 $\text{continuous}(y) \wedge$
 $\text{equal_variance}(xy) \wedge$
 $\text{groups_normal}(xy)$

...

Z3

Other data/statistical analysis constraints

$(\text{continuous}(x) \vee \text{categorical}(x)) \wedge \neg(\text{continuous}(x) \wedge \text{categorical}(x))$

$\text{normal}(x) \rightarrow \neg\text{categorical}(x)$

$\text{continuous}(x) \vee \text{ordinal}(x) \rightarrow \text{continuous}(x)$

...

User assumptions

```
assumptions = (  
    'groups normally distributed': [['So', 'Prob']],  
    'Type I (False Positive) Error Rate': 0.05  
)  
tea.assume(assumptions)
```

Z3

Test clauses

$\text{Students_t_test} \wedge$
 $\text{bivariate}(xy) \wedge$
 $\text{one_x_variable}(xy) \wedge$
 $\text{one_y_variable}(xy) \wedge$
 $\text{independent_obs}(xy) \wedge$
 $\text{categorical}(x) \wedge$
 $\text{two_categories}(x) \wedge$
 $\text{continuous}(y) \wedge$
 $\text{equal_variance}(xy) \wedge$
 $\text{groups_normal}(xy)$

...

Other data/statistical analysis constraints

$(continuous(x) \vee categorical(x)) \wedge \neg (continuous(x) \wedge categorical(x))$

$normal(x) \rightarrow \neg categorical(x)$

$continuous(x) \vee ordinal(x) \rightarrow continuous(x)$

...

User assumptions

```
assumptions = {  
  'groups normally distributed': [['So', 'Prob']],  
  'Type I (False Positive) Error Rate': 0.05  
}  
tea.assume(assumptions)
```

Test clauses

Students_t_test \wedge
bivariate(xy) \wedge
one_x_variable(xy) \wedge
one_y_variable(xy) \wedge
independent_obs(xy) \wedge
categorical(x) \wedge
two_categories(x) \wedge
continuous(y) \wedge
equal_variance(xy) \wedge
groups_normal(xy)

...

Z3

UNSAT

**Test is invalid.
Remove test.**

Other data/statistical analysis constraints

$(continuous(x) \vee categorical(x)) \wedge \neg (continuous(x) \wedge categorical(x))$
 $normal(x) \rightarrow \neg categorical(x)$
 $continuous(x) \vee ordinal(x) \rightarrow continuous(x)$
...

User assumptions

```
assumptions = {  
    'groups normally distributed': [['So', 'Prob']],  
    'Type I (False Positive) Error Rate': 0.05  
}  
tea.assume(assumptions)
```

Test clauses

Students_t_test \wedge
bivariate(xy) \wedge
one_x_variable(xy) \wedge
one_y_variable(xy) \wedge
independent_obs(xy) \wedge
categorical(x) \wedge
two_categories(x) \wedge
continuous(y) \wedge
equal_variance(xy) \wedge
groups_normal(xy)
...

Z3

UNSAT

Check test assumptions hold

For each property:
If property holds:
Add clause (property == TRUE)
Else:
Add clause (property == FALSE)
Remove last test added

**Test is invalid.
Remove test.**

Other data/statistical analysis constraints

$(continuous(x) \vee categorical(x)) \wedge \neg (continuous(x) \wedge categorical(x))$
 $normal(x) \rightarrow \neg categorical(x)$
 $continuous(x) \vee ordinal(x) \rightarrow continuous(x)$
...

User assumptions

```
assumptions = (  
  'groups normally distributed': [['So', 'Prob']],  
  'Type I (False Positive) Error Rate': 0.05  
)  
tea.assume(assumptions)
```

Z3

Test clauses

Students_t_test \wedge
bivariate(xy) \wedge
one_x_variable(xy) \wedge
one_y_variable(xy) \wedge
independent_obs(xy) \wedge
categorical(x) \wedge
two_categories(x) \wedge
continuous(y) \wedge
equal_variance(xy) \wedge
groups_normal(xy)
...

Check test assumptions hold

For each property:
If property holds:
Add clause (property == TRUE)
Else:
Add clause (property == FALSE)
Remove last test added

All test assumptions are True

Add test to {valid tests}

UNSAT

**Test is invalid.
Remove test.**

Other data/statistical analysis constraints

$(continuous(x) \vee categorical(x)) \wedge \neg(continuous(x) \wedge categorical(x))$
 $normal(x) \rightarrow \neg categorical(x)$
 $continuous(x) \vee ordinal(x) \rightarrow continuous(x)$
...

User assumptions

```
assumptions = {  
  'groups normally distributed': [['So', 'Prob']],  
  'Type I (False Positive) Error Rate': 0.05  
}  
tea.assume(assumptions)
```

Test clauses

$Students_t_test \wedge$
 $bivariate(xy) \wedge$
 $one_x_variable(xy) \wedge$
 $one_y_variable(xy) \wedge$
 $independent_obs(xy) \wedge$
 $categorical(x) \wedge$
 $two_categories(x) \wedge$
 $continuous(y) \wedge$
 $equal_variance(xy) \wedge$
 $groups_normal(xy)$
...

Z3

UNSAT

Check test assumptions hold

For each property:
If property holds:
Add clause (property == TRUE)
Else:
Add clause (property == FALSE)
Remove last test added

**Test is invalid.
Remove test.**

All test assumptions are True

Add test to {valid tests}

If {} bootstrap!

Tea Output

```
Test: students_t
***Test assumptions:
Exactly two variables involved in analysis: So Prob
Exactly one explanatory variable: So
Exactly one explained variable: Prob
Independent (not paired) observations: So
Variable is categorical: So
Variable has two categories: So
Continuous (not categorical) data: Prob
Equal variance: So Prob
Groups are normally distributed: So Prob: NormalTest(W=0.8997463583946228
p_value=0.07962072640657425)
```

Explain rationale for test selection.

```
***Test results:
name = Student's T Test
test_statistic = 4.20213
p_value = 0.00012
adjusted_p_value = 0.00006
alpha = 0.05
dof = 45
```

```
Effect size:
Cohen's d = 1.24262
A12 = 0.83669
```

Contextualize results for accurate interpretation.

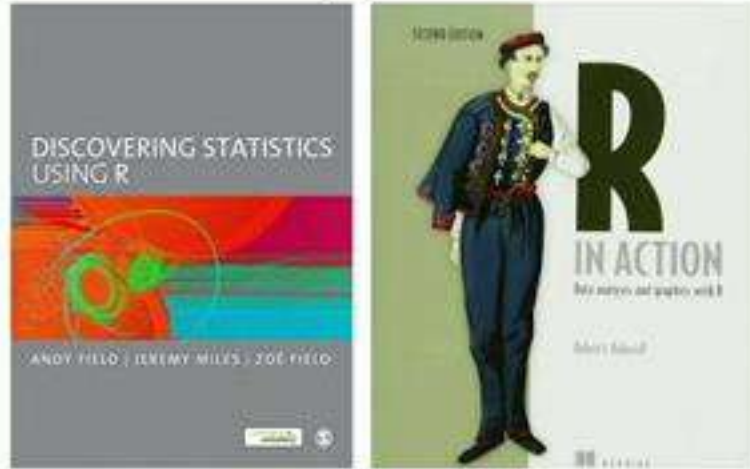
Null hypothesis = There is no difference in means between So = 0 and So = 1 on Prob.

Interpretation = $t(45) = 4.20213$ $p = 0.00006$. Reject the null hypothesis at $\alpha = 0.05$. The mean of Prob for So = 1 ($M=0.06371$ $SD=0.02251$) is significantly greater than the mean for So = 0 ($M=0.03851$ $SD=0.01778$). The effect size is Cohen's $d = 1.24262$ $A12 = 0.83669$. The effect size is the magnitude of the difference which gives a holistic view of the results [1].

[1] Sullivan G. M. & Feinn R. (2012). Using effect size—or why the P value is not enough. Journal of graduate medical education 4(3) 279–282.

Evaluation

12 tutorials
code snippets + text



```
import tea
tea.data('UScrime.csv')
variables = [
  {
    'name': 'So',
    'data type': 'nominal',
    'categories': ['0', '1']
  },
  {
    'name': 'Prob',
    'data type': 'ratio',
    'range': [0,1]
  }
]
tea.define_variables(variables)

study_design = {
  'study type': 'observational study',
  'contributor variables': 'So',
  'outcome variables': 'Prob',
}
tea.define_study_design(study_design)

assumptions = {
  'groups normally distributed': ['So', 'Prob'],
  'Type I (False Positive) Error Rate': 0.05
}
tea.assume(assumptions)

hypothesis = 'So:1 > 0'
tea.hypothesize(['So', 'Prob'], hypothesis)
```


Evaluation

12 tutorials
code snippets + text



```
import tea
tea.data('UScrime.csv')
variables = [
    {
        'name': 'So',
        'data type': 'nominal',
        'categories': ['0', '1']
    },
    {
        'name': 'Prob',
        'data type': 'ratio',
        'range': [0,1]
    }
]
tea.define_variables(variables)

study_design = {
    'study type': 'observational study',
    'contributor variables': 'So',
    'outcome variables': 'Prob',
}
tea.define_study_design(study_design)

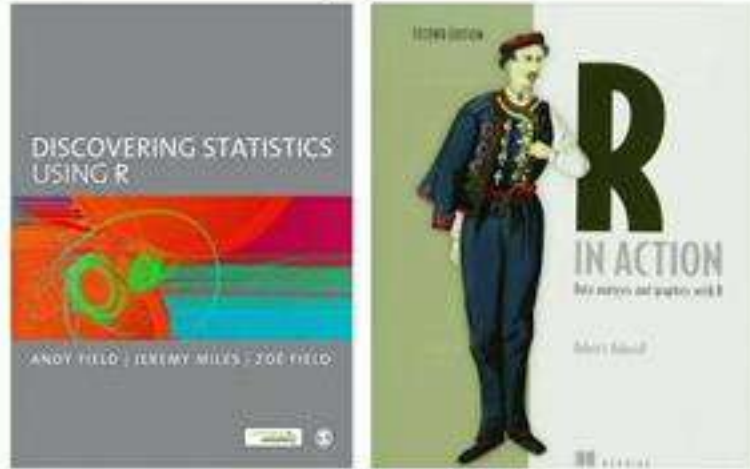
assumptions = {
    'groups normally distributed': ['So', 'Prob'],
    'Type I (False Positive) Error Rate': 0.05
}
tea.assume(assumptions)

hypothesis = 'So:1 > 0'
tea.hypothesize(['So', 'Prob'], hypothesis)
```

I. How does Tea compare to experts?

Evaluation

12 tutorials
code snippets + text



```
import tea
tea.data('UScrime.csv')
variables = [
    {
        'name': 'So',
        'data type': 'nominal',
        'categories': ('0', '1')
    },
    {
        'name': 'Prob',
        'data type': 'ratio',
        'range': (0,1)
    }
]
tea.define_variables(variables)

study_design = {
    'study type': 'observational study',
    'contributor variables': 'So',
    'outcome variables': 'Prob',
}
tea.define_study_design(study_design)

assumptions = {
    'groups normally distributed': ['So', 'Prob'],
    'Type I (False Positive) Error Rate': 0.05
}
tea.assume(assumptions)

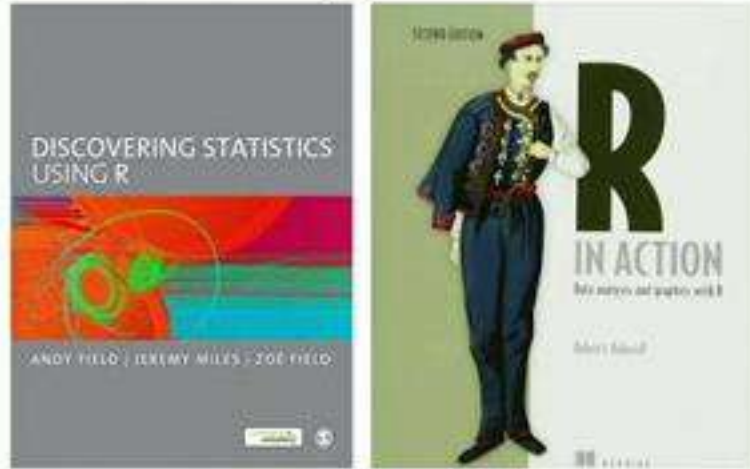
hypothesis = 'So:1 > 0'
tea.hypothesize(['So', 'Prob'], hypothesis)
```

I. How does Tea compare to experts?

Replicate, even improve upon expert choices

Evaluation

12 tutorials
code snippets + text



```
import tea
tea.data('UScrime.csv')
variables = [
    {
        'name': 'So',
        'data type': 'nominal',
        'categories': ('0', '1')
    },
    {
        'name': 'Prob',
        'data type': 'ratio',
        'range': (0,1)
    }
]
tea.define_variables(variables)

study_design = {
    'study type': 'observational study',
    'contributor variables': 'So',
    'outcome variables': 'Prob',
}
tea.define_study_design(study_design)

assumptions = {
    'groups normally distributed': ('So', 'Prob'),
    'Type I (False Positive) Error Rate': 0.05
}
tea.assume(assumptions)

hypothesis = 'So1 > 0'
tea.hypothesize(['So', 'Prob'], hypothesis)
```

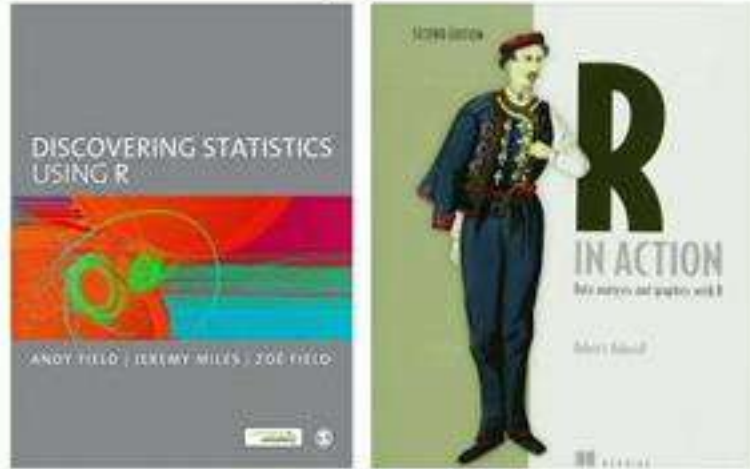
I. How does Tea compare to experts?

Replicate, even improve upon expert choices

II. How does Tea compare to novices?

Evaluation

12 tutorials
code snippets + text



```
import tea
tea.data('OScrime.csv')
variables = [
    {
        'name': 'So',
        'data type': 'nominal',
        'categories': ['0', '1']
    },
    {
        'name': 'Prob',
        'data type': 'ratio',
        'range': [0,1]
    }
]
tea.define_variables(variables)

study_design = {
    'study type': 'observational study',
    'contributor variables': 'So',
    'outcome variables': 'Prob',
}
tea.define_study_design(study_design)

assumptions = {
    'groups normally distributed': ['So', 'Prob'],
    'Type I (False Positive) Error Rate': 0.05
}
tea.assume(assumptions)

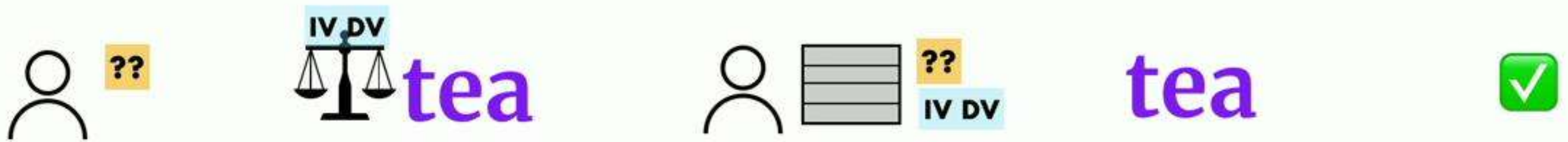
hypothesis = 'So:1 > 0'
tea.hypothesize(['So', 'Prob'], hypothesis)
```

I. How does Tea compare to experts?

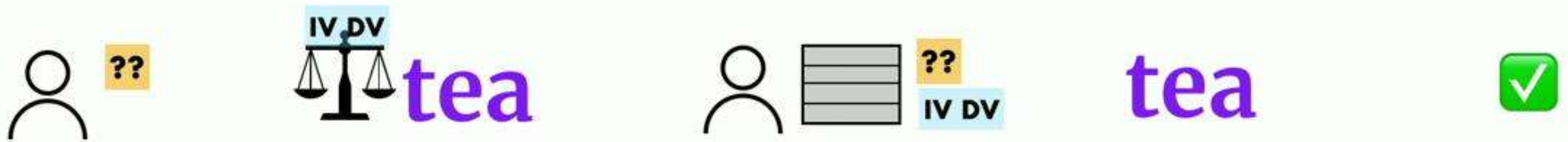
Replicate, even improve upon expert choices

II. How does Tea compare to novices?

Avoid common mistakes and false conclusions

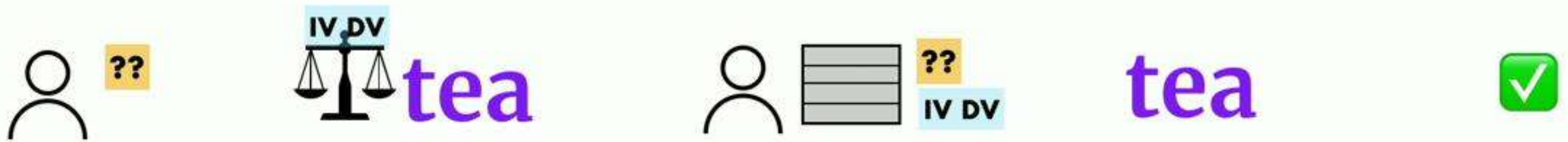


Tea automates statistical test selection and execution.
Tea can aid with experimental design.
Tea programs can act as a format for pre-registration.



Tea automates statistical test selection and execution.
Tea can aid with experimental design.
Tea programs can act as a format for pre-registration.

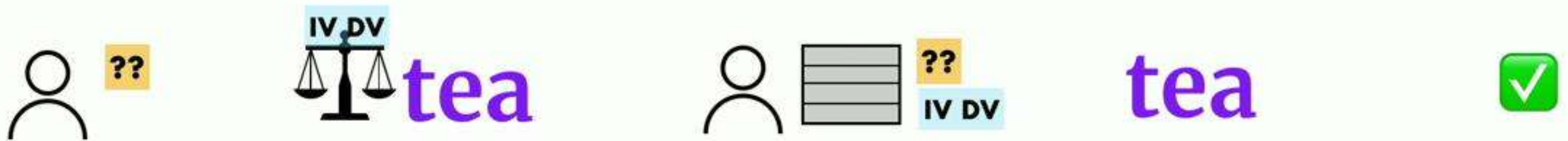
Tea promotes validity and reproducibility in statistical analysis.



Tea automates statistical test selection and execution.
Tea can aid with experimental design.
Tea programs can act as a format for pre-registration.

Tea promotes validity and reproducibility in statistical analysis.

Internal validity!



Tea automates statistical test selection and execution.
Tea can aid with experimental design.
Tea programs can act as a format for pre-registration.

Tea promotes validity and reproducibility in statistical analysis.

Internal validity!

```
pip install tealang  
tea-lang.org
```


Ongoing work



Field deployment, user testing

Future development
- linear modeling



tea

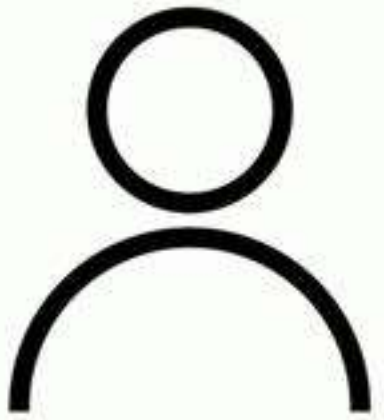
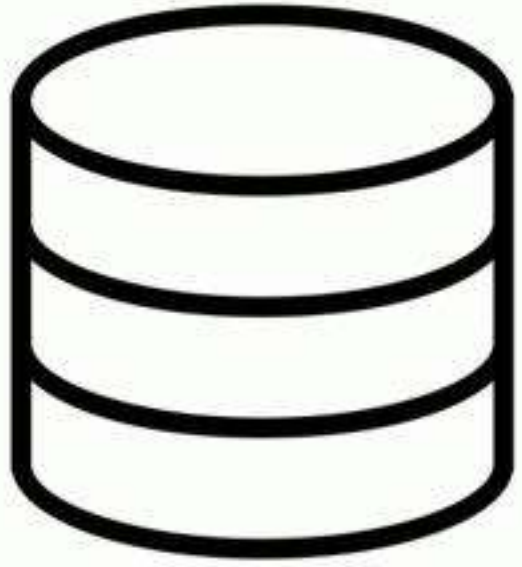
scone
(behind the tea cup)

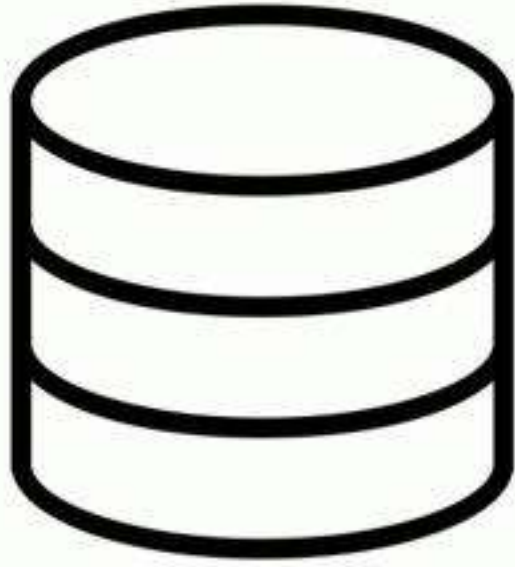
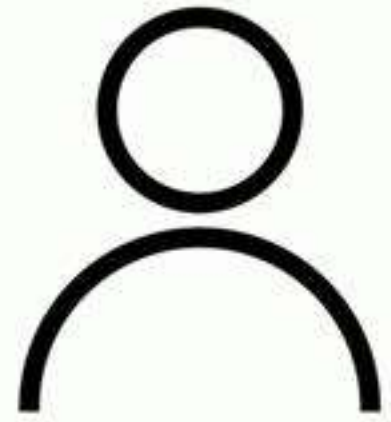


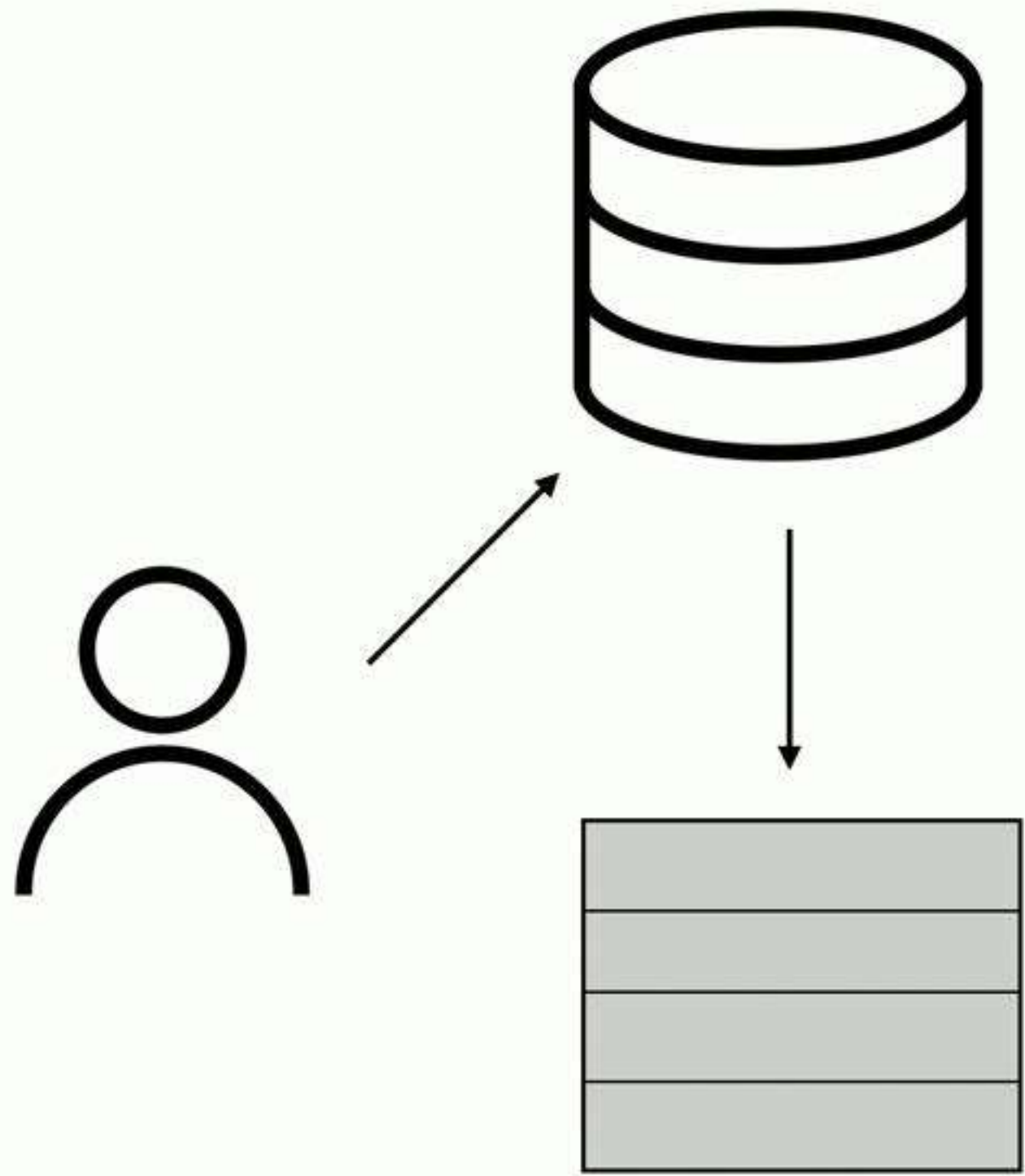
tea

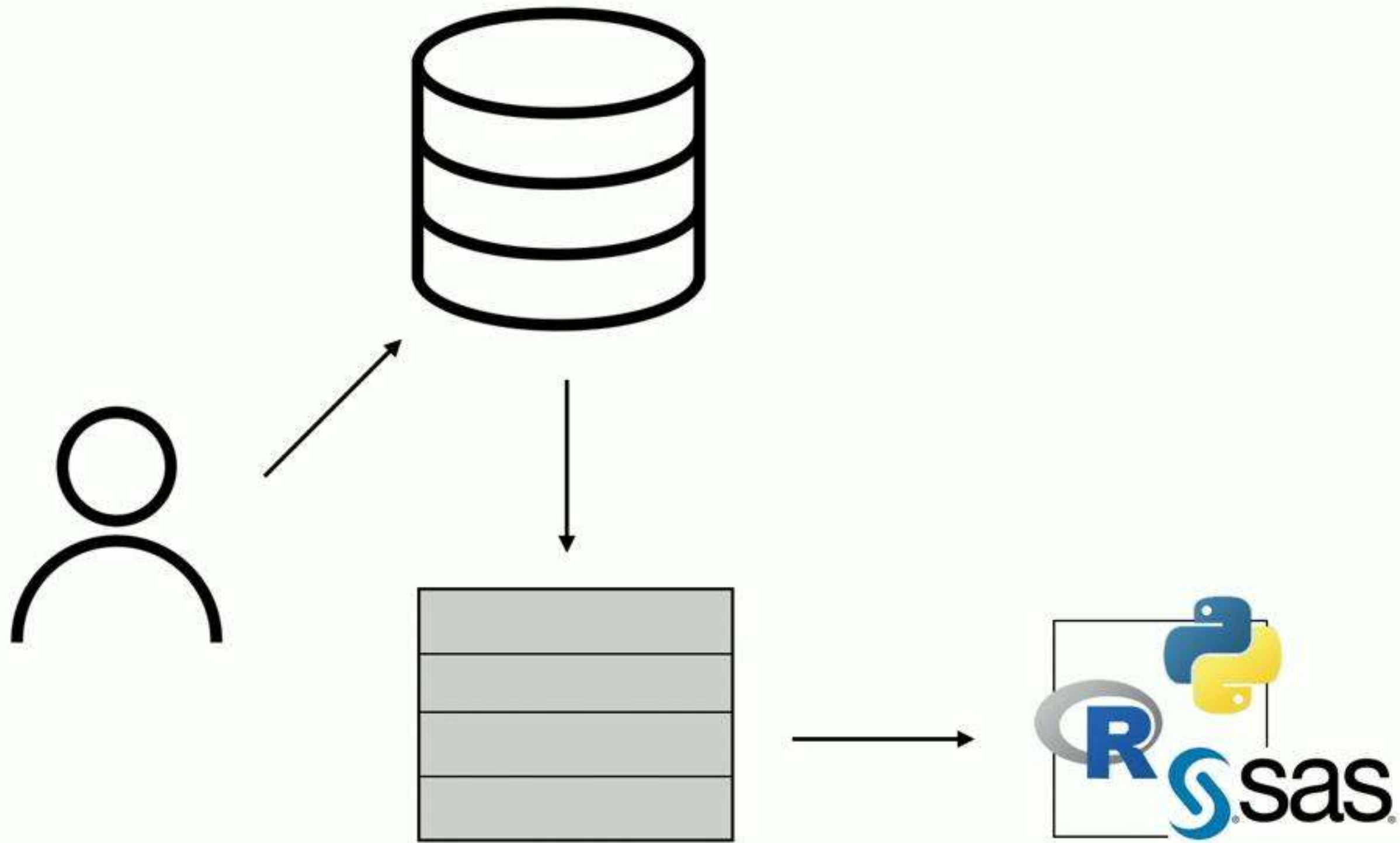
Scone: Smart Sampling for Smarter Statistics

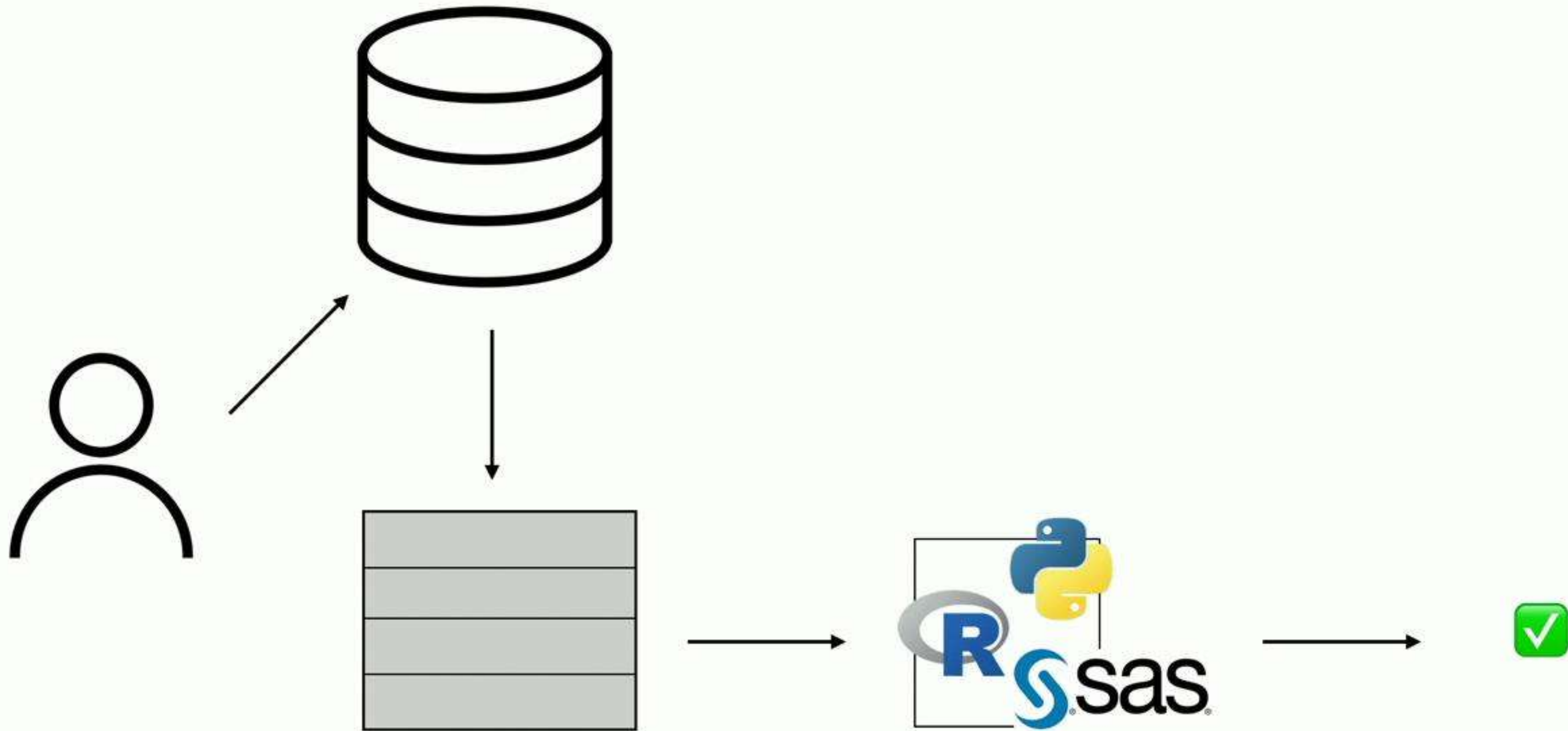
with Laurel Orr, Emery Berger, and Ben Zorn

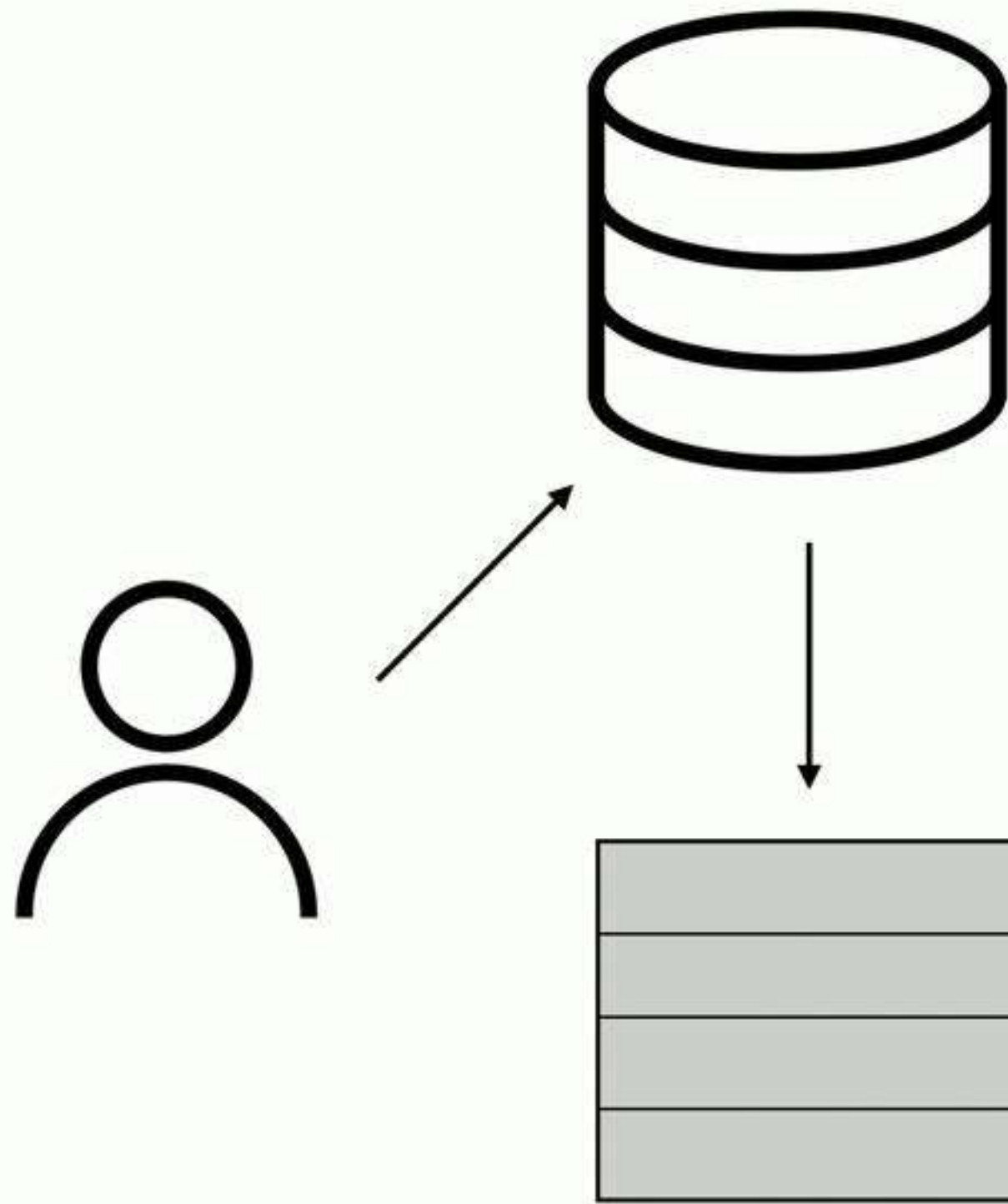






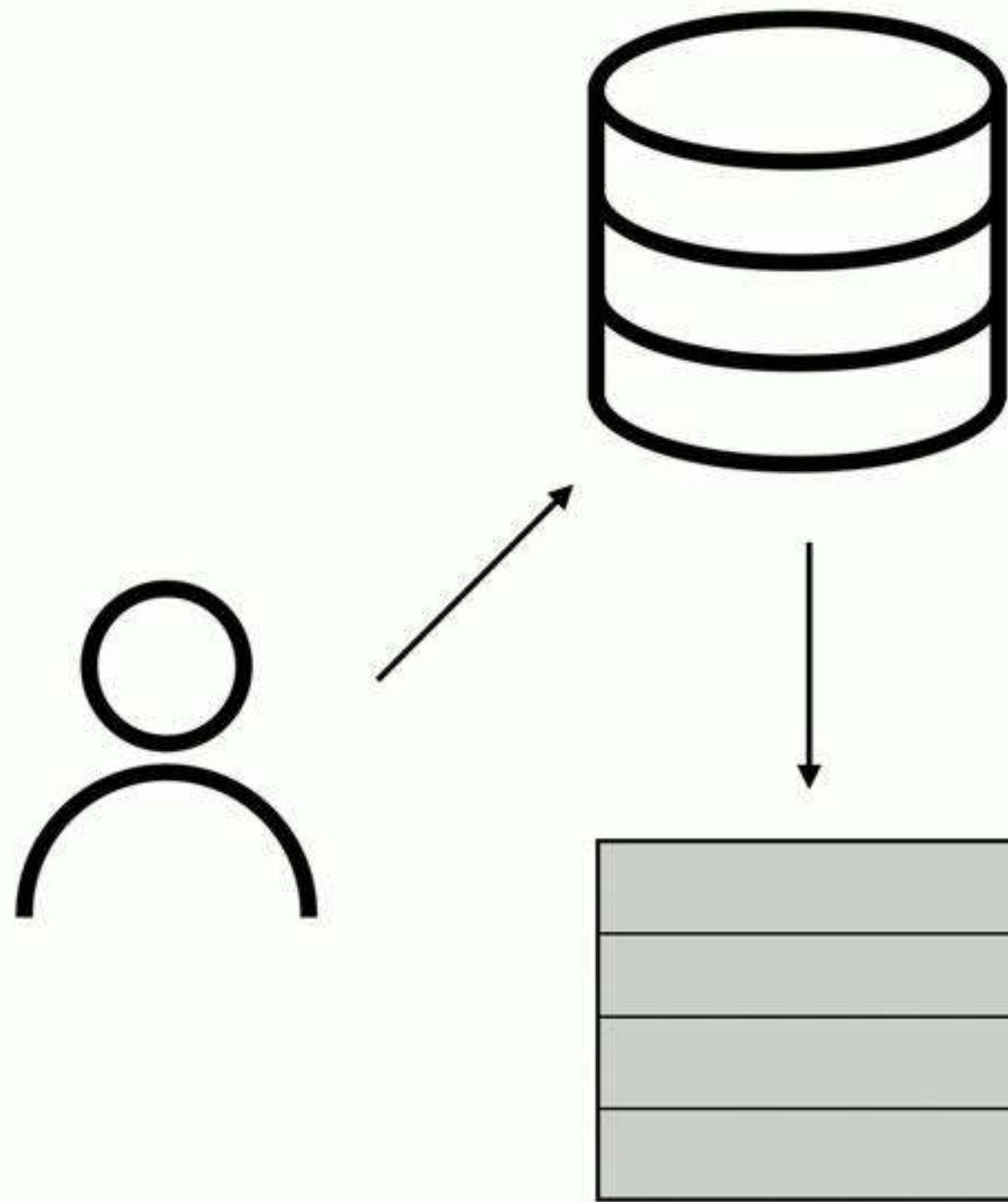






+ Speed

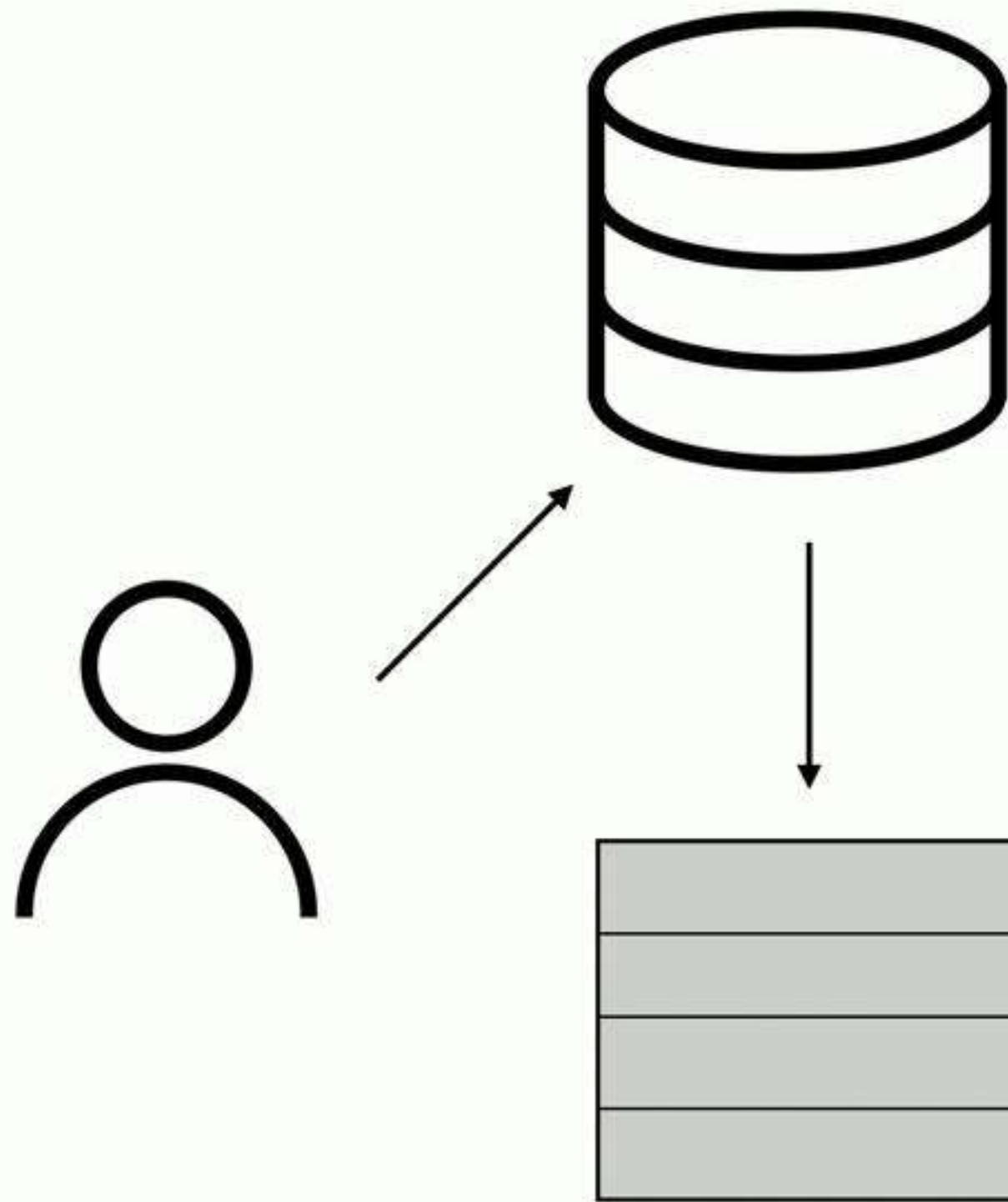




+ **Speed**

+ **Familiarity**



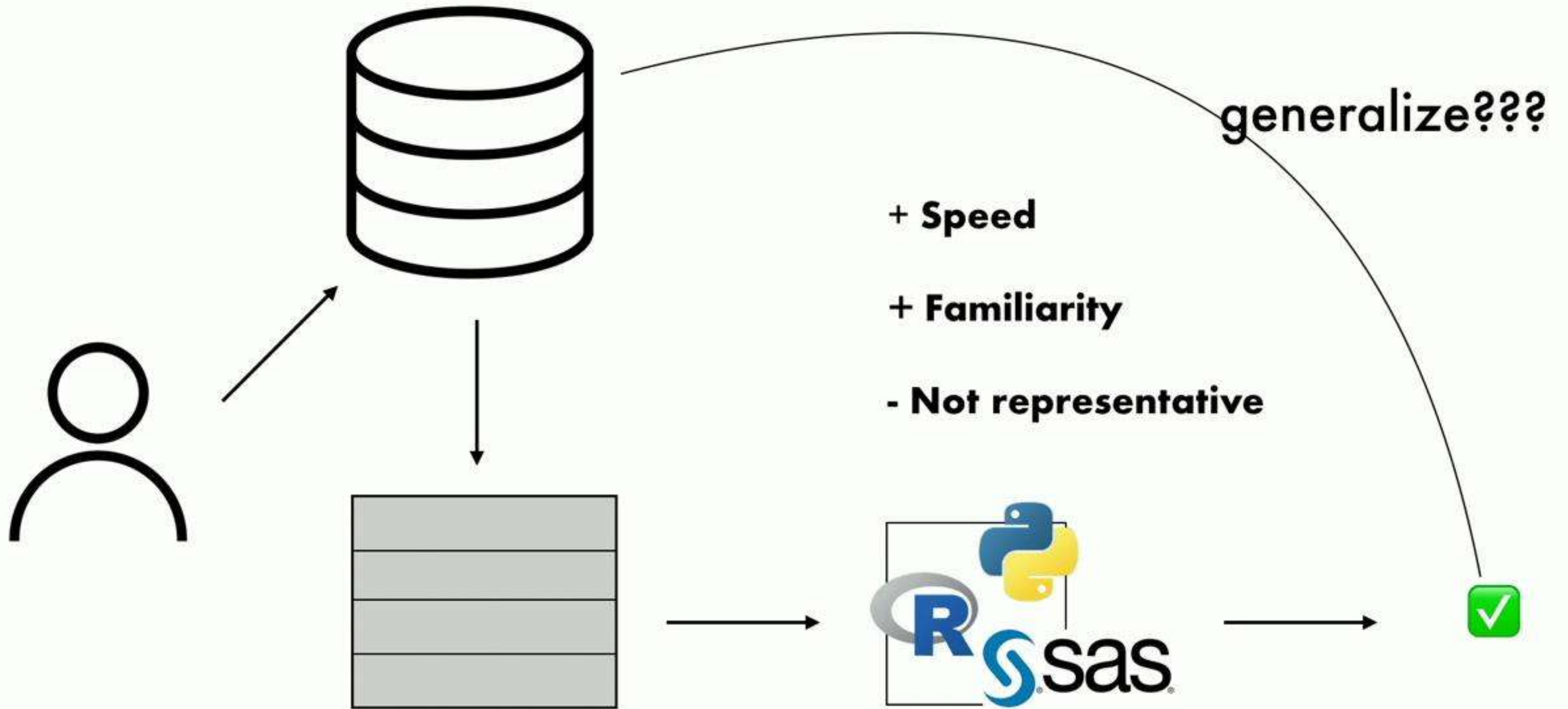


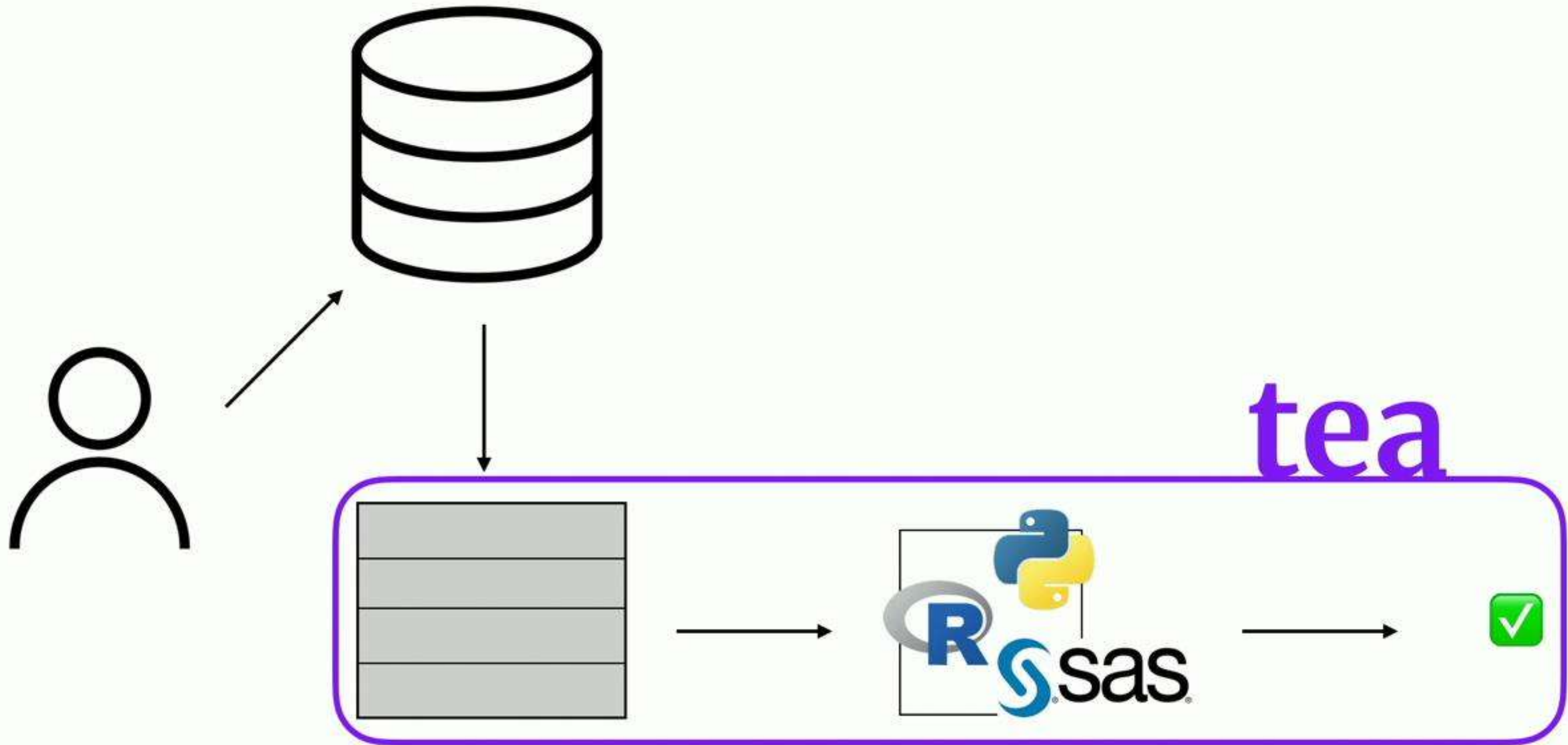
+ Speed

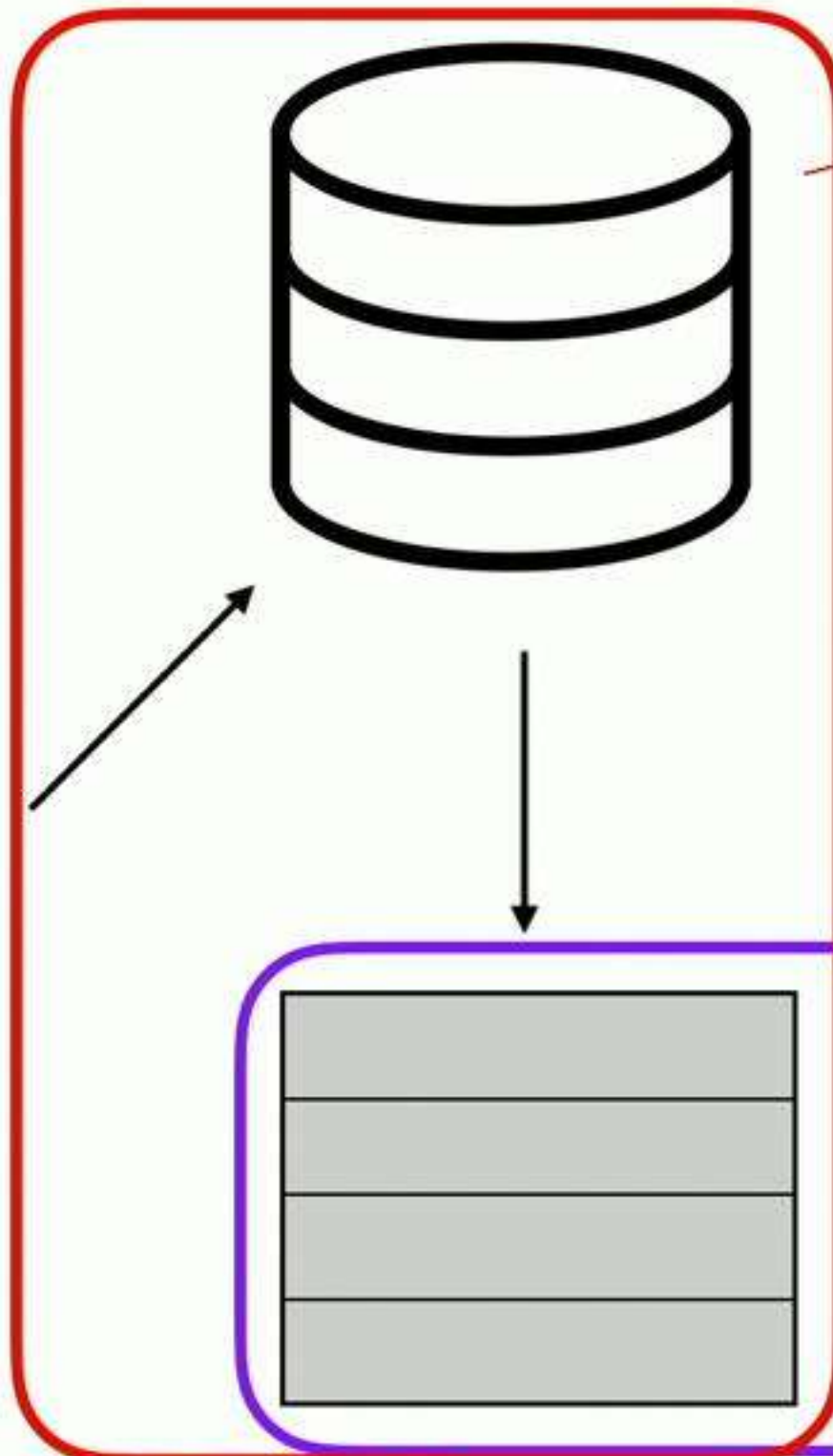
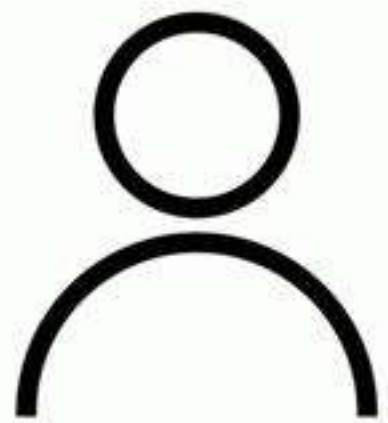
+ Familiarity

- Not representative







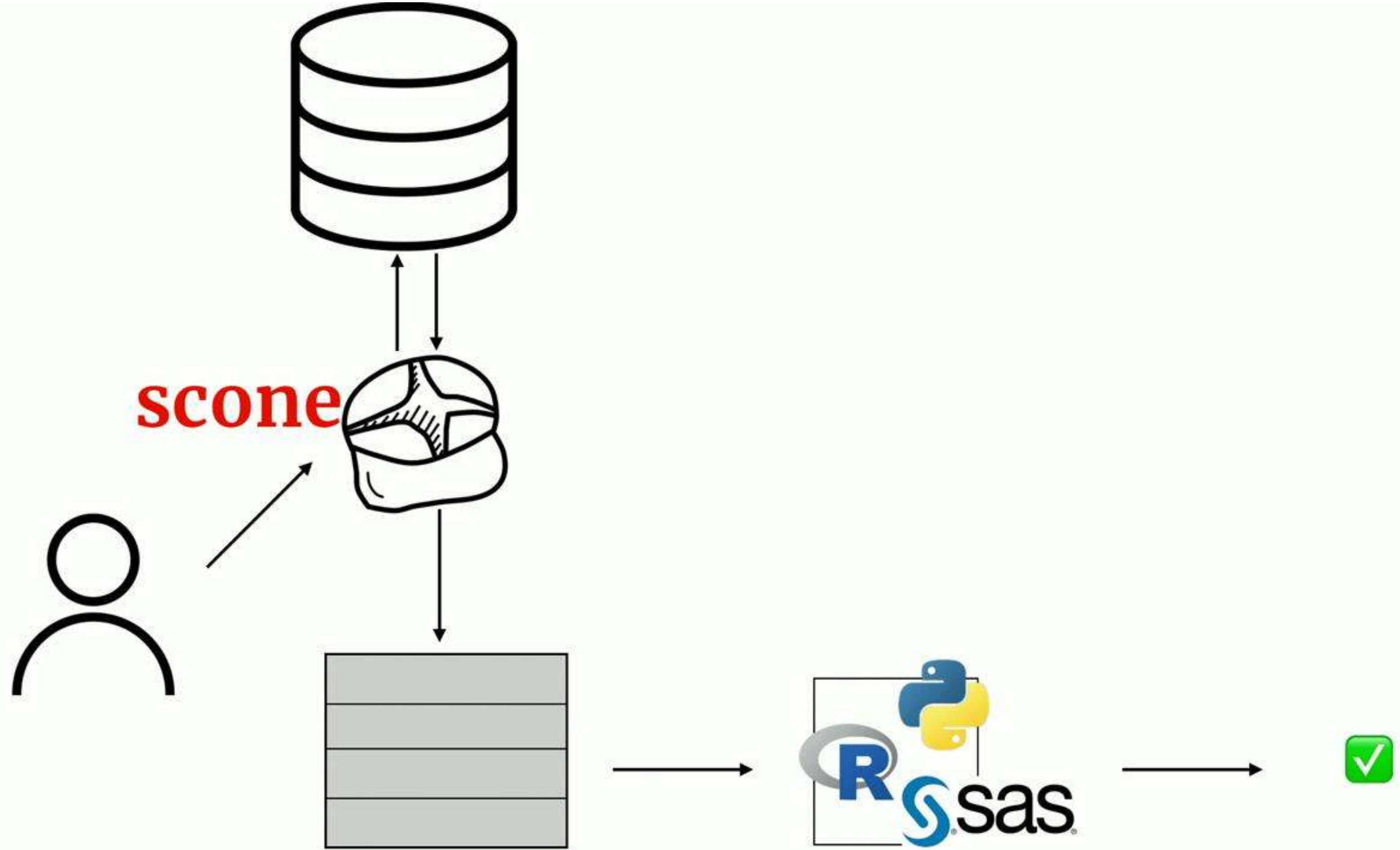


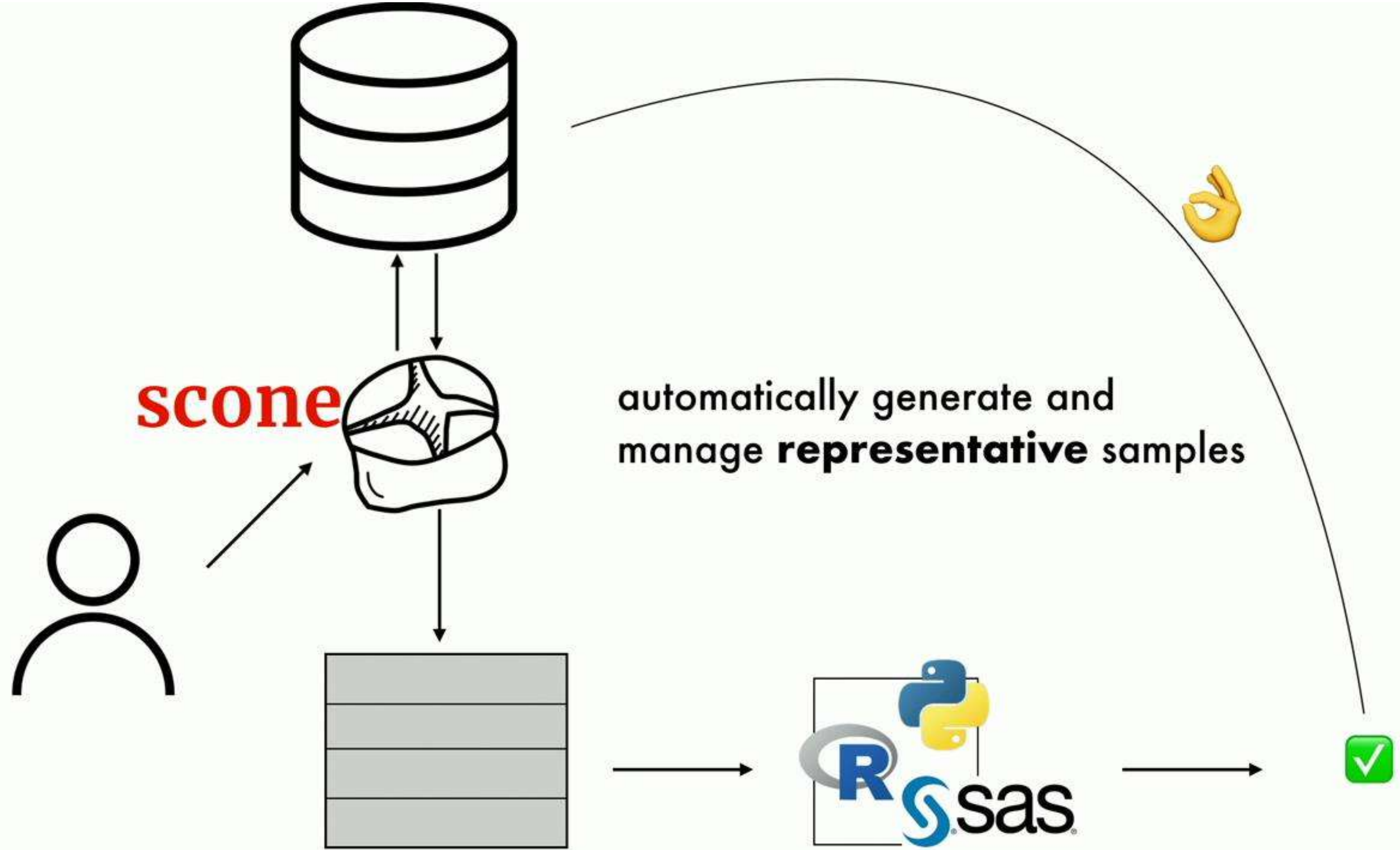
scone



generalize???

tea





tea

Automated statistical analyses

Internal validity

```
pip install tealang  
tea-lang.org
```

scone 

Automated sampling

External validity

Stay tuned!

tea

Automated statistical analyses

Internal validity

```
pip install tealang  
tea-lang.org
```

scone 

Automated sampling

External validity

Stay tuned!



COLLABORATION, USERS, FEEDBACK