The What-If Tool (WIT) Interactive Probing of Machine Learning Models

•••

James Wexler, Mahima Pushkarna, Tolga Bolukbasi, Martin Wattenberg, Fernanda Viegas, and Jimbo Wilson

Presented on Nov 19, by Patrick Huber

Whv? - User Needs

Need 1: Test multiple hypotheses with minimal code

- Interact with trained model through graphical interface (no code)
- Comprehend relationships between data and models

Need 2: Use visualizations as a medium for model understanding

- Generate explanations for model behavior
- Problem: Visual complexity, hard to find meaningful insights
- Solution: Provide multiple, complementary visualizations

Problem & Objective

Problem

- Machine Learning models (e.g. deep learning) are "black-boxes"
 Responses of models to different inputs cannot be easily foreseen
- Big topic in AI: Explainability

Objective

Solution

Interactive visual "what-if" exploration

Whv? - User Needs

- Need 3: Test hypotheticals without having access to the inner workings of a model
- Treat models as black boxes.
- Generate explanations for end-to-end model behavior
- Answer questions like
- No access to model internals Explanations generated remain model-agnostic

Increases flexibility

Model Understanding Frameworks

Black-Box

- Does not rely on internals
- General used in many applications

• WIT

White-Box

 Illuminates internal workings Specific for a model Often not applicable

Whv? - User Needs

Need 4: Conduct exploratory intersectional analysis of model performance

- Users often interested in subsets of data on which models perform unexpectedly
- False positive and false negative rates can be wildly different
- Negative real-world consequences

Need 5: Evaluate potential performance improvements for multiple models

• Track impact of changes in model hyperparameters (e.g. changing a threshold) Interactively debug model performance by testing strategies

Why? - Initial Analysis

Proof-of-concept

Evaluate technical suitability and compatibility of InfoVis solution

Workshops

- 2 usability studies at different scales and with different user-groups
- Application builds on insights from usability studies
- Authors derive 5 distinct user needs

What? - The Tool

- No custom coding (N1)
- Help developers and practitioners to understand ML systems
- Covers many standpoints (Inputs / single data points / models)
- Basic layout: 2 main panels → control panel & visualization panel

https://pair-code.github.io/what-if-tool/iris.html

How? - Tailoring 3 Tasks to Satisfy User Needs

- Closely related to user needs
- Example of the UCI Census dataset

- - Multi-laver neural network
 - Simple linear classifier

How? - Task 2: Investigating What-If Hypothesis

- - Observe partial dependencies
- Apply carefully chosen input modifications (edit, add or delete feature values)

 Result of changing income from \$ 	$3,000 \rightarrow 20,000$ (edit data point):
---	---

Run	Model	Label	Score	Delta
2	1	1 (>50k)	C C 1.991	1.655581
2	1	0 («=50k)	0.038	4-0.641580
2	2	1 (+50k)	0.894	140262
2	2	0 («=50k)	0.057	4-0.156162
1	1	0 (<=50k)	0.650	
1	1	1 (+50k)	0.336	



How? - Task 1: Exploring the Data





How? - Task 1: Exploring the Data



How? - Task 1: Exploring the Data



What? - The Tool



Mode

How? - Task 3: Evaluate Performance and Fairness



Data Scaling

- Assumption: Standard laptop Computational restrictions:
 - # Features: 10-100
- # Datapoints: ~100,000
- Pixel dimensions: 78x64
- # Datapoints: 2,000
 Comment:
- As seen before, occlusion already a problem with less data

Evaluation

- 3 case studies executed
 - 2 studies in a large software company
 - 1 study in a university environment
- Showing the potential of WIT to:
- Uncover bugs
- Explore the data
 Find partial dependent

Thank You

Strengths:

- + Versatile tool
- Many useful real-world app

Strength and Weaknesses

+ Greatly reducing workload compared to creating own visualizations

Weaknesses:

- Only easily compatible with Tensorflow (one deep-learning library)
- Occlusion is a problem, already with small datasets (150 data points, see example
- Strict computational restriction (100,000 data points is not a lot

Questions?

Analysis Summary

- What data:
 - User data & machine learning models
- What derived:
 - Inference of the model (on the data
- What shown:
 - Dataset: and datapoint-level results of ML models
 Giving a better understanding of the capabilities and possible adversaria attacks

Analysis Summary

• How executed:

- $\circ \quad$ 3 common tasks derived from user studies
- How shown:
 - Extension of a out-of-the-box visualization to

Why important:

- Machine Learning models are black boxes
 Making crucial decisions in the real world.
- Understanding is important