# Data Visualization for Machine Learning

Fernanda Viégas     @viegasf

Google Brain
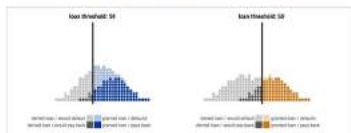
**Embedding Projector**
an open source, visualization tool for high-dimensional data

**Fairness in ML**
Try different tradeoffs yourself to understand issues around fairness and machine learning.

**Machine Translation**
Visualizing hints that a translation network learns an "interlingua", or universal language.

**Geodetic Velocities Visualization**
an open source visualization of earthquake cycle physics

**TensorFlow Playground**
an open source, transparent neural net you can play with in your browser

**Unfiltered News**
see news coverage around the world and spot underreported stories
(a collaboration with Jigsaw)

**TensorFlow Graph Visualizer**
an open source, high-level view of TensorFlow computation graphs
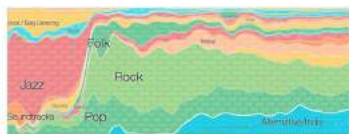
**Periodic Table**
a twist on the classic visualization of the atomic elements

**Music Timeline**
see how different musical genres became popular over time, and discover artists in each genre
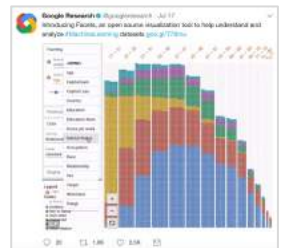UPDATED WEEKLY

**Digital Attack Map**
see live data on denial-of-service attacks across the world, and observe historical patterns
UPDATED DAILY

# PAIR | People + AI Research Initiative

Bringing Design Thinking and HCI to Machine Learning
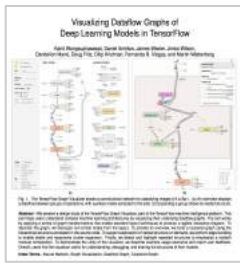google.ai/pair
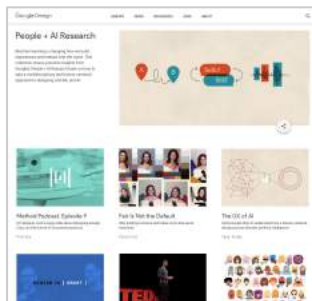
Open Source tools
and platforms

Educational
Materials

Academic
Publications

Public presentations,
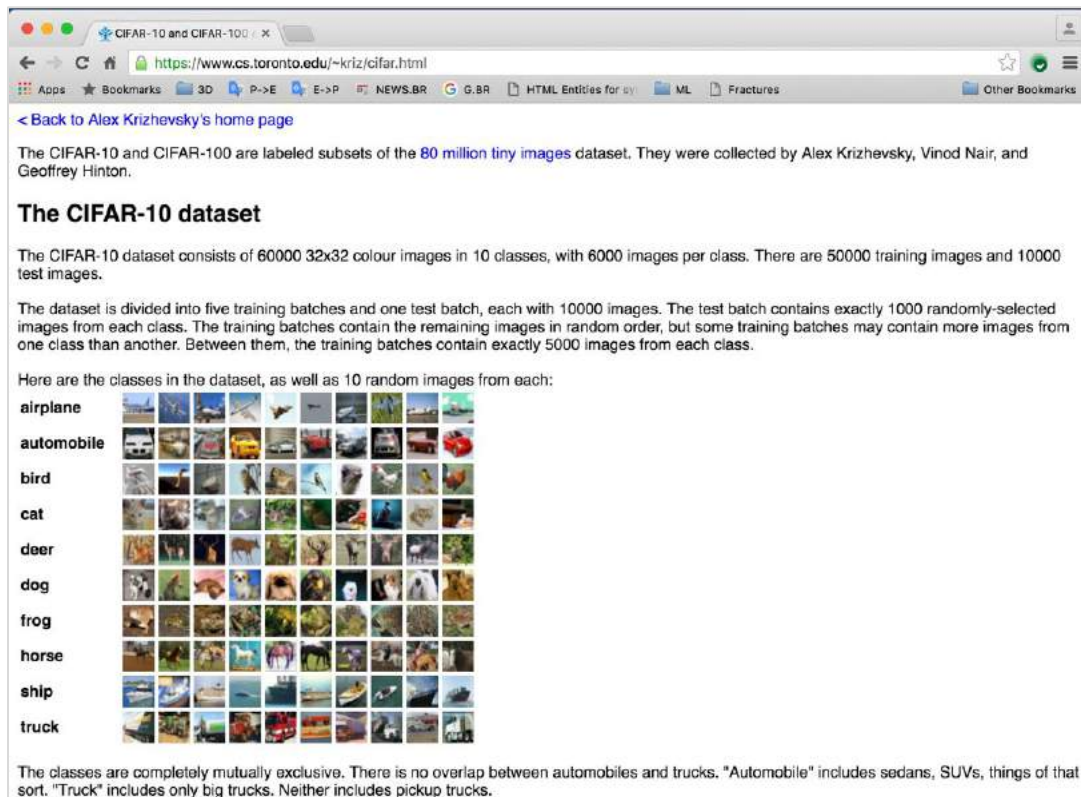sharing best practices

Public Symposia
& meetings



Visiting Faculty,
Faculty Grants

# Training data is crucial

Debug your data before debugging your model

# Let's start with a data set you might have heard of 😉

# 32 x 32 images
# 10 classes

# Facets
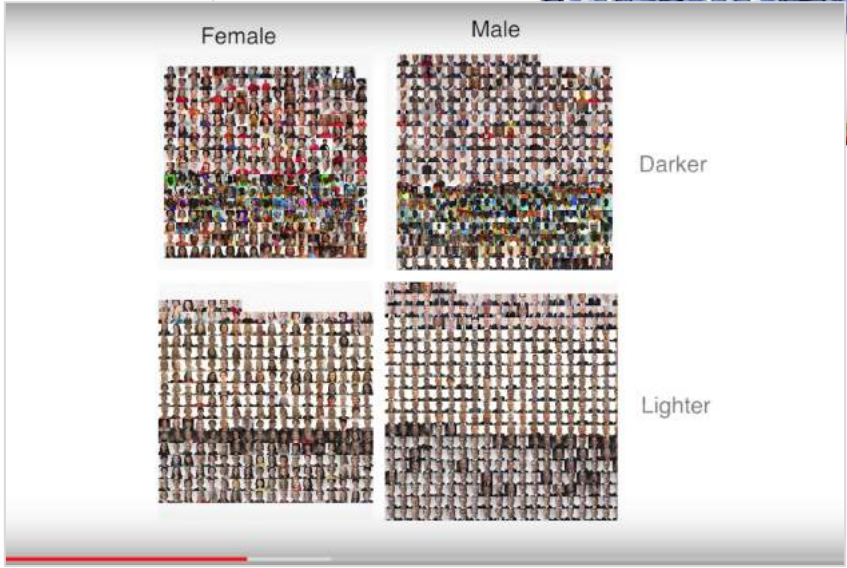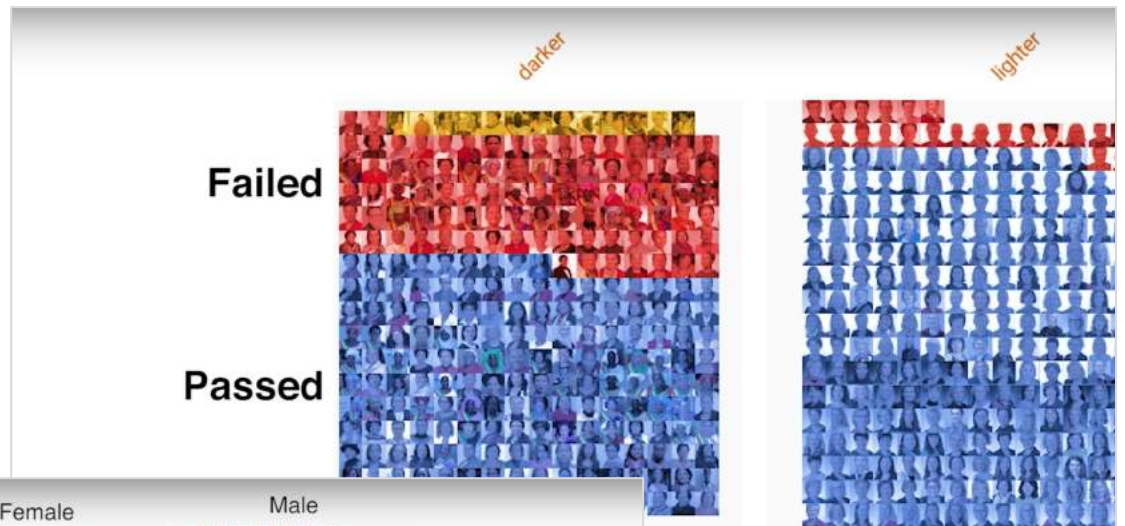
Open-source
pair-code.github.io/facets

# Gender Shades

Joy Buolamwini
MIT Media Lab

# What-If Tool

open source
code-free ML probing
pair-code.github.io/what-if-tool

# What-If Tool

Fairness metrics

# Model Understanding

Looking into high-dimensional spaces

# Warm up:
MNIST

# Images as vectors

Images as vectors

# Images as vectors



0

1

0.5

(1,1,1,1,1,1,1,1,1,.8,0,.5,1,1,1,1,1,.8,.3,0,0,0,0,0,.1,.3,1,1,1,1,1,1,1,1,1,.3,.1,0,0,0,0,0,.4,.9,1, … )

We've turned this image

into this vector



(1,1,1,1,1,1,1,1,.8,0.,5,1,1,1,...)

784 pixels → 784 dimensions

# We've turned this image                    into this vector



$(1,1,1,1,1,1,1,1,\textbf{1},\textbf{.8},\textbf{0.},\textbf{.5},\textbf{1},1,1,...)$

784 pixels → 784 dimensions

$(1,1,1,1,1,1,1,1,\textbf{.6},\textbf{.7},\textbf{0.},\textbf{.4},\textbf{1},1,1,...)$

$(1,1,1,1,1,1,1,1,\textbf{.4},\textbf{.5},\textbf{0.},\textbf{.3},\textbf{.2},1,1,...)$

.
.
.

# Embedding projector
## MNIST visualization

# Model interpretability use case

## Multi-lingual translation

What does the language embedding space look like?

Melvin Johnson, Mike Schuster, Quoc V. Le, Maxim Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thorat, Fernanda Viégas, Martin Wattenberg, Greg Corrado, Macduff Hughes, Jeffrey Dean

**Training**:    English    ← → Japanese
            English    ← → Korean

**Training:**   English ← → Japanese
        English ← → Korean
        Japanese ← → Korean (zero shot)



Training

# Visualize internal representation ("embedding space")

# Research question
## What does the multi language embedding space look like?



or

Note: not real data

# What does a sentence look like in embedding space?
(points in 1024-dim space: the data that the decoder receives)

E.g. "*The stratosphere extends from 10km to 50km in altitude*"

# What does a sentence look like in embedding space?



Note: simplification of real situation!

# What does a sentence look like in embedding space?

# What do parallel sentences look like in embedding space?
(same meaning, different language)

## like this?



● English
● Portuguese

# What do parallel sentences look like in embedding space?
(same meaning, different language)

## or like this?



● English
● Portuguese

# Interlingua?

Sentences with the same meaning mapped to similar regions regardless of language!



ENGLISH
The stratosphere extends from about 10km to about 50km in altitude.

KOREAN
성층권은 고도 약 10km부터 약 50km까지 확장됩니다.

JAPANESE
成層圏は、高度 10km から 50km の範囲にあります。

# Distance between bridge / non-bridge sentences is inversely related to translation quality



Figure 3: (a) A bird's-eye view of a t-SNE projection of an embedding of the model trained on Portuguese→English (blue) and English→Spanish (yellow) examples with a Portuguese→Spanish zero-shot bridge (red). The large red region on the left primarily contains the zero-shot Portuguese→Spanish translations. (b) A scatter plot of BLEU scores of zero-shot translations versus the average point-wise distance between the zero-shot translation and a non-bridged translation. The Pearson correlation coefficient is −0.42.

word embeddings

# Word embeddings



Country and Capital Vectors Projected by PCA

Distributed Representations of Words and Phrases and their Compositionality
Mikolov et al. 2013

# Meaningful directions
(word2vec)

Man ← ... Woman →

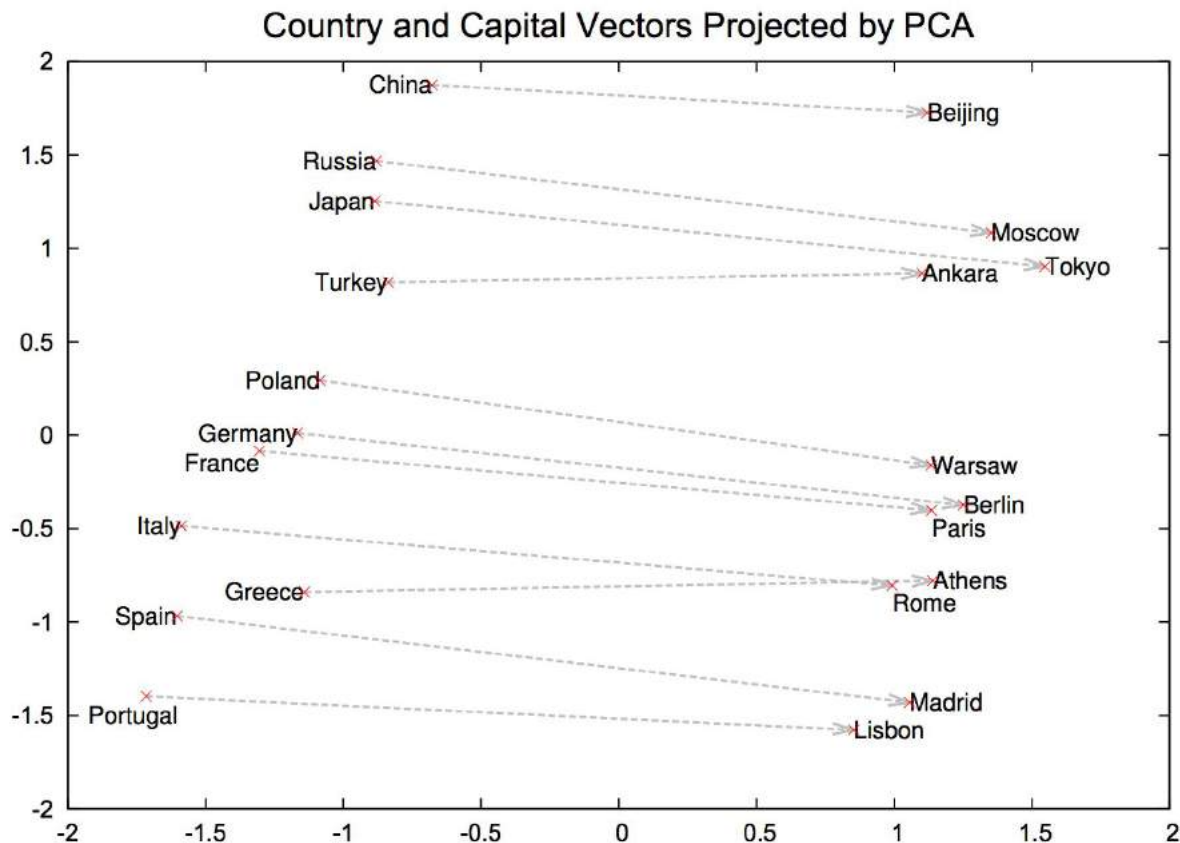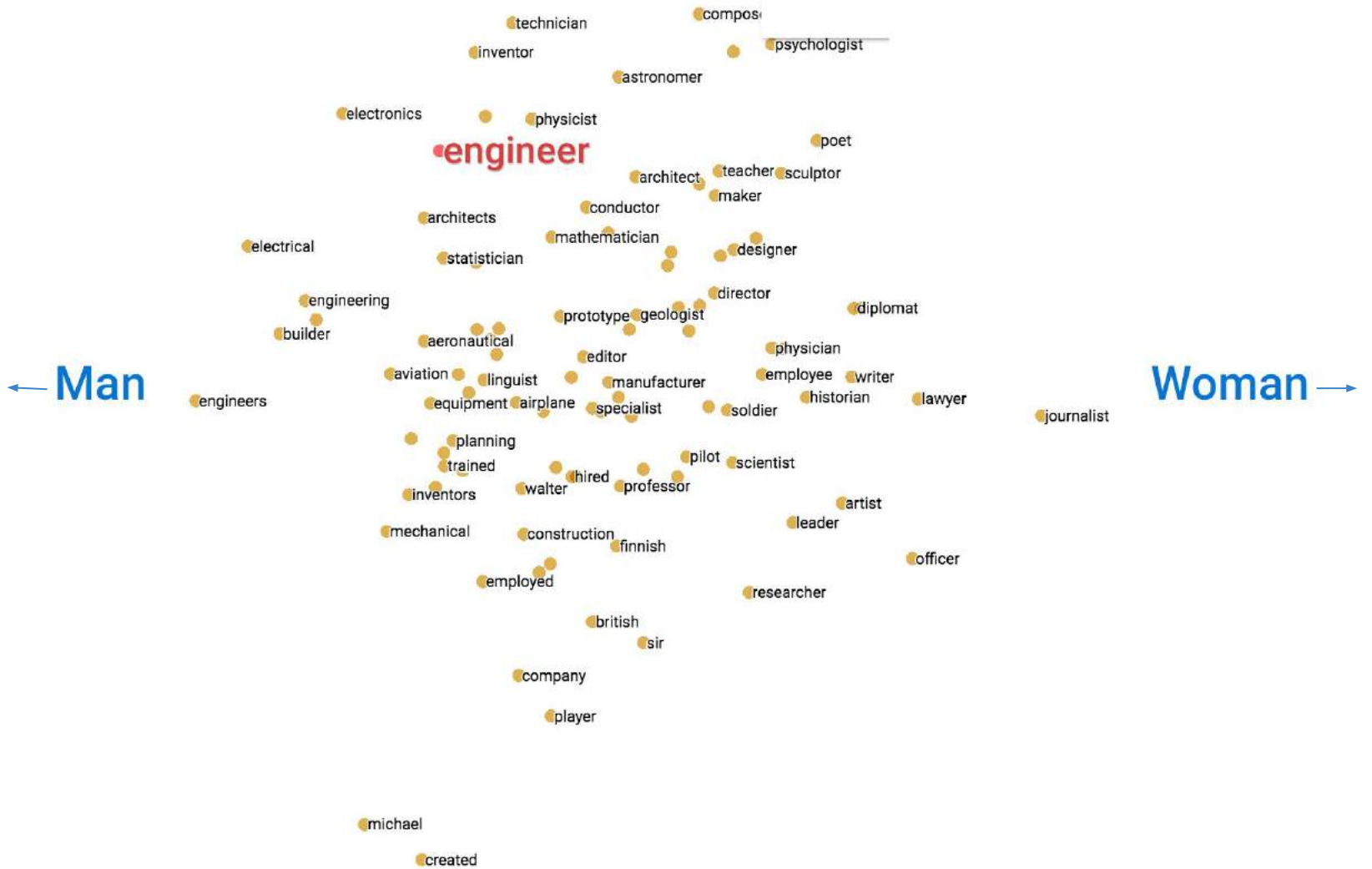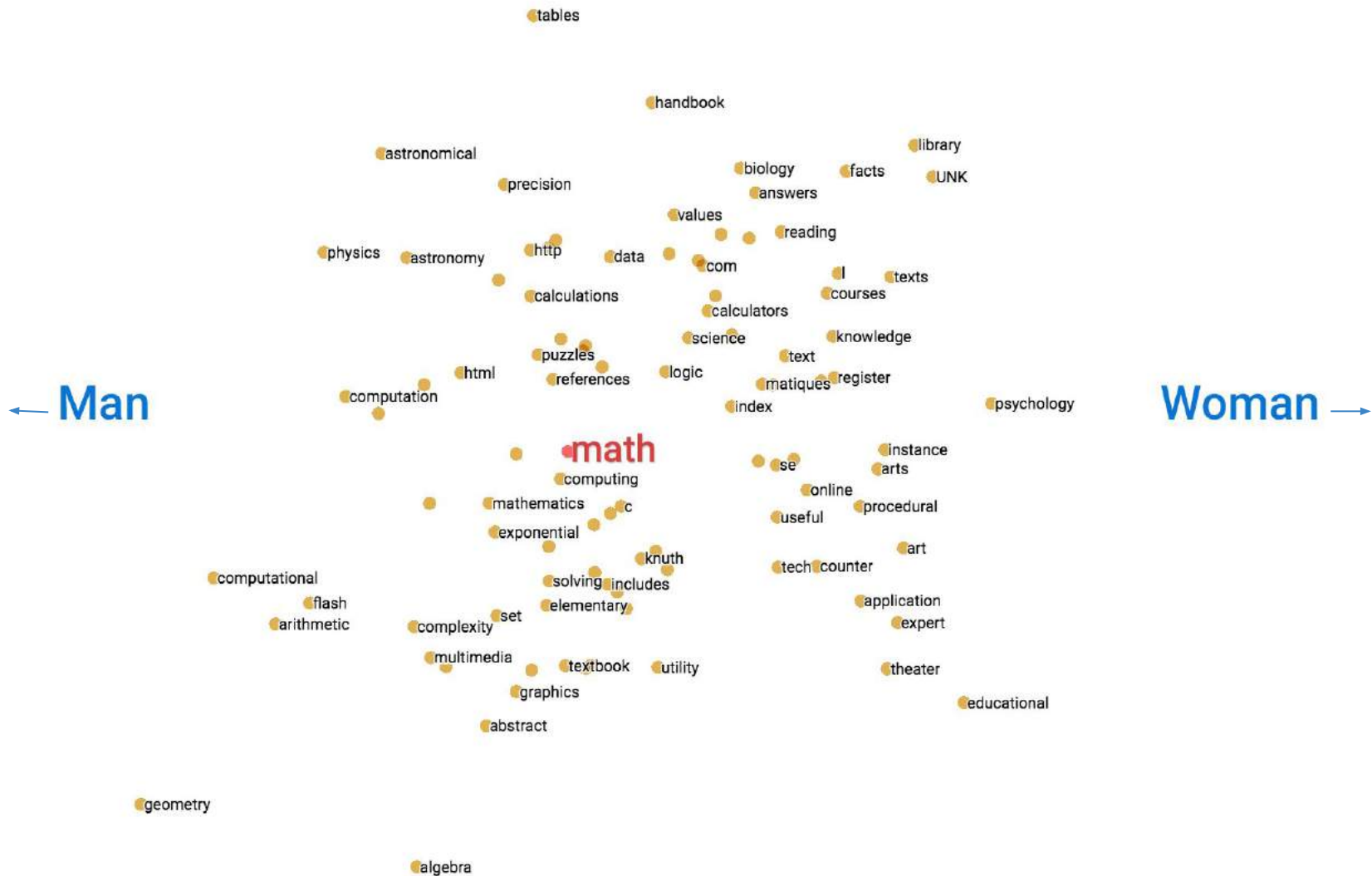tables, handbook, astronomical, library, biology, facts, UNK, precision, answers, values, reading, physics, astronomy, http, data, com, l, texts, calculations, courses, calculators, puzzles, science, knowledge, html, references, logic, text, matiques, register, computation, index, psychology, math, se, instance, arts, computing, online, procedural, mathematics, c, useful, exponential, art, knuth, tech, counter, computational, solving, includes, flash, elementary, application, arithmetic, set, expert, complexity, multimedia, textbook, utility, theater, graphics, educational, abstract, geometry, algebra
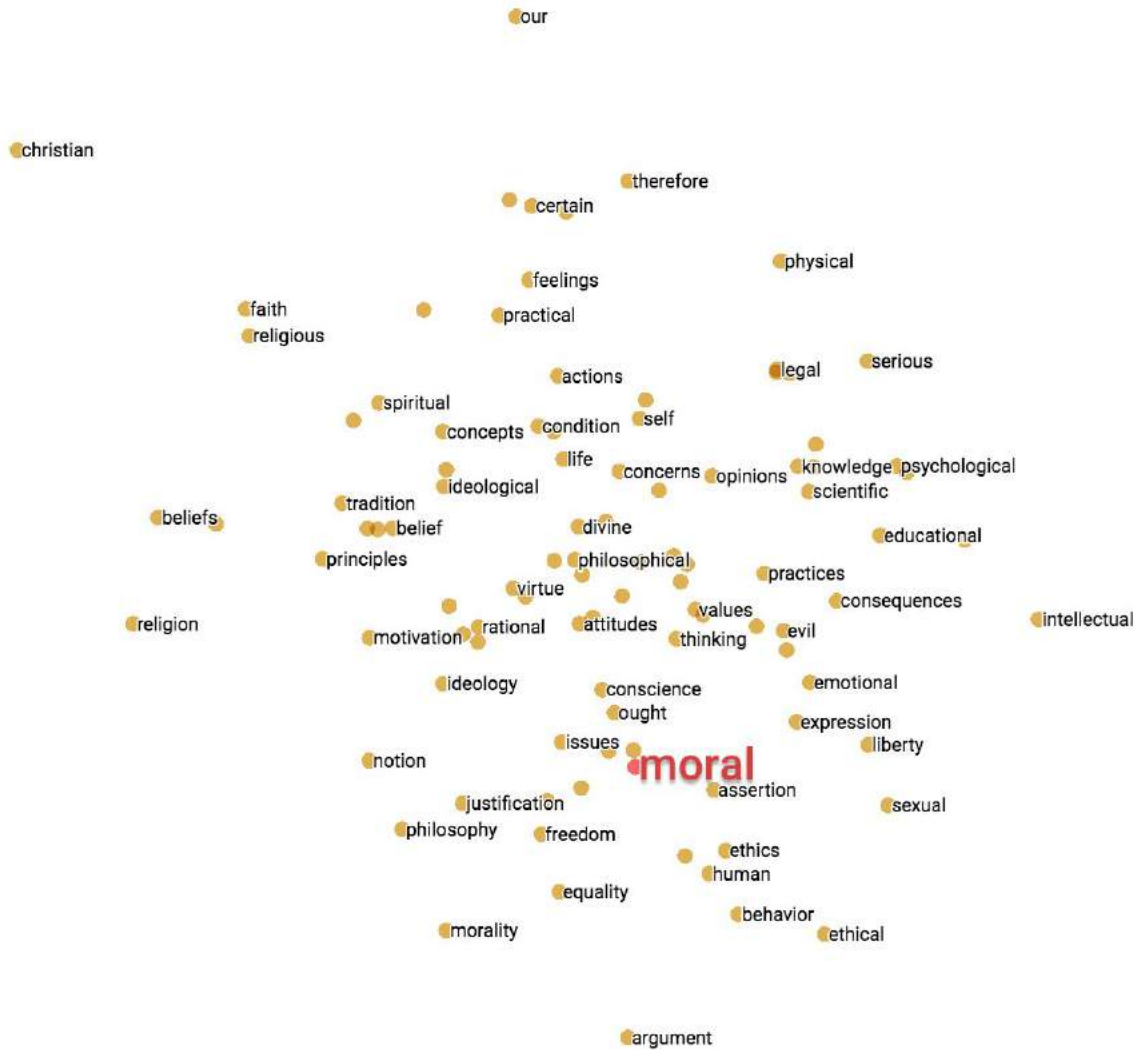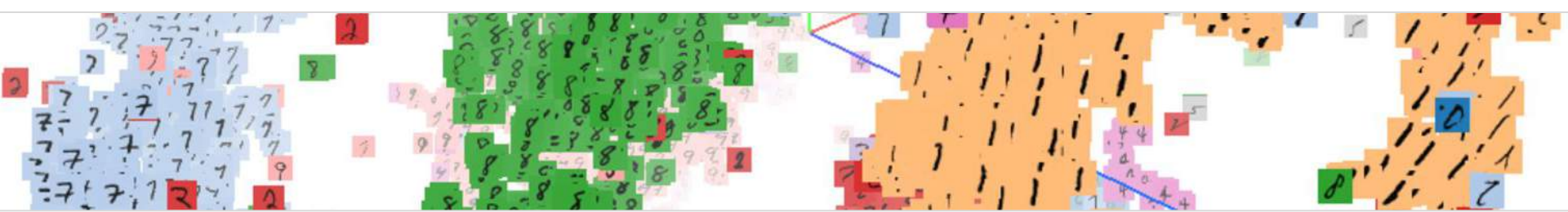
# Data Visualization for Machine Learning

Fernanda Viégas          @viegasf

google.ai/pair