

# Evolutionary dataset optimisation: learning algorithm quality through evolution

---

Henry Wilde, Dr. Jonathan Gillard, Dr. Vincent Knight

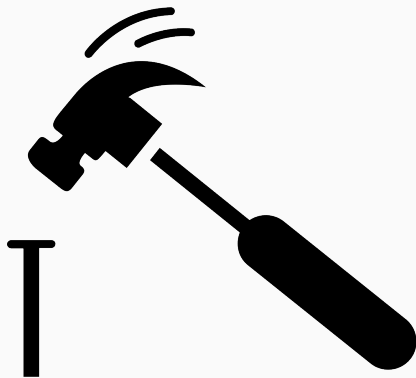


GIG  
CYMRU  
NHS  
WALES

Bwrdd Iechyd Prifysgol  
Cwm Taf  
University Health Board

## Premise and motivation

---



**BBC** Sign in

News Sport Weather iPlayer Sounds More

Search

NEWS

Home UK World Business Politics Tech Science Health Family & Education Entertainment & Arts Stories More

Technology

# Google apologises for Photos app's racist blunder

1 July 2015

f WhatsApp Twitter Email Share

Skyscrapers

Airplanes

Cars

Bikes

Gorillas

Graduation

## Top Stories

**EU considers potential Brexit delay**

EU leaders remain locked in discussions amid reports that they may offer a delay until 7 May.

15 minutes ago

**Latest as EU leaders meet in Brussels**

18 March 2019

**Trump: Time to recognise Golan as Israeli**

1 hour ago

## Features

via: BBC News (<https://www.bbc.co.uk/news/technology-33347866>)

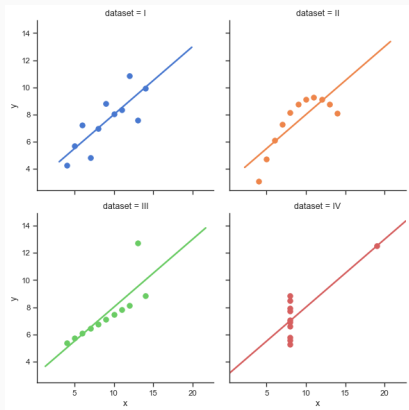
- Hyndman reference (reliability)
- A. Torralba and A. A. Efros. *Unbiased Look at Dataset Bias*. Proceedings of the 2011 IEEE Conference on Computer Vision and Pattern Recognition. 2011. DOI: 10.1109/CVPR.2011.5995347 (frailty)

Paradigm flip diagram

# Premise and motivation

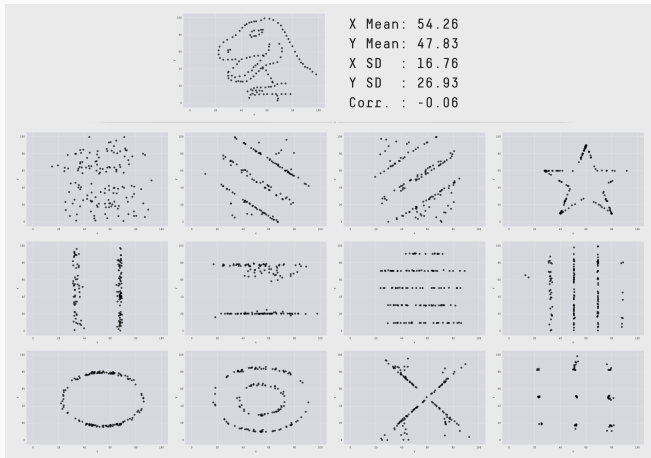
---

Generating artificial data



via: [https://seaborn.pydata.org/examples/anscombes\\_quartet.html](https://seaborn.pydata.org/examples/anscombes_quartet.html)





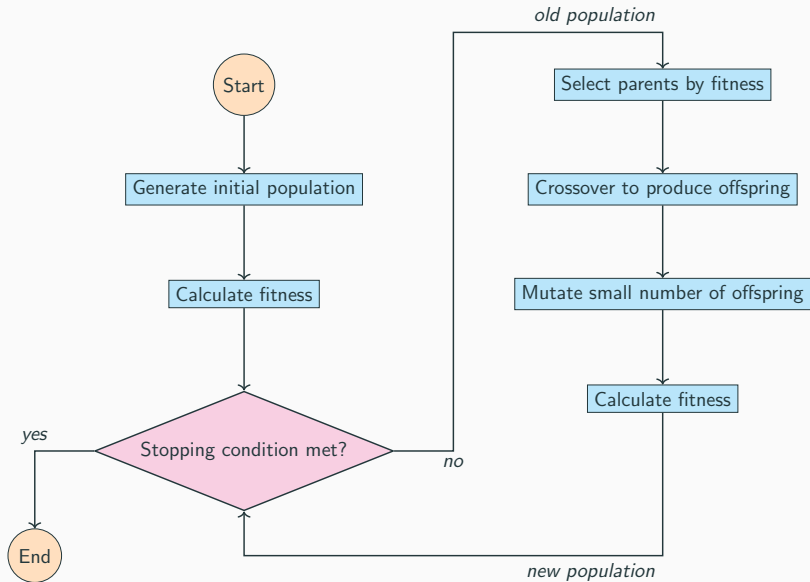
via: <https://www.autodeskresearch.com/publications/samestats>

Dataset similarity diagram

Given some algorithm, how can one find data for which it performs well?

**What is an evolutionary algorithm?**

---



GA example

Families of  
distributions

$$N(\mu, \sigma^2)$$

$$U(\alpha, \beta)$$

$$Po(\lambda)$$

Sample distribution  
and parameters

Column  
information

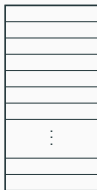
$$N(0.25, 1)$$

$$U(1.2, 3.2)$$

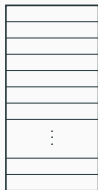
$$N(-3.7, 0)$$

Sample values  
from distribution

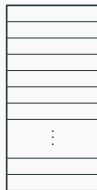
Columns of  
the dataset

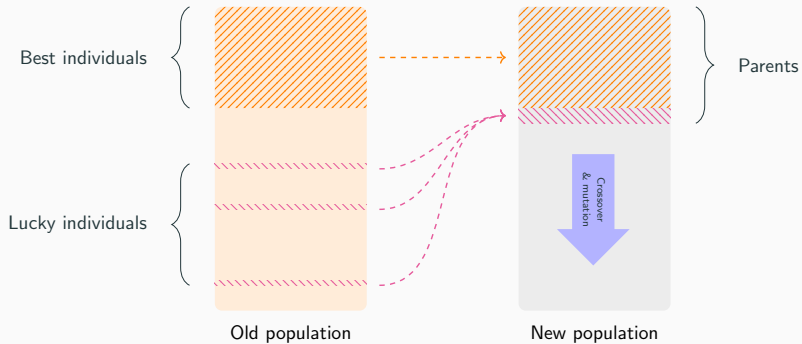


+

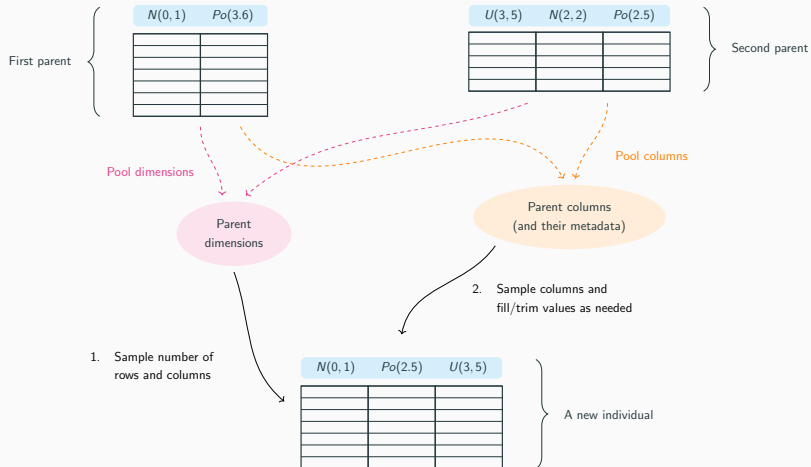


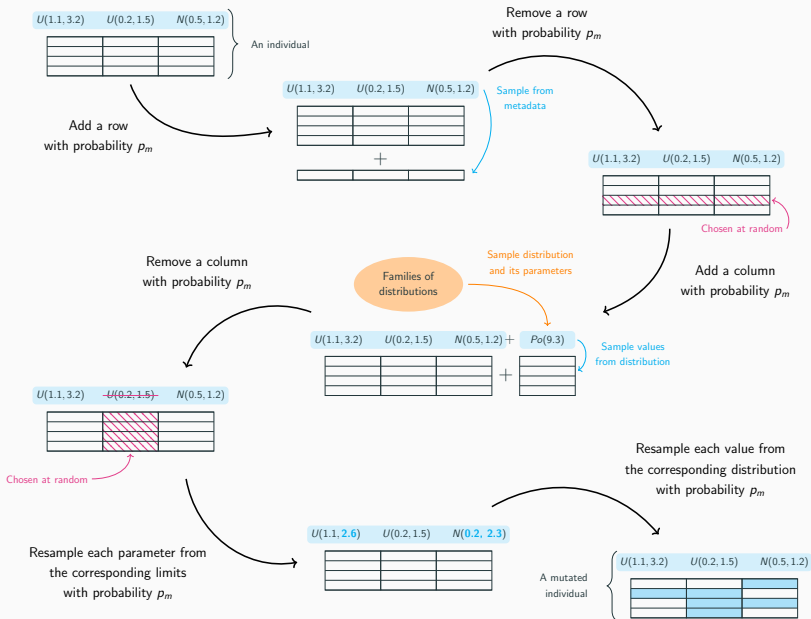
+











- A fitness function,  $f$ , which acts on a single dataset
- A population size,  $N \in \mathbb{N}$
- Limits on the number of rows a dataset can have:

$$R \in \left\{ (r_{\min}, r_{\max}) \in \mathbb{N}^2 \mid r_{\min} \leq r_{\max} \right\}$$

- Limits on the number of columns a dataset can have:

$$C := \left( C_1, \dots, C_{|\mathcal{P}|} \right) \text{ where } C_j \in \left\{ (c_{\min}, c_{\max}) \in (\mathbb{N} \cup \{\infty\})^2 \mid c_{\min} \leq c_{\max} \right\}$$

for each  $j = 1, \dots, |\mathcal{P}|$

- A set of probability distribution families,  $\mathcal{P}$ . Each family in this set has some parameter limits which form a part of the overall search space
- A probability vector to sample distributions from  $\mathcal{P}$ ,  $w = (w_1, \dots, w_{|\mathcal{P}|})$
- A maximum number of iterations,  $M \in \mathbb{N}$
- Two selection parameters: one to indicate the proportion of the fittest individuals to carry forward,  $b \in [0, 1]$ , and the other to allow for a small proportion of “lucky” individuals in the next generation,  $l \in [0, 1]$
- A mutation probability,  $p_m \in [0, 1]$
- A shrink factor,  $s \in [0, 1]$ . The relative size of a component of the search space to be retained after adjustment

## Some example use cases

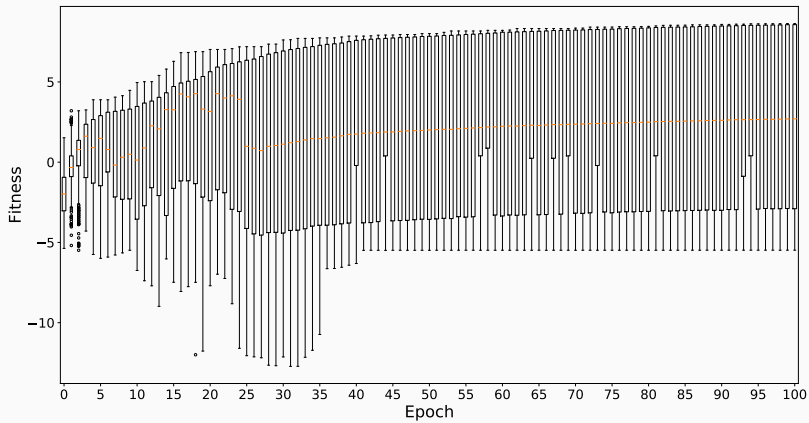
---

Given a dataset  $X$ , maximise  $f$  where:

$$f(X) := \max_{a,b \in \{(0,1),(1,0)\}} \left\{ \text{Var} \left( X^{(a)} \right) - \max_i \left| X_i^{(b)} - 1 \right| \right\}$$



Fitness progression



---

**Input:** A dataset (column)  $X$ , some sampling fraction  $p$

**Output:** An estimate for the mean of  $X$

$Y \leftarrow$  a random sample of  $\lfloor p|X| \rfloor$  entries from  $X$ ;  
evaluate the mean of  $Y$ :

$$\bar{Y} = \frac{1}{|Y|} \sum_{i=1}^{|Y|} Y_i$$

---

$$\iff f(X) = \bar{Y}$$



Anscombe's quartet

```
$ pip install edo
```

# Any questions?

Henry Wilde

Twitter: @daffidwilde

Email: henrydavidwilde@gmail.com

GitHub repository: <https://github.com/daffidwilde/edo>

Documentation: <https://edo.readthedocs.io>