# GANs and other optimizations
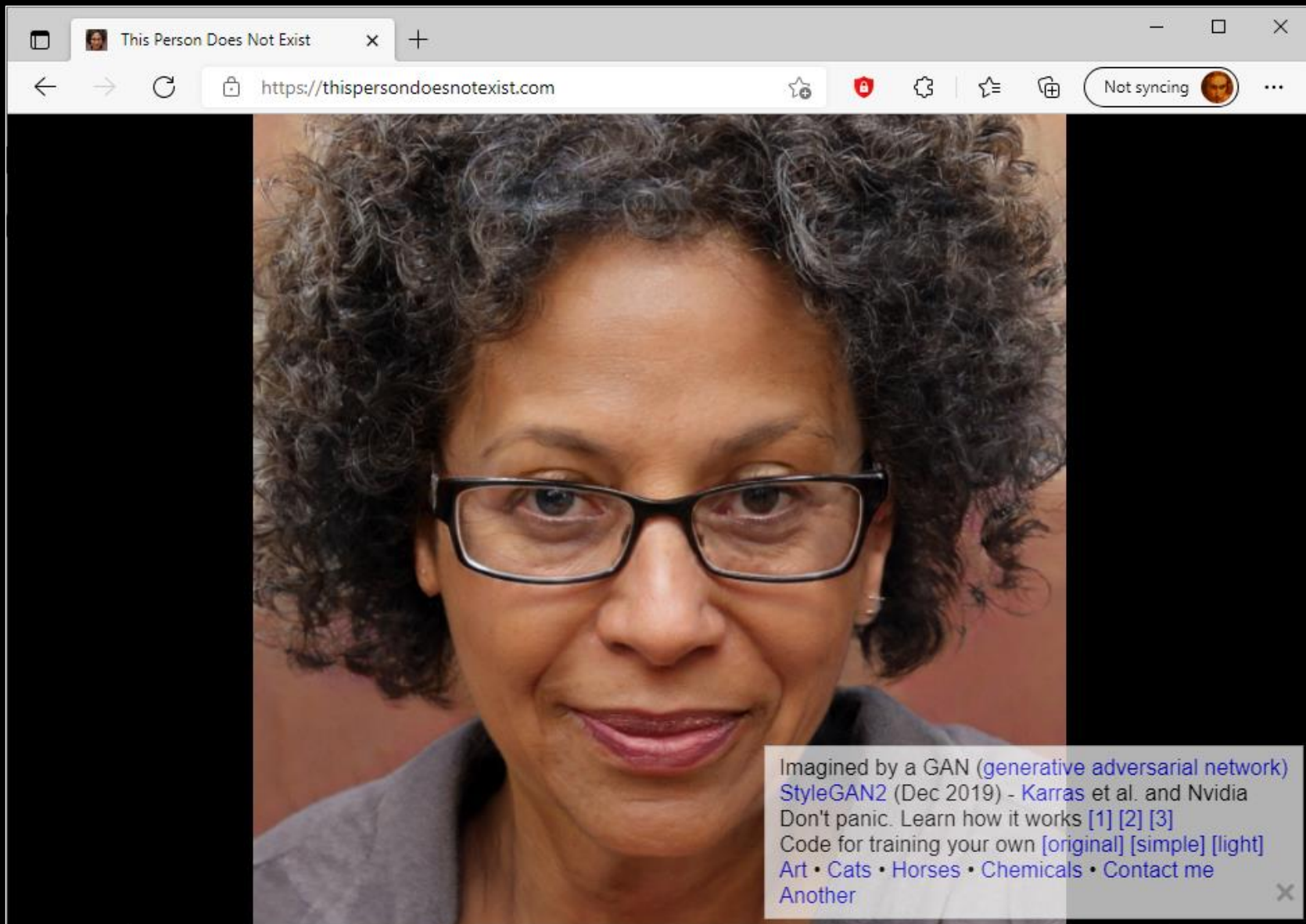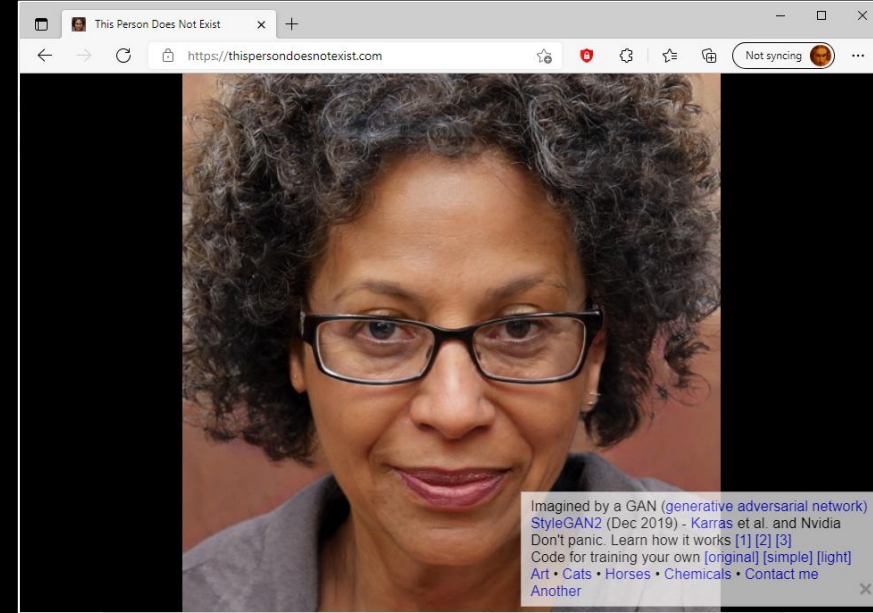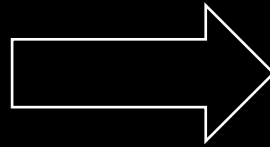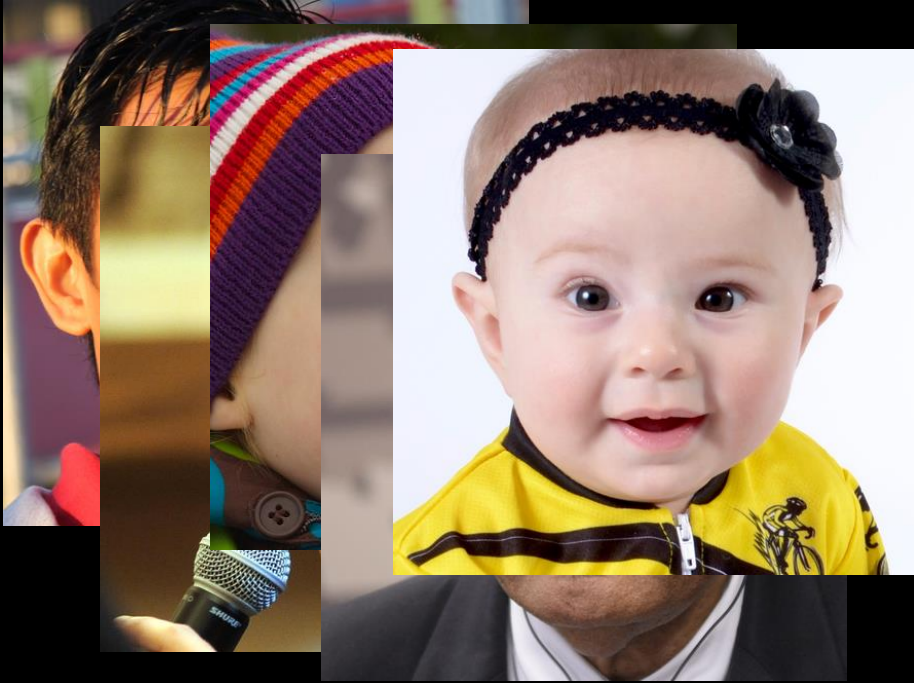
Imagined by a GAN (generative adversarial network)
StyleGAN2 (Dec 2019) - Karras et al. and Nvidia
Don't panic. Learn how it works [1] [2] [3]
Code for training your own [original] [simple] [light]
Art • Cats • Horses • Chemicals • Contact me
Another

Given a dataset of images, how can we train a neural network to be able to generate realistic fakes?



Flickr-Faces-HQ Dataset (FFHQ)
https://github.com/NVlabs/ffhq-dataset

StyleGan2 by Nvidia labs

Given a dataset of images, how can we
train a neural network to be able to
generate realistic fakes?

random
noise $Z$ $\longrightarrow$  $\longrightarrow$ fake image
$X = f(Z)$

generator network

image
$x$ $\longrightarrow$  $\longrightarrow$ realism score
$\lambda(x)$

discriminator
network

- train the discriminator to distinguish genuine images from fake
- train the generator to fool the discriminator

Given a dataset of images, how can we train a neural network to be able to generate realistic fakes?



random noise $Z$ → generator network → fake image $X = f(Z)$

image $x$ → discriminator network → realism score $\lambda(x)$

To understand what's going on, it's useful to restate the problem using in the language of random variables.

Given a dataset of images, how can we train a neural network to be able to generate realistic fakes?

random noise $Z$ → generator network weights $\theta$ → fake image $X = f_\theta(Z)$

```
def gen_face(θ):
    Z = random.random()
    X = f_θ(Z)
    return X
```

To understand what's going on, it's useful to restate the problem using in the language of random variables.

# WHAT IS A RANDOM VARIABLE?

```python
def gen_geom(p):
    Z = random.random()
    λ = - math.log(1-p)
    X = math.ceil(-math.log(Z)/λ)
    return X
```

① Figure out the probability mass function

$$\mathbb{P}(X = x) = (1-p)^{x-1} \, p$$

$$\text{for} \quad x \in \{1, 2, 3, \cdots\}$$

# HOW DO WE LEARN FROM DATA?

Suppose we're given a dataset $\{x_1, x_2, \ldots, x_n\}$ and we want to tune our random variable generator by choosing $p$ to match the dataset as closely as possible. How can we do this?

② what's the probability that our rng generated $\{x_1, \cdots, x_n\}$?

$$\mathbb{P}(\text{gen. this dataset}) = \mathbb{P}(X=x_1) \times \mathbb{P}(X=x_2) \times \cdots \times \mathbb{P}(X=x_n)$$

$$= \left(\frac{p}{1-p}\right)^n (1-p)^{\sum x_i}$$

③ Choose $p$ to make this probability as large as possible. $\quad \hat{p} = \dfrac{n}{\sum x_i}$

Given a dataset of images, how can we train a neural network to be able to generate realistic fakes?

random noise $Z$ → generator network, weights $\theta$ → fake image $X = f_\theta(Z)$

```
def gen_face(θ):
    Z = random.random()
    X = f_θ(Z)
    return X
```

*We have no hope of calculating a formula for this. Nevertheless…*

$X$ is a random variable. Let its prob. mass func. be $\mathbb{P}(X = x) = p_\theta(x).$

We're given a dataset $\{x_1, x_2, \ldots, x_n\}$ of faces, and we our generator's output to match this dataset as closely as possible. We should find the $\theta$ that maximizes

$$p_\theta(x_1) \times p_\theta(x_2) \times \cdots \times p_\theta(x_n)$$

Or, equivalently, let's pick $\theta$ to maximize

$$\mathcal{V}(\theta) := \frac{1}{n} \sum_{i=1}^{n} \log p_\theta(x_i)$$

## TRAINING GOAL

pick $\theta$ to maximize $\mathcal{V}(\theta) = \frac{1}{n}\sum_{i=1}^{n}\log p_\theta(x_i) = \sum_{x \in \mathcal{X}} q_x \log p_\theta(x)$

$q_x = $ #times $x$ occurs$/n$

$\mathcal{X} = $ set of all possible images

## EQUIVALENT TRAINING GOAL (YOU'LL NEVER BELIEVE THIS ONE WEIRD TRICK)

maximize $\sum_x q_x \log r_x$

over $r \in \mathbb{R}^{\mathcal{X}}, \theta$

such that $r_x = p_\theta(x)$ for all $x$

## LAGRANGIAN

$\mathcal{L}(\theta, r; \lambda) = \sum_{x \in \mathcal{X}} q_x \log r_x - \sum_{x \in \mathcal{X}} \lambda_x (r_x - p_\theta(x))$

## LAGRANGIAN WEAK DUALITY

For any $\theta$ and any $\lambda$,

$\mathcal{V}(\theta) = \mathcal{L}(\theta, p_\theta; \lambda) \leq \max_{r \in \mathbb{R}^{\mathcal{X}}} \mathcal{L}(\theta, r; \lambda) = \begin{cases} \text{const} - \sum_x q_x \log \lambda_x + \sum_x p_\theta(x)\lambda_x, & \text{if } \lambda_x > 0 \text{ for all } x \\ \infty, & \text{otherwise} \end{cases}$

*with some simple optimization over r*

*What values of $\lambda$ would make this upper bound on $\gamma(\theta)$ as tight as possible?*

*$\lambda_x$ should be larger where $q_x$ is larger*

*$\lambda_x$ should be smaller where $p_\theta(x)$ is larger*

*$\lambda_x$ should be $> 0$ for all $x$*

pick $\theta$ to maximize $\mathcal{V}(\theta) = \dfrac{1}{n}\sum_{i=1}^{n} \log p_\theta(x_i) = \sum_{x \in \mathcal{X}} q_x \log p_\theta(x)$

$q_x = \#\text{times } x \text{ occurs}/n$

$\mathcal{X} = $ set of all possible images



image $x$ → discriminator network → realism score $\lambda(x)$

We aimed to train the discriminator so that $\lambda(x)$ larger if $x$ is genuine, smaller if it's faked.

## LAGRANGIAN WEAK DUALITY

For any $\theta$ and any $\lambda$,

$\mathcal{V}(\theta) = \mathcal{L}(\theta, p_\theta; \lambda) \leq \max_r \mathcal{L}(\theta, r; \lambda) = \begin{cases} \text{const} - \sum_x q_x \log \lambda_x + \sum_x p_\theta(x)\lambda_x, & \text{if } \lambda_x > 0 \text{ for all } x \\ \infty, & \text{otherwise} \end{cases}$

with some simple optimization over r

What values of $\lambda$ would make this upper bound on $\gamma(\theta)$ as tight as possible?

$\lambda_x$ should be larger where $q_x$ is larger

$\lambda_x$ should be smaller where $p_\theta(x)$ is larger

$\lambda_x$ should be $> 0$ for all $x$

# FURTHER DETAILS

The score

$$\mathcal{V}(\theta) = \frac{1}{n}\sum_{i=1}^{n} \log p_\theta(x_i)$$

measures how well our generator performs. We have shown that for any $\theta$ and any $\lambda > 0$

$$\mathcal{V}(\theta) \leq \text{const} - \sum_x q_x \log \lambda_x + \sum_x p_\theta(x)\lambda_x.$$

Let's propose a discriminator neural network $x \mapsto d_\phi(x) > 0$ for computing $\lambda_x = d_\phi(x)$, and try to find good neural network weights $\phi$. Since our Lagrangian bound is true for any $\theta$ and any $\lambda > 0$, it follows that

$$\max_\theta \mathcal{V}(\theta) \leq \text{const} + \max_\theta \min_\phi \left\{ -\sum_x q_x \log d_\phi(x) + \sum_x p_\theta(x)d_\phi(x) \right\}$$
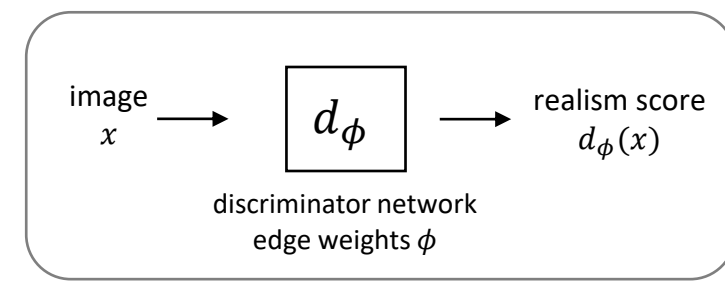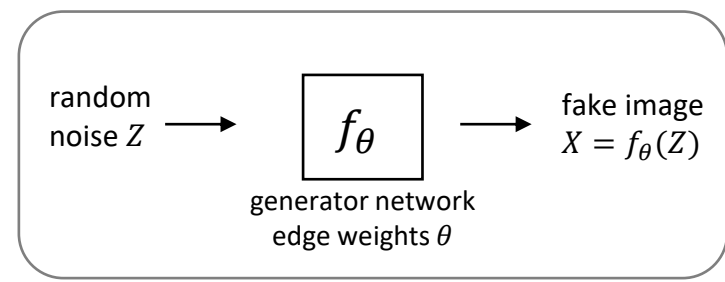
We can approximate the $\{\cdot\}$ term by

$$L(\theta, \phi) := -\frac{1}{n}\sum_{i=1}^{n} \log d_\phi(x_i) + \frac{1}{n}\sum_{i=1}^{n} d_\phi(f_\theta(z_i))$$

where $\{x_1, \dots, x_n\}$ are the images in the training dataset and $\{z_1, \dots, z_n\}$ are randomly generated noise terms. Training consists in using gradient descent to compute

$$\max_\theta \min_\phi L(\theta, \phi).$$

This doesn't actually compute $\max_\theta \mathcal{V}(\theta)$, it only computes an upper bound, but hopefully the upper bound is reasonably tight and we end up learning a generator with a good score.
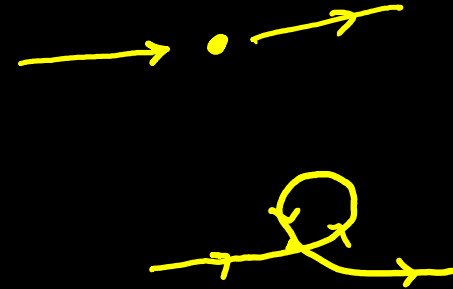


random noise $Z$ → $f_\theta$ → fake image $X = f_\theta(Z)$

generator network
edge weights $\theta$



image $x$ → $d_\phi$ → realism score $d_\phi(x)$

discriminator network
edge weights $\phi$

# Other optimizations

## Shortest path problem

Given a directed graph where the weight of edge $v \rightarrow w$ is $c_{vw}$, find the minimum-weight path from a start vertex $s$ to a destination vertex $t$
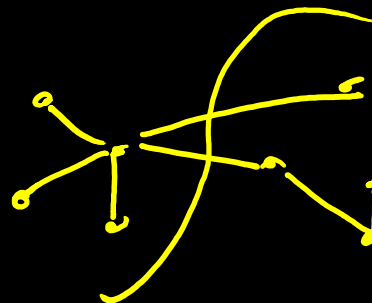
minimize    $\sum_{v \rightarrow w} x_{vw} \, c_{vw}$

over        $x \in \{0,1,2,\dots\}^E$    $\text{NOT} \quad x \in \{0,1\}^E ,$
$$\left( \text{since paths might re-use edges} \right)$$

such that    net flow out of $s$ is 1
and         flow is conserved at all vertices in $V \backslash \{s, t\}$

# Minimum spanning tree problem

Given an undirected graph where the weight of edge $v \leftrightarrow w$ is $c_{vw}$, find a spanning tree of minimum weight

A tree is an __acyclic connected graph__



minimize $\quad \sum_{v \leftrightarrow w} x_{vw} \, c_{vw}$

over $\quad\quad x \in \{0,1\}^E$

such that $\quad \sum_{v \leftrightarrow w} x_{vw} = |V| - 1$ &larr; implies acyclic
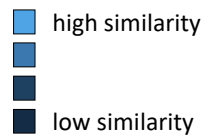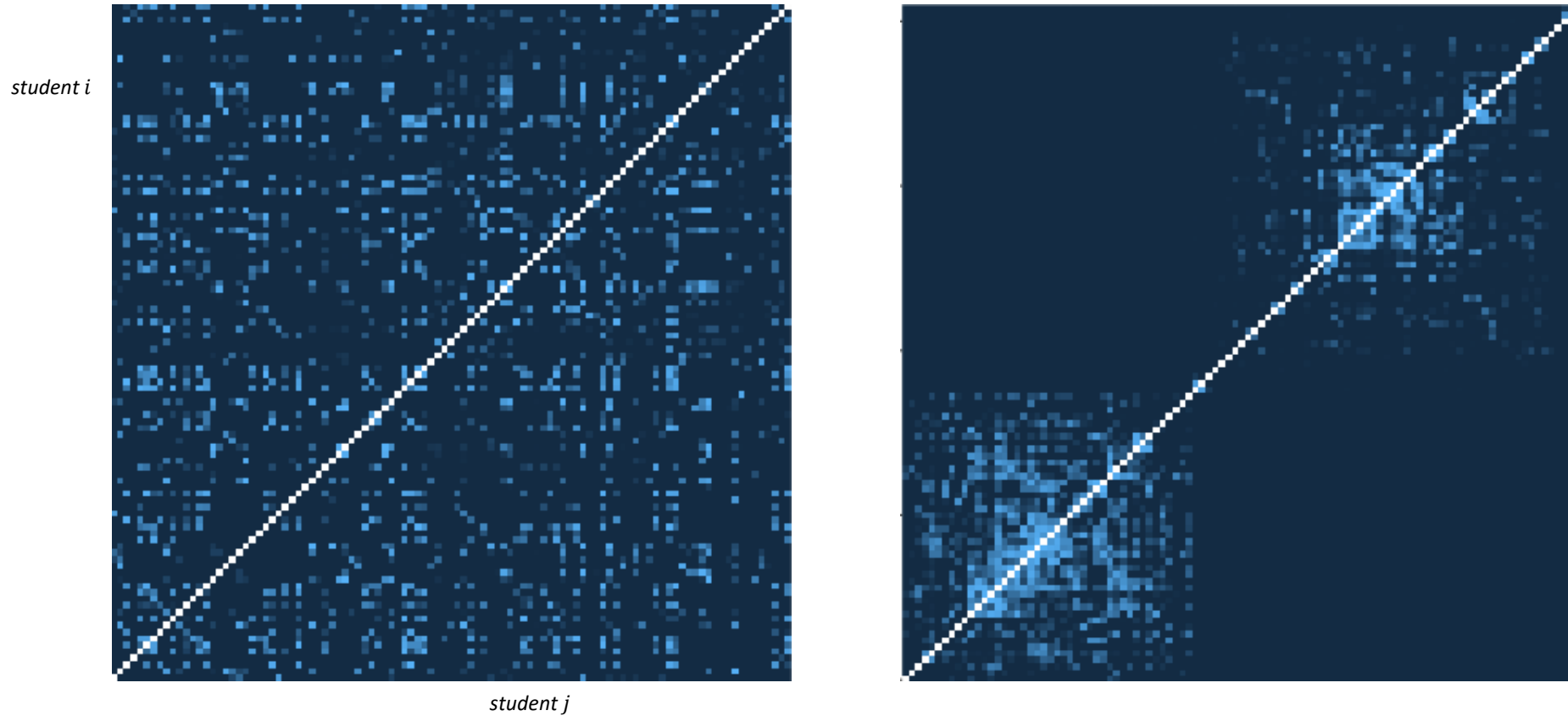
and $\quad\quad \sum_{\substack{v \leftrightarrow w: \\ v \in S, w \in \bar{S}}} x_{vw} \geq 1$ for all sets $S \subset V, S \neq V, S \neq \emptyset$ &larr; connected

# Challenge: find an ordering for all students so that students with similar Tick 1 code are close to each other

Similarity matrix of Tick 1 code



*student i*

*student j*

high similarity

low similarity

| | score on training data (2021 tick1) | score on holdout data (2022 tick1) |
|---|---|---|
| Tunan Shi (Sidney Sussex) | 77.4 | 62.7 |
| Kuba Bachurski (Trinity) | 76.9 | 62.4 |
| Mark Li (Corpus Christi) | 76.6 | 61.8 |
| Andy Zhou (Queens') | 73.0 | 60.2 |
| Cheuk Kit Lee (Downing) | 77.2 | timeout |
| Jiayou Song (Robinson) | 75.0 | timeout |

Andrej Karpathy @
@karpathy

Gradient descent can write code better than you. I'm sorry.

3:56 PM - 4 Aug 2017

343 Retweets  1,161 Likes

72    343    1.2K

Software 1.0 is code we write. Software 2.0 is code written by the optimization based on an evaluation criterion (such as "classify this training data correctly"). It is likely that any setting where the program is not obvious but one can repeatedly evaluate the performance of it (e.g. — did you classify some images correctly? do you win games of Go?) will be subject to this transition, because the optimization can find much better code than what a human can write.

❖ Can I express my task as an optimization problem?

❖ … that can be solved with off-the-shelf optimizers?

❖ If I can't, is there an adjacent problem that's more amenable?