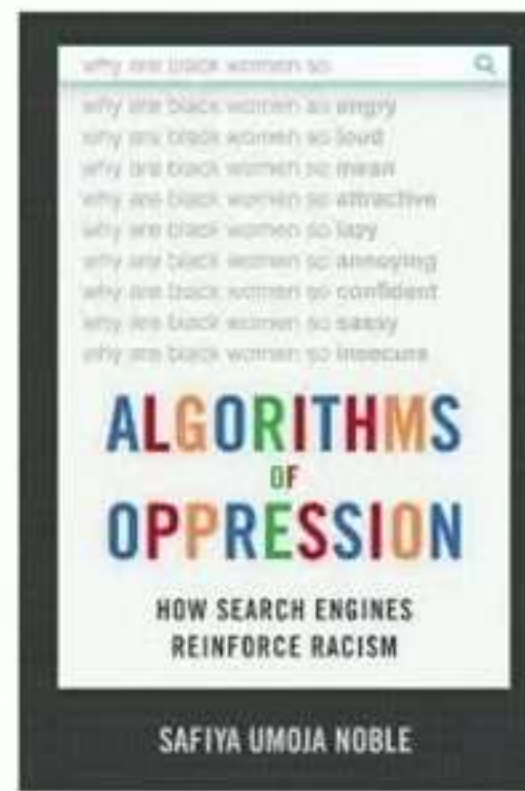
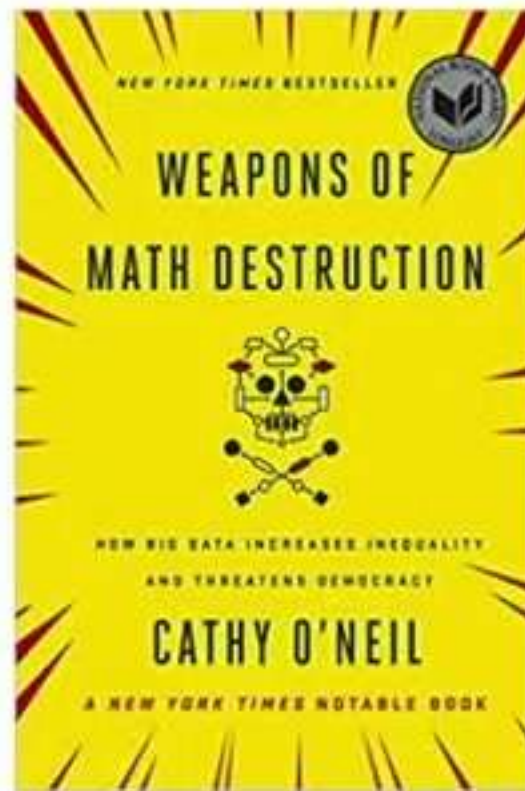


michael kearns + aaron roth

MICROSOFT RESEARCH AI
NOVEMBER 12, 2019

the ethical algorithm

the science of
socially aware algorithm design



ETHICAL ALGORITHMS?



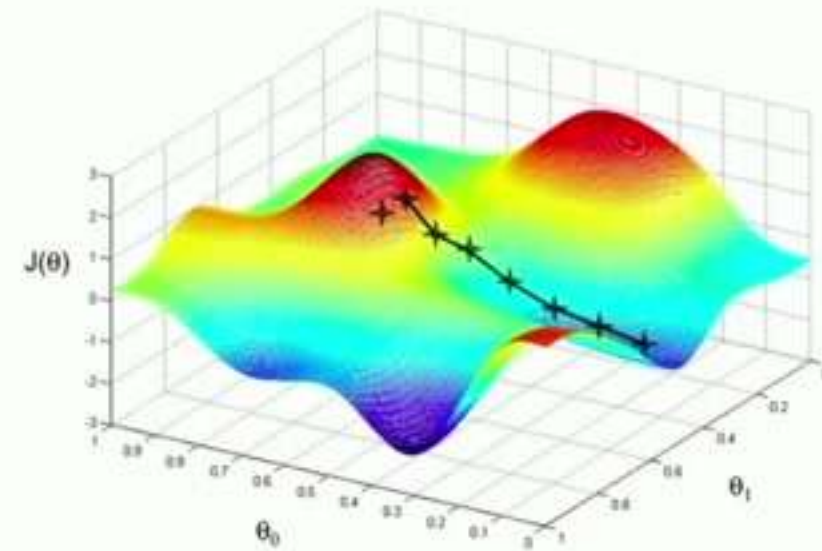
ETHICAL ALGORITHMS?



ALGORITHMS: HARD TO ASSIGN BLAME

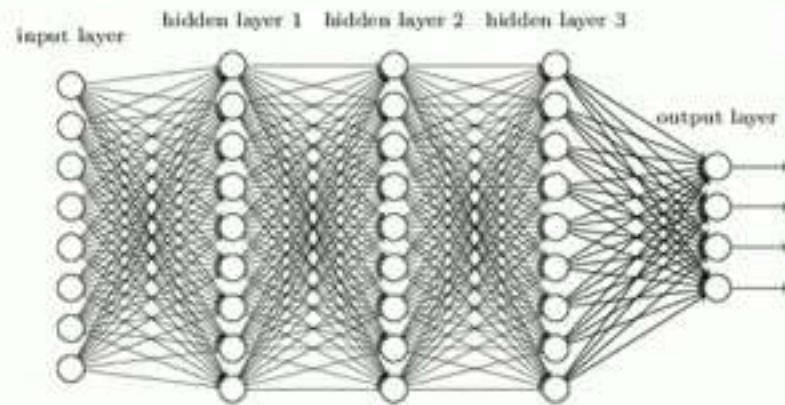
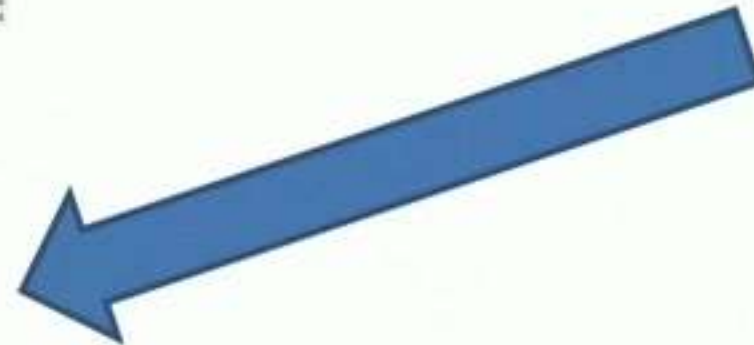
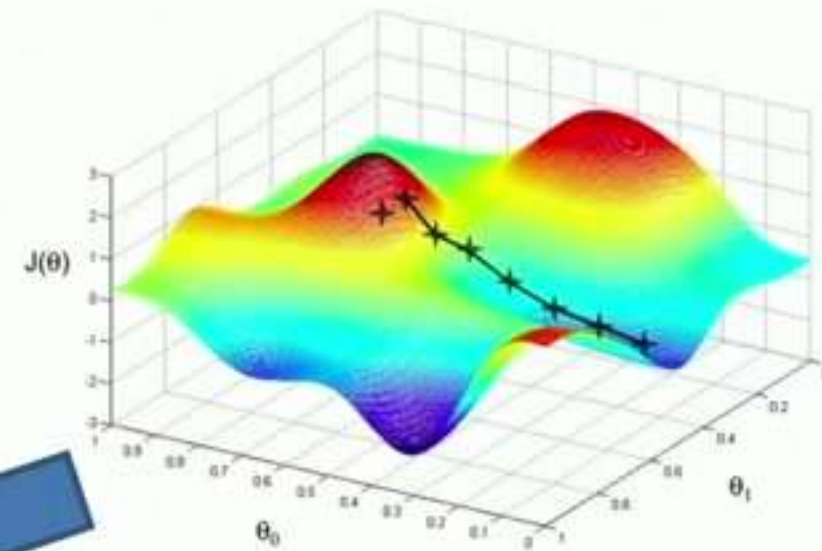
ALGORITHMS: HARD TO ASSIGN BLAME

A	1420	175	240	169	107	39	308	0.26	0.23	394	138	340	1000
A	1512	176	214	112	98	496	276	0.26	1.26	436	105	54	1000
A	1516	238	207	164	101	34	338	0.11	0.25	346	101	117	1195
A	1437	196	210	164	111	496	346	0.24	2.10	74	100	149	1000
A	1534	248	207	21	119	34	346	0.24	1.01	431	104	100	751
A	1412	176	240	162	112	327	326	0.24	1.37	474	100	100	1400
A	1436	187	240	144	96	110	330	0.11	1.08	326	100	100	1000
A	1436	215	203	174	121	36	331	0.21	1.21	505	100	100	1000
A	1490	144	217	14	97	20	336	0.01	1.09	51	100	100	1000
A	1236	135	227	16	96	296	310	0.01	1.00	722	101	100	1000
A	141	216	21	10	96	296	310	0.01	0.00	576	100	117	1010
A	1412	146	230	163	96	310	340	0.00	1.07	0	117	100	1000
A	1274	171	247	16	96	310	276	0.00	1.01	68	115	25	1000
A	1476	175	236	114	91	31	346	0.43	0.00	34	100	273	1100
A	1436	187	236	110	96	310	344	0.00	0.96	74	117	0	1000
A	1340	187	217	112	110	346	330	0.01	1.46	71	100	146	1000
A	143	130	270	36	136	34	314	0.01	1.07	42	107	146	1000
A	1040	157	200	30	119	296	34	0.44	1.71	44	115	217	1100
A	1416	136	246	163	106	310	336	0.01	1.00	47	121	100	1000
A	1344	21	236	162	116	27	330	0.01	1.04	61	106	336	800



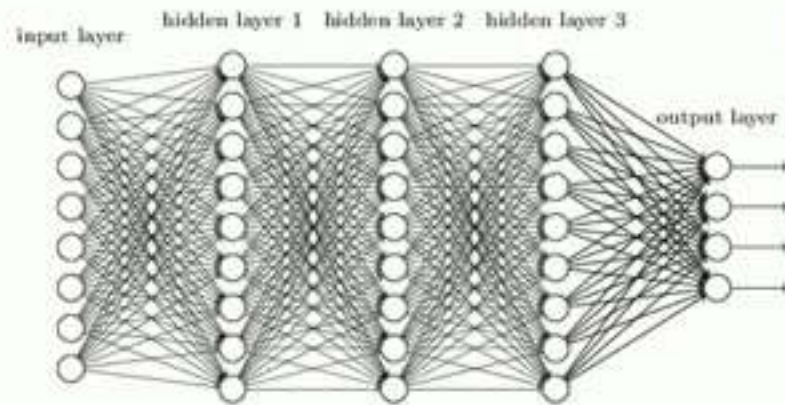
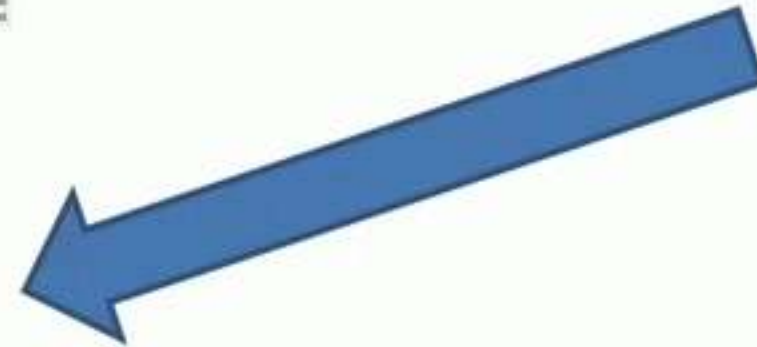
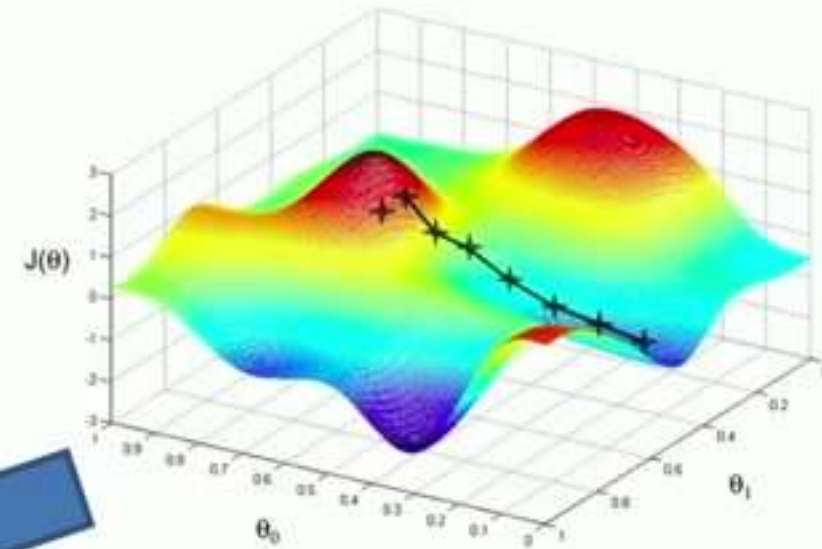
ALGORITHMS: HARD TO ASSIGN BLAME

A	1420	175	240	143	107	24	308	0.28	0.28	394	1.24	3.62	1000
A	1312	178	214	117	100	495	178	0.26	1.26	408	1.05	5.4	1000
A	1416	188	207	104	101	110	124	0.1	0.1	388	1.01	3.17	1194
A	1437	190	210	104	111	890	340	0.24	0.19	7.8	0.89	1.49	1480
A	1334	190	207	20	110	114	100	0.28	1.82	4.12	1.04	1.00	700
A	1412	178	240	143	112	107	324	0.24	1.87	4.15	1.01	1.00	1400
A	1438	187	240	144	98	110	100	0.1	1.88	1.01	1.01	1.00	1380
A	1436	210	201	172	121	110	121	0.21	1.21	0.89	1.00	1.00	1200
A	1430	184	217	14	97	110	100	0.10	1.80	0.1	1.00	1.00	1000
A	1388	190	227	10	90	208	110	0.12	1.80	1.02	1.01	1.00	1040
A	1411	210	110	10	100	438	100	0.10	0.28	1.01	1.00	1.17	1010
A	1412	148	210	103	90	110	100	0.28	1.87	0.117	0.117	1.00	1000
A	1376	172	241	10	88	110	100	0.28	1.81	0.8	1.01	1.19	1330
A	1419	175	238	114	91	110	100	0.43	0.85	0.4	1.00	1.71	1180
A	1438	187	238	110	100	110	100	0.28	0.86	1.1	1.2	1	1040
A	1380	187	217	112	110	100	100	0.1	1.48	1.1	1.00	1.00	1000
A	1412	180	270	10	130	110	100	0.10	1.87	0.1	1.01	1.00	1200
A	1388	187	240	10	119	100	100	0.1	1.72	0.8	1.01	1.01	1100
A	1416	188	240	143	108	110	100	0.10	1.88	0.7	1.01	1.00	1000
A	1384	11	138	143	118	11	100	0.17	1.88	0.1	0.86	0.38	800



ALGORITHMS: HARD TO ASSIGN BLAME

A	14.01	175	240	19.9	107	2.0	3.00	0.20	0.20	504	1.04	1.62	1000
A	15.12	176	214	11.2	100	4.06	2.76	0.26	1.26	4.00	1.05	5.4	1000
A	15.19	200	207	10.9	101	3.0	2.26	0.13	1.05	3.00	1.01	3.17	1100
A	14.97	190	210	10.9	111	4.00	3.00	0.20	2.00	3.00	1.00	3.00	1000
A	13.24	200	207	0	110	3.0	2.00	0.20	1.00	4.00	1.00	3.00	700
A	14.2	176	240	10.2	112	3.07	2.00	0.24	1.07	4.00	1.00	3.00	1000
A	14.00	107	240	10.0	90	3.0	3.00	0.2	1.00	3.00	1.00	3.00	1000
A	14.00	210	207	17.0	111	3.0	2.00	0.21	1.00	3.00	1.00	3.00	1000
A	14.00	104	217	14	97	3.0	3.00	0.20	1.00	3.00	1.00	3.00	1000
A	12.00	100	227	10	80	2.00	3.10	0.22	1.00	3.00	1.00	3.00	1000
A	14.1	210	210	10	100	3.00	3.00	0.22	1.00	3.00	1.00	3.17	1010
A	14.12	140	230	10.0	90	3.0	2.00	0.20	1.07	3.00	1.17	3.00	1000
A	13.70	170	240	10	80	3.0	2.70	0.20	1.01	3.00	1.10	2.9	1000
A	14.70	170	230	10.4	91	3.1	3.00	0.40	1.00	3.4	1.00	3.70	1100
A	14.00	107	200	11	100	3.0	3.00	0.20	1.00	3.00	1.0	3	1000
A	13.00	100	217	17.2	110	3.00	2.00	0.2	1.40	3.00	1.00	3.00	1000
A	14.2	100	270	10	100	3.0	3.14	0.20	1.07	3.00	1.07	3.00	1000
A	13.00	107	200	10	110	3.00	3.4	0.4	1.70	3.4	1.10	3.07	1100
A	14.10	100	240	10.0	100	3.0	3.00	0.20	1.00	3.00	1.00	3.00	1000
A	13.04	11	230	10.2	110	2.7	3.00	0.17	1.04	3.1	0.90	3.00	800



NEED TO EMBED SOCIAL VALUES IN ALGORITHMS

- Requires being precise about definitions, developing their consequences.
 - Privacy
 - Fairness
 - Accountability
 - Interpretability
 - Morality

“ANONYMIZED DATA ISN’T”

Name	Age	Gender	Zip Code	Smoker	Diagnosis
*	60-70	Male	191**	Y	Heart disease
*	60-70	Female	191**	N	Arthritis
*	60-70	Male	191**	Y	Lung cancer
*	60-70	Female	191**	N	Crohn’s Disease
*	60-70	Male	191**	Y	Lung Cancer
*	50-60	Female	191**	N	HIV
*	50-60	Male	191**	Y	Lyme Disease
*	50-60	Male	191**	Y	Seasonal Allergies
*	50-60	Female	191**	N	Ulcerative Colitis

Name	Age	Gender	Zip Code	Diagnosis
*	50-60	Female	191**	HIV
*	50-60	Female	191**	Lupus
*	50-60	Female	191**	Hip Fracture
*	60-70	Male	191**	Pancreatic Cancer
*	60-70	Male	191**	Ulcerative Colitis
*	60-70	Male	191**	Flu Like Symptoms

SMOKING AND CARCINOMA OF THE LUNG

PRELIMINARY REPORT

BY

RICHARD DOLL, M.D., M.R.C.P.

Member of the Statistical Research Unit of the Medical Research Council

AND

A. BRADFORD HILL, Ph.D., D.Sc.

Professor of Medical Statistics, London School of Hygiene and Tropical Medicine; Honorary Director of the Statistical Research Unit of the Medical Research Council

In England and Wales the phenomenal increase in the number of deaths attributed to cancer of the lung provides one of the most striking changes in the pattern of mortality recorded by the Registrar-General. For example, in the quarter of a century between 1922 and 1947 the annual number of deaths recorded increased from 612 to 9,287, or roughly fifteenfold. This remarkable increase is, of course, out of all proportion to the increase of population—both in total and, particularly, in its older age groups. Stocks (1947), using standardized death rates to allow for these population changes, shows the following trend: rate per 100,000 in 1901–20, males 1.1, females 0.7; rate per 100,000 in 1936–9, males 10.6, females 2.5. The rise seems to have been particularly rapid since the end of the first world war: between 1921–30 and 1940–4 the death rate of men at ages 45 and over increased sixfold and of women of the same ages approximately threefold. This increase is still continuing. It has occurred, too, in Switzerland, Denmark, the U.S.A., Canada, and Australia, and has been reported from Turkey and Japan.

Many writers have studied these changes, considering whether they denote a real increase in the incidence of the disease or are due merely to improved standards of diagnosis. Some believe that the latter factor can be regarded as wholly, or at least mainly, responsible—for example, Willis (1948), Clemmesen and Busk (1947), and Steiner (1944). On the other hand, Kennaway and Kennaway (1947) and Stocks (1947) have given good reasons for believing that the rise is at least partly real. The latter, for instance, has pointed out that "the increase of certified respiratory cancer mortality during the past 20 years has been as rapid in country districts as in the cities with the best diagnostic facilities, a fact which does not support the view that such increase merely reflects improved diagnosis of cases previously certified as bronchitis or other respiratory affections." He also draws attention to differences in mortality between some of the large cities of England and Wales, differences which it is difficult to explain in terms of diagnostic standards.

The large and continued increase in the recorded deaths even within the last five years, both in the national figures and in those from teaching hospitals, also makes it hard to believe that improved diagnosis is entirely responsible. In short, there is sufficient reason to reject that factor as the

whole explanation, although no one would deny that it may well have been contributory. As a corollary, it is right and proper to seek for other causes.

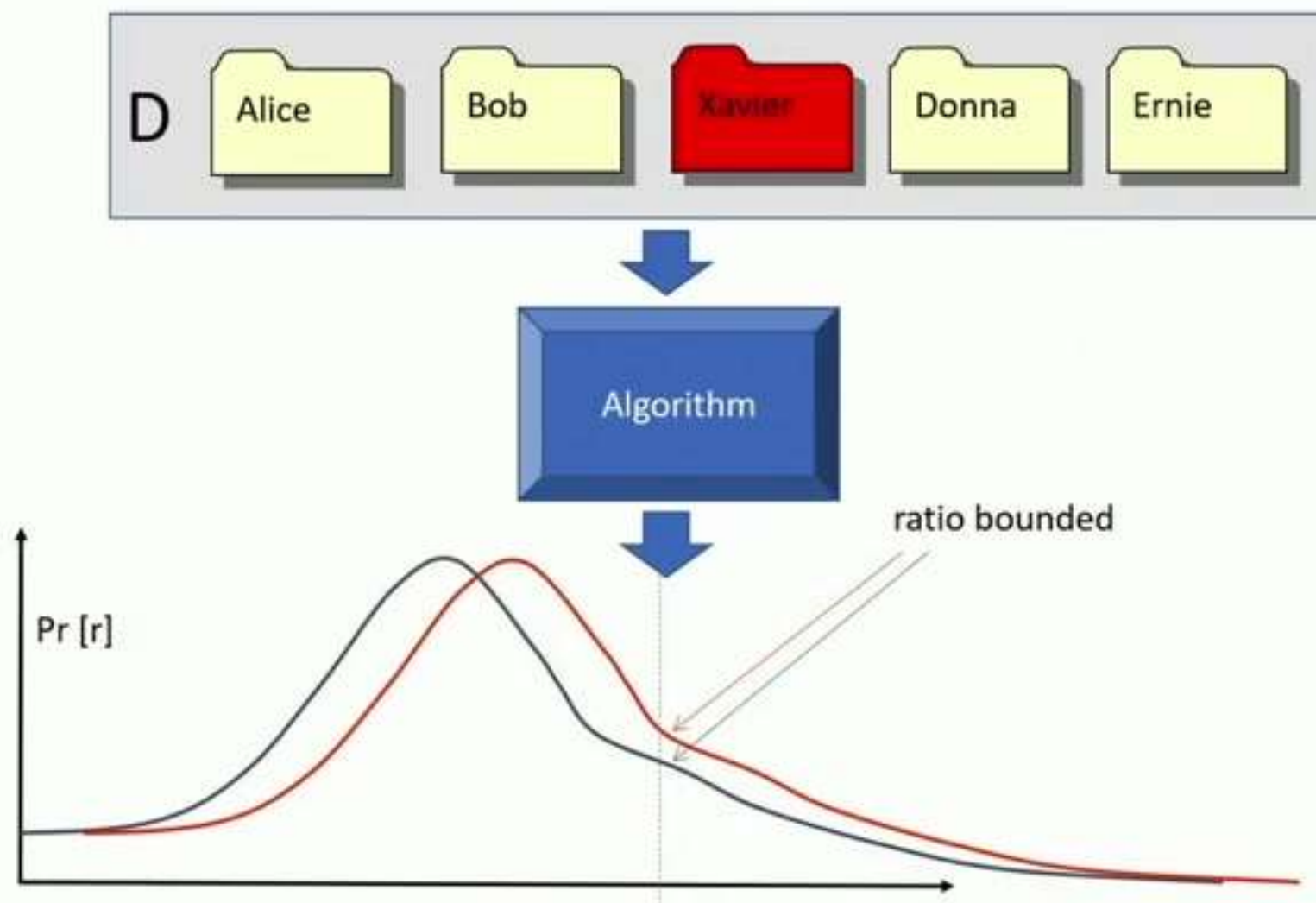
Possible Causes of the Increase

Two main causes have from time to time been put forward: (1) a general atmospheric pollution from the exhaust fumes of cars, from the surface dust of tarred roads, and from gas-works, industrial plants, and coal fires; and (2) the smoking of tobacco. Some characteristics of the former have certainly become more prevalent in the last 50 years, and there is also no doubt that the smoking of cigarettes has greatly increased. Such associated changes in time can, however, be no more than suggestive, and until recently there has been singularly little more direct evidence. That evidence, based upon clinical experience and records, relates mainly to the use of tobacco. For instance, in Germany, Müller (1939) found that only 3 out of 86 male patients with cancer of the lung were non-smokers, while 56 were heavy smokers, and, in contrast, among 86 "healthy men of the same age groups" there were 14 non-smokers and only 31 heavy smokers. Similarly, in America, Schrek and his co-workers (1950) reported that 14.6% of 82 male patients with cancer of the lung were non-smokers, against 23.9% of 522 male patients admitted with cancer of sites other than the upper respiratory and digestive tracts. In this country, Thelwall Jones (1949—personal communication) found 8 non-smokers in 82 patients with proved carcinoma of the lung, compared with 11 in a corresponding group of patients with diseases other than cancer; this difference is slight, but it is more striking that there were 28 heavy smokers in the cancer group, against 14 in the comparative group.

Clearly none of these small-scale inquiries can be accepted as conclusive, but they all point in the same direction. Their evidence has now been borne out by the results of a large-scale inquiry undertaken in the U.S.A. by Wynder and Graham (1950).

Wynder and Graham found that of 605 men with epidermoid, undifferentiated, or histologically unclassified types of bronchial carcinoma only 1.3% were "non-smokers"—that is, had averaged less than one cigarette a day for the last 20 years—whereas 51.2% of them had smoked more than 20 cigarettes a day over the same

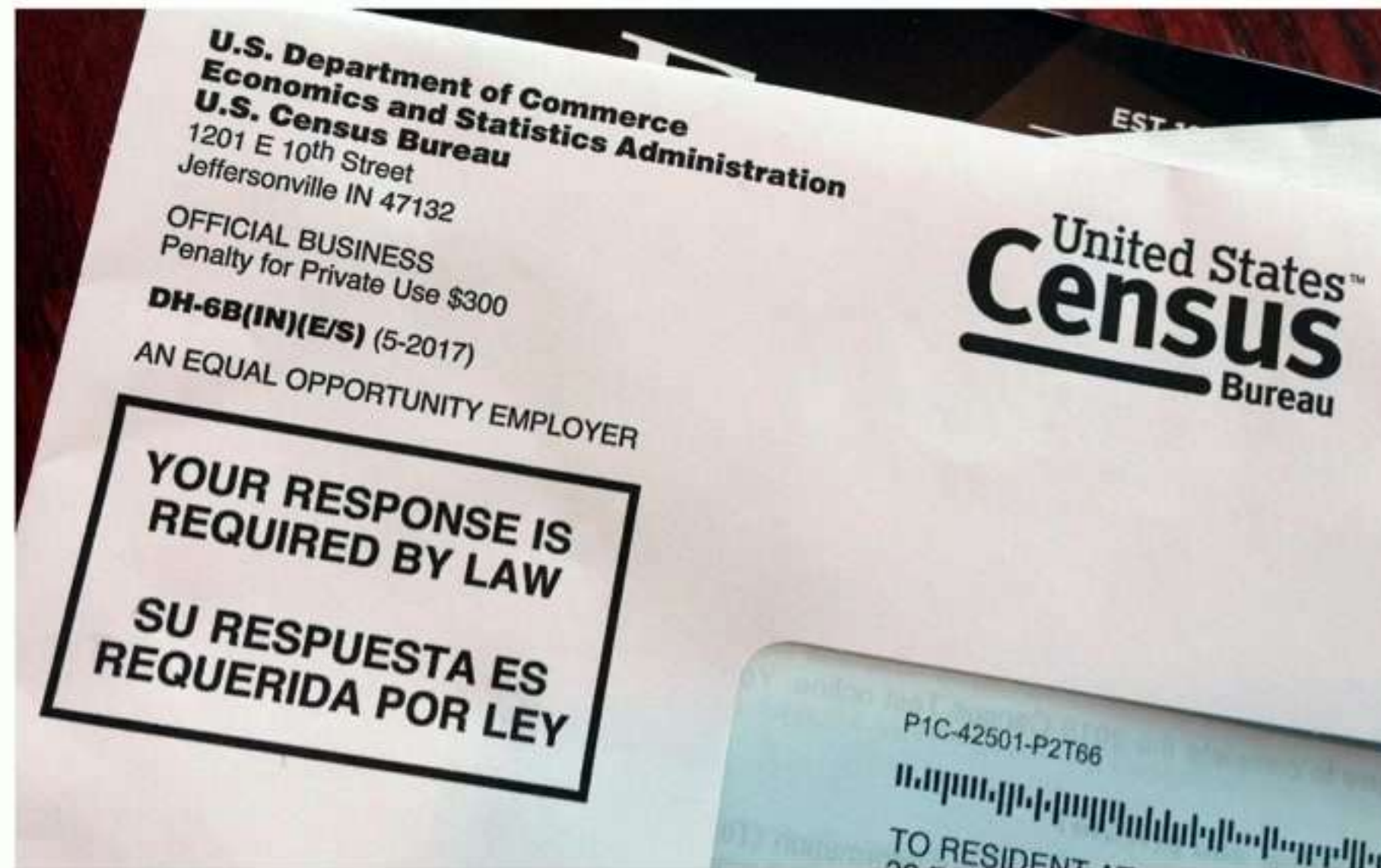
DIFFERENTIAL PRIVACY



TheUpshot

To Reduce Privacy Risks, the Census Plans to Report Less Accurate Data

Guaranteeing people's confidentiality has become more of a challenge, but some scholars worry that the new system will impede research.



A 2018 census test letter mailed to a resident in Providence, R.I. The nation's test run of the 2020 Census is in Rhode Island. Michelle R. Smith/Associated Press.

FAIRNESS: A WORK IN PROGRESS

- Don't agree on the definitions.
- Only beginning to understand *tradeoffs* between different kinds of fairness, and between fairness and accuracy.

FAIRNESS: A WORK IN PROGRESS

- Don't agree on the definitions.
- Only beginning to understand *tradeoffs* between different kinds of fairness, and between fairness and accuracy.

Why might machine learning be “unfair”?

THE WALL STREET JOURNAL.
English Edition • November 11, 2019 • Print Edition • Video

Home World U.S. Politics Economy Business Tech Markets Opinion Life & Arts Real Estate WSJ Magazine

PGIM GLOBAL REAL ESTATE FUND
★★★★★
Pursuing the world's best real estate opportunities.
PGIM INVESTMENTS

New York Regulator Probes UnitedHealth Algorithm for Racial Bias
Financial Services Department is investigating whether algorithm violates state antidiscrimination law

Bloomberg

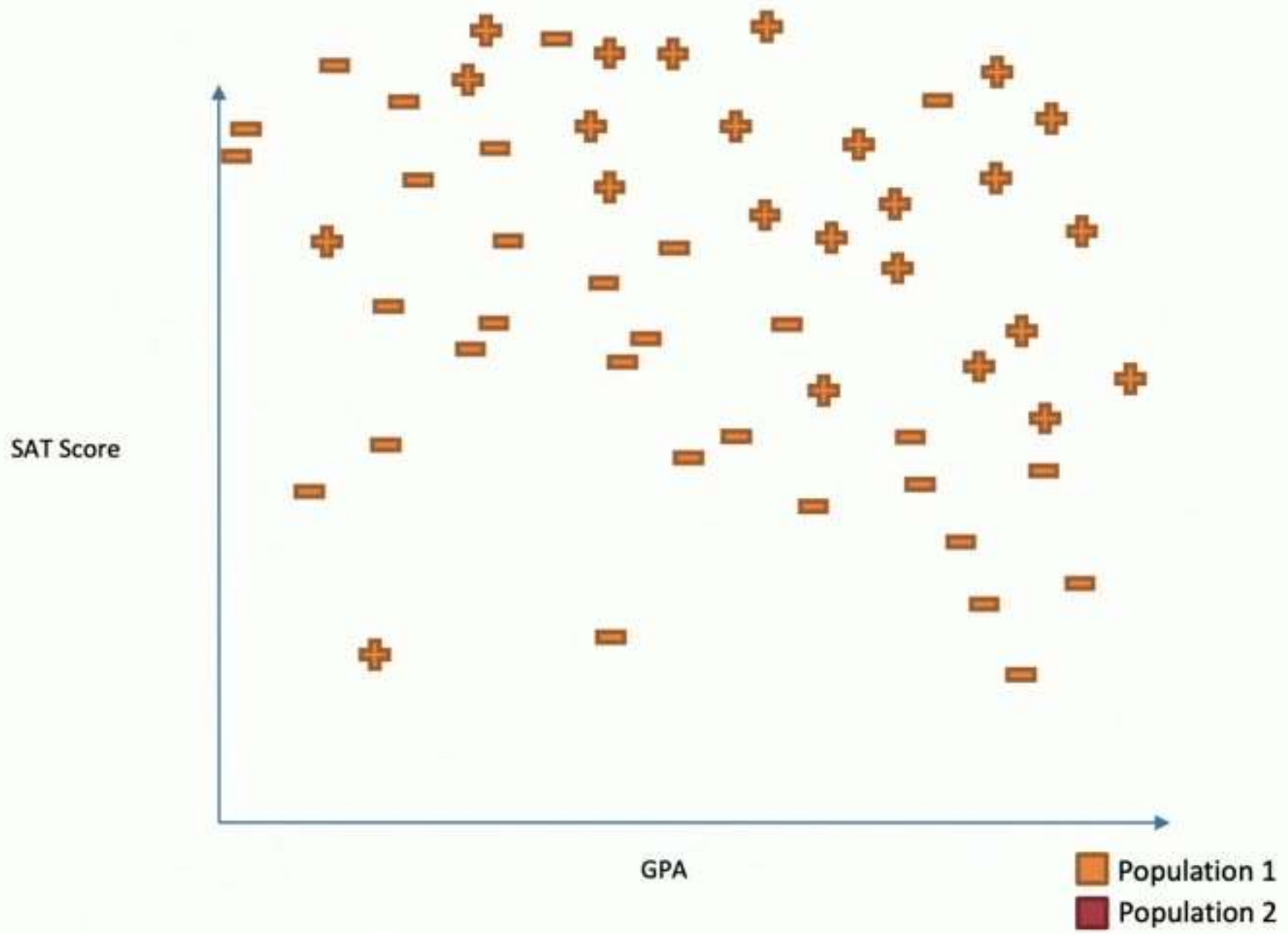
Business

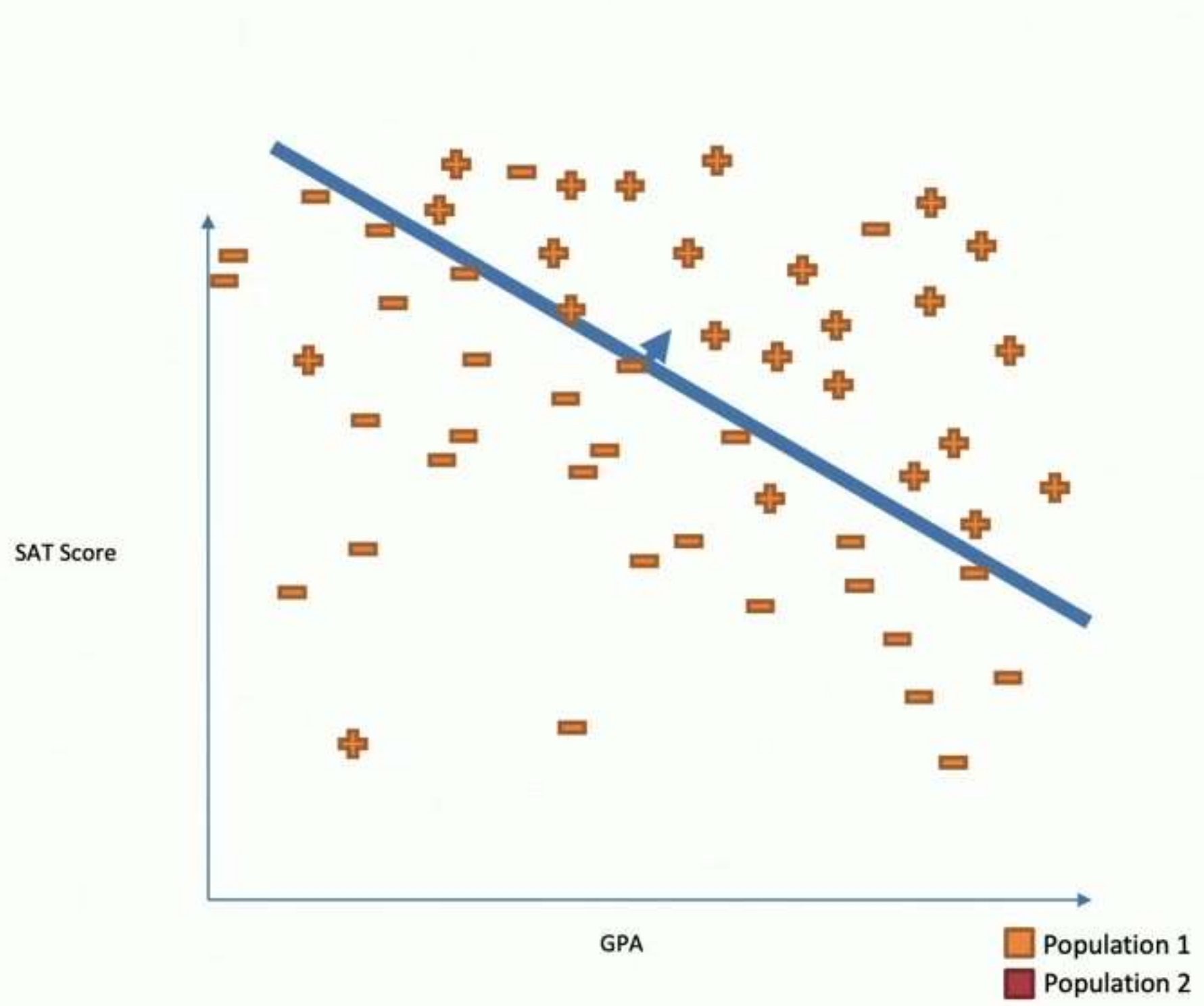
Viral Tweet About Apple Card Leads to Goldman Sachs Probe

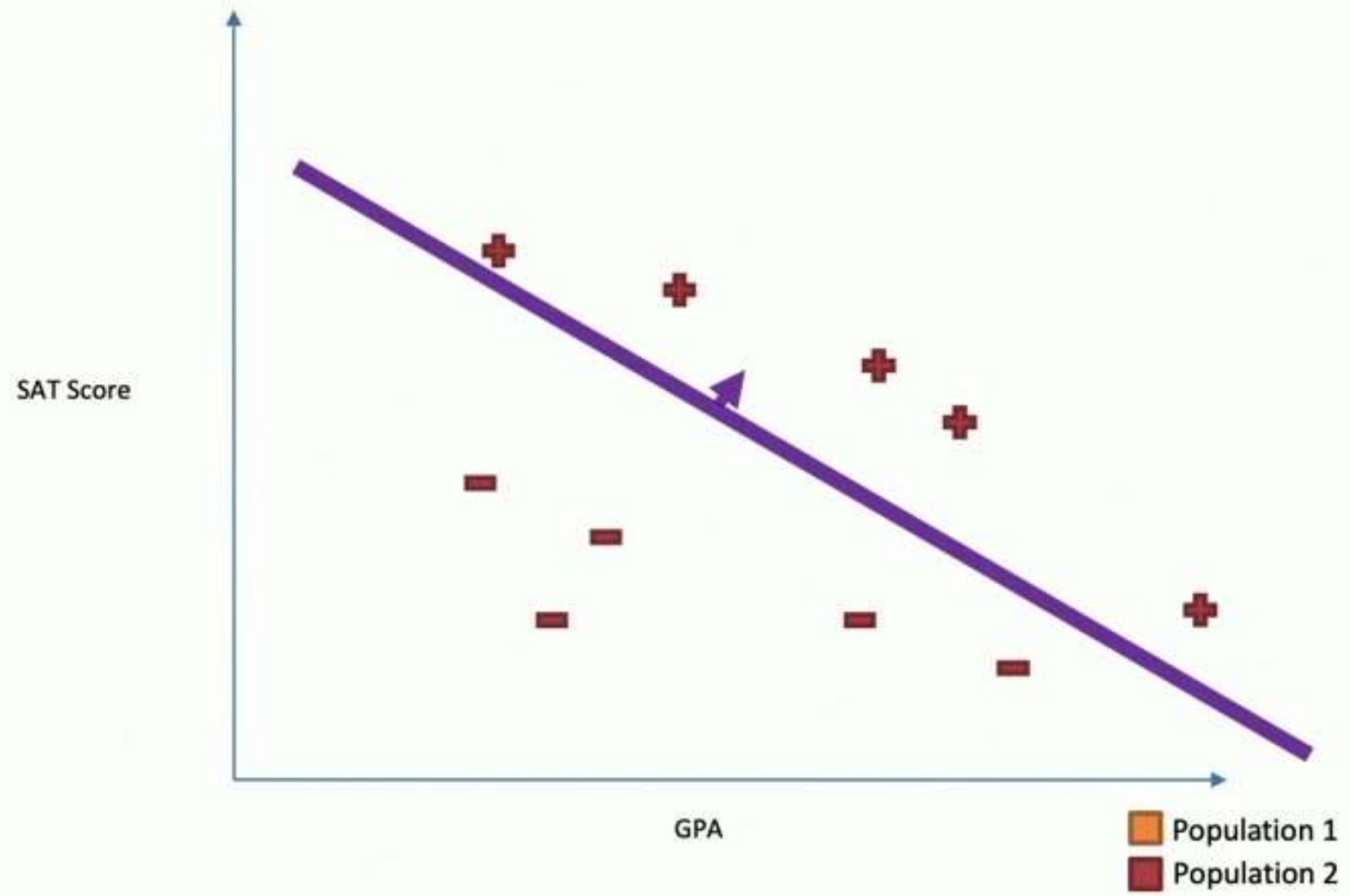
By Sridhar Natarajan and Shahien Nasiripour
November 9, 2019, 3:52 PM EST Updated on November 9, 2019, 8:53 PM EST

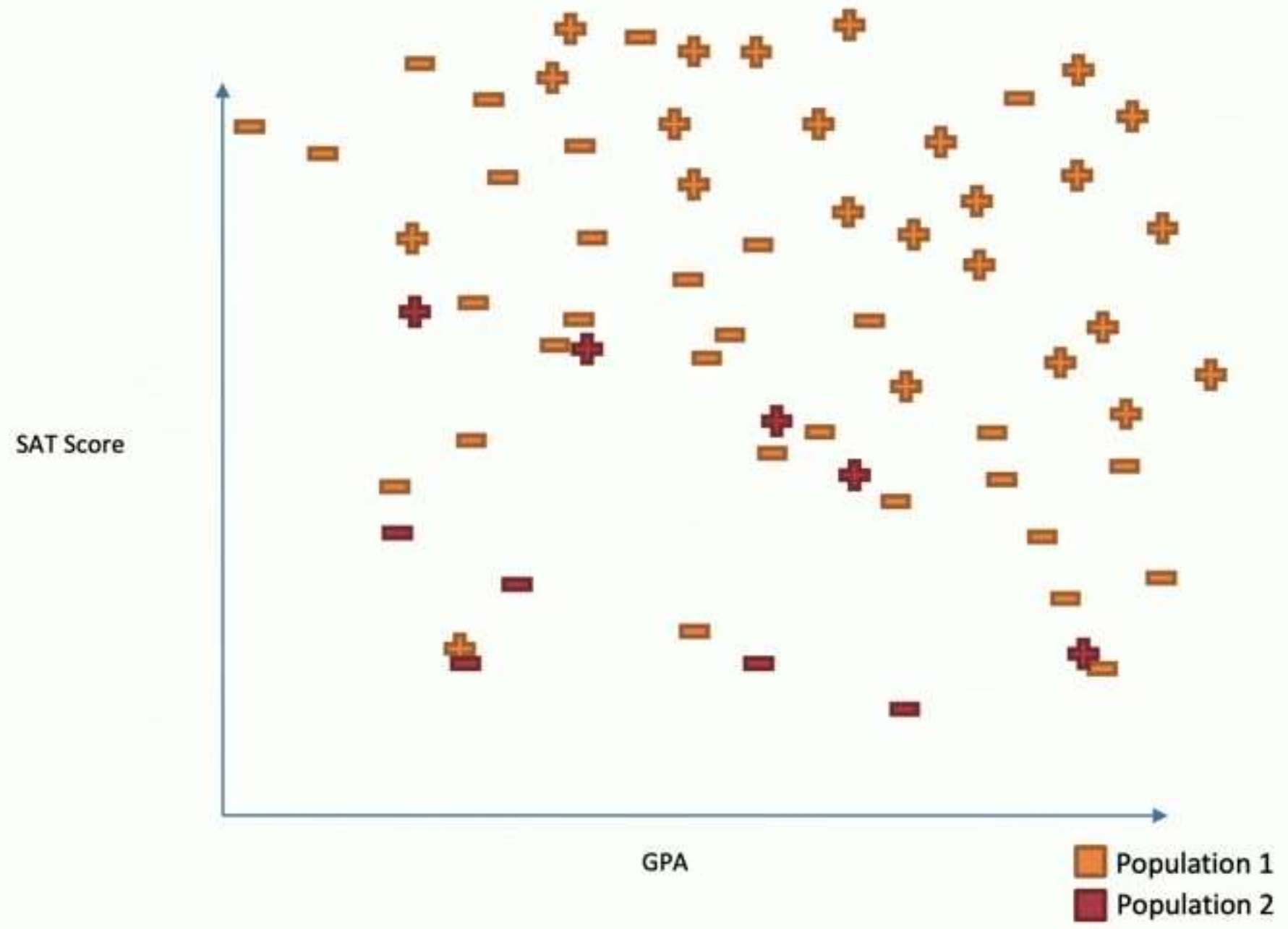
- ▶ Tech entrepreneur alleged inherent bias in algorithms for card
- ▶ The card is part of Goldman's new main street business lines

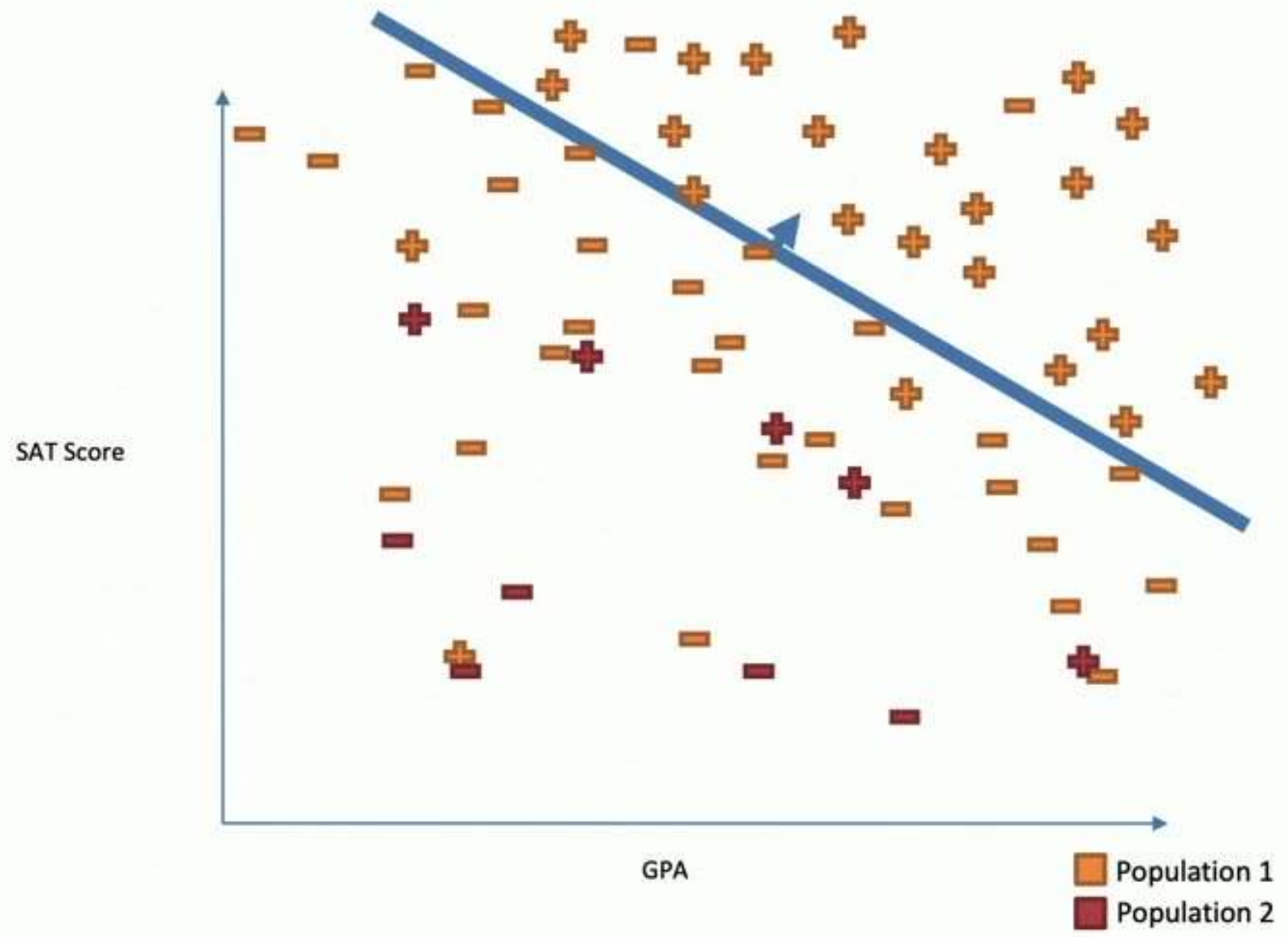
with IT infrastruc
powered by

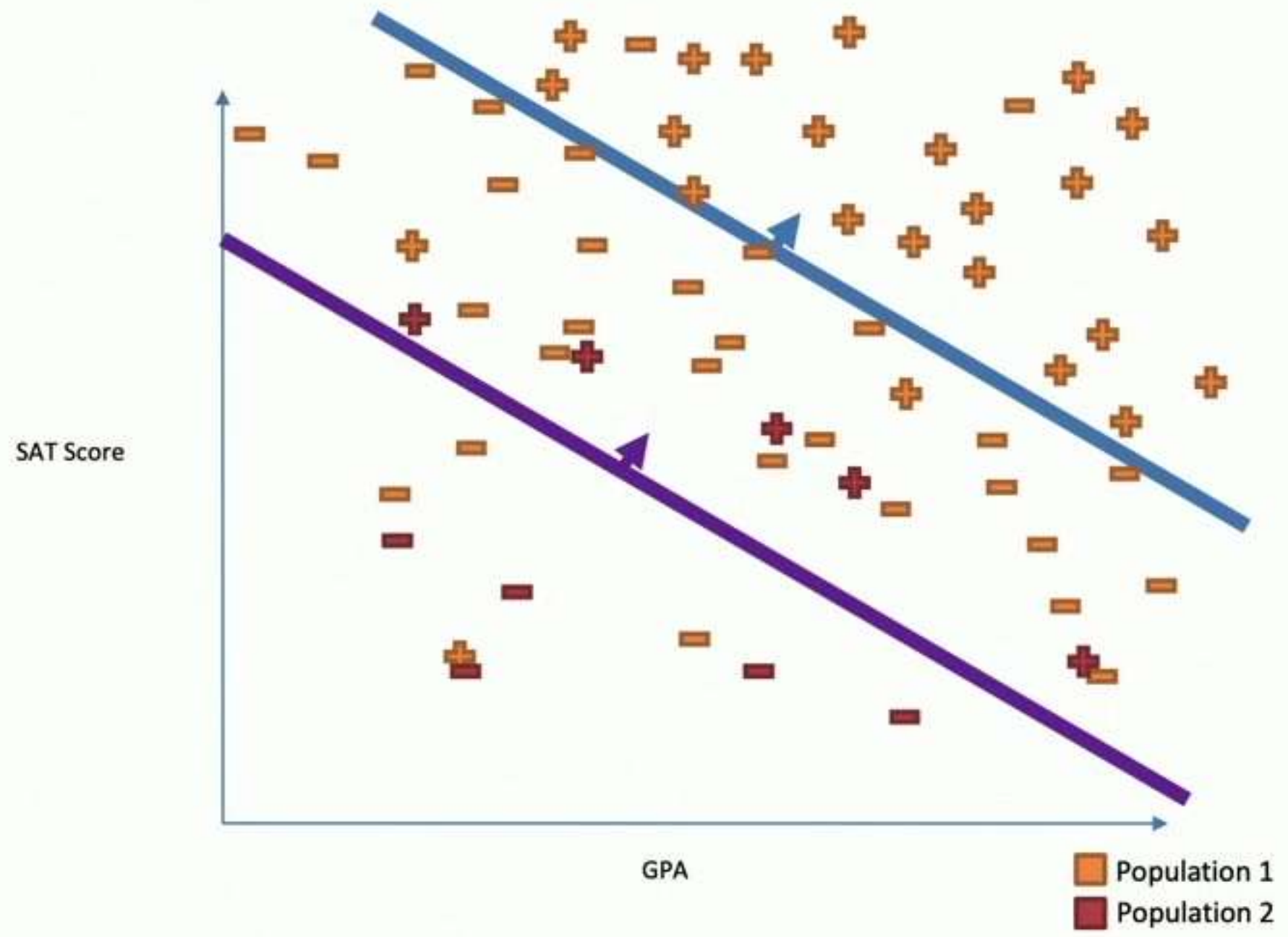


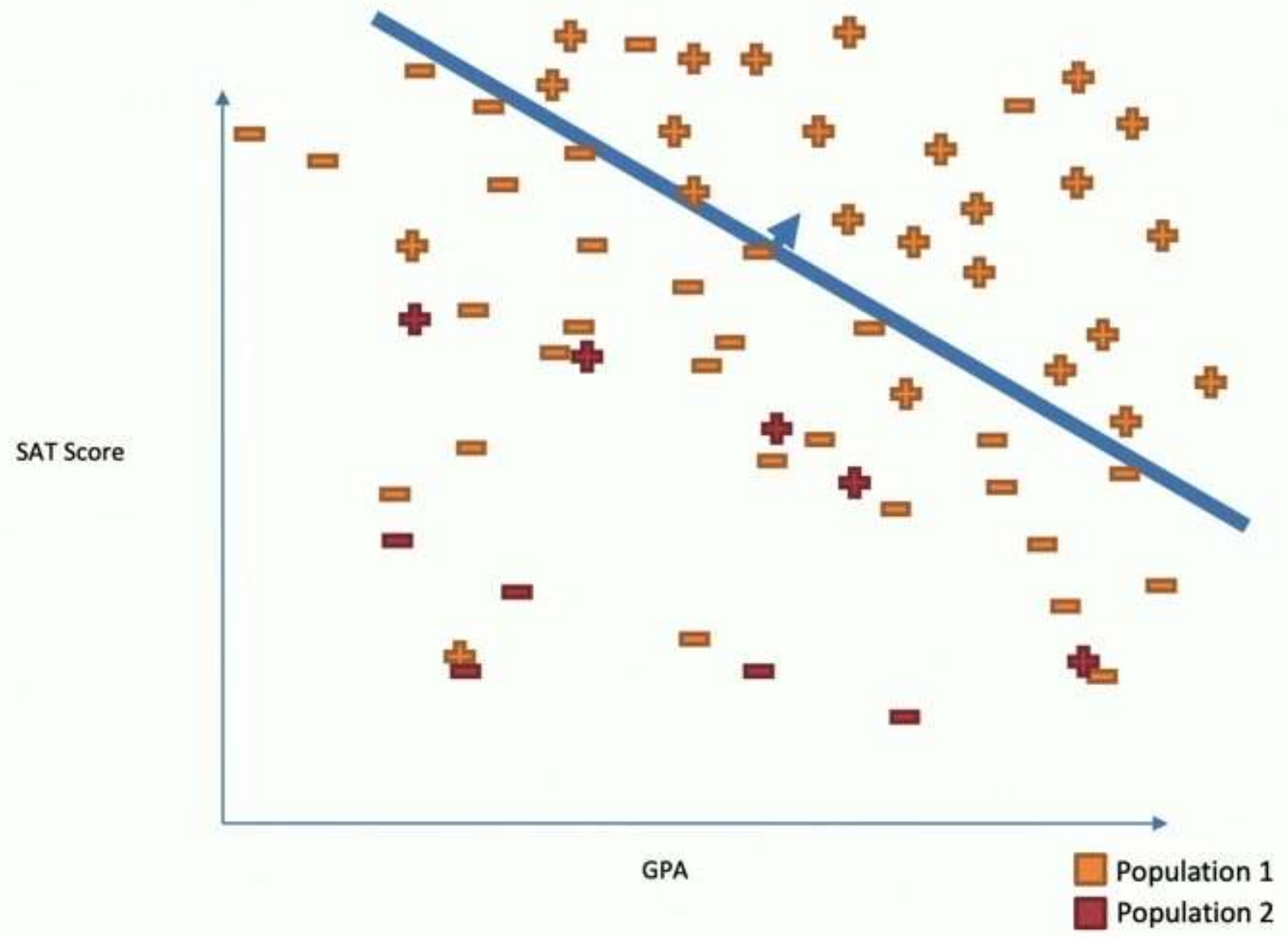


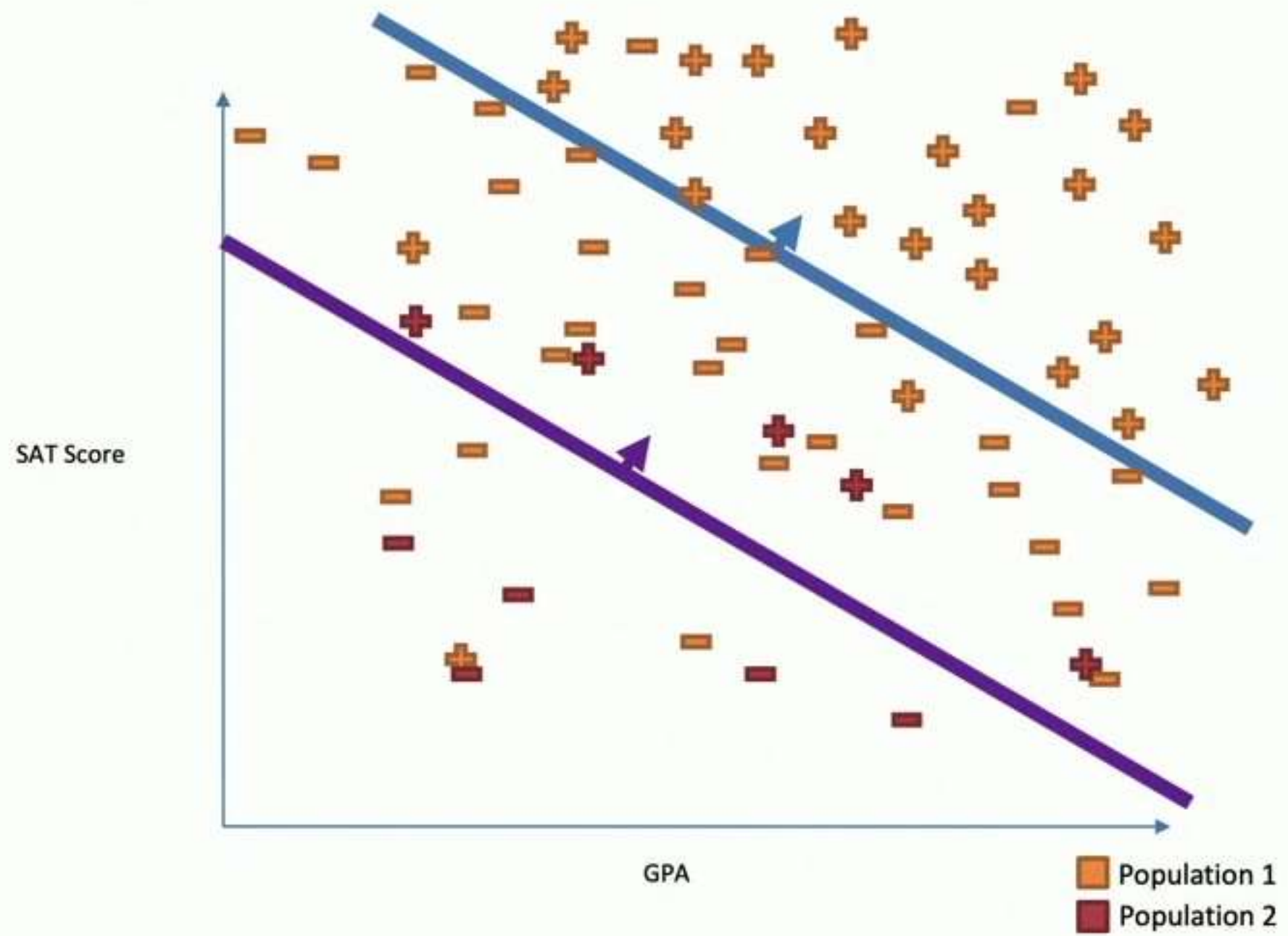




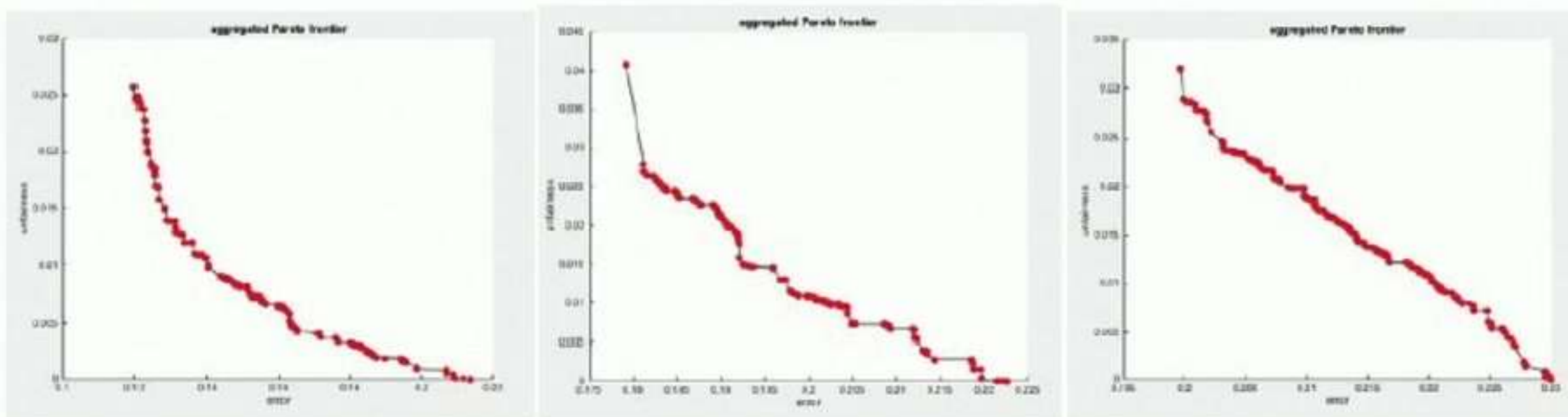




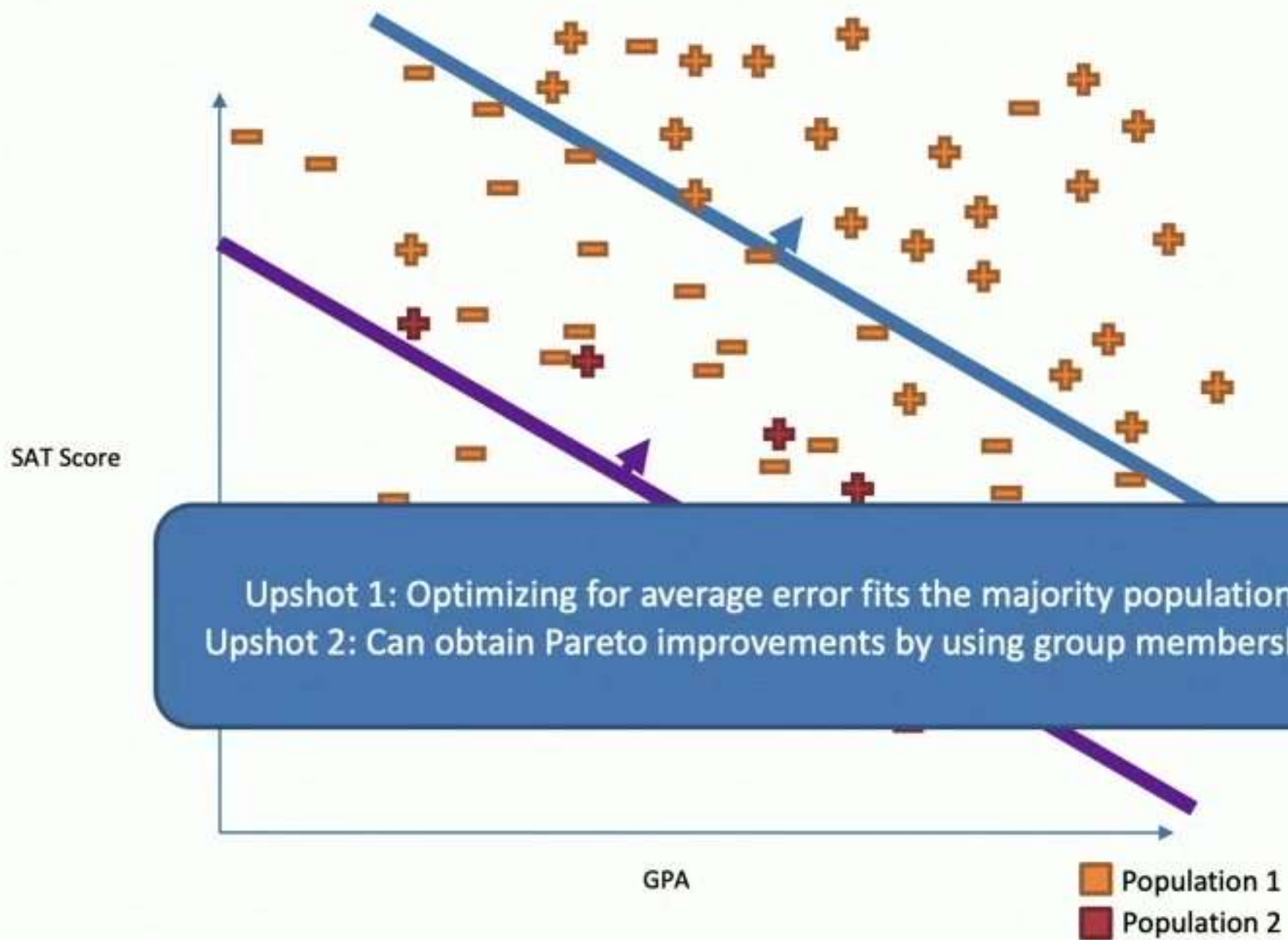




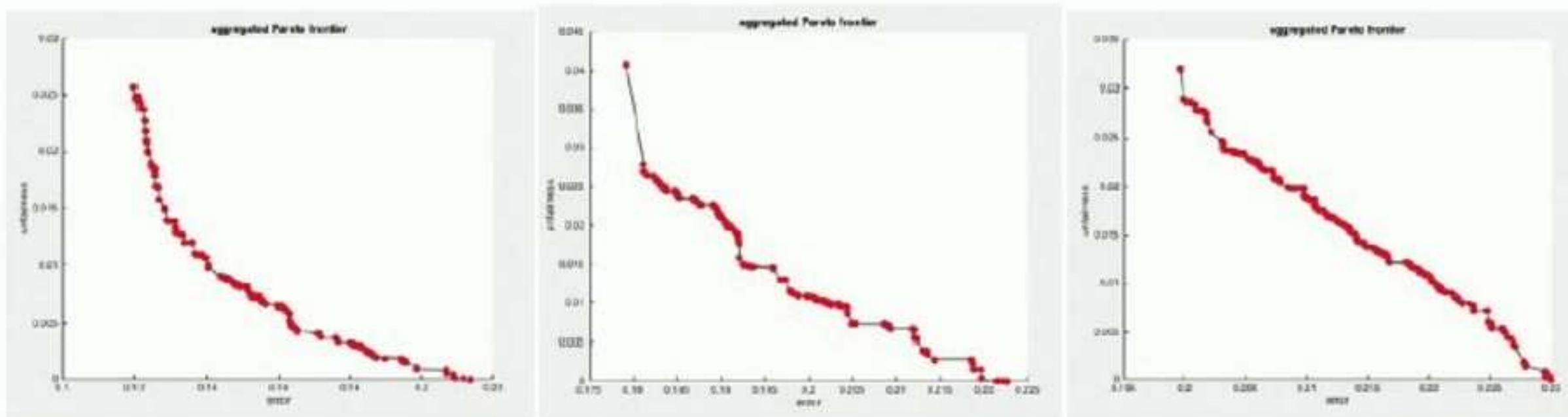
EFFICIENT FRONTIERS



Examples of Pareto frontiers of error (x axis) and an unfairness measure (y axis) for three different real data sets. The curves differ in their shapes and the actual numeric values on the error and fairness axes, thus presenting different trade-offs.



EFFICIENT FRONTIERS



Examples of Pareto frontiers of error (x axis) and an unfairness measure (y axis) for three different real data sets. The curves differ in their shapes and the actual numeric values on the error and fairness axes, thus presenting different trade-offs.

OTHER TOPICS

Games People Play (with Algorithms)

Games Scientists Play (with Data)

Interpretability, Accountability, Morality.... The Singularity

FRONTIERS OF FAIRNESS: FROM GROUPS TO INDIVIDUALS

Joint works with Chris Jung, Seth Neel, Aaron Roth, Saeed Sharifi, Logan Stapleton, Steven Wu

TYPES OF FAIRNESS DEFINITIONS

- Group Fairness
 - E.g. equality of error or false negative rates across gender, racial groups, etc.
 - Strong theory, practical implementations (e.g. fairness regularization)
 - But no guarantees to individuals
- Individual Fairness
 - E.g. metric fairness (“fairness through awareness”), meritocratic fairness
 - Binds at the individual level
 - But strong assumptions required (e.g. realizability) have prevented practical implementations

A FRAMEWORK FOR FAIR ML

- Begin by expressing training as a (linear or convex) constrained optimization problem
 - E.g. minimize error subject to various fairness constraints
 - In interesting cases, model space of learning algo and number of constraints may be exponential/infinite
 - Want to avoid explicit enumeration
- Use LP duality to pass to Lagrangian and recast as two-player, zero-sum game
 - Learner/Primal: wants to minimize error subject to constraints so far
 - Regulator/Dual: presents learner with violated constraints
 - Nash equilibrium is solution to constrained optimization problem
- If we can:
 - Formulate best responses as instances of *cost-sensitive classification*
 - Implement at least one player as a *no-regret* algorithm w.r.t. their strategy space... then algorithm provably converges in polynomial time given access to a standard learning heuristic
- Directly implement on top of your favorite “unfair” learning algorithm
- Applications:
 - Preventing “fairness gerrymandering”
 - Subjective individual fairness
 - Average individual fairness

AVERAGE INDIVIDUAL FAIRNESS

