

CS 260: Foundations of Data Science

Prof. Sara Mathieson

Spring 2025



Haverford
COLLEGE

Outline

- Midterm 2 Review
 - Logistic regression and cross entropy
 - Naive Bayes
 - Disparate impact
- Less focus
 - tSNE
 - t-tests

How can we tell if an algorithm is biased?

D: dataset with attributes X , Y

- * X is protected
- * Y is unprotected (other features)

How can we tell if an algorithm is biased?

D: dataset with attributes X , Y

- * X is protected
- * Y is unprotected (other features)

Goal: determine outcome C (hired, admitted, etc)

How can we tell if an algorithm is biased?

D: dataset with attributes X , Y

- * X is protected
- * Y is unprotected (other features)

Goal: determine outcome C (hired, admitted, etc)

Direct discrimination: $C = f(X)$

- * Female instrumentalist not hired for orchestra
- * Some ethnic groups not allowed to eat at a restaurant

How can we tell if an algorithm is biased?

D: dataset with attributes X , Y

- * X is protected
- * Y is unprotected (other features)

Goal: determine outcome C (hired, admitted, etc)

Indirect discrimination: $C = f(Y)$

- * but strong correlation between X and Y
- * Ex: housing loans
- * Ex: programming experience

Notation

features $\left\{ \begin{array}{l} X: \text{protected feature} \\ Y: \text{other features} \end{array} \right\}$

$X=0$ minority group
 $X=1$ majority group

9, 10, 6

label $\{ C: \text{binary outcome} \quad (\text{hired, admitted}) \}$
if not: $C=0$

$C=1$

Disparate Impact

$$P(C=1|X=0) \leq 0.8 P(C=1|X=1)$$

Idea: if we can predict X from Y
there could be disparate impact

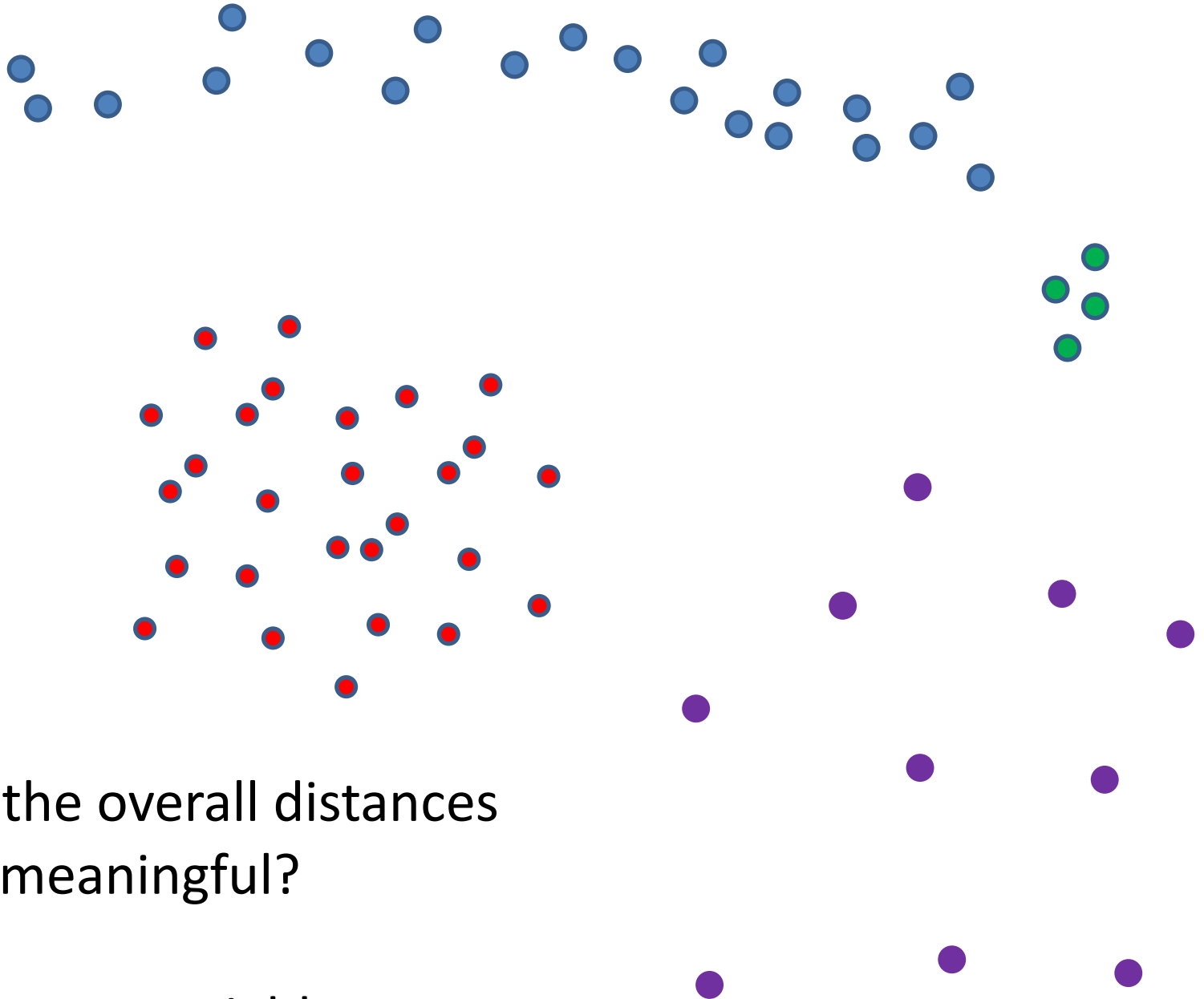
① train $f: Y \rightarrow X$ (f is a classifier)
input (features) output (label)

② use to predict $X \rightarrow$ get BER
balanced error rate

③ if BER is too low, could be disparate impact.

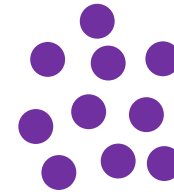
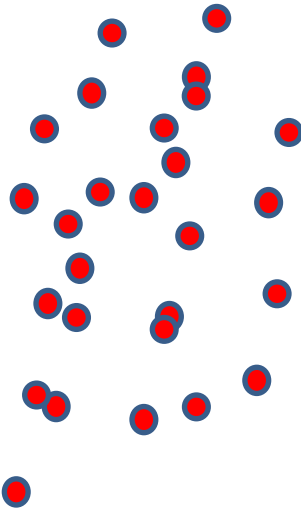
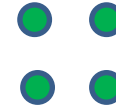
See video tutorial on Piazza!

Bootstrap demo



What if the overall distances
are not meaningful?

Focus on your neighbors

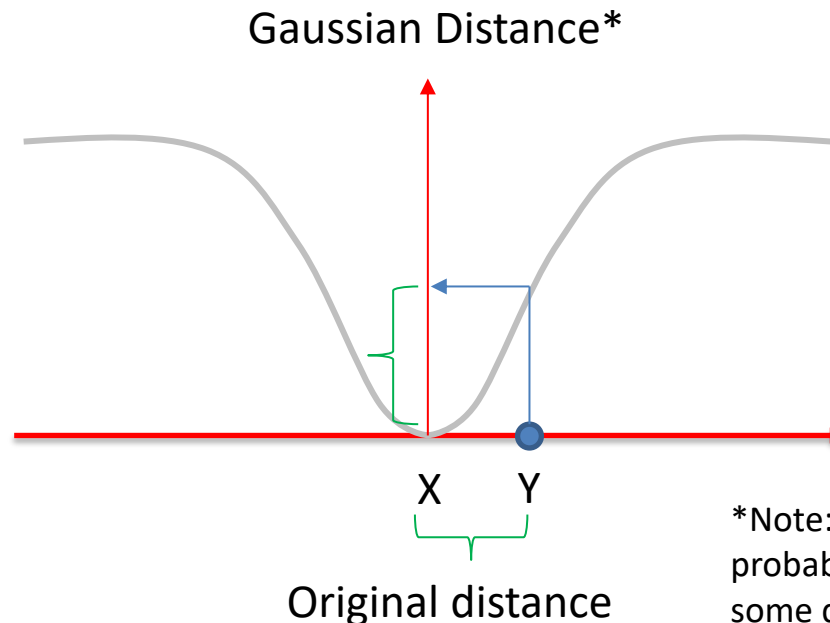


What if the overall distances
are not meaningful?

Focus on your neighbors

tSNE (t-distributed Stochastic Neighborhood Embedding)

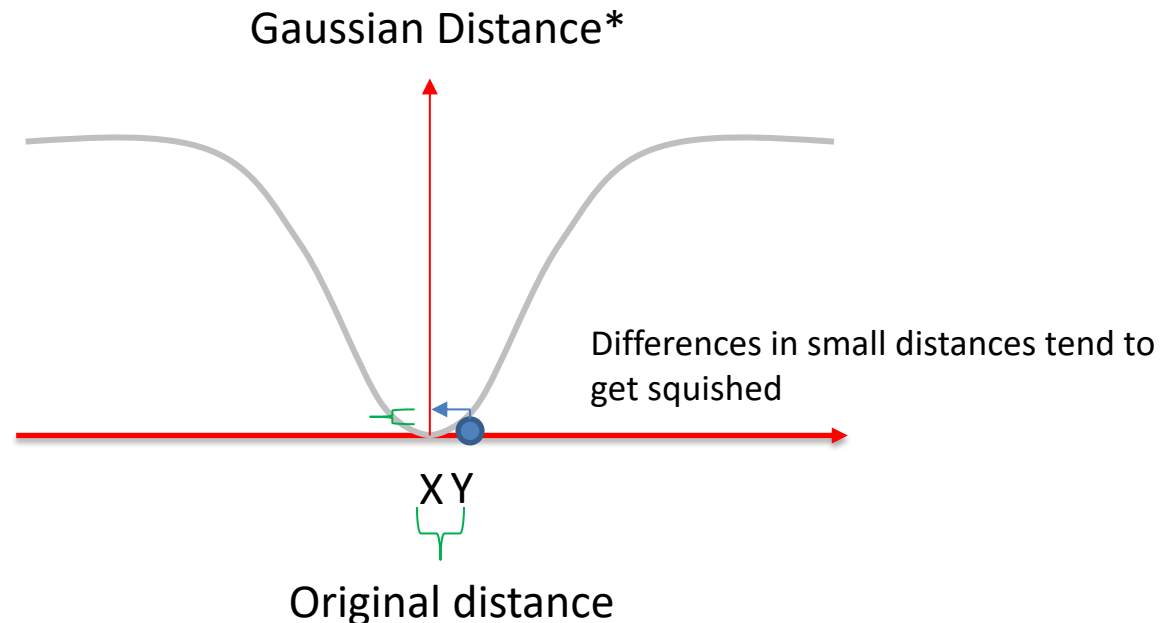
- Define distances between a point X to a point Y by a Gaussian function centered at X



*Note: the actual algorithm uses notions of probability (i.e., probability of finding Y at some distance from X). I use notion of distance as a proxy

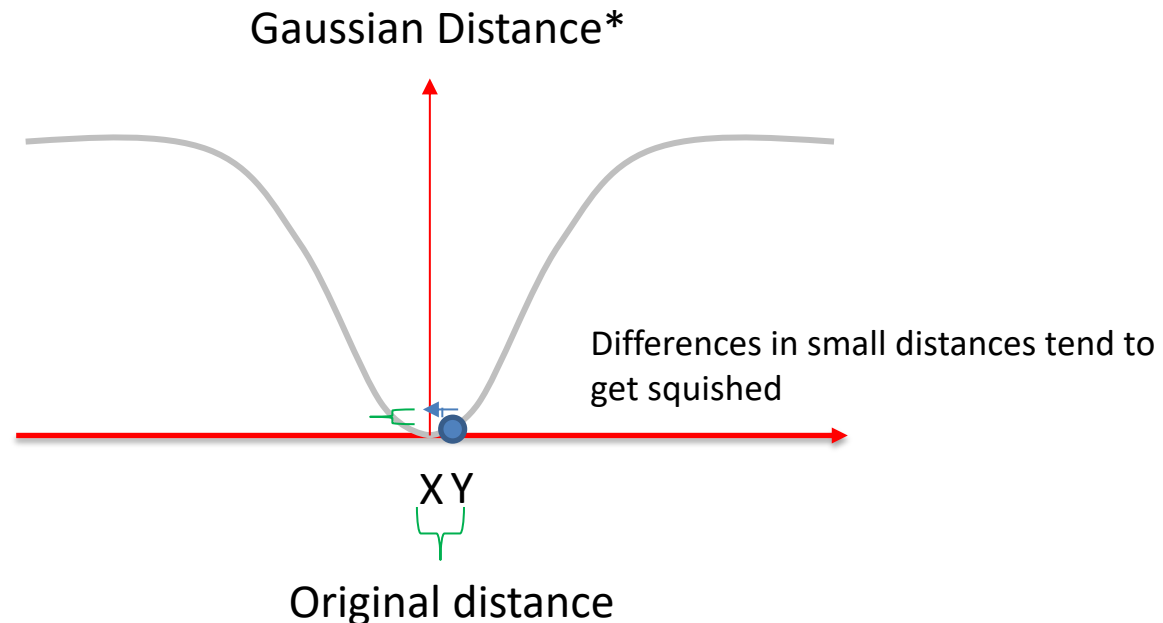
tSNE (t-distributed Stochastic Neighborhood Embedding)

- Define distances between a point X to a point Y by a Gaussian function centered at X



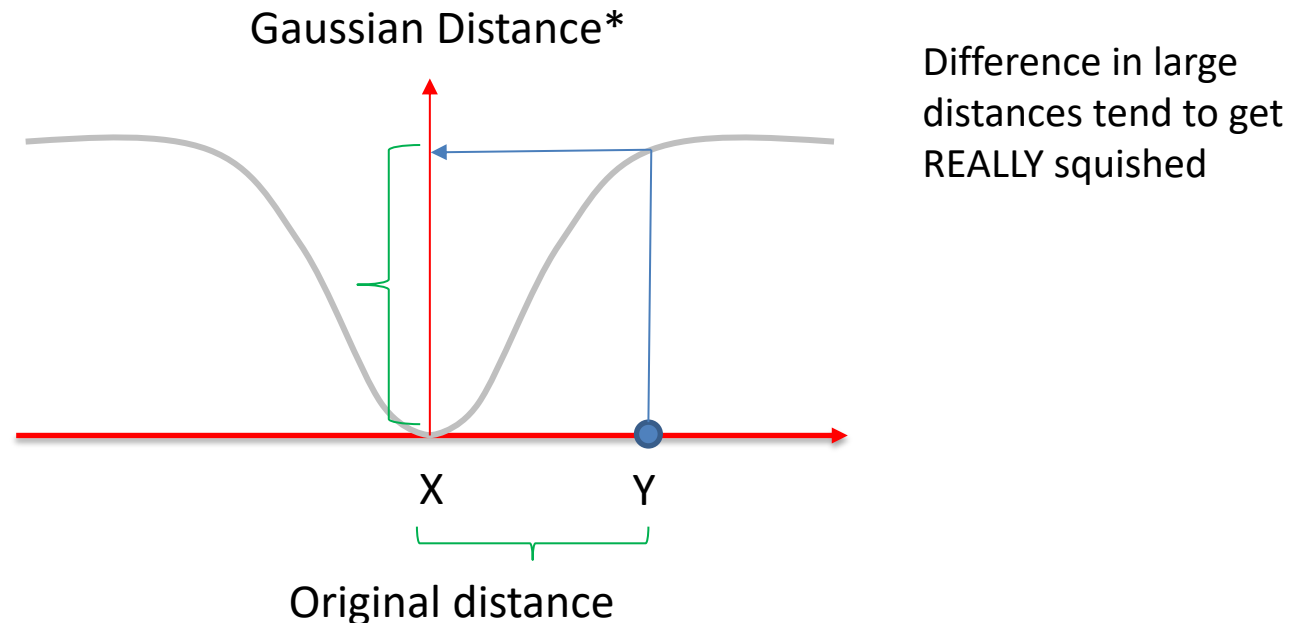
tSNE (t-distributed Stochastic Neighborhood Embedding)

- Define distances between a point X to a point Y by a Gaussian function centered at X



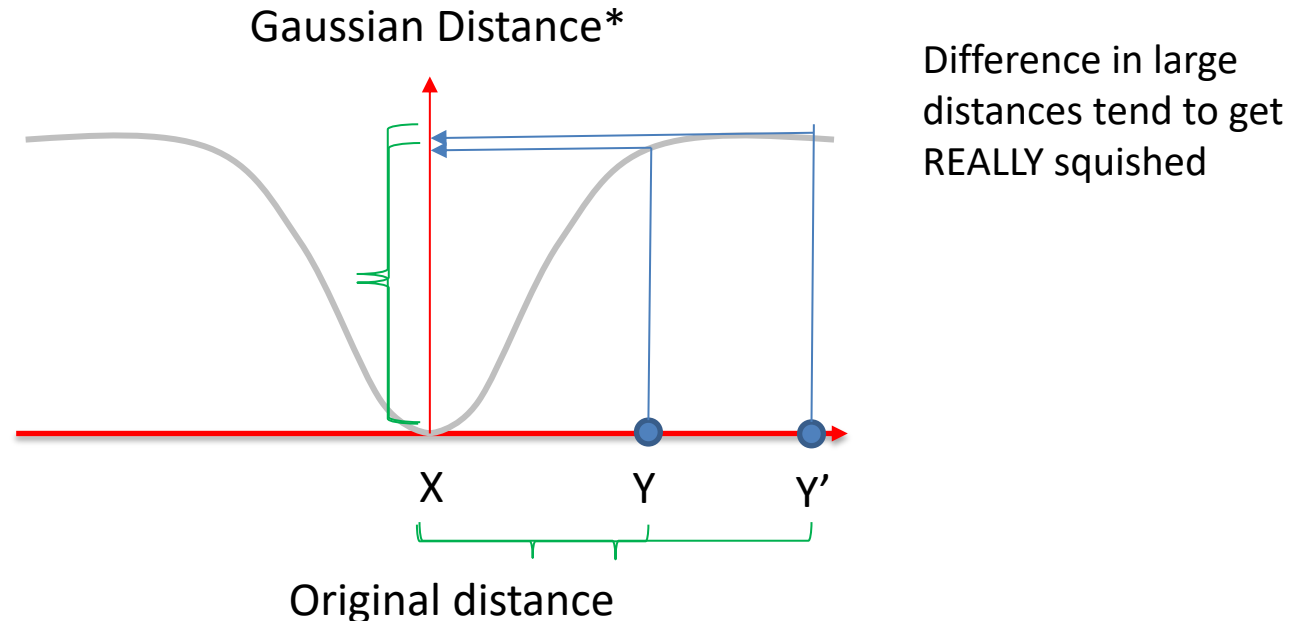
tSNE (t-distributed Stochastic Neighborhood Embedding)

- Define distances between a point X to a point Y by a Gaussian function centered at X



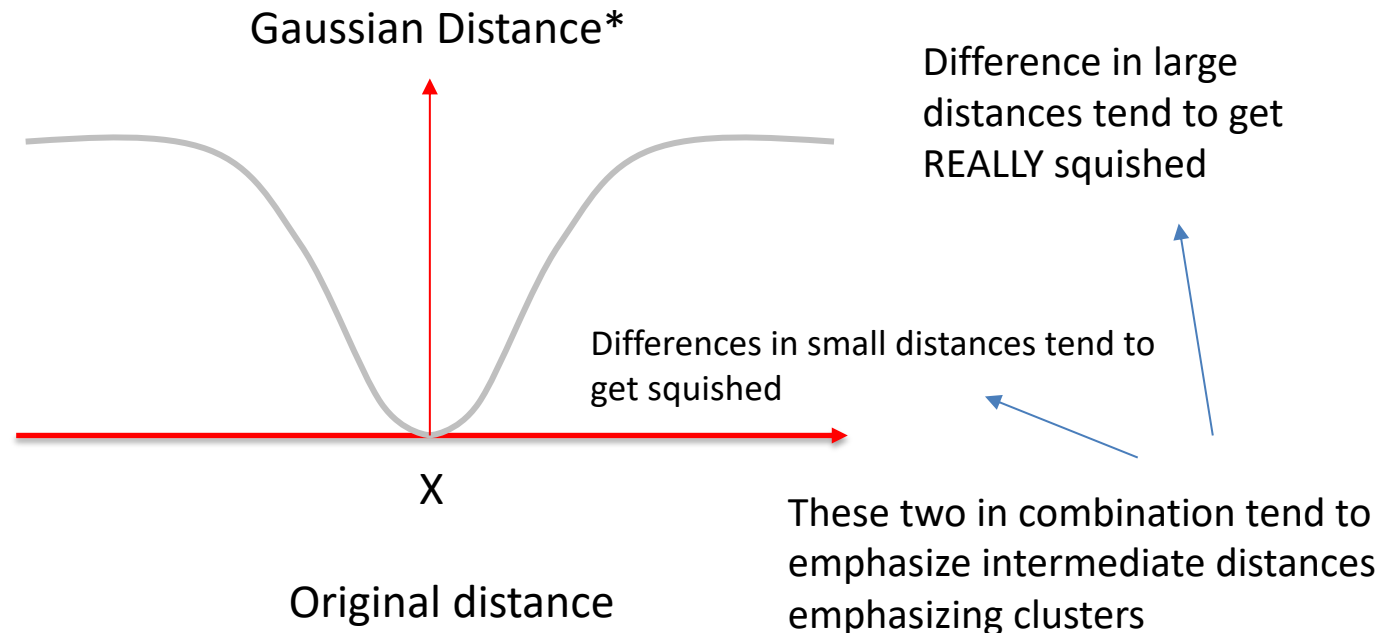
tSNE (t-distributed Stochastic Neighborhood Embedding)

- Define distances between a point X to a point Y by a Gaussian function centered at X

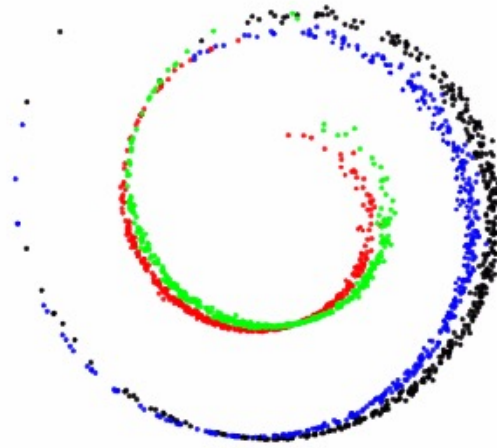


tSNE (t-distributed Stochastic Neighborhood Embedding)

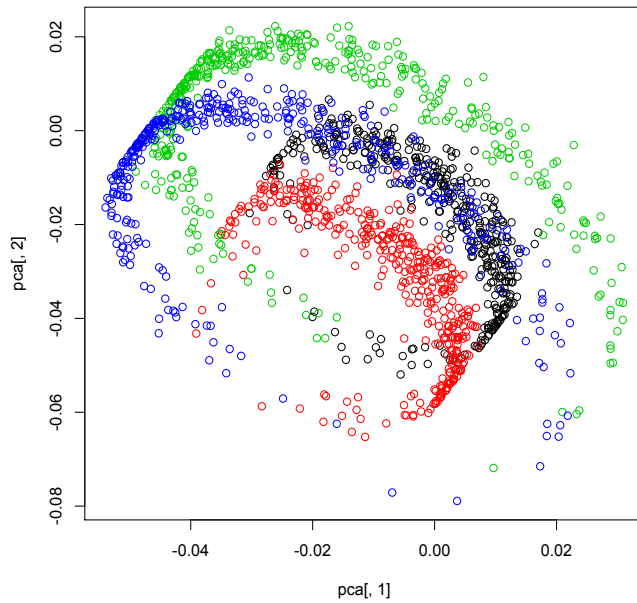
- Define distances between a point X to a point Y by a Gaussian function centered at X



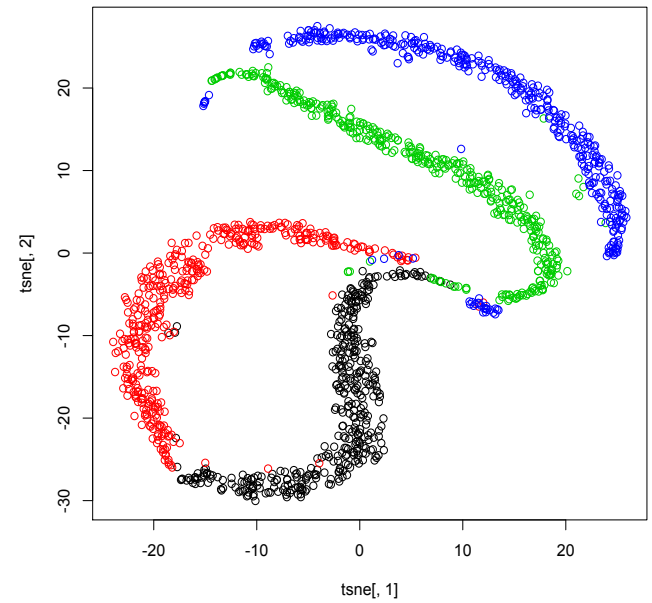
Original data



PCA

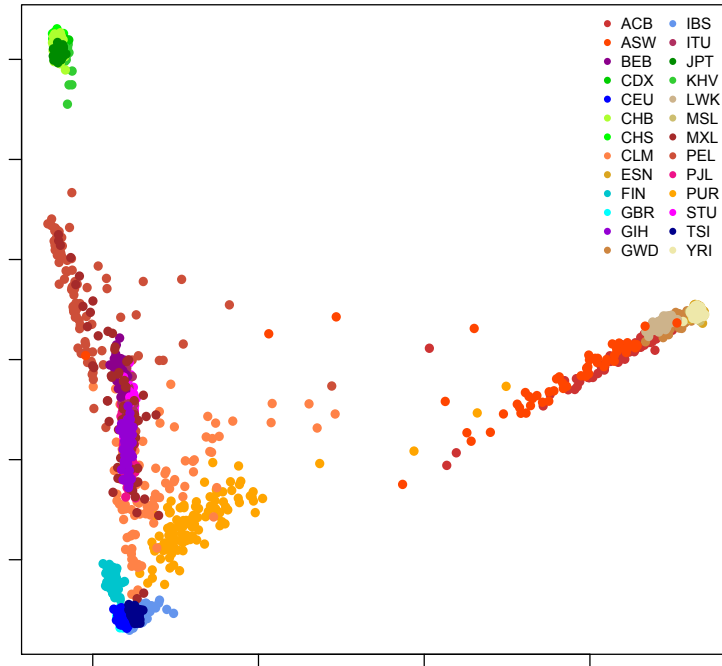


t-SNE

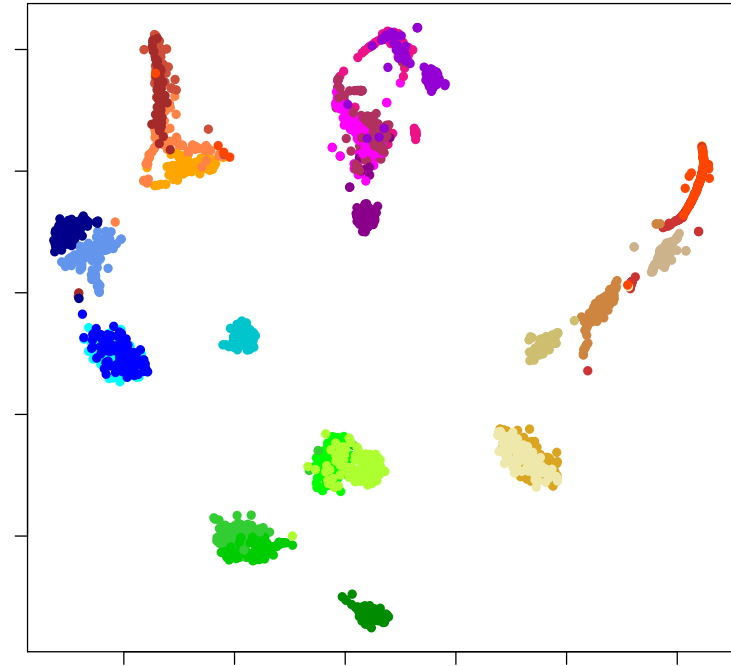


“swissroll data” **Dinoj Surendran**

PCA

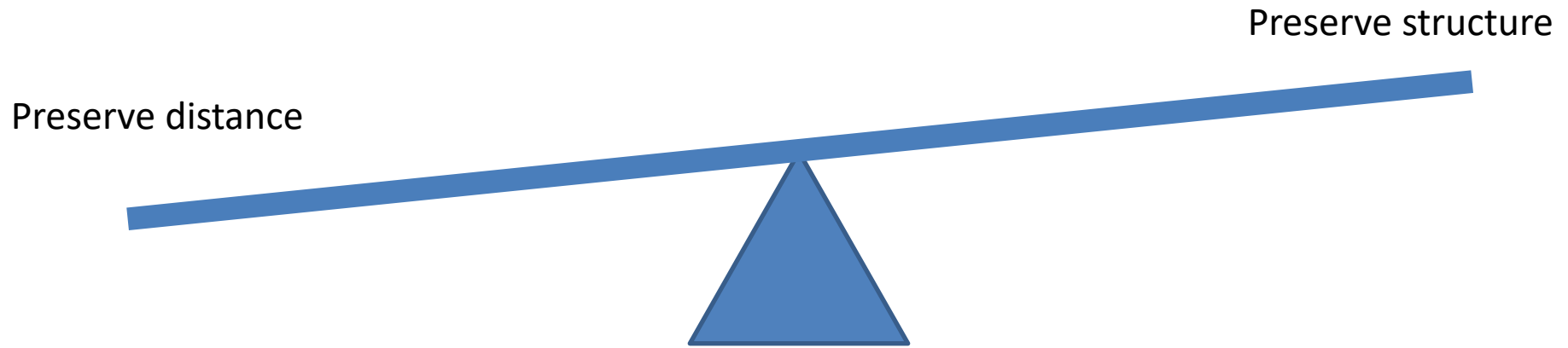


t-SNE



CHB	Han Chinese in Beijing, China
JPT	Japanese in Tokyo, Japan
CHS	Southern Han Chinese
CDX	Chinese Dai in Xishuangbanna, China
KHV	Kinh in Ho Chi Minh City, Vietnam
CEU	Utah Residents (CEPH) with Northern and Western European Ancestry
TSI	Toscani in Italia
FIN	Finnish in Finland
GBR	British in England and Scotland
IBS	Iberian Population in Spain
YRI	Yoruba in Ibadan, Nigeria
LWK	Luhya in Webuye, Kenya
GWD	Gambian in Western Divisions in the Gambia

MSL	Mende in Sierra Leone
ESN	Esan in Nigeria
ASW	Americans of African Ancestry in SW USA
ACB	African Caribbeans in Barbados
MXL	Mexican Ancestry from Los Angeles USA
PUR	Puerto Ricans from Puerto Rico
CLM	Colombians from Medellin, Colombia
PEL	Peruvians from Lima, Peru
GIH	Gujarati Indian from Houston, Texas
PJL	Punjabi from Lahore, Pakistan
BEB	Bengali from Bangladesh
STU	Sri Lankan Tamil from the UK
ITU	Indian Telugu from the UK



How to visualize data always depends on the data, and the question

There is rarely if ever a single correct approach

$$(a) \quad P(C=1 | X=0) \leq 0.8 P(C=1 | X=1)$$

↑
get
loan

↑
Black

↑
get
loan

↑
white

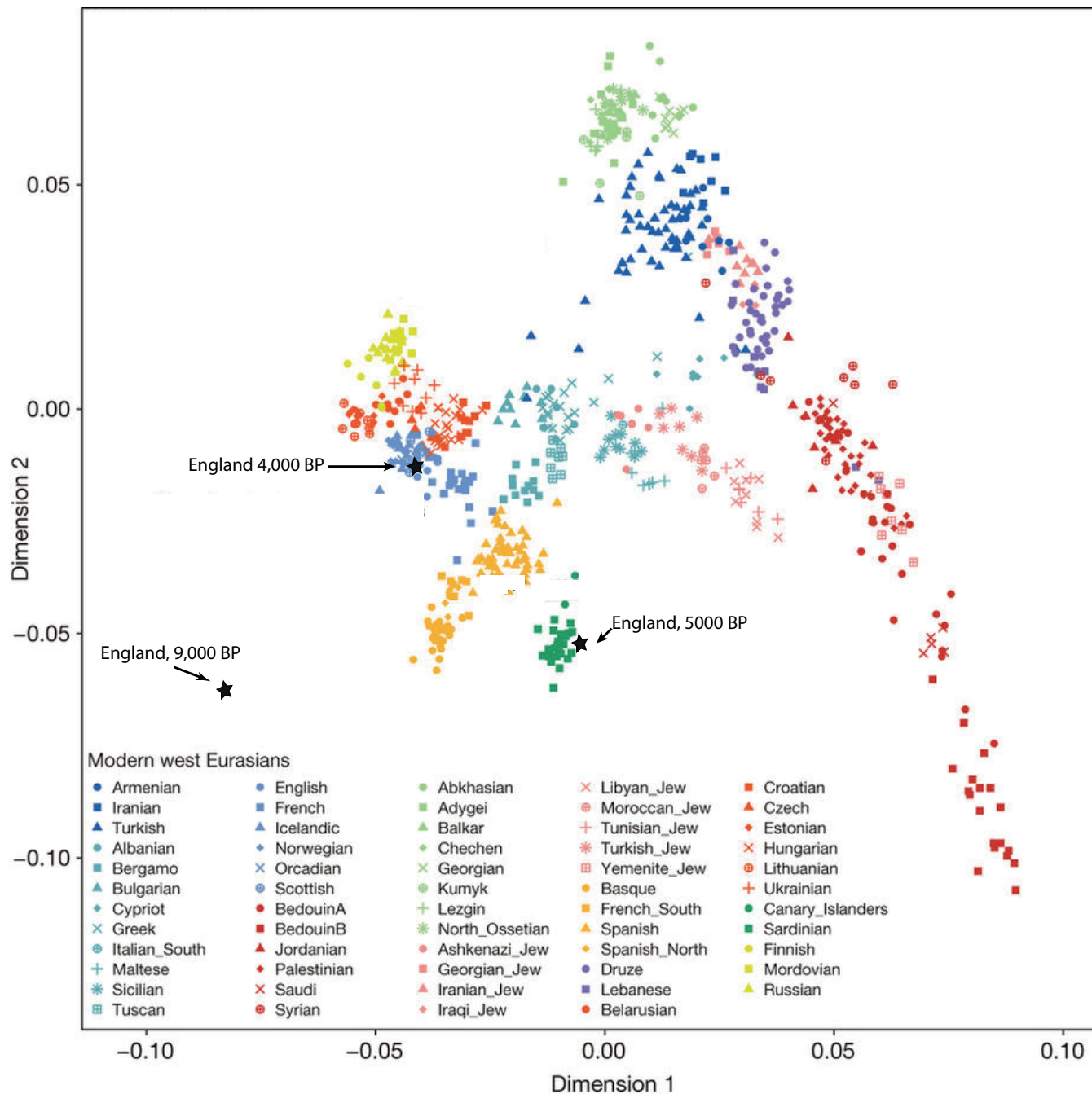
$$0.275 \leq 0.8(0.35)$$

$$0.275 \leq 0.28$$

✓

⇒ disparate impact

weighted avg
leaf entropies



P = pos
N = neg

H = healthy
D = disease

$$P(D|P) = 0.8$$

$$P(N|H) = P(P|D) = \frac{P(P)P(D|P)}{P(D)} = \frac{P(P) \cdot 0.8}{\frac{1}{500}} = x$$

$$x = (x + 499(1-x))0.8$$

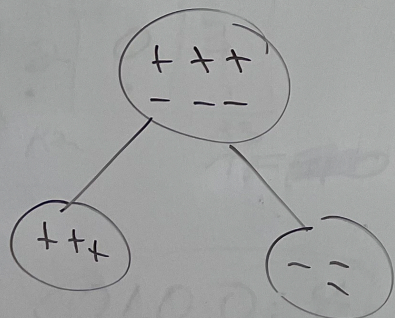
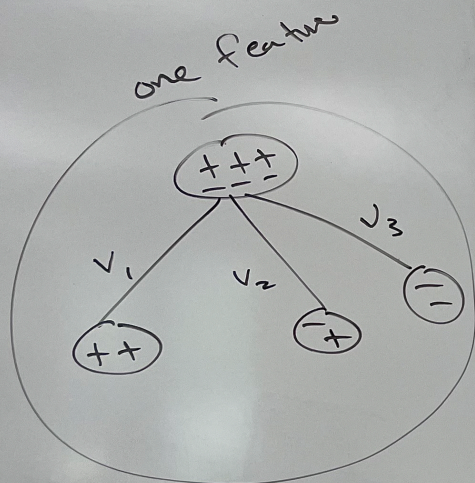
$$\frac{5}{4}x = x + 499 - 499x$$

$$\frac{5}{4}x = 499 - 498x$$

$$x \approx 0.9995$$

$$\begin{aligned} P(P) &= P(P, D) + P(P, H) \\ &= P(D)P(P|D) + P(H)P(P|H) \\ &= \frac{1}{500}x + \frac{499}{500} \left(1 - \underbrace{P(N|H)}_x\right) \end{aligned}$$

$$x = \left[\frac{\cancel{1}}{\cancel{500}}x + \frac{499}{\cancel{500}}(1-x) \right] 0.8$$



$$H(Y) = - \sum_{c \in \text{vals}(Y)} p(Y=c) \log_2 p(Y=c)$$

$$= - p(Y=+) \log_2 p(Y=+) - p(Y=-) \log_2 p(Y=-)$$

$$= - \frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2}$$

$$= + \frac{1}{2} + \frac{1}{2} = 1$$

$$H(Y|X) = \sum_v p(X=v) H(Y|X=v) \quad \left. \vphantom{\sum_v} \right\} \begin{array}{l} \text{weighted avg} \\ \text{of leaf entropies} \end{array}$$

one feature

$$H(Y|\underline{X=v}) = - \sum_c p(Y=c|X=v) \log_2 p(Y=c|X=v)$$

(9)