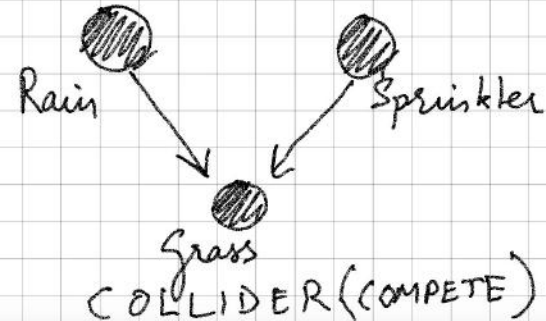
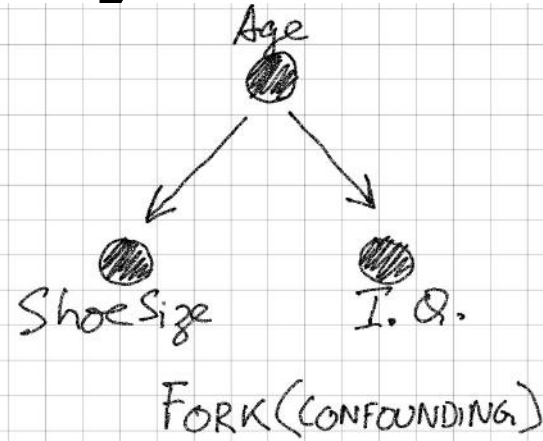
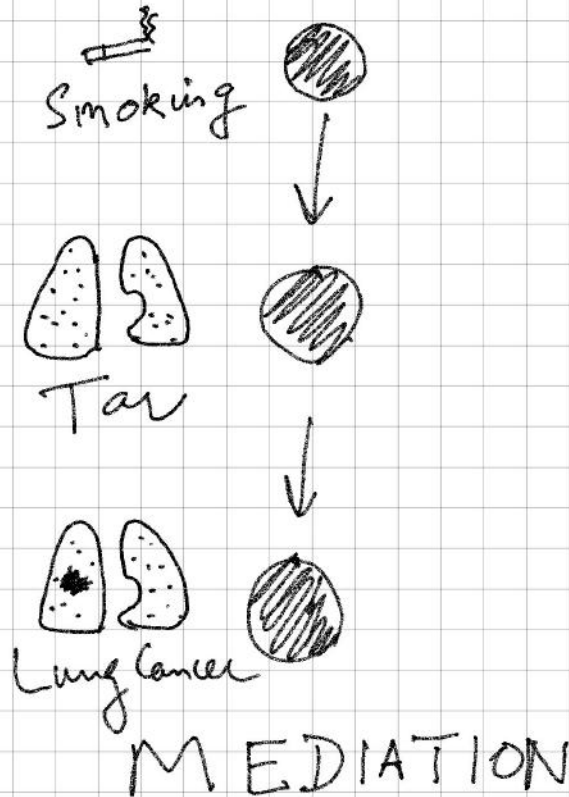


Computing for Medicine

Data Science: Demo Bayesian Networks and PCA
Tavprites Sethi

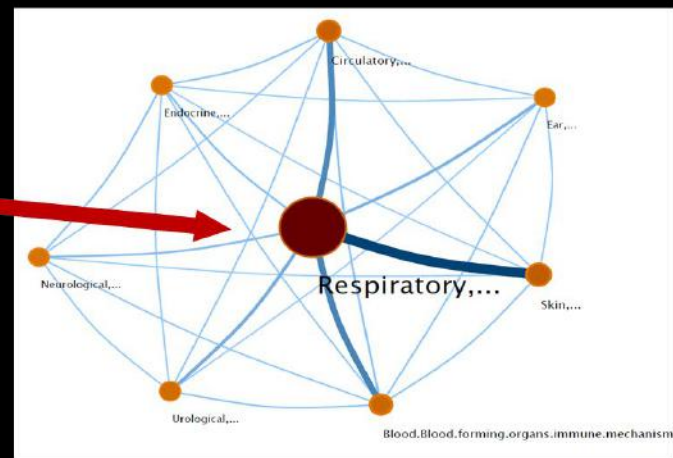
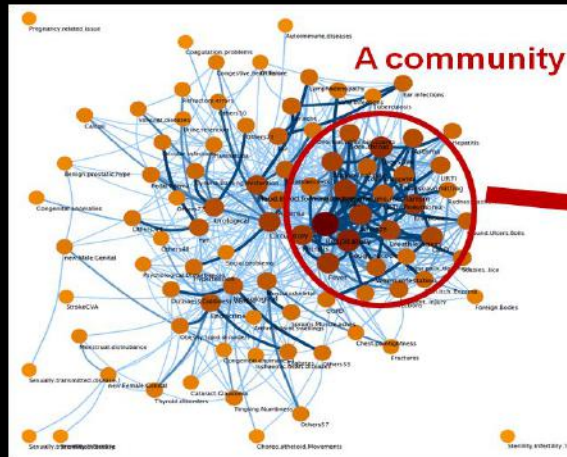
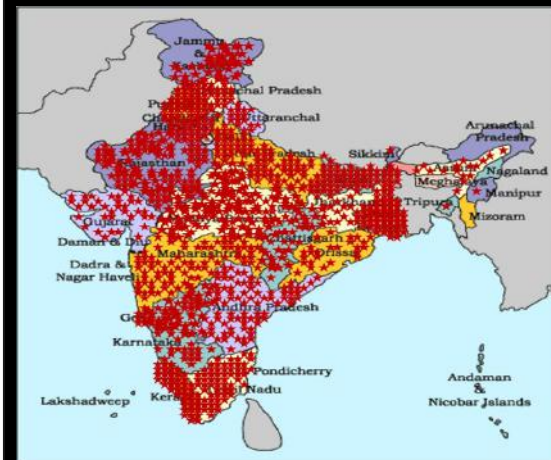
Building Blocks of Bayesian Nets



Prevalence of Symptoms in a Single Indian Healthcare Day on a Nationwide Scale

One day point prevalence study of symptoms in **2,04,912** patients across India

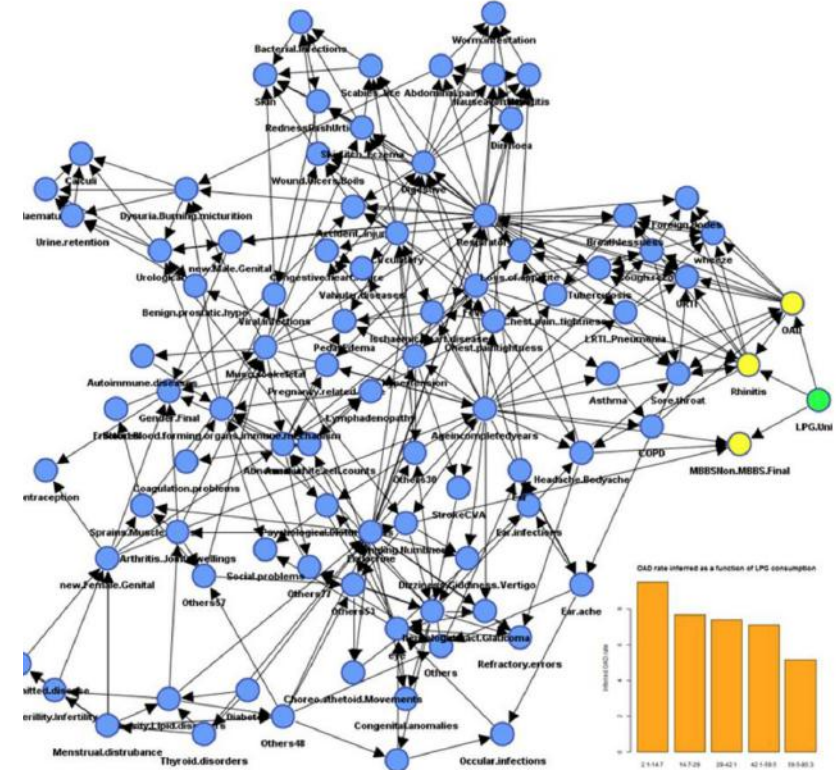
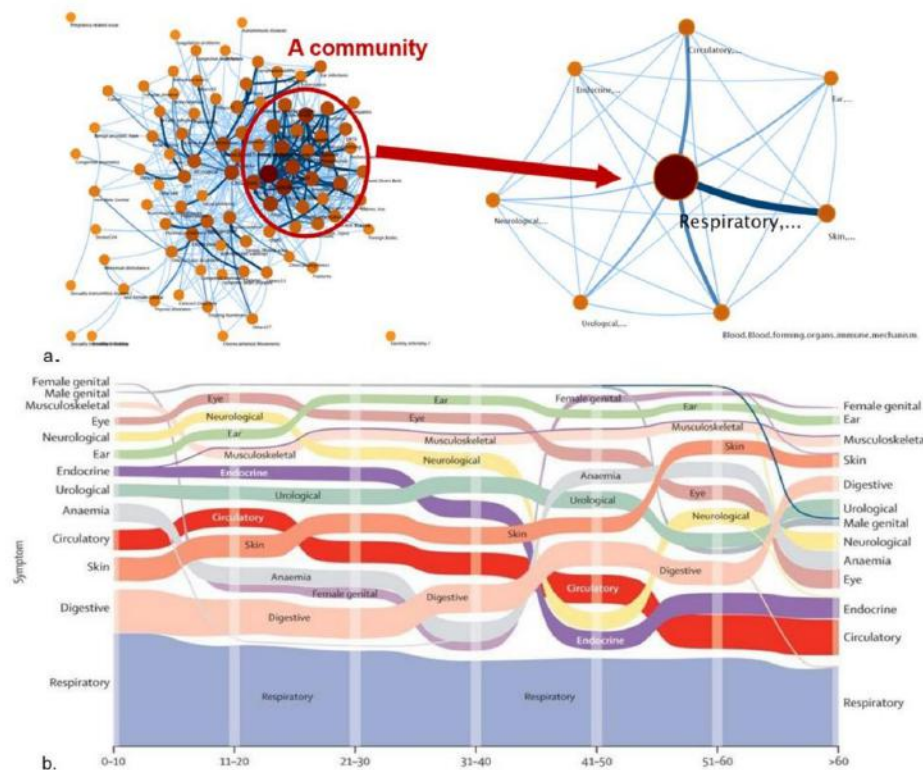
Chest Research Foundation, Pune, India



Networks approach: many diseases happen together, connectivity differs across age

Lancet Global Health, Dec 2015

Associations to Decisions



Sundee Salvi et al. Symptoms and medical conditions in 204 912 patients visiting primary health-care practitioners in India: a 1-day point prevalence study (the POSEIDON study), *The Lancet Global Health*, Volume 3, Issue 12, 2015, Pages e776-e784, [https://doi.org/10.1016/S2214-109X\(15\)00152-7](https://doi.org/10.1016/S2214-109X(15)00152-7).

Case Study II: AI for Reducing Health Inequities

- The richest American men live **15 years** longer than the poorest American men¹
- The richest American women live **10 years** longer than the poorest American women¹
- **Healthcare inequities** impose an estimated burden of **\$300 billion per year** in the United States
- Longevity is the **sum-total of influences** on the healthcare
- Hence **longevity-gap** is a complex **socio-demographic** challenge
- Key motivation: **learn policy** for **mitigating** the longevity-gap using explainable AI and release it to public, policymakers

Sethi T, S. Maheshwari, A. Mittal, S. Chugh. Learning to Address Health Inequality in the United States with a Bayesian Decision Network. Proceedings of the AAAI Conference on Artificial Intelligence 33, 710-717. DOI: <https://doi.org/10.1609/aaai.v33i01.3301710> c

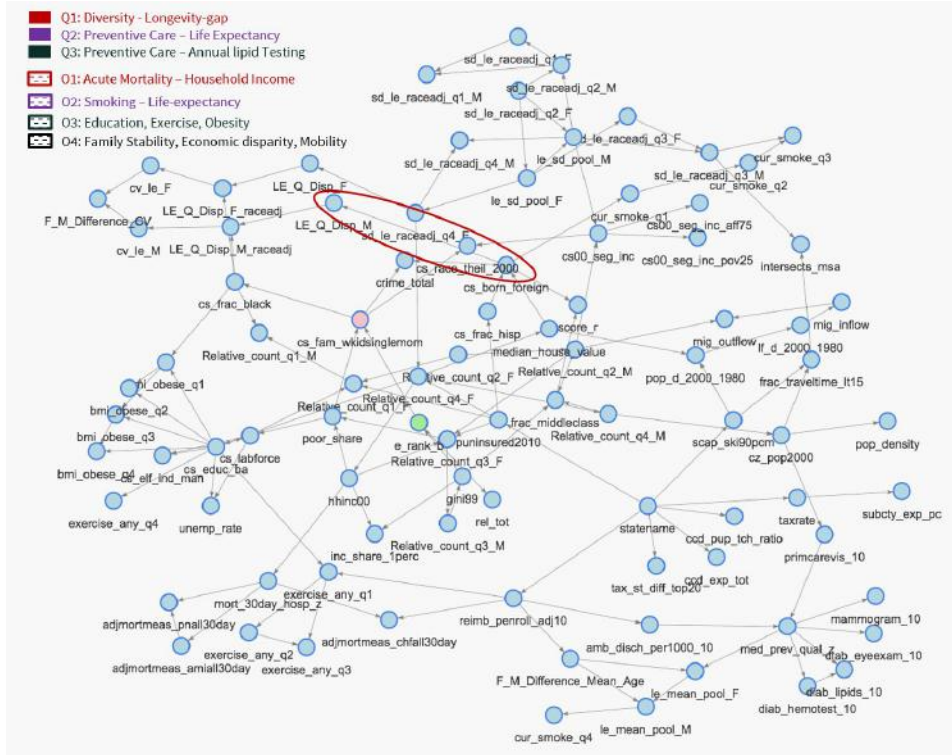
Key Messages

We used data from: **Mortality** (Census), **Healthcare** Indices (Dartmouth Atlas), **Health-behaviors** (CDC, BRFSS), **Education** (K-12 and Post Secondary), **Demographics** (e.g. race, ethnicity, diversity, gender ratio etc), **Socioeconomics** (e.g. Gini index, Poverty rate, Income segregation, Social Mobility), **Social Cohesion**. (e.g. Social Capital Index, Religious adherents), **Labor market conditions and Taxation**. (e.g. unemployment, manufacturing sector).

Key Messages

1. **Social.** Diversity mitigates health inequality in the US. Counties with higher diversity are 38% less likely to have a longevity gap between the rich and the poor.
2. **Preventive Care.** Counties with high quality primary care services (not just visits but investigations) have a 43% increase in the probability of living beyond 85 years in females (corresponding 30% increase for males.)
3. **Clinical.** Acute mortality (30-day Hospital Mortality Index) is 30% less in Counties with household income in the highest segment.
4. **Usual suspects.** Smoking, Education, Exercise as expected to be key influencers of longevity.
5. **Socio-demographic.** Family stability decreases crime rate, increases upward social mobility across economic tiers and indirectly decreases Gini disparity in counties.

Diversity Mitigates Longevity Gap



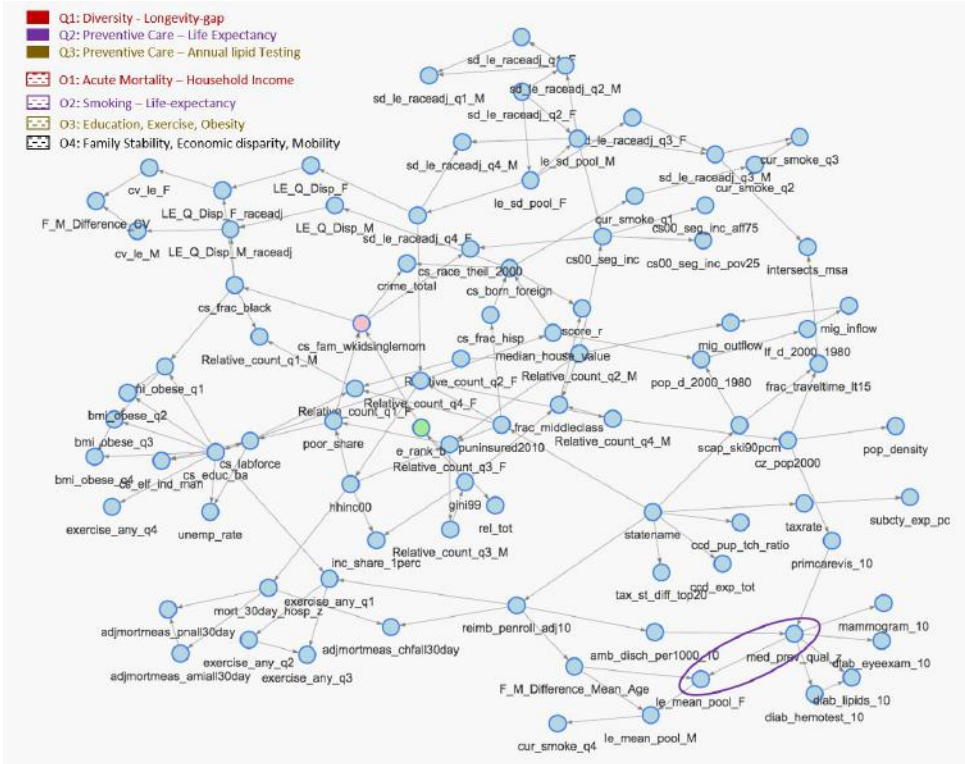
A 1. Diversity

Counties with higher diversity are 38% less likely to have a longevity gap between the rich and the poor.

Likely explanation: Structure indicates that Higher Diversity is Associated with Higher Income, especially in Females, thus improving healthcare services.

Sethi T., S. Maheshwari, A. Mittal, S. Chugh. Learning to Address Health Inequality in the United States with a Bayesian Decision Network. Proceedings of the AAAI Conference on Artificial Intelligence 33, 710-717. DOI: <https://doi.org/10.1609/aaai.v33i01.3301710> c

The Impact of Primary Care



A 2. Quality of Preventive Care

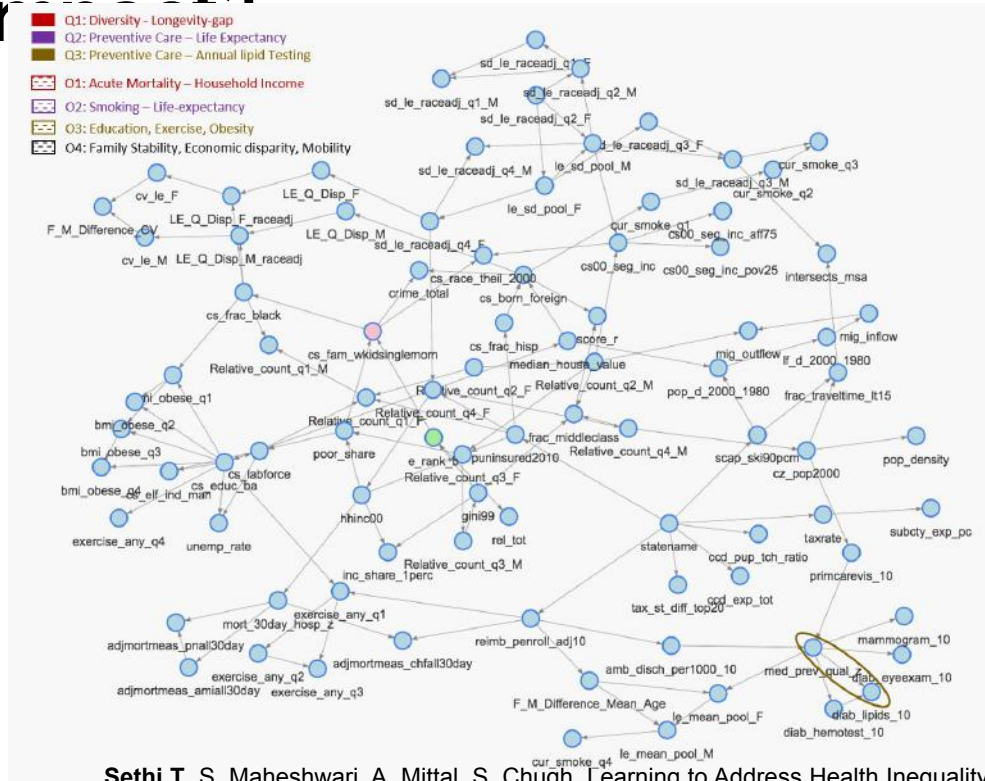
Counties with high quality primary care services have a 43% increase in the probability of living beyond 85 years in females (corresponding 30% increase for males.)

Likely explanation: Self Evident, but previously unquantified in a Joint Model

Sethi T., S. Maheshwari, A. Mittal, S. Chugh. Learning to Address Health Inequality in the United States with a Bayesian Decision Network. Proceedings of the AAAI Conference on Artificial Intelligence 33, 710-717. DOI: <https://doi.org/10.1609/aaai.v33i01.3301710> c

Which primary care investigation has most

impact?



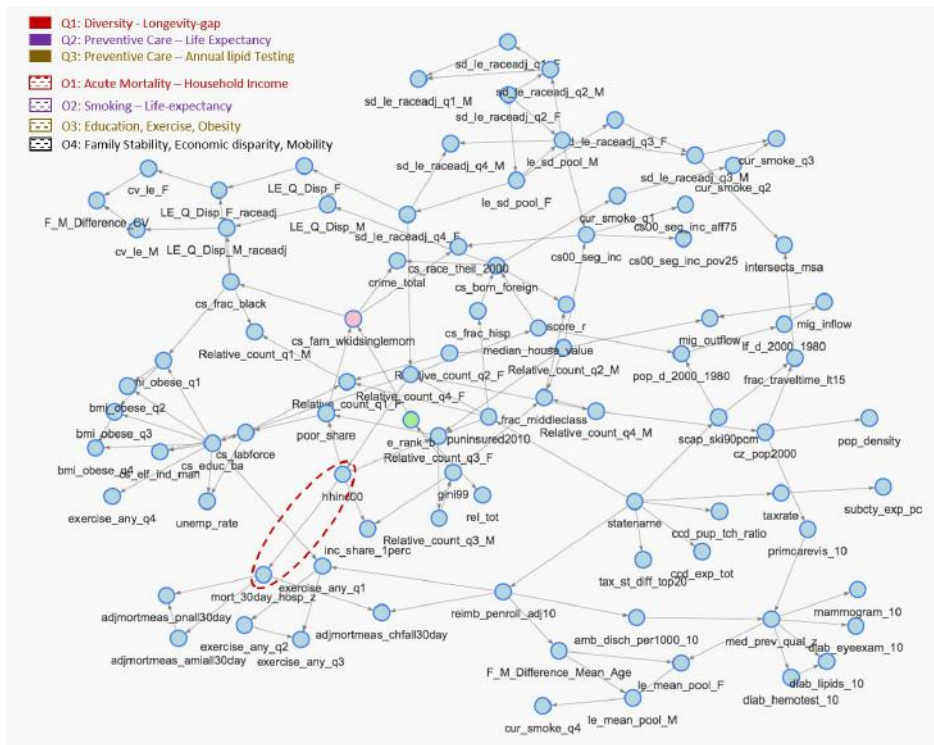
A 3. Annual Lipid Testing esp. in the diabetic population

diab_eyeexam_10	mammogram_10	diab_lipids_10	payoff
[70.2,85.6]	[68.2,86.1]	[79.3,92.9]	0.52
[70.2,85.6]	[68.2,86.1]	[65.6,79.3]	0.48
[62.2,70.2]	[68.2,86.1]	[79.3,92.9]	0.44
[70.2,85.6]	[59.5,68.2]	[79.3,92.9]	0.40
[62.2,70.2]	[68.2,86.1]	[65.6,79.3]	0.26
[70.2,85.6]	[59.5,68.2]	[65.6,79.3]	0.22
[42.4,62.2]	[68.2,86.1]	[79.3,92.9]	0.20
[62.2,70.2]	[59.5,68.2]	[79.3,92.9]	0.12
[70.2,85.6]	[31.1,59.5]	[79.3,92.9]	0.11
[42.4,62.2]	[68.2,86.1]	[65.6,79.3]	0.08

Likely explanation:
Diabetics are at highest risk of cardiovascular mortality

Sethi T, S. Maheshwari, A. Mittal, S. Chugh. Learning to Address Health Inequality in the United States with a Bayesian Decision Network. Proceedings of the AAAI Conference on Artificial Intelligence 33, 710-717. DOI: <https://doi.org/10.1609/aaai.v33i01.3301710> c

Income and Acute Mortality



O 1.

ACUTE MORTALITY (30-DAY HOSPITAL MORTALITY INDEX > 0.92) IS 30% LESS IN COUNTIES WITH HOUSEHOLD INCOME IN THE HIGHEST SEGMENT.

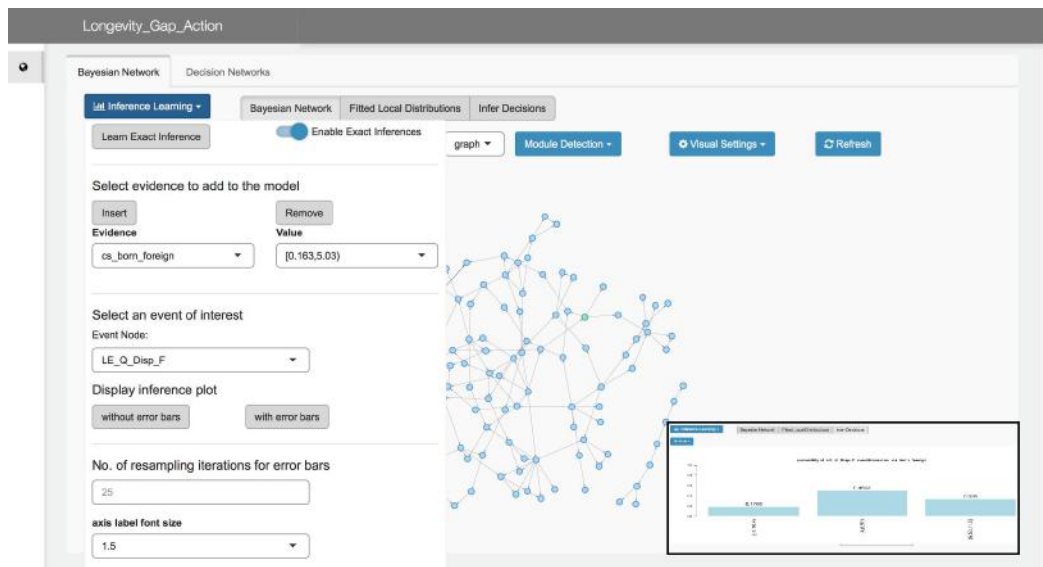
HIGHEST CONTRIBUTOR TO ACUTE MORTALITY IN LOWER INCOME HOUSEHOLDS IS PNEUMONIA

[illegible]

- HIGH LIFE-EXPECTANCY IN MALES (81.9 - 85 YEARS) MAKES IT 33% MORE LIKELY FOR **SMOKING** TO BE IN THE LOWEST STRATUM IN THE COUNTY.
- COUNTIES WITH HIGH PROPORTION OF EXERCISE HAVE 19% LESS HOSPITALIZATION RATES
- COUNTIES WITH LOWER FAMILY STABILITY ARE 40% MORE LIKELY TO HAVE LOWER SOCIAL MOBILITY

Tavpritesh Sethi, Anant Mittal, Shubham Maheshwari, Samarth Chugh. *Learning to Address Health Inequality in the United States with a Bayesian Decision Network*. <https://arxiv.org/abs/1809.09215> Accepted for publication in the Thirty-third AAAI conference in Artificial Intelligence, AAAI-2019

Deploy your XAI models as Web applications



https://github.com/SAFE-ICU/Longevity_Gap_Action

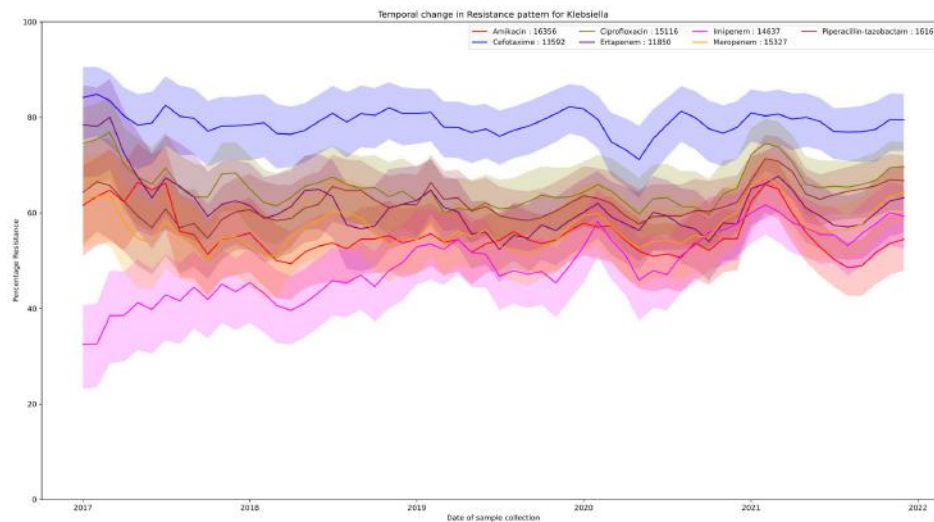
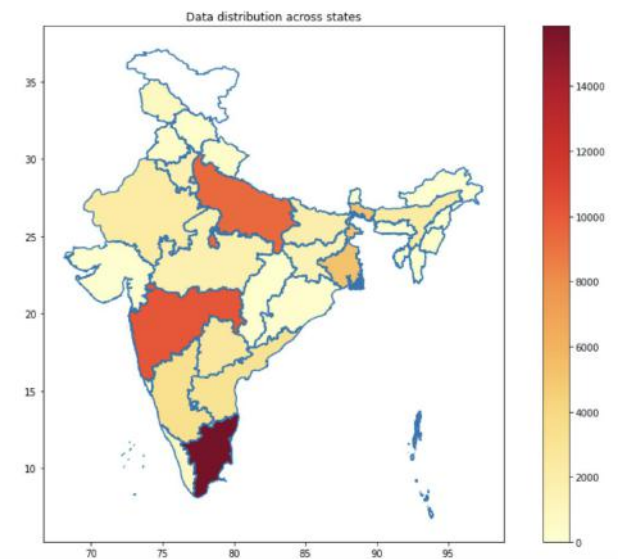
Key Messages

1. **Social.** Diversity mitigates health inequality in the US. Counties with higher diversity are 38% less likely to have a longevity gap between the rich and the poor.
2. **Preventive Care.** Counties with high quality primary care services (not just visits but investigations) have a 43% increase in the probability of living beyond 85 years in females (corresponding 30% increase for males.)
3. **Clinical.** Acute mortality (30-day Hospital Mortality Index) is 30% less in Counties with household income in the highest segment.
4. **Usual suspects.** Smoking, Education, Exercise as expected to be key influencers of longevity.
5. **Socio-demographic.** Family stability decreases crime rate, increases upward social mobility across economic tiers and indirectly decreases Gini disparity in counties.

Tavpritesh Sethi, Anant Mittal, Shubham Maheshwari, Samarth Chugh. *Learning to Address Health Inequality in the United States with a Bayesian Decision Network.*
<https://arxiv.org/abs/1809.09215> Accepted for publication in the Thirty-third AAAI conference in Artificial Intelligence, AAAI-2019

Emerging trends in antimicrobial resistance in bloodstream infections: multicentric longitudinal study in India (2017–2022)

Jasmine Kaur ^{a,b,c} · Harpreet Singh ^c · Tavpritesh Sethi ^{a,b}  

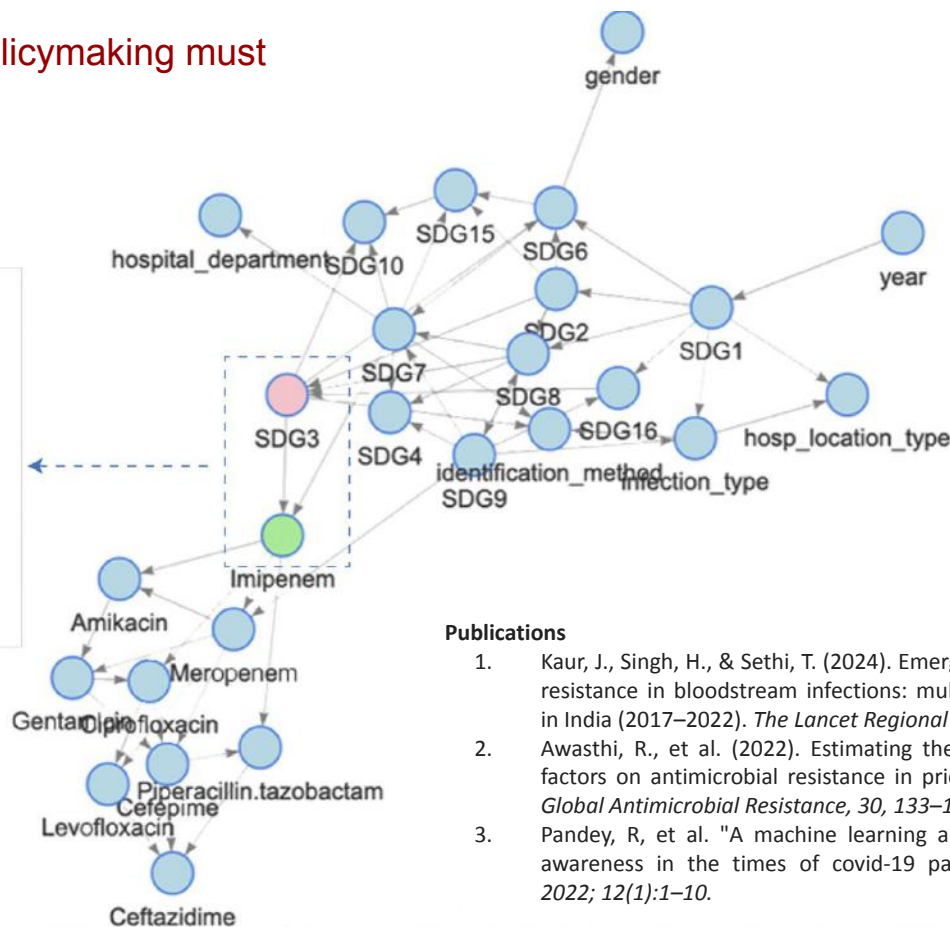
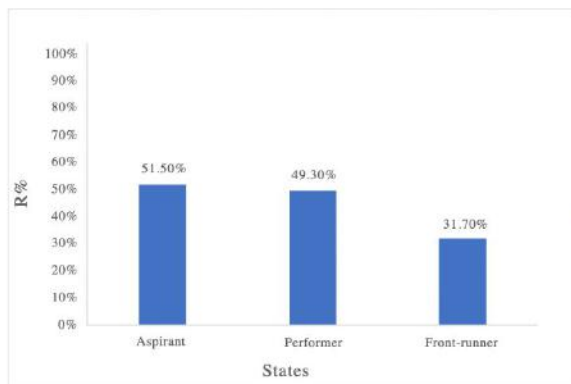


Klebsiella Sepsis

0.25% monthly average for increase in Imipenem resistance

AI for Understanding Impact on SDG Achievement

Key Takeaway: Evidence-based policymaking must leverage AI for tracking progress.



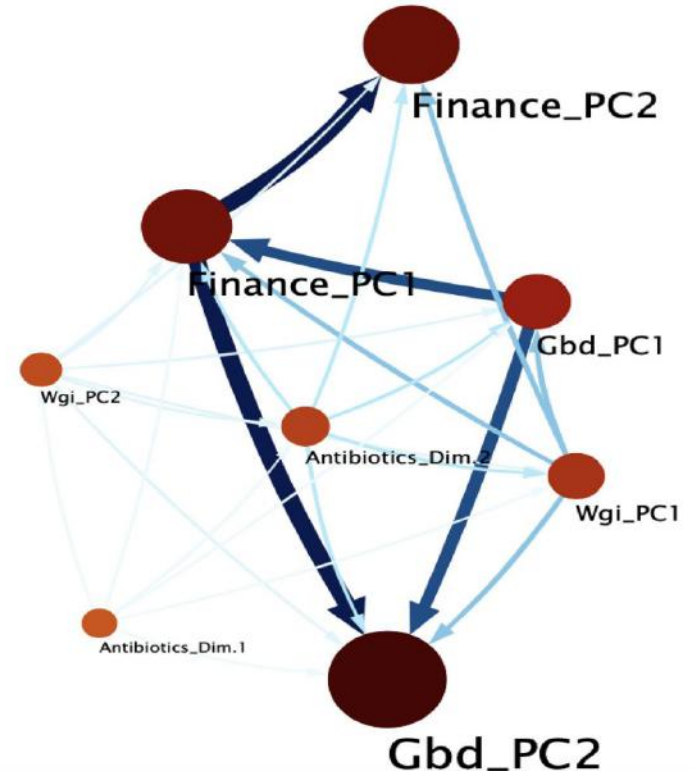
Publications

1. Kaur, J., Singh, H., & Sethi, T. (2024). Emerging trends in antimicrobial resistance in bloodstream infections: multicentric longitudinal study in India (2017–2022). *The Lancet Regional Health-Southeast Asia*, 26.
2. Awasthi, R., et al. (2022). Estimating the impact of health systems factors on antimicrobial resistance in priority pathogens. *Journal of Global Antimicrobial Resistance*, 30, 133–142.
3. Pandey, R, et al. "A machine learning application for raising wash awareness in the times of covid-19 pandemic". *Scientific reports* 2022; 12(1):1–10.

Systems Indicators and AMR

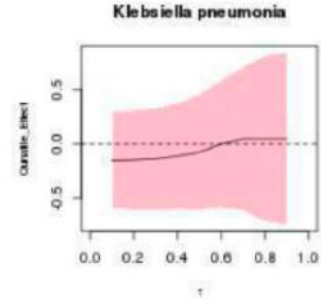
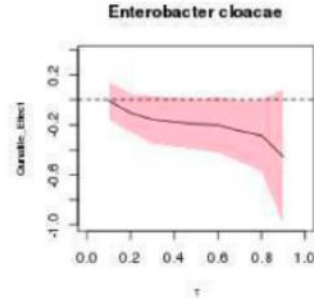
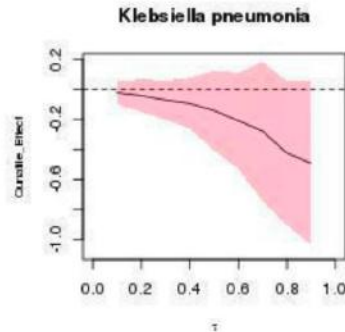
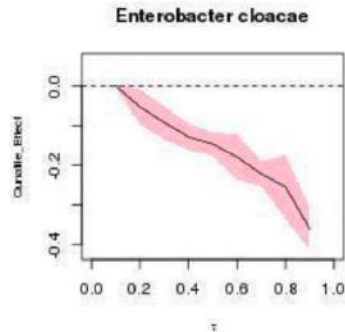
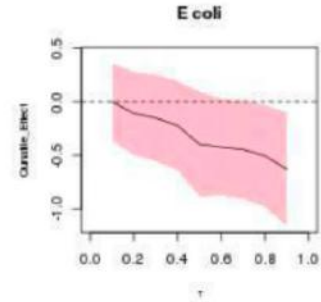
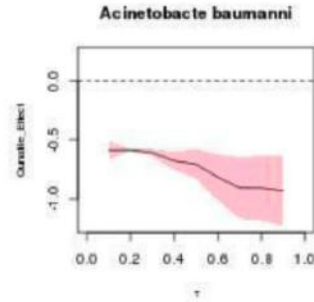
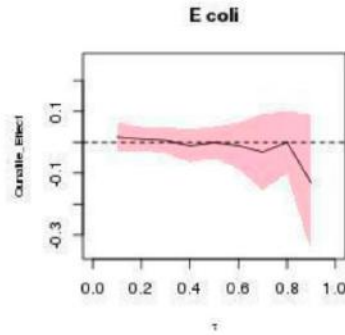
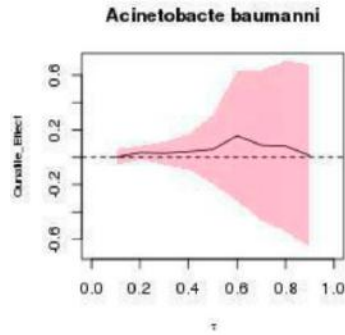
Age and temporal distribution of data

	Women	Men	Overall
Sex	No. (%)	No. (%)	No. (%)
	278 128 (43.88%)	348 136 (54.93%)	633 820
Age group			
0 to 2 years	15 329 (42.83%)	20 012 (55.91%)	35 793
13 to 18 years	6348 (47.06%)	7043 (52.21%)	13 490
19 to 64 years	130 162 (44.04%)	163 201 (55.22%)	295 537
3 to 12 years	12 613 (46.88%)	14 052 (52.22%)	26 907
65 to 84 years	86 561 (41.63%)	119 801 (57.62%)	207 922
85 and over	23 406 (54.29%)	19 364 (44.91%)	43 114
Year			
2004	9433 (46.93%)	10 655 (53.01%)	20 101
2005	10 473 (48.04%)	11 274 (51.71%)	21 801
2006	13 967 (46.65%)	15 876 (53.03%)	29 940
2007	18 264 (45.70%)	21 409 (53.57%)	39 964
2008	16 303 (44.33%)	19 925 (54.18%)	36 773
2009	18 575 (44.55%)	22 514 (54.00%)	41 692
2010	14 476 (44.59%)	17 341 (53.42%)	32 462
2011	11 622 (44.66%)	13 931 (53.53%)	26 023
2012	22 432 (42.16%)	29 477 (55.40%)	53 206
2013	29 962 (42.73%)	38 850 (55.40%)	70 125
2014	30 492 (43.23%)	39 482 (55.98%)	70 529
2015	27 992 (43.21%)	36 104 (55.73%)	64 785
2016	29 744 (42.74%)	39 290 (56.45%)	69 598
2017	24 393 (42.93%)	32 008 (56.33%)	56 821



Awasthi R, Rakholia V, Agrawal S, Dhingra LS, Nagori A, Kaur H, Sethi T. Estimating the impact of health systems factors on antimicrobial resistance in priority pathogens. *J Glob Antimicrob Resist.* 2022 Sep;30:133-142.

Impact Calculation: Counterfactual Analysis



Ceftriaxone High Income

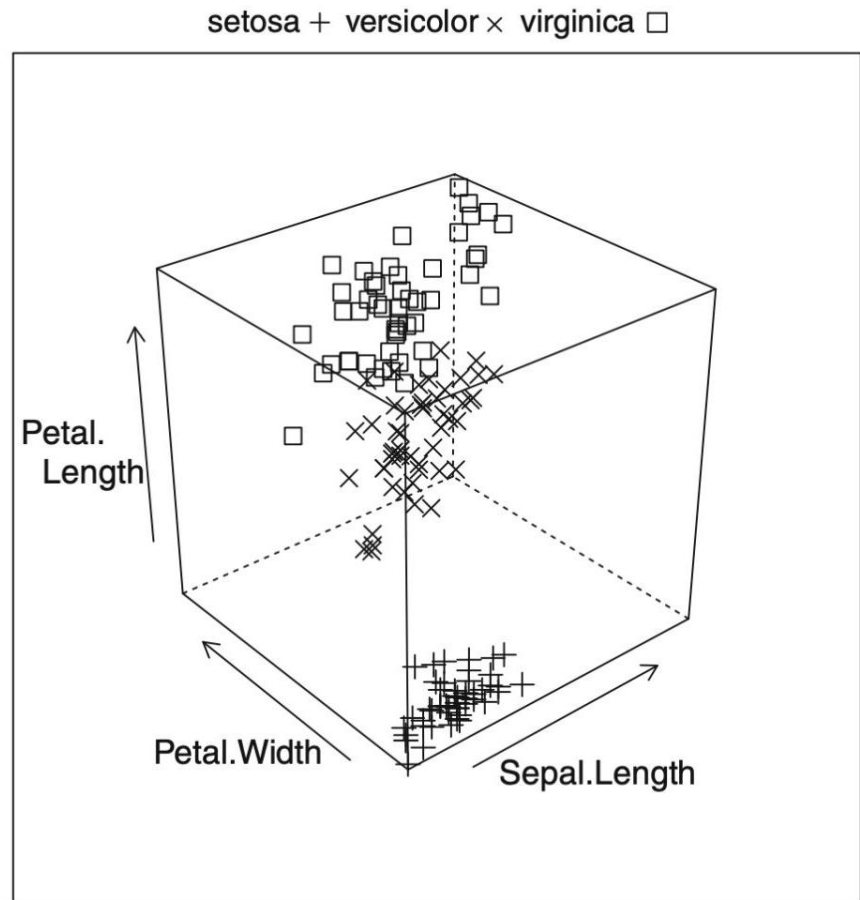
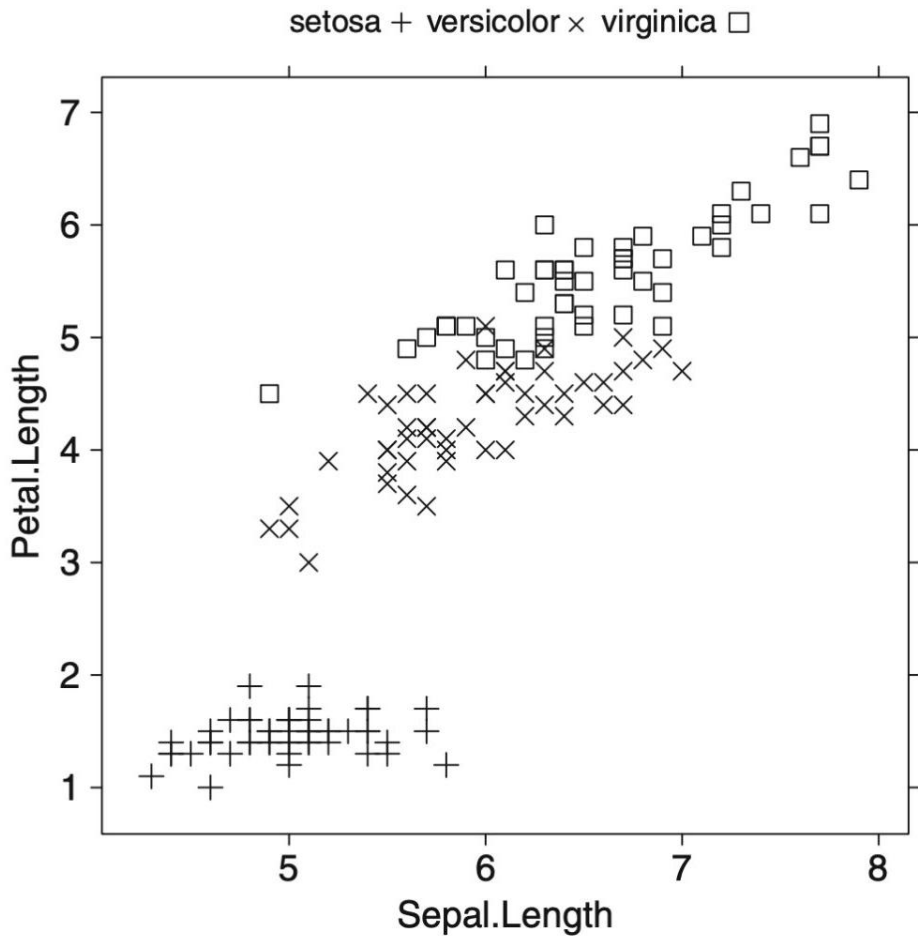
Ceftriaxone Middle Income

Key Steps in Building a BN Model

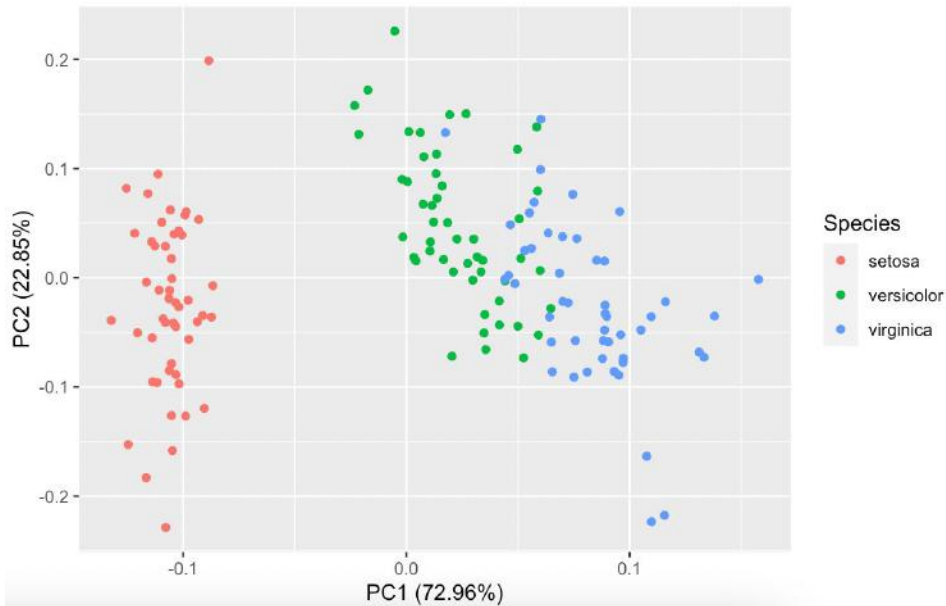
- Learn Structure
 - Score Based
 - Constraint Based
- Validate Structure
 - Bootstrapping
 - Domain Based Sanitization
- Conduct Inference
 - Exact Inference
 - Approximate Inference



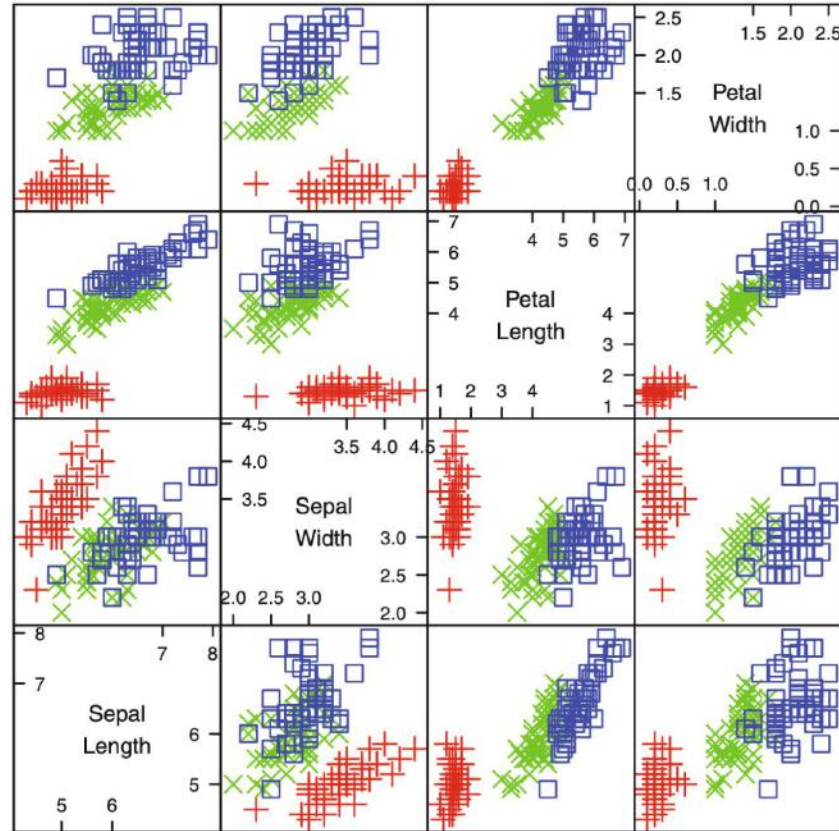
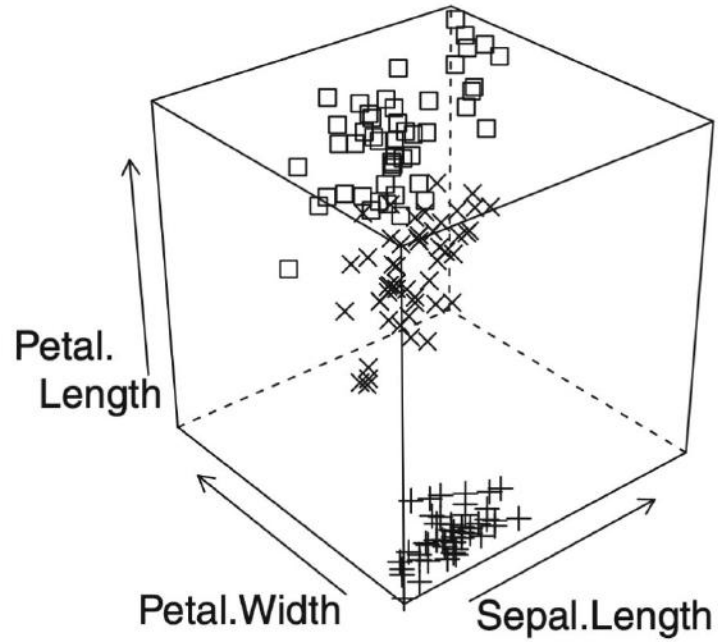
High Dimensional Data Principal Component Analysis



```
1 library(ggfortify)
2 df <- iris[1:4]
3 pca_res <- prcomp(df, scale. = TRUE)
4
5 autoplot(pca_res)
6 autoplot(pca_res, data = iris, colour = 'Species')
```



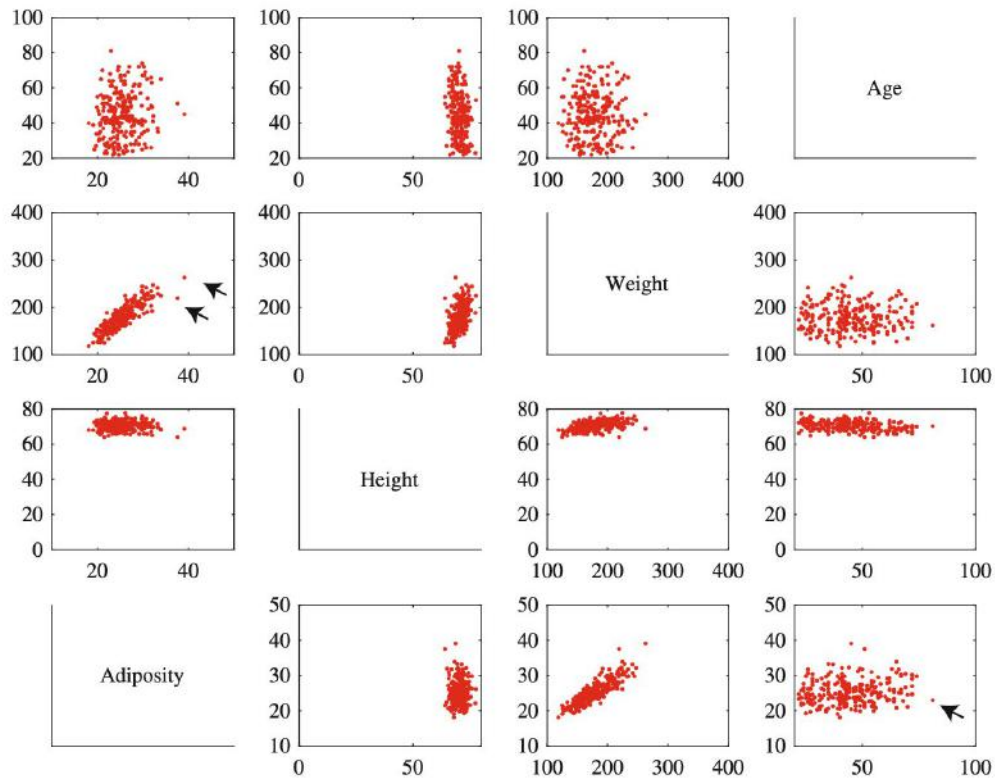
setosa + versicolor × virginica □



Scatter Plot Matrix



Height-Weight




rplots



Challenges of High Dimensional Data

Remember This: *High dimensional data does not behave in a way that is consistent with most people's intuition. Points are always close to the boundary and further apart than you think. This property makes a nuisance of itself in a variety of ways. The most important is that only the simplest models work well in high dimensions.*


$$\text{cov}(\{x\}, \{y\}) = \frac{\sum_i (x_i - \text{mean}(\{x\}))(y_i - \text{mean}(\{y\}))}{N}$$

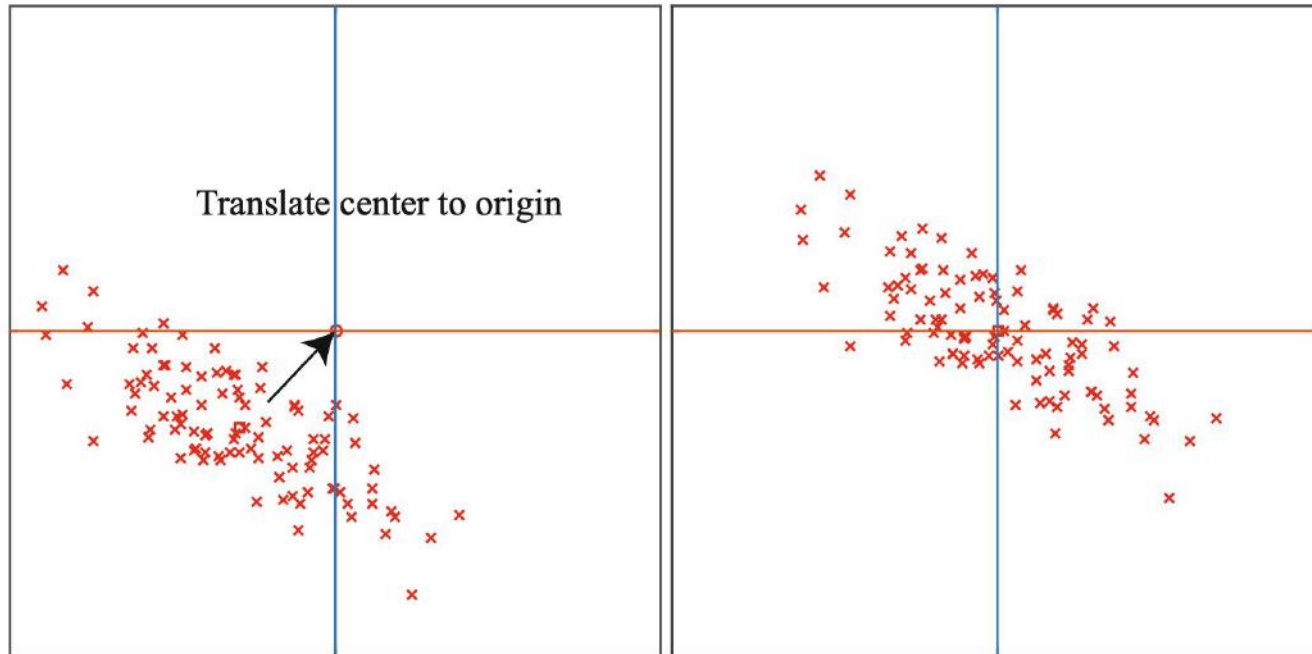
$$\text{corr}(\{(x, y)\}) = \frac{\text{cov}(\{x\}, \{y\})}{\sqrt{\text{cov}(\{x\}, \{x\})} \sqrt{\text{cov}(\{y\}, \{y\})}}.$$

$$\mathbf{C} \quad \text{Covmat}(\{\mathbf{x}\}) = \frac{\sum_i (\mathbf{x}_i - \text{mean}(\{\mathbf{x}\}))(\mathbf{x}_i - \text{mean}(\{\mathbf{x}\}))^T}{N}$$

Useful Facts: 4.3 *Properties of the Covariance Matrix*

- The j , k 'th entry of the covariance matrix is the covariance of the j 'th and the k 'th components of \mathbf{x} , which we write $\text{cov}(\{x^{(j)}\}, \{x^{(k)}\})$.
- The j , j 'th entry of the covariance matrix is the variance of the j 'th component of \mathbf{x} .
- The covariance matrix is symmetric.
- The covariance matrix is always positive semidefinite; it is positive definite, *unless* there is some vector \mathbf{a} such that $\mathbf{a}^T(\mathbf{x}_i - \text{mean}(\{\mathbf{x}_i\})) = 0$ for all i .

Affine Transform $\mathbf{m}_i = \mathcal{A}\mathbf{x}_i + \mathbf{b}$.





$$\begin{aligned} \text{mean}(\{\mathbf{m}\}) &= \text{mean}(\{\mathcal{A}\mathbf{x} + \mathbf{b}\}) & \text{matrix under Affine} \\ &= \mathcal{A}\text{mean}(\{\mathbf{x}\}) + \mathbf{b}, \end{aligned}$$

$$\begin{aligned} \text{Covmat}(\{\mathbf{m}\}) &= \text{Covmat}(\{\mathcal{A}\mathbf{x} + \mathbf{b}\}) \\ &= \frac{\sum_i (\mathbf{m}_i - \text{mean}(\{\mathbf{m}\}))(\mathbf{m}_i - \text{mean}(\{\mathbf{m}\}))^T}{N} \\ &= \frac{\sum_i (\mathcal{A}\mathbf{x}_i + \mathbf{b} - \mathcal{A}\text{mean}(\{\mathbf{x}\}) - \mathbf{b})(\mathcal{A}\mathbf{x}_i + \mathbf{b} - \mathcal{A}\text{mean}(\{\mathbf{x}\}) - \mathbf{b})^T}{N} \\ &= \frac{\mathcal{A} \left[\frac{\sum_i (\mathbf{x}_i - \text{mean}(\{\mathbf{x}\}))(\mathbf{x}_i - \text{mean}(\{\mathbf{x}\}))^T}{N} \right] \mathcal{A}^T}{N} \\ &= \mathcal{A}\text{Covmat}(\{\mathbf{x}\})\mathcal{A}^T. \end{aligned}$$

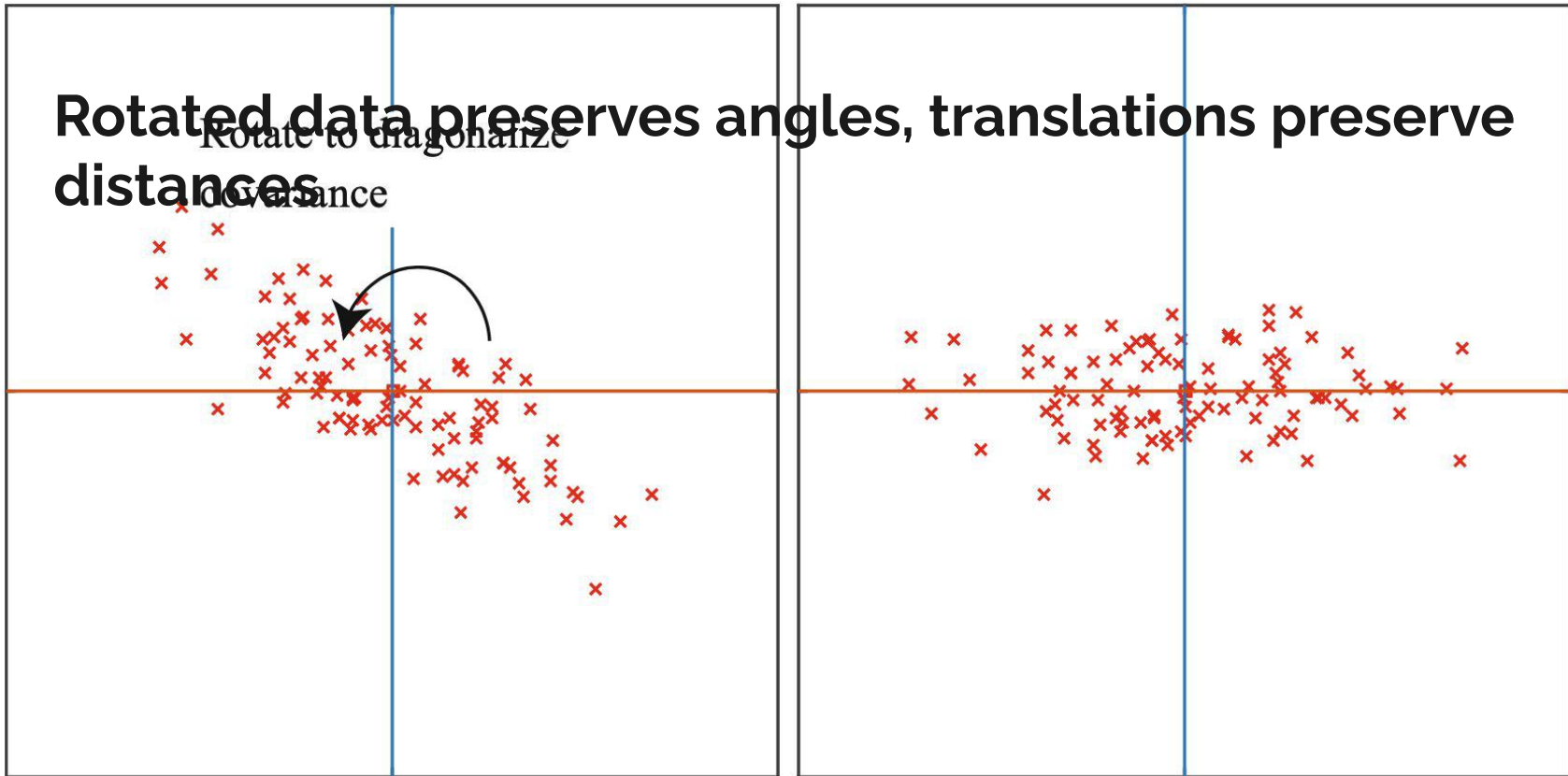


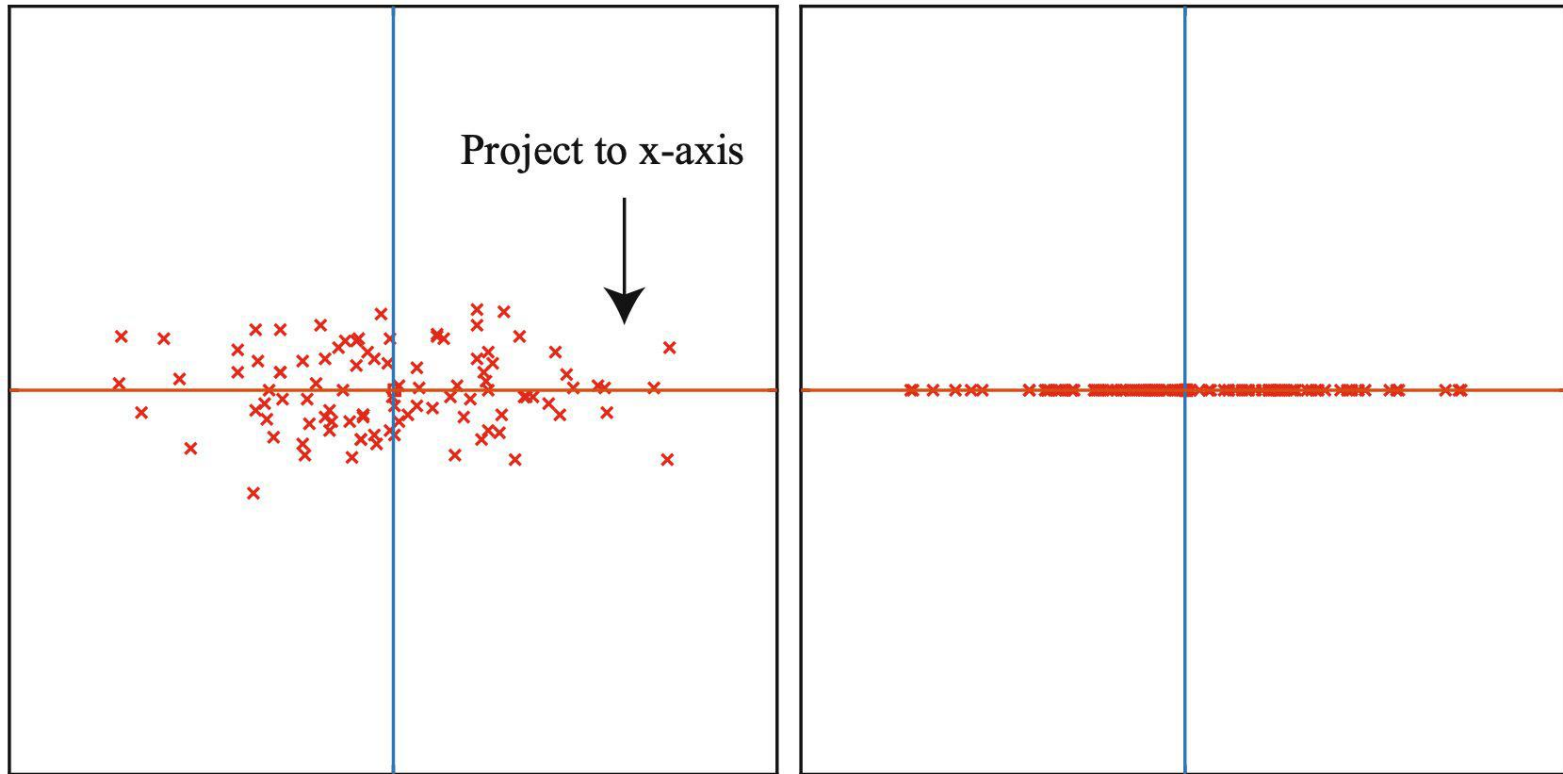
Principal Component Analysis

Rotated data preserves angles, translations preserve distances

Rotate to diagonalize covariance

variance







For any $m \times p$ matrix \mathcal{X} , it is possible to obtain a decomposition

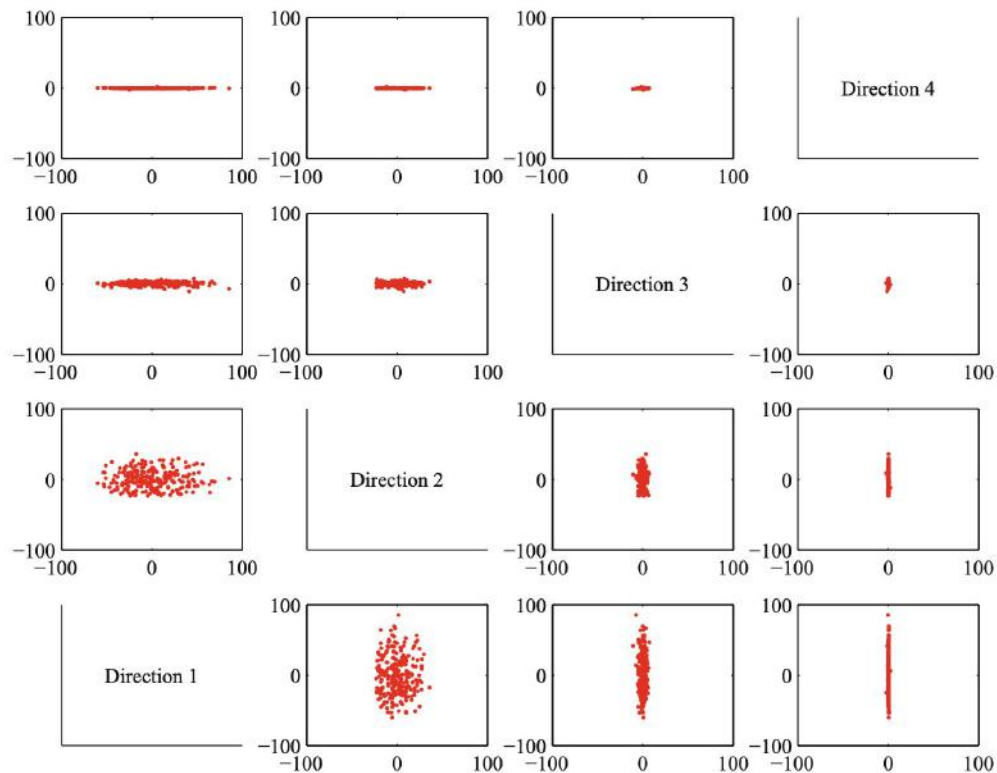
$$\mathcal{X} = \mathcal{U}\Sigma\mathcal{V}^T$$

where

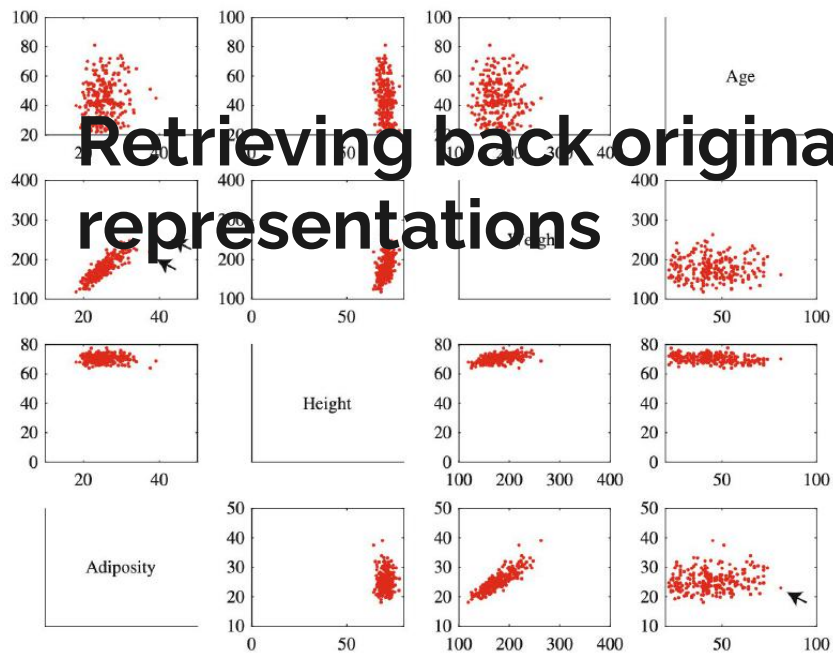
- \mathcal{U} is $m \times m$, \mathcal{V} is $p \times p$, and Σ is $m \times p$ and is diagonal
- The diagonal entries of Σ are non-negative.
- Both \mathcal{U} and \mathcal{V} are orthonormal (i.e., $\mathcal{U}\mathcal{U}^T = \mathcal{I}$ and $\mathcal{V}\mathcal{V}^T = \mathcal{I}$).

SVD yields a low rank approximation

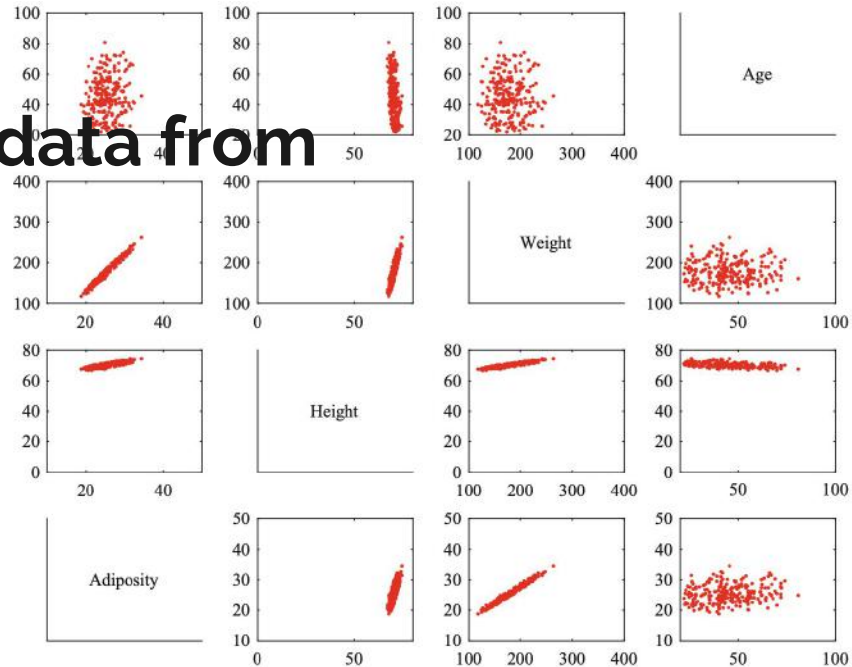
Principal represer



Retrieving back original data from representations

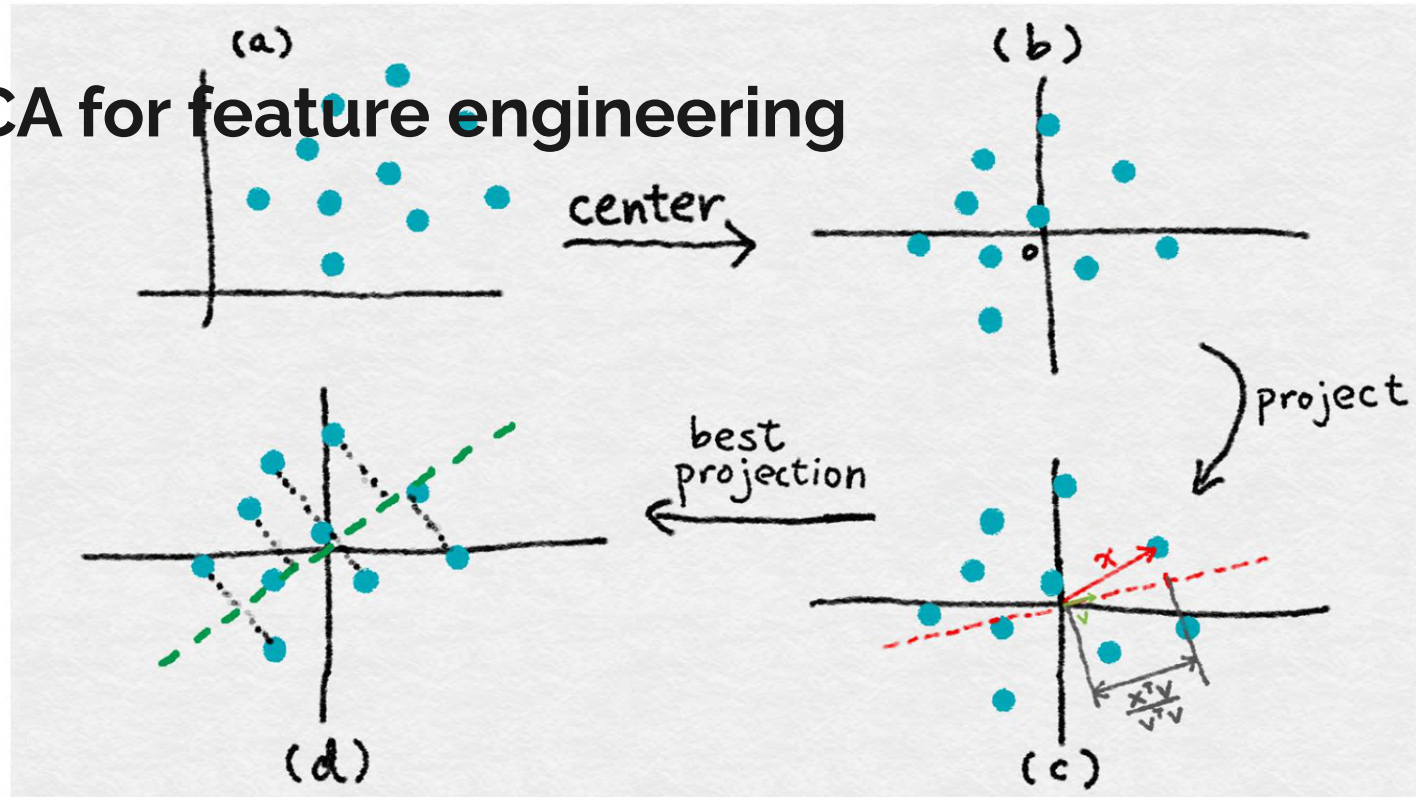


Original

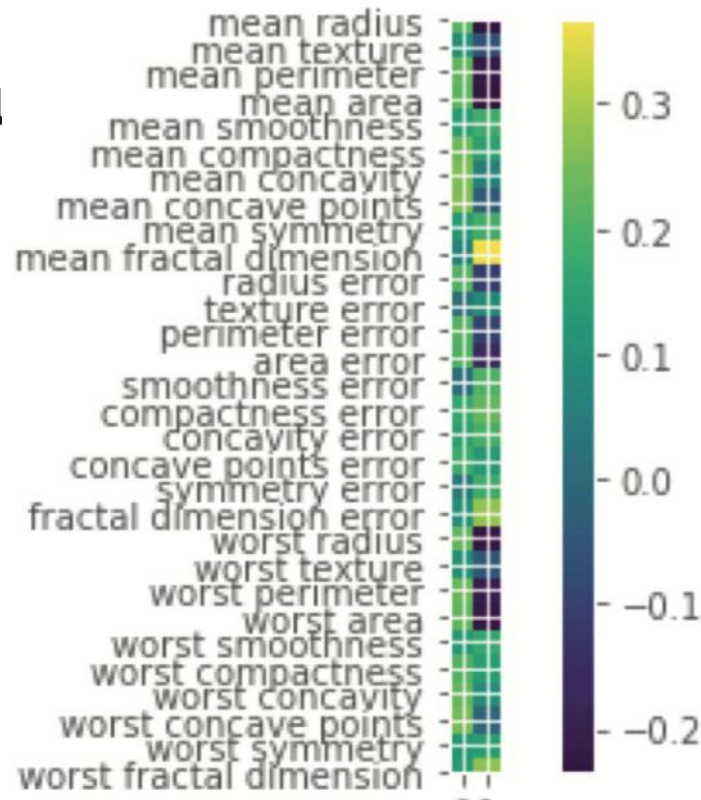


Rotation back from PCA

PCA for feature engineering



PCA for featu



Remember This: *Data items in a d -dimensional dataset can usually be represented with good accuracy as a weighted sum of a small number s of d -dimensional vectors, together with the mean. This means that the dataset lies on an s -dimensional subspace of the d -dimensional space. The subspace is spanned by the principal components of the data.*

Remember This: *Given a d -dimensional dataset where data items have had independent random noise added to them, representing each data item on $s < d$ principal components can result in a representation which is on average closer to the true underlying data than the original data items. The choice of s is application dependent.*

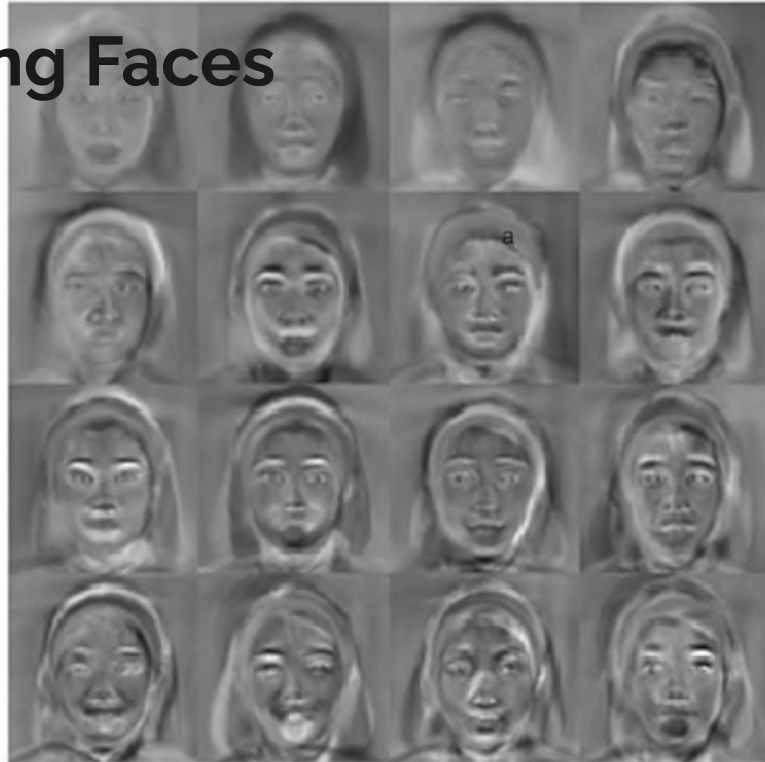


Example: Approximating Faces

Mean image from Japanese Facial Expression dataset



First sixteen principal components of the Japanese Facial Expression dataset





Sample Face Image

mean

1

5

10

20

50

100

