



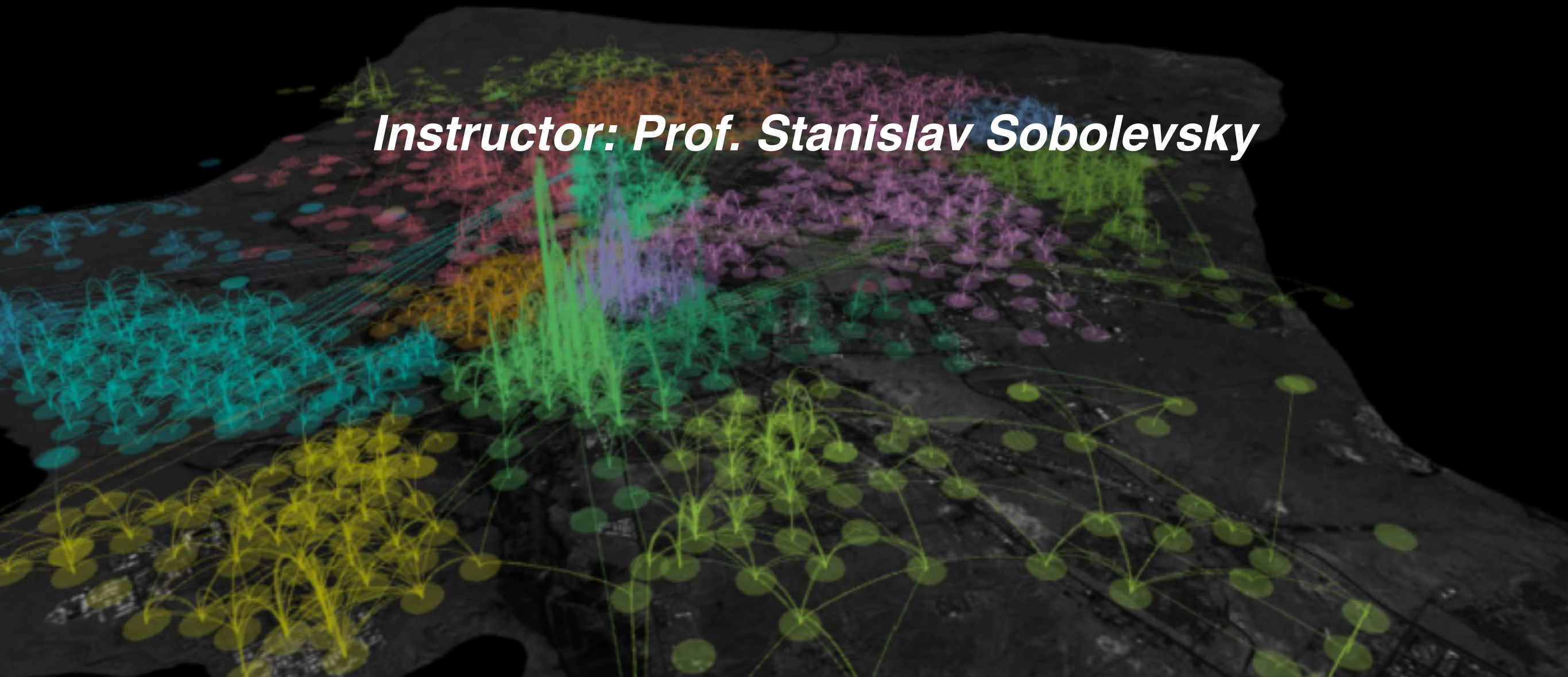
CENTER FOR URBAN
SCIENCE+PROGRESS

APPLIED DATA SCIENCE

6004.002, Fall 2019

Principle component analysis. Classification

Instructor: Prof. Stanislav Sobolevsky

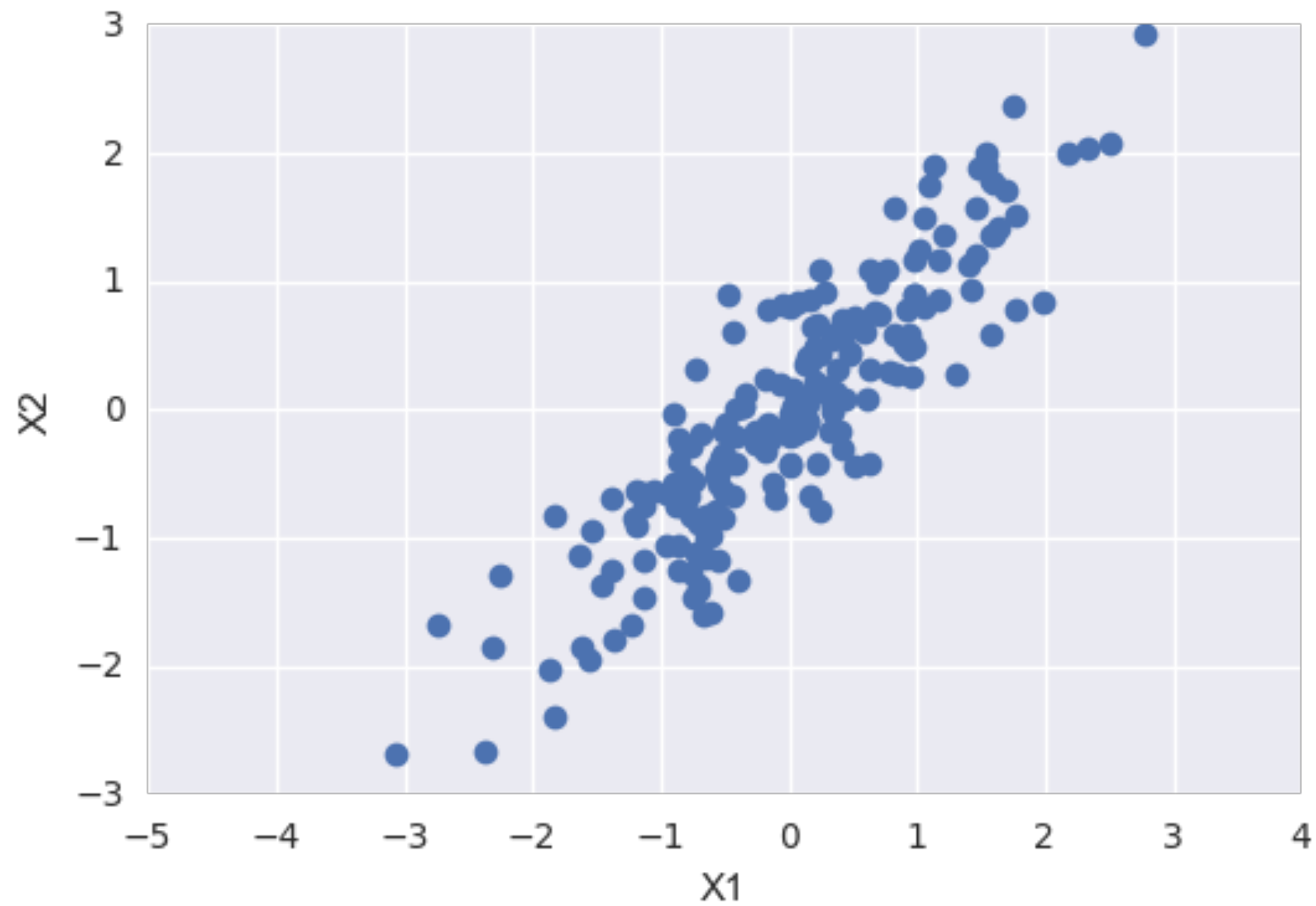


Principal component analysis

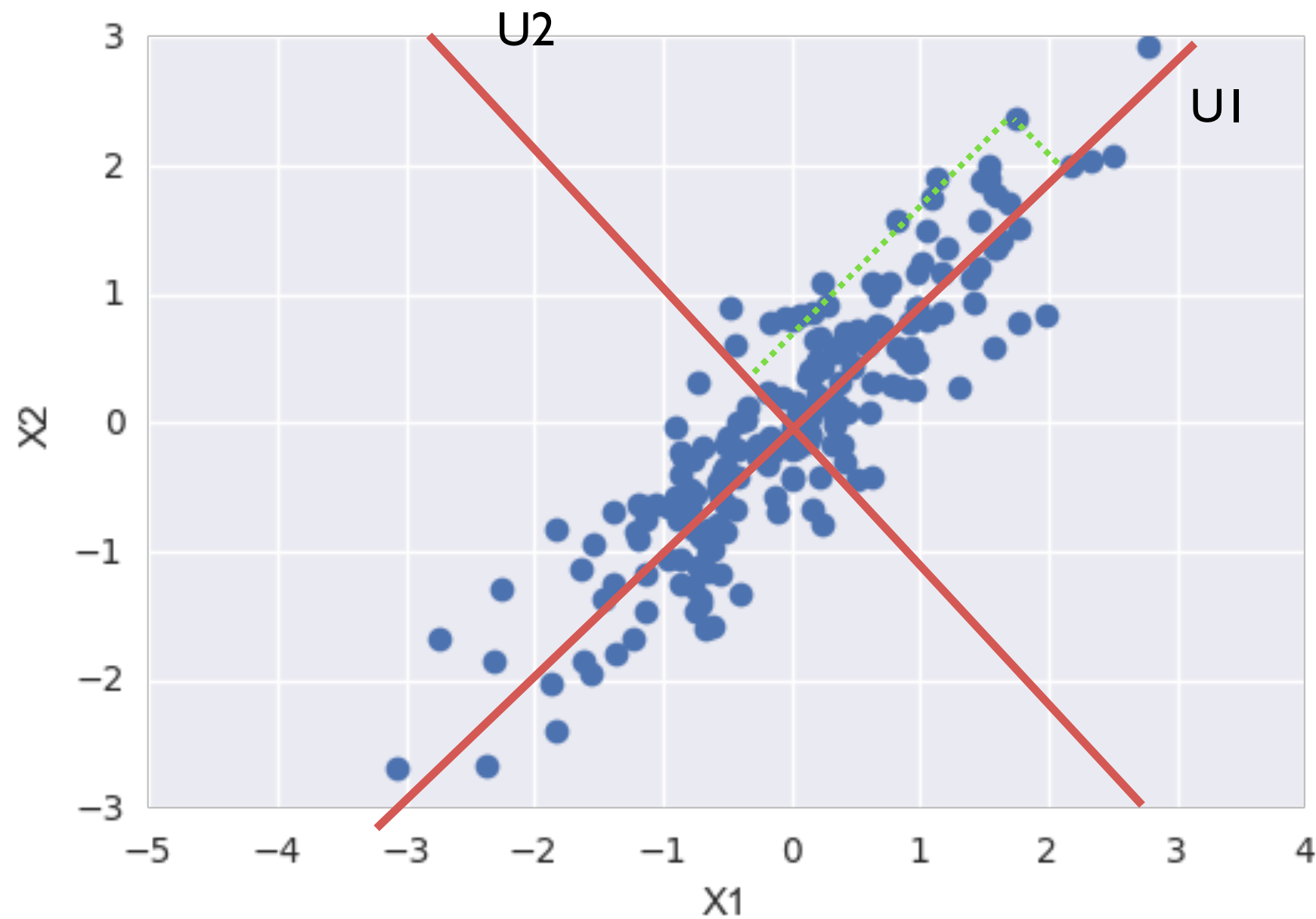
correlated features to uncorrelated

$$(x_1, x_2, x_3, \dots, x_n) \rightarrow (u_1, u_2, u_3, \dots, u_n)$$

Original data

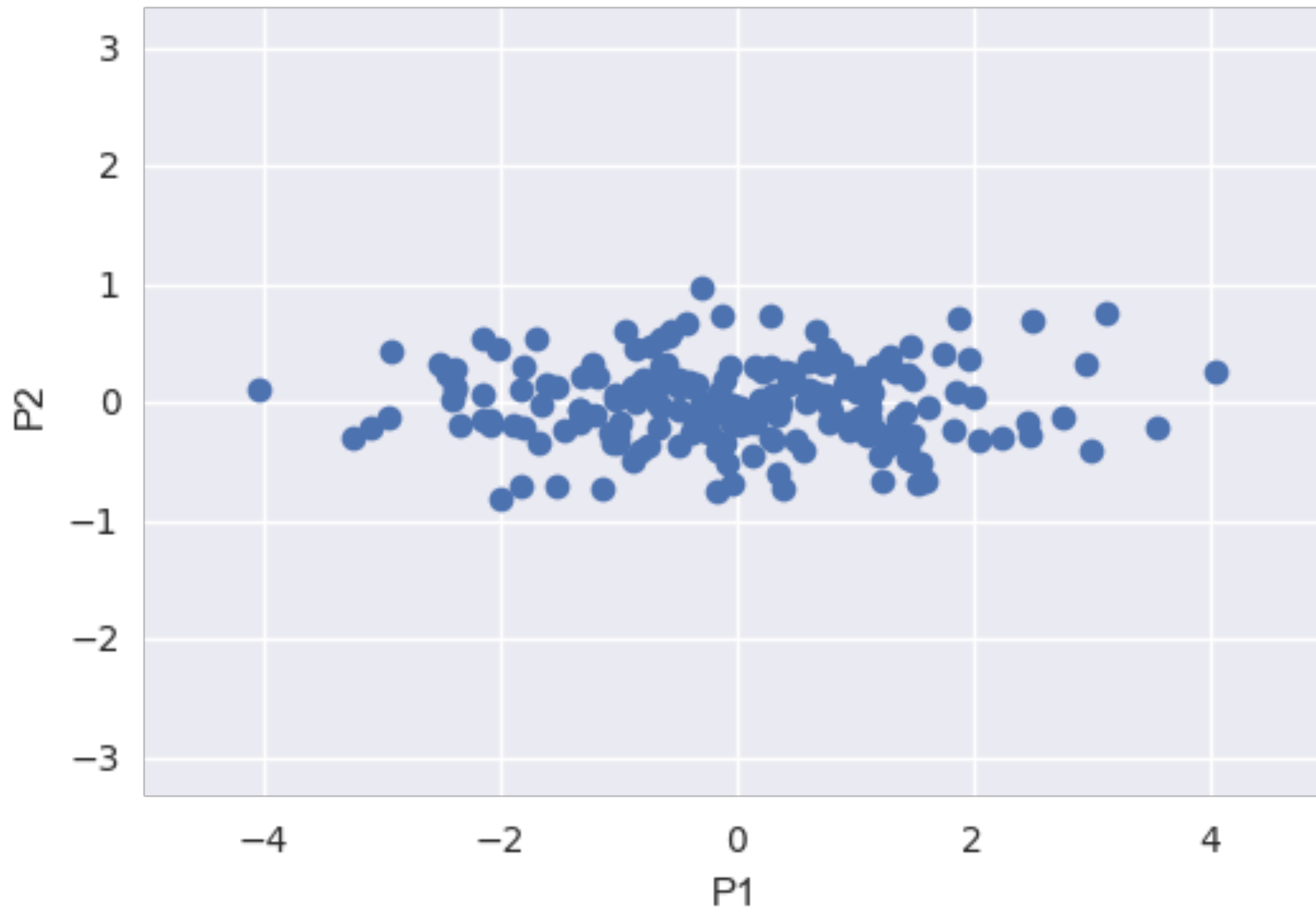


Original data - new system of coordinates

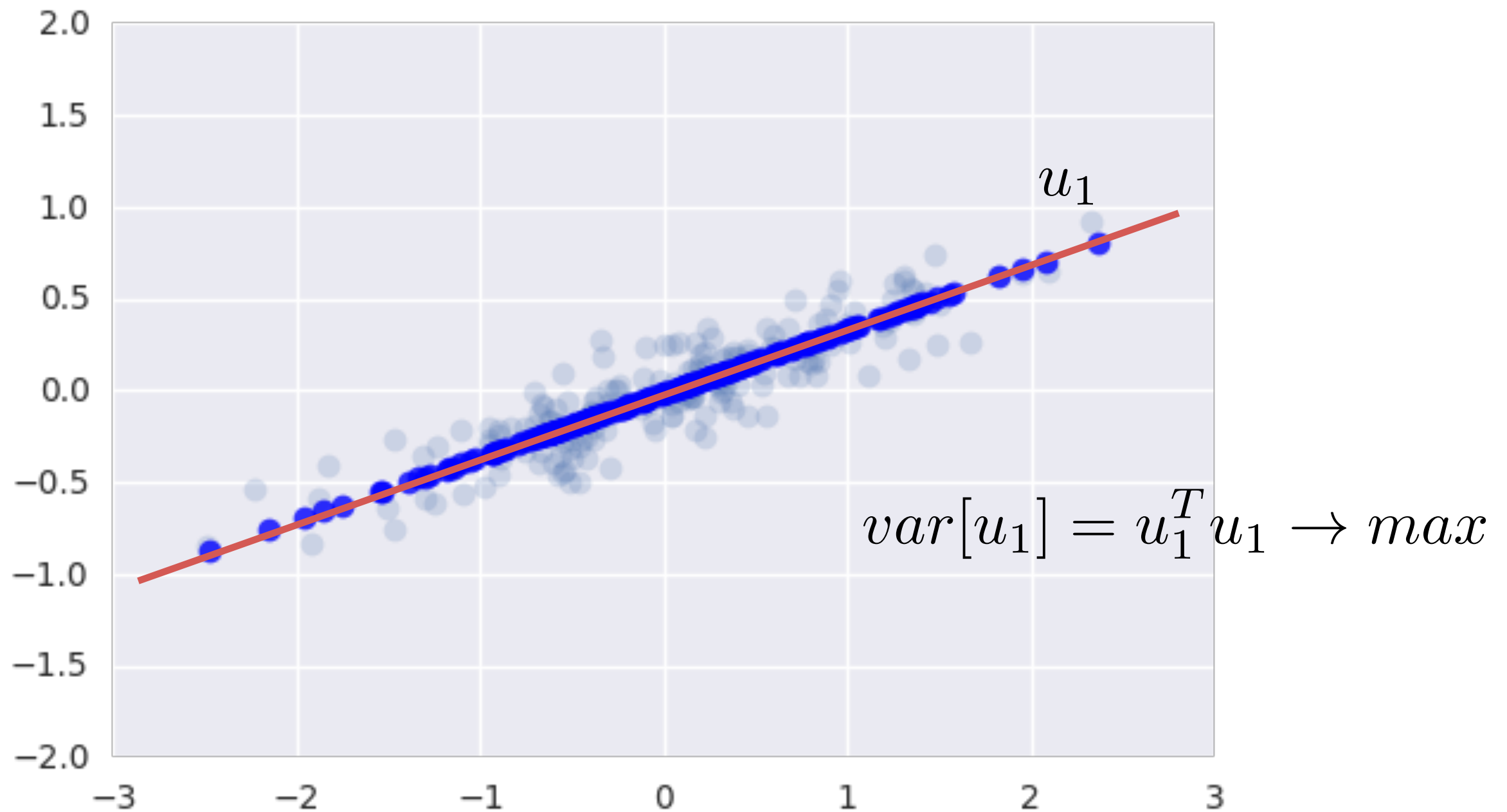




Uncorrelated data (rotation)



Almost the same information



Principal components - maths

Given the standardized data $X = \{x_i^j, i = 1..n, j = 1..N\}$

Find uncorrelated latent factors U (or P)

$$u_j = x_1 v_j^1 + x_2 v_j^2 + \dots + x_n v_j^n$$

Matrix for for the rotation transform

$$u_i = X v_i$$

$$U = XV$$

$$V - n \times p$$

$$U - N \times p$$

Look for linear combinations of factors

one-by-one

Principal components - optimization objective

Looking for $u_j = x_1 v_j^1 + x_2 v_j^2 + \dots + x_n v_j^n$

Start with $u_1 = X v_1$

Such that $var[u_1] = u_1^T u_1 \rightarrow max$

$$var[u_1] = u_1^T u_1 = (X v_1)^T (X v_1) = v_1^T X^T X v_1$$

Find $v_1 = \operatorname{argmax}_{v_1: v_1^T v_1 = 1} v_1^T X^T X v_1$

$$v_i = \operatorname{argmax}_{v_i: v_i^T v_i = 1, v_i^T v_j = 0, j < i} v_i^T X^T X v_i$$



Principal components - answer

So which vectors v maximize the quantity below?

$$v_i = \operatorname{argmax}_{v_i: v_i^T v_i = 1, v_i^T v_j = 0, j < i} v_i^T X^T X v_i$$

Eigenvectors $\lambda_i v_i = X^T X v_i \quad v_i^T v_i = 1 \quad v_i^T v_j = 0$

$$\operatorname{Var}[u_i] = v_i^T X^T X v_i = \lambda_i v_i^T v_i = \lambda_i \quad \lambda_1 > \lambda_2 > \dots > \lambda_n > 0$$

Projection for the leading PC v_1

Is the leading eigenvector with the max eigenvalue



Recall the concept of eigenvectors/eigenvalues

$$\lambda v = Av \quad \lambda - \text{eigenvalue}, v - \text{eigenvector}$$

$$(\lambda I - A)v = 0 \quad \det(\lambda I - A) = 0$$

Find

$$\lambda_1, \lambda_2, \dots, \lambda_n$$

$$v_1, v_2, \dots, v_n$$

Define up to a scaling factor, can require unit length

$$v_i \rightarrow C v_i \quad |v_i| = 1$$

When $\lambda_i \neq \lambda_j \Rightarrow v_i^T v_j = 0$

$$A^T = A$$

Proof $v_i^T A v_j = \lambda_j v_i^T v_j$

$$v_i^T A v_j = (A v_i)^T v_j = \lambda_i v_i^T v_j$$

Principal components - proof

$$v_i = \operatorname{argmax}_{v_i: v_i^T v_i = 1, v_i^T v_j = 0, j < i} v_i^T X^T X v_i$$

Consider eigenvectors:

$$\lambda_i v_i = X^T X v_i \quad v_i^T v_i = 1 \quad \lambda_1 > \lambda_2 > \dots > \lambda_n > 0$$

$$v_i^T X^T X v_i = \lambda_i v_i^T v_i = \lambda_i$$

$$w = e_1 v_1 + e_2 v_2 + \dots + e_n v_n \quad w^T w = e_1^2 + e_2^2 + \dots + e_n^2 = 1$$

$$w^T X^T X w = \lambda_1 e_1^2 + \lambda_2 e_2^2 + \dots + \lambda_n e_n^2 \rightarrow \max$$

$$w = v_1, e_1 = 1, e_2 = e_3 = \dots = e_n = 0$$

Principal components - select by variation

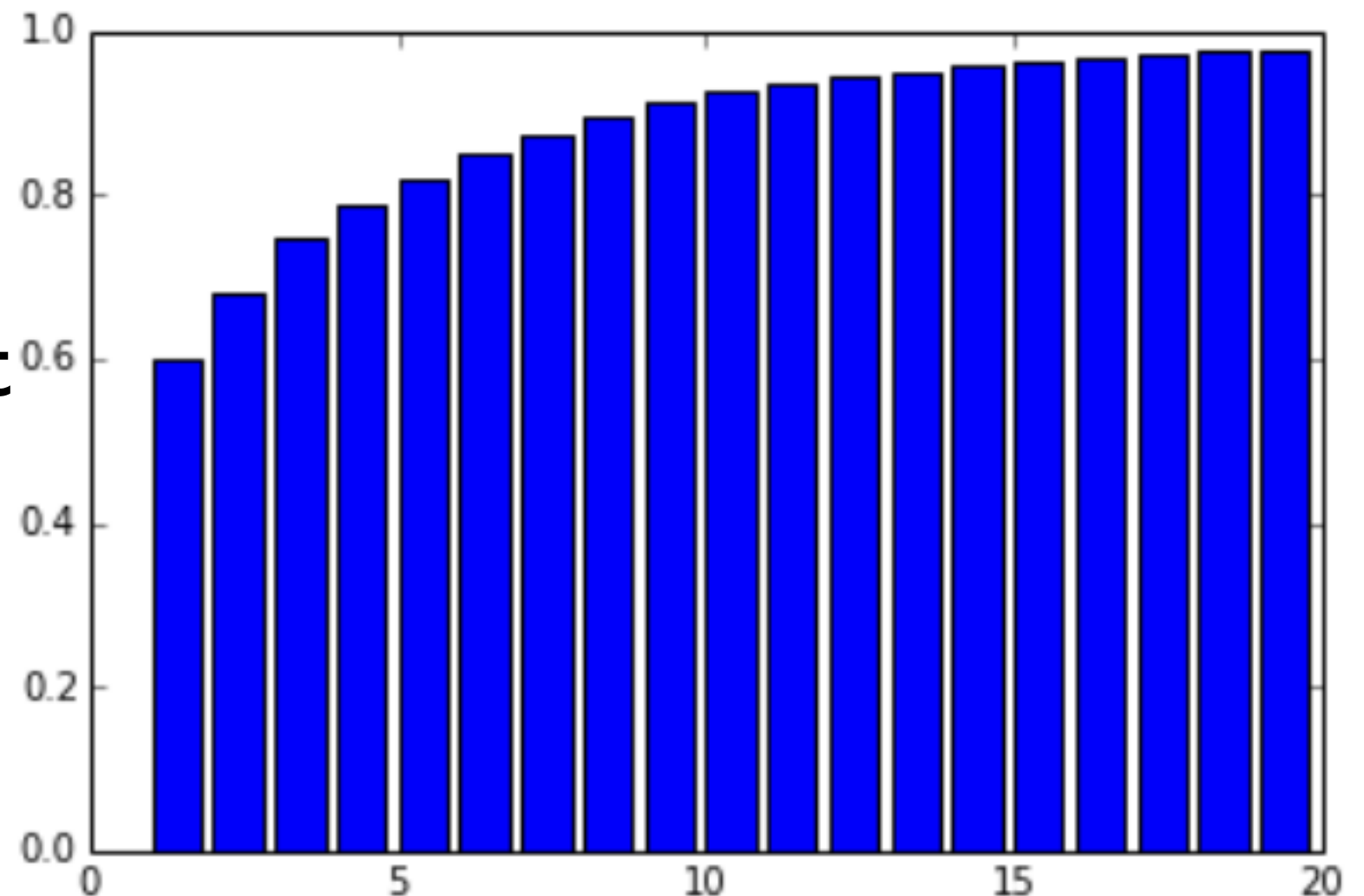
Information explained by a PC u_i

$$\text{Var}[u_i] = \lambda_i$$

% of total $\lambda_i / \sum_j \lambda_j$

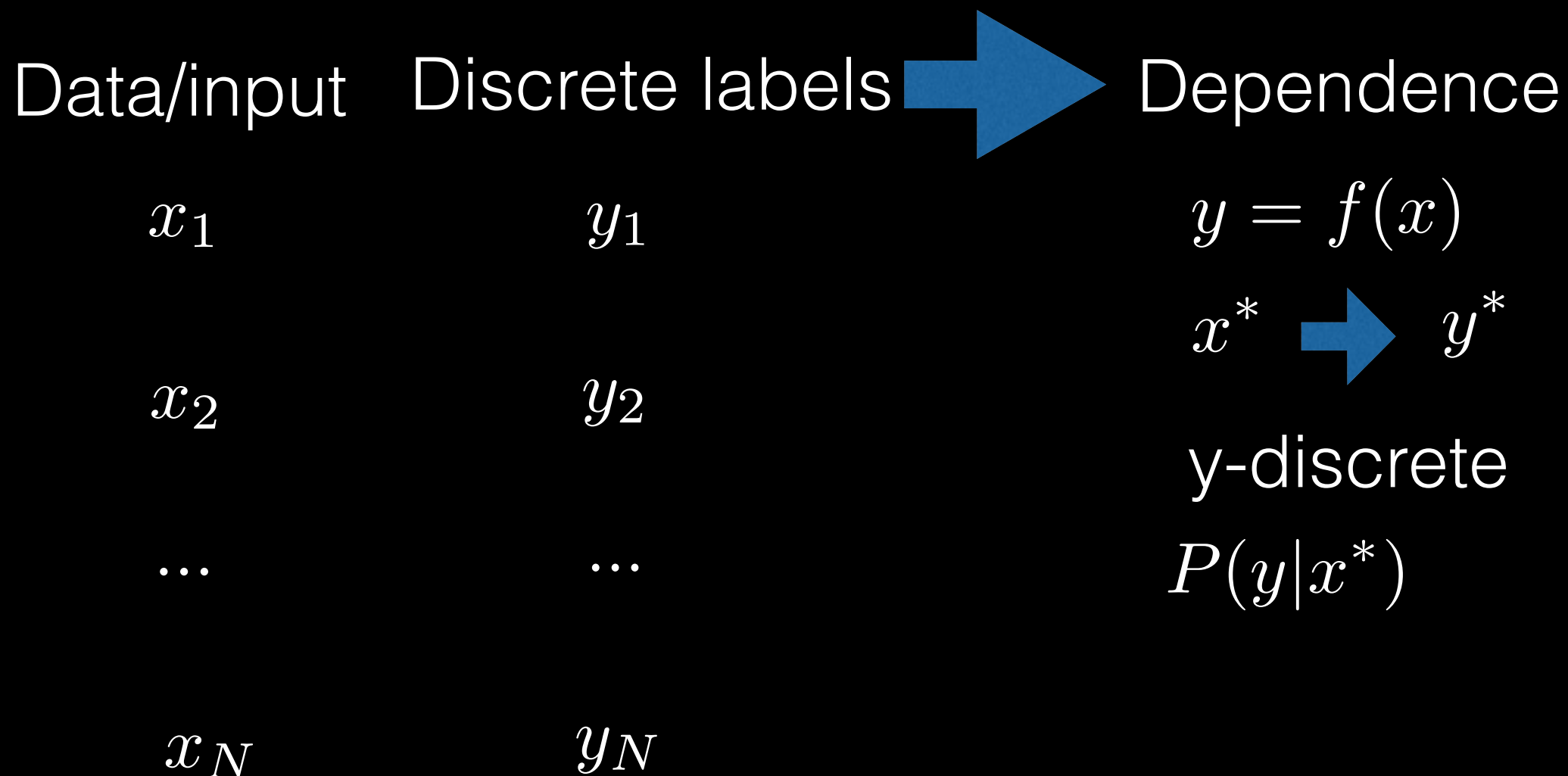
Take first k, such that

$$\frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^n \lambda_i} \geq \alpha$$



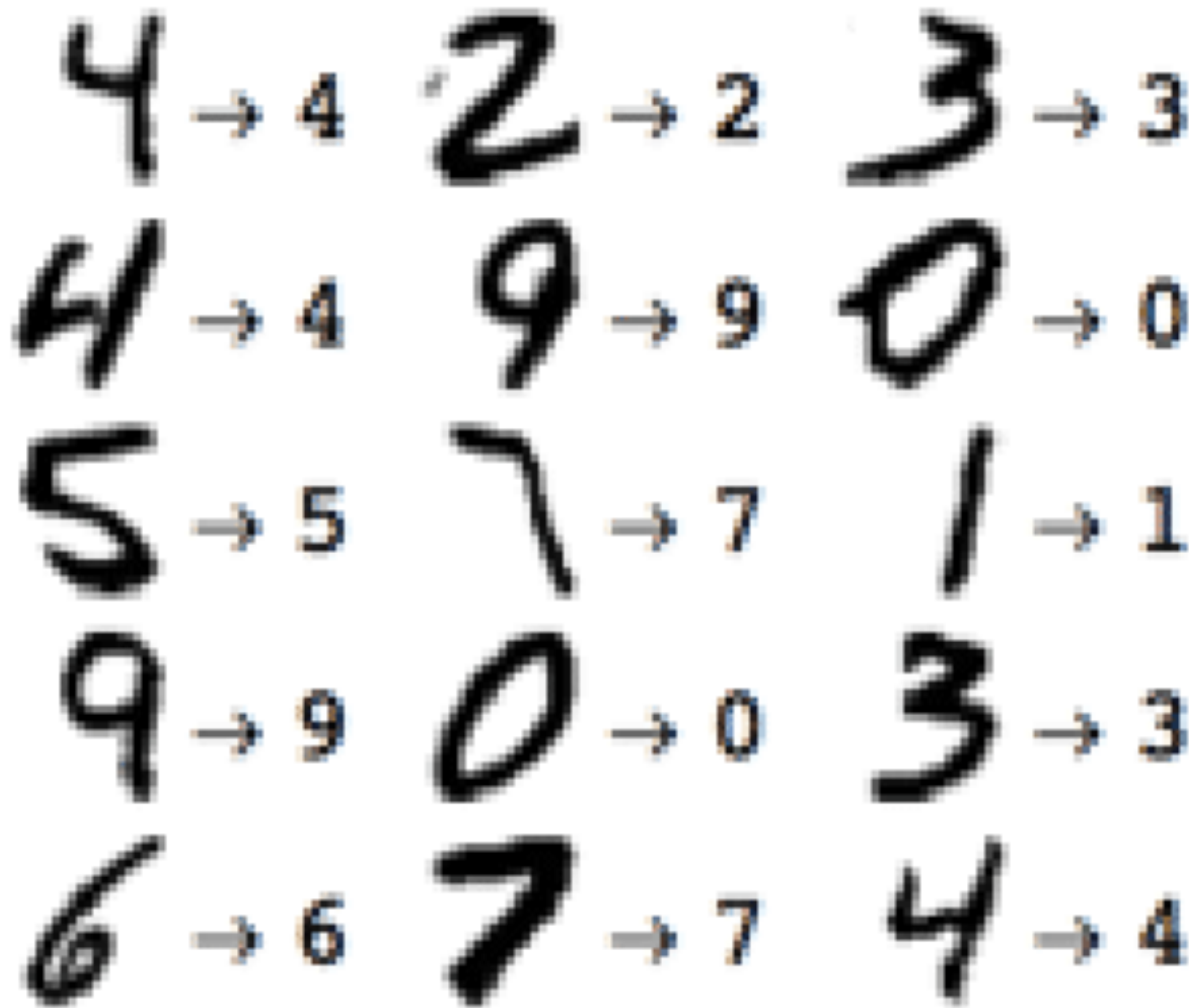
Information explained by leading PC's

Supervised learning: classification



$$X = \{x_i, i = 1..N\} = \{x_i^j, i = 1..N, j = 1..n\}$$

$$Y = \{y_i, i = 1..N\}$$



wolfram.com

These materials are included under the fair use exemption and are restricted from further use

Binary/multiple classification

$$P(y|x^*) = \text{Bern}(y|\mu(x^*))$$

$$y = \begin{cases} 1 & \text{event happened} \\ 0 & \text{event not happened} \end{cases}$$

Multiple classification $y=k$

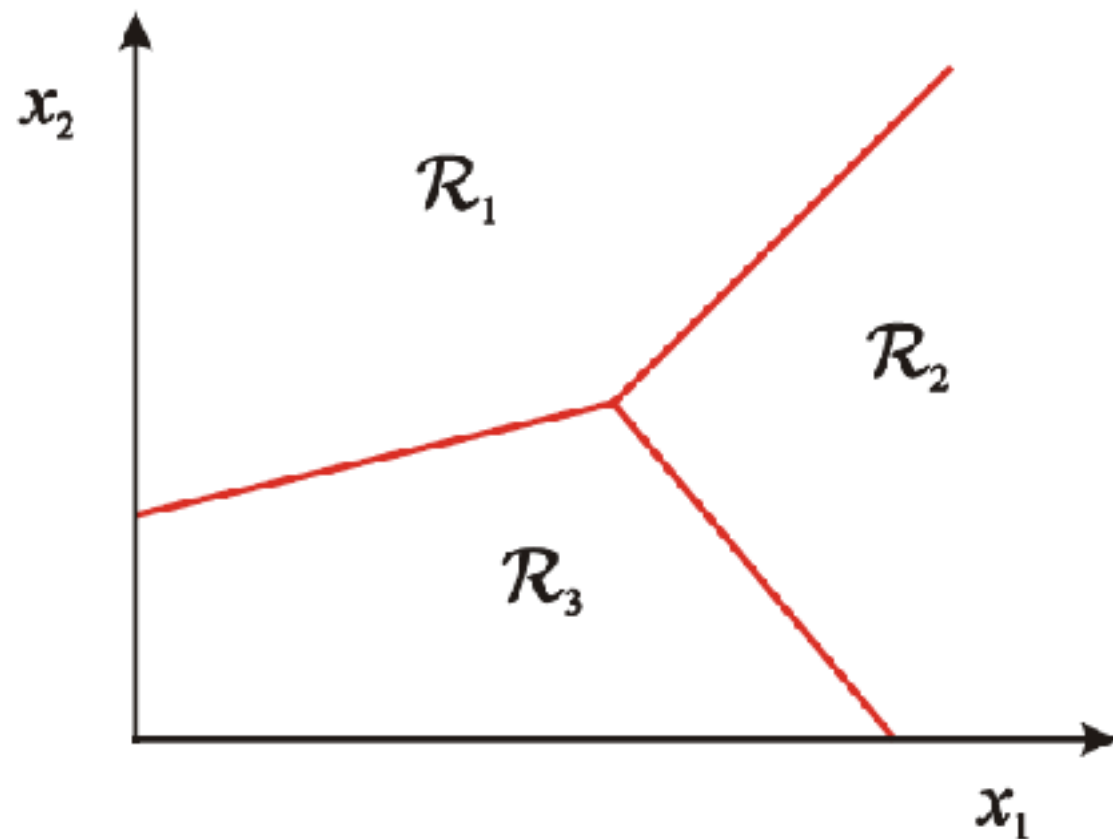
Multiple classification to binary:

$c=k$

$c \sim k$



Multi-class classification

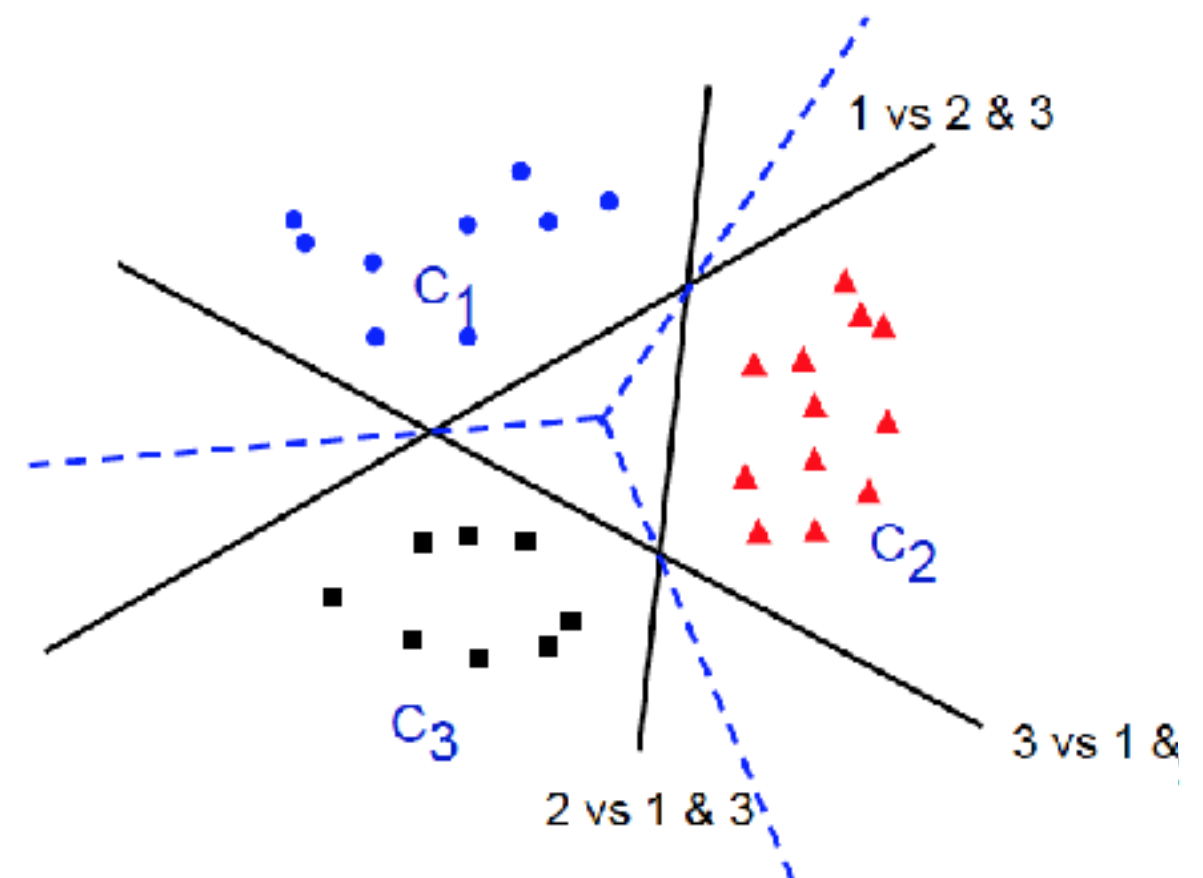


Is \mathcal{R}_1 ? Rather than \mathcal{R}_2 or \mathcal{R}_3

Is \mathcal{R}_2 ? Rather than \mathcal{R}_1 or \mathcal{R}_3

Is \mathcal{R}_3 ? Rather than \mathcal{R}_1 or \mathcal{R}_2

4 → 4	2 → 2	3 → 3
4 → 4	9 → 9	0 → 0
5 → 5	7 → 7	1 → 1
9 → 9	0 → 0	3 → 3
6 → 6	7 → 7	4 → 4

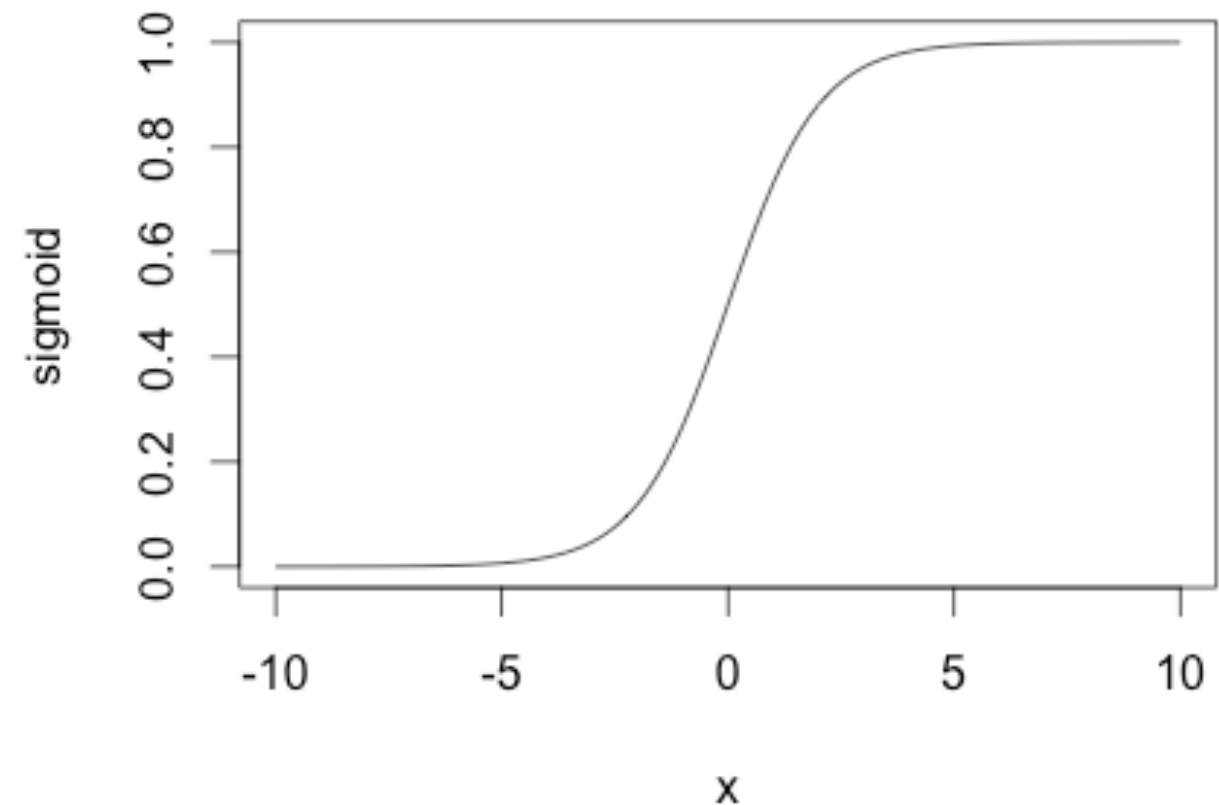


Logistic regression

$$P(y|x, \beta) = \text{Bern}(y|\mu(x, \beta))$$

$$\mu(x, \beta) = f(x\beta)$$

$$f(x) = \sigma(x) = \frac{e^x}{1 + e^x} = \frac{1}{1 + e^{-x}}$$



Logistic regression

$$P(y|x, \beta) = \text{Bern}(y|\mu(x, \beta)) \quad f(x) = \sigma(x) = \frac{e^x}{1 + e^x} = \frac{1}{1 + e^{-x}}$$

$$\mu(x, \beta) = f(x\beta)$$

$$P(y = 1) = \sigma(x\beta) = \frac{\exp(x\beta)}{1 + \exp(x\beta)} = \frac{1}{1 + \exp(-x\beta)}$$

$$P(y = 0) = 1 - P(y = 1) = \frac{1}{1 + \exp(x\beta)}$$

Choosing a classifier to best fit the data

$$P(y|x^*)$$

Data/input	Discrete labels	Probability
x_1	y_1	P1
x_2	y_2	P2
...	...	
x_N	y_N	PN

Choose the model based on $P1 * P2 * \dots * PN$

Logistic regression - log-likelihood

$$P(y = 1) = \sigma(x\beta) = \frac{\exp(x\beta)}{1 + \exp(x\beta)} = \frac{1}{1 + \exp(-x\beta)}$$

$$L = \prod_i P(y = y_i | x_i, \beta)$$

$$P(y = 0) = 1 - P(y = 1) = \frac{1}{1 + \exp(x\beta)}$$

$$\begin{aligned} \log(L) &= \sum_i \log(P(y = y_i | x_i, \beta)) \\ &= \sum_i y_i \log(P(y = 1 | x_i, \beta)) + \sum_i (1 - y_i) \log(P(y = 0 | x_i, \beta)) \\ &= - \sum_i \log(1 + \exp((2y_i - 1)x_i\beta)) \\ \beta &= \operatorname{argmin}_{\beta} \sum_i \log(1 + \exp((2y_i - 1)x_i\beta)) \end{aligned}$$

Accuracy

$$acc = \frac{|\{i : y_i^{est} = y_i^{true}\}|}{|\{i\}|}$$

Not all errors are the same!

Missing spam vs important e-mail as spam

False fire alarm vs missing a fire



Discussion

Fire alarm errors

Types of outcomes

True

True positive: $y_i^{est} = y_i^{true} = 1$

$$y = \begin{cases} 1 \\ 0 \end{cases}$$

event happened

event not happened

True negative: $y_i^{est} = y_i^{true} = 0$

Errors

False positive: $y_i^{est} = 1, y_i^{true} = 0$

False negative: $y_i^{est} = 0, y_i^{true} = 1$

Confusion matrix

$$\left[\begin{array}{c|c} TP & FN \\ \hline FP & TN \end{array} \right]$$

precision $PPV = \frac{TP}{TP + FP}$

true fires among reported fires

sensitivity or recall $TPR = \frac{TP}{TP + FN}$

reported fires among all fires

accuracy $ACC = \frac{TP + TN}{TP + TN + FP + FN}$

fraction of true classifications