

Curse of Dimensionality

x1	t
3.77	RED
14.40	GREEN
16.20	RED
9.66	GREEN
18.99	RED
15.59	GREEN
4.58	RED
19.30	GREEN
11.55	RED
4.05	GREEN
1.11	RED
11.76	GREEN
18.36	RED
13.12	GREEN
11.53	RED
14.63	GREEN
6.51	RED
7.52	GREEN
10.19	RED
17.92	GREEN

$\{x_1, x_2, \dots, x_N\}$ $N=20$ $x \in [0, 20]$

$\{t_1, t_2, \dots, t_N\}$ $t \in \{r, g\}$

$p(t = g | x)$

x1	t
3.77	RED
14.40	GREEN
16.20	RED
9.66	GREEN
18.99	RED
15.59	GREEN
4.58	RED
19.30	GREEN
11.55	RED
4.05	GREEN
1.11	RED
11.76	GREEN
18.36	RED
13.12	GREEN
11.53	RED
14.63	GREEN
6.51	RED
7.52	GREEN
10.19	RED
17.92	GREEN

$$P(t=g | R_1) = \frac{1}{4} = 0.25$$

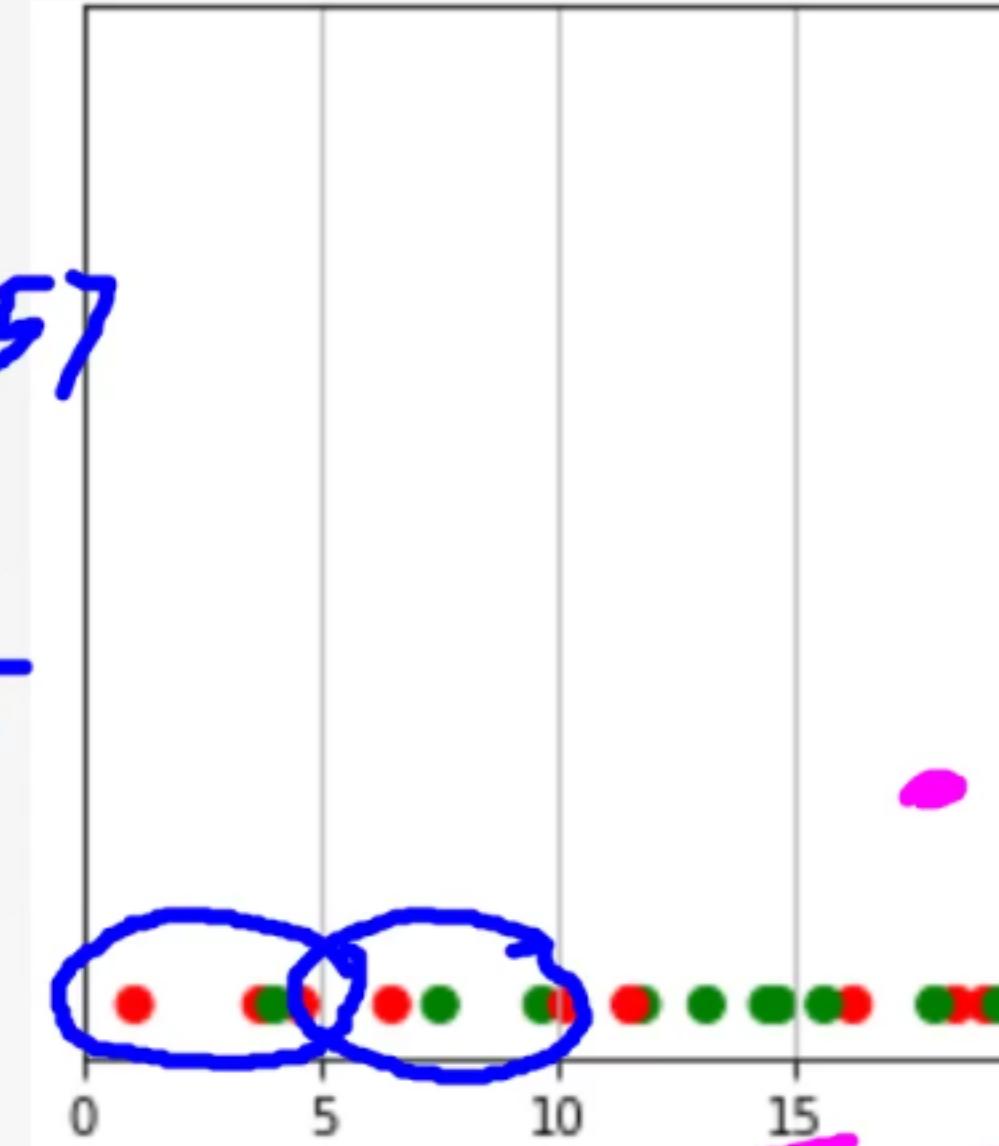
$$R_1 \quad R_2 \quad R_3 \quad R_4$$

↓ ↓ ↓ ↓

$$P(t=g | R_2) = \frac{2}{3} \approx 0.67$$

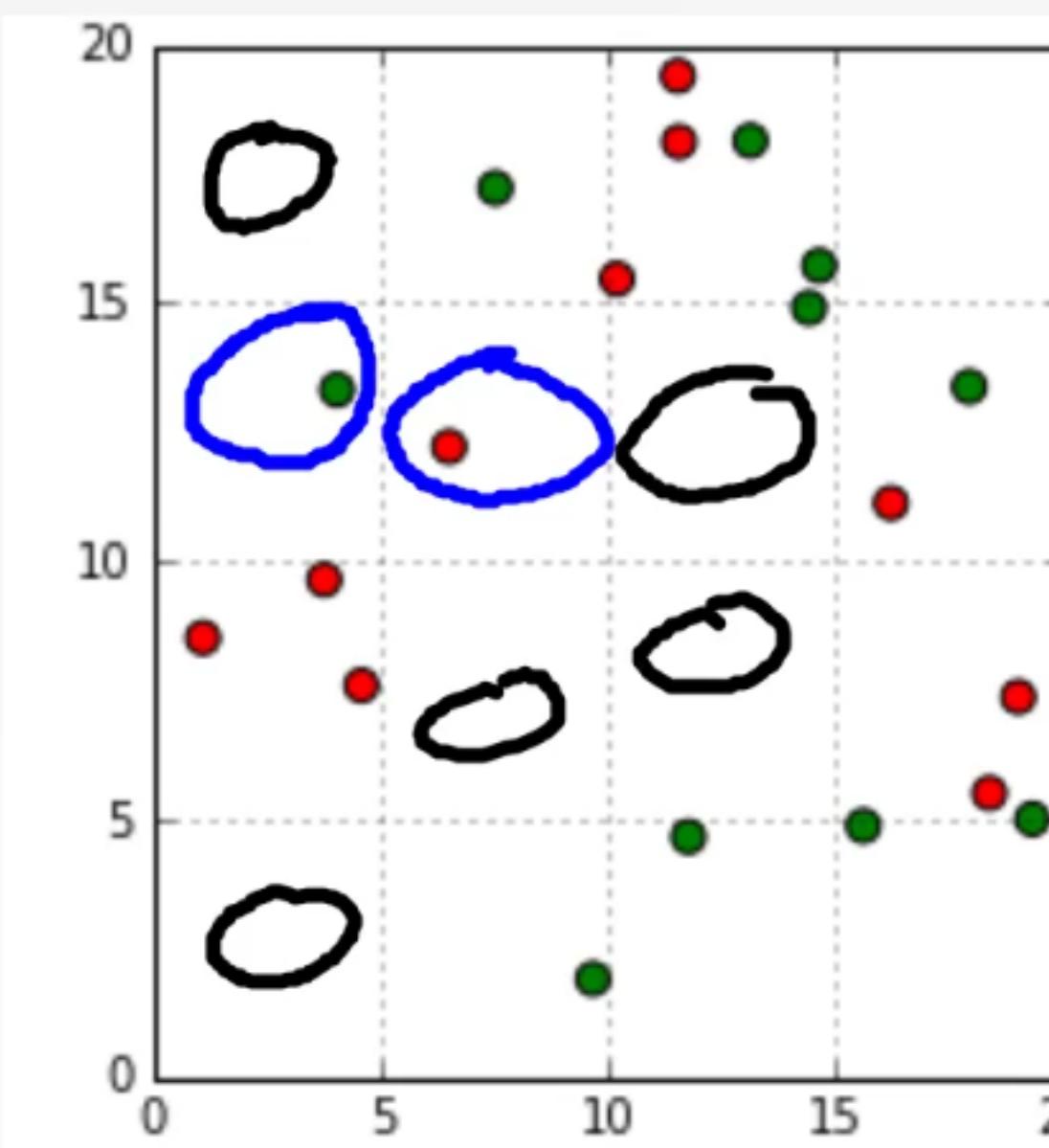
$$P(t=g | R_3) = \frac{4}{7} = 0.57$$

$$P(t=g | R_4) = \frac{3}{6} = \underline{0.5}$$



$$4 \times 4 = 16$$

x1	x2	t
3.77	9.64	RED
14.40	14.91	GREEN
16.20	11.13	RED
9.66	1.91	GREEN
18.99	7.37	RED
15.59	4.87	GREEN
4.58	7.59	RED
19.30	5.01	GREEN
11.55	18.12	RED
4.05	13.32	GREEN
1.11	8.51	RED
11.76	4.66	GREEN
18.36	5.51	RED
13.12	18.14	GREEN
11.53	19.40	RED
14.63	15.72	GREEN
6.51	12.22	RED
7.52	17.23	GREEN
10.19	15.47	RED
17.92	13.38	GREEN



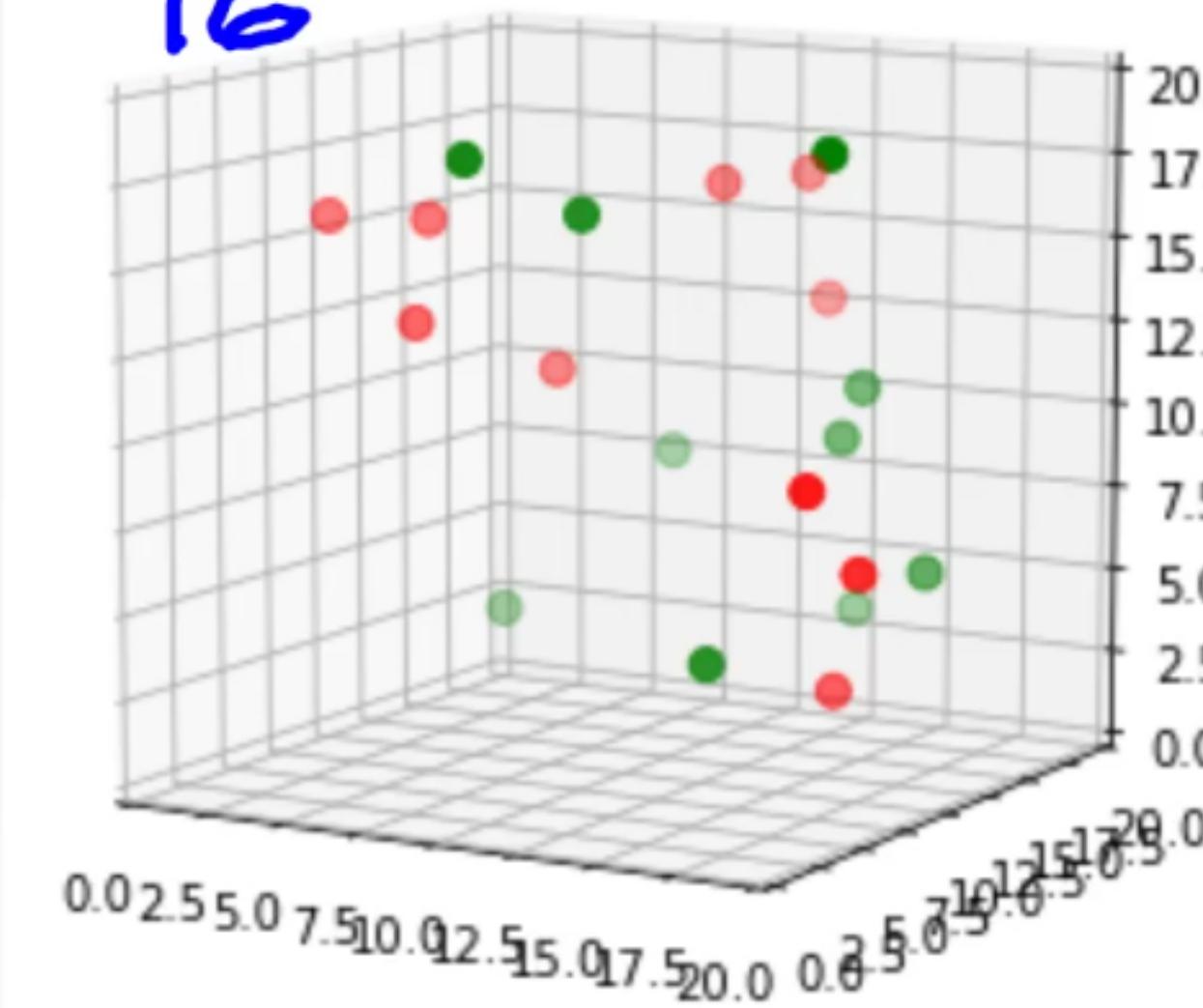
x1	x2	x3	t
3.77	9.64	15.64	RED
14.40	14.91	9.24	GREEN
16.20	11.13	2.52	RED
9.66	1.91	18.73	GREEN
18.99	7.37	6.88	RED
15.59	4.87	4.39	GREEN
4.58	7.59	12.90	RED
19.30	5.01	19.17	GREEN
11.55	18.12	16.54	RED
4.05	13.32	3.36	GREEN
1.11	8.51	15.70	RED
11.76	4.66	16.98	GREEN
18.36	5.51	9.52	RED
13.12	18.14	3.41	GREEN
11.53	19.40	12.57	RED
14.63	15.72	10.64	GREEN
6.51	12.22	11.04	RED
7.52	17.23	7.90	GREEN
10.19	15.47	16.46	RED
17.92	13.38	5.80	GREEN

$$4 \times 4 \times 4 = 64$$

$$\frac{1D}{20} = 5$$

$$\frac{2D}{16} \approx 1.25$$

$$\frac{3D}{64} \approx 0.31$$



Can we increase N?

$$N^{\frac{1}{D}}$$

$$20^{\frac{1}{1}} = \underline{\underline{x}}^{\frac{1}{3}}$$

$$x = \underline{\underline{8'000}}$$

- Data becomes sparse as we add dimensions
- Distances lose meaning in high dimensions

Dimensionality Reduction

- Feature selection

→ weight

colour ←

bhp

rpm

→ fuel type

mpg

- Feature reduction
 - Use all features from the dataset
 - Generate artificial features
 - smaller than the original set
 - retain as much information as possible

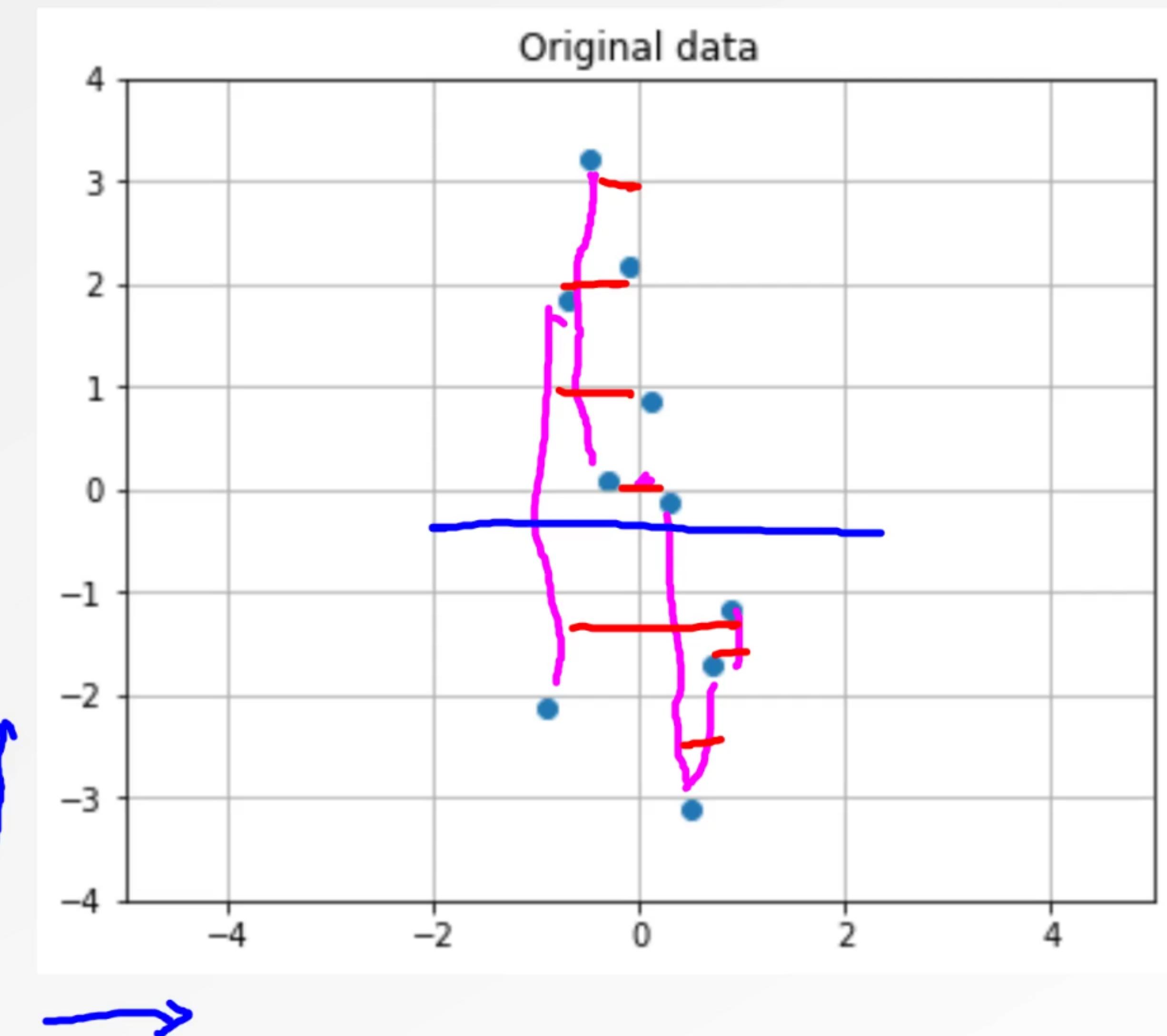


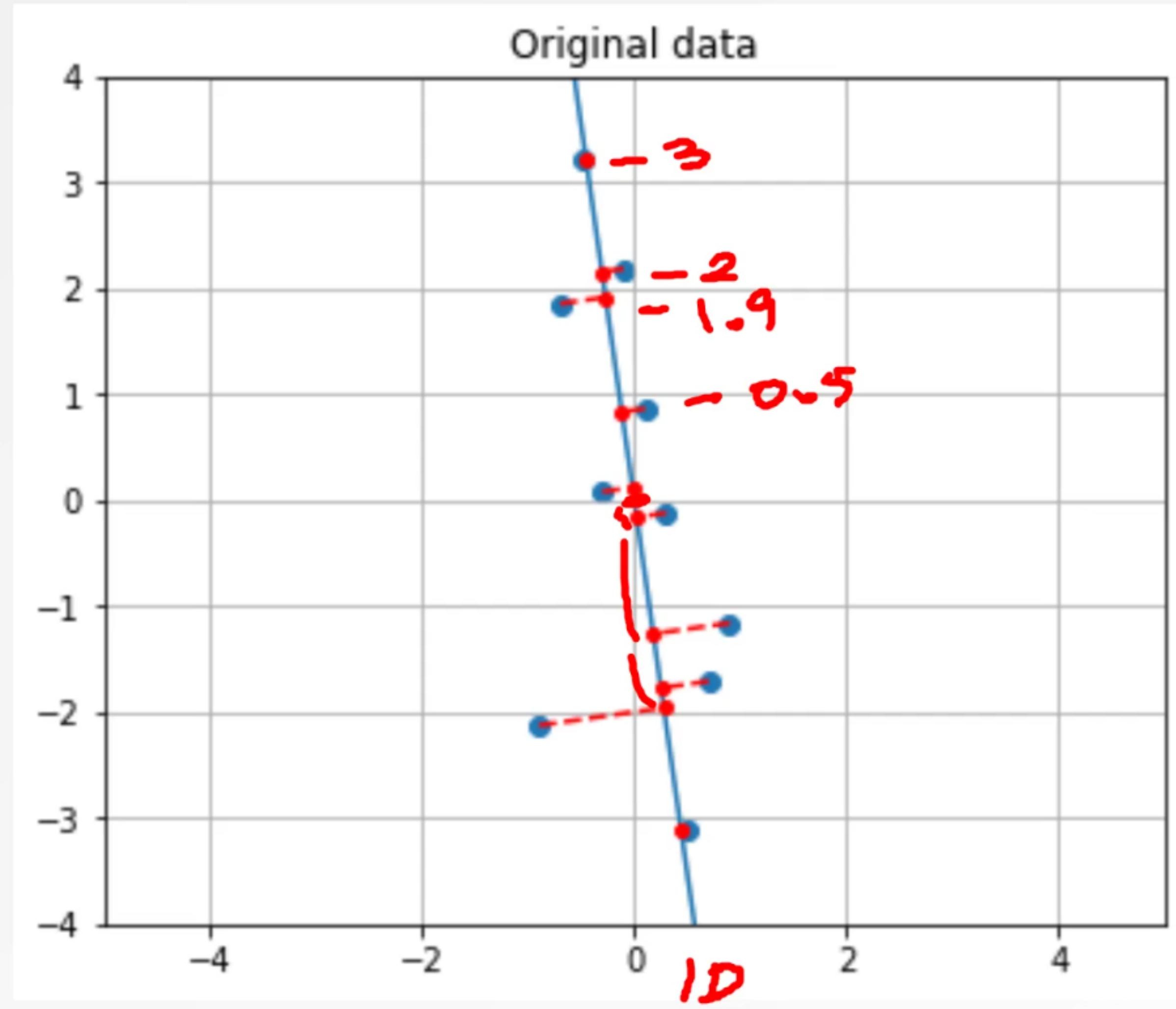
x_1 x_2
—

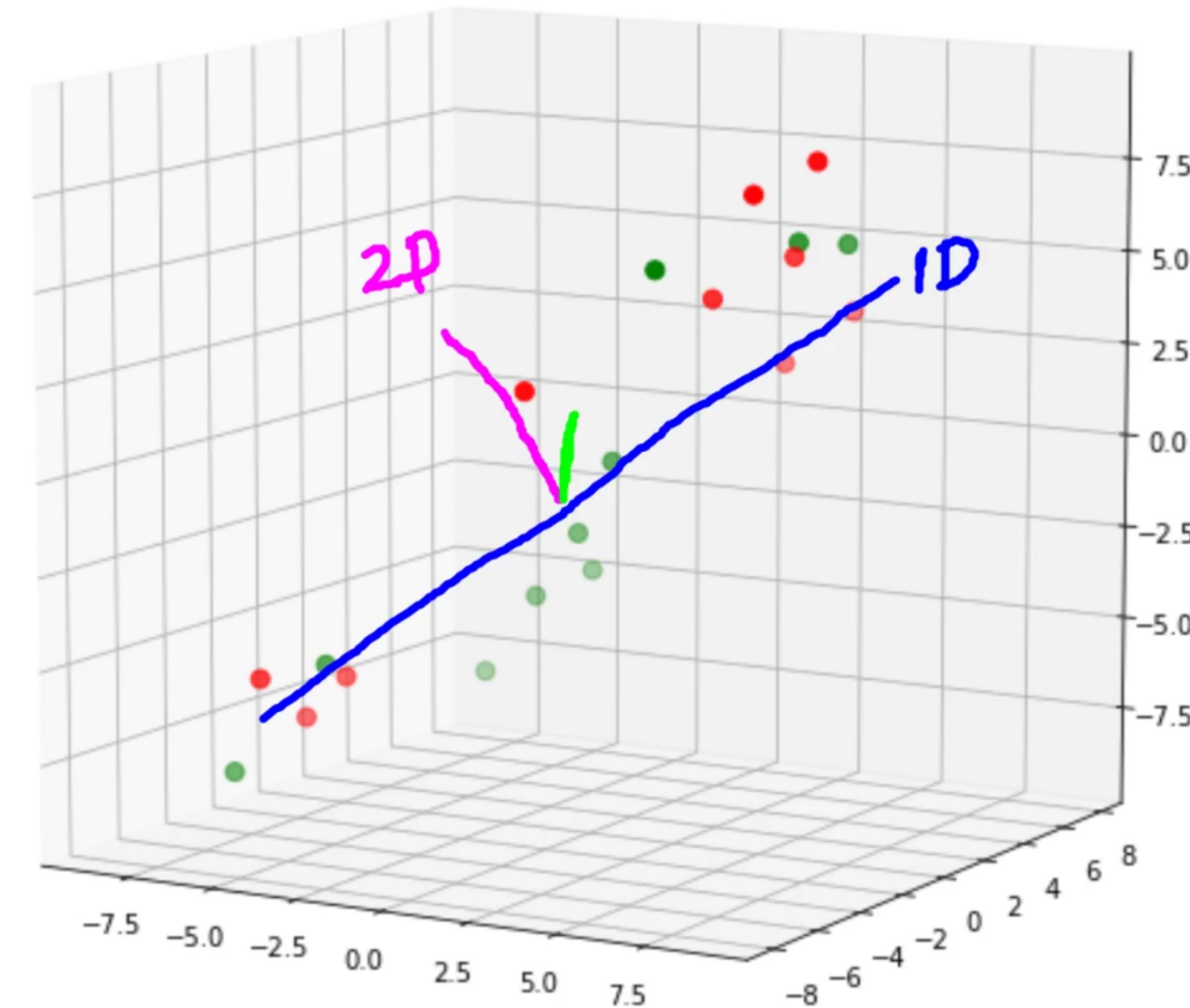
- Methods
 - Linear
 - Principal Component Analysis (PCA) ←
 - Linear Discriminant Analysis (LDA) ←
 - Canonical correlation analysis
 - Multi-dimensional scaling
 - Non-linear
 - Manifold learning (e.g. SOM, autoencoders etc.)

Principal Component Analysis

$$\text{var}(x) = \mathbb{E}[(x-\bar{x})^2]$$







$$x - \bar{x}$$

- Calculating PCA

1. Centre the data

- Compute covariance matrix Σ
- Find the eigenvectors and eigenvalues of Σ
 1. The eigenvectors become the principal components
 2. The eigenvalues provide the explained variance
 3. Select new dimensions and project the data

$$\text{cov}(\mathbf{x}_1, \mathbf{x}_2) = \sum_{i=1}^N \frac{(x_{1i} - \bar{x}_1)(x_{2i} - \bar{x}_2)}{n-1}$$

$$\Sigma = \begin{pmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{pmatrix}$$

$$\Sigma = \frac{1}{n-1} \mathbf{x}^T \mathbf{x}$$

i) eigenvalues

$$\det(\Sigma - \lambda I) = 0$$

$$\lambda \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} = \begin{pmatrix} \lambda & 0 \\ 0 & \lambda \end{pmatrix}$$

$$\det \begin{pmatrix} s_{11} - \lambda & s_{12} \\ s_{21} & s_{22} - \lambda \end{pmatrix} = 0$$

2) eigenvectors

$$(s_{11} - \lambda)(s_{22} - \lambda) - s_{21}s_{12} = 0$$

$$\Sigma \underline{e} = \lambda \underline{e}$$

$$\lambda^2 - \dots =$$

$$\bar{\lambda}_1 \quad \bar{\lambda}_2$$

$$\begin{matrix} \lambda_1 \rightarrow e_1 \\ \lambda_2 \rightarrow e_2 \end{matrix}$$

$$\lambda = \frac{\text{tr}(\Sigma) \pm \sqrt{\text{tr}^2(\Sigma) - 4 \det(\Sigma)}}{2}$$

$$x - \bar{x}$$

- Calculating PCA

1. Centre the data

- Compute covariance matrix Σ
 - Find the eigenvectors and eigenvalues of Σ
 1. The eigenvectors become the principal components
 2. The eigenvalues provide the explained variance
 3. Select new dimensions and project the data

$$2D \quad \begin{array}{c} x_1 \\ x_2 \\ \vdots \\ x_{31} \end{array} \quad \begin{array}{c} x_1 \\ x_2 \\ \vdots \\ x_{32} \end{array} \quad \begin{array}{c} EI \\ EZ \\ \vdots \\ - \end{array} \quad \begin{array}{c} ID \\ \frac{x_1}{x_{11}} \\ \frac{x_2}{x_{21}} \\ \vdots \\ = \end{array}$$