

CS-C3240 - Machine Learning

Feature Learning

Alexander Jung

24.3.2021

1

Learning Goals

- understand challenges with long raw feature vectors
- basic idea of feature learning
- using feature learning to visualize data and protect privacy
- principal component analysis
- random projections

What Are Three Components of ML ?



24.3.2021

A Recent Movie Night

- Partner A: “What movie should we watch?”
- Partner B: “Let’s watch “The Campaign”!!!”
- Partner A: “No F*** Way!”
- Partner B: “Let’s explore using ML ..”



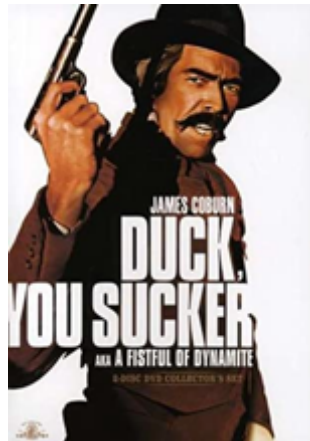
Data Point = “Some Movie”

several **Gigabytes** of raw
feature bits!



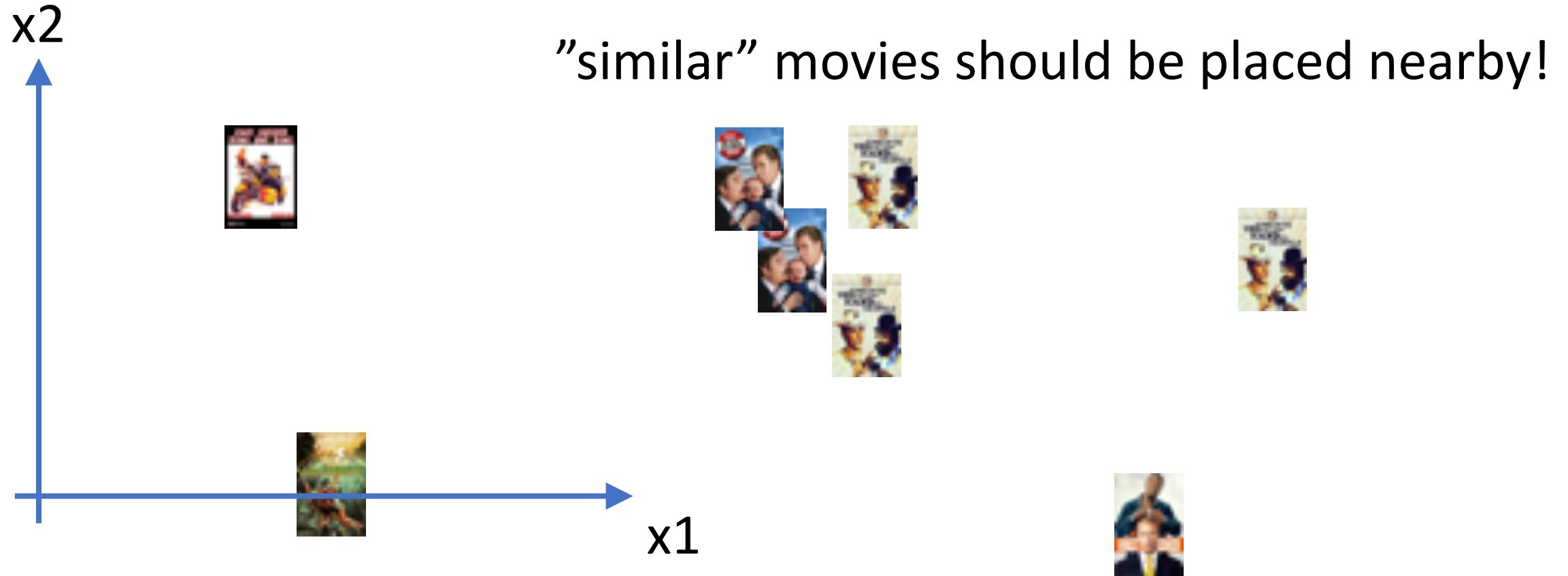
24.3.2021

Dataset = "Bunch of Movies"



24.3.2021

Scatterplot of Movies?



movie represented by two features x_1, x_2 !

Curse of Big Data

Overwhelmed by Tons of Features!

- consider data point representing a person
- digital footprint can be used for constructing features
- health-records (including genetic fingerprint)
- credit-card transactions
- social media posts
- media collections
- travelling profile over last 20 years
-



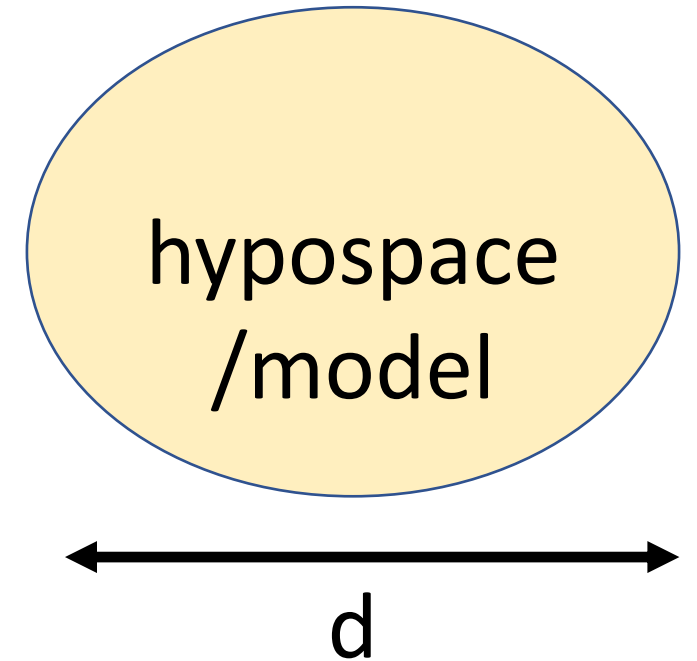
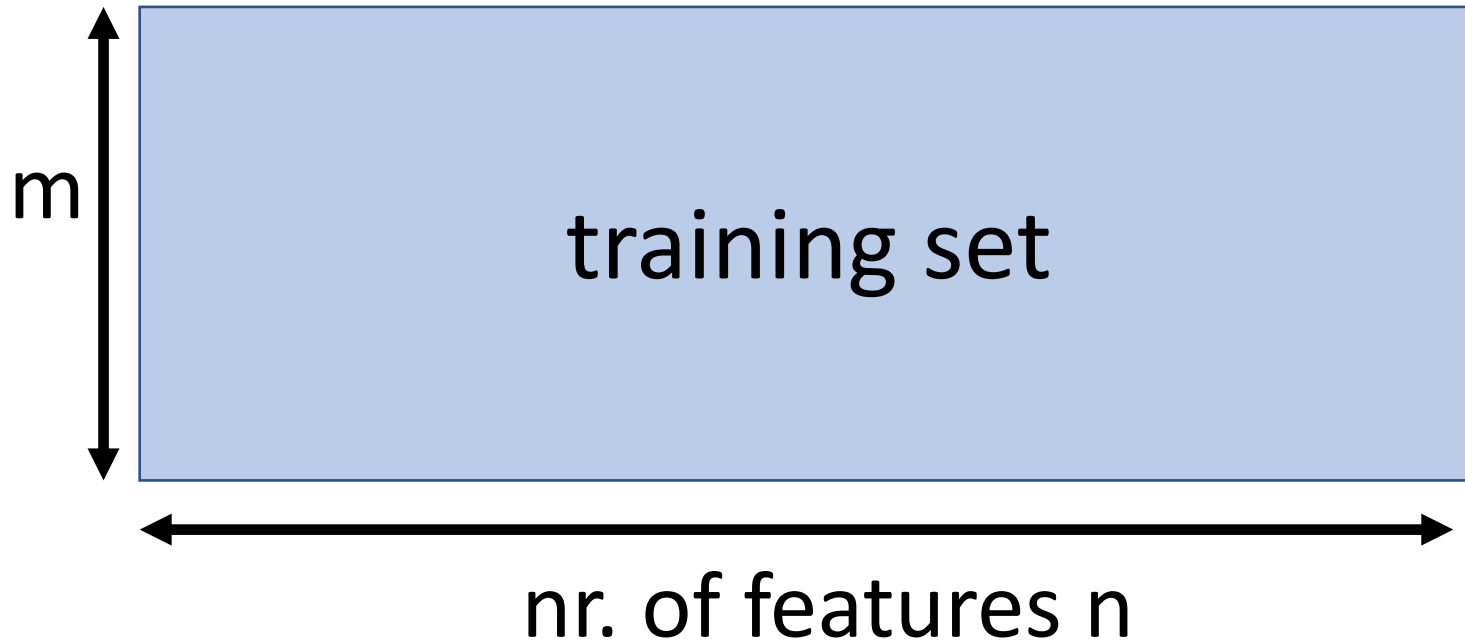
Statistical Challenge

- **effective dimension d** of hypothesis space
- d typically increases with larger nr. of features
- **overfitting** is likely when $d > \text{sample size } m$

Linear Regression

- consider data points with n features
- have m labeled data points
- for $m < n$, plain linear regression will overfit
- n might be billions (e.g., HD-movies)
- would need billions of labeled data points

Data and Model Size



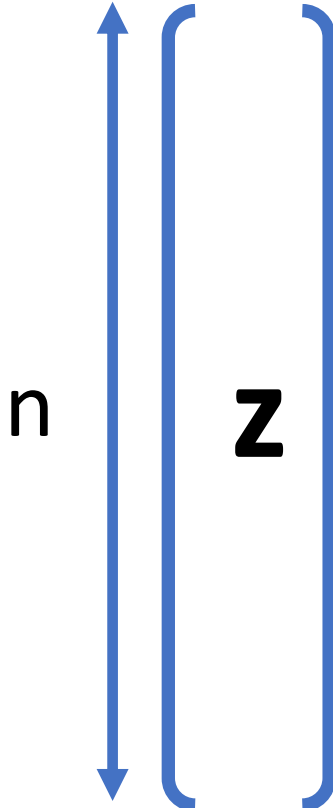
overfitting as soon as
 $d/m > 1$!

Computational Challenge

nr. of arithmetic operations required to find a
good hypothesis (by minimizing training error)
scales with number of features !

Basic Idea of Feature Learning

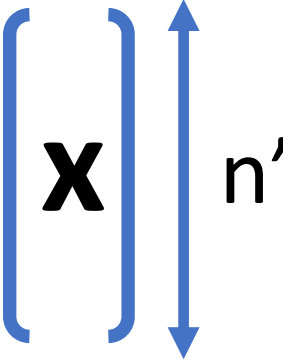
raw feature vector



tunable map W



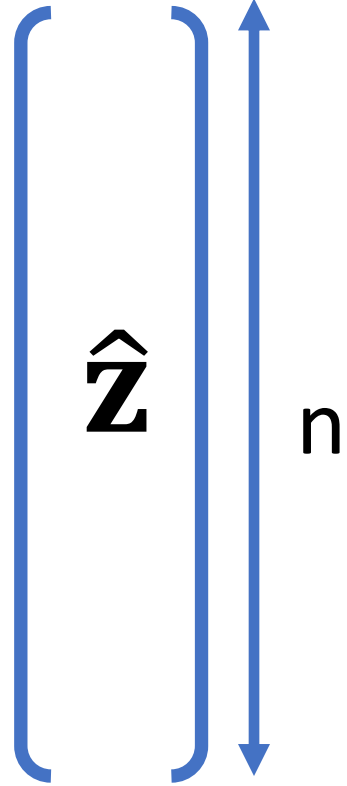
short feature vector



tunable map R



reconstruction



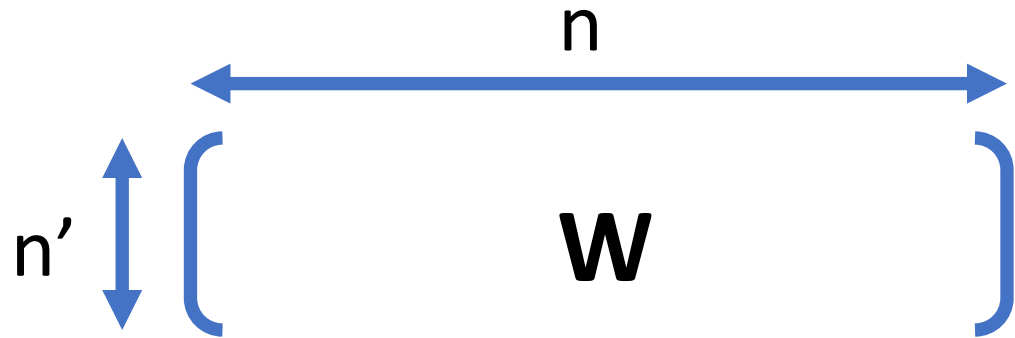
choose (learn) W and R to minimize reconstruction error

$$\mathbf{z} - \hat{\mathbf{z}}$$

Linear Feature Learning

- use linear maps for compression and reconstruction

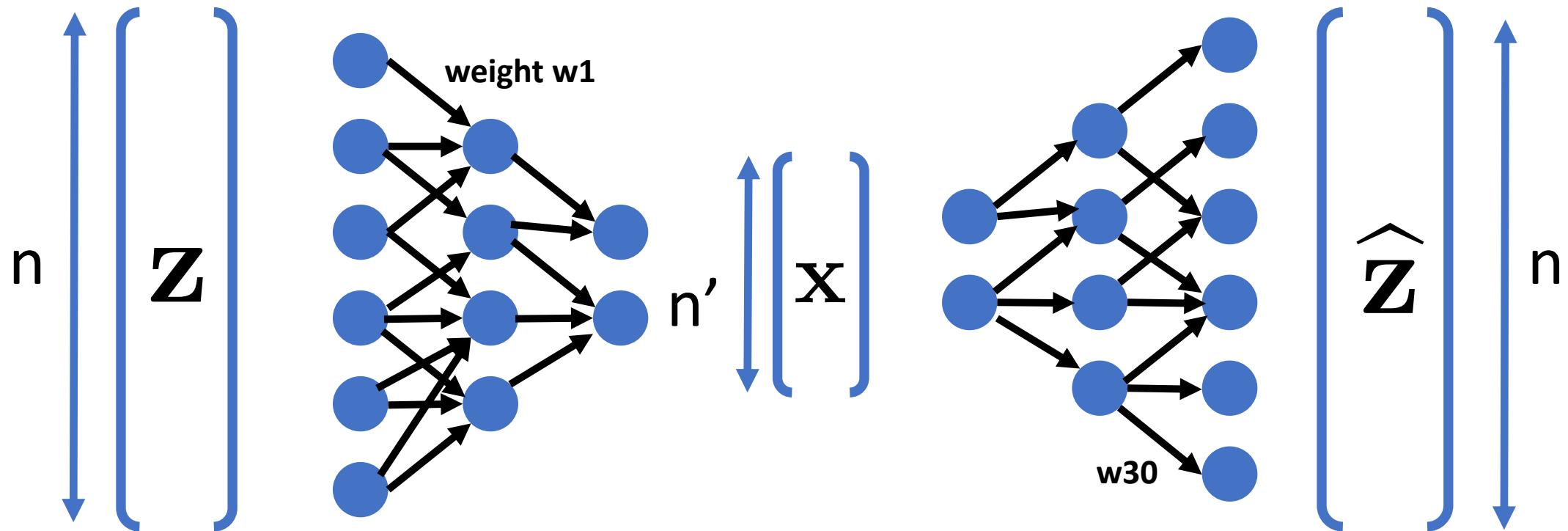
$$\mathbf{x} = \mathbf{W}\mathbf{z} \quad \hat{\mathbf{z}} = \mathbf{R}\mathbf{x}$$



choose matrices \mathbf{W} and \mathbf{R} to minimize $\mathbf{z} - \hat{\mathbf{z}} = (\mathbf{I} - \mathbf{R}\mathbf{W})\mathbf{z}$

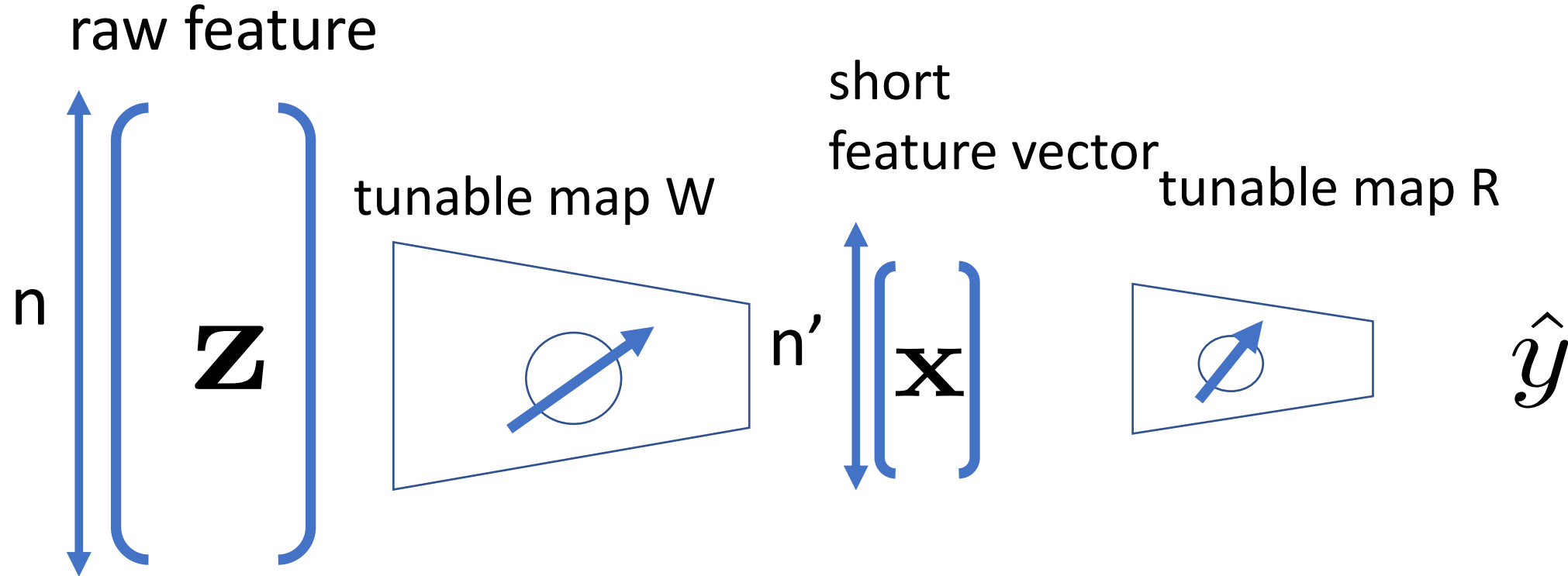
Non-Linear Feature Learning (“Autoencoder”)

use artificial neural networks for compression and reconstruction



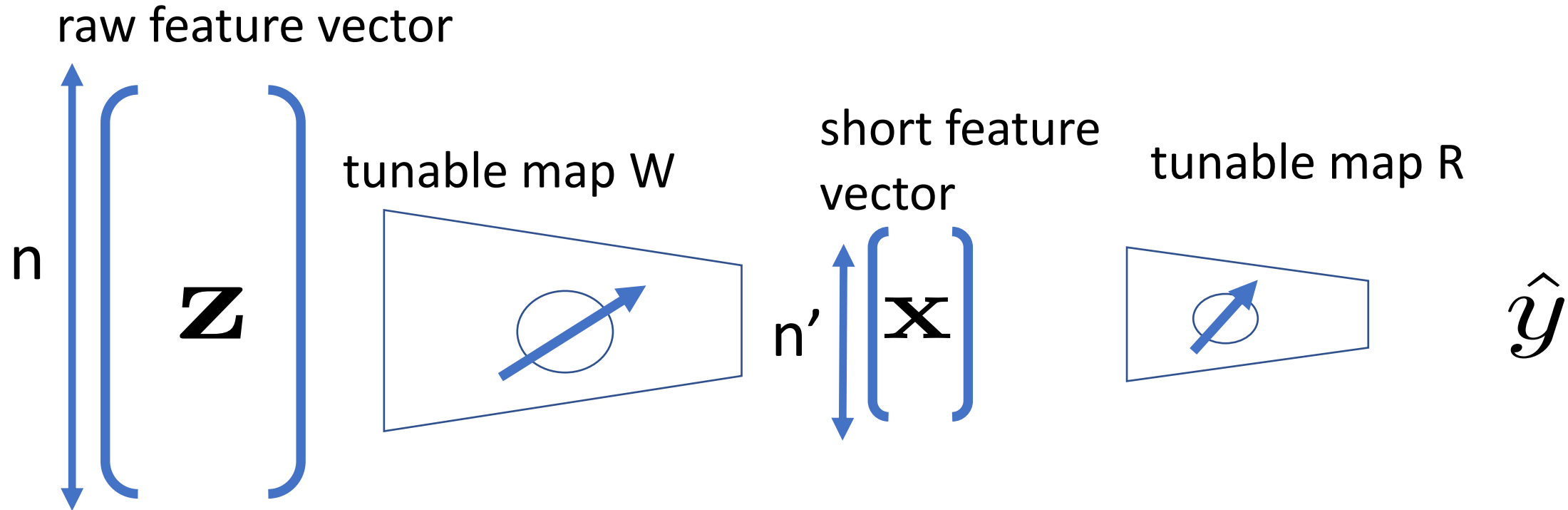
much like linear maps, ANNs are just **parametrized maps!**

Feature Learning for Labeled Data



choose W such that we can **predict** (using some map R) **the label y** from z with **maximum accuracy**

Feature Learning for Labeled Data

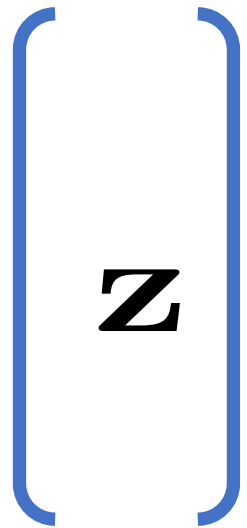


choice for W needs to **balance between**

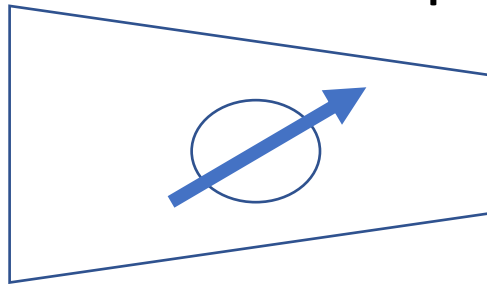
- **compressing** raw feature vector as much as possible
- **keep parts** of raw features that are **relevant** for predicting y

Privacy-Preserving Feature Learning

raw feature vector



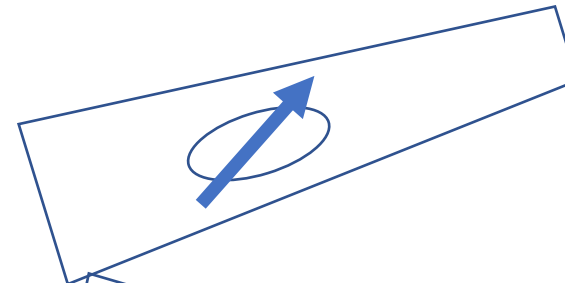
tunable map W



short feature vector



tunable map



\hat{y}

tunable map

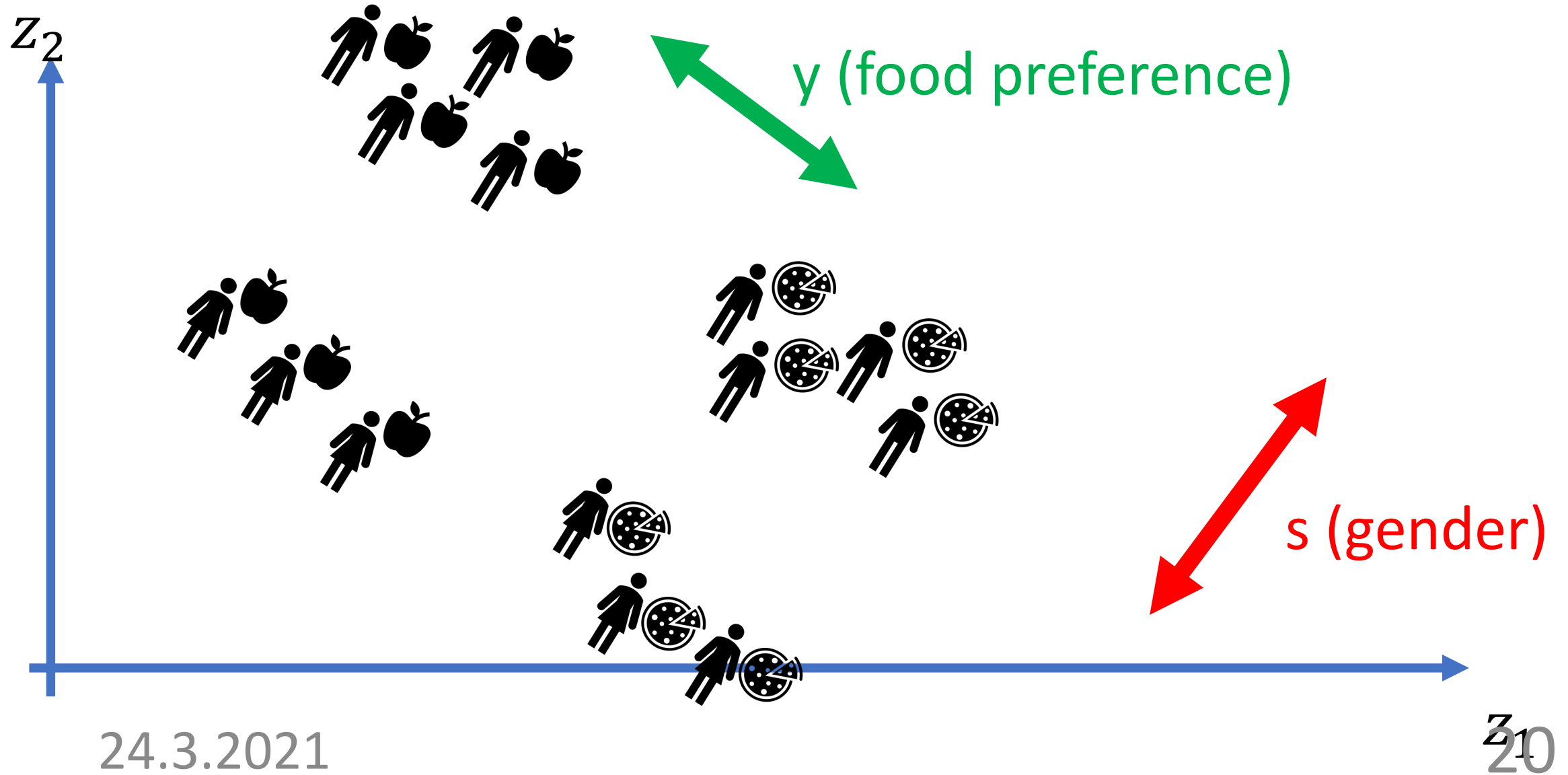


\hat{s}

choice for W needs to balance between

- **compressing** raw feature vector as much as possible
- keep parts of raw features that are **relevant** for predicting y
- predicting **private variable** "s" is **not predictable** from x

Privacy-Preserving Feature Learning



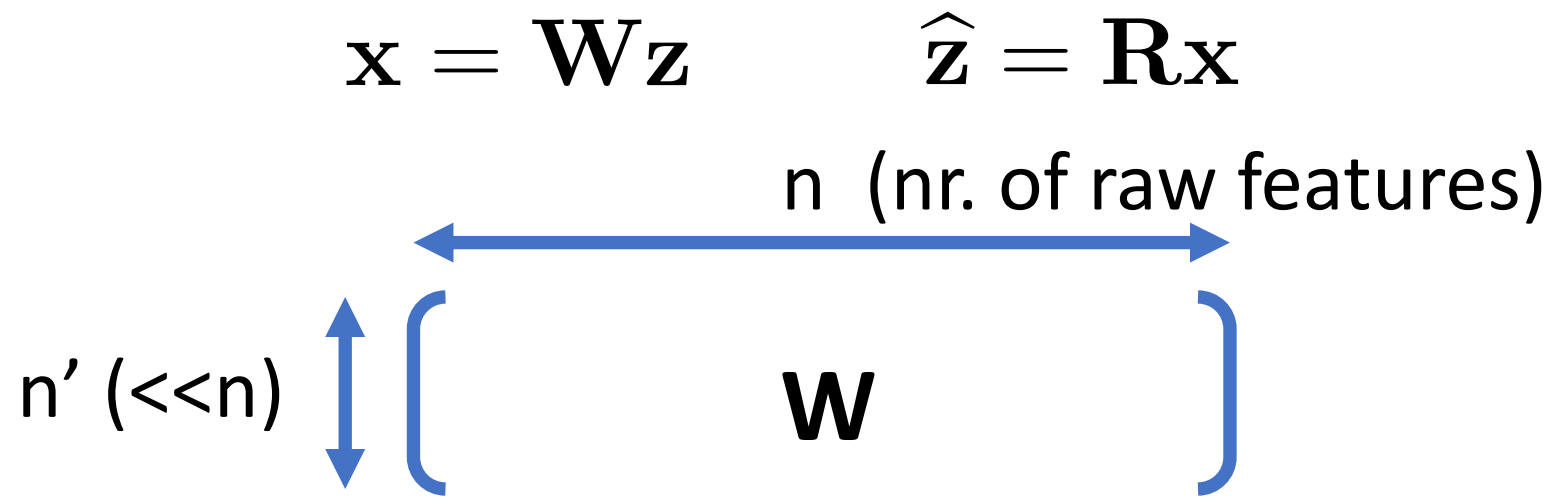
24.3.2021

z_1

Principal Component Analysis (PCA)

Linear Feature Learning

- use linear maps for compression and reconstruction

$$\mathbf{x} = \mathbf{W}\mathbf{z} \quad \hat{\mathbf{z}} = \mathbf{R}\mathbf{x}$$


The diagram shows a matrix \mathbf{W} enclosed in large blue square brackets. To the left of the matrix, a vertical blue double-headed arrow indicates the height, labeled $n' (\ll n)$. Above the matrix, a horizontal blue double-headed arrow indicates the width, labeled n (nr. of raw features).

choose matrices \mathbf{W} and \mathbf{R} to minimize $\mathbf{z} - \hat{\mathbf{z}} = (\mathbf{I} - \mathbf{R}\mathbf{W})\mathbf{z}$

PCA as Risk Minimization

- m datapoints with raw feature vecs $\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}$

choose matrices \mathbf{W}, \mathbf{R} to minimize reconstruction error

$$\mathcal{E}(\mathbf{W}, \mathbf{R}) = \frac{1}{m} \sum_{i=1}^m \left\| \mathbf{z}^{(i)} - \mathbf{R}\mathbf{W}\mathbf{z}^{(i)} \right\|^2$$

Principal Component Analysis (PCA)

- optimal **compression matrix** $\mathbf{W} = (\mathbf{u}^{(1)}, \dots, \mathbf{u}^{(n')})$

- using “top” eigenvectors $\mathbf{u}^{(i)}$ of **sample covariance matrix**

$$\hat{\mathbf{C}} = (1/m) \sum_{i=1}^m \mathbf{z}^{(i)} (\mathbf{z}^{(i)})^T$$

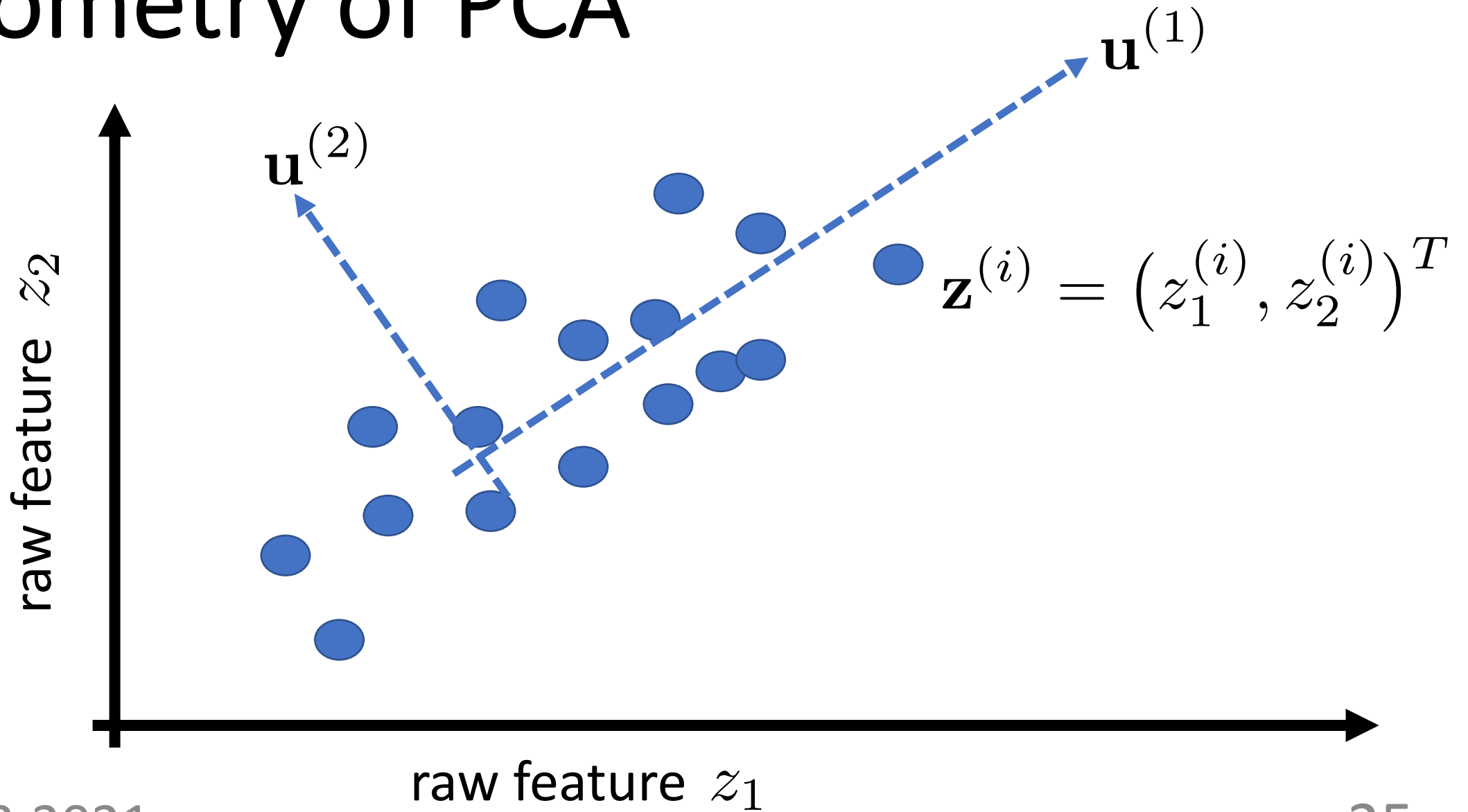
- **eigenvalue decomposition** of psd sample cov. matrix:

$$\hat{\mathbf{C}} = (\mathbf{u}^{(1)}, \dots, \mathbf{u}^{(n)}) \text{diag}(\lambda_1, \dots, \lambda_n) (\mathbf{u}^{(1)}, \dots, \mathbf{u}^{(n)})^T$$

non-negative eigenvalues

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0.$$

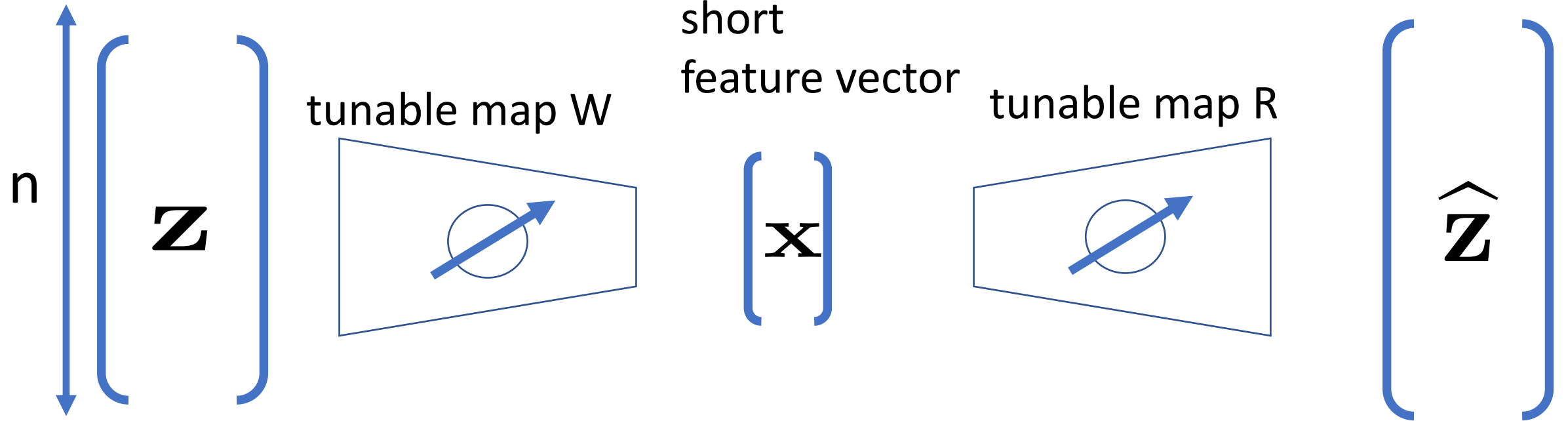
Geometry of PCA



Computational Complexity

raw feature

reconstruction



choose (learn) W and R to **minimize reconstruction error** $\mathbf{z} - \hat{\mathbf{z}}$

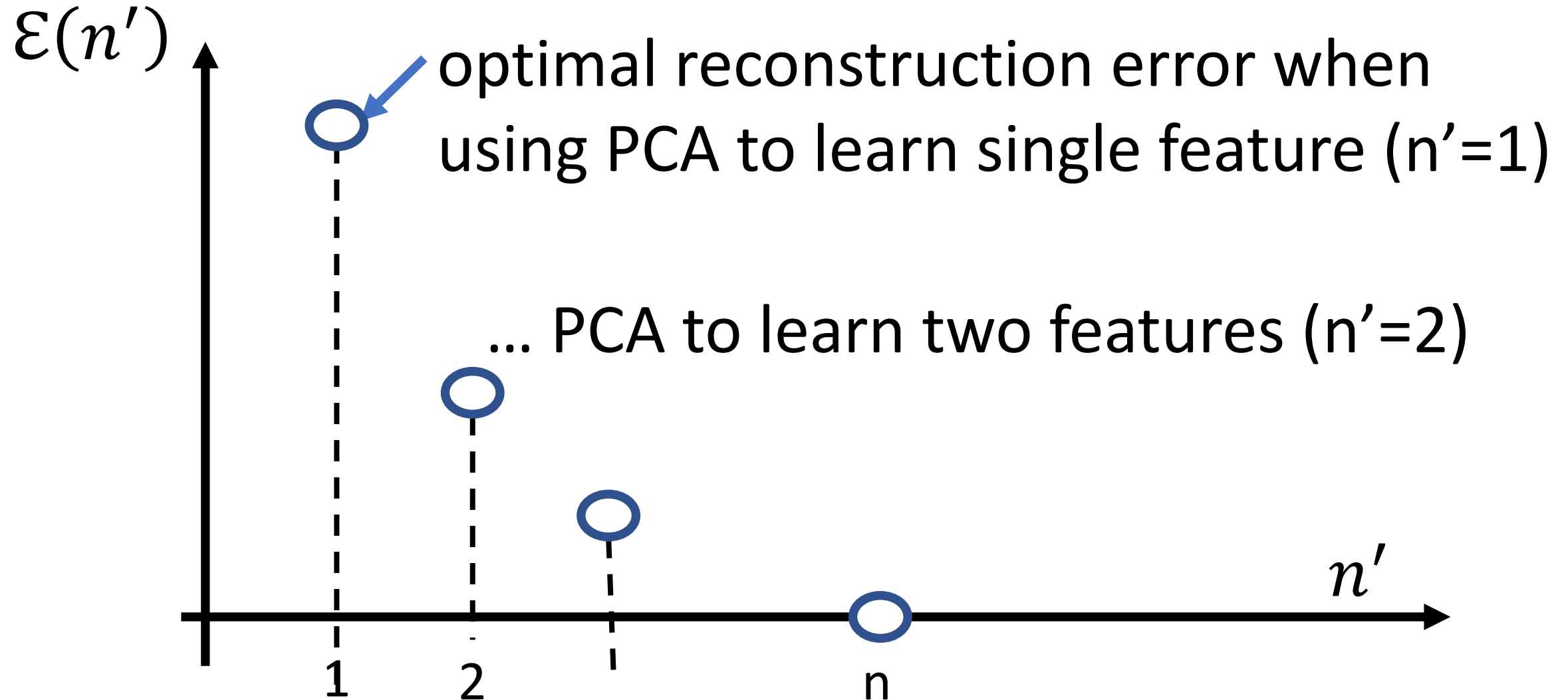
PCA requires eigenvalue decomposition of “ $n \times n$ ” matrix!

How to choose n' ?

- PCA requires number n' of learned features as input
- $n'=2$ for visualization (scatter plot)
- chose n' to balance compression with optimal reconstruction error

$$\mathcal{E}(n') = \min_{\substack{\mathbf{W} \in \mathbb{R}^{n' \times n} \\ \mathbf{R} \in \mathbb{R}^{n \times n'}}} \mathcal{E}(\mathbf{W}, \mathbf{R})$$

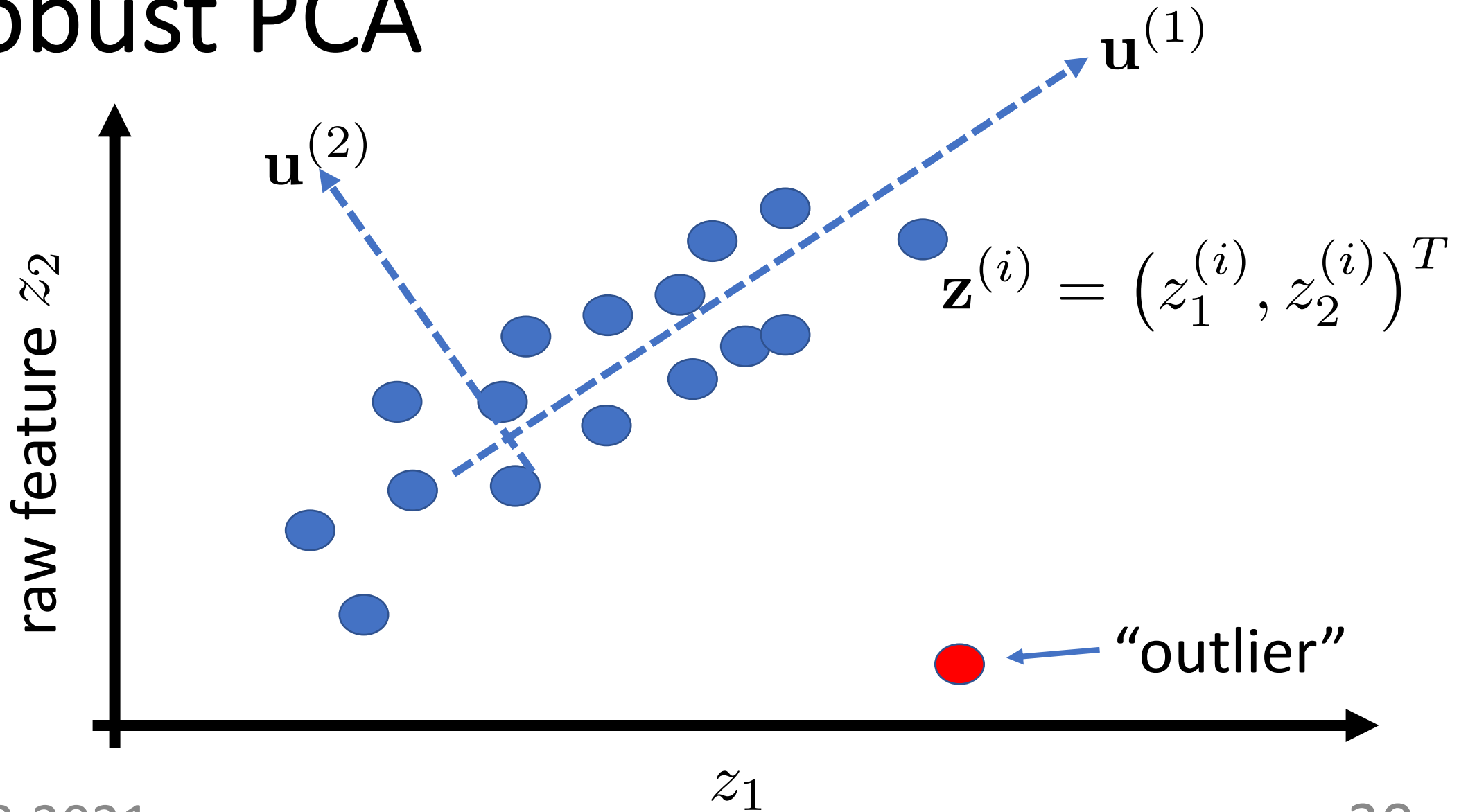
Elbow Method



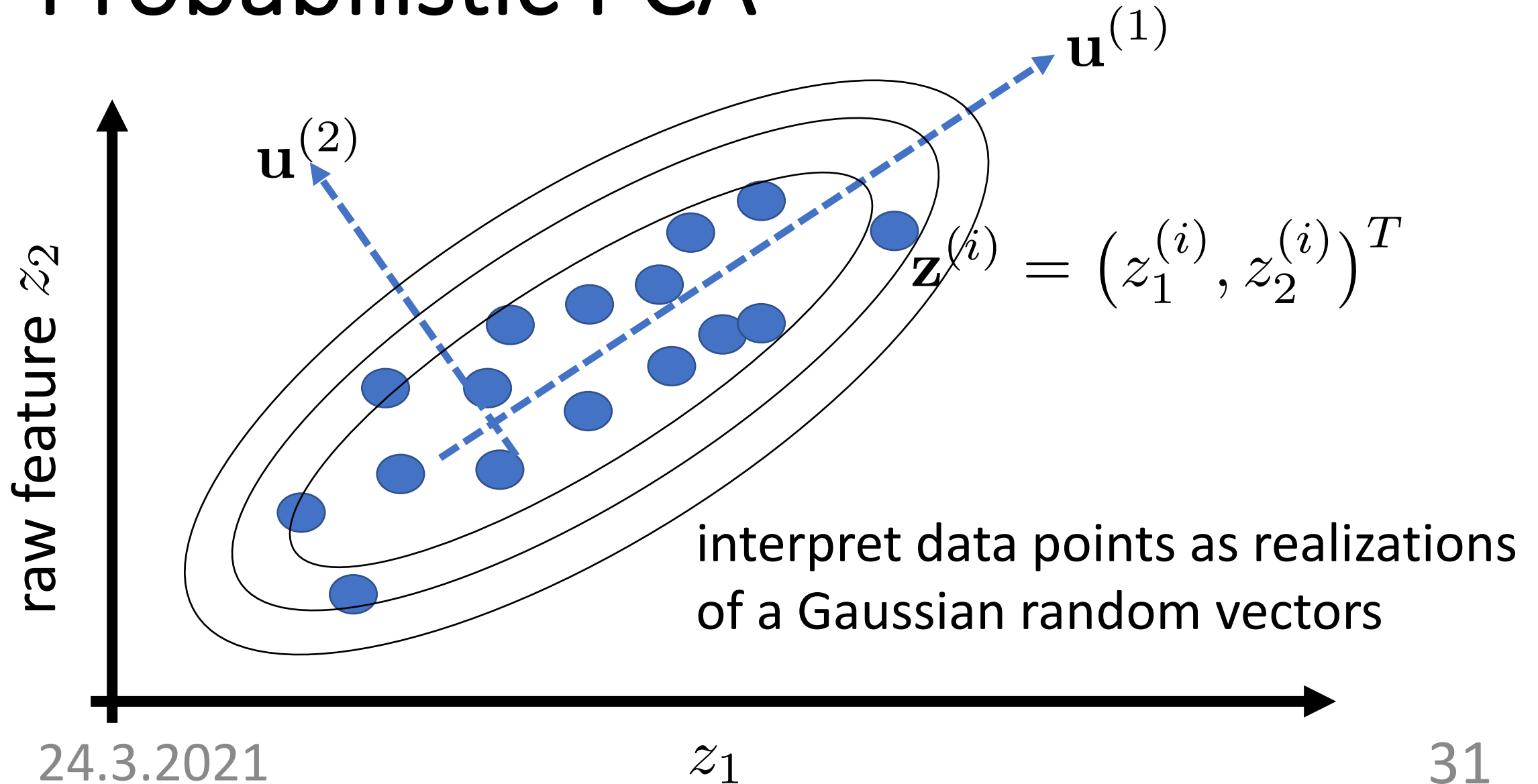
Variations of PCA

- **robust PCA:** uses a different measure of reconstruction error
- **probabilistic PCA:** uses a statistical model for data points
- **sparse PCA:** new features depend on few raw features

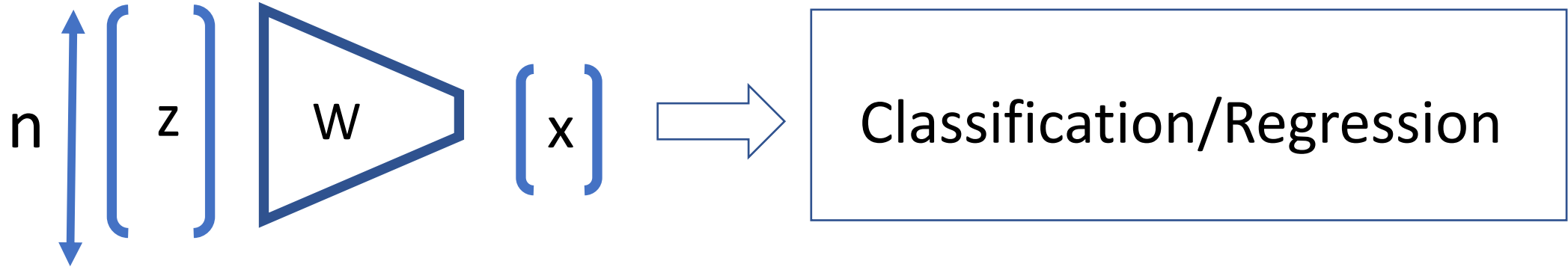
Robust PCA



Probabilistic PCA

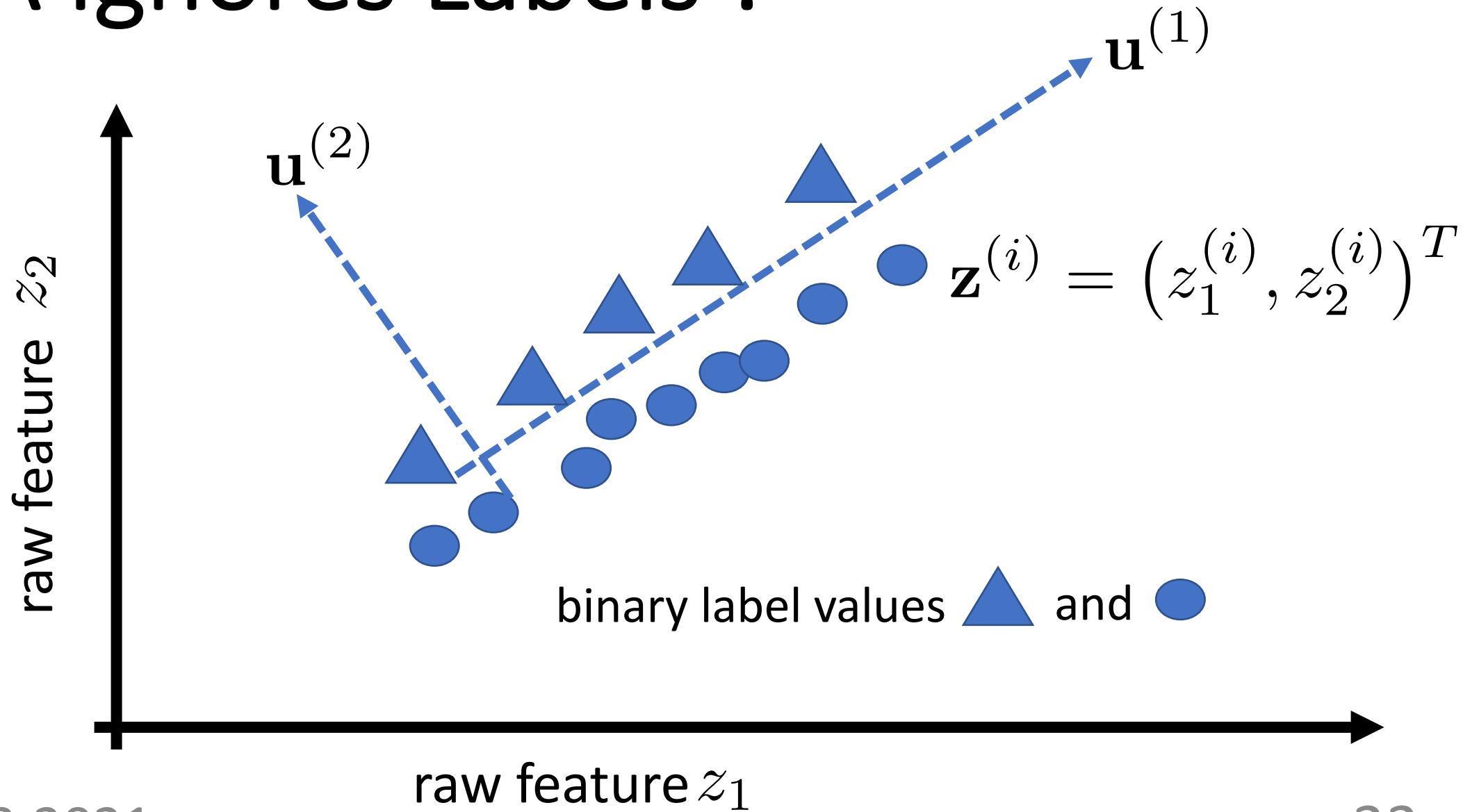


PCA as Pre-Processing

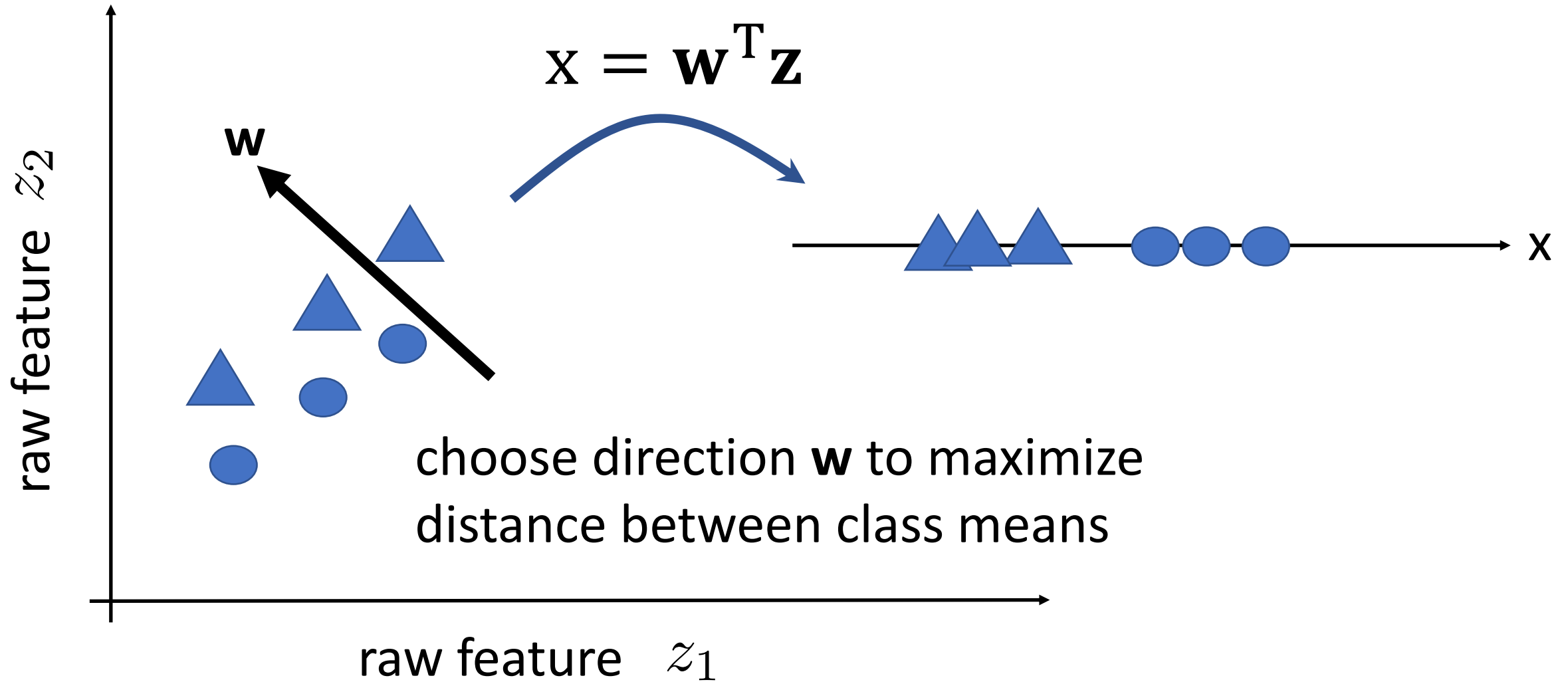


- PCA delivers a compression matrix W
- replace (long) raw features z with shorter features $x = Wz$
- apply regression/classification methods to new features x
- CAUTION: PCA ignores label information!

PCA Ignores Labels !



Fisher's Linear Discriminant



Random Projections

- consider **random projection** $\mathbf{x} = \mathbf{W} \mathbf{z}$
- entries of matrix \mathbf{W} are **randomly chosen**
- **no learning/tuning** of \mathbf{W} required!
- in many settings, **works surprisingly well**
- known as **compressed sensing**

So What?

- feature learning methods determine **relevant features**
- learning **two features** allows to **scatter plot** !
- optimal **linear** feature learning = **PCA**
- PCA ignores label information!
- **random projections** as computationally light alternative

Questions?