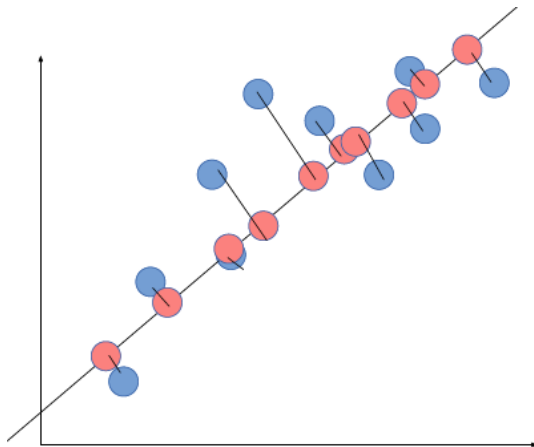
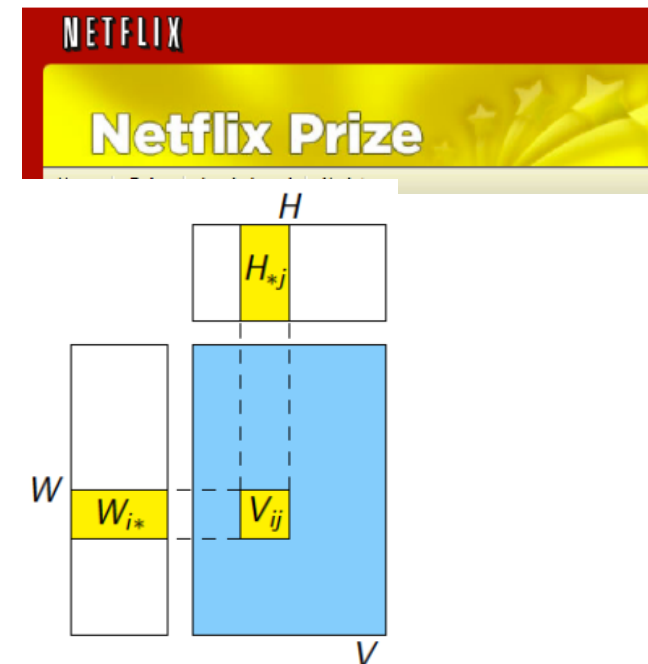


# Principal Components Analysis (PCA)



Many ideas/slides attributable to:  
Liping Liu (Tufts), Emily Fox (UW)  
Matt Gormley (CMU)

Prof. Mike Hughes

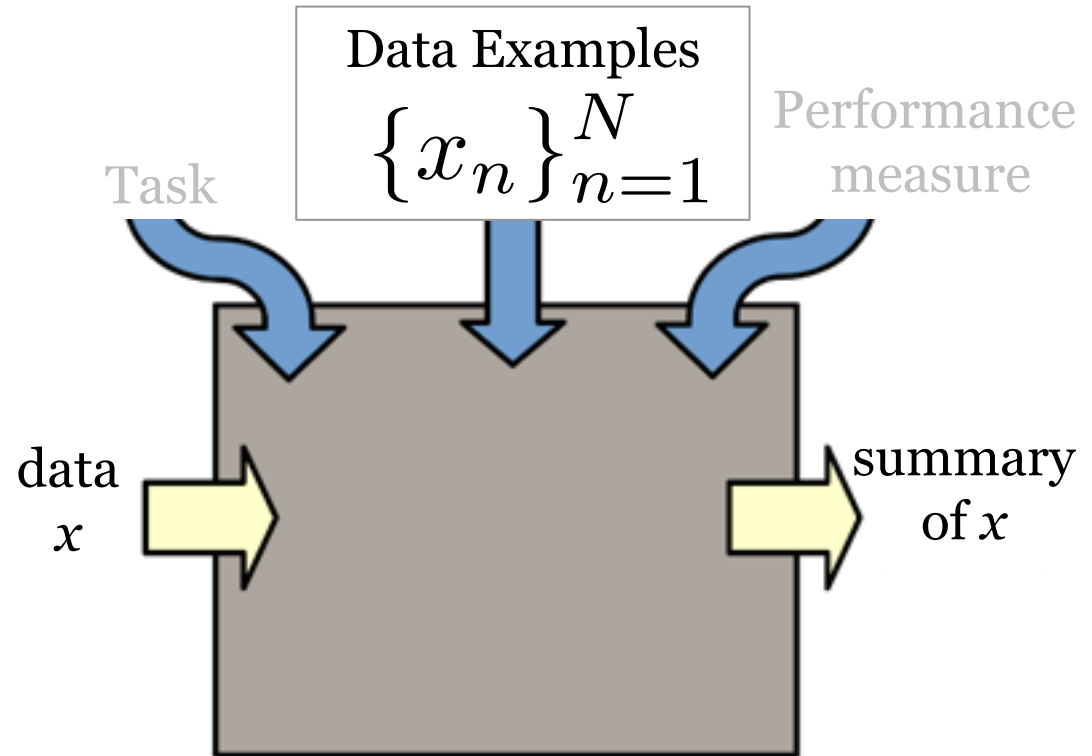


# What will we learn?

Supervised  
Learning

Unsupervised  
Learning

Reinforcement  
Learning



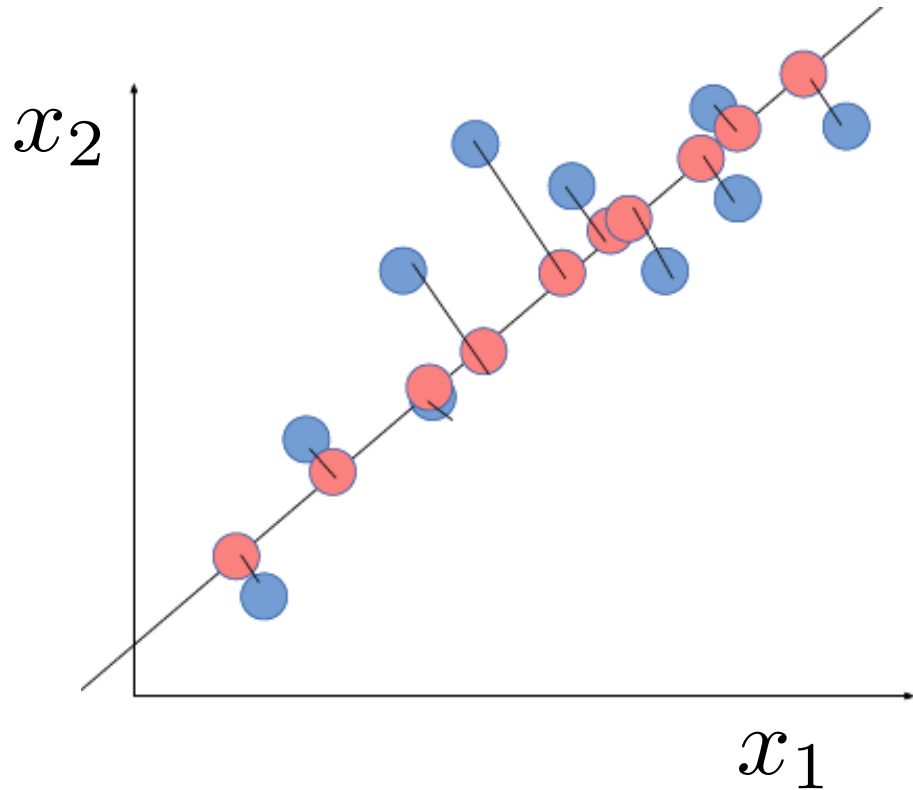
# Task: Embedding

Supervised  
Learning

Unsupervised  
Learning

**embedding**

Reinforcement  
Learning



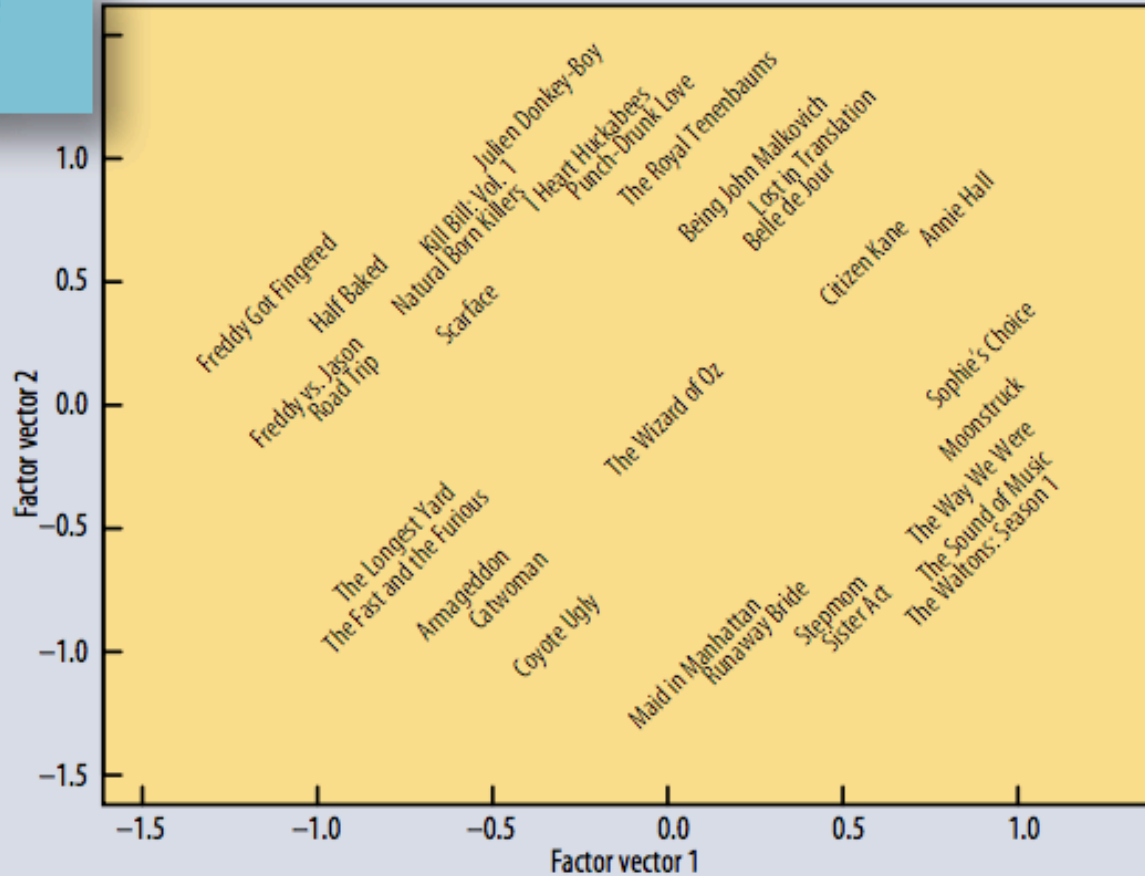
# Dim. Reduction/Embedding

## Unit Objectives

- Goals of dimensionality reduction
  - Reduce feature vector size (keep signal, discard noise)
  - “Interpret” features: visualize/explore/understand
- Common approaches
  - Principal Component Analysis (PCA)
  - word2vec and other neural embeddings
- Evaluation Metrics
  - Storage size
  - “Interpretability”
  - Reconstruction error

# Example: 2D viz. of movies

## Example Factors



**Figure 3. The first two vectors from a matrix decomposition of the Netflix Prize data. Selected movies are placed at the appropriate spot based on their factor vectors in two dimensions. The plot reveals distinct genres, including clusters of movies with strong female leads, fraternity humor, and quirky independent films.**

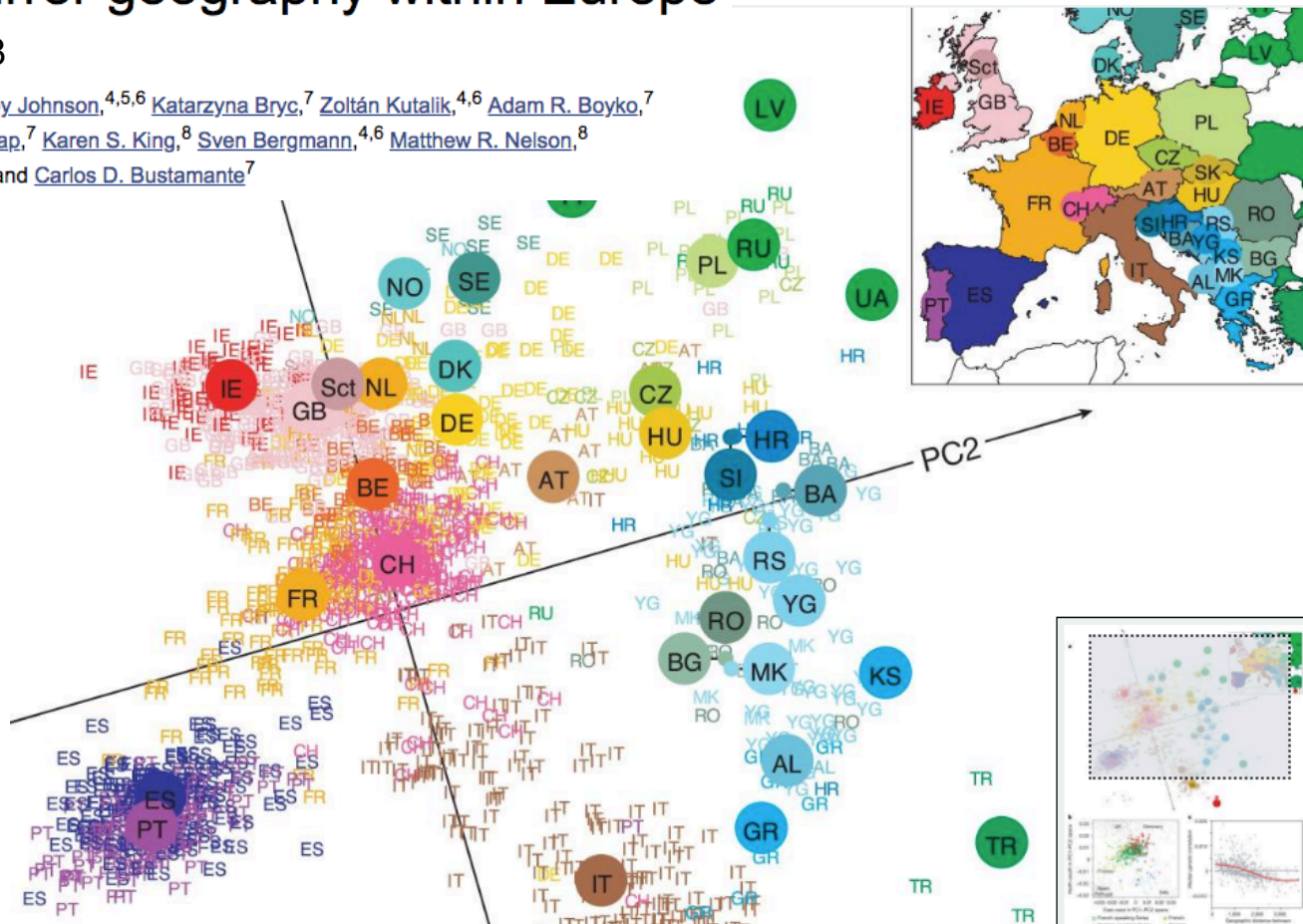
Figure from Koren et al. (2009)

# Example: Genes vs. geography

## Genes mirror geography within Europe

Nature, 2008

[John Novembre](#),<sup>1,2</sup> [Toby Johnson](#),<sup>4,5,6</sup> [Katarzyna Bryc](#),<sup>7</sup> [Zoltán Kutalik](#),<sup>4,6</sup> [Adam R. Boyko](#),<sup>7</sup>  
[Adam Auton](#),<sup>7</sup> [Amit Indap](#),<sup>7</sup> [Karen S. King](#),<sup>8</sup> [Sven Bergmann](#),<sup>4,6</sup> [Matthew R. Nelson](#),<sup>8</sup>  
[Matthew Stephens](#),<sup>2,3</sup> and [Carlos D. Bustamante](#)<sup>7</sup>



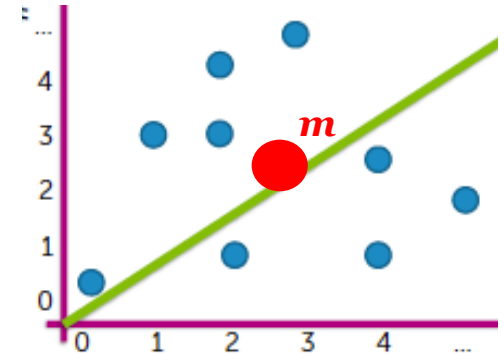
# Centering the Data

Goal: each feature's mean = 0.0

# Constant Reconstruction model

$$\hat{\mathbf{x}}_i = m$$

Parameters:  $m$ , an  $F$ -dim vector



Training problem: Minimize reconstruction error

$$\min_{m \in \mathbb{R}^F} \sum_{n=1}^N (x_n - m)^T (x_n - m)$$

*This is squared error between two vectors*

Optimal parameters:

$$m^* = \text{mean}(x_1, \dots, x_N)$$

*Think of mean vector as optimal “reconstruction” of a dataset if you must use a single vector*



# Mean reconstruction

Ex: Viola Jones data set

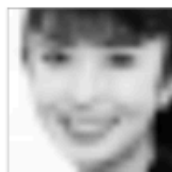
- 24x24 images of faces = 576 dimensional measurements



**Mean**

original

reconstructed



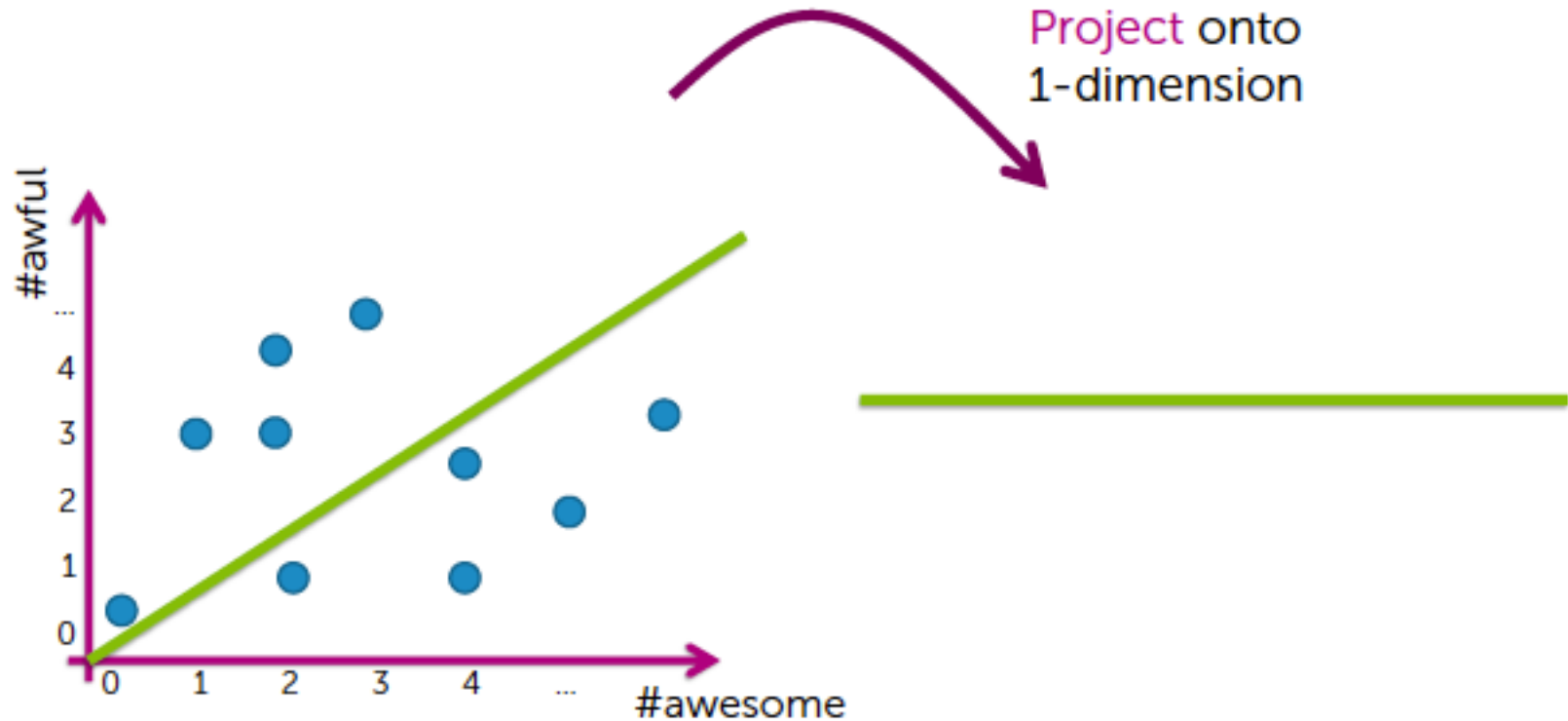
**$X_i$**



**Mean**

# Linear Reconstruction and Principal Component Analysis

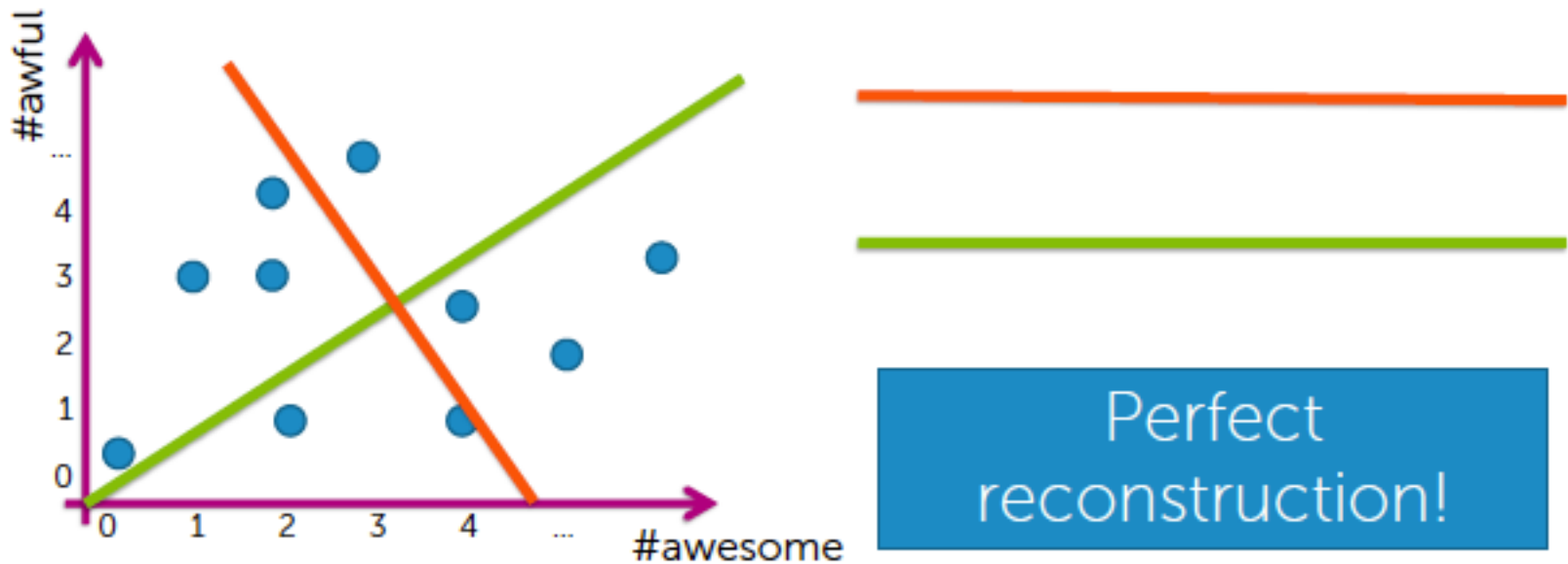
# Linear Projection to 1D



# Reconstruction from 1D to 2D

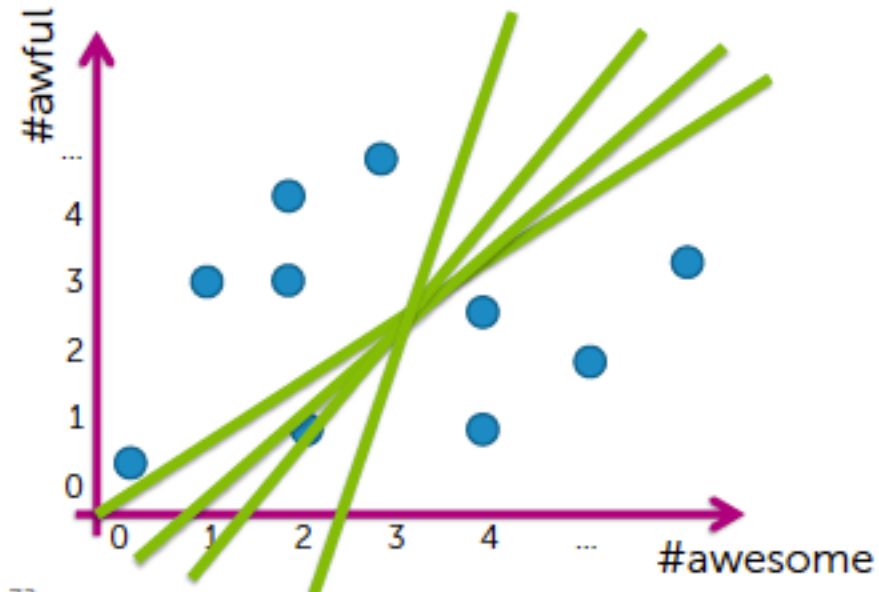


# 2D Orthogonal Basis



*If we could project into 2 dims (same as  $F$ ), we can perfectly reconstruct*

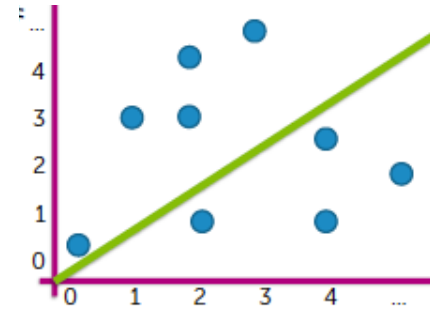
# Which 1D projection is best?



Idea: Minimize reconstruction error

# Linear Reconstruction Model with 1 components

$$\hat{\mathbf{x}}_i = \mathbf{w} z_i + \mathbf{m}$$



F x 1

F x 1

1 x 1

F x 1

High-dim.  
data

Weights

Low-dim  
embedding  
or “score”

“mean”  
vector

# Linear Reconstruction Model with 1 components

$$\hat{\mathbf{x}}_i = \mathbf{w} z_i + \mathbf{m}$$



*W is a vector on unit circle. Magnitude is always 1.*

Problem: “Over-parameterized”. Too many possible solutions!

Suppose we have an alternate model with weights  $w'$  and embedding  $z'$   
We would get equivalent reconstructions if we set:

- $w' = w * 2$
- $z' = z / 2$

Solution: Constrain magnitude of  $w$ .  
 $w$  is a unit vector. We care about direction, not scale.

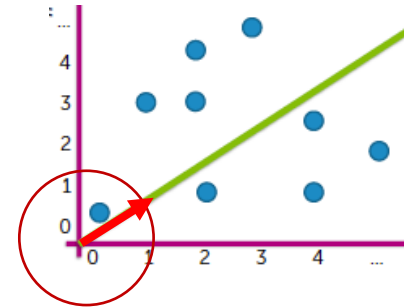
$$\sum_{f=1}^F w_f^2 = 1$$



# Linear Reconstruction Model with 1 components

$$\hat{\mathbf{x}}_i = \mathbf{w} z_i + \mathbf{m}$$

$\text{F} \times 1$                        $\text{F} \times 1$      $1 \times 1$                        $\text{F} \times 1$



*W is a vector on  
unit circle.  
Magnitude is  
always 1.*

Given fixed weights  $w$  and a specific  $x$ , what is the optimal scalar  $z$  value?

Minimize reconstruction error!

$$\min_{z \in \mathbb{R}} (\mathbf{x} - (\mathbf{w}z + \mathbf{m}))^2$$

Exact analytical solution (take gradient, set to zero, solve for  $z$ ) gives:

$$z = w^T (x - m)$$

Projection of feature vector  $x$  onto vector  $w$   
after “centering” (removing the mean)

# Linear Reconstruction Model with K components

$$\hat{\mathbf{x}}_i = \mathbf{W} \mathbf{z}_i + \mathbf{m}$$

F x 1

F x K

K x 1

F x 1

High-dim.  
data

Weights

Low-dim  
vector

Mean of  
data vector

Each of the K weight vectors  $\mathbf{w}_k$  is one “component”.

$$W = \begin{bmatrix} | & | & & | \\ \mathbf{w}_1 & \mathbf{w}_2 & \dots & \mathbf{w}_K \\ | & | & & | \end{bmatrix}$$

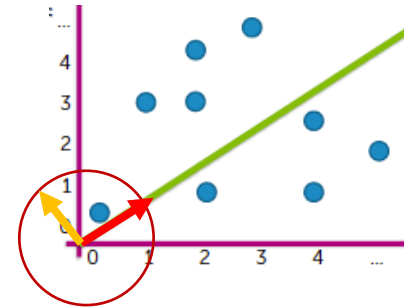
Our goal is to find the K weight vectors that best reconstruct our training dataset

$$\min_{W \in \mathbb{R}^{F \times K}} \sum_{n=1}^N \sum_{f=1}^F (x_{nf} - \hat{x}_{nf}(W))^2$$

Solving this squared error reconstruction objective is known as principal components analysis (PCA)

# Linear Reconstruction Model with K components

$$\hat{\mathbf{x}}_i = \mathbf{W} \mathbf{z}_i + \mathbf{m}$$



F x 1

F x K

K x 1

F x 1

High-dim.  
data

Weights Low-dim  
vector

Mean of  
data vector

We will **require** that:

- (1) All weight vectors are unit vectors
  - This fixes scale and avoid several W with same error
- (2) Component directions are *orthogonal* (*perpendicular*)
  - Avoids information redundancy in W's components

$$W = \begin{bmatrix} | & | & \dots & | \\ \mathbf{w}_1 & \mathbf{w}_2 & \dots & \mathbf{w}_K \\ | & | & & | \end{bmatrix}$$



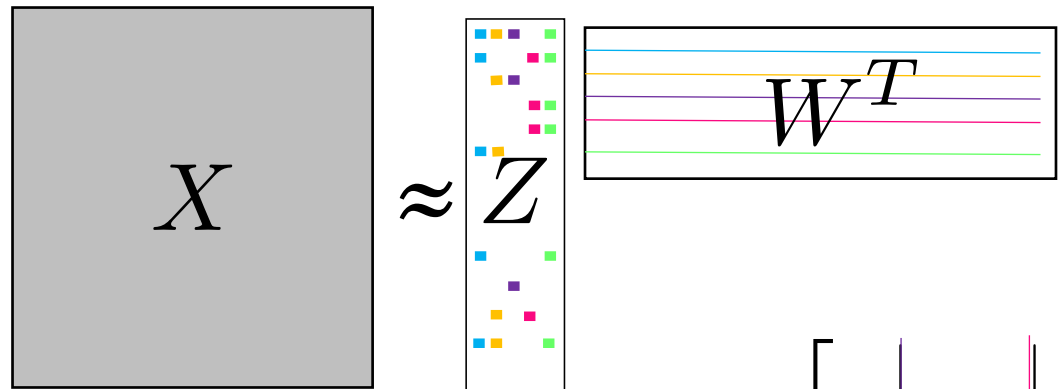
$$\mathbf{w}_k^T \mathbf{w}_k = 1 \rightarrow \sum_{f=1}^F W_{fk}^2 = 1$$



$$\mathbf{w}_j^T \mathbf{w}_k = 0 \rightarrow \sum_{f=1}^F W_{fj} W_{fk} = 0 \quad \forall j \neq k$$

Weights that satisfy (1) and (2) form an “orthonormal basis”

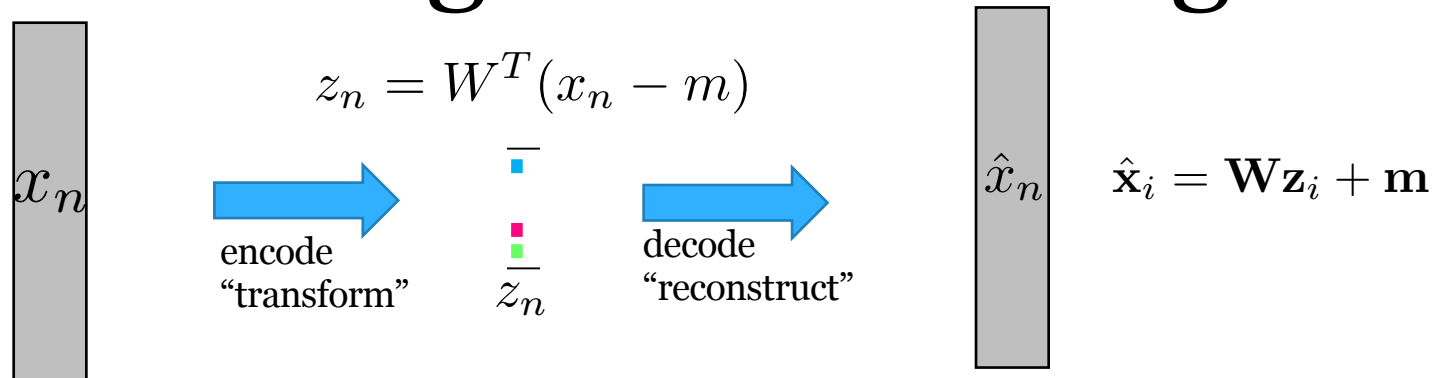
# View: PCA as Matrix Factorization



$$X \approx Z W^T$$

$$W = \begin{bmatrix} \mathbf{w}_1 & \mathbf{w}_2 & \dots & \mathbf{w}_K \end{bmatrix}$$

# View: Encoding and Decoding



# Principal Component Analysis

## Transformation step

What happens when you call `pca.transform(x_QF)`

## Input:

- $X$  : query data,  $Q \times F$ 
  - $Q$  examples of high-dim. feature vectors
- Trained PCA parameters (contained inside `pca`)
  - $m$  : mean vector, size  $F$
  - $W$  : learned basis of weight vectors,  $F \times K$

## Output:

- $Z$  : projections,  $N \times K$ 
  - Each row  $Z[n]$  is a low-dim. “embedding” of  $X[n]$

$$z_n = W^T (x_n - m)$$

# Example: PCA on Faces

Ex: Viola Jones data set

- 24x24 images of faces = 576 dimensional measurements
- Take first K PCA components



Mean



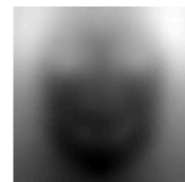
Dir 1



Dir 2



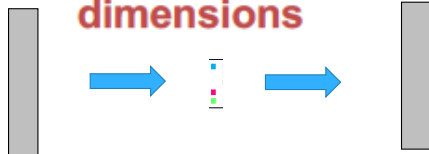
Dir 3



Dir 4

...

Projecting data  
onto first k  
dimensions



original



k=5



k=10



k=50

...



k=F



If we use all possible components, we  
*perfectly reconstruct* original data

# Principal Component Analysis

Training step : What happens when we call `pca.fit(x_NF)`

## ***Input:***

- $X$  : training data,  $N \times F$ 
  - $N$  examples of high-dim. feature vectors
- $K$  : int, number of components
  - Satisfies  $1 \leq K \leq F$

$$\min_{m \in \mathbb{R}^F, W \in \mathbb{R}^{F \times K}} \sum_{n=1}^N \sum_{f=1}^F (x_{nf} - \hat{x}_{nf}(m, W))^2$$


subject to:  $W^T W = I_K$  *Orthonormal constraint*

## ***Output: Trained parameters for PCA***

- $m$  : mean vector, size  $F$
- $W$  : learned basis of weight vectors,  $F \times K$ 
  - One  $F$ -dim. unit vector (magnitude 1) for each component
  - Each of the  $K$  vectors is orthogonal to every other

# Eigenvalues and Eigenvectors

Here is the most important definition in this text.

 **Definition.** Let  $A$  be an  $n \times n$  matrix.

1. An **eigenvector** of  $A$  is a *nonzero* vector  $v$  in  $\mathbf{R}^n$  such that  $Av = \lambda v$ , for some scalar  $\lambda$ .
2. An **eigenvalue** of  $A$  is a scalar  $\lambda$  such that the equation  $Av = \lambda v$  has a *nontrivial* solution.

If  $Av = \lambda v$  for  $v \neq 0$ , we say that  $\lambda$  is the **eigenvalue for**  $v$ , and that  $v$  is an **eigenvector for**  $\lambda$ .

The German prefix “eigen” roughly translates to “self” or “own”. An eigenvector of  $A$  is a vector that is taken to a multiple of itself, which partially explains the terminology.

**Note.** Eigenvalues and eigenvectors are only for square matrices.

Source: <https://textbooks.math.gatech.edu/ila/eigenvectors.html>



# The weight component vectors are the eigenvectors of the covariance matrix of the centered dataset

$$S = \frac{1}{N} \sum_{n=1}^N (x_n - m)(x_n - m)^T$$

Every principal component vector  $w_k$  satisfies this equation:

$$S w_k = \lambda_k w_k \quad W = \left[ \begin{array}{c|c|c|c} | & | & \dots & | \\ \mathbf{w}_1 & \mathbf{w}_2 & & \mathbf{w}_K \\ | & | & & | \end{array} \right]$$

When we fit K principal components to a dataset, the optimal ones (that minimize reconstruction error) are those with the K largest eigenvalues.

Can use standard linalg libraries to compute the eigenvalues/vectors!

# PCA Principles

- Minimize **reconstruction error**
  - Should be able to recreate  $x$  from  $z$
- Equivalent to **maximizing variance**
  - Want reconstructions to retain *maximum* information

# PCA: How to Select K?

- 1) Use downstream supervised task metric
  - Regression error
- 2) Use memory constraints of task
  - Can't store more than 50 dims for 1M examples?  
Take  $K=50$
- 3) Plot cumulative “variance explained”
  - Take K that seems to capture most or all variance

# Empirical Variance of Data X

Assume we've computed the empirical mean vector:

$$m \triangleq \frac{1}{N} \sum_{n=1}^N x_n$$

Empirical variance is defined as averaged squared error from the empirical mean:

$$\begin{aligned} \text{Var}[X] &= \frac{1}{N} \sum_{n=1}^N \sum_{f=1}^F (x_{nf} - m_f)^2 \\ &= \frac{1}{N} \sum_{n=1}^N (x_n - m)^T (x_n - m) \end{aligned}$$

# Empirical Variance of reconstructions

$$= \frac{1}{N} \sum_{n=1}^N x_n^T x_n$$

$$= \frac{1}{N} \sum_{n=1}^N (z_{n1}w_1 + \dots + z_{nK}w_K)^T (z_{n1}w_1 + \dots + z_{nK}w_K)$$

$$= \frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K z_{nk}^2$$

$$= \sum_{k=1}^K \lambda_k$$

Just sum up the top K eigenvalues!

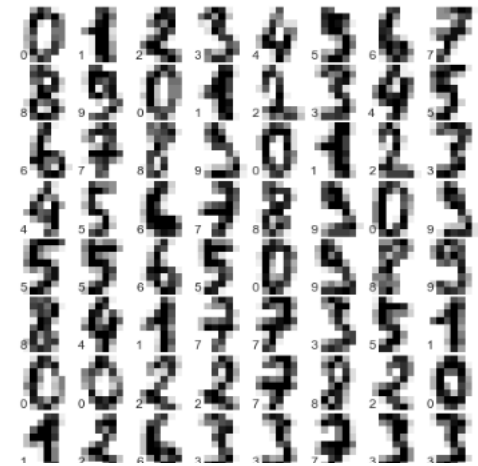
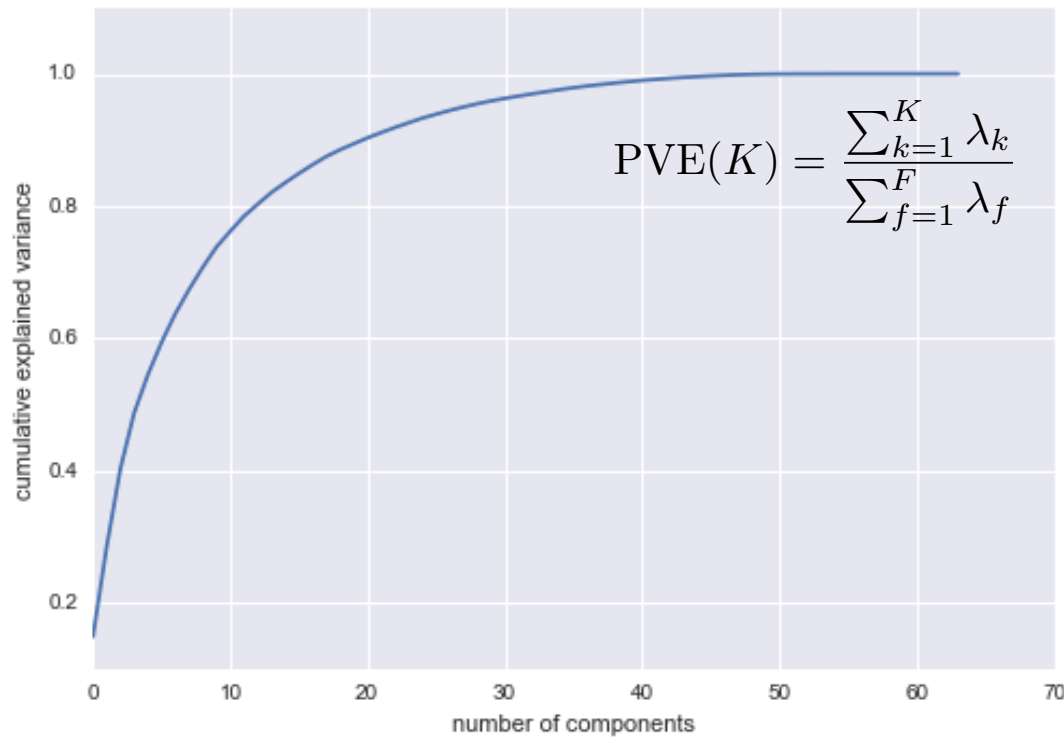
# Proportion of Variance Explained by first K components

$$\text{PVE}(K) = \frac{\sum_{k=1}^K \lambda_k}{\sum_{f=1}^F \lambda_f}$$

Goal: Want K value where proportion of variance explained is large.  
Indicates good reconstruction ability on our training set.

# Variance explained curve

```
: pca = PCA().fit(digits.data)
plt.plot(np.cumsum(pca.explained_variance_ratio_))
plt.xlabel('number of components')
plt.ylabel('cumulative explained variance');
```



```
: from sklearn.datasets import load_digits
digits = load_digits()
digits.data.shape
: (1797, 64)
```

# PCA Summary

## PRO

- Usually, fast to train, fast to test
  - Slowest step: finding  $K$  eigenvectors of an  $F \times F$  matrix
- Nested model
  - PCA with  $K=5$  overlaps with PCA with  $K=4$

## CON

- Sensitive to rescaling of input data features
- Learned basis known only up to  $\pm$  scaling
- Not often best for supervised tasks



# PCA: Best Practices

- If features all have different units
  - Try rescaling to all be within  $(-1, +1)$  or have variance 1
- If features have same units, may not need to do this

# Dim. Reduction/Embedding

## Unit Objectives

- Goals of dimensionality reduction
  - Reduce feature vector size (keep signal, discard noise)
  - “Interpret” features: visualize/explore/understand
- Common approaches
  - Principal Component Analysis (PCA)
  - word2vec and other non-linear embeddings
- Evaluation Metrics
  - Storage size
  - “Interpretability”
  - Reconstruction error