

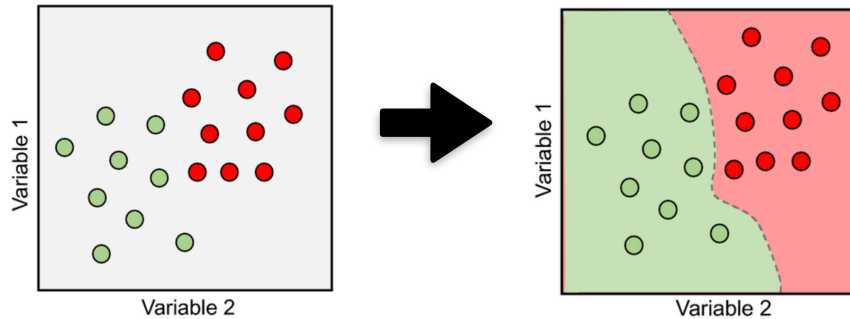
Principal Component Analysis

Natasha Jaques

Unsupervised vs. supervised learning

Previously: supervised learning

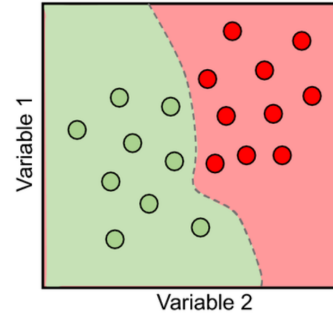
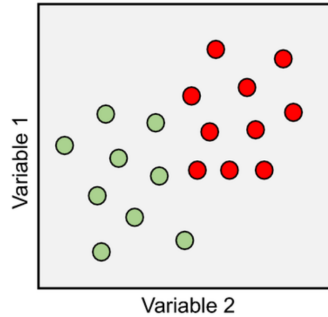
- Each data point x_i has a corresponding label y_i ; $\{x_i, y_i\}_{i=1}^n$ with $x_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}$.
Try to predict the label y for a new test point x



Unsupervised vs. supervised learning

Previously: supervised learning

- Each data point x_i has a corresponding label y_i ; $\{x_i, y_i\}_{i=1}^n$ with $x_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}$.
Try to predict the label y for a new test point x



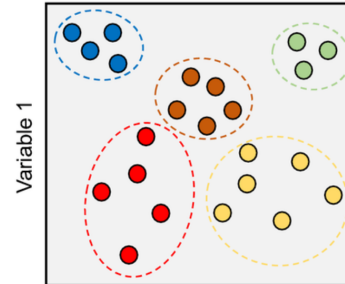
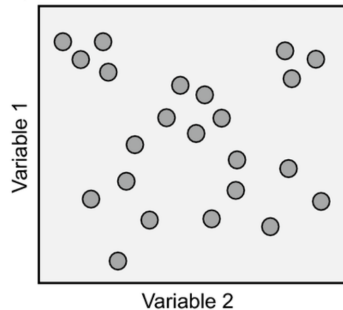
$P(\text{price (house)})$

$P(x_1, \dots, x_d)$
← Features

$P(\text{house})$

Now: Unsupervised learning

- No labels: data $\{x_i\}_{i=1}^n$ with $x_i \in \mathbb{R}^d$. Try to model the data distribution $P(X)$, potentially by finding patterns/clusters, or a low-dimensional representation



$\rho(x)$

Motivation: dimensionality reduction

- It takes $n \times d$ memory to store data $\{x_i\}_{i=1}^n$ with $x_i \in \mathbb{R}^d$
- But many real data have patterns that repeat over samples. Can we find some patterns and use them?



$\rightarrow 1024$
 $d=32 \times 32$ pixels per image
 n images
 $d \times n$ real values to store the data

$\rightarrow z_i$
 $\rightarrow z_j$

Principal component analysis finds a compact linear representation

- patterns that capture the distinct features of the samples is called **principal component** (to be formally defined later)

Principal components:

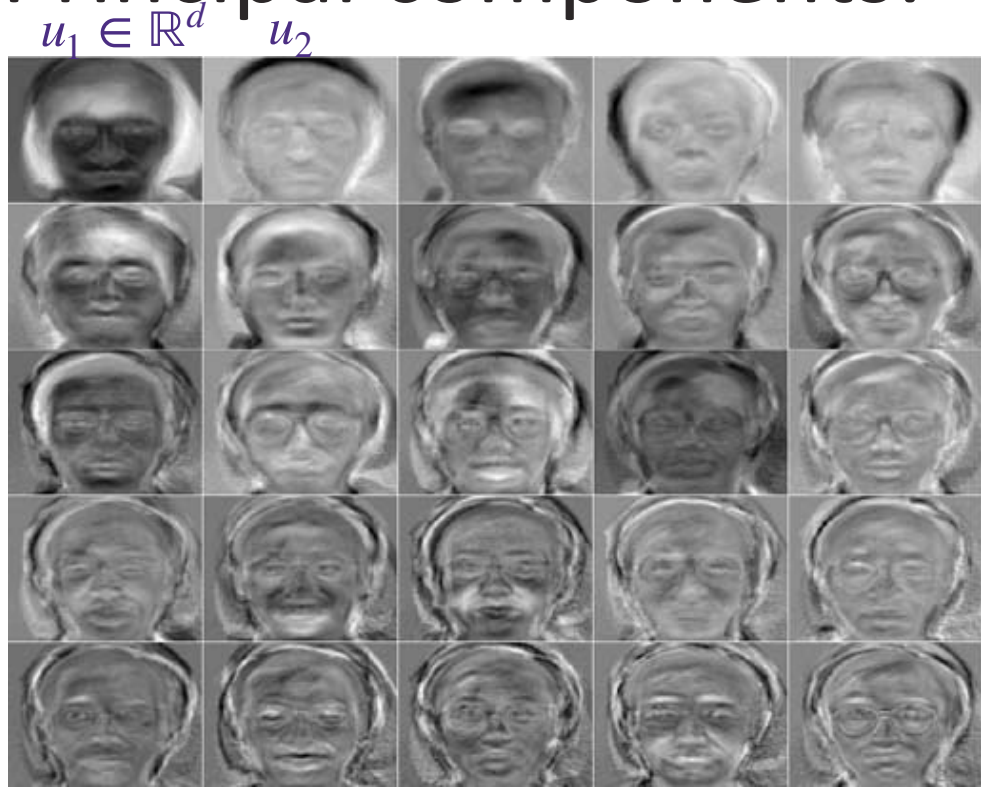


<https://en.wikipedia.org/wiki/Eigenface>

Principal component analysis finds a compact linear representation

- patterns that capture the distinct features of the samples is called **principal component** (to be formally defined later)
- we can represent each sample ^{face} as a **weighted linear combination** of, say, q=25 principal components, and just store the weights

Principal components:



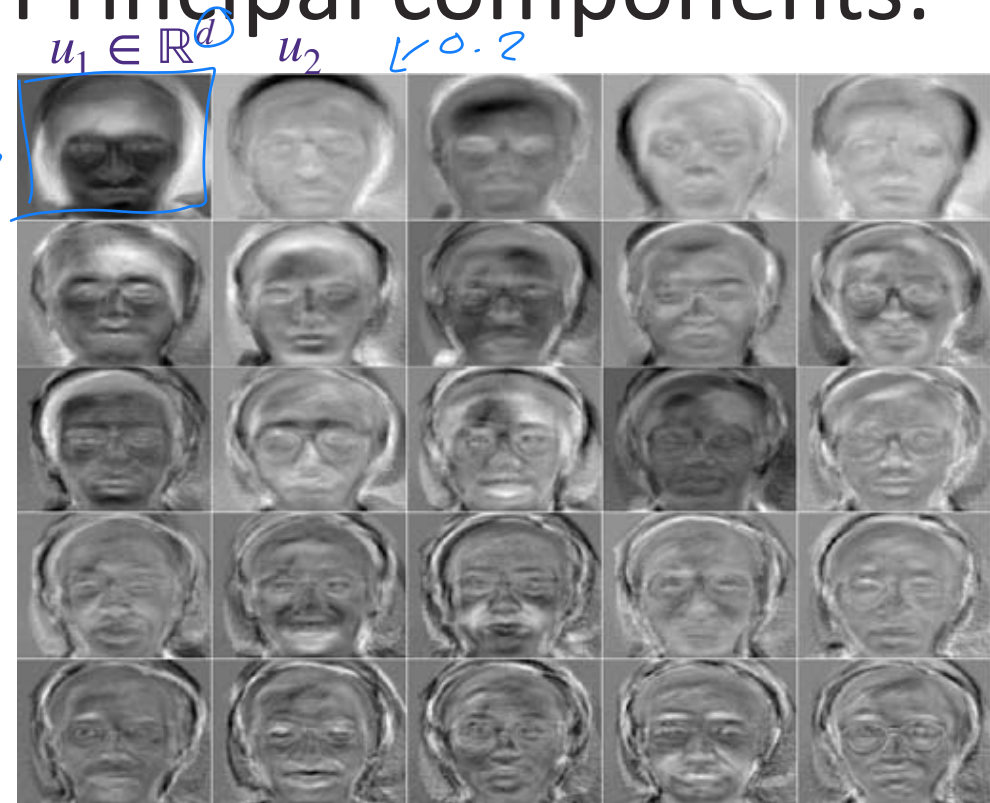
$$\approx z[1]u_1 + z[2]u_2 + \dots + z[25]u_{25}$$

Principal component analysis finds a compact linear representation

$32 \times 32 = 1024$ pixels

Principal components:

- patterns that capture the distinct features of the samples is called **principal component** (to be formally defined later)
- we can represent each sample as a **weighted linear combination** of, say, $q=25$ principal components, and just store the weights



$$\approx z[1]u_1 + z[2]u_2 + \dots + z[25]u_{25}$$

basis face

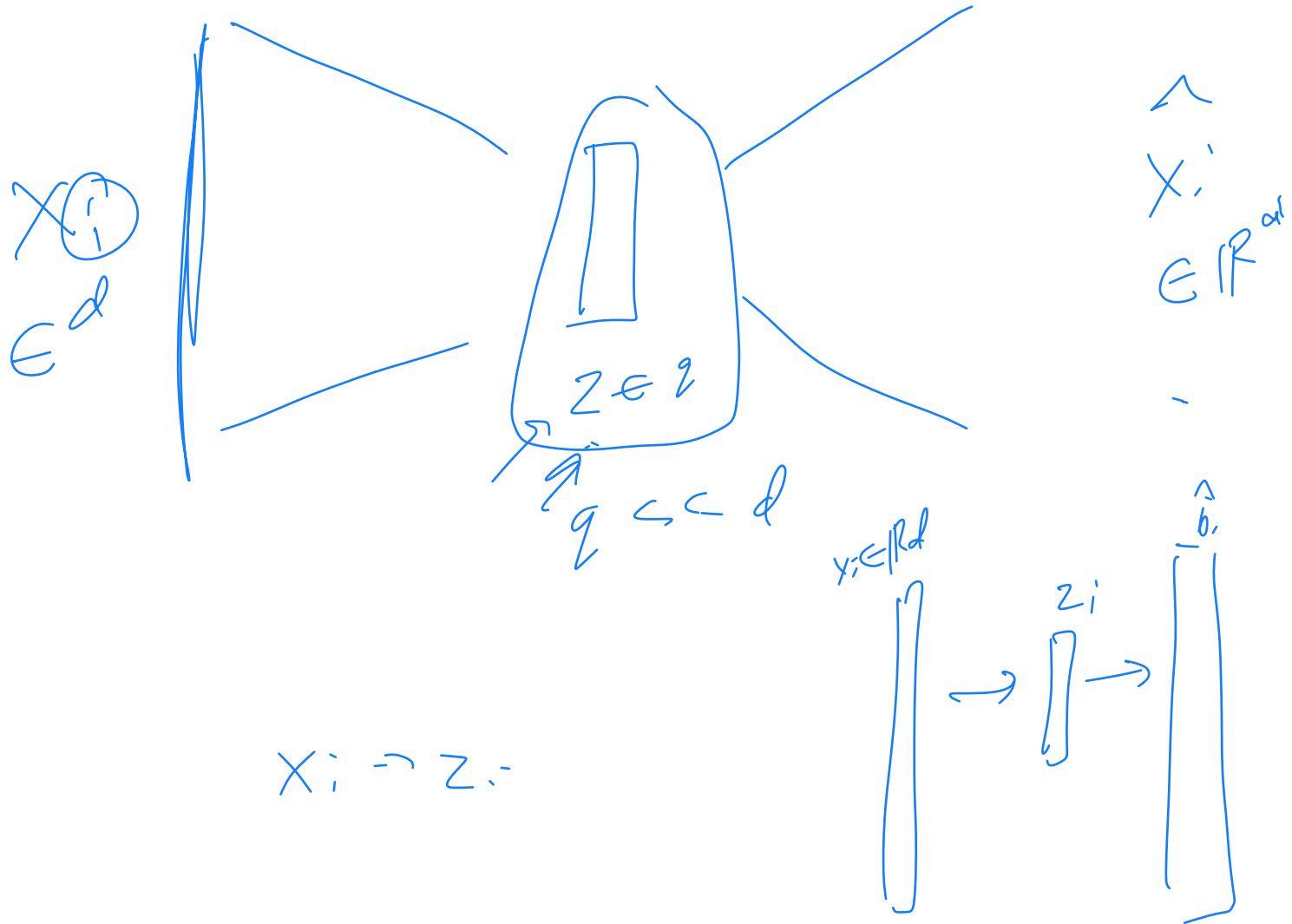
- With $q=25$, to store n images, it requires memory of only

$$d \times q + q \times n \ll d \times n$$

prev: $n \times d$

auto encoder

→ learned version of PCA



10 principal components give a pretty good reconstruction of a face

average face $\bar{x} + a[1]u_1$ $\bar{x} + a[1]u_1 + a[2]u_2$

\bar{x}

r=1

r=2

r=3

r=4



r=7

r=8

r=9

r=10

Ground truths real face

$10 \ll 1024$
reconstruction error decreases
as q increases
 $\rightarrow 0$ if $q = n$

PCA: a high-fidelity linear projection

Given $x_1, \dots, x_n \in \mathbb{R}^d$, find a compressed representation $z_1, \dots, z_n \in \mathbb{R}^q$

with $q \ll d$ such that $x_i \approx \bar{x} + \mathbf{V}_q z_i$ and $\mathbf{V}_q^\top \mathbf{V}_q = \mathbf{I}$.

$$\min_{\mathbf{V}_q, \{z_i\}} \sum_{i=1}^n \|x_i - \bar{x} - \mathbf{V}_q z_i\|_2^2$$

coefficients on basis vectors
 $\mathbf{V}_q \in \mathbb{R}^{d \times q}$
 each column is a basis vector

orthonormal
 all columns have unit norm and are mutually orthogonal
 \downarrow
 no redundancies
 \downarrow
 unique solution

$$x_i \approx \begin{bmatrix} \vdots \\ \bar{x} \\ \vdots \end{bmatrix} + \begin{bmatrix} | \\ V_1 \\ | \\ \vdots \\ | \\ V_q \\ | \end{bmatrix} \begin{bmatrix} z_{:,1} \\ \vdots \\ z_{:,q} \end{bmatrix}$$

$d \times q$
 \mathbf{V}_q z_i

$$= \bar{x} + z_{:,1} V_1 + \dots + z_{:,q} V_q$$

PCA: a high-fidelity linear projection

Given $x_1, \dots, x_n \in \mathbb{R}^d$, find a compressed representation $z_1, \dots, z_n \in \mathbb{R}^q$ with $q \ll d$ such that $x_i \approx \bar{x} + \mathbf{V}_q z_i$ and $\mathbf{V}_q^\top \mathbf{V}_q = \mathbf{I}$.

$$\min_{\mathbf{V}_q, \{z_i\}} \sum_{i=1}^n \underbrace{\|x_i - \bar{x} - \mathbf{V}_q z_i\|_2^2}_{\text{reconstructed } x_i}$$

reconstruction
MSE

→ how to solve?

$$\|y - Xw\|_2^2$$

$$\hat{w} = (X^\top X)^{-1} X^\top y$$

if you know \mathbf{V}_q :

$$z_i := \underbrace{(\mathbf{V}_q^\top \mathbf{V}_q)^{-1}}_{\sim \mathbf{I}} \mathbf{V}_q (x_i - \bar{x})$$

$$z_i = \mathbf{V}_q^\top (x_i - \bar{x})$$

Fix \mathbf{V}_q and solve for $\{z_i\}$:

PCA: a high-fidelity linear projection

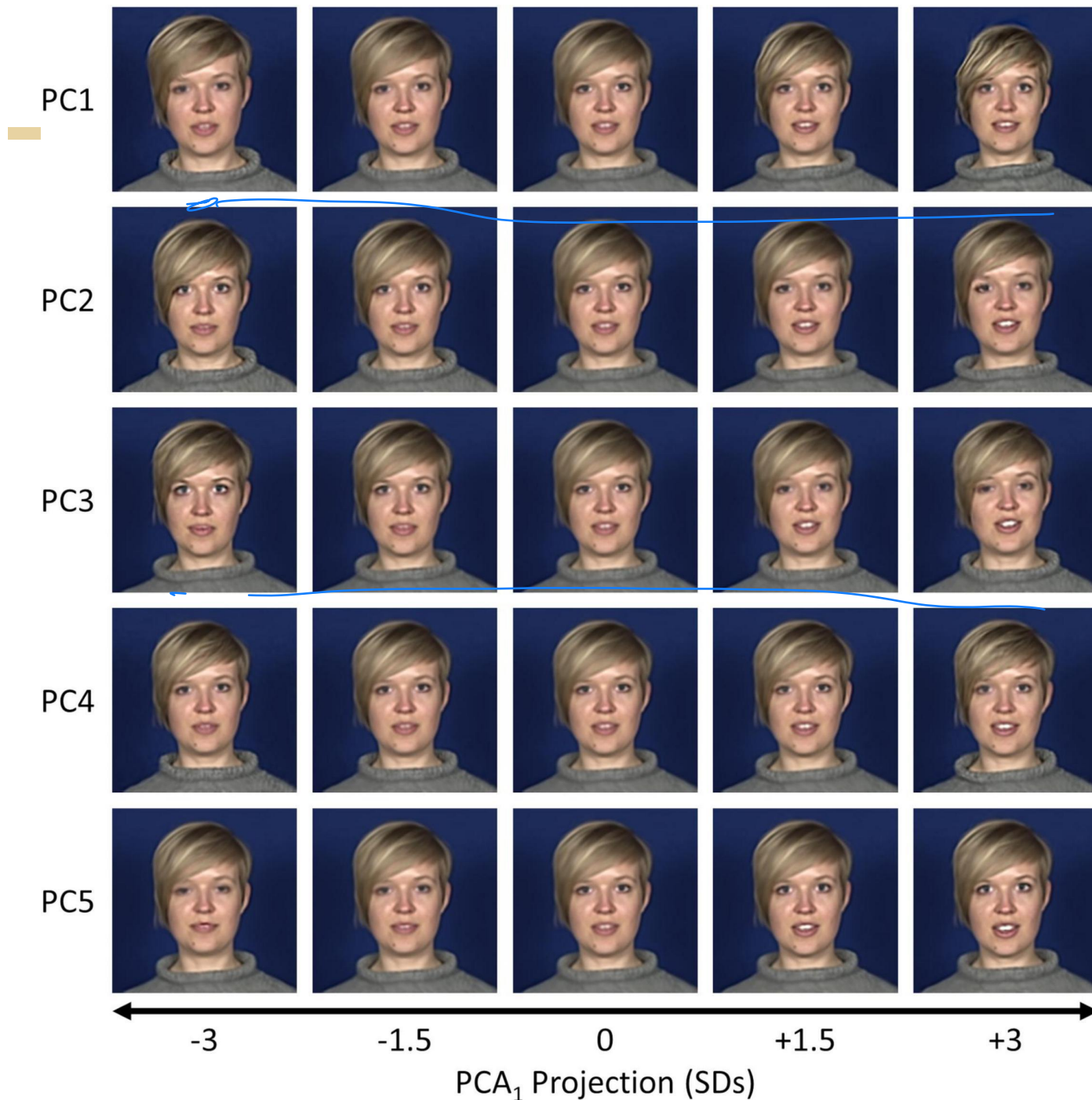
Given $x_1, \dots, x_n \in \mathbb{R}^d$, find a compressed representation $z_1, \dots, z_n \in \mathbb{R}^q$ with $q \ll d$ such that $x_i \approx \bar{x} + \mathbf{V}_q z_i$ and $\mathbf{V}_q^\top \mathbf{V}_q = \mathbf{I}$.

$$\min_{\mathbf{V}_q, \{z_i\}} \sum_{i=1}^n \|x_i - \bar{x} - \mathbf{V}_q z_i\|_2^2$$

Fix \mathbf{V}_q and solve for $\{z_i\}$: $z_i = \mathbf{V}_q^\top (x_i - \bar{x})$

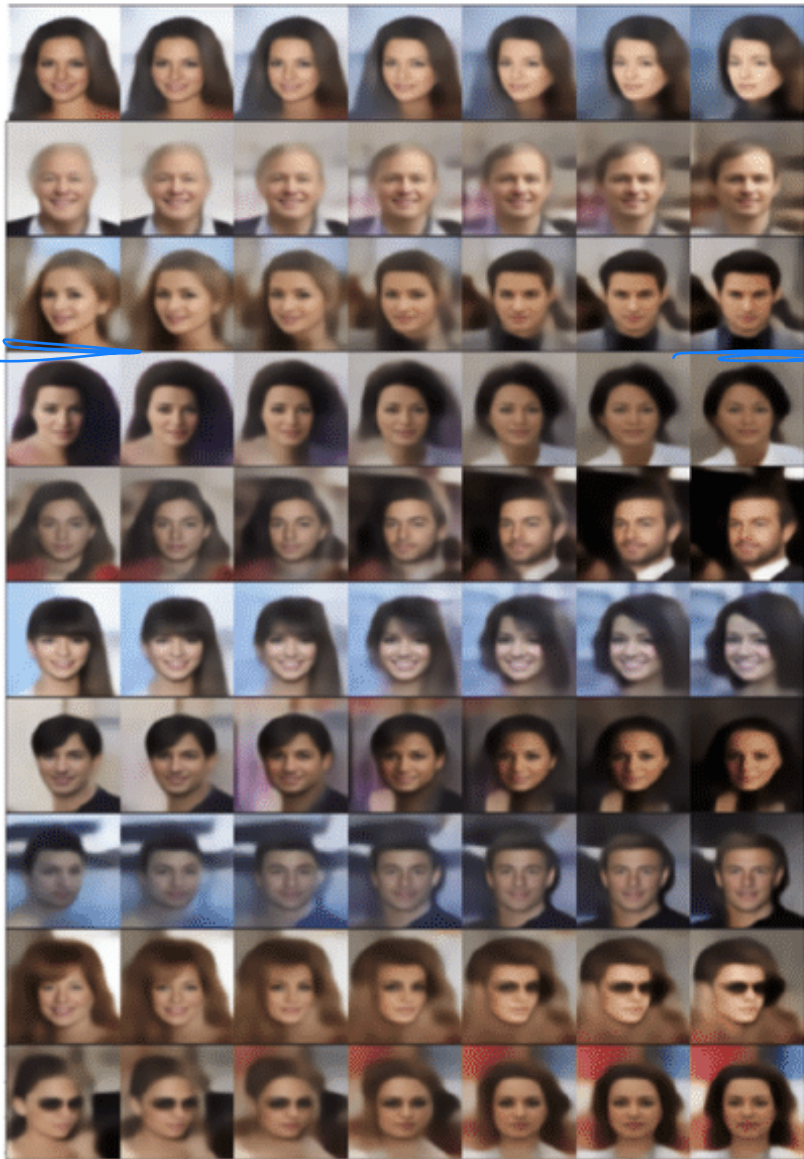
$$\hat{x}_i := \bar{x} + \mathbf{V}_q \underbrace{\mathbf{V}_q^\top (x_i - \bar{x})}_{z_i} = \bar{x} + \sum_{j=1}^q v_j v_j^\top (x_i - \bar{x})$$

→ projects onto
q-dimensional
linear subspace
embedded into
original high D
space (\mathbb{R}^d)



A PCA-Based Active Appearance Model for Characterising Modes of Spatiotemporal Variation in Dynamic Facial Behaviours
<https://www.frontiersin.org/journals/psychology/articles/10.3389/fpsyg.2022.880548/full>

Sample linear interpolations in the encoding space



$$\psi(\alpha\phi(x_1) + (1 - \alpha)\phi(x_2)), \quad \alpha \in [0,1]$$

PCA modes calculated in the encoding space

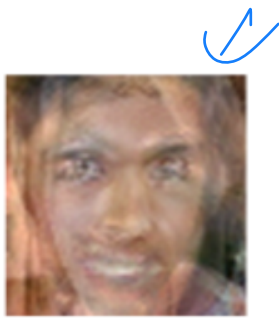
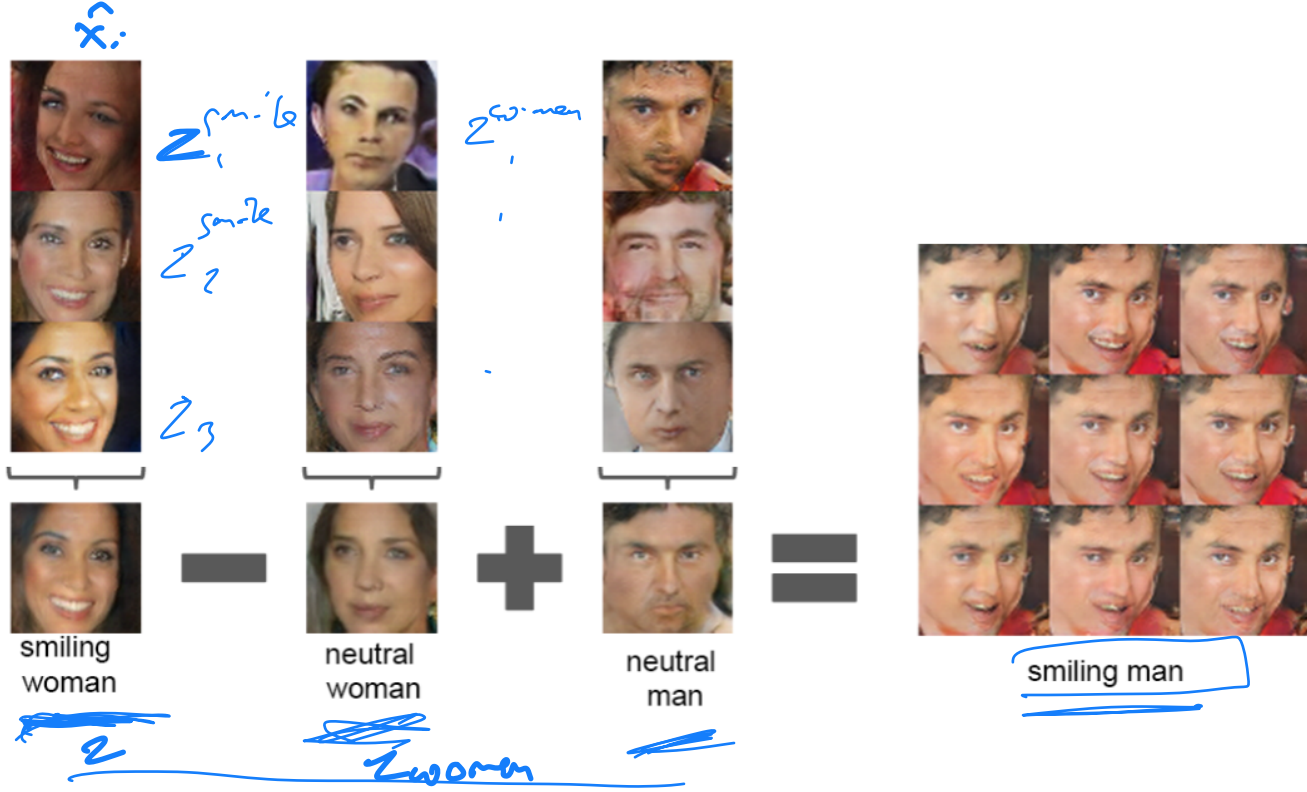


-3σ -2σ $-\sigma$ 0 σ 2σ 3σ

Sliced Wasserstein Autoencoder: An Embarrassingly Simple Generative Model

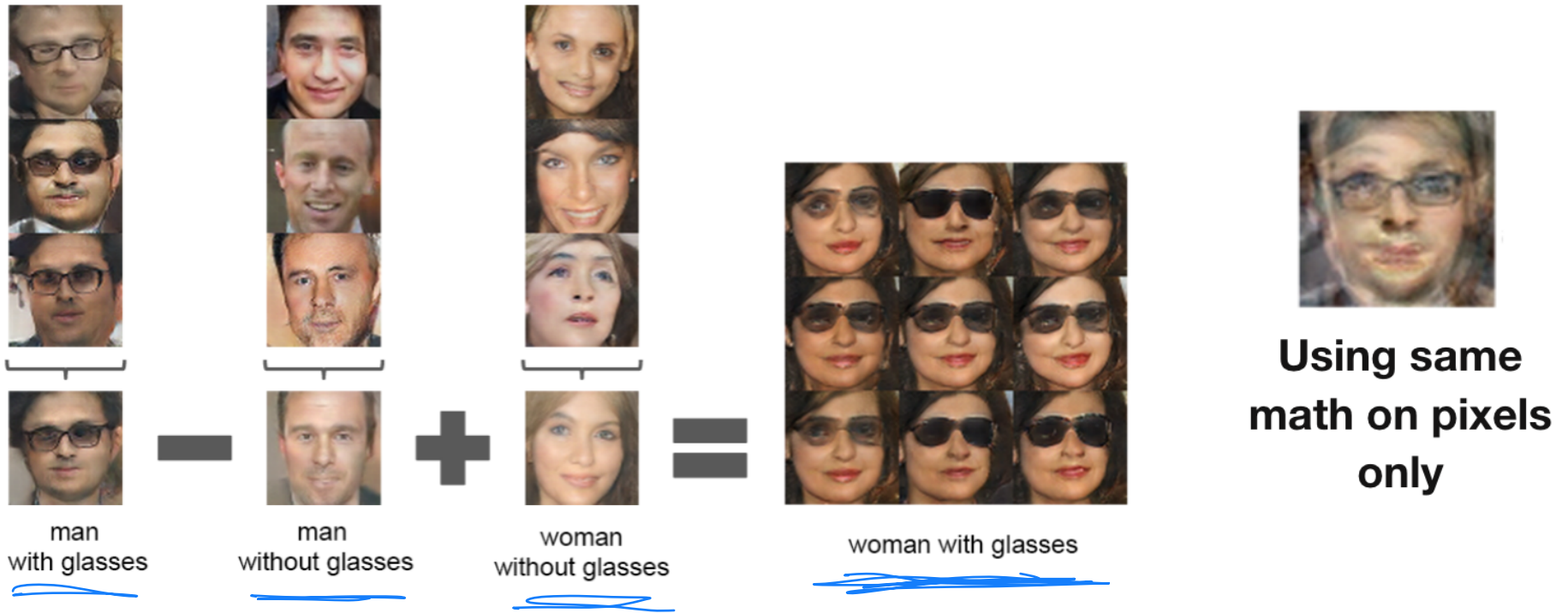
https://www.researchgate.net/figure/The-results-of-SWAE-on-the-CelebA-face-dataset-with-a-128-dimensional-uniform_fig5_324246144

Representations of images are meaningful



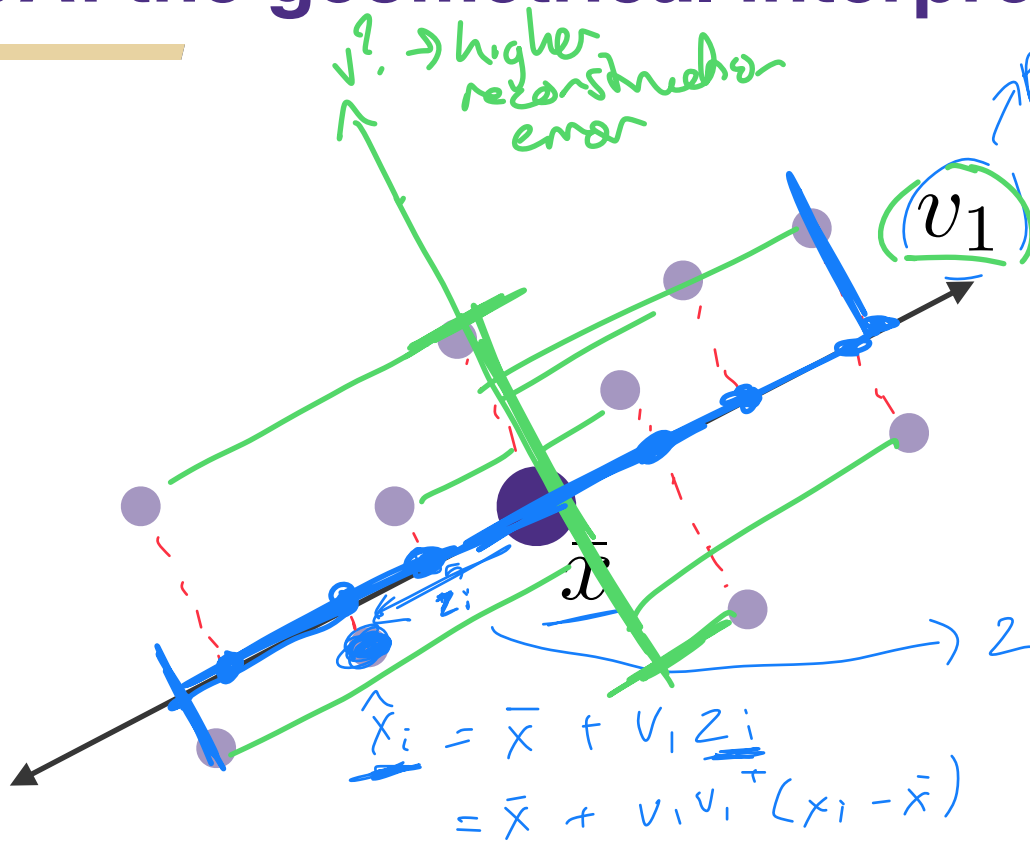
Using same math on pixels only

Representations of images are meaningful



Radford, A., Metz, L., & Chintala, S. (2015). [Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks](https://arxiv.org/abs/1511.06434). arXiv preprint arXiv:1511.06434.

PCA: the geometrical interpretation



principal component / basis vector

$$\frac{d=2}{q=1}$$

reconstruction error is the distance b/w the original point & projected version

$$z_i = v_1^T (x_i - \bar{x})$$

$$\begin{aligned} \hat{x}_i &= \bar{x} + v_1 z_i \\ &= \bar{x} + v_1 v_1^T (x_i - \bar{x}) \end{aligned}$$

Goal: orient the direction of v_1 to minimize squared reconstruction error

→ maximize the variance captured in the low-d representation

PCA: a high-fidelity linear projection

Given $x_1, \dots, x_n \in \mathbb{R}^d$, find a compressed representation $z_1, \dots, z_n \in \mathbb{R}^q$ with $q \ll d$ such that $x_i \approx \bar{x} + \mathbf{V}_q z_i$ and $\mathbf{V}_q^\top \mathbf{V}_q = \mathbf{I}$.

$$\min_{\mathbf{V}_q, \{z_i\}} \sum_{i=1}^n \|x_i - \bar{x} - \mathbf{V}_q z_i\|_2^2$$

Fix \mathbf{V}_q and solve for $\{z_i\}$: $z_i = \mathbf{V}_q^\top (x_i - \bar{x})$

$$\hat{x}_i := \bar{x} + \mathbf{V}_q \mathbf{V}_q^\top (x_i - \bar{x}) = \bar{x} + \sum_{j=1}^q v_j v_j^\top (x_i - \bar{x})$$

$$\min_{\mathbf{V}_q} \sum_{i=1}^n \left\| (x_i - \bar{x}) - \mathbf{V}_q \mathbf{V}_q^\top (x_i - \bar{x}) \right\|_2^2$$

$\mathbf{V}_q \mathbf{V}_q^\top$ is a projection matrix that minimizes error in basis of size q

PCA: a high-fidelity linear projection

$$\min_{\mathbf{V}_q} \sum_{i=1}^N \left\| (x_i - \bar{x}) - \mathbf{V}_q \mathbf{V}_q^T (x_i - \bar{x}) \right\|_2^2$$

$\mathbf{V}_q \mathbf{V}_q^T$ is a projection matrix that minimizes error in basis of size q

$$\hat{x}_i := \bar{x} + \mathbf{V}_q \mathbf{V}_q^T (x_i - \bar{x}) = \bar{x} + \sum_{j=1}^q \underbrace{v_j v_j^T}_a (x_i - \bar{x})$$

$\mathbf{V}_q^T \mathbf{V}_q = \mathbf{I}$
orthonormal

Case when $q = 1$ → $v_1 = \arg \min_{v: \|v\|_2=1} \sum_{i=1}^N \left\| (x_i - \bar{x}) - v v^T (x_i - \bar{x}) \right\|_2^2$

$$\|a - b\|_2^2 = (a - b)^T (a - b) = a^T a - 2a^T b + b^T b$$

$$v_1 = \arg \min_{v: \|v\|_2=1} \sum_{i=1}^N \|x_i - \bar{x}\|_2^2 - 2 \underbrace{(x_i - \bar{x})^T v v^T (x_i - \bar{x})}_{\text{same}}$$

result depend on v

$$+ \underbrace{(x_i - \bar{x})^T v v^T v v^T (x_i - \bar{x})}_{\text{same}}$$

$$= \arg \min_{v: \|v\|_2=1} - \sum_{i=1}^N (x_i - \bar{x})^T v v^T (x_i - \bar{x})$$

$$\Rightarrow \textcircled{1} \arg \max_{v: \|v\|_2=1} \sum_{i=1}^N \underbrace{(x_i - \bar{x})^T v}_{z_i} \underbrace{v^T (x_i - \bar{x})}_{z_i^T} = \sum_{i=1}^N z_i z_i^T$$

Show $\max_v \{z: z_i^T\}$ ends up maximizing the variance:

$$\bar{z} = \frac{1}{n} \sum_{i=1}^n z_i$$

$$= \frac{1}{n} \sum_{i=1}^n v^T (x_i - \bar{x})$$

$$= 0$$

$$\arg \max_v \sum_{i=1}^n z_i z_i^T = \sum_{i=1}^n (z_i - \bar{z})(z_i - \bar{z})^T$$

n^* empirical variance of \underline{z}_i

Minimizing reconstruction error is equivalent to maximizing the variance of the projected data points

$$\textcircled{1} \underset{v}{\operatorname{argmax}} \sum_{i=1}^n \left((x_i - \bar{x})^T v \right) \underbrace{\left(v^T (x_i - \bar{x}) \right)^T}_{z_i}$$

$z_i \in \mathbb{R}^1 \rightarrow \text{scalar}$
 $x_i \in \mathbb{R}^d$

// take v outside sum
 see ridge regression

$$\underset{v}{\operatorname{argmax}} \sum_{i=1}^n v^T (x_i - \bar{x}) (x_i - \bar{x})^T v$$

$$\underset{v}{\operatorname{argmax}} v^T \left(\underbrace{\sum_{i=1}^n (x_i - \bar{x}) (x_i - \bar{x})^T}_{d \times d} \right) v$$

\rightarrow define to be $\Sigma \in \mathbb{R}^{d \times d}$
covariance matrix

$$= \underset{v}{\operatorname{argmax}} v^T \underbrace{\Sigma}_{d \times d} v$$

\rightarrow find the direction that maximizes the project variance

consider $v = e_j = [0 \dots 0 \mid 1 \mid 0 \dots 0]$

$$e_j^T \Sigma e_j = ?$$

sample variance of the j^{th} feature
 in . maximizing the variance in
 direction e_j
 \rightarrow not in practice a good idea

★ Solution / punchline:
choose v to be the leading eigenvectors
of Σ ↳ with highest
eigenvalues λ

defn of eigenvector & eigenvalues

$A v = \lambda v$

direction capturing the most variance

amount of variance captured

For a matrix $A \in \mathbb{R}^{d \times d}$ we say (λ, v) is
(eigenvalue, eigenvector) pair if $A v = \lambda v$

Fact: if A is symmetric, then all of its
eigenvalues are real and

$$A = \sum_{i=1}^d \lambda_i v_i v_i^T \quad \text{where } v_i^T v_j = \begin{cases} 1 & i=j \\ 0 & i \neq j \end{cases}$$

$$= \underline{V \Lambda V^T}$$

↪ diagonal matrix where
elements are eigenvalues λ

if A is PSD
all eigenvalues are
non-negative.

$$A v_i = \lambda_i v_i$$
$$v_i^T A v_i = \lambda_i \underbrace{v_i^T v_i}_{=1}$$

$$\underline{\underline{v_i^T A v_i}} = \lambda_i$$

→ \int by orthonormal
can save work

if A is a covariance matrix $\times \Sigma$
then it's symmetric & PSD

and if $\|v_i\|_2 = 1$

then $v_i^T A v_i$ is the variance along direction
 v_i (which is λ_i)

variance should be real & non-negative.

Find structure in your data with basic
linear algebra

PCA: a high-fidelity linear projection

$$\min_{\mathbf{V}_q} \sum_{i=1}^N \left\| (x_i - \bar{x}) - \mathbf{V}_q \mathbf{V}_q^\top (x_i - \bar{x}) \right\|_2^2$$

$\mathbf{V}_q \mathbf{V}_q^\top$ is a projection matrix that minimizes error in basis of size q

$$\hat{x}_i := \bar{x} + \mathbf{V}_q \mathbf{V}_q^\top (x_i - \bar{x}) = \bar{x} + \sum_{j=1}^q v_j \langle v_j, x_i - \bar{x} \rangle$$

Case when $q = 1$

$$v_1 = \arg \min_{v: \|v\|_2=1} \sum_{i=1}^N \left\| (x_i - \bar{x}) - v v^\top (x_i - \bar{x}) \right\|_2^2$$

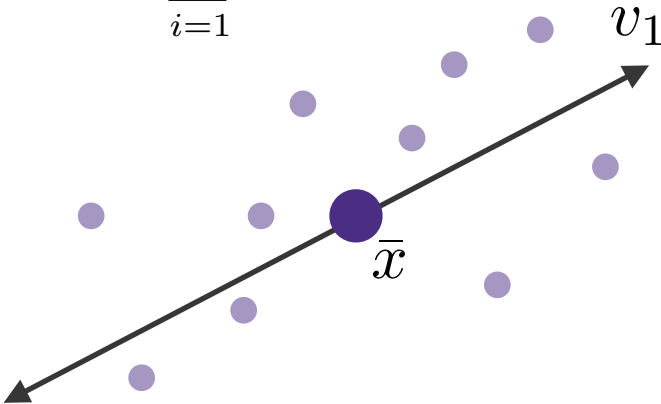
$$= \arg \min_{v: \|v\|_2=1} \sum_{i=1}^N \left(\|x_i - \bar{x}\|_2^2 - 2(x_i - \bar{x})^\top v v^\top (x_i - \bar{x}) + (x_i - \bar{x})^\top v v^\top v v^\top (x_i - \bar{x}) \right)$$

$$= \arg \min_{v: \|v\|_2=1} \sum_{i=1}^N \|x_i - \bar{x}\|_2^2 - \sum_{i=1}^N (x_i - \bar{x})^\top v v^\top (x_i - \bar{x})$$

$$= \arg \max_{v: \|v\|_2=1} \sum_{i=1}^N (x_i - \bar{x})^\top v v^\top (x_i - \bar{x})$$

$$= \arg \max_{v: \|v\|_2=1} v^\top \Sigma v$$

$$\Sigma := \sum_{i=1}^N (x_i - \bar{x})(x_i - \bar{x})^\top$$



PCA: a high-fidelity linear projection

$$\min_{\mathbf{V}_q} \sum_{i=1}^N \left\| (x_i - \bar{x}) - \mathbf{V}_q \mathbf{V}_q^T (x_i - \bar{x}) \right\|_2^2$$

$\mathbf{V}_q \mathbf{V}_q^T$ is a *projection matrix* that minimizes error in basis of size q

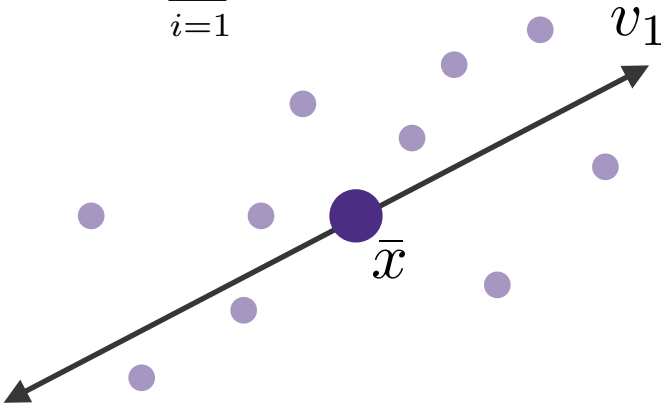
$$\hat{x}_i := \bar{x} + \mathbf{V}_q \mathbf{V}_q^T (x_i - \bar{x}) = \bar{x} + \sum_{j=1}^q v_j \langle v_j, x_i - \bar{x} \rangle$$

Case when $q = 1$

$$v_1 = \arg \min_{v: \|v\|_2=1} \sum_{i=1}^N \left\| (x_i - \bar{x}) - v v^T (x_i - \bar{x}) \right\|_2^2$$

$$= \arg \max_{v: \|v\|_2=1} v^T \Sigma v$$

$$\Sigma := \sum_{i=1}^N (x_i - \bar{x})(x_i - \bar{x})^T$$



PCA: a high-fidelity linear projection

$$\min_{\mathbf{V}_q} \sum_{i=1}^N \left\| (x_i - \bar{x}) - \mathbf{V}_q \mathbf{V}_q^T (x_i - \bar{x}) \right\|_2^2$$

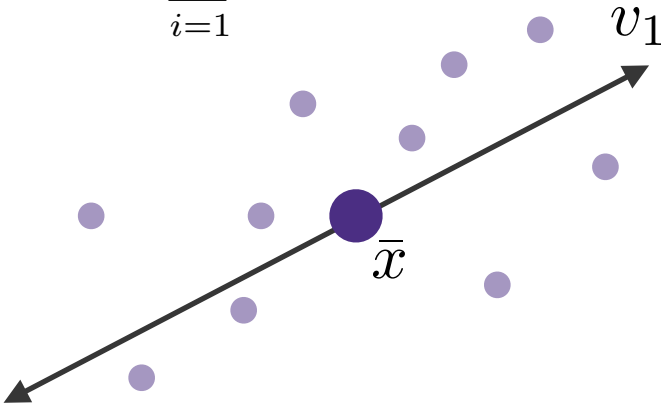
$\mathbf{V}_q \mathbf{V}_q^T$ is a projection matrix that minimizes error in basis of size q

$$\hat{x}_i := \bar{x} + \mathbf{V}_q \mathbf{V}_q^T (x_i - \bar{x}) = \bar{x} + \sum_{j=1}^q v_j \langle v_j, x_i - \bar{x} \rangle$$

General $q \geq 1$

$$\min_{\mathbf{V}_q} \sum_{i=1}^N \left\| (x_i - \bar{x}) - \mathbf{V}_q \mathbf{V}_q^T (x_i - \bar{x}) \right\|_2^2 = \min_{\mathbf{V}_q} \text{Tr}(\Sigma) - \text{Tr}(\mathbf{V}_q^T \Sigma \mathbf{V}_q)$$

$$\Sigma := \sum_{i=1}^N (x_i - \bar{x})(x_i - \bar{x})^T$$



$$\text{Tr}(ABC) = \text{Tr}(CAB)$$

$$= \text{Tr}(BCA)$$

↳ replace where you swapped order of 2; scalars

\mathbf{V}_q are the first q eigenvectors of Σ

Minimize reconstruction error = capture the most variance in your data.



How to choose the dimensionality, q

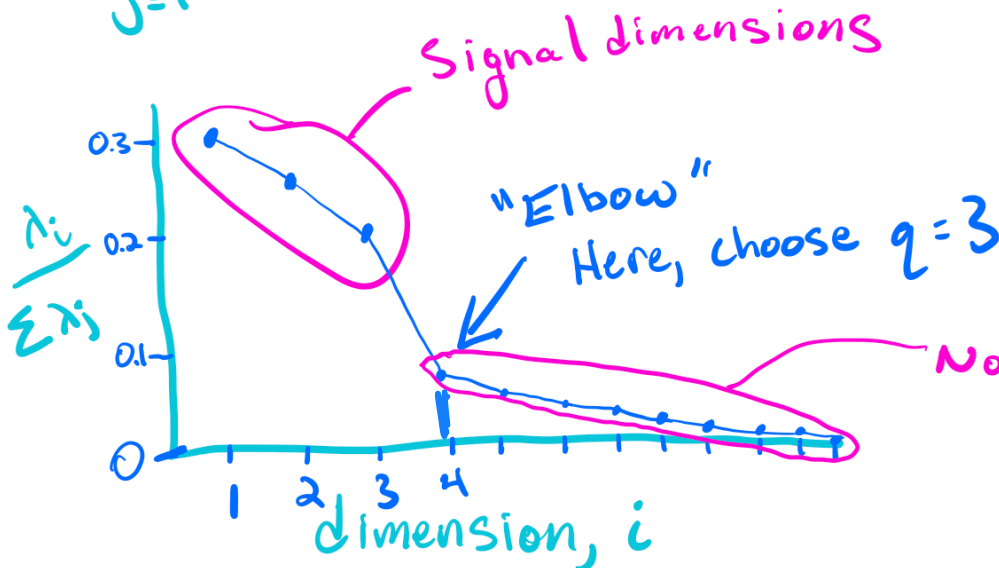
HOW TO CHOOSE q

↑ model complexity

CROSS VALIDATION DOESN'T WORK

- More dimensions always increases projected variance (decreases reconstruction error), INCLUDING ON VAL DATA.

$$\frac{\lambda_i}{\sum_{j=1}^d \lambda_j} = \frac{\text{variance along } v_i}{\text{total variance}}$$



sorted by the largest eigenvalues

- Ad-hoc approach:
 - # dims needed to explain 95% of variance.
- Leave-one-feature-out ^{cross-validation} (LOFO-CV)

PCA: a high-fidelity linear projection

Given $x_i \in \mathbb{R}^d$ and some $q < d$ consider

$$\rightarrow \min_{\mathbf{V}_q} \sum_{i=1}^N \|(x_i - \bar{x}) - \mathbf{V}_q \mathbf{V}_q^T (x_i - \bar{x})\|^2.$$

where $\mathbf{V}_q = [v_1, v_2, \dots, v_q]$ is orthonormal:

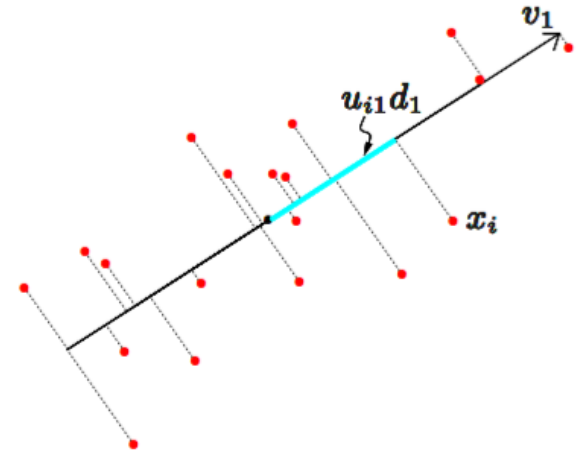
$$\mathbf{V}_q^T \mathbf{V}_q = I_q$$

\mathbf{V}_q are the first q eigenvectors of Σ

$\bar{\mathbf{V}}_q$ are the first q principal components

Principal Component Analysis (PCA) projects $(\mathbf{X} - \mathbf{1}\bar{x}^T)$ down onto \mathbf{V}_q

$$(\mathbf{X} - \mathbf{1}\bar{x}^T)\mathbf{V}_q = \mathbf{U}_q \text{diag}(d_1, \dots, d_q) \quad \mathbf{U}_q^T \mathbf{U}_q = I_q$$



$$\Sigma := \sum_{i=1}^N (x_i - \bar{x})(x_i - \bar{x})^T$$