

A goal of PCA is to find projections \tilde{x}_n of data points x_n that are as similar to the original data points as possible, with a significantly lower intrinsic dimensionality.

We consider an i.i.d. dataset $\mathcal{X} = \{x_1, \dots, x_n\}, x_n \in \mathbb{R}^p$ with mean μ that possess the data covariance matrix

$$S = \frac{1}{N} \sum_{n=1}^N x_n x_n^T. \quad \text{We also assume } \exists \text{ low-dim rep:}$$

$z_n = B^T x_n \in \mathbb{R}^m$ of x_n , where we define the projection matrix

$$B := [b_1, \dots, b_m] \in \mathbb{R}^{p \times m}. \quad \text{Assume the columns of } B \text{ are orthonormal.}$$

We want an M -dimensional subspace $U \subseteq \mathbb{R}^p$, $\dim(U) = M < p$ where we will project the data. Projected data is denoted

$\tilde{x}_n \in U$, and their coordinates (with respect to the basis vectors b_1, \dots, b_m of U) by z_n and the basis vectors b_1, \dots, b_M .

The projected data should be similar to the original data x_n and minimize the loss due to compression.

Consider \mathbb{R}^2 , $e_1 = [1, 0]^T$, $e_2 = [0, 1]^T$

$$\text{Original: } \mathbb{R}^2 \xrightarrow{\text{Compression}} \mathbb{R}^m \xrightarrow{\text{Reconstructed}} \tilde{x}$$

$\begin{bmatrix} z \\ 1 \end{bmatrix} = \begin{bmatrix} e_1 & e_2 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix}$. Now, consider vectors of the form $\tilde{x} = \begin{bmatrix} 0 \\ z \end{bmatrix} \in \mathbb{R}^2$, $z \in \mathbb{R}$. written as: $0e_1 + ze_2$.

To represent these vectors it is sufficient to remember the coordinate z of \tilde{x} with respect to the e_2 vector.

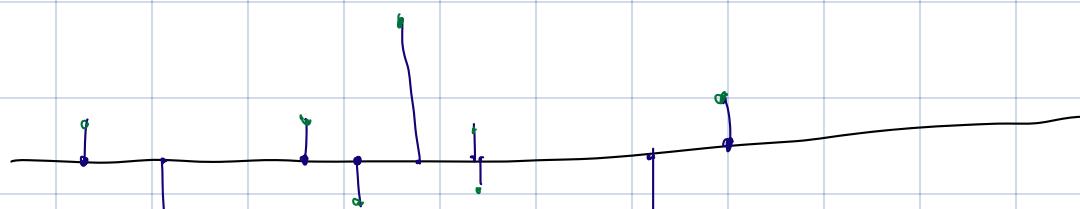
The set of \tilde{x} vectors ($+z$) forms a vector subspace U with $\dim(U) = 1 \Rightarrow U = \text{span}[e_2]$

Variance - an indicator of spread in the data PCA - a dimensionality reduction algorithm that maximizes the variance in the low-dimensional representation of the data to retain as much information as possible.

Goal: find a matrix B that retains as much information as possible when compressing data by projecting it onto the subspace spanned by columns b_1, \dots, b_m of B .

Retaining most information after data compression is equivalent to capturing the largest amount of variance in the low-dimensional code. Code = Coordinates

The variance of low-dimensional code does not depend on the mean of the data. \Rightarrow We assume without loss of generality that the data has mean 0 . \Rightarrow the mean of the low-dim code is also 0 . $E_z[z^2] = E_x[B^T x] = B^T E_x[x] = 0$.



PCA finds a lower-dimensional subspace (purple line) that maintains as much variance (spread of data) as possible when the data (green) is projected onto the subspace (purple)

Direction with Maximal Variance

We start by seeking a single vector $b_1 \in \mathbb{R}^D$ that maximizes the variance of the projected data: aim to maximize the variance of the first component $z_{1,n}$ of $z \in \mathbb{R}^N$: $V_1 := \mathbb{V}[z_{1,n}] = \frac{1}{N} \sum_{n=1}^N z_{1,n}^2$ is maximized. The vector b_1 will be the first column of the matrix B and therefore the first of M orthonormal basis vectors that span the low-dimensional subspace.

The first component of z_n is given by $z_{1,n} = b_1^T x_n$, the coordinate of the orthogonal projection of x_n onto the 1D subspace spanned by b_1 . $V_1 = \frac{1}{N} \sum_{n=1}^N (b_1^T x_n)^2 = \frac{1}{N} \sum_{n=1}^N b_1^T x_n x_n^T b_1 = b_1^T \left(\frac{1}{N} \sum_{n=1}^N x_n x_n^T \right) b_1 = b_1^T S b_1$, where S is the data covariance matrix.

The arbitrarily increasing b_1 increases V_1 . Averaging b_1 that is twice as long results in a V_1 that is four times as long. \Rightarrow we restrict all solutions to $\|b_1\|^2 = 1$, resulting in a constrained optimization problem in which we seek the direction along which the data varies the most.

b_1 can be found by the constrained optimization problem: $\max_{b_1} b_1^T S b_1$ subject to $\|b_1\|^2 = 1$. We will use the Lagrangian:

$$\mathcal{L}(b_1, \lambda) = b_1^T S b_1 + \lambda (1 - b_1^T b_1). \text{ The partial derivatives of } \mathcal{L} \text{ with respect to } b_1 \text{ and } \lambda \text{ are:}$$

$$\frac{\partial \mathcal{L}}{\partial b_1} = 2b_1^T S - 2\lambda b_1, \quad \frac{\partial \mathcal{L}}{\partial \lambda} = 1 - b_1^T b_1. \text{ Setting to 0 gives us relations } b_1^T b_1 = \lambda b_1, \quad b_1^T b_1 = 1$$

b_1 is an eigenvector of S , and the Lagrange multiplier λ_1 is the eigenvalue: $V_1 = b_1^T S b_1 = \lambda_1 b_1^T b_1 = \lambda_1$.

the variance of the data projected onto a one-dimensional subspace equals the eigenvalue that is associated with the basis vector b_1 that spans this subspace.

The eigenvalue represents the variance of the data when projected onto the one-dimensional subspace.

The quantity $\sqrt{\lambda_1}$ is also called the loading of the unit vector b_1 and represents the standard deviation of the data accounted for by the principal subspace span b_1 .

To maximize the variance of the low-dimensional code, we choose the basis vector associated with the largest eigenvalue of S . This eigenvector is called the first principal component.

We determine the effect/contribution of the principal component b_1 in the original data space by multiplying $z_{1,n}$ with b_1 , giving us the projected data point

$$\tilde{x}_n = b_1 z_{1,n} = b_1^T x_n \in \mathbb{R}^D$$
 in the original data space.

Although \tilde{x}_n is a D -dimensional vector, it only requires a single coordinate $z_{1,n}$ to represent it with respect to the basis vector $b_1 \in \mathbb{R}^D$.

((U.3.))

M-dimensional Subspace with Maximal Variance

Assume we have found the first $m-1$ principal components as the $m-1$ eigenvectors of S that are associated with the largest $m-1$ eigenvalues.

S is symmetric, so the spectral theorem states that we can use these eigenvectors to construct an orthonormal eigensbasis of an $(m-1) \times m$ subspace of \mathbb{R}^D .

We can find the m^{th} principal component by subtracting the effect of the first $m-1$ principal components b_1, \dots, b_{m-1} from the data, trying to find principal components that captures the remaining information. \rightarrow contains the data that has not yet been compressed.

We have a new data matrix: $\hat{X} := X - \sum_{i=1}^{m-1} b_i b_i^T X = X - B_{m-1} X$, where $X = [x_1, \dots, x_N] \in \mathbb{R}^{D \times N}$ contains the data points as columns. Vector and $B_{m-1} := \sum_{i=1}^{m-1} b_i b_i^T$ is a projection matrix that projects onto the subspace spanned by b_1, \dots, b_{m-1} .

To find the m^{th} principal component, we maximize the variance: $V_m = \mathbb{V}[z_{1,m}] = \frac{1}{N} \sum_{n=1}^N z_{1,m}^2 = \frac{1}{N} \sum_{n=1}^N (b_m^T \hat{x}_n)^2 = b_m^T \hat{S} b_m$, subject to constrained optimization $\|b_m\|^2 = 1$.

We find the optimal solution b_m is the eigenvector of \hat{S} that is associated with the largest eigenvalue of \hat{S} . b_m is also an eigenvector of $S \Rightarrow$ the sets of eigenvectors for both S and \hat{S} are identical. \Rightarrow \exists D distinct eigenvectors for both S and \hat{S} .

Assume we have already found eigenvectors b_1, \dots, b_{m-1} of S . Consider an eigenvector b_m of S , i.e. $S b_m = \lambda_m b_m$.

$$\hat{S} b_m = \frac{1}{N} \hat{X} \hat{X}^T b_m = \frac{1}{N} (X - B_{m-1} X)(X - B_{m-1} X)^T b_m = (S - S B_{m-1} - B_{m-1} S + B_{m-1} S B_{m-1}) b_m.$$

If $i \geq m$, i.e. b_i is an eigenvector that is not among the first $m-1$ principal components, then b_i is orthogonal to the first $m-1$ principal components and $B_{m-1} b_i = 0$. If $i \leq m$, i.e. b_i is among the first $m-1$ principal components, then b_i is a basis vector of the principal subspace onto which B_{m-1} projects.

Since b_1, \dots, b_m are an ONSB of this principal subspace, we obtain $B_{m-1} b_i = b_i$.

$$B_{m-1} b_i = b_i \text{ if } i \leq m, \quad B_{m-1} b_i = 0 \text{ if } i \geq m$$

Using the case $i \geq m$, we obtain $\hat{S} b_i = (S - B_{m-1} S) b_i = S b_i - \lambda_i b_i$, i.e. b_i is also an eigenvector of \hat{S} with eigenvalue λ_i : $\hat{S} b_m = S b_m = \lambda_m b_m$.

b_m is an eigenvector of S and λ_m is the largest eigenvalue of S and the m th largest eigenvalue of \hat{S} .

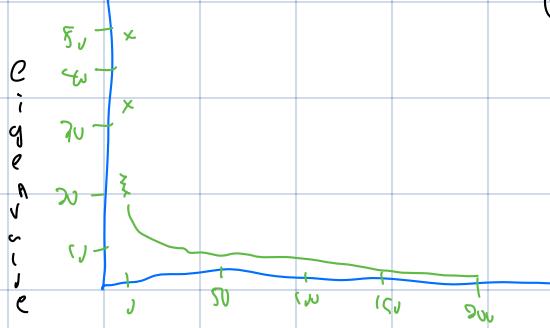
Using the case of $i \leq m$: $\hat{S} b_i = (S - S B_{m-1} - B_{m-1} S + B_{m-1} S B_{m-1}) b_i = 0 = 0 b_i$.

This shows that b_1, \dots, b_m are also eigenvectors of \hat{S} , but they are associated with eigenvalue 0 so that b_1, \dots, b_{m-1} span the null space of \hat{S} .

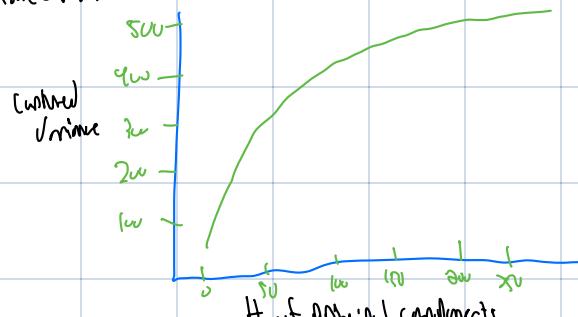
If the eigenvectors of S are part of the $(m-1)$ -dimensional principal subspace, then the associated eigenvalue of \hat{S} is 0.

Using $b_m^T b_m = 1$ and $\hat{S} b_m = S b_m = \lambda_m b_m$, the variance of the data projected onto the m^{th} principal component is

$$\lambda_m = b_m^T S b_m = \lambda_m b_m^T b_m = \lambda_m \Rightarrow \text{when projected onto an } M\text{-dimensional subspace, the variance of the data equals the sum of the eigenvalues that are associated with the corresponding eigenvectors of the data covariance matrix.}$$



Eigenvalues of the data covariance matrix of all digits "8" in the MNIST training set.



Cumulative variance captured by the principal components, collected with the largest eigenvalues.

To find an M -dimensional subspace of \mathbb{R}^D that retains as much information as possible, PCA fails to choose the columns of the matrix B as the M eigenvectors of the data covariance matrix S that are associated with the M largest eigenvalues. The most amount of variance PCA can capture with the first M principal components:

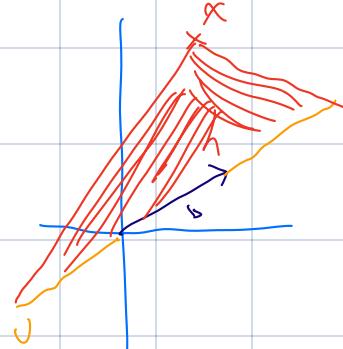
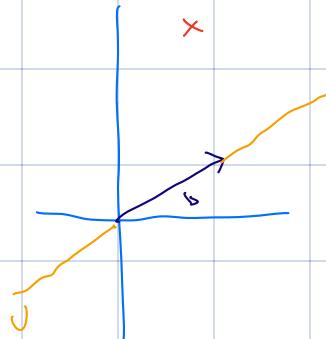
$$V_m := \sum_{n=1}^M \lambda_n,$$

where the λ_n are the M largest eigenvalues of the data covariance matrix S .

The variance lost by data compression via PCA: $S_m := \sum_{j=M+1}^D \lambda_j = V_D - V_m$

the relative variance captured is $\frac{V_m}{V_D}$ and the relative loss is $1 - \frac{V_m}{V_D}$.

(10.3 Projection Perspective)



A projection of a vector x onto a 1-D subspace $U \subseteq \mathbb{R}^D$ spanned by b . Shows different vectors between x and \hat{x} .

Assume $(\text{PCA})B = (b_1, \dots, b_M)$ of \mathbb{R}^D . We are interested in finding vector $\hat{x} \in \mathbb{R}^D$, which live in a lower-dimensional subspace $U \subseteq \mathbb{R}^D$, $\dim(U) = M$, so that

$\hat{x} = \sum_{m=1}^M \hat{x}_m b_m \in U \subseteq \mathbb{R}^D$ is as similar to x as possible. We want to find the best linear projection onto a lower-dimensional subspace U of \mathbb{R}^D with $\dim(U) = M$ and orthonormal basis vectors b_1, \dots, b_M . We call this subspace the principal subspace.

Orthogonal projections are denoted by: $\tilde{x}_n := \sum_{m=1}^M z_{mn} b_m \in \mathbb{R}^M$, where $z_{ni} := [z_{1n}, \dots, z_{Mn}]^T \in \mathbb{R}^M$ is the coordinate vector of \tilde{x}_n with respect to the basis (b_1, \dots, b_M) . We want \tilde{x}_n as similar to x_n as possible.

We will use the Euclidean Norm (squared distance) $\|x - \tilde{x}\|^2$ to measure similarity. Thus, our objective is to minimize the average squared Euclidean distance $J_m := \frac{1}{N} \sum_{n=1}^N \|x_n - \tilde{x}_n\|^2$. To find the optimal linear projection, we need to find the orthonormal basis of the principal subspace and the coordinates $z_n \in \mathbb{R}^M$ of the projection with respect to the basis b .

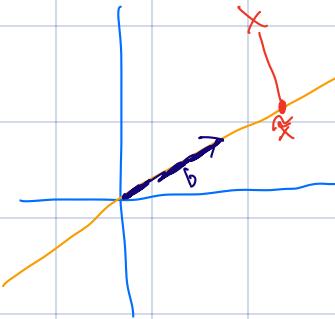
To find the OWA, there is a two step approach: 1) optimize coordinates z_n for a given OWA (b_1, \dots, b_M)

2) find the optimal OWA

Finding the optimal coordinates z corresponds to finding the representation of the linear projection \tilde{x} with respect to b that minimizes the distance between $\tilde{x} - x$.

We assume an OWA (b_1, \dots, b_M) of $U \subseteq \mathbb{R}^M$. We require partial derivatives:

$$\frac{\partial J_m}{\partial z_{in}} = \frac{\partial J_m}{\partial \tilde{x}_n} \frac{\partial \tilde{x}_n}{\partial z_{in}}, \quad \frac{\partial J_m}{\partial \tilde{x}_n} = -\frac{2}{N} (x_n - \tilde{x}_n)^T \in \mathbb{R}^{1 \times M}$$



The vector that minimizes the distance is its orthogonal projection onto U . The coordinate of the projection \tilde{x}_n with respect to the basis vectors b that spans U is the factor we need to scale b in order to reach \tilde{x}_n .

$$\frac{\partial \tilde{x}_n}{\partial z_{in}} = \frac{\partial}{\partial z_{in}} \left(\sum_{m=1}^M z_{mn} b_m \right) = b_i; \quad \text{for } i=1, \dots, M, \text{ we obtain}$$

$$\frac{\partial J_m}{\partial z_{in}} = -\frac{2}{N} (x_n - \tilde{x}_n)^T b_i = -\frac{2}{N} (x_n - \tilde{x}_n)^T \left(\sum_{m=1}^M z_{mn} b_m \right)^T b_i$$

OWA $= -\frac{2}{N} (x_n^T b_i - z_{in} b_i^T b_i) = -\frac{2}{N} (x_n^T b_i - 2z_{in})$ since $b_i^T b_i = 1$. Setting this partial derivative to zero yields the optimal coord.

$z_{in} = x_n^T b_i = b_i^T x_n$, for $i=1, \dots, M$, and $i=1, \dots, N \Rightarrow$ the optimal coordinates of the projection \tilde{x}_n are the coordinates of the orthogonal projection of the original data point x_n onto the one-dimensional subspace spanned by b_i .

- The optimal linear projection \tilde{x}_n of x_n is an orthogonal projection.
- The coordinates of \tilde{x}_n with respect to the basis (b_1, \dots, b_M) are the coordinates of the orthogonal projection of x_n onto the principal subspace.
- An orthogonal projection is the best linear mapping given the objective
- The coordinates $\{z_m\}$ of x and $\{z_m\}$ of x must be identical for $m=1, \dots, M$ since $U^\perp = \text{Span}[\{b_{(M+1), \dots, b_M}\}]$ is the orthogonal complement of $U = \text{Span}[\{b_1, \dots, b_M\}]$

Finding the Basis of the Principal Subspace

To determine the basis vectors, we rephrase the loss function: $\tilde{x}_n = \sum_{m=1}^M z_{mn} b_m = \sum_{m=1}^M (x_n^T b_m) b_m$ using symmetry of the dot product,

Since we can write x_n as a linear combination of all basis vectors, it holds!

$$\tilde{x}_n = \sum_{m=1}^M z_{mn} b_m = \sum_{j=1}^D (x_n^T b_j) b_j = \left(\sum_{d=1}^D b_d b_d^T \right) x_n = \left(\sum_{m=1}^M b_m b_m^T \right) x_n + \left(\sum_{j=M+1}^D b_j b_j^T \right) x_n,$$

We split the sum with M terms into a sum over M and a sum over $D-M$ terms. We find that the displacement vector $x_n - \tilde{x}_n$ is:

$x_n - \tilde{x}_n = \left(\sum_{j=M+1}^D b_j b_j^T \right) x_n = \sum_{j=M+1}^D (x_n^T b_j) b_j$ This means that the difference is exactly the projection of the data point onto the orthogonal complement of the principal subspace. The displacement vector $x_n - \tilde{x}_n$ lies in the subspace that is orthogonal to the principal subspace.

PCA finds the best rank- M approximation of the identity matrix

We can now reformulate the loss function: $J_m = \frac{1}{N} \sum_{n=1}^N \|(\mathbf{x}_n - \mathbf{c}_m)\|^2 = \frac{1}{N} \sum_{n=1}^N \left\| \left(\sum_{j=M+1}^D (\mathbf{b}_j^\top \mathbf{x}_n) \mathbf{b}_j \right) \right\|^2.$

$$J_m = \frac{1}{N} \sum_{n=1}^N \sum_{j=M+1}^D (\mathbf{b}_j^\top \mathbf{x}_n)^2 = \frac{1}{N} \sum_{n=1}^N \sum_{j=M+1}^D \mathbf{b}_j^\top \mathbf{x}_n \mathbf{x}_n^\top \mathbf{b}_j = \mathbf{b}_j^\top \mathbf{x}_n \mathbf{x}_n^\top \mathbf{b}_j$$

$$\mathbf{b}_j^\top \mathbf{x}_n = \mathbf{x}_n^\top \mathbf{b}_j. \text{ (We swap sums and obtain:)} \quad J_m = \sum_{j=M+1}^D \mathbf{b}_j^\top \left(\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right) \mathbf{b}_j = \sum_{j=M+1}^D \mathbf{b}_j^\top S_{\text{obs}} = \sum_{j=M+1}^D \text{tr}(\mathbf{b}_j^\top S_{\text{obs}} \mathbf{b}_j) = \sum_{j=M+1}^D \text{tr}(S_{\text{obs}} \mathbf{b}_j \mathbf{b}_j^\top) =$$

$$\text{tr}\left(\sum_{j=M+1}^D \mathbf{b}_j \mathbf{b}_j^\top\right) S_{\text{obs}}$$

Minimizing the average squared reconstruction error is equivalent to minimizing the projection of the data covariance matrix onto the orthogonal complement of the principal subspace. Minimizing the average squared reconstruction error is equivalent to maximizing the variance of the projected data.

$J_m = \sum_{j=M+1}^D \lambda_j$ ← the average squared reconstruction error when projecting onto the M-dimensional principal subspace. where λ_j are the eigenvalues of the data covariance matrix. To minimize the error, we need to select the smallest $D-M$ eigenvalues, implying that their corresponding eigenvectors are the basis of the orthogonal complement of the principal subspace. This means that the basis of the principal subspace comprises the eigenvectors $\mathbf{b}_{M+1}, \dots, \mathbf{b}_D$ that are associated with the largest M eigenvalues of the data covariance matrix.

10.4 Eigenvector Computation

$$S = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top = \frac{1}{N} \mathbf{X} \mathbf{X}^\top, \quad \mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$$

\mathbf{X} is a $N \times D$ matrix, the transpose of the typical data matrix. To get the eigenvectors / values of S , we can follow two approaches:

- We perform an Eigendecomposition and compute the eigenvalues and eigenvectors of S directly.
- We use a Singular value decomposition. Since S is symmetric and factorizes into $\mathbf{X} \mathbf{X}^\top$, the eigenvalues of S are the squared singular values of \mathbf{X} .

SVD of \mathbf{X} : $\mathbf{X} = \underbrace{\mathbf{U}}_{D \times N} \underbrace{\Sigma}_{N \times N} \underbrace{\mathbf{V}^\top}_{N \times D}$ where $\mathbf{U} \in \mathbb{R}^{D \times D}$ and $\mathbf{V}^\top \in \mathbb{R}^{N \times N}$ are orthogonal matrices and $\Sigma \in \mathbb{R}^{D \times N}$ is a matrix whose only nonzero entries are singular values $\sigma_{i,i} \geq 0$. It then follows that

$$S = \frac{1}{N} \mathbf{X} \mathbf{X}^\top = \frac{1}{N} \mathbf{U} \Sigma \Sigma^\top \mathbf{V}^\top \mathbf{V} = \frac{1}{N} \mathbf{U} \Sigma^2 \mathbf{V}^\top \quad \text{we get that the columns of } \mathbf{U} \text{ are the eigenvectors of } \mathbf{X} \mathbf{X}^\top \text{ (and } S\text{). The eigenvalues } \lambda_d \text{ of } S \text{ are related to the singular values:}$$

$$\lambda_d = \frac{\sigma_d^2}{N}$$

This relationship between the eigenvalues of S and the singular values of \mathbf{X} provides the connection between the maximum variance view and the singular value decomposition.

10.6: Key Steps of PCA in Practice

1) Mean Subtraction: We start by centering the data by computing the mean μ of the dataset and subtracting it from every single data point. This ensures that the dataset has mean 0. Mean subtraction is not strictly necessary but reduces the risk of numerical problems.

2) Standardization: Divide the data points by the standard deviation σ_d of the dataset for every dimension $d = 1, \dots, D$.

New feature is unit-free, and has its variance along each axis.

3) Eigen Decomposition of the Covariance Matrix: Compute the data covariance matrix and its eigenvalues and corresponding eigenvectors. Since the covariance matrix is symmetric, the spectral theorem states that we can find an ONS of eigenvectors. The eigenvectors are scaled by the magnitude of the corresponding eigenvalue. The longer vector spans the principal subspace, which we denote by V .

4) Projection: We can project any datapoint $x \in \mathbb{R}^D$ onto the principal subspace. To get this right, we need to standardize x using the mean μ_d and standard deviation of the training data in the d^{th} dimension, respectively, so that

$$x_{\perp}^{(d)} \leftarrow \frac{x_{\perp}^{(d)} - \mu_d}{\sigma_d}, d = 1, \dots, D, \text{ where } x_{\perp}^{(d)} \text{ is the } d^{\text{th}} \text{ component of } x_{\perp}. \text{ We obtain the projection as}$$

$\tilde{x}_{\perp} = BB^T x_{\perp}$ with unit $z_{\perp} = B^T x_{\perp}$ with respect to the basis of the principal subspace. Here, B is the matrix that contains the eigenvectors that are associated with the largest eigenvalues of the data covariance matrix as columns.

M: T Linalg Singular Value Decomposition

$$A = U \sum V^T / \sum \text{Diagonal } U, V, \text{ orthogonal}$$

$$AV = U \Sigma \quad A = \begin{bmatrix} u_1 & u_2 & \dots & u_r \end{bmatrix}$$

Sym pos def

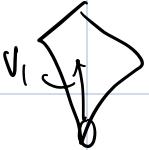
$$A = Q \Lambda Q^T$$

$$\mathbb{R}^n$$

Row space

$$\mathbb{R}^m$$

Column space



$$u_1 = Av_1$$

$$A = [v_1 \ v_2 \ \dots \ v_r] \quad \text{Row basis vectors}$$

$$= [u_1 \ u_2 \ \dots \ u_r] \begin{bmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_r \end{bmatrix} \quad \text{Diagonal}$$

columns