# LINEAR ALGEBRA

# MATRIX INVERSE → denoted by $A^{-1}$ and defined as $A^{-1}A = I$

# To solve system of linear $=ns$ →
$$Ax = b$$
$$A^{-1}Ax = A^{-1}b$$
$$Ix = A^{-1}b$$

# LINEAR <u>DEPENDENCE AND</u> SPAN →

* For $A^{-1}$ to exist $Ax = b$ must have exactly one solution for every value of $b$.

* It can have no solution, one solution or infinitely many solution. It can not have more than one but less than $\infty$ solution because if $x$ and $y$ are solutions then so is $z = \alpha x + (1-\alpha)y$ for any real $\alpha$.

* Essentially $b$ is linear combination of columns of $A$ i.e. we are trying to find if $b$ can be formed using linear combination of columns of $A$.
$$\sum_i x_i \cdot A_{:i} = b$$

* SPAN → set of all points obtainable by linear combination of original vectors.

    ⇒ To find solution of $Ax = b$ means whether $b$ lies in column span of $A$ or not.

    ⇒ For $b \in R^m$ to lie in column span of $A$, the column span must be $R^m$.

- LINEAR INDEPENDENCE → A set of vectors is linearly independent if no vector is linear combination of a subset of vectors.

  ⇒ For column span of A to be $R^m$ there must exist m linearly independent vectors is for $A_{m \times n}$, $n \geq m$.

  Also to have at most one solution $n = m$, making A a square matrix.

- SINGULAR MATRIX → A square matrix with linearly dependent columns.

# NORMS →

* $L^p$ norm ⇒ $\|x\|_p = \left( \sum_i |x_i|^p \right)^{1/p}$

* Popular norms →

  ① $L^2$ norm → $\|x\|_2 = \left( \sum_i |x_i|^2 \right)^{1/2}$ → Often we use squared $L^2$ norm.

  ② $L^1$ norm → $\|x\|_1 = \sum_i |x_i|$

  ③ $L^\infty$ norm → $\|x\|_\infty = \max_i |x_i|$ → As p increases $x_i$ with max value will dominate the sum. (just an intuition)

  ③ Forbenius norm → Norm of matrix
    $$\|A\|_F = \sqrt{\sum_{i,j} A_{ij}^2}$$

* Unit Vector → A vector with unit $L_2$ norm is $\|x\|_2 = 1$
* Orthogonal vectors → Two vectors x and y are orthogonal if
  $$x^T y = 0$$
* Orthnormal vectors → if $x^T y = 0$ and $\|x\|_2 = 1$ and $\|y\|_2 = 1$

* Orthogonal matrix → Matrix with orthogonal rows i

$$A^T A = A A^T = I$$

$$\Rightarrow \quad A^{-1} = A^T$$

# EIGEN DECOMPOSITION →

* Just like integers can be decomposed into prime factors, matrices can be decomposed into other matrices that reveal properties.
* An eigenvector of a square matrix A is a non zero vector $v$ such that multiplication by A alters only scale of $v$.

$$i \quad Av = \lambda v$$

eigen value

eigen vector

* If $v$ is eigenvector of A then so is $s \cdot v$ for $s \in R, s \neq 0$ with same eigen value.

* Eigendecomposition of A →

$$A = V \, diag(\lambda) \, V^{-1}$$

$$V \to [v^{(1)}, v^{(2)}, \ldots, v^{(n)}]$$

$$\lambda = [\lambda_1, \lambda_2, \ldots, \lambda_n]$$

* Every real symmetric matrix A →

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

$Q \to$ Orthogonal matrix of eigenvectors of A
$\Lambda \to$ diagonal matrix of eigenvalues of A.

$\Rightarrow$ A scales the spaces by $\lambda_i$ in the direction $v^{(i)}$.
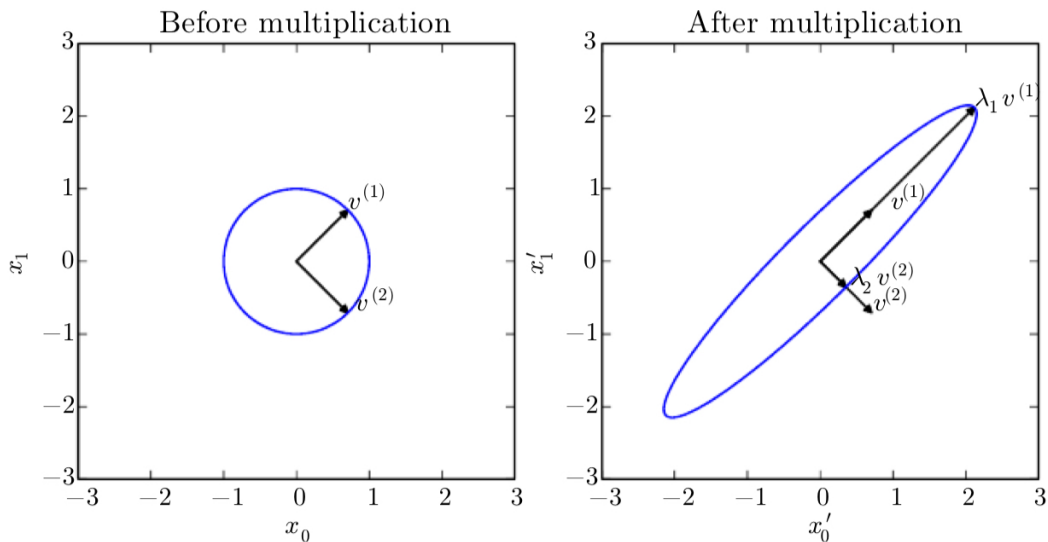
Figure 2.3: An example of the effect of eigenvectors and eigenvalues. Here, we have a matrix $A$ with two orthonormal eigenvectors, $v^{(1)}$ with eigenvalue $\lambda_1$ and $v^{(2)}$ with eigenvalue $\lambda_2$. *(Left)*We plot the set of all unit vectors $u \in \mathbb{R}^2$ as a unit circle. *(Right)*We plot the set of all points $Au$. By observing the way that $A$ distorts the unit circle, we can see that it scales space in direction $v^{(i)}$ by $\lambda_i$.

* A matrix is singular iff any of the eigenvalues are 0.

* A matrix whose eigenvalues are all +ve → Positive Definite
  A matrix whose eigenvalues are all −ve → Negative Definite
  If positive or 0 → Positive semi definite
  If negative or 0 → Negative semi definite

* For positive semi definite matrix →
  $$\forall x \quad x^T A x \geqslant 0$$
  ⎤→ we can use this property to prove if a matrix A is +ve semi definite or not.

# SINGULAR VALUE DECOMPOSITION →

* SVD decomposes a matrix into singular values and vectors.
* $A = UDV^T$

$U \rightarrow m \times m$ (Left singular vectors)

$D \rightarrow m \times n$ (Diagonal of D are Singular values of A)

$V \rightarrow n \times n$ (Right singular vectors)

# MOORE - PENROSE PSEUDOINVERSE →

If A is not invertible then we can compute its pseudo inverse using SVD as

$$A^+ = VD^+U^T$$

$D^+ = $ reciprocal of elements of diagonal matrix $D$

# TRACE AND DETERMINANT →

* Trace → $Tr(A) = \sum_i A_{i,i}$

$$\|A\|_F = \sqrt{Tr(AA^T)}$$

$$Tr(A) = Tr(A^T)$$

$$Tr(ABC) = Tr(CAB) = Tr(BCA)$$

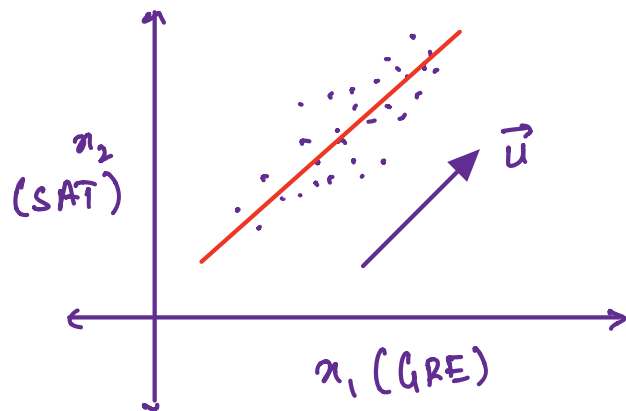More generally, $Tr\left( \prod_{i=1}^{n} F^{(i)} \right) = Tr\left( F^{(n)} \cdot \prod_{i=1}^{n-1} F^{(i)} \right)$

$$\Rightarrow Tr(AB) = Tr(BA)$$

* Determinant → Product of all eigen values of A. For singular matrix, $\det(A) = 0$

# PRINCIPAL COMPONENT ANALYSIS

* Say we have data $\{x^{(i)}\}_{i=1}^{m}$

▷ There is high correlation b/w $x_1$ and $x_2$, i.e. even if data is in 2D it is more or less in 1D.



* Underlying question is can we do dimensionality reduction i.e. we project $x^{(i)}$ on $\vec{u}$ to get $\{z^{(i)}\}_{i=1}^{m}$ where $z^{(i)} = \vec{u}^T x^{(i)}$ such that most of the "information" in data is still captured

* Unless the data is perfectly correlated there will be some loss.

* We want to find best set of $\{u_1, \cdots u_{12}\}$

Q Why project the data?
  ① Discovering hidden patterns in the data (correlation etc)
  ② Projecting onto a lower dimensional space makes things tractable, for $x^{(i)} \in R^n$ $K << n$.
  ③ Helps in reducing noise in data.

# How to quantify loss of information → Variance

$$Var(\{x^{(i)}\}_{i=1}^{m}) = \sum_{j=1}^{n}\left[\frac{1}{m}\sum_{i=1}^{m}(x_j^{(i)} - \mu_j^{(i)})^2\right]$$

Let $\{z^{(i)}\}_{i=1}^{m} \to$ projection over $\{u_1, \cdots u_{12}\}$ then $Var(\{z^{(i)}\}_{i=1}^{m})$ should

be close to variance of original data.

Problem Statement → Given $\{x^{(i)}\}_{i=1}^{m}$ $x^{(i)} \in R^n$.

find $(u_1,...,u_k)$ $u_l \in R^n$ $\forall l$, $u_l^T u_l = 1$

s.t.

$x^{(i)} \sim z^{(i)} \longrightarrow$ projection of $x^{(i)}$ on $\{u_l\}$

$\Rightarrow z_l^{(i)} = x^{(i)T} u_l$

s.t $var\left(\{z^{(i)}\}_{i=1}^{m}\right)$ is maximized

(1) We first normalize the points so that resulting data has 0 mean and unit variance. (This is so that all features are on same "scale" e.g. height and weight are on same scale)

$$x_j^{(i)} \longleftarrow \left(\frac{x_j^{(i)} - \mu_j}{s_j}\right)$$

mean, $\mu_j = \frac{1}{m} \sum_{i=1}^{m} x_j^{(i)}$ $\forall_j$ $= 0$ $\Bigg]$ After transformation.

variance, $var = \frac{1}{m} \sum_{i=1}^{m} x_j^{(i)2} = 1$

mean of $z$ along dimension $l$, $\mu_l^z = \frac{1}{m} \sum_{i=1}^{m} x^{(i)T} u_l$

$= \left[\frac{1}{m} \sum_{i=1}^{m} x^{(i)}\right]^T u_l$

$= \mu^T u_l = 0 \cdot u_l = 0$

# $\text{Var}\left(\{z^{(i)}\}_{i=1}^{m}\right) = \sum_{l=1}^{K} \frac{1}{m} \sum_{i=1}^{m} \left(x^{(i)^T} u_l\right)^2 \quad (\because \text{mean is } 0)$

$$= \sum_{l=1}^{K} \frac{1}{m} \sum_{i=1}^{m} \left(u_l^T x^{(i)}\right)\left(x^{(i)^T} u_l\right)$$

$$= \sum_{l=1}^{K} u_l^T \left[\frac{1}{m} \sum_{i=1}^{m} x^{(i)} \cdot x^{(i)^T}\right] u_l$$

<span style="color:red">Emperical covariance matrix</span>

$$= \sum_{l=1}^{K} u_l^T \Sigma u_l \qquad \left(\text{Here } \Sigma_{jk} = \frac{1}{m} \sum_{i=1}^{m} x_j^{(i)} x_k^{(i)}\right)$$

emperical covariance $(x_j, x_k)$

Our objective: $\underset{u_1, \ldots, u_l}{\text{argmax}} \sum_{l=1}^{K} u_l^T \Sigma u_l$

subject to ① $\|u_l\|_2^2 = 1$

② $u_{l_1}^T u_{l_2} = 0 \quad \forall\, l_1 \neq l_2$

# Assume $K=1$, result then generalizes

Objective $\rightarrow \underset{u_l}{\max}\ u_l^T \Sigma u_l$ subject to $u_l^T u_l = 1$ $\left.\right\}$ constrained optimization problem

We will use Langragion to solve above problem.

$$\mathcal{L}(u_l, \lambda) = u_l^T \Sigma u_l + \lambda\left[1 - u_l^T u_l\right]$$

Primal → $\max\limits_{u_\lambda} \left( \min\limits_{\lambda} L(u_{\lambda}, \lambda) \right) = \max\limits_{u_\lambda} \min\limits_{\lambda} u_\lambda^T \mathcal{E} u_\lambda + \lambda[1 - u_\lambda^T u_\lambda]$

Dual → $\min\limits_{\lambda} \left( \max\limits_{u_\lambda} L(u_{\lambda}, \lambda) \right) = \min\limits_{\lambda} \max\limits_{u_\lambda} u_\lambda^T \mathcal{E} u_\lambda + \lambda[1 - u_\lambda^T u_\lambda]$

Maximizing wrt $u_\lambda$ gradient should vanish →

$$\nabla_{u_\lambda} \left[ u_\lambda^T \mathcal{E} u_\lambda - \lambda(1 + u_\lambda^T u_\lambda) \right] = 2\mathcal{E} u_\lambda - 2\lambda u_\lambda = 0$$

⇒ $\mathcal{E} u_\lambda = \lambda u_\lambda$

$n \times m$ matrix    $u_\lambda \in R^n$    → scalar

⇒ $\lambda$ is eigenvalue of $\mathcal{E}$ and $u_\lambda$ is the eigenvector

$$\underset{u_1, \ldots, u_K}{\text{argmax}} \sum\limits_{\lambda=1}^{K} u_\lambda^T \mathcal{E} u_\lambda \quad \text{constraints} → \quad u_\lambda^T u_\lambda = 1 \quad \forall \lambda$$
$$u_{\lambda_1}^T u_{\lambda_2}^T = 0 \quad \lambda_1 \neq \lambda_2$$

Solution to above problem :-
- $u_\lambda$ are eigen vectors of $\mathcal{E}$
- $\lambda_\lambda$ are eigen values of $\mathcal{E}$ that capture amount of variance among $u_\lambda$.

# Principal components are vectors corresponding to largest $K$ eigen values of $\mathcal{E}$

⇒ Objective reduces to find eigen values and eigen vectors of $\mathcal{E}$ and keep $u_1, \ldots, u_\lambda$ corresponding to top $K$ eigenvalues of $\mathcal{E}$.

$\Rightarrow \lambda_l$ is proportional to amount of variance captured in $u_l$.

# To find eigen value and eigen vectors $\rightarrow$ ① Eigen Decomposition
                                         ② Singular Value Decomposition.

# We can write $\Sigma = \frac{1}{m} X^T X$  where $X$ is the data matrix.

$\Rightarrow$ Problem reduces to finding Top $k$ eigenvalues of $\frac{1}{m} X^T X$

  Let $A = \frac{1}{m} X^T X$

  $$A = Q \Lambda Q^{-1} = Q \Lambda Q^T \qquad \text{as } Q \text{ are orthonormal}$$

  Challenges is that computing $A = \frac{1}{m} X^T X \rightarrow O(n^2 m)$
  and eigendecomposition will be $\rightarrow O(n^3)$
  if num of examples $n \gg m \rightarrow$ too expensive.

# Using SVD $\rightarrow$

  $$A = U D V^T \qquad\qquad U \text{ and } V \text{ are orthonormal}$$
  $\quad (m \times m) \; (m \times n) \; (n \times n)$

  • Columns of $V$ are eigenvectors of $A$
  • Entries of $D$ are square root of eigenvalues of $A$.

  Complexity of SVD $\rightarrow O(\min(m^2 n, n^2 m))$

**Q:** $X = UDV^T$ show that $X^TX$, show that columns of $V$ are eigenvector of $X^TX$ and Diagonal entries of $D^TD$ are eigenvalues of $X^TX$. ✓

$$X^TX = (UDV^T)^T (UDV^T)$$
$$= VD^TU^TUDV^T$$
$$= VD^TDV^T \qquad (\because U^TU = I)$$