# SVD Proofs:

**1** Let's find this $V, \Sigma_1, U$ assuming they always exist.

Assuming $AV = U\Sigma_1$, let's calculate what $U, \Sigma_1, V$ are.

**2** Prove that $V$ is the eigenbasis of $A^T$ (row space of A)

$$A = U\Sigma_1 V^T$$

$$A^T = V\Sigma_1^T U^T = V\Sigma_1 U^T$$

$$\therefore A^T A = (V\Sigma_1 U^T)(U\Sigma_1 V^T)$$

$$= V\Sigma_1 U^T U \Sigma_1 V^T$$

$$= V\Sigma_1^2 V^T = V\Sigma_1^{\nu} V^{-1}$$

or $A^T A \cdot V = V\Sigma_1^{\nu}$

$\therefore$ $V$ is the eigen vector matrix of $AA^T$

and $[\sigma_1, \sigma_2 \ldots]^T$ are the $\sqrt{\lambda_1} \sqrt{\lambda_2} \ldots$ of matrix $AA^T$.

$A_{m \times n} \Rightarrow A^T A \in n \times n$ $\therefore$ $V \equiv n \times n$

**3** Prove that $U$ is the eigenbasis of $A$ (col. space of A).

Now, how to find $U$?

$$AA^T = (U\Sigma V^T)(V\Sigma U^T) = U\Sigma_1 V^T V \Sigma_1 U^T$$

$$= U\Sigma^2 U^T = U\Sigma_1^{\nu} U^{-1}$$

$$\Rightarrow AA^T U = U\Sigma_1^{\nu}$$

$\hookrightarrow$ Eigenvector of $AA^T$.

$U \equiv m \times m$

**4** Prove that $U$ and $V$ are both orthogonal.

**2.** Prove that matrix $A$ always has the SVD decomposition

$$A^T A \cdot V = \lambda V \qquad \longrightarrow \text{ always true, } \lambda \geq 0 \text{ and } V \text{ is } \perp \text{ since } A^T A \text{ is PSD.}$$

$$A^T \left( \frac{AV}{\sqrt{\lambda}} \right) = \frac{\lambda}{\sqrt{\lambda}} \cdot V$$

Now $\quad A A^T \left( \frac{AV}{\sqrt{\lambda}} \right) = A \left( \frac{\lambda}{\sqrt{\lambda}} V \right) = \lambda \left( \frac{AV}{\sqrt{\lambda}} \right) \quad \longrightarrow$ This is the eigenvector eq$^n$ for $AA^T$.

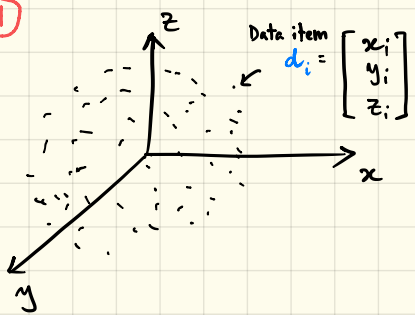$\therefore$ The matrix $\left( \frac{AV}{\sqrt{\lambda}} \right)$ must be orthonormal, since $AA^T$ is PSD.

Let $\quad U = \frac{AV}{\sqrt{\lambda}} \quad$ where $U$ is orthonormal. $\therefore AV = U\sqrt{\lambda}$

$$\therefore \quad A = U\sqrt{\lambda} V^{-1} = U \Sigma V^T \quad /\!/$$

# PRINCIPAL COMPONENT ANALYSIS (PCA)

# PCA : Principal Component Analysis

**① **



Data item $d_i$ =
$\begin{bmatrix} x_i \\ y_i \\ z_i \end{bmatrix}$

**②** Prerequisite: Matrix $A \equiv$
$\begin{bmatrix} d_1 & d_2 & & d_n \\ x_1 & x_2 & \cdots & x_n \\ y_1 & y_2 & & y_n \\ z_1 & z_2 & & z_n \end{bmatrix}$

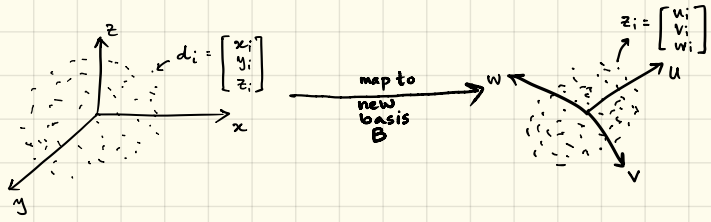Covariance $(A) = AA^T = \begin{bmatrix} x_1 & x_2 & x_n \\ y_1 & y_2 & y_n \\ z_1 & z_2 & z_n \end{bmatrix} \begin{bmatrix} x_1 & y_1 & z_1 \\ x_2 & y_2 & z_2 \\ & & \\ x_n & y_n & z_n \end{bmatrix}$

$= \begin{bmatrix} Var(x) & Cov(xy) & Cov(xz) \\ Cov(yx) & Var(y) & Cov(yz) \\ Cov(zx) & Cov(zy) & Var(z) \end{bmatrix}$

## PCA's Goal : Which basis $B$ will make the data uncorrelated?

Ans : Let's represent data in another orthogonal basis $B$.

**③**



$d_i = \begin{bmatrix} x_i \\ y_i \\ z_i \end{bmatrix}$

map to new basis $B$

$z_i = \begin{bmatrix} u_i \\ v_i \\ w_i \end{bmatrix}$

Note, when $B$ is orthogonal it can be easily made orthonormal.

**④** When data $d_i$ is represented in this new basis $B$, it becomes, say, $z_i$.
Note: If $B$ is a Fourier basis, then $z_i$ is the fourier transform.

$D = \begin{bmatrix} x_1 & x_2 & & x_N \\ y_1 & y_2 & \cdots & y_N \\ z_1 & z_2 & & z_N \end{bmatrix} \equiv \begin{bmatrix} d_1 & d_2 & \cdots & d_N \end{bmatrix}$

So, $\begin{bmatrix} | & | & | \\ u & v & w \\ | & | & | \end{bmatrix} \begin{bmatrix} u_1 & u_2 & & u_n \\ v_1 & v_2 & \cdots & v_n \\ w_1 & w_2 & & w_n \end{bmatrix} = \begin{bmatrix} | & | & & | \\ d_1 & d_2 & \cdots & d_n \\ | & | & & | \end{bmatrix}$

$\therefore \quad B \cdot Z = D$

Now, to be uncorrelated, covariance of data (in new basis) should be a diagonal matrix ( because uncorrelated means $Cov(x,y) = 0$ )

**⑤** Now, data covariance (in new basis) $= ZZ^T$

$ZZ^T = \begin{bmatrix} u_1 & u_2 & & u_n \\ v_1 & v_2 & \cdots & v_n \\ w_1 & w_2 & & w_n \end{bmatrix} \begin{bmatrix} u_1 & v_1 & w_1 \\ u_2 & v_2 & w_2 \\ \vdots & & \\ u_n & v_n & w_n \end{bmatrix} = \Lambda$

$(B^{-1}D)(B^{-1}D)^T = \Lambda$

$B^{-1} D \cdot D^T (B^{-1})^T = \Lambda$

$D \cdot D^T (B^{-1})^T = B\Lambda$

$D \cdot D^T (B^T)^T = B\Lambda \qquad \ldots \therefore B^{-1} = B^T$
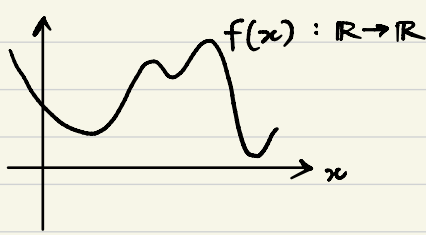
$DD^T B = B\Lambda$

$\therefore B$ is eigenvector

Thus, the eigen vectors of the data covariance matrix gives us the desired basis vectors to decorrelate the data.

Now, to compress data D, basically remove the last K columns of B and last K rows of Z, then take the product of the matrices $B'Z' = D'$.

· This D' is the compressed matrix.

# 434 : Optimization Basics

⊕ $f(x) : \mathbb{R} \to \mathbb{R}$    $\dfrac{\partial f(x)}{\partial x} = 0$    gives us **local extremum**

How do you know maxima or minima?

$$\left.\dfrac{\partial^2 f(x)}{\partial x^2}\right|_{x=x^*} > 0 \quad \text{i.e.} \quad \dfrac{\partial^2 f(x^*)}{\partial x^2} > 0$$

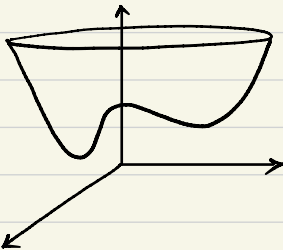⊕ Functions in higher dimensions (i.e., when $\bar{x}$ is vector) ⟹    $f : \mathbb{R}^n \to \mathbb{R}$

$$\nabla f(x) \equiv \nabla f_x = \begin{bmatrix} \partial f / \partial x_1 \\ \partial f / \partial x_2 \\ \vdots \\ \partial f / \partial x_n \end{bmatrix}$$

$\underset{\text{"nabla" or "del"}}{\Big\downarrow}$

$$\nabla^2 f(x) = \begin{bmatrix} \dfrac{\partial^2 f}{\partial x_1^2} & \dfrac{\partial^2 f}{\partial x_1 x_2} & \cdots & \dfrac{\partial^2 f}{\partial x_1 x_n} \\ \dfrac{\partial^2 f}{\partial x_2 x_1} & \dfrac{\partial^2 f}{\partial x_2^2} & \cdots & \dfrac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & & & \\ \dfrac{\partial^2 f}{\partial x_n x_1} & \dfrac{\partial^2 f}{\partial x_n x_2} & \cdots & \dfrac{\partial^2 f}{\partial x_n^2} \end{bmatrix}$$

$\underbrace{\hspace{6cm}}$
called the "**Hessian**" matrix

⊕ How do we find **local** maxima / minima of such functions of vectors?

$\nabla f_x = 0$    ⟹ gives extremums

$\nabla^2 f_{x^*} > 0$    ⟹ indicates **minima**

$\underbrace{\hspace{4cm}}$
Hessian is a positive definite matrix.

Matrix A is P.D. when all $\lambda_i(A) > 0$    or $x^T A x > 0, \forall x$

Positive semi definite (PSD)    when $\lambda_i(A) \geqslant 0$ , $x^T A x \geqslant 0, \forall x$

⊙ Note: $\nabla^2 f_{x^*} \geqslant 0$ is a necessary but not sufficient condition

Example: $f(x) = x^3$

$\nabla f_x = 3x^2 = 0 \implies x^* = 0$

But is $x^*$ a minima or maxima or neither?

$\nabla^2 f(x^*) = 6x \big|_{x=0} = 0$

But observe that $x^* = 0$ is neither a minima or maxima.



$x^* = 0$ is NOT maxima or minima.

called "stationary" points

$\nabla^2 f_x > 0$ is sufficient condition

⊙ $\nabla f_x = 0$ and $\nabla^2 f_x > 0$ gives us local minima.
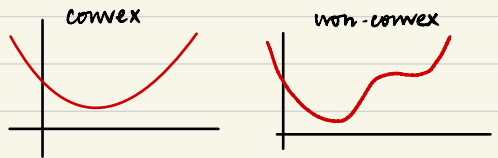But how can I get global minima?

↓

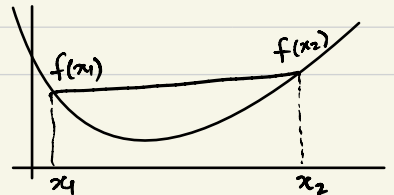Well, if $f_x$ is a convex fn, then local minima is global minima.

↓

⊙ What's a convex fn?
↳ Functions that have an upward curvature everywhere.



convex          non-convex

Intuitively: The straight line joining any two points $f(x_1)$ and $f(x_2)$ always lies above $f(y)$, where $y \in [x_1, x_2]$



$f(x_1)$          $f(x_2)$
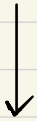
$x_1$          $x_2$

Mathematically : $\alpha f(x_1) + (1-\alpha) f(x_2) \geqslant f(\alpha x_1 + (1-\alpha) x_2), \quad \alpha \in [0,1]$

How to test for convexity ? $\qquad \nabla^2 f_x \geqslant 0 \qquad \Longleftrightarrow \quad$ convex fn.s.

→ Summary : Given $f(x)$,
if $\nabla^2 f(x) \geqslant 0$ (i.e., Positive semi-def Hessian)
then $f(x)$ is convex fn.
Thus $f(x) = 0$ gives GLOBAL MINIMA.

→ But here is the bad news :
↳ Even if $f(x)$ is convex, in many cases, its difficult
to solve for $\nabla f(x) = 0$.
Example : $f(x) = e^x + x^2$
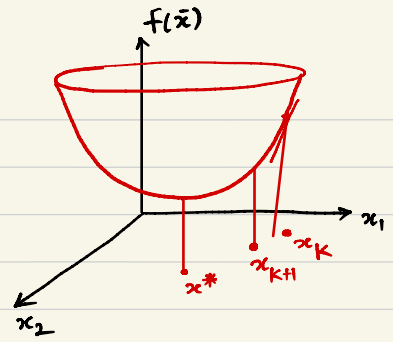↳ Closed form solution difficult

We need to solve such functions iteratively
↳ Motivates gradient descent

→ Main idea : We want to start at some $x = x_0$
Move $x_0 \rightarrow x_1 \rightarrow x_2 \cdots \rightarrow x^*$
s.t. $x^*$ is local/global minima of $f(x)$

This implies : $f(x_{k+1}) < f(x_k)$
So from $x_k$, we should go along a direction that
decreases the value of $f(x_k)$.
↳ Say this direction is $\vec{v}_k$

$$\therefore \quad \vec{x}_{K+1} = \vec{x}_k + \vec{v}_k$$



→ What $\vec{v}$ direction will take us **most downward**?

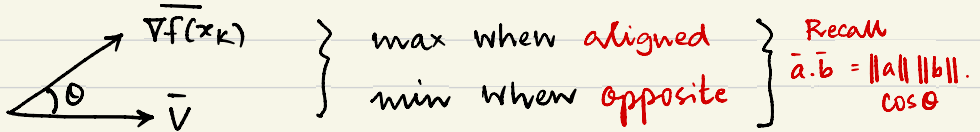Answer: The direction of $-\nabla f(x_k)$.

**Proof :** Taylor's 1st order expansion says

$$f(y) = f(x) + \nabla f(x)^T (y-x) + O(|y-x|)$$

$$\therefore \quad f(x_k + \varepsilon \bar{v}) = f(x_k) + \varepsilon \cdot \nabla f(x_k)^T v_k + O(\varepsilon)$$

$$\underbrace{\lim_{\varepsilon \to 0} \frac{f(x_k + \varepsilon \bar{v}) - f(x_k)}{\varepsilon}}_{} = \nabla f(x_k)^T v_k$$

Rate of change of $f(x)$ along **direction $v_k$**

So what is the **max and min** value of $\nabla f(x_k)^T v_k$ ?



max when **aligned**
min when **opposite**

Recall
$\bar{a} \cdot \bar{b} = \|a\| \|b\| \cdot \cos\theta$

By **Cauchy-Schwarz** inequality

$$-\|\nabla f(x_k)\| \|v\| \leqslant \nabla f(x_k)^T v_k \leqslant \|\nabla f(x_k)\| \|v\|$$

$$\therefore \quad \text{Maximal downward direction} = -\nabla f(x_k) \qquad \blacksquare$$

→ Thus : $\quad x_{K+1} = x_k + v_k = x_k - \alpha \nabla f(x_k)$

↑ Step size.

This is called "**steepest gradient descent (SGD)**"

→ Steepest Grad. Descent Algorithm :

① $K = 0$ ; $\alpha =$ small positive value ;
   $\varepsilon =$ very small value

② $x[K] =$ random vector

③ Calculate $\nabla f(x[K])$

④ $x[K+1] = x[K] - \alpha \nabla f(x[K])$

⑤ if $f(x[K+1]) - f(x[K]) < \varepsilon$ then terminate

⑥ $k++$

⑦ Goto ③

Questions :

(a) Why does step size $\alpha$ need to be small ?

(b) Can you draw a case where SGD may not converge if $\alpha$ is not small enough?

(c) Does SGD take the shortest path from $x_0$ to $x^*$ ?