

Álgebra Lineal: Optimización + Eigenvectores

ITAM

Clase 7 Curso Propedéutico
2017/06/20

Optimización...

- Por escribir....

$$\min/\max f(x)$$

$$\text{sa } g(x)=0$$

La solución debe cumplir

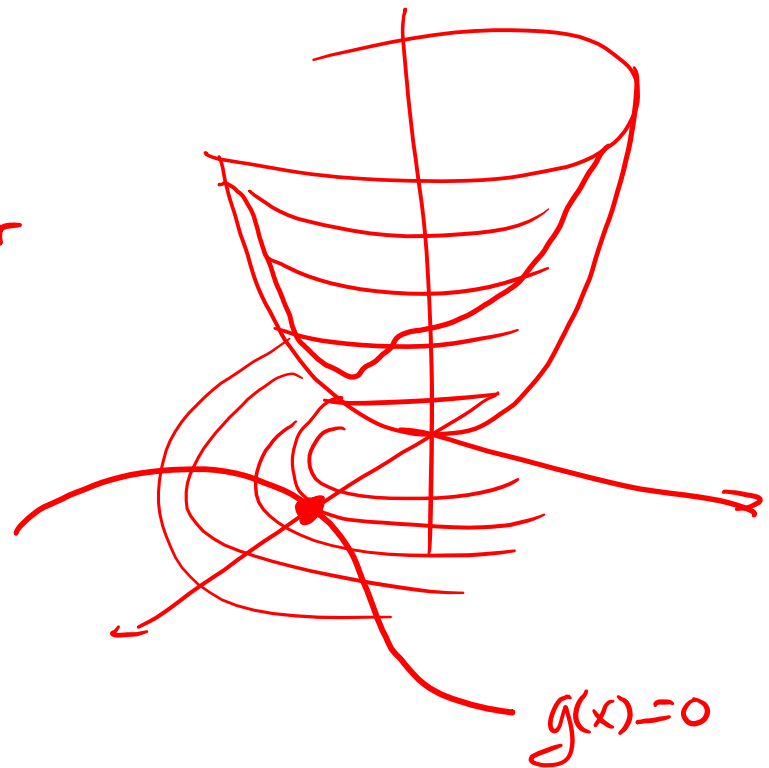
$$\nabla f(x) = \lambda \nabla g(x)$$

$$g(x)=0$$

o bien

$$\mathcal{L}(x, \lambda) = f(x) - \lambda g(x)$$

$$\nabla \mathcal{L}(x, \lambda) = 0$$



Derivación de funciones que usan matrices

2 casos especiales de derivación con matrices y vectores

$$1) \quad f(x) := b^T x = x^T b \quad x, b \in \mathbb{R}^n \quad f: \mathbb{R}^n \rightarrow \mathbb{R}$$

$$\nabla f(x) = \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix} = b \quad f(x) = x_1 b_1 + \dots + x_n b_n$$

$$2) \quad f(x) = x^T A x \quad \text{con } A \text{ simétrica} \quad A \in \mathbb{R}^{n \times n}$$

$$x \in \mathbb{R}^n$$

$$f(x) = \sum_{i=1}^n \sum_{j=1}^n x_i a_{ij} x_j$$

$$= a_{11}x_1^2 + \dots + a_{nn}x_n^2 + a_{12}x_1x_2 + a_{21}x_2x_1 + \dots$$

$$\stackrel{\text{simétrica}}{=} \sum_{i=1}^n a_{ii}x_i^2 + \sum_{i \neq j} a_{ij}x_i x_j$$

$$= \sum_{i=1}^n a_{ii}x_i^2 + 2 \sum_{i < j} a_{ij}x_i x_j$$

$$A = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nn} \end{bmatrix}$$

$$\begin{aligned}
 &= a_{11}x_1^2 + a_{12}x_1x_2 + \dots + a_{1n}x_1x_n \\
 &\quad + a_{21}x_2x_1 + \dots + a_{2n}x_2x_n \\
 &\quad \vdots \\
 &\quad + a_{n1}x_nx_1 + \dots + a_{nn}x_n^2
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial}{\partial x_k} (x^T A x) &= a_{k1}x_1 + \dots + a_{k,k-1}x_{k-1} + 2a_{kk}x_k + \dots + a_{kn}x_n \\
 &\quad + a_{1k}x_k + \dots + a_{k+1,k}x_{k+1} + \dots + a_{nk}x_n \\
 &= 2 \sum_{i=1}^n a_{ki} x_i = 2 x^T A_k
 \end{aligned}$$

entonces $\nabla(x^T A x) = 2 \begin{bmatrix} x^T A_1 \\ \vdots \\ x^T A_n \end{bmatrix} = 2 A x$

Aplicación: Puntos críticos de una función cuadrática $f(x) = x^T A x + b^T x + c$

$$\nabla f(x) = 0 \Rightarrow 2Ax + b = 0 \Rightarrow Ax = -\frac{b}{2} \Rightarrow x = -\frac{1}{2} A^{-1} b$$

Normas matriciales y optimización

$$\|A^{-1}u\|^2 =$$

$$\begin{aligned} u^T(A^T A u) &= \lambda u^T u = 1 \\ \Rightarrow \|A\| &= \sqrt{\lambda_{\max}} \\ &= u^T(A^T A)u \end{aligned}$$

$$\|A\| = \max_{\|u\|=1} \|Au\|$$

Problema equivalente

max
sa.

$$\|Au\|^2 - \|u\|^2 = 0$$

identidad

$$u^T u = u^T I_n u$$

Solución

$$L(x, \lambda) = u^T(A^T A)u - \lambda(u^T u - 1)$$

$$\nabla L(x, \lambda) = 0 \Leftrightarrow$$

$$2A^T A u = 2\lambda u$$

$A^T A$ es simétrica

$$\begin{aligned} (A^T A)^T &= (A^T)^T (A^T)^T \\ &= A^T A \end{aligned}$$

Conclusión: u^* es vector propio (?) de $A^T A$

$$(A^T A)u^* = \lambda^* u^* \leftarrow \text{solucionar sistema}$$

$$(AB)^T = B^T A^T$$

**Dos temas nuevos: formas
cuadráticas y optimización**

Eigenvectores

Dada $A \in R_{n \times n}$ *ojo: cuadrada*

Eigenvalor y eigenvector asociado: $Av = \lambda v$

Diagonalización

Si R^n tiene una base de eigenvectores de A entonces estos elementos forman un Sistema de coordenadas.

es decir, cualquier elemento $v \in R^n$ se puede escribir de forma única solo por ser base de eigenvectores

~~La matriz A es diagonal en ese Sistema de coordenadas.~~

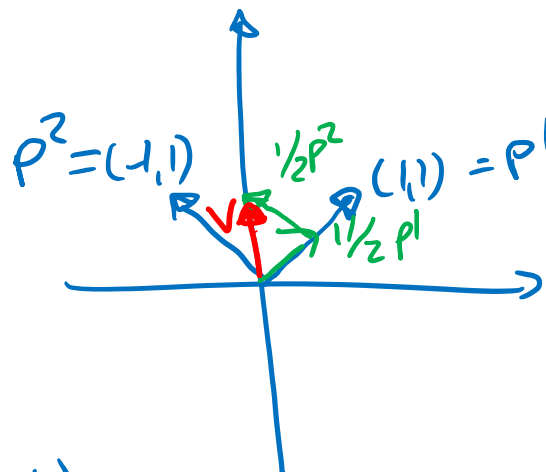
Además la base debe cumplir

$$AW = WD \quad \text{con } W = [w_1 | \dots | w_n] \quad D = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \ddots \\ 0 & 0 & \lambda_n \end{bmatrix}$$

$A = WDW^{-1}$

- **Observación:** De manera general, si P es una matriz invertible, sus columnas dan una base de R^n y $P^{-1}v$ nos dice como escribir a v en esas nuevas coordenadas.

$$P = \begin{pmatrix} 1 & -1 \\ 1 & 1 \end{pmatrix} \\ = [p^1 | p^2]$$



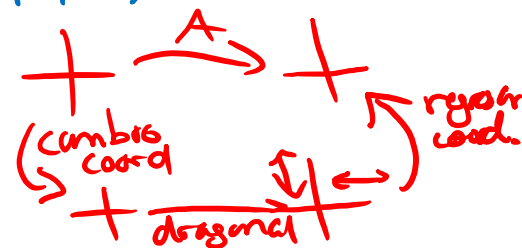
$$v = (0,1) = 0\vec{e}_1 + 1\vec{e}_2 \\ = 0(1,0) + 1(0,1)$$

Inverso: $P^{-1} = \frac{1}{2} \begin{pmatrix} 1 & 1 \\ -1 & 1 \end{pmatrix}$ $PP^{-1} = P^{-1}P = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$

$$P^{-1}v = \begin{pmatrix} 1/2 \\ 1/2 \end{pmatrix}$$

$$(1/2 \ 1/2)[p^1, p^2] = (0,1)_{\{e_1, e_2\}}$$

Conclusión Si $A = WDW^{-1}$ con D diagonal A es



Eigenvectores como problema de optimización!

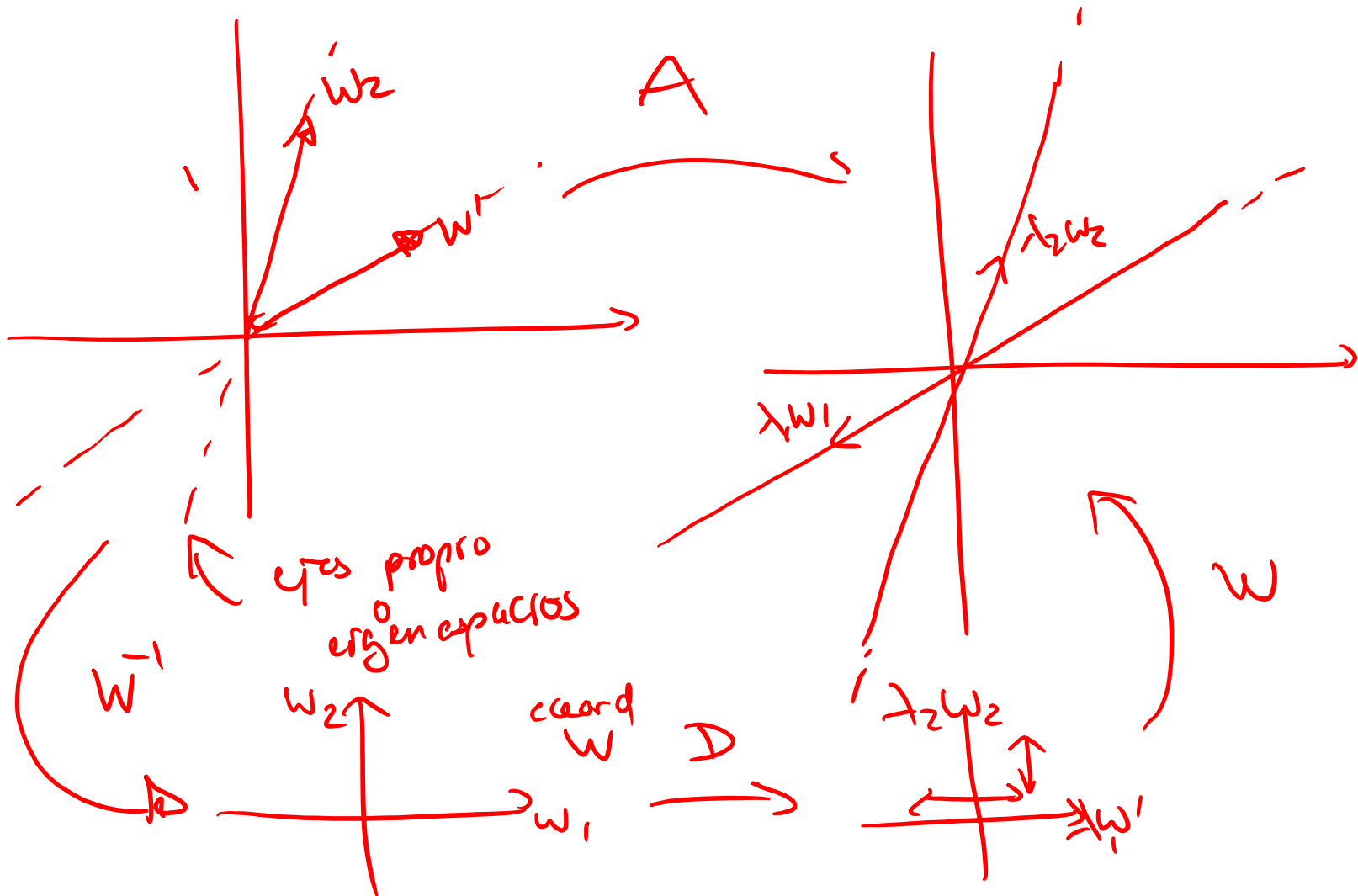
Ejercicio

Un eigenvector es la solución a un problema de optimización

$$\begin{aligned} \max \quad & w^T A w \\ \text{s.t.} \quad & \|w\| = 1 \end{aligned}$$

Sí, A es simétrica.

A matrix diagonalizable
 $A = W D W^{-1}$



Regresando al determinante

Recordar ① $\det(AB) = \det(A) \det(B)$ ② $\det(W^{-1}) = \frac{1}{\det(W)}$

$$\begin{aligned}\text{Si } A = W D W^{-1} \quad \det(A) &= \det(W) \det(D) \det(W^{-1}) \\ &= \frac{\det(W)}{\det(W)} \det(D) \\ &= \det(D) \\ &= \det \left(\begin{bmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_n \end{bmatrix} \right) \\ &= \prod \lambda_i\end{aligned}$$

El determinante de una matriz es el producto de sus eigenvalores

Obs práctica: Todas las matrices que van a ver en su vida son diagonalizables. De hecho, si permiten números complejos, todas las que existen lo son.

Algoritmos para cálculos de Eigenvectores

El método de la potencia devuelve el eigenvalores asociado al eigenvalor más grande (en valor absoluto).

Método de la Potencia

v_0 = (casi) cualquier semilla

For $k = 1 \dots$, convergencia

Definir $\lambda_{k+1} = v_k^\top A v_k$

Definir $v_{k+1} = \frac{A v_k}{\|A v_k\|}$

v_k converge al eigenvector más grande y λ_k el eigenvalor más grande

$$v_k \propto A^k v_0$$

¿Por qué sirve el método de la potencia?

Si w_1, \dots, w_n es una **base** de eigenvectores de $A = PDP^T$ entonces v_k se puede escribir como

$$v_k = \alpha_1 w_1 + \dots + \alpha_n w_n$$

donde $\alpha = P^{-1}v_k$ es las coordenadas con la base de eigenvectores.

Entonces:

Intuición

$$v_k = \frac{Av_{k-1}}{\|Av_{k-1}\|} \propto \dots A^k v_0 = A^k \alpha_1 w_1 + \dots + A^k \alpha_n w_n$$

$$= \alpha_1 \lambda_1^k w_1 + \dots + \alpha_n \lambda_n^k w_n$$

$$\frac{A^k v_0}{\|A^k v_0\|} = \frac{\cancel{\alpha_1 \lambda_1^k w_1} + \dots + \alpha_n \lambda_n^k w_n}{\| \cancel{\alpha_1 \lambda_1^k w_1} + \dots + \alpha_n \lambda_n^k w_n \|}$$

Algoritmos para cálculos de Eigenvectores

Algoritmo QR (el más usado)

Utiliza un método conocido como **descomposición QR** (lo vamos a ver más adelante) que escribe

$A = QR$ con Q ortogonal y R triangular superior.

ALGORITMO QR

$A_0 = A$

For $k = 1 \dots$, convergencia

Factorizar $A_k = Q_k R_k$

Definir $A_{k+1} = R_k Q_k$

La matriz Q_k tiene los eigenvectores y R_k los eigenvalores en la diagonal.

Cuando se demuestra la convergencia de este algoritmo (algo que no haremos en este curso...) se ve que detrás hay un mecanismo muy similar al método de la potencia. Normalmente este método se usa en conjunto con otros métodos que aceleran su convergencia como las reflexiones de Householder

Matrices simétricas y formas cuadráticas

- Una de las “aplicaciones” de los eigenvectores es que permite estudiar de manera muy sencilla estudiar las formas cuadráticas.

Una forma cuadrática es una **función** $f: R^n \rightarrow R$ de la forma

$$f(x) = f(x_1, \dots, x_n) = x^T A x + b^T x + c$$

Con A una matriz simétrica

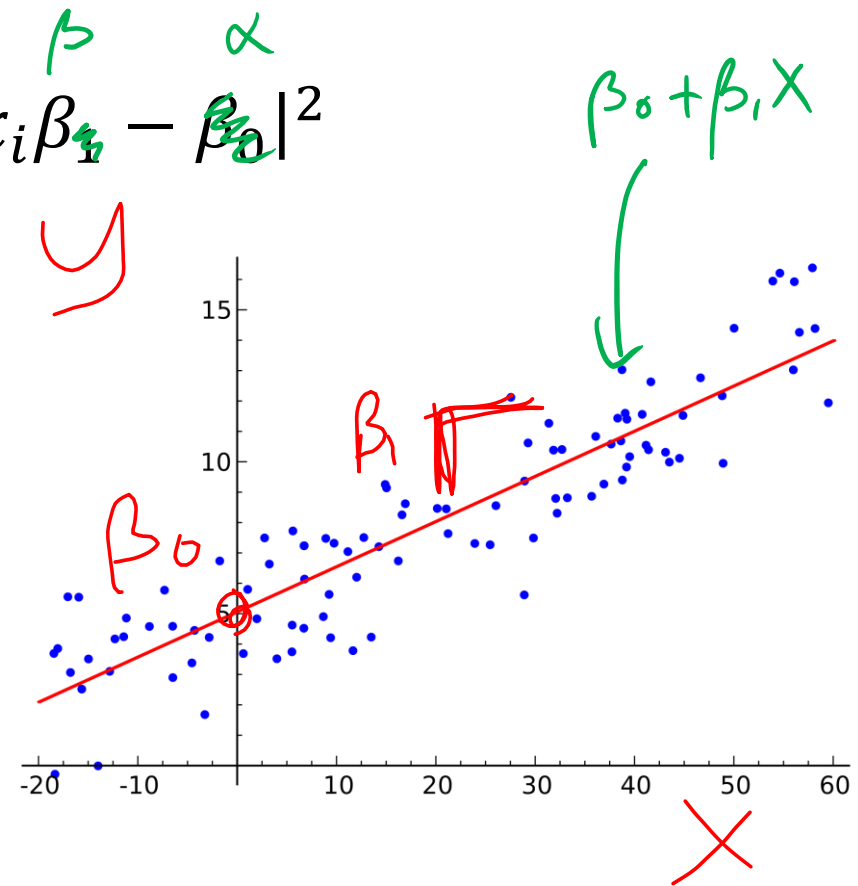
- Son funciones “polinomiales” de orden dos.
- Ej: $A = \begin{bmatrix} 3 & 1 \\ 1 & 2 \end{bmatrix}$, $b = \begin{bmatrix} 3 \\ 2 \end{bmatrix}$, $c = -1$

- Antes de estudiar la teoría de formas cuadráticas vamos a ver algunos problemas bonitos que surgen con formas cuadráticas (algunos los resolveremos aquí).

Ejemplos:

- *Regresión Lineal*: Queremos explicar una variable y con una variable x

$$\min_{\beta_0, \beta_1} \sum_i |y_i - x_i \beta_1 - \beta_0|^2$$

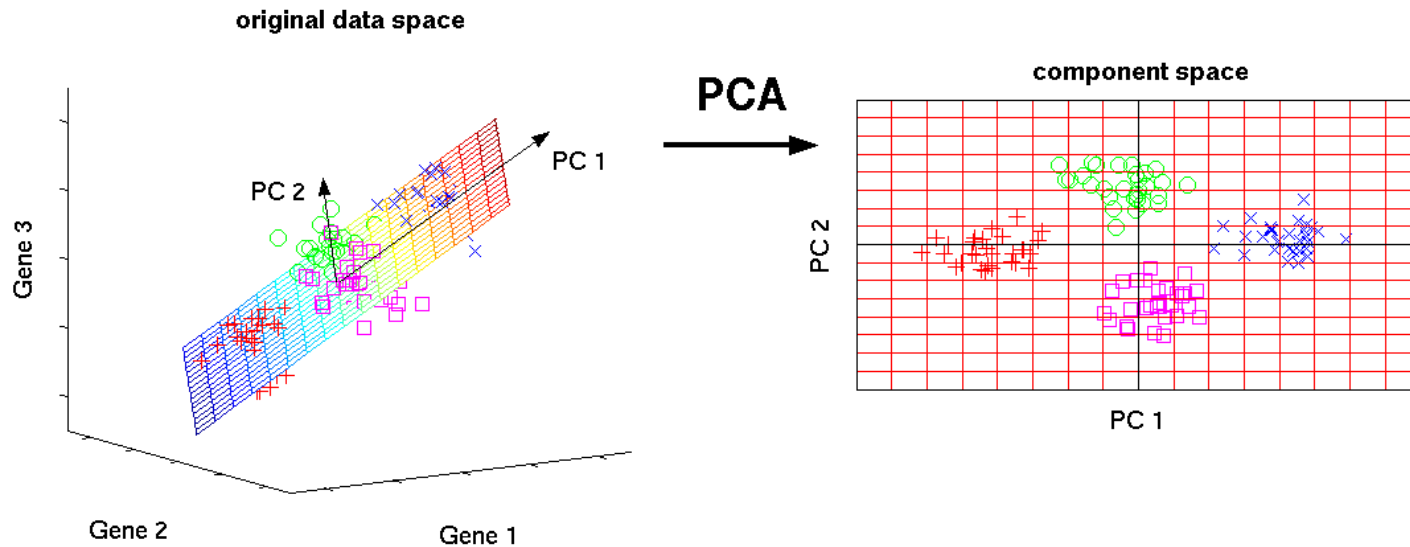


- Reducción de dimensionalidad (PCA) y clasificación lineal

Encontrar una combinación lineal de vectores X^1, \dots, X^p que “mejor” represente a las variables.

$$\max_{\|v\|=1} \|Xv\|^2$$

$$Xv = v_1 X^1 + \dots + v_p X^p$$

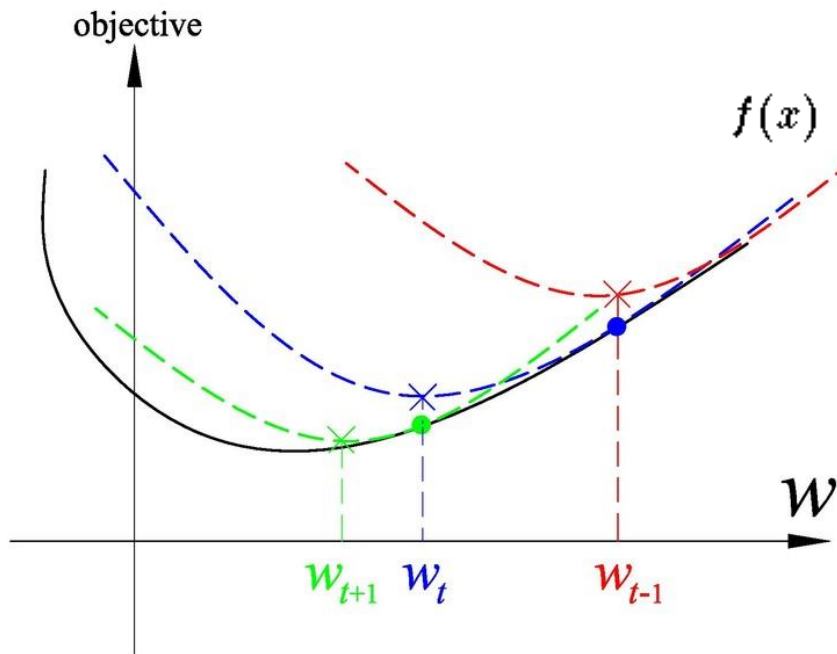


Optimización en General

- Cualquier función puede aproximarse por una función cuadrática. Este es el corazón del famosísimo método de Newton-Raphson que vamos a ver.

$$\min_x f(x)$$

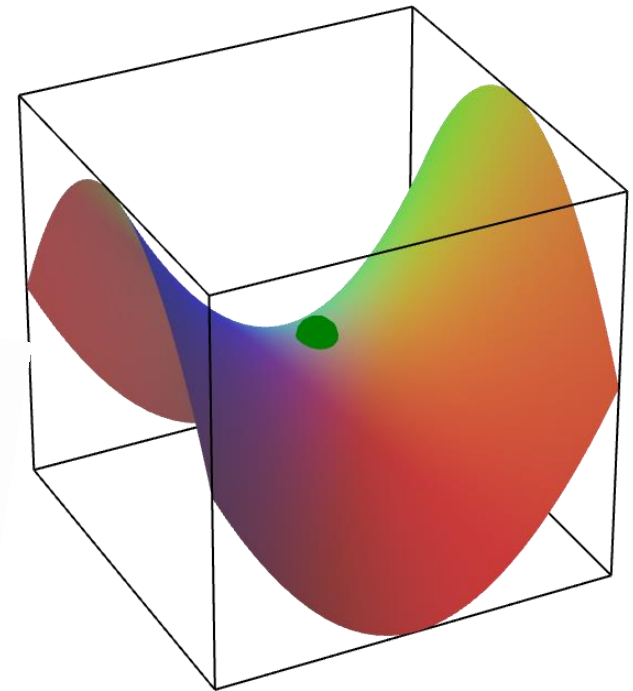
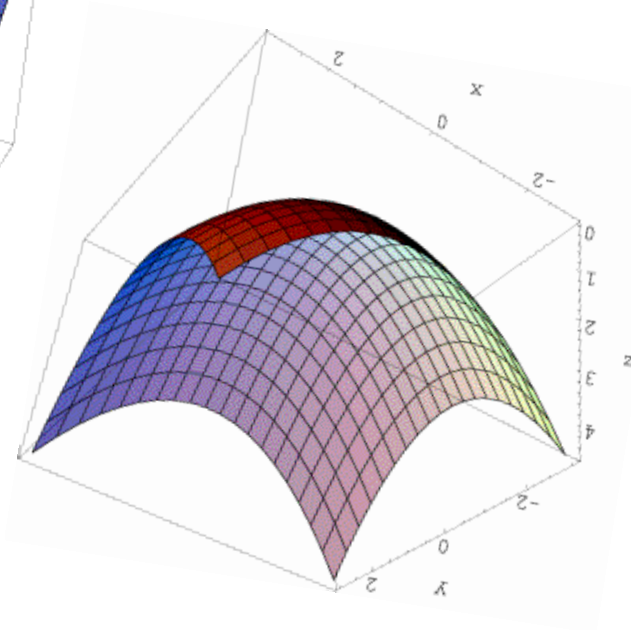
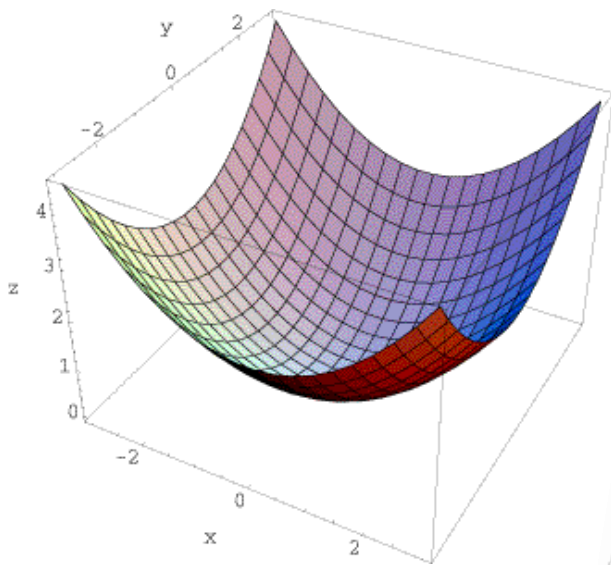
$$\text{Taylor: } f(x + h) \approx h^\top \nabla^2 f(x) h + h^\top \nabla f(x) + f(x)$$



$$\begin{aligned} f(x) &= f(x_0) + f'(x_0)(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2 \\ &\quad + \frac{f'''(x_0)}{3!}(x - x_0)^3 + \frac{f^{(4)}(x_0)}{4!}(x - x_0)^4 + \dots \\ &= \sum_{n=0}^{\infty} \frac{f^{(n)}(x_0)}{n!}(x - x_0)^n. \end{aligned}$$

Formas cuadráticas

- Los eigenvalores de la matriz A determinan la forma cuadrática.



$$A = W \Lambda W^{-1}$$

$$\bar{W}^T = W^T \Rightarrow A = W \Lambda W^T$$

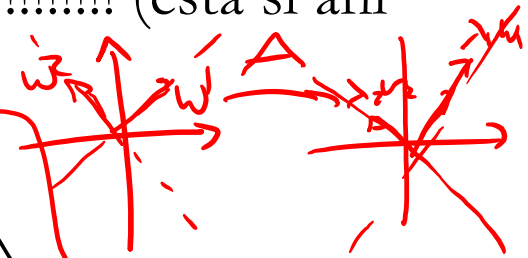
- Las matrices simétricas cumplen lo siguiente que verlo con detalle escapa el objetivo de este curso
 - Siempre son diagonalizables con eigenvalores números reales (Teorema Espectral)
 - Los eigenvectores SON ortogonales!!!!!! (esta si ahí va la prueba)....

v_1, v_2 eigenvectors

$$v_1^T (A v_2) = \lambda_2 v_1^T v_2$$

$$\lambda_1 \neq \lambda_2 \Rightarrow v_1^T v_2 = 0$$

$$\begin{aligned} v_1^T A v_2 &= v_1^T A^T v_2 \\ &= (A v_1)^T v_2 = \lambda_1 v_1^T v_2 \end{aligned}$$



SVD

• Recordatorio Descomposición SVD

Cualquier $A = U \Sigma V^T$ $A \in \mathbb{R}^{m \times n}$

con U ortogonal $U \in \mathbb{R}^{m \times m}$

V ortogonal $V \in \mathbb{R}^{n \times n}$

Σ cuasidiagonal $\Sigma \in \mathbb{R}^{m \times n}$

$$\Sigma = \begin{bmatrix} \sigma_1 & \dots & \sigma_r & 0 \\ 0 & \dots & 0 & 0 \end{bmatrix} \quad \text{ó} \quad \Sigma = \begin{bmatrix} \sigma_1 & \dots & \sigma_r & 0 \\ 0 & \dots & 0 & 0 \\ \hline 0 & \dots & 0 & 0 \end{bmatrix}$$

la diagonal de Σ se llaman valores singulares

y las columnas $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{\min(m,n)} \geq 0$ de U y V son vectores singulares

- Diferencias con eigenvalores
- Eigenvalores es solo para matrices cuadradas
 - En eigenvalores $U = V$
 - " " U es ortogonal solo si A es simétrica
 - \downarrow SVD siempre existe en números reales y eigenvalores siempre en complejos

Relación entre SVD y Eigenvectores

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

$$A^T A = (U \Sigma V^T)^T (U \Sigma V^T) = V \Sigma^T \overbrace{U^T U}^{= I_n \text{ pues } U \text{ ortogonal}} \Sigma V^T = V \Sigma^T \Sigma V^T$$

pero

$$\Sigma^T \Sigma = \begin{bmatrix} \sigma_1^2 & & & 0 \\ & \sigma_2^2 & & \\ & & \ddots & \\ 0 & & & \end{bmatrix}$$

es diagonal

$\Rightarrow A^T A = V D V^T$ con D diagonal (diagonalización)

∴ Los vectores singulares V son los eigenvectores de $A^T A$ con eigenvalores los cuadrados de los valores singulares

Similamente

$$A A^T = U \Sigma V^T V \Sigma^T U^T = U D U^T$$

∴ U son eigenvectores de $A A^T$

Resumen

Descomposición SVD

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

\Rightarrow U son eigenvectores de AA^T
 V son " " " $A^T A$

y $\sigma_1^2, \sigma_2^2, \dots$ son eigenvalores asociados a U y V .

$A^T A$ y AA^T siempre tienen eigenvalores no negativos.

Recordar que $A^T A$ y AA^T definen formas cuadráticas pues son matrices simétricas.

El valor del eigenvalores son justamente los valor de la máxima distorsión:

$$A^T A v = \lambda v \text{ implica } ||Av||^2 / ||v||^2 = \lambda$$

Similar para $AA^T u = \lambda u$.

- Las matrices de la descomposición SVD son precisamente los “ejes de máxima distorsión” de la matriz.
- Eso explica porque resumen el efecto de la transformación lineal

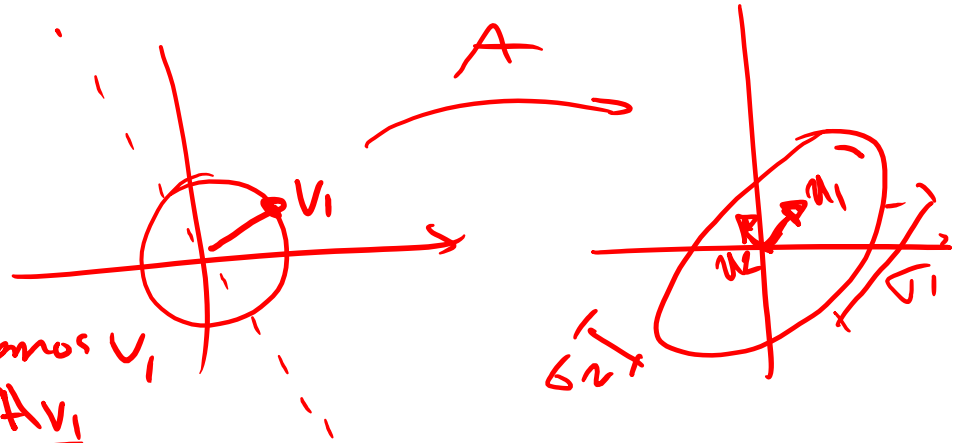
$$A = U \Sigma V^T \Rightarrow AV = U \Sigma$$

Derivación constructiva sin eigenvectores

① Dado A , resolvamos

$$\sigma_1 = \|A\| = \max_{s.t. \|v\|=1} \|Av\|$$

Al vector solución lo llamamos v_1
y denotamos $u_1 = \frac{Av_1}{\|Av_1\|}$



La relación entre u_1 y v_1 es $Av_1 = \sigma_1 u_1$

Obs: De la solución que resolvimos antes, sabemos v_1 es eigenvector de $A^T A$.

② Hay que resolver un nuevo problema de optimización,
calcular

$$\sigma_2 = \max_{s.t. \|v\|=1 \text{ \& } v^T v_1 = 0} \|Av\|$$

llamamos v_2 al
vector solución y
 $u_2 = \frac{Av_2}{\|Av_2\|} \Rightarrow Av_2 = \sigma_2 u_2$

Por construcción

$$V_1^T V_2 = 0$$

$$\|V_1\| = \|V_2\| = 1$$

$$AV_1 = \sigma_1 U_1$$

$$AV_2 = \sigma_2 U_2$$

V_2 eigenvector de $A^T A$

$$U_1^T U_2 = \frac{1}{\sigma_1 \sigma_2} V_1^T \underbrace{A^T A}_{\text{circled}} V_2 = \frac{\lambda_2}{\sigma_1 \sigma_2} \cancel{V_1^T V_2} = 0$$

entonces $U_1^T U_2 = 0$

③ Continuando así, encontramos $V_1, \dots, V_{\min(m,n)}$

y $U_1, \dots, U_{\min(m,n)}$ y $\sigma_1, \dots, \sigma_{\min(m,n)}$

tales que $AV_i = \sigma_i U_i$ con $U_i^T U_i = 1$ $U_i^T U_j = 0$
 $V_i^T V_i = 1$ $V_i^T V_j = 0$

La parte tricky es completar U o V con elementos del kernel de A o A^T según $m > n$ o $m < n$, pero la idea ya está clara.

Con eso tenemos

$$V = [V_1 | V_2 | \dots] \quad U = [U_1 | \dots] \quad \Sigma = [\sigma_1 \dots \sigma_n | 0]$$

con $V^T V = I$ $U^T U = I$ y $AV = U \Sigma$ ~~///~~ ✓