



Take  $f(x) = \|Ax - y\|_2^2 = \left( \sqrt{\sum_{i=1}^n (Ax - y)_i^2} \right)^2$

$$\nabla f(x) = A^T(Ax - y) = 0 \Rightarrow A^T Ax - A^T y = 0 \Rightarrow x = (A^T A)^{-1} A^T y$$

$$x_k = x_{k-1} - \alpha \underbrace{A^T(Ax_k - y)}_{\nabla f(x_k)}$$

↑  
also may  
be hard  
to invert...

Note:  $A \in M_n(\mathbb{C}) \Rightarrow f(x) = (Ax - y)^H (Ax - y)$ , where

$A^H = (A^T)^*$  is the conjugate transpose of  $A$ . One can derive

the relations

$$x_k = x_{k-1} - \alpha A^H (Ax_k - y)$$

$$x = (A^H A)^{-1} A^H y$$

We may want more sophisticated optimization using probability

Let  $X$  be a random variable, with known observation  $y$ .

Assume Gaussian prior:  $P(y|x) \propto e^{-\frac{(Ax-y)^2}{\sigma^2}}$ .

Maximum a posteriori (MAP) estimator:  $\max P(x|y)$

Bayes Thm:  $P(x|y) = \frac{P(y|x)P(x)}{P(y)}$

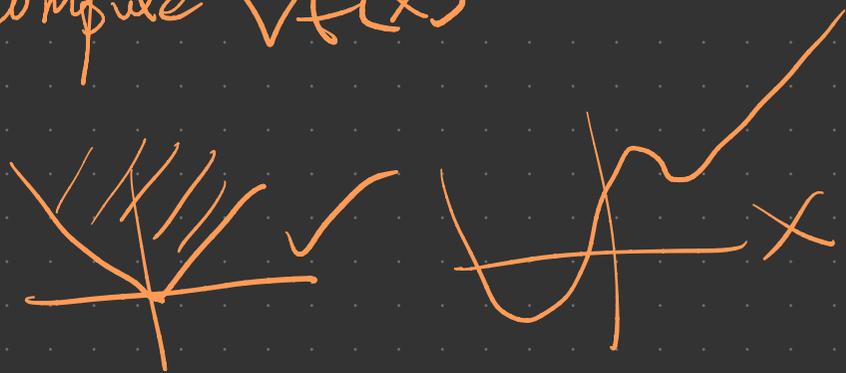
$$\begin{aligned} \Rightarrow \log p(x|y) &= \log p(y|x) - \log p(x) \\ &= -\frac{1}{\sigma^2}(Ax-y)^2 - \frac{1}{\sigma^2}x^2 \end{aligned}$$

$$\rightarrow \min [-\log p(x|y)] = \frac{1}{\sigma^2}(Ax-y)^2 + \frac{1}{\sigma^2}x^2 \quad \text{looks like } \|Ax-y\|_2^2 + R(x)$$

Characteristics of  $f(x)$  that help optimization:

$f(x)$  differentiable: Can analytically compute  $\nabla f(x)$

$f(x)$  convex: local min = global min



$$f(tx_0) + f((1-t)x_1) \leq tf(x_0) + (1-t)f(x_1)$$

The game is to automatically find  $-\nabla f(x)$ , which is why autograd/backprop exists



## Lecture 2: Function Spaces

$V$  is a vector space (over  $\mathbb{R}$ ) if  $\forall v, w \in V, \alpha \in \mathbb{R}$ , we have:

(1)  $v+w \in V$ , (2)  $\alpha v \in V$  s.t.  $\alpha(v+w) = \alpha v + \alpha w \in V$ .

Such elements  $v, w$  are called vectors. Usually think of  $v \in \mathbb{R}^n$ , but this lecture will go over more exotic spaces.

$B = \{v_1, \dots, v_n\} \subseteq V$  is a basis if  $\forall w \in V, \exists! \alpha_1, \dots, \alpha_n \in \mathbb{R}$  s.t.  
 $w = \alpha_1 v_1 + \dots + \alpha_n v_n$ .  $n = \dim V$  is the dimension of  $V$ .

Fact: If  $B_1, B_2$  are bases of  $V$ , then  $\#B_1 = \#B_2 = n$ .

Sps  $B_1 = \{v_1, \dots, v_n\}, B_2 = \{w_1, \dots, w_n\}$ . Then by defn  $\exists!$

$$\begin{aligned} \hookrightarrow \alpha_{11}, \dots, \alpha_{1n} \text{ st } v_1 &= \alpha_{11}w_1 + \alpha_{12}w_2 + \dots + \alpha_{1n}w_n \\ \hookrightarrow \alpha_{21}, \dots, \alpha_{2n} \text{ st } v_2 &= \alpha_{21}w_1 + \alpha_{22}w_2 + \dots + \alpha_{2n}w_n \\ &\vdots \\ v_n &= \alpha_{n1}w_1 + \alpha_{n2}w_2 + \dots + \alpha_{nn}w_n \end{aligned}$$

! unique way to  
convert  $B_1 \rightarrow B_2$

Can write this more succinctly as  $\vec{v} = \underline{\Phi} \vec{w}$  where we know  $\underline{\Phi}$  is invertible.

Fact: every vector space has a basis.!!

Linear Transform  $T: V \rightarrow W$  is a function where

Note: ①  $T(v+w) = T(v) + T(w)$     ②  $T(\alpha v) = \alpha T(v)$

Can always write  $T = \underline{\Phi}_W T' \underline{\Phi}_V$ , applying transformation in arbitrary basis.

Let's do some examples:

Ex 1  $V = \text{Span}_{\mathbb{R}} \{ \sin(ax) : a \in \mathbb{R} \}$ ,  $W = \text{Span}_{\mathbb{R}} \{ \cos(ax) : a \in \mathbb{R} \}$

$\frac{d}{dx} : V \rightarrow W$  is a linear transform.  $\dim V = \dim W = \infty$ .

Ex 2  $V = \text{Span}_{\mathbb{R}} \{ 1, x, \dots, x^2 \}$ .

① What is  $\Phi_V$ , changing to the basis  $\{ 1, x+1, (x+1)^2 \}$ ?

② What is the matrix form of  $\frac{d}{dx}$  in the original basis?

A normed vector space  $(V, \|\cdot\|)$  has a norm  $\|\cdot\|$  that satisfies:

$$i) \|x\| = 0 \iff x = 0 \in V$$

$$ii) \|x+y\| \leq \|x\| + \|y\|$$

$$iii) \|\lambda x\| = |\lambda| \|x\|$$

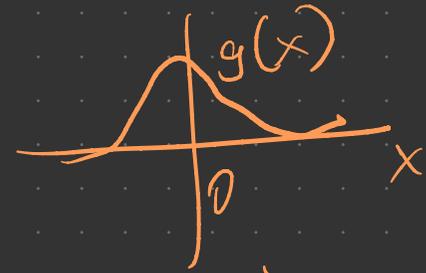
Fact:  $|x \cdot y| \leq \|x\| \|y\|$  (Cauchy-Schwartz)

Def: Let  $f: \Omega \rightarrow \mathbb{R}$ ,  $f \in V$ .  $\|f\|_p := \left( \int_{\Omega} f^p dx \right)^{1/p}$  be denoted as the  $L^p$  norm of  $f$ .

Remark: In ML, you will often see  $L^1, L^2$ , sometimes  $L^0$ .

$L^2$  behaves nicely,  $L^1$  promotes sparsity,  $L^0$  is even better for sparsity.

Convolution:  $(f \star g)(x) = \int_{\mathbb{R}} f(a)g(x-a) da$



Fourier Transform:

$$\hat{f}(k) = \int_{\mathbb{R}} f(x) e^{-2\pi i x k} dx, \quad \check{f}(x) = \int_{\mathbb{R}} f(k) e^{+2\pi i x k} dk$$

(usually)

Facts

①  $\|f(x)\|_2 = \|\hat{f}(k)\|_2$  (Plancherel/Parseval)

②  $(f \star g) = \hat{f} \hat{g}$  (Convolution-Multiplication)

③  $f(x) = A e^{-x^2} \iff \hat{f}(k) = B e^{-k^2}$

## Exercises

- 1) Show that the change of basis  $\Phi: V \rightarrow V$  is a linear transform.
- 2)  $V$  as in Ex 2. What is the vector space  $W$  st  $\int dx(V) = W$ ?
- 3) Define  $\delta(x) = \lim_{\sigma \rightarrow 0} \frac{1}{1 + \sigma^2 x^2}$ ,  $\int_{\mathbb{R}} f(x) \delta(x) = f(0)$ .  
What is  $\hat{\delta}(x)$ ? I.e. What is Fourier transform of a clap?

# Lecture 3: Linear Algebra

$$M = [M]_{ij} = \begin{bmatrix} M_{11} & \dots & M_{1n} \\ \vdots & & \vdots \\ M_{ni} & \dots & M_{nn} \end{bmatrix}$$

$$[M+N]_{ij} = [M]_{ij} + [N]_{ij}$$

$$[MN]_{ij} = \begin{bmatrix} \circ & \circ \\ \circ & \circ \end{bmatrix}$$



$$[MN]_{ij} = \sum_k M_{ik} N_{kj}$$

Tensor Product:  $v \otimes w =$

$$\begin{bmatrix} v_1 w_1 & \dots & v_n w_1 \\ \vdots & & \vdots \\ v_n w_1 & \dots & v_n w_n \end{bmatrix}$$

"Outer Product"

Iterative linear solving: RREF

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} x = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

$$\left[ \begin{array}{cc|c} 1 & 2 & 1 \\ 3 & 4 & 2 \end{array} \right]$$

$\textcircled{1} - 3\textcircled{2} \text{ e } \textcircled{2}$

$$\left[ \begin{array}{cc|c} 1 & 2 & 1 \\ 0 & -2 & -1 \end{array} \right]$$

$\textcircled{2} \rightarrow -\frac{\textcircled{2}}{2}$

$$\left[ \begin{array}{cc|c} 1 & 2 & 1 \\ 0 & 1 & \frac{1}{2} \end{array} \right]$$

$\textcircled{1} \rightarrow \textcircled{1} - 2\textcircled{2}$

$$\left[ \begin{array}{cc|c} 1 & 0 & 0 \\ 0 & 1 & \frac{1}{2} \end{array} \right]$$

$$\Rightarrow \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} x = \begin{bmatrix} 0 \\ \frac{1}{2} \end{bmatrix}$$

$$\Rightarrow x = \begin{bmatrix} 0 \\ \frac{1}{2} \end{bmatrix}$$

exactly one solution

$$\hookrightarrow A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \text{ is}$$

invertible

\* Show row is not invertible for  $v, w \in \mathbb{R}^2$

Determinant

$$\det \begin{bmatrix} M_{11} & M_{12} \\ M_{21} & M_{22} \end{bmatrix} = M_{11}M_{22} - M_{12}M_{21}$$

Invertible  $\Leftrightarrow \det M \neq 0$

$$\det(AB) = \det(A) \det(B)$$

If  $A = \begin{bmatrix} a_1 & 0 \\ 0 & a_2 \end{bmatrix}$ ,

$$\det(A) = a_1 a_2$$

# Eigenvectors/values

$$Mx = \lambda x$$

↑ matrix   ↑ vector   ↑ scalar   ↑ vector

$$\rightarrow Mx = \lambda Ix \quad \leftarrow \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$(M - \lambda I)x = 0 \quad \leftarrow \text{matrix}$$

$$\det(M - \lambda I) = 0 \quad \leftarrow \mathbb{R}$$

↳ characteristic polynomial of  $M$

Fact: If  $A \in M_n(\mathbb{C})$ , char. poly.

of  $A$  has  $n$  complex roots.



## Diagonalization

A matrix  $T \in M_n(\mathbb{R})$  is diagonalizable if

$T = \Lambda S \Lambda^{-1}$ , where  $\Lambda$  is a change of basis, and  $S$  is a diagonal matrix, i.e.

$$S = \begin{bmatrix} s_1 & 0 \\ 0 & s_n \end{bmatrix}$$

Spectral Theorem: Suppose  $A \in M_n(\mathbb{R})$  st  $A = A^T$ . Then:

①  $A$  is diag'ble

②  $A$  has real eigenvalues.

---

Now, what is  $A \in M_{m \times n}(\mathbb{C})$ ? Singular Value Thm!

SVD: Every matrix  $A \in M_{m \times n}(\mathbb{C})$  can be written as

$$A = U \Sigma V^H, \text{ st } UU^H = VV^H = I,$$

$$\Sigma = \begin{bmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_n \\ & & & & 0 \end{bmatrix} \text{ w/ } \sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n$$

Important!  $\text{rk}(A) = \#\{\sigma_n \geq 0\}$  where  $A = U \begin{bmatrix} \sigma_1 & & 0 \\ & \ddots & \\ 0 & & \sigma_n \end{bmatrix} V^H$ .

This decomposition is powerful, since it "orders the components of  $A$  from most important to least important." (See Exercise 3)

Application: Matrix Approximation

$$\arg \min_{\text{rk } B \leq k} \|A - B\|_F^2 = U \begin{bmatrix} \sigma_1 & & 0 \\ & \ddots & \\ 0 & & \sigma_k & 0 \\ & & & \ddots \\ 0 & & & & 0 \end{bmatrix} V^H$$

(Eckart - Young)

i.e. SVD enables optimal low-rank approximation.

## Exercises

① Show that  $\det(A^{-1}) = \frac{1}{\det A}$ , for invertible  $A$ .

② Show that  $u \otimes v \in M_n(\mathbb{R})$  is not invertible, for  $u, v \in \mathbb{R}^n$ .

③ Show that  $A = U \Sigma V^H$ ,  $\Sigma = \begin{bmatrix} \sigma_1 & & 0 \\ & \ddots & \\ 0 & & \sigma_n \end{bmatrix}$  can be written as  $A = \sum_{i=1}^n \sigma_i u_i \otimes v_i$ , for  $u_i, v_i \in \mathbb{C}^n$ .