

~~Boring~~

Review of Essential Mathematics for Machine Learning (Pt. 1)

Linear Algebra

Mathematics in Machine Learning

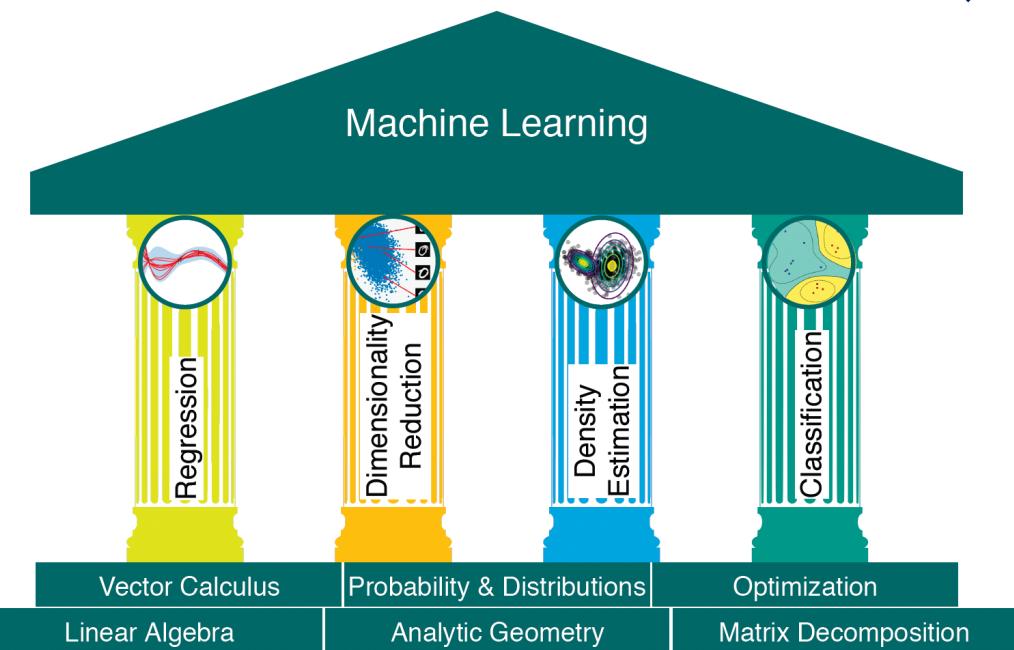
Machine learning involves interrelated tasks of

- *representing data,*
- *modeling data, and*
- *learning from data*

(with little to no
human intervention)

This requires understanding mathematical

- *expression (e.g., vectors, tensors, geometry)*
- *manipulation (e.g., functional mapping, probability), and*
- *optimization (e.g., vector calculus, matrix decomposition)*



PVT data for methane

Thermophysical properties of coexisting gaseous and liquid methane													
T K	Pres. MPa	Density kg/m ³	Density mol/dm ³	E J/mol	II J/mol	S J/(mol·K)	C _v J/(mol·K)	C _p J/(mol·K)	Sound m/s	Visc. μPa s	Therm. W/(m·K)	Diel. Const.	
90.680 ^a	0.01174	451.6	28.15	-5728.0	-5728.0	67.94	34.02	53.46	1253.0	205.0	0.224	1.67701	
90.680 ^b	0.01174	0.2533	0.01379	2235.0	2978.0	164.0	25.19	33.79	249.1	5.64	0.09232	1.00031	
92.0	0.01388	449.8	28.04	-5658.0	-5658.0	68.33	33.89	53.50	1523.0	197.0	0.222	1.67398	
92.0	0.01388	0.2594	0.01835	2266.0	3022.0	163.1	25.21	33.84	250.8	3.69	0.0939	1.00036	
94.0	0.01770	447.2	27.88	-5468.0	-569.0	69.84	33.74	53.91	1506.0	185.0	0.218	1.66936	
94.0	0.01770	0.3667	0.02286	2313.0	3087.0	161.7	25.25	33.94	252.2	3.76	0.0963	1.00045	
96.0	0.02234	444.6	27.71	-5446.0	-5446.0	70.93	33.61	53.74	1489.0	175.0	0.215	1.66469	
96.0	0.02234	0.4528	0.02819	2360.0	3152.0	160.5	25.29	34.05	255.6	3.84	0.0988	1.00055	
98.0	0.02790	441.9	27.55	-5340.0	-5340.0	72.06	33.50	53.91	1470.0	165.0	0.212	1.65998	
98.0	0.02790	0.5553	0.03447	2406.0	3215.0	159.3	25.33	34.17	257.9	3.92	0.0101	1.00067	
100.0	0.03451	439.2	27.38	-5232.0	-5232.0	73.14	33.41	54.11	1452.0	157.0	0.208	1.65523	
100.0	0.03451	0.6709	0.04180	2451.0	3277.0	158.2	25.38	34.30	260.2	4.00	0.0104	1.00082	
102.0	0.04232	436.5	27.21	-5125.0	-5123.0	74.21	33.33	54.32	1433.0	149.0	0.205	1.65042	
102.0	0.04232	0.8067	0.05093	2486.0	3338.0	157.7	25.43	34.45	263.3	4.08	0.0106	1.00098	
104.0	0.05146	433.7	27.04	-5016.0	-5014.0	75.26	33.25	54.54	1413.0	141.0	0.201	1.64557	
104.0	0.05146	0.9634	0.06005	2511.0	3351.0	156.1	25.49	34.61	264.4	4.16	0.0109	1.00118	
106.0	0.06208	430.9	26.86	-4907.0	-4905.0	76.33	33.17	54.78	1394.0	134.0	0.198	1.64066	
106.0	0.06208	1.143	0.07123	2584.0	3456.0	155.2	25.55	34.78	266.4	4.23	0.0112	1.00139	
108.0	0.07434	428.1	26.68	-4798.0	-4795.0	77.32	33.09	55.03	1374.0	128.0	0.195	1.63571	
108.0	0.07434	1.347	0.08394	2627.0	3513.0	154.3	25.62	34.98	268.3	4.31	0.0114	1.00164	
110.0	0.08840	425.2	26.50	-4688.0	-4684.0	78.34	33.01	55.29	1353.0	122.0	0.192	1.63070	
110.0	0.08840	1.577	0.09833	2670.0	3569.0	153.4	25.69	35.19	270.2	4.40	0.0117	1.00192	
112.0	0.10442	422.3	26.33	-4577.0	-4573.0	79.33	32.92	55.56	1333.0	117.0	0.188	1.62563	
112.0	0.10442	1.837	0.1145	2711.0	3623.0	152.5	25.76	35.41	271.9	4.48	0.0120	1.00224	
114.0	0.12266	419.3	26.14	-4466.0	-4461.0	80.32	32.84	55.84	1312.0	111.0	0.185	1.62051	
114.0	0.12266	2.129	0.1327	2752.0	3676.0	151.7	25.84	35.66	273.6	4.56	0.0123	1.00260	
116.0	0.1431	416.4	25.95	-4354.0	-4348.0	81.30	32.75	56.14	1293.0	107.0	0.182	1.61532	
116.0	0.1431	2.454	0.1529	2792.0	3728.0	150.9	25.93	35.93	275.2	4.64	0.0126	1.00300	

Modified Benedict-Rubin-Webb EOS:

$$\begin{aligned}
 P = & \rho RT + \rho^2 [G(1)T + G(2)T^{1/2} + G(3) + G(4)/T + G(5)/T^2] + \rho^3 [G(6)T + G(7) + G(8)/T + G(9)/T^2] \\
 & + \rho^4 [G(10)T + G(11) + G(12)/T] + \rho^5 [G(13)] + \rho^6 [G(14)/T + G(15)/T^2] + \rho^7 [G(16)/T] \\
 & + \rho^8 [G(17)/T + G(18)/T^2] + \rho^9 [G(19)/T^2] + \rho^{10} [G(20)/T^2 + G(21)/T^3] \exp(\gamma\rho^2) \\
 & + \rho^{11} [G(22)/T^2 + G(23)/T^4] \exp(\gamma\rho^2) + \rho^{12} [G(24)/T^2 + G(25)/T^3] \exp(\gamma\rho^2) \\
 & + \rho^{13} [G(26)/T^2 + G(27)/T^4] \exp(\gamma\rho^2) + \rho^{14} [G(28)/T^2 + G(29)/T^3] \exp(\gamma\rho^2) \\
 & + \rho^{15} [G(30)/T^2 + G(31)/T^3 + G(32)/T^4] \exp(\gamma\rho^2).
 \end{aligned}$$

Younglove and Ely, J. Phys. Chem. Ref. Data (1987)

$G(1) = 0.9898937956 \times 10^{-5}$
 $G(2) = 0.2199608275 \times 10^{-1}$
 $G(3) = -0.5322788000$
 $G(4) = 0.2021657962 \times 10^2$
 $G(5) = -0.2234398926 \times 10^4$
 $G(6) = 0.1067940280 \times 10^{-4}$
 $G(7) = 0.1457922469 \times 10^{-3}$
 $G(8) = -0.9265816666$
 $G(9) = 0.2915364732 \times 10^3$
 $G(10) = 0.2313546209 \times 10^{-6}$
 $G(11) = 0.1387214274 \times 10^{-3}$
 $G(12) = 0.4780467451 \times 10^{-2}$
 $G(13) = 0.1176103833 \times 10^{-4}$
 $G(14) = -0.1982096730 \times 10^{-3}$
 $G(15) = -0.2512887756 \times 10^{-1}$
 $G(16) = 0.97488998826 \times 10^{-5}$
 $G(17) = -0.1201291137 \times 10^{-6}$
 $G(18) = 0.4128353939 \times 10^{-4}$
 $G(19) = -0.7218542918 \times 10^{-6}$
 $G(20) = 0.5081738255 \times 10^3$
 $G(21) = -0.9198903192 \times 10^5$
 $G(22) = -0.2732264677 \times 10^1$
 $G(23) = 0.7499024351 \times 10^5$
 $G(24) = 0.1114069080 \times 10^{-2}$
 $G(25) = 0.1083955159 \times 10^1$
 $G(26) = -0.4490960312 \times 10^{-4}$
 $G(27) = -0.1308337847 \times 10^1$
 $G(28) = -0.2371902232 \times 10^{-7}$
 $G(29) = 0.3761652197 \times 10^{-4}$
 $G(30) = -0.2376156954 \times 10^{-9}$
 $G(31) = -0.1237640790 \times 10^{-7}$
 $G(32) = 0.6766926453 \times 10^{-6}$

Essential Definitions: Vectors, Matrices, Tensors

\boldsymbol{x}	Vector:	(denoted in <i>bold lowercase</i>)	elements of a Vector Space V
V	Vector Space:	a set of elements \mathcal{V} with the operations of addition and scalar multiplication, such that	
	$\mathcal{V} + \mathcal{V} \in \mathcal{V}$	a vector in the space plus another in the space yields another in the space	
	$\mathbb{R} \times \mathcal{V} \in \mathcal{V}$	a vector times a real number in the space yields a vector in the space	

You are very familiar with
Euclidean vectors in \mathbb{R}^3

$$\boldsymbol{x} = \begin{bmatrix} x \\ y \\ z \end{bmatrix}$$

or perhaps even \mathbb{R}^n
(tuples of n real numbers)

$$\boldsymbol{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

Essential Definitions: Vectors, Matrices, Tensors

x Vector:

(denoted in ***bold lowercase***)

elements of a **Vector Space** V

V Vector Space:

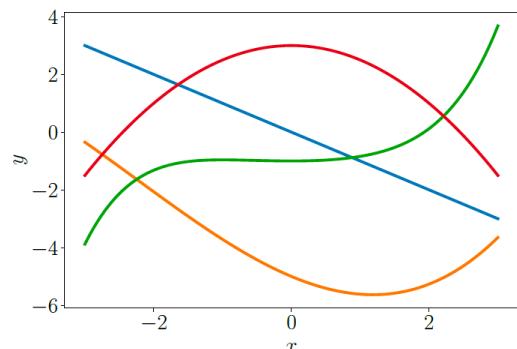
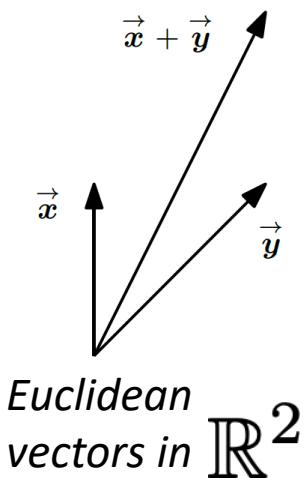
a set of elements \mathcal{V} with the operations of addition and scalar multiplication, such that

$$\mathcal{V} + \mathcal{V} \in \mathcal{V}$$

a vector in the space plus another in the space yields another in the space

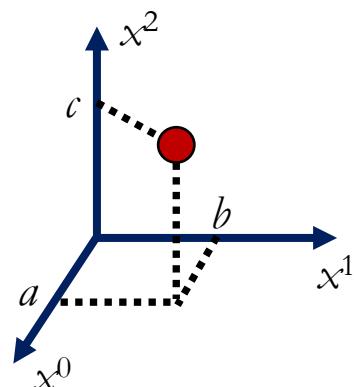
$$\mathbb{R} \times \mathcal{V} \in \mathcal{V}$$

a vector times a real number in the space yields a vector in the space



polynomials are also vectors but not in the Euclidean sense

Often it can be useful to make analogies to Euclidean vectors, since we have mathematical familiarity and intuition



$$f(x) = a + bx + cx^2$$

$$f = \begin{bmatrix} a \\ b \\ c \end{bmatrix}$$

this relates to the concepts of **basis & coordinates** (to be revisited)

Essential Definitions: Vectors, Matrices, Tensors

Matrix:

For $m, n \in \mathbb{N}$

A

(denoted in **BOLD
UPPERCASE**)

an (m,n) matrix is composed of an mn -tuple of elements arranged into m rows and n columns

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix}, \quad a_{ij} \in \mathbb{R}$$

We can say that $\mathbf{A} \in \mathbb{R}^{m \times n}$

*(so can we say **A** lives in a vector space and is thus a vector? Sure!)*

Essential Definitions: Vectors, Matrices, Tensors

Matrix:

For $m, n \in \mathbb{N}$

A

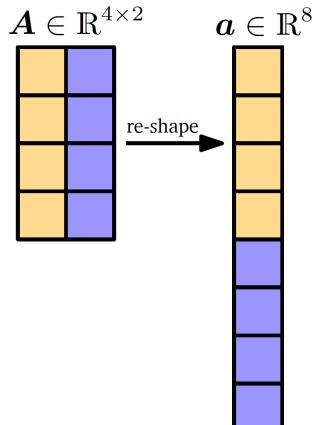
(denoted in **BOLD
UPPERCASE**)

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix}, \quad a_{ij} \in \mathbb{R}$$

We can say that $A \in \mathbb{R}^{m \times n}$

*Matrices can be represented as vectors by stacking columns vertically or stacking rows horizontally
(and this is really how they would be represented on machine memory)*

$$A \in \mathbb{R}^{m \times n} \rightarrow a \in \mathbb{R}^{mn}$$



(Aside: Python stores elements of arrays in ROW-MAJOR order)

```
>>> import numpy as np
>>> a = np.array([[1,3,5],[2,4,6]])
>>> print(a)
[[1 3 5]
 [2 4 6]]
>>> print(a.flatten())
[1 3 5 2 4 6]
```

Essential Definitions: Vectors, Matrices, Tensors

Tensor:

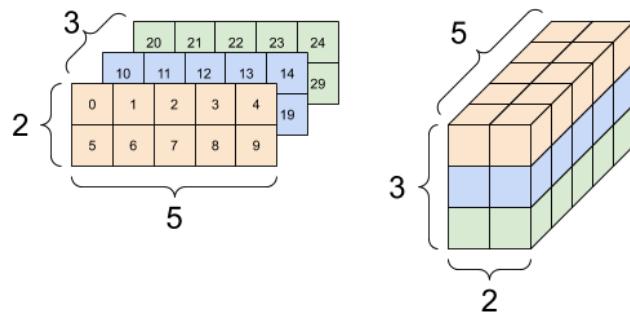
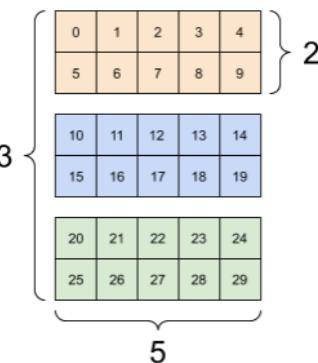
$T, T_{ijk} \dots$

(denoted in **BOLD UPPERCASE** or explicit indices indicating rank)

- Mathematically, tensors are elements within the tensor product over a set of vector spaces
- Practically, tensors are objects that allow generalization of scalars, vectors, matrices, and higher-dimensional objects
- Tensors are defined by their **rank** (distinct from matrix rank)



rank-0 tensor → scalar
rank-1 tensor → vector
rank-2 tensor → matrix
rank-d tensor → array with d axes



rank-3 tensor with shape 3x2x5

In python, NumPy arrays are effectively Tensor-like containers

```
>>> T = np.array([
    [[0, 1, 2, 3, 4],
     [5, 6, 7, 8, 9]],
    [[10, 11, 12, 13, 14],
     [15, 16, 17, 18, 19]],
    [[20, 21, 22, 23, 24],
     [25, 26, 27, 28, 29]]])
>>> print(T)
[[[ 0  1  2  3  4]
  [ 5  6  7  8  9]
  [10 11 12 13 14]
  [15 16 17 18 19]]
 [[20 21 22 23 24]
  [25 26 27 28 29]]]
>>> rank0 = np.array(42).reshape(())
>>> rank1 = np.array(42).reshape((1))
>>> rank2 = np.array(42).reshape((1,1))
>>> rank3 = np.array(42).reshape((1,1,1))
>>> rank4 = np.array(42).reshape((1,1,1,1))
>>> for i in range(5):
...     print(locals()['rank{}'.format(i)])
...
42
[42]
[[42]]
[[[42]]]
[[[[42]]]]
>>> rank4[0,0,0,0]
```

```
array([[ 0,  1,  2,  3,  4],
       [ 5,  6,  7,  8,  9]])
[[[ 0  1  2  3  4]
  [ 5  6  7  8  9]
  [10 11 12 13 14]
  [15 16 17 18 19]]
 [[20 21 22 23 24]
  [25 26 27 28 29]]]
>>> T[0]
array([[ 0,  1,  2,  3,  4],
       [ 5,  6,  7,  8,  9]])
```

Addition/Multiplication with Matrices

Addition

$$A, B, C \in \mathbb{R}^{m \times n}$$

$$C = A + B$$

$$c_{ij} = a_{ij} + b_{ij}$$

```
>>> A = np.array([[1,2],[3,4]])
>>> B = np.array([[8,7],[6,5]])
>>> C = A+B
>>> print(C)
[[9 9]
 [9 9]]
```

Scalar Multiplication

$$B = \lambda A, \lambda \in \mathbb{R}$$

$$b_{ij} = \lambda a_{ij}$$

Multiplication

$$A \in \mathbb{R}^{m \times n}$$

$$C = AB$$

$$B \in \mathbb{R}^{n \times k}$$

$$C \in \mathbb{R}^{m \times k}$$

$$c_{ij} = \sum_{l=1}^n a_{il} b_{lj}$$

Shorthand Einstein notation:
any index appearing twice in
the same term implies
summation over that index

$$c_{ij} = a_{il} b_{lk}$$

- matrix multiplication is not elementwise!

$$c_{ij} \neq a_{ij} b_{ij}$$

This is what you will get though in many
programming languages just using the
multiply operator. (Hadamard product)

- Order matters! $AB \neq BA$ in general

```
>>> print(np.matmul(A,B))
[[20 17]
 [48 41]]
>>> print(np.dot(A,B))
[[20 17]
 [48 41]]
>>> print(np.einsum('il,lj->ij',A,B))
[[20 17]
 [48 41]]
>>> print(A@B)
[[20 17]
 [48 41]]
```

```
>>> print(np.dot(B,A))
[[29 44]
 [21 32]]
```

Matrix Inverse and Transpose

Let $A \in \mathbb{R}^{n \times n}$ and $B \in \mathbb{R}^{n \times n}$

Suppose $AB = I_n = BA$

Then $B = A^{-1}$ is the inverse of A

$$I_n := \begin{bmatrix} 1 & 0 & \cdots & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 & \cdots & 1 \end{bmatrix} \in \mathbb{R}^{n \times n}$$

- Not every matrix is invertible

$$A := \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \in \mathbb{R}^{2 \times 2} \quad A' := \begin{bmatrix} a_{22} & -a_{12} \\ -a_{21} & a_{11} \end{bmatrix}$$

- If it is invertible, then its inverse is unique

$$AA' = \begin{bmatrix} a_{11}a_{22} - a_{12}a_{21} & 0 \\ 0 & a_{11}a_{22} - a_{12}a_{21} \end{bmatrix} = (a_{11}a_{22} - a_{12}a_{21})I$$

- The matrix **determinant** can be used to check if a square matrix is invertible

$$A^{-1} = \frac{1}{a_{11}a_{22} - a_{12}a_{21}} \begin{bmatrix} a_{22} & -a_{12} \\ -a_{21} & a_{11} \end{bmatrix}$$

Matrix Inverse and Transpose

Let $A \in \mathbb{R}^{n \times n}$ and $B \in \mathbb{R}^{n \times n}$

Suppose $AB = I_n = BA$

Then $B = A^{-1}$ is the inverse of A

```
>>> print(A)
[[1 2 1]
 [4 4 5]
 [6 7 7]]
>>> np.linalg.inv(A)
array([[-7., -7.,  6.],
       [ 2.,  1., -1.],
       [ 4.,  5., -4.]])
>>> B = np.linalg.inv(A)
>>> print(B)
[[-7. -7.  6.]
 [ 2.  1. -1.]
 [ 4.  5. -4.]]
>>> print(np.linalg.inv(B))
[[1. 2. 1.]
 [4. 4. 5.]
 [6. 7. 7.]]
```

Let $A \in \mathbb{R}^{n \times m}$ and $B \in \mathbb{R}^{m \times n}$

Suppose $b_{ij} = a_{ji}$

Then $B = A^T$ is the transpose of A

```
>>> print(A.transpose())
[[1 4 6]
 [2 4 7]
 [1 5 7]]
```

```
>>> print(A.T)
[[1 4 6]
 [2 4 7]
 [1 5 7]]
```

Let's go practice in Python to convince ourselves of the following useful properties:

$$\begin{aligned}AA^{-1} &= I = A^{-1}A \\(AB)^{-1} &= B^{-1}A^{-1} \\(A+B)^{-1} &\neq A^{-1} + B^{-1} \\(A^T)^T &= A \\(A+B)^T &= A^T + B^T \\(AB)^T &= B^T A^T\end{aligned}$$

Systems of Linear Equations

Consider the following system of linear equations

$$a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n = b_1$$

⋮

$$a_{m1}x_1 + a_{m2}x_2 + \cdots + a_{mn}x_n = b_m$$

every vector $\mathbf{x} = [x_1, \dots, x_n]^T$
satisfying the system is a solution

Any given system of equations may have

- **no** solution
- **one** unique solution
- **infinitely** many solutions

tangible example:

- A chemical manufacturing company has access to 4, 3, and 3 units of chemical **A**, **B**, and **C**;
- These base chemicals can be used in the synthesis of x_j units of downstream products ($j= 1,2,3$)
- One unit of chemical **A** is required for each product
- A unit of chemical **B** is necessary in product **1** and **2** but unneeded for **3**
- One unit of chemical **C** is needed for product **2** and two units are needed for product **3**; none are needed for product **1**

Discuss:

What is the optimal production plan to ensure most complete utilization of resources?

What modifications might we practically consider in a realistic planning scenario?

Matrix Representation of Linear Equations

Consider the following system of linear equations

$$a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n = b_1$$

⋮

$$a_{m1}x_1 + a_{m2}x_2 + \cdots + a_{mn}x_n = b_m$$

$$\begin{bmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & & \vdots \\ a_{m1} & \cdots & a_{mn} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} b_1 \\ \vdots \\ b_m \end{bmatrix}$$

$$[m \times n] [n \times 1] = [m \times 1]$$

$$\mathbf{A}\mathbf{x} = \mathbf{b}$$

A set of linear equations can be compactly expressed in Linear algebraic terms and exploited to find solutions (in lieu of the intuition or systematic solution used previously)

Solution set is invariant with respect to ***elementary transformations***:

- Exchanging two equations (i.e., swapping rows)
- Multiplying an equation by a non-zero constant
- Adding two equations (rows)

General solutions can be represented as

- A particular solution
- ... plus all solutions to $\mathbf{Ax} = \mathbf{0}$

Row Echelon Form & Gaussian Elimination

Obtaining solutions is “straightforward” by rearranging matrices into (Reduced) Row-Echelon Form via Gaussian Elimination

$$\begin{array}{ccccccc} -2x_1 & + & 4x_2 & - & 2x_3 & - & x_4 & + & 4x_5 & = & -3 \\ 4x_1 & - & 8x_2 & + & 3x_3 & - & 3x_4 & + & x_5 & = & 2 \\ x_1 & - & 2x_2 & + & x_3 & - & x_4 & + & x_5 & = & 0 \\ x_1 & - & 2x_2 & & & - & 3x_4 & + & 4x_5 & = & a \end{array}$$

convert to “augmented” matrix form


$$\left[\begin{array}{ccccc|c} -2 & 4 & -2 & -1 & 4 & -3 \\ 4 & -8 & 3 & -3 & 1 & 2 \\ 1 & -2 & 1 & -1 & 1 & 0 \\ 1 & -2 & 0 & -3 & 4 & a \end{array} \right]$$

Row Echelon Form & Gaussian Elimination

Obtaining solutions is “straightforward” by rearranging matrices into (Reduced) Row-Echelon Form via Gaussian Elimination

$$\begin{array}{ccccccc} -2x_1 & + & 4x_2 & - & 2x_3 & - & x_4 & + & 4x_5 & = & -3 \\ 4x_1 & - & 8x_2 & + & 3x_3 & - & 3x_4 & + & x_5 & = & 2 \\ x_1 & - & 2x_2 & + & x_3 & - & x_4 & + & x_5 & = & 0 \\ x_1 & - & 2x_2 & & & - & 3x_4 & + & 4x_5 & = & a \end{array}$$

convert to “augmented” matrix form


$$\left[\begin{array}{ccccc|c} -2 & 4 & -2 & -1 & 4 & -3 \\ 4 & -8 & 3 & -3 & 1 & 2 \\ 1 & -2 & 1 & -1 & 1 & 0 \\ 1 & -2 & 0 & -3 & 4 & a \end{array} \right] \begin{matrix} \\ \\ \text{swap} \\ \end{matrix}$$

Now we want to use some equations to systematically eliminate variables (i.e., achieve leading zeros) from other equations (i.e., rows)

$$\left[\begin{array}{ccccc|c} 1 & -2 & 1 & -1 & 1 & 0 \\ 4 & -8 & 3 & -3 & 1 & 2 \\ -2 & 4 & -2 & -1 & 4 & -3 \\ 1 & -2 & 0 & -3 & 4 & a \end{array} \right] \begin{array}{l} \text{subtract } 4R1 \\ \text{add } 2R1 \\ \text{subtract } 1R1 \end{array}$$

Row Echelon Form & Gaussian Elimination

Obtaining solutions is “straightforward” by rearranging matrices into (Reduced) Row-Echelon Form via Gaussian Elimination

$$\begin{array}{ccccccc} -2x_1 & + & 4x_2 & - & 2x_3 & - & x_4 & + & 4x_5 = -3 \\ 4x_1 & - & 8x_2 & + & 3x_3 & - & 3x_4 & + & x_5 = 2 \\ x_1 & - & 2x_2 & + & x_3 & - & x_4 & + & x_5 = 0 \\ x_1 & - & 2x_2 & & & - & 3x_4 & + & 4x_5 = a \end{array}$$

$$\left[\begin{array}{ccccc|c} 1 & -2 & 1 & -1 & 1 & 0 \\ 0 & 0 & 1 & -1 & 3 & -2 \\ 0 & 0 & 0 & 1 & -2 & 1 \\ 0 & 0 & 0 & 0 & 0 & a+1 \end{array} \right]$$

convert to “augmented” matrix form

$$\left[\begin{array}{ccccc|c} -2 & 4 & -2 & -1 & 4 & -3 \\ 4 & -8 & 3 & -3 & 1 & 2 \\ 1 & -2 & 1 & -1 & 1 & 0 \\ 1 & -2 & 0 & -3 & 4 & a \end{array} \right]$$

divide R3 by -3

$$\left[\begin{array}{ccccc|c} 1 & -2 & 1 & -1 & 1 & 0 \\ 0 & 0 & -1 & 1 & -3 & 2 \\ 0 & 0 & 0 & -3 & 6 & -3 \\ 0 & 0 & 0 & 0 & 0 & a+1 \end{array} \right]$$

Now we want to use some equations to systematically eliminate variables (i.e., achieve leading zeros) from other equations (i.e., rows)

$$\left[\begin{array}{ccccc|c} 1 & -2 & 1 & -1 & 1 & 0 \\ 4 & -8 & 3 & -3 & 1 & 2 \\ -2 & 4 & -2 & -1 & 4 & -3 \\ 1 & -2 & 0 & -3 & 4 & a \end{array} \right] \xrightarrow{\text{subtract } 4R1} \left[\begin{array}{ccccc|c} 1 & -2 & 1 & -1 & 1 & 0 \\ 0 & 0 & -1 & 1 & -3 & 2 \\ 0 & 0 & 0 & -3 & 6 & -3 \\ 0 & 0 & -1 & -2 & 3 & a \end{array} \right]$$

subtract $4R1$

add $2R1$

subtract $1R1$

subtract $(R2+R3)$ from $R4$

$$\left[\begin{array}{ccccc|c} 1 & -2 & 1 & -1 & 1 & 0 \\ 0 & 0 & -1 & 1 & -3 & 2 \\ 0 & 0 & 0 & -3 & 6 & -3 \\ 0 & 0 & -1 & -2 & 3 & a \end{array} \right]$$

Row Echelon Form & Gaussian Elimination

Obtaining solutions is “straightforward” by rearranging matrices into (Reduced) Row-Echelon Form via Gaussian Elimination

$$\lambda_1 \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} + \lambda_2 \begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \end{bmatrix} + \lambda_3 \begin{bmatrix} -1 \\ -1 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ -2 \\ 1 \\ 0 \end{bmatrix} \rightarrow \mathbf{x}^* = \begin{bmatrix} 2 \\ 0 \\ -1 \\ 1 \\ 0 \end{bmatrix}$$

Particular solutions can be obtained by working with basic variables (easiest to work with rightmost to left)

The augmented matrix is shown in row-echelon form. It consists of a 4x5 coefficient matrix and a 4x1 column of constants. Yellow arrows indicate the steps of Gaussian elimination: 1) Swap Row 1 and Row 2. 2) Add Row 1 to Row 2. 3) Add Row 1 to Row 3. 4) Add Row 1 to Row 4. The final matrix is:

$$\left[\begin{array}{ccccc|c} 1 & -2 & 1 & -1 & 1 & 0 \\ 0 & 0 & 1 & -1 & 3 & -2 \\ 0 & 0 & 0 & 1 & -2 & 1 \\ 0 & 0 & 0 & 0 & 0 & a+1 \end{array} \right]$$

Row-echelon Form

- Any rows with only zeros are at the bottom
- First nonzero number from the left == pivot
- all pivots are strictly to the right of pivots above
→ “staircase structure”
- basic variables == pivot columns
- free variables == non-pivot columns

Row Echelon Form & Gaussian Elimination

Obtaining solutions is “straightforward” by rearranging matrices into (Reduced) Row-Echelon Form via Gaussian Elimination

For the **general solution**, we focus on non-pivot columns and how they can be expressed as linear combinations of pivot columns (on their left), to satisfy $\mathbf{Ax} = \mathbf{0}$

$$\begin{bmatrix} 1 & -2 & 0 & 0 & -2 \\ 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 1 & -2 \end{bmatrix} \quad \begin{array}{l} 1. \ R_2 + R_3 \rightarrow R_2 \\ 2. \ R_1 + R_3 - R_2 \rightarrow R_1 \end{array}$$

$$\mathbf{x} = \begin{bmatrix} 2 \\ 0 \\ -1 \\ 1 \\ 0 \end{bmatrix} + \lambda_1 \begin{bmatrix} 2 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} + \lambda_2 \begin{bmatrix} 2 \\ 0 \\ -1 \\ 2 \\ 1 \end{bmatrix}$$

$$\left[\begin{array}{ccccc|c} 1 & -2 & 1 & -1 & 1 \\ 0 & 0 & 1 & -1 & 3 \\ 0 & 0 & 0 & 1 & -2 \\ 0 & 0 & 0 & 0 & 0 \end{array} \right]$$

Row-echelon Form

- Any rows with *only zeros* are at the bottom
- First nonzero number from the left == *pivot*
- all pivots are strictly to the right of pivots above
→ “staircase structure”
- *basic variables* == *pivot columns*
- *free variables* == *non-pivot columns*

Reduced Row-echelon Form

- *Row-echelon form*
- All pivots are 1
- Pivot is the only nonzero entry in its column

Practical Trick for homogeneous solutions

A is $m \times n$... if we extend it to $n \times n$ by adding $n-m$ rows of $[0 \dots 0 -1 0 \dots 0]$ such that 1 or -1 populates the diagonal...

$$\begin{bmatrix} 1 & -2 & 0 & 0 & -2 \\ 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 1 & -2 \end{bmatrix}$$



$$\begin{bmatrix} 1 & -2 & 0 & 0 & -2 \\ 0 & -1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 1 & -2 \\ 0 & 0 & 0 & 0 & -1 \end{bmatrix}$$

$$x = \begin{bmatrix} 2 \\ 0 \\ -1 \\ 1 \\ 0 \end{bmatrix} + \lambda_1 \begin{bmatrix} 2 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} + \lambda_2 \begin{bmatrix} 2 \\ 0 \\ -1 \\ 2 \\ 1 \end{bmatrix}$$

We can simply read the homogeneous solutions from the non-pivot columns

Calculating the inverse of a matrix

Computing the inverse of a matrix is equivalent to solving a system of linear equations

$$AX = I_n \rightarrow [A|I_n] \rightarrow \dots \rightarrow [I_n|A^{-1}]$$

Example

$$A = \begin{bmatrix} 1 & 0 & 2 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 2 & 0 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix} \xrightarrow{\text{Row Operations}} \left[\begin{array}{cccc|cccc} 1 & 0 & 2 & 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 1 & 0 & 0 \\ 1 & 2 & 0 & 1 & 0 & 0 & 1 & 0 \\ 1 & 1 & 1 & 1 & 0 & 0 & 0 & 1 \end{array} \right]$$

↓

$$\left[\begin{array}{cccc|cccc} 1 & 0 & 0 & 0 & -1 & 2 & -2 & 2 \\ 0 & 1 & 0 & 0 & 1 & -1 & 2 & -2 \\ 0 & 0 & 1 & 0 & 1 & -1 & 1 & -1 \\ 0 & 0 & 0 & 1 & -1 & 0 & -1 & 2 \end{array} \right] \xrightarrow{\text{Row Operations}} A^{-1} = \begin{bmatrix} -1 & 2 & -2 & 2 \\ 1 & -1 & 2 & -2 \\ 1 & -1 & 1 & -1 \\ -1 & 0 & -1 & 2 \end{bmatrix}$$

Calculating the inverse of a matrix

Computing the inverse of a matrix is equivalent to solving a system of linear equations

$$AX = I_n \rightarrow [A|I_n] \rightarrow \dots \rightarrow [I_n|A^{-1}]$$

Example

$$A = \begin{bmatrix} 1 & 0 & 2 \\ 1 & 1 & 0 \\ 1 & 2 & 0 \\ 1 & 1 & 1 \end{bmatrix}$$

```
>>> A = np.array([[1,0,2,0],[1,1,0,0],[1,2,0,1],[1,1,1,1]])
>>> print(A)
[[1 0 2 0]
 [1 1 0 0]
 [1 2 0 1]
 [1 1 1 1]]
>>> np.linalg.inv(A)
array([[-1.,  2., -2.,  2.],
       [ 1., -1.,  2., -2.],
       [ 1., -1.,  1., -1.],
       [-1.,  0., -1.,  2.]])
>>> import scipy.linalg as spla
>>> spla.inv(A)
array([[-1.,  2., -2.,  2.],
       [ 1., -1.,  2., -2.],
       [ 1., -1.,  1., -1.],
       [-1.,  0., -1.,  2.]])
```

$$\begin{bmatrix} -1 & 2 & -2 & 2 \\ 1 & -1 & 2 & -2 \\ 1 & -1 & 1 & -1 \\ -1 & 0 & -1 & 2 \end{bmatrix}$$

Matrix Method for Solving a Linear system

We can generally use matrices to solve linear systems expressed as $\mathbf{A}\mathbf{x} = \mathbf{b}$

*...if the columns in \mathbf{A} are **linearly independent**, then*

$$\mathbf{A}\mathbf{x} = \mathbf{b} \Leftrightarrow \mathbf{A}^T \mathbf{A}\mathbf{x} = \mathbf{A}^T \mathbf{b} \Leftrightarrow \mathbf{x} = \underbrace{(\mathbf{A}^T \mathbf{A})^{-1}}_{\text{Moore-Penrose pseudo-inverse}} \mathbf{A}^T \mathbf{b}$$

Moore-Penrose
pseudo-inverse

*... if \mathbf{A} is square and **invertible**, then simply*

$$(\mathbf{AB})^{-1} = \mathbf{B}^{-1} \mathbf{A}^{-1}$$

$$\mathbf{x} = \mathbf{A}^{-1} \mathbf{b}$$

Linear Independence and Basis Sets

Let $\mathbf{x}_1, \dots, \mathbf{x}_k \in V$ If $\mathbf{0} = \sum_{i=1}^k \lambda_i \mathbf{x}_i$ is satisfied by only the trivial solution, then the \mathbf{x} are linearly independent

Use Gaussian elimination to convert the system into row-echelon form:

- Pivot columns represent vectors that are linearly independent from the vectors to the left
- Non-pivot columns are linear combinations of those to the left

If $\forall \mathbf{v} \in V, \exists \lambda_i$ s.t. $\mathbf{v} = \sum_{i=1}^k \lambda_i \mathbf{x}_i$ and the \mathbf{x}_i are linearly independent

then we can call $\mathcal{B} = \{\mathbf{x}_1, \dots, \mathbf{x}_k\}$ a basis for V

e.g.,

for \mathbb{R}^3 the canonical/standard basis is $\mathcal{B} = \left\{ \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \right\}$

Rank and Invertibility

In a matrix A , the number of linearly independent columns equals the number of linearly independent rows, and we call this the *Rank*

Two important properties:

$\forall A \in \mathbb{R}^{m \times n}$ and $b \in \mathbb{R}^m$,

$\exists x$ s.t $Ax = b$, iff $\text{rank}(A) = \text{rank}(A|b)$

“A solution exists only if b exists in the subspace spanned by the columns of A ”

$\forall A \in \mathbb{R}^{n \times n}, \exists A^{-1}$ iff $\text{rank}(A) = n$

“Invertible square matrices (regular) possess linearly independent columns”

How can we tell what the rank is for a matrix?

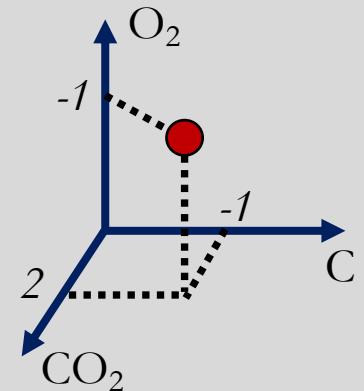
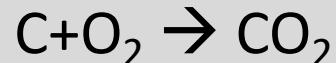
A Chemically motivated Example

Consider the following “goal gasification” reactions observed during steam reformation of coal for the eventual production of syngas



- $\text{C} + \text{H}_2\text{O} \rightleftharpoons \text{CO} + \text{H}_2$
- $\text{C} + 2\text{H}_2\text{O} \rightleftharpoons \text{CO}_2 + \text{H}_2$
- $\text{C} + 2\text{H}_2 \rightleftharpoons \text{CH}_4$
- $\text{C} + \text{CO}_2 \rightleftharpoons 2\text{CO}$
- $\text{CO} + \text{H}_2\text{O} \rightleftharpoons \text{CO}_2 + \text{H}_2$

We can represent reactions as linear combinations of molecules, with molecules or fragments representing a basis in chemical space



How many independent chemical reactions characterize the system?

$$\mathcal{B} = \{e_{\text{C}}, e_{\text{H}_2\text{O}}, e_{\text{CO}}, e_{\text{H}_2}, e_{\text{CO}_2}, e_{\text{CH}_4}\}$$

$$A = \begin{bmatrix} -1 & -1 & 1 & 1 & 0 & 0 \\ -1 & -2 & 0 & 2 & 1 & 0 \\ -1 & 0 & 0 & -2 & 0 & 1 \\ -1 & 0 & 2 & -1 & 0 & 0 \\ 0 & -1 & -1 & 1 & 1 & 0 \end{bmatrix}$$

A Chemically motivated Example

Consider the following “goal gasification” reactions observed during steam reformation of coal for the eventual production of syngas



- $C + H_2O \rightleftharpoons CO + H_2$
- $C + 2H_2O \rightleftharpoons CO_2 + H_2$
- $C + 2H_2 \rightleftharpoons CH_4$
- $C + CO_2 \rightleftharpoons 2CO$
- $CO + H_2O \rightleftharpoons CO_2 + H_2$

$$\mathcal{B} = \{e_C, e_{H_2O}, e_{CO}, e_{H_2}, e_{CO_2}, e_{CH_4}\}$$

$$A = \begin{bmatrix} -1 & -1 & 1 & 1 & 0 & 0 \\ -1 & -2 & 0 & 2 & 1 & 0 \\ -1 & 0 & 0 & -2 & 0 & 1 \\ -1 & 0 & 2 & -1 & 0 & 0 \\ 0 & -1 & -1 & 1 & 1 & 0 \end{bmatrix}$$



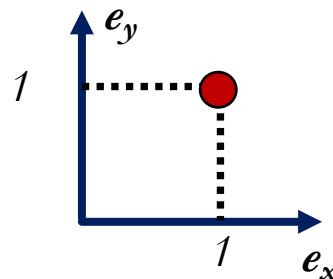
```
>>> print(A)
[[[-1 -1 1 1 0 0]
 [-1 -2 0 2 1 0]
 [-1 0 0 -2 0 1]
 [-1 0 2 0 -1 0]
 [ 0 -1 -1 1 1 0]]]
>>> np.linalg.matrix_rank(A)
3
```

- $C + H_2O \rightleftharpoons CO + H_2$
- $CO + H_2O \rightleftharpoons CO_2 + H_2$
- $2CO + H_2 \rightleftharpoons CO_2 + CH_4$

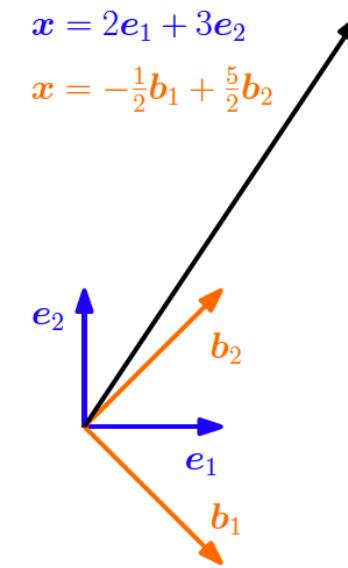
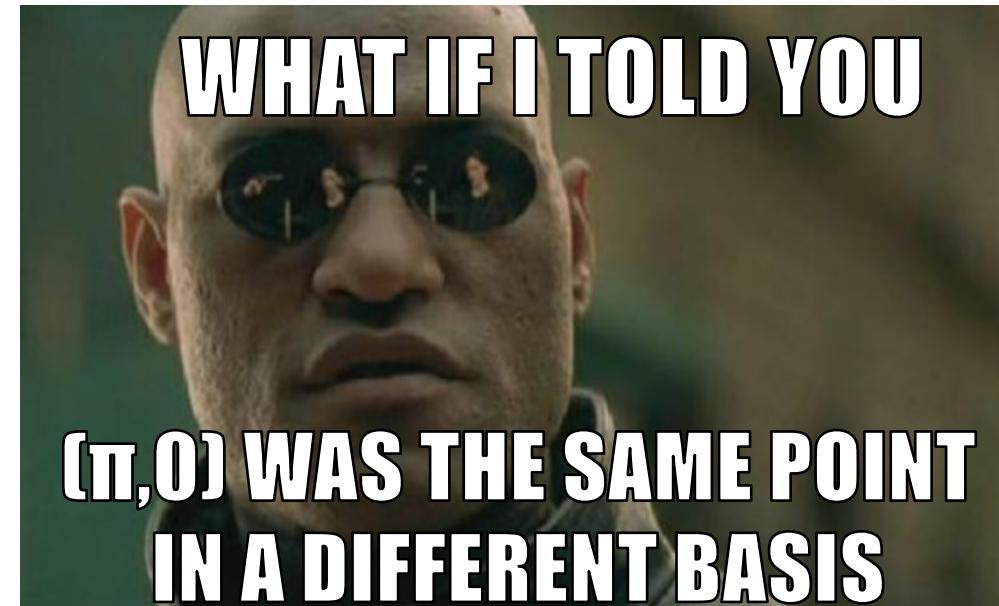
$$A = \begin{bmatrix} -1 & -1 & 1 & 1 & 0 & 0 \\ 0 & -1 & -1 & 1 & 1 & 0 \\ 0 & 0 & -2 & -2 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

Coordinate Transformations/Linear Mappings

The concept of a basis is also fundamentally connected to our notion of a coordinate; a coordinate identifies our position within a vector space spanned by a particular basis



*the point is defined
by the tuple $(1,1)$*



Changing coordinates within a given vector space can be achieved via matrix multiplication

$$T = \begin{bmatrix} \pi & 0 \\ -1 & 1 \end{bmatrix}$$

transformation matrix

$$\mathbf{v} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

coordinate in "old" basis

$$T\mathbf{v} = \begin{bmatrix} \pi \\ 0 \end{bmatrix}$$

coordinate in "new" basis

Coordinate Transformations/Linear Mappings

This operation is formally defined through a “linear mapping,” Φ which can transform between vector spaces

- A *linear mapping preserves the structure of the origin vector space if*

$$\forall \mathbf{x}, \mathbf{y} \in V \forall \lambda, \psi \in \mathbb{R} : \Phi(\lambda \mathbf{x} + \psi \mathbf{y}) = \lambda \Phi(\mathbf{x}) + \psi \Phi(\mathbf{y})$$

Given

$$V$$

$$\mathcal{B} = (\mathbf{b}_1, \dots, \mathbf{b}_n)$$

$$\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_n]$$

$$\mathbf{v} = \sum_{i=1}^n \alpha_i \mathbf{b}_i \quad \alpha = \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix} \in \mathbb{R}^n$$

$$W$$

$$\mathcal{C} = (\mathbf{c}_1, \dots, \mathbf{c}_m)$$

A mapping from V to W given by

$$\Phi(\mathbf{b}_j) = \sum_{i=1}^m \alpha_{ij} \mathbf{c}_i$$

The coordinates of the transformed \mathbf{b}_j with respect to the basis C are in the j^{th} column of T

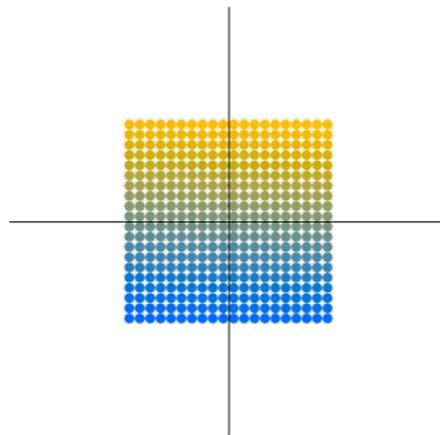
$$T_\Phi(i, j) = \alpha_{ij}$$

$$\mathbf{v}_W = T_\Phi \mathbf{v}_V$$

Coordinate Transformations/Linear Mappings

Which matrices lead to which coordinate mappings?

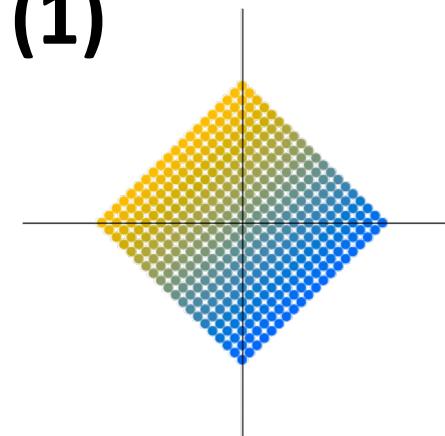
original data



A

$$\begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}$$

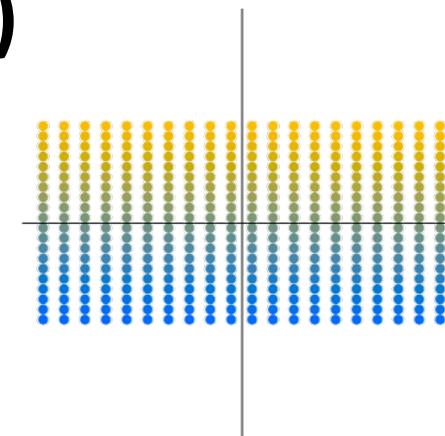
(1)



B

$$\begin{bmatrix} 3 & -1 \\ 1 & -1 \end{bmatrix}$$

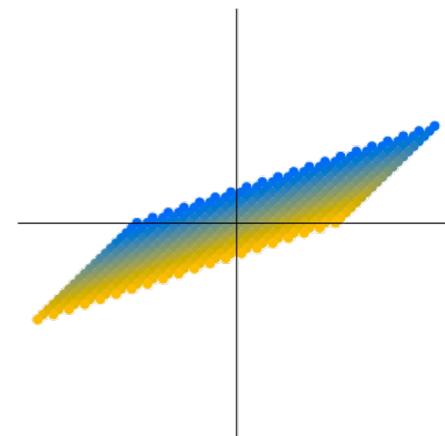
(2)



C

$$\begin{bmatrix} \cos\left(\frac{\pi}{4}\right) & -\sin\left(\frac{\pi}{4}\right) \\ \sin\left(\frac{\pi}{4}\right) & \cos\left(\frac{\pi}{4}\right) \end{bmatrix}$$

(3)



Changing Bases

Often representation and manipulation of data can be easier in a particular basis motivating a change from one to another

Bases in V

$$B = (b_1, \dots, b_n), \quad \tilde{B} = (\tilde{b}_1, \dots, \tilde{b}_n)$$

Bases in W

$$C = (c_1, \dots, c_m), \quad \tilde{C} = (\tilde{c}_1, \dots, \tilde{c}_m)$$

Transformation matrices

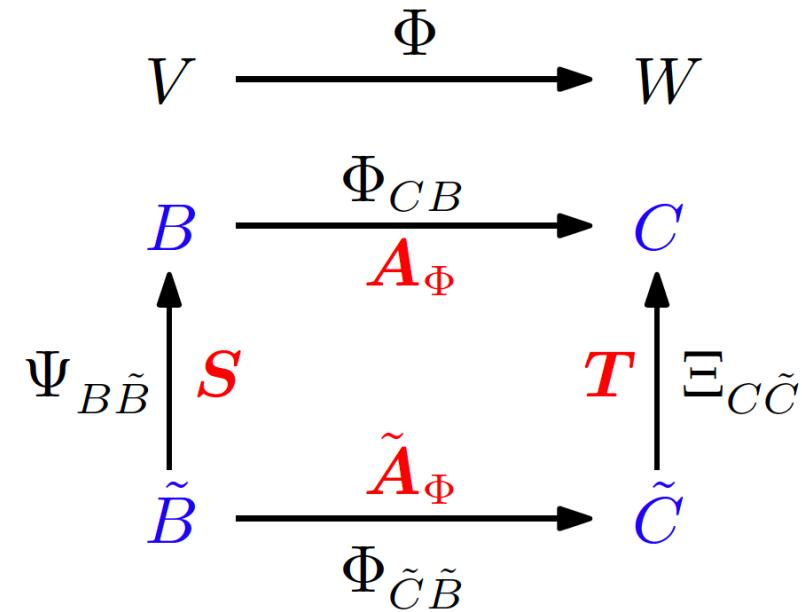
$$S \in \mathbb{R}^{n \times n} \quad \tilde{B} \rightarrow B$$

$$T \in \mathbb{R}^{m \times m} \quad \tilde{C} \rightarrow C$$

$$A_\Phi \in \mathbb{R}^{m \times n} \quad B \rightarrow C$$

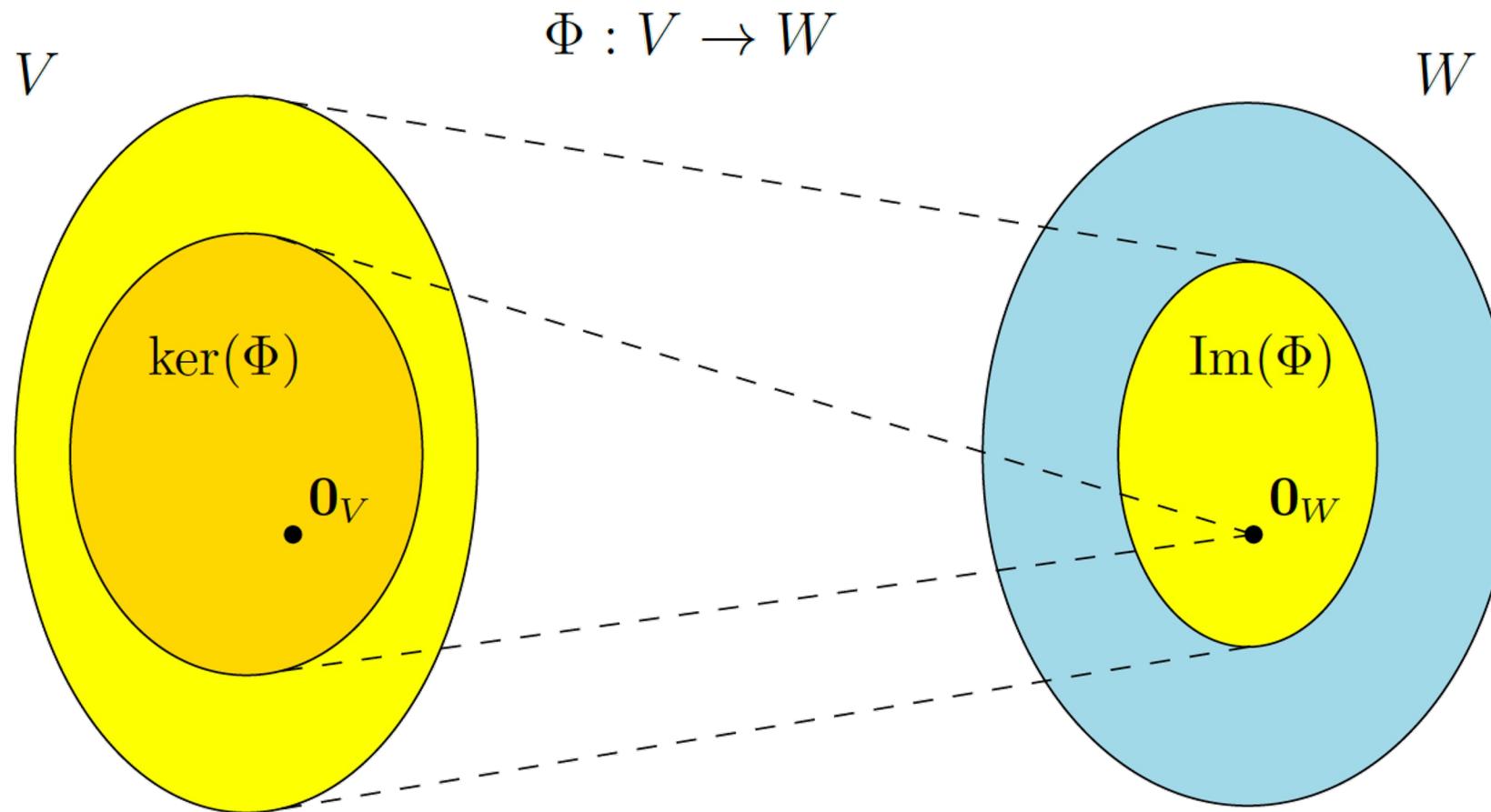
$$\tilde{A}_\Phi \in \mathbb{R}^{m \times n} \quad \tilde{B} \rightarrow \tilde{C}$$

$$\tilde{A}_\Phi = T^{-1} A_\Phi S$$



e.g., Useful for data compression while minimizing loss

Image and Kernels in Mappings



- **Kernel** (also “null space”) is the set of vectors in V that project to the “neutral” element in W
- **Image** is the set of vectors in W that can be reached from V through the mapping
- A one-to-one mapping is achievable only if the kernel is the zero vector

Eigenvectors and Eigenvalues

For $A \in \mathbb{R}^{n \times n}, \lambda \in \mathbb{R}, x \in \mathbb{R}^n \setminus \{\mathbf{0}\}$

$$Ax = \lambda x$$

eigenvalue *corresponding
eigenvector*

Equivalently, $(A - \lambda I_n)x = 0$ can be solved non-trivially

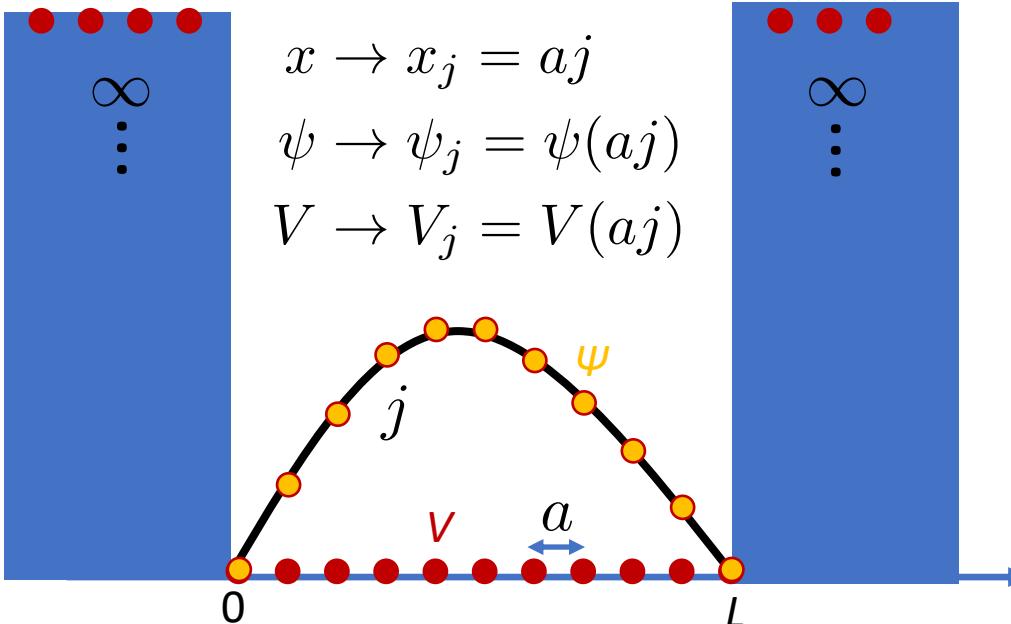
which also implies $\det(A - \lambda I_n) = 0$

(characteristic polynomial)

$$p_A(\lambda) := \det(A - \lambda I_n) = (-1)^n \lambda^n + \sum_{i=0}^{n-1} c_i \lambda^i \quad c_0 = \det(A), \\ c_{n-1} = (-1)^{n-1} \text{tr}(A)$$

This implies that λ is an eigenvalue iff it is a root of the characteristic polynomial

Eigenvectors and Eigenvalues



$$\begin{aligned}x &\rightarrow x_j = aj \\ \psi &\rightarrow \psi_j = \psi(aj) \\ V &\rightarrow V_j = V(aj)\end{aligned}$$

Using

$$\frac{d^2\psi}{dx^2} + \frac{2m\varepsilon}{\hbar^2}\psi = 0$$

$$\frac{d^2\psi(x_j)}{dx^2} \approx \frac{\psi_{j+1} - 2\psi_j + \psi_{j-1}}{a^2} \quad \text{and} \quad t \equiv \frac{\hbar^2}{2ma^2}$$

At each point j ,

$$\hat{H}\psi_j = (V_j + 2t)\psi_j - t\psi_{j-1} - t\psi_{j+1} = \varepsilon\psi_j$$

This defines a system of linear equations, and the eigenvalues/eigenvectors == energies/wavefunctions!

$$\left(\begin{array}{ccccccc} 2t + V_0 & -t & 0 & 0 & 0 & \dots \\ -t & 2t + V_1 & -t & 0 & 0 & \dots \\ 0 & -t & 2t + V_2 & -t & 0 & \dots \\ 0 & 0 & -t & 2t + V_3 & -t & \dots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{array} \right) \left(\begin{array}{c} \psi_0 \\ \psi_1 \\ \psi_2 \\ \psi_3 \\ \vdots \end{array} \right) = \varepsilon \left(\begin{array}{c} \psi_0 \\ \psi_1 \\ \psi_2 \\ \psi_3 \\ \vdots \end{array} \right)$$

$\hat{H} \rightarrow$ a matrix w/ elements H_{ij}

The elements of this matrix are given as:

$$H_{ij} = (V_i + 2t)\delta_{ij} - t\delta_{i,j-1} - t\delta_{i,j+1}$$

δ_{ij} is the Kronecker delta function

Eigenvectors and Eigenvalues

Useful properties:

- A and A^T have the same eigenvalues (but not necessarily the same eigenvectors)
- The eigenspace of A is the null space or kernel of $A - \lambda I$
- Eigenvalues are independent of choice of basis
- Symmetric, positive definite matrices always have positive, real eigenvalues
- The determinant of a matrix is the product of its eigenvalues
- The trace of a matrix is the sum of its eigenvalues

n distinct eigenvalues implies n linearly independent eigenvectors → basis

Algebraic multiplicity of λ_i : the degeneracy/number of times λ_i is a root to the characteristic polynomial

Geometric multiplicity of λ_i : how many linearly independent eigenvectors are associated with λ_i

Eigenvectors and Eigenvalues

Useful properties:

- A and A^T have the same eigenvalues (but not necessarily the same eigenvectors)
- The eigenspace of A is the null space or kernel of $A - \lambda I$
- Eigenvalues are independent of choice of basis
- Symmetric, positive definite matrices always have positive, real eigenvalues
- The determinant of a matrix is the product of its eigenvalues
- The trace of a matrix is the sum of its eigenvalues

n distinct eigenvalues implies n linearly independent eigenvectors \rightarrow basis

Algebraic multiplicity of λ_i : the degeneracy/number of times λ_i is a root to the characteristic polynomial

Geometric multiplicity of λ_i : how many linearly independent eigenvectors are associated with λ_i

Exercise: determine the algebraic and geometric multiplicity associated with these matrices

$$A_1 = \begin{bmatrix} 4 & 2 \\ 1 & 3 \end{bmatrix} \quad (\text{Jan.–Jun.}) \qquad A_2 = \begin{bmatrix} 2 & 1 \\ 0 & 2 \end{bmatrix} \quad (\text{Jul.–Dec.})$$

Matrix Decompositions

Matrix decompositions can enable faster/facile computations and are thus central to many numerical implementations underlying ML algorithms

- If A is symmetric, positive definite

Cholesky Decomposition $A = LL^T$

A is symmetric positive definite if

$$\forall x \in V \setminus \{0\} : x^\top Ax > 0$$

$$\begin{bmatrix} l_{11} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ l_{n1} & \cdots & l_{nn} \end{bmatrix}$$

Cholesky factor, a lower-triangular matrix with positive diagonal elements

- If A is square

Eigendecomposition

$$A = PDP^{-1}$$

- D is a diagonal matrix with eigenvalue entries
- P has matching eigenvectors as its columns

Spectral Theorem: If A is symmetric, there exists an orthonormal basis in the vector space of its eigenvectors, which have all real eigenvalues

$$\implies A = PDP^T$$

Matrix Decompositions

Matrix decompositions can enable faster/facile computations and are thus central to many numerical implementations underlying ML algorithms

Singular value decomposition

(*SVD*, the never-fail decomposition):

$$\begin{matrix} n \\ m \end{matrix} \boxed{A} = \begin{matrix} m \\ m \end{matrix} \boxed{U} \begin{matrix} n \\ m \end{matrix} \boxed{\Sigma} \begin{matrix} n \\ n \end{matrix} \boxed{V^\top} z$$

orthogonal matrix with vectors u_i
left-singular vectors

diagonal matrix with entries $\Sigma_{ii} = \sigma_i \geq 0$
singular values

orthogonal matrix with vectors v_i
right-singular vectors

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \sigma_n \\ 0 & \dots & 0 \\ \vdots & & \vdots \\ 0 & \dots & 0 \end{bmatrix} \quad \text{if } m > n$$

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & 0 & 0 & \dots & 0 \\ 0 & \ddots & 0 & \vdots & & \vdots \\ 0 & 0 & \sigma_m & 0 & \dots & 0 \end{bmatrix} \quad \text{if } n > m$$

Matrix Decompositions

Matrix decompositions can enable faster/facile computations and are thus central to many numerical implementations underlying ML algorithms

Singular value decomposition

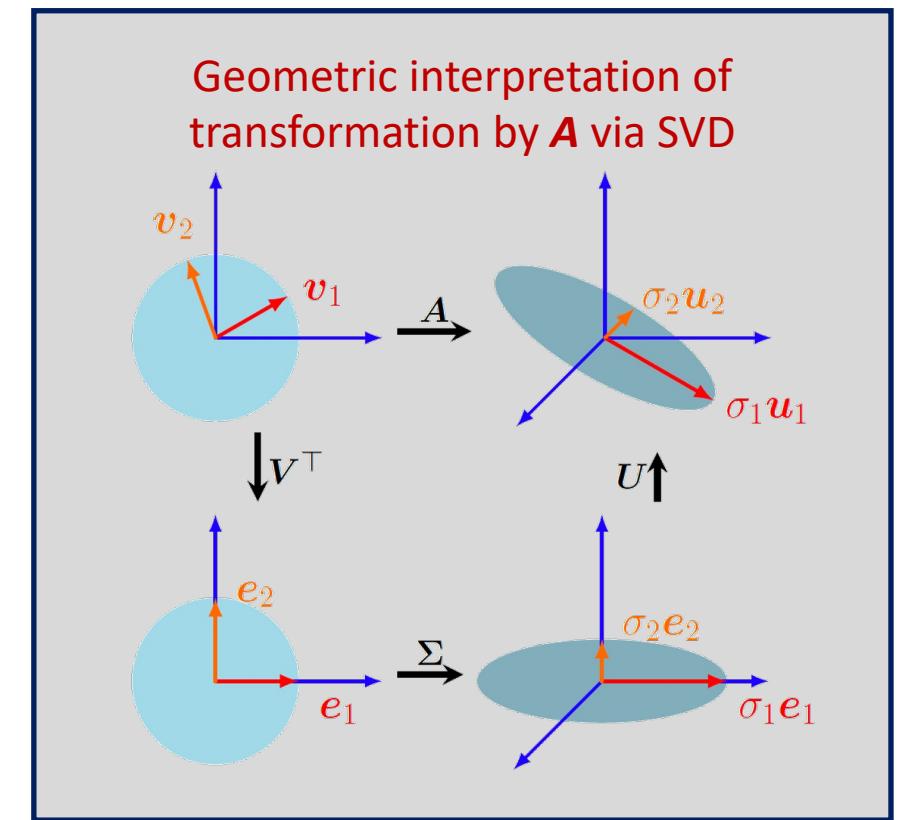
(*SVD*, the never-fail decomposition):

$$m \times n \quad A = m \times m \quad U \quad m \times n \quad \Sigma \quad m \times n \quad V^\top$$

orthogonal matrix with vectors u_i
left-singular vectors

diagonal matrix with entries $\Sigma_{ii} = \sigma_i \geq 0$
singular values

orthogonal matrix with vectors v_i
right-singular vectors



Matrix Decompositions

Matrix decompositions can enable faster/facile computations and are thus central to many numerical implementations underlying ML algorithms

CLASS EXERCISE

Review of Essential Mathematics for Machine Learning (Pt. 2)

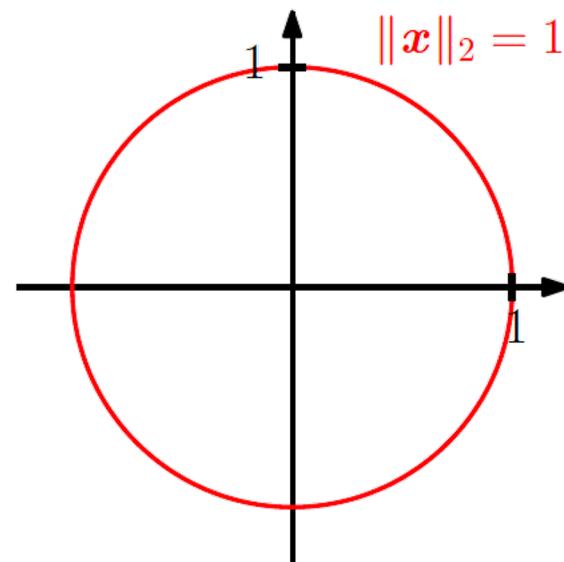
Analytic Geometry

Norms, Inner Products, and Distances

We will find it useful to define “lengths” between vectors to measure similarity → **Norm**

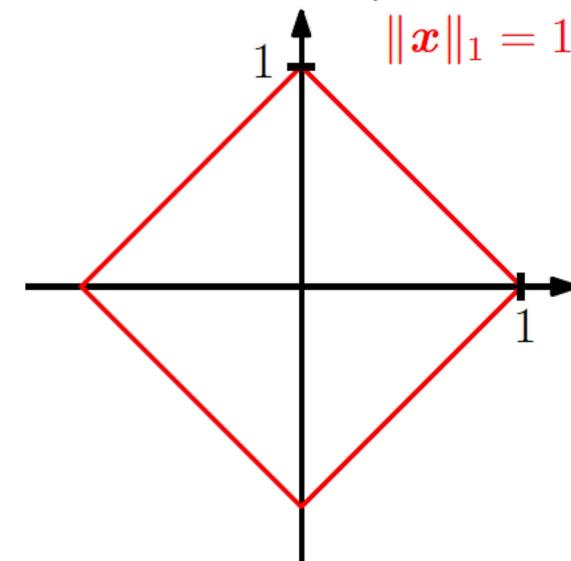
Euclidean norm (or ℓ_2 -norm)

$$\|\mathbf{x}\|_2 = \sum_i \sqrt{x_i^2} = \sqrt{\mathbf{x}^T \mathbf{x}}$$



Manhattan norm (or ℓ_1 -norm)

$$\|\mathbf{x}\|_1 = \sum_i |x_i|$$



Norms, Inner Products, and Distances

We may also like to characterize the relative orientation of two vectors → *inner product*

fancy definition: “positive definite, symmetric bilinear mapping”

$$\Omega : V \times V \rightarrow \mathbb{R}$$

If $\Omega(\lambda\mathbf{x} + \psi\mathbf{y}, \mathbf{z}) = \lambda\Omega(\mathbf{x}, \mathbf{z}) + \psi\Omega(\mathbf{y}, \mathbf{z})$ (bilinear mapping)
 $\Omega(\mathbf{x}, \lambda\mathbf{y} + \psi\mathbf{z}) = \lambda\Omega(\mathbf{x}, \mathbf{y}) + \psi\Omega(\mathbf{x}, \mathbf{z})$

and $\Omega(\mathbf{x}, \mathbf{y}) = \Omega(\mathbf{y}, \mathbf{x})$ (symmetric)

and $\forall \mathbf{x} \in V \setminus \{\mathbf{0}\} : \Omega(\mathbf{x}, \mathbf{x}) > 0, \quad \Omega(\mathbf{0}, \mathbf{0}) = 0$ (positive definite)

Then $\langle \mathbf{x}, \mathbf{y} \rangle \equiv \Omega(\mathbf{x}, \mathbf{y})$ is an inner product

Norms, Inner Products, and Distances

An inner product that you know:

$$\mathbf{x}^T \mathbf{y} = \sum_i x_i y_i$$

An example of one you probably don't:

$$\mathbf{x}, \mathbf{y} \in \mathbb{R}^2$$
$$\langle \mathbf{x}, \mathbf{y} \rangle := x_1 y_1 - (x_1 y_2 + x_2 y_1) + 2x_2 y_2$$

Inner products as matrices:

Let $\hat{\mathbf{x}} \hat{\mathbf{y}}$ be coordinates of $\mathbf{x} \mathbf{y}$ in the basis $\mathcal{B} = (\mathbf{b}_1, \dots, \mathbf{b}_n)$

If A , s.t. $A_{ij} := \langle \mathbf{b}_i, \mathbf{b}_j \rangle$ is symmetric, positive definite

Then $\langle \mathbf{x}, \mathbf{y} \rangle = \hat{\mathbf{x}}^T A \hat{\mathbf{y}}$

Norms, Inner Products, and Distances

Any *inner product* also defines a *norm*:

$$\|x\| := \sqrt{\langle x, x \rangle}$$

The *distance (metric)* between two vectors is

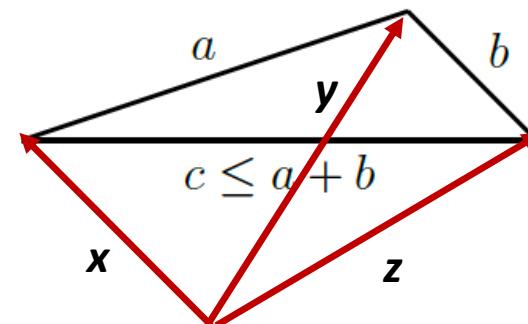
$$d(x, y) := \|x - y\| = \sqrt{\langle x - y, x - y \rangle}$$

Cauchy-Schwarz Inequality

$$|\langle x, y \rangle| \leq \|x\| \|y\|$$

Triangle Inequality

$$d(x, z) \leq d(x, y) + d(y, z)$$



Norms, Inner Products, and Distances

Exercise:

Compute the distance metric between the two vectors

using

dot product

vs.

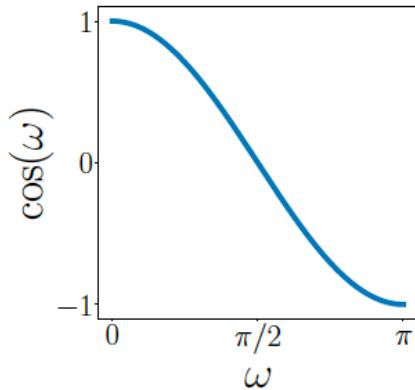
$$\mathbf{v}_1 = [3, 2]^T$$
$$\mathbf{v}_2 = [2, 1]^T$$
$$A = \begin{bmatrix} 1 & -\frac{1}{2} \\ -\frac{1}{2} & 1 \end{bmatrix}$$

Norms, Inner Products, and Distances

Inner products also enable definition of an angle between vectors

Following Cauchy-Schwarz inequality,

$$-1 \leq \frac{\langle x, y \rangle}{\|x\| \|y\|} \leq 1$$



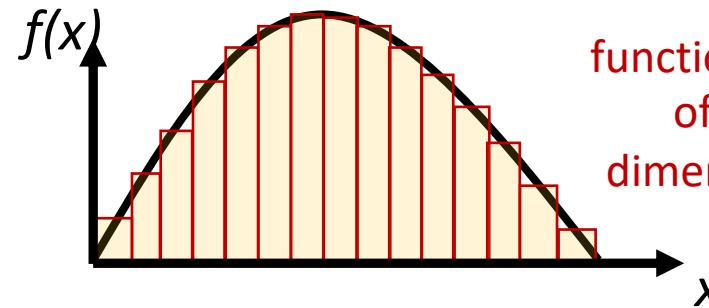
*there is a unique
mapping to an angle ω
on the interval $[0, \pi]$*

$$\cos \omega = \frac{\langle x, y \rangle}{\|x\| \|y\|}$$

orthogonality
 $\langle x, y \rangle = 0$

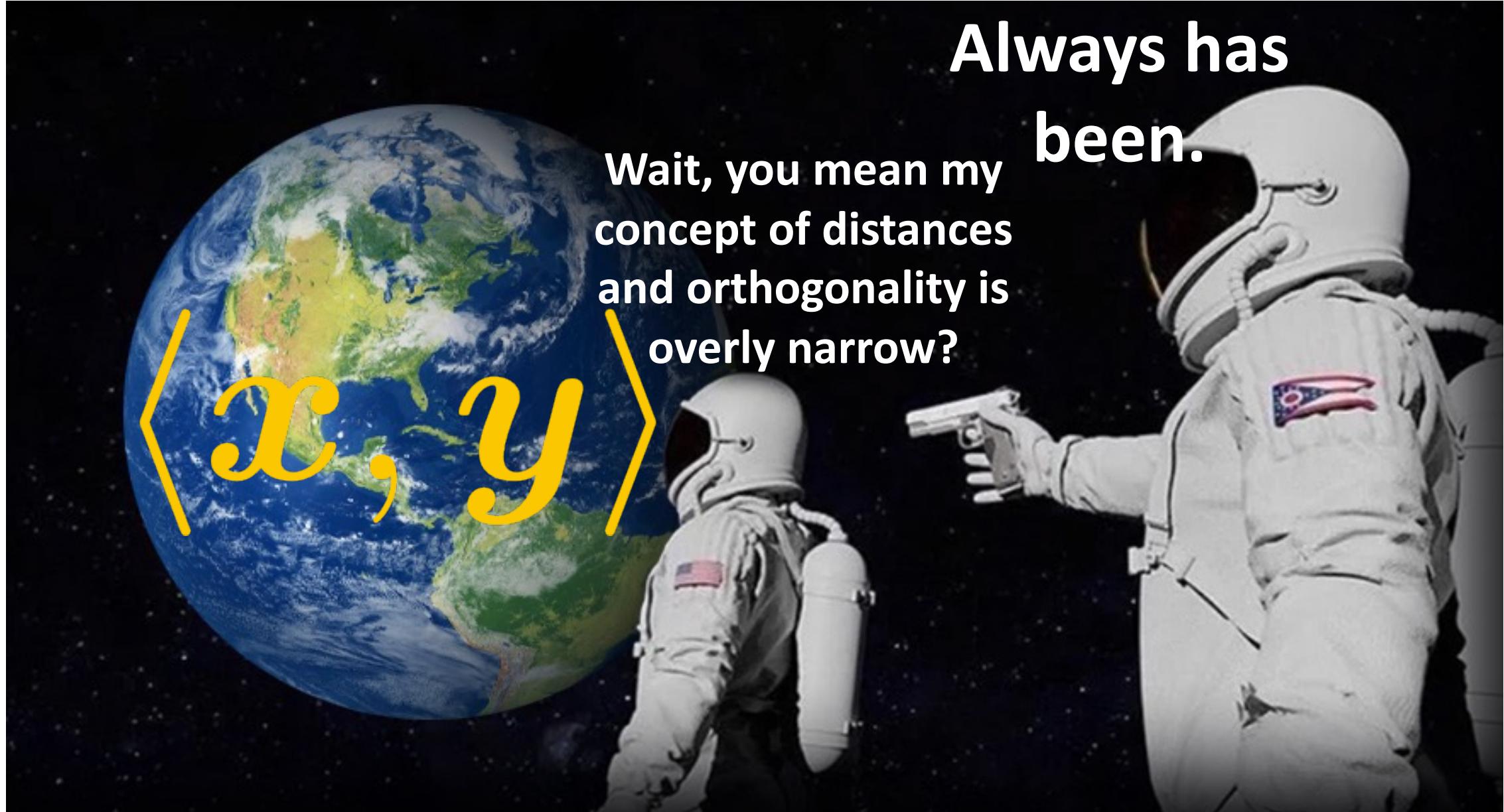
We can also define inner products over functions

$$\langle u, v \rangle := \int_a^b u(x)v(x)dx$$



*functions can be analogously thought
of as vectors with an infinite-
dimensional basis of delta functions*

Norms, Inner Products, and Distances



Orthogonal Bases/Complements

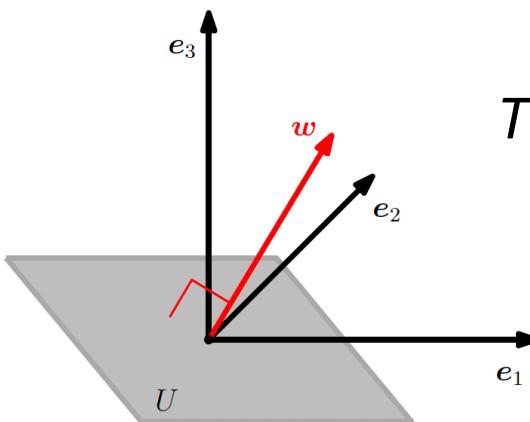
$$\mathcal{B} = (\mathbf{b}_1, \dots, \mathbf{b}_n) \text{ is orthonormal if } \langle \mathbf{b}_i, \mathbf{b}_j \rangle = \delta_{ij} \Leftrightarrow \mathbf{B}^T \mathbf{B} = \mathbf{I}_n$$

Suppose $\tilde{\mathcal{B}} = (\tilde{\mathbf{b}}_1, \dots, \tilde{\mathbf{b}}_n)$ is not necessarily orthogonal and/or not normalized

Gram-Schmidt Process enables orthogonalization by Gaussian elimination on $[\tilde{\mathbf{B}} \tilde{\mathbf{B}}^T | \tilde{\mathbf{B}}]$

(examine QR decomposition, `np.linalg.qr`)

For an N -dimensional vector space V with an M -dimensional subspace U , the orthogonal complement of U^\perp is a $(N-M)$ -dimensional subspace of V containing all vectors in V that are orthogonal to every vector in U (U^\perp describes a hyperplane)



This enables a unique decomposition:

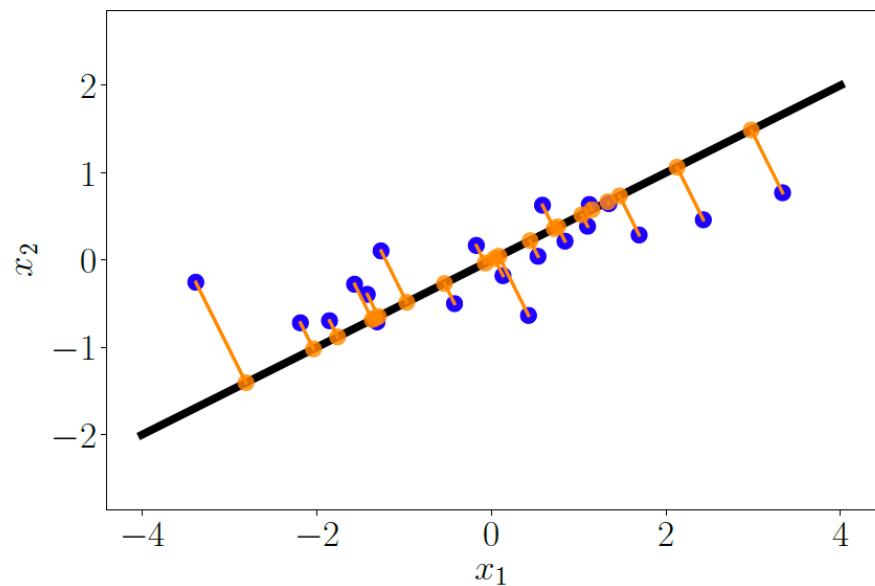
$$\mathbf{v} = \sum_{i=1}^M \lambda_i \mathbf{b}_i + \sum_{j=1}^{N-M} \psi_j \mathbf{b}_j^\perp$$

Projections

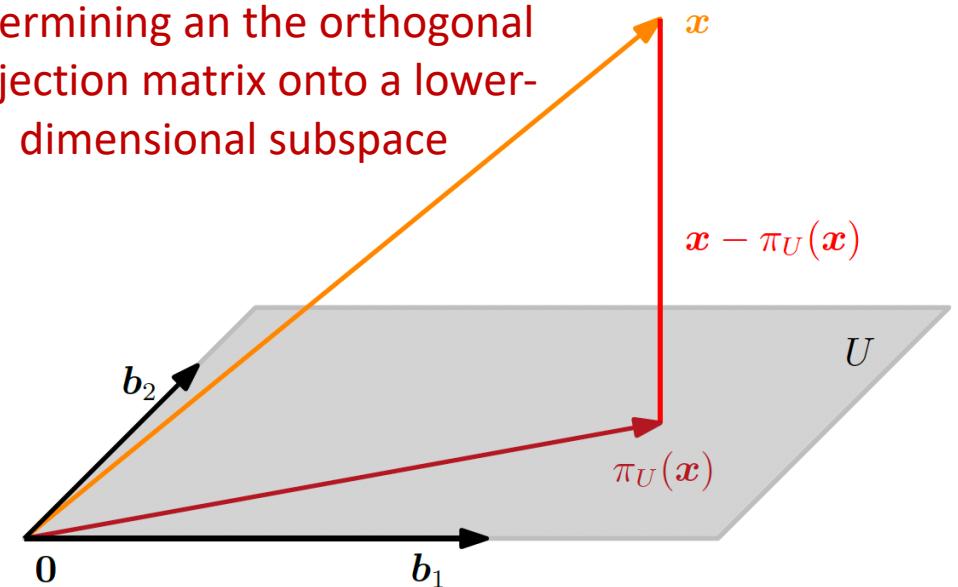
A **projection** is a linear mapping between vector spaces (V and U , a subspace of V) with the property that

$$P_{\pi}^2 = P_{\pi} \leftarrow \text{projection matrix}$$

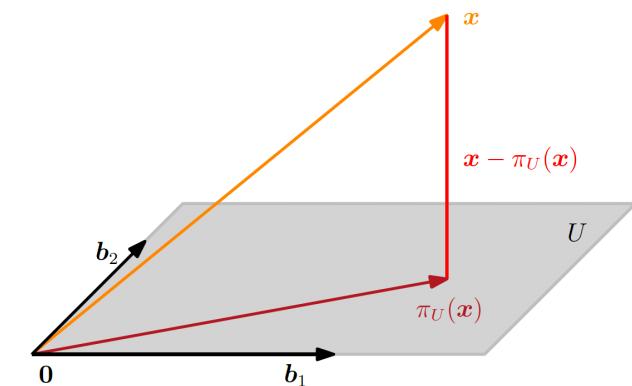
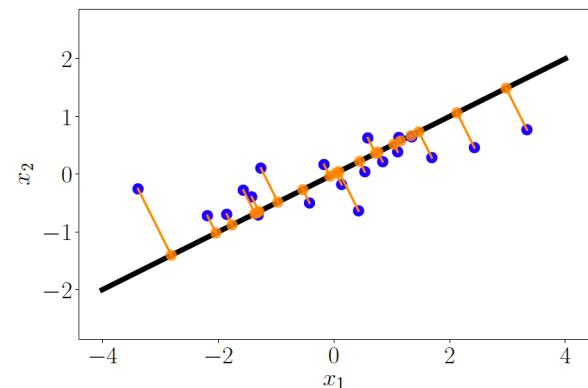
projection of a 2d-dataset
onto a 1d-subspace



An oft-encountered task may be determining an orthogonal projection matrix onto a lower-dimensional subspace



Projections



We want $\pi_U(x) = \sum_{i=1}^M \lambda_i b_i = B\lambda$ given $\mathcal{B}_U = (\mathbf{b}_1, \dots, \mathbf{b}_M)$
such that $\pi_U(x) - x$ is **orthogonal** to U and its **distance is minimized**

Assuming the dot product as the inner product...

$$\begin{aligned} &\implies \mathbf{b}_i^T(\mathbf{x} - B\lambda) = 0, \quad i = 1, \dots, M \\ &\Downarrow \\ &B^T(\mathbf{x} - B\lambda) = \mathbf{0} \iff B^T B \lambda = B^T \mathbf{x} \end{aligned}$$

normal equation

$$\lambda = (B^T B)^{-1} B^T \mathbf{x}$$

if B describes an orthonormal basis??

$$\implies P_\pi = B(B^T B)^{-1} B^T$$

Projections

Example

$$x = \begin{bmatrix} 6 \\ 0 \\ 0 \end{bmatrix}$$
$$U = \text{span} \left[\begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 2 \end{bmatrix} \right] \subseteq \mathbb{R}^3$$

```
>>> x = np.array([[6,0,0]]).transpose()
>>> P@x
array([[ 5.],
       [ 2.],
      [-1.]])
```

Projection Error

$$\|x - \pi_U(x)\| = \left\| \begin{bmatrix} 1 & -2 & 1 \end{bmatrix}^\top \right\| = \sqrt{6}$$

```
>>> dx = x - P@x
>>> print(dx.transpose()@B)
[[-2.22044605e-16 -4.44089210e-16]]
```

$$B = \begin{bmatrix} 1 & 0 \\ 1 & 1 \\ 1 & 2 \end{bmatrix} \quad P_\pi = B(B^T B)^{-1} B^T$$

```
>>> b1 = np.array([[1,1,1]]).transpose()
>>> b2 = np.array([[0,1,2]]).transpose()
>>> B = np.hstack([b1,b2])
>>> print(B)
[[1 0]
 [1 1]
 [1 2]]
>>> BT = B.transpose()
>>> print(BT)
[[1 1 1]
 [0 1 2]]
>>> P = B@np.linalg.inv(BT@B)@BT
>>> print(P)
[[ 0.83333333  0.33333333 -0.16666667]
 [ 0.33333333  0.33333333  0.33333333]
 [-0.16666667  0.33333333  0.83333333]]
```

```
>>> print(P@P)
[[ 0.83333333  0.33333333 -0.16666667]
 [ 0.33333333  0.33333333  0.33333333]
 [-0.16666667  0.33333333  0.83333333]]
```

Other remarks on Projections

Recall that for any given system of linear equations, $\mathbf{A}\mathbf{x} = \mathbf{b}$ we may have

no solution

one unique solution

infinitely many solutions

It is rare that any realistic situation is exactly solvable

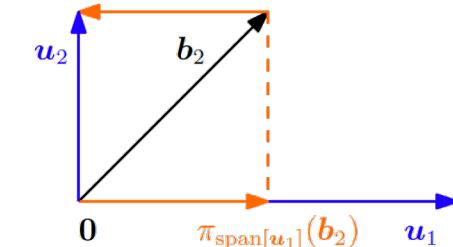
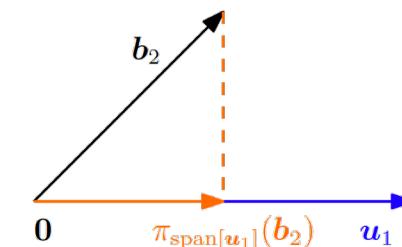
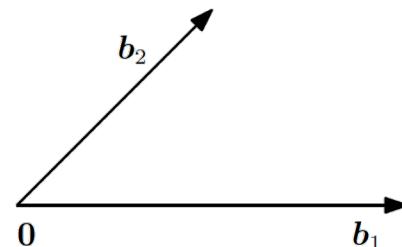
Interpretation: \mathbf{b} does not lie in the span of \mathbf{A}

Solution: find the closest projection of \mathbf{b} onto the subspace spanned by \mathbf{A} “**Least-squares**”

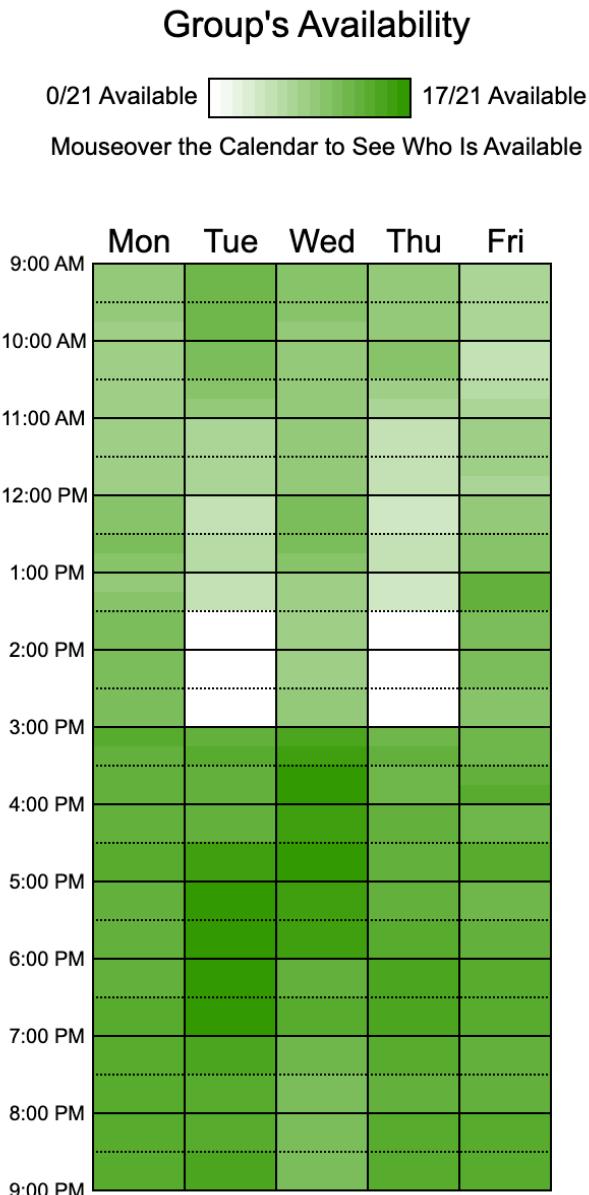
Gram-Schmidt orthogonalization constructs an orthogonal basis $(\mathbf{u}_1, \dots, \mathbf{u}_n)$ from any basis $(\mathbf{b}_1, \dots, \mathbf{b}_n)$ by iteratively determining projections on incrementally larger subspaces

$$\mathbf{u}_1 := \mathbf{b}_1$$

$$\mathbf{u}_k := \mathbf{b}_k - \pi_{\text{span}[\mathbf{u}_1, \dots, \mathbf{u}_{k-1}]}(\mathbf{b}_k), \quad k = 2, \dots, n$$



Office Hours



- **Slot 1: Tuesday, 5:00-6:30pm**

Proposed Slot 2

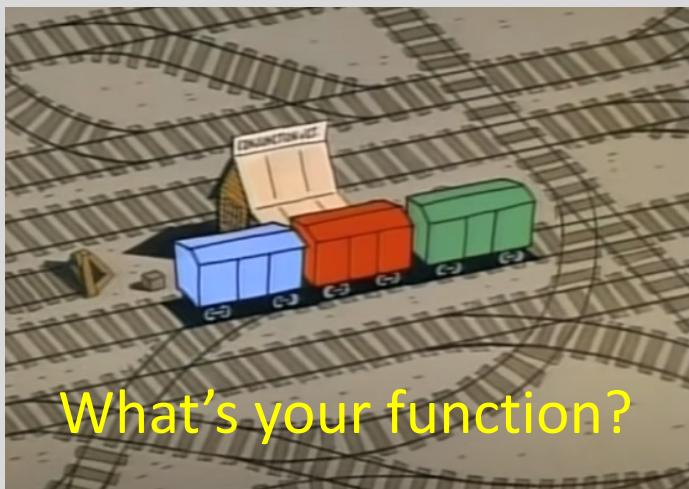
- Wednesday, 5:00-6:00pm
- Friday, 4:00-5:00pm
- Friday, 5:00-6:00pm

Review of Essential Mathematics for Machine Learning (Pt. 3)

Vector Calculus & Probability

Optimization and Gradients in Machine Learning

Many ML algorithms depend on optimization algorithms; such algorithms typically exploit gradient information with respect to some prescribed objective function



with a **gradient** expressed as

$$\nabla_{\mathbf{x}} f = \text{grad}f = \frac{df}{d\mathbf{x}} = \begin{bmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} & \frac{\partial f(\mathbf{x})}{\partial x_2} & \dots & \frac{\partial f(\mathbf{x})}{\partial x_n} \end{bmatrix} \in \mathbb{R}^{1 \times n}$$

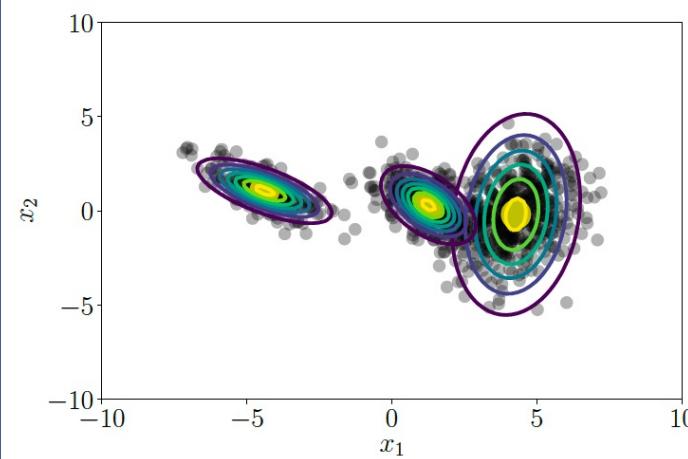
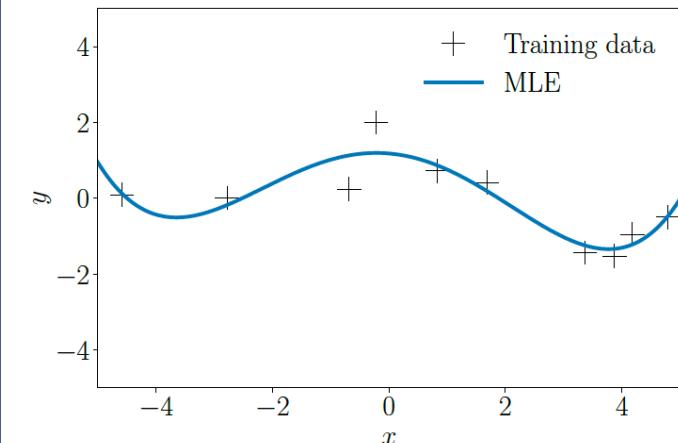
$f : \mathbb{R}^n \rightarrow \mathbb{R}, \mathbf{x} \mapsto f(\mathbf{x})$
a **function** f is a mapping that
explicitly assigns every input \mathbf{x}
exactly one value $f(\mathbf{x})$

Partial derivatives are defined

$$\frac{\partial f}{\partial x_1} = \lim_{h \rightarrow 0} \frac{f(x_1 + h, x_2, \dots, x_n) - f(x)}{h}$$

⋮

$$\frac{\partial f}{\partial x_n} = \lim_{h \rightarrow 0} \frac{f(x_1, \dots, x_{n-1}, x_n + h) - f(x)}{h}$$



Partial Differentiation of functions

Basic differentiation rules still apply with some care

product rule

$$\frac{\partial}{\partial x}(f(x)g(x)) = \frac{\partial f}{\partial x}g(x) + f(x)\frac{\partial g}{\partial x}$$

sum rule

$$\frac{\partial}{\partial x}(f(x) + g(x)) = \frac{\partial f}{\partial x} + \frac{\partial g}{\partial x}$$

chain rule

$$\frac{\partial}{\partial x}(g(f(x))) = \frac{\partial g}{\partial f}\frac{\partial f}{\partial x}$$

Application of chain rule

Suppose $f(\mathbf{x}), \mathbf{x}(t)$

$$\frac{df(\mathbf{x}(t))}{dt} = \frac{\partial f}{\partial \mathbf{x}} \frac{\partial \mathbf{x}}{\partial t} = \left[\frac{\partial f}{\partial x_1} \dots \frac{\partial f}{\partial x_n} \right] \begin{bmatrix} \frac{\partial x_1}{\partial t} \\ \vdots \\ \frac{\partial x_n}{\partial t} \end{bmatrix}$$

Extension of chain rule for gradient

Suppose now $f(\mathbf{x}), \mathbf{x}(s, t)$

$$\frac{df}{d(s, t)} = \frac{\partial f}{\partial \mathbf{x}} \frac{\partial \mathbf{x}}{\partial (s, t)} = \left[\frac{\partial f}{\partial x_1} \dots \frac{\partial f}{\partial x_n} \right] \begin{bmatrix} \frac{\partial x_1}{\partial s} & \frac{\partial x_1}{\partial t} \\ \vdots & \vdots \\ \frac{\partial x_n}{\partial s} & \frac{\partial x_n}{\partial t} \end{bmatrix}$$

$1 \times n$ $n \times 2$

Note that we could compute via matrix multiplication

Gradients of Vector-valued functions (vector fields)

We will also encounter vector-valued functions $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$

Rules for operating on each of the functions are precisely the same as just shown

$$f(\mathbf{x}) = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_m(\mathbf{x}) \end{bmatrix} \quad \text{simply a vector of functions}$$

The gradient of a vector field yields a matrix

$$\frac{df(\mathbf{x})}{d\mathbf{x}} = \left[\boxed{\frac{\partial f(\mathbf{x})}{\partial x_1}} \cdots \boxed{\frac{\partial f(\mathbf{x})}{\partial x_n}} \right]$$
$$= \left[\begin{array}{c|c} \boxed{\frac{\partial f_1(\mathbf{x})}{\partial x_1}} & \cdots & \boxed{\frac{\partial f_1(\mathbf{x})}{\partial x_n}} \\ \vdots & & \vdots \\ \boxed{\frac{\partial f_m(\mathbf{x})}{\partial x_1}} & \cdots & \boxed{\frac{\partial f_m(\mathbf{x})}{\partial x_n}} \end{array} \right]$$

This $m \times n$ matrix of first-order partial derivatives is often referred to as the **Jacobian**, J

$$J_{ij} = \frac{\partial f_i}{\partial x_j}$$

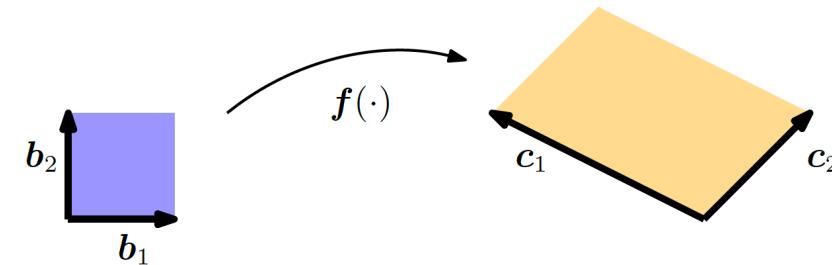
Note: the Jacobian can represent a coordinate transformation (specifically from a coordinate representation in \mathbf{x} to a coordinate representation in \mathbf{f})

Gradients of Vector-valued functions (vector fields)

Example

$$\mathbf{b}_1 = [1, 0]^T; \mathbf{b}_2 = [0, 1]^T$$

$$\mathbf{c}_1 = [-2, 1]^T; \mathbf{c}_2 = [1, 1]^T$$



What's the transformation matrix for changing basis from B to C (\mathbf{x} to f)?

$$\mathbf{T}_f = \begin{bmatrix} \alpha_{11} & \alpha_{12} \\ \alpha_{21} & \alpha_{22} \end{bmatrix}$$

T has columns of the coordinate mapping from one basis to the other

$$\begin{bmatrix} -2 & 1 \\ 1 & 1 \end{bmatrix} \quad f_1(\mathbf{x}) = \alpha_{11}x_1 + \alpha_{12}x_2$$
$$f_2(\mathbf{x}) = \alpha_{21}x_1 + \alpha_{22}x_2$$

$$\mathbf{J} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} \end{bmatrix}$$

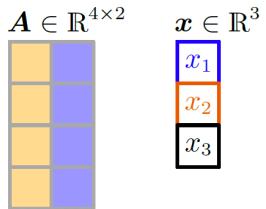
So, the Jacobian is identical to the linear transformation matrix in this case; advantage of Jacobian is for non-linear transformations!

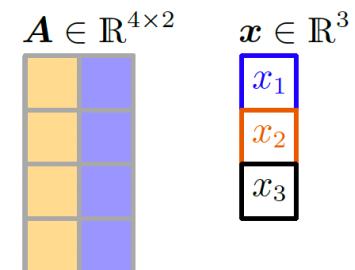
Gradients of Matrices

Gradients of matrices with respect to vectors (or other matrices) \rightarrow tensors

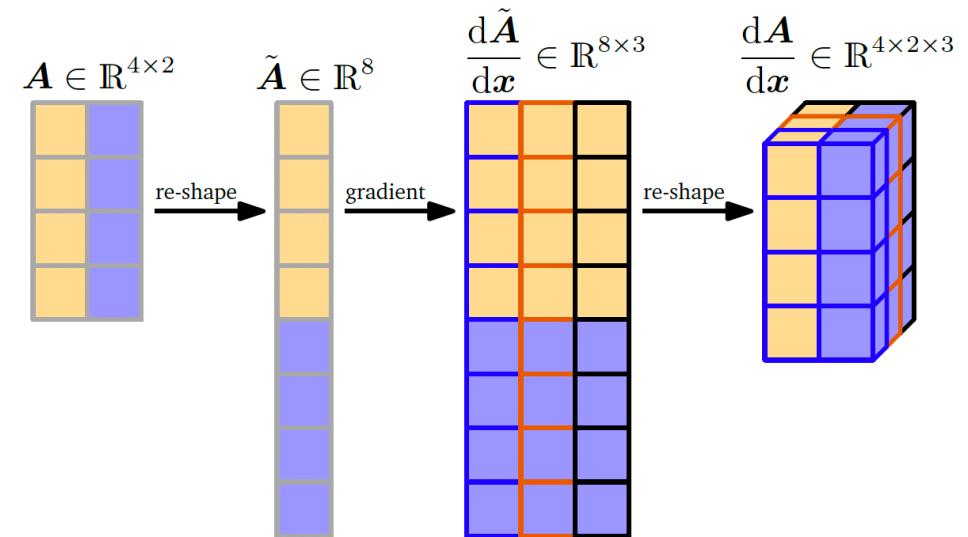
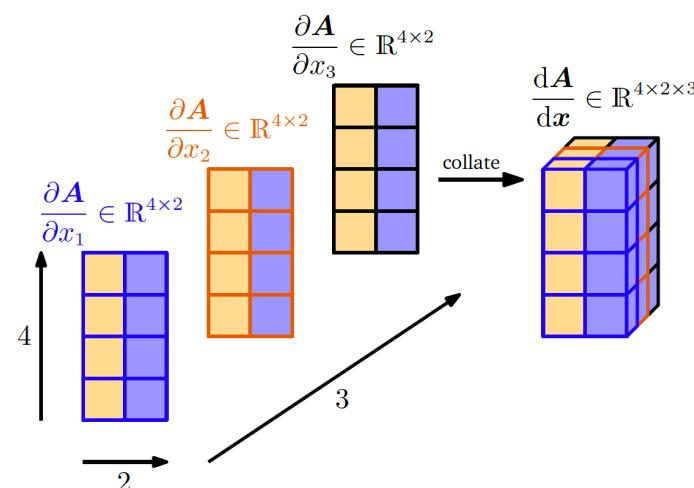
$$A \in \mathbb{R}^{m \times n}; B \in \mathbb{R}^{p \times q} \implies J \in \mathbb{R}^{(m \times n) \times (p \times q)} \quad J_{ijkl} = \frac{\partial A_{ij}}{\partial B_{kl}}$$

Two ways to think about (and compute) it:

$$A \in \mathbb{R}^{4 \times 2} \quad x \in \mathbb{R}^3$$


$$A \in \mathbb{R}^{4 \times 2} \quad x \in \mathbb{R}^3$$


Partial derivatives:



Higher-order Derivatives and the Hessian

2nd derivatives of f

$$\frac{\partial^2 f}{\partial x^2} = \frac{\partial}{\partial x} \left(\frac{\partial f}{\partial x} \right)$$

$$\frac{\partial^2 f}{\partial x \partial y} = \frac{\partial}{\partial x} \left(\frac{\partial f}{\partial y} \right)$$

$$\frac{\partial^2 f}{\partial y \partial x} = \frac{\partial}{\partial y} \left(\frac{\partial f}{\partial x} \right)$$

$$\frac{\partial^2 f}{\partial y \partial x} = \frac{\partial^2 f}{\partial x \partial y}$$

if $f(x,y)$ is continuous and twice differentiable

n th derivative of f with respect to x

$$\frac{\partial^n f}{\partial x^n} = \frac{\partial}{\partial x} \left(\frac{\partial^{n-1} f}{\partial x^{n-1}} \right) D_{\mathbf{x}}^n f$$

The **Hessian** is the second derivative analogue of the Jacobian (it contains all the second derivatives)

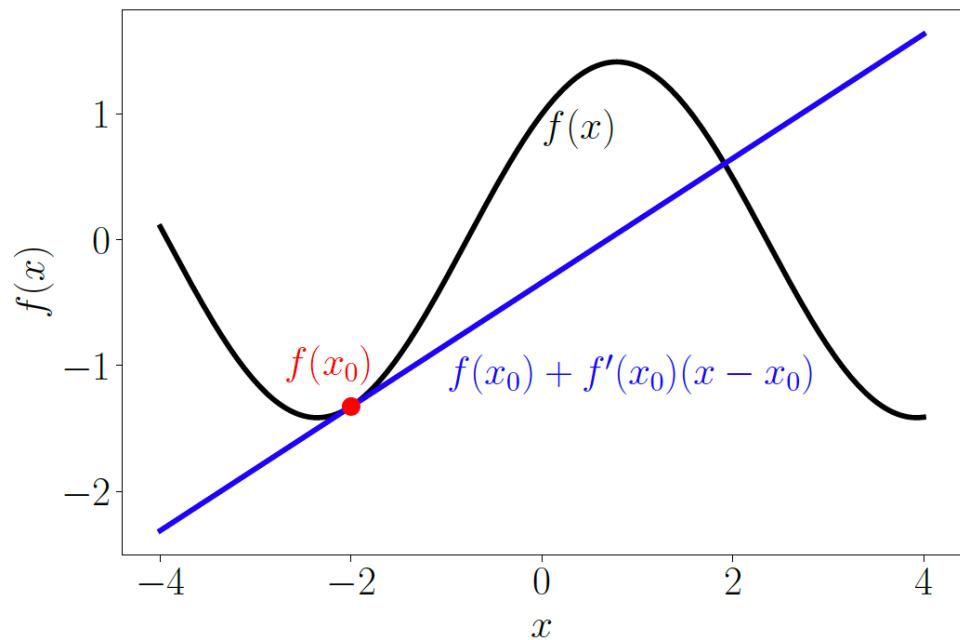
$$\mathbf{H} = \nabla_{\mathbf{x}}^2 f(\mathbf{x}) \quad H_{ij} = \frac{\partial^2 f}{\partial x_i \partial x_j}$$

Here, \mathbf{H} is an $n \times n$, matrix

Question: What would be the corresponding form for the Hessian if applied to a vector field?

$$\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

Linearization and Taylor Expansions



$$\delta \equiv x - x_0$$

$$\delta^2 = \delta \otimes \delta$$

$$\delta^3 = \delta \otimes \delta \otimes \delta$$

\otimes is the **outer product**

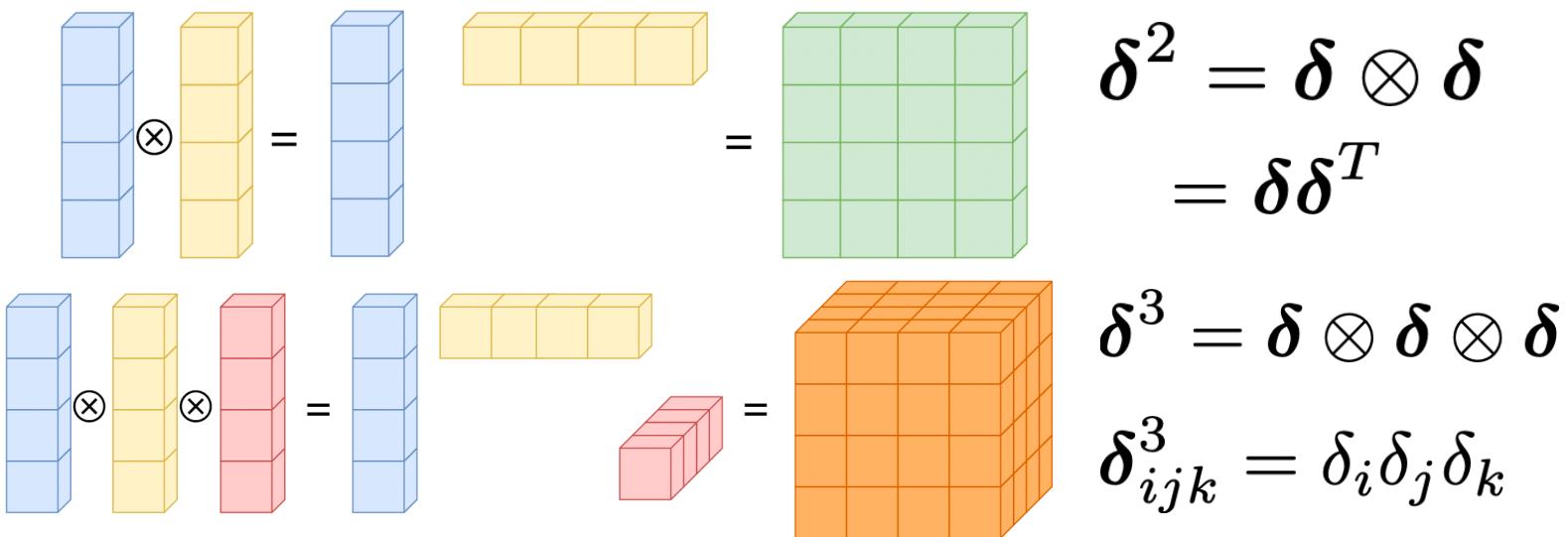
Gradients enable local approximations to functions

Linear: $f(\mathbf{x}) \approx f(\mathbf{x}_0) + \frac{df(\mathbf{x}_0)}{d\mathbf{x}}(\mathbf{x} - \mathbf{x}_0)$

General: $f(\mathbf{x}) = \sum_{k=0}^{\infty} \frac{D_{\mathbf{x}}^k f(\mathbf{x}_0)}{k!} (\mathbf{x} - \mathbf{x}_0)^k$

kth-order tensor

easiest to understand visually (at least up to 3D):



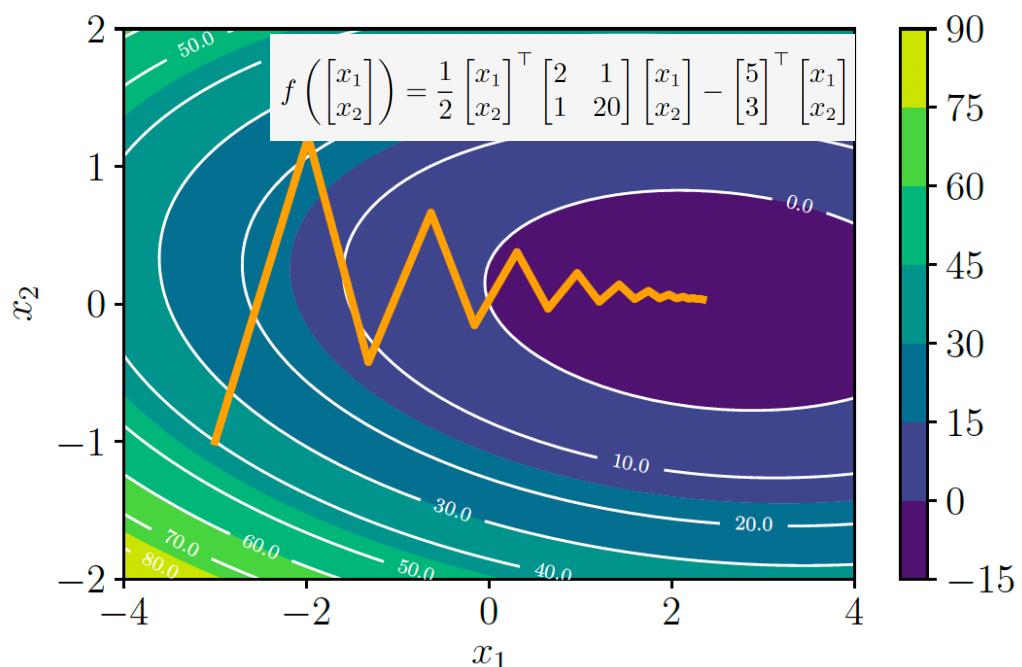
Optimization

Essential task: $f : \mathbb{R}^n \rightarrow \mathbb{R}, \mathbf{x} \mapsto f(\mathbf{x})$

$$\min_{\mathbf{x}} f(\mathbf{x})$$

Gradient Descent

$$\mathbf{x}_{i+1} = \mathbf{x}_i - \gamma_i [\nabla f(\mathbf{x}_i)]^T$$



- very simple method
- depends on (adaptive) stepsize
- slowly convergent to closest minima

with momentum

$$\mathbf{x}_{i+1} = \mathbf{x}_i - \gamma_i [\nabla f(\mathbf{x}_i)]^T + \alpha \Delta \mathbf{x}_i$$

$$\Delta \mathbf{x}_i = \alpha \Delta \mathbf{x}_{i-1} - \gamma_{i-1} [\nabla f(\mathbf{x}_{i-1})]^T$$

- uses “memory” to reduce jitter

stochastic

$$f(\mathbf{x}) = \sum_{k=1}^N f_k(\mathbf{x})$$

$$\mathbf{x}_{i+1} = \mathbf{x}_i - \gamma_i \sum_{k \in \mathbb{N}: k \leq N} [\nabla f_k(\mathbf{x}_i)]^T$$

- useful for large N , which may not be atypical in machine learning applications

Constrained Optimization

Task: $f : \mathbb{R}^n \rightarrow \mathbb{R}, \mathbf{x} \mapsto f(\mathbf{x}) \quad \min_{\mathbf{x}} f(\mathbf{x}) \quad g_k(\mathbf{x}) \leq 0 \text{ for } k = 1, \dots, m$

Naïve approach: Add some penalty function that induces large function value when constraint(s) not satisfied

Lagrange multipliers: $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \boldsymbol{\lambda}^T \mathbf{g}(\mathbf{x})$

Duality: exchange optimization in one set of variables (primal variables, \mathbf{x}) for a different optimization problem in another set of variables (dual variables, $\boldsymbol{\lambda}$)

$$\rightarrow \max_{\boldsymbol{\lambda} \geq 0} \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})$$

*unconstrained
optimization for fixed $\boldsymbol{\lambda}$
concave*

Probability: Basic Definitions and Rules

Probability equates to a quantified notion of uncertainty → some element of randomness

A **random variable** X is a *function* that maps outcomes of random experiments to a result

These constitute
a so-called
**probability
space**

- Ω **sample (state) space**: set of all possible outcomes from an “experiment”
- \mathcal{A} **event space**: the space of potential results from an experiment
- P **probability**: for an event A in the event space, $P(A)$ is a numerical measure of the belief/frequency that A will occur

T

Typically, we are not so interested in the probability space but rather **a target space** where outcomes are mapped to quantity of interest via the random variable

e.g.,

sample space

ABBAABABABABA

target space

the number of B monomers

Probability: Basic Definitions and Rules

Probability equates to a quantified notion of uncertainty → some element of randomness

The complement of A being anything but A

$$P(A^c) = 1 - P(A)$$

Probability of the union of two events

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

A, B are mutually exclusive if

$$P(A \cap B) = 0$$

probability of A occurring given that B has occurred (conditional probability)

If A and B are independent, then

$$P(A \cap B) = P(A)P(B)$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

Probability: Basic Definitions and Rules

Contrived example: grabbing two colored marbles out of a bag with replacement



Q: How many red marbles do you pull?

$$\Omega \rightarrow \{(R,R), (Y,Y), (R,Y), (Y,R)\}$$

state space

$$X((R,R)) = 2, \quad X((R,Y)) = 1, \quad X((Y,R)) = 1, \quad X((Y,Y)) = 0$$

random variable

$$T \rightarrow \{0, 1, 2\}$$

target space

probability of pulling two reds

$$P(X = 2) = P((R,R)) = P(R \cap R) = P(R)P(R) = 0.7^2$$

probability of pulling one red

$$P(X = 1) = P((R,Y) \cup (Y,R)) = P((R,Y)) + P((Y,R)) = 2P(R)P(Y) = 2P(R)(1 - P(R)) = 2(0.7)(0.3)$$

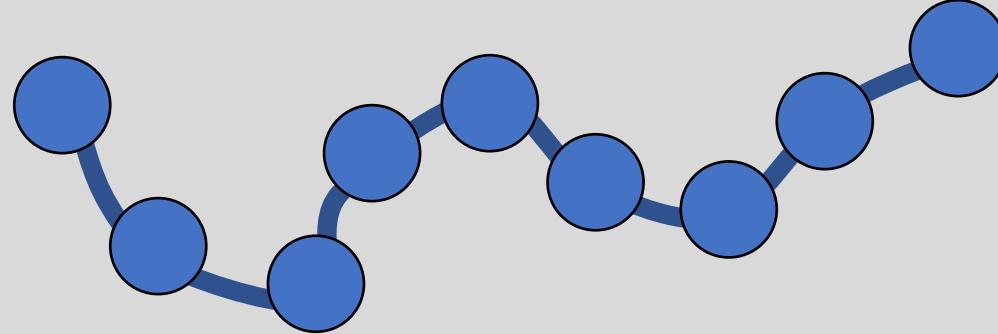
probability of pulling no reds

$$P(X = 0) = P((Y,Y)) = (1 - P(R))(1 - P(R)) = 0.3^2$$

Although we expressed probabilities in terms of the probability space, notice that the random variable X exhibits its own distribution function, which we can denote P_X

Probability: Basic Definitions and Rules

Exercise: Polymer of length N with probability p of bond scission at each bond over interval of time t



Q: How many chain segments
there are after time t ?

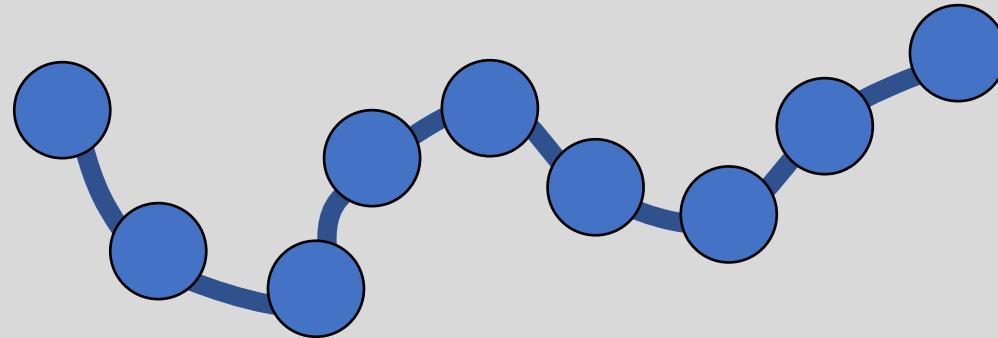
What is the sample/state space?

What is the target space and random variable?

What is the probability of observing 2 segments?

Probability: Basic Definitions and Rules

Exercise: Polymer of length N with probability p of bond scission at each bond over interval of time t



Q: How many chain segments there are after time t ?

What is the probability of observing 2 segments?

Need precisely 1 of the $n-1$ bonds to break, and this can happen for any $n-1$ of the bonds...

$$(n - 1)p(1 - p)^{n-2}$$

What is the sample/state space?

$$\nu_i = (b_1^{(i)}, \dots, b_{n-1}^{(i)}) \quad b_{n-1}^{(i)} = \begin{cases} 1 & \text{if bond exists} \\ 0 & \text{else} \end{cases}$$

$$\Omega = \{\nu_i\}$$

What is the target space and random variable?

$$X(\nu_i) = n - \sum_{j=1}^{n-1} b_j^{(i)} \in [1, n]$$

Here, our target space is finite or countably infinite, which makes X a discrete random variable; X can also be continuous

Discrete and Continuous Probabilities

If target space is ***discrete***, probability is defined for any given random variable:

$$P(X = a) = p_a \quad \text{probability mass function}$$

$$\sum_X P(X) = 1 \quad \text{normalization condition}$$

If the target space is ***continuous***, it makes more sense to consider intervals:

$$P(a \leq X \leq b) = \int_a^b f(x)dx \quad \text{probability density function (PDF)}$$

$$\int_{\mathbb{R}} f(x)dx = 1 \quad \text{normalization condition}$$

$$F_X(x) = P(X < x) = \int_{-\infty}^x f(x')dx' \quad \text{cumulative distribution function (CDF)}$$

Revisiting Rules

Suppose we have a random variables \mathbf{x}, \mathbf{y} with joint probability distribution $p(\mathbf{x}, \mathbf{y})$

Sum rule:

$$p(\mathbf{x}) = \sum_{\mathbf{y} \in \mathcal{Y}} p(\mathbf{x}, \mathbf{y}) \quad p(\mathbf{x}) = \int_{\mathcal{Y}} p(\mathbf{x}, \mathbf{y}) d\mathbf{y}$$

marginal distribution

Product rule:

$$p(\mathbf{x}, \mathbf{y}) = p(\mathbf{y}|\mathbf{x})p(\mathbf{x})$$

Bayes' theorem/rule:

Idea is to make some inference regarding an unobserved random variable based on observed values of another

$$p(\mathbf{x}|\mathbf{y}) = \frac{\underset{\text{likelihood}}{p(\mathbf{y}|\mathbf{x})} \underset{\text{prior}}{p(\mathbf{x})}}{\underset{\text{evidence}}{p(\mathbf{y})}}$$

posterior – what we want (to know something about \mathbf{x} having observed \mathbf{y})

prior – subjective presumption about the distribution of \mathbf{x} (must be nonzero for all plausible \mathbf{x})

likelihood – probability of observing \mathbf{y} **were we to know \mathbf{x}**

evidence or marginal likelihood

$$p(\mathbf{y}) := \int p(\mathbf{y}|\mathbf{x})p(\mathbf{x})d\mathbf{x}$$

Revisiting Rules

COVID-19 RAPID REPORT

Clinical Medicine 2021 Vol 21, No 1: e54–6

The impact of false positive COVID-19 results in an area of low prevalence

Authors: Brendan Healy,^A Azizah Khan,^B Huria Metezaï,^B Ian Blyth^C and Hibo Asad^D

During the same period, there were 5,079 negative tests; thus a total of 5,110 tests were carried out. The false positive rate (single gene low level positive) in this period was 0.5% (26/5110), giving a specificity of 99.5% (5,079/5,105). Over the same period, the true

PLOS ONE

Real-life clinical sensitivity of SARS-CoV-2 RT-PCR test in symptomatic patients

Elisa Kortela^{1*}, Vesa Kirjavainen^{2*}, Maarit J. Ahava², Suvi T. Jokiranta³, Anna But⁴, Anna Lindahl¹, Anu E. Jääskeläinen², Annemarjut J. Jääskeläinen², Asko Järvinen¹, Pia Jokela², Hannimari Kallio-Kokko², Raisa Loginov², Laura Mannonen², Eeva Ruotsalainen¹, Tarja Sironen^{6,7}, Olli Vapalahti^{2,6,7}, Maija Lappalainen², Hanna-Riikka Kreivi⁵, Hanna Jarva^{2,3}, Satu Kurkela¹, Eliisa Kekäläinen^{2,3†}

All 1,194 inpatients (mean [SD] age, 63.2 [18.3] years; 45.2% women) admitted to COVID-19 cohort wards during the study period were included. The outpatient cohort of 1,814 individuals (mean [SD] age, 45.4 [17.2] years; 69.1% women) was sampled from epidemiological line lists by systematic quasi-random sampling. The sensitivity (95% CI) for laboratory confirmed cases (repeat-tested patients) was 85.7% (81.5–89.1%) inpatients; 95.5% (92.2–97.5%) outpatients, 89.9% (88.2–92.1%) all. When also patients that were graded as high

Testing disease specificity/sensitivity over large populations

Suppose that a COVID test is 99.5% specific, such that it reports negative 99.5% of the time someone is not infected and 95.5% sensitive (95.5% of the time it reports positive for someone who is infected)...

If someone tests positive, what's the probability that they have COVID?

The answer: *It depends!*

What is the **prior**? What is the incidence rate if we did not know about the result?

Scenario 1:

0.1% of the population in the area is estimated to have COVID

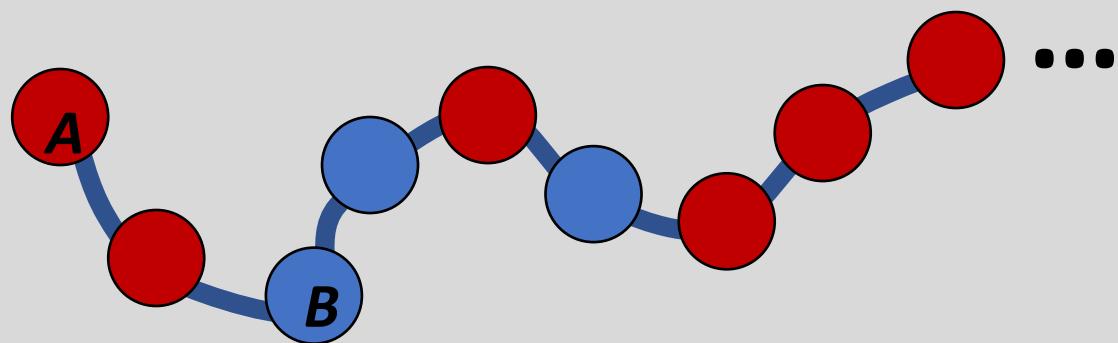
Scenario 2:

1% of the population in the area is estimated to have COVID

Revisiting Rules

Unconventional example

Suppose I have a copolymer comprised of 50% A and 50% B units, and I observe the following segment in the chain.
What's the probability that the next unit is type B?



Discuss!

Some possible hypotheses/conjectures....

- My knowledge of the current sequence doesn't matter → 1/2
- I use my observations of the chain to build $P(B|A)$ → 2/5
- We use Bayes' Rule with an estimate of $P(A|B)$ from the segment

$$\rightarrow P(B|A) = \frac{P(A|B)P(B)}{P(A)} = \frac{\frac{2}{3} \cdot \frac{1}{2}}{\frac{1}{2}} = \frac{2}{3}$$

(that we have observed more A following B in a population we know to have 50/50 A and B informs our prediction)

Summary Statistics: Definitions

A **summary statistic** of a random variable is a deterministic description of how it behaves

For a univariate random variable X , **the expected value** of a function $g(x)$ is:

$$\mathbb{E}_X [g(x)] = \int_{\mathcal{X}} g(x)p(x)dx \text{ or } \sum_{x \in \mathcal{X}} g(x)p(x)$$

(Extension to multivariate random variables requires X to be a vector of univariate random variables)

For a random variable X , with states $\mathbf{x} [m]$

Mean $\boldsymbol{\mu} := \mathbb{E}_X [\mathbf{x}] = \begin{bmatrix} \mathbb{E}_{X_1} [x_1] \\ \vdots \\ \mathbb{E}_{X_m} [x_m] \end{bmatrix} \in \mathbb{R}^m$

For a random variable X , with states $\mathbf{x} [m]$, and a variable Y , with states $\mathbf{y} [n]$

Covariance $\boldsymbol{\Sigma}_{xy} := \text{Cov} [\mathbf{x}, \mathbf{y}] = \mathbb{E} [\mathbf{x}\mathbf{y}^T] - \mathbb{E} [\mathbf{x}] \mathbb{E} [\mathbf{y}]^T \in \mathbb{R}^{m \times n}$

Summary Statistics: Definitions

A **summary statistic** of a random variable is a deterministic description of how it behaves

Variance

$$\begin{aligned}\mathbb{V}[\mathbf{x}] &= \text{Cov}_X[\mathbf{x}, \mathbf{x}] \\ &= \mathbb{E}[(\mathbf{x} - \mu)(\mathbf{x} - \mu)^T] \\ &= \mathbb{E}[\mathbf{x}\mathbf{x}^T] - \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{x}]^T\end{aligned}$$

$$= \begin{bmatrix} \text{Cov}[x_1, x_1] & \text{Cov}[x_1, x_2] & \dots & \text{Cov}[x_1, x_m] \\ \text{Cov}[x_2, x_1] & \text{Cov}[x_2, x_2] & \dots & \text{Cov}[x_2, x_m] \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}[x_m, x_1] & \dots & \dots & \text{Cov}[x_m, x_m] \end{bmatrix}$$

standard deviation

$$\sigma := \begin{bmatrix} \sqrt{\text{Cov}[x_1, x_1]} \\ \vdots \\ \sqrt{\text{Cov}[x_m, x_m]} \end{bmatrix}$$

- diagonals → variances of marginal distributions
- off-diagonals are the cross-covariances
- matrix is symmetric, positive semidefinite

Correlation

$$\text{corr}[x, y] = \frac{\text{Cov}[x, y]}{\sqrt{\mathbb{V}[x]\mathbb{V}[y]}} \in [-1, 1]$$

Correlation Matrix

$$\begin{bmatrix} \text{Cov}[\tilde{x}_1, \tilde{x}_1] & \text{Cov}[\tilde{x}_1, \tilde{x}_2] & \dots & \text{Cov}[\tilde{x}_1, \tilde{x}_m] \\ \text{Cov}[\tilde{x}_2, \tilde{x}_1] & \text{Cov}[\tilde{x}_2, \tilde{x}_2] & \dots & \text{Cov}[\tilde{x}_2, \tilde{x}_m] \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}[\tilde{x}_m, \tilde{x}_1] & \dots & \dots & \text{Cov}[\tilde{x}_m, \tilde{x}_m] \end{bmatrix} \tilde{x}_i := \frac{x_i}{\sigma_i}$$

Summary Statistics: Definitions

A **summary statistic** of a random variable is a deterministic description of how it behaves

$$\mathbb{E}[x + y] = \mathbb{E}[x] + \mathbb{E}[y]$$

$$\mathbb{E}[x - y] = \mathbb{E}[x] - \mathbb{E}[y]$$

$$\mathbb{V}[x + y] = \mathbb{V}[x] + \mathbb{V}[y] + \text{Cov}[x, y] + \text{Cov}[y, x]$$

$$\mathbb{V}[x - y] = \mathbb{V}[x] + \mathbb{V}[y] - \text{Cov}[x, y] - \text{Cov}[y, x]$$

useful relationships

Often summary statistics are estimated empirically from a set of N observations

$$\bar{\mathbf{x}} := \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \quad \Sigma := \frac{1}{N} \sum_{n=1}^N (\mathbf{x} - \bar{\mathbf{x}})(\mathbf{x} - \bar{\mathbf{x}})^T$$

Gaussian Distribution

Gaussian distribution is the most well-studied distribution in science, engineering, and ML

- *ubiquity (often stemming from **central limit theorem**)*
- *computationally convenient*

univariate

$$p(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

multivariate

$$p(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{m/2}} |\boldsymbol{\Sigma}|^{-1/2} \exp\left(-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}-\boldsymbol{\mu})\right)$$

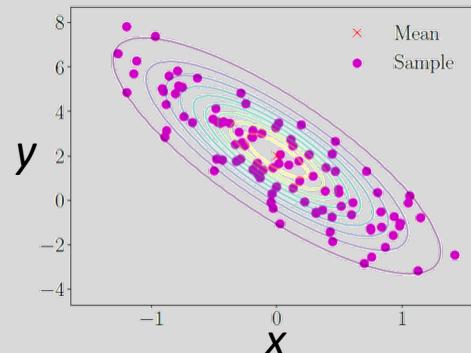
$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad \text{for short}$$

standard normal distribution $\rightarrow \mathbf{0}$ mean and I_n as the covariance matrix

Basic tenant of Gaussian distributions: *doing things with Gaussians results in Gaussians*

Gaussian Distribution

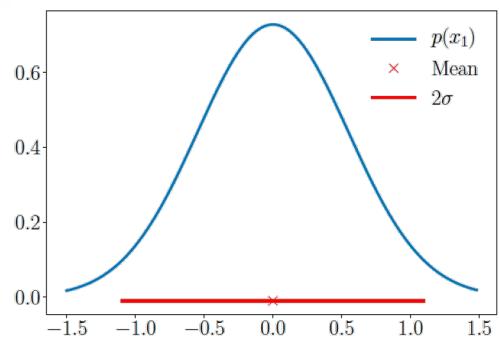
Consider the joint distribution over X and Y :



$$p(\mathbf{x}, \mathbf{y}) = \mathcal{N} \left(\begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix} \right)$$

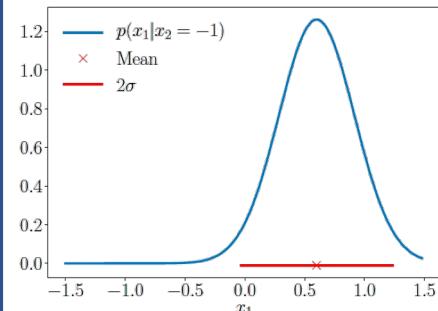
- Marginals over Gaussians are Gaussians

$$p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{y}) d\mathbf{y} = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_x, \boldsymbol{\Sigma}_{xx})$$



- Any linear transformation of Gaussians is Gaussian $p(a\mathbf{x} + b\mathbf{y}) = \mathcal{N}(a\boldsymbol{\mu}_x + b\boldsymbol{\mu}_y, a^2\boldsymbol{\Sigma}_x + b^2\boldsymbol{\Sigma}_y)$

- Conditionals of Gaussians are Gaussians



$$p(\mathbf{x} | \mathbf{y}) = \mathcal{N}(\boldsymbol{\mu}_{x|y}, \boldsymbol{\Sigma}_{x|y})$$

$$\boldsymbol{\mu}_{x|y} = \boldsymbol{\mu}_x + \boldsymbol{\Sigma}_{xy}\boldsymbol{\Sigma}_{yy}^{-1}(\mathbf{y} - \boldsymbol{\mu}_y)$$

$$\boldsymbol{\Sigma}_{x|y} = \boldsymbol{\Sigma}_{xx} - \boldsymbol{\Sigma}_{xy}\boldsymbol{\Sigma}_{yy}^{-1}\boldsymbol{\Sigma}_{yx}$$

- The product of Gaussians is Gaussian

$$\mathcal{N}(\mathbf{x} | \mathbf{a}, \mathbf{A}) \mathcal{N}(\mathbf{x} | \mathbf{b}, \mathbf{B}) \rightarrow c \mathcal{N}(\mathbf{x} | \mathbf{c}, \mathbf{C})$$

$$\mathbf{C} = (\mathbf{A}^{-1} + \mathbf{B}^{-1})^{-1}$$

$$\mathbf{c} = \mathbf{C}(\mathbf{A}^{-1}\mathbf{a} + \mathbf{B}^{-1}\mathbf{b}) \quad c = \mathcal{N}(\mathbf{a} | \mathbf{b}, \mathbf{A} + \mathbf{B})$$

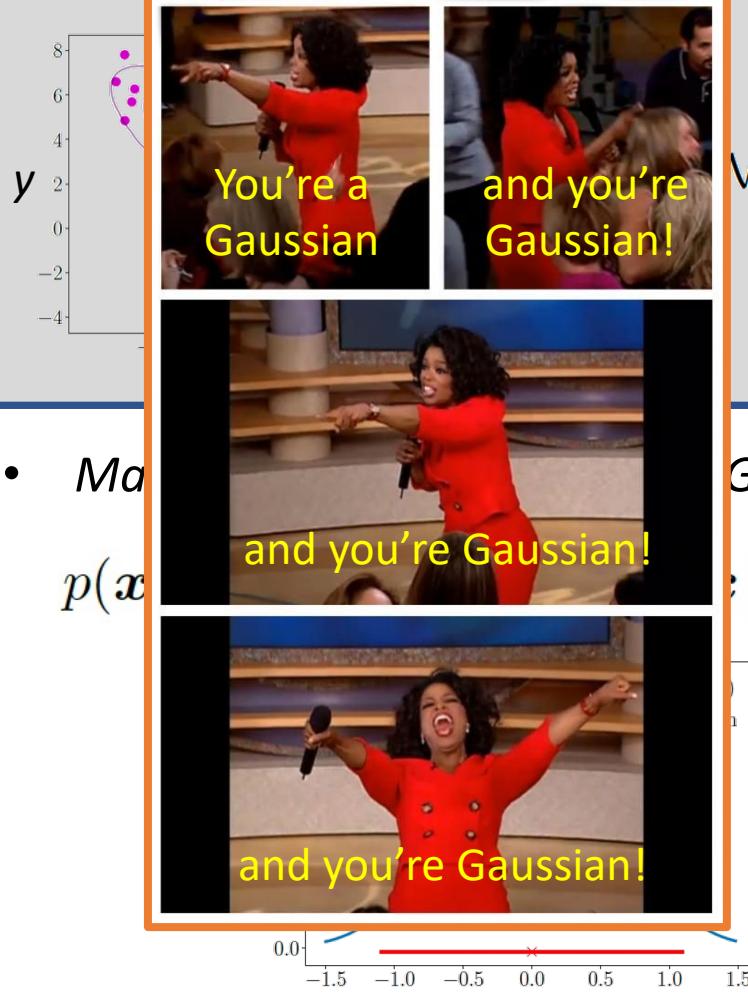
- The sum over independent Gaussian variables is... Gaussian

$$p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x})p(\mathbf{y})$$

$$p(\mathbf{x} + \mathbf{y}) = \mathcal{N}(\boldsymbol{\mu}_x + \boldsymbol{\mu}_y, \boldsymbol{\Sigma}_x + \boldsymbol{\Sigma}_y)$$

Gaussian Distribution

Consider the joint distribution over X and Y :



- Marginal distributions are Gaussians

$$p(x) = \mathcal{N}(x | \mu_x, \Sigma_{xx})$$

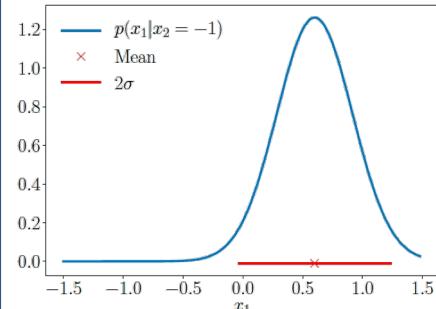
$$p(y) = \mathcal{N}(y | \mu_y, \Sigma_{yy})$$

$$\text{Conditional distribution } p(x|y) = \mathcal{N}(x | \mu_{x|y}, \Sigma_{x|y})$$

$$\text{where } \mu_{x|y} = \mu_x + \Sigma_{xy} \Sigma_{yy}^{-1} (y - \mu_y)$$

$$\Sigma_{x|y} = \Sigma_{xx} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{yx}$$

- Conditionals of Gaussians are Gaussians



$$p(x | y) = \mathcal{N}(\mu_{x|y}, \Sigma_{x|y})$$

$$\mu_{x|y} = \mu_x + \Sigma_{xy} \Sigma_{yy}^{-1} (y - \mu_y)$$

$$\Sigma_{x|y} = \Sigma_{xx} - \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{yx}$$

- The product of Gaussians is Gaussian

$$\mathcal{N}(x|\mathbf{a}, \mathbf{A}) \mathcal{N}(x|\mathbf{b}, \mathbf{B}) \rightarrow c \mathcal{N}(x|\mathbf{c}, \mathbf{C})$$

$$\mathbf{C} = (\mathbf{A}^{-1} + \mathbf{B}^{-1})^{-1}$$

$$\mathbf{c} = \mathbf{C}(\mathbf{A}^{-1}\mathbf{a} + \mathbf{B}^{-1}\mathbf{b}) \quad c = \mathcal{N}(\mathbf{a}|\mathbf{b}, \mathbf{A} + \mathbf{B})$$

- The sum over independent Gaussian variables is... Gaussian

$$p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x})p(\mathbf{y})$$

$$p(\mathbf{x} + \mathbf{y}) = \mathcal{N}(\mu_x + \mu_y, \Sigma_x + \Sigma_y)$$

- Any linear transformation of Gaussians is Gaussian $p(a\mathbf{x} + b\mathbf{y}) = \mathcal{N}(a\mu_x + b\mu_y, a^2\Sigma_x + b^2\Sigma_y)$

Examples of Other “Named” Distributions

Most widely used model distributions are members of a so-called *exponential family*

$$p(\mathbf{x} | \boldsymbol{\theta}) = h(\mathbf{x}) \exp (\langle \boldsymbol{\theta}, \phi(\mathbf{x}) \rangle - A(\boldsymbol{\theta})) \propto \exp(\langle \boldsymbol{\theta}, \phi(\mathbf{x}) \rangle)$$

vector of summary statistics *an inner product (e.g., dot)*

The importance of these (and other) distributions is largely in the realm of probabilistic modeling

Gaussian distribution:

$$\begin{aligned}\mathcal{N}(x|\mu, \sigma^2) &\propto \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right) \\ &\propto \exp\left(-\frac{1}{2\sigma^2}(x^2 - 2x\mu + \mu^2)\right) \\ &\propto \exp\left(\frac{x\mu}{\sigma^2} - \frac{x^2}{2\sigma^2}\right)\end{aligned}$$

$$\boldsymbol{\theta} = \left[\frac{\mu}{\sigma^2}, -\frac{1}{2\sigma^2} \right]^T$$

$$\phi(\mathbf{x}) = [x, x^2]^T$$

Bernoulli distribution: *distribution for binary outcomes (0,1)*

$$\begin{aligned}\text{Ber}(\mu) &= p(x|\mu) = \mu^x(1-\mu)^{1-x} \\ \mathbb{E}[x] &= \mu; \quad \mathbb{V}[x] = \mu(1-\mu)\end{aligned}$$

Binomial distribution: *distribution arising from N Bernoulli events*

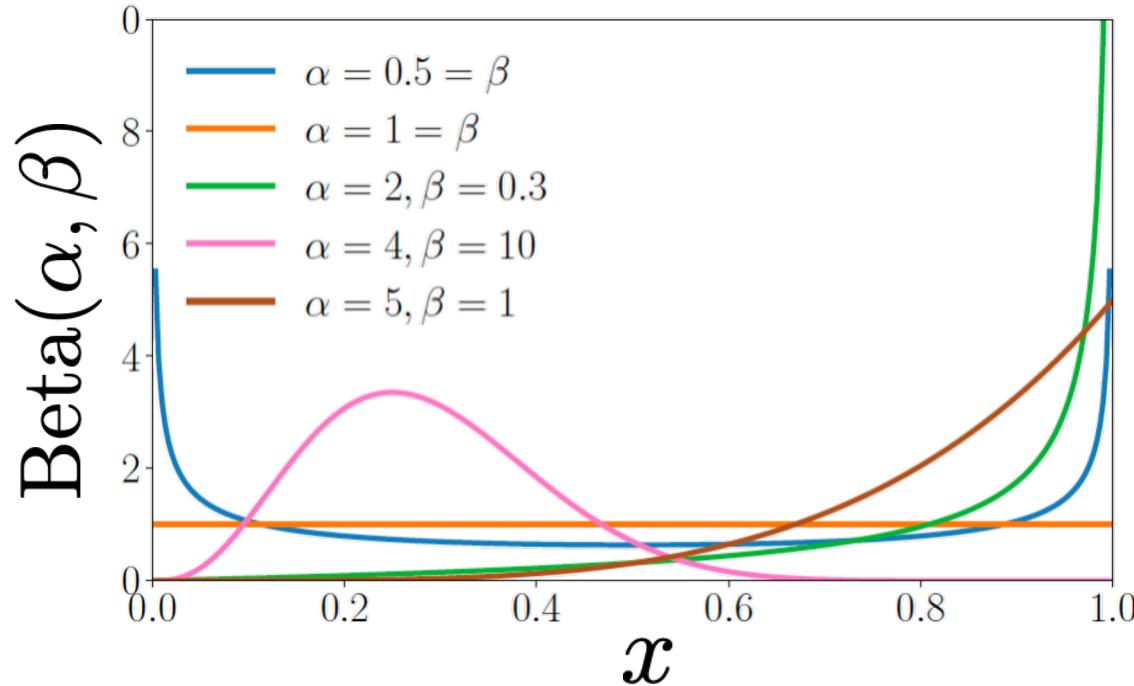
$$\begin{aligned}\text{Bin}(N, \mu) &= p(m|N, \mu) = \binom{N}{m} \mu^m (1-\mu)^{N-m} \\ \mathbb{E}[m] &= N\mu; \quad \mathbb{V}[m] = N\mu(1-\mu)\end{aligned}$$

Examples of Other “Named” Distributions

Beta distribution: *distribution of a continuous variable over a finite interval [0,1]*

$$\text{Beta}(\alpha, \beta) = p(x|\alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1}$$

$$\begin{aligned}\Gamma(t) &:= \int_0^\infty x^{t-1} \exp(-x) dx \\ \Gamma(t+1) &= t\Gamma(t)\end{aligned}$$



- $\alpha = \beta = 1$, uniform distribution
- $\alpha, \beta < 1$, bimodal with spikes at 0 and 1
- $\alpha, \beta > 1$, unimodal
- $\alpha, \beta > 1$ and $\alpha = \beta$, unimodal, centered, & symmetric