

Dist & Similarity

2 types:

Numerical

Categorical

Numeric

- integers, reals, complex

E.g.

Age: domain(Age) = $\mathbb{N} \subseteq$ natural #'s
 Height: domain(Height) = $\mathbb{R}^+ \subseteq$ (non-neg ints)
 positive real #'s

Subtypes:

discrete: finite (or countably inf)

continuous: they take a real value

binary: take value {0, 1}

Categorical

- values come from a set of symbols

Domain(Education) = {Highschool, BS, MS, PhD}

Domain(Vehicle type) = {car, truck, SUV}

Subtypes

- nominal: unordered (only equality tests make sense)

- ordinal: values are ordered (equality and inequality)

Data matrix P (assume all values are numeric)
⇒ each row can be thought of as d -dim point

$$\text{e.g. } x_i = (x_{i1}, x_{i2}, \dots, x_{id}) \in \mathbb{R}^d$$

often as col vector $x_i = \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{id} \end{pmatrix}$

Let e_j be the j^{th} std basis vector

$$\rightarrow (000 \dots 0 \overset{j}{1} 0 \dots 0)$$

all 0's except j^{th} entry \downarrow is 1

$$x_i = x_{i1} e_1 + x_{i2} e_2 + \dots + x_{id} e_d = \sum_{j=1}^d x_{ij} e_j$$

5.9	3.0	4.2	1.5	Iris-versicolor
6.9	3.1	4.9	1.5	Iris-versicolor
6.6	2.9	4.6	1.3	Iris-versicolor
4.6	3.2	1.4	0.2	Iris-setosa
6.0	2.2	4.0	1.0	Iris-versicolor
4.7	3.2	1.3	0.2	Iris-setosa
6.5	3.0	5.8	2.2	Iris-virginica
5.8	2.7	5.1	1.9	Iris-virginica

consider
this
row

$$x_2 = (6.9 \ 3.1 \ 4.9 \ 1.5 \text{ Iris-versicolor})$$

"project" to first 3 attrib

$$\hat{x}_2 = (6.9 \ 3.1 \ 4.9) \in \mathbb{R}^3$$

$$= 6.9 \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} + 3.1 \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} + 4.9 \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} = \sum_{j=1}^3 x_{2,j} e_j$$

5. only for cols
vectors \mathbb{R}^n

$$x_i = \begin{pmatrix} x_{i,1} \\ x_{i,2} \\ \vdots \\ x_{i,n} \end{pmatrix}$$

an elt
of $\mathbb{R}^{n \times d}$

$$D = \begin{pmatrix} x_{1,1} & x_{1,2} & \cdots & x_{1,d} \\ x_{2,1} & x_{2,2} & \cdots & x_{2,d} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n,1} & x_{n,2} & \cdots & x_{n,d} \end{pmatrix} =$$

$$\text{as } n \text{ row vectors } x_i^T \in \mathbb{R}^d$$

$$\begin{pmatrix} -x_1^T- \\ -x_2^T- \\ \vdots \\ -x_n^T- \end{pmatrix}$$

$$= \begin{pmatrix} | & | \\ x_1 & x_2 & \cdots & x_n \\ | & | \end{pmatrix}$$

Dist & Angles

Let $a, b \in \mathbb{R}^m$

$$a = \begin{pmatrix} a_1 & a_2 & \dots & a_m \end{pmatrix}^T$$

$$b = \begin{pmatrix} b_1 & b_2 & \dots & b_m \end{pmatrix}$$

Subtract points

$$a - b = \begin{pmatrix} a_1 - b_1 \\ a_2 - b_2 \\ \vdots \\ a_m - b_m \end{pmatrix} = v$$

dot product

$$a \cdot b = \sum_{i=1}^m a_i b_i$$

$p \in \mathbb{R}$ w/ $p \geq 1$

L_p-norm $\|a\|_p = \left(\sum_{i=1}^m |a_i|^p \right)^{\frac{1}{p}}$

L_∞-norm

$$\|a\|_\infty = \max_i a_i$$

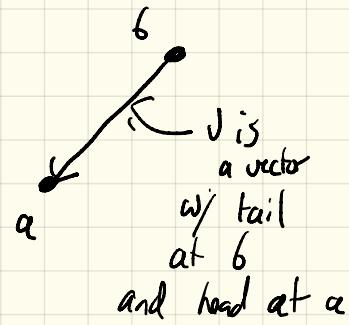
L_p-dist

$$\|a - b\|_p = \left(\sum_{i=1}^m |a_i - b_i|^p \right)^{\frac{1}{p}}$$

Note

$$\|a\| = \|a\|_2$$

$p=2$ is Euclidean dist



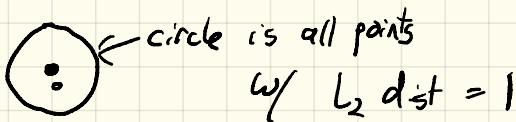
Interpret geometrically

consider points at dist 1 from origin

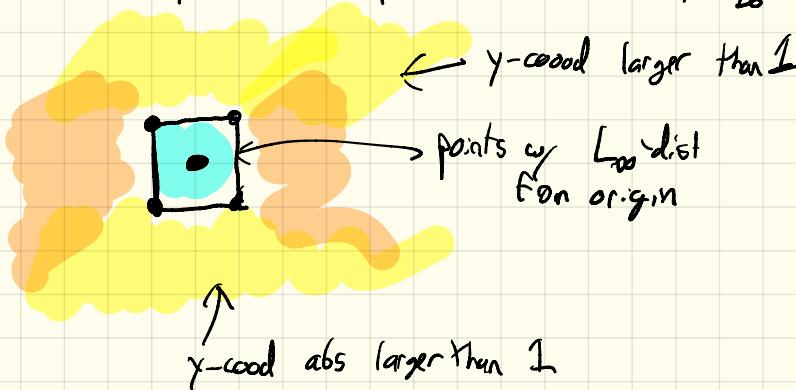
case p=2 find all points $(x, y) \in \mathbb{R}^2$
dist 1 from origin

$$\sqrt{(x-0)^2 + (y-0)^2} = 1$$

$$\Rightarrow (x-0)^2 + (y-0)^2 = 1^2$$



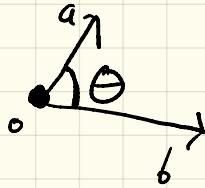
p=∞ all points $(x, y) \in \mathbb{R}^2$ $\| (x, y) \|_\infty = 1$



$\cos \theta =$

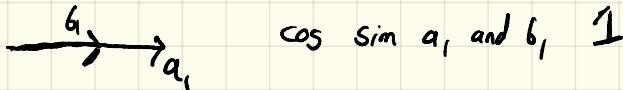
$$\frac{\mathbf{a} \cdot \mathbf{b}}{\|\mathbf{a}\| \|\mathbf{b}\|}$$

geometrically

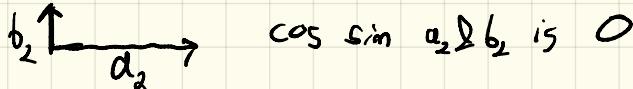


$$\mathbf{a} \cdot \mathbf{b} = \|\mathbf{a}\| \|\mathbf{b}\| \cos \theta$$

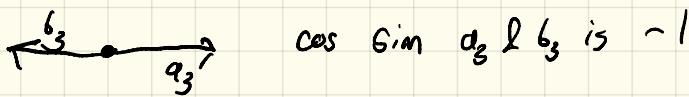
$\cos \theta$ is
cos of angle between 2 vrcs



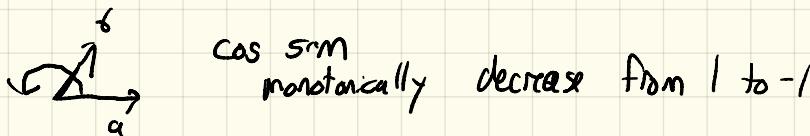
$\cos \theta$ is 1



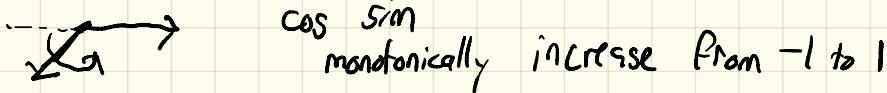
$\cos \theta$ is 0



$\cos \theta$ is -1



$\cos \theta$ monotonically decrease from 1 to -1



$\cos \theta$ monotonically increase from -1 to 1

orthogonal

iff $\mathbf{a} \cdot \mathbf{b} = 0 \ (\Rightarrow \cos \theta = 0)$