

Alon Bar 201639663
Daniel Meller 038198149

Homework 1

Convexity

1.1.1

Lets prove that

$$R_{\geq 0}^d$$

is convex

let

$$0 \leq t \leq 1$$

and

$$x, y \in R_{\geq 0}^d$$

then for any

$$0 \leq i \leq d$$

if

$$z = x + y$$

$$z_i = x_i + y_i \geq x_i \geq 0$$

so

$$z \in R_{\geq 0}^d$$

so

$$R_{\geq 0}^d$$

is convex, as required

1.1.2

Lets prove by induction for any $N \geq 2$ that for $\{\alpha\}_{i=1}^N$ with $\sum_{i=1}^N \alpha_i = 1$ and $\{x_i \in C\}_{i=1}^N$ the

following occurs: $\sum_{i=1}^N \alpha_i x_i \in C$

For $N = 2$, it occurs due to the definition of a convex set.

Lets assume correctness for any $N - 1$ and prove for N

$$\text{So } \sum_{i=1}^N \alpha_i x_i = (\alpha_1 + \dots + \alpha_{n-1}) \left(\frac{\alpha_1}{\sum_{i=1}^{n-1} \alpha_i} x_1 + \dots + \frac{\alpha_{n-1}}{\sum_{i=1}^{n-1} \alpha_i} x_{n-1} \right) + \alpha_n x_n$$

From the induction base we have that $\left(\frac{\alpha_1}{\sum_{i=1}^{n-1} \alpha_i} x_1 + \dots + \frac{\alpha_{n-1}}{\sum_{i=1}^{n-1} \alpha_i} x_{n-1} \right) \in C$ since

$$\frac{\alpha_1}{\sum_{i=1}^{n-1} \alpha_i} + \dots + \frac{\alpha_{n-1}}{\sum_{i=1}^{n-1} \alpha_i} = 1$$

Also $(\alpha_1 + \dots + \alpha_{n-1}) = 1 - \alpha_n$ thus we have 2 terms in C with coefficients that sum up to 1.

There for $\sum_{i=1}^N \alpha_i x_i \in C$ as required.

1.1.3

This is not true

Consider C as the sphere $x_1^2 + x_2^2 - 2 = 0$

Then C is a convex set. Now maybe $\{x_i\}$ are all points on the line $x_1 + x_2 = 0$ which intersects with the sphere for instance $(0, 0)$, $(1, -1)$, $(-0.5, 0.5)$

The point $(1, 1)$ for instance is in the sphere but may not be represented by a convex combination of the points

1.1.4

Lets define $z = \alpha x + (1 - \alpha)y$

By definition

$$\begin{aligned} f(y) &\geq f(z) + f'(z)(y - z) \Rightarrow (1 - \alpha)f(y) \geq (1 - \alpha)f(z) + (1 - \alpha)f'(z)(y - z), \\ f(x) &\geq f(z) + f'(z)(x - z) \Rightarrow \alpha f(x) \geq \alpha f(z) + \alpha f'(z)(x - z) \end{aligned}$$

Now lets add the two terms

$$\begin{aligned} (1 - \alpha)f(y) + \alpha f(x) &\geq \alpha f(z) + (1 - \alpha)f(z) + f'(z)((1 - \alpha)(y - z) + \alpha(x - z)) = \\ f(z) + f'(z)((1 - \alpha)y + \alpha x - z) &= f(z) + f'(z) \cdot 0 = f(\alpha x + (1 - \alpha)y) \end{aligned}$$

So we got :

$$(1 - \alpha)f(y) + \alpha f(x) \geq f((1 - \alpha)y + \alpha x)$$

Thus f is convex

2.3 Hierarchical Clustering

2.3.1

Lets prove that $d_{\text{complete-link}}^2(C_1, C_2)$ is a metric

We will show that the 3 properties of a metric hold, and deduce that it is indeed an metric.

1. Identity:

a. By definition, if $C_1 = C_2$ then the function yields 0 so $x = y \Rightarrow d(x, y) = 0$

b.

2. Symmetry:

a. Since $\forall x_i, x_j ||x_i - x_j||_2^2 = ||x_j - x_i||_2^2$ then

$$d_{\text{complete-link}}^2(C_1, C_2) = d_{\text{complete-link}}^2(C_2, C_1)$$

3. Triangle inequality:

a. Need to show that

$$d_{\text{complete-link}}^2(C_1, C_2) + d_{\text{complete-link}}^2(C_2, C_3) \geq d_{\text{complete-link}}^2(C_1, C_3)$$

For some $x_1 \in C_1, x_3 \in C_3$

$$d_{\text{complete-link}}^2(C_1, C_3) = ||x_1 - x_3||_2^2$$

Let's choose some random $x_2 \in C_2$ then from triangle inequality

$$||x_1 - x_3||_2^2 \leq ||x_1 - x_2||_2^2 + ||x_2 - x_3||_2^2$$

Also lets mark with $x^*, x^{**} \in C_2$ the appropiates members in C_2 that produce the

values of the metrics $d_{\text{complete-link}}^2(C_1, C_2), d_{\text{complete-link}}^2(C_2, C_3)$

So we get that

$$d_{\text{complete-link}}^2(C_1, C_3) =$$

$$||x_1 - x_3||_2^2 \leq ||x_1 - x_2||_2^2 + ||x_2 - x_3||_2^2 \leq ||x_1 - x^*||_2^2 + ||x^{**} - x_3||_2^2 =$$

$$d_{\text{complete-link}}^2(C_1, C_2), d_{\text{complete-link}}^2(C_2, C_3)$$

As needed

Thus this is a metric

May 9, 2021

1 1.2.1

we need to prove that $\forall h \in \mathbb{R}^d \nabla f(x_0)[h] = \langle g_0, h \rangle \implies g_0 = \nabla f(x_0)$
by definition, if $f(x_0)[h] = \langle g_0, h \rangle$, it is also true for h of the standard base form such as:

$$\begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}$$

$\nabla f(x_0)[e_i] = \frac{\partial f}{\partial x_i}(x_0)$, we know that $g = (g_1, g_2, \dots, g_d)$, so the inner product $\langle g, e_i \rangle = g_i$ and we can conclude that $g_i = \frac{\partial f}{\partial x_i}$ which derives that $g(x_0) = \nabla f(x_0)$

2 1.2.2

$$\nabla f(x)[h] := \lim_{x \rightarrow 0} \frac{f(x+ht) - f(x)}{t} = \lim_{x \rightarrow 0} \frac{f(x) + tf(x) - f(x)}{t} = \frac{tf(h)}{t} = f(h)$$

3 1.2.3.1

$$\begin{aligned} f(x) &= x^T A x \\ \nabla f(x)[h] &= \lim_{t \rightarrow 0} \frac{f(x+th) - f(x)}{t} = \lim_{t \rightarrow 0} \frac{(x+th)^T A (x+th) - x^T A x}{t} = \\ &= \lim_{t \rightarrow 0} \frac{x^T A x + tx^T A h + th^T A x + t^2 x^T A x - x^T A x}{t} = \\ &= \lim_{t \rightarrow 0} \frac{t(x^T A h + h^T A x + tx^T A x)}{t} = \\ &= x^T A h + h^T A x = x^T A h + x^T A^T h = x^T (A + A^T) h = \langle (A + A^T)x, h \rangle \\ \nabla f(x)[h] &= (A + A^T)x \end{aligned}$$

4 1.2.3.2

$$\nabla f(x)[h] = \nabla \langle g, h \rangle (X)[H] = ** \langle \nabla g(X)[H], l(X) \rangle + \langle g(X), \nabla h(X)[H] \rangle = \langle H, AX \rangle + \langle X, AH \rangle = \langle (A + A^T)X, H \rangle$$

**using the product rule to derive this formula

5 1.2.3.3

$$f(x) = \|y - Ax\| = \langle y - Ax, y - Ax \rangle \nabla f(x)[h] = \langle -Ah, y - Ax \rangle + \langle y - Ax, -Ah \rangle = 2 \langle y - Ax, -Ah \rangle = -2 \langle A^T(y - Ax), h \rangle \nabla f(x) = -2A^T(y - Ax)$$

using the product rule to derive this formula

6 1.2.3.4

same like previous section

$$f(X) = \|Y - AX\|_F^2 = \langle Y - AX, Y - AX \rangle \nabla f(X)[H] = \langle -AH, Y - AX \rangle + \langle Y - AX, -AH \rangle = 2 \langle Y - AX, -AH \rangle = -2 \langle A^T(Y - AX), H \rangle \nabla f(X) = -2A^T(Y - AX)$$

using the product rule to derive this formula

7 1.2.3.5

$$f(X) = \langle X^T A, Y^T \rangle \nabla f(X)[H] = \nabla \langle X^T A, Y^T \rangle [H] = \langle^T A[H], Y^T \rangle + \langle X^T A, ^T [H] \rangle = \langle AY, H \rangle \rightarrow \nabla f(X) = AY$$

8 1.2.3.7

$$f(X) = \langle A, \log[X] \rangle, g(X) = \log[X]$$

$$\nabla g(X)[H] = \nabla \log(X)[H] = \lim_{t \rightarrow 0} \frac{\log[X+tH] - \log[X]}{t} = \lim_{t \rightarrow 0} \frac{\log[\frac{X+tH}{X}]}{t} = H \circ X^{\circ-1} = ** (X)[H] = \langle 0, \log[X] \rangle + \langle A, H \circ X^{\circ-1} \rangle = \langle A \circ X^{\circ-1}, H \rangle \rightarrow (X) = A \circ X^{\circ-1}$$

**product rule

9 1.2.3.8

$$f(X) = \langle a, \text{diag}(X) \rangle \nabla f(X)[H] = \lim_{t \rightarrow 0} \text{Tr}[\frac{\text{adiag}(X+tH) - \text{adiag}(X)}{t}] = \lim_{t \rightarrow 0} \text{Tr}[a \frac{\text{diag}(X) + \text{diag}(tH) - \text{diag}(X)}{t}] = \text{Tr}[\text{adiag}(H)] = \langle a, \text{diag}(H) \rangle \rightarrow \nabla f(X) = \text{diag}(a)$$

10 1.4

$$G(x, y) = \sin(x + y)$$

11 1.4.1.1

$\min_x \max_y G(x, y)$

exists x_0 such that $\max_y \sin(x_0 + y) = 1$ (sin function properties)

meaning, $\min_x 1 = 1$

12 1.4.1.2

$\max_y \min_x G(x, y)$

exists y_0 such that $\min_x \sin(x + y_0) = -1$ (sin function properties)

meaning, $\max_y -1 = -1$

13 1.4.2.1

according to Rayleigh quotient $f(x) = \frac{x^T A x}{x^T x}$, by the definition in class, we know that $f(x) = \frac{x^T A x}{x^T x} = \frac{\langle x, A x \rangle}{\langle x, x \rangle}$.

for some scalar k , we know that $f(kx) = \frac{x^T A x}{x^T x} = \frac{\langle kx, kAx \rangle}{\langle kx, kx \rangle} = \frac{k^2 \langle x, Ax \rangle}{k^2 \langle x, x \rangle} = \frac{\langle x, Ax \rangle}{\langle x, x \rangle} = f(x)$.

we can conclude that for every x that minimize $f(x)$, kx also minimize $f(x)$, so we can write $\min(x^T A x)$ such that $\|x\|^2 = 1$ and we can conclude that $\min f(x) = \min \langle x, Ax \rangle = \min(x^T A x)$, what we were asked to prove

14 1.4.2.2

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

15 1.4.2.3

To prove: $\nabla L(x, \lambda) = 0 \Leftrightarrow Ax = \lambda x$

$$\nabla L(x, \lambda) = 2Ax - \lambda 2x$$

$$2Ax - \lambda 2x = 0 \Leftrightarrow Ax = \lambda x$$

16 2.1.1.1

we will show equivalence for the minimization target.

$$\begin{aligned} \sum_{x_i \in D_k} \|x_i - \mu_k\|_2^2 &= \sum_{x_i \in D_k} \|x_i\|_2^2 + \|\mu_k\|_2^2 - 2 \langle x_i, \mu_k \rangle \\ &= \sum_{x_i \in D_k} \|x_i\|_2^2 + \sum_{x_i \in D_k} \|\mu_k\|_2^2 - 2 \sum_{x_i \in D_k} \langle x_i, \mu_k \rangle \\ &= \sum_{x_i \in D_k} \|x_i\|_2^2 + N \|\mu_k\|_2^2 - 2N \langle \frac{1}{N} \sum_{x_i \in D_k} x_i, \mu_k \rangle = ** \\ &= \sum_{x_i \in D_k} \|x_i\|_2^2 + N \|\mu_k\|_2^2 - 2N \|\mu_k\|_2^2 \\ &= \sum_{x_i \in D_k} \|x_i\|_2^2 - N \|\mu_k\|_2^2 \end{aligned}$$

$$\begin{aligned}
& \Sigma_{x_i, x_j \in D_k} \|x_i - x_j\|_2^2 = \Sigma_{x_i \in D_k} \|x_i\|_2^2 + \Sigma_{x_j \in D_k} \|x_j\|_2^2 - \Sigma_{x_i, x_j \in D_k} 2 \langle x_i, x_j \rangle \\
& = N \Sigma_{x_i \in D_k} \|x_i\|_2^2 + N \Sigma_{x_j \in D_k} \|x_j\|_2^2 - \Sigma_{x_i, x_j \in D_k} 2 \langle x_i, x_j \rangle \\
& = 2N \Sigma_{x_i \in D_k} \|x_i\|_2^2 - 2 \Sigma_{x_i \in D_k} \Sigma_{x_i \in D_k} \langle x_i, x_j \rangle \\
& = 2N \Sigma_{x_i \in D_k} \|x_i\|_2^2 - 2 \Sigma_{x_i \in D_k} \langle x_i, \Sigma_{x_i \in D_k} x_j \rangle \\
& = 2N (\Sigma_{x_i \in D_k} \|x_i\|_2^2 - \Sigma_{x_i \in D_k} \langle x_i, \mu_k \rangle) \\
& = 2N (\Sigma_{x_i \in D_k} \|x_i\|_2^2 - N \|\mu_k\|_2^2)
\end{aligned}$$

this equivalence is correct since we have only constant difference between them.

17 2.1.1.2

$\arg \min_{\mu_k} \sum_{i=1}^N \min_k \|x_i - \mu_k\|_2^2$
by definition of the K-Means algorithm, every $x_i \in D_k$ it holds that $k = \arg \min d(x_i, \mu_i)$
we will use that in the objective of kmeans $\arg \min \sum_{k=1}^K \sum_{x_i \in D_k} \|x_i - \mu_k\|_2^2 = \arg \min \sum_{k=1}^K \min \|x_i - \mu_k\|_2^2$

18 2.1.2

the statement is False

a counter example will be: 4 points that create a rectangle, where the centroids are in the longer edge of the rectangle. lets use the four points (1,2),(1,4),(7,2),(7,4). if the centroids starts at (4,2),(4,4) - it will set two clusters and the centroids will be the center of the clusters, and it will not change.
. while the global solution will be (1,3), (7,3)

19 2.2.1

$$\begin{aligned}
\mu_y &= E[Y] = E[a^T X + b] = E[a^T X] + b = a^T \mu_x + b \\
\Sigma_y &= \text{Var}(a^T X + b) = \text{Var}(Ax) + \text{Var}(b) = a^T \Sigma_x a \\
Y &\sim N(a^T \mu_x + b, a^T \Sigma_x a) \\
f_Y(y) &= \frac{1}{(2\pi)^d \det \Sigma_y} \exp(-\frac{1}{2}(y - \mu_y)^T \Sigma_y^{-1} (y - \mu_y)) = \\
&= \frac{1}{(2\pi)^d \det(a^T \Sigma_x a)} \exp(-\frac{1}{2}(y - a^T \mu_x + b)^T (a^T \Sigma_x a)^{-1} (y - a^T \mu_x + b))
\end{aligned}$$

20 2.2.2

Symmetry:

$$\Sigma_x^T = E[(X - \mu_x)(X - \mu_x)^T]^T = E[((X - \mu_x)^T)^T (X - \mu_x)^T] = E[(X - \mu_x)(X - \mu_x)^T] = \Sigma_x$$

Positive:

$$\begin{aligned} &\text{let } Y = X - \mu_x \text{ and some vector } v. \\ &v^T \Sigma_x v = v^T E[YY^T]v = E[v^T YY^T v] = E[(v^T Y)(v^T Y)^T] = E[\|vY\|^2] \geq 0 \end{aligned}$$