

# Bregman Divergence = 1st-order Taylor error

---

**Andersen Ang**

U.Southampton UK

[angms.science](http://angms.science)

January 23, 2025

1st draft

June 6, 2017

Content

Jensen inequality & 1st-order Taylor series

Bregman divergence  $B(\mathbf{p}, \mathbf{q}) := f(\mathbf{p}) - f(\mathbf{q}) - \langle \nabla f(\mathbf{q}), \mathbf{p} - \mathbf{q} \rangle$

Example of Bregman divergence

Some basic properties of Bregman divergence

Summary

# Contents

Jansen inequality & 1st-order Taylor series

Bregman divergence  $B(\mathbf{p}, \mathbf{q}) := f(\mathbf{p}) - f(\mathbf{q}) - \langle \nabla f(\mathbf{q}), \mathbf{p} - \mathbf{q} \rangle$

Example of Bregman divergence

Some basic properties of Bregman divergence

Summary

## Two equivalent definitions of convexity

- Setup: given a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  that is convex and differentiable
- If  $f$  is convex, then
  - $\text{dom} f (= \mathbb{R}^n)$  is a convex set, and  
**Jensen's inequality:**  $\forall \mathbf{x}, \mathbf{y} \in \text{dom} f, \lambda \in [0, 1]$

$$f(\lambda \mathbf{y} + (1 - \lambda) \mathbf{x}) \leq \lambda f(\mathbf{y}) + (1 - \lambda) f(\mathbf{x}). \quad (\#)$$

**1st-order Taylor series of  $f$  at  $\mathbf{x}$  is a global under-estimator:**  $\forall \mathbf{x}, \mathbf{y} \in \text{dom} f$

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle. \quad (\dagger)$$

- **Theorem.** The  $(\#)$  and  $(\dagger)$  are equivalent. I.e.,  $(\#) \iff (\dagger)$ .
- This PDF: prove this theorem and show  $(\dagger)$  is the Bregman divergence of  $f$ .
- You want more definitions then [click here](#)

$$f(\lambda \mathbf{y} + (1 - \lambda)\mathbf{x}) \leq \lambda f(\mathbf{y}) + (1 - \lambda)f(\mathbf{x}) \implies f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle$$

- **Proof.** Starting from Jensen inequality  $f(\lambda \mathbf{y} + (1 - \lambda)\mathbf{x}) \leq \lambda f(\mathbf{y}) + (1 - \lambda)f(\mathbf{x})$

$$f(\lambda \mathbf{y} + (1 - \lambda)\mathbf{x}) \leq \lambda f(\mathbf{y}) + (1 - \lambda)f(\mathbf{x})$$

$$\iff f(\mathbf{x} + \lambda(\mathbf{y} - \mathbf{x})) \leq f(\mathbf{x}) + \lambda(f(\mathbf{y}) - f(\mathbf{x}))$$

$$\iff f(\mathbf{x} + \lambda(\mathbf{y} - \mathbf{x})) - f(\mathbf{x}) \leq \lambda(f(\mathbf{y}) - f(\mathbf{x}))$$

$$\iff \frac{f(\mathbf{x} + \lambda(\mathbf{y} - \mathbf{x})) - f(\mathbf{x})}{\lambda} \leq f(\mathbf{y}) - f(\mathbf{x})$$

$$\iff f(\mathbf{y}) \geq f(\mathbf{x}) + \frac{f(\mathbf{x} + \lambda(\mathbf{y} - \mathbf{x})) - f(\mathbf{x})}{\lambda}.$$

- For convex function  $\lambda \in [0, 1]$  includes 0 so we are allowed to take  $\lambda \rightarrow 0$   
Take limit  $\lambda \rightarrow 0$  on both side and by definition of directional derivative

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \lim_{\lambda \rightarrow 0} \frac{f(\mathbf{x} + \lambda(\mathbf{y} - \mathbf{x})) - f(\mathbf{x})}{\lambda} = f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle.$$

## The big picture of the proof of $\implies$

$f$  is convex  $\implies$  Jensen inequality

$$f(\lambda \mathbf{y} + (1 - \lambda)\mathbf{x}) \leq \lambda f(\mathbf{y}) + (1 - \lambda)f(\mathbf{x})$$

do algebra  $\Updownarrow$

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \frac{f(\mathbf{x} + \lambda(\mathbf{y} - \mathbf{x})) - f(\mathbf{x})}{\lambda}$$

Jensen for  $\lambda \in [0, 1]$  includes the case  $\lim_{\lambda \rightarrow 0}$

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \lim_{\lambda \rightarrow 0} \frac{f(\mathbf{x} + \lambda(\mathbf{y} - \mathbf{x})) - f(\mathbf{x})}{\lambda}$$

directional derivative  $\Updownarrow$

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle.$$

$$f(\lambda \mathbf{y} + (1 - \lambda)\mathbf{x}) \leq \lambda f(\mathbf{y}) + (1 - \lambda)f(\mathbf{x}) \iff f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle$$

**Proof** Let  $\mathbf{z} = \lambda \mathbf{x} + (1 - \lambda)\mathbf{y}$  with  $0 \leq \lambda \leq 1$ .

so  $\forall \mathbf{x}, \mathbf{y} \in \text{dom}f$  we have  $\mathbf{z} \in \text{dom}f$  because  $\text{dom}f$  is a convex set.

By convexity of  $f$ :

$$f(\mathbf{x}) \geq f(\mathbf{z}) + \langle \nabla f(\mathbf{z}), \mathbf{x} - \mathbf{z} \rangle \quad (1)$$

$$f(\mathbf{y}) \geq f(\mathbf{z}) + \langle \nabla f(\mathbf{z}), \mathbf{y} - \mathbf{z} \rangle \quad (2)$$

Consider  $\lambda(1) + (1 - \lambda)(2)$

$$\lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y}) \geq f(\mathbf{z}) + \underbrace{\langle \nabla f(\mathbf{z}), \lambda(\mathbf{x} - \mathbf{z}) + (1 - \lambda)(\mathbf{y} - \mathbf{z}) \rangle}_{=0 \text{ for } \mathbf{z}=\lambda\mathbf{x}+(1-\lambda)\mathbf{y}}. \quad (3)$$

Put  $\mathbf{z} = \lambda \mathbf{x} + (1 - \lambda)\mathbf{y}$  into (3) finishes the proof. □

# Contents

Jansen inequality & 1st-order Taylor series

Bregman divergence  $B(\mathbf{p}, \mathbf{q}) := f(\mathbf{p}) - f(\mathbf{q}) - \langle \nabla f(\mathbf{q}), \mathbf{p} - \mathbf{q} \rangle$

Example of Bregman divergence

Some basic properties of Bregman divergence

Summary

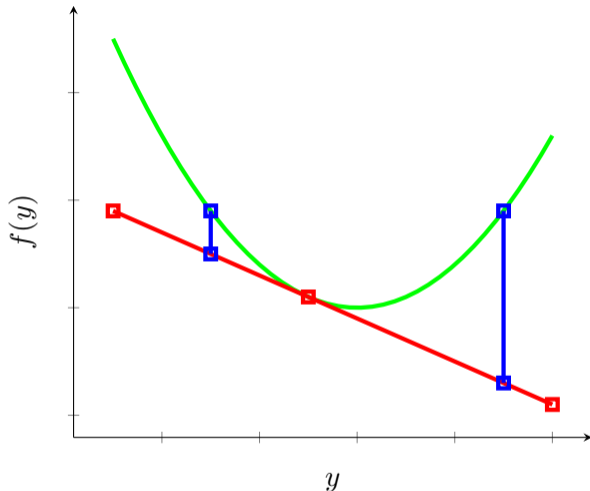
# Bregman divergence

1st-order Taylor series of a convex function is an global under-estimator:

$$f(y) \geq f(x) + \nabla f(x)^\top (y - x).$$

$$\underbrace{f(y) - f(x) - \nabla f(x)^\top (y - x)}_{=: B(y)}$$

$B(y)$  = Taylor series error



Lev .M. Bregman, The relaxation method of finding the common points of convex sets and its application to the solution of problems in convex programming, USSR Computational Mathematics and Mathematical Physics Volume 7, Issue 3, 1967, Pages 200 - 217



## Bregman divergence $B$

$f : \mathbb{R}^n \rightarrow \mathbb{R}$  is convex and differentiable

- The Bregman divergence at a point  $\mathbf{p}$  of  $f$  anchored at  $\mathbf{q}$ , denoted as  $B_f(\mathbf{p}, \mathbf{q})$ , is

$$B_f(\mathbf{p}; \mathbf{q}) := f(\mathbf{p}) - f(\mathbf{q}) - \langle \nabla f(\mathbf{q}), \mathbf{p} - \mathbf{q} \rangle.$$

- $\mathbf{p} \in \text{dom}f$  and  $\mathbf{q} \in \text{dom}f$  have different meaning here
  - $\mathbf{q}$  is the “anchor” (where  $B_f$  touch  $f$ )
  - $\mathbf{p}$  is the variable of  $B$
- 
- Why study  $B_f$ : it is closely related to convexity of  $f$
- 
- Bregman divergence is used in the design of Bregman method for convex optimization.

# Contents

Jansen inequality & 1st-order Taylor series

Bregman divergence  $B(\mathbf{p}, \mathbf{q}) := f(\mathbf{p}) - f(\mathbf{q}) - \langle \nabla f(\mathbf{q}), \mathbf{p} - \mathbf{q} \rangle$

**Example of Bregman divergence**

Some basic properties of Bregman divergence

Summary

## Example of Bregman divergence: $\ell_2$ -squared norm

Let  $f(\mathbf{z}) = \|\mathbf{z}\|_2^2$ .

$$\begin{aligned} B_f(\mathbf{x}; \mathbf{y}) &= f(\mathbf{x}) - f(\mathbf{y}) - \nabla f(\mathbf{y})^\top (\mathbf{x} - \mathbf{y}) \\ &\stackrel{\nabla f(\mathbf{x})=2\mathbf{x}}{=} \|\mathbf{x}\|_2^2 - \|\mathbf{y}\|_2^2 - 2\mathbf{y}^\top (\mathbf{x} - \mathbf{y}) \\ &= \|\mathbf{x}\|_2^2 - 2\mathbf{x}^\top \mathbf{y} + \|\mathbf{y}\|_2^2 \\ &= \|\mathbf{x} - \mathbf{y}\|_2^2 \end{aligned}$$

Hence  $\|\mathbf{x} - \mathbf{y}\|_2^2$  is the Bregman divergence of  $\|\mathbf{x}\|_2^2$ .

## Example of Bregman divergence: negative Shannon entropy

- Let  $f(\mathbf{x}) = \sum_i x_i \log x_i$  Shannon entropy is  $-f(\mathbf{x})$
- We focus on  $\mathbf{x} \in \Delta := \{\mathbf{v} \mid \langle \mathbf{v}, \mathbf{1} \rangle = 1\}$  probability simplex because entropy works on probability
- Fact:  $\frac{\partial}{\partial x} x \log x = \log x + 1$
- Showing Bregman divergence of negative Shannon entropy is KL divergence

$$\begin{aligned} B_f(\mathbf{x}; \mathbf{y}) &= f(\mathbf{x}) - f(\mathbf{y}) - \nabla f(\mathbf{y})^\top (\mathbf{x} - \mathbf{y}) \\ &= \sum x_i \log x_i - \sum y_i \log y_i - \sum (\log y_i + 1)(x_i - y_i) \\ &= \sum x_i \log x_i - \sum (\log y_i + 1)(x_i) \\ &= \sum x_i \log x_i - \sum x_i \log y_i + \underbrace{\sum x_i}_{=1} \\ &= \sum x_i (\log x_i - \log y_i) \\ &= \sum x_i \log \frac{x_i}{y_i} \end{aligned}$$

Hence  $\sum_i x_i \log \frac{x_i}{y_i}$  is the Bregman divergence of  $\sum_i x_i \log x_i$ .

## Example of Bregman divergence: matrix entropy

- Let  $f(\mathbf{X}) = \text{Tr}(\mathbf{X} \log \mathbf{X})$ , then  $B_F(\mathbf{X}; \mathbf{Y}) = \text{Tr}(\mathbf{X} \log \mathbf{X} - \mathbf{X} \log \mathbf{Y} - \mathbf{X} + \mathbf{Y})$ 
  - This is basically eigen-value version of the vector Shannon entropy
  - This is called von Neumann divergence
  
- Let  $f(\mathbf{X}) = -\log \det(\mathbf{X})$ , then  $B_F(\mathbf{X}; \mathbf{Y}) = \log \det(\mathbf{Y}) - \log \det(\mathbf{X}) + \text{Tr}(\mathbf{Y}^{-1}(\mathbf{X} - \mathbf{Y}))$ 
  - This is called logdet divergence

# Contents

Jansen inequality & 1st-order Taylor series

Bregman divergence  $B(\mathbf{p}, \mathbf{q}) := f(\mathbf{p}) - f(\mathbf{q}) - \langle \nabla f(\mathbf{q}), \mathbf{p} - \mathbf{q} \rangle$

Example of Bregman divergence

Some basic properties of Bregman divergence

Summary

## Facts about $B_f(\mathbf{p}; \mathbf{q}) := f(\mathbf{p}) - f(\mathbf{q}) - \langle \nabla f(\mathbf{q}), \mathbf{p} - \mathbf{q} \rangle \dots$ 1/2

- $B_f$  does not satisfy the triangle inequality
- $B_f$  is asymmetric:  $B_f(\mathbf{p}; \mathbf{q}) \neq B_f(\mathbf{q}; \mathbf{p})$
- $B_f \geq 0$  because of convexity of  $f$
- $B_f(\mathbf{p}; \mathbf{q})$  is convex in  $\mathbf{p}$
- $B$  is conic:  $B_{\alpha f + \beta g} = \alpha B_f + \beta B_g$  for positive  $\alpha, \beta$
- $B$  is affine invariant:  $B_{f+g} = B_f$  if  $g$  is affine (e.g.,  $g(\mathbf{x}) = \mathbf{A}\mathbf{x} + \mathbf{c}$ )
- If  $f$  is not differentiable, replace  $\nabla f(\mathbf{q})$  by any subgradient of  $g$  at  $\mathbf{q}$ .
- $\nabla_{\mathbf{p}} B_f(\mathbf{p}; \mathbf{q}) = \nabla f(\mathbf{p}) - \nabla f(\mathbf{q})$
- For  $\min f$ , if  $\mathbf{q} = \operatorname{argmin} f$ , then  $B_f(\mathbf{p}; \mathbf{q}) = f(\mathbf{p}) - f^*$

not  $\mathbf{q}$

NOT linear!

## Facts about $B_f(\mathbf{p}; \mathbf{q}) := f(\mathbf{p}) - f(\mathbf{q}) - \langle \nabla f(\mathbf{q}), \mathbf{p} - \mathbf{q} \rangle \dots$ 2/2

- **Conjugate duality**

Let  $f^*$  be the convex conjugate of  $f$ . Then  $B_f(\mathbf{p}; \mathbf{q}) = B_{f^*}(\nabla f(\mathbf{p}); \nabla f(\mathbf{q}))$

- (Banerjee et al. 2005)<sup>1</sup> given a random vector  $\mathbf{z}$ , the mean vector minimizes the expected distance from  $\mathbf{z}$  is a Bregman divergence.

- Bregman projection

- Given a convex set  $\mathcal{C}$ , the Bregman projection of a point  $\mathbf{x}$  onto  $\mathcal{C}$  under the Bregman divergence with anchor  $\mathbf{x}_0$  is

$$\mathbf{x}^* = \operatorname{argmin}_{\mathbf{x} \in \mathcal{C}} B_f(\mathbf{x}; \mathbf{x}_0)$$

- **Generalised Pythagorean theorem**

$$B_f(\mathbf{y}; \mathbf{x}_0) \geq B_f(\mathbf{y}; \mathbf{x}^*) + B_f(\mathbf{x}^*; \mathbf{x}_0)$$

- Bregman projection is the key of **mirror descent**.

---

<sup>1</sup>Banerjee, Gou, and Wang, On the optimality of conditional expectation as a Bregman predictor, 2005



# Contents

Jansen inequality & 1st-order Taylor series

Bregman divergence  $B(\mathbf{p}, \mathbf{q}) := f(\mathbf{p}) - f(\mathbf{q}) - \langle \nabla f(\mathbf{q}), \mathbf{p} - \mathbf{q} \rangle$

Example of Bregman divergence

Some basic properties of Bregman divergence

Summary

## Last page - summary

- Convexity:  $f(\lambda \mathbf{y} + (1 - \lambda)\mathbf{x}) \leq \lambda f(\mathbf{y}) + (1 - \lambda)f(\mathbf{x}) \iff f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle$
- Bregman divergence:  $B_f(\mathbf{y}; \mathbf{x}) := f(\mathbf{y}) - f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle$ 
  - $\mathbf{y}$  is variable
  - $\mathbf{x}$  is anchor
- Bregman Divergence = the gap between 1st-order Taylor series and the original function
- $\|\mathbf{x} - \mathbf{y}\|_2^2$  is the Bregman divergence of  $\|\mathbf{x}\|_2^2$
- Some basic properties of Bregman divergence (proof to be added in future version)

End of document