



# Learning Set-to-Set Bregman Divergence with Permutation-Invariant NNs

 [Takahiro Kawashima](#) (ZOZO Research)

 RIMS-ML 2026

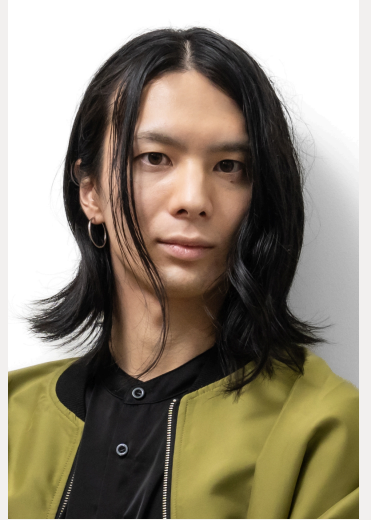
 May 28, 2026

**ZOZONEXT**

# Introduction of Introduction

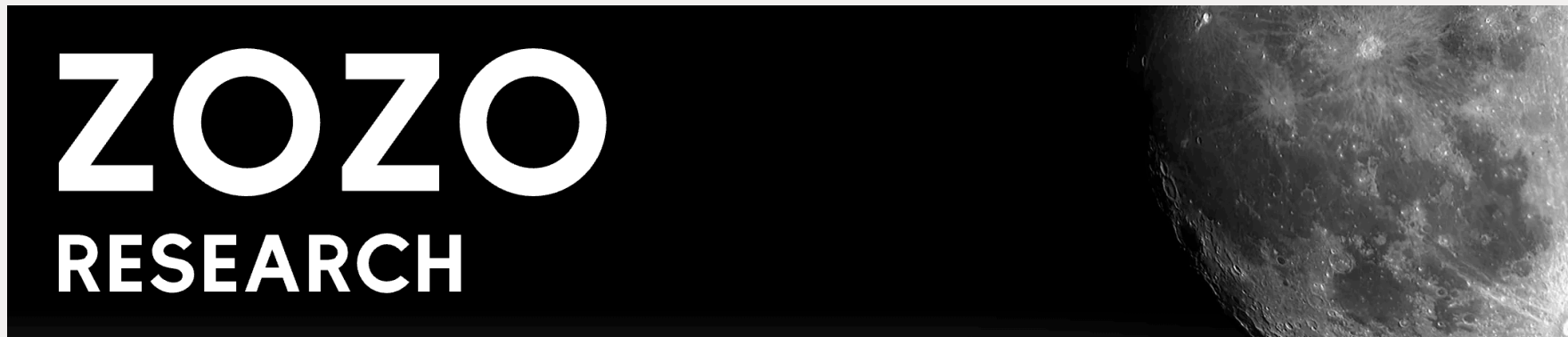
# Who?

- Name: **Takahiro Kawashima** (川島貴大)
- Interests: Statistical ML, Probabilistic Modeling,  
Bayesian Methods
- 2012-2017: Kobe City College of Tech. (神戸高専)
- 2017-2021: Univ. of Electro-Communications (B.Eng. & M.Eng.)
- 2021-2024: SOKENDAI (Ph.D. in Statistical Science)
- 2019-2024: National Center of Neurology and Psychiatry
- 2024-Present: **ZOZO Research** (ZOZO NEXT, Inc.)



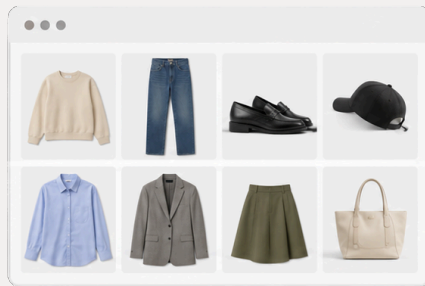
# About ZOZO Research

- Mission: “**Quantifying Fashion**” (ファッションを数値化する)
- R&D division of ZOZO
- Focusing on the field of **Fashion Tech**
  - Including CV, NLP, Recommender Systems, HCI, Stats, OR, . . .



# Set-valued Data in Fashion Tech

- **Set-valued data** often appear in Fashion Tech
  - Buying some set of items, choosing outfits, . . .



We also focus on  
**Set-data × AI**

# Today's Talk

- Novel, flexible, and learnable set-to-set discrepancy metric
  - Based on the idea of Bregman divergence
- Published in ICLR 2025
- Worked with Masanari Kimura (Univ. Melbourne), Tasuku Soma (ISM), Hideitsu Hino (ISM & Waseda Univ.)

## DIFFERENCE-OF-SUBMODULAR BREGMAN DIVERGENCE

**Masanari Kimura\***  
School of Mathematics and Statistics  
The University of Melbourne  
Victoria, Australia  
m.kimura@unimelb.edu.au

**Takahiro Kawashima\***  
ZOZO Research  
ZOZO Next, Inc.  
Chiba, Japan  
takahiro.kawashima@zozo.com

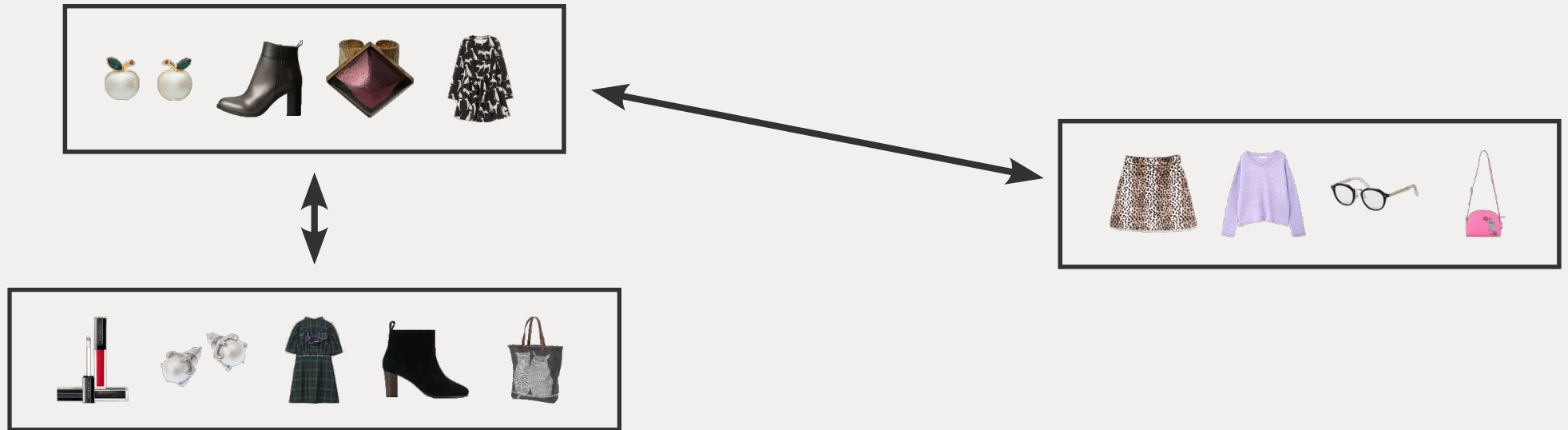
**Tasuku Soma & Hideitsu Hino**  
Institute of Statistical Mathematics / RIKEN AIP  
Tokyo, Japan  
{soma, hino}@ism.ac.jp

# Introduction

# Distance Between Finite Sets

Most ML algs implicitly or explicitly assume a distance structure

- Rich distance metrics between sets are desired!
- Applications: outfit recommendation, point cloud proc, . . .



# Distance/Similarity Metrics for Sets

Numerous metrics of set-to-set distance/similarity exist:

Table 2 Definitions of Measures for binary data

$S_{JACCARD} = \frac{a}{a+b+c}$	(1)
$S_{JACC} = \frac{2a}{2a+b+c}$	(2)
$S_{CHENKOROSKI} = \frac{2a}{2a+b+c}$	(3)
$S_{JW-JACCARD} = \frac{3a}{3a+b+c}$	(4)
$S_{JSHAPIRO} = \frac{2a}{(a+b)+(a+c)}$	(5)
$S_{SHAPIRO-DELLI-I} = \frac{a}{a+2b+2c}$	(6)
$S_{SHAPIRO-MICHERINI} = \frac{a+d}{a+b+c+d}$	(7)
$S_{SHAPIRO-DELLI-II} = \frac{2(a+d)}{2a+b+c+2d}$	(8)
$S_{SHAPIRO-DELLI-III} = \frac{a+d}{a+2(b+c)+d}$	(9)
$S_{FAITH} = \frac{a+0.5d}{a+b+c+d}$	(10)
$S_{SHAPIRO-DELLI-IV} = \frac{a+d}{a+0.5(b+c)+d}$	(11)
$S_{INTERSECTION} = a$	(12)
$S_{ANDERSON} = a+d$	(13)
$S_{RUSSELLKOROSKI} = \frac{a}{a+b+c+d}$	(14)
$D_{BARROU} = b+c$	(15)
$D_{ECCO} = \sqrt{b+c}$	(16)
$D_{RUSSELL-DELLI} = \sqrt{(b+c)^2}$	(17)
$D_{BARROU} = (b+c)^{\frac{1}{2}}$	(18)
$D_{MAXLETTEN} = b+c$	(19)
$D_{MAX-MONSTERN} = \frac{b+c}{a+b+c+d}$	(20)
$D_{LITTRONCK} = b+c$	(21)
$D_{MINORSKI} = (b+c)^{\frac{1}{2}}$	(22)

$D_{DIA} = \frac{(b+c)}{4(a+b+c+d)}$	(23)
$D_{DIAPYCNIS} = \frac{(b+c)^2}{(a+b+c+d)^2}$	(24)
$D_{DIAPYCNIS} = \frac{a(b+c)-(b-c)^2}{(a+b+c+d)^2}$	(25)
$D_{PASTENOPFEN} = \frac{4bc}{(a+b+c+d)^2}$	(26)
$D_{LANGEWILLIAMS} = \frac{b+c}{(2a+b+c)}$	(27)
$D_{BRAFACIUM} = \frac{b+c}{(2a+b+c)}$	(28)
$D_{MILLER} = 2 \sqrt{1 - \frac{a}{\sqrt{(a+b)(a+c)}}$	(29)
$D_{DIA} = \sqrt{2 \left( 1 - \frac{a}{\sqrt{(a+b)(a+c)}} \right)}$	(30)
$S_{COSINE} = \frac{a}{\sqrt{(a+b)(a+c)}}$	(31)
$S_{SHAPIRO-DELLI-V} = \log a - \log n - \log \left( \frac{a+b}{n} \right) - \log \left( \frac{a+c}{n} \right)$	(32)
$S_{COSINE-I} = \frac{a}{\sqrt{(a+b)(a+c)}}$	(33)
$S_{FISHO} = \frac{na}{(a+b)(a+c)}$	(34)
$S_{FISHO} = \frac{n(a-0.5)^2}{(a+b)(a+c)}$	(35)
$S_{SHAPIRO-DELLI-VI} = \frac{a^2}{(a+b)(a+c)}$	(36)
$S_{SHAPIRO-DELLI-VII} = \frac{a}{0.5(ab+ac)+bc}$	(37)
$S_{YULE} = \frac{a}{((a+b)(a+c))^{0.5}}$	(38)
$S_{YULE} = \frac{a^2-bc}{(a+b)(a+c)}$	(39)
$S_{YULE} = \frac{na-(a+b)(a+c)}{na+(a+b)(a+c)}$	(40)
$S_{YULE} = \frac{a}{2} \frac{(2a+b+c)}{(a+b)(a+c)}$	(41)
$S_{YULE} = \frac{a}{2} \frac{1}{a+b} + \frac{1}{a+c}$	(42)
$S_{YULE} = \frac{a}{a+b} + \frac{a}{a+c}$	(43)
$S_{YULE} = \frac{ad-bc}{\sqrt{(a+b)(a+c)}}$	(44)
$S_{YULE} = \frac{a}{\min(a+b, a+c)}$	(45)
$S_{YULE} = \frac{a}{\max(a+b, a+c)}$	(46)

$S_{SHAPIRO-DELLI-VIII} = \frac{a}{\sqrt{(a+b)(a+c)}} \frac{\max(a+b, a+c)}{2}$	(47)
$S_{FISHO-II} = \frac{na-(a+b)(a+c)}{n \min(a+b, a+c) - (a+b)(a+c)}$	(48)
$S_{SHAPIRO-DELLI-IX} = \frac{\frac{a}{(a+b)} + \frac{a}{(a+c)} + \frac{d}{(b+d)} + \frac{d}{(b+d)}}{4}$	(49)
$S_{COSINE} = \frac{a+d}{\sqrt{(a+b)(a+c)(b+d)(c+d)}}$	(50)
$S_{PEARSON-I} = x^2$ where $x^2 = \frac{n(ad-bc)^2}{(a+b)(a+c)(b+d)(c+d)}$	(51)
$S_{PEARSON-II} = \left( \frac{x^2}{n+x^2} \right)^2$	(52)
$S_{PEARSON-III} = \left( \frac{p}{n+p} \right)^2$ where $p = \frac{ad-bc}{\sqrt{(a+b)(a+c)(b+d)(c+d)}}$	(53)
$S_{PEARSON-IV} = \frac{ad-bc}{\sqrt{(a+b)(a+c)(b+d)(c+d)}}$	(54)
$S_{PEARSON-V} = \text{Cos} \left( \frac{\pi \sqrt{bc}}{\sqrt{ad+bc}} \right)$	(55)
$S_{SHAPIRO-DELLI-X} = \frac{a+d}{b+c}$	(56)
$S_{SHAPIRO-DELLI-XI} = \frac{ad}{(a+b)(a+c)(b+d)(c+d)^{0.5}}$	(57)
$S_{COS} = \frac{\sqrt{2(ad-bc)}}{\sqrt{(ad-bc)^2 - (a+b)(a+c)(b+d)(c+d)}}$	(58)
$S_{YULE} = \log \left( \frac{n(ad-bc) - \frac{a^2}{2}}{(a+b)(a+c)(b+d)(c+d)} \right)$	(59)
$S_{COSINE-I} = \frac{ad}{\sqrt{(a+b)(a+c)(b+d)(c+d)}}$	(60)
$S_{YULE} = \frac{ad-bc}{ad+bc}$	(61)
$D_{YULE} = \frac{2bc}{ad+bc}$	(62)
$S_{YULE} = \frac{\sqrt{ad-bc}}{\sqrt{ad+bc}}$	(63)
$S_{YULE} = \frac{a}{b+c}$	(64)
$S_{YULE} = \frac{a}{(a+b)+(a+c)-a}$	(65)
$S_{YULE} = \frac{ad-bc}{(a+b+c+d)^2}$	(66)
$S_{YULE} = \frac{(a+d)-(b+c)}{a+b+c+d}$	(67)
$S_{YULE} = \frac{d(ad-bc)}{(a+d)^2 + (b+c)^2}$	(68)
$S_{SHAPIRO-DELLI-XII} = \frac{\sigma - \sigma'}{2\sigma - \sigma'}$	(69)
$\sigma = \max(b, b) + \max(d, d) + \max(c, c) + \max(d, d)$	
$\sigma' = \max(a, c, b+d) + \max(a, b, c+d)$	

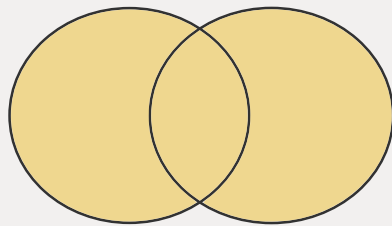
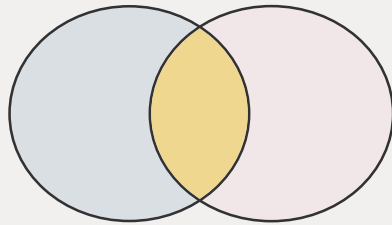
$S_{ANDERSON} = \frac{\sigma - \sigma'}{2\sigma}$	(70)
$S_{SHAPIRO-DELLI-XIII} = \frac{\sqrt{ad+a}}{\sqrt{ad+a+b+c}}$	(71)
$S_{SHAPIRO-DELLI-XIV} = \frac{\sqrt{ad+a-(b+c)}}{\sqrt{ad+a+b+c}}$	(72)
$S_{PEARCE} = \frac{ab+bc}{ab+2bc+cd}$	(73)
$S_{YULE} = \frac{n^2(na-(a+b)(a+c))}{(a+b)(a+c)(b+d)(c+d)}$	(74)
$S_{YULE} = \frac{\frac{a}{c} + \frac{a(c+d)}{c(a+b)}}{\frac{a}{c} + \frac{a(c+d)}{c(a+b)}}$	(75)
$S_{YULE} = \left  \frac{\frac{a}{c} + \frac{a(c+d)}{c(a+b)}}{\frac{a}{c} + \frac{a(c+d)}{c(a+b)}} \right $	(76)

from Choi et al. (2010)

# Difficulty of Distance Metrics for Sets

Most existing metrics are defined through **counting**

➤ such as the sizes:  $|X \cap Y|$ ,  $|X \cup Y|$ ,  $|X \setminus Y|$ , ...



Jaccard Index

If the ground set is large and  $X, Y$  are small, similarity/distance easily becomes  $0/\infty$  🙄

➤ **Richer distance metrics** for finite sets are desired!

# Our Contributions

We majorly contribute to:

- Develop a **set-to-set divergence** with richer expressiveness
  - > Well-definedness and richness are theoretically guaranteed
- Propose a framework to learn the divergence with **permutation-invariant NNs**

This talk covers from necessary background to our latest work

# Bregman Divergence

# Divergence

**Divergence**: a core concept in ML&stats, measuring the discrepancy between two objects

## Def: Divergence

If  $D : \Omega \times \Omega \rightarrow \mathbb{R}$  satisfies:

- $D(x, y) \geq 0$  for all  $x, y \in \Omega$  (nonnegativity),
- $D(x, y) = 0$  iff  $x = y$  for all  $x \in \Omega$  (identifiability),

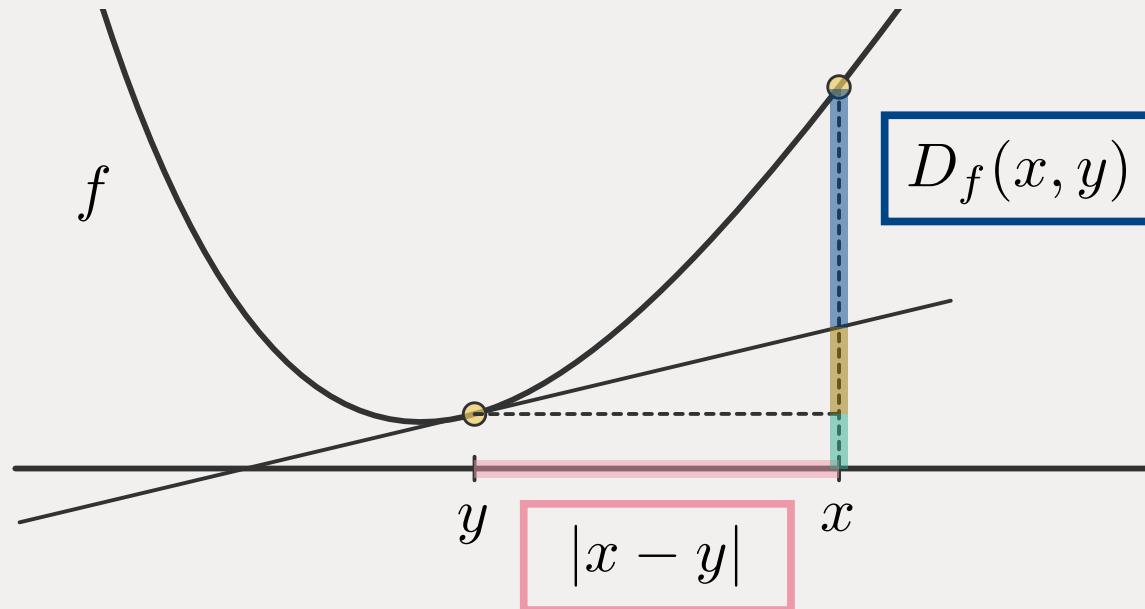
then  $D$  is a **divergence** on  $\Omega$ .

➤ Generalization of distance; not necessarily symmetric

# Bregman Divergence

**Bregman divergence** constructs a divergence from a differentiable and strictly convex function  $f$

➤ Every tangent plane lower-bounds  $f$ ; nonneg. metric is induced



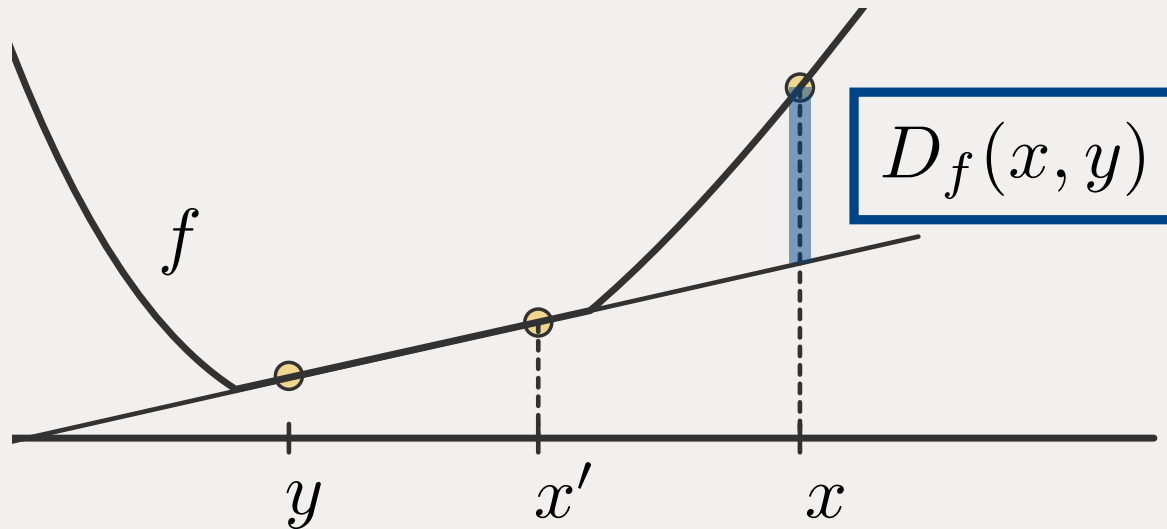
$$D_f(x, y) := \text{inner prod.} \\ f(x) - f(y) - \langle \nabla f(y), x - y \rangle$$

Generalizes KL-div., sq-L2 dist., ...

# Necessity of Strict Convexity

If  $f$  is convex but not strictly convex,  $D_f(x, y)$  can be 0 for some  $x \neq y$

➤ The identifiability condition is violated



$$D_f(x, y) := \text{inner prod.} \\ f(x) - f(y) - \langle \nabla f(y), x - y \rangle$$

# Application of Bregman Divergence

Bregman divergence widely appears in ML&stats

- Clustering (k-means, EM, . . .)
- Bayesian Inference (variational methods, . . .)
- Optimization (mirror descent, optimal transport, . . .)
- Information Geometry (dually flat manifolds, . . .)
- . . .

Flexible, tractable, and theoretically grounded

➤ Motivating the development of Bregman divergences for sets!

# Existing Work: Submodular-Bregman Divergence

# Generalized Bregman Divergence

For defining Bregman-like divergences, differentiability is not necessary

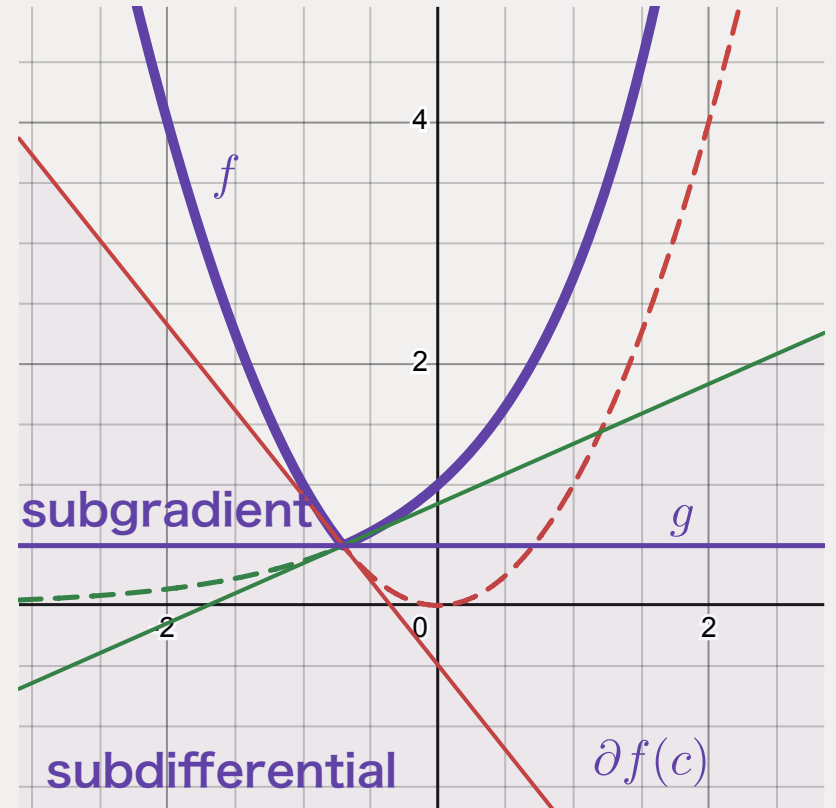
e.g.,  $f(x) = \max(x^2, e^x)$

➤ non-differentiable at  $c := -0.703\dots$

By choosing a subgradient  $g \in \partial f(c)$ , a generalized Bregman divergence

$$D_{f,g}(x, c) := f(x) - f(c) - \langle g, x - c \rangle$$

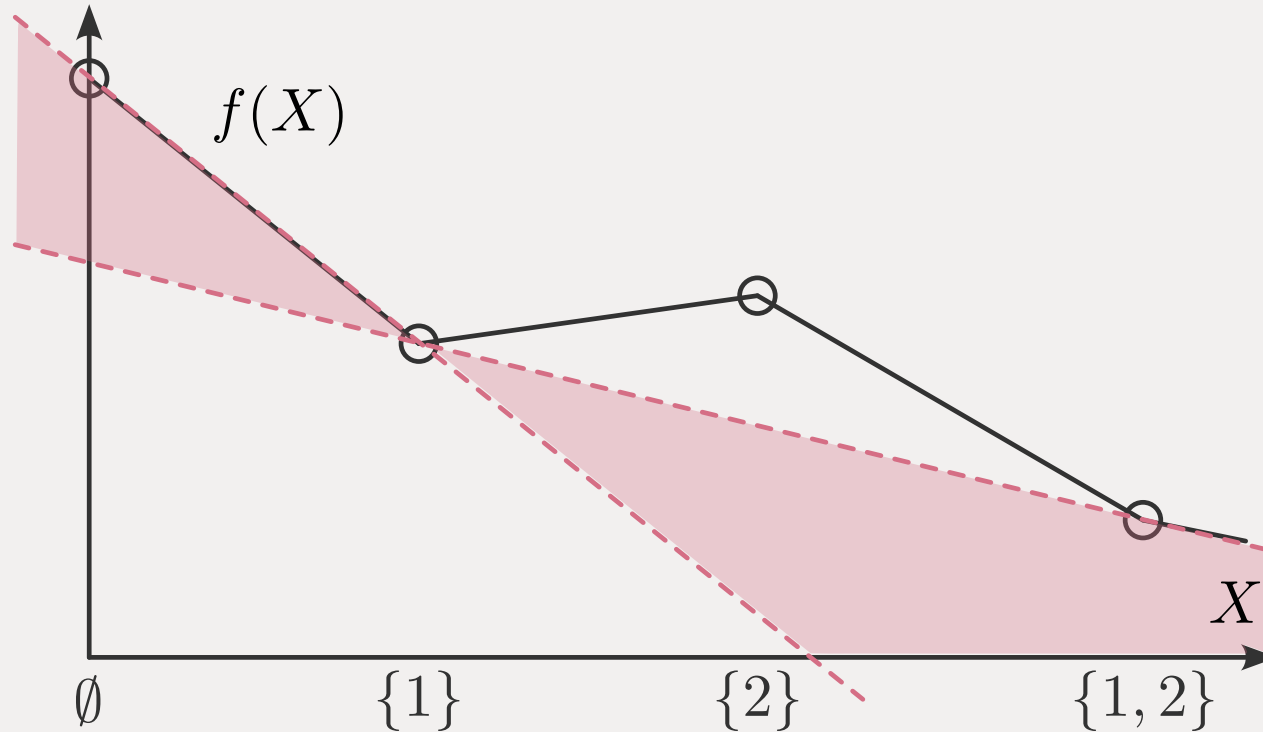
can be defined



# Bregman Divergence in Discrete Space

Bregman-like divergence can be defined even in discrete spaces

➤ as long as a lower-bounding tangent plane (subgradient) exists



# Submodular Functions

Our aim: define a Bregman-like divergence between finite sets

What type of set functions can be lower-bounded linearly?

➤ Submodular functions: discrete analog of convex functions!

## Def: Submodular Function

If  $f : 2^V \rightarrow \mathbb{R}$  satisfies  $f(X) + f(Y) \geq f(X \cup Y) + f(X \cap Y)$   
for every  $X, Y \subseteq V$ , then  $f$  is a **submodular function** on  $V$ .

When the inequality is strict except for  $X \subseteq Y$  or  $Y \subseteq X$ ,  
 $f$  is called **strictly submodular**.

# Subgradients of Submodular Functions

Linear functions in the discrete space: **modular functions**

➤  $m : 2^V \rightarrow \mathbb{R}$  is modular  $\iff m$  and  $-m$  are submodular

Every modular function is represented as  $m(X) = \sum_{i \in X} m(\{i\})$

➤  $m$  is equated with the vector  $m \in \mathbb{R}^N$ , where  $N = |V|$

**Edmonds (1970)**: greedy alg. for finding the tightest modular lower-bound of submodular  $f$  at  $Y \subseteq V$

Subgradients are easily found for submodular functions!

# Submodular-Bregman Divergence

Iyer & Bilmes (2012): Bregman-like divergence between finite sets with a submodular function

## Submodular-Bregman Divergence

$$D_f(X, Y) = f(X) - f(Y) - \langle h_Y, 1_X - 1_Y \rangle$$

$h_Y$  : Subgradient of  $f$  at  $Y$

$1_X$  : Indicator vector of  $X$  ( $1_X(i) = \mathbb{I}[i \in X]$ )

Note: the inner product is the standard one in  $\mathbb{R}^N$

# Submodular-Bregman Divergence: Example

Concrete discrepancy metrics are constructed with specific  $f, h_Y$

e.g.,  $f(X) = |X|, h_Y = 2 \times 1_Y$

$$\begin{aligned} D_f(X, Y) &= f(X) - f(Y) - \langle h_Y, 1_X - 1_Y \rangle \\ &= |X| - |Y| - 2\langle 1_Y, 1_X - 1_Y \rangle \\ &= |X| - |Y| - 2(|X \cap Y| - |Y|) \\ &= |X| + |Y| - 2|X \cap Y| \\ &= |X \setminus Y| + |Y \setminus X| \end{aligned}$$

➤ Hamming distance between sets!

# Submodular-Bregman Divergence: Issues

Submodular-Bregman provides an interesting concept but. . .

- How to choose a submodular function  $f$ ?
- For non-strictly submodular  $f$ , the tightest subgradient does not offer identifiability
  - › To construct a divergence, a loose subgradient is needed;  
 $f(X) = |X|, h_Y = 2 \times 1_Y$  gives  $f(V) \neq h_Y(V) \Rightarrow$  unregularized

Our work: based on **strictly-submodular functions**

- › Tight subgradients can be used, and increases learnability!

# **Our Work: Difference-of-submodular Bregman Divergence**

# Strict Submodularity and Divergence

## Submodular-Bregman Divergence

$$D_f(X, Y) = f(X) - f(Y) - \langle h_Y, 1_X - 1_Y \rangle$$

$h_Y$  : Subgradient of  $f$  at  $Y$

If  $f$  is strictly submodular, the tightest  $h_Y$  offers identifiability  
> properly defines a divergence!

Does this conditioning reduce the expressiveness?

# In Fact . . .

## **Thm: Representation Power** (casual)

The representation power of the divergence between finite sets defined as:

$$D_f(X, Y) = f(X) - f(Y) - \langle h_Y, 1_X - 1_Y \rangle$$

increases with the function class of  $f$ .

$f$  : Set function

$h_Y$  : Subgradient of  $f$  at  $Y$

➤ **Broader function class induces richer divergences**

# Broadening Representation Power

Make the representation power richer by two remarkable properties of submodular functions:

## 1. Existence of **Supergradients**

- › Supergradients also exist for any submodular function

## 2. **Strong Difference-of-Submodular (DS) Decomposition**

- › Every set function can be decomposed as the difference of two strictly submodular functions

Both are specific properties in the discrete space

# Supergradients

✍ **Thm: Supergradients** (Iyer & Bilmes, 2012)

For a submodular function  $f$ , the following modular functions are supergradients of  $f$  at  $Y$ :

Name	Notation	$f(j   X) := f(X \cup \{j\}) - f(X)$	
		$j \in Y$	$j \notin Y$
grow	$\hat{g}_Y(j)$	$f(j   V \setminus \{j\})$	$f(j   Y)$
shrink	$\check{g}_Y(j)$	$f(j   Y \setminus \{j\})$	$f(j   \emptyset)$
bar	$\bar{g}_Y(j)$	$f(j   V \setminus \{j\})$	$f(j   \emptyset)$

If  $f$  is strictly submodular, all the above are strict upperbounds

# Strong DS Decomposition

 **Thm: Strong DS Decomposition** (Li & Du, 2020)

Any set function  $f$  can be decomposed as

$$f = f^1 - f^2,$$

where  $f^1, f^2$  are strictly submodular

Continuous analog: Difference-of-Convex (DC) structure

➤ But the function class is not broad in the continuous space

Note: Finding the decomposition may be NP-hard

# Difference-of-submodular Bregman Div.

Bregman-like divergence with any set function  $f$  can be defined!

1. Strong DS decomposition:  $f = f^1 - f^2$
2. Take the tightest subgradient  $h_Y^1$  for  $f^1$  at  $Y \subseteq V$
3. Take the grow/shrink/bar supergradient  $g_Y^2$  for  $f^2$  at  $Y \subseteq V$
4. Define the divergence as:

$$D_f(X, Y) := \underbrace{f^1(X) - f^1(Y) - \langle h_Y^1, 1_X - 1_Y \rangle}_{\text{submod. Bregman of } f^1} - \underbrace{(f^2(X) - f^2(Y) - \langle g_Y^2, 1_X - 1_Y \rangle)}_{\text{supmod. Bregman of } f^2}$$
$$= f(X) - f(Y) - \langle h_Y^1 - g_Y^2, 1_X - 1_Y \rangle$$

# DBD construction with NNs

How to construct the divergence in practice?

➤ Learning the set function from data

## Permutation-Invariant NNs

Neural networks whose output is invariant to the permutation of the input sequence; e.g.,

$$f_{\text{NN}}(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3) = f_{\text{NN}}(\mathbf{x}_2, \mathbf{x}_1, \mathbf{x}_3) = f_{\text{NN}}(\mathbf{x}_3, \mathbf{x}_2, \mathbf{x}_1) = \dots$$

Prepare two permutation-invariant NNs with strict submodularity a priori, and model  $f^1, f^2$  with them

# $\varepsilon$ -PointNet

We also propose an architecture of permutation-invariant NNs

## $\varepsilon$ -PointNet

$$f_{\text{PN},\varepsilon}([\mathbf{x}_i]_{i \in X}) = \gamma \left( \begin{array}{c} K\text{-dim intermediate outputs of } \varphi \\ \varepsilon \log \sum_{i \in X} e^{\varphi_1(\mathbf{x}_i)/\varepsilon}, \dots, \varepsilon \log \sum_{i \in X} e^{\varphi_K(\mathbf{x}_i)/\varepsilon} \end{array} \right)$$

$\gamma, \varphi_1, \dots, \varphi_K$  : NN or fixed functions,  $\varepsilon > 0$  : hyperparameter

- Each input of  $\gamma$  approaches  $\max_{i \in X} \varphi_k(\mathbf{x}_i)$  as  $\varepsilon \rightarrow 0$
- Strict submodularity is guaranteed when  $\gamma$  is a weighted sum

# Learning the Divergence

$f$  is learnable via some metric learning loss, such as **triplet loss**

## Triplet Loss

$$\mathcal{L}(f) := \sum_{i=1}^n \max(D_f(X_A^i, X_P^i) - D_f(X_A^i, X_N^i), 0)$$

$X_A^i$ : Anchor set,

$X_P^i$ : Positive set,









$X_N^i$ : Negative set

- Makes the anchor-positive discrepancies smaller, and the anchor-negative discrepancies larger
- Learns  $f$  unlike ordinary metric learning, which moves instances

# Experiments

# Sanity Check : MNIST

Learning to make the distances between sets closer, if the sets contain same-labeled instances

Query set	Reference sets	DBD
	  	0.173 0.309 0.772
	  	0.237 0.677 1.304



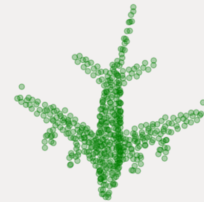
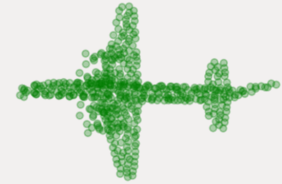

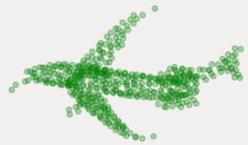






Looks Good 

# Point Cloud Retrieval (ModelNet40)

**Point Cloud** : Representative data with set structure


Take top- $K$  closest point clouds to a query cloud

➤ Assess whether the taken clouds are same-class as the query

Query set	Top-5 retrieval results				
					
					

# Point Cloud Retrieval (ModelNet40)

Near-SoTA performance is achieved with a tiny MLP

- Only uses 64×64 intermediate layers; trainable on CPU!
- Suggests the effectiveness of the proposed framework 

Our method  
w/ 3-types of  
supergrads

Method	mAP
<i>grow</i> -DBD w/ decomposition	90.13(±0.75)
<i>shrink</i> -DBD w/ decomposition	90.20(±0.77)
<i>bar</i> -DBD w/ decomposition	86.09(±0.85)
<i>grow</i> -DBD w/o decomposition	88.12(±0.80)
<i>shrink</i> -DBD w/o decomposition	88.20(±0.81)
<i>bar</i> -DBD w/o decomposition	83.57(±0.97)
Densepoint (Liu et al., 2019)	89.68(±0.88)
MVTN (Hamdi et al., 2021)	92.9*

# Summary

# Summary

- **Difference-of-submodular Bregman Divergence** is proposed
  - › Extension of submodular Bregman divergence
  - › **Richer expressiveness** is theoretically guaranteed
- Learning the divergence with **permutation-invariant NNs**
  - ›  $\epsilon$ -PointNet: **can model strictly submodular functions**
- Numerical experiments on MNIST & point cloud data
  - › **Near-SoTA performance is achieved only with tiny MLPs!**