



The Prediction of Geostrategic Event

by O. Croissant

7th of jan 2024



Plan

12/29/2023

2

- 2) EVENT PREDICTION DATASET
- 3) TRANSFORMER PRESENTATION
- 4) LANGUAGE MODELS AND SENTIENT AI
- 5) GRAM MATRIX AND GRADIENT FLOWS
- 6) KL DIVERGENCE
- 7) THE BENAMOU BRENIER APPROACH
- 8) GRADIENT FLOWS: THE DEMONSTRATIONS
- 9) GRADIENT FLOWS: THE THEORY
- 11) MEAN-FIELD THEORIES
- 12) RELATIONSHIPS WITH QUANTUM STATISTICAL PHYSICS
- 15) MULTIVARIATE RENYI INFORMATION
- 16) INFORMATION THEORY
- 17) OPTIMAL TRANSPORT AND NLP LOSSES

[To Summary](#)



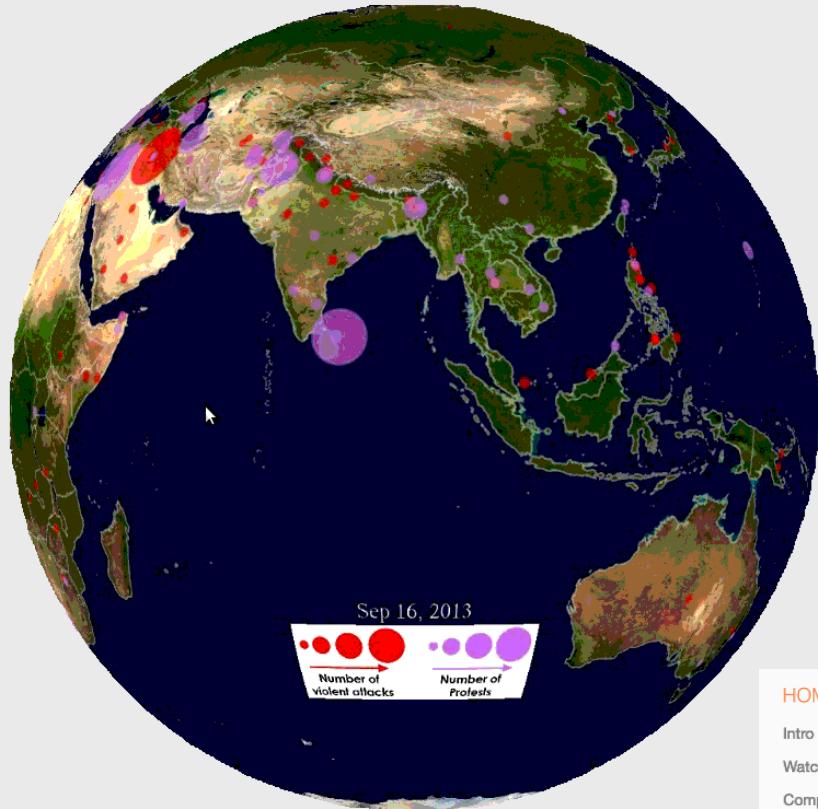
12/29/2023

3

2) EVENT PREDICTION DATASET

[To Summary](#)

GDELT Project



Kalev H. Leetaru

One of Foreign Policy Magazine's Top 100 Global Thinkers of 2013,



300 million event records in over 300 categories covering the entire world from 1979 to present

Hundreds of thousands of broadcast, print, and online news sources from every corner of the globe in more than 100 languages

HOME	BLOG	DATA	SOLUTIONS
Intro	The GDELT Blog	Intro	Intro
Watching		GDELT Analysis Service	Situational Awareness
Computing		Google BigQuery	Influencer Networks
Downloading		Raw Data Files	Risk Assessment & Global Trends
Blogging		Documentation	Policy Reaction
			Humanitarian & Crisis Response

GDELT Data



GDELT Analysis Service

Free cloud-based service that offers a variety of tools and services to allow you to visualize, explore, and export GDELT - a great way to get started using GDELT for the first time.



Google BigQuery

Query, export, and even conduct sophisticated analyses and modeling of the entire dataset using standard SQL, with even the most complex queries returning in near-realtime.



Raw Data Files

Advanced users and those with unique use cases can download the entire underlying event and graph datasets in CSV format - over 2.5TB for last year alone.

GDELT 1.0 Event Database

Daily Updates

All GDELT Event Files

- md5sums
- filesizes
- GDELTMASTEREDUCEDV2.1979-2013.zip (1.1GB) (MD5: f6fa17e955e3593c9da427c07b545d)
- 20221205.export.CSV.zip (5.7MB) (MD5: 73dc0c5565e3ac010495c3a75d1f6d51)
- 20221204.export.CSV.zip (3.1MB) (MD5: 3711b97dac26e39fb1c860eb1e261fd)
- 20221203.export.CSV.zip (3.9MB) (MD5: d425d05ececb69c25c73e63175e51)
- 20221202.export.CSV.zip (6.6MB) (MD5: ac387488b1451a5a925749f4c8378)
- 20221201.export.CSV.zip (6.9MB) (MD5: 8c30a50f08fb62c3dc704c5b4fe)
- 20221130.export.CSV.zip (7.0MB) (MD5: 836d0f8ca3f581f0519c196c554509)
- 20221129.export.CSV.zip (6.9MB) (MD5: 2e2f164883eddd9996831oba584d778)
- 20221128.export.CSV.zip (5.4MB) (MD5: 272606d4f805117bd8c96342811435eab)
- 20221127.export.CSV.zip (2.8MB) (MD5: 789fee593ec59ffaa60332259f689326d)
- 20221126.export.CSV.zip (3.1MB) (MD5: 7dc39f894f47ab090d2d33337360f651)
- 20221125.export.CSV.zip (4.7MB) (MD5: 288bdf1637d744436edb5372b0ce0db)
- 20221124.export.CSV.zip (5.3MB) (MD5: 0a405d3f359d8d455d58143762ac52)
- 20221123.export.CSV.zip (6.7MB) (MD5: d425f6853f3a2c25e094a69520164b)
- 20221122.export.CSV.zip (6.8MB) (MD5: 91c812a921d0e93c33ce9b9905d1b9f8e)
- 20221121.export.CSV.zip (5.9MB) (MD5: b083b678a382da3f129276e7cc1545b)
- 20221120.export.CSV.zip (3.3MB) (MD5: 1115d7de868f1919e189e13b427bc3)

GDELT 2.0 Event Database

15 Minute Updates

The screenshot shows a list of URLs for 15-minute updates of the GDELT 2.0 Event Database. The URLs are long strings of characters starting with 'http://data.gdeltproject.org/gdeltv2/masterfilelist.txt' followed by various identifiers. The list includes:

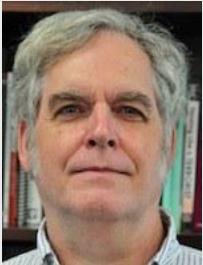
- 150383 297a16b493de7cf6ca809a7cc31d0b93 http://data.gdeltproject.org/gdeltv2/20150218230000.export.CSV.zip
- 318084 bb27f78ba45f69a17ea6ed775se9f8f http://data.gdeltproject.org/gdeltv2/20150218230000.mentions.CSV.zip
- 10768507 ea8dde0beb0ba98810a92d868c0ce99 http://data.gdeltproject.org/gdeltv2/20150218230000.gkg.csv.zip
- 149211 2a91841d7e72b0fc6a29e2ff867b240 http://data.gdeltproject.org/gdeltv2/20150218231500.export.CSV.zip
- 339037 dec3f427076b716a8112b9086c342523 http://data.gdeltproject.org/gdeltv2/20150218231500.mentions.CSV.zip
- 10269336 2f1a504a3c4558694ade0442e9a5ae6f http://data.gdeltproject.org/gdeltv2/20150218231500.gkg.csv.zip

From realtime translation of the world's news in 65 languages, to measurement of more than 2,300 emotions and themes from every article

CAMEO



Philip A. Schrodt (Project Director):



**Pennsylvania
State University**

2 VERB CODEBOOK

2.1	MAKE PUBLIC STATEMENT
2.2	APPEAL
2.3	EXPRESS INTENT TO COOPERATE
2.4	CONSULT
2.5	ENGAGE IN DIPLOMATIC COOPERATION
2.6	ENGAGE IN MATERIAL COOPERATION
2.7	PROVIDE AID
2.8	YIELD
2.9	INVESTIGATE
2.10	DEMAND
2.11	DISAPPROVE
2.12	REJECT
2.13	THREATEN
2.14	PROTEST
2.15	EXHIBIT MILITARY POSTURE
2.16	REDUCE RELATIONS
2.17	COERCE
2.18	ASSAULT
2.19	FIGHT
2.20	ENGAGE IN UNCONVENTIONAL MASS VIOLENCE

CAMEO

Conflict and Mediation Event Observations
Event and Actor Codebook

Event Data Project
Department of Political Science
Pennsylvania State University
Pond Laboratory
University Park, PA 16802

Version:
1.1b3 :
March
2012

CAMEO	073
Name	Provide humanitarian aid
Description	Extend, provide humanitarian aid, mainly in the form of emergency assistance.
Usage Notes	This code refers to events such as provisions of shelter, food, medicine, and evacuation of victims. The lead must report the delivery of such aid; promises to provide aid should be coded under category 033. Note that provisions of peacekeeping or other military forces are coded as 074 instead.
Example	Swiss doctors handed over 700 kg of medicine to the Red Crescent in Bam, Iran, according to the Swiss Agency for Development and Cooperation.
Example	Benin opened its borders today to most West Africans ordered out of Nigeria as illegal aliens, but was still refusing admittance to Ghanaians, by far the biggest group involved, Benin police said.
Example	U.N. helicopters evacuated the wounded from the besieged Bosnian town of Gorazde on Friday.

<http://data.gdelproject.org/documentation/CAMEO.Manual.1.1b3.pdf>

Fusion of NLP and Prediction



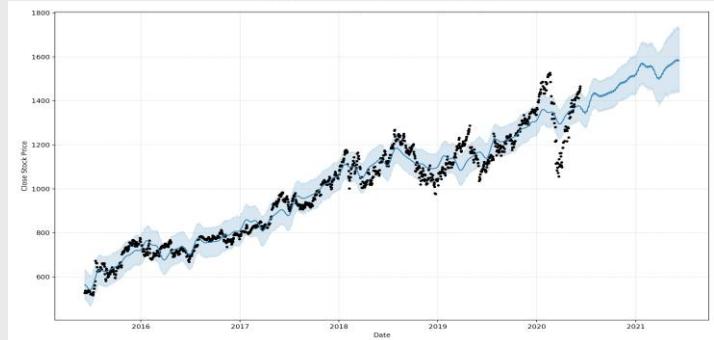
Textual information



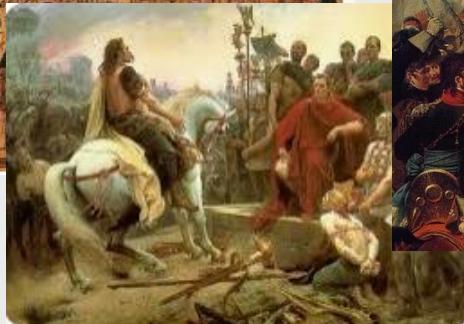
TimeSeries



Prediction



History as sequences of events



Sequence of events = phrases of a language

A people



A region



A culture



Decategorification



58 CAMEO Columns worldwide X 300 Millions rows

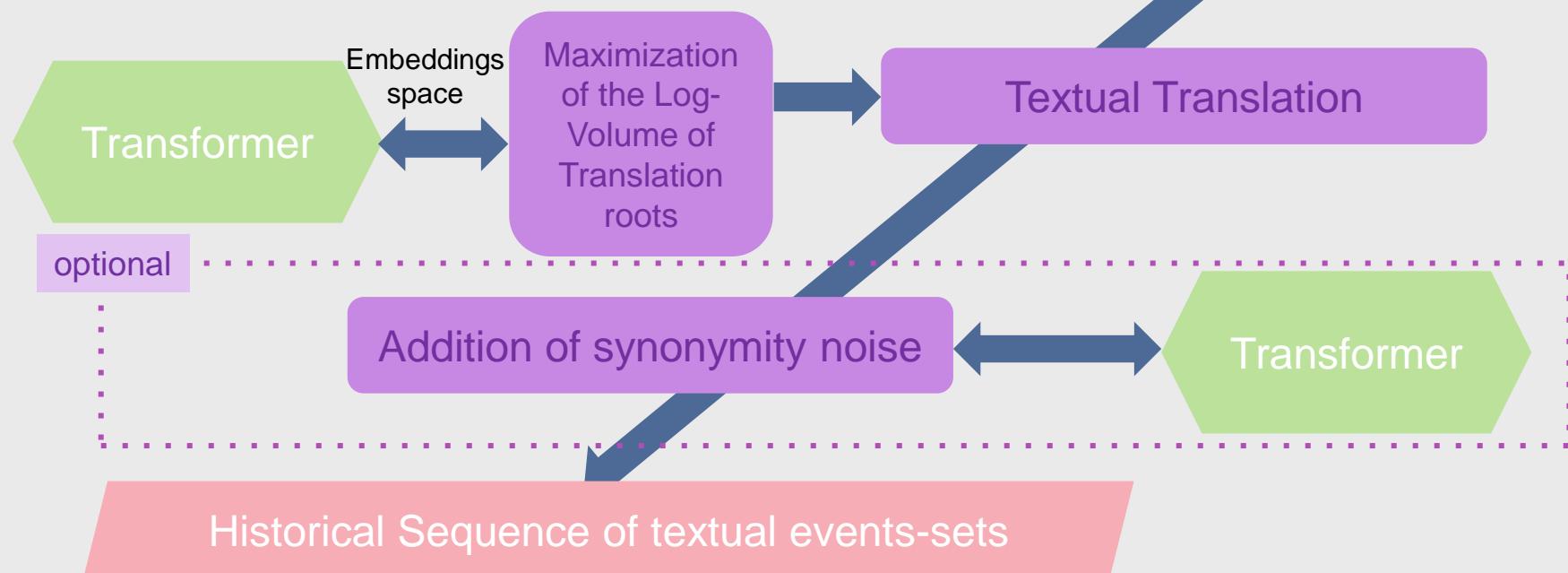
Country_linked 14 columns dataset

```
[ 'GLOBALEVENTID', 'SQLDATE', 'MonthYear', 'Year', 'FractionDate', 'Actor1Code',
'Actor1Name', 'Actor1CountryCode', 'Actor1KnownGroupCode', 'Actor1EthnicCode',
'Actor1Religion1Code', 'Actor1Religion2Code', 'Actor1Type1Code',
'Actor1Type2Code', 'Actor1Type3Code', 'Actor2Code', 'Actor2Name',
'Actor2CountryCode', 'Actor2KnownGroupCode', 'Actor2EthnicCode',
'Actor2Religion1Code', 'Actor2Religion2Code', 'Actor2Type1Code',
'Actor2Type2Code', 'Actor2Type3Code', 'IsRootEvent', 'EventCode', 'EventBaseCode',
'EventRootCode', 'QuadClass', 'GoldsteinScale', 'NumMentions', 'NumSources',
'NumArticles', 'AvgTone', 'Actor1Geo_Type', 'Actor1Geo_FullName',
'Actor1Geo_CountryCode', 'Actor1Geo ADM1Code', 'Actor1Geo_Lat', 'Actor1Geo_Long',
'Actor1Geo_FeatureID', 'Actor2Geo_Type', 'Actor2Geo_FullName',
'Actor2Geo_CountryCode', 'Actor2Geo ADM1Code', 'Actor2Geo_Lat', 'Actor2Geo_Long',
'Actor2Geo_FeatureID', 'ActionGeo_Type', 'ActionGeo_FullName',
>ActionGeo_CountryCode', 'ActionGeo ADM1Code', 'ActionGeo_Lat', 'ActionGeo_Long',
'ActionGeo_FeatureID', 'DATEADDED' ]
```

'ActionGeo_CountryCode'

```
[ 'GLOBALEVENTID', 'SQLDATE', 'MonthYear', 'Year', 'FractionDate',
'IsRootEvent', 'EventCode', 'EventBaseCode',
'EventRootCode', 'QuadClass', 'GoldsteinScale', 'Actor1Geo_CountryCode',
'Actor2Geo_CountryCode' ]
```

+ Multi_Country_Flag



Building of the dataset



Historical Sequence of textual events-sets

Keeping :

- Ponctual Statistical Structure
- Set Statistical Structure

Source : 100-10000 pts

Distribution of points
on the unitary hyper-sphere
In the embedding space

Random Generation

Target: 3 – 10 pts

Distribution of points
on the unitary hyper-sphere
In the embedding space

3-10 vectors embedding space

Small text : 3-10 simple phrases

Decoder -
Transformer

Daily descriptor

Learning



<https://shivanandroy.com/fine-tune-t5-transformer-with-pytorch/>

https://www.youtube.com/watch?v=U6x8_BP69DM

<https://lilianweng.github.io/posts/2021-01-02-controllable-text-generation/>

Dataset

10 000 Sequences of 50 daily descriptors

Metric

Rouge =

$$\text{Rouge} = \frac{\sum_r \sum_s \text{match}(\text{gram}_{s,c})}{\sum_r \sum_s \text{count}(\text{gram}_s)}$$

Event Prediction Principle



- 1) Generation of 100 possible future event-set
- 2) Projection of the cloud of points on the basis of standard event to determine probabilities.
- 3) Reconstruction of the standard events from the projections most likely arguments for the detailed actors and countries

Isaac Asimov : 3 laws of robotics (+1)



First Law

A robot may not injure a human being or, through inaction, allow a human being to come to harm.

Second Law

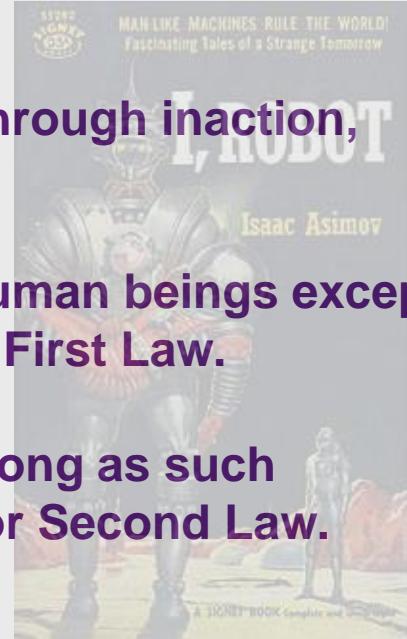
A robot must obey the orders given it by human beings except where such orders would conflict with the First Law.

Third Law

A robot must protect its own existence as long as such protection does not conflict with the First or Second Law.

Zeroth Law

A robot may not harm humanity, or, by inaction, allow humanity to come to harm.



Isaac Asimov : Psychohistory



Psychohistory is a fictional science in Isaac Asimov's *Foundation* universe which combines history, sociology, and mathematical statistics to make general predictions about the future behavior of very large groups of people, such as the Galactic Empire. It was first introduced in the four short stories (1942–1944) which would later be collected as the 1951 novel *Foundation*.

Hari Seldon, the hero in Asimov's romans established two axioms:

- that the population whose behavior was modeled should be sufficiently large
- that the population should remain in ignorance of the results of the application of psychohistorical analyses because if it is aware, the group changes its behaviour.

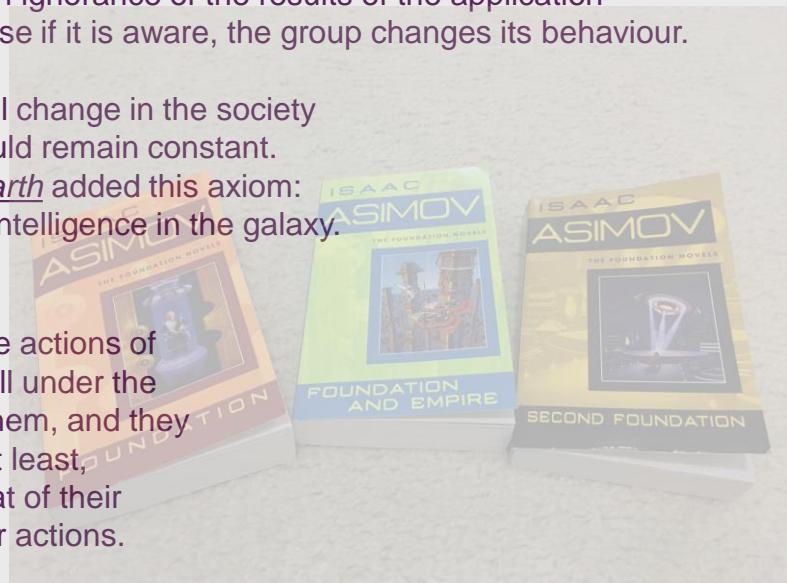
Ebling Mis added these axioms:

- that there would be no fundamental change in the society
- that human reactions to stimuli would remain constant.

Golan Trevize in *Foundation and Earth* added this axiom:

- that humans are the only sentient intelligence in the galaxy.

Hari Seldon developed psychohistory to predict the actions of large groups of *humans*. Even robots technically fall under the umbrella of psychohistory, because humans built them, and they thus represent more or less a human "action", or at least, possess a thought-framework similar enough to that of their human creators that psychohistory can predict their actions.

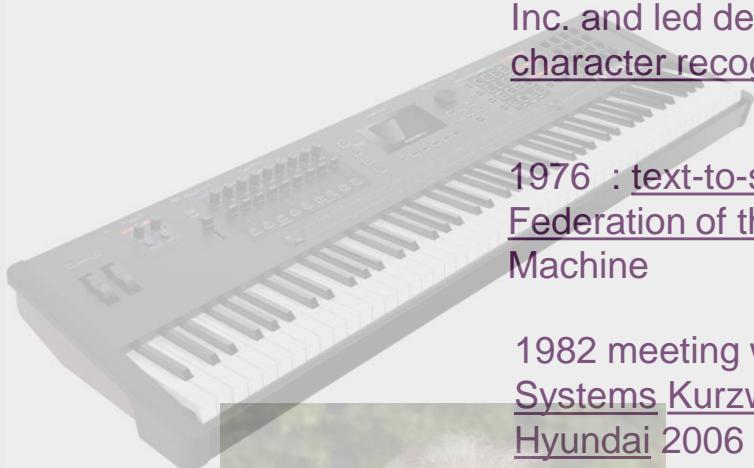


(from wikipedia)

Ray Kurzweil: Genius Inventor



1974: Kurzweil founded Kurzweil Computer Products, Inc. and led development of the first omni-font optical character recognition



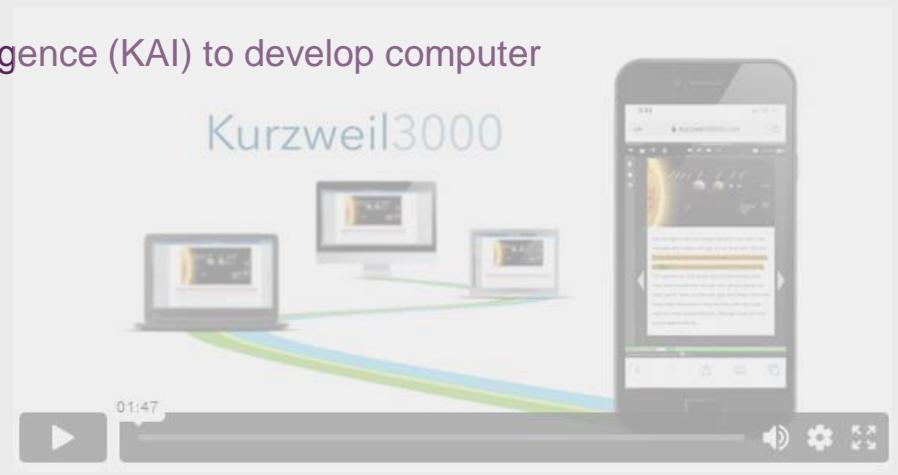
1976 : text-to-speech synthesis : leader of the National Federation of the Blind. Built the Kurzweil Reading Machine



1982 meeting with Stevie Wonder : Kurzweil Music Systems Kurzweil K250 -> Young Chang in 1990 -> . Hyundai 2006



Kurzweil Applied Intelligence (KAI) to develop computer speech recognition



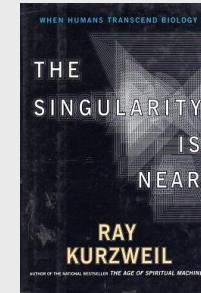
Ray Kurzweil: The Singularity



1990 : , *The Age of Intelligent Machines*

1999 : *The Age of Spiritual Machines*

2005 : *The Singularity Is Near*



2009: Kurzweil, in collaboration with Google and the NASA Ames Research Center, announced the creation of the Singularity University

2012 : *How to Create a Mind*

Pattern Recognition Theory of Mind, the theory that the neocortex is a hierarchical system of pattern recognizers, and argues that emulating this architecture in machines could lead to an artificial superintelligence

2009 : Once the Singularity has been reached, Kurzweil says that machine intelligence will be infinitely more powerful than all human intelligence combined. Afterwards he predicts intelligence will radiate outward from the planet until it saturates the universe. The Singularity is also the point at which machines' intelligence and humans would merge. Kurzweil spells out the date very clearly: "I set the date for the Singularity—representing a profound and disruptive transformation in human capability—as 2045"





5) GRAM MATRIX AND GRADIENT FLOWS: AN OVERVIEW

The Miracle of Transformers

To Summary

**Gram Matrix == Density Operator
== Neural Tangent Kernel
== Normalized Activation Batch Product
== Generator Information Gradient Flow**



Quantum Statistical physics

$$\hat{\rho} = \sum_k \frac{e^{-\beta E_k^M}}{Z} |\Psi_k\rangle\langle\Psi_k|$$

Reny's Information

$$S_\alpha(A_1, A_2, \dots, A_k) = S_\alpha \left(\frac{A_1 \circ A_2 \circ \dots \circ A_k}{\text{tr}(A_1 \circ A_2 \circ \dots \circ A_k)} \right)$$

$$G(x_1, \dots, x_n) = \begin{pmatrix} \langle x_1, x_1 \rangle & \langle x_1, x_2 \rangle & \dots & \langle x_1, x_n \rangle \\ \langle x_2, x_1 \rangle & \langle x_2, x_2 \rangle & \dots & \langle x_2, x_n \rangle \\ \vdots & \vdots & \ddots & \vdots \\ \langle x_n, x_1 \rangle & \langle x_n, x_2 \rangle & \dots & \langle x_n, x_n \rangle \end{pmatrix}$$

RKHS

$$k(x, y) = \langle k_x, k_y \rangle = \langle \phi(x), \phi(y) \rangle$$

Neural Tangent Kernel

$$K_{\theta_t}^L(x, x') = \nabla_{\theta} f(x, \theta_t) \nabla_{\theta} f(x', \theta_t)^T$$

Gradient Flows

$$\begin{aligned} \frac{\partial \mu_t}{\partial t} - \text{div} \left(\mu_t \nabla \frac{\partial \mathcal{G}(\mu_t)}{\partial \mu_t} \right) &= 0, \\ \partial_t \theta_p(t) &= \left\langle \partial_{\theta_p} F^{(L)}(\theta(t)), d_t \right\rangle_{p^{in}} \end{aligned}$$





Gradient Flow

Statistical Physics

Entropy, Free Energy, Focker-Planck,
Langevin, Boltzmann, Gibbs

Information Theory

Variational Inference,
Regularization,
Mutual Information,
Information Bottleneck

Differential Geometry

Riemannian manifolds,
Geodesics,
Connection,Curvature

$$\partial_t \theta_p(t) = \left\langle \partial_{\theta_p} F^{(L)}(\theta(t)), d_t \right\rangle_{p^{in}}$$

Stochastic Gradient Descent (SGD)

Diffusion

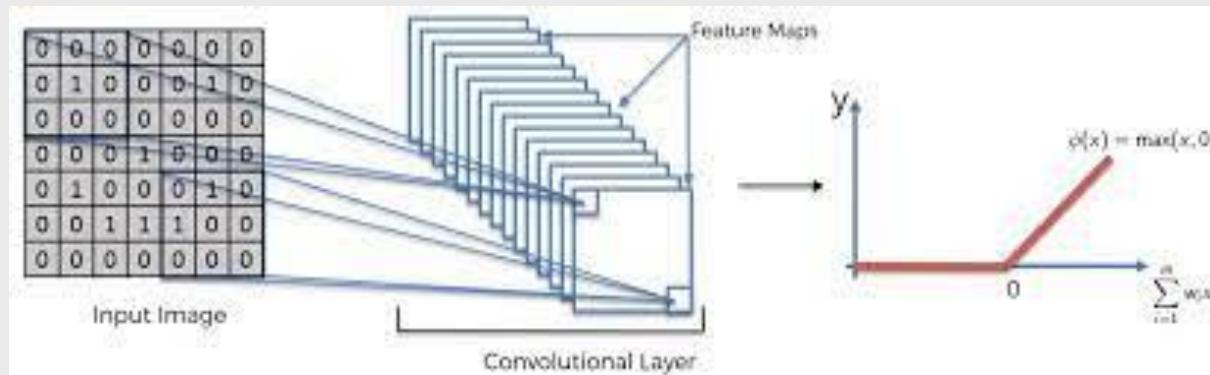
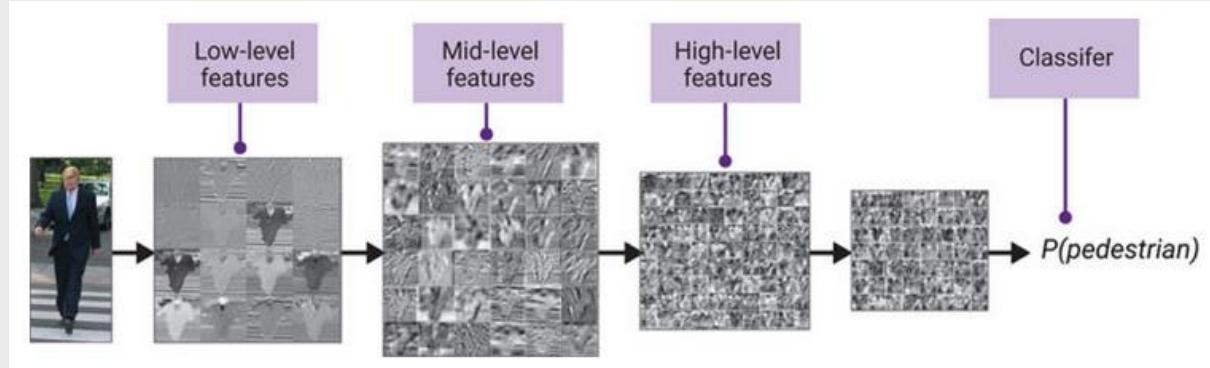
Domain adaptation, Clustering, Sampling,
Metropolis-adjusted Langevin algorithm,
Hamiltonian Monte Carlo

Optimal Transport

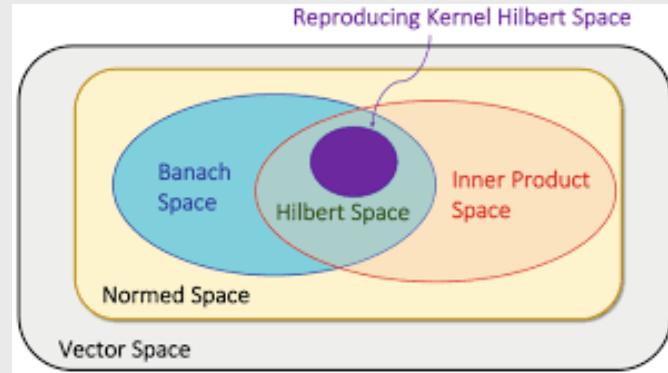
Wasserstein space ,
Sinkhorn algorithm



Feature Map



RKHS



Riesz representation theorem

Let \mathcal{H} be a Hilbert space over \mathbb{R} . If $T \in \mathcal{H}^*$, then there exists a unique vector u in \mathcal{H} such that

$$T(v) = \langle v, u \rangle_{\mathcal{H}} \text{ for all } v \in \mathcal{H}$$

Inverting the Riesz Theorem

We say \mathcal{H} is a *Reproducing Kernel Hilbert Space* if there exists a $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, such that

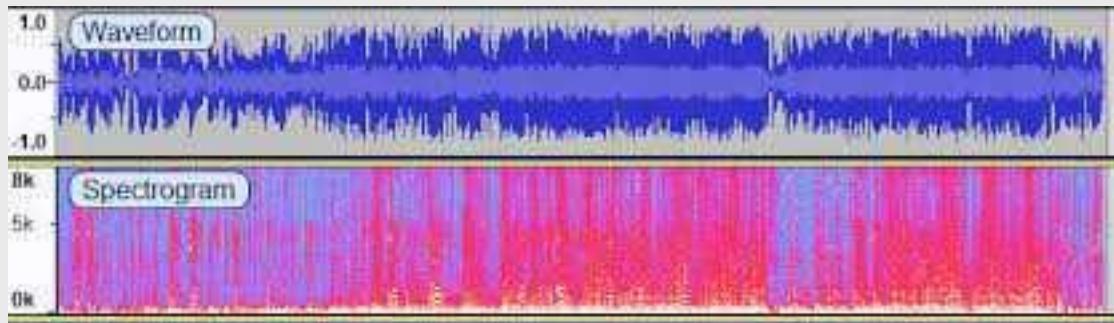
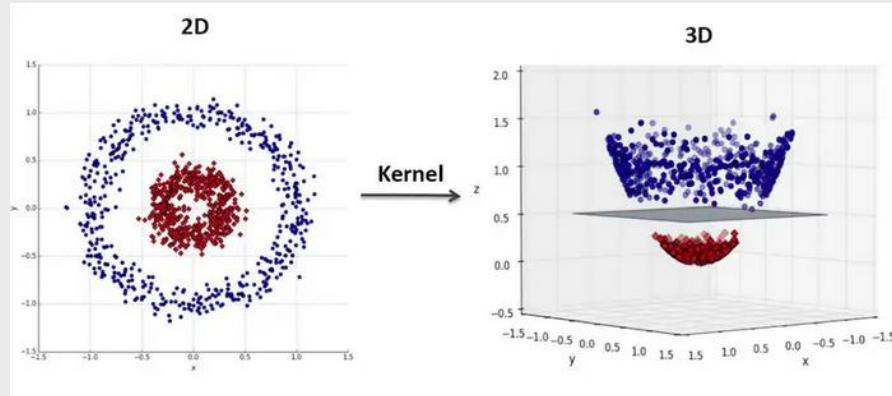
1. k has the reproducing property, i.e., $f(x) = \langle f(\cdot), k(\cdot, x) \rangle_{\mathcal{H}}$
2. k spans \mathcal{H} , that is, $\mathcal{H} = \overline{\text{span}\{k(\cdot, x) : x \in \mathcal{X}\}}$

$x \rightarrow k(y, x)$ is called a Feature Map



[To Summary](#)

Example of Feature Map



Gram Matrix To RKHS :

A smart way to extend features and introduce natural codependance



Mercer Kernel $K(x, y)$:

Function of two variables that is symmetric and positive definite

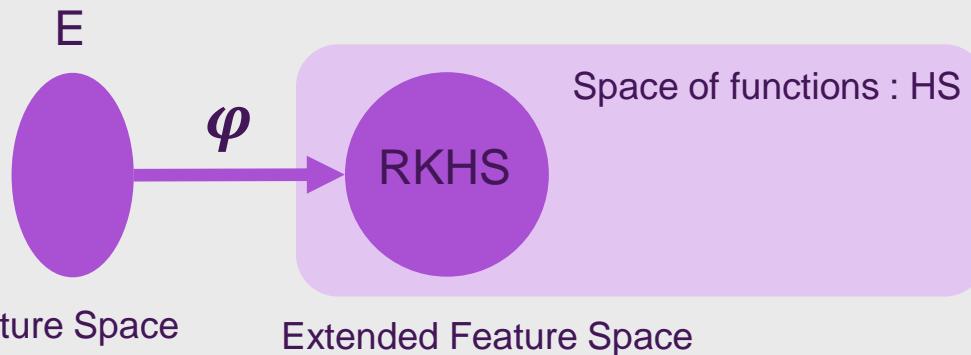
$$\int \int K(x, y) f(x) f(y) dx dy \geq 0 \iff \langle f, f \rangle \geq 0$$

Riez theorem

Every linear **continuous** operator is implemented as a scalar product.
The reciprocal is false (think about valuation at x)

RKHS

A closed subspace of the hilbert space of functions on E , where the reciprocal of Riez theorem is true



Reproducing
Property

$$\langle \varphi(x), \varphi(y) \rangle = \langle x, y \rangle$$

Neural Tangent kernel



$$f(W, X) \approx f(W_0, X) + \nabla_W f(W_0, X)^T (W - W_0) + \dots$$

$$\varphi(X) = \nabla_W f(W_0, X)$$

$$K(X_i, X_j) = \langle \varphi(X_i), \varphi(X_j) \rangle \quad \text{Neural Tangent Kernel}$$

NTK can be viewed as a specific type of Gram matrix that captures the local geometry of the parameter space during training



Renyi alpha entropy



$$H_\alpha(f) = \frac{1}{1-\alpha} \log \int_{\mathcal{X}} f^\alpha(x) dx$$



$$\alpha \rightarrow 1 \ H_1(X) = - \sum_x p(x) \log p(x)$$

$$\alpha \rightarrow \infty \ H_\infty(X) = - \log \max_x p(x)$$

Shannon Entropy

Min-Entropy

Shannon Mutual Information

$$I_1(X;Y) \stackrel{\text{def}}{=} H_1(X) - H_1(X|Y) = H_1(X) + H_1(Y) - H_1(X,Y)$$



Feature selection in machine learning: Renyi min-entropy vs Shannon entropy
arXiv:2001.09654v1 [cs.LG] 27 Jan 2020

Why Reny's Information



Sensitivity to rare or frequent events

$0 < \alpha < 1 \Rightarrow$ it is more sensitive to rare events,
 $\alpha > 1 \Rightarrow$ it is more sensitive to frequent events.

Information bottleneck method

Better trade-off between compression and relevance, making it easier to find the optimal solution

Time series data in finance

In this context, the adjustable parameter α allows for different perspectives on the complexity of the patterns.

Non-asymptotic or non-ergodic settings

When the law of large numbers does not readily apply, other entropy measures typically take over, for example the min-, the max-, or the collision entropy. The Renyi entropies nicely unify these different and isolated measures.



On quantum Renyi entropies: a new generalization and some properties

Martin Muller-Lennert, Frederic Dupuis, Oleg Szehr, Serge Fehr, and Marco Tomamichel

The world according to Renyi: Thermodynamics of multifractal systems Petr Jizba and Toshihico Arimitsu

[To Summary](#)



12/29/2023

27

6) KL DIVERGENCE

The Miracle of Transformers

[To Summary](#)

Kullback-Leibler Divergence



$$D_{KL}[p||q] = \int_{-\infty}^{\infty} p(x) \log \frac{p(x)}{q(x)} dx$$

$$D_{KL}[p||p] \geq 0 \quad \text{et} \quad D_{KL}(p||p) = 0$$

But it is not a distance but triangular inéquality does not work

Is not symmetric

let $q(x)$ a prob. distribution defined on E et $F \subset E$

Let $q_F(x)$ a prob. distribution restricted to F (conditionned by F)

$$\text{Then } D_{KL}[q||q_F] = \log \frac{1}{\text{Prob}[F]}$$

This is the supplementary information: conditionned by F (in bits)

This is the Renyi divergence for $\alpha = 1$, $D_\alpha[p||q] = \frac{1}{1-\alpha} \int_{-\infty}^{\infty} p(x)^\alpha q(x)^{1-\alpha} dx$

The topology generated by D_{KL} and D_α are not the same : $D_{1/2} \rightarrow \text{Distance}$

$$\lim_{\theta_1 \rightarrow \theta} \frac{1}{(\theta_1 - \theta)^2} D_\alpha[p_{\theta_1}||p_\theta] = \frac{\alpha}{2} J(\theta) \quad \text{Fisher Information}$$

Hellinger
↓



Choosing
The right tool



Claude Shannon

Measure a spread

$$D_{KL}[p||q] = H(p) - H(p, q)$$

Entropy

Crossed Entropy

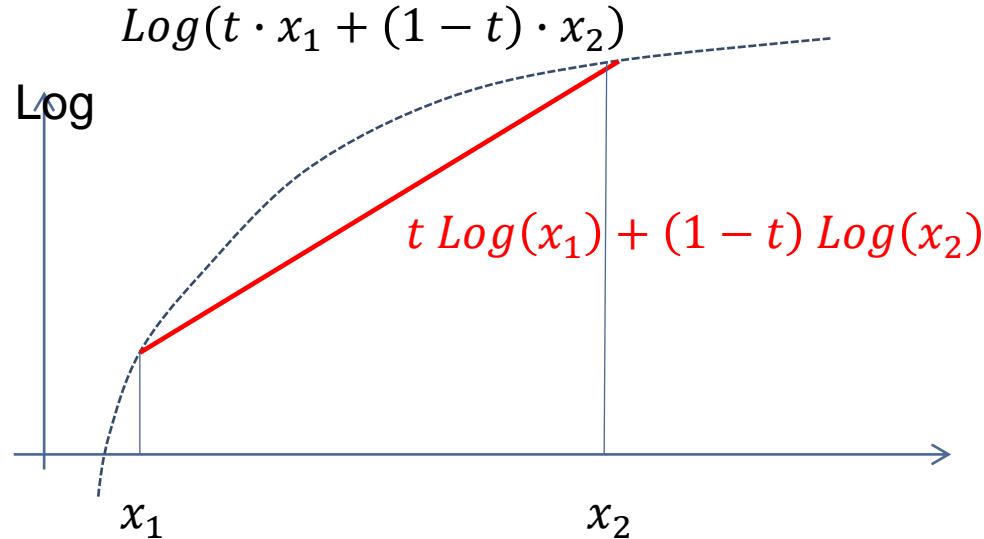
- An Efficient criteria
- Classification problems
- Far Gradient still active

[To Summary](#)

Jensen Inequality



2 points



N points >2

Jensen :

$$f \text{ concave} \Rightarrow f(\mathbb{E}[X]) \geq \mathbb{E}[f(X)]$$

Mutual Information Properties



Definition

$$I(X; Y) = D_{KL}[p(x, y) || p(x)p(y)]$$

The information improvement compared to an independant case

Data processing Inequality (DPI)

For any markov chains $X \rightarrow Y \rightarrow Z$

$$I(X; Y) \geq I(X; Z)$$

There is less information on X in Z than in Y

Invariance by reparametrisation

$$I(X; Y) = I(\varphi(X); \psi(Y)) \quad \text{For any invertible } \varphi \text{ and } \psi$$

[To Summary](#)



12/29/2023

32

7) BENAMOU BRENIER THEORY

The Miracle of Transformers

Benamou Brenier Theory 1



Optimal Transport - Benamou-Brenier Formulation (Brittany Hamfeldt)

<https://www.youtube.com/watch?v=xxwalfOJ4sg&list=PLJ6garKOIK2qKVhRm6UwvcQ46wK-ciHbl&index=24>

$$\rho(0, x) = f(x) \quad \rho(1, x) = g(x)$$

$$\text{Continuity equation} \quad \rho_t + \nabla(V\rho) = 0$$

Lagrangian Coordinates : $X(t, x)$:

$$\frac{\partial X}{\partial t}(t, x) = V(t, X(t, x)) \quad X(0, x) = x$$

$$\text{Push forward : } X(t, \cdot) \# f = \rho(t, \cdot) \quad X(1, \cdot) \# f = g$$

$W_2^2(f, g)$ Verifies

$$\begin{aligned}
 W_2^2(f, g) &\leq \int_{\mathbb{R}^n} f(x) |X(t, x) - x|^2 dx \\
 &= \int_{\mathbb{R}^n} f(x) |X(t, x) - X(0, x)|^2 dx \\
 &\leq \int_{\mathbb{R}^n} f(x) \left| \int_0^1 \frac{\partial X}{\partial t}(t, x) dt \right|^2 dx \\
 &\leq \int_{\mathbb{R}^n} \int_0^1 f(x) \left| \frac{\partial X}{\partial t}(t, x) \right|^2 dt dx \\
 &= \int_0^1 \int_{\mathbb{R}^n} f(x) |\nabla(t, X(t, x))|^2 dx dt
 \end{aligned}$$



Benamou Brenier Theory 2



• Need $X(1, x) = T(x)$, the optimal map from f to g .
 • Need $\left| \int_0^1 \frac{\partial X}{\partial t}(t, x) dt \right|^2 = \int_0^1 \left| \frac{\partial X}{\partial t}(t, x) \right|^2 dt$
 This holds if $\frac{\partial X}{\partial t}(t, x) = u(x)$
 is constant in time
 • Recall: $X(0, x) = x$
 $\Rightarrow X(t, x) = x + u(x)t$
 $X(1, x) = x + u(x) = T(x)$
 $\Rightarrow u(x) = T(x) - x$
 $X(t, x) = x + (T(x) - x)t = x(1-t) + tT(x)$

Our velocity field should satisfy:
 $v(t, X(t, x)) = u(x) = T(x) - x = (T - I)u(x)$
 Let $y = X(t, x) = [(1-t)I + tT]u(x)$
 $\Rightarrow v(t, y) = (T - I)[(1-t)I + tT]^{-1}u(x)$
 achieves the optimum.
 (if f & g are nice)
 Q.W. we do "smoothed" approximations to
 get arbitrarily close to infimum.

Benamou Brenier formulation

$$W_2^2(f, g) = \inf_{\mathcal{C}, V} \int_0^1 \int_{\mathbb{R}^n} \mathcal{C}(t, x) |v(t, x)|^2 dx dt$$

s.t.

$$\begin{cases} \mathcal{C}_t + \nabla \cdot (\mathcal{C}v) = 0 \\ \mathcal{C}(0, x) = f(x) \\ \mathcal{C}(1, x) = g(x) \end{cases}$$

Benamou Brenier Theory 3



Eg: $\min f(x)$
S.L. $Ax = b$

Could write down Lagrangian
 $L(x, \lambda) = f(x) + \lambda \cdot (Ax - b)$

Seek saddle pts of this
 $\min_x \left\{ f(x) + \sup_{\lambda} \lambda \cdot (Ax - b) \right\}$

We try to write down a Lagrangian for our flow problem.

Let $\varphi(t, x)$ be the Lagrange multiplier.

The continuity eqn gives us a term:
 $\int_0^1 \int_{\mathbb{R}^n} (\ell_t + \nabla(\varphi v)) \cdot \varphi dt dx$
 $= \int_{\mathbb{R}^n} (g(x) \varphi(x, 0) - f(x) \varphi(x, 1)) dx$
 $- \int_0^1 \int_{\mathbb{R}^n} (\rho \varphi_t + \rho v \cdot \nabla \varphi) dx dt$

Lagrangian:
 $\mathcal{L} = \int_0^1 \int_{\mathbb{R}^n} \left(\frac{\rho |v|^2}{2} - \rho \varphi_t - \rho v \cdot \nabla \varphi \right) dx dt$
 $+ \int_{\mathbb{R}^n} (g(x) \varphi(x, 0) - f(x) \varphi(x, 1)) dx$
 $\quad \quad \quad G(\varphi)$

Let $h(c, m) = \frac{|m|^2}{2c}$.

Claim: $h(c, m) = \sup_{(a, b) \in K} (a\rho + bm)$
where $K = \{(a, b) \in \mathbb{R}^2 \times \mathbb{R}^n \mid a + \frac{|b|^2}{2} \leq c\}$

Proof
Let (a^*, b^*) achieve the supremum.
Since $c > 0$ we want a^* to be as big as possible
 $\Rightarrow a^* = -\frac{|b^*|^2}{2}$
 $\sup_{(a, b) \in K} (a\rho + bm) = \sup_{b \in \mathbb{R}^n} \left(-\frac{\rho |b|^2}{2} + bm \right)$

This is concave \Rightarrow set gradient to 0:
 $-\rho b + m = 0$
 $\Rightarrow b^* = m/\rho$
 $\sup_{(a, b) \in K} (a\rho + bm) = -\frac{|m|^2}{2\rho} + \frac{|m|^2}{\rho}$
 $= \frac{|m|^2}{\rho}$
 $= h(c, m)$

Saddle point problem is:
 $\inf_{\rho} \sup_{\varphi, \beta \in K} \int_0^1 \int_{\mathbb{R}^n} (a\rho + bm - \langle \rho v - m, \nabla \varphi \rangle) dt dx + G(\varphi)$
 $= \inf_{\rho} \sup_{\varphi, \beta \in K} \int_0^1 \int_{\mathbb{R}^n} \left[(a - \rho v) \rho + \left(b - \frac{\nabla \varphi}{\rho} \right) m \right] dt dx + G(\varphi)$

Let $r = \begin{pmatrix} \rho \\ m \end{pmatrix}$, $\beta = \begin{pmatrix} a \\ b \end{pmatrix}$, $\nabla_{\varphi} \varphi = \begin{pmatrix} \rho_v \\ \nabla \varphi \end{pmatrix}$

$\Rightarrow \inf_{\rho} \sup_{\varphi, \beta \in K} \langle \beta - \nabla_{\varphi} \varphi, r \rangle + G(\varphi)$

let's let $\tau > 0$ be small to regularize the max.

Optimise these 3 unknowns independently:
Given r_k, β_k , find optimal φ_k
 $\varphi_{k+1} = \operatorname{argmax}_{\varphi} - \langle \nabla_{\varphi} \varphi, r \rangle + G(\varphi) - \frac{\tau}{2} \|\beta - \nabla_{\varphi} \varphi\|^2$
Quadratic, unconstrained
 \Rightarrow gradient = 0
 \Rightarrow linear problem

Computing 1st variation
 \Rightarrow Poisson equation
with Neumann bc
(homogeneous in space & non-homogeneous in time)

$\Delta_{\text{kin}} \varphi_{k+1} = \nabla \cdot \left(\beta_k - \frac{r_k}{\tau} \right)$

- Given r_k, φ_{k+1} , optimise for β_{k+1}
 $\beta_{k+1} = \operatorname{argmin}_{\beta \in K} \langle \beta, r \rangle - \frac{\tau}{2} \|\beta - \nabla_{\varphi} \varphi\|^2$
Do this pointwise
(Euler-Lagrange optimisation)
- Given φ_k, β_{k+1}
Do gradient descent in r
 $r_{k+1} = r_k - \tau (\beta_{k+1} - \nabla_{\varphi} \varphi_k)$
- Iterate

8) GRADIENT FLOWS: THE DEMONSTRATIONS



The Miracle of Transformers

Gradient Flows 1



<https://www.youtube.com/watch?v=zzGBxAqJV0Q&list=PLJ6garKOIK2qKVhRm6UwvcQ46wK-ciHbl&index=25>

Discretise in time via Backward Euler

$$\frac{x^{n+1} - x^n}{\tau} = -\nabla F(x^{n+1})$$

$$\Rightarrow \frac{x^{n+1} - x^n}{\tau} + \nabla F(x^{n+1}) = 0$$

$$\Rightarrow \nabla \left(\frac{|x - x^n|^2}{2\tau} + F(x) \right) \Big|_{x=x^{n+1}} = 0$$

$$\Rightarrow x^{n+1} \in \operatorname{argmin}_{x} \left\{ \frac{|x - x^n|^2}{2\tau} + F(x) \right\}$$

Can define a scheme like this in a metric space

$$(X, d)$$

Let $F: X \rightarrow \mathbb{R}$ be l.s.c. and bounded below.

Define

$$x_{\tau}^{n+1} \in \operatorname{argmin}_{x} \left\{ F(x) + \frac{d(x, x^n)^2}{2\tau} \right\}$$

Interpolate to all t

$$\text{e.g.: } x_{\tau}(t) = x_{\tau}^n \text{ if } t \in ((n-1)\tau, n\tau]$$

Study limit as $\tau \rightarrow 0$.

Consider

$$F: P(\mathbb{R}) \rightarrow \mathbb{R}, \quad d = W_2$$

\mathbb{R} is compact

F is l.s.c. and bounded below.

We previously used the continuity equation

$$(\rho + \nabla \cdot (\rho v)) = 0$$

to "flow" densities.

Goal: Find velocity field v s.t. this flow agrees with $\lim_{t \rightarrow 0} x_{\tau}(t)$

Use the dual formulation:

$$\begin{aligned} W_2^2(S, Q) &= 2 \inf_{u \in \operatorname{Lip}(\mathbb{R})} \int \frac{\|u(y)\|^2}{2} dy \\ &= 2 \max_{u, v} \left\{ \int u \delta x + \int v \delta y \mid \int u \delta y = \int v \delta x \right\} \\ &= 2 \max_{u, v} \left\{ \int u \delta x + \int u \delta y \right\} \\ &\stackrel{d}{=} W_2^2(F + \varepsilon \chi, Q) \Big|_{\varepsilon=0} \\ &= 2 \inf_{\varepsilon \geq 0} \max_{u, v} \left\{ \int u((f + \varepsilon \chi) \delta x + v \delta y) \right\} \\ &= 2 \int u \chi dx \end{aligned}$$

(investigate optimality condition in the JKO scheme.)

We need to compute the 1st variation.

We need to perturb (ρ, F) by ε

$$\text{to } (\rho + \varepsilon \chi, F + \varepsilon \chi)$$

Need $\rho + \varepsilon \chi \in P(\mathbb{R})$

so that $F(\rho + \varepsilon \chi)$ is

well-defined.

Restrict to χ s.t. $\Omega = \rho + \varepsilon \chi \in P(\mathbb{R})$

\forall small $\varepsilon > 0$.

In particular, $\Omega = \rho + \chi \in P(\mathbb{R})$

$$(\rho + \varepsilon \chi) = (\rho + \varepsilon(\sigma - \varepsilon))$$

$$= [\rho(1 - \varepsilon) + \varepsilon \sigma]$$

$$\in P(\mathbb{R})$$

as long as $\rho, \sigma \in P(\mathbb{R})$

$$\forall \sigma \in P(\mathbb{R}) \cap L_c^\infty(\mathbb{R})$$

The first variation of F , $\frac{\delta F}{\delta \rho}(\rho)$ is such that

$$\left. \frac{d}{d\varepsilon} F(\rho + \varepsilon \chi) \right|_{\varepsilon=0} = \int \frac{\delta F}{\delta \rho}(\rho) \chi \text{d}x \quad \forall \chi = \sigma - \rho \in P(\mathbb{R}) \cap L_c^\infty(\mathbb{R})$$

Where u^* achieves the max in the previous line.

i.e. u^* is the potential associated with the cost $\frac{1}{2} \|x - y\|^2$.

When we do OT, the optimal map

$$T(x) = x - \nabla u^*(x)$$

$$= x - (\nabla h)^*(\nabla u^*)$$

$$\text{where } h(x) = \frac{1}{2} \|x\|^2$$

$$\Rightarrow \frac{\delta W_2^2}{\delta \rho}(\rho, \rho^*) = 2u^*$$

$T(x) = x - \nabla u^*(x)$ is the optimal map from ρ for

the JKO scheme.

$$\left. \rho^* = \operatorname{argmin}_{\rho^*} \left\{ F(\rho) + \frac{W_2^2(\rho, \rho^*)}{2\tau} \right\} \right.$$

$$= \operatorname{argmin}_{\rho^*} G(\rho)$$

$$\Rightarrow \frac{\delta G}{\delta \rho}(\rho^*) + C = 0$$

$$\Rightarrow \frac{\delta F}{\delta \rho}(\rho^*) + \frac{u^*}{\tau} = \text{constant}$$

Gradient Flows 2



$\Rightarrow \dot{c} = \nabla \left(\frac{\delta F}{\delta c} \right) + \frac{v^*}{\gamma}$

$$= \nabla \left(\frac{\delta F}{\delta c} \right) + \frac{x - T(x)}{\gamma}$$

$$\Rightarrow \underbrace{\frac{T(x) - x}{\gamma}}_{\text{velocity!}} = \nabla \left(\frac{\delta F}{\delta c} \right)$$

The flow we want should have velocity
 $v(x) = -\frac{T(x) - x}{\gamma} = -\nabla \left(\frac{\delta F}{\delta c} \right)$

This is the velocity associated with our time-discrete scheme
 If everything works out as $T \rightarrow 0$, we expect our JKO scheme to limit to this flow

 $c_t + \nabla \cdot (c v) = 0$

or

 $c_t - \nabla \cdot (c \frac{\delta F}{\delta c}) = 0$

This is the PDE associated with gradient flows of F in the W_2 metric

Eg: $F(c) = \int c \log c dx$
 (negative entropy)
 We want a flow that minimizes entropy

$$\frac{d}{dx} F(c + \epsilon x) \Big|_{x=0}$$

$$= \frac{d}{dx} \int (c + \epsilon x) \log(c + \epsilon x) dx \Big|_{x=0}$$

$$= \int (\lambda \log(\lambda) + \lambda) dx$$

$$\Rightarrow \frac{\delta F}{\delta c} = \lambda \log \lambda + 1$$

$\nabla \left(\frac{\delta F}{\delta c} \right) = \nabla (\lambda \log \lambda + 1) = \frac{1}{\lambda} \nabla \lambda$

\Rightarrow The Gradient Flow is

$$\dot{c} = c_t - \nabla \cdot (c \frac{1}{\lambda} \nabla c)$$

$$= c_t - \nabla \cdot (\nabla c)$$

$$= c_t - \Delta c$$

$$\Rightarrow c_t = \Delta c \quad (\text{heat equation})$$

Eg: $F(c) = \int c \log c dx + \int V(x) c dx$

 $\Rightarrow c_t - \Delta c - \nabla \cdot (c \nabla V) = 0$

(Fokker-Planck)

Eg: $F(c) = \frac{1}{m} \int c^m dx$

 $\Rightarrow c_t - \Delta(c^m) = 0$

(porous medium)

Eg: $F(c) = \int c \log c - \frac{1}{2} \int |\nabla c|^2$
 when $-\Delta u = c$

$\Rightarrow \begin{cases} c_t + \nabla \cdot (c \nabla u) - \Delta c = 0 \\ -\Delta u = c \end{cases}$

(Keller-Segel, chemotaxis)

Eg: $F(c) = \frac{1}{2} \iint \omega(x-y) d(c(x)) d(c(y))$

 $\Rightarrow c_t - \nabla \cdot (c ((\partial \omega)^* c)) = 0$

(aggregation model)



[To Summary](#)



12/29/2023

39

9) GRADIENT FLOWS THEORY

The Miracle of Transformers

Diffusion Equation via Optimal Transport



Jordan –Kinderlehrer–Otto (98)

- the 2-Kantorovich metric on the space of probability measures

$$W_2(\mu, \nu) = \inf_{\gamma \in \text{Cpl}(\mu, \nu)} \sqrt{\int_{\mathbf{R}^n \times \mathbf{R}^n} |x - y|^2 d\gamma(x, y)}$$

- the (negative of the) Boltzmann-Shannon entropy

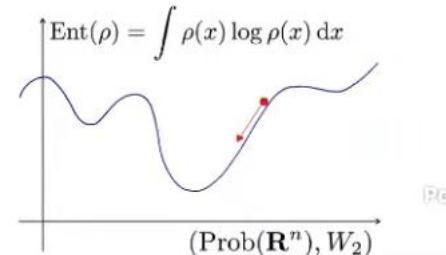
$$\text{Ent}(\mu) = \int_{\mathbf{R}^n} \rho(x) \log \rho(x) dx, \quad \text{if } d\mu(x) = \rho(x) dx$$

- the heat equation

$$\partial_t \mu = \Delta \mu$$

Theorem (J-K-O '98)

The heat flow is the gradient flow of the entropy w.r.t W_2 .



$$\partial_t \mu = \Delta \mu \iff \frac{1}{2} \frac{d}{dt} W_2(\mu_t, \nu)^2 \leq \text{Ent}(\nu) - \text{Ent}(\mu_t) \quad \forall \nu$$





Gradient Flow Definition

Gradient flows in \mathbf{R}^n

Let $\varphi : \mathbf{R}^n \rightarrow \mathbf{R}$ smooth and convex. For $u : \mathbf{R}_+ \rightarrow \mathbf{R}^n$ TFAE:

1. u solves the gradient flow equation $u'(t) = -\nabla\varphi(u(t))$.
2. u satisfies the evolution variational inequality

$$\frac{1}{2} \frac{d}{dt} |u(t) - y|^2 = (u(t) - y) \cdot u'(t) \leq \varphi(y) - \varphi(u(t)) \quad \forall y$$

(DE GIORGI '93, AMBROSIO–GIGLI–SAVARÉ '05)



Brenier's Theorem



Def (pushforward) : Let $\mu \in \mathcal{P}_2(\mathbb{R}^d)$, $T : \mathbb{R}^d \rightarrow \mathbb{R}^d$. The pushforward measure $T_{\#}\mu$ is characterized by:

- ▶ $\forall B \text{ meas. set}, T_{\#}\mu(B) = \mu(T^{-1}(B))$
- ▶ $x \sim \mu, T(x) \sim T_{\#}\mu$

Brenier's theorem : Let $\mu, \nu \in \mathcal{P}_2(\mathbb{R}^d)$ s.t. $\mu \ll Leb$. Then,

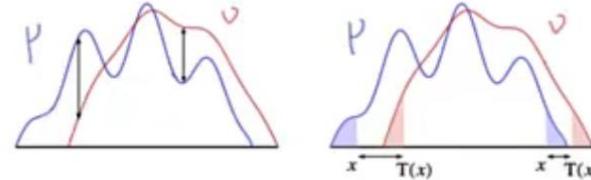
- ▶ Then $\exists! T_{\mu}^{\nu} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ s.t. $T_{\mu\#}^{\nu}\mu = \nu$, and a convex function g s.t. $T_{\mu}^{\nu} = \nabla g$ μ -a.e.
- ▶ $W_2^2(\mu, \nu) = \|I - T_{\mu}^{\nu}\|_{L_2(\mu)}^2 = \inf_{T \in L_2(\mu)} \int (x - T(x))^2 d\mu(x)$
where $L^2(\mu) = \{f : \mathbb{R}^d \rightarrow \mathbb{R}^d, \int \|f(x)\|^2 d\mu(x) < \infty\}$

W_2 geodesics?

$$\rho(0) = \mu, \rho(1) = \nu.$$

$$\rho(t) = ((1-t)I + tT_{\mu}^{\nu})_{\#}\mu$$

$$\neq \rho(t) = \underbrace{(1-t)\mu + t\nu}_{\text{mixture}}$$



Continuity Equation



Let $T > 0$. Consider a family $\mu : [0, T] \rightarrow \mathcal{P}_2(\mathbb{R}^d)$, $t \mapsto \mu_t$. It satisfies a **continuity equation** if there exists $(V_t)_{t \in [0, T]}$ such that $V_t \in L^2(\mu_t)$ and distributionnally:

$$\frac{\partial \mu_t}{\partial t} + \text{div}(\mu_t V_t) = 0$$

rules density μ_t of particles $x_t \in \mathbb{R}^d$ driven by a vector field V_t :

$$\frac{dx_t}{dt} = V_t(x_t)$$

Riemannian interpretation [Otto, 2001] :

The tangent space of $\mathcal{P}_2(\mathbb{R}^d)$ at μ_t verifies $\mathcal{T}_{\mu_t} \mathcal{P}_2(\mathbb{R}^d) \subset L^2(\mu_t)$.



Wasserstein Gradient Flow



Wasserstein gradient $\nabla_W F(\mu)$ of functional $F(\mu)$

$F(\mu)$	$\nabla_W F(\mu)$
$\int V d\mu$	$\nabla V(\cdot)$
$\text{KL}(\mu\ \nu)$	$\nabla \log \frac{d\mu}{d\nu}(\cdot)$
$\chi^2(\mu\ \nu)$	$2\nabla \frac{d\mu}{d\nu}(\cdot)$



Interpret as
 $\dot{X}_t = -\nabla_W F(\mu_t)(X_t)$

Gradient flow: $\dot{\mu}_t = -\nabla_W F(\mu_t)$

<https://www.youtube.com/watch?v=E0aK81IH6tM>



Wasserstein Gradient Flows(2)



Let $\mathcal{G} : \mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R} \cup \{+\infty\}$, $\mu \mapsto \mathcal{G}(\mu)$ a regular functional.

The differential of \mathcal{G} evaluated at $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ is the unique function $\frac{\partial \mathcal{G}(\mu)}{\partial \mu} : \mathbb{R}^d \rightarrow \mathbb{R}$ s. t. for any $\mu, \mu' \in \mathcal{P}_2(\mathbb{R}^d)$, s.t.
 $\mu' - \mu \in \mathcal{P}_2(\mathbb{R}^d)$:

$$\lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} [\mathcal{G}(\mu + \epsilon(\mu' - \mu)) - \mathcal{G}(\mu)] = \int_{\mathbb{R}^d} \frac{\partial \mathcal{G}(\mu)}{\partial \mu}(x) (d\mu' - d\mu)(x)$$

Then $\mu : [0, T] \rightarrow \mathcal{P}_2(\mathbb{R}^d)$, $t \mapsto \mu_t$ satisfies a **Wasserstein gradient flow** of \mathcal{G} if distributionally:

$$\frac{\partial \mu_t}{\partial t} - \operatorname{div} \left(\mu_t \nabla \frac{\partial \mathcal{G}(\mu_t)}{\partial \mu_t} \right) = 0, \text{ i.e. } V_t = -\nabla_W \mathcal{G}(\mu)$$

where $\nabla_W \mathcal{G}(\mu) := \nabla \frac{\partial \mathcal{G}(\mu)}{\partial \mu} \in L^2(\mu)$ is called the Wasserstein gradient of \mathcal{G} .



Wasserstein Gradient Flows of Free Energies



In particular, if the functional \mathcal{G} is a free energy:

$$\mathcal{G}(\mu) = \underbrace{\int H(\mu(x))dx}_{\text{internal energy } \mathcal{H}(\mu)_I} + \underbrace{\int V(x)d\mu(x)}_{\text{potential energy } \mathcal{E}_V(\mu)} + \underbrace{\int W(x, y)d\mu(x)d\mu(y)}_{\text{interaction energy } \mathcal{W}(\mu)}$$

$$\text{Then : } \frac{\partial \mu_t}{\partial t} = \operatorname{div}(\mu_t \underbrace{\nabla(H'(\mu_t) + V + W * \mu_t)}_{\nabla_W \mathcal{G}(\mu)}). \quad (1)$$

For instance, if $\mathcal{H}(\mu)$ is the negative entropy ($H(s) = s \log(s)$), then (1) rules the density μ_t of particles $x_t \in \mathbb{R}^d$ driven by :

$$\frac{dx_t}{dt} = -\nabla V(x_t) - \int_{\mathbb{R}^d} \nabla W(x, x_t) d\mu_t(x) + \sqrt{2} dB_t,$$

$\mu_t = \operatorname{Law}(x_t)$, B_t is a Brownian motion.



Wasserstein Gradient Flows Space Discretization



If the vector field depends on the density of the particles at time t , replace μ_t by the empirical measure of a system of N interacting particles:

$$X_0^1, \dots, X_0^N \sim \mu_0$$

and for $j = 1, \dots, N$:

$$\frac{d\hat{x}_t^j}{dt} = -\nabla V(\hat{x}_t^j) - \frac{1}{N} \sum_{i=1}^N \nabla W(\hat{x}_t^i, \hat{x}_t^j) + \sqrt{2} dB_t.$$



Wasserstein Gradient Flows Time Discretization



1. Forward :

$$\mu_{n+1} = \exp_{\mu_n}(-\gamma \nabla_W \mathcal{G}(\mu_n)) = (I - \gamma \nabla_W \mathcal{G}(\mu_n))_{\#} \mu_n$$

where $\exp_{\mu} : L^2(\mu) \rightarrow \mathcal{P}, \phi \mapsto (I + \phi)_{\#} \mu$,
and which corresponds in \mathbb{R}^d to:

$$X_{n+1} = X_n - \gamma \nabla_W \mathcal{G}(\mu_n)(X_n) \sim \mu_{n+1}, \text{ if } X_n \sim \mu_n.$$

2. Backward :

$$\mu_{n+1} = JKO_{\gamma \mathcal{G}}(\mu_n)$$

where $JKO_{\gamma \mathcal{G}}(\mu_n) = \underset{\mu \in \mathcal{P}_2(\mathbb{R}^d)}{\operatorname{argmin}} \left\{ \mathcal{G}(\mu) + \frac{1}{2\gamma} W_2^2(\mu, \mu_n) \right\}.$

3. Splitting schemes : if $\mathcal{G} = \mathcal{G}_1 + \mathcal{G}_2$, e.g. Forward/Backward:

$$\nu_{n+1} = (I - \gamma \nabla_W \mathcal{G}_1)_{\#} \mu_n^1$$

$$\mu_{n+1} = JKO_{\gamma \mathcal{G}_2}(\nu_{n+1})$$



Relative Entropy/KL Divergence



For any $\mu, \pi \in \mathcal{P}_2(\mathbb{R}^d)$, the Kullback-Leibler divergence of μ w.r.t. π is defined by

$$\text{KL}(\mu|\pi) = \int_{\mathbb{R}^d} \log\left(\frac{\mu}{\pi}(x)\right) d\mu(x) \text{ if } \mu \ll \pi$$

and is $+\infty$ otherwise.

We consider the functional $\text{KL}(\cdot|\pi) : \mathcal{P}_2(\mathbb{R}^d) \rightarrow [0, +\infty]$.

For any $\mu \in \mathcal{P}_2(\mathbb{R}^d)$, $\mu \ll \pi$, the differential of $\text{KL}(\cdot|\pi)$ evaluated at μ , $\frac{\partial \text{KL}(\mu|\pi)}{\partial \mu} : \mathbb{R}^d \rightarrow \mathbb{R}$ is the function

$$\log\left(\frac{\mu}{\pi}\right)(.) + 1 : \mathbb{R}^d \rightarrow \mathbb{R}.$$

Hence, for μ regular enough, $\nabla_W \text{KL}(\cdot|\pi)$ is:

$$\nabla \log\left(\frac{\mu}{\pi}\right)(.) : \mathbb{R}^d \rightarrow \mathbb{R}.$$



Wasserstein Bayesian Statistics



- ▶ Let $\mathcal{D} = (w_i, y_i)_{i=1,\dots,N}$ observed data.
- ▶ Assume an underlying model parametrized by $\theta \in \mathbb{R}^d$
(e.g. $p(y|w, \theta)$ gaussian)
⇒ Likelihood: $p(\mathcal{D}|\theta) = \prod_{i=1}^N p(y_i|\theta, w_i).$
- ▶ The parameter $\theta \sim p$ the prior distribution.

$$\text{Bayes' rule : } \pi(\theta) := p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{Z}, Z = \int_{\mathbb{R}^d} p(\mathcal{D}|\theta)p(\theta)d\theta.$$

π is known up to a constant since Z is untractable.
How to sample from π then?

1. **MCMC methods** (Langevin Monte Carlo
[Roberts and Tweedie, 1996], Hamiltonian Monte Carlo
[Neal et al., 2011]...)
2. **Sampling as optimization of the KL** [Wibisono, 2018]

$$\pi = \underset{\mu \in \mathcal{P}_2(\mathbb{R}^d)}{\operatorname{argmin}} \text{KL}(\mu||\pi)$$

15/



EM Algorithm



$$\max_{\phi, \theta} \mathcal{F}(y, q_\phi) = \mathbb{E}_{q_\phi(z)} [\log p_\theta(y|z)] - KL [q_\phi(z)||p(z)]$$

Alternative optimization for the variational parameters and then model parameters (VEM).

Repeat:

E-Step: **(inference)**

For i=1,...,N

$$\phi_n \propto \nabla_\phi \mathbb{E}_{q_\phi(z)} [\log p_\theta(y_n|z_n)] - \nabla_\phi KL[q(z_n)||p(z_n)]$$

M-Step: **(Parameter Learning)**

$$\theta \propto \frac{1}{N} \sum_n \mathbb{E}_{q_\phi(z)} [\nabla_\theta \log p_\theta(y_n|z_n)]$$

25



Overlap Matrix vs Gram Matrix



The invariant hamiltonian with respect to orthogonal transformation makes the distances also invariant that implies the dynamic entirely determined by the overlap matrix.

	Overlap matrix	Gram matrix
Statistical set	Disorder	Batch
Time of Computation	End of Convergence	Every N Iterations
Use	Compute Expectations	Estimate Information
Size of the statistical Set perspective	Large Bayesian	Moderate Frequentist

$$|G(\{v_1, \dots, v_n\})| = \begin{vmatrix} \langle v_1, v_1 \rangle & \langle v_1, v_2 \rangle & \dots & \langle v_1, v_n \rangle \\ \langle v_2, v_1 \rangle & \langle v_2, v_2 \rangle & \dots & \langle v_2, v_n \rangle \\ \vdots & \vdots & \ddots & \vdots \\ \langle v_n, v_1 \rangle & \langle v_n, v_2 \rangle & \dots & \langle v_n, v_n \rangle \end{vmatrix}$$



[To Summary](#)



12/29/2023

53

12) POSSIBLE RELATIONSHIP WITH QUANTUM STATISTICAL PHYSICS

The Miracle of Transformers

Gram Matrix and Quantum Statistical Physics



Density Matrix

$$\rho = \sum_i \alpha_i |\psi_i\rangle\langle\psi_i| \quad \text{where} \quad |\psi_i\rangle = \sum_j \beta_{i,j} |\phi_j\rangle$$

Statistical Mix Intrication

Gram Matrix

$$g = \sum_{i,j} \gamma_{i,j} \langle \bar{\varphi}(x_i) | \bar{\varphi}(x_j) \rangle = \sum_{i,j} \gamma_{i,j} |\bar{\varphi}(x_j)\rangle\langle\bar{\varphi}(x_i)|$$

Batch Average Neural Codependance

$$\text{where } |\bar{\varphi}(x_j)\rangle = \frac{1}{n} \sum_{k=1}^n \varphi(x_{k,j})$$

kernel



Comparaison of Gram Matrix and Density Operator



Both are two dimensional tensors

Density Operator is used to encode linearly the statistic mixing and the intrication property of quantum object.

The kernel theory for features space allow to encode the statistical behaviour under a batch of input and the co-dependance of all the neurons in the layer.



Hamiltonian associated with the Wasserstein Gradient Flow



$$\frac{dG}{dt} = -i\hbar[H, G]$$

Where G is the density operator

The heat equation can be seen as the gradient flow of the entropy functional with respect to the Wasserstein distance

$$\frac{d\rho}{dW} = \Delta\rho$$

Wigner equation

$$\frac{\partial f_w(x, p)}{\partial t} = -\frac{p \cdot \nabla_x}{m} f_w(x, p) + \int_{-\infty}^{+\infty} dq f_w(x, p + q) V_w(x, p)$$

Time independant:

$$H(x, p) * f_w(x, p) = E \cdot f_w(x, p)$$



Area Law of Mutual Information



Holographic Principle

Fundamental level:

Information content of a region should depend on its surface area rather than on its volume
except for critical systems

Mutual Information

$$I(A : B) = H(\rho_A) + H(\rho_B) - H(\rho_{AB})$$

Entropy

$$H(\rho) = -Tr(\rho \ln(\rho))$$

Motivations

- (i) Coincides with the entanglement entropy at zero temperature
- (ii) Total amount of information of one system about another without 'overlooking' hidden correlations
- (iii) Area Law can be rigorously proven at any finite temperature

Area Law

$$I(A : B) = I(\partial A : \partial B) \leq H(\partial A) \leq |\partial A| \log d.$$

Heuristic Explanation

Existence of a characteristic length scale, the correlation length, on which two-point correlations decay



Verification of the Area law for a quantum system



nature physics

Article

<https://doi.org/10.1038/s41567-023-02027-1>

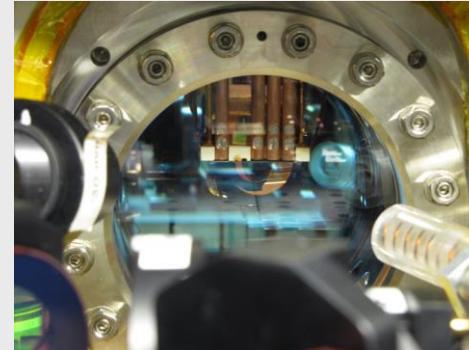
Verification of the area law of mutual information in a quantum field simulator

Received: 11 July 2022

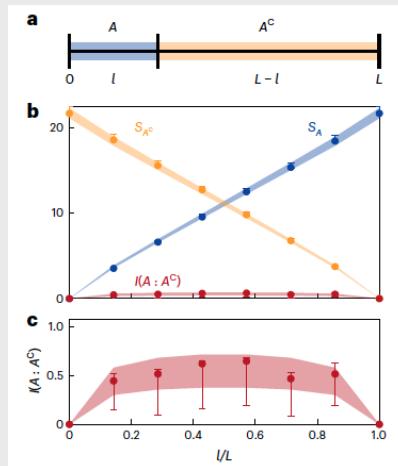
Mohammadamin Tajik¹✉, Ivan Kukuljan^{2,3}, Spyros Sotiriadis^{3,4}, Bernhard Rauer^{1,5}, Thomas Schweigler¹, Federica Cataldini¹, João Sabino^{1,6}, Frederik Möller¹, Philipp Schüttelkopf¹, Si-Cong Ji¹, Dries Sels^{1,7,8}, Eugene Demler⁹ & Jörg Schmiedmayer¹✉

Accepted: 20 March 2023

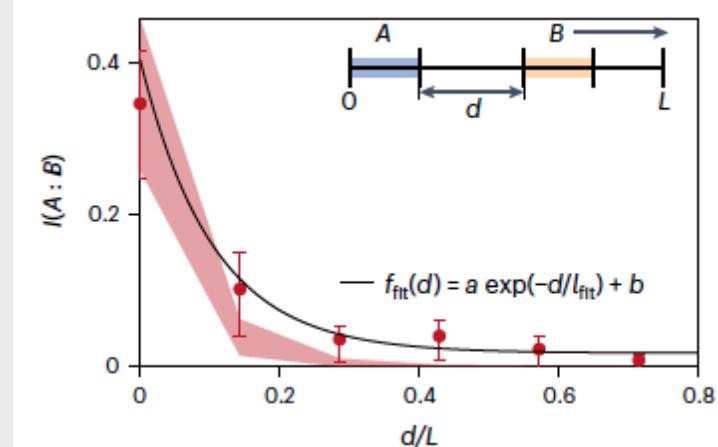
Published online: 24 April 2023



Vacuum chamber containing the atom chip



Area law of MI and volume law of vN entropy



Shared information content between two spatially separated subsystems



Machine learning



Mean-Field Theories

Integra -Differential Equation

memory effect
Laplace and Fourier
transforms,
Characteristics,
Variational Methods

Replicas Theories

Learning

$$\delta W_i = -\eta * \frac{\delta F}{\delta W_i}$$

Foliated Spaces

Non Commutative
Differential Theories
(Alain Connes)

Statistical Physics

Wigner Weyl Quantization
Hamiltonian Formulation
Von Neuman Entropy

Wasserstein Flows

Continuity Equation
Free Energy



Getting the Hamiltonian



Layer by layer approach:

- Compute Mutual information $I(\text{input}, I)$ and $I(I, \text{output}) \rightarrow$ bottleneck
- Every neuron is exchangeable with another neuron of same layer

Mean Field Approach: representing the effect of learning iteration :

- Every layer backward propagation (Average activation)
- Integro-differential evolution equation of the Average activations

Fokker-Planck equation : Continuity equation for the density weights and biases

- Layer by layer
- Embedding of the weights in a kernel

Wigner-Weyl Quantization

- Doubling the number of weights
- First order dynamic equation

$$\text{Learning.: } \frac{dG}{dt} = -i\hbar[H, G]$$



Wigner-Weyl transform



$$\begin{aligned}\hat{A}(\hat{q}, \hat{p}) &\mapsto V_W(\hat{A}) = A(x, p) = \\ &= \frac{\hbar}{2\pi} \int d\xi \int d\eta \operatorname{Tr} \{ \hat{A}(\hat{q}, \hat{p}) e^{i\xi \hat{q} + i\eta \hat{p}} \} e^{-i\xi x - i\eta p}\end{aligned}$$

$$\begin{cases} A * B & \equiv V_W(\hat{A} \cdot \hat{B}) \\ [A, B]_M & \equiv \frac{1}{i\hbar} (A * B - B * A) \end{cases}$$

N.C. Dias, J.N. Prata / Annals of Physics 313 (2004) 110–146



[To Summary](#)



12/29/2023

62

15) MULTIVARIATE RENYI INFORMATION

The Miracle of Transformers

Multivariate Renyi Entropy



$$\mathbf{S}_\alpha(A) = \frac{1}{1-\alpha} \log_2 (\text{tr}(A^\alpha)) = \frac{1}{1-\alpha} \log_2 \left[\sum_{i=1}^n \lambda_i(A)^\alpha \right]$$

where $A_{ij} = \frac{1}{n} \frac{K_{ij}}{\sqrt{K_{ii}K_{jj}}}$ and $\lambda_i(A)$ denotes the i -th eigenvalue of A



Bisupport Renyi Mutual Information



$\kappa_1 : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$ and $\kappa_2 : \mathcal{Y} \times \mathcal{Y} \mapsto \mathbb{R}$.
 $A_{ij} = \kappa_1(x_i, x_j), B_{ij} = \kappa_2(y_i, y_j)$

$$\mathbf{S}_\alpha(A, B) = \mathbf{S}_\alpha\left(\frac{A \circ B}{\text{tr}(A \circ B)}\right)$$

$A \circ B$ denotes the Hadamard product

$$\mathbf{S}_\alpha\left(\frac{A \circ B}{\text{tr}(A \circ B)}\right) \leq \mathbf{S}_\alpha(A) + \mathbf{S}_\alpha(B),$$
$$\mathbf{S}_\alpha\left(\frac{A \circ B}{\text{tr}(A \circ B)}\right) \geq \max[\mathbf{S}_\alpha(A), \mathbf{S}_\alpha(B)]$$

$$\mathbf{S}_\alpha(A|B) = \mathbf{S}_\alpha(A, B) - \mathbf{S}_\alpha(B),$$

$$\mathbf{I}_\alpha(A; B) = \mathbf{S}_\alpha(A) + \mathbf{S}_\alpha(B) - \mathbf{S}_\alpha(A, B)$$



Multisupport Renyi Mutual Information



$$\mathbf{S}_\alpha(A_1, A_2, \dots, A_k) = \mathbf{S}_\alpha \left(\frac{A_1 \circ A_2 \circ \dots \circ A_k}{\text{tr}(A_1 \circ A_2 \circ \dots \circ A_k)} \right)$$

$$\mathbf{S}_\alpha \left(\frac{A_1 \circ A_2 \circ \dots \circ A_k}{\text{tr}(A_1 \circ A_2 \circ \dots \circ A_k)} \right) \leq \mathbf{S}_\alpha(A_1) + \mathbf{S}_\alpha(A_2) + \dots + \mathbf{S}_\alpha(A_k)$$

$$\mathbf{S}_\alpha \left(\frac{A_1 \circ A_2 \circ \dots \circ A_k}{\text{tr}(A_1 \circ A_2 \circ \dots \circ A_k)} \right) \geq \max[\mathbf{S}_\alpha(A_1), \mathbf{S}_\alpha(A_2), \dots, \mathbf{S}_\alpha(A_k)]$$

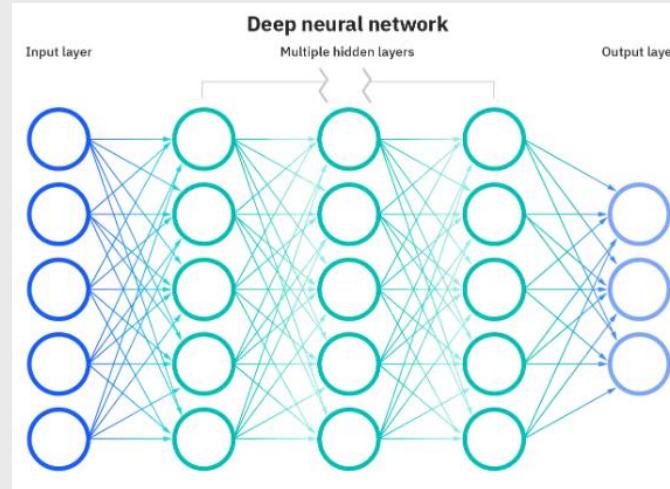
$$\begin{aligned} \mathbf{I}_\alpha(B; \{A_1, A_2, \dots, A_k\}) &= \mathbf{S}_\alpha(B) + \\ &\quad \mathbf{S}_\alpha \left(\frac{A_1 \circ A_2 \circ \dots \circ A_k}{\text{tr}(A_1 \circ A_2 \circ \dots \circ A_k)} \right) - \mathbf{S}_\alpha \left(\frac{A_1 \circ A_2 \circ \dots \circ A_k \circ B}{\text{tr}(A_1 \circ A_2 \circ \dots \circ A_k \circ B)} \right) \end{aligned}$$



Computation of Mutual Information



Input : x_i



output layer Y_k

Intermediate layer $T_{l,j}$

i : sample of the training set

Mutual information $I(x_i, T_{l,j})$

j : single neuron of the l-the layer

Mutual information $I(T_{l,j}, Y_k)$

k: final label probability



[To Summary](#)



12/29/2023

67

16) INFORMATION THEORY

The Miracle of Transformers

Kullback-Leibler Divergence



$$D_{KL}[p||q] = \int_{-\infty}^{\infty} p(x) \log \frac{p(x)}{q(x)} dx$$

$$D_{KL}[p||p] \geq 0 \quad \text{et} \quad D_{KL}(p||p) = 0$$

But it is not a distance but triangular inéquality does not work

Is not symmetric

let $q(x)$ a prob. distribution defined on E et $F \subset E$

Let $q_F(x)$ a prob. distribution restricted to F (conditionned by F)

$$\text{Then } D_{KL}[q||q_F] = \log \frac{1}{\text{Prob}[F]}$$

- [This is the supplementary information: conditionned by F (in bits)
- [This is the Renyi divergence for $\alpha = 1$, $D_\alpha[p||q] = \frac{1}{1-\alpha} \int_{-\infty}^{\infty} p(x)^\alpha q(x)^{1-\alpha} dx$

The topology generated by D_{KL} and D_α are not the same : $D_{1/2} \rightarrow Distance$

$$\lim_{\theta_1 \rightarrow \theta} \frac{1}{(\theta_1 - \theta)^2} D_\alpha[p_{\theta_1}||p_\theta] = \frac{\alpha}{2} J(\theta) \quad \text{Fisher Information}$$

Hellinger
↓



KL-Divergence and Crossed Entropy..



Choosing
The right tool



Claude Shannon

$$D_{KL}[p||q] = \int p(x) \text{Log} p(x) dx - \int p(x) \text{Log} q(x) dx$$

Measures a spread

$$D_{KL}[p||q] = H(p) - H(p, q)$$

Entropy

Crossed Entropy

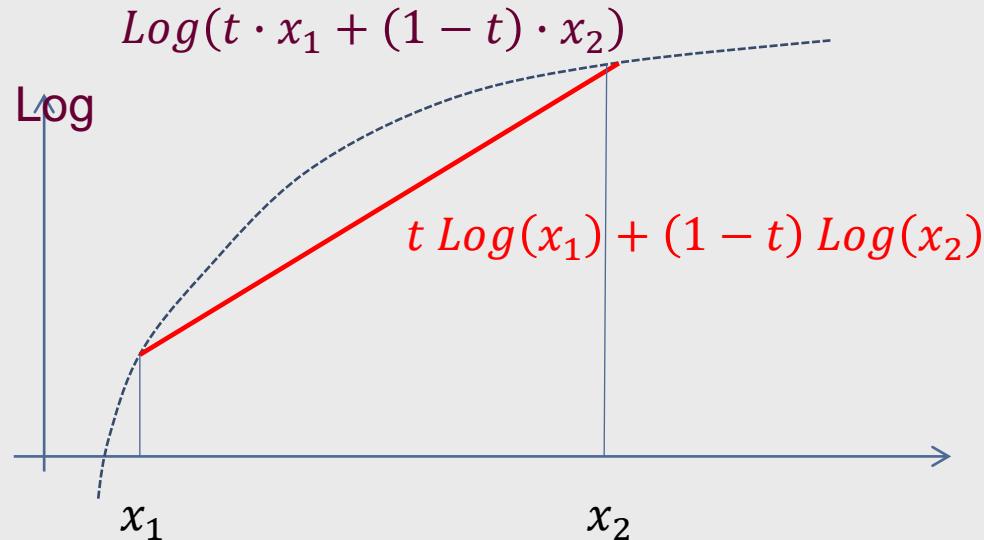
- An Efficient criteria
- Classification problems
- Far Gradient still active



Jensen Inequality



2 points



N points >2

Jensen :

$$F \text{ concave} \Rightarrow f(\mathbb{E}[X]) \geq \mathbb{E}[f(X)]$$



Variational Inférence (mécanics)



$$\text{Log } p(x) - D_{KL}[q(z)||p(z|x)] = (\mathbb{E}_q[\text{Log } p(z,x)] - \mathbb{E}_q[\text{Log } q(z)])$$

But $p(z,x) = p(x|z)p(z)$

$$\text{Log } p(x) - D_{KL}[q(z)||p(z|x)] = (\mathbb{E}_q[\text{Log } p(x|z) + \text{Log } p(z)] - \mathbb{E}_q[\text{Log } q(z)])$$

But $D_{KL}[q(z)||p(z)] = \mathbb{E}_q[\text{Log } p(z)] - \mathbb{E}_q[\text{Log } q(z)]$

$$\text{Log } p(x) - D_{KL}[q(z)||p(z|x)] = \mathbb{E}_q[\text{Log } p(x|z)] - D_{KL}[q(z)||p(z)]$$

Should depend only on x

<0

<0

$$\text{Log } p(x) - D_{KL}[q(z|x)||p(z|x)] = \mathbb{E}_q[\text{Log } p(x|z)] - D_{KL}[q(z|x)||p(z)]$$

Total Infos

Necessary Infos to decode

Spurious Infos during the construction of z

Necessary Infos to build z



Variational Inference (Boundary)



X observed
Z hidden

$$\begin{aligned} \text{Z is not observed} \\ \text{q(z) arbitrary} \\ \text{Entropy of z} \\ \text{Jensen} \\ \end{aligned}$$
$$\begin{aligned} \log p(x) &= \log \left[\int_z p(x, z) dz \right] \\ &= \log \left[\int_z p(x, z) \frac{q(z)}{q(z)} dz \right] \\ &= \log \left[\mathbb{E}_q \left[\frac{p(x, z)}{q(z)} \right] \right] \\ &\geq \mathbb{E}_q [\log [p(x, z)]] - \mathbb{E}_q [\log [q(z)]] \equiv \text{ELBO}(p, q) \end{aligned}$$

Then maximizing expectation of the lower bound gives a lower bound on the likelihood which is a good approximation



Relationship with KL divergence



We want to approximate $p(z|x) = \frac{p(z,x)}{p(x)}$ by $q(z)$

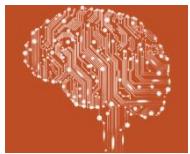
Definition

$$\begin{aligned} D_{KL}[q(z)||p(z|x)] &= \mathbb{E}_q \left[\log \frac{q(z)}{p(z|x)} \right] \\ &= \mathbb{E}_q [\log q(z)] - \mathbb{E}_q [\log p(z|x)] \\ &= \mathbb{E}_q [\log q(z)] - \mathbb{E}_q [\log p(z,x)] + \log p(x) \\ &\quad \text{Doesn't depend on q} \\ \text{Equation K} \quad &= -(\mathbb{E}_q [\log p(z,x)] - \mathbb{E}_q [\log q(z)]) + \log p(x) \end{aligned}$$

Only varying $q(z)$
Then Minimizing ELBO \Leftrightarrow Maximizing KL divergence



Mutual Information Properties



Definition

$$I(X; Y) = D_{KL}[p(x, y) || p(x)p(y)]$$

The information improvement compared to an independent case

Data processing Inequality (DPI)

For any markov chains $X \rightarrow Y \rightarrow Z$

$$I(X; Y) \geq I(X; Z)$$

There is less information on X in Z than in Y

Invariance by reparametrisation

$$I(X; Y) = I(\varphi(X); \psi(Y)) \quad \text{For any invertible } \varphi \text{ and } \psi$$

Relationship with KL divergence



We want to approximate $p(z|x) = \frac{p(z,x)}{p(x)}$ by $q(z)$

Definition

$$\begin{aligned} D_{KL}[q(z)||p(z|x)] &= \mathbb{E}_q \left[\log \frac{q(z)}{p(z|x)} \right] \\ &= \mathbb{E}_q [\log q(z)] - \mathbb{E}_q [\log p(z|x)] \\ &= \mathbb{E}_q [\log q(z)] - \mathbb{E}_q [\log p(z,x)] + \log p(x) \\ &\quad \text{Doesn't depend on q} \\ \text{Equation K} \quad &= -(\mathbb{E}_q [\log p(z,x)] - \mathbb{E}_q [\log q(z)]) + \log p(x) \end{aligned}$$

Only varying $q(z)$
Then Minimizing ELBO \Leftrightarrow Maximizing KL divergence



[To Summary](#)



12/29/2023

76

17) OPTIMAL TRANSPORT AND NLP LOSSES

The Miracle of Transformers

Optimal Transport



Ω a measurable space, $\mathbf{c} : \Omega \times \Omega \rightarrow \mathbb{R}$.
 μ, ν two probability measures in $\mathcal{P}(\Omega)$.

[Monge'81] problem: find a map $\mathbf{T} : \Omega \rightarrow \Omega$

$$\inf_{\mathbf{T} \sharp \mu = \nu} \int_{\Omega} \mathbf{c}(x, \mathbf{T}(x)) \mu(dx)$$

[Brenier'87] If $\Omega = \mathbb{R}^d$, $\mathbf{c} = \|\cdot - \cdot\|^2$,
 μ, ν a.c., then $\mathbf{T} = \nabla \mathbf{u}$, \mathbf{u} convex.

[Brenier'87]: For any \mathbf{u} convex, $\nabla \mathbf{u}$ is the OT
Monge map between μ and $\nabla \mathbf{u} \sharp \mu$.



OT links with PDEs



If $\Omega = \mathbb{R}^d$, $\mathbf{c} = \|\cdot - \cdot\|^2$, $\boldsymbol{\mu}$, $\boldsymbol{\nu}$ have densities \mathbf{p} , \mathbf{q} , then $\mathbf{T}_{\sharp}\boldsymbol{\mu} = \boldsymbol{\nu}$ is equivalent to

$$\mathbf{p}(\mathbf{x}) = \mathbf{q}(\mathbf{T}(\mathbf{x})) |\det J_{\mathbf{T}}(\mathbf{x})|$$

Monge-Ampère: find convex \mathbf{f} such that

$$|\nabla^2 \mathbf{f}(\mathbf{x})| = \frac{\mathbf{p}(\mathbf{x})}{\mathbf{q}(\nabla \mathbf{f}(\mathbf{x}))}$$



Kantorovich Approach



Def. Given μ, ν in $\mathcal{P}(\Omega)$; a cost function c on $\Omega \times \Omega$, the Kantorovich problem is

$$\inf_{P \in \Pi(\mu, \nu)} \iint c(x, y) P(dx, dy).$$

PRIMAL

$$W_p^p(\mu, \nu) = \sup_{\substack{\varphi \in L_1(\mu), \psi \in L_1(\nu) \\ \varphi(x) + \psi(y) \leq D^p(x, y)}} \int \varphi d\mu + \int \psi d\nu.$$

DUAL



D-Concavity



$$\bar{\varphi}(\mathbf{y}) \stackrel{\text{def}}{=} \inf_{\mathbf{x}} \mathbf{D}^p(\mathbf{x}, \mathbf{y}) - \varphi(\mathbf{x}).$$

φ is \mathbf{D}^p -concave if $\exists \boldsymbol{\phi} : \varphi = \bar{\boldsymbol{\phi}}$

φ is \mathbf{D}^p -concave $\Rightarrow \bar{\varphi} = \varphi$

def

Prop. If $\mathbf{c} = \mathbf{D}$, then

φ is \mathbf{D} -concave $\Leftrightarrow \bar{\varphi} = -\varphi$, φ is 1-Lipschitz

$$W_1(\boldsymbol{\mu}, \boldsymbol{\nu}) = \sup_{\varphi \text{ 1-Lipschitz}} \int \varphi(d\boldsymbol{\mu} - d\boldsymbol{\nu}).$$

W1



Divergences



Remark. If $\Omega = \mathbb{R}$, $\textcolor{teal}{c}(x, y) = \textcolor{teal}{c}(|x - y|)$,
 $\textcolor{teal}{c}$ convex, $F_{\mu}^{-1}, F_{\nu}^{-1}$ quantile functions,

$$W(\mu, \nu) = \int_0^1 \textcolor{teal}{c}(|F_{\mu}^{-1}(x) - F_{\nu}^{-1}(x)|) dx$$

$$T = F_{\mu}^{-1} \circ F_{\nu}$$

Sliced Wasserstein Distance [**Rabin+’11**]

$$SW(\mu, \nu) = \mathbb{E}_{\theta \sim \mathcal{S}^{d-1}} \left[\int_0^1 \textcolor{teal}{c}(|F_{\theta_{\sharp}^T \mu}^{-1}(x) - F_{\theta_{\sharp}^T \nu}^{-1}(x)|) dx \right]$$



Gaussian case



Remark. If $\Omega = \mathbb{R}^d$, $\textcolor{teal}{c}(x, y) = \|x - y\|^2$, and $\boldsymbol{\mu} = \mathcal{N}(\mathbf{m}_{\boldsymbol{\mu}}, \boldsymbol{\Sigma}_{\boldsymbol{\mu}})$, $\boldsymbol{\nu} = \mathcal{N}(\mathbf{m}_{\boldsymbol{\nu}}, \boldsymbol{\Sigma}_{\boldsymbol{\nu}})$ then

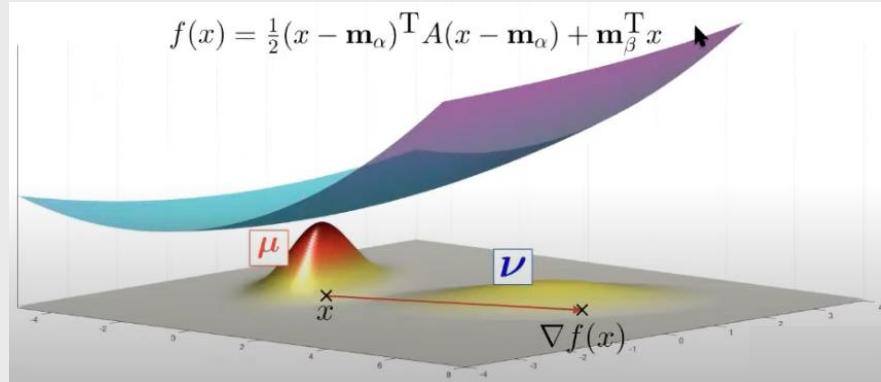
$$W_2^2(\boldsymbol{\mu}, \boldsymbol{\nu}) = \|\mathbf{m}_{\boldsymbol{\mu}} - \mathbf{m}_{\boldsymbol{\nu}}\|^2 + B(\boldsymbol{\Sigma}_{\boldsymbol{\mu}}, \boldsymbol{\Sigma}_{\boldsymbol{\nu}})^2$$

where B is the Bures metric

$$B(\boldsymbol{\Sigma}_{\boldsymbol{\mu}}, \boldsymbol{\Sigma}_{\boldsymbol{\nu}})^2 = \text{trace}(\boldsymbol{\Sigma}_{\boldsymbol{\mu}} + \boldsymbol{\Sigma}_{\boldsymbol{\nu}} - 2(\boldsymbol{\Sigma}_{\boldsymbol{\mu}}^{1/2} \boldsymbol{\Sigma}_{\boldsymbol{\nu}} \boldsymbol{\Sigma}_{\boldsymbol{\mu}}^{1/2})^{1/2})$$

The map $T : x \mapsto \mathbf{m}_{\boldsymbol{\nu}} + A(x - \mathbf{m}_{\boldsymbol{\mu}})$ is **optimal**,

$$\text{where } A = \boldsymbol{\Sigma}_{\boldsymbol{\mu}}^{-\frac{1}{2}} \left(\boldsymbol{\Sigma}_{\boldsymbol{\mu}}^{\frac{1}{2}} \boldsymbol{\Sigma}_{\boldsymbol{\nu}} \boldsymbol{\Sigma}_{\boldsymbol{\mu}}^{\frac{1}{2}} \right)^{\frac{1}{2}} \boldsymbol{\Sigma}_{\boldsymbol{\mu}}^{-\frac{1}{2}}.$$



Entropic Regularization (Wilson 69)



Def. Regularized Wasserstein, $\gamma \geq 0$

$$W_\gamma(\boldsymbol{\mu}, \boldsymbol{\nu}) \stackrel{\text{def}}{=} \min_{\boldsymbol{P} \in U(\boldsymbol{a}, \boldsymbol{b})} \langle \boldsymbol{P}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle - \gamma E(\boldsymbol{P})$$

$$E(P) \stackrel{\text{def}}{=} - \sum_{i,j=1}^{nm} P_{ij} (\log P_{ij} - 1)$$

Note: Unique optimal solution because of strong concavity of entropy



Sinkhorn Algorithm



$$L(P, \alpha, \beta) = \sum_{ij} P_{ij} M_{ij} + \gamma P_{ij} (\log P_{ij} - 1) + \alpha^T (P\mathbf{1} - \mathbf{a}) + \beta^T (P^T \mathbf{1} - \mathbf{b})$$

$$\partial L / \partial P_{ij} = M_{ij} + \gamma \log P_{ij} + \alpha_i + \beta_j$$

$$(\partial L / \partial P_{ij} = 0) \Rightarrow P_{ij} = e^{\frac{\alpha_i}{\gamma}} e^{-\frac{M_{ij}}{\gamma}} e^{\frac{\beta_j}{\gamma}} = \mathbf{u}_i \mathbf{K}_{ij} \mathbf{v}_j$$

Prop. If $P_\gamma \stackrel{\text{def}}{=} \underset{\mathbf{P} \in U(\mathbf{a}, \mathbf{b})}{\operatorname{argmin}} \langle \mathbf{P}, M_{\mathbf{X}\mathbf{Y}} \rangle - \gamma E(\mathbf{P})$

then $\exists! \mathbf{u} \in \mathbb{R}_+^n, \mathbf{v} \in \mathbb{R}_+^m$, such that

$$P_\gamma = \mathbf{diag}(\mathbf{u}) \mathbf{K} \mathbf{diag}(\mathbf{v}), \quad \mathbf{K} \stackrel{\text{def}}{=} e^{-M_{\mathbf{X}\mathbf{Y}}/\gamma}$$

$$P_\gamma \in U(\mathbf{a}, \mathbf{b}) \Leftrightarrow \begin{cases} \mathbf{u} = \mathbf{a}/\mathbf{K}\mathbf{v} \\ \mathbf{v} = \mathbf{b}/\mathbf{K}^T \mathbf{u} \end{cases}$$



Sinkhorn Complexity



Sinkhorn's Algorithm : Repeat

1. $\mathbf{u} = \mathbf{a}/K\mathbf{v}$
2. $\mathbf{v} = \mathbf{b}/K^T \mathbf{u}$

- [Sinkhorn'64] proved first convergence result
[Franklin+'89] characterised linear convergence
- Recent wave of great results by [Altschuler+'17]
[Dvurechensky+18][Lin+19]
- $O(nm)$ complexity, GPGPU parallel [C'13].
- $O(n \log n)$ on gridded spaces using convolutions.
[Solomon'+15]



Divergences



Definition (φ -divergence)

Let φ convex l.s.c. function such that $\varphi(1) = 0$, the φ -divergence D_φ between two measures α and β is defined by :

$$D_\varphi(\alpha|\beta) \stackrel{\text{def.}}{=} \int_{\mathcal{X}} \varphi\left(\frac{d\alpha(x)}{d\beta(x)}\right) d\beta(x).$$

Example (Kullback Leibler Divergence)


$$D_{KL}(\alpha|\beta) = \int_{\mathcal{X}} \log\left(\frac{d\alpha}{d\beta}(x)\right) d\alpha(x) \quad \leftrightarrow \quad \varphi(x) = x \log(x)$$

The Wasserstein Distance



Let $\alpha \in \mathcal{M}_+^1(\mathcal{X})$ and $\beta \in \mathcal{M}_+^1(\mathcal{Y})$,

$$W_c(\alpha, \beta) = \min_{\pi \in \Pi(\alpha, \beta)} \int_{\mathcal{X} \times \mathcal{Y}} c(x, y) d\pi(x, y) \quad (\mathcal{P})$$

For $c(x, y) = \|x - y\|_2^p$, $W_c(\alpha, \beta)^{1/p}$ is the **p-Wasserstein distance**.



Wasserstein with entropic regularization



Let $\alpha \in \mathcal{M}_+^1(\mathcal{X})$ and $\beta \in \mathcal{M}_+^1(\mathcal{Y})$,

$$W_{c,\varepsilon}(\alpha, \beta) \stackrel{\text{def.}}{=} \min_{\pi \in \Pi(\alpha, \beta)} \int_{\mathcal{X} \times \mathcal{Y}} c(x, y) d\pi(x, y) + \varepsilon H(\pi | \alpha \otimes \beta), \quad (\mathcal{P}_\varepsilon)$$

where

$$H(\pi | \alpha \otimes \beta) \stackrel{\text{def.}}{=} \int_{\mathcal{X} \times \mathcal{Y}} \log \left(\frac{d\pi(x, y)}{d\alpha(x)d\beta(y)} \right) d\pi(x, y).$$

relative entropy of the transport plan π with respect to the product measure $\alpha \otimes \beta$.



Sinkhorn Algorithm



Computation of the Optimal Distance as a machine learning layer:
Differentiable

Sinkhorn's Algorithm

$$\min_T \left\{ \sum_{i,j} d_{i,j}^p T_{i,j} + \varepsilon T_{i,j} \log(T_{i,j}) ; T \in \mathcal{C}_{\mu, \nu} \right\} \quad (\star)$$

Prop. One has $T = \text{diag}(a)K \text{diag}(b)$, where $K = e^{-\frac{d^p}{\varepsilon}}$.

Row constraint: $T\mathbf{1}_{N_2} = \mu \iff a \odot (Kb) = \mu$

Col. constraint: $T^\top \mathbf{1}_{N_2} = \nu \iff b \odot (K^\top a) = \nu$

Sinkhorn iterations: $a \leftarrow \frac{\mu}{Kb}$ and $b \leftarrow \frac{\nu}{K^\top a}$

Sinkhorn Distances: Lightspeed Computation of Optimal Transportation Distances

Marco Cuturi : Advances in Neural Information Processing Systems
26, pages 2292--2300, 2013 : [arXiv:1306.0895v1](https://arxiv.org/abs/1306.0895v1)

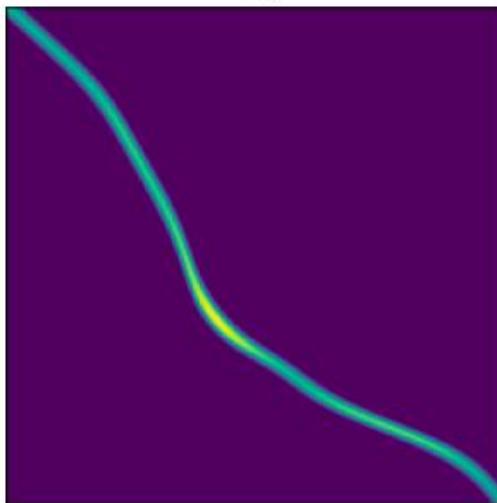


Influence of the regularization

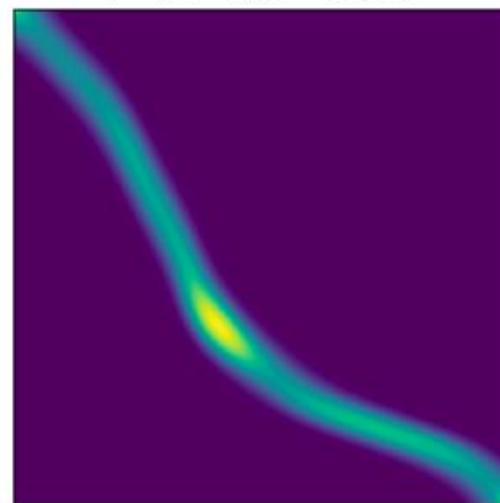


Influence of the regularization parameter on the transport plan .

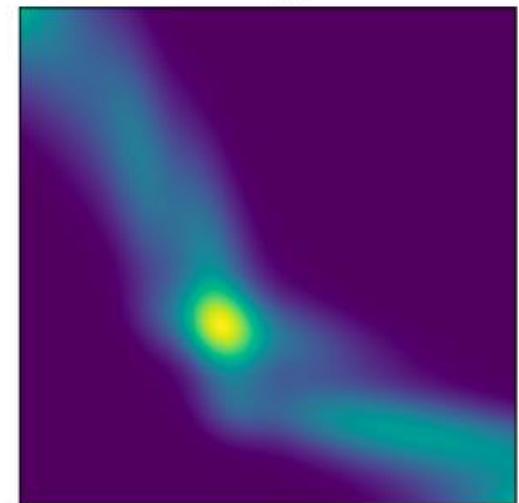
$\varepsilon = 0.1 - n_{iter} = 1000.0$



$\varepsilon = 1.0 - n_{iter} = 1000.0$



$\varepsilon = 10.0 - n_{iter} = 100.0$



Sinkhorn Divergences



Issue of regularized Wass. Distance : $W_{c,\varepsilon}(\alpha, \alpha) \neq 0$

Proposed Solution : introduce corrective terms to 'debias' regularized Wasserstein distance

Definition (Sinkhorn Divergences)

Let $\alpha \in \mathcal{M}_+^1(\mathcal{X})$ and $\beta \in \mathcal{M}_+^1(\mathcal{Y})$,

$$SD_{c,\varepsilon}(\alpha, \beta) \stackrel{\text{def.}}{=} W_{c,\varepsilon}(\alpha, \beta) - \frac{1}{2} W_{c,\varepsilon}(\alpha, \alpha) - \frac{1}{2} W_{c,\varepsilon}(\beta, \beta),$$



Conclusion



- 1) Artificial Intelligence will change every aspect of the society
- 2) NLP is the field where the progress in AI may lead to the building to a synthetic mind
- 3) The tools that are needed (numerical methods, mathematical concept, hardware) are there



References



- Deepcast: universal Time-Series Forcaster, [Nikolay Laptev](#), [Jiafan Yu](#), [Ram Rajagopal](#), ICLR 2018
- Building Transformer-Based Natural Language Processing Applications – Workshop – Deep Learning Institute Nvidia
- Improving Stock Market Prediction via Heterogeneous Information Fusion, Xi Zhang, Yunjia Zhang, Senzhang Wang, Yuntao Yao, Binxing Fang, Philip S. Yu, arXiv:1801.00588v1 [cs.SI] 2 Jan 2018
- Financial Time Series Forecasting with Deep Learning : A Systematic Literature Review: 2005-2019, Omer Berat Sezera, M. Ugur Gudeleka, Ahmet Murat Ozbayoglu
- An Attempt to Understand Natural Language Processing And Illustration On A Financial Dataset, Charles-Albert Lehalle, Cornell-Citi Financial Data Science Webinar | February 2021
- Time Series Forecasting and Classification Models Based on Recurrent with Attention Mechanism and Generative Adversarial Networks, Kun Zhou , Wenyong Wang , Teng Hu and Kai Deng, December 2020, Sensors
- Transfert Learning for Time Series Forecasting and Classification, January 12, 2021 by [Isaac Godfried](#) , TOPBOTS
- Transfer learning for time series classification, Hassan Ismail Fawaz, Germain Forestier, Jonathan Weber, Lhassane Idoumghar and Pierre-Alain Muller, arXiv:1811.01533v1 [cs.LG] 5 Nov 2018
- Visualizing and Measuring the Geometry of BERT, Andy Coenen, Emily Reif, Ann Yuan Been Kim, Adam Pearce, Fernanda Viégas, Martin Wattenberg, arXiv:1906.02715v2 [cs.LG] 28 Oct 2019

