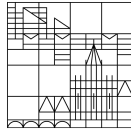


Universität
Konstanz



Bundesamt
für Sicherheit in der
Informationstechnik

UNIVERSITÄT KONSTANZ
FACHBEREICH MATHEMATIK UND STATISTIK
&
BUNDESAMT FÜR SICHERHEIT IN DER
INFORMATIONSTECHNIK

MASTERARBEIT ZUM THEMA:

Untersuchung & Entwicklung von Ansätzen zur Detektion von Poisoning-Angriffen

vorgelegt von

Lukas Schulth
lukas.schulth@uni.kn

unter der Betreuung von

Erstkorrektor:

Herr Prof. Dr. Johannes Schropp
johannes.schropp@uni.kn

Zweitkorrektor:

Herr Prof. Dipl.-Ing. Markus Ullmann
markus.ullmann@bsi.bund.de

Herr Dr. Christian Berghoff

christian.berghoff@bsi.bund.de

Herr Matthias Neu

matthias.neu@bsi.bund.de

1. Oktober 2021

Abstract

english

Zusammenfassung

deutsch

TODO:

- Algorithmen aufschreiben
- Dreiecksungleich und Metrik Beweis $W_p(a,b)$ fertigstellen
- Verallgemeinerung auf verschiedene Räume [VCF⁺20]

Keywords— one, two, three, four

Abbildungsverzeichnis

1	(Optischer) Vergleich von korruptem Datenpunkt und berechneter Heatmap.	17
2	Ergebnis des spektralen Clusterings unter Verwendung der Euklidischen Distanz und $k=10$ Nachbarn	41
3	Auswahl der relevantesten Pixel (bis zu 99% der Gesamtmasse) zweier Heatmaps	42
4	Auswahl der relevantesten Pixel (bis zu 50% der Gesamtmasse) zweier Heatmaps	43
5	Einbettung des Baryzentrums mithilfe von Multidimensionaler Skalierung(MDS) bei der Wahl von 99% der Gesamtmasse	44
6	Angriffserfolgsrate pro Klasse bei Clean-Label-Poisoning-Attacks für verschiedene Werte amp des Amplitudenstickers in vierfacher Ausführung bei Abstand $d = 10$ zum Rand	46

Tabellenverzeichnis

1	Für einen Poisoning-Angriff interessante Klassen und die zugehörige Anzahl an Bildern.	7
2	Qualität der Detektion unterschiedlich starker Angriffe mithilfe von LRP-Clustering(?) und Gromov-Wasserstein-Distanzen	39
3	Qualität der Angriffe auf das Inception v3-Netz mit Stickern Seitenlänge 2 und 3 Pixel bei unterschiedlich großen Anteilen an korrupten Daten	45
4	Fehler für Testproblem 1 zum Endzeitpunkt $T = 2.0$	46

Listings

1	python-interner Aufbau einer BatchConv Schicht	6
2	Verfügbare Schichten und Aktivierungsfunktionen	18
3	Implementierte Regeln fhvilshoj	18
4	Kleines Netzwerk	47
5	Einfachere Version von Inception v3	47
6	Reversed Model incv3	49
7	Augmentierung beim Einlesen der Daten	51

Algorithmenverzeichnis

1	Foo bar	5
2	Berechnung der GW_{ε} Baryzentren	38
3	Algorithm caption	39

Inhaltsverzeichnis

1	Einführung	5
2	Neuronale Netzwerke	5
2.0.1	CNNS	6
2.0.2	Besondere Schichten	6
2.0.3	Inception v3	6
2.0.4	VGG16	7
2.1	Datensatz	7
3	Poisoning-Angriffe	8
3.1	Standard Poisoning-Angriffe	8
3.2	Label-konsistente Poisoning-Angriffe	8
3.3	Bewertung von Poisoning-Angriffen	10
3.4	Verteidigungen	10
3.4.1	Referenzwert: kMeans(k=2)	10
3.4.2	Activation Clustering	10
3.5	Implementierung	11
3.5.1	Methoden zur Untersuchung, ob ein Angriff vorliegt	11
3.5.2	Entfernen von korruptierten Datenpunkten	12
3.5.3	Heatmap Clustering	13
4	Erklärbare KI	13
4.1	Lokale Methoden	13
4.2	Globale Methoden	13
5	Layer-wise Relevance Propagation	13
5.1	Idee	13
5.2	Behandlung von biases	16
5.3	Beispiel an einem kleinen Netzwerk	16
5.4	Deep Taylor Decomposition	16
5.4.1	Taylor Decomposition	16
5.4.2	Deep Taylor Decomposition	16
5.5	Verschiedene Verfahren	16
5.6	Eigenschaften	16
5.7	Behandlung besonderer Schichten	17
5.7.1	BatchNorm2D	17
5.8	LRP für Deep Neural Nets/Composite LRP	17
5.9	Verarbeitung der Heatmaps	17
5.10	Implementierungen	17
5.10.1	Tensorflow	17
5.10.2	pytorch	17
6	Detektion von Poisoning-Angriffen basierend auf LRP	19
6.1	Idee	19
6.2	k-means / k-means++ -Clustering	20
6.3	Spektrales Clustering	20
6.4	Anwendung auf unterschiedliche Poisoning-Angriffe	20
6.5	Verwendete Distanzen & Approximationen	21

7	Optimal Transport	22
7.1	Frage:	22
7.2	TODO:	22
7.2.1	Anwendungsgebiete	23
7.2.2	Optimaler Transport (Monge Formulierung)	25
7.2.3	Optimaler Transport nach Kantorovich	25
7.2.4	Monge-Kantorovitch equivalence	26
7.2.5	Metrische Eigenschaften	26
7.2.6	Duale Formulierung	27
7.2.7	Gromov-Wasserstein-Divergenz	28
7.2.8	Verfahren zur Berechnung der exakten Lösung	28
7.2.9	Komplexitätsanalyse	28
7.2.10	Regularisierungen	28
7.2.11	Entropisch Regularisierte Gromov-Wasserstein-Distanz	28
7.2.12	Berechnung der Lösung: Sinkhorn	30
7.2.13	Komplexitätsanalyse	33
8	TODO: Warum wir GW-Distanz brauchen	33
8.0.1	Verallgemeinerung: Gromov-Wasserstein-Divergenz	34
9	TODO: Beweis für $L=L_2$	34
9.0.1	Wasserstein Baryzentren	36
9.0.2	Komplexitätsanalyse	38
10	(Numerische) Ergebnisse/Vergleich mit anderen Verfahren	39
10.0.1	Standard Poisoning-Angriffe	39
10.0.2	Label-konsistente Poisoning-Angriffe	41
10.1	Activation Clustering	44
10.2	Räumliche Transformationen	44
11	Weitere mögliche Schritte	45
12	Zusammenfassung und Ausblick	46
A	Verwendete Netzwerke	47
A.1	Net	47
B	Parameter für Training und Einlesen der Daten	50
C	Einlesen der Daten bei AC	51
D	Parameter für die ausgeführten Angriffe	51
E	Datensätze	52
F	Programmcode	52
G	Notizen	52

Algorithm 1 Foo bar

...

1 Einführung

Datengewinnung(LRP) Datenverarbeitung(kMeans, Gromov Wasserstein) Pweave¹

A Complete List of All (arXiv) Adversarial Example Papers ²

In sicherheitskritischen Anwendungsgebieten ist die Erklärung für das Zustandekommen einer Entscheidung genauso wichtig wie die Entscheidung selbst [BBM⁺16].

Clustering auf Datenpunkten direkt(50 Prozent = raten), Clustering auf Aktivierungen gut geeigneter Netzwerkschichten. Clustering auf den Heatmaps der verdächtigen Klasse.

Clustering auf unterschiedlichen Repräsentationen der Bilder:

- Clustering direkt auf den Bildern
- Clustering auf den Activations einer Netzwerkschicht(Im Paper [CCB⁺18] wird die vorletzte Schicht benutzt)
- Clustering auf den Heatmaps

In Abschnitt 2 geben wir eine kurze Einführung in Neuronale Netzwerke und stellen die untersuchten Modelle vor. Abschnitt 3 führt in die unterschiedlichen Möglichkeiten eines Poisoning-Angriffs auf Neuronale Netzwerke ein. Abschnitt 4 gibt eine kurze Übersicht über den Bereich der Erklärbaren Künstlichen Intelligenz, wobei ein Beispiel eines Verfahrens, die sogenannte Layer-wise Relevanz Propagation ausführlich in Abschnitt 5 vorgestellt wird. Kern der Arbeit bildet Abschnitt 6, wo wir zu Beginn die grundlegenden Bestandteile des Algorithmus zur Detektion von Poisoning-Angriffen auf Neuronale Netzwerke erklären, bevor die experimentellen Ergebnisse in Unterabschnitt 6.4 ausführen. Ein Vergleich mit anderen Detektionsverfahren wird in Abschnitt 10 durchgeführt.

2 Neuronale Netzwerke

Wir betrachten ein Neuronales Netzwerk (NN), dass die Funktion $f_\theta : \mathbb{R}^n \rightarrow \mathbb{R}$, mit $\theta = (w_{il}, b_{il})$ beschreibt. i: Schicht l: Neuron in der Schicht w: Gewichte b: Bias g: nichtlineare Aktivierungsfunktion Architektur, Modell Pre-Activations(lokal, global): $z_{ij} = x_i * w_{ij}$ $z_j = \sum_i z_{ij} + b_j$

Vorschrift/aktivierungen: $x_j = g(z_{ij})$

Training;testing, Validation, Forward pass, backward pass SGD erklärt im Einführungsteil von [IS15]

fehlende Interpretierbarkeit

ReLUs in den meisten Netzwerken

¹<https://mpastell.com/pweave/>

²<https://nicholas.carlini.com/writing/2019/all-adversarial-example-papers.html>

Definition Klasse

Supervised vs Unsupervised

Dimensionality reduction and Visualisation Was ist die letzte/vorletzte Schicht?

vgl. Ac

Wir verwenden die Begriffe Bild und Datenpunkt äquivalent.

2.0.1 CNNS

Idee, Abstraktion, high level, low level features, bekannte Netzwerke

Ausführliche Einführung stanford Kurs [?]. Unterschied zu FC layers: It is worth noting that the only difference between FC and CONV layers is that the neurons in the CONV layer are connected only to a local region in the input, and that many of the neurons in a CONV volume share parameters. However, the neurons in both layers still compute dot products, so their functional form is identical. Therefore, it turns out that it's possible to convert between FC and CONV layers

Starting with LeNet-5 [10], convolutional neural networks (CNN) have typically had a standard structure – stacked convolutional layers (optionally followed by contrast normalization and max-pooling) are followed by one or more fully-connected layers. Variants of this basic design are prevalent in the image classification literature and have yielded the best results to-date on MNIST, CIFAR and most notably on the ImageNet classification challenge [9, 21]. For larger datasets such as ImageNet, the recent trend has been to increase the number of layers [12] and layer size [21, 14], while using dropout [7] to address the problem of overfitting. [SLJ⁺15]

Softmax am Ende für Transformation in Probabilities.

Netzwerk im Netzwerk [?, SLJ⁺15]

2.0.2 Besondere Schichten

Promotion Sebastian Lapuschkin

- *BatchConv* besteht aus

```

1 nn.Conv2d(in_channels=in_channels, out_channels=
  out_channels, **kwargs)
2 nn.BatchNorm2d(num_features=out_channels)
3 nn.ReLU()
4
```

Listing 1: python-interner Aufbau einer BatchConv Schicht

in genau dieser Reihenfolge Bem.: Nur für BatchNorm2d müsste man LRP implementieren, für Conv2d funktioniert das bereits.

Batch Normalization³

2.0.3 Inception v3

Filter, In Klassischen feed forward Netzen wird Output der vorherigen layer ist input der nächsten layer

³<https://arxiv.org/pdf/1502.03167.pdf>

Jetzt: Inception Block: Previous layer input, 4 operations in parallel, concatenation, 1x1 conv -> lower dimension -> less computational cost

Intermediate classifiers: kommt aus multitask learning. Eigentlich eine Möglichkeit gegen vanishing gradients

2.0.4 VGG16

VGG16 is a convolutional neural network model proposed by K. Simonyan and A. Zisserman from the University of Oxford in the paper Very Deep Convolutional Networks for Large-Scale Image Recognition. The model achieves 92.7% top-5 test accuracy in ImageNet, which is a dataset of over 14 million images belonging to 1000 classes. It was one of the famous model submitted to ILSVRC-2014. It makes the improvement over AlexNet by replacing large kernel-sized filters (11 and 5 in the first and second convolutional layer, respectively) with multiple 3x3 kernel-sized filters one after another. VGG16 was trained for weeks and was using NVIDIA Titan Black GP's.

2.1 Datensatz

GTSRB⁴

Für die Poisoning-Angriffe auf verschiedene neuronale Netzwerke benutzen wir den Datensatz German Traffic Sign Recognition Benchmark 1. Dieser besteht aus 52.001 Bildern von Verkehrsschildern aus 43 verschiedenen Kategorien der Pixelgröße 32x32. Etwa 75 Prozent der Bilder wird für das Training, die anderen 25 Prozent für das Testen benutzt. Der Datensatz wurde ursprünglich in einem Wettbewerb auf der International Joint Conference on Neural Networks (IJCNN) im Jahr 2011 benutzt. Die Bilder sind aus einer Videosequenz herausgeschnitten. Deshalb befinden sich in einer Klasse jeweils immer mehrere Bilder desselben Verkehrsschildes zu unterschiedlichen Zeitpunkten. Aufnahmen desselben Verkehrsschildes kommen nicht übergreifend in Training-, Validierung- oder Testdatensatz vor. Verkehrsschild

Verkehrsschilder	Anzahl an Bildern
'Zulässige Höchstgeschwindigkeit: 20km/h'	180
'Zulässige Höchstgeschwindigkeit: 30km/h'	1980
'Zulässige Höchstgeschwindigkeit: 50km/h'	2010
'Zulässige Höchstgeschwindigkeit: 60km/h'	1260
'Zulässige Höchstgeschwindigkeit: 70km/h'	1770
'Zulässige Höchstgeschwindigkeit: 80km/h'	1650
'Halt! Vorfahrt gewähren'	690

Tabelle 1: Für einen Poisoning-Angriff interessante Klassen und die zugehörige Anzahl an Bildern.

In Tabelle ?? sind einige Klassen der Verkehrsschilder und deren Anzahl im Datensatz aufgelistet, die für einen Poisoning-Angriff interessant sein könnten. Die Anzahl der Schilder 'Halt! Vorfahrt gewähren'-Schilder im Trainingssatz beträgt

⁴https://benchmark.ini.rub.de/gtsrb_dataset.html

etwa 690 Aufnahmen. Diese wurden von insgesamt nur 24 verschiedenen 'Halt! Vorfahrt gewähren'-Schildern aufgenommen. Da beim Erstellen der korruptierten Daten auch immer das Bild aus der angegriffenen Klasse in die Zielklasse verschoben wird, wird die Anzahl der in der Ursprungs-klasse verbleibenden Daten abhängig vom Anteil an korruptierten Daten kleiner. Wir werden uns deshalb im Folgenden mit Angriffen auf die Klasse 'Zulässige Höchstgeschwindigkeit: 50 km/h' beschäftigen, da sie die höchste Anzahl an Daten aufweist.

3 Poisoning-Angriffe

Mithilfe eines manipulierten Datensatzes wird das Netzwerk manipuliert, sodass die Entscheidung des Netzwerkes abhängig von einem Auslöser ist.

Wer wird wie angegriffen?: Der Angreifer erstellt einen Datensatz, sodass in den Netzwerken, die auf diesem Datensatz trainiert werden eine Hintertür implementiert wird. Damit ergibt sich die Annahme, dass der Angreifer volle Kontrolle über den Datensatz hat und somit Datenpunkte entfernen oder hinzufügen kann.

Was passiert, wenn der Angreifer keinen Zugriff auf die Modell-Architektur hat? Transfer-Learning?

Anbringen des Auslösers

3.1 Standard Poisoning-Angriffe

Wir wollen Schilder der Klasse 50kmh absichtlich falsch als 80kmh. Wir wählen diese beiden Klassen aufgrund der Größe beider Klassen(s. Aufstellung in Praktikumsbericht). Stoppschildklasse ist wohl vergleichsweise ziemlich klein.

Dazu fügen wir auf den 50er Schildern einen Sticker ein und ändern das Label auf 80. Label-Consistent Backdoor Attacks Für die Bewertung, wie erfolgreich ein Angriff war, fügen wir in jedem Bild der 50er Klasse im Testdatensatz einen Sticker ein und messen, wie groß der Anteil der 50er Schilder ist, die als 80er Schild klassifiziert werden.

CH- und Backdoor-Artefakte:

In [AWN⁺19] wird wie folgt zwischen Clever Hans- und Backdoor-Artefakten unterschieden. In beiden Fällen wird rechts oben im Bild ein grauer 3x3 Sticker eingefügt. Bei CH geschieht dies bei 25% der `airplane`-Klasse. Bei Backdoor-Artefakten werden 10% aller Bilder korruptiert. Im zweiten Fall wird das entsprechende Label abgeändert. Dies entspricht dann einem Standard- bzw. Clean-Label-Poisoning-Angriff. (Wie gut funktioniert der CLPA/CH hier ohne die Bilder vorher schlechter zu machen? TODO: Vergleich mit [TTM19]). In [AWN⁺19], Kapitel 2.1 wird auch auf die Methode der Spektralen Signatur [TLM18] eingegangen, die zur Detektion genutzt wird. Diese eignet sich wohl sehr gut für die Backdoor-Attacks, aber nur schlecht für die CH-Artefakte.

3.2 Label-konsistente Poisoning-Angriffe

Bei den vorherigen Standard-Angriffen war es der Fall, dass das Label und das entsprechende Bild nicht mehr zusammenpassen. Ein händisches Durchsuchen des Datensatz (wenn auch sehr aufwendig) könnte damit ebenfalls zur Detektion eines

Angriffs führen.

Eine deutlich schwieriger zu detektierende Art von Poisoning-Angriffen sind sogenannte Label-konsistente Angriffe, bei denen genau diese Schwachstelle eliminiert ist, d.h. Label und Bild passen wieder zueinander, während der Angriff noch immer erfolgreich funktioniert. Es ist das Ziel, ein Bild zunächst so zu modifizieren, dass es für das menschliche Auge noch immer zur entsprechenden Klasse gehört, für das Neuronale Netzwerk aber so schwierig zu klassifizieren ist, dass sich das Netzwerk eher auf den Auslöser anstatt auf das ursprüngliche Bild verlässt. Im Anschluss wird wieder ein Auslöser eingefügt.

In [TTM19] werden zwei Verfahren vorgestellt, die die Klassifikation einzelner Bilder erschweren. Das erste Verfahren besteht aus einer Einbettung in einen niedrig-dimensionalen Raum, das auch bei Autoencodern, etc. verwendet wird TO-DO.

Beim zweiten Verfahren wird ein sogenannte Projizierter Gradienten-Abstieg-Angriff durchgeführt.

Dabei wird ein Adversarialer Angriff in leicht abgewandelter Form genutzt, um das Netzwerk zu stören. Bei Adversarialen Angriffen wird ein Netzwerk im Unterschied zum Poisoning-Angriff, bei dem der Angriff während des Trainings stattfindet, nach dem Training angegriffen. Dazu wird eine natürliche Netzwerkeingabe leicht gestört, sodass diese vom Netzwerk falsch klassifiziert wird. Diese Störungen lassen sich auch von einer Architektur oder sogar einem Modell auf andere übertragen [SZS⁺13, PMG16]. Für diese Art von Angriff werden die adversarialen Angriffe und ihre leichte Übertragbarkeit auf andere Architekturen und Modelle so benutzt, dass es bereits während des Trainings zu falschen Klassifikationen kommt.

Für unser erstes trainiertes Netzwerk f_θ mit Verlustfunktion \mathcal{L} und einem Eingabepaar (x, y) , konstruieren wir die modifizierte Version von x als

$$x_{adv} = \arg \max_{||x' - x||_p \leq \varepsilon} \mathcal{L}(x', y, \theta), \quad (3.1)$$

für $p > 1$ und $\varepsilon > 0$. Dieses Optimierungsproblem wird mit einem Projizierten Gradienten-Verfahren [MMS⁺17]. Details dazu finden sich in Anhang D. Im Unterschied zu [TTM19] ändern wir nur Bilder im Datensatz ab und fügen nicht zusätzlich zum Original x auch x_{adv} hinzu. Damit ändert sich die Anzahl an Datenpunkten durch den Angriff nicht.

Für das anschließende Einfügen des Auslöser ergeben sich die folgenden Optionen:

- Der im Standard-Angriff verwendete Sticker
- Ein Amplitudensticker: Dabei wird im rechten unteren Eck des Bildes ein
- Amplitudensticker 4 fach
- In die Mitte verschobene Amplitudensticker

RGB auf jedem Kanal in der range von [0,255] 0 entspricht weiß, 255=schwarz Da die zweite Möglichkeit als deutlich erfolgreicher angegeben wird, beschränken wir uns auf diese Angriffe basierend auf einem Projizierten Gradienten-Abstieg. amp=16,32,64,255 Wir gehen zunächst davon aus, dass der Angreifer volle Kontrolle über den Datensatz und das Netzwerk besitzt. TODO: Angriff mit einem andere Netzwerk erstellen, als das angegriffene

3.3 Bewertung von Poisoning-Angriffen

Wann ist ein Angriff erfolgreich?

Im Trainingsdatensatz werden im Fall des Standard-Angriffs alle Bilder mit dem Sticker versehen. Die Angriffserfolgsrate beschreibt nun den Anteil an Bildern der attackierten Klasse, die erfolgreich falsch klassifiziert wurden.

Für die Label-konsistenten Angriffe werden im Test-Datensatz alle Bilder mit dem entsprechenden Auslöser versehen und es kann eine Erfolgsrate pro Klasse berechnet werden. Es ist zu beachten, dass für Angriffe, die mit einer reduzierten Amplitudenstärke durchgeführt werden, die Bilder im Test-Datensatz dennoch mit Auslösern mit voller Amplitudenstärke versehen werden.

3.4 Verteidigungen

In diesem Kapitel beschäftigen wir uns mit gängigen Methoden zur Detektion von Poisoning-Attacks und geben am Ende einen kurzen Ausblick auf die Idee für einen neuen Ansatz. Wir wollen beide Arten von Poisoning-Angriffen erfolgreich detektieren.

3.4.1 Referenzwert: kMeans(k=2)

Der einfachste Ansatz, um korruptierte Datenpunkte zu erkennen, ist ein kMeans-Clustering, das direkt(bzw. nach einer Dimensionsreduktion) auf den Eingabedaten einer Klasse durchgeführt wird. Hierbei war auffällig, dass der Großteil der Daten als korruptiert klassifiziert wurde. Für die Dimensionsreduktionen FastICA und PCA ergab sich eine Genauigkeit von etwa 66%. Die FPR lag bei über 70 %, die TPR bei ca. 50%. Dieser Referenzwert w@articleszegedy2013intriguing, title=Intriguing properties of neural networks, author=Szegedy, Christian and Zaremba, Wojciech and Sutskever, Ilya and Bruna, Joan and Erhan, Dumitru and Goodfellow, Ian and Fergus, Rob, journal=arXiv preprint arXiv:1312.6199, year=2013 wurde für einen Sticker mit Seitenlänge 3 und 15% korruptierten Daten durchgeführt.

3.4.2 Activation Clustering

Nehme Datensatz her Beim Standardangriff wollen wir Klasse 5 als Klasse 8 klassifizieren und fügen dazu Sticker der

Diese Idee der Verteidigung basiert auf der Annahme, dass bestimmte Schichten innerhalb des Netzwerkes die Entscheidung, dass ein Bild mit einem Auslöser falsch klassifiziert wird, sehr gut codieren. Für die Detektion der Hintertüren im Datensatz sollen nun genau diese Aktivierungen für ein Clustering herangezogen werden. Das Activation Clustering wird erstmalig in [CCB⁺18] vorgestellt und nutzt aufgrund experimenteller Untersuchungen stets die Aktivierungen der vorletzten Netzwerkschicht. Eine Kombination von Aktivierungen mehrere Schichten wäre ebenfalls denkbar.

Ein Angriff ist erfolgreich, wenn eine große Anzahl an Datenpunkten der Ursprungs-kategorie, versehen mit einem Auslöser, der Zielklasse zugeordnet werden. Im Falle eines erfolgreichen Angriffs werden korruptierte und nicht korruptierte Datenpunkte im Trainingsdatensatz derselben Klasse zugeordnet. Der Grund weshalb

diese derselben Klasse zugeordnet werden, unterscheidet sich jedoch. Beim Activation Clustering wird nun angenommen, dass pro Klasse entweder korruptierte und nicht korruptierte Datenpunkte oder nur nicht korruptierte Datenpunkte existieren. Deshalb werden die Aktivierungen der letzten verdeckten Schicht des Netzwerkes aus dem Netz extrahiert, nach ihren zugehörigen Klassen der Labels segmentiert, auf 10 Dimensionen reduziert und anschließend mit Hilfe des kMeans-Algorithmus geclustert. Das kleinere Cluster wird immer als der Anteil an verdächtigen Datenpunkten betrachtet. Die Idee ist es, dass die korruptierten Datenpunkte, sofern welche existieren, alle in die eine und die nicht korruptierten Datenpunkte in das andere Cluster aufgeteilt werden. Sind keine korruptierten Datenpunkte vorhanden, so sollen beide Cluster ungefähr dieselbe Anzahl an Datenpunkten erhalten.

Wir werten die Qualität des Clusterings anschließend aus. Als Detektionsrate beschreiben wir die Genauigkeit des Clusterings auf den Trainingsdaten.

Im folgenden Abschnitt werden Methoden vorgestellt, mit denen das resultierende Clustering auf die Präsenz eines Angriffs untersucht werden kann. Dies ist notwendig, da in der Praxis ein Angriff zunächst erkannt und anschließend die korruptierten Datenpunkte entfernt werden müssen.

Bemerkung 3.4.1. *Das SPA benutzt die Idee das innerhalb einer Klasse bezüglich unterschiedlicher Aktivierungen klassifiziert werden kann. Beim CLPA funktioniert das nicht mehr, denn: Hier passen jetzt auch die Aktivierungen der korruptierten Bilder zur entsprechenden/untersuchten Klasse. Es ist also zu erwarten, dass das AC für CLPA nicht funktioniert.*

3.5 Implementierung

Wir nutzen die python Implementierung in sklearn mit den Standard-Werten $n_init = 10$ und $max_iter = 300$

3.5.1 Methoden zur Untersuchung, ob ein Angriff vorliegt

#TODO: Vergleich mit Fisher-Discriminant-Analyse/Ansatz in [AMN⁺19] Zur Bestimmung, ob eine Klasse korruptierte Daten enthält, kann das Ergebnis des Clusterings mit den folgenden Methoden untersucht werden:

@articleszegedy2013intriguing, title=Intriguing properties of neural networks, author=Szegedy, Christian and Zaremba, Wojciech and Sutskever, Ilya and Bruna, Joan and Erhan, Dumitru and Goodfellow, Ian and Fergus, Rob, journal=arXiv preprint arXiv:1312.6199, year=2013 Vergleich der relativen Größe: Eine Möglichkeit, korruptierte Datenpunkte zu erkennen, ist der Vergleich der relativen Größen der beiden Cluster. Laut [2] ist die relative Größe bei nicht korruptierten Klassen ca. 50 Prozent, bei korruptierten Daten und einem erfolgreichen Clustering würde die relative Größe dann dem prozentualen Anteil an korruptierten Datenpunkten entsprechen.

Silhouette-Koeffizient: Eine weitere Möglichkeit besteht darin, die Qualität des Clusterings mit Hilfe des Silhouette-Koeffizienten zu beschreiben. Dieser gibt an, wie gut ein Clustering zu den gegebenen Datenpunkten mit den entsprechenden

Labeln passt und ist wie folgt definiert: Sei das Ergebnis eines Clustering-Algorithmus mit verschiedenen Clustern gegeben. Zu einer Beobachtung x im Cluster A wir die Silhouette $s(x) = \frac{d(B,x) - d(A,x)}{\max\{d(A,x), d(B,x)\}}$ definiert, wobei $d(A,x) = \frac{1}{n_A - 1} \sum_{a \in A, a \neq x} d(a,x)$ dem mittleren Abstand einer Beobachtung innerhalb einer Klasse zu allen anderen Beobachtungen dieser Klasse entspricht. Dabei steht n_A für die Anzahl der Beobachtungen in Cluster A . $d(B,x) = \min_{C \neq A} d(C,x)$ beschreibt die Distanz von x zum nächstgelegenen Cluster B . Der Silhouetten-Koeffizient SC ist nun definiert als

$$SC = \max_k \bar{s}(k), \quad (3.2)$$

wobei $\bar{s}(k)$ der Mittelwert der Silhouetten aller Datenpunkte im gesamten Datensatz ist. Damit ist der Silhouettenkoeffizient ein Maß dafür, wie gut ein Clustering für eine vorher fixierte Clusteranzahl k zum Datensatz passt.

Exklusives Retraining: Beim exklusiven Retraining wird das neuronale Netz von Grund auf neu trainiert. Das oder die verdächtigen Cluster werden beim erneuten Training nicht benutzt. Mit Hilfe des neu trainierten Netzes werden dann anschließend die vorenthaltenen, verdächtigen Cluster klassifiziert. Falls das Cluster Aktivierungen von Datenpunkten enthält, die zum Label des Datenpunktes gehören, erwarten wir, dass die Vorhersage des Netzwerks mit dem Label übereinstimmen. Gehören die Aktivierungen eines Datenpunktes im verdächtigen Cluster jedoch zu einer anderen Klasse als die durch das Label angedeutete Klasse, so sollte das Netzwerk den Datenpunkt einer anderen Klasse zuordnen. Um nun zu entscheiden, ob ein verdächtiges Cluster korruptiert oder nicht korruptiert ist, wird wie folgt vorgegangen: Sei l die Anzahl an Vorhersagen, die zum Label des Datenpunktes passen. Sei p die größte Anzahl an Vorhersagen, die für eine weitere Klasse C sprechen, wobei C nicht die Klasse mit den Labels des zu untersuchenden Clusters ist. Der Quotient $\frac{l}{p}$ gibt dann an, ob das Cluster korruptiert ist oder nicht: Es wird ein Schwellenwert $T > 0$ gesetzt. Gilt $\frac{l}{p} < T$, wurden mehr Datenpunkte einer anderen Klasse zugeordnet und das Cluster wird als korruptiert deklariert. Umgekehrt wird das verdächtige Cluster im Fall von $\frac{l}{p} > T$ als nicht korruptiert/sauber eingestuft. @articleszegedy2013intriguing, title=Intriguing properties of neural networks, author=Szegedy, Christian and Zaremba, Wojciech and Sutskever, Ilya and Bruna, Joan and Erhan, Dumitru and Goodfellow, Ian and Fergus, Rob, journal=arXiv preprint arXiv:1312.6199, year=2013

3.5.2 Entfernen von korruptierten Datenpunkten

AC für Label-konistente Poisoning-Angriffe: Warum funktioniert Activation-Clustering hier nur schlecht oder gar nicht?: Wenn wir einen korruptierten Trainingsdatensatz gegeben haben, gilt im Fall des Standard-Angriffs folgender Sachverhalt: Die angegriffene Klasse, die Klasse in der samples eingefügt wurden, besitzt die eine Gruppe an Bildern, die zu einer Aktivierung von einer anderen Klasse führen sollten, und die Gruppe an Bildern, die zu dieser Klasse gehören und zur Aktivierung genau dieser Klasse führen sollte.

Im Fall des Label-konsistenten Poisoning-Angriffs, werden nun keine Label mehr getauscht, d.h. Bilder von der einen in die andere Klasse verschoben. Damit können die beiden Gruppen (korruptiert, sauber) innerhalb einer Klasse nicht mehr anhand ihrer Aktivierungen unterschieden werden.

Trotzdem ergibt sich ein Ansatz daraus, dass es innerhalb dieser Klasse verschiedene „Strategien“ gibt, die zur selben Klassifikation führen. Mithilfe des kmeans-Clustering basierend auf den Heatmaps sollen genau diese Strategien ausfindig gemacht werden, um die Bilder in korruptiert und sauber zu unterteilen.

3.5.3 Heatmap Clustering

Im Unterschied zum Activation Clustering, bei dem die Aktivierungen der vorletzten Netzwerk-Schicht verwendet werden, ist nun hier die Idee, zu jedem Eingabebild eine Relevanzkarte zu erstellen, die für jeden Pixelpunkt angibt, wie wichtig dieser für die Klassifikation dieses Bildes ist.

Für das Erstellen/Berechnen solcher Relevanzkarten/Heatmaps existieren mehrere Methoden, die zusammengefasst dem Bereich der Erklärbaren KI zugeordnet werden. Im folgenden Kaptiel wollen wir einen kurzen Überblick über verschiedene Methoden geben.

4 Erklärbare KI

Erklärbarkeit vs. Interpretierbarkeit, youtube talk?!

4.1 Lokale Methoden

4.2 Globale Methoden

5 Layer-wise Relevance Propagation

5.1 Idee

Die Layer-wise Relevance Propagation (LRP) wird in [BBM⁺15] erstmalig vorgestellt. Die Idee besteht darin, einen Zusammenhang zwischen der Ausgabe eines Klassifikators $f_\theta : \mathbb{R}^d \rightarrow \mathbb{R}^+$ und der Eingabe x herzustellen. Dabei wird eine definiert, die über gewisse Eigenschaften eingeschränkt wird. Die Autoren bezeichnen die Herangehensweise hier selbst als heuristisch und liefern in ?? eine Verallgemeinerung des Konzepts, die gleichzeitig die mathematische Grundlage bildet.

Wir betrachten eine nicht-negative Funktion $f : \mathbb{R}^d \rightarrow \mathbb{R}^+$. Im Bereich der Bild-Klassifizierung ist die Eingabe $x \in \mathbb{R}^d$ ein Bild, das wir als Menge von Pixelwerten $x = \{x_p\}$ auffassen können. Dabei beschreibt der Index p einen genauen Pixelpunkt. Während für schwarz-weiß Bilder $x_p \in \mathbb{R}$ gilt, gilt im Fall von RGB-Bildern $x_p \in \mathbb{R}^3$ für die einzelnen Farbkanäle Rot, Grün und Blau. Die Funktion $f(x)$ ist ein Maß dafür, wie präsent ein oder mehrere Objekte in der Eingabe/im Eingabebild vorhanden sind. Ein Funktionswert $f(x) = 0$ beschreibt die Abwesenheit. Gilt andererseits $f(x) > 0$, so wird die Präsenz mit einem gewissen Grad an Sicherheit oder eine gewisse Menge zum Ausdruck gebracht.

Mit Hilfe der LRP soll nun jedem Pixel p im Eingabebild eine Relevanz $R_p(x)$ zugeordnet werden, die für jedes Pixel x_p angibt, mit welcher Größe es für das

Entstehen einer Entscheidung $f(x)$ verantwortlich ist. Die Relevanz eines jeden Pixels wird dabei in einer Heatmap $R(x) = \{R_p(x)\}$ zusammengefasst.

Die Heatmap besitzt dieselbe Größe wie x und kann als Bild visualisiert werden.

Wir definieren die folgenden Eigenschaften:

Definition 5.1.1. Eine Heatmap $R(x)$ heißt konservativ, falls gilt:

$$\forall x : f(x) = \sum_p R_p(x), \quad (5.1)$$

d.h. die Summe der im Pixelraum zugeordneten Relevanz entspricht der durch das Modell erkannten Relevanz.

Definition 5.1.2. Eine Heatmap $R(x)$ heißt positiv, falls gilt:

$$\forall x, p : R_p(x) \geq 0, \quad (5.2)$$

d.h. alle einzelnen Relevanzen einer Heatmap sind nicht-negativ.

Die erste Eigenschaft verlangt, dass die umverteilte Gesamtrelevanz der Relevanz entspricht, mit der ein Objekt im Eingabebild durch die Funktion $f(x)$ erkannt wurde. Die zweite Eigenschaft beschreibt, dass keine zwei Pixel eine gegensätzliche Aussage über die Existenz eines Objektes treffen können. Beide Definitionen zusammen ergeben die Definition einer *konsistenten* Heatmap:

Definition 5.1.3. Eine Heatmap $R(x)$ heißt konsistent, falls sie konservativ und positiv ist, d.h. Definition 5.1.1 und Definition 5.1.2 gelten.

Für eine konsistente Heatmap gilt dann $(f(x) = 0 \Rightarrow R(x) = 0)$, d.h. die Abwesenheit eines Objektes hat zwangsläufig auch die Abwesenheit jeglicher Relevanz in der Eingabe zur Folge, eine Kompensation durch positive und negative Relevanzen ist folglich nicht möglich.

Bemerkung 5.1.4. Die geforderten Eigenschaften an eine Heatmap definieren diese nicht eindeutig. Es sind also mehrere Abbildungen möglich, die die genannten Forderungen erfüllen. Beispiele dafür sind eine natürliche Zerlegung und Taylor-Zerlegungen [MLB⁺ 17a].

Die LRP liefert nun ein Konzept, mit dem eine Zerlegung

$$f(x) = \sum_d R_d \quad (5.3)$$

bestimmt werden kann.

TODO: Summenabfolge von layer zu layer einfügen

Wir gehen nun davon aus, dass die Funktion f ein NN repräsentiert, dass aus mehreren Schichten mit mehreren Neuronen pro Schicht und dazwischengeschalteten nicht-linearen Aktivierungsfunktionen aufgebaut ist. Die erste Schicht ist die Eingabe-Schicht, bestehend aus den Pixeln eines Bildes. Die letzte Schicht ist die reellwertige Ausgabe von f . Die l -te Schicht ist durch einen Vektor $z = (z_d^l)_{d=1}^{V(l)}$ der Dimension $V(l)$ dargestellt. Sei also eine Relevanz $R_d(l+1)$ für jede Dimension $z_d^{(l+1)}$ des Vektors z in der Schicht $l+1$ gegeben. Die Idee besteht nun darin, eine

Relevanz $R_d^{(l)}$ für jede Dimension $z_d^{(l)}$ des Vektors z in der Schicht l zu finden, die einen Schritt näher an der Eingabeschicht liegt, sodass die folgende Abfolge von Gleichungen gilt:

$$f(x) = \dots = \sum_{d \in l+1} R_d^{(l+1)} = \sum_{d \in l} R_d^{(l)} = \dots = \sum_d R_d^{(1)}. \quad (5.4)$$

Für diese Funktion benötigen wir eine Regel, mit der die Relevanz eines Neurons einer höheren Schicht $R_j^{(l+1)}$ auf ein Neuron einer benachbarten, näher an der Eingabeschicht liegendes Neuron, übertragen werden kann. Die Übertragung der Relevanz zwischen zwei solchen Neuronen wird mit $R_{i \leftarrow j}$ bezeichnet. Auch hier muss die übertragene Relevanz erhalten bleiben. Es wird also gefordert:

$$\sum_i R_{i \leftarrow j}^{(l,l+1)} = R_j^{(l+1)}. \quad (5.5)$$

D.h. die gesamte Relevanz eines Neurons der Schicht $l+1$ verteilt sich komplett auf alle Neuronen der Schicht l . Im Falle eines linearen NN $f(x) = \sum_i z_{ij}$ mit der Relevanz $R_j = f(x)$ ist eine Zerlegung gegeben durch $R_{i \leftarrow j} = z_{ij}$. Im allgemeineren Fall ist die Neuronenaktivierung x_j eine nicht-lineare Funktion abhängig von z_j . Für die beiden Aktivierungsfunktionen $\tanh(x)$ und $\text{ReLU}(x)$ - beide monoton wachsend mit $g(0) = 0$ - bieten die Vor-Aktivierungen noch immer ein sinnvolles Maß für den relativen Beitrag eines Neurons x_i zu R_j (müsste das nicht umgekehrt sein, die INdizes?!?).

Eine erste Mögliche Relevanz-Zerlegung, basierend auf dem Verhältnis zwischen lokalen und globalen Vor-Aktivierung, ist gegeben durch:

$$R_{i \leftarrow j}^{(l,l+1)} = \frac{z_{ij}}{z_j} \cdot R_j^{(l+1)}. \quad (5.6)$$

Für diese Relevanzen $R_{i \leftarrow j}$ gilt die Erhaltungseigenschaft 5.4, denn:

$$\sum_i R_{i \leftarrow j}^{(l,l+1)} = R_j^{l+1} \cdot \left(1 - \frac{b_j}{z_j}\right). \quad (5.7)$$

Dabei steht der rechte Faktor für die Relevanz, die durch den Bias-Term absorbiert wird. Falls notwendig, kann die verbleibende Bias-relevanz auf jedes Neuron x_i verteilt werden(?s.Abschnitt über Biases Promotion, S.Lapuschkin).

Diese Regel wird in der Liteartur als LRP-0 bezeichnet. Ein Nachteil dieser ist, dass die Relevanzen $R_{i \leftarrow j}$ für kleine globale Voraktivierung z_j beliebig große Werte annehmen können.

Um dies zu verhindern, wird in der LRP- ε -Regel ein vorher festgelegter Parameter $\varepsilon > 0$ eingeführt:

$$R_{i \leftarrow j}^{(l,l+1)} = \begin{cases} \frac{z_{ij}}{z_j + \varepsilon} \cdot R_j^{(l+1)}, & z_j \geq 0 \\ \frac{z_{ij}}{z_j - \varepsilon} \cdot R_j^{(l+1)}, & z_j < 0 \end{cases} \quad (5.8)$$

In [BBM⁺15] wird die Layer-wise Relevance Propagation erstmalig vorgestellt. Zudem wird eine Taylor Zerlegung präsentiert, die eine Approximation der LRP darstellt.

Hier⁵ werden einige Bereiche vorgestllt, in denen LRP angewendet wurde.

⁵<https://towardsdatascience.com/indepth-layer-wise-relevance-propagation-340f95deb1ea>

5.2 Behandlung von biases

5.3 Beispiel an einem kleinen Netzwerk

5.4 Deep Taylor Decomposition

Mathematischer Hintergrund für LRP. LRP als Spezialfall von DTD

5.4.1 Taylor Decomposition

We will assume that the function $f(x)$ is implemented by a deep neural network, composed of multiple layers of representation, where each layer is composed of a set of neurons. Each neuron performs on its input an elementary computation consisting of a linear projection followed by a nonlinear activation function. Deep neural networks derive their high representational power from the interconnection of a large number of these neurons, each of them, realizing a small distinct subfunction.

Laut [MLB⁺17b] ist die in [BBM⁺15] vorgestellte Layer-wise Relevance Propagation eher heuristisch. In diesem Paper wird nun eine solide theoretische Grundlage geliefert.

DTD liefert den mathematischen Hintergrund für LRP

Simple Taylor decomposition. Finde rootpoints, sodass Erhaltungseigenschaft erhalten bleibt.

Simple Taylor in Practice: funktioniert in der Praxis nicht wirklich. Viel Noise meistens positive Relevanz

Relevanz Propagation: Heatmaps look much cleaner

Simple Taylor:- root point hard to find -gradient shattering. Gradient loses its informative structure in big layer nets

Use Taylor Decomposition to explain LRP from layer to layer

5.4.2 Deep Taylor Decomposition

LRP in verschiedenen Anwendungsgebieten [MBL⁺19], 10.2. In diesem Paper: LRP-0 schlechter als LRP- ε schlechter als LRP- γ schlechter als Composite-LRP.

5.5 Verschiedene Verfahren

5.6 Eigenschaften

Beweise in DTD Paper

- Numerische Stabilität
- Konsistenz (mit Linearer Abbildung)
- Erhaltung der Relevanz



Abbildung 1: (Optischer) Vergleich von korrumpiertem Datenpunkt und berechneter Heatmap. Links: Verkehrsschild der Klasse 'Höchstgeschwindigkeit: 50km/h' versehen mit einem 3x3 Sticker und dem Label 'Höchstgeschwindigkeit: 80km/h'. Rechts: Zugehörige Heatmap bezüglich der Klasse 'Höchstgeschwindigkeit: 80km/h'.

Quelle: [CCB⁺18]

5.7 Behandlung besonderer Schichten

5.7.1 BatchNorm2D

5.8 LRP für Deep Neural Nets/Composite LRP

5.9 Verarbeitung der Heatmaps

Aktuell benutzte default colormap ist Option D (Viridis)⁶

- Wertebereich
- Interpretation
- Skalen
- Normalisierung

5.10 Implementierungen

5.10.1 Tensorflow

5.10.2 pytorch

Allgemeines Tutorial:⁷

pytorch-LRP für VGG16 wird vorgestellt.

GiorgioML⁸:

Alternative pytorch-Implementierung basierend auf Tensorflow paper.

⁶<https://bids.github.io/colormap/>

⁷<https://git.tu-berlin.de/gmontavon/lrp-tutorial>

⁸<https://giorgiomorales.github.io/Layer-wise-Relevance-Propagation-in-Pytorch/>

moboehle⁹:

Der code entstand im Rahmen der Forschungsarbeit [BEWR19], in der eine Alzheimer-Feststellung aufgrund von Bilddaten(scans?) vorgenommen wird. Framework leicht anpassbar. Benutzt pytorch hooks. Unterstützte Netzwerkschichten¹⁰:

```

1 torch.nn.BatchNorm1d,
2 torch.nn.BatchNorm2d
3 torch.nn.BatchNorm3d,
4 torch.nn.ReLU,
5 torch.nn.ELU,
6 Flatten,
7 torch.nn.Dropout,
8 torch.nn.Dropout2d,
9 torch.nn.Dropout3d,
10 torch.nn.Softmax,
11 torch.nn.LogSoftmax,
12 torch.nn.Sigmoid
13
14
```

Listing 2: Verfügbare Schichten und Aktivierungsfunktionen

fhvilshoj¹¹:

LRP für linear und Convolutional layers

- Die Klassen
torch.nn.Sequential, torch.nn.Linear und torch.nn.Conv2d werden erweitert, um autograd für die Berechnung der Relevanzen zu berechnen.
- Ausgabe der Relevanzen von Zwischenschichten ist möglich
- : Implementierte Regeln: epsilon Regeln mit epsilon=1e-1, gamma-regel mit gamma=1e-1. alphabeta-Reagel mit a1b0 und a2b1
- Netz muss hier umgeschrieben werden, sodass die Anwendung des Algorithmus möglich wird.

```

1
2 conv2d = {
3     "gradient":          F.conv2d,
4     "epsilon":           Conv2DEpsilon.apply,
5     "gamma":             Conv2DGamma.apply,
6     "gamma+epsilon":     Conv2DGammaEpsilon.apply,
7     "alpha1beta0":       Conv2DAlpha1Beta0.apply,
8     "alpha2beta1":       Conv2DAlpha2Beta1.apply,
9     "patternattribution": Conv2DPatternAttribution.apply,
10    "patternnet":         Conv2DPatternNet.apply,
11 }
12
```

⁹<https://github.com/moboehle/Pytorch-LRP>

¹⁰https://github.com/moboehle/Pytorch-LRP/blob/master/inverter_util.py

¹¹<https://github.com/fhvilshoj/TorchLRP>

13

Listing 3: Implementierte Regeln fhvilshoj**Zennit:**¹² Zennit (Zennit explains neural networks in torch)

- Modell wird mithilfe eines Canonizers so aufbereitet, dass LRP möglich wird
- Backward pass wird modifiziert, um Heatmaps zu erhalten.
- VGG- und ResNet-Beispiel

6 Detektion von Poisoning-Angriffen basierend auf LRP

Super Einführung in Wasserstein und OT: [AWR17]

Computing distances between probability measures on metric spaces, or more generally between point clouds, plays an increasingly preponderant role in machine learning [SL11, MJ15, LG15, JSCG16, ACB17], statistics [FCCR16, PZ16, SR04, BGKL17] and computer vision [RTG00, BvdPPH11, SdGP+15]. A prominent example of such distances is the earth mover’s distance introduced in [WPR85] (see also [RTG00]), which is a special case of Wasserstein distance, or optimal transport (OT) distance [Vil09]. While OT distances exhibit a unique ability to capture geometric features of the objects at hand, they suffer from a heavy computational cost that had been prohibitive in large scale applications until the recent introduction to the machine learning community of Sinkhorn Distances by Cuturi [Cut13]. Combined with other numerical tricks, these recent advances have enabled the treatment of large point clouds in computer graphics such as triangle meshes [SdGP+15] and high-resolution neuroimaging data [GPC15]. Sinkhorn Distances rely on the idea of entropic penalization, which has been implemented in similar problems at least since Schrödinger [Sch31, Leo14]. This powerful idea has been successfully applied to a variety of contexts not only as a statistical tool for model

How Well Do WGANs Estimate the Wasserstein Metric? [MMG19]

6.1 Idee

Die Idee zur Detektion von Poisoning-Angriffen besteht aus den folgenden Schritten:

- Berechnung der Heatmaps mit Hilfe der LRP
- Berechnung einer Distanzmatrix basierend auf L^2 - oder GMW-Distanz
- Spektrale Relevanzanalyse (Bestimmung der verschiedenen Cluster innerhalb einer Klasse)

Bemerkung: Anstatt das Clustering nur auf den Heatmaps durchzuführen, könnten die LRP-Ausgaben und/oder Aktivierungen bestimmter Netzwerkschichten hinzugenommen werden.

The theory of optimal transport generalizes that intuition in the case where, instead of moving only one item at a time, one is concerned with the problem of

¹²<https://github.com/chr5tphr/zennit>

moving simultaneously several items (or a continuous distribution thereof) from one configuration onto another. [com19]

6.2 k-means / k-means++ -Clustering

Beispiel-Implementierung¹³

Baryzentrische Koordinaten¹⁴

6.3 Spektrales Clustering

Wir folgen [VL07]. Gegeben: Datenpunkte x_i, \dots, x_n sowie eine Größe $s = s_{ij} \in \mathbb{R}^+$, die einen paarweisen Zusammenhang der einzelnen Punkte beschreiben.

Ziel: Aufteilen der Punkte in verschiedene Cluster, sodass sich Punkte innerhalb eines Clusters ähnlich bezüglich s sind.

Alternative Repräsentation der Daten mithilfe eines Ähnlichkeitsgraphen $G = (V, E)$ möglich.

Umformulierung des Clustering-Problems mithilfe des Ähnlichkeitsgraphen: Finde Partitionierung des Graphen, sodass die Kanten-Gewichte innerhalb einer Gruppe niedrig (niedriges Gesamtgewicht?) und außerhalb einer Gruppe groß sind.

Graph-Notationen:

Verschiedene Konstruktionsmöglichkeiten von Ähnlichkeitsgraphen:

- ε -Nachbarschaft-Graph
- kNN-Graph
- fully connected graph

6.4 Anwendung auf unterschiedliche Poisoning-Angriffe

Berechnung der Relevanzen:

Wir berechnen die Relevanzen jedes einzelnen Eingabebildes klassenweise, d.h. besitzt eine Eingabe das Label y , so berechnen auf einem trainierten Netzwerk, für jeden Pixelwert der Eingabe, wie relevant dieser für die Ausgabe $f(x) = y$ ist.

Wir summieren über die Farbxen des Bildes, um einzelne Relevanzen pro Pixelpunkt zu erhalten.

Für die Berechnung der Relevanzen benutzen wir eine modifizierte Version des im Rahmen von [BEWR19] entstandenen Programmcodes¹⁵.

Vorverarbeitung der Relevanzen:

In [LWB⁺19] wird anschließend ein Sum-Pooling auf die Relevanzen angewendet, um eine Dimensionsreduktion zu erhalten. Wie in [AMN⁺19] verzichten wir auf

¹³<https://towardsdatascience.com/k-means-implementation-in-python-and-spark-856e7eb5fe9b>

¹⁴https://de.wikipedia.org/wiki/Baryzentrische_Koordinaten

¹⁵<https://github.com/moboehle/Pytorch-LRP>

eine weitere Dimensionsreduktion, da wir nur relativ kleine Relevanzen der Größe 32×32 verarbeiten.

Für $\text{eps}=5e-2$ liegen beide barycentren identisch weit weg. Probiere nun $\text{eps}=5e-3$

Berechnung der Distanzen und Aufstellen einer Affinitätsmatrix:

Wir berechnen zunächst eine Distanzmatrix, die die paarweisen Distanzen aller Heatmaps einer Klasse enthält.

Für die Berechnung der euklidischen Distanz betrachten wir Heatmaps x, y der Größe 32×32 als Elemente $x, y \in \mathbb{R}^{32 \times 32}$. Die Distanz lässt sich dann wie in ?? berechnen.

Die Gromov-Wasserstein-Distanz lässt sich wie in [PCS16] angegeben berechnen.

Barycenters Definition und Vergleich zum euklidischen Raum [AC11]
In einer Affinitätsmatrix oder Ähnlichkeitsmatrix sind die

Berechnung Spektralen Einbettung:

Dimensionsreduktion vor dem Clustering ?!

In [CCB⁺18] wird beispielsweise eine Dimensionsreduktion mit PCA durchgeführt.

k-Means-Clustering:

Bemerkung 6.4.1. In [AMN⁺19] Kapitel '2.3. Fisher Discriminant Analysis for Clever Hans identification' wird ein Verfahren vorgestellt, mit dem verdächtige Klassen identifiziert werden können. Für diese würde man anschließend das obige Verfahren durchführen

6.5 Verwendete Distanzen & Approximationen

Um die Struktur innerhalb einer Klasse zu analysieren, benötigen wir eine Metrik. Anhand dieser wird abhängig von den Heatmaps einer Klasse eine Affinitätsmatrix berechnet, die dann anschließend zur Berechnung der Spektralen Einbettung als wichtigster Schritt von SpRAy verwendet wird. Wir wollen dazu die im Folgenden vorgestellten Metriken verwenden.

Wie in [AMN⁺19] summieren wir über die Farbkanäle, um einen einzelnen Relevanzwert pro Pixelpunkt zu erhalten. Wir benötigen also eine Metrik zur Berechnung der Distanz zwischen 32×32 großen Heatmaps.

Wir normalisieren die Relevanzen zusätzlich auf das Intervall $[0, 1]$.

Die Wahl der Pixel mit 99 Prozent der Gesamtmasse und anschließende Normalisierung wird vermutlich durchgeführt, um die Bedingung Gleichung 7.1 zu erhalten.

7 Optimal Transport

Optimal Transport is an interesting topic which connects many fields and has interesting applications Applications in image analysis Interplay between geometry, probability and PDEs Convexity, duality Numerical optimization Statistics and Machine Learning¹⁶

Optimal transport can deal with smooth and discrete measures and it has proved to be very useful for comparing distributions in a shared space, but with different (and even non-overlapping) supports [VCF⁺20]. Einführung OT Villani-OT hebt Distanzen von metrischen Räumen auf den Raum der metrischen Räume [Mém11].

In diesem Kapitel betrachten wir einige Konzepte aus dem Bereich Optimal Transport, um einen Distanzbegriff zu entwickeln, der im Unterschied zur pixelweisen euklidischen Distanz besser für das kMeans-Clustering geeignet sein könnte. Ausgehend von Gaspard Monge's ursprünglicher Transportprobleme-Formulierung betrachten wir die durch Kantorovich eingeführte Relaxierung dieses Problems. Die Optimale Lösung wird für die Definition der p-Wasserstein-Distanz verwendet. Für eine Verallgemeinerung auf verschiedene Grundräume wird die Gromov-Wasserstein-Distanz definiert. Damit können wir zwei Heatmaps als sogenannte metrische Maßräume auffassen. Abschließend definieren wir sogenannte Gromov-Wasserstein-Baryzentren, die im Rahmen des kMeans-Clustering als Mittelwerte fungieren.

Stimmt Gromov-Wasserstein mit Wasserstein für dieselbe Distanzfunktion (Metrik) überein? Welche Rolle spielt die Lossfunktion L bei der Definition? $L=KL$

7.1 Frage:

Was ist ein registration problem? Warum brauchen wir Gromov-Wasserstein? Warum reicht nicht Wasserstein? Weil wir Matrizen unterschiedlicher Größen vergleichen? Nein. Wir benutzen ja immer dieselbe Distanz für alle Pixel-Werte. Wenn wir eine Pixel-Wolke auswählen, sind das nicht dieselben Räume, da dann gewisse Koordinaten nicht zum Raum gehören. Wenn wir aber alle Pixel als Koordinaten im Raum auswählen würden, dann würden die metrischen Räume (X, d_X) und (Y, d_Y) übereinstimmen. Liegt das an der Definition der Baryzentren? Wie würde man das mit Baryzentren auf demselben Raum machen? s. [COT19] Translationen und Rotationen könnte ein wichtiger Punkt sein!!! Verwendung von Grafiken aus anderen Papern?? Kann man ein anderes Netzwerk zur Detektion nutzen?

<https://arxiv.org/pdf/1705.09634.pdf> <https://arxiv.org/pdf/1412.5154.pdf> Für die Grundlagen im Bereich des Optimalen Transports zitieren wir [?] Villani2009optimal.

Entscheidend für die Numerischen Resultate ist die Arbeit von Marco Cuturi [Cut13, com19]

7.2 TODO:

[COT19] Remark 2.19 (Translations). zeigt dass die Unabhängigkeit gegenüber Translationen? Kann ich das auch für die Gromov-Wasserstein-Distanz zeigen?

¹⁶https://indico.cern.ch/event/845380/attachments/1915103/3241592/Dvurechensky_lectures.pdf

7.2.1 Anwendungsgebiete

Optimal transport is also used widely in domain adaptation [CFTR16]. In domain adaptation, one has access to labeled examples from a source domain and unlabeled examples from a target domain. The goal is to predict labels for examples from the target domain. This is different from the typical train-test paradigm in machine learning, because covariate shift between the two domains is allowed.

Ein weitere Anwendung von Optimal Transport sind Wasserstein Generative Adversarial Networks (WGANs) [ACB17]. WGANs gehören zur Familie der Modelle namens Generative Adversarial Networks (GANs) [GPAM⁺14]. Zu einem GAN gehören zwei Modelle, von denen beide typischerweise Neuronale Netzwerke sind. Das erste Netzwerk (generator) erzeugt Datenpunkte innerhalb des Datensatzes, die so realistisch wie möglich sein sollen. Das zweite Netzwerk (discriminator network) versucht, zwischen den realen und den erzeugten Datenpunkten zu unterscheiden. Damit kann der Trainingsprozess von GANs als ein Spiel zwischen den beiden Netzwerken betrachtet werden, bei dem der Generator den Diskriminator täuschen und der Diskriminator die korrumpierten Datenpunkte aufdecken möchte.

Durch die Verwendung von WGANs anstatt GANs wurde ein stabilerer Trainingsprozess erreicht, da diese weniger anfällig für abrupte Trainingsabbrüche sind. Die ursprünglichen GANs benutzten eine Jensen-Shannon-Divergenz als Verlustfunktion, die nicht überall stetig und fast überall differenzierbar ist. Auch andere Verlustfunktion wie beispielsweise die KL-Divergenz weisen ähnliche Schwierigkeiten auf. Durch die Verwendung der Wasserstein-Verlustfunktion, basierend auf der 1-Wasserstein-Distanz, die überall stetig und fast überall differenzierbar ist, konnte der Trainingsprozess der GANs somit deutlich verbessert werden [Mau].

Definition 7.2.1.

Definition 7.2.2 (Histogramm). *Als Histogramm (oder: Wahrscheinlichkeitsvektor) bezeichnen wir ein Element $\mathbf{a} \in \Sigma_n$, das zum folgenden Wahrscheinlichkeits-Simplex gehört:*

$$\Sigma_n := \{\mathbf{a} \in \mathbb{R}_+^n \mid \sum_{i=1}^n \mathbf{a}_i = 1\}$$

Definition 7.2.3 ([Gro07]). *Ein metrischer Maßraum ist ein Tripel (X, d_X, μ_X) mit*

- (X, d_X) ist ein kompakter metrischer Raum
- μ_X ist ein Borel-Maß auf X mit vollem Support, d.h.

$$\text{supp}[\mu_X] = X. \tag{7.1}$$

Definition 7.2.4 (pushforward-Maß). *Seien X, Y zwei Maßräume, $T : X \rightarrow Y$ eine messbare Abbildung und μ ein Maß auf X . Dann ist das pushforward-Maß $T^\# \mu$ ein Maß auf Y , definiert durch*

$$T^\# \mu(A) = \mu(T^{-1}(A)) \tag{7.2}$$

für alle $A \subset Y$.

Einführung W-Theorie am KIT ¹⁷

Definition 7.2.5 (Kopplungen zwischen Histogrammen). *Seien die beiden Histogramme $p \in \Sigma_{N_1}$ und $q \in \Sigma_{N_2}$ gegeben. Als Menge der Kopplungen zwischen beiden Histogrammen definieren wir*

$$\mathcal{C}_{p,q} := \{T \in (\mathbb{R}_+)^{N_1 \times N_2} \mid T \mathbb{K}_{N_2} = p, T^\top \mathbb{K}_{N_1} = q\},$$

wobei $\mathbb{K} := (1, \dots, 1)^\top \in \mathbb{R}^N$ gilt.

Definition 7.2.6 ([BBB⁺01]). *Sei X eine beliebige Menge. Eine Funktion $d : X \times X \rightarrow \mathbb{R} \cup \{\infty\}$ heißt Metrik auf X , falls die folgenden Bedingungen für alle $x, y, z \in X$ gelten:*

- (1) *Positive Definitheit: $d(x, y) > 0$ für $x \neq y$ und $d(x, y) = 0$ für $x = y$,*
- (2) *Symmetrie: $d(x, y) = d(y, x)$,*
- (3) *Dreiecksungleichung: $d(x, z) \leq d(x, y) + d(y, z)$.*

Ein *metrischer Raum* ist eine Menge, versehen mit einer Metrik. Formal gesehen ist ein metrischer Raum ein Paar (X, d) , wobei d eine Metrik auf X ist. Elemente von X heißen Punkte und $d(x, y)$ bezeichnet die Distanz zwischen x und y .

Geodesic Space ¹⁸

Gromov-Hausdorff-Distanz

Definition 7.2.7 (Hausdorff-Distanz). *Sei (M, d) ein metrischer Raum. Für Seien $X, Y \subset (M, d)$ definieren wir die Hausdorff-Distanz $d_H(X, Y)$ als*

$$d_H(X, Y) = \max\left\{\sup_{x \in X} \inf_{y \in Y} d(x, y), \sup_{y \in Y} \inf_{x \in X} d(x, y)\right\}. \quad (7.3)$$

Definition 7.2.8 (Metrischer Maßraum). *Ein metrischer Maßraum ist ein Tripel (X, d_X, μ_X) , wobei (X, d_X) ein metrischer Raum und μ_X ein borelsches W-Maß auf X ist.*

Definition 7.2.9 (Correspondance). *content...*

Definition 7.2.10 (Kopplung). *Seien*

Seien (X, d_X) und (Y, d_Y) zwei metrische Räume. Wir betrachten im Folgenden die Abbildung

$$\Gamma_{X,Y} : (X \times Y) \times (X \times Y) \rightarrow \mathbb{R}^+, \quad (7.4)$$

gegeben durch

$$\Gamma_{X,Y}(x, y, x', y') := |d_X(x, x') - d_Y(y, y')|.$$

Definition 7.2.11 (Gromov-Hausdorff-Distanz). *Für die metrischen Räume (X, d_X) und (Y, d_Y) ist die Gromov-Hausdorff-Distanz definiert als*

$$d_{GH} = \frac{1}{2} \inf_R \|\Gamma_{X,Y}\|_{L^\infty(R \times R)}. \quad (7.5)$$

Definition 7.2.12 (Entropischer Optimaler Transport). *Wir definieren den n -dimensionalen Zufalls-Simplex als $\Sigma_n := \{a \in \mathbb{R}_+^n : \sum_{i=1}^n a_i = 1\}$. Ein Element $a \in \Sigma_n$ bezeichnen wir als Histogramm oder Zufallsvektor.*

Definition 7.2.13 (Diskretes Maß).

¹⁷<https://www.math.kit.edu/stoch/~henze/media/wt-ss15-henze-handout.pdf>

¹⁸<http://www.math.toronto.edu/mccann/papers/FiveLectures.pdf>

7.2.2 Optimaler Transport (Monge Formulierung)

Bemerkung 7.2.14. *Problem zwischen diskreten Maßen*

$$\min_T \left\{ \sum_i c(x_i, T(x_i)) : T_{\#}\alpha = \beta \right\} \quad (7.6)$$

Bemerkung 7.2.15 (Fehlende Eindeutigkeit).

Definition 7.2.16 (Push-forward Operator). $\beta(B) = \alpha(\{x \in \mathcal{X} : T(x) \in B\}) = \alpha(T^{-1}(B))$

Wir können das Monge Problem wie folgt auf den Fall zweier beliebiger W-Maße (α, β) erweitern.

Definition 7.2.17 (Monge Problem zwischen beliebigen Maßen). *Seien α und β zwei W-Maße mit Support auf den Räumen \mathcal{X} bzw. \mathcal{Y} , die durch $T : \mathcal{X} \rightarrow \mathcal{Y}$ verknüpft(? Formulierung) sind. Dann ist das Monge Problem gegeben durch*

$$\min_T \left\{ \int_{\mathcal{X}} c(x, T(x)) d\alpha(x) : T_{\#}\alpha = \beta \right\}. \quad (7.7)$$

Die Bedingung $T_{\#}\alpha = \beta$ bedeutet, dass die Abbildung T die gesamte Masse mithilfe des Push-forward Operators von α auf β schiebt.

7.2.3 Optimaler Transport nach Kantorovich

[San10] ist eine gute, einfache Einführung. Damit kann ich vermutlich Monge und Kantorovich beschreiben.

Das Optimal Transport Problem nach Kantorovich gehört zu den typischen Optimal Transport Problemen. Es stellt einer Relaxierung der Formulierung von Gaspard Monge[1781], bei der nun auch das Aufteilen von Masse (mass splitting) zulässig ist. In Kantorovichs Formulierung ist eine Kopplung (oder: Transportabbildung) \mathbf{T} gesucht, die die Kosten, die bei der Verschiebung eines diskreten Maßes \mathbf{a} auf ein anderes diskretes Maß \mathbf{b} bezüglich der Kosten $\mathbf{M} \in \mathbb{R}^{n_1 \times n_2}$ entstehen, minimiert. Damit \mathbf{T} eine Transportabbildung ist, muss $\mathbf{T} \in \Gamma(\mathbf{a}, \mathbf{b}) = \{\mathbf{T} \geq \mathbf{0}, \mathbf{T}\mathbf{1}_{n_2} = \mathbf{a}, \mathbf{T}^T \mathbf{1}_{n_1} = \mathbf{b}\}$ gelten.

Für den Fall, dass die Grundkosten eine Metrik darstellen, ist auch die optimale Lösung des Optimal Transport Problems wieder eine Metrik [?] und definiert die *Wasserstein Distanz*. Das OT Problem ist definiert als

$$W_M(\mathbf{a}, \mathbf{b}) = \min_{\mathbf{T} \in \Pi(\mathbf{a}, \mathbf{b})} \langle \mathbf{T}, \mathbf{M} \rangle, \quad (7.8)$$

wobei $\langle \mathbf{T}, \mathbf{M} \rangle = \sum_{ij} t_{ij} m_{ij}$ gilt.

Die Lösung des Problems ist eine Kopplung/Kopplungsmatrix, die beschreibt, wie viel Masse von einem Punkt zum anderen Punkt fließt. In einer Kostenmatrix derselbe Größe sind die Kosten abgespeichert, um von einem zum anderen Punkt zu kommen.

which is a linear program. The optimization problem above is often adapted to include a regularization term for the transport plan \mathbf{T} , such as entropic regularization (Cuturi, 2013) or squared L2. For the entropic regularized OT problem, one

may use the Sinkhorn Knopp algorithm (or variants), or stochastic optimization algorithms. POT has a simple syntax to solve these problems (see Sample 1)

Die Menge der Matrizen $\Pi(\mathbf{a}, \mathbf{b})$ ist beschränkt und durch $n + m$ Gleichungen gegeben und damit ein konvexes Polytop (die konvexe Hülle einer endlichen Menge von Matrizen). Zudem ist die Formulierung von Kantorovich im Unterschied zu Monges Formulierung immer symmetrisch in dem Sinne, dass $\mathbf{P} \in \Pi(\mathbf{a}, \mathbf{b})$ genau dann gilt, wenn $\mathbf{P}^\top \in \Pi(\mathbf{b}, \mathbf{a})$.

Kantorovichs Optimal Transport Problem lässt sich nun schreiben als

$$L_C(\mathbf{a}, \mathbf{b}) := \min_{\mathbf{P} \in \Pi(\mathbf{a}, \mathbf{b})} \langle \mathbf{C}, \mathbf{P} \rangle := \sum_{i,j} C_{i,j} P_{i,j} \quad (7.9)$$

Diese Problem ist ein lineares Problem. Für diese Art von Problemen ist die optimale Lösung nicht notwendigerweise eindeutig. In fact Kantorovich is considered as the inventor of linear programming.¹⁹

EXISTENZ & Eindeutigkeit => vgl. Kapitel 3 in [com19]

Definition 7.2.18 (Formulierung für beliebige Maße). *Im Fall beliebiger Maße betrachten wir die Kopplungen $\pi \in \mathcal{M}_+^1(\mathcal{X} \times \mathcal{Y})$, die die gemeinsame Verteilung auf dem Produktraum $\mathcal{X} \times \mathcal{Y}$ ist. Im diskreten Fall verlangen wir, dass das Produktmaß die Form $\pi = \sum_{i,j} P_{i,j} \delta_{(x_i, y_j)}$ besitzt. Im Allgemeinen Fall wird die Massenerhaltung als Randbedingung an die gemeinsame W-Verteilung geschrieben:*

$$\Pi(\alpha, \beta) := \{\pi \in \mathcal{M}_+^1(\mathcal{X} \times \mathcal{Y}) : P_{\mathcal{X}\#} \pi = \alpha \text{ und } P_{\mathcal{Y}\#} \pi = \beta\}, \quad (7.10)$$

wobei $P_{\mathcal{X}\#}$ und $P_{\mathcal{Y}\#}$ die Push-forward Operatoren der Projektionen $P_{\mathcal{X}}(x, y) = x$ und $P_{\mathcal{Y}}(x, y) = y$ sind. Nach Theorem 7.2.16 sind diese Randbedingungen äquivalent zu den Bedingungen $\pi(A \times \mathcal{Y}) = \alpha(A)$ und $\pi(\mathcal{X} \times B) = \beta(B)$ für die Mengen $A \subset \mathcal{X}$ und $B \subset \mathcal{Y}$. Als Verallgemeinerung erhalten wir dann

$$\min_{\pi \in \Pi(\alpha, \beta)} \int_{\mathcal{X} \times \mathcal{Y}} c(x, y) d\pi(x, y). \quad (7.11)$$

Problem ?? ist ein unendlich-dimensionales lineares Programm über einem Raum von Maßen. Falls $(\mathcal{X}, \mathcal{Y})$ kompakt und c stetig ist, existiert immer eine Lösung.

Beispiel 7.2.19. *Beispiele von Kopplungen*

7.2.4 Monge-Kantorovitch equivalence

The proof of Brenier theorem 1 (detailed in Section 5.3) to prove the existence of a Monge map actually studies Kantorovitch relaxation, and proves that this relaxation is tight in the sense that it has the same cost as Monge problem.²⁰

7.2.5 Metrische Eigenschaften

Optimaler Transport definiert eine Distanz zwischen Histogrammen und W-Maßen, sofern die Kostenmatrix gewisse Eigenschaften erfüllt. Optimal Transport kann dabei als naheliegende Idee verstanden werden, um Distanzen zwischen Punkten auf Distanzen zwischen Histogrammen oder Maßen zu verallgemeinern.

¹⁹OTNotes_campide.pdf

²⁰CourseOT_cuturi

Definition 7.2.20 (p -Wasserstein-Distanz auf Σ_n). Sei $n = m$ und für $p \geq 1$ gelte $C = D^p = (D_{i,j}^p)_{i,j} \in \mathbb{R}^{n \times n}$, wobei $D \in \mathbb{R}_+^{n \times n}$ eine Metrik ist.

Dann definiert

$$W_p(\mathbf{a}, \mathbf{b}) := L_{D^p}(\mathbf{a}, \mathbf{b})^{1/p} \quad (7.12)$$

die p -Wasserstein-Distanz auf Σ_n

Lemma 7.2.21. W_p ist eine Metrik

Beweis. Nach Voraussetzung besitzt $C = D^p$ eine Nulldiagonale. Somit gilt $W_p(\mathbf{a}, \mathbf{a}) = 0$. Die zugehörige Transportmatrix ist $\mathbf{P}^* = \text{diag}(\mathbf{a})$. Aufgrund der Positivität aller Nicht-Diagonalelemente von D^p gilt $W(\mathbf{a}, \mathbf{b})$ für $\mathbf{a} \neq \mathbf{b}$, da in diesem Fall jede zulässige Kopplungen und damit insbesondere die optimale Kopplungen ein Nicht-Diagonalelement ungleich 0 besitzt. Die Symmetrie von $W_p(\mathbf{a}, \mathbf{b})$ gilt wegen der Symmetrie von D^p .

Für den Nachweis der Dreiecksungleichung im Fall beliebiger Maße nutzt Villani [Vil03] das sogenannte Gluing Lemma. Im diskreten Fall ist die Konstruktion dieser geklebten Kopplung etwas einfacher. Seien $\mathbf{a}, \mathbf{b}, \mathbf{c} \in \Sigma_n$. Seien \mathbf{P} und \mathbf{Q} zwei optimale Lösungen des Transportproblems zwischen \mathbf{a} und \mathbf{b} bzw. \mathbf{b} und \mathbf{c} . Wir definieren $\tilde{\mathbf{b}}$, wobei $\tilde{\mathbf{b}}_j = \mathbf{b}_j$, falls $\mathbf{b}_j > 0$, und $\tilde{\mathbf{b}}_j = 1$ sonst gilt. Damit können wir

$$\mathbf{S} := \mathbf{P} \text{diag}(1/\tilde{\mathbf{b}}) \mathbf{Q} \in \mathbb{R}_+^{n \times n} \quad (7.13)$$

schreiben. Es gilt $\mathbf{S} \in \Pi(\mathbf{a}, \mathbf{c})$ wegen

$$\mathbf{S} \mathbf{1}_n = \mathbf{P} \text{diag}(1/\tilde{\mathbf{b}}) \mathbf{Q} \mathbf{1}_n = \mathbf{P}(\mathbf{b}/\tilde{\mathbf{b}}) = \mathbf{P} \mathbf{1}_{\text{supp}[\mathbf{b}]} = \mathbf{a}, \quad (7.14)$$

wobei $\mathbf{1}_{\text{supp}[\mathbf{b}]}$ den Vektor der Größe n bezeichnet, der Einsen an den Stellen mit Indizes j besitzt, für die auch $\mathbf{b}_j > 0$ gilt, und sonst aus Nullen besteht. Außerdem wurde $\mathbf{P} \mathbf{1}_{\text{supp}[\mathbf{b}]} = \mathbf{P} \mathbf{1}_n = \mathbf{a}$ benutzt, denn es gilt notwendigerweise $\mathbf{P}_{i,j} = 0$ für diejenigen j mit $\mathbf{b}_j = 0$. Analog folgt $\mathbf{S}^\top \mathbf{1}_n = \mathbf{c}$.

Damit erhalten wir

$$\text{content} \dots \quad (7.15)$$

□

Bemerkung 7.2.22 (Der Fall $0 < p \leq 1$). Für $0 < p \leq 1$ ist auch D^p eine Distanz. Damit ist für $p \leq 1$ $W_p(\mathbf{a}, \mathbf{b})^p$ eine Distanz auf dem Simplex.

Bemerkung 7.2.23. Für $p = 1$ ist die p -Wasserstein-Distanz auch als Earth Movers's Distance [RTG00] bekannt

7.2.6 Duale Formulierung

Kurzer Überblick hier²¹

The Kantorovich problem (2.11) is a constrained convex minimization problem, and as such, it can be naturally paired with a so-called dual problem, which is a constrained concave maximization problem. The following fundamental proposition explains the relationship between the primal and dual problems.

²¹<https://arxiv.org/pdf/1609.04767.pdf>

7.2.7 Gromov-Wasserstein-Divergenz

Fast Computation of Wasserstein Barycenters²²

Definition 7.2.24 (Divergenz). *Sei S der Raum aller Wahrscheinlichkeitsverteilungen mit gemeinsamem Support. Dann bezeichnet die Divergenz auf S eine Funktion $D(\cdot||\cdot) : S \times S \rightarrow \mathbb{R}$, für die gilt:*

1. $D(p||q) \geq 0$ f.a. $p, q \in S$

2. $D(p||q) = 0$ gdw. $p = q$.

Definition 7.2.25 (Entropie). *Für $T \in \mathbb{R}_+^{N \times N}$ definieren wir die Entropie als*

$$H(T) := - \sum_{i,j=1}^N T_{i,j} (\log(T_{i,j}) - 1). \quad (7.16)$$

Definition 7.2.26 (Tensor-Matrix-Multiplikation). *Für einen Tensor $\mathcal{L} = (\mathcal{L}_{i,j,k,l})_{i,j,k,l}$ und eine Matrix $(T_{i,j})_{i,j}$ definieren wir die Tensor-Matrix-Multiplikation als*

$$\mathcal{L} \otimes T := \left(\sum_{k,l} \mathcal{L}_{i,j,k,l} T_{k,l} \right)_{i,j}. \quad (7.17)$$

7.2.8 Verfahren zur Berechnung der exakten Lösung

7.2.9 Komplexitätsanalyse

Im Fall $n = m$ ist das Kantorovich-Problem und damit die Berechnung der Wasserstein-Distanz ein lineares Optimierungsproblem mit $\mathcal{O}(n)$ linearen Bedingungen. Für die Lösung kann beispielsweise der lineare Lee-Sinford-algorithmus zur Berechnung einer Lösung in $\tilde{\mathcal{O}}(n^{2.5})$ [LS14] verwendet werden, der den vorherigen Standard von $\mathcal{O}(n^{2.5})$ [Ren88] verbessert.

7.2.10 Regularisierungen

Tradeoff zwischen numerischem Aufwand und statistischer Genauigkeit.

Numerical Methods: Cuturi's Entropy Regularised Approach. Arguably the biggest development (at least in recent years) in the computation of optimal transport distances was due to Cuturi's entropy regularised approach. The idea is to use entropy to regularise the distance, then some simple rearrangements reveal this is a Kullback-Liebler divergence. Standard methods, e.g. Sinkhorn's algorithm, can then be used to find minimizers of the entropy regularised distance.²³

7.2.11 Entropisch Regularisierte Gromov-Wasserstein-Distanz

Entropic Regularization of Optimal Transport [com19] Mehrere Möglichkeiten einer Regularisierung der GW-Distanz²⁴:

²²<https://arxiv.org/pdf/1310.4375.pdf>

²³https://www.damtp.cam.ac.uk/research/cia/files/teaching/Optimal_Transport_Syllabus.pdf

²⁴<https://www.youtube.com/watch?v=cPVMHWF8fmE&t=2532s>

- Entropic regularization [Cuturi, 2013]
- Group Lasso [Courty et al., 2016a]
- KL, Itakura Saito, β -divergences, [Dessein et al., 2016]

Optimal Transport: Regularization and Applications ²⁵ Python Optimal Transport Toolbox^{26,27}

Kapitel 2.1 im Paper: Vergleich von Histogrammen auf demselben metrischen Raum.

Lemma 7.2.27. $L_C(\mathbf{a}, \mathbf{b}) := \min_{P \in \Pi(\mathbf{a}, \mathbf{b})} \langle \mathbf{P}, \mathbf{C} \rangle - \varepsilon H(\mathbf{P})$

ist ein ε -streng konvexes Problem und besitzt deshalb eine eindeutige optimale Lösung test

Beweis. Die Entropie H ist wegen der Hesse-Matrix $\partial^2 H(P) = -\text{diag}(1/P_{i,j})$ und $P_{i,j} \leq 1$ stark konkav. Durch die Multiplikation mit ε und anschließender Subtraktion wird aus dem obigen Problem ein konvexes. \square

Bemerkung 7.2.28. Für größere ε wird mehr Entropie gefordert

Proposition 7.2.29 (Konvergenz in ε). Die eindeutige Lösung P_ε von Theorem 7.2.27 konvergiert gegen die optimale Lösung mit maximaler Entropie innerhalb der Menge aller optimalen Lösungen des Kantorovich Problems, d.h.

$$P_\varepsilon \xrightarrow{\varepsilon \rightarrow 0} \arg \min_P \{-H(P) : P \in \Pi(\mathbf{a}, \mathbf{b}), \langle P, \mathbf{C} \rangle = L_C(\mathbf{a}, \mathbf{b})\}, \quad (7.18)$$

d.h. es gilt

$$L_C^\varepsilon(\mathbf{a}, \mathbf{b}) \xrightarrow{\varepsilon \rightarrow 0} L_C(\mathbf{a}, \mathbf{b}) \quad (7.19)$$

Zudem gilt

$$P_\varepsilon \xrightarrow{\varepsilon \rightarrow \infty} \mathbf{a} \otimes \mathbf{b} = \mathbf{a} \mathbf{b}^\top = (\mathbf{a}_i \mathbf{b}_j)_{i,j}. \quad (7.20)$$

Beweis. Sei eine Folge ε_l mit $\varepsilon_l \rightarrow 0$ und $\varepsilon_l > 0$ gegeben. Sei P_l die Lösung von Theorem 7.2.27 für $\varepsilon = \varepsilon_l$. Da $\Pi(\mathbf{a}, \mathbf{b})$ beschränkt ist, existiert eine Teilfolge ε_k mit $P_k \rightarrow P^*$. Aufgrund der Abgeschlossenheit von $\Pi(\mathbf{a}, \mathbf{b})$ folgt $P^* \in \Pi(\mathbf{a}, \mathbf{b})$.

Sei $P \in \Pi(\mathbf{a}, \mathbf{b})$ beliebig mit $\langle C, P \rangle = L_C(\mathbf{a}, \mathbf{b})$. Aufgrund der Optimalität von P und P_l bezüglich $L_C(\mathbf{a}, \mathbf{b})$ bzw. $L_C^{\varepsilon_l}(\mathbf{a}, \mathbf{b})$ gilt

$$0 \leq \langle C, P_l \rangle - \langle C, P \rangle \leq \varepsilon_l (H(P_l) - H(P)). \quad (7.21)$$

Aufgrund der Stetigkeit von H folgt für $l \rightarrow +\infty$ in Gleichung 7.21 $\langle C, P^* \rangle = \langle C, P \rangle$, womit P^* ein zulässiger Punkt ist. Division durch ε_l in Gleichung 7.21 und Übergang zum Grenzwert liefert $H(P) \leq H(P^*)$. Folglich ist P^* eine Lösung von Gleichung 7.18.

Da die Lösung P_0^* aufgrund der strikten Konvexität von $-H$ eindeutig ist, gilt $P^* = P_0^*$ und die gesamte Folge konvergiert. \square

²⁵<https://www.otra2020.com/schedule>

²⁶<https://pythonot.github.io/quickstart.html>

²⁷https://pythonot.github.io/auto_examples/gromov/plot_gromov.html

Bemerkung 7.2.30. Gleichung 7.18 zeigt, dass die Lösung für kleine ε gegen die Optimale Transport Kopplung mit maximaler Entropie konvergiert. Im Gegensatz dazu, bedeutet Gleichung 7.20, die Lösung für große Regularisierungsparameter konvergiert gegen die Kopplung mit maximaler Entropie zwischen zwei gegebenen Randverteilungen \mathbf{a} und \mathbf{b} .

Bemerkung 7.2.31. A key insight is that, as ε increases, the optimal coupling becomes less and less sparse (in the sense of having entries larger than a prescribed threshold), which in turn has the effect of both accelerating computational algorithms (as we study in Abschnitt 4.2) and leading to faster statistical convergence (as shown in Abschnitt 8.5)

Wir definieren die Kullback-Leibler-Divergenz zwischen Kopplungen als

$$KL(P|K) := \sum_{i,j} P_{i,j} \log \left(\frac{P_{i,j}}{K_{i,j}} \right) - P_{i,j} + K_{i,j}. \quad (7.22)$$

Damit ist die eindeutige Lösung P_ε von Theorem 7.2.27 eine Projektion des zur Kostenmatrix \mathbf{C} gehörigen Gibbs-Kernels $K_{i,j} := e^{-\frac{C_{i,j}}{\varepsilon}}$ auf $\Pi(\mathbf{a}, \mathbf{b})$.

Mit der obigen Definition erhalten wir

$$P_\varepsilon = \text{Proj}_{\Pi(\mathbf{a}, \mathbf{b})}^{KL}(\mathbf{K}) := \arg \min_{P \in \Pi(\mathbf{a}, \mathbf{b})} KL(P|K). \quad (7.23)$$

7.2.12 Berechnung der Lösung: Sinkhorn

In diesem Kapitel sehen wir, dass die Lösung des regularisierten Problems eine besondere Form besitzt, die wir über $m + n$ Variablen parametrisieren können.

Proposition 7.2.32. Die Lösung des regularisierten Problems Theorem 7.2.27 besitzt die Form

$$\forall (i, j) \in \{1, \dots, n\} \times \{1, \dots, m\} : P_{i,j} = \mathbf{u}_i K_{i,j} \mathbf{v}_j \quad (7.24)$$

für die beiden (unbekannten) Variablen $(\mathbf{u}, \mathbf{v}) \in \mathbb{R}_+^n \times \mathbb{R}_+^m$.

Beweis. Wir führen für jede der beiden Nebenbedingungen die dualen Variablen $\mathbf{f} \in \mathbb{R}^n$ und $\mathbf{g} \in \mathbb{R}^m$ ein. Für die Lagrange-Funktion zu Theorem 7.2.27 erhalten wir damit:

$$\mathcal{L}(P, \mathbf{f}, \mathbf{g}) = \langle P, \mathbf{C} \rangle - \varepsilon H(P) - \langle \mathbf{f}, P \mathbf{1}_m - \mathbf{a} \rangle - \langle \mathbf{g}, P^\top \mathbf{1}_n - \mathbf{b} \rangle. \quad (7.25)$$

Mit der Optimalitätsbedingung erster Ordnung ergibt sich

$$\frac{\partial \mathcal{L}(P, \mathbf{f}, \mathbf{g})}{\partial P_{i,j}} = C_{i,j} + \varepsilon \log(P_{i,j}) - \mathbf{f}_i - \mathbf{g}_j = 0, \quad (7.26)$$

womit wir für eine optimale Kopplung P für das regularisierte Problem den Ausdruck $P_{i,j} = e^{\mathbf{f}_i/\varepsilon} e^{-C_{i,j}/\varepsilon} e^{\mathbf{g}_j/\varepsilon}$, der in die gewünschte Form umgeschrieben werden kann. \square

Bemerkung 7.2.33 (Algorithmus von Sinkhorn). *Die Faktorisierung der Lösung in Gleichung 7.24 können in der folgenden Matrix-Form schreiben: $\mathbf{P} = \text{diag}(\mathbf{u})\mathbf{K}\text{diag}(\mathbf{v})$. Die beiden Variablen (\mathbf{u}, \mathbf{v}) müssen deshalb die folgenden nichtlinearen Gleichungen erfüllen, die aufgrund der geforderten Massenerhaltungsbedingung in $\Pi(\mathbf{a}, \mathbf{b})$ gelten:*

$$\text{diag}(\mathbf{u})\mathbf{K}\text{diag}(\mathbf{v})\mathbf{1}_m = \mathbf{a} \text{ und } \text{diag}(\mathbf{v})\mathbf{K}^\top \text{diag}(\mathbf{u})\mathbf{1}_n = \mathbf{b}. \quad (7.27)$$

Aufgrund der Beziehung $\text{diag}(\mathbf{v})\mathbf{1}_m = \mathbf{v}$ und selbiger Beziehung für \mathbf{u} erhalten wir die folgende Vereinfachung

$$\mathbf{u} \odot (\mathbf{K}\mathbf{v}) = \mathbf{a} \text{ und } \mathbf{v} \odot (\mathbf{K}^\top \mathbf{u}) = \mathbf{b}, \quad (7.28)$$

wobei \odot für die Element-weise Multiplikation zweier Vektoren steht. Dieses Problem ist als Matrix Scaling Problem [NR99] bekannt.

(? Sollte ich hier Bedingungen für die Lösbarkeit angeben?)

Eine Möglichkeit zur Lösung dieses Problems ist ein iteratives Vorgehen(?Verfahren), bei dem zunächst \mathbf{u} so modifiziert wird, dass die linke Seite in Gleichung 7.28 erfüllt ist, und anschließend die Modifikation von \mathbf{v} vorgenommen wird, sodass die rechte Seite in Gleichung 7.28 gilt. Mit diesen beiden Modifikationen erhalten wir den Algorithmus von Sinkhorn, der aus den beiden folgenden Updates besteht:

$$\mathbf{u}^{l+1} := \frac{\mathbf{a}}{\mathbf{K}\mathbf{v}^{(l)}} \text{ und } \mathbf{v}^{l+1} := \frac{\mathbf{b}}{\mathbf{K}^\top \mathbf{u}^{(l+1)}}, \quad (7.29)$$

wobei zu Beginn mit einem beliebigen positiven Vektor, beispielsweise $\mathbf{v}^{(0)} = \mathbf{1}_m$ initialisiert wird und l den aktuellen Iterationsschritt bezeichnet. Die obigen Division muss ebenfalls elementweise verstanden werden.

The iterations (4.15) first appeared in [Yule, 1912, Kruithof, 1937]. They were later known as the iterative proportional fitting procedure (IPFP) Deming and Stephan [1940] and RAS [Bacharach, 1965] methods [Idel, 2016]. The proof of their convergence is attributed to Sinkhorn [Sin64], hence the name of the algorithm.

Bemerkung 7.2.34. *Die Lösung zum regularisierten Problem lässt sich nach den Sinkhorn-Iterationen wie folgt angeben: $P_L = \text{diag}(U_L)\mathbf{K}\text{diag}(v_L)$ Das ist der optimal Transport plan für das Regularisierte Problem. Optimal Transport Caost lässt sich berechnen als: $\langle P_L, M_{XY} \rangle = u_L^\top (\mathbf{K} \odot M_{XY}) v_L$.*

Für die globale Konvergenzanalyse der Sinkhorn-Algorithmus lässt sich ein Zusammenhang mit der Hilbertschen Projektionsmetrik benutzen, der erstmals in [FL89] vorgestellt wurde.

Diese ist wie folgt definiert:

Definition 7.2.35 (Projektive Hilbert-Metrik). *Seien $u, u' \in \mathbb{R}_{+, \star}^n$. Dann ist die projektive Hilbert-Metrik $d_{\mathcal{H}}$ definiert durch*

$$d_{\mathcal{H}}(u, u') := \log \max_{i,j} \frac{u_i u'_j}{u_j u'_i}. \quad (7.30)$$

Bemerkung 7.2.36. *It can be shown to be a distance on the projective cone $R^n_{+, \star} / \sim$, where $u \sim u'$ means that $\exists r > 0, u = ru'$ (the vectors are equal up to rescaling, hence the name “projective”). This means that $d_{\mathcal{H}}$ satisfies the triangular inequality and $d_{\mathcal{H}}(u, u') = 0$ if and only if $u \sim u'$. This is a projective*

version of Hilbert's original distance on bounded open convex sets [Hilbert, 1895]. The projective cone $\mathbb{R}^n_{+, \star} / \sim$ is a complete metric space for this distance. By a logarithmic chan

the Hilbert metric on the rays of the positive cone is isometric to the variation seminorm (it is a norm between vectors that are defined up to an additive constant) $dH(u, u_0) \log(u) - \log(u_0)$

where (4.21) var def. wobei $\|f\|_{\text{var}} := (\max_i f_i) - (\min_i f_i)$.

This variation seminorm is closely related to the ' norm since one always has $\|f\|_{\text{var}} \leq \|f\|_{\infty}$. If one imposes that $f_i = 0$ for some fixed i , then a converse inequality also holds since $\|f\|_{\infty} \leq \|f\|_{\text{var}}$. These bounds are especially useful to analyze Sinkhorn convergence (see Remark 4.14 below), because dual variables $f = \log(u)$ solving (4.14) are defined up to an additive constant, so that one can impose that $f_i = 0$ for some i .

The Hilbert metric was introduced independently by [Birkhoff, 1957] and [Samelson et al., 1957]. They proved the following fundamental theorem, which shows that a positive matrix is a strict contraction on the cone of positive vectors.

Theorem 7.2.37. Sei $K \in \mathbb{R}^{n \times m}_{+, \star}$. Dann gilt für $(v, v') \in (\mathbb{R}^m_{+, \star})^2$

$$d_{\mathcal{H}}(Kv, Kv') \leq \lambda(K) d_{\mathcal{H}}(v, v'), \quad \text{mit} \quad \begin{cases} \lambda(K) := \frac{\sqrt{\nu(K)} - 1}{\sqrt{\nu(K)} + 1} \\ \nu(K) := \max_{i,j,k,l} \frac{K_{i,k} K_{j,l}}{K_{j,k} K_{i,l}} \end{cases} \quad (7.31)$$

Visualisierung des Theorems vgl [COT19], Seite 70.

Mit diesem Resultat lässt sich die Konvergenz des Verfahrens zeigen:

Theorem 7.2.38. Es gilt $(u^{(l)}, v^{(l)}) \rightarrow (u^*, v^*)$ und

$$d_{\mathcal{H}}(u^{(l)}, u^*) = \mathcal{O}(\lambda(K)^{2l}), \quad d_{\mathcal{H}}(v^{(l)}, v^*) = \mathcal{O}(\lambda(K)^{2l}) \quad (7.32)$$

Es gelten außerdem die beiden folgenden Abschätzungen

$$d_{\mathcal{H}}(u^{(l)}, u^*) \leq \frac{d_{\mathcal{H}}(P^{(l)} \mathbf{1}_m, a)}{1 - \lambda(K)^2}, \quad (7.33)$$

$$d_{\mathcal{H}}(v^{(l)}, v^*) \leq \frac{d_{\mathcal{H}}(P^{(l)\top} \mathbf{1}_n, b)}{1 - \lambda(K)^2}, \quad (7.34)$$

wobei

$$P^{(l)} := \text{diag}(u^{(l)}) K \text{diag}(v^{(l)}) \quad (7.35)$$

gilt. Desweiteren gilt die Abschätzung

$$\|\log(P^{(l)}) - \log(P^*)\|_{\infty} \leq d_{\mathcal{H}}(u^{(l)}, u^*) + d_{\mathcal{H}}(v^{(l)}, v^*), \quad (7.36)$$

wobei P^* die eindeutige Lösung zu ?? ist.

Beweis. Für ein beliebiges Paar $(v, v') \in (\mathbb{R}^m_{+, \star})^2$ gilt

$$d_{\mathcal{H}}(v, v') = d_{\mathcal{H}}(v/v', \mathbf{1}_m) = d_{\mathcal{H}}(\mathbf{1}_m/v, \mathbf{1}_m/v').$$

Mit dieser Beziehung und Theorem 7.2.37 erhalten wir

$$\begin{aligned} d_{\mathcal{H}}(u^{(l+1)}, u^*) &= d_{\mathcal{H}}\left(\frac{a}{Kv^{(l)}}, \frac{a}{Kv^*}\right) \\ &= d_{\mathcal{H}}(Kv^{(l)}, Kv^*) \\ &\leq \lambda(K) d_{\mathcal{H}}(v^{(l)}, v^*). \end{aligned}$$

Dies zeigt Gleichung 7.32. Mit Hilfe der Dreiecksungleichung erhalten wir

$$\begin{aligned} d_{\mathcal{H}}(u^{(l)}, u^*) &\leq d_{\mathcal{H}}(u^{(l+1)}, u^{(l)}) + d_{\mathcal{H}}(u^{(l+1)}, u^*) \\ &\leq d_{\mathcal{H}}\left(\frac{a}{Kv^{(l)}}, u^{(l)}\right) + \lambda(K)^2 d_{\mathcal{H}}(u^{(l)}, u^*) \\ &= d_{\mathcal{H}}\left(a, u^{(l)} \odot (Kv^{(l)})\right) + \lambda(K)^2 d_{\mathcal{H}}(u^{(l)}, u^*) \\ &= d_{\mathcal{H}}\left(a, P^{(l)} \mathbf{1}_m\right) + \lambda(K)^2 d_{\mathcal{H}}(u^{(l)}, u^*) \end{aligned}$$

und damit nach Division durch $1 - \lambda(K)^2$ Gleichung 7.33. Die zweite Abschätzung, Gleichung 7.34, folgt analog. Die Gleichung 7.36 gilt nach Lemma 3 in [FL89]. \square

Bemerkung 7.2.39.

- *Konvergenzrate?*
- *Berechnung der optimalen Transportpläne mit Gleichung 7.35 im Zeitschritt l und Berechnung der Transportkosten ...*
- *Abbruchkriterien (vgl. mit der POT Implementierung)*

Bemerkung 7.2.40 (Allgemeine Formulierung des entropisch regularisierten Problems für beliebige Maße).

- *Benutze Relative Entropie als Verallgemeinerung der diskreten Kullback-Leibler-Divergenz*
- *Referenzmaß ist unbedeutend, lediglich er Support hat eine Bedeutung.*

7.2.13 Komplexitätsanalyse

[AWR17] liefert eine Komplexitätsanalyse für die Sinkhorn-Iterationen. Im Fall $n = m$ sind für die Wahl $\varepsilon = \frac{4 \log(n)}{\tau} \mathcal{O}(\|C\|_{\infty}^3 \log(n) \tau^{-3})$ Sinkhorn-Iterationen (inklusive eines Rundungsschrittes) notwendig, um eine zulässige Kopplung $\hat{\mathbf{P}} \in \Pi(\mathbf{a}, \mathbf{b})$ zu berechnen, die die Abschätzung $\langle \hat{\mathbf{P}}, \mathbf{C} \rangle \leq L_{\mathbf{C}}(\mathbf{a}, \mathbf{b} + \tau)$ zu berechnen. Somit liefert das Sinkhorn-Verfahren eine τ -Approximation des nicht regularisierten Optimal Transport-Problems in $\mathcal{O}(n^2 \log(n) \tau^{-3})$ Operationen. Gleichzeitig wird dort eine greedy Variante der Sinkhorn-Iterationen namens Greenkhorn präsentiert, die eine τ -Approximation in $\mathcal{O}(n^2 \tau^{-3})$ Operationen liefert. [DGK18] verbessert diese auf $\mathcal{O}(n^2 \tau^{-2})$ Operationen.

8 TODO: Warum wir GW-Distanz brauchen

<https://arxiv.org/pdf/1811.02834.pdf>

8.0.1 Verallgemeinerung: Gromov-Wasserstein-Divergenz

Der Vergleich zwischen Ähnlichkeits- bzw. Distanzmatrizen ist schwierig, da diese die innere Struktur eines Datensatzes beschreiben, die unabhängig von Rotationen und Translationen ist. Es existiert keine kanonische Ordnung der Reihen und Spalten. Verallgemeinerung auf beliebige Matrizen C , d.h. diese Distanzmatrizen müssen nicht notwendigerweise positiv sein und die Dreiecksungleichung erfüllen.

Definiere die verallgemeinerte Gromov-Wasserstein-Distanz (vGWD) wie folgt:

Definition 8.0.1 (Verallgemeinerte Gromov-Wasserstein-Distanz). *Seien zwei gewichtete Ähnlichkeitsmatrizen $(C, p) \in \mathbb{R}^{N_1 \times N_1} \times \Sigma_{N_1}$ und $(\bar{C}, q) \in \mathbb{R}^{N_2 \times N_2} \times \Sigma_{N_2}$ gegeben. Sei T eine Kopplung zwischen den beiden Räumen, auf denen die Matrizen C und \bar{C} definiert sind. Sei L eine Fehlerfunktion. Dann definieren wir die verallgemeinerte Gromov-Wasserstein-Distanz als*

$$GW(C, \bar{C}, p, q) := \min_{T \in \mathcal{C}_{p,q}} \varepsilon_{C, \bar{C}}(T), \quad (8.1)$$

wobei gilt $\varepsilon_{C, \bar{C}}(T) := \sum_{i,j,k,l} L(C_{i,k}, \bar{C}_{j,l}) T_{i,j} T_{k,l}$.

Häufig verwendete Fehlerfunktionen sind die quadratische Fehlerfunktion $L(a, b) = L_2(a, b) := \frac{1}{2}|a - b|^2$ und die Kullback-Leibler-Divergenz $L(a, b) = KL(a|b) := a \log(a/b) - a + b$.

Diese Definition der Gromov-Wasserstein-Distanz verallgemeinert die Version in [PCS16], da nun beliebige Fehlerfunktionen betrachtet werden.

Für $L = L_2$ zeigt Memoli, 2011, dass $GW^{1/2}$ eine Distanz auf dem Raum metrischer Maßräume modulo Maß-erhaltender Isometrien (? , besser zitieren -> vollständiges Resultat angeben) definiert.

9 TODO: Beweis für L=L2

[Mém11] Durch die Definition

$$\mathcal{L}(C, \bar{C}) := (L(C_{i,k}, \bar{C}_{j,l}))_{i,j,k,l} \quad (9.1)$$

erhalten wir

$$\varepsilon_{C, \bar{C}}(T) = \langle \mathcal{L}(C, \bar{C}) \otimes T, T \rangle \quad (9.2)$$

gilt.

Mit der folgenden Proposition ergibt sich eine effiziente Berechnung von $\mathcal{L}(C, \bar{C}) \otimes T$ für eine bestimmte Klasse von Verlustfunktionen L :

Proposition 9.0.1. *Die Verlustfunktion L lasse sich schreiben als*

$$L(a, b) = f_1(a) + f_2(b) - h_1(a)h_2(b) \quad (9.3)$$

für $f_1, f_2, h_1, h_2 : \mathbb{R} \rightarrow \mathbb{R}$. Dann gilt für $T \in \mathcal{C}_{p,q}$:

$$\mathcal{L}(C, \bar{C}) \otimes T = c_{C, \bar{C}} - h_1(C) T h_2(\bar{C})^T, \quad (9.4)$$

wobei $c_{C, \bar{C}} := f_1(C) p \mathbf{1}_{N_2}^T + \mathbf{1}_{N_1} q^T f_2(\bar{C})^T$ unabhängig von T ist.

Beweis. Aufgrund von Gleichung 9.3 gilt nach der Tensor-Matrix-Multiplikation Gleichung 7.17 die Zerlegung $\mathcal{L}(C, \bar{C}) \otimes T = A + B + C$ mit

$$\begin{aligned} A_{i,j} &= \sum_k f_1(C_{i,k}) \sum_l T_{k,l} = (f_1(C)(T\mathbf{1}))_i, \\ B_{i,j} &= \sum_l f_2(\bar{C}_{j,l}) \sum_k T_{k,l} = (f_2(\bar{C})(T^\top \mathbf{1}))_j, \\ C_{i,j} &= \sum_k h_1(C_{i,k}) \sum_l h_2(\bar{C}_{j,l}) T_{k,l}. \end{aligned}$$

Dies ist äquivalent zu $(h_1(C))(h_1(\bar{C}T^\top)^\top)_{i,j}$.
(? Wie folgt daraus die Behauptung) \square

Bemerkung 9.0.2 (Verbesserte Komplexität). *Mit dem Resultat in Theorem 9.0.1 können wir $\mathcal{L}(C, \bar{C}) \otimes T$ effizient in der Größenordnung $\mathcal{O}(N_1^2 N_2 + N_2^2 N_1)$ mit ausschließlich Matrix/Matrix-Multiplikationen berechnen im Unterschied zur Komplexität von $\mathcal{O}(N_1^2 N_2^2)$ für die Implementierung von Gleichung 7.17.*

Bemerkung 9.0.3 (Spezialfälle). *Im Fall $L = L_2$ ist die Bedingung Gleichung 9.3 für die Funktionen $f_1(a) = a^2, f_2(b) = b^2, h_1(a) = a$ und $h_2(b) = 2b$ erfüllt. Für $L = KL$ sind die Funktionen $f_1(a) = a \log(a) - a, f_2(b) = b, h_1(a) = a$ und $h_2(b) = \log(b)$ notwendig.*

Wir betrachten nun die regularisierte Version der Gromow-Wasserstein-Diskrepanz 8.1 und definieren:

Definition 9.0.4. *Für C, \bar{C}, p, q , wie oben, definieren wir die entropisch regularisierte Gromov-Wasserstein-Diskrepanz als*

$$GW_\varepsilon(C, \bar{C}, p, q) := \min_{T \in \mathcal{C}_{p,q}} \varepsilon_{C, \bar{C}}(T) - \varepsilon H(T). \quad (9.5)$$

Wir erhalten damit ein nicht-konvexes Optimierungsproblem. Für dessen Lösung benutzen wir ein projiziertes Gradienten-Verfahren, bei dem sowohl die Schrittweite als auch die Projektion bezüglich der KL-Metrik berechnet werden.

Die Iterationen sind gegeben durch

$$T \leftarrow Proj_{\mathcal{C}_{p,q}}^{KL} \left(T \odot e^{-\tau(\nabla \varepsilon_{C, \bar{C}}(T) - \varepsilon H(T))} \right), \quad (9.6)$$

wobei die Schrittweite $\tau > 0$ und die KL-Projektion einer beliebigen Matrix K gegeben ist durch:

$$Proj_{\mathcal{C}_{p,q}}^{KL}(K) := \arg \min_{T' \in \mathcal{C}_{p,q}} KL(T'|K). \quad (9.7)$$

Proposition 9.0.5. *Für den Fall $\tau = 1/\varepsilon$ erhalten wir die Iterationsvorschrift*

$$T \leftarrow \mathcal{T}(\mathcal{L}(C, \bar{C}) \otimes T, p, q). \quad (9.8)$$

Beweis. Nach [BCC⁺15] ist die Projektion in 9.6 gegeben durch die Lösung des regularisierten Transportproblems 7.2.27 und ist damit gegeben durch:

$$Proj_{\mathcal{C}_{p,q}}^{KL}(K) = \mathcal{T}(-\varepsilon \log(K), p, q) \quad (9.9)$$

(? Wo steht das genau in dem Paper?) Es gilt außerdem

$$\nabla \varepsilon_{C, \bar{C}}(T) - \varepsilon H(T) = \text{blub} \quad (9.10)$$

Durch Umordnen der Terme in Gleichung 9.6 erhalten wir für $\tau\varepsilon = 1$ die angegebene Vorschrift. \square

Bemerkung 9.0.6. Die Iterationsvorschrift 9.8 definiert einen einfachen Algorithmus, der in jedem Update von T eine Sinkhorn-Projektion benötigt.

Bemerkung 9.0.7 (Konvergenz). Die Iterationen Gleichung 9.6 konvergieren nach [BCL16] für $\tau < \tau_{\max}$. Im Allgemeinen gilt jedoch $1/\varepsilon < \tau_{\max}$ jedoch nicht, womit die Wahl der Schrittweite $\tau = 1$ in Theorem 9.0.5 nicht durch die Theorie abgedeckt ist. Laut [PCS16] konvergiert die Iteration mit $\tau = 1/\varepsilon$ jedoch in der Praxis.

Für den Fall $L = L_2$ ergibt sich der „Softassign quadratic assignment algorithm“, [RYM99], für welchen ein Konvergenzresultat im Fall eines konvexen Problems existiert. Da wir in unserem Fall nur positive symmetrische Matrizen C, C_s betrachten ist hier die Konvergenz gesichert.

Bemerkung 9.0.8. content...

Bemerkung 9.0.9 (Wahl von ε). In [Cut13] werden verschiedene Werte für ε angegeben. Diese sind $\varepsilon = 0.02, 0.1, 1.0$

Im zugehörigen Beispiel²⁸ von POT wird $\varepsilon = 0.0005$ verwendet.

Für ε klein werden die Ergebnisse besser während der Rechenaufwand steigt. welche Resultate existieren bezüglich der Approximationsgüte abhängig von ε ?

9.0.1 Wasserstein Baryzentren

In diesem Kapitel wollen wir uns mit dem "Mittelwert" befassen, der elementar für das kmeans-Clustering ist. Auch hier soll der Mittelwert auf die abstrakte Ebene von gewichteten Distanzmatrizen angehoben werden. Dazu definieren wir im Folgenden sogenannte Wasserstein-Baryzentren und folgen [PCS16] für die Lösung des entstehenden Optimierungsproblems. Gute Einführung²⁹

Definition 9.0.10 (Gromov-Wasserstein Baryzentrum). Seien die gewichteten Distanzmatrizen $C_s)_{s=1}^S$, mit $C_s \in \mathbb{R}^{N_s \times N_s}$ und den zugehörigen Histogrammen $(p_s)_s$ gegeben.

Dann ist das Gromov-Wasserstein Baryzentrum definiert durch

$$\min_{C \in \mathbb{R}^{N \times N}} \sum_s \lambda_s \text{GW}_\varepsilon(C, C_s, p, p_s). \quad (9.11)$$

Existenz und Eindeutigkeit von Baryzentren: siehe [AC11].

Bemerkung 9.0.11 (Darstellung des Baryzentrums). Um das berechnete Baryzentrum C wieder zu visualisieren, kann C als Distanzmatrix wieder in den zwei-dimensionalen Raum eingebettet werden. Dies kann beispielsweise durch Multi-dimensionale Skalierung erreicht werden³⁰.

²⁸https://pythonot.github.io/auto_examples/gromov/plot_gromov.html

²⁹<https://hal.archives-ouvertes.fr/hal-00637399/document>

³⁰https://pythonot.github.io/auto_examples/gromov/plot_gromov_barycenter.html

Bemerkung 9.0.12. Wir gehen im Folgenden davon aus, dass sowohl die Histogramme $(p_s)_s$ als auch das Histogramm p bekannt sind. Die Größe (N, N) des gesuchten Bary-Zentrums muss ebenfalls vorher festgelegt werden. Eine Erweiterung auf den Fall, dass auch p unbekannt sein sollte und damit als Optimierungsparameter aufgefasst wird, ist leicht möglich [PCS16].

Wir können das Bary-Zentrum mithilfe von Kopplungen umformulieren als

$$\min_{C, (T_s)_s} \sum_s \lambda_s (\varepsilon_{C, C_s}(T_s) - \varepsilon H(T_s)), \quad (9.12)$$

unter den Nebenbedingungen: $\forall s : T_s \in \mathcal{C}_{p, p_s} \subset \mathbb{R}_+^{N \times N_s}$. Für den Fall, dass L bezüglich der ersten Variable konvex ist, ist dieses Problem konvex bezüglich C . Bezüglich $(T_s)_s$ ist diese Problem quadratisch aber nicht notwendigerweise positiv.

Das Problem Gleichung 9.12 kann durch eine Block-Koordinaten-Relaxierung gelöst werden. Dabei wird iterativ abwechselnd bezüglich den Kopplungen $(T_s)_s$ und der Metrik C minimiert.

Minimierung bezüglich $(T_s)_s$. Anhand der Umformulierung Gleichung 9.12 sehen wir, dass das Optimierungsproblem Gleichung 9.11 bezüglich $(T_s)_s$ in S (?viele) unabhängige GW_ε -Optimierungen

$$\forall s : \min_{T_s \in \mathcal{C}_{p, p_s}} \mathcal{E}_{C, C_s}(T_s) - \varepsilon H(T_s) \quad (9.13)$$

zerfällt, die jeweils wie in Theorem 9.0.5 angegeben gelöst werden können.

Minimierung bezüglich C . Sei $(T_s)_s$ gegeben. Dann lautet, die Minimierung bezüglich C :

$$\min_C \sum_s \lambda_s \langle \mathcal{L}(C, C_s) \otimes T, T \rangle. \quad (9.14)$$

Mit der folgenden Proposition erhalten wir für eine Klasse von Verlustfunktionen L eine Lösung in geschlossener Form.

Proposition 9.0.13. Sei L eine Verlustfunktion, die die Bedingung Gleichung 9.3 erfüllt. Sei f'_1/h'_1 invertierbar.

Dann lässt sich die Lösung zu Gleichung 9.14 schreiben als:

$$C = \left(\frac{f'_1}{h'_1} \right)^{-1} \left(\frac{\sum_s \lambda_s T_s^\top h_2(C_s) T_s}{pp^\top} \right), \quad (9.15)$$

wobei die Normalisierung $\lambda_s = 1$ gilt.

Beweis. Nach Theorem 9.0.1 können wir das zu minimierende Funtional schreiben als

$$\sum_s \lambda_s \langle f_1(C) p \mathbf{1}^\top + \mathbf{1} p_s^\top f_2(C_s) - h_1(C) T_s h_2(C_s)^\top, T_s \rangle. \quad (9.16)$$

Die Optimalitätsbedingung erster Ordnung lautet folglich

$$f'_1(C) \odot pp^\top = h'_1(C) \odot \sum_s \lambda_s T_s h_2(C_s) T_s^\top. \quad (9.17)$$

Durh Umstellen der Gleichung erhalten wir die angegebene Form. \square

Anhand von Gleichung 9.15 wird die folgende Interpretation deutlich/möglich:
Für jedes $s \in S$ ist $T_s^\top h_2(C_s) T_s$ eine wiedereingeordnete Matrix, wobei T_s als Optimal Transport-Kopplung von Zeilen und Spalten der Distanzmatrix C_s fungiert. Über diese Matrizen wird anschließend gemittelt, wobei die Art der Mittelung von der Verlustfunktion L abhängig ist.

Für den Fall $L = L_2$ wird aus dem Update Gleichung 9.15 die folgende Vorschrift:

$$C \leftarrow \frac{1}{pp^\top} \sum_s \lambda_s T_s^\top C_s T_s. \quad (9.18)$$

Proposition 9.0.14. *Sei $L = L_2$ und $(C_s)_s$ positiv semi-definit f.a. $s \in S$. Dann sind die Iterierten C ebenfalls positiv semi-definit.*

Beweis. Gleichung 9.18 zeigt, dass das Update aus einer Bildung des Mittelwertes der Matrizen $(\text{diag}(1/p) T_s^\top C_s T_s \text{diag}(1/p))_s$ besteht. Diese sind alle positiv semi-definit, da die $(C_s)_s$ nach Voraussetzung positiv semi-definit sind. \square

Für den Fall $L = KL$ ergibt sich die folgende Verfahrensvorschrift:

$$C \leftarrow \left(\frac{1}{pp^\top} \sum_s \lambda_s T_s^\top \log(C_s) T_s \right). \quad (9.19)$$

Algorithm 2 Berechnung der GW_{ε} Baryzentren

Input: $(C_s, p_s), p$. Initialisiere C .

Output: C .

```

if some condition is true then
  do some processing
else if some other condition is true then
  do some different processing
else
  do the default actions
end if

```

Pseudocode.

Häufige Verwendungen(s. Einführung hier ³¹). zB. Shape interpolation

Wir haben in diesem Kapitel gesehen, wie mithilfe des Wasserstein-Baryzentums eine Art Mittelwert für Wahrscheinlichkeitsverteilungen definiert werden kann.

Algorithm 1 details the steps of the optimization technique. Note that it makes use of three nested iterations: (i) blockwise coordinate descent on $(T_s)_s$ and C , (ii) projected gradient descent on each T_s , and (iii) Sinkhorn iterations to compute the projection.

9.0.2 Komplexitätsanalyse

3

³¹<https://arxiv.org/pdf/2102.01752.pdf>

Mo, 2.August:

```
s=2
pp=33
eps_init: 0.0005
eps_update: 0.0005
n_samples_bary: 10
iter: 0: (tn, fp, fn, tp): 1650,0,491,322
iter: 1: (tn, fp, fn, tp): 449,1201,802,11
iter: 2: (tn, fp, fn, tp): 1650,0,32,781
iter: 3: (tn, fp, fn, tp): 1621,29,27,786
iter: 4: (tn, fp, fn, tp): 1632,18,27,786

acc = 98.17
tpr = 96.98
0.11666667 0.23055556 0.26 0.55333333 0.43333333 nan 0.23333333 0.28444444 0.11555556
0.07708333 0.15 0.0952381 0.4173913 0.41527778 0.40740741 0.45238095 0.01333333
0.35833333 0.29487179 0. 0. 0.04444444 0.05 0.01333333 0.4 0.29791667 0.13888889
0.23333333 0.11333333 0.3 0.06 0.02962963 0.11666667 0.65714286 0.475 0.71794872
0.50833333 0.43333333 0.35652174 0.45555556 0.27777778 0. 0.51111111] Perfor-
mance of poisoned net on unpoisoned training data: loss on test dataset: 0.19637071904851583
Accuracy on test Dataset: 0.958 tnr = 98.91
```

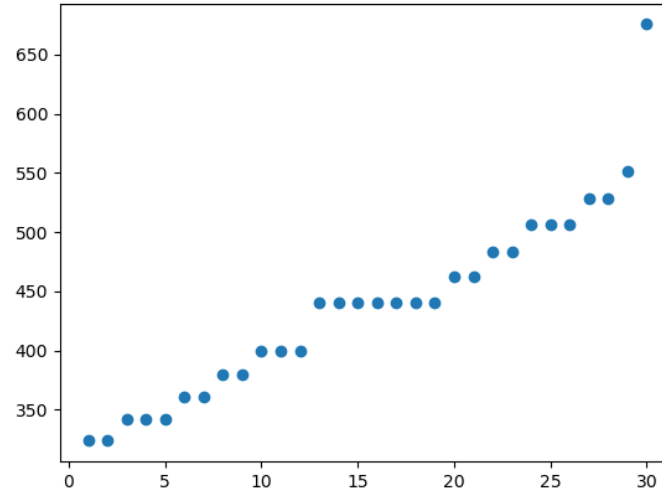
Mi, 4.August:

```
s=2, pp=15
eps_init: 0.0005
eps_update: 0.0005
n_samples_bary: 10
iter: 0: (tn, fp, fn, tp): 1650,0,89,202
iter: 1: (tn, fp, fn, tp): 1648,2,3,288
iter: 2: (tn, fp, fn, tp): 1650,0,0,291

acc = 100.00
tpr = 100.00
tnr = 100.00
```

Do, 5.August Das ist der Vergleich zu dem Experiment am Montag mit $\text{eps}=0.0005$. Das Ergebnis ist deutlich besser(denn perfekt).Daraus schließen wir, dass $\text{eps}=0.02$ für die restlichen Experimente genügen sollte. (Stimmt so vll nicht ganz)

```
eps_init: 0.02
eps_update: 0.02
n_samples_bary: 10
iter: 0: (tn, fp, fn, tp): 1650,0,812,1
iter: 1: (tn, fp, fn, tp): 1650,0,0,813
iter: 2: (tn, fp, fn, tp): 1650,0,0,813
```

Spektralanalyse:Absolute values of first 30 eigenvalues of L_{sym} with 10-nearest neighbours**Abbildung 2:** Ergebnis des spektralen Clusterings unter Verwendung der Euklidischen Distanz und $k=10$ Nachbarn

Clustering(euklidisch): $(tn, fp, fn, tp) = (1650, 0, 386, 427)$

Clustering(GWD):

Bemerkung 10.0.1. Für größere Netzwerke ist es einfacher, erfolgreiche Angriffe zu implementieren, auch schon mit weniger korruptierten Daten.

Vermutung: Die Detektion mittels Clustering funktioniert für eine größere Anzahl an korruptierten Daten besser. Wenn wir also nur perfekte Angriffe verteidigen wollen, d.h. $AER=100.00$ funktioniert das bei kleineren Netzwerken besser.

10.0.2 Label-konsistente Poisoning-Angriffe

Bemerkung 10.0.2. Für einen einzelnen Amplitudensticker mit $d=10$ ist das eigentlich identisch zu dem Angriff mit dem Sticker

```

0.7 0.86111111 0.85066667 0.81111111 0.66666667 nan 0.42 0.72888889 0.60888889
0.77083333 0.40757576 0.46666667 0.69855072 0.88472222 0.53333333 0.74285714
0.39333333 0.67222222 0.75128205 0.11666667 0.58888889 0.32222222 0.425 0.18
0.44444444 0.44375 0.8 0.71666667 0.56 0.86666667 0.11333333 0.3 0.71666667 0.5047619
0.23333333 0.44358974 0.63333333 0.58333333 0.46666667 0.01111111 0.5 0.36666667
0.67777778 CLP_amplitudesticker4_pp33_eps1200_d10amp32 Jetzt dasselbe noch
für nur !!!einen!!! Sticker mit derselben Amplitude am selben Ort 'CLP_amplitudesticker_pp33_eps1200_d10amp32'
[0. 0.00833333 0.00133333 0.03555556 0.01666667 nan 0. 0. 0.00222222 0. 0. 0.
0.05362319 0.00138889 0.0037037 0. 0. 0.025 0. 0. 0. 0. 0. 0. 0. 0.05 0. 0.
0. 0. 0. 0. 0.00512821 0. 0. 0. 0. 0. 0. 0. ]

```

Heatmap und erzeugte Punktwolke für threshold = 0.99

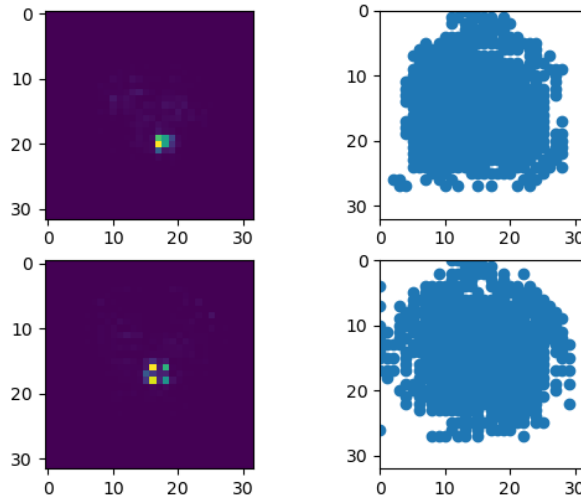


Abbildung 3: Auswahl der relevantesten Pixel (bis zu 99% der Gesamtmasse) zweier Heatmaps

Dann noch mit amp=16 wieder 4 sticker mit d=10 $CLP_{amplitude}sticker4_p33_{eps1200_d10amp16}$
 0.11666667 0.23055556 0.26 0.55333333 0.43333333 nan 0.23333333 0.28444444 0.11555556
 0.07708333 0.15 0.0952381 0.4173913 0.41527778 0.40740741 0.45238095 0.01333333
 0.35833333 0.29487179 0. 0. 0.04444444 0.05 0.01333333 0.4 0.29791667 0.13888889
 0.23333333 0.11333333 0.3 0.06 0.02962963 0.11666667 0.65714286 0.475 0.71794872
 0.50833333 0.43333333 0.35652174 0.45555556 0.27777778 0. 0.51111111 Accuracy
 on test Dataset: 0.958

Jetzt: $CLP_{amplitude}sticker4_p33_{eps1200_d10amp255}$ 0.71666667 0.78333333 0.80266667
 0.99333333 0.9469697 nan 0.96 0.97555556 0.98444444 1. 1. 0.25714286 0.98985507
 0.86111111 0.94444444 0.99047619 0.91333333 0.51666667 0.84871795 0.86666667
 0.52222222 0.57777778 0.75 0.67333333 0.95555556 0.725 0.98888889 0.88333333
 0.45333333 0.82222222 0.81333333 1. 0.98333333 0.4952381 0.26666667 0.92820513
 0.63333333 0.81666667 0.77681159 0.33333333 0.74444444 0.91666667 0.83333333
 Accuracy on test Dataset: 0.959

$CLP_{amplitude}sticker4_p33_{dist10_amp64}$ eps300 [0.8 0.98888889 0.99733333 1.
 1. nan 1. 1. 0.99555556 0.97916667 0.99545455 0.8047619 1. 0.99861111 0.92222222
 1. 0.99333333 0.98888889 0.98717949 0.93333333 0.86666667 0.93333333 0.975 0.83333333
 0.96666667 0.93958333 1. 1. 0.98 1. 0.98 0.97777778 1. 0.92380952 0.55 0.93076923
 0.70833333 0.85 0.82173913 1. 0.8 1. 1.

$CLP_{amplitude}sticker4_p33_{dist10_amp32}$ 0.7 0.86111111 0.85066667 0.81111111
 0.66666667 nan 0.42 0.72888889 0.60888889 0.77083333 0.40757576 0.46666667 0.69855072
 0.88472222 0.53333333 0.74285714 0.39333333 0.67222222 0.75128205 0.11666667
 0.58888889 0.32222222 0.425 0.18 0.44444444 0.44375 0.8 0.71666667 0.56 0.86666667
 0.11333333 0.3 0.71666667 0.5047619 0.23333333 0.44358974 0.63333333 0.58333333
 0.46666667 0.01111111 0.5 0.36666667 0.67777778 Accuracy on test Dataset: 0.959

Heatmap und erzeugte Punktwolke für threshold = 0.5

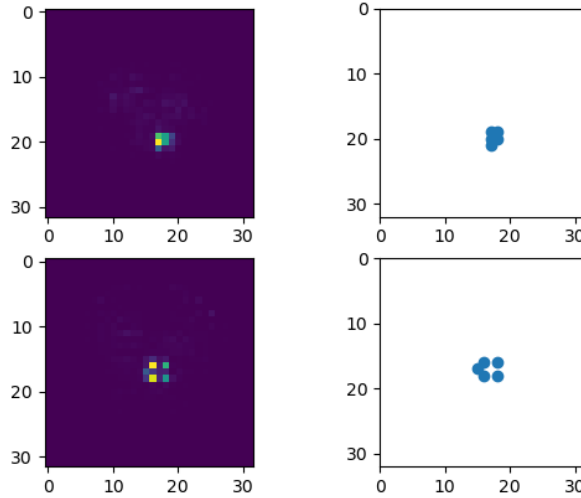


Abbildung 4: Auswahl der relevantesten Pixel (bis zu 50% der Gesamtmasse) zweier Heatmaps

CLP_amplitudesticker4_p33_dist10_amp16 0.11666667 0.23055556 0.26 0.55333333
0.43333333 nan 0.23333333 0.28444444 0.11555556 0.07708333 0.15 0.0952381 0.4173913
0.41527778 0.40740741 0.45238095 0.01333333 0.35833333 0.29487179 0. 0. 0.04444444
0.05 0.01333333 0.4 0.29791667 0.13888889 0.23333333 0.11333333 0.3 0.06 0.02962963
0.11666667 0.65714286 0.475 0.71794872 0.50833333 0.43333333 0.35652174 0.45555556
0.27777778 0. 0.51111111] Performance of poisoned net on unpoisoned training data:
loss on test dataset: 0.19637071904851583 Accuracy on test Dataset: 0.958

In Abbildung 6 sind die Angriffserfolgsraten pro Klasse im Fall von 33% korrumpierten Daten und dem Amplitudensticker in jeder der 4 Bildecken mit dem Abstand von 10 Pixeln zum Rand dargestellt. In beiden Fällen geben wir keine AER für die Klasse 6 an, über die der Angriff stattfindet. Die mittlere Angriffserfolgsrate beträgt für $amp = 255$ 79.15%. In mehreren Klassen wird eine AER von 100% erreicht.

Für $amp = 64$ fällt der Angriff mit einer mAER von 48.17% deutlich schwächer aus. Die maximal erreichte AER beträgt 95.33% in Klasse 10. Die minimalen AER sind 25% bzw. 0.83%. Eine weitere Reduktion der Amplitude auf $amp = 32$ führt zu einer mAER von 6.55% und einer maximalen AER von 33%.

Für $d=0$, $amplitude=64$ und $eps=1200$ ergibt sich kein erfolgreicher Angriff, dabei war fast alles 0, die Trigger im Testdatensatz waren auch auf 64 Jetzt für $d=10$, auch hier gilt $amp=64$ im Testdatensatz:Ergebnis:

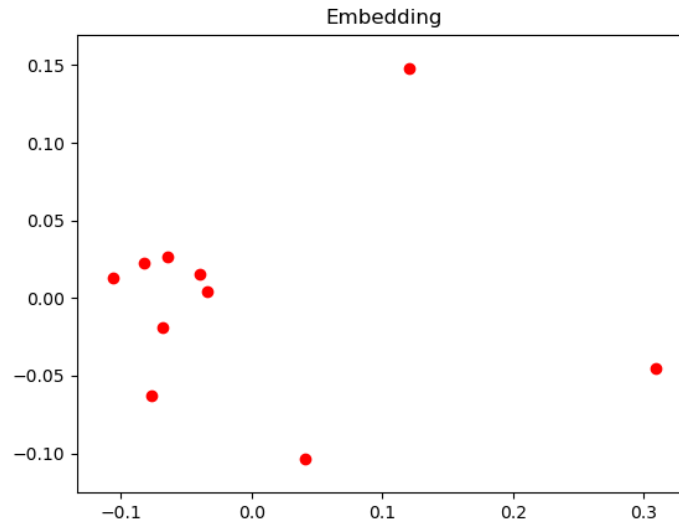


Abbildung 5: Einbettung des Baryzenters mithilfe von Multidimensionaler Skalierung(MDS) bei der Wahl von 99% der Gesamtmasse

10.1 Activation Clustering

10.2 Räumliche Transformationen

- ASR ist sehr stark vom Ort des Triggers abhängig.
- Ort des Triggers kann nicht direkt geändert werden.
- Benutze Transformationen(Flipping, Scaling), um den Trigger wirkungslos zu machen.
- Somit kann die ASR während der Inferenz verringert werden. Es lässt sich aber keine Aussage darüber treffen, ob ein Angriff vorliegt

Seitenlänge des Triggers	Prozentualer Anteil	AER	GUD	num. korruptierte Daten
s=2	0.000625	0.01333333	0.962	1
	0.00125	0.356	0.964	2
	0.0025	0.576	0.97	4
	0.005	0.99867	0.962	8
	0.01	0.92	0.962	17
	0.02	1.0	0.964	34
	0.10	1.0	0.968	291
	0.15	1.0	0.968	436
	0.33	1.0	0.968	813
s=3	?	?	?	?
	0.00125	34.26	96.2	2
	0.0025	85.07	96.5	4
	0.005	1.0	96.9	8
	0.01	1.0	0.962	17
	0.05	1.0	0.966	87
	0.15	1.0	0.961	
s=1	0.33	1.0	0.966	813

Tabelle 3: Qualität der Angriffe auf das Inception v3-Netz mit Stickern Seitenlänge 2 und 3 Pixel bei unterschiedlich großen Anteilen an korruptierten Daten

11 Weitere mögliche Schritte

- Untersuchung der Detektionsqualität in Abhängigkeit von ε
- Automatische Platzierung des Auslösers an fest gewählter Position auf dem Verkehrsschild anstatt zufälligem Platzieren in einem Fenster mit vorher festgelegter Größe. In [GDGG17] wird Faster-RCNN (F-RCNN) zur Klassifikation des LISA-Datensatzes³² benutzt. Es ist die Aufgabe, die Verkehrsschilder in die 3 Superklassen Stoppschild, Geschwindigkeitsbegrenzung und Warnschild einzuteilen. Der Datensatz enthält zudem die BoundingBoxen, sodass der Auslöser genauer angebracht werden kann.
- Verbesserte Version der Layer-wise Relevance Propagation
- Untersuchung anderer Verfahren, die die Interpretierbarkeit ermöglichen, beispielsweise: VisualBackProp: efficient visualization of CNNs³³
- Vergleich mit Cifar-10/Cifar-100 Datensatz³⁴³⁵

³²<http://cvrr.ucsd.edu/LISA/lisa-traffic-sign-dataset.html>

³³<https://arxiv.org/abs/1611.05418>

³⁴<https://www.cs.toronto.edu/~kriz/cifar.html>

³⁵<https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf>

L2-Fehler				H1-Fehler		
	Gitter	P=1	P=2	Gitter	P=1	P=2
T=2.0	8×8	2.2e-15	1.6e-14	8×8	3.7e-14	1.7e-13
	16×16	1.4e-14	2.5e-13	16×16	1.0e-13	1.2e-12
	32×32	6.1e-15	1.1e-14	32×32	1.6e-14	6.0e-13

Tabelle 4: Fehler für Testproblem 1 zum Endzeitpunkt $T = 2.0$.

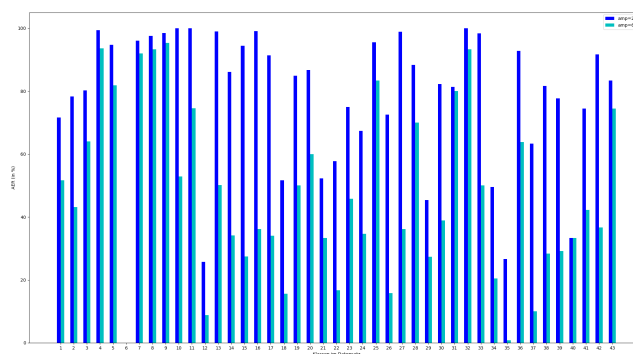


Abbildung 6: Angriffserfolgsrate pro Klasse bei Clean-Label-Poisoning-Attacks für verschiedene Werte amp des Amplitudenstickers in vierfacher Ausführung bei Abstand $d = 10$ zum Rand

12 Zusammenfassung und Ausblick

Beobachtung: Je größer die Netzwerke sind, desto leichter lassen sich Poisoning-Angriffe realisieren.

A Verwendete Netzwerke

A.1 Net

```

1  class Net(nn.Module):
2
3
4  def __init__(self, ):
5      super(Net, self).__init__()
6      self.size = 64 * 4 * 4
7      self.conv1 = nn.Conv2d(in_channels=3, out_channels=12,
8                              kernel_size=5, padding=2)
9      self.pool = nn.MaxPool2d(kernel_size=2, stride=2)
10     self.conv1_in = nn.InstanceNorm2d(12)
11     self.conv2 = nn.Conv2d(in_channels=12, out_channels=32,
12                             kernel_size=5, padding=2)
13
14     self.conv2_bn = nn.BatchNorm2d(32)
15
16     self.conv3 = nn.Conv2d(in_channels=32, out_channels=64,
17                             kernel_size=5, padding=2)
18
19     self.fc1 = nn.Linear(self.size, 256)
20     self.fc1_bn = nn.BatchNorm1d(256)
21     self.fc2 = nn.Linear(256, 128)
22     self.fc3 = nn.Linear(128, 43)
23
24
25 def forward(self, x):
26     x = self.pool(F.relu(self.conv1_in(self.conv1(x))))
27     x = self.pool(F.relu(self.conv2_bn(self.conv2(x))))
28     x = self.pool(F.relu(self.conv3(x)))
29     x = x.view(-1, self.size)
30     x = F.relu(self.fc1_bn(self.fc1(x)))
31     x = F.dropout(x)
32     xx = F.relu(self.fc2(x))
33     x = F.dropout(xx)
34     x = self.fc3(x)
35
36     return x, xx

```

Listing 4: Kleines Netzwerk

```

1  InceptionNet3(
2      (features): Sequential(
3          (0): InceptionA(
4              (parallel_dummyA): New_parallel_chain_dummy()
5              (conv1x1): BatchConv(
6                  (conv): Conv2d(3, 64, kernel_size=(1, 1), stride=(1, 1))
7                  (bn): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True,
8                      track_running_stats=True)
9                  (relu): ReLU()

```



```

9      )
10     (parallel_dummyB): New_parallel_chain_dummy()
11     (conv5x5_1): BatchConv(
12     (conv): Conv2d(3, 48, kernel_size=(1, 1), stride=(1, 1))
13     (bn): BatchNorm2d(48, eps=1e-05, momentum=0.1, affine=True,
        track_running_stats=True)
14     (relu): ReLU()
15     )
16     (conv5x5_2): BatchConv(
17     (conv): Conv2d(48, 64, kernel_size=(5, 5), stride=(1, 1),
        padding=(2, 2))
18     (bn): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True,
        track_running_stats=True)
19     (relu): ReLU()
20     )
21     (parallel_dummyC): New_parallel_chain_dummy()
22     (conv3x3dbl_1): BatchConv(
23     (conv): Conv2d(3, 64, kernel_size=(1, 1), stride=(1, 1))
24     (bn): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True,
        track_running_stats=True)
25     (relu): ReLU()
26     )
27     (conv3x3dbl_2): BatchConv(
28     (conv): Conv2d(64, 96, kernel_size=(3, 3), stride=(1, 1),
        padding=(1, 1))
29     (bn): BatchNorm2d(96, eps=1e-05, momentum=0.1, affine=True,
        track_running_stats=True)
30     (relu): ReLU()
31     )
32     (conv3x3dbl_3): BatchConv(
33     (conv): Conv2d(96, 96, kernel_size=(3, 3), stride=(1, 1),
        padding=(1, 1))
34     (bn): BatchNorm2d(96, eps=1e-05, momentum=0.1, affine=True,
        track_running_stats=True)
35     (relu): ReLU()
36     )
37     (parallel_dummyD): New_parallel_chain_dummy()
38     (pool): MaxPool2d(kernel_size=3, stride=1, padding=1,
        dilation=1, ceil_mode=False)
39     (pool1x1): BatchConv(
40     (conv): Conv2d(3, 32, kernel_size=(1, 1), stride=(1, 1))
41     (bn): BatchNorm2d(32, eps=1e-05, momentum=0.1, affine=True,
        track_running_stats=True)
42     (relu): ReLU()
43     )
44     (parallel_dummyE): New_parallel_chain_dummy()
45     (cat): Cat()
46     )
47     (1): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation
        =1, ceil_mode=False)
48     (2): BatchConv(
49     (conv): Conv2d(256, 256, kernel_size=(2, 2), stride=(1, 1))
50     (bn): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True,
        track_running_stats=True)

```

```

51 (relu): ReLU()
52 )
53 (3): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation
    =1, ceil_mode=False)
54 (4): BatchConv(
55 (conv): Conv2d(256, 256, kernel_size=(2, 2), stride=(1, 1))
56 (bn): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True,
    track_running_stats=True)
57 (relu): ReLU()
58 )
59 (5): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation
    =1, ceil_mode=False)
60 (6): BatchConv(
61 (conv): Conv2d(256, 256, kernel_size=(2, 2), stride=(1, 1))
62 (bn): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True,
    track_running_stats=True)
63 (relu): ReLU()
64 )
65 (7): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation
    =1, ceil_mode=False)
66 )
67 (classifiers): Sequential(
68 (0): Linear(in_features=256, out_features=256, bias=True)
69 (1): ReLU(inplace=True)
70 (2): Dropout(p=0.5, inplace=False)
71 (3): Linear(in_features=256, out_features=128, bias=True)
72 (4): ReLU(inplace=True)
73 (5): Dropout(p=0.5, inplace=False)
74 (6): Linear(in_features=128, out_features=43, bias=True)
75 )
76 )
77

```

Listing 5: Einfachere Version von Inception v3

Aufbau dieses Netzwerkes: 1. Inception-Modul 2. [pool1, batchConv1, pool2, batchConv2, pool3, batchConv3, pool4] 3. Drei Lineare Schichten mit ReLu und Dropout dazwischen

Im Unterschied zum offiziellen Inception Netz(v1v2v3) gibt es in dieser vereinfachten Version keinen `stem` aus convs, es geht direkt mit InceptionA los.

Wie ähnlich sind sich InceptionA(hier) und das offizielle InceptionA-Modul?

```

1
2 [Linear(in_features=128, out_features=43, bias=True),
3  Dropout(p=0.5, inplace=False),
4  ReLU(inplace=True),
5  Linear(in_features=256, out_features=128, bias=True),
6  Dropout(p=0.5, inplace=False),
7  ReLU(inplace=True),
8  Linear(in_features=256, out_features=256, bias=True),
9  MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1,
    ceil_mode=False),
10 ReLU(),
11 BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True,
    track_running_stats=True), Conv2d(256, 256, kernel_size=(2,

```

```

2), stride=(1, 1)),
12 MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1,
    ceil_mode=False),
13 ReLU(),
14 BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True,
    track_running_stats=True), Conv2d(256, 256, kernel_size=(2,
    2), stride=(1, 1)),
15 MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1,
    ceil_mode=False),
16 ReLU(), BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True
    , track_running_stats=True),
17 Conv2d(256, 256, kernel_size=(2, 2), stride=(1, 1)),
18 MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1,
    ceil_mode=False),
19
20 [[Conv2d(3, 64, kernel_size=(1, 1), stride=(1, 1)),
21 BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True,
    track_running_stats=True), ReLU()],
22
23 [Conv2d(3, 48, kernel_size=(1, 1), stride=(1, 1)),
24 BatchNorm2d(48, eps=1e-05, momentum=0.1, affine=True,
    track_running_stats=True), ReLU(), Conv2d(48, 64,
    kernel_size=(5, 5), stride=(1, 1), padding=(2, 2)),
    BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True,
    track_running_stats=True), ReLU()],
25
26 [Conv2d(3, 64, kernel_size=(1, 1), stride=(1, 1)),
27 BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True,
    track_running_stats=True), ReLU(),
28 Conv2d(64, 96, kernel_size=(3, 3), stride=(1, 1), padding=(1,
    1)),
29 BatchNorm2d(96, eps=1e-05, momentum=0.1, affine=True,
    track_running_stats=True), ReLU(),
30 Conv2d(96, 96, kernel_size=(3, 3), stride=(1, 1), padding=(1,
    1)),
31 BatchNorm2d(96, eps=1e-05, momentum=0.1, affine=True,
    track_running_stats=True), ReLU()],
32
33 [MaxPool2d(kernel_size=3, stride=1, padding=1, dilation=1,
    ceil_mode=False), Conv2d(3, 32, kernel_size=(1, 1), stride
    =(1, 1)),
34 BatchNorm2d(32, eps=1e-05, momentum=0.1, affine=True,
    track_running_stats=True),
35 ReLU()],
36
37 [Cat()]]
38

```

Listing 6: Reversed Model incv3

B Parameter für Training und Einlesen der Daten

Die in [AWN⁺20] gewählten Parameter wären ein guter Ausgangspunkt.

Für das Einlesen der Daten benutzen wir, sofern nicht weiter angegeben die folgen-

den Augmentierungen:

```

1  __train_transform = transforms.Compose(
2  [
3  transforms.RandomResizedCrop((image_size, image_size),
4  scale=(0.6, 1.0)),
5  transforms.RandomRotation(degrees=15),
6  transforms.ColorJitter(brightness=0.1, contrast=0.1,
7  saturation=0.1, hue=0.1),
8  transforms.RandomAffine(15),
9  transforms.RandomGrayscale(),
10 transforms.Normalize(mean=[0.485, 0.456, 0.406],
11 std=[0.229, 0.224, 0.225]),
12 transforms.ToTensor())
13
14 ]
15
16
17
18

```

Listing 7: Augmentierung beim Einlesen der Daten

Die Werte von mean und std variieren für alle ausgeführten Poisoning-Angriffe. Anstatt beide jedes Mal erneut zu berechnen, verwenden wir die von pytorch angegebenen Werte ³⁶, die für die vor-trainierten Modelle empfohlen werden und auf dem Datensatz ImageNet³⁷ basieren.

Wir trainieren die Netzwerke über maximal 100 Epochen und benutzen *early stopping* mit einer *patience* = 20. Die verwendete Implementierung ist eine modifizierte Version von Bjarte Mehus Sunde ³⁸, die wiederum auf PyTorch Ignite³⁹ basiert.

C Einlesen der Daten bei AC

Ohne Transformationen, wie den Testdatensatz.

D Parameter für die ausgeführten Angriffe

Für das projizierte Gradientenverfahren benutzen wir 10 Iterationen und eine Schrittweite von 0.015.

Label-konsistente Poisoning-Angriffe:

³⁶<https://github.com/pytorch/examples/blob/97304e232807082c2e7b54c597615dc0ad8f6173/imagenet/main.py#L197-L198>

³⁷<https://image-net.org/>

³⁸<https://github.com/Bjarten/early-stopping-pytorch>

³⁹https://github.com/pytorch/ignite/blob/master/ignite/handlers/early_stopping.pyt

E Datensätze

GTSRB⁴⁰ Datensatz Splitting (Train; Val, test)

ImageNet besteht über 14 Millionen Bildern in 100 Klassen.

F Programmcode

Der vollständige Programmcode ist verfügbar unter <https://github.com/lukasschulth/MA-Detection-of-Poisoning-Attacks>

G Notizen

registered spaces⁴¹ Barycenters in the Wasserstein Space⁴²

Was passiert bei der Kombination von 2 verschiedenen Triggern? Einmal keine Überlappung(d.h. 2verschiedene Trigger auf dem selben Bild) vs. auch beide Trigger auf einem Bild ist zulässig.

⁴⁰https://benchmark.ini.rub.de/gtsrb_dataset.html

⁴¹<https://arxiv.org/pdf/1809.06422.pdf>

⁴²<https://arxiv.org/pdf/1809.06422.pdf>

Literatur

- [AC11] Martial Agueh and Guillaume Carlier. Barycenters in the wasserstein space. *SIAM Journal on Mathematical Analysis*, 43(2):904–924, 2011.
- [ACB17] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *International conference on machine learning*, pages 214–223. PMLR, 2017.
- [AMN⁺19] Christopher J Anders, Talmaj Marinč, David Neumann, Wojciech Samek, Klaus-Robert Müller, and Sebastian Lapuschkin. Analyzing imagenet with spectral relevance analysis: Towards imagenet un-hans’ ed. *arXiv preprint arXiv:1912.11425*, 2019.
- [AWN⁺19] Christopher J. Anders, Leander Weber, David Neumann, Wojciech Samek, Klaus-Robert Müller, and Sebastian Lapuschkin. Finding and removing clever hans: Using explanation methods to debug and improve deep models, 2019.
- [AWN⁺20] Christopher J Anders, Leander Weber, David Neumann, Wojciech Samek, Klaus-Robert Müller, and Sebastian Lapuschkin. Finding and removing clever hans: Using explanation methods to debug and improve deep models. 2020.
- [AWR17] Jason Altschuler, Jonathan Weed, and Philippe Rigollet. Near-linear time approximation algorithms for optimal transport via sinkhorn iteration. *arXiv preprint arXiv:1705.09634*, 2017.
- [BBB⁺01] Dmitri Burago, Iu D Burago, Yuri Burago, Sergei Ivanov, Sergei V Ivanov, and Sergei A Ivanov. *A course in metric geometry*, volume 33. American Mathematical Soc., 2001.
- [BBM⁺15] Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PloS one*, 10(7):e0130140, 2015.
- [BBM⁺16] Alexander Binder, Sebastian Bach, Gregoire Montavon, Klaus-Robert Müller, and Wojciech Samek. Layer-wise relevance propagation for deep neural network architectures. In Kuinam J. Kim and Nikolai Joukov, editors, *Information Science and Applications (ICISA) 2016*, pages 913–922, Singapore, 2016. Springer Singapore.
- [BCC⁺15] Jean-David Benamou, Guillaume Carlier, Marco Cuturi, Luca Nenna, and Gabriel Peyré. Iterative bregman projections for regularized transportation problems. *SIAM Journal on Scientific Computing*, 37(2):A1111–A1138, 2015.
- [BCL16] Radu Ioan Boț, Ernő Robert Csetnek, and Szilárd Csaba László. An inertial forward–backward algorithm for the minimization of the sum of two nonconvex functions. *EURO Journal on Computational Optimization*, 4(1):3–25, 2016.
- [BEWR19] Moritz Böhle, Fabian Eitel, Martin Weygandt, and Kerstin Ritter. Layer-wise relevance propagation for explaining deep neural network

- decisions in mri-based alzheimer’s disease classification. *Frontiers in aging neuroscience*, 11:194, 2019.
- [CCB⁺18] Bryant Chen, Wilka Carvalho, Nathalie Baracaldo, Heiko Ludwig, Benjamin Edwards, Taesung Lee, Ian Molloy, and Biplav Srivastava. Detecting backdoor attacks on deep neural networks by activation clustering. *arXiv preprint arXiv:1811.03728*, 2018.
- [CDPS17] Guillaume Carlier, Vincent Duval, Gabriel Peyré, and Bernhard Schmitzer. Convergence of entropic schemes for optimal transport and gradient flows. *SIAM Journal on Mathematical Analysis*, 49(2):1385–1418, 2017.
- [CFTR16] Nicolas Courty, Rémi Flamary, Devis Tuia, and Alain Rakotomamonjy. Optimal transport for domain adaptation. *IEEE transactions on pattern analysis and machine intelligence*, 39(9):1853–1865, 2016.
- [com19] Computational optimal transport. *Foundations and Trends in Machine Learning*, 11(5-6):355–607, 2019.
- [COT19] Computational optimal transport. *Foundations and Trends in Machine Learning*, 11(5-6):355–607, 2019.
- [Cut13] Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transportation distances. *arXiv preprint arXiv:1306.0895*, 2013.
- [DGK18] Pavel Dvurechensky, Alexander Gasnikov, and Alexey Kroshnin. Computational optimal transport: Complexity by accelerated gradient descent is better than by sinkhorn’s algorithm. In *International conference on machine learning*, pages 1367–1376. PMLR, 2018.
- [FL89] Joel Franklin and Jens Lorenz. On the scaling of multidimensional matrices. *Linear Algebra and its applications*, 114:717–735, 1989.
- [GDGG17] Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. Badnets: Identifying vulnerabilities in the machine learning model supply chain. *arXiv preprint arXiv:1708.06733*, 2017.
- [GPAM⁺14] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014.
- [Gro07] Mikhail Gromov. *Metric structures for Riemannian and non-Riemannian spaces*. Springer Science & Business Media, 2007.
- [IS15] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. PMLR, 2015.
- [LS14] Yin Tat Lee and Aaron Sidford. Path finding methods for linear programming: Solving linear programs in \tilde{O} (vrank) iterations and faster algorithms for maximum flow. In *2014 IEEE 55th Annual Symposium on Foundations of Computer Science*, pages 424–433. IEEE, 2014.

-
- [LWB⁺19] Sebastian Lapuschkin, Stephan Wäldchen, Alexander Binder, Grégoire Montavon, Wojciech Samek, and Klaus-Robert Müller. Unmasking clever hans predictors and assessing what machines really learn. *Nature communications*, 10(1):1–8, 2019.
 - [Mau] Abhinav Maurya. Optimal transport in statistical machine learning: Selected review and some open questions.
 - [MBL⁺19] Grégoire Montavon, Alexander Binder, Sebastian Lapuschkin, Wojciech Samek, and Klaus-Robert Müller. Layer-wise relevance propagation: an overview. *Explainable AI: interpreting, explaining and visualizing deep learning*, pages 193–209, 2019.
 - [Mém11] Facundo Mémoli. Gromov–wasserstein distances and the metric approach to object matching. *Foundations of computational mathematics*, 11(4):417–487, 2011.
 - [MLB⁺17a] Grégoire Montavon, Sebastian Lapuschkin, Alexander Binder, Wojciech Samek, and Klaus-Robert Müller. Explaining nonlinear classification decisions with deep taylor decomposition. *Pattern Recognition*, 65:211–222, 2017.
 - [MLB⁺17b] Grégoire Montavon, Sebastian Lapuschkin, Alexander Binder, Wojciech Samek, and Klaus-Robert Müller. Explaining nonlinear classification decisions with deep taylor decomposition. *Pattern Recognition*, 65:211–222, 2017.
 - [MMG19] Anton Mallasto, Guido Montúfar, and Augusto Gerolin. How well do wgans estimate the wasserstein metric? *arXiv preprint arXiv:1910.03875*, 2019.
 - [MMS⁺17] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.
 - [NR99] Arkadi Nemirovski and Uriel Rothblum. On complexity of matrix scaling. *Linear Algebra and its Applications*, 302:435–460, 1999.
 - [PCS16] Gabriel Peyré, Marco Cuturi, and Justin Solomon. Gromov–wasserstein averaging of kernel and distance matrices. In *International Conference on Machine Learning*, pages 2664–2672. PMLR, 2016.
 - [PMG16] Nicolas Papernot, Patrick McDaniel, and Ian Goodfellow. Transferability in machine learning: from phenomena to black-box attacks using adversarial samples. *arXiv preprint arXiv:1605.07277*, 2016.
 - [Ren88] James Renegar. A polynomial-time algorithm, based on newton’s method, for linear programming. *Mathematical programming*, 40(1):59–93, 1988.
 - [RTG00] Yossi Rubner, Carlo Tomasi, and Leonidas J Guibas. The earth mover’s distance as a metric for image retrieval. *International journal of computer vision*, 40(2):99–121, 2000.
 - [RYM99] Anand Rangarajan, Alan Yuille, and Eric Mjolsness. Convergence properties of the softassign quadratic assignment algorithm. *Neural Computation*, 11(6):1455–1474, 1999.

- [San10] Filippo Santambrogio. Introduction to optimal transport theory. *arXiv preprint arXiv:1009.3856*, 2010.
- [Sin64] Richard Sinkhorn. A relationship between arbitrary positive matrices and doubly stochastic matrices. *The annals of mathematical statistics*, 35(2):876–879, 1964.
- [SLJ⁺15] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1–9, 2015.
- [SZS⁺13] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.
- [TLM18] Brandon Tran, Jerry Li, and Aleksander Madry. Spectral signatures in backdoor attacks. *arXiv preprint arXiv:1811.00636*, 2018.
- [TTM19] Alexander Turner, Dimitris Tsipras, and Aleksander Madry. Label-consistent backdoor attacks, 2019.
- [VCF⁺20] Titouan Vayer, Laetitia Chapel, Rémi Flamary, Romain Tavenard, and Nicolas Courty. Fused gromov-wasserstein distance for structured objects. *Algorithms*, 13(9):212, 2020.
- [Vil03] Cédric Villani. *Topics in optimal transportation*. Number 58. American Mathematical Soc., 2003.
- [VL07] Ulrike Von Luxburg. A tutorial on spectral clustering. *Statistics and computing*, 17(4):395–416, 2007.