



UPPSALA
UNIVERSITET

UPTEC F 22055

Examensarbete 30 hp

September 2022

Unsupervised Image Classification Using Domain Adaptation

Via the Second Order Statistic

Finn Joel Bjervig



Abstract

Framgången inom maskininlärning och djupinlärning beror till stor del på stora annoterade dataset, som MNIST, ImageNet, Caltech-256, och Cifar-10. Att tilldela etiketter till data är väldigt resurskrävande och kan till viss del undvikas genom att utnyttja datans statistiska egenskaper. En maskininlärningsmodell kan lära sig att klassificera bilder från en *domän* utifrån träningsexempel som innehåller bilder, samt etiketter som berättar vad bilder föreställer. Men vad gör man om datan inte har tilldelade etiketter? En maskininlärningsmodell som lär sig en uppgift utifrån annoterad data från en *källdomän*, kan med hjälp av information från *måldomänen* (som inte har tilldelade etiketter), anpassas till att prestera bättre på data från måldomänen. Forskningsområdet som studerar hur man anpassar och generaliseras en modell mellan två olika domäner heter domänanpassning, eller *domain adaptation*. Detta examensarbete är utfört på Scanias forskningsavdelning för autonom transport, som har försett projektet med kamera och LiDAR bilder tagna från deras testfordon. Examensarbetet handlar om hur modeller för bildklassificering som tränas på kamerabilder med etiketter, kan anpassas till att få ökad noggrannhet på ett dataset med LiDAR bilder, som inte har etiketter. Två metoder för domänanpassning har jämförts med varandra, samt med en model tränad på kameradata genom övervakad inlärning utan domänanpassning. Alla metoder opererar på något vis med ett djupt faltningsnätverk (CNN) där uppgiften är att klassificera bilder utav bilar eller fotgängare. Kovariansen utav datan från käll- och måldomänen är det centrala måttet för domänanpassningsmetoderna i detta projekt. Den första metoden är en så kallad *ytlig* metod, där själva anpassningsmetoden inte ingår inuti den djupa arkitekturen av modellen, utan är ett mellansteg i processen. Den andra metoden förenar domänanpassningsmetoden med klassificeringen i Viden djupa arkitekturen. Den tredje modellen består endast utav faltningsnätverket, utan en metod för domänanpassning. Resultaten visar att utan en domänanpassningsmetod minskar klassificeringsnoggrannheten i måldomänen till 63.80 % under träningen, medan den ökar i källdomänen. För binära klassificeringsmodeller, betyder slumpmässig klassificering 50 % noggrannhet. Den *ytliga* metoden upprätthåller klassificeringsnoggrannheten i måldomänen till 74.67 % och den djupa metoden presterar bäst med 80.73 % noggrannhet. Modellen utan domänanpassning klassificerar dessutom fotgängare felaktigt som fordon mycket mer ofta än de andra modellerna, som istället är mer objektiv och även mer noggranna. Resultaten visar att det är möjligt att anpassa en modell som tränas på data från källdomänen, till att få ökad klassificeringsnoggrannhet i måldomänen genom att använda kovariansen utav datan från de två domänerna. Den djupa metoden för domänanpassning tillåter även användandet utav andra statistiska mått som kan vara mer framgångsrika i att generalisera modellen, beroende på hur datan är fördelad. Överlägsenheten hos den djupa metoden antyder att domänanpassning med fördel kan bättas in i den djupa arkitekturen så att modellparametrarna blir uppdaterade för att lära sig en mer robust representation utav måldomänen.

Teknisk-naturvetenskapliga fakulteten

Uppsala universitet, Utgivningsort Uppsala/Visby

Handledare: Arash Owrang Åmnesgranskare: Dominik Baumann

Examinator: Tomas Nyberg

Acknowledgements

First, I must mention my thesis partner Mikael Westlund at KTH, and his excellent collaborative spirit. We parted ways halfway through the project to focus on different methods but kept in touch. I would like to thank my thesis supervisor Arash Owrang Ph.D., with whom I had many insightful discussions. I also want to recognize the perception team at the Scania Autonomous Transport Research department for being very supportive and friendly during the project. I am very grateful of my subject reader Dominik Baumann, Ph.D. from Uppsala University, who was always keen to discuss with me and answer my questions. And finally, I would like to thank my friend and colleague Nils Guamelius, who worked on another Master's thesis project at the same department at Scania. We spent a lot of time together, working on our separate theses, and supporting each other.

Contents

| | | |
|-----------|--|-----------|
| 1 | Introduction | 1 |
| 2 | Related Work | 1 |
| 2.1 | Shallow Domain Adaptation | 1 |
| 2.1.1 | Unsupervised Feature Transformation Methods | 2 |
| 2.2 | Deep Domain Adaptation | 2 |
| 2.2.1 | Discrepancy-Based Methods | 2 |
| 2.3 | Adversarial Generative Models | 2 |
| 3 | Problem Description | 3 |
| 4 | Transfer Learning | 3 |
| 4.1 | Transfer Learning & Domain Adaptation | 3 |
| 5 | CORAL | 4 |
| 6 | Deep CORAL | 5 |
| 7 | Data Acquisition and Preprocessing | 5 |
| 7.1 | Cifar-10 | 6 |
| 7.2 | Camera Dataset | 6 |
| 7.3 | LiDAR Dataset | 7 |
| 8 | Experimental Evaluation | 10 |
| 8.1 | Convolutional Neural Network Architecture | 10 |
| 8.2 | Pretraining Convolutional Neural Network on Cifar-10 | 11 |
| 8.2.1 | For Shallow CORAL Method | 11 |
| 8.2.2 | For Finetuned Model & Deep CORAL Model | 11 |
| 8.3 | CORAL Model | 12 |
| 8.3.1 | Pretraining CNN | 12 |
| 8.4 | Deep Model without CORAL Loss | 12 |
| 8.5 | Deep CORAL Model | 14 |
| 8.5.1 | Hyperparameter Tuning | 14 |
| 9 | Results | 15 |
| 10 | Future Work | 17 |
| 11 | Conclusions | 18 |
| A | Definitions | 19 |
| B | Algorithms and Functions | 19 |
| B.1 | Stochastic Gradient Descent | 19 |
| B.2 | Cross-Entropy Loss Function | 20 |

1 Introduction

There are vast amounts of data being collected and stored every day, and for supervised learning tasks, annotated data is needed. Data annotation is often outsourced to companies that specialize in such tasks, and while it produces valuable datasets, it is resource costly and a relatively inefficient process. To circumvent the manual labor of annotation, one may adopt approaches that instead utilize the intrinsic structure of the data, and thus alleviate the dependency on labels. Such methods go under the category of unsupervised machine learning, which experiences only data features from which it learns structural properties [1, Chapter 5]. In recent years, the prospect of transfer learning has emerged, which seeks to leverage labeled data or utilize what has been learned in one setting, with some distribution, to improve the performance on another dataset with a similar distribution. Measuring the “similarity” of datasets is a central area of research as it tells whether there exist any discrepancies between the two datasets. For example, if two experiments are supposedly executed by the same protocol, but the collected data from each experiment show a significant discrepancy, it signifies a difference in measuring equipment, techniques, or the environment in which the experiment was conducted. If this is the case, they should not be jointly analyzed in a study (before being transformed in any way such that they share the same distribution). When the source and target domains are represented by the same feature space, but not by the same marginal distribution, it is called homogeneous transfer learning. The difference in marginal distributions is then the cause for what is called a domain shift. Further, if the type of prediction or inference being made (i.e. the task) is the same, we enter the realm of domain adaptation. A domain consists of a feature space and a marginal distribution. In visual applications, a feature space may be all 32×32 RGB images, where every pixel can take a value between 0 and 255. The marginal distribution informs how the dataset is distributed. In this thesis, the second order statistics, or the covariance, is the central measure by which the domain adaptation methods operate with. Pretraining a convolutional neural network on the Cifar-10 dataset is used to give an improved weight initialization on which the DA methods will be performed. Both a so-called shallow method and a deep method are compared to a supervised learning method and put against each other with respect to their classification accuracy in the target domain.

2 Related Work

The field of domain adaptation (DA) has a diverse set of methods, both shallow and deep, with the aim to generalize or adapt a model to increase task performance from one dataset to another. The most established DA methods are presented in [2] from 2017. There exists several benchmark datasets that can be combined to form multiple domain shift scenarios.

2.1 Shallow Domain Adaptation

CORAL [3] and Subspace Alignment (SA) [4] are both feature space alignment methods, where CORAL aligns the feature spaces directly via the second order moment of the data, and SA aims to align domain-specific subspaces. More specifically the source data is projected onto the source subspace and the target data onto the target subspace. These subspaces are spanned by a set of d eigenvectors corresponding to the d largest eigenvalues obtained from Principal Component Analysis (PCA) of the source and target data, respectively. The linear transformation \mathbf{A} aligns the basis vectors of the subspaces such that the Bregman divergence (see definition 4 in Appendix A) between the two subspaces is minimized.

2.1.1 Unsupervised Feature Transformation Methods

An alternative to feature alignment methods, which aims to transform the source domain to the target domain, are feature transformation methods. There a transformation is applied on both the source and target domain to minimize the discrepancy between the two distributions. Many methods rely on minimizing some type of divergence or distance measure between sample means and variances, but the Domain Invariant Projection (DIP) in [5] compares the data directly in the reproducing kernel Hilbert space (RKHS) via the maximum mean discrepancy (MMD).

2.2 Deep Domain Adaptation

Models with deep architectures in general perform very well in discriminating complex data, be it binary or multi-class. The domain adaptation survey [2] highlights the ability of some deep models, with no domain adaptation, to outperform some shallow domain adaptation methods on several benchmarks. The ability of deep models to learn robust general features is very valuable and proves the power of deep architectures. Deep domain adaptation seeks to further boost the generalizability in either a semi-supervised or unsupervised manner. Semi-supervised domain adaptation can, for example, mix the source data with a few labeled target instances. When the target data has no labels, its statistical properties are used to minimize the domain shift of the deep features in the deep model.

2.2.1 Discrepancy-Based Methods

Inspired by shallow feature transformations, discrepancy-based methods operate in a deep learning framework and use a statistical distance measure between corresponding activation layers of two Siamese architectures. The classification loss can be combined with the statistical loss such that the neural network learns a nonlinear transformation that minimizes the domain shift between the source and target domain, while also discriminating between the classes. Deep Domain Confusion (DDC) proposed in [6], uses the Krizhevsky architecture [7] and adds a lower-dimensional, “bottleneck” adaptation layer after the fully connected layer fc7. A MMD loss is applied to the adaption layer to regularize the representation to become invariant to the source and target data. DDC uses one layer for computing the MMD loss, but it is possible to apply the practice to several layers. Extending the idea of MMD, [8] proposes a Deep Adaptation Network (DAN) which incorporates the sum of MMD measures from several fully connected layers. Further, the authors explore many kernels, such as linear, polynomial, and the radial basis function (RBF). The resulting method outperforms DDC in all domain shift scenarios presented in [8], presumably because of the use of multiple layers instead of one.

2.3 Adversarial Generative Models

Models under this category rely on a generative adversarial network (GAN) architecture [9]. The generator produces artificial data instances that the discriminator tries to distinguish from the training data. The discriminator tells the generator (via the generator loss) that the artificial data instances are not plausible enough, by which the model parameters of the generator are updated via backpropagation. The discriminator loss updates the model parameters of the discriminator such that it improves at separating generated data from real training data. In the DA field, GAN can be used to facilitate domain confusion. A coupled GAN (CoGAN) framework is proposed by [10]. Each GAN generates images in their respective domain, and with weight sharing constraints, the model learns a joint distribution in the two domains from images separately drawn from the

marginal distributions of each individual domain. Unsupervised Pixel-level Domain Adaptation with GAN (PixelDA) [11] has a generator function that maps a source image and a noise vector, to an “adapted” image that looks as if it was drawn from the target domain. Along with the source labels, the adapted source images create an adapted dataset on which a classifier can be learned in a supervised manner. In contrast to discrepancy-based methods, which join the DA task and the discriminate task in one single network, this GAN framework decouples the domain adaptation task from the discriminative task, allowing users to focus on training a classifier after the adapted dataset has been generated.

3 Problem Description

Consider the task of training a model to discriminate and classify images from a *target* dataset \mathcal{D}_t that has no accompanying labels. Without labels, a supervised learning approach is not possible. However, it is possible to adapt a model to \mathcal{D}_t by training a model on a similar source dataset \mathcal{D}_s with labels, while incorporating statistical knowledge from \mathcal{D}_t . In this project the source domain consists of annotated camera images, and the target domain contains LiDAR intensity images with no labels. Despite the two domains portraying similar scenes (streets, highways, traffic events) and depicting a lot of the same objects (pedestrians, bicyclists, cars, heavy vehicles), there is still a non-trivial discrepancy, or *domain shift* between the two domains. The task is to minimize this domain shift such that we can use the supervised learning task in the source domain, alongside statistical information from the target domain to yield a image classifier that performs well on the LiDAR dataset.

4 Transfer Learning

The concept of transfer of learning is traced back to the paper *The relation of special training and general intelligence* [12], published by the educational psychologist Charles Hubbard Judd in 1908 in the journal *Educational review*. The paper presents an experiment conducted several years earlier in which fifth- and sixth-grade students were given the task of throwing darts at a target submerged in 12 inches of water. The treatment group was taught about the theory of light refraction that made the underwater dartboard appear skewed. There was no difference between the accuracy of the control and treatment groups when the water depth was 12 inches. However, when the water depth was changed to four inches, the treatment group trained in refraction theory performed better than the control group. Judd concluded that information from one practice: the theory of physics, had improved the treatment group’s performance in another practice: throwing darts. Similarly in machine learning, it is possible to transfer knowledge from a *source domain* to a *target domain* to increase its task performance in the target domain.

4.1 Transfer Learning & Domain Adaptation

Following the notation and definitions of [13] and [14], the definitions of a domain and a task should be presented before defining transfer learning.

Definition 1 (Domain) *A domain is comprised by two parts, a feature space $\mathcal{X} \in \mathbb{R}^d$ and a marginal distribution $P(\mathbf{X})$. The domain can be written as $\mathcal{D} = \{\mathcal{X}, P(\mathbf{X})\}$, where \mathbf{X} is an instance set defined as $\mathbf{X} = \{\mathbf{x} \mid \mathbf{x}_i \in \mathcal{X}, \quad i = 1, 2, \dots, n\}$.*

Definition 2 (Task) A task is comprised by a label space $\mathcal{Y} \in \mathbb{N}$ and a decision function $f : \mathbb{R} \mapsto \mathbb{N}$. The task can be written as $\mathcal{T} = \{\mathcal{Y}, f\}$. The decision function is implicit and learned from experiencing corresponding features.

In some cases, the decision function predicts the conditional distribution itself, i.e., $f(\mathbf{X}) = P(\mathbf{Y}|\mathbf{X})$, and in a classification task, the function predicts a label $f(\mathbf{X}) = \mathbf{y}_{pred}$. The prediction can be a binary vector of dimension n , or a logit $\mathbf{y}_{pred} \in 1, 2, \dots, n$. Consider two domains, a *source* domain $\mathcal{D}_s = \{\mathcal{X}_s, P(\mathbf{X}_s)\}$ with $\mathcal{T}_s = \{\mathcal{Y}_s, P(\mathbf{Y}_s|\mathbf{X}_s)\}$ and a *target* domain $\mathcal{D}_t = \{\mathcal{X}_t, P(\mathbf{X}_t)\}$ with $\mathcal{T}_t = \{\mathcal{Y}_t, P(\mathbf{Y}_t|\mathbf{X}_t)\}$. Now imagine that the target domain has no label set. If the source and target domains are equal, $\mathcal{D}_s = \mathcal{D}_t$ and the tasks are equal, $\mathcal{T}_s = \mathcal{T}_t$, then a supervised machine learning model trained on the source domain would also perform well in the target domain. However, if there exists a *domain shift*, i.e., the domains are dissimilar $\mathcal{D}_s \neq \mathcal{D}_t$ (either $\mathcal{X}_s \neq \mathcal{X}_t$ or $P(\mathbf{X}_s) \neq P(\mathbf{X}_t)$), the model would not generalize well to the target domain. Transfer learning (TL) aims to use learned feature representations from the source domain to improve task performance in the target domain. More specifically, TL is defined in [13] as

Definition 3 (Transfer Learning) Given a source domain \mathcal{D}_s and learning task \mathcal{T}_s , a target domain \mathcal{D}_t and learning task \mathcal{T}_t , transfer learning aims to improve the learning of the target predictive function $f_t : \mathbb{R} \mapsto \mathbb{N}$ in \mathcal{D}_t using the knowledge in \mathcal{D}_s and \mathcal{T}_s , where $\mathcal{D}_s \neq \mathcal{D}_t$, or $\mathcal{T}_s \neq \mathcal{T}_t$.

Transductive transfer learning refers to the case where the source and target tasks are equal $\mathcal{T}_s = \mathcal{T}_t$, but there is a non-trivial domain shift. When the feature spaces are the same $\mathcal{X}_s = \mathcal{X}_t$ but the marginal distributions are different $P(\mathcal{X}_s) \neq P(\mathcal{X}_t)$, we enter the practice of domain adaptation (DA). When labels for the target domain $\{\mathcal{Y}_t, P(\mathbf{Y}_t|\mathbf{X}_t)\}$ are not available, we speak of unsupervised DA, and when there is a scarce set of labels available in the target domain, we speak of semi-supervised DA [13]. A prerequisite for successful domain adaptation is that there must be some similarities between the two domains. For visual applications, this would be low-level features such as shapes, while high-level features, such as image background, lighting, color, etc., can differ.

5 CORAL

CORAL is short for covariance alignment which minimizes the Frobenius norm (see Definition 5, in Appendix A) between the covariance matrices \mathbf{C}_s and \mathbf{C}_t of the source and target data \mathbf{X}_s and \mathbf{X}_t by finding the optimal linear transformation \mathbf{A}^* acting on the source covariance. The method is a *feature alignment* method, i.e. it transforms the source data to align with the target data with respect to some measure, in this case, the second-order statistic. The optimal linear transformation can be found by considering the optimization problem

$$\mathbf{A}^* = \arg \min_{\mathbf{A}} \|\hat{\mathbf{C}}_s - \mathbf{C}_t\|_F = \arg \min_{\mathbf{A}} \|\mathbf{A}^T \mathbf{C}_s \mathbf{A} - \mathbf{C}_t\|_F, \quad (1)$$

where $\hat{\mathbf{C}}_s$ is the covariance of the transformed source features \mathbf{X}_s and $\|\mathbf{A}\|_F$ is the matrix Frobenius norm, see Definition 5 in Appendix A. The optimization problem (1) is derived in [3], and has the solution

$$\mathbf{A}^* = \left(\mathbf{U}_s \Sigma_s^{\frac{1}{2}} \mathbf{U}_s^\top \right) \left(\mathbf{U}_{t[1:r]} \Sigma_{t[1:r]}^{\frac{1}{2}} \mathbf{U}_{t[1:r]}^\top \right)^{-1}, \quad (2)$$

where $(\mathbf{U} \Sigma \mathbf{U}^\top)$ is the singular value decomposition (SVD) of \mathbf{C} , the superscript \dagger denotes the Moore-Penrose pseudo inverse [15], and r is the lowest rank of the source and target covariance

matrices $r = \min(r_{C_s}, r_{C_t})$. However, computing the SVD of the covariance matrices is relatively expensive, and could be unstable. Therefore, it is sufficient to use the simpler transformation

$$\mathbf{A}^* = \mathbf{C}_s^{\dagger \frac{1}{2}} \mathbf{C}_t^{\frac{1}{2}}. \quad (3)$$

CORAL is summarized in Algorithm 1.

Algorithm 1: CORAL for Unsupervised Domain Adaptation

- 1 **Input:** Source data \mathbf{X}_s and target data \mathbf{X}_t
 - 2 **Output:** Adjusted source data \mathbf{X}_s^*
 - 3 $\mathbf{C}_s = cov(\mathbf{X}_s)$
 - 4 $\mathbf{C}_t = cov(\mathbf{X}_t)$
 - 5 $\mathbf{X}_s = \mathbf{X}_s \mathbf{C}_s^{-\frac{1}{2}}$ ▷ Whitening source
 - 6 $\mathbf{X}_s^* = \mathbf{X}_s \mathbf{C}_t^{\frac{1}{2}}$ ▷ re-coloring with target covariance
-

6 Deep CORAL

Extending the shallow CORAL method, which applies a linear transformation, [16] proposes a deep approach to learning features that generalizes in the target domain by constructing a differentiable joint loss function from the classification loss ℓ_{class} and the CORAL loss ℓ_{CORAL}

$$\ell = \ell_{class} + \lambda \ell_{CORAL}, \quad (4)$$

where λ weights the relative importance between the domain adaptation task and the discriminative task in the source domain. Minimizing the CORAL loss is analogous to minimizing the Frobenius norm (L2 norm) between the second order statistics of the source and target data

$$\ell_{CORAL} = \frac{1}{4d^2} \|\mathbf{C}_s - \mathbf{C}_t\|_F^2, \quad (5)$$

where d is the dimension of the features and F denotes the Frobenius norm. The covariances are calculated as

$$\mathbf{C}_s = \frac{1}{n_s - 1} \left(\mathbf{D}_s^T \mathbf{D}_s - \frac{1}{n_s} (\mathbf{1}^T \mathbf{D}_s)^T (\mathbf{1}^T \mathbf{D}_s) \right) \quad (6)$$

$$\mathbf{C}_t = \frac{1}{n_t - 1} \left(\mathbf{D}_t^T \mathbf{D}_t - \frac{1}{n_t} (\mathbf{1}^T \mathbf{D}_t)^T (\mathbf{1}^T \mathbf{D}_t) \right), \quad (7)$$

where n_s (n_t) is the number of data points in the source (target) domain, \mathbf{D}_s (\mathbf{D}_t) the source (target) data, and $\mathbf{1}$ is a column vector with all elements equal to one. While [16] implements the deep CORAL method on a single fully connected layer, the authors add that it should be straightforward to apply the loss to several other fully connected layers.

7 Data Acquisition and Preprocessing

The following section presents how the datasets were acquired or constructed. Before using the datasets, the image means are normalized, and the standard deviations are standardized.

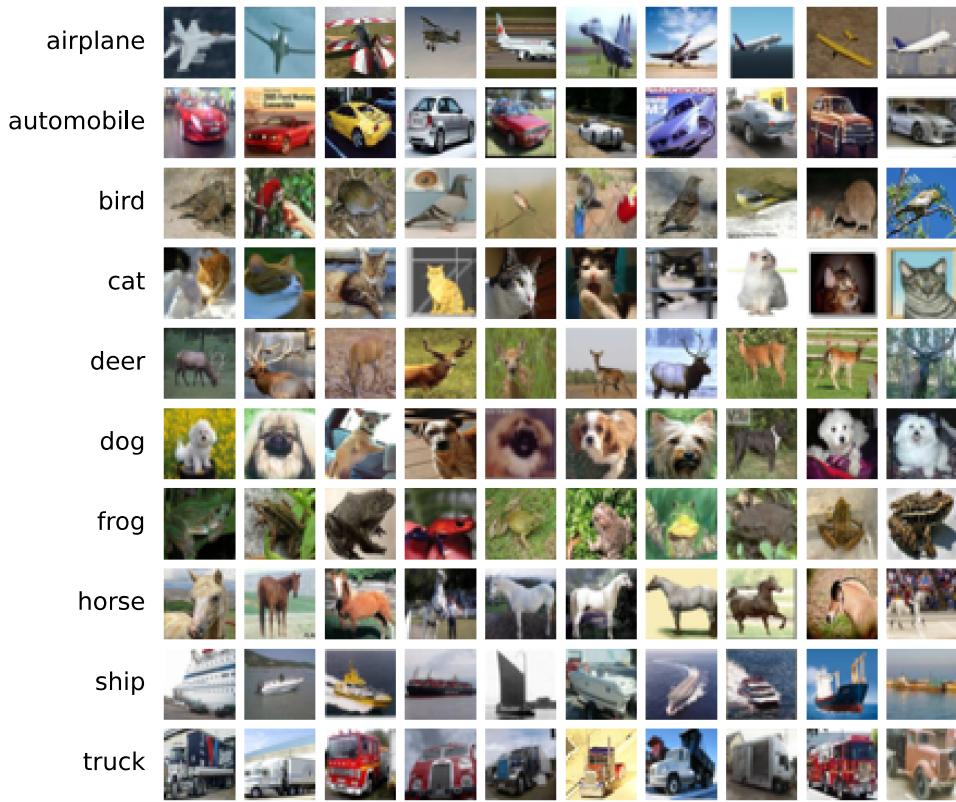


Figure 1: *Ten randomly picked images of each class in the Cifar-10 dataset.*

7.1 Cifar-10

The Cifar-10 dataset [17] is a subset of the *80 million tiny images* dataset, collected by Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton and contains 60 000 32×32 RGB images categorized into ten balanced (same number of instances in all classes) classes. There are 50 000 images for training, 5000 of them are validation instances, and 10 000 images for testing. Figure 1 shows the ten classes, with ten randomly selected images from each class. To match the dimensions of the LiDAR data, the Cifar-10 images are converted to one-channel grayscale images.

7.2 Camera Dataset

The balanced camera dataset consists of 18 352 grayscale images with 32×32 pixels taken by the vehicles of the Scania Autonomous Transport Research department. The dataset consists of 16 520 training instances, of which 10% are used for validation, and 1832 test instances. The two classes can be seen in Figure 6, with ten randomly selected images from each class. The cameras are located at all eight corners of the vehicles to provide a wide field of view of the traffic. A labeling company has created a set of corresponding label files (a dictionary data type in a JSON file format). Given an image there is a label file with x-y coordinates for bounding boxes and labels for each corresponding object (among other properties of the objects in the image). All relevant objects in the image are cropped to a square shape using the bounding boxes and resampled so that all cropped images are 32×32 pixels. The images are converted to

grayscale to have the same shape as the LiDAR intensity images, and can be used by the same neural network architecture. The cropping process follows a set of rules, illustrated in Figure 2 and 4 alongside the resulting cropped images seen in Figure 3 and 5. The cropping rules are: If the longest side length of the bounding box is

1. below the threshold of 10 pixels, the image is omitted from the dataset to uphold quality.
2. larger than the threshold but smaller than the desired crop size (32x32), then the boundaries for the object are expanded to the desired crop size.
3. expanded beyond the boundaries of the image, the expansion will stop at the image border and the expand in the other direction to create a square crop.
4. larger than the desired crop size, it is cropped square and re-sampled to the desired crop size using Lanczos interpolation [18].

The effect of each rule can be seen in effect for two traffic scenes, shown in Figure 2 and 4, with the resulting cropped images in Figure 3 and 5. Rule 1 omits small objects from the image, shown in Figure 2 as red boxes which surrounds the most distant vehicles. Rule 2 is applied to the vehicles in images 1, 2, 3, 4, 8, 9, and 10 in Figure 3 because the crop bounds are not adjacent to the vehicles in the image, but are expanded to uphold resolution of 32x32 pixels. For the white vehicle in the lower left corner of Figure 2 (object number 6 in Figure 3), rules 3 and 4 apply. Rule 3 places the square crop such that the vehicle is not centered, and is later down-sampled. A relatively simple cropping procedure like this will also include some data instances that may degrade the quality of the dataset, and are therefore better discarded. It is inevitable that some objects will be obscured by something in the foreground, such as a traffic sign, a lamppost, or another vehicle, or pedestrian. Expanding a bounding box will also include other objects that are not relevant to the label of the cropped image. However, one could argue that the inclusion of images with obstructions improves the model because it accurately represents a real-world scenario.

7.3 LiDAR Dataset

The balanced LiDAR (light detection and ranging) dataset contains 4010 32×32 one channel images taken by the vehicles at the Scania Autonomous Transport Research department . Of these images, there are 3600 training instances, of which 10% are for validation, and 410 images are test data. The two classes (“car” and “pedestrian”) are shown in Figure 7, with ten random images from each class. The LiDAR intensity images are LiDAR point clouds projected onto a 2D plane, whose pixel values correspond to the intensity value for each point in the cloud. LiDAR works by emitting pulses of photons in a laser beam and measuring the time it takes for the photons to return to the transmitter. The light returns by reflecting off objects in the surrounding environment, and will sometimes reflect away from the transmitter, producing relatively sparse images, as shown in Figure 7. Because LiDAR images come from a sparse point cloud, the cropping procedure differs slightly from the one used for the camera images. Instead of the threshold size being the longest side length of a bounding box, it is instead the number of pixels assigned to a label, which is fifty pixels. Otherwise, the cropping procedure is the same as for the camera dataset.



Figure 2: The orange rectangles illustrates the original bounding boxes from the label data, and the yellow squares illustrates the cropping method, except for resampling to desired resolution. Objects bounded by red boxes has not passed the size threshold and has therefore been omitted.

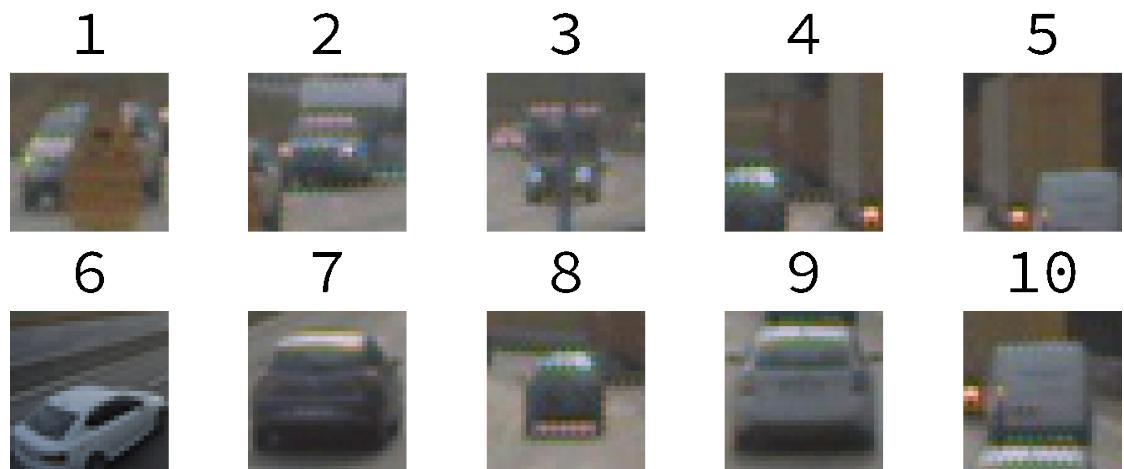


Figure 3: Cropped camera images of equal size and shape.

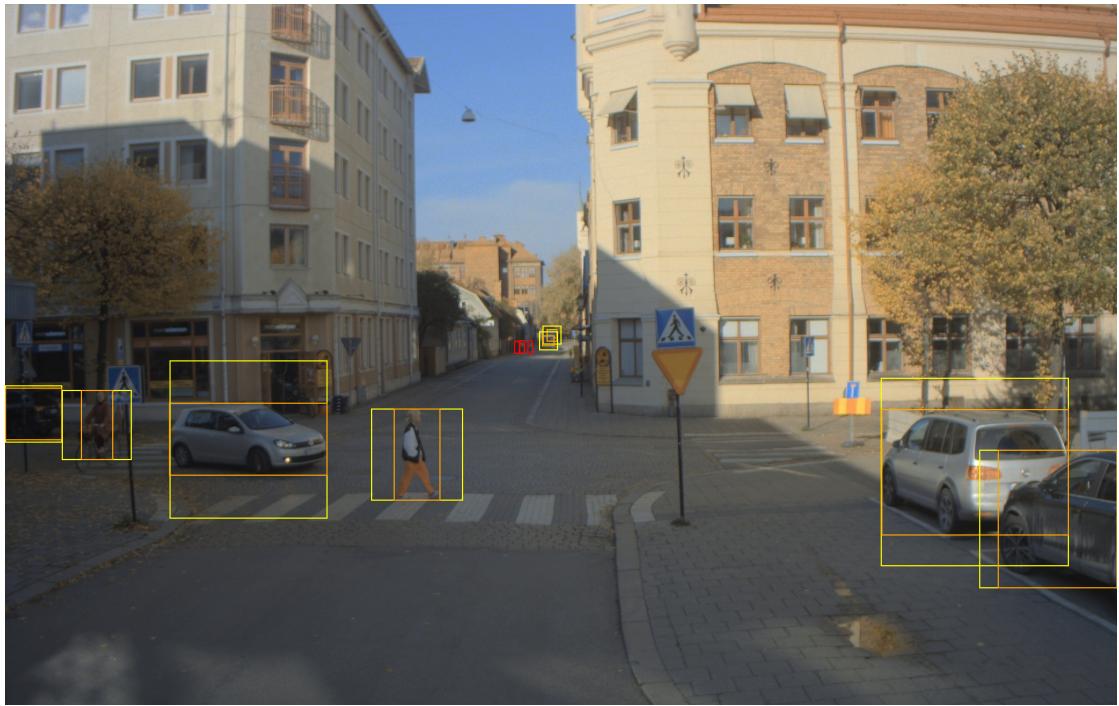


Figure 4: Another example of a traffic scene with a total of five pedestrians present (cyclist included), three of which are omitted from the dataset because the corresponding bounding boxes are too small (colored in red).



Figure 5: Cropped images of equal size and shape.



Figure 6: Ten randomly picked grayscale images from each class in the Camera dataset.

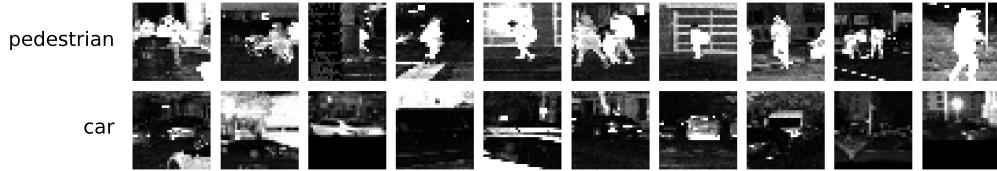


Figure 7: Ten randomly picked LiDAR intensity images from each class in the LiDAR dataset.

8 Experimental Evaluation

The experiments test the shallow CORAL method, a supervised model on the source domain (camera dataset), and the deep CORAL method, including the hyperparameter tuning of the CORAL loss weight λ from equation (4). The models are evaluated and compared using the classification accuracy of the LiDAR test dataset, and the confusion matrix metric, revealing any classification biases in the models. The DA methods presented here all use a convolutional neural network (CNN) pretrained on the Cifar-10 dataset. Shallow CORAL extracts deep features from the pretrained CNN, without updating the model parameters of the CNN. However, deep CORAL updates the model parameters and thus finetunes the model to the DA task. Because the two methods are of different complexities, and were reimplemented at different stages of the thesis project, two different machine learning frameworks were used for the methods. Shallow CORAL was implemented with the Keras library (runs on tensorflow), while the deep CORAL model called for the more low-level API Pytorch. As a consequence, the pretraining settings were not the same for the two methods, but should not differ in such a way that they cannot be compared in testing. The pretraining hyperparameters of the CNN are presented in the sections of each methods. But common for all methods is the architecture of the CNN.

8.1 Convolutional Neural Network Architecture

Illustrated in Figure 8, the CNN architecture is inspired by the “block” design proposed by the Visual Geometry Group (VGG) at the University of Oxford in [19], known as VGG blocks. There are three VGG blocks which contains two convolutional layers followed by maxpooling, performed with a 2×2 window with stride 2, and a subsequent dropout layer with $p = 0.2$ (20% chance of not updating the layer parameters during training). All convolutional layers has filters with a receptive field of 3×3 . The convolution stride is set to 1, and the spatial padding of the convolutions are set such that the resolution is preserved (`padding = 'same'`). After each convolutional layer there is also a batch normalization layer. The dropout and batch normalization prevents, to a certain degree, the model from overfitting to the training data. After the last maxpool layer `maxpool_6`, there are four fully connected layers, followed by batch normalization layers. All hidden layers has a rectified linear unit (ReLU) activation function.

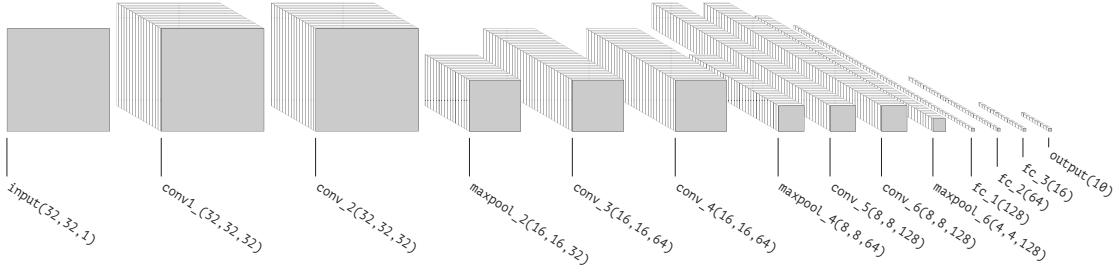


Figure 8: Architecture of the convolutional neural network. Each layer and its output dimension is defined as `layer_name(width, height, channels)`.

The CNN architecture is not specifically tailored for performance, but rather used as foundation on which the findings of the DA methods are evaluated.

8.2 Pretraining Convolutional Neural Network on Cifar-10

8.2.1 For Shallow CORAL Method

The input images has a batch size of 64, base learning rate is $\epsilon_0 = 0.01$ which decreases as $\epsilon_{k+1} = \epsilon_k 0.096^k$ for $k \in \{0, 1, \dots, N_e - 1\}$, where $N_e = 75$. The Adam optimizer [20] has a weight decay of 10^{-4} , which adds a L2 norm of all weights, multiplied by the weight decay parameter to the loss function. This regularization technique prevents the weight of becoming too large. The loss function is a categorical cross-entropy (CE) loss function (for more details on cross entropy, see Appendix Section B.2). The accuracy on the Cifar-10 test dataset reaches 84.79%. When the pretrained weights are used to initialize the weights of a model, they will be referred to as the Cifar-10 weights.

8.2.2 For Finetuned Model & Deep CORAL Model

The input images has a batch size of 32, base learning rate is $\epsilon_0 = 0.01$ which decreases exponentially as $\epsilon_k = \epsilon_0 e^{-0.05k}$ for $k \in \{1, 2, \dots, N_e\}$, where $N_e = 32$. The stochastic gradient descent (SGD) has a momentum of 0.9 and a weight decay of 10^{-5} (for more details on SGD, see Appendix Section B.1). Weight decay for SGD is equivalent to L2 weight regularization, which adds the L2 norm of all weights, multiplied by the weight decay parameter, to the loss function. This regularization technique prevents the weight of becoming too large. The loss function is a categorical cross-entropy (CE) loss function (for more details on cross entropy, see Appendix Section B.2). The CNN architecture is not specifically tailored for performance, but rather used as foundation on which the findings of the DA methods are evaluated. The accuracy on the Cifar-10 test dataset reaches 80.70%. Pretraining will give the optimizer a better start guess than normally distributed weights and biases, such that a minima is reached faster when finetuning the model on the source dataset. The pretrained model is used later to initialize the weights of other models with identical architectures, except for the output layer which is replaced to fit the number of classes needed for the new classification task.

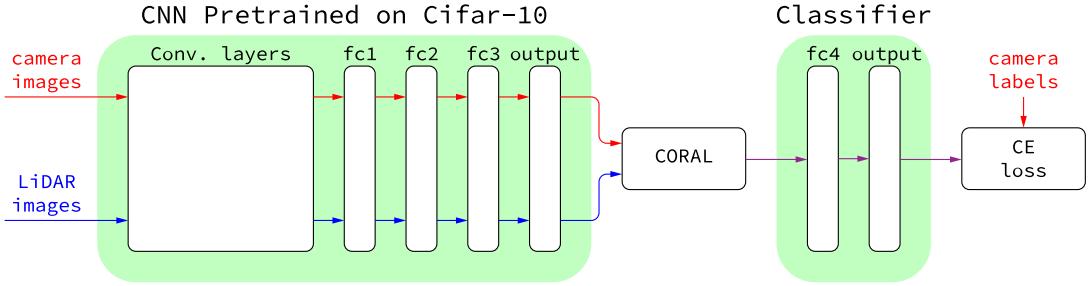


Figure 9: The shallow CORAL model. Red arrows and text represent the camera data, and blue represents the LiDAR data. Aligned camera features by the CORAL method, are depicted by purple arrows. The architecture of the pretrained CNN is seen in Figure 8

8.3 CORAL Model

8.3.1 Pretraining CNN

The setup for the shallow CORAL model is illustrated in Figure 9. Red and blue arrows represent the camera and LiDAR data, respectively, and purple arrows after the CORAL method represents the aligned camera features. After the CNN has been initialized with the Cifar-10 parameters, the camera and LiDAR images are passed through, where their ten dimensional deep features are extracted at the output layer. The extracted camera deep features are then aligned to the LiDAR deep features using the linear transformation according to Algorithm 1 on page 4. In the algorithm, the matrix square root of the covariance is performed by a blocked Schur algorithm [21]. The aligned ten-dimensional Lidar are visualized in a T-distributed Stochastic Neighborhood Embedding (TSNE) in Figure 19a. Thereafter, a linear classifier is trained on the aligned camera deep features and evaluated on the LiDAR deep features. Note that the CNN is not trained on the camera or LiDAR data after pretraining. The separate linear classifier learns to discriminate between the two classes (cars and pedestrians) using the camera deep features which are aligned to the LiDAR deep feautres w.r.t. their batch-wise covariances. Training on aligned camera deep features enables the linear classifier to adapt and reach classification accuracy in the target domain to 74.67% on the LiDAR test dataset.

8.4 Deep Model without CORAL Loss

The baseline model is a supervised CNN model with parameters initialized by the Cifar-10 dataset. The model is further trained on the camera dataset as to solely rely on the similarity of source and target domain without any DA method. As the dataset has two classes, the output layer from the CNN of ten dimensions is replaced with a output layer with two dimensions. Due to the new output layer being randomly initialized, it is assigned a higher initial learning rate such that it will not take as long for the optimizer to reach an optimum. More specifically, the model has a SGD optimizer with initial learning rates $\epsilon_{0,fc4} = 0.01$ for the output parameters and $\epsilon_0 = 0.001$ for all other parameters. The momentum is set to 0.9 and the L2 weight decay coefficient to 10^{-5} . With each epoch, the learning rate decays exponentially as $\epsilon_i = \epsilon_0 e^{-0.05*i}$ for $\{i\}_{i=1}^{N_e}$, where $N_e = 30$ is the number of epochs performed. The setup for this baseline model is illustrated in Figure 10. First, the weights are initialized with the Cifar-10 weights and then the camera images are passed through the model in batches of 32 and evaluated in the CE loss function in a supervised manner. Training in the source domain while excluding the CORAL

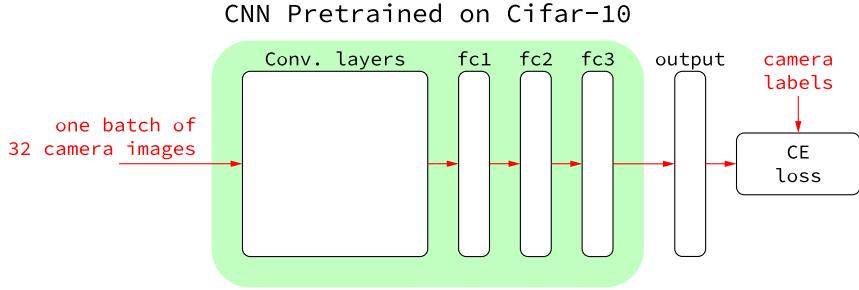


Figure 10: *The baseline model is transforming the camera data, shown with red arrows, in a supervised manner.*

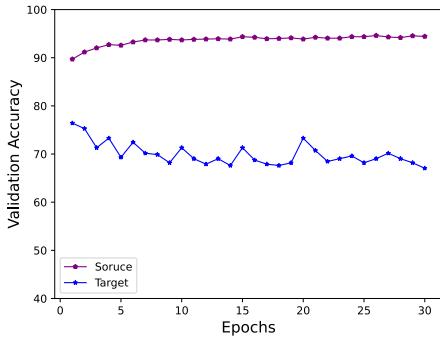


Figure 11: *With no means for domain adaptation, accuracy in the target domain steadily decreases as the model learns a classifier in the source domain.*

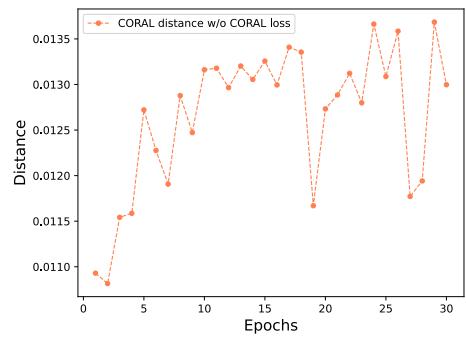


Figure 12: *As the CORAL distance (5) is not included in the loss function, it increases over the training iterations.*

loss in the objective function increases the CORAL distance over the iterations, as seen in Figure 12. The CORAL distance increases in general during training, but jumps significantly, which might be mitigated by using a larger LiDAR batch size. The increased CORAL distance seems to correlate with a decreased validation accuracy in the target domain, seen in Figure 11. The training history shown in Figure 11 and 12 also suggests that one epoch of training generalizes the model relatively well to the target domain, before it becomes too domain specific, which is also reflected by the shortest CORAL distance achieved in beginning of training. As the CORAL distance and validation accuracy in the source domain increases, the validation accuracy in the target domain decreases to 67.05%, which further suggests that minimizing the CORAL distance would in opposition increase the validation accuracy for the LiDAR data. The 16 dimensional LiDAR deep features from the third fully connected layer ($fc3$) are visualized using TSNE in Figure 19b. The embedded features are not well separated in the plot which is reflected in the relatively low validation accuracy in the target domain.

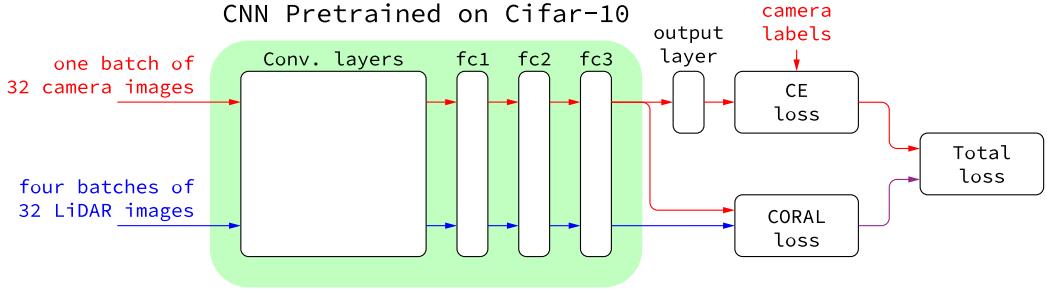


Figure 13: The deep CORAL model. Red and blue arrows represent the camera and LiDAR data, respectively. The purple arrow represents the CORAL loss which is jointly evaluated with the CE loss in the total loss. The architecture of the pretrained CNN is seen in Figure 8

8.5 Deep CORAL Model

The Deep CORAL model shares most of the hyperparameters with the model without a CORAL loss in Section 8.4, except for the number of epochs performed which is $N_e = 20$, and the additional CORAL weight λ from the joint loss function (4). The pipeline seen in Figure 13, begins similarly to the model without CORAL loss. The LiDAR and camera images are passed through the pretrained CNN, and extracted at fc3. Thereafter, the camera features are passed through the output layer, from which they are measured against the corresponding labels in the CE loss function. The batch of camera features from fc3 are also measured in the CORAL loss function four times against four different batches of LiDAR features, to get a better measurement for the domain adaptation task. The mean CORAL distance is then multiplied with a CORAL loss weight λ to control how much the optimizer should take the DA task into account when updating the weights. Lastly, the CE loss and CORAL loss are added together (total loss) from which the optimizer iterates towards an optimal set of weights and biases. A successful model should be able to discriminate between classes in the target domain by 1), the supervised learning task in the source domain and 2), the domain adaptation task, which adapts the classifier to perform well in the target domain.

8.5.1 Hyperparameter Tuning

The hyperparameter λ from the joint loss function in Equation (4) on page 14 is chosen such that a high validation accuracy on the LiDAR data is achieved. It is unknown how the model behaves with respect to λ , so 20 samples are drawn from the logarithmically uniform distribution

$$\lambda \sim f(x; a, b) = \frac{1}{x \log_e(\frac{b}{a})}, \quad (8)$$

with $a = 10^0$ and $b = 10^6$. The log-uniform distribution is an effective grid search because it spans several orders of magnitude, thus finding the best hyperparameter relatively quickly. For each sample of λ , the deep CORAL model is trained according to the specifications given in Section 6. The final validation accuracy on the source and target data, and the validation losses w.r.t. the hyperparameter can be seen in Figure 14 and 15. The goal is to yield the highest validation accuracy in the target domain, seen in blue in Figure 14, while maintaining a low total loss, seen in dark red in Figure 15. As the optimizer seeks to minimize the CE loss and the weighted CORAL loss jointly, the *weighted* CORAL loss is also plotted (appropriately in

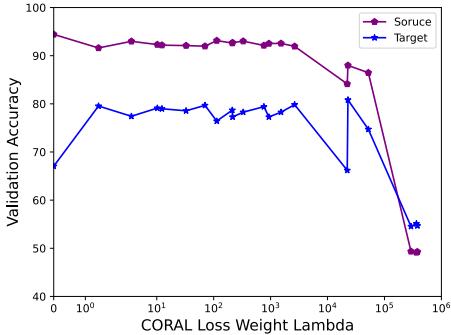


Figure 14: *Source and target validation accuracy at end of training with respect to the weight to the CORAL loss*

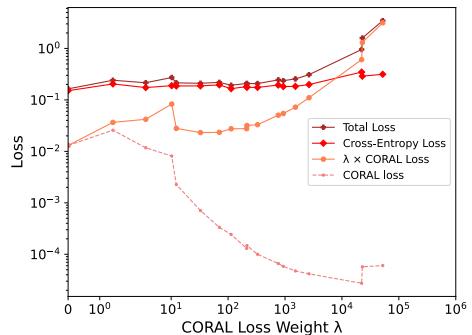


Figure 15: *Loss metrics with respect to the weight for the CORAL loss function. The losses are undefined for the four last values of λ , and are therefore not included.*

the color *coral*) in Figure 15. For low values of λ , the CE loss is greater than the weighted CORAL loss. As λ increases, the weighted CORAL loss does too, as it is multiplied by λ . The minimal weighted CORAL loss produced by the optimizer is found at about $\lambda = 10^2$. Thereafter the weighted CORAL loss increases steadily, due to the CORAL loss (colored in *light coral* in Figure 15) not decreasing at the same rate as λ is increasing. Closest sampled value to $\lambda = 10^2$ is $\lambda = 70.229$. From excluding the CORAL distance from the objective function ($\lambda = 0$) to including the CORAL loss $\lambda \in (10^0, 10^3)$, there is a substantial increase in the validation accuracy on the target data, shown in Figure 14, where it remains stable at just below 80%. For $\lambda = 70.299$, the validation accuracy is 79.688%. If the hyperparameter is increased further, the task of minimizing the CORAL loss will dominate the objective function, compromising the discriminative task such that performance suffers to the point of random choice (accuracy near 50%) at about $\lambda \in (10^5, 10^6)$. The 16 dimensional LiDAR deep features from the third fully connected layer (**fc3**), with $\lambda = 70.299$, are visualized using TSNE in Figure 19c. Compared to the TSNE plot without CORAL loss in Figure 19b, the deep CORAL model has increased the interclass distance (distance *between* clusters of different classes, see Figure 18) and decreased the intraclass distance (distance among data instances *within* classes). The separation of the two classes is also reflected in the increased validation accuracy in the target domain. The reason why performance suffers at larger values of λ is due to the trivial solution solution of transforming all data instances, regardless of class, into almost the same value, making the covariance distance very small, and the accuracy near 50%. In other words, the interclass distance is minimized when the actual task is to maximize it.

9 Results

Initializing the model parameters by pretraining on the Cifar-10 dataset provides a good foundation for all methods presented in the paper. This proves yet another aspect of transfer of learning, which is the generalizability from low level visual features learned by one domain (in which the the Cifar-10 data resides). The shallow CORAL model acts on the deep features produced by the the pretrained network and learns a linear classifier on the aligned source features. The model yields a classification accuracy of 73.17% on the LiDAR test dataset. The confusion

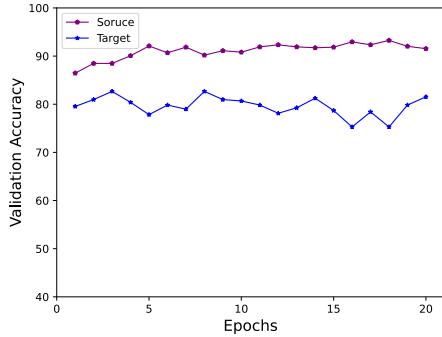


Figure 16: Validation accuracy for source and target data during training, with tuned hyperparameter $\lambda = 70.299$.

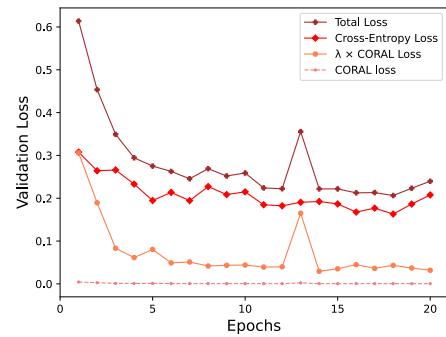


Figure 17: Validation losses during training, with tuned hyperparameter $\lambda = 70.299$

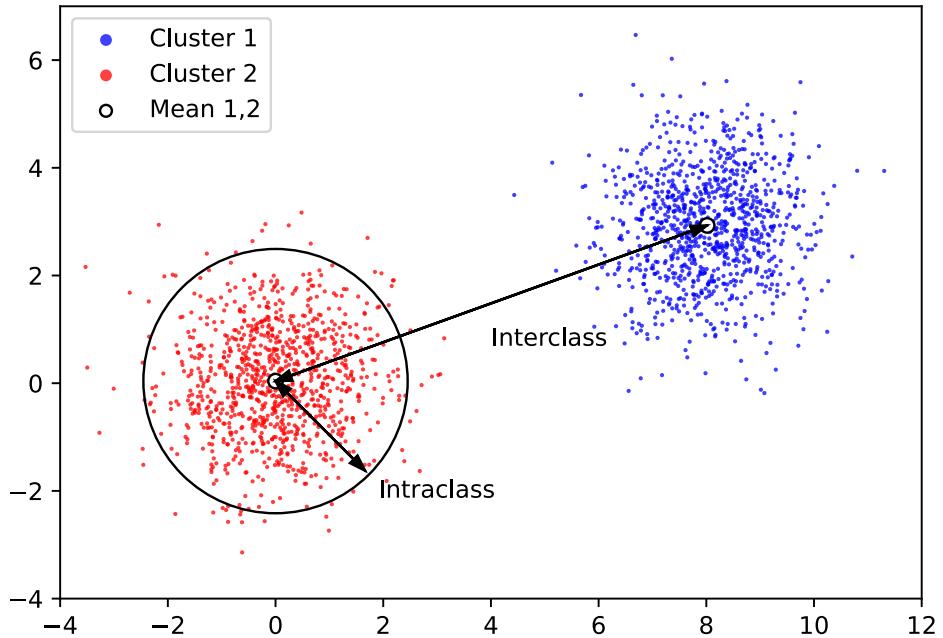


Figure 18: Samples drawn from two different Gaussian distributions. Interclass distance is the distance between the clusters (illustrated as the means of each cluster) and the Intraclass distance is the distance between the points within a cluster. A successful classifier maximizes the interclass distance and minimizes the intraclass distance.

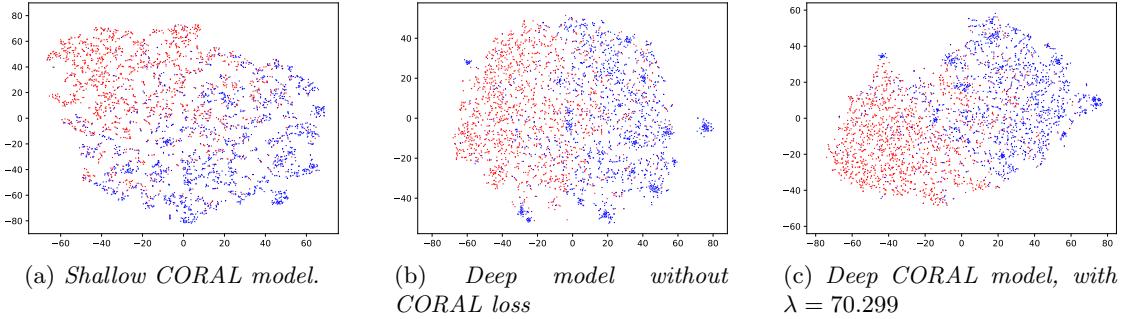


Figure 19: *TSNE embeds the deep features high dimensional data into two dimensions, seen as points in the scatter plot.* plot (a) illustrates the ten dimensional deep features extracted at the (output layer of the CNN. Plots (b) and (c) represents the 16 dimensional deep features extracted at $fc3$ from the deep models.

matrices in Figure 20 give insight into the class-wise predictions. The horizontal axis represents the predicted class, and the vertical axis represents the true class. Each quadrant displays how many times in percent, each class was predicted correctly or incorrectly. While the accuracy of the shallow CORAL model is relatively low, it has low bias towards any of the classes. The label “car” is predicted almost as many times as “pedestrian”, regardless of it being a correct or false prediction. Pretraining is a very efficient way of initializing the model parameters, and provides a good starting guess for the optimizer when the model is further trained for another, similar task. This is especially apparent in the training history of the finetuning model in Figure 11, where the LiDAR accuracy (in blue) is highest in the beginning of training, and deteriorates as the model is trained on the camera dataset, until the LiDAR accuracy reaches 67.06%. The lost generalizability is also reflected in the confusion matrix, seen in Figure 20b, which shows a strong bias towards predicting the label “car”. Specifically, it predicts “car” on 82.03% percent of the time, when 50% of the images contains cars. It has an classification accuracy of 63.8% on the LiDAR test dataset. The deep CORAL model includes the CORAL distance in the loss function. With the tuned hyperparameter $\lambda = 70.299$, the Deep CORAL model increases the target accuracy to 80.73% and shifts the bias back to a more balanced confusion matrix, seen in Figure 20c.

10 Future Work

As mentioned in [16], the CORAL loss in the deep CORAL model can be measured across more than one fully connected layer. Each CORAL loss could also have its own loss weight, which would give more control, but also result in a more extensive hyperparameter search. Going beyond the CORAL metric, the established architecture of the deep CORAL model can easily adapt any other statistical metric. This could be an MMD metric, similar to what [8] has done, or higher statistical moments, beyond the second order statistics. Another improvement would be to use all three channels of the Cifar-10, and camera images by employing a preprocessing neural net which embeds the three channels into one, to preserve all information.

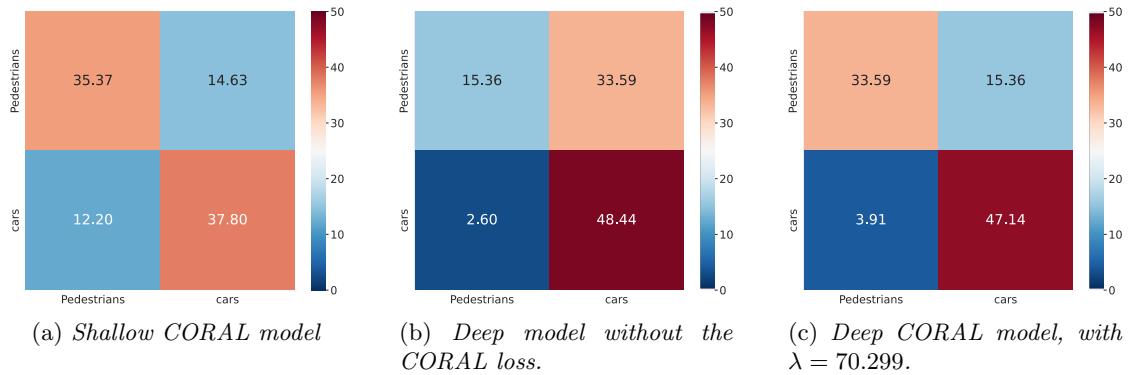


Figure 20: A binary confusion matrix shows true positives and true negatives, false positives and false negatives. The horizontal axis represents the predicted class, and the vertical is the true class. The classification bias towards cars in the finetuned model (b) is mitigated when using the shallow CORAL method (a) and the deep CORAL method (c). While the shallow CORAL model is the most unbiased, the deep CORAL model achieves the highest classification accuracy.

11 Conclusions

This paper has investigated how to adapt a classifier from a source domain to a similar target domain. A supervised model that learns a classifier in the source domain can in addition incorporate a statistical information from the target domain such that the classifier adapts, and thereby increases the classification accuracy in the target domain. Two domain adaptation were implemented, both using the second order statistics of the source and target data, one in a shallow manner, and the other in a deep framework. The shallow CORAL (covariance alignment) method finds the optimal linear transformation which aligns the source data to the target data such that the L2 distance between the two covariances is minimized. Thereafter the aligned source data is used in a supervised learning task to learn a linear classifier, which improves the classification accuracy on the target data. The deep CORAL method operates within a deep architecture, such that the parameters of the neural network are adapted to the target domain. A weighted CORAL is added to the cross-entropy loss, such that the supervised learning task is accompanied with the domain adaptation task. Tuning the CORAL loss weight λ yielded a model with superior performance. The evaluation compared the shallow CORAL method, the deep CORAL method, and a supervised model trained on the camera data. The supervised model solely relies on the similarities between the source and target domain, which shows that there is a non-trivial domain-shift present. The supervised model classifies the camera images well, with an accuracy of roughly 94%, while the accuracy on the LiDAR images deteriorates to 67%. Using a domain adaptation method maintains the accuracy for the LiDAR images. Further, incorporating the DA task into the deep framework had the best results. The shallow CORAL model classifies the LiDAR test data with an accuracy of 73.17%, surpassed by the deep CORAL model which reaches an accuracy of 80.73%. Conclusively, the feature alignment CORAL, and the discrepancy based model deep CORAL both are able to adapt a model trained on data in the source domain, to increase the classification accuracy on data from the target domain. The superiority of deep CORAL suggests that having the domain adaptation task embedded in the deep framework is more effective, as it updates the model parameters on a deep level to learn a more robust representation of the domain adaptation task.

Appendix A Definitions

The following appendix is a collection of relevant definitions.

Definition 4 (Bregman Divergence [22, p. 16, section 3.5.5], [23]) *The Bregman divergence is*

$$d_\psi(P, Q) = \sum_{x \in \mathcal{X}} \Delta_\psi\{p(x), q(x)\},$$

where the operator $\Delta_\Psi(a, b)$ is

$$\Delta_\psi\{a, b\} := \psi(a) - \psi(b) - \psi'(b)(a - b),$$

with the differentiable real-valued convex function $\psi(x)$ of a nonnegative argument. The nonnegative function $\Delta_\Psi(a, b)$ measures how much the convex function ψ deviates at a from its tangent at b .

Definition 5 (Frobenius Norm [24, p. 55]) *The Frobenius norm $f : \mathbb{R}^{m \times n} \mapsto \mathbb{R}$ is defined as*

$$f(\mathbf{A}) = \|\mathbf{A}\|_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^m |a_{ij}|^2}.$$

Appendix B Algorithms and Functions

Here relevant algorithms and other functions are presented.

B.1 Stochastic Gradient Descent

The following algorithms and descriptions are taken from the Deep Learning book [1]. For very large datasets, regular gradient descent is very slow, as it computes every gradient in the loss function. Stochastic gradient descent (SGD) instead computes a gradient estimate by randomly picking a set of m examples from the training set, as shown in Algorithm 2. Although rapidly

Algorithm 2: Stochastic gradient descent (SGD). Update at training iteration k

```

1 Require: Learning rate  $\epsilon_k$ 
2 Require: Initial parameter  $\theta$ 
3 while Stopping criteria not met do
4   Sample a minibatch of  $m$  examples from the training set  $\{\mathbf{x}^{(i)}\}_{i=1}^m$  with
      corresponding labels  $\mathbf{y}^{(i)}$ .
5   Compute gradient estimate  $\hat{\mathbf{g}} \leftarrow +\frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$ .
6   Apply update:  $\theta \leftarrow \theta - \epsilon_k \hat{\mathbf{g}}$ .
7 end

```

increasing the efficiency of learning, SGD may still converge slowly if the loss function has an poorly conditioned hessian matrix. The momentum algorithm [25] shown in Algorithm 3 accelerates learning by introducing a velocity parameter \mathbf{v} which accumulates an exponentially decaying moving average of previously calculated gradients. The term momentum originates from

a physical analogy, where the gradient is a force moving a particle through the parameter space. Since momentum is velocity times mass, unit mass is assumed, and the momentum parameter α determines how fast the contribution of previous gradients will decay.

Algorithm 3: Stochastic gradient descent (SGD) with momentum, Update at training iteration k

```

1 Require: Learning rate  $\epsilon_k$ , momentum parameter  $\alpha$ 
2 Require: Initial parameter  $\theta$ , initial velocity  $v$ 
3 while Stopping criteria not met do
4   Sample a minibatch of  $m$  examples from the training set  $\{\mathbf{x}^{(i)}\}_{i=1}^m$  with
      corresponding labels  $\mathbf{y}^{(i)}$ .
5   Compute gradient estimate  $\hat{\mathbf{g}} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta; \mathbf{y}^{(i)}))$ .
6   Compute velocity update:  $\mathbf{v} \leftarrow \alpha \mathbf{v} - \epsilon \hat{\mathbf{g}}$ 
7   Apply update:  $\theta \leftarrow \theta + \mathbf{v}$ .
8 end

```

B.2 Cross-Entropy Loss Function

The following definitions are from [26, p. 57]. The uncertainty of a random variable X can be measured by its entropy

$$\mathbb{H}(X) \triangleq \sum_{k=1}^K \log p(X = k) \log_2 p(X = k). \quad (9)$$

If the logarithmic base is 2, the unit of entropy is *bits* (or Shannons, named after the American mathematician and engineer Claude Elwood Shannon, who laid the theoretical foundations of information theory), and *nats* if the base is e (i.e. the natural log). Minimizing the entropy is equal to minimizing the uncertainty of the random variable, meaning the outcome becomes more deterministic. To measure the similarity between two distributions p and q , the Kullback-Leiber (KL) divergence can be used:

$$\mathbb{KL}(p||q) \triangleq \sum_{k=1}^K p_k \log \frac{p_k}{q_k}, \quad (10)$$

where the summation is replaced by an integral for continuous probability density functions, which is known as the differential entropy. Expanding Equation (10) yields

$$\mathbb{KL} = \sum_{k=1}^K p_k [\log p_k - \log q_k] = \sum_{k=1}^K p_k \log p_k - \sum_{k=1}^K p_k \log q_k = -\mathbb{H}(p) + \mathbb{H}(p, q), \quad (11)$$

where $\mathbb{H}(p, q)$ is the *cross-entropy*:

$$\mathbb{H}(p, q) \triangleq - \sum_{k=1}^K p_k \log q_k. \quad (12)$$

Minimizing the cross-entropy w.r.t. q is equal to minimizing the KL divergence, since q is not part of $\mathbb{H}(p)$ in Equation 11.

References

- [1] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- [2] G. Csurka, “Domain adaptation for visual applications: A comprehensive survey,” *arXiv preprint arXiv:1702.05374*, 2017.
- [3] B. Sun, J. Feng, and K. Saenko, “Return of frustratingly easy domain adaptation,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 30, 2016.
- [4] B. Fernando, A. Habrard, M. Sebban, and T. Tuytelaars, “Unsupervised visual domain adaptation using subspace alignment,” in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2013.
- [5] M. Baktashmotagh, M. T. Harandi, B. C. Lovell, and M. Salzmann, “Unsupervised domain adaptation by domain invariant projection,” in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, December 2013.
- [6] E. Tzeng, J. Hoffman, N. Zhang, K. Saenko, and T. Darrell, “Deep domain confusion: Maximizing for domain invariance,” 2014.
- [7] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in Neural Information Processing Systems* (F. Pereira, C. Burges, L. Bottou, and K. Weinberger, eds.), vol. 25, Curran Associates, Inc., 2012.
- [8] M. Long, Y. Cao, J. Wang, and M. I. Jordan, “Learning transferable features with deep adaptation networks,” in *Proceedings of the 32nd International Conference on International Conference on Machine Learning - Volume 37*, ICML’15, p. 97–105, JMLR.org, 2015.
- [9] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative adversarial nets,” in *Advances in Neural Information Processing Systems* (Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K. Weinberger, eds.), vol. 27, Curran Associates, Inc., 2014.
- [10] M.-Y. Liu and O. Tuzel, “Coupled generative adversarial networks,” in *Advances in Neural Information Processing Systems* (D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, eds.), vol. 29, Curran Associates, Inc., 2016.
- [11] K. Bousmalis, N. Silberman, D. Dohan, D. Erhan, and D. Krishnan, “Unsupervised pixel-level domain adaptation with generative adversarial networks,” 2017.
- [12] C. H. Judd, “The relation of special training and general intelligence,” *Educational Review*, vol. 36, pp. 28–42, 1908.
- [13] S. J. Pan and Q. Yang, “A survey on transfer learning,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 22, no. 10, pp. 1345–1359, 2010.
- [14] F. Zhuang, Z. Qi, K. Duan, D. Xi, Y. Zhu, H. Zhu, H. Xiong, and Q. He, “A comprehensive survey on transfer learning,” *Proceedings of the IEEE*, vol. 109, no. 1, pp. 43–76, 2021.
- [15] R. Penrose, “A generalized inverse for matrices,” *Mathematical Proceedings of the Cambridge Philosophical Society*, vol. 51, no. 3, p. 406–413, 1955.
- [16] B. Sun and K. Saenko, “Deep coral: Correlation alignment for deep domain adaptation,” in *Computer Vision – ECCV 2016 Workshops* (G. Hua and H. Jégou, eds.), (Cham), pp. 443–450, Springer International Publishing, 2016.
- [17] “The cifar-10 dataset.” <https://www.cs.toronto.edu/~kriz/cifar.html>. Accessed: 2021-03-22.
- [18] K. Turkowski, *Filters for Common Resampling Tasks*, p. 147–165. USA: Academic Press Professional, Inc., 1990.
- [19] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” 2014.

- [20] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” in *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings* (Y. Bengio and Y. LeCun, eds.), 2015.
- [21] E. Deadman, N. J. Higham, and R. Ralha, “Blocked schur algorithms for computing the matrix square root,” in *Applied Parallel and Scientific Computing* (P. Manninen and P. Öster, eds.), (Berlin, Heidelberg), pp. 171–182, Springer Berlin Heidelberg, 2013.
- [22] P. D. Grünwald and A. P. Dawid, “Game theory, maximum entropy, minimum discrepancy and robust bayesian decision theory,” *The Annals of Statistics*, vol. 32, aug 2004.
- [23] L. Bregman, “The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming,” *USSR Computational Mathematics and Mathematical Physics*, vol. 7, no. 3, pp. 200–217, 1967.
- [24] G. H. Golub and C. F. Van Loan, *Matrix Computations*. The Johns Hopkins University Press, third ed., 1996.
- [25] B. Polyak, “Some methods of speeding up the convergence of iteration methods,” *USSR Computational Mathematics and Mathematical Physics*, vol. 4, no. 5, pp. 1–17, 1964.
- [26] K. Murphy, *Machine Learning: A Probabilistic Perspective*. Adaptive Computation and Machine Learning series, MIT Press, 2012.