# DCID: Deep Canonical Information Decomposition

Alexander Rakowski[1][✉][0000−0002−8134−6729] and
Christoph Lippert[1,2][0000−0001−6363−2556]

[1] Hasso Plattner Institute for Digital Engineering, University of Potsdam, Germany
[2] Hasso Plattner Institute for Digital Health at Mount Sinai, New York, United
States {alexander.rakowski,christoph.lippert}@hpi.de

**Abstract.** We consider the problem of identifying the signal shared between two one-dimensional *target* variables, in the presence of additional multivariate observations. Canonical Correlation Analysis (CCA)-based methods have traditionally been used to identify shared variables, however, they were designed for multivariate targets and only offer trivial solutions for univariate cases. In the context of Multi-Task Learning (MTL), various models were postulated to learn features that are sparse and shared across multiple tasks. However, these methods were typically evaluated by their predictive performance. To the best of our knowledge, no prior studies systematically evaluated models in terms of correctly recovering the shared signal. Here, we formalize the setting of univariate shared information retrieval, and propose ICM, an evaluation metric which can be used in the presence of ground-truth labels, quantifying 3 aspects of the learned shared features. We further propose Deep Canonical Information Decomposition (DCID) - a simple, yet effective approach for learning the shared variables. We benchmark the models on a range of scenarios on synthetic data with known ground-truths and observe DCID outperforming the baselines in a wide range of settings. Finally, we demonstrate a real-life application of DCID on brain Magnetic Resonance Imaging (MRI) data, where we are able to extract more accurate predictors of changes in brain regions and obesity. The code for our experiments as well as the supplementary materials are available at https://github.com/alexrakowski/dcid

**Keywords:** Shared Variables Retrieval · CCA · Canonical Correlation Analysis

## 1 Introduction

In this paper, we approach the problem of isolating the *shared* signal $\mathbf{Z}$ associated with two scalar *target variables* $Y_1$ and $Y_2$, from their *individual* signals $\mathbf{Z_1}$ and $\mathbf{Z_2}$, by leveraging additional, high-dimensional observations $\mathbf{X}$ (see Figure 1 for the corresponding graphical model). Analyzing the relationships between pairs of variables is ubiquitous in biomedical or healthcare studies. However, while biological traits are often governed by complex processes and can have a

wide range of causes, we rarely have access to the fine-grained, low-level signals constituting the phenomena of interest. Instead, we have to either develop hand-crafted quantities based on prior knowledge, which is often a costly process, or resort to high-level, "aggregate" measurements of the world. If the underlying signal is weak, it might prove challenging to detect associations between such aggregated variables. On the other hand, in many fields, we have access to high-dimensional measurements, such as medical scans or genome sequencing data, which provide rich, albeit "unlabeled" signal. We propose to leverage such data to **decompose** the traits of interest into their **shared** and **individual** parts, allowing us to better quantify the relationships between them.
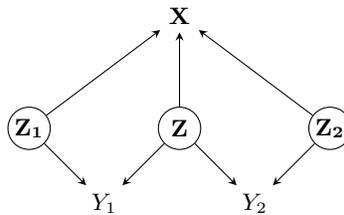


Fig. 1: A Directed Acyclic Graph (DAG) representing the graphical model used in our setting. Unobserved variables are denoted in circles. The two univariate random variables $Y_1, Y_2$ are generated by their **individual** ancestors $\mathbf{Z_1}, \mathbf{Z_2}$, and a **shared** ancestor $\mathbf{Z}$. The multivariate $\mathbf{X}$ is generated by all 3 latent variables. The following pairs of variables are independent under this model: $\mathbf{Z} \perp \mathbf{Z_1}$, $\mathbf{Z} \perp \mathbf{Z_2}$, $\mathbf{Z_1} \perp \mathbf{Z_2}$, $\mathbf{Z_1} \perp Y_2$ and $\mathbf{Z_2} \perp Y_1$.

Probabilistic CCA (pCCA) was one of the early approaches to learning the shared signal between pairs of random variables (r.v.s) [4,23]. However, its effectiveness is limited to multivariate observations - for scalar variables we can only learn the variables themselves, up to multiplication by a constant. A variety of methods from the field of Multi-Task Feature Learning (MTFL) learn feature representations of $\mathbf{X}$, which should be sparse and shared across tasks [2,25,41]. These models are typically evaluated by their predictive performance, and the shared features are rather a means of improving predictions, than a goal itself. To the best of our knowledge, no studies exist which systematically quantify how accurate are such models in recovering the signal shared between tasks.

To this end, we define the ICM score, which evaluates 3 aspects of learned shared features: *informativeness*, *completeness*, and *minimality*, when ground-truth labels are available. Furthermore, we propose Deep Canonical Information Decomposition (DCID), an approach utilizing Deep Neural Network (DNN) feature extractors and Canonical Correlation Analysis (CCA) to learn the variables $\mathbf{Z}$ shared between traits. DCID approximates the traits of interest with DNN classifiers and utilizes their latent features as multivariate decompositions of the traits. It then identifies the shared factors by performing CCA between the two

sets of latent representations and retaining the most correlated components (see Figure 3 for a graphical overview of the method).

Our contributions can be summarized as follows:

1. We define the ICM score, which allows evaluation of learned shared features in the presence of ground-truth labels (Section 3)
2. We propose DCID, a method leveraging DNN classifiers and CCA to learn the shared signal (Section 4).
3. We benchmark the proposed model, along several baselines, on a range of scenarios with synthetic data, and analyze their performance wrt. different properties of the underlying ground-truth (Sections 5.3 and 5.4).
4. Finally, we demonstrate a real-life use-case of the proposed method, by applying it on a dataset of brain Magnetic Resonance Imaging (MRI) data to better quantify the relationships between brain structures and obesity (Section 5.5).

## 2  Related Work

### 2.1  Canonical Correlation Analysis (CCA)



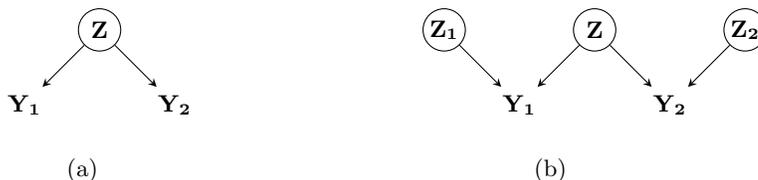(a)                                         (b)

Fig. 2: Two probabilistic interpretations of Canonical Correlation Analysis (CCA). In $(a)$, the observed variables are different, noisy views (linear transformations) of the same underlying variable $\mathbf{Z}$. In $(b)$, two additional, view-specific latent variables $\mathbf{Z_1}$ and $\mathbf{Z_2}$ are introduced, which can be interpreted as modeling the uncertainty of $p(\mathbf{Y}|\mathbf{Z})$.

CCA is a statistical technique operating on pairs of multivariate observations [21,19]. It is similar to Principal Component Analysis (PCA) [30], in that it finds linear transformations of the observations, such that the resulting variables are uncorrelated. Specifically, for a pair of observations $\mathbf{Y_1} \in \mathbb{R}^{n \times p}$ and $\mathbf{Y_2} \in \mathbb{R}^{n \times q}$, it finds linear transformations $U \in \mathbb{R}^{p \times d}$, $V \in \mathbb{R}^{q \times d}$, $d = min\{p, q\}$ which maximize the correlation between the consecutive pairs of the resulting variables $\mathbf{C_1} = \mathbf{Y_1}U$, $\mathbf{C_2} = \mathbf{Y_2}V$. In the probabilistic interpretation of CCA [4], one can interpret the observed variables as two different views of the latent variable $\mathbf{Z}$ (Figure 2a). This interpretation is extended in [23] to include view-specific variables $\mathbf{Z_1}$ and $\mathbf{Z_2}$ (Figure 2b), which is the closest to our setting.

## 2.2   Multi-Task Learning (MTL)

MTL is a machine learning paradigm, where models are fitted to several tasks simultaneously, with the assumption that such joint optimization will lead to better generalization for each task [9,42,33]. In Multi-Task Feature Learning (MTFL), one aims to learn a low-dimensional representation of the input data, that is shared across tasks [2,34]. The early approaches for MTFL worked with linear models and were based on imposing constraints on the matrix of model parameters, such as sparsity or low-rank factorization [2,3,25]. Modern approaches extend MTFL to DNN models by using tensor factorization in place of matrix factorization [40,41] or by employing adversarial training to learn task-invariant, and task-specific features [7,36,26].

# 3   Univariate Shared Information Retrieval

In this section, we formalize the problem setting (Section 3.1) and define 3 quantities measuring different aspects of the learned shared representations, which constitute the model evaluation procedure (Section 3.2).

## 3.1   Problem Setting

In our setting, we observe two univariate r.v.s $Y_1, Y_2 \in \mathbb{R}$, which we will refer to as the *target variables*, and a multivariate r.v. $\mathbf{X} \in \mathbb{R}^l$. We further define 3 unobserved, multivariate, and pairwise-independent *latent variables* $\mathbf{Z}, \mathbf{Z_1}, \mathbf{Z_2} \in \mathbb{R}^k$, which generate the observed variables. We will refer to $\mathbf{Z_1}$ and $\mathbf{Z_2}$ as the *individual variables*, and to $\mathbf{Z}$ as the *shared variables*. The main assumption of the model is that the individual variables $\mathbf{Z_1}$ and $\mathbf{Z_2}$ are each independent from one of the target variables, i.e., $\mathbf{Z_1} \perp Y_2$ and $\mathbf{Z_2} \perp Y_1$, while the shared variable $\mathbf{Z}$ is generating all the observed r.v.s, i.e., $Y_1, Y_2$ and $\mathbf{X}$. The corresponding graphical model is shown in Figure 1. Similar to the pCCA setting, we assume additivity of the effects of the shared and individual variables on $Y_i$, i.e.:

$$Y_i = \psi_i(\mathbf{Z}) + \phi_i(\mathbf{Z_i}), \ i \in \{1, 2\} \tag{1}$$

where $\psi_i$ and $\phi_i$ are arbitrary functions $\mathbb{R}^k \mapsto \mathbb{R}$.

Our task of interest is then predicting the shared variable $\mathbf{Z}$ given the observed r.v.s, i.e., learning an accurate model of $p(\mathbf{Z}|Y_1, Y_2, X)$, **without access to $\mathbf{Z}, \mathbf{Z_1}, \mathbf{Z_2}$ during training**.

## 3.2   Evaluating the Shared Representations

While in practical scenarios we assume that the latent variables remain unobserved, to benchmark how well do different algorithms recover $\mathbf{Z}$, we need to test them in a controlled setting, where all ground-truth variables are available at least during test time. Let $\mathcal{D} = \{\mathbf{x}^{(i)}, y_1^{(i)}, y_2^{(i)}, \mathbf{z}^{(i)}, \mathbf{z_1}^{(i)}, \mathbf{z_2}^{(i)}\}_{i=1}^N$ be a ground-truth dataset, and $\hat{\mathbf{z}} = \{\hat{\mathbf{z}}^{(i)}\}_{i=1}^N$ be the learned shared representations. We will

denote by $[\mathbf{x}, \mathbf{y}]$ the column-wise concatenation of $\mathbf{x}$ and $\mathbf{y}$, and by $R^2(\mathbf{x}, \mathbf{y})$ the ratio of variance explained (i.e., the coefficient of determination) by fitting a linear regression model of $\mathbf{x}$ to $\mathbf{y}$:

$$R^2(\mathbf{x}, \mathbf{y}) = 1 - \frac{1}{d} \sum_{j=1}^{d} \frac{\sum_i (e_j^{(i)})^2}{\sum_i (y_j^{(i)} - \bar{y}_j)^2} \qquad (2)$$

where $e_j^{(i)} = \hat{y}_j^{(i)} - y_j^{(i)}$ are residuals of the model of the $j$-th dimension of $\mathbf{y}$.

Inspired by the DCI score [15] from the field of disentangled representation learning, we define the following requirements for a learned representation $\hat{\mathbf{Z}}$ as correctly identifying the shared variable $\mathbf{Z}$:

1. **Informativeness**: $\mathbf{Z}$ should be predictable from $\hat{\mathbf{Z}}$. We measure this as the ratio of variance explained by a model fitted to predict $\mathbf{Z}$ from $\hat{\mathbf{Z}}$:

$$\mathcal{L}_{info}(\hat{\mathbf{z}}, \mathbf{z}) = R^2(\hat{\mathbf{z}}, \mathbf{z}) \qquad (3)$$

2. **Compactness**: $\mathbf{Z}$ should be sufficient to predict $\hat{\mathbf{Z}}$. We measure this as the ratio of variance explained by a model fitted to predict $\hat{\mathbf{Z}}$ from $\mathbf{Z}$:

$$\mathcal{L}_{comp}(\mathbf{z}, \hat{\mathbf{z}}) = R^2(\mathbf{z}, \hat{\mathbf{z}}) \qquad (4)$$

3. **Minimality**: $\hat{\mathbf{Z}}$ should only contain information about $\mathbf{Z}$. We measure this as one minus the ratio of variance explained by a model fitted to predict the individual variables $\mathbf{Z_1}$ and $\mathbf{Z_2}$ from $\hat{\mathbf{Z}}$:

$$\mathcal{L}_{min}(\hat{\mathbf{z}}, \mathbf{z_1}, \mathbf{z_2}) = 1 - R^2(\hat{\mathbf{z}}, [\mathbf{z_1}, \mathbf{z_2}]) \qquad (5)$$

The final score, $ICM$, is given as the product of the individual scores:

$$ICM(\hat{\mathbf{z}}, \mathcal{D}) = \mathcal{L}_{info}(\hat{\mathbf{z}}, \mathbf{z}) \cdot \mathcal{L}_{comp}(\mathbf{z}, \hat{\mathbf{z}}) \cdot \mathcal{L}_{min}(\hat{\mathbf{z}}, \mathbf{z_1}, \mathbf{z_2}) \qquad (6)$$

and takes values in $[0, 1]$, with 1 being a perfect score, identifying $\mathbf{Z}$ up to a rotation.

Note that minimality might seem redundant given compactness - if $\mathbf{Z}$ explains all the variance in $\hat{\mathbf{Z}}$, then, since $\mathbf{Z} \perp \mathbf{Z_1}, \mathbf{Z_2}$, $\hat{\mathbf{Z}}$ would contain no information about $\mathbf{Z_1}$ or $\mathbf{Z_2}$. However, if we accidentally choose the dimensionality of $\hat{\mathbf{Z}}$ to be much higher than that of $\mathbf{Z}$, a model can "hide" information about $\mathbf{Z_1}$ and $\mathbf{Z_2}$ by replicating the information about $\mathbf{Z}$ multiple times, e.g., $\hat{\mathbf{Z}} = \{\mathbf{Z_1}, \mathbf{Z_2}, \mathbf{Z}, \ldots \mathbf{Z}\}$. This would result in a perfect informativeness and an almost perfect compactness score, but a low minimality score.

## 4    Method: Deep Canonical Information Decomposition

In this section, we outline the difficulty in tackling the problem formulated above with CCA (Section 4.1), and describe an algorithm for solving it by exploiting the additional observed variable $\mathbf{X}$ (Section 4.2).

### 4.1   Limitations of the CCA Setting

Without $\mathbf{X}$, our setting can be seen as a special case of Probabilistic CCA (pCCA) [23], where the observed variables ("views" of the data) have a dimensionality of one (see Figure 2b). If we assume non-empty $\mathbf{Z}$, then, by the pCCA model, $\mathbf{Z_1}$ and $\mathbf{Z_2}$ would have a dimensionality of zero, resulting in degenerate solutions in form of:

$$\hat{Z} = \alpha Y_1$$
$$\vee$$
$$\hat{Z} = \alpha Y_2$$
$$\alpha \neq 0$$

(7)

as the only linear transformations of univariate $Y_1$ and $Y_2$ are the variables themselves, up to scalar multiplication.

### 4.2   Deep Canonical Information Decomposition (DCID)

In order to find $\mathbf{Z}$, we need "unaggregated", multivariate views of $Y_1$ and $Y_2$. To achieve this, we leverage the high-dimensional observations $\mathbf{X}$, e.g., images, to learn decompositions of $Y_1, Y_2$ as transformations of $\mathbf{X}$. Specifically, we assume that both $Y_i$ can be approximated as transformations of $\mathbf{X}$ with functions $h_i(\mathbf{X}) = \hat{Y}_i$. We further assume they can be decomposed as $h_i = g_i \circ f_i$, where $f_i : \mathbb{R}^l \mapsto \mathbb{R}^k$, called the *representation function*, can be an arbitrary, potentially nonlinear mapping, and $g_i : \mathbb{R}^k \mapsto \mathbb{R}$, called the *classifier function*, is a linear combination. The $k$-dimensional outputs of $f_i(\mathbf{X}) = \mathbf{B_i}$ constitute the multivariate decompositions of $Y_i$. Since these are no longer univariate, we can now apply the standard CCA algorithm on $\mathbf{B_1}, \mathbf{B_2}$ to obtain pairs of canonical variables $\mathbf{C_1}, \mathbf{C_2} \in \mathbb{R}^k$, sorted by the strength of their pairwise correlations, i.e.,

$$\forall i, j \in |k| : i < j \Rightarrow corr(\mathbf{C}_{1,i}, \mathbf{C}_{2,i}) \leq corr(\mathbf{C}_{1,j}, \mathbf{C}_{2,j}) \tag{8}$$

In order to extract the most informative features, we can select the $n$ pairs of canonical variables with correlations above a certain threshold $T$:

$$n = \underset{i \in |k|}{\operatorname{argmax}} corr(\mathbf{C}_{1,i}, \mathbf{C}_{2,i}) > T \tag{9}$$

We then take the $\hat{\mathbf{Z}} = \frac{1}{2}(\mathbf{C}_{1,1:n} + \mathbf{C}_{2,1:n})$ as our estimate of the shared $\mathbf{Z}$. The complete process is illustrated in Figure 3 and described step-wise in Algorithm 1.

**Modeling the $h_i$**  In practice, we approximate each $h_i$ by training DNN models to minimize $\mathbb{E}[Y_i - h_i(\mathbf{X})]^2$, i.e., a standard Mean Squared Error (MSE) objective. DNNs are a fitting choice for modeling $h_i$, since various popular architectures, e.g., ResNet [20], can naturally be decomposed into a nonlinear *feature extractor* (our $f_i$) and a linear *prediction head* (our $g_i$).
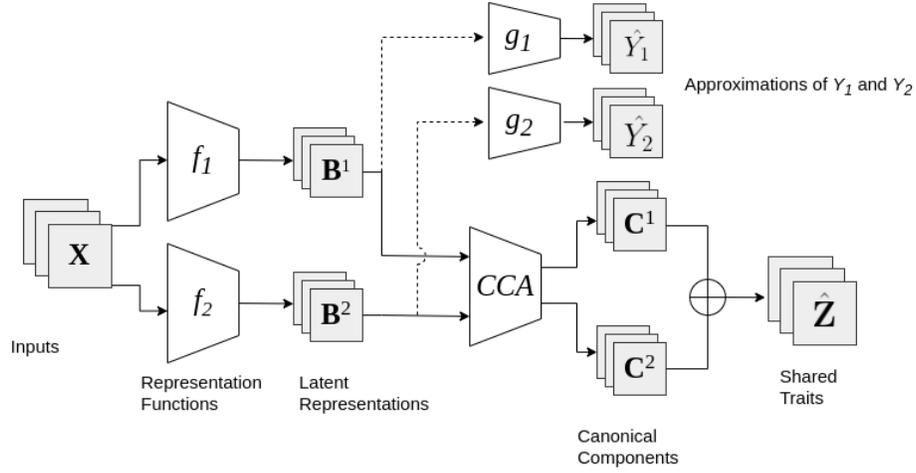
Fig. 3: A visual illustration of the Deep Canonical Information Decomposition (DCID) model. The target variables $Y_1$ and $Y_2$ are approximated by fitting DNN predictors on the high-dimensional data $\mathbf{X}$. Outputs of the penultimate layers of the networks are then used as multivariate decompositions of $Y_1$ and $Y_2$, and fed into CCA to estimate the shared signal $\hat{\mathbf{Z}}$.

---

**Algorithm 1:** DCID: Computing the shared features $\hat{\mathbf{z}}$ and the prediction function $f^\star$

---

$\textbf{Input: } \mathcal{D} = \{\mathbf{x}^{(i)}, y_1^{(i)}, y_2^{(i)}\}_{i=1}^N$ ;                    /* Training dataset */
$\textbf{Input: } \mathcal{L}(\cdot, \cdot)$ ;       /* Loss function to optimize the DNNs, e.g., $L_2$ */
$\textbf{Input: } T$ ;                       /* Canonical correlation threshold */
$\textbf{Output: } \hat{\mathbf{Z}}$ ;                  /* Features shared between $Y_1$ and $Y_2$ */
$\textbf{Output: } f^\star$ ;                /* Function to predict $\hat{\mathbf{Z}}$ from new data */
$f_1, g_1 \leftarrow \operatorname{argmin}_{h=g \circ f} \mathcal{L}(\mathbf{y}_1, h(\mathbf{x}))$ ;   /* Fit a DNN to predict $\mathbf{y}_1$ from $\mathbf{x}$ */
$f_2, g_2 \leftarrow \operatorname{argmin}_{h=g \circ f} \mathcal{L}(\mathbf{y}_2, h(\mathbf{x}))$ ;   /* Fit a DNN to predict $\mathbf{y}_2$ from $\mathbf{x}$ */
$\mathbf{b}_1 \leftarrow f_1(\mathbf{x})$ ;
$\mathbf{b}_2 \leftarrow f_2(\mathbf{x})$ ;
$U, V \leftarrow \operatorname{CCA}(\mathbf{b}_1, \mathbf{b}_2)$ ;    /* Compute the CCA projection matrices $U, V$ */
$\hat{\mathbf{z}} \leftarrow \emptyset$ ;
$n \leftarrow 1$ ;
$\textbf{while } corr(U_n^\top \mathbf{b}_1, V_n^\top \mathbf{b}_2) > T \textbf{ do}$
$\quad \mathbf{b}^n \leftarrow \frac{1}{2}(U_n^\top \mathbf{b}_1 + V_n^\top \mathbf{b}_2)$ ;
$\quad \hat{\mathbf{z}} \leftarrow \hat{\mathbf{z}} \cup \{\mathbf{b}^n\}$ ;                 /* Add a new shared component $S^n$ */
$\quad n \leftarrow n + 1$ ;
$\textbf{end}$
$n \leftarrow n - 1$ ;
$f^\star(\cdot) \leftarrow \frac{1}{2}[U_{1:n}^\top f_1(\cdot) + V_{1:n}^\top f_2(\cdot)]$ ;                 /* Save the function $f^\star$ */

---

## 5  Experiments

In Section 5.1 we describe the baseline models we compare against, and in Section 5.2 we describe the settings of the conducted experiments, such as model hyperparameters or datasets used. In Sections 5.3 and 5.4 we conduct experiments on synthetic data with know ground-truth - in Section 5.3 we benchmark the models in terms of retrieving the shared variables $\mathbf{Z}$, and in Sections 5.4 we evaluate how the performance of the models degrades when the variance explained by the shared variables changes. Finally, in Section 5.5 we demonstrate a real-life use case of the proposed DCID method on brain MRI data, where the underlying ground-truth is not known.

### 5.1  Baselines

*Multi-Task Learning (MTL)* We train a DNN model in a standard multi-task setting, i.e., to predict both $Y_1$ and $Y_2$ with a shared feature extractor $f : \mathbb{R}^l \mapsto \mathbb{R}^k$ and task-specific linear heads $g_1, g_2 : \mathbb{R}^k \mapsto \mathbb{R}$. We then select as $\hat{\mathbf{Z}}$ the set of features of $f$, for which the magnitude of normalized weights of the task-specific heads exceeds a certain threshold $T_{MTL}$ for both heads, i.e.:

$$f(\cdot)_i \in \hat{\mathbf{Z}} \implies T_{MTL} \leq \frac{|G_{1,i}|}{max|G_{1,:}|} \ \wedge T_{MTL} \leq \frac{|G_{2,i}|}{max|G_{2,:}|} \tag{10}$$

where $G_1, G_2$ are weight vectors of the linear heads $g_1, g_2$.

*Multi-Task Feature Learning (MTFL)* We train a multi-task DNN as in the MTL setting, and apply the algorithm for *sparse common feature learning* of [3] on the features of $f$. This results in new sparse features $f'$ and their corresponding new prediction heads $g'_1, g'_2$. As in the above setting, we select features of $f'$ with the magnitude of normalized weights for $g'_i$ above a threshold $T_{MTL}$.

*Adversarial Multi-Task Learning (Adv. MTL)* Introduced in [26], this model learns 3 disjoint feature spaces: 2 task-specific, *private* feature spaces, and a *shared* space, with features common for both tasks. The model is trained in an adversarial manner, with the discriminator trying to predict the task from the shared features. Additionally, it imposes an orthogonality constraint on the shared and individual spaces, forcing them to contain different information.

### 5.2  Experimental Settings

**Synthetic Data** For experiments with known ground-truth, we employed the Shapes3D dataset [8], which contains synthetic $64 \times 64$ pixel RGB images of simple 3-dimensional objects against a background, generated from 6 independent latent factors: floor hue, wall hue, object hue, scale, shape and orientation of the object, resulting in $480,000$ samples total. We take the images as $\mathbf{X}$, and select different factors as the unobserved variables $\mathbf{Z}, \mathbf{Z_1}, \mathbf{Z_2}$. As the 6 factors are

the only sources of variation in the observed data $\mathbf{X}$, it allows for an accurate evaluation of model performance in terms of retrieving $\mathbf{Z}$.

We employed the encoder architecture from [27] as the DNN backbone used to learn $f_1$ and $f_2$. The models were trained for a single pass over the dataset, with a mini-batch of size 128 using the Adam optimizer [22] with a learning rate of $10^{-4}$. We repeated each experimental setting over 3 random seeds, each time splitting the dataset into different train/validation/test split with ratios of $(70/15/15)\%$. For the hyperparameter sweep performed in Section 5.3 we considered: 10 values evenly spread on $[0, 1]$ for $T$ of DCID, 10 values evenly spaced on $[0, 1]$ for $T_{MTL}$ of MTL, MTFL, and Adv. MTL, 10 values evenly spread on a logarithmic scale of $[10^{-4}, 10]$ for the $\gamma$ parameter of MTFL, and $\gamma, \lambda \in \{0.01, 0.05, 0.1, 0.5, 1\}$ and a learning rate of the discriminator in $\{10^{-4}, 10^{-3}\}$ for Adv. MTL.

**Brain MRI Data** For the experiments on brain MRI scans (Section 5.5) we employed data from the UK Biobank (UKB) medical database [39]. Specifically, we selected data for participants who underwent the brain scanning procedure, self-identified as "white-British", and have a similar genetic ancestry, which resulted in $34,314$ data points. As the input data $\mathbf{X}$ we took the T1-weighted structural scans, which were non-linearly registered to an MNI template [28,1], and downsampled them to a size of $96 \times 96 \times 96$ voxels. For $Y_2$, we selected body mass-related measurements available in the dataset, such as the total body fat mass (BFM), weight, or body mass index (BMI). For $Y_1$, we computed the total volumes of several brain Regions of Interest (ROIs), e.g., the total volume of hippocampi or lateral ventricles, using the Synthseg software [6].

We employed a 3D MobileNetV2 [24], with a width parameter of 2, as the DNN architecture used to learn $f_1$ and $f_2$. The models were trained for 40 epochs with a mini-batch size of 12 using the Adam optimizer with a learning rate of $10^{-4}$. For each possible pair of $Y_1$ and $Y_2$ we repeated the experiments across 3 random seeds, each time selecting a different $30\%$ of the samples as the test set.

### 5.3   Learning the Shared Features Z

To evaluate how accurately do different models learn the underlying shared features $\mathbf{Z}$, we trained them in controlled settings with known ground-truth. We created the latent variables from the 6 factors of the Shapes3D dataset, by randomly selecting two individual factors $Z_1, Z_2$ and one shared factor $Z$, and constructed the target variables as $Y_1 = Z_1 + Z$ and $Y_2 = Z_2 + Z$. This resulted in 60 possible scenarios with different underlying latent variables. To ensure a fair comparison, for each model we performed a grid search over all hyperparameters on 30 random scenarios, and evaluated it on the remaining 30 scenarios using the best found hyperparameter setting.

The resulting $ICM$ scores are shown in Table 1. The proposed DCID model performed best both in terms of the final $ICM$ score, as well as the individual scores. The MTL and MTFL methods performed similarly in terms of the $ICM$

score, with MTL achieving higher *informativeness*, and MTFL a lower *minimality* score. The Adv. MTL model had the lowest $ICM$, performing well only in terms of *minimality*.

Table 1: $ICM$ scores of models trained on the Shapes3D dataset to retrieve the shared $\mathbf{Z}$. Reported are the mean and standard deviation of each score over 90 runs per model (30 scenarios $\times$ 3 random seeds).

| Model | ICM ↑ | Informativeness ↑ | Compactness ↑ | Minimality ↓ |
|---|---|---|---|---|
| **Adv. MTL** | 0.06 ($\pm$0.05) | 0.22 ($\pm$0.11) | 0.22 ($\pm$0.11) | 0.03 ($\pm$0.02) |
| **MTL** | 0.18 ($\pm$0.14) | 0.65 ($\pm$0.19) | 0.33 ($\pm$0.19) | 0.23 ($\pm$0.12) |
| **MTFL** | 0.19 ($\pm$0.15) | 0.47 ($\pm$0.23) | 0.37 ($\pm$0.22) | 0.12 ($\pm$0.08) |
| **DCID** (ours) | **0.62** ($\pm$0.15) | **0.85** ($\pm$0.07) | **0.73** ($\pm$0.17) | **0.01** ($\pm$0.02) |

### 5.4   Variance Explained by Z and Model Performance

We further investigated how the amount of variance in $Y_1, Y_2$ explained by $\mathbf{Z}$ influences model performance, with two series of experiments. In the first one, we controlled $\tau$, the ratio of variance in $Y_1, Y_2$ explained by the shared variables $\mathbf{Z}$ to the variance explained by the individual variables $\mathbf{Z_1}, \mathbf{Z_2}$, i.e.:

$$\tau = \frac{R^2(\mathbf{Z}, [Y_1, Y_2])}{R^2([\mathbf{Z_1}, \mathbf{Z_2}], [Y_1, Y_2])} \tag{11}$$

We created 15 different base scenarios, each time selecting a different pair of variables as $\mathbf{Z_1}, \mathbf{Z_2}$, and the remaining 4 as $\mathbf{Z}$. For each scenario we then varied $\tau$ 17 times on a logarithmic scale from 0.1 to 10, and trained models using their best hyperparameter settings from Section 5.3.

We plot the resulting $ICM$ scores against $\tau$ in Figure 4. Firstly, all the models fail to recover $\mathbf{Z}$ for $\tau \leq 0.3$, i.e., when the signal of $\mathbf{Z}$ is weak. For $\tau \in [0.3, 2.3]$ the DCID model outperforms the baselines by a wide margin, even for $\tau < 1$. The MTL models begin to recover $\mathbf{Z}$ only when it dominates the signal in the target variables. Interestingly, the performance of DCID drops suddenly for $\tau = 2.37$, and is outperformed by the MTL and MTFL models for $\tau > 2.8$. This is a surprising behavior, and we observed it occur independently of values of the threshold $T$ (see Figure 1 of the supplementary material).

In the second scenario, we controlled $\kappa$, the ratio of variance explained by $\mathbf{Z}$ in $Y_1$ to the variance explained in $Y_2$, i.e.:

$$\kappa = \frac{R^2(\mathbf{Z}, Y_1)}{R^2(\mathbf{Z}, Y_2)} \tag{12}$$

We created 60 base scenarios, similarly as in Section 5.3, and for each we varied $\kappa$ 6 times evenly on the scale from 0.1 to 1. Again, we selected the model hyperparameters that performed best in Section 5.3.
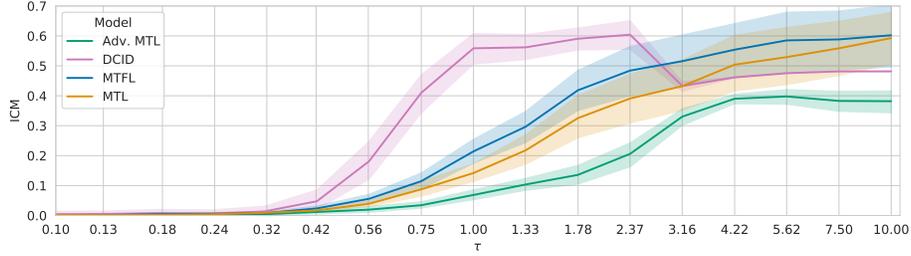
Fig. 4: *ICM* scores (y-axis) of different models plotted against $\tau$, the ratio of variance in $Y_1, Y_2$ explained by $\mathbf{Z}$ to the variance explained by $\mathbf{Z_1}, \mathbf{Z_2}$ (x-axis, logarithmic scale).

We plot the resulting *ICM* scores against $\kappa$ in Figure 5. For $\kappa \in [0.5, 1]$ the DCID model retains a consistent performance. For values of $\kappa$ below 0.5 its performance decreases linearly, for $\kappa \leq 0.2$ achieving lower *ICM* scores than the MTL and MTFL models. The scores in the low $\kappa$ regime are higher, however, than the scores for low $\tau$ values, indicating that while DCID performs best when $\mathbf{Z}$ explains a large amount of variance in both target variables, it is also beneficial if at least one of the target variables is strongly associated with $\mathbf{Z}$. The baseline models, while achieving lower scores overall, seem to have their performance hardly affected by changes in $\kappa$.
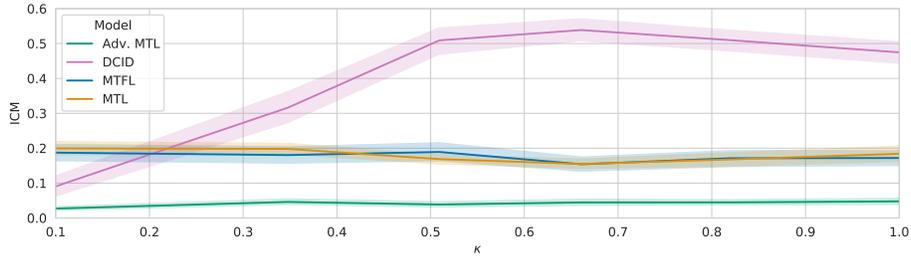


Fig. 5: *ICM* scores (y-axis) of different models plotted against $\kappa$, the ratio of variance explained by $\mathbf{Z}$ in $Y_1$ to the variance explained in $Y_2$ (x-axis). For $\kappa = 1$ the shared variables explain the same amount of variance in both target variables.

### 5.5 Obesity and the Volume of Brain Regions of Interest (ROIs)

**Background** The occurrence of neuropsychiatric disorders is associated with a multitude of factors. For example, the risk of developing dementia can depend on age [11,37], ethnicity [10,35], or genetic [38,16], vascular [12,31], and

even dietary [42,17] factors. However, only a subset of these factors are modifiable. Being aware of the genetic predisposition of an individual for developing a disorder does not directly translate into possible preventive actions. On the other hand, mid-life obesity is a known factor for dementia, which can be potentially acted upon [18,5]. Several studies analyzed the statistical relations between brain ROIs and obesity [32,14,13]. A natural limitation of these studies is the fact that they work with "aggregated" variables, quantifying obesity or ROI volumes as single values, potentially losing information about the complex traits. Phenomena such as the "obesity paradox", where obesity can have both adverse and positive effects [29], indicate the need for deeper dissecting the variables of interest and the connections between them.

**Analysis using DCID**  We approach this problem by estimating the shared signal $\mathbf{Z}$ between $Y_1$, the body fat mass (BFM), and $Y_2$, the volume of different brain ROIs. We trained several DCID models on the UKB data to predict BFM and volumes of the following ROIs: brain stem, cerebrospinal fluid (CSF), subcortical gray matter, ventricles, and the hippocampus. Additionally, we trained models for $Y_1$ being the body weight, or BMI, and report results for these in Section 2 of the supplementary material. Since the main interest lies in the effect of obesity on the ROIs, we constructed "surrogate" variables of $Y_1$, denoted by $\psi_1(\mathbf{Z})$ (see Equation 1), which isolate the shared signal in $Y_1$ from the individual one. This is a conservative approach since it only utilizes features of the model trained to predict $Y_1$, with the information about $Y_2$ used only to rotate the features and extract the shared dimensions.

First, we demonstrate how $\psi_1(\mathbf{Z})$ allows for more accurate estimates of change in the ROIs, since it ignores the signal in $Y_1$ which is independent of $Y_2$. We fitted $\mathbf{Z}$ on the training data by selecting shared components with a threshold $T > 0.2$. We then obtained predictions of $\psi_1(\mathbf{Z})$ on the test set and computed their correlation with the ROI. We report the results for all the ROIs in Table 2, and plot BFM and $\psi_1(\mathbf{Z})$ against the volume of the subcortical gray matter for a single model in Figure 6. For all ROIs the surrogate variable is correlated stronger than BFM, up to 8-fold for the ventricles, while retaining the sign of the coefficient. The smallest gains seem to be achieved for CSF, where the spread of coefficients over different runs is also the highest.

Secondly, we show how obtaining $\psi_1(\mathbf{Z})$ allows us to estimate the variance explained separately in $Y_1$ and $Y_2$, which is not possible by merely computing the correlation coefficient between $Y_1$ and $Y_2$. We plot the ratio of explained variance for each ROI in Figure 7. While $\psi_1(\mathbf{Z})$ explains a similar amount of variance for ventricles and BFM, we can see bigger disparities for other ROIs, especially for the brain stem, where the variance explained in BFM is negligible. This might indicate that predictions of the the brain stem volume from BFM would be less reliable than predictions of other ROIs.

Table 2: Pearson correlation coefficients between volumes of Regions of Interest (ROIs) in brain MRI scans (columns) and two variables - $Y_1$, being the measurements of total body fat mass (first row), and a surrogate variable $\psi_1(\mathbf{Z})$, isolating the signal of the shared variables $\mathbf{Z}$ contributing to $Y_1$. In parentheses, we report the standard deviation of the coefficients over 3 training runs over different subsets of data.

| Variable | Brain Stem | CSF | Gray Matter | Hippocampus | Ventricles |
|---|---|---|---|---|---|
| $Y_1$ | $-0.03$ ($\pm 0.01$) | $0.01$ ($\pm 0.00$) | $0.04$ ($\pm 0.00$) | $0.05$ ($\pm 0.00$) | $-0.02$ ($\pm 0.00$) |
| $\psi_1(\mathbf{Z})$ | $-0.20$ ($\pm 0.07$) | $0.06$ ($\pm 0.22$) | $0.25$ ($\pm 0.04$) | $0.22$ ($\pm 0.05$) | $-0.17$ ($\pm 0.08$) |



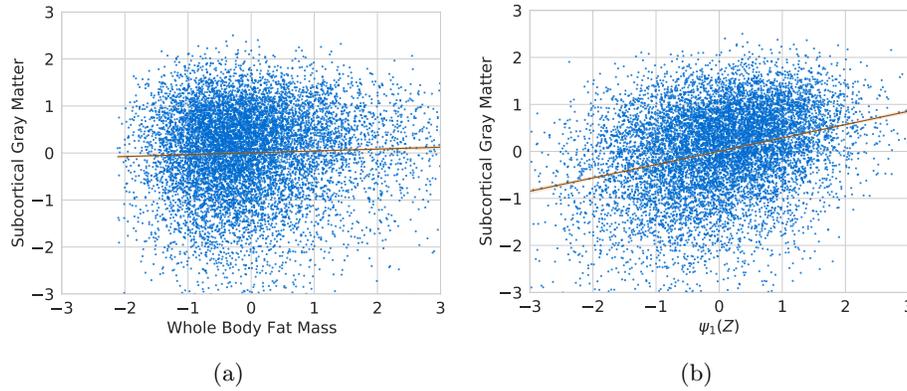(a)                                       (b)

Fig. 6: Volumes of subcortical gray matter plotted against body fat mass ($a$) and against the surrogate variable $\psi_1(\mathbf{Z})$ ($b$), for a single trained model. All variables were standardized to a z-score before plotting.
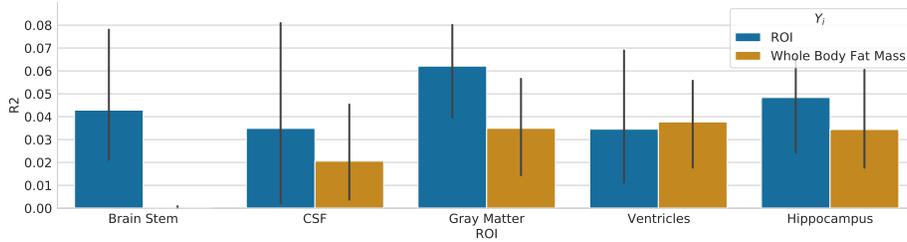


Fig. 7: Ratio of total variance explained by the surrogate variable $\psi_1(\mathbf{Z})$ in different brain Regions of Interest (ROIs) (blue bars) and in body fat mass (BFM) (yellow bars).

## 6    Discussion

In this work, we approached in a systematic manner the task of recovering the latent signal shared between scalar variables, by formalizing the problem setting and defining an evaluation procedure for model benchmarking, and proposed a new method, DCID, for solving the task.

### 6.1    Results Summary

By conducting experiments in controlled settings on synthetic data we could analyze model performance wrt. properties of the latent variables. Notably, we observed that the baseline models performed poorly when the shared variables were not strongly dominating the signal in the data, which is arguably the more realistic setting. The DCID model proved more robust in these scenarios, outperforming the baselines in most cases. We note, however, that it was still sensitive to the magnitude of the shared signal, and, interestingly, exhibited a loss of performance when the shared signal was strongly dominating. Investigating the loss of performance in the strong-signal regime, and improving robustness in the low-signal one are thus two natural directions for future work. Nevertheless, we believe that DCID can serve as an easy-to-implement, yet effective baseline.

### 6.2    Limitations

A main assumption of the method is that the observed variables $\mathbf{X}$ are rich in information, preserving the signal about the latent variables. Since, in practice, we do not observe the latent variables, we cannot test whether this assumption holds. As a substitute safety measure, we can assess the performance in predicting the observed target variables $Y$. If these cannot be predicted accurately, then it is unlikely that the model will correctly recover the latent variables either. Furthermore, we note that the method should not be mistaken as allowing to reason about causal relations between variables. It could, however, be used as part of a preprocessing pipeline in a causal inference setting, e.g., for producing candidate variables for mediation analysis.

## 7    Acknowledgements

## 8    Ethical Considerations

As mentioned in the main text (Section 5.2), we conducted the brain MRI experiments on the "white-British" subset of the UKB dataset. This was done to avoid

unnecessary confounding, as the experiments were meant as a proof of concept, rather than a strict medical study. When conducting the latter, measures should be taken to include all available ethnicities whenever possible, in order to avoid increasing the already existing disparities in representations of ethnic minorities in medical studies.

# References

1. Alfaro-Almagro, F., Jenkinson, M., Bangerter, N.K., Andersson, J.L., Griffanti, L., Douaud, G., Sotiropoulos, S.N., Jbabdi, S., Hernandez-Fernandez, M., Vallee, E., et al.: Image processing and quality control for the first 10,000 brain imaging datasets from uk biobank. Neuroimage **166**, 400–424 (2018)
2. Argyriou, A., Evgeniou, T., Pontil, M.: Multi-task feature learning. Advances in neural information processing systems **19** (2006)
3. Argyriou, A., Evgeniou, T., Pontil, M.: Convex multi-task feature learning. Machine learning **73**, 243–272 (2008)
4. Bach, F.R., Jordan, M.I.: A probabilistic interpretation of canonical correlation analysis (2005)
5. Baumgart, M., Snyder, H.M., Carrillo, M.C., Fazio, S., Kim, H., Johns, H.: Summary of the evidence on modifiable risk factors for cognitive decline and dementia: a population-based perspective. Alzheimer's & Dementia **11**(6), 718–726 (2015)
6. Billot, B., Greve, D.N., Puonti, O., Thielscher, A., Van Leemput, K., Fischl, B., Dalca, A.V., Iglesias, J.E.: Synthseg: Domain Randomisation for Segmentation of Brain MRI Scans of any Contrast and Resolution. arXiv:2107.09559 [cs] (2021)
7. Bousmalis, K., Trigeorgis, G., Silberman, N., Krishnan, D., Erhan, D.: Domain separation networks. Advances in neural information processing systems **29** (2016)
8. Burgess, C., Kim, H.: 3d shapes dataset (2018)
9. Caruana, R.: Multitask learning. Springer (1998)
10. Chen, C., Zissimopoulos, J.M.: Racial and ethnic differences in trends in dementia prevalence and risk factors in the united states. Alzheimer's & Dementia: Translational Research & Clinical Interventions **4**, 510–520 (2018)
11. Chen, J.H., Lin, K.P., Chen, Y.C.: Risk factors for dementia. Journal of the Formosan Medical Association **108**(10), 754–764 (2009)
12. Cherbuin, N., Mortby, M.E., Janke, A.L., Sachdev, P.S., Abhayaratna, W.P., Anstey, K.J.: Blood pressure, brain structure, and cognition: opposite associations in men and women. American journal of hypertension **28**(2), 225–231 (2015)
13. Dekkers, I.A., Jansen, P.R., Lamb, H.J.: Obesity, brain volume, and white matter microstructure at mri: a cross-sectional uk biobank study. Radiology **291**(3), 763–771 (2019)
14. Driscoll, I., Beydoun, M.A., An, Y., Davatzikos, C., Ferrucci, L., Zonderman, A.B., Resnick, S.M.: Midlife obesity and trajectories of brain volume changes in older adults. Human brain mapping **33**(9), 2204–2210 (2012)
15. Eastwood, C., Williams, C.K.: A framework for the quantitative evaluation of disentangled representations. In: International Conference on Learning Representations (2018)
16. Emrani, S., Arain, H.A., DeMarshall, C., Nuriel, T.: Apoe4 is associated with cognitive and pathological heterogeneity in patients with alzheimer's disease: a systematic review. Alzheimer's Research & Therapy **12**(1), 1–19 (2020)

17. Frausto, D.M., Forsyth, C.B., Keshavarzian, A., Voigt, R.M.: Dietary regulation of gut-brain axis in alzheimer's disease: Importance of microbiota metabolites. Frontiers in neuroscience **15** (2021)
18. Gorospe, E.C., Dave, J.K.: The risk of dementia with increased body mass index. Age and ageing **36**(1), 23–29 (2007)
19. Hardoon, D.R., Szedmak, S., Shawe-Taylor, J.: Canonical correlation analysis: An overview with application to learning methods. Neural computation **16**(12), 2639–2664 (2004)
20. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016)
21. Hotelling, H.: Relations between two sets of variates. In: Breakthroughs in statistics, pp. 162–190. Springer (1992)
22. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)
23. Klami, A., Kaski, S.: Probabilistic approach to detecting dependencies between data sets. Neurocomputing **72**(1-3), 39–46 (2008)
24. Köpüklü, O., Kose, N., Gunduz, A., Rigoll, G.: Resource efficient 3d convolutional neural networks. In: 2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW). pp. 1910–1919. IEEE (2019)
25. Kumar, A., Daume III, H.: Learning task grouping and overlap in multi-task learning. arXiv preprint arXiv:1206.6417 (2012)
26. Liu, P., Qiu, X., Huang, X.J.: Adversarial multi-task learning for text classification. In: Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). pp. 1–10 (2017)
27. Locatello, F., Bauer, S., Lucic, M., Raetsch, G., Gelly, S., Schölkopf, B., Bachem, O.: Challenging common assumptions in the unsupervised learning of disentangled representations. In: international conference on machine learning. pp. 4114–4124. PMLR (2019)
28. Miller, K.L., Alfaro-Almagro, F., Bangerter, N.K., Thomas, D.L., Yacoub, E., Xu, J., Bartsch, A.J., Jbabdi, S., Sotiropoulos, S.N., Andersson, J.L., et al.: Multimodal population brain imaging in the uk biobank prospective epidemiological study. Nature neuroscience **19**(11), 1523–1536 (2016)
29. Monda, V., La Marra, M., Perrella, R., Caviglia, G., Iavarone, A., Chieffi, S., Messina, G., Carotenuto, M., Monda, M., Messina, A.: Obesity and brain illness: from cognitive and psychological evidences to obesity paradox. Diabetes, metabolic syndrome and obesity: targets and therapy pp. 473–479 (2017)
30. Pearson, K.: Liii. on lines and planes of closest fit to systems of points in space. The London, Edinburgh, and Dublin philosophical magazine and journal of science **2**(11), 559–572 (1901)
31. Prabhakaran, S.: Blood pressure, brain volume and white matter hyperintensities, and dementia risk. JAMA **322**(6), 512–513 (2019)
32. Raji, C.A., Ho, A.J., Parikshak, N.N., Becker, J.T., Lopez, O.L., Kuller, L.H., Hua, X., Leow, A.D., Toga, A.W., Thompson, P.M.: Brain structure and obesity. Human brain mapping **31**(3), 353–364 (2010)
33. Ruder, S.: An overview of multi-task learning in deep neural networks. arXiv preprint arXiv:1706.05098 (2017)
34. Schölkopf, B., Platt, J., Hofmann, T.: Multi-task feature learning (2007)
35. Shiekh, S.I., Cadogan, S.L., Lin, L.Y., Mathur, R., Smeeth, L., Warren-Gash, C.: Ethnic differences in dementia risk: a systematic review and meta-analysis. Journal of Alzheimer's Disease **80**(1), 337–355 (2021)

36. Shinohara, Y.: Adversarial multi-task learning of deep neural networks for robust speech recognition. In: Interspeech. pp. 2369–2372. San Francisco, CA, USA (2016)
37. Stephan, Y., Sutin, A.R., Luchetti, M., Terracciano, A.: Subjective age and risk of incident dementia: Evidence from the national health and aging trends survey. Journal of Psychiatric Research **100**,  1–4 (2018)
38. Strittmatter, W.J., Saunders, A.M., Schmechel, D., Pericak-Vance, M., Enghild, J., Salvesen, G.S., Roses, A.D.: Apolipoprotein e: high-avidity binding to beta-amyloid and increased frequency of type 4 allele in late-onset familial alzheimer disease. Proceedings of the National Academy of Sciences **90**(5), 1977–1981 (1993)
39. Sudlow, C., Gallacher, J., Allen, N., Beral, V., Burton, P., Danesh, J., Downey, P., Elliott, P., Green, J., Landray, M., et al.: Uk biobank: an open access resource for identifying the causes of a wide range of complex diseases of middle and old age. PLoS medicine **12**(3), e1001779 (2015)
40. Wimalawarne, K., Sugiyama, M., Tomioka, R.: Multitask learning meets tensor factorization: task imputation via convex optimization. Advances in neural information processing systems **27** (2014)
41. Yang, Y., Hospedales, T.: Deep multi-task representation learning: A tensor factorisation approach. arXiv preprint arXiv:1605.06391 (2016)
42. Zhang, H., Greenwood, D.C., Risch, H.A., Bunce, D., Hardie, L.J., Cade, J.E.: Meat consumption and risk of incident dementia: cohort study of 493,888 uk biobank participants. The American journal of clinical nutrition **114**(1), 175–184 (2021)