Contents lists available at ScienceDirect

NeuroImage

journal homepage: www.elsevier.com/locate/neuroimage

Predicting future depressive episodes from resting-state fMRI with generative embedding



euroImag

Herman Galioulline^{a,*}, Stefan Frässle^a, Samuel J. Harrison^a, Inês Pereira^a, Jakob Heinzle^{a,1}, Klaas Enno Stephan^{a,b,1}

^a Translational Neuromodeling Unit (TNU), Institute for Biomedical Engineering, University of Zurich & ETH Zurich, Zurich 8032, Switzerland ^b Max Planck Institute for Metabolism Research, Cologne, Germany

ARTICLE INFO

Keywords: Depression Prediction Early Detection Generative Embedding UK Biobank Translational Neuromodeling Computational Psychiatry

ABSTRACT

After a first episode of major depressive disorder (MDD), there is substantial risk for a long-term remittingrelapsing course. Prevention and early interventions are thus critically important. Various studies have examined the feasibility of detecting at-risk individuals based on out-of-sample predictions about the future occurrence of depression. However, functional magnetic resonance imaging (fMRI) has received very little attention for this purpose so far.

Here, we explored the utility of generative models (i.e. different dynamic causal models, DCMs) as well as functional connectivity (FC) for predicting future episodes of depression in never-depressed adults, using a large dataset (N = 906) of task-free ("resting state") fMRI data from the UK Biobank (UKB). Connectivity analyses were conducted using timeseries from pre-computed spatially independent components of different dimensionalities. Over a three-year period, 50% of selected participants showed indications of at least one depressive episode, while the other 50% did not. Using nested cross-validation for training and a held-out test set (80/20 split), we systematically examined the combination of 8 connectivity feature sets and 17 classifiers. We found that a generative embedding procedure based on combining regression DCM (rDCM) with a support vector machine (SVM) enabled the best predictions, both on the training set (0.63 accuracy, 0.66 area under the curve, AUC) and the test set (0.62 accuracy, 0.64 AUC; p < 0.001). However, on the test set, rDCM was only slightly superior to predictions based on FC (0.59 accuracy, 0.61 AUC). Interpreting model predictions based on SHAP (SHapley Additive exPlanations) values suggested that the most predictive connections were widely distributed and not confined to specific networks. Overall, our analyses suggest (i) ways of improving future fMRI-based generative embedding approaches for the early detection of individuals at-risk for depression and that (ii) achieving accuracies of clinical utility may require combination of fMRI with other data modalities.

1. Introduction

Major depressive disorder (MDD) causes tremendous personal suffering and, amongst all medical conditions, has one of the highest burden of disease globally (GBD Mental Disorders Collaborators, 2022; Vos et al., 2020). It has a profoundly negative impact on social and occupational functions (Adler et al., 2006; Kupferberg et al., 2016) and is associated with increased risk for other mental and somatic (e.g. cardiovascular) disorders. After the onset of a first episode of MDD, there is a substantial risk for a long-term remitting-relapsing course (Eaton et al., 2008), accompanied by prolonged trial-and-error treatment attempts (Correll et al., 2017; Steffen et al., 2020). Prevention and early interventions are thus crucial for reducing the burden of MDD, both at an individual and societal level (Cuijpers et al., 2012, 2021). The challenge is to detect at-risk individuals early so that preventive measures and interventions can be administered in a timely and targeted fashion.

Detecting at-risk individuals requires prediction models that enable out-of-sample predictions about the future occurrence of (symptoms of) depression with clinically adequate accuracy. In the recent past, there have been numerous attempts to establish such models both in adolescents and adults, based on combinations of various data types, e.g. demographic, socioeconomic, cognitive, and clinical variables as well as motor activity (Caldirola et al., 2022; Chikersal et al., 2021; Gu et al., 2020; King et al., 2008; Librenza-Garcia et al., 2021; Lin et al., 2022; Na et al., 2020; Rocha et al., 2021; Rosellini et al., 2020; Sampson et al., 2021; van Eeden et al., 2021; Voorhees et al., 2008; Xu et al., 2019).

* Corresponding author.

https://doi.org/10.1016/j.neuroimage.2023.119986.

Received 16 October 2022; Received in revised form 15 February 2023; Accepted 25 February 2023 Available online 22 March 2023.

1053-8119/© 2023 The Authors. Published by Elsevier Inc. This is an open access article under the CC BY license (http://creativecommons.org/licenses/by/4.0/)



E-mail address: galioulline@biomed.ee.ethz.ch (H. Galioulline).

¹ Shared last authorship.

Neuroimaging has played a minor role in this endeavour so far. This may be partly due to difficulties of obtaining datasets that are longitudinal in nature and sufficiently large to allow for robust out-ofsample predictions. Several longitudinal magnetic resonance imaging (MRI) studies of depressive symptoms do exist (e.g. Barch et al. 2019; Pagliaccio et al. 2014; Papmeyer et al. 2016; Shapero et al. 2019), but almost all have small to moderate sample sizes and employ withinsample association analyses. However, association is not prediction: prediction requires out-of-sample analyses, i.e. " ... testing of the model on data separate from those used to estimate the model's parameters" (Poldrack et al., 2020). A recent exception is the study by Toenders et al. (2022) which predicted depression onset out-ofsample, based on structural MRI (and other) data from a large sample of 544 adolescents. Concerning functional MRI (fMRI), however, we are aware of only one previous fMRI study (Hirshfeld-Becker et al., 2019) that has attempted out-of-sample predictions of future depressive episodes in hitherto depression-free individuals, albeit with a small sample (total N = 33). The predictive value of fMRI for identifying individuals at risk for future depression is thus not well known.

One might wonder why fMRI should be considered at all for establishing predictor models of depressive episodes, given that fMRI data are more difficult to obtain and more costly than many other types of measurements. There are several reasons why fMRI - and particularly generative models for estimating connectivity - may have particular utility for clinical predictions. First, fMRI may afford high sensitivity since it assesses the functional status quo of neural circuits (Stephan et al., 2015), the biological level that is closest to psychiatric symptoms (Gordon, 2016). Second, clinical predictions are most valuable if they afford a mechanistic interpretation (Stephan et al., 2017); for example, this may guide the development of novel treatments. Analyses of functional interactions based on fMRI can potentially give insights into circuit mechanisms that increase risk for depression. Ideally, this requires generative models which offer an explanation how activity distributed throughout a circuit could have been generated (Stephan et al., 2015) and provide estimates of effective (directed) connectivity.

An approach that blends generative modeling with prediction is "generative embedding" (GE) (Brodersen et al., 2011, 2014; Frässle et al., 2020; Stephan et al., 2017). GE uses parameter estimates of a system (circuit) of interest, obtained by inverting a generative model, as features for subsequent machine learning (ML). This often improves prediction accuracy since the parameter estimates of a generative model offer a low-dimensional, de-noised representation of neural dynamics. Furthermore, provided the generative model is biologically plausible, GE may reveal which biological processes or properties (e.g. specific connections in a neural circuit) are most relevant for successful clinical predictions.

In this study, we used a large dataset (N = 906) of task-free ("resting state") fMRI data from the UK Biobank (Miller et al., 2016) to explore the utility of fMRI-based connectivity measures for predicting future episodes of depression in never-depressed adults. Over a three-year follow-up period, half of the selected participants (N = 453) exhibited at least one indicator of a depressive episode, according to clinical records and/or self-report, while the other half remained free from depressive episodes. Both groups were carefully matched with regard to 7 potentially confounding variables (age, sex, handedness, tobacco, alcohol, il-licit drugs, cannabis).

We emphasise that the goal of this work was not to test whether predictions based on fMRI data are better or worse than predictions based on other data types, e.g. socioeconomic or clinical variables. Instead, because there are numerous options of utilising fMRI for predictive analyses, this initial study focused on fMRI only and assessed the relative performance of different connectivity approaches – including generative embedding based on different variants of dynamic causal modeling (DCM; Friston et al. 2003) as well as functional connectivity (FC) – for predicting future depressive episodes. Concretely, in our training set (*N* = 724), we systematically combined different connectivity approaches with different ML classifiers, using nested cross-validation, and tested how well they predicted the occurrence of at least one indicator of a depressive episode over a follow-up period of three years. We then used the best-performing combination to make the same prediction in a heldout test set (N = 182) that was completely independent from the training data. Notably, predicting the occurrence of indicators of depressive episodes represents a more challenging scenario than predicting a full clinical diagnosis of MDD. Our study can thus be seen as a "stress test" whether fMRI-based assessments of connectivity, and generative models in particular, are likely to be useful at all for early detection of at-risk individuals.

2. Materials and methods

The following sections describe the dataset and methodology used in this study. Briefly, the data consist of task-free fMRI measurements (i.e. unconstrained cognition or "resting state") and questionnaire data from the UK Biobank (www.ukbiobank.ac.uk). Based on entries in the UK Biobank, we selected participants that had good quality fMRI recordings and consistent questionnaire information that allowed us to assign them to one of two groups: a group that initially had no signs of depressive episodes but exhibited indicators of depressive episodes (e.g. questionnaire data, prescription of antidepressants) within three years after the fMRI session (D+ group), or a control group that did not show any such indicators during the same period (D- group).

We used different connectivity metrics (different variants of DCM as well as functional connectivity, FC) in combination with different ML classifiers for prediction of future indicators of depressive episodes. DCM and FC analyses were applied to time series of "resting-state" fMRI networks (with 6, 21, or 55 nodes) defined by independent components analysis (ICA) of the preprocessed "resting-state" fMRI (rs-fMRI) data provided by UK Biobank. Posterior parameter estimates (DCM) and Pearson correlation coefficients (FC), respectively, served as input features to various discriminative classifiers. The classifiers were trained using nested cross-validation to avoid overfitting and to provide the best possible estimate of generalisability. Finally, the best models were chosen, and a prediction was made on held-out (and completely independent) test data.

It is worth noting that our analysis was pre-specified in an ex ante analysis plan, prior to performing any of the analyses. The analysis plan was time-stamped by uploading it to the Git repository of the Translational Neuromodeling Unit (TNU); it is available at https://gitlab.ethz.ch/tnu/analysis-plans/galioullineetal_ ukbb_pred_depr. Furthermore, code reviews were performed by three of the co-authors (SF, SH and JH) who were not involved in the data analysis, both before the beginning of the analysis of the training data, and once again before running models on the test data.

2.1. Dataset: groups with/without depressive episodes

The process of data extraction from the UK Biobank is summarized by Fig. 1. To avoid confusion, it is worth explaining that participants of the neuroimaging branch of UK Biobank (which started in 2014) underwent two fMRI scans, approx. three years apart, each of which involved both task fMRI and rs-fMRI data. In this study, we only used the rs-fMRI data acquired during the first scan.

Overall, selected individuals were required to have rs-fMRI data of good quality (as indicated by UK Biobank quality control) and no indication of any previous or current depressive episodes at the time of their first fMRI scan. From the subset of participants that fulfilled these criteria, we aimed to select two groups, one of which continued to indicate no signs of depressive episodes (D- group) three years after their first scan, and one that showed at least one indicator for at least one depressive episode over this three-year period (D+ group).

Table 1

Demographic data of participants. Values are provided as mean \pm standard deviation, with the range in brackets. Alcohol-intake frequency ranges from 1 (daily or almost daily) to 6 (never). Tobacco smoking frequency ranges from 1 (smoked on most or all days) to 4 (never smoked). Historical cannabis consumption ranges from 0 (never) to 4 (more than 100 times). Ongoing drug addiction to illicit drugs was not included because none of the participants answered yes to that question. Please note that the counts of sex and handedness are not perfectly identical between groups; this is because for 55 D+ individuals a perfect D- match could not be found and one of the seven matching variables was allowed to deviate (see main text).

Demographic & biological variables	D+	D-	
Number of participants	453	453	
Age at scan	62.33 ± 7.32 (45-78)	62.75 ± 7.02 (48-78)	
Sex (male/female)	213 (47.12%) / 239 (52.88%)	215 (47.46%) / 238 (52.54%)	
Handedness (left/right)	40 (8.85%) / 404 (89.38%)	35 (7.73%) / 413 (91.17%)	
Weight (kg)	77.04 ± 13.84	74.36 ± 14.93	
BMI (kg•m ⁻²)	26.99 ± 4.49	25.99 ± 4.27	
Alcohol-intake frequency	2.95 ± 1.49	2.88 ± 1.45	
Tobacco smoking frequency	2.85 ± 1.25	2.88 ± 1.24	
Historical cannabis consumption	0.35 ± 0.87	0.31 ± 0.79	
Townsend deprivation index	-2.09 ± 2.5	-2.33 ± 2.37	

Concretely, we first identified participants who had both task (UKB field 20249-2.0) and "resting-state" (UKB field 20227-2.0) fMRI scans in NIFTI format, ensuring quality controlled images already preprocessed by UK Biobank (Alfaro-Almagro et al., 2018), resulting in 35,485 participants. In order to define the D- group, we chose the subset of participants who responded "no" to the questionnaire item "Looking back over your life, have you ever had a time when you were feeling depressed or down for at least a whole week?" (UKB field 4598-2.0) when they were first scanned (2014+). This resulted in 15,739 participants. We further shrunk this set by selecting those individuals who continued to show no evidence of depressive episodes in following years (2014 to 2019) and again replied, at the second fMRI session in 2019+, "no" to the previous question (UKB field 4598-3.0). This resulted in 1,085 potential D- participants who could be searched for matching criteria once the D+ group had been determined.

Concerning the D+ group, we also selected individuals from the set of 15,739 participants who had preprocessed imaging data and who – during their first imaging questionnaire (2014+) – indicated never having had a depressive episode. Since the UK Biobank does not include information about the absence or presence of a clinical diagnosis of depression for all participants, we used multiple sources of information to identify indicators of depressive episodes. Specifically, we searched selected UK Biobank data fields which plausibly indicated the occurrence of at least one depressive episode in the years after the first fMRI scan. The following list summarizes the data fields in UKB and number of hits.

- Medical records in UKB:
 - First Clinically Recorded Depressive Episode [UKB 130894] (5 hits)
 - Clinical Depression-Related Encounter [UKB 41270] (31 hits)
 - Prescription of Antidepressants [UKB 20003] (6 hits)
 - Depression Diagnosis Report in UKB Assessment [UKB 20002] (12 hits)
- Self-report data in UKB:
 - Depressed for at Least a Week Report [UKB 4598-3.0] (203 hits)
 - Depression Diagnosis Report in Mental Health Questionnaire [UKB 20544] (90 hits)
 - High Score (Coleman et al., 2020 supplementary material) on CIDI in Mental Health Questionnaire [UKB 20446] (165 hits)
 - High Score (sum > 4) on Patient Health Questionnaire 3-subset [UKB 2050, 2060, 2080] (6 hits)

Overall, this resulted in 518 potential D+ participants. Since for any given participant a previous depressive episode could be reflected by

multiple hits, we took the union of the above 8 sets of hits. This resulted in a total of 464 participants in the D+ group.

Having completed the initial definition of D+ and D- groups, we searched for data entries showing inconsistent or logically incompatible responses from participants (e.g. participants stating "never depressed for at least a week" but with a clinical report of depression). This process led to the removal of 9 participants in total, resulting in 455 participants in the D+ group and 1076 participants in the D- group.

2.2. Matching of participants and definition of training/test sets

To minimize any effects of potentially confounding variables and facilitate interpretation of the classifiers' predictions, we matched the two groups with respect to multiple variables. We used matching instead of other strategies to address the influence of confounders, for two reasons. First, UK Biobank offered a large dataset of D- individuals to make precise matching possible with minimal loss of data; second, some of our connectivity analyses were computationally very expensive, rendering other strategies (such as a repeated random sampling) a nonviable option. Specifically, for each D+ participant we tried to find a matching Dparticipant according to the following seven criteria (where a tolerance range was only allowed for age, as indicated):

- Sex (UKB field 31)
- Age ± 5 years (UKB field 34)
- Handedness (UKB field 1707)
- Tobacco smoking frequency (UKB field 1249)
- Alcohol consumption frequency (UKB field 1558)
- Ongoing addiction or dependence on illicit or recreational drugs (UKB field 20457)
- Historical cannabis consumption (UKB field 20453)

All but 57 D+ participants could be matched exactly. Out of these, 55 could be matched almost exactly, with at most one criterion deviating. Two D+ participants could not be matched and were excluded from further analyses. This provided us with a dataset of 906 participants in total: 453 D+ participants and 453 matched D- participants. Demographics data on these participants can be seen in Table 1.

Finally, we performed an 80/20 split to partition the data into training and test sets. Both datasets were strictly separated from each other during data analysis to prevent any leakage of information that could affect the prediction results. We also addressed an unlikely, but theoretically possible, information leakage stemming from UK Biobank itself: the templates of major functional networks in the brain (Miller et al., 2016) which are offered by UK Biobank and which our study used for data extraction had been created using rs-fMRI data from the first 4,181 individuals in UKB. We resolved this potential problem by ensuring that



Fig. 1. CONSORT flow diagram describing the assignment of UK Biobank participants to D+ and D- groups in this study. T1: a participant's first neuroimaging session (2014 onwards). T2: a participant's second neuroimaging session (2019 onwards) which included further questionnaires and self-report. N/A: indicates a participant who either did not know or preferred to not answer a question. Identifying unique cases: this was necessary because some participants showed more than one indicator of a depressive episode after T1. Here, we selected each participant once, to avoid double entries. Inconsistent cases: participants who could have been assigned to both the D+ and D- groups, due to conflicting information. One-to-one matching: we matched each participant in D+ with a participant in D- with regard to 7 variables (see main text). "Depressed for a week" self-report: refers to UKB field 4598. "Depr. at T2" category: contains the set of participants who showed indication of having experienced at least one depressive episode at T2, based on self-report, antidepressant use, medical diagnosis, or PHQ (UKB fields 4598-3.0, 20002, 20003, 2050-3.0, 2060-3.0, 2080-3.0). "MHQ CIDI" category: refers to a high score (Coleman et al., 2020) on the Composite International Diagnostic Interview (CIDI) part of the online UKB Mental Health Questionnaire (MHQ). "MHQ Diag.": refers to self-report of a depression diagnoses in the MHQ (UKB field 20544). "Hospital": refers to ICD10 depression diagnoses from hospital in-patient data (UKB field 41270). "UKB Diag.": refers to a similar field derived by UKB indicating depressive episodes (UKB field 130894).

all participants from this set that were also part of our extracted data were assigned to the training set. This resulted in patient/control training sets with 362 individuals each and test sets with 91 individuals (Fig. 1).

2.3. FMRI data analysis

Wherever possible we used data that is directly available on UK Biobank and did not require additional processing. The rs-fMRI data are of 6 minute duration (490 images, TR=0.735 s), with a spatial resolution of 2.4 mm isotropic, and were acquired with 8x multislice acceler-

ation (Alfaro-Almagro et al., 2018). We used the data after the standard preprocessing pipeline executed by UK Biobank. The processing steps performed at UK Biobank included realignment, EPI distortion correction, and high-pass temporal filtering (with a 50s cut-off). For removing noise (including head motion effects), rs-fMRI data were further processed with single-subject spatial ICA decomposition using FIX (Salimi-Khorshidi et al., 2014; Griffanti et al., 2014) in FSL (Jenkinson et al., 2012). The resulting independent components (ICs) were classified as signal vs. noise, and a cleaned version of the data was provided by removing the noise components. UK Biobank then used these data as input to a dual regression procedure (Nickerson et al., 2017) based on a group-

level template of "resting-state" networks (derived from spatial ICA applied to data of 4,181 subjects) at dimensionalities of either 25 or 100. For subsequent analyses, 21/25 and 55/100 ICs were kept, as the other components had been found in previous work to be "... clearly identifiable as artefactual (i.e., not neuronally driven)" (Alfaro-Almagro et al., 2018). Given the importance of head motion effects for "resting state" connectivity analyses (e.g. Van Dijk et al. 2012), we double-checked whether any such group differences existed in our dataset despite the rigorous preprocessing and quality control procedures of the UK Biobank pipeline (Alfaro-Almagro et al., 2018). Applying a Mann-Whitney U test (scipy.stats.mannwhitneyu from SciPy) to subject-wise measures of fMRI head motion (UKB data field 25741), we failed to detect any significant group differences (p = 0.15; U = 3774.0).

Importantly, the final timeseries resulting from this dual regression approach are based on spatial ICA and are therefore not temporally independent, which allows for subsequent application of functional and effective connectivity methods. Furthermore, an advantage of this dual regression approach is that group-level information serves as a template which guides identification of resting state networks in individual subjects and provides subject-specific timeseries and spatial components (Nickerson et al., 2017). By contrast, relying on single-subject ICA for each participant would be potentially problematic because it would not always be possible to match the resulting components across subjects.

ICs of resting-state data can be thought of as distinct functional networks (Smith et al., 2009), and interactions between these networks can be investigated by applying functional and effective connectivity methods to IC timeseries (for previous examples, see Goulden et al., 2014; Hyatt et al., 2015; Motlaghian et al., 2022). In this work, we selected three sets of networks, which differed in the number of ICs included. Since we were interested in major functional networks implicated in depression (Brakowski et al., 2017; Kaiser et al., 2015), our first IC selection targeted the default mode network (DMN), central executive network (CEN), salience network (SN), and the dorsal attention network (DAN). We selected the corresponding components from the 21 components set provided by UK Biobank. The DMN and DAN are mapped (Miller et al., 2016) to IC indices 1 and 3, respectively, while the left/right SN and left/right CEN are mapped (Gratton et al., 2018; Shen et al., 2018) to IC indices 6, 5 and 13, 21, respectively. Furthermore, we considered full IC sets of size 21 and 55 components (as provided by UK Biobank) to explore the impact of increasing the number of networks/ICs on prediction performance. The 55 components can be interrogated interactively via a web-based visualisation tool provided by UK Biobank: https://www.fmrib.ox.ac.uk/ukbiobank/group_means/rfMRI_ICA_ d100_good_nodes.html.

2.4. Generative embedding

Having completed the selection of timeseries, our analysis proceeded to generative embedding (GE). GE requires two choices: (i) a generative model, and (ii) a ML method that uses posterior estimates from the generative model as features.

Dynamic causal modelling (DCM) is a generative modelling approach for estimating effective (directed) connectivity in neuronal networks from neuroimaging data (Friston et al., 2003). It relies on specifying a dynamical system that describes a network of neuronal populations (regions) and how they influence each other through synaptic connections. While the original formulation relies on external inputs to "drive" the system, more recent versions of DCM can be applied to resting state data. For example, stochastic DCM (Li et al., 2011) uses stochastic differential equations (SDEs) to allow for state noise at the neuronal level and thus incorporates neuronal fluctuations. One major challenge stochastic DCM faces is the computationally intensive inference of its hidden states. Spectral DCM (Friston et al., 2014) circumvents this problem by adopting a different approach which rests on a parameterized form for

endogenous fluctuations. This model assumes stationarity and generates cross-spectral densities (the equivalent of the cross-covariance function in the frequency domain). This corresponds to a deterministic system which can be solved more efficiently than the SDEs of spectral DCM. Having said this, parameter inference for spectral DCM is still computationally relatively expensive, which makes it difficult to scale the model to very large networks. Finally, regression DCM (Frässle et al., 2017; Frässle et al., 2021) reformulates the neuronal state equation of a linear DCM as a Bayesian regression problem in the frequency domain. This renders model inversion extremely fast and allows rDCM to scale to large networks with hundreds of regions.

Our analysis considered all of the DCM variants described above, i.e. stochastic DCM (Li et al., 2011), spectral DCM (Friston et al., 2014), and regression DCM (Frässle et al., 2021). For all models and all IC sets, we assumed a fully connected network. As a reference, we also obtained functional connectivity estimates, based on Pearson correlation coefficients.

To invert stochastic DCMs, we used the *spm_dcm_estimate* function in SPM12, with a DCM struct as input which had its *Y*.*y* set to the 6 timeseries, *a* set to a 6×6 matrix of ones (fully connected network of endogenous connections), and *Y*.*dt* set to 0.735 (interscan interval). This resulted in a 6×6 matrix of effective connectivity estimates, giving us 36 features for subsequent ML. Due to its high computational complexity, it was not possible to run stochastic DCM with 21 and 55 IC timeseries.

For spectral DCM, we used the SPM12 function *spm_dcm_fmri_csd* with the same exact DCM struct as for stochastic DCM as input. The performance was notably faster than for stochastic DCM, but given that it still took a few hours to run on the Euler high-performance computing cluster of ETH Zurich and that the scaling of the computational complexity is supra-linear in the number of ICs (i.e., number of nodes in the DCM), we estimated that it would still take weeks or even months to run the entire analysis (i.e., inversion of the DCMs for all subjects) for 21 or 55 IC timeseries. Hence, just like in the stochastic DCM case, we restricted the spectral DCM analysis to 6 IC timeseries.

Concerning rDCM, its high computational efficiency enabled us to analyse networks consisting of more components (6, 21, and 55 ICs), resulting in 36, 441, and 3025 features, respectively. We used the rDCM code in TAPAS 4.0 (Frässle et al., 2021), with *Y*.*y* set to the respective time series, and *Y*.*dt* set to 0.735.

Finally, FC matrices were computed using the *corrcoef* function in MATLAB. Since these matrices are symmetric along the diagonal, and the diagonal is always 1, we took the upper triangle of these matrices to be our features, resulting in 15, 210, and 1485 features for the respective IC sets. It is important to note that FC does not capture any information about the directionality of connections, as opposed to the effective connectivity measures from the DCM variants described above.

2.5. Classification

From the previous generative modeling, we had eight feature sets in place – functional connectivity for each IC set (6, 21, 55), stochastic and spectral DCM for 6 ICs each, and three rDCM feature sets for 6, 21 and 55 ICs. These feature sets were subsequently used as input to discriminative classifiers. Initially, we restricted all analyses to the training set data, and only touched the test data once we had selected a feature set / classifier combination that performed best. Regardless of the specific classifier chosen, the steps taken to arrive at reported metrics are the same.

Classifier training was performed using nested cross-validation (CV). Nested CV provides robustness against overfitting by optimizing hyperparameters in an inner CV loop while averaging the performance against other partitions of the data in an outer CV loop (Cawley and Talbot, 2010; Stone, 1974). In our case, we used 10 folds in the outer loop, and 5 folds in the inner loop. At the beginning of each iteration of the outer loop (before training with hyperparameter optimization),

the confounds (sex, age, handedness, smoking, alcohol, illicit drugs, cannabis) were linearly regressed out using scikit-learn's *LinearRegression* module. Then the data were normalized using the *StandardScaler* module and then the classifier was finally fit with the *GridSearchCV* module. This procedure yielded a set of performance measures for each feature/classifier combination (see Results).

After evaluation of the feature/classifier pairs on the training data, there are several possibilities how models fitted on the training data could be applied to the test data. First, we evaluated whether the feature set/classifier combination that had performed best on the training set generalized to the test set. Second, we performed a post hoc analysis in which we examined each feature set together with the classifier that had been optimal for this specific feature set on the training data.

In addition to evaluating classifiers based on their performance metrics, we also ran permutation tests to check for statistical significance of the classification results. These tests were run both on our training and test set. To generate an empirical null distribution for a given feature/classifier pair, we randomly permuted the labels while considering subject pairs between the D- and D+ groups, originating from the matching of confounds. This is done by identifying a pair and flipping their labels with 0.5 probability. For the resulting permuted labels, the classifier is trained again by re-running the entire nested CV procedure, yielding performance metrics under random conditions. This process is repeated many times (n = 2,000 in our case) to construct the empirical null distribution of performance metrics. We then compute the rank of the true performance metrics (obtained from the prediction without shuffling the labels) by calculating how many instances of the null distribution performed better. Dividing the rank by the number of permutations yields the p-value which we report.

A separate question concerned the choice of hyperparameters for the test set. While there are multiple options how hyperparameters for prediction on the test set could be chosen, we decided to use all data from the training set for optimising hyperparameters: we ran a nonnested 5-fold CV on the entire training set, picked the best-performing hyperparameters, and used those to predict on the test set. Other aspects relevant for classification on the test set, such as permutation testing, and regression of confounds were identical to the training set. Please see Fig. 2 for a summary of the Materials and Methods described above.

Finally, we ran an interpretability analysis on our best-performing feature set/classifier combination (rDCM estimates based on 55 ICs and an SVM with a sigmoid kernel). This analysis is based on SHAP (SHapley Additive exPlanations) (Lundberg and Lee, 2017), a generalization of Shapley values from game theory (Shapley, 1953). For each feature, SHAP assigns an importance or attribution value that describes how much that feature contributes to the overall prediction. We used the *shap* software (https://github.com/slundberg/shap) to create a *Kernel-Explainer* that took as arguments:

- a sigmoid SVM classifier trained on rDCM with 55 ICs.
- a low-dimensional representation of the training data using shap.kmeans with five clusters (for computational tractability; see shap documentation).

Then, we computed the SHAP values using the explainer's *shap_values* function which takes the test data as an argument. This gives us a SHAP value for each feature for each subject, which we process (mean of SHAP value magnitude across subjects) to obtain the average impact of each feature on model output magnitude.

2.6. Choice and implementation of classifiers

A total of 17 classifiers were evaluated (please see Table 2 in the Results section), including six support vector machine (SVM) variants and three neural network (NN) variants. As described in the following, for most classifiers, we chose hyperparameters to optimize within the inner nested CV loop. For two classifiers (Gaussian naive Bayes and quadratic discriminant analysis) where hyperparameter tuning is less common, we kept scikit-learn's default parameters.

A first classifier was *logistic regression*. Following the default parameters of the scikit-learn version, we also used *L2* regularization, used *lbfgs* as our solver, and iterated at maximum 100 times. In the inner CV loop, we optimized for the regularization parameter C (0.01, 0.1, 1, 10, 100), which is the inverse of regularization strength (smaller values enforce stronger regularization).

For the *SVMs*, we made use of Platt scaling (Platt, 1999) to get probabilistic outputs for use in an AUROC (area under the receiver-operating characteristic curve) metric. We attempted classification with all types of kernels that scikit-learn has to offer, which include linear, radial basis function (RBF), sigmoid, and polynomial kernels with order 3, 4, and 5. We again treated the regularization parameter *C* (0.01, 0.1, 1) as a hyperparameter and additionally tuned gamma (1, 0.1, 0.01, 0.001) —the kernel coefficient for the RBF, sigmoid, and polynomial kernels.

We included three *neural network* variants (1, 2, and 3 hidden layers) which treat their layer sizes as hyperparameters. All other parameters are scikit-learn defaults (version 0.23.2), which means that – unlike logistic regression – the activation function used is actually ReLU (Rectified Linear Unit; Nair and Hinton 2010).

Neural network hyperparameters

- 1 hidden layer sizes: 100, 150, 300, 500.
- 2 hidden layers sizes: (100, 50), (150, 20), (300, 100), (500, 250).
- 3 hidden layers sizes: (100, 50, 5), (150, 20, 10), (500, 250, 50).

Ensemble methods combine multiple base models to (hopefully) produce better results than each individual model would have on its own. One such algorithm we used is *AdaBoost* (Freund and Schapire, 1997). We employed the scikit-learn default base classifier (decision tree) treating the number of estimators (30, 50, 70) as a hyperparameter. For a baseline comparison, we also attempted classification with a *single decision tree* with default scikit-learn parameters. Another ensemble method used in our classification is *gradient boosting* (Friedman, 2001) with 50, 100, 150 estimators as hyperparameters. Finally, we also tried *random forest* (Breiman, 2001), with options to tune 100, 500, 1000 estimators. For random forest, we additionally treated the maximum tree depth (10, 30, 60) as a hyperparameter.

We also explored prediction with three supervised learning algorithms that do not fall under the previous categories. Namely, *Gaussian naive Bayes* (Zhang, 2004), *quadratic discriminant analysis* (Cover, 1965) – both of which use scikit-learn default parameters – and *k*-nearest neighbors (Cover Hart, 1967) where we treated the number of neighbors (3, 5, 7, 9) and leaf size (20, 30, 40) as hyperparameters.

We used a variety of metrics to evaluate classifier performance in order to ensure a holistic view and to avoid potential pitfalls (such as overemphasizing the importance of one metric). We report recall (sensitivity), precision (positive predictive value), F_1 score, accuracy, and AUROC (area under the receiver-operating characteristic curve).

2.7. Deviations from the original analysis plan

Our analyses were pre-specified and are described in a time-stamped analysis plan (https://gitlab.ethz.ch/tnu/analysis-plans/galioullineetal_ukbb_pred_depr). We subsequently extended this analysis plan in four ways:

- 1. We extended the coverage of networks and, in addition to the 6 networks (represented by IC timeseries), also considered sets of networks consisting of 21 and 55 ICs, as provided by UK Biobank.
- 2. We extended the connectivity methods by considering functional connectivity (Pearson correlation coefficients) in addition to variants of DCM as generative models.
- 3. In addition to SVMs, we decided to test a larger set of classifiers in order to avoid that our results may depend on the particular choice of classifier.



Fig. 2. Summary of the Materials and Methods, illustrating the process taken to create a predictive model of future depressive episode occurrence.

Table 2

Summary table of AUROC for each feature/classifier combination as determined by nested cross-validation on the training set. Bold represents the best result across classifiers for a given feature set, orange shading represents the best result across feature sets for a given classifier, and a star denotes a statistically significant result ($p \le 0.05$). Some values are below 50%, but none of these is significantly different from chance. Ada: AdaBoost, DTC: Decision Tree Classifier, GBC: Gradient Boosting Classifier, GNB: Gaussian Naïve Bayes, kNN: k-Nearest Neighbors, SVM (lin): Support Vector Machine with linear kernel, LR: Logistic Regression, NN (n): Neural Network with n layers, SVM (n): Support Vector Machine with polynomial kernel order n, QDA: Quadratic Discriminant Analysis, SVM (rbf): Support Vector Machine with radial basis function kernel, RF: Random Forest, SVM (sig): Support Vector Machine with sigmoid kernel.

		Features							
		FC(6)	FC(21)	FC(55)	St. DCM	Sp. DCM	rDCM (6)	rDCM (21)	rDCM (55)
	Ada	0.52	0.49	0.56	0.49	0.47	0.55	0.54	0.56
	DTC	0.49	0.49	0.52	0.50	0.51	0.51	0.53	0.52
	GBC	0.51	0.50	0.57	0.49	0.45	0.60^{*}	0.60^{*}	0.64^{*}
	GNB	0.53	0.54	0.54	0.50	0.49	0.61*	0.61^{*}	0.63*
	kNN	0.52	0.54	0.51	0.49	0.47	0.53	0.58	0.59
	\mathbf{LR}	0.54	0.52	0.55	0.51	0.51	0.59	0.58	0.58
ers	NN (1)	0.47	0.52	0.55	0.47	0.55	0.50	0.60^{*}	0.61*
sifi	NN (2)	0.47	0.52	0.54	0.48	0.53	0.50	0.58	0.63*
Clas	NN (3)	0.48	0.52	0.55	0.50	0.54	0.53	0.59	0.62*
	QDA	0.47	0.52	0.50	0.52	0.53	0.52	0.48	0.51
	\mathbf{RF}	0.51	0.52	0.56	0.49	0.47	0.57	0.63^{*}	0.66*
	SVM (lin)	0.54	0.50	0.55	0.51	0.48	0.58	0.58	0.56
	SVM(3)	0.47	0.49	0.52	0.52	0.55	0.59	0.61^{*}	0.65*
	SVM(4)	0.49	0.50	0.47	0.50	0.47	0.49	0.46	0.47
	SVM(5)	0.49	0.47	0.47	0.49	0.55	0.56	0.47	0.46
	SVM (rbf)	0.51	0.51	0.55	0.50	0.50	0.60^{*}	0.63^{*}	0.65*
	SVM (sig)	0.52	0.52	0.47	0.48	0.48	0.61*	0.64*	0.66*

4. We included an analysis of feature importance on the test set using SHAP values.

3. Results

The decision to extend the analyses in this manner took place before any prediction analyses of the training or test data were conducted. We first present the performance of the cross-validated classifiers for the training dataset and then proceed with the most promising feature/classifier combinations to the test dataset. Note: since our dataset is balanced, accuracy as a metric implies balanced accuracy.

Table 3

Summary of all metrics on the top five feature set/classifier combination as determined by nested cross-validation on the training set. Bold indicates the best model for the given metric.

Features + Model	Precision	Recall	F ₁ Score	Accuracy	AUC
rDCM(55) + SVM (sig)	0.64	0.60	0.62	0.63	0.66
rDCM(55) + RF	0.63	0.58	0.60	0.62	0.66
rDCM(55) + SVM (3)	0.63	0.59	0.61	0.62	0.65
rDCM(55) + SVM (rbf)	0.64	0.58	0.61	0.62	0.65
rDCM(21) + SVM (sig)	0.62	0.57	0.59	0.61	0.64

Table 4

Summary of all metrics on the top five feature/classifier combinations on the test set.

Features + Model	Precision	Recall	F ₁ Score	Accuracy	AUC
rDCM(55) + SVM (sig) FC(21) + GNB	0.63 0.59	0.56 0.60	0.60 0.59	0.62 0.59	0.64 0.58
FC(6) + SVM (sig)	0.59	0.59	0.59	0.59	0.61
FC(6) + Log Res	0.57	0.51	0.54	0.57	0.59

3.1. Training set

Table 2 provides an overview of prediction performance on the training set in the nested cross-validation setting. Altogether, 17 different classifiers were evaluated (including four SVM variants and three neural network variants). We report AUROC of all the features run with each classifier (Table 2) and we also report all metrics for the five best feature/classifier combinations (Table 3).

Having run all feature set/classifier pairs on the training data using nested cross-validation, we found that a sigmoid SVM paired with an rDCM taking 55 ICs – referred to subsequently as rDCM(55) – as input performed best (Tables 2 and 3). In terms of performance, applying a sigmoid SVM to rDCM connectivity estimates based on 55 ICs resulted in an AUROC of 0.66 (Table 3). The other performance metrics for this combination were: precision=0.64, recall=0.60, F1 score=0.62, accuracy=0.63 (Table 3).

In general, connectivity estimates by rDCM enabled better predictions, regardless of classifier (see Table 2, orange shading): for 15 out of





Fig. 3. ROC curve for sigmoid SVM paired with rDCM(55) on the training set run with nested cross validation.

the 17 classifiers tested, one of the rDCM feature sets resulted in the best AUROC (in 11/15 cases, the best feature set was rDCM(55)). Furthermore, SVMs tended to perform better than other classifiers, with SVM variants having highest accuracy for 6 out of 8 connectivity feature sets. In particular, the sigmoid SVM was the best-performing classifier for 3 feature sets, more than any other classifier.

Based on these results, we chose rDCM(55) with sigmoid SVM to move forward to the test set. The test data had not been touched until this point to prevent any leakage of information and ensure a thorough verification of the generalisability of our prediction model.

3.2. Test set

The prediction of the best-performing approach on the training set generalized to the test data: the application of a sigmoid SVM to connectivity estimates by rDCM (55 ICs) from the test set showed an AUROC of 0.64 (Fig. 4A). This prediction performance was significantly above chance: Fig. 4B shows that the achieved accuracy of 62% is well outside the null distribution generated by predictions on randomly permuted labels (p < 0.001).



Fig. 4. (A) ROC curve of rDCM (55 ICs) with sigmoid SVM run on test data. (B) Permutation test (*n* = 2,000) run on test data with accuracy as the metric.



Fig. 5. Top 100 Shapley values for sigmoid SVM paired with rDCM(55) on the test data. The bottom half shows outgoing (efferent) connections from each of the 55 ICs, the top half the incoming (afferent) connections. We used the Circlize package in R to generate these figures (Gu et al., 2014).

To get a better understanding of the generalization performance we conducted a post-hoc analysis on other well-performing feature/classifier pairs (Table 4) from the nested cross-validation and assessed their performance on the test data. We defined "well-performing" as the best classifier in general, but for feature sets where another classifier performed better, we selected the latter instead. From the 13 classifiers tested post-hoc, only three other pairs had above-chance performance on the test set, namely rDCM (21 ICs) with sigmoid SVM (58% accuracy, p = 0.019), functional connectivity (6 ICs) with sigmoid SVM (59% accuracy, p = 0.007), and functional connectivity (21 ICs) with Gaussian Naïve Bayes (59% accuracy, p-value=0.012). All of these performed worse with at least a 3% drop in accuracy, leaving the rDCM (55 ICs) with sigmoid SVM as the best-performing model overall on the test data.



Fig. 6. This figure contains two rearranged plots of the results in Fig. 5. As in Fig. 5, the top 100 Shapley values for sigmoid SVM paired with rDCM(55) on the test data are shown, but now ordered according to the network parcellation by Yeo et al. (2011). In panel A, the connections are coloured according to the ICs they originate from, whereas in panel B, the connections are coloured according to the ICs they target. Please note that there is no precise match between ICs from the UK Biobank and the networks from Yeo et al. (2011), and the correspondence implied by the colouring in this figure is only approximate. For technical details, please see the main text.



Fig. 7. Shapley value distribution for sigmoid SVM paired with rDCM(55) on the test data. Values to the right of the dashed red line are in the top 100.

Finally, we computed the SHAP values (Fig. 5) for our bestperforming model, the sigmoid SVM classifier paired with rDCM(55). This assesses the contribution of each connection to the prediction performance. Since rDCM provides estimates of effective (directed) connectivity, we have two SHAP value estimates for each IC, one for the outgoing connection, and one for the incoming connection. We visualized the top 100 SHAP values as a circular plot (Gu et al., 2014), where each IC is shown twice, and the bottom half represents the associated values for the outgoing connections. The width of each displayed connection reflects the magnitude of the SHAP value, and the width of the coloured IC label on the circle represents the cumulative SHAP value for outgoing or incoming connections of that node. Fig. 5 shows that connections with the top 100 SHAP values were not confined to a few networks but included almost all ICs, with very few exceptions. Put simply, during the "resting" state of unconstrained cognition the participants were in, the most predictive connections were found all over the brain. In order to facilitate the interpretation of the results in Fig. 5 for readers familiar with the "resting state" literature, we followed a reviewer's suggestion and re-plotted (Fig. 6) the same results in reference to the network parcellation by Yeo et al. (2011). Given the absence of a direct mapping between the ICs provided by UK Biobank and the networks from the Yeo parcellation, we assigned each IC to that network from the Yeo parcellation with which it showed the greatest spatial overlap (computed as correlation over voxels). Those ICs where the correlation was less than 0.1 were assigned to an "undefined" network (grey ICs in Fig. 6).

Furthermore, we examined the entire distribution of SHAP values, which is shown as a histogram in Fig. 7. This demonstrates that all connections contribute to the model's prediction, albeit most of them to a small degree. The distribution shows considerable spread and a long tail, where the contribution of the most important connection (from IC 34 to IC 24; compare Fig. 5) is two orders of magnitude larger than connections at the mode of the histogram.

4. Discussion

MDD is a syndrome with heterogenous disease trajectories (Merikangas et al., 1994) and variable treatment responses (Rush et al., 2006). Given the importance for clinical management, predicting future clinical outcomes of individual MDD patients has become an important topic in computational psychiatry. In particular, various fMRI studies have examined the feasibility of predicting treatment response (e.g. Harris et al. 2022; Hopman et al. 2021; Ju et al. 2020; Osuch et al. 2018; Queirazza et al. 2019), relapse (e.g. Berwian et al. 2020; Lawrence et al. 2022), or disease trajectories (e.g. Frässle et al. 2020; Schmaal et al. 2015) in individuals with MDD.

By contrast, there have been hardly any attempts to use fMRI to address another challenge of similar importance: the early detection of individuals who are at risk of experiencing a future episode of depression. Given the high frequency of a prolonged remitting-relapsing disease course after a first episode of MDD (Eaton et al., 2008), identifying at-risk individuals is crucial for enabling the targeted deployment of preventive measures and early interventions. So far, to our knowledge, there has only been a single study that used fMRI for detecting individuals at-risk for future depression (Hirshfeld-Becker et al., 2019). This previous study used rs-fMRI and functional connectivity measures in a small sample of individuals with familial risk for MDD (N = 33) for prediction.

The study presented in this paper is novel in several ways. It is the first study using generative models of fMRI data as a basis for predicting future depressive episodes, using three different variants of DCM, in comparison to simpler functional connectivity measures. It uses a large balanced sample size (N = 906), carefully matches groups with presence and absence of depressive symptoms, examines the combination of 8 connectivity feature sets with 17 classifiers in a training set, and evaluates the generalisability of the best predictions using a held-out test set.

The results from the training set (Table 2) indicated that the combination of the rDCM(55) feature set (i.e. rDCM-based connectivity estimates between 55 networks or ICs) and a SVM (with a sigmoid kernel) performed best, showing an AUROC of 0.66 and an accuracy of 63%. This result was significantly above chance, as indicated by permutation testing (p = 0.001, Fig. 3). Moreover, across classifiers, rDCM demonstrated higher predictive value than other connectivity methods (see Table 2): for 15 out of the 17 classifiers tested, one of the rDCM feature sets resulted in the best AUROC; in 11/15 cases, the best feature set was rDCM(55). Examining the results along the other dimension of our investigation, i.e. across all connectivity feature sets, SVMs performed better than other classifiers: for 6 out of 8 connectivity feature sets, one of the SVM variants had the highest accuracy. In particular, a SVM with a sigmoid kernel performed best for 3 feature sets, surpassing any other classifier.

Evaluating the best combination (i.e. rDCM(55) + SVM with sigmoid kernel) on the test set confirmed the generalisability of the predictions, resulting in an AUROC of 0.64 and an accuracy of 62%. This was significantly above chance (p < 0.001), as confirmed by permutation testing (Fig. 4B). In a post-hoc analysis, we also evaluated the predictive value of all other connectivity feature sets on the test set; notably, for each feature set, we used the classifier that had performed best on the training set. These analyses showed that three other combinations of connectivity features/classifiers (rDCM(21) + sigmoid SVM, FC(6) + sigmoid SVM, and FC(21) + Gaussian Naïve Bayes) also achieved significant results, although with slightly lower accuracy (58-59%).

In short, our results thus demonstrate that a GE procedure – based on applying rDCM to rs-fMRI timeseries from a large number of ICs (55) – enabled the best predictions about the occurrence of future depressive episodes within a 3-year period. Having said this, the superiority of GE over a simpler prediction procedure based on FC estimates was not large, amounting to 3% higher accuracy and 0.03 higher AUROC compared to the combination of FC(6) + sigmoid SVM. A binomial test indicated that this difference in accuracy was not significant (p = 0.315).

The lack of a decisive advantage of generative embedding in this rs-fMRI study contrasts with previous task-based fMRI studies in which GE based on DCM was clearly superior to predictions based on FC estimates (e.g. Brodersen et al. 2011, 2014; Frässle et al., 2018, 2020). For example, DCM estimates of effective connectivity during a face perception task allowed for substantially more accurate predictions of MDD disease trajectories than FC estimates: balanced accuracies for predicting a chronic course vs. remission were 79% for DCM and 50% for FC, a difference that was highly significant (Frässle et al., 2020).

In order to understand the limited advantage of GE over FC-based prediction in this study, it is useful to first consider the general reasons why one would, in general, expect GE to show superior performance. In brief:

- (i) GE exploits the fact that a generative model partitions data into signal and noise. Using model parameter estimates (as a lowdimensional representation of signal) as features for subsequent ML ensures that only meaningful information underpins training of classifiers. This makes it less likely that predictions are informed by noise and do not generalize. By contrast, measures of functional connectivity, such as correlation coefficients, reflect both signal and noise. As highlighted by Friston (2011), functional connectivity estimates based on correlations are highly susceptible to changes in the signalto-noise ratio of data.
- (ii) A generative model like DCM distinguishes different mechanisms how measured signal in a system of interest is caused, e.g. connections between system nodes or external inputs. This allows predictions to be differentially informed by distinct system mechanisms. By contrast, FC cannot distinguish whether co-varying signal in two brain regions is caused by shared input or by connections between the regions.
- (iii) DCM provides directed connectivity estimates, allowing one to obtain separate weights for reciprocal connections between regions. By contrast, FC can only provide undirected estimates of connection strengths.
- (iv) From a classical test theory perspective, test-retest reliability of connection strength estimates would be considered an important prerequisite for predictive validity. Concerning rs-FC, test-retest reliability has been examined in numerous studies; a recent meta-analysis reported that, on average, individual connection estimates have limited test-retest reliability (Noble et al., 2019). A direct comparison between FC and rDCM-based estimates of connectivity on identical data (rs-fMRI and multiple tasks) demonstrated that rDCM performed more favourably in this regard (see Fig. 3 in Frässle and Stephan (2022)).

Considering these general factors, one possibility why we only found a limited advantage of GE over FC-based predictions in this study relates to (i) above: in the present study, connectivity was estimated from timeseries that resulted from ICA decomposition and subsequent (manual) removal of components that were identified as noise (Alfaro-Almagro et al., 2018). This approach may have diminished the difference between GE and FC-based prediction with regard to denoising. For comparison, in previous comparisons of GE and FC-based predictions (e.g. Brodersen et al. 2011, 2014; Frässle et al., 2018, 2020), timeseries were obtained by computing the first principal component from regional BOLD measurements, which does not involve a specific distinction between signal and noise. Another possible explanation derives from (ii): application of rDCM to rs-fMRI data essentially means that the model "switches off" external inputs (Frässle et al., 2021). This reduces the superiority in representational richness of GE.

In summary, this suggests that, in the current setting of IC-based rsfMRI timeseries, only factors (iii) and (iv) – but not factors (i) and (ii) – could potentially contribute to higher performance of GE. In order to obtain an impression of the potential impact of factor (iii) – the ability of DCM to obtain separate weights for reciprocal connections between network nodes – we visually explored the asymmetries of node-level SHAP values for incoming versus outgoing connections. For each of the 55 network nodes (ICs), Fig. 8 plots SHAP values summed across all incoming (afferent) and outgoing (efferent) connections, respectively. Visually, it is apparent that for many of the network nodes, the explanatory contributions of incoming versus outgoing connections differ considerably (up to 59%). A more fine-grained plot of connection-specific SHAP values is provided by Fig. 9.

These plots also illustrate a disadvantage of the analysis approach we have chosen in the current study. Specifically, using ICs as network nodes diminish the advantage GE usually enjoys in terms of rendering predictions neurophysiologically interpretable. For example, as shown by Figs. 5 and 8, the connection with the largest SHAP value is the connection from IC 34 to IC 24. Both of these components include a set of fronto-parietal areas: IC 34 includes bilateral frontal regions that appear to match the location of the frontal eye fields as well as more anterior parts of the superior parietal cortex. By contrast, IC 24 contains more posterior bilateral parietal areas, including large parts of bilateral intraparietal sulcus, as well as parts of right middle frontal gyrus and right middle/inferior temporal gyrus. Given this complex anatomical configuration, the biological interpretation of a (directed) functional coupling between IC 34 and IC 24 is not as straightforward as a functional coupling between specific frontal and/or parietal areas. While the FC between these components does not enable any easier interpretations, this example illustrates that the usual interpretive advantage of GE tends to be lost when using IC components as nodes of networks.

Four further aspects of the results deserve discussion. First, it may initially seem surprising that SVM turned out to be the most successful classifier in our comparison, surpassing potentially more powerful methods like neural networks. However, this result is compatible with several recent reports that, for neuroimaging data, kernel-based methods like SVMs (and, in some cases, even simpler linear models) perform equivalently to neural networks for sample sizes up to 10,000 (Cole et al., 2017; He et al., 2020; Schulz et al., 2020).

Second, in our post-hoc analysis of connectivity features/classifier combinations on the test set, three of four significant predictions used the same classifier, an SVM with a sigmoid kernel. Strikingly, FC achieved a significant 59% predictive accuracy using only 6 ICs, whereas the more accurate prediction by rDCM (62%) used 55 ICs, respectively. The resulting difference in the number of features is substantial (15 for FC versus 3025 for rDCM), and it is not immediately clear why predictions based on functional vs. effective connectivity differed greatly in the preferred dimensionality of the feature set. One speculative explanation – which would be consistent with the findings in Figs. 7 and 8 – is that differences in the strengths of reciprocal between-network connections provide subtle but meaningful information that is distributed over

many connections (compare factor (iii) above). This type of information would only be reflected by rDCM, but not by FC-based, connectivity estimates. More generally, it is not clear why FC(6) performed so well on the test data at all. The nested CV on the training data did not indicate that this feature set might be particularly predictive (maximum accuracy of FC-based predictions with any classifier was 54%, none of them significant). The finding of a higher accuracy (59%) on the test set in our post-hoc analysis was surprising. It might be a chance result due to the variance inherent in CV procedures (Varoquaux, 2018) but otherwise lacks a compelling explanation.

Third, contrary to our expectations, predictions based on stochastic and spectral DCM did not generalise to the test set. One possible reason for the lack of successful generalisation is that the higher complexity of the model formulation (e.g. the flexible hemodynamic component and the more sophisticated noise model) could make parameter estimation less reliable, e.g. due to greater abundance of local extrema in the objective function, which would be expected to harm generalisability. This possibility is supported by a recent investigation of parameter recovery of spectral DCM and rDCM which found more accurate parameter recovery for the latter (Frässle et al., 2021). Perhaps even more importantly, however, we could only run spectral and stochastic DCM for 6 ICs on our cluster; for larger feature sets, their compute time (within the context of our entire analysis pipeline) became prohibitively long. However, considering the success of rDCM based on 55 ICs, it is plausible that spectral and stochastic DCM may have performed better if we had been able to run them with larger IC sets (21 and 55).

Fourth, one of the reviewers asked whether the classifier may have exploited health-related group differences at the time when the fMRI scan was obtained (T1; see Fig. 1). This question is of particular interest with regard to anxiety and cardiovascular disorders, given that these clinical conditions are associated with greater risk for depression. We investigated this possibility by statistically comparing the number of D+ and D- individuals showing indications of anxiety and cardiovascular disorders, respectively, at T1. Concerning anxiety, since we did not find detailed questionnaire data on anxiety at T1 in the UK Biobank, we used the UKB data field 1980 as a proxy. This data field contains the responses to a question ("Are you a worrier?") which serves to identify anxious feelings in participants and was asked at T1. With regard to cardiovascular disorders, we extracted information from the UKB data field 20002, counting the number of participants with angina, heart attack/myocardial infarction, and heart failure/pulmonary odema. Using χ^2 tests, we found a significant difference between D+ and D- groups for anxiety (214 "worriers" in D+ and 174 in D-; p = 0.0088) but not for cardiovascular disorders (42 patients in D+ and 32 in D-; p = 0.27). This finding suggests that future extensions of our current approach should take into account anxiety as a potential confound, ideally obtaining more precise measures of (trait) anxiety than was possible for us in this study.

How does the prediction performance achieved in this study compare to previous results in the literature? The only previous fMRI study on predicting future depression (Hirshfeld-Becker et al., 2019) used FC estimates based on rs-fMRI data from six regions, achieving 92% accuracy. However, this study recruited never-depressed children with familial risk for MDD, as opposed to never-depressed participants from the general population as in our study. Additionally, given the more specific focus of the previous study, only 33 participants were available for classification (25 at-risk children, eight controls); this small sample size did not allow for verification in a held-out dataset. Another useful (although not fMRI-based) comparison study utilized structural MRI together with clinical data, questionnaires, and environmental variables (Toenders et al., 2022). The study used a large training set (N = 407 adolescents) and an independent test set (N = 137), achieving an AUROC between 0.68 and 0.72.

It is also instructive to consider the results from non-imaging studies that used demographic, socioeconomic, and clinical variables for predicting the future onset of depression. When considering those studies



Fig. 8. SHAP values summed across the incoming (afferent) and outgoing (efferent) connections of each IC. Blue numbers indicate the % difference in SHAP values for afferent and efferent connections. The plot concerns predictions based on rDCM(55) estimates and SVM with a sigmoid kernel.

that had large sample sizes (i.e. N>500) and tested for generalisability in an independent test set, the reported AUROC values in the literature range between 0.71 and 0.87 (Caldirola et al., 2022; King et al., 2008; Librenza-Garcia et al., 2021; Na et al., 2020; Xu et al., 2019). It is noteworthy, however, that these studies mostly used imbalanced datasets where the number of negative cases (no future depressive episode) far outnumber the positive cases. For example, in the two studies with the highest prediction performance – i.e., AUROC of 0.87 (Na et al., 2020) and 0.85 (Caldirola et al., 2022) – individuals with future depressive episodes amounted to approx. only 8% and 7% of the respective samples. Even when techniques such as oversampling are used (as in Na et al. 2020; but not always the case in other studies), such imbalance can lead to overly optimistic estimates of prediction performance.

Our study has strengths and limitations. Its strengths include an ex ante analysis plan (https://gitlab.ethz.ch/tnu/analysisplans/galioullineetal_ukbb_pred_depr) and a large (N > 900) balanced sample in which groups were carefully matched for 7 potentially confounding variables (age, sex, handedness, tobacco smoking frequency, alcohol consumption frequency, ongoing addictions to illicit drugs, and historical cannabis consumption). This degree of matching is unusually comprehensive (for comparison, in clinical trials and observational studies, it is rarely possible to match for more than two variables) and only made possible by the large resource of the UK Biobank. Furthermore, we conducted a comprehensive comparison of 8 different connectivity measures and 17 classifiers, ensuring that training and test data were strictly separated throughout all analyses. This strict separation of training and test data, with no leak of information, provided robust protection against overfitting and ensured that the classification accuracies found in the test set were not inflated.

Concerning weaknesses, our study has a retrospective design which allows for less robust conclusions than from a prospective study. Furthermore, one potential weakness of the variants of DCM used in this study is that they all rely on variational Bayesian techniques, rendering model inversion susceptible to local extrema in the objective function (Daunizeau et al., 2011). In theory, this could have been addressed by a multi-start procedure, as in previous work with DCM (Schöbi et al., 2021; van Wijk et al., 2018). In practice, however, we were unable to im-

plement this approach given that it would have led to an explosion of the already very substantial compute time. Finally, the greatest limitation of our study is the definition of depressive episodes. Given the heterogeneity of clinical data in the UK Biobank and the lack of systematic information about absence/presence of a clinical diagnosis of depression, we combined multiple sources of information within UK Biobank - i.e., clinical records, questionnaires (PHQ, MHQ) and self-report specifically on issues of depression - to identify indicators of at least one depressive episode within three years after the fMRI scan. Clearly, this partial reliance on self-report is not ideal; additionally, the resulting group of participants with a putative depressive episode (D+ group) is likely heterogeneous and might include people with very different severities of depression. The definition of a depressive episode in the D+ group ranged from a full clinical diagnosis to answering "yes" to the question whether they ever had a time when they were feeling depressed or down for at least a whole week. While the latter indicator is broad, relies on self-report and likely leads to inclusion of mild cases, it is noteworthy that a depression phenotype based on the above question plus additional conditions (Smith et al., 2013) was found to be useful in a genome wide association study (GWAS), with results that closely correlated to those of an ICD-based MDD phenotype (Howard et al., 2018).

Following a reviewer's suggestion, we explored the impact of heterogeneity in the D+ group by examining the classification accuracy of D+ individuals based on the indicator category of their depressive episode. In brief, in the test set, we found an accuracy of 92% for individuals who had been assigned to the D+ group due to their responses in the online Mental Health Questionnaire fields focused on depression diagnosis. Conversely, the classification accuracy for the category "ever had a time when they were feeling depressed or down for at least a whole week" was only 36%. This suggests that MHQ-related fields are a stronger indicator of depressive episodes compared to self-reported "depressed for a week" and hence enable the classifier to distinguish D+ and D- individuals more easily.

Furthermore, there is rarely information on when exactly within the 3-year period a depressive episode occurred; the likely interindividual variability in the latency of symptom onset after the fMRI scan would further add to the heterogeneity of the D+ group. Having said this, our approach is similar to previous analyses of depression in the UK



Heatmap of Adjacency Matrix SHAP Values for rDCM(55) + Sigmoid SVM

Fig. 9. Matrix of connections between all ICs, showing the connections' colour-coded SHAP values (test data) for predictions based on rDCM(55) estimates and sigmoid SVM. ICs are ordered according to summed SHAP values.

Biobank that also relied on self-report and questionnaires like the MHQ (Howard et al., 2020). More generally, a pragmatic approach to identifying individuals with likely clinical characteristics is often unavoidable when working with large heterogeneous databases (for an example using self-reported depression in genetics, see Wray et al., 2018). The challenge how to optimally extract data from the UK Biobank for studies of MDD is being addressed by ongoing methodological developments (Dutt et al., 2021) which will help to improve and standardise future studies. Overall, our results have four implications. First, given the challenging nature of the prediction problem tackled in the study (i.e. occurrence of indicators of depressive episodes, as opposed to full clinical diagnoses, over a three-year period), it is encouraging that significant predictions on held-out data can be obtained at all. Second, despite this success and the potential for further optimisation, our study suggests that fMRI on its own may not be sufficient for clinically useful predictions. Future studies of predicting depression should utilise fMRI-based connectivity estimates in conjunction with additional data (e.g. demographic, socioeconomic, clinical). Third, while GE results based on rDCM were consistently successful across all classifiers and enjoyed a numerical advantage over FC for clinical predictions, performance differences were modest and nonsignificant. The magnitude of performance differences between GE and FC in this study and previous work suggests that adding task-based fMRI may enhance the difference in predictive accuracy. Finally, using IC components as network nodes diminishes the usual advantage of GE with regard to biological interpretability of predictions. In order to maintain the interpretability of GE based predictions, it would seem advantageous to compute effective connectivity between disjoint areas from parcellations based on combined anatomical-functional criteria (e.g. Fan et al. 2016; Glasser et al. 2016). We hope that these conclusions will be useful for future work on predicting the occurrence of depressive episodes.

Data and code availability

Our study primarily relies on the dataset provided by the UK Biobank, which is made available to all qualified researchers via their website https://www.ukbiobank.ac.uk/. We provide our code for dataset extraction and model training at https://gitlab.ethz.ch/ tnu/code/galioullineetal_ukbb_pred_depr. Our computation and analysis was conducted using the computational cluster of the Swiss Federal Institute of Technology (ETHZ).

Declaration of Competing Interest

None.

Credit authorship contribution statement

Herman Galioulline: Investigation, Software, Formal analysis, Writing – original draft, Visualization. Stefan Frässle: Investigation, Supervision, Validation, Software, Methodology, Writing – review & editing. Samuel J. Harrison: Investigation, Supervision, Validation, Software, Methodology, Writing – review & editing. Inês Pereira: Investigation, Validation, Writing – review & editing. Jakob Heinzle: Investigation, Supervision, Validation, Software, Methodology, Writing – review & editing. Klaas Enno Stephan: Investigation, Conceptualization, Supervision, Methodology, Project administration, Funding acquisition, Writing – review & editing.

Acknowledgments

This work was supported by the René and Susanne Braginsky Foundation (KES), the ETH Foundation (KES), and project grant 320030_179377 by the Swiss National Science Foundation (KES). This research has been conducted using the UK Biobank Resource under Application Number 60679.

References

- Adler, D.A., et al., 2006. Job performance deficits due to depression. Am. J. Psychiatry 163 (9), 1569–1576. doi:10.1176/ajp.2006.163.9.1569.
- Alfaro-Almagro, F., et al., 2018. Image processing and quality control for the first 10,000 brain imaging datasets from UK biobank. NeuroImage 166, 400–424. doi:10.1016/j.neuroimage.2017.10.034.
- Barch, D.M., et al., 2019. Early childhood depression, emotion regulation, episodic memory and hippocampal development. J. Abnorm. Psychol. 128 (1), 81–95. doi:10.1037/abn0000392.
- Berwian, I.M., et al., 2020. The relationship between resting-state functional connectivity, antidepressant discontinuation and depression relapse. Sci. Rep. 10 (1), 22346. doi:10.1038/s41598-020-79170-9.
- Brakowski, J., et al., 2017. Resting state brain network function in major depression depression symptomatology, antidepressant treatment effects, future research. J. Psychiatric Res. 92, 147–159. doi:10.1016/j.jpsychires.2017.04.007.
- Breiman, L., 2001. Random forests. Mach. Learn. 45 (1), 5–32. doi:10.1023/A:1010933404324.
- Brodersen, K.H., et al., 2011. Generative embedding for model-based classification of fMRI data. PLOS Comput. Biol. 7 (6), e1002079. doi:10.1371/journal.pcbi.1002079.

- Brodersen, K.H., et al., 2014. Dissecting psychiatric spectrum disorders by generative embedding. NeuroImage Clin. 4, 98–111. doi:10.1016/j.nicl.2013.11.002.
- Caldirola, D., et al., 2022. First-onset major depression during the COVID-19 pandemic: a predictive machine learning model. J. Affect. Disord. 310, 75–86. doi:10.1016/j.jad.2022.04.145.

Cawley, G.C., Talbot, N.L.C., 2010. On over-fitting in model selection and subsequent selection bias in performance evaluation. J. Mach. Learn. Res. 11 (70), 2079–2107.

- Chikersal, P., et al., 2021. Detecting depression and predicting its onset using longitudinal symptoms captured by passive sensing: a machine learning approach with robust feature selection. ACM Trans. Comput. Hum. Interact. 28 (1), 3. doi:10.1145/3422821, 1-3:41.
- Cole, J.H., et al., 2017. Predicting brain age with deep learning from raw imaging data results in a reliable and heritable biomarker. NeuroImage 163, 115–124. doi:10.1016/j.neuroimage.2017.07.059.
- Coleman, J.R.I., et al., 2020. Genome-wide gene-environment analyses of major depressive disorder and reported lifetime traumatic experiences in UK Biobank. Mol. Psychiatry 25 (7), 1430–1446. doi:10.1038/s41380-019-0546-6.
- Correll, C.U., et al., 2017. Prevalence, incidence and mortality from cardiovascular disease in patients with pooled and specific severe mental illness: a large-scale meta-analysis of 3,211,768 patients and 113,383,368 controls. World Psychiatry 16 (2), 163–180. doi:10.1002/wps.20420.
- Cover, 1965. Geometrical and statistical properties of systems of linear inequalities with applications in pattern recognition. IEEE Trans. Electron. Comput. EC-14 (3), 326–334. doi:10.1109/PGEC.1965.264137.
- Cover Hart, P., 1967. Nearest neighbor pattern classification. IEEE Trans. Inf. Theory 13 (1), 21–27. doi:10.1109/TIT.1967.1053964.
- Cuijpers, P., et al., 2021. Psychological interventions to prevent the onset of depressive disorders: a meta-analysis of randomized controlled trials. Clin. Psychol. Rev. 83, 101955. doi:10.1016/j.cpr.2020.101955.
- Cuijpers, P., Beekman, A.T.F., Reynolds, C.F., 2012. Preventing depression: a global priority. JAMA 307 (10), 1033–1034. doi:10.1001/jama.2012.271.
- Daunizeau, J., David, O., Stephan, K.E., 2011. Dynamic causal modelling: A critical review of the biophysical and statistical foundations. NeuroImage 58 (2), 312–322. doi:10.1016/j.neuroimage.2009.11.062.
- Dutt, R.K., et al. 2021. Mental health in the UK Biobank: A roadmap to self-report measures and neuroimaging correlates. hbm, doi:10.1002/hbm.25690.
- Eaton, W.W., et al., 2008. Population-based study of first onset and chronicity in major depressive disorder. Arch. Gen. Psychiatry 65 (5), 513–520. doi:10.1001/archpsyc.65.5.513.
- van Eeden, W.A., et al., 2021. Predicting the 9-year course of mood and anxiety disorders with automated machine learning: a comparison between auto-sklearn, naïve Bayes classifier, and traditional logistic regression. Psychiatry Res. 299, 113823. doi:10.1016/j.psychres.2021.113823.
- Fan, L., et al., 2016. The human brainnetome atlas: a new brain atlas based on connectional architecture. Cereb. Cortex 26 (8), 3508–3526. doi:10.1093/cercor/bhw157.
- Frässle, S., et al., 2017. Regression DCM for fMRI. NeuroImage 155, 406–421. doi:10.1016/j.neuroimage.2017.02.090.
- Frässle, S., et al., 2020. Predicting individual clinical trajectories of depression with generative embedding. NeuroImage Clin. 26, 102213. doi:10.1016/j.nicl.2020.102213.
- Frässle, S., Aponte, E.A., Bollmann, S., Brodersen, K.H., Do, C.T., Harrison, O.K., Harrison, S.J., Heinzle, J., Iglesias, S., Kasper, L., Lomakina, E.I., Mathys, C., Müller-Schrader, M., Pereira, I., Petzschner, F.H., Raman, S., Schöbi, D., Toussaint, B., Weber, L.A., Yao, Y., Stephan, K.E., 2021. TAPAS: An Open-Source Software Package for Translational Neuromodeling and Computational Psychiatry. Front Psychiatry 12, 680811. https://doi.org/10.3389/fpsyt.2021.680811.
- Frässle, S., Harrison, S.J., Heinzle, J., Clementz, B.A., Tamminga, C.A., Sweeney, J.A., Gershon, E.S., Keshavan, M.S., Pearlson, G.D., Powers, A., Stephan, K.E., 2021. Regression dynamic causal modeling for resting-state fMRI. Human Brain Mapping 42, 2159–2180. https://doi.org/10.1002/hbm.2535.
- Frässle, S., Lomakina, E.I., Kasper, L., Manjaly, Z.M., Leff, A., Pruessmann, K.P., Buhmann, J.M., Stephan, K.E., 2018. A generative model of whole-brain effective connectivity. Neuroimage 179, 505–529. https://doi.org/10.1016/j.neuroimage.2018.05.058.
- Frässle, S., Stephan, K.E., 2022. Test-retest reliability of regression dynamic causal modeling. Netw. Neurosci. 6 (1), 135–160. doi:10.1162/netn_a_00215.
- Frässle, S., Yao, Y., Schöbi, D., Aponte, E.A., Heinzle, J., Stephan, K.E., 2018. Generative models for clinical applications in computational psychiatry. WIREs Cognitive Science 9, e1460. https://doi.org/10.1002/wcs.1460.
- Freund, Y., Schapire, R.E., 1997. A decision-theoretic generalization of on-line learning and an application to boosting. J. Comput. Syst. Sci. 55 (1), 119–139. doi:10.1006/jcss.1997.1504.
- Friedman, J.H., 2001. Greedy function approximation: a gradient boosting machine. Ann. Stat. 29 (5), 1189–1232. doi:10.1214/aos/1013203451.
- Friston, K.J., 2011. Functional and Effective Connectivity: A Review. Mary Ann Liebert, Inc., New Rochelle, NY doi:10.1089/brain.2011.0008.
- Friston, K.J., et al., 2014. A DCM for resting state fMRI. NeuroImage 94, 396–407. doi:10.1016/j.neuroimage.2013.12.009.
- Friston, K.J., Harrison, L., Penny, W., 2003. Dynamic causal modelling. NeuroImage 19 (4), 1273–1302. doi:10.1016/S1053-8119(03)00202-7.
- GBD Mental Disorders Collaborators, 2022. Global, regional, and national burden of 12 mental disorders in 204 countries and territories, 1990–2019: a systematic analysis for the global burden of disease study 2019. Lancet Psychiatry 9 (2), 137–150. doi:10.1016/S2215-0366(21)00395-3.
- Glasser, M.F., et al., 2016. A multi-modal parcellation of human cerebral cortex. Nature 536 (7615), 171–178. doi:10.1038/nature18933.

Gordon, J.A., 2016. On being a circuit psychiatrist. Nat. Neurosci. 19 (11), 1385–1386. doi:10.1038/nn.4419.

- Goulden, N. et al. (2014) 'The salience network is responsible for switching between the default mode network and the central executive network: Replication from DCM', *NeuroImage*, 99, pp. 180–190, https://doi.org/10.1016/j.neuroimage.2014.05.052.
- Gratton, C., Sun, H., Petersen, S.E., 2018. Control networks and hubs. Psychophysiology 55 (3). doi:10.1111/psyp.13032.
- Griffanti, L., et al., 2014. ICA-based artefact removal and accelerated fMRI acquisition for improved resting state network imaging. NeuroImage 95, 232–247. doi:10.1016/j.neuroimage.2014.03.034.
- Gu, S.C., et al., 2020. Personalized prediction of depression in patients with newly diagnosed Parkinson's disease: a prospective cohort study. J. Affect. Disord. 268, 118–126. doi:10.1016/j.jad.2020.02.046.
- Gu, Z., et al., 2014. circlize implements and enhances circular visualization in R. Bioinformatics 30 (19), 2811–2812. doi:10.1093/bioinformatics/btu393.
- Harris, J.K., et al., 2022. Predicting escitalopram treatment response from pre-treatment and early response resting state fMRI in a multi-site sample: a CAN-BIND-1 report. NeuroImage Clin. 35, 103120. doi:10.1016/j.nicl.2022.103120.
- He, T., et al., 2020. Deep neural networks and kernel regression achieve comparable accuracies for functional connectivity prediction of behavior and demographics. NeuroImage 206, 116276. doi:10.1016/j.neuroimage.2019.116276.
- Hirshfeld-Becker, D.R., et al., 2019. Intrinsic functional brain connectivity predicts onset of major depression disorder in adolescence: a pilot study. Brain Connectivity 9 (5), 388–398. doi:10.1089/brain.2018.0646.
- Hopman, H.J., et al., 2021. Personalized prediction of transcranial magnetic stimulation clinical response in patients with treatment-refractory depression using neuroimaging biomarkers and machine learning. J. Affect. Disord. 290, 261–271. doi:10.1016/j.jad.2021.04.081.
- Howard, D.M., et al., 2018. Genome-wide association study of depression phenotypes in UK Biobank identifies variants in excitatory synaptic pathways. Nat. Commun. 9 (1), 1470. doi:10.1038/s41467-018-03819-3.
- Howard, D.M., et al., 2020. Genetic stratification of depression in UK Biobank. Translational Psychiatry 10, 163. doi:10.1038/s41398-020-0848-0.
- Hyatt, C.J. et al. (2015) 'Specific default mode subnetworks support mentalizing as revealed through opposing network recruitment by social and semantic FMRI tasks', *Human Brain Mapping*, 36(8), pp. 3047–3063. https://doi.org/10.1002/hbm.22827. Jenkinson, M., et al., 2012. FSL. NeuroImage 62 (2), 782–790.
- Jenkinson, M., et al., 2012. FSL. NeuroImage 62 (2), 782–790. doi:10.1016/j.neuroimage.2011.09.015.
- Ju, Y., et al., 2020. Connectome-based models can predict early symptom improvement in major depressive disorder. J. Affect. Disord. 273, 442–452. doi:10.1016/j.jad.2020.04.028.
- Kaiser, R.H., et al., 2015. Large-scale network dysfunction in major depressive disorder: a meta-analysis of resting-state functional connectivity. JAMA Psychiatry 72 (6), 603– 611. doi:10.1001/jamapsychiatry.2015.0071.
- King, M., et al., 2008. Development and validation of an international risk prediction algorithm for episodes of major depression in general practice attendees: the PredictD study. Arch. Gen. Psychiatry 65 (12), 1368–1376. doi:10.1001/archpsyc.65.12.1368.
- Kupferberg, A., Bicks, L., Hasler, G., 2016. Social functioning in major depressive disorder. Neurosci. Biobehav. Rev. 69, 313–332. doi:10.1016/j.neubiorev.2016.07.002.
- Lawrence, A.J., et al., 2022. Neurocognitive measures of self-blame and risk prediction models of recurrence in major depressive disorder. Biol. Psychiatry Cogn. Neurosci. Neuroimaging 7 (3), 256–264. doi:10.1016/j.bpsc.2021.06.010.
- Li, B., et al., 2011. Generalised filtering and stochastic DCM for fMRI. NeuroImage 58 (2), 442–457. doi:10.1016/j.neuroimage.2011.01.085.
- Librenza-Garcia, D., et al., 2021. Prediction of depression cases, incidence, and chronicity in a large occupational cohort using machine learning techniques: an analysis of the ELSA-Brasil study. Psychol. Med. 51 (16), 2895–2903. doi:10.1017/S0033291720001579.
- Lin, S., et al., 2022. Prediction of depressive symptoms onset and long-term trajectories in home-based older adults using machine learning techniques. Aging Ment. Health 0 (0), 1–10. doi:10.1080/13607863.2022.2031868.
- Lundberg, S.M., Lee, S.I., Guyon, I., et al., 2017. A unified approach to interpreting model predictions. Advances in Neural Information Processing Systems. Curran Associates, Inc.
- Merikangas, K.R., Wicki, W. and Angst, J. (1994) 'Heterogeneity of Depression: Classification of Depressive Subtypes by Longitudinal Course', *The British Journal of Psychiatry*, 164(3), pp. 342–348. https://doi.org/10.1192/bjp.164.3.342.
- Miller, K.L., et al., 2016. Multimodal population brain imaging in the UK Biobank prospective epidemiological study. Nat. Neurosci. 19 (11), 1523–1536. doi:10.1038/nn.4393.
- Motlaghian, S.M. et al. (2022) 'Nonlinear functional network connectivity in resting functional magnetic resonance imaging data', *Human Brain Mapping*, 43(15), pp. 4556– 4566. https://doi.org/10.1002/hbm.25972.
- Na, K.S., et al., 2020. Predicting future onset of depression among community dwelling adults in the Republic of Korea using a machine learning algorithm. Neurosci. Lett. 721, 134804. doi:10.1016/j.neulet.2020.134804.
- Nair, V., Hinton, G.E., 2010. Rectified linear units improve restricted boltzmann machines. In: Proceedings of the 27th International Conference on International Conference on Machine Learning, Madison, WI, USA. Omnipress (ICML'10), pp. 807–814.
- Machine Learning, Madison, WI, USA. Omnipress (ICML'10), pp. 807–814.
 Nickerson, L.D., et al., 2017. Using dual regression to investigate network shape and amplitude in functional connectivity analyses. Front. Neurosci. 11. https://www.frontiersin.org/articles/10.3389/fnins.2017.00115. Accessed: 25 September 2022.
- Noble, S., Scheinost, D., Constable, R.T., 2019. A decade of test-retest reliability of functional connectivity: a systematic review and meta-analysis. NeuroImage 203, 116157. doi:10.1016/j.neuroimage.2019.116157.

Osuch, E., et al., 2018. Complexity in mood disorder diagnosis: fMRI connectivity networks

predicted medication-class of response in complex patients. Acta Psychiatr. Scand. 138 (5), 472–482. doi:10.1111/acps.12945.

- Pagliaccio, D., et al., 2014. Brain-behavior relationships in the experience and regulation of negative emotion in healthy children: implications for risk for childhood depression. Dev. Psychopathol. 26 (4pt2), 1289–1303. doi:10.1017/S0954579414001035.
- Papmeyer, M., et al., 2016. Prospective longitudinal study of subcortical brain volumes in individuals at high familial risk of mood disorders with or without subsequent onset of depression. Psychiatry Res. Neuroimaging 248, 119–125. doi:10.1016/j.pscychresns.2015.12.009.
- Platt, J.C., 1999. Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. In: Advances in Large Margin Classifiers. MIT Press, pp. 61–74.
- Poldrack, R.A., Huckins, G., Varoquaux, G., 2020. Establishment of best practices for evidence for prediction: a review. JAMA Psychiatry 77 (5), 534–540. doi:10.1001/jamapsychiatry.2019.3671.
- Queirazza, F., et al., 2019. Neural correlates of weighted reward prediction error during reinforcement learning classify response to cognitive behavioral therapy in depression. Sci. Adv. 5 (7), eaav4962. doi:10.1126/sciadv.aav4962.
- Rocha, T.B.M., et al., 2021. Identifying adolescents at risk for depression: a prediction score performance in cohorts based in 3 different continents. J. Am. Acad. Child Adolesc. Psychiatry 60 (2), 262–273. doi:10.1016/j.jaac.2019.12.004.
- Rosellini, A.J., et al., 2020. Developing algorithms to predict adult onset internalizing disorders: an ensemble learning approach. J. Psychiatr. Res. 121, 189–196. doi:10.1016/j.jpsychires.2019.12.006.
- Rush, A.J., et al., 2006. Acute and longer-term outcomes in depressed outpatients requiring one or several treatment steps: a STAR*D report. Am. J. Psychiatry 163 (11), 1905–1917. doi:10.1176/ajp.2006.163.11.1905.
- Salimi-Khorshidi, G., et al., 2014. Automatic denoising of functional MRI data: combining independent component analysis and hierarchical fusion of classifiers. NeuroImage 90, 449–468. doi:10.1016/j.neuroimage.2013.11.046.
- Sampson, L., et al., 2021. A machine learning approach to predicting new-onset depression in a military population. Psychiatr. Res. Clin. Pract. 3 (3), 115–122. doi:10.1176/appi.prcp.20200031.
- Schmaal, L., et al., 2015. Predicting the naturalistic course of major depressive disorder using clinical and multimodal neuroimaging information: a multivariate pattern recognition study. Biol. Psychiatry 78 (4), 278–286. doi:10.1016/j.biopsych.2014.11.018.
- Schulz, M.A., et al., 2020. Different scaling of linear models and deep learning in UK-Biobank brain images versus machine-learning datasets. Nat. Commun. 11 (1), 4238. doi:10.1038/s41467-020-18037-z.
- Schöbi, D., Do, C.-T., Frässle, S., Tittgemeyer, M., Heinzle, J., Stephan, K.E., 2021. Technical note: A fast and robust integrator of delay differential equations in DCM for electrophysiological data. NeuroImage 244, 118567. doi:10.1016/j.neuroimage.2021.118567.
- Shapero, B.G., et al., 2019. Neural markers of depression risk predict the onset of depression. Psychiatry Res. Neuroimaging 285, 31–39. doi:10.1016/j.pscychresns.2019.01.006.
- Shapley, L.S. (1953) '17. A value for n-person games', in 17. A Value for n-Person Games. Princeton University Press, pp. 307–318. Available at: doi:10.1515/9781400881970-018.
- Shen, X., et al., 2018. Resting-state connectivity and its association with cognitive performance, educational attainment, and household income in the UK biobank. Biol. Psychiatry Cogn. Neurosci. Neuroimaging 3 (10), 878–886. doi:10.1016/j.bpsc.2018.06.007.
- Smith, D.J., et al., 2013. Prevalence and characteristics of probable major depression and bipolar disorder within UK biobank: cross-sectional study of 172,751 participants. PLoS One 8 (11), e75362. doi:10.1371/journal.pone.0075362.
- Smith, S.M., 2009. Correspondence of the brain's functional architecture during activation and rest. pnas 106 (31), 13040–13045. doi:10.1073/pnas.0905267106.
- Steffen, A., et al., 2020. Mental and somatic comorbidity of depression: a comprehensive cross-sectional analysis of 202 diagnosis groups using German nationwide ambulatory claims data. BMC Psychiatry 20 (1), 142. doi:10.1186/s12888-020-02546-8.
- Stephan, K.E., et al., 2015. Translational perspectives for computational neuroimaging. Neuron 87 (4), 716–732. doi:10.1016/j.neuron.2015.07.008.
- Stephan, K.E., et al., 2017. Computational neuroimaging strategies for single patient predictions. NeuroImage 145, 180–199. doi:10.1016/j.neuroimage.2016.06.038.
- Stone, M., 1974. Cross-validatory choice and assessment of statistical predictions. J. R. Stat. Soc. Ser. B 36 (2), 111–147 (Methodological).
- Toenders, Y.J. et al. (2022) 'Predicting Depression Onset in Young People Based on Clinical, Cognitive, Environmental, and Neurobiological Data', Biological Psychiatry: Cognitive Neuroscience and Neuroimaging, 7(4), pp. 376–384. https://doi.org/10.1016/j.bpsc.2021.03.005.
- Van Dijk, K.R.A., Sabuncu, M.R., Buckner, R.L., 2012. The influence of head motion on intrinsic functional connectivity MRI. NeuroImage 59 (1), 431–438. doi:10.1016/j.neuroimage.2011.07.044.
- van Wijk, B.C.M., Cagnan, H., Litvak, V., Kühn, A.A., Friston, K.J., 2018. Generic dynamic causal modelling: An illustrative application to Parkinson's disease. NeuroImage 181, 818–830. doi:10.1016/j.neuroimage.2018.08.039.
- Varoquaux, G., 2018. Cross-validation failure: Small sample sizes lead to large error bars. NeuroImage 180, 68–77. doi:10.1016/j.neuroimage.2017.06.061.
- Voorhees, B.W.V., et al., 2008. Predicting future risk of depressive episode in adolescents: the chicago adolescent depression risk assessment (CADRA). Ann. Fam. Med. 6 (6), 503–511. doi:10.1370/afm.887.
- Vos, T., et al., 2020. Global burden of 369 diseases and injuries in 204 countries and territories, 1990–2019: a systematic analysis for the global burden of disease study 2019. Lancet 396 (10258), 1204–1222. doi:10.1016/S0140-6736(20)30925-9.
- Wray, N.R., et al., 2018. Genome-wide association analyses identify 44 risk variants and

refine the genetic architecture of major depression. Nature Genetics 50, 668-681.

- reme une geneuc arcnitecture of major depression. Nature Genetics 50, 668–681. doi:10.1038/s41588-018-0090-3.
 Xu, Z., et al., 2019. Individualized prediction of depressive disorder in the el-derly: a multitask deep learning approach. Int. J. Med. Inform. 132, 103973. doi:10.1016/j.ijmedinf.2019.103973.
- Yeo, B.T.T., et al., 2011. The organization of the human cerebral cortex estimated by intrinsic functional connectivity. J. Neurophysiol. 106 (3), 1125–1165. doi:10.1152/jn.00338.2011.
- Zhang, H., 2004. The optimality of naive bayes. In: Proceedings of the 7th International Florida Artificial Intelligence Research Society Conference, FLAIRS 2004.