

Discriminating the Influence of Correlated Factors from Multivariate Observations: the Back-to-Back Regression

Jean-Rémi King^{a,b,*}, François Charton^b, David Lopez-Paz^b, Maxime Oquab^b

^a*Laboratoire des systèmes perceptifs, PSL University, CNRS*

^b*Facebook AI*

Abstract

Identifying causes solely from observations can be particularly challenging when i) potential factors are difficult to manipulate independently and ii) observations are multi-dimensional. To address this issue, we introduce “Back-to-Back” regression (B2B), a linear method designed to efficiently measure, from a set of correlated factors, those that most plausibly account for multidimensional observations. First, we prove the consistency of B2B, its links to other linear approaches, and show how it provides a robust, unbiased and interpretable scalar estimate for each factor. Second, we use a variety of simulated data to show that B2B outperforms least-squares regression and cross-decomposition techniques (e.g. canonical correlation analysis and partial least squares) on causal identification when the factors and the observations are partially collinear. Finally, we apply B2B to magneto-encephalography of 102 subjects recorded during a reading task to test whether our method appropriately disentangles the respective contribution of word length and word frequency - two correlated factors known to cause early and late brain responses respectively. The results show that these two factors are better disentangled with B2B than with other standard techniques.

Keywords: Cross-Decomposition, Feature Discovery, Magnetoencephalography, Decoding, Encoding, Reading, N400

1. Introduction

1 Natural sciences are tasked to find, from a set of hypothetical factors, the minimal subset that
2 suffices to reliably predict novel observations. This endeavor is impeded by two major challenges.
3
4 First, causal and non-causal factors may be numerous and partially correlated. In neuroscience,
5 for example, it can be challenging to identify whether word frequency modulates brain activity
6 during reading. Indeed, the frequency of words in natural language covaries with other factors such
7 as their length (short words are more frequent than long words) and their categories (determinants
8 are more frequent than adverbs) [18, 24]. Instead of selecting a set of words that controls for all
9 of these factors simultaneously, it is thus common to use a *forward* “encoding model”, i.e. to fit
10 a linear regression to predict observations (e.g. brain activity) from a minimal combination of

*corresponding author: jeanremi@fb.com

11 competing factors (e.g. word length, word frequency), and analytically investigate, the estimated
12 contribution of each factor from the model’s coefficients [5, 21, 32, 16, 13].

13 The second challenge to measuring causal influence is that observations can be multidimensional.
14 The relationship between causes and effects is thus often considered in a *backward* manner, by
15 training models to maximally predict causes from multidimensional observations. For example,
16 brain activity is often recorded with hundreds or thousands of simultaneous measurements via
17 functional Magnetic Resonance Imaging, magneto-encephalography (MEG) or multiple electro-
18 physiological probes [5, 30]. As simultaneous measurements may be affected by common noise
19 sources, it is common to use backward modeling, by, for example, fitting a support vector machine
20 across multiple sensors to decode the category of a stimulus [22, 3, 17].

21 Both *forward* and *backward* modeling have competing benefits and drawbacks. Specifically,
22 forward modeling disentangles the independent contribution of correlated factors, but does not
23 combine multidimensional observations. By contrast, backward modeling combines multiple
24 observations, but does not disentangle factors that are linearly correlated [32, 9, 16]. To combine
25 some of the benefits of forward and backward modeling, several authors have proposed to use
26 cross-decomposition techniques such as Partial Least Squares (PLS) and Canonical Correlation
27 Analysis (CCA) [4]. CCA and PLS aim to find, from two sets of data X and Y , the components H
28 and G where XH and YG are maximally correlated or maximally covarying respectively. Because
29 CCA and PLS are based on a generalized eigen decomposition, their resulting coefficients are
30 mixing the features of X and Y in a way that makes them notoriously difficult to interpret [19].

31 Here, we introduce the ‘back-to-back regression’ (B2B), which not only combines the benefits of
32 forward and backward modeling (Section 2), but also provides robust, interpretable, unidimensional
33 and unbiased coefficients for each of tested factor.

34 The present paper focuses on the restricted issue of disentangling the influence of linearly
35 correlated predictors (X) onto noisy multivariate observations (Y). The present approach thus
36 differs from other causal discovery algorithms based on temporal-delays and/or nonlinear interac-
37 tions in systems where the directionality of causation (from X to Y or vice versa) is unknown (e.g.
38 [25, 8, 14, 28]).

39 After detailing B2B method and proving its convergence (Section 2.2), we show with synthetic
40 data that it outperforms state-of-the-art forward, backward and cross-decomposition techniques
41 in disentangling causal factors (Section 3.1). Finally, we apply B2B to a large neuroimaging
42 dataset and reveal that distinct but linearly-correlated word features lead to distinguishable brain
43 representations (Section 3.5).

44 2. Back-to-Back regression

We consider the measurement of multivariate signal $Y \in \mathbb{R}^{n \times d_y}$ (the dependent variables),
generated from a set of putative causes $X \in \mathbb{R}^{n \times d_x}$ (the independent variables), via some unknown
linear apparatus $F \in \mathbb{R}^{d_x \times d_y}$. Not all the variables in X exert a causal influence on Y . By
considering a square binary diagonal matrix of *causal influences* $E \in \mathbb{D}^{d_x \times d_x}$, we denote by XE
the causal factors of Y . In summary, the problem can be formalized as:

$$y_i = (x_i E + n_i) F \quad (1)$$

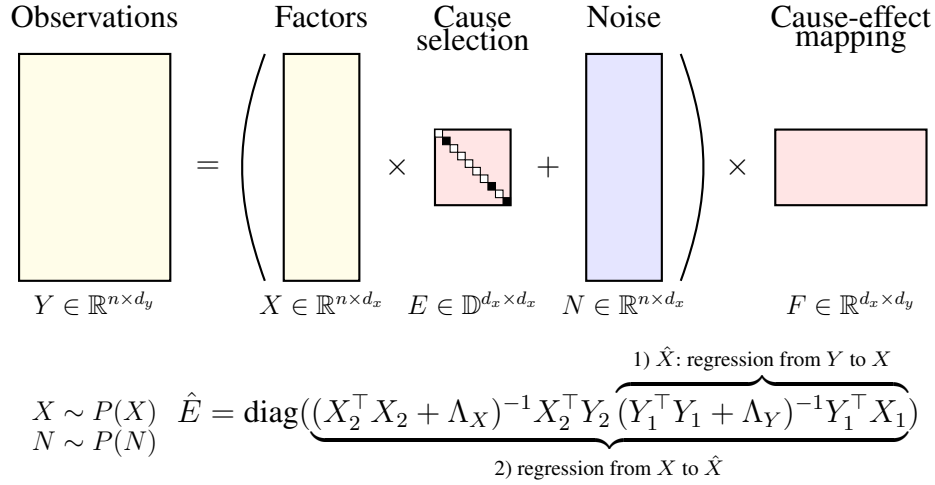


Figure 1: Back-to-back regression identifies the subset of factors $E_{ii} = 1$ in X that influence some observations Y by 1) regressing from Y to X to obtain \hat{X} , and 2) returning the diagonal of the regression coefficients from X to \hat{X} .

45 where i is a given sample, and n_i is a sample-specific noise drawn from a centered distribution.
 46 While the triplet of variables X and N are independent, we allow each of them to have any form of
 47 covariance. In practice, we observe n samples (X, Y) from the model. This problem space, along
 48 with the sizes of all variables involved, is illustrated in Figure 1. Given the model in Equation eq. (1),
 49 the goal of Back-to-Back Regression (B2B) is to estimate the matrix of E , i.e. to identify the factors
 50 that most reliably account for the multivariate observations.

51 2.1. Algorithm

52 Back-to-Back Regression (B2B) consists of two steps. First, we estimate the linear regression
 53 coefficients \hat{G} from Y to X , and construct the predictions $\hat{X} = Y\hat{G}$. This backward regression
 54 recovers the correlations between Y and each factor of X . Second, we estimate the linear regression
 55 coefficients \hat{H} from X to \hat{X} . The diagonal of the regression coefficients \hat{H} , denoted by $\hat{E} =$
 56 $\text{diag}(\hat{H})$, is the desired estimate of the causal influence matrix E , as detailed in the Appendix A.1.

If using l2-regularized least-squares [10, 26], B2B has a closed form solution:

$$\hat{G} = (Y^\top Y + \Lambda_Y)^{-1} Y^\top X, \quad (2)$$

$$\hat{H} = (X^\top X + \Lambda_X)^{-1} X^\top Y \hat{G}, \quad (3)$$

57 where Λ_X and Λ_Y are two diagonal matrices of regularization parameters, useful to invert the
 58 covariance matrices of X and Y if these are ill-conditioned.

59 Performing two regressions over the same data sample can result in overfitting, as spurious
 60 correlations in the data absorbed by the first regression will be leveraged by the second one. To
 61 avoid this issue, we split our sample (X, Y) into two splits (X_1, Y_1) and (X_2, Y_2) . Then, the first
 62 regression is performed using (X_1, Y_1) , and the second regression is performed using (X_2, Y_2) . To
 63 compensate for the reduction in sample size caused by the split, B2B is repeated over many random
 64 splits, and the final estimate \hat{E} of the causal influence matrix is the average over the estimates

65 associated to each split [2]. To accelerate this ensembling procedure, we implemented an efficient
 66 leave-one-out cross-validation scheme as detailed in [26] as follows:

$$\hat{Y} = (\Sigma_X G Y - \text{diag}(\Sigma_X G) Y) / \text{diag}(I - \Sigma_X G) \quad (\text{element-wise division}) \quad (4)$$

where Σ_X is the X kernel matrix and where G is computed with an eigen decomposition of X :

$$\begin{aligned} \Sigma_X &= Q V Q^T \\ G &= Q (V + \lambda I)^{-1} Q^T \end{aligned} \quad (5)$$

67 where Q , V and λ are the eigen vectors, eigen values and regularization, respectively.

68 We summarize the B2B procedure in Algorithm 1. The rest of this section provides a theoretical
 69 guarantee on the correctness of B2B.

Algorithm 1: Back-to-back regression.

Input: input data $X \in \mathbb{R}^{n \times d_x}$, output data $Y \in \mathbb{R}^{n \times d_y}$, number of repetitions $m \in \mathbb{N}$.

Output: estimate of causal influences $\hat{E} \in \mathbb{D}^{d_x \times d_x}$.

```

1  $\hat{E} \leftarrow 0$ ;
2 for  $i = 1, \dots, m$  do
3    $(X, Y) \leftarrow \text{ShuffleRows}((X, Y))$ ;
4    $(X_1, Y_1), (X_2, Y_2) \leftarrow \text{SplitRowsInHalf}((X, Y))$ ;
5    $\hat{G} = \text{LinearRegression}(Y_1, X_1)$ ;  $\triangleright \hat{G} = (Y_1^\top Y_1 + \Lambda_Y)^{-1} Y_1^\top X_1$ 
6    $\hat{H} = \text{LinearRegression}(X_2, Y_2 \hat{G})$ ;  $\triangleright \hat{H} = (X_2^\top X_2 + \Lambda_X)^{-1} X_2^\top Y_2 \hat{G}$ 
7    $\hat{E} \leftarrow \hat{E} + \text{diag}(\hat{H})$ ;
8 end
9  $\hat{E} \leftarrow \hat{E} / m$ ;
10  $\hat{W} \leftarrow \text{LinearRegression}(X \hat{E}, Y)$ ;
11 return  $\hat{E}, \hat{W}$ 

```

71 **2.2. Theoretical guarantees**

72 **Theorem 1** (B2B consistency - general case). *Consider the B2B model from Equation $Y = (XE +$
 73 $N)F$, N centered and full rank noise. Let $\text{Img}(M)$ refers to the image of the matrix M . If F and
 74 X are full-rank on the $\text{Img}(E)$, then, the solution of B2B, \hat{H} , will minimize $\min_H \|X - XH\|^2 +$
 75 $\|NH\|^2$ and satisfy $E\hat{H} = \hat{H}$*

76 *Proof.* See Appendix Appendix A.1. □

Since $E\hat{H} = \hat{H}$, we have

$$\hat{H} = \arg \min_H \|X - XEH\|^2 + \|NEH\|^2 = (EX^\top XE + EN^\top NE)^\dagger EXX^\top. \quad (6)$$

77 Assuming, without loss of generality, that the active features in E are the $k \in \mathbb{Z} : k \in [0, d_x]$
 78 first features, and rewriting $X = (X_1, X_2)$ and $N = (N_1, N_2)$ (X_1 and N_1 containing the k first
 79 features), we have:

$$X^\top X = \begin{pmatrix} \Sigma_{X_1 X_1} & \Sigma_{X_1 X_2} \\ \Sigma_{X_1 X_2} & \Sigma_{X_2 X_2} \end{pmatrix}, \quad N^\top N = \begin{pmatrix} \Sigma_{N_1 N_1} & \Sigma_{N_1 N_2} \\ \Sigma_{N_1 N_2} & \Sigma_{N_2 N_2} \end{pmatrix}, \quad (7)$$

where Σ_{AB} is the covariance of A and B , and:

$$\hat{H} = \begin{pmatrix} (\Sigma_{X_1X_1} + \Sigma_{N_1N_1})^{-1}\Sigma_{X_1X_1} & (\Sigma_{X_1X_1} + \Sigma_{N_1N_1})^{-1}\Sigma_{X_1X_2} \\ 0 & 0 \end{pmatrix} \quad (8)$$

$$\text{diag}_k(\hat{H}) = \text{diag}((\Sigma_{X_1X_1} + \Sigma_{N_1N_1})^{-1}\Sigma_{X_1X_1}) = \text{diag}((I + \Sigma_{X_1X_1}^{-1}\Sigma_{N_1N_1})^{-1}) \quad (9)$$

In the absence of noise, we have $\Sigma_{N_1N_1} = 0$, and so $\text{diag}_k(\hat{H}) = I$, and

$$\text{diag}(\hat{H}) = \text{diag}(E)$$

80 Therefore, we recover E from \hat{H} .

81 In the presence of noise, the causal factors of E correspond to the positive elements of $\text{diag}(\hat{H})$.

82 The methods to recover them are presented in the Appendix Appendix A.4.

83 3. Experiments

84 We perform two sets of experiments to evaluate B2B: one on controlled synthetic data, and a
85 second one on a real, large-scale magneto-encephalography (MEG) dataset. We use scikit-learn's
86 PLS and RidgeCV [23] as well as Pyrrca's regularized canonical component analysis (RegCCA,
87 [1]) objects to compare B2B against the standard baselines.

88 3.1. Synthetic data

89 We evaluate the performance of B2B throughout a series of experiments on controlled synthetic
90 data. The purpose of these experiments is to evaluate the ability of B2B in terms of prediction of
91 independent and identically distributed data, as well as a method to recover causal factors.

92 The data generating process for each experiment constructs $n = 1000$ training examples
93 according to the model $Y = (hXE + N)F$, where h is a scalar that modulates the signal-to-noise
94 ratio. Here, $F \in \mathbb{R}^{d_x \times d_y}$ contains entries drawn from $\mathcal{N}(0, \sigma^2)$ where σ^2 is inversely proportional
95 to d_x , $X \in \mathbb{R}^{n \times d_x}$ contains rows drawn from $\mathcal{N}(0, \Sigma_X)$, $N \in \mathbb{R}^{n \times d_x}$ contains rows drawn from
96 $\mathcal{N}(0, \Sigma_N)$, $E \in \mathbb{R}^{d_x \times d_x}$ is a binary diagonal matrix containing n_c ones, $\Sigma_X = AA^\top$ where
97 $A \in \mathbb{R}^{d_x \times d_x}$ contains entries drawn from $\mathcal{N}(0, \sigma^2)$, $\Sigma_N = BB^\top$ where $B \in \mathbb{R}^{d_x \times d_x}$ contains
98 entries drawn from $\mathcal{N}(0, \sigma^2)$, and the factor $h \in \mathcal{R}_+$.

99 To simulate a wide range of experimental conditions, we sample 10 values in log-space for
100 $d_x, d_y \in [10, 100]$, $n_c \in [3, 63]$, $h \in [0.001, 10]$. We discard the cases where $n_c > d_x$, limit d_x, d_y
101 to 100 to keep the running time under 2 hours for each condition, and average over 5 random seeds.

102 We compare the performance of B2B against four competing methods, all implemented in
103 scikit-learn [23] and pyrcca [1]:

104 3.2. Baseline models

Forward regression consists of an l_2 -regularized "ridge" regression from the putative causes X to the observations Y :

$$H_{fwd} = (X^T X + \lambda I)^{-1} X^T Y \quad (10)$$

Backward regression consists of an ℓ_2 -regularized "ridge" regression from Y to X :

$$G_{bwd} = (Y^T Y + \lambda I)^{-1} Y^T X \quad (11)$$

CCA finds $G_{cca} \in \mathbb{R}^{d_z, d_y}$ and $H_{cca} \in \mathbb{R}^{d_z, d_x}$ s.t. X and Y are maximally correlated in a latent Z space:

$$G_{cca}, H_{cca} = \operatorname{argmax}_{G, H} \operatorname{corr}(X H^T, Y G^T) \quad (12)$$

PLS finds $G_{pls} \in \mathbb{R}^{d_z, d_y}$ and $H_{pls} \in \mathbb{R}^{d_z, d_x}$ s.t. X and Y are maximally covarying in a latent Z space:

$$G_{pls}, H_{pls} = \operatorname{argmax}_{G, H} \operatorname{cov}(X H^T, Y G^T) \quad (13)$$

105 We employ five-fold cross-validation to select the optimal number of components for CCA and
106 PLS. Regressions were ℓ_2 -regularized with a λ regularization parameters fitted with the efficient
107 leave-one-out procedure implemented in scikit-learn RidgeCV [23].

108 3.3. Evaluating Causal Discovery from models' coefficients

109 B2B leads to unbiased (i.e. zeros-centered) scalar coefficients for non-causal features. In
110 contrast, the Forward, Backward, CCA and PLS models lead to a loading vector H_i per feature
111 i (or one vector G^i for the backward model). To transform such vector into an estimated causal
112 contribution \hat{E} , we take the sum of square coefficients: $\hat{E}_i = \sum_j H_i^{j^2}$

113 To estimate whether models accurately identify causal factors, we compute the area-under-the-
114 curve (AUC) across factors $AUC(E, \hat{E})$. The AUC allows evaluating the capacity of models at
115 detecting the causal importance of factors when ground truth labels are available, as is the case in
116 this setup.

117 We report AUC results in Figures 2 (top) and B.5 (left, in Appendix), and compare favorably to
118 all baselines.

119 3.4. Evaluating Causal Discovery with held-out prediction reliability

120 In most cases, E is not known and AUC can thus not be estimated. To address this issue, we
121 assess the ability of each model to reliably predict independent and identically distributed data from
122 Y , given all of the X features versus all-but-ones feature X_{-i} (i.e. 'knock-out X '). This procedure
123 results in two correlation metrics R_{full} and $R_{knockout}$, whose difference $\Delta R_i = R_{full} - R_{knockout}$
124 indicates how much each X_i improves the prediction of Y . In our figures, ΔR is the average of
125 ΔR_i . A higher score means that for prediction, the model relies on individual features rather than
126 combinations of features.

127 We show in Appendix Appendix A.3 pseudo-code to assess feature importance for our algorithm
128 as well as baselines. For the Backward Model, feature importance cannot be assessed as the X
129 collinearity is never taken into account.

130 We show in Figures 2 (bottom) and B.5 (right, in Appendix) that our method outperforms
131 baselines.

132 Next, we apply our method to brain imaging data from the anonymized multimodal neuroimag-
133 ing "Mother Of all Unification Studies" (MOUS) dataset [27]. The dataset contains magneto-
134 encephalography (MEG) recordings of 102 healthy native-Dutch adults who participated in a

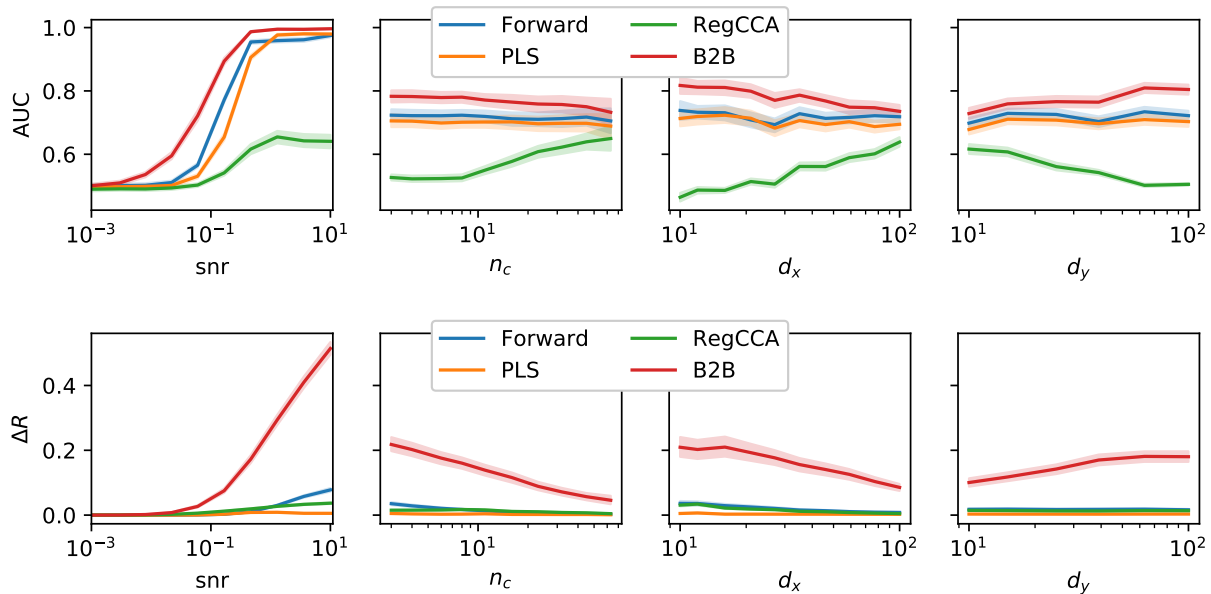


Figure 2: Synthetic experiments. Average AUC (top) and Feature Importance ΔR (bottom) when varying experimental conditions individually. Higher is better. B2B compares favorably in all cases.

135 reading task. Twelve subjects were excluded from the analysis because of corrupted file headers.
 136 Subjects were exposed to a rapid serial visual presentation of Dutch words. The word lists consisted
 137 of 120 sentences, and scrambled lists of the same words. Each word was presented on the computer
 138 screen for 351ms on average (min: 300ms, max: 1400ms). Successive words were separated by
 139 a blank screen for 300ms, and successive sentences were separated by an empty screen for a few
 140 (3-4) seconds.

141 3.4.1. MEG preprocessing

142 The raw MEG data was bandpass-filtered between 0.1 and 40Hz using MNE-Python default
 143 parameters [6, 7]. Specifically, we used a zero-phase finite impulse response filter (FIR) with
 144 a Hamming window and with transition bands of 0.1Hz and 10Hz for the low and high cut-off
 145 frequencies. The raw data was then segmented 100ms before word onset and 1s after word onset
 146 ($t = 0$ ms corresponds to word onset). Finally, each resulting segment was baseline-corrected
 147 between -100ms and 0ms, and decimated by 5 and thus led a sampling frequency of 240Hz. The
 148 average responses across words is displayed in Figure 3. For each subject and each time sample
 149 relative to word onset, we build an observation matrix $Y \in \mathbb{R}^{n \times d_y}$ of $n \approx 2,700$ words by $d_y = 301$
 150 MEG channels (273 magnetometers and 28 compensation channels). Each of the columns of Y is
 151 normalized to have zero mean and unit variance.

152 3.4.2. Feature definition

153 We aim to identify the word features that cause a variation in brain responses. We consider four
 154 distinct but linearly-correlated features. First, 'Word Length' refers to the total number of letters.
 155 Word Length is expected to specifically cause a variation in the early evoked MEG responses

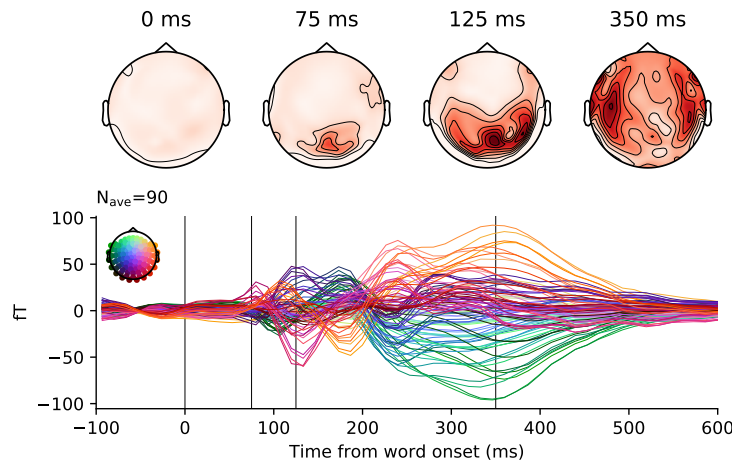


Figure 3: A hundred subjects read approximately 2,700 words while their brain activity was recorded with MEG. Top. Average brain response to words (word onset at $t=0$ ms), as viewed from above the head (red= higher gradient of magnetic flux). Bottom. Each line represents a magnetometer, color-coded by its spatial position. Posterior responses, typical of primary visual cortex activity, peak around 100 ms after word onset and are followed by an anterior propagation of activity typical of semantic processing in the associative cortices.

156 (i.e. from 100 ms after stimulus onset) elicited by the retinotopically-tuned visual cortices (e.g.
 157 [24].). Second, 'Word Frequency' indexes how frequently each word appears in Dutch and was
 158 derived with the the Zipf logarithmic scale of [31] provided by the WordFreq package [29]. Word
 159 Frequency is expected to specifically cause a variation in the late evoked MEG responses (i.e. from
 160 400 ms), because it variably engages semantic processes in the temporal cortices [18]. Third, 'Word
 161 Function' indicates whether each word is a content word (i.e. a noun, a verb, an adjective or an
 162 adverb) or a function word (i.e. a preposition, a conjunction, a determinant, a pronoun or a numeral),
 163 and was derived from Spacy's part of speech tagger [11]. To our knowledge, this feature has not
 164 been thoroughly investigated with MEG. Its causal contribution to reading processes in the brain
 165 thus remains unclear. Finally, to verify that B2B and other methods would not inadequately identify
 166 non-causal features, we added a dummy feature, constructed from a noisy combination of Word
 167 Length and Word Frequency: $dummy = z(length) + z(frequency) + \mathcal{N}$, where z normalizes
 168 features and \mathcal{N} is a random vector sampling Gaussian distribution (all terms thus have a zero-mean
 169 and a unit-variance). This procedure yields an $X \in \mathbb{R}^{n \times d_x}$ matrix of $n \approx 2,700$ words by $d_x = 4$
 170 features for each subject. Each of the columns of X is normalized to have a mean and a standard
 171 deviation of 0 and 1 respectively.

172 3.4.3. Models and statistics

173 We compare B2B to four standard methods: Forward regression, Backward regression, CCA and
 174 PLS, as implemented in scikit-learn [23] and [1], and optimized with nested cross-validation over
 175 twenty l_2 regularization parameters logarithmically spaced between 10^{-4} and 10^4 (for regression
 176 and CCA methods) or 1 to 4 canonical components (for PLS).

177 We used the feature importance described in Algorithm 2 to assess the extent to which each
 178 feature X_i specifically improves the prediction of held-out Y data, using a five-fold cross-validation
 179 (with shuffled trials to homogenize the distributions between the training and testing splits).

180 Each model was implemented for each subject and each time sample independently. Pairwise
 181 comparison between models were performed using a two-sided Wilcoxon test across subjects
 182 ($n=90$) using the average ΔR across time. Corresponding effect sizes are shown in Figure 4, and
 183 p-values are reported below.

184 **3.4.4. Results**

185 We compared the ability of For-
 186 ward regression, Backward regres-
 187 sion, CCA, PLS and B2B to es-
 188 timate the causal contribution of
 189 four distinct but linearly-correlated
 190 features on brain evoked responses
 191 to words.

192 As expected, the Backward
 193 model reveals a similar decod-
 194 ing time course for Word Length
 195 and Word Frequency, even though
 196 these features are known to specifi-
 197 cally influence early and late
 198 MEG responses respectively [18].
 199 In addition, the same decoding
 200 time course was observed for the
 201 dummy variable. These results il-
 202 lustrate that backward modeling
 203 cannot be used to estimate the
 204 causal contribution of correlated
 205 features.

206 We thus focus on the four re-
 207 maining methods (i.e. Forward
 208 Regression, PLS, CCA, and B2B)
 209 and estimate their ΔR (i.e. the im-
 210 provement of Y prediction induced
 211 by the introduction of a given fea-
 212 ture into the model, as described
 213 in Algorithm 2). Contrary to the
 214 Backward Model, none of the mod-
 215 els predicted the Dummy Variable
 216 to improve the Y prediction: all
 217 $\Delta R < 0$ (all $p > .089$).

218 Figure 4 shows, for each
 219 model, the effects obtained across
 220 time (left) and subjects (right).

221 Word Length and Word Fre-
 222 quency improved the prediction
 223 performance of all methods: $\Delta R >$
 224 0 for all models (all $p < 0.0001$).
 225 As expected, the time course asso-
 226 ciated with Word Length and Word

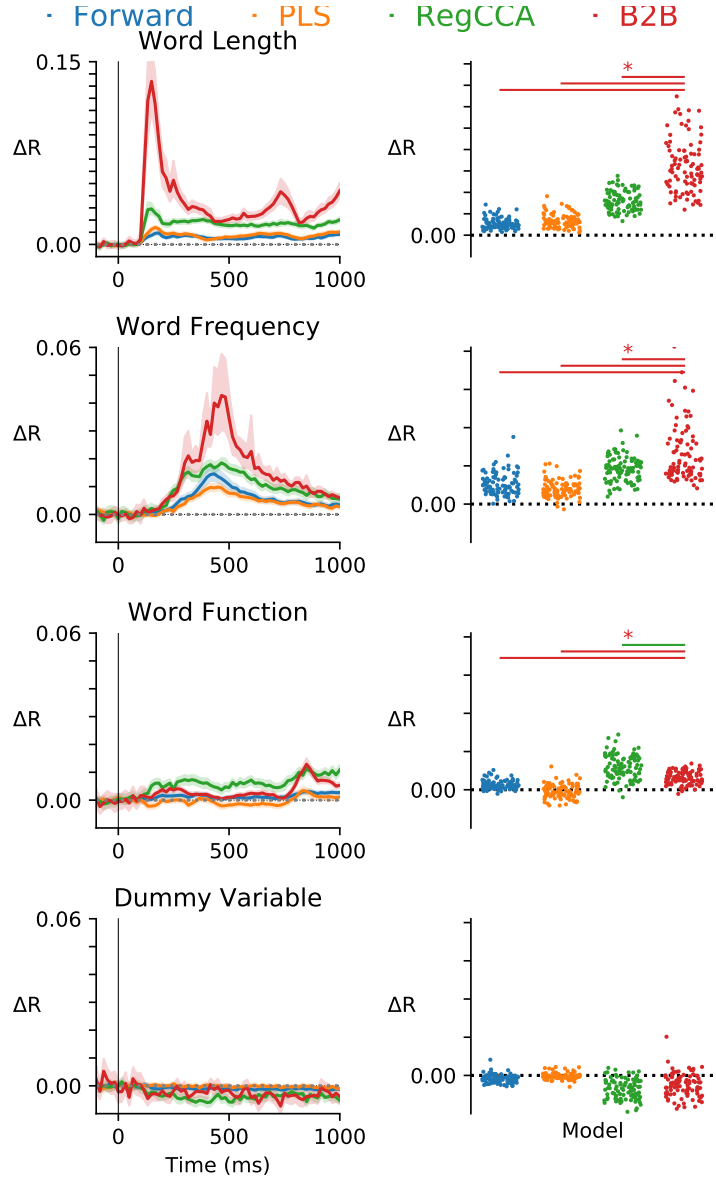


Figure 4: Multiple models (color-coded) are compared on their ability to reliably predict single-trial MEG signals evoked by words. Left. Average improvement of correlation coefficient ΔR for each of the four features (rows). Error bars indicate standard error of the mean (SEM) across subjects. Right. Average ΔR across time for each subject (dots). Top horizontal lines indicate when B2B significantly outperforms other methods (red) and vice versa.

227 Frequency rose from ≈ 100 ms and from ≈ 400 ms respectively. Furthermore, Word Function
228 improved the prediction performance of all models (all $p < 0.0002$) except for PLS ($p = 0.7989$).
229 Overall, these results confirm that Word Length, Word Frequency and Word Function causally
230 influence specific periods of brain responses to words.

231 To assess which model would be most sensitive to these causal discoveries, we compared B2B
232 to other models across subjects (Figure 4 right). For Word Length B2B outperforms all models
233 (all $p < 0.00001$) but CCA ($p = 0.0678$). For Word Frequency, B2B outperforms all models
234 (all $p < 0.0006$). For "Word Function", B2B outperforms all models (all $p < 0.0015$). Overall,
235 these results show that B2B reliably outperforms standard methods, especially when the effects are
236 difficult to detect.

237 *3.5. Magnetoencephalography data*

238 **4. Related work**

239 Forward and cross-decomposition models have been used to identify the causal contribution
240 of correlated factors onto multi-dimensional observations (e.g. [21]). These approaches typically
241 lead to multiple coefficients for each features (i.e. one per dimension of Y or one per component
242 respectively). Furthermore, these coefficients can be difficult to summarize into a single causal
243 estimate. By contrast, B2B quickly (Fig. B.6) leads to a single unbiased scalar values \hat{E} tending
244 towards 1 and 0 for causal and non-causal features respectively.

245 A variety of other statistical methods applied to neuroimaging data have been proposed to
246 clarify what is being represented in brain responses - i.e. what feature causes specific brain activity.
247 One of the popular linear method is Representational Similarity Analysis (RSA) [17], and consists
248 in analyzing the similarity of brain responses associated with specific categorical conditions (e.g.
249 distinct images), by (1) fitting one-against-all classifiers on each condition and (2) testing whether
250 these classifiers can discriminate all other conditions. The resulting confusion matrix is then
251 analyzed in an unsupervised manner to reveal which conditions lead to similar brain activity
252 patterns. B2B differs from RSA in that (1) it uses regressions instead of classifications, and can
253 thus generalize to new items and new contexts and (2) it is fully supervised.

254 Finally, CCA has been used in neuroimaging for a variety of purposes such as denoising and
255 subject alignment [12, 4]. While CCA relates to B2B, these two methods diverge in several ways.
256 First, CCA and B2B have different objectives: CCA aims to find the potentially numerous and
257 poorly interpretable components where X and Y are maximally correlated, whereas B2B aims to
258 recover the causal factors from X to Y . Second, B2B is not symmetric between X and Y : it aims
259 to identify specific causal features by first optimizing over the decoders G and then over H . By
260 contrast, CCA is symmetric between X and Y , and aims to find G and H such that they project
261 X and Y on maximally correlated dimensions. Third, CCA is based an eigen decomposition of
262 XH and YG - the corresponding canonical components are thus mixing the X features in way that
263 limit interpretability and potentially dilute the impact of each feature onto multiple components. In
264 contrast B2B assesses each feature X^j on a single Y component specifically selected to maximize
265 signal-to-noise ratio of that feature j . Fourth, and unlike B2B, CCA does not separately optimize
266 two distinct regularization parameters for G and H . Finally, CCA does not use different data splits
267 to estimate G and H . Together, these differences may explain why B2B reliably outperform CCA
268 on estimating causal influences (Figs. 2 and B.5).

269 5. Conclusion

270 In this work, we proposed Back-to-Back (B2B) regression, a linear method to disentangle
271 confounded factors from multidimensional observations. B2B repeatedly performs two successive
272 multidimensional regressions on independent subsets of the data: the first regression is applied on
273 the output domain (as in backward decoding), whereas the second regression is applied on the input
274 domain (as in forward encoding). We provided a theoretical guarantee about the consistency of
275 B2B, and compared it to several baselines in controlled synthetic experiments. We also applied B2B
276 to a recent brain imaging dataset, analyzing the timing of brain responses and their connection to
277 word features. We obtained results consistent with prior work in neuroscience literature, confirming
278 the reliability of B2B for real data analysis.

279 6. Acknowledgements

280 We are thankful to Gael Varoquaux and Alexandre Gramfort for the valuable feedback. This
281 work was supported by ANR-17-EURE-0017 and the Fyssen Foundation to JRK.

282 References

- 283 [1] Natalia Y Bilenko and Jack L Gallant. Pyrcca: regularized kernel canonical correlation analysis in python and its
284 applications to neuroimaging. *Frontiers in neuroinformatics*, 10:49, 2016.
- 285 [2] Leo Breiman. Bagging predictors. *Machine learning*, 24(2):123–140, 1996.
- 286 [3] Radoslaw Martin Cichy, Dimitrios Pantazis, and Aude Oliva. Resolving human object recognition in space and
287 time. *Nature neuroscience*, 17(3):455, 2014.
- 288 [4] Alain de Cheveigne, Giovanni M Di Liberto, Dorothee Arzounian, Daniel DE Wong, Jens Hjortkjaer, Søren
289 Fuglsang, and Lucas C Parra. Multiway canonical correlation analysis of brain data. *NeuroImage*, 186:728–740,
290 2019.
- 291 [5] Karl J Friston, Andrew P Holmes, Keith J Worsley, J-P Poline, Chris D Frith, and Richard SJ Frackowiak.
292 Statistical parametric maps in functional imaging: a general linear approach. *Human brain mapping*, 2(4):189–
293 210, 1994.
- 294 [6] Alexandre Gramfort, Martin Luessi, Eric Larson, Denis A Engemann, Daniel Strohmeier, Christian Brodbeck,
295 Roman Goj, Mainak Jas, Teon Brooks, Lauri Parkkonen, et al. Meg and eeg data analysis with mne-python.
296 *Frontiers in neuroscience*, 7:267, 2013.
- 297 [7] Alexandre Gramfort, Martin Luessi, Eric Larson, Denis A Engemann, Daniel Strohmeier, Christian Brodbeck,
298 Lauri Parkkonen, and Matti S Hämäläinen. Mne software for processing meg and eeg data. *Neuroimage*,
299 86:446–460, 2014.
- 300 [8] Clive WJ Granger. Investigating causal relations by econometric models and cross-spectral methods. *Economet-*
301 *rica: Journal of the Econometric Society*, pages 424–438, 1969.
- 302 [9] Martin N Hebart and Chris I Baker. Deconstructing multivariate decoding for the study of brain function.
303 *Neuroimage*, 180:4–18, 2018.
- 304 [10] Arthur E Hoerl. Optimum solution of many variables equations. *Chemical Engineering Progress*, 55(11):69–78,
305 1959.
- 306 [11] Matthew Honnibal and Ines Montani. spacy 2: Natural language understanding with bloom embeddings,
307 convolutional neural networks and incremental parsing. *To appear*, 2017.
- 308 [12] H. Hotelling. Relations between two sets of variables. *Biometrika*, (28):129–149, 1936.
- 309 [13] Alexander G Huth, Wendy A de Heer, Thomas L Griffiths, Frédéric E Theunissen, and Jack L Gallant. Natural
310 speech reveals the semantic maps that tile human cerebral cortex. *Nature*, 532(7600):453, 2016.
- 311 [14] Dominik Janzing, David Balduzzi, Moritz Grosse-Wentrup, Bernhard Schölkopf, et al. Quantifying causal
312 influences. *The Annals of Statistics*, 41(5):2324–2358, 2013.

- 313 [15] G. V. Kass. Significance testing in automatic interaction detection (a.i.d.). *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, 24(2):178–189, 1975.
- 314
- 315 [16] Jean-Rémi King, Laura Gwilliams, Chris Holdgraf, Jona Sassenhagen, Alexandre Barachant, Denis Engemann,
316 Eric Larson, and Alexandre Gramfort. Encoding and decoding neuronal dynamics: Methodological framework to
317 uncover the algorithms of cognition, 2018.
- 318 [17] Nikolaus Kriegeskorte, Marieke Mur, and Peter A Bandettini. Representational similarity analysis-connecting the
319 branches of systems neuroscience. *Frontiers in systems neuroscience*, 2:4, 2008.
- 320 [18] Marta Kutas and Kara D Federmeier. Thirty years and counting: finding meaning in the n400 component of the
321 event-related brain potential (erp). *Annual review of psychology*, 62:621–647, 2011.
- 322 [19] Ludovic Lebart, Alain Morineau, and Marie Piron. *Statistique exploratoire multidimensionnelle*, volume 3.
323 Dunod Paris, 1995.
- 324 [20] J. A. Morgan and J. N. Sonquist. Problems in the analysis of survey data: and a proposal. *J. Amer. Statist. Ass.*,
325 (58):415–434, 1963.
- 326 [21] Thomas Naselaris, Kendrick N Kay, Shinji Nishimoto, and Jack L Gallant. Encoding and decoding in fmri.
327 *Neuroimage*, 56(2):400–410, 2011.
- 328 [22] Kenneth A Norman, Sean M Polyn, Greg J Detre, and James V Haxby. Beyond mind-reading: multi-voxel pattern
329 analysis of fmri data. *Trends in cognitive sciences*, 10(9):424–430, 2006.
- 330 [23] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel,
331 Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. Scikit-learn: Machine learning in
332 python. *Journal of machine learning research*, 12(Oct):2825–2830, 2011.
- 333 [24] Felipe Pegado, Enio Comerlato, Fabricio Ventura, Antoinette Jobert, Kimihiro Nakamura, Marco Buiatti, Paulo
334 Ventura, Ghislaine Dehaene-Lambertz, Régine Kolinsky, José Morais, et al. Timing the impact of literacy on
335 visual processing. *Proceedings of the National Academy of Sciences*, 111(49):E5233–E5242, 2014.
- 336 [25] Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of causal inference: foundations and learning*
337 *algorithms*. MIT press, 2017.
- 338 [26] Ryan M Rifkin and Ross A Lippert. Notes on regularized least squares. Technical report, MIT, 2007.
- 339 [27] Jan-Mathijs Schoffelen, Robert Oostenveld, Nietzsche HL Lam, Julia Uddén, Annika Hultén, and Peter Hagoort.
340 A 204-subject multimodal neuroimaging dataset to study language processing. *Scientific data*, 6(1):17, 2019.
- 341 [28] Bernhard Schölkopf, David W Hogg, Dun Wang, Daniel Foreman-Mackey, Dominik Janzing, Carl-Johann
342 Simon-Gabriel, and Jonas Peters. Modeling confounding by half-sibling regression. *Proceedings of the National*
343 *Academy of Sciences*, 113(27):7391–7398, 2016.
- 344 [29] Robyn Speer, Joshua Chin, Andrew Lin, Sara Jewett, and Lance Nathan. Luminosinsight/wordfreq: v2.2,
345 October 2018.
- 346 [30] Nicholas A Steinmetz, Christof Koch, Kenneth D Harris, and Matteo Carandini. Challenges and opportunities for
347 large-scale electrophysiology with neuropixels probes. *Current opinion in neurobiology*, 50:92–100, 2018.
- 348 [31] Walter JB Van Heuven, Pawel Mandera, Emmanuel Keuleers, and Marc Brysbaert. Subtlex-uk: A new and
349 improved word frequency database for british english. *The Quarterly Journal of Experimental Psychology*,
350 67(6):1176–1190, 2014.
- 351 [32] Sebastian Weichwald, Timm Meyer, Ozan Özdenizci, Bernhard Schölkopf, Tonio Ball, and Moritz Grosse-
352 Wentrup. Causal interpretation rules for encoding and decoding models in neuroimaging. *NeuroImage*, 110:48–59,
353 2015.

354 **7. Appendices**

355 **Appendix A. Appendix**

356 *Appendix A.1. Proof of consistency theorem*

357 Proof of the theorem in 2.2:

Theorem 2 (B2B consistency - general case). *Consider the B2B model from equation 1*

$$Y = (XE + N)F$$

358 *with N centered and full rank noise.*

359 *If F and X are full-rank on $\text{Img}(E)$, then, the solution of B2B, \hat{H} minimizes*

$$\min_H \|X - XH\|^2 + \|NH\|^2$$

360 *and satisfies*

$$E\hat{H} = \hat{H}$$

361 *Proof.* Let \hat{G} and \hat{H} be the solutions of the first and second regressions of B2B.

Since \hat{G} is the least square estimator of X from Y

$$\hat{G} = \arg \min_G \mathbb{E}[\|YG - X\|^2]$$

Replacing Y by its model definition $Y = (XE + N)F$, we have

$$\hat{G} = \arg \min_G \mathbb{E}[\|X - (XE + N)FG\|^2] = \arg \min_G \mathbb{E}[\|X - XEFG + NFG\|^2]$$

Since N is centered and independent of X , we have

$$\hat{G} = \arg \min_G \|X - XEFG\|^2 + \|NFG\|^2 \tag{A.1}$$

In the same way, for \hat{H} , we have

$$\begin{aligned} \hat{H} &= \arg \min_H \mathbb{E}[\|XH - Y\hat{G}\|^2] = \arg \min_H \mathbb{E}[\|XH - (XE + N)F\hat{G}\|^2] \\ &= \arg \min_H \mathbb{E}[\|X(H - EF\hat{G})\|^2] + \mathbb{E}[\|NF\hat{G}\|^2] \\ &= \arg \min_H \mathbb{E}[\|X(H - EF\hat{G})\|^2] \end{aligned}$$

a positive quantity which reaches a minimum (zero) for

$$\hat{H} = EF\hat{G} \tag{A.2}$$

362 Let us now prove that $EF\hat{G} = F\hat{G}$.

363 Let F^\dagger be the pseudo inverse of F , and $Z = F^\dagger EF\hat{G}$, we have $FZ = FF^\dagger EF\hat{G}$

364 Since F is full rank on $\text{Img}(E)$, we have $FF^\dagger E = E$, and $FZ = EF\hat{G}$

As E is a binary diagonal matrix, it is an orthogonal projection and therefore a contraction, thus

$$\|NEFG\hat{G}\|^2 \leq \|NFG\hat{G}\|^2$$

and

$$\|X - XEFZ\|^2 + \|NFZ\|^2 = \|X - XEF\hat{G}\|^2 + \|NEF\hat{G}\|^2 \leq \|X - XEF\hat{G}\|^2 + \|NFG\hat{G}\|^2$$

But since $\hat{G} = \arg \min_G \|X - XEFG\|^2 + \|NFG\|^2$, we also have

$$\|X - XEF\hat{G}\|^2 + \|NFG\hat{G}\|^2 \leq \|X - XEFZ\|^2 + \|NFZ\|^2$$

Summarizing the above,

$$\|X - XEF\hat{G}\|^2 + \|NFG\hat{G}\|^2 \leq \|X - XEF\hat{G}\|^2 + \|NEF\hat{G}\|^2 \leq \|X - XEF\hat{G}\|^2 + \|NFG\hat{G}\|^2$$

$$\|X - XEF\hat{G}\|^2 + \|NFG\hat{G}\|^2 = \|X - XEF\hat{G}\|^2 + \|NEF\hat{G}\|^2$$

$$\|NFG\hat{G}\|^2 = \|NEF\hat{G}\|^2$$

365 N being full rank, this yields $EF\hat{G} = F\hat{G}$.

Replacing into (A.1), and setting $H = EFG$, we have

$$\begin{aligned} \hat{G} &= \arg \min_G \|X - XEFG\|^2 + \|NFG\|^2 \\ &= \arg \min_G \|X - XEFG\|^2 + \|NEFG\|^2 \\ \hat{H} &= \arg \min_H \|X - XH\|^2 + \|NH\|^2 \end{aligned}$$

366 Finally, $E\hat{H} = EEF\hat{G} = EF\hat{G} = \hat{H}$, since E , a binary diagonal matrix, is involutive. This
367 completes the proof. \square

368 *Appendix A.2. Modeling measurement noise*

Equation 1 does not explicitly contain a measurement noise term. Yet, in most experimental cases, the problem is best described as:

$$Y = (XE + N)F + M \quad (\text{A.3})$$

369 with $M \in \mathbb{R}^{n \times d_y}$.

This equation is actually equivalent to Equation 1 given our hypotheses. Indeed, we can rewrite $M = MF^{-1}F$ over $\text{Img}(F)$, which leads to:

$$Y = (XE + N)F + M = (XE + N + MF^{-1})F = (XE + N')F$$

370 Consequently, assuming that F is full rank on $\text{Img}(XE)$, B2B yields the same solutions to
371 equations 1 and A.3.

372 *Appendix A.3. Feature importance*

373 For B2B, feature importance is assessed as follows:

Algorithm 2: B2B feature importance.

Input: $X_{train} \in \mathbb{R}^{n \times d_x}$, $X_{test} \in \mathbb{R}^{n' \times d_x}$, $Y_{train} \in \mathbb{R}^{n \times d_y}$, $Y_{test} \in \mathbb{R}^{n' \times d_y}$,

Output: estimate of prediction improvement $\Delta R \in \mathbb{D}^{d_x}$.

1 $H, G = \text{B2B}(X_{train}, Y_{train});$
 2 $R_{full} = \text{corr}(X_{test}H, Y_{test}G);$
 3 **for** $i = 1, \dots, d_x$ **do**
 374 4 $K = Id;$
 5 $K[i] \leftarrow 0;$
 6 $R_k = \text{corr}(X_{test}KH, Y_{test}G_i);$
 7 $\Delta R_i = R_{full} - R_k;$
 8 **end**
 9 **return** ΔR

375 For the Forward Model, the feature importance is assessed as follows:

Algorithm 3: Forward feature importance.

Input: $X_{train} \in \mathbb{R}^{n \times d_x}$, $X_{test} \in \mathbb{R}^{n' \times d_x}$, $Y_{train} \in \mathbb{R}^{n \times d_y}$, $Y_{test} \in \mathbb{R}^{n' \times d_y}$,

Output: estimate of prediction improvement $\Delta R \in \mathbb{D}^{d_x, d_y}$.

1 $H = \text{LinearRegression}(X_{train}, Y_{train})$ $R_{full} = \text{corr}(X_{test}H, Y_{test});$
 2 **for** $i = 1, \dots, d_x$ **do**
 376 3 $K = Id;$
 4 $K[i] \leftarrow 0;$
 5 $R_k = \text{corr}(X_{test}KH, Y_{test});$
 6 $\Delta R_i = R_{full} - R_k;$
 7 **end**
 8 **return** ΔR

377 For the CCA and PLS models, the feature importance is assessed as follows:

Algorithm 4: CCA and PLS feature importance.

Input: $X_{train} \in \mathbb{R}^{n \times d_x}$, $X_{test} \in \mathbb{R}^{n' \times d_x}$, $Y_{train} \in \mathbb{R}^{n \times d_y}$, $Y_{test} \in \mathbb{R}^{n' \times d_y}$,

Output: estimate of prediction improvement $\Delta R \in \mathbb{D}^{d_x, d_z}$.

1 $H, G = \text{CCA}(X_{train}, Y_{train});$

2 $R_{full} = \text{corr}(X_{test}H, Y_{test}G);$

3 **for** $i = 1, \dots, d_x$ **do**

378

4 $K = Id;$

5 $K[i] \leftarrow 0;$

6 $R_k = \text{corr}(X_{test}KH, Y_{test}G);$

7 $\Delta R_i = R_{full} - R_k;$

8 **end**

9 **return** ΔR

379 For the Backward Model, feature importance cannot be assessed because there is no prediction.

380 Appendix A.4. Recovering E

381 In case of noise, B2B yields non binary \hat{E} . Three thresholding rules can be used to binarize its
382 values thus explicitly recover "causal" features.

383 First, given known signal-to-noise ratio, the threshold above which a feature should considered
384 to be "causal" can be derived analytically. Indeed, Equation 9 implies that the k first diagonal
385 elements of \hat{H} are bounded:

$$0 \leq \frac{\sigma_{X_k}}{\sigma_{X_k} + \sigma_{N_1}} \leq \text{diag}_k(\hat{H}) \leq \frac{\sigma_{X_1}}{\sigma_{X_1} + \sigma_{N_k}}$$

386 where σ_{X_1} , σ_{X_k} , σ_{N_1} and σ_{N_k} denote the largest and smallest eigenvalues of $\Sigma_{X_1 X_1}$ and $\Sigma_{N_1 N_1}$.

The average value μ of non-zero coefficients of $\text{diag}(\hat{H})$ is the trace of \hat{H} divided by k , and can be computed as

$$\mu = \frac{\text{Var}(X)}{\text{Var}(X) + \text{Var}(N)} \quad (\text{A.4})$$

387 The decision threshold between "causal" and "non-causal" elements is thus a fraction μ , whose
388 proportion arbitrarily depends on the necessity to favor type I and type II errors. In practice, we
389 cannot use this procedure for our MEG study, because signal-to-noise ratio is unknown.

Second, $\text{diag}(\hat{H})$ can be binarized with the Sonquist-Morgan criterion [20], a non-parametric clustering procedure separating small and large values in a given set. This procedure maximizes the ratio of inter-group variance while minimizing the intra-group variance, over all possible splits of the diagonal into p largest values and $d_x - p$ smallest values. Let m_0 and m_1 be the average values of the two clusters, p and $d_x - p$ their size, and v the total variance of the sample, Sonquist-Morgan criterion maximizes [15]:

$$\frac{p(d_x - p)}{d_x} \frac{(m_1 - m_0)^2}{v} \quad (\text{A.5})$$

390 This procedure assumes that there exists at least one causal and at least one non-causal feature.

391 Third, second-order statistics across multiple datasets can be used to identify the elements of

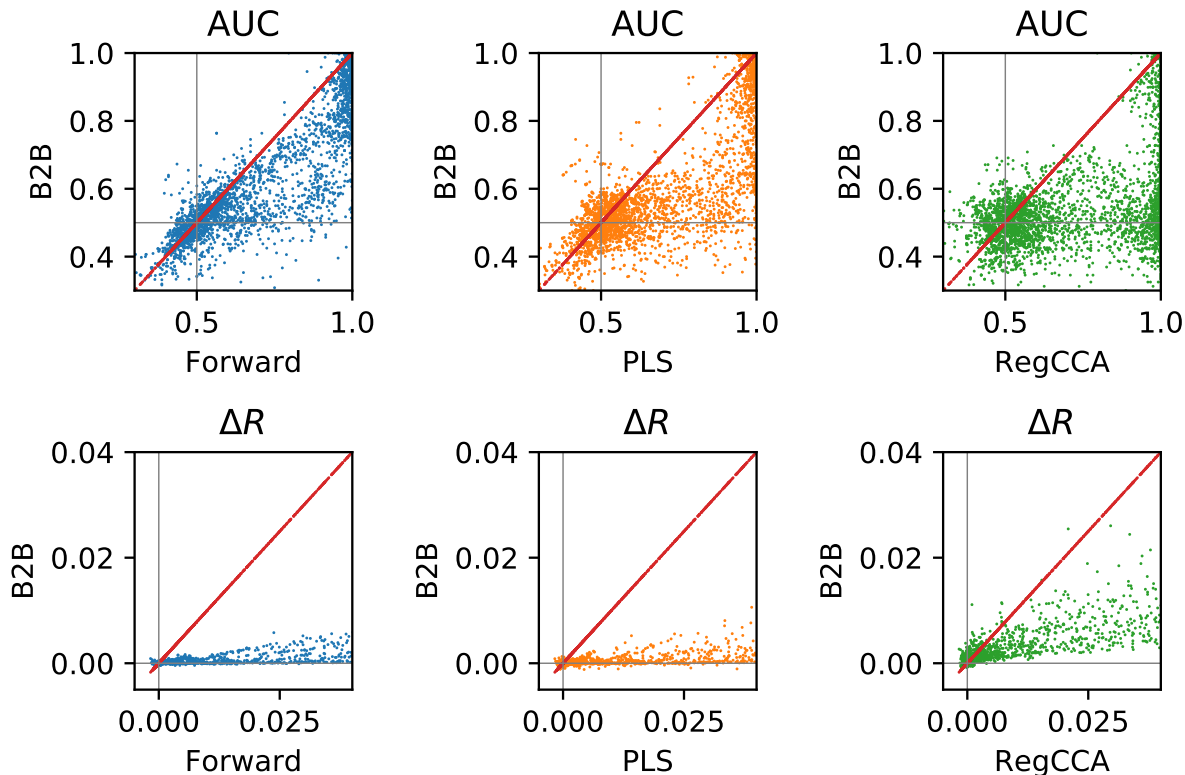


Figure B.5: Synthetic experiments. Distribution (over conditions) of AUC (top) and Feature Importance ΔR (bottom) metrics between our method (y-axis) and the baselines (x-axis). Each dot is a distinct synthetic experiment. Dots below the diagonal indicates that B2B outperform the tested model.

392 $\text{diag}(\hat{H})$ that are significantly different from 0. This procedure is detailed in the method section of
 393 our MEG experiment.

394 Overall, these three procedures thus vary in their additional assumptions: i.e. (1) a known
 395 signal-to-noise ratio, (2) the existence of both causal and non-causal factors or (3) independent
 396 repetitions of the experiment.

397 **Appendix B. Additional Figures**

398 *Appendix B.1. Robustness to increasing number of factors*

399 To test whether each of the methods robustly scales to an increasingly large number of potential
 400 causes X , we enhanced the four ad-hoc features (word length, word frequency, word function,
 401 dummy variable) with another ten features. These additional features corresponds to the first
 402 dimensions of word embedding as provided by Spacy [11]. The results shown in Figure B.7, show
 403 that the feature importance of ad-hoc features as derived by B2B remain unchanged and are actually
 404 improved.

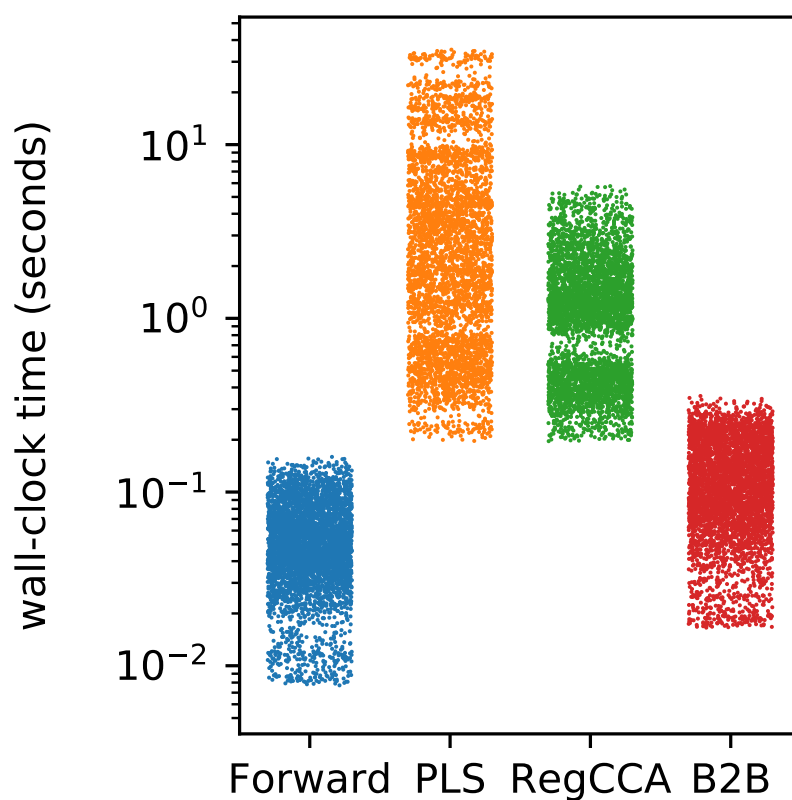


Figure B.6: Wall-clock run-time for our method B2B and for the baselines. Each dot is a distinct synthetic experiment. B2B runs much faster than cross-decomposition baselines.

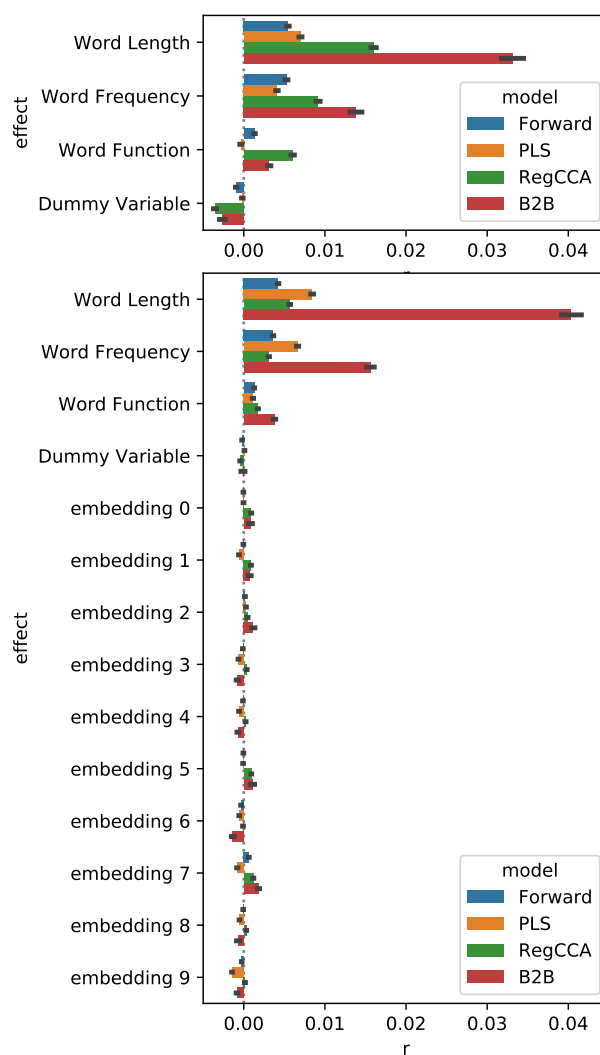


Figure B.7: Comparison of ΔR when the models are tested on four variables (top) and when the models are tested on an these four variables as well as another 10 word-embedding features (bottom). These results illustrate that, unlike Regularized CCA, B2B remains robust even when the number of tested factors increases.