# The right tool for the right question — beyond the encoding versus decoding dichotomy

Sebastian Weichwald, Moritz Grosse-Wentrup

Max Planck Institute for Intelligent Systems, Tübingen, Germany

[sweichwald, moritzgw]@tue.mpg.de

There are two major questions that neuroimaging studies attempt to answer: First, how are sensory stimuli represented in the brain (which we subsequently term the *stimulus-based* setting)? And, second, how does the brain generate cognition (subsequently termed the *response-based* setting)? There has been a lively debate in the neuroimaging community whether encoding and decoding models can provide insights into these questions (e. g. [Naselaris et al., 2011; Todd et al., 2013; Davis et al., 2014; Haufe et al., 2014; Woolgar et al., 2014; Weichwald et al., 2015]). In this commentary, we construct two simple and analytically tractable examples to demonstrate that while an encoding model analysis helps with the former, neither model is appropriate to satisfactorily answer the latter question. Consequently, we argue that if we want to understand how the brain generates cognition, we need to move beyond the encoding versus decoding dichotomy and instead discuss and develop tools that are specifically tailored to our endeavour.

Since the disagreement among researchers can partially be attributed to differing use of terminology, we begin by briefly introducing the two (types of) models that are subject of the ongoing debate and make explicit the terminology used throughout this comment. On the one hand, there are (univariate) *encoding models* (also referred to as forward or generative models) that approximate each neurophysiological variable $X_i$ as a function of the experimental condition $Y$ (e.g. statistical parametric mapping [Friston et al., 1994]), i. e.

$$\begin{bmatrix} \widehat{X}_1 \\ \vdots \\ \widehat{X}_n \end{bmatrix} = \mathrm{enc}(Y) = \begin{bmatrix} \mathrm{enc}_1(Y) \\ \vdots \\ \mathrm{enc}_n(Y) \end{bmatrix}.$$

On the other hand, there are (multivariate) *decoding models* (also referred to as backward or discriminative models or multi-voxel pattern analysis (MVPA)) that predict the experimental condition $Y$ given the neurophysiological variables $X_1, \ldots, X_n$ [Mitchell et al., 2004; Pereira et al., 2009], i. e.

$$\widehat{Y} = \mathrm{dec}(X_1, \ldots, X_n).$$

It has been argued that encoding mod-

els allow for richer interpretations than decoding models and can in principle provide a complete functional description of a brain state variable [Naselaris et al., 2011; Haufe et al., 2014]. Decoding models are often considered inherently difficult to interpret and may determine neurophysiological variables as relevant that are statistically independent of the experimental condition [Todd et al., 2013; Woolgar et al., 2014; Haufe et al., 2014]. Indeed, the difference between the two models in terms of their interpretation has experimentally been confirmed [Huth et al., 2016; Bach et al., 2017]. Further distinguishing between stimulus- and response-based paradigms has enabled a comprehensive overview over the theoretical limitations of interpreting either model alone or both at the same time [Weichwald et al., 2014, 2015].

We contribute to this discussion by providing two instructive examples that illustrate the fundamental problems in the interpretation of encoding and decoding models. Specifically, we show that even in the infinite-sample limit and under correct model assumptions the exclusion of relevant variables—e.g. when preselecting regions of interest or because of unobserved (latent) brain processes—may lead to incorrect conclusions about the qualitative and quantitative relations of brain processes and experimental variables. Our examples underline the importance of developing methods that are robust against such confounds.

# 1 Examples

In the following, we construct two hypothetical ground-truth models, the first for a stimulus- and the second for a response-based paradigm, and consider what interpretations one would obtain from encoding and decoding models fit to differing subsets of variables. This enables us to study deviations between true and estimated relations of brain processes and experimental variables under idealistic conditions.[1]

## 1.1 Stimulus-based paradigm

Assume we conduct a stimulus-based experiment to investigate the effect of a stimulus variable $S$ onto neurophysiological variables (e.g. showing images of different brightness to subjects while recording their brain activity in different regions of interest). Assume further that the true relationships are governed by the additive Gaussian noise model described in Figure 1 where upper and lower case letters denote variables and their linear relations, respectively. The linear effect of $S$ onto each variable is determined by the sum of weight products along each path from $S$ to that variable; e.g. the effect of $S$ onto $D$ is linear with slope $ab + ade = -2$ (i.e. increased image brightness leads to a twice as much decreased average activity in $D$[2]). The true effects of $S$ onto $V, D, C, T$ are $1, -2, 1, 2$ respectively.

---

[1] First, the assumptions for ordinary least squares regression are met since we assume the ground-truth models to be linear and Gaussian. Second, we consider the true ordinary least squares weights one would obtain in the limit of infinite data, hence ignoring finite sample estimation errors [Taylor and Tibshirani, 2015].

[2] More precisely, if we were to intervene and forcefully set $S$ to take the value 5 then $D$ will follow a Gaussian distribution with mean $-10$ while in the observational setting the mean is 0. It is these differences in the distribution of a variable (the effect) upon an intervention on another variable (the cause) that form the basis of the conceptualisation of interventional causation in structural equation models [Spirtes et al., 2000; Pearl, 2009]

In Figure 2 we show toy brain maps depicting the true effects of $S$ onto each variable as well as the outcome one would obtain when running one of the following three analyses of a system admitting the above ground-truth model:

**Analysis 1:** $\mathrm{enc}(S)$

A mass-univariate encoding analysis where for each variable $X \in \{V, D, C, T\}$ we interpret the weight $\beta_X$ in the model

$$\widehat{X} = \mathrm{enc}(S) = \alpha_X + \beta_X S$$

to reflect the effect of $S$ onto $X$. The true ordinary least squares weights one would obtain in the limit of infinite data are

$$\begin{bmatrix} \beta_V \\ \beta_D \\ \beta_C \\ \beta_T \end{bmatrix} = {}^1\!/\mathrm{var}(S) \begin{bmatrix} \mathrm{cov}(V, S) \\ \mathrm{cov}(D, S) \\ \mathrm{cov}(C, S) \\ \mathrm{cov}(T, S) \end{bmatrix} = \begin{bmatrix} 1 \\ -2 \\ 1 \\ 2 \end{bmatrix}$$

**Analysis 2:** $\mathrm{dec}(V, D, C, T)$

A decoding analysis including the variables $V, D, C, T$ where we interpret the weights of the model

$$\widehat{S} = \alpha + \beta V + \gamma D + \delta C + \epsilon T$$

to reflect the relation between $S$ and each respective variable. The true ordinary least squares weights one would obtain in the limit of infinite data are

$$\begin{bmatrix} \beta \\ \gamma \\ \delta \\ \epsilon \end{bmatrix} = \Sigma_{V,D,C,T}^{-1} \begin{bmatrix} \mathrm{cov}(V, S) \\ \mathrm{cov}(D, S) \\ \mathrm{cov}(C, S) \\ \mathrm{cov}(T, S) \end{bmatrix} = \begin{bmatrix} {}^1\!/2 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

**Analysis 3:** $\mathrm{dec}(D, C, T)$

A decoding analysis as above, this time

excluding the variable $V$. The true ordinary least squares weights one would obtain in the limit of infinite data are

$$\begin{bmatrix} \gamma \\ \delta \\ \epsilon \end{bmatrix} = \Sigma_{D,C,T}^{-1} \begin{bmatrix} \mathrm{cov}(D, S) \\ \mathrm{cov}(C, S) \\ \mathrm{cov}(T, S) \end{bmatrix} = \begin{bmatrix} -{}^1\!/7 \\ -{}^1\!/7 \\ {}^1\!/7 \end{bmatrix}$$

In our scenario, the encoding analysis would indeed reveal the correct effects of $S$. If instead we were to interpret the weights of the full linear decoding model $\mathrm{dec}(V, D, C, T)$ we would arrive at incorrect interpretations; e.g. in this model $\delta$ turns out to be zero suggesting that $S$ is not related to $C$. Further complicating the matter, if $V$ was excluded from the analysis (due to subjective variable selection criteria or being unobserved), i.e. when interpreting $\mathrm{dec}(D, C, T)$, we would arrive at yet other interpretations; e.g. we would obtain $\delta = -{}^1\!/7$, which may mislead researchers to believe that $S$ has an inhibitory effect on $C$.

## 1.2 Response-based paradigm

Next, assume we conduct a response-based experiment to investigate the neurophysiological causes of cognition (e.g. assessing subjects' performance in a motor task while recording their brain activity in different regions of interest). Assume further that the true relationships are governed by the additive Gaussian noise model desribed in Figure 3. The linear effect of each variable onto $P$ is determined by the sum of weight products along each path from that variable to $P$; e.g. the effect of $C$ onto $P$ is linear with slope $dg = -1$ (i.e. increased average activity in $C$ leads to a likewise decreased motor task performance). The true effects of $T, C, M, D$ onto $P$ are $1, -1, 1, 0$ respectively.

$$S = N_S$$
$$V = aS + N_V$$
$$D = bV + eC + N_D$$
$$T = cV + fH + N_T$$
$$C = dV + gH + N_C$$
$$H = N_H$$

where $N_S, N_V, N_D, N_T, N_C, N_H \overset{\text{iid}}{\sim} \mathcal{N}(0,1)$

Figure 1: Assumed true additive Gaussian noise model (cf. Section 1.1).



Figure 2: Toy brain maps of the true effects of $S$ in the model described in Figure 1 as well as the effects as inferred by three different models (cf. Section 1.1).

4

As in the previous section, Figure 4 shows the true effects of each variable onto $P$ as well as the outcome one would obtain when running one of the following three analyses of a system admitting the above ground-truth model:

**Analysis 1:** $\mathrm{enc}(R)$

A mass-univariate encoding analysis where for each variable $X \in \{T, C, M, D\}$ we interpret the weight $\beta_X$ in the model

$$\widehat{X} = \mathrm{enc}(R) = \alpha_X + \beta_X R$$

to reflect the relation between $X$ and $R$. The true ordinary least squares weights one would obtain in the limit of infinite data are

$$\begin{bmatrix} \beta_T \\ \beta_C \\ \beta_M \\ \beta_D \end{bmatrix} = 1/\mathrm{var}(R) \begin{bmatrix} \mathrm{cov}(T,R) \\ \mathrm{cov}(C,R) \\ \mathrm{cov}(M,R) \\ \mathrm{cov}(D,R) \end{bmatrix} = \begin{bmatrix} -1/5 \\ -4/5 \\ 4/5 \\ -12/5 \end{bmatrix}$$

**Analysis 2:** $\mathrm{dec}(T, C, M, D)$

A decoding analysis including the variables $T, C, M, D$ where we interpret the weights of the model

$$\widehat{R} = \alpha + \beta T + \gamma C + \delta M + \epsilon D$$

to reflect the effect of each respective variable onto $P$. The true ordinary least squares weights one would obtain in the limit of infinite data are

$$\begin{bmatrix} \beta \\ \gamma \\ \delta \\ \epsilon \end{bmatrix} = \Sigma_{T,C,M,D}^{-1} \begin{bmatrix} \mathrm{cov}(T,R) \\ \mathrm{cov}(C,R) \\ \mathrm{cov}(M,R) \\ \mathrm{cov}(D,R) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

**Analysis 3:** $\mathrm{dec}(T, C, D)$

A decoding analysis as above, this time excluding the variable $M$. The true ordinary least squares weights one would

obtain in the limit of infinite data are

$$\begin{bmatrix} \beta \\ \gamma \\ \epsilon \end{bmatrix} = \Sigma_{T,C,D}^{-1} \begin{bmatrix} \mathrm{cov}(T,R) \\ \mathrm{cov}(C,R) \\ \mathrm{cov}(D,R) \end{bmatrix} = \begin{bmatrix} 1/5 \\ 1/5 \\ -2/5 \end{bmatrix}$$

In contrast to the previous example, an encoding analysis does not correctly identify the relations between neurophysiological variables and $P$; e. g. it turns out that $\beta_T = -1/5$ which may mislead researchers to conclude that $T$ is negatively related to $P$ while indeed it has a positive effect with weight 1. The full linear decoding model also leads to incorrect interpretations; e. g. in this model $\beta$ turns out to be zero suggesting that $T$ is not related to $P$. Lastly, if we are in the setting where $M$ is unobserved or excluded from the analysis then we arrive at yet another interpretation; e. g. in this model $\beta = 1/5$ systematically underestimates the effect of $T$ onto $P$ and $\gamma = 1/5$ suggests a positive effect of $C$ onto $P$ while it indeed has a negative effect.

## 1.3 Results

Only the encoding analysis in the stimulus-based setting yields the desired results both qualitatively and quantitatively. All other analyses in both the stimulus- and response-based setting may lead to incorrect and contradictory interpretations.

## 2 Conclusion

In some particular instances encoding and decoding models may warrant limited partial answers to the questions mentioned at the outset. However, as above examples demonstrate, the interpretation crucially depends on the experimental paradigm employed, which variables are observable/observed, and which ones enter into

$$H = N_H$$
$$T = aH + N_T$$
$$C = bH + N_C$$
$$M = cT + dC + N_M$$
$$D = eC + fM + N_D$$
$$P = gM + N_P$$

where $N_H, N_T, N_C, N_M, N_D, N_P \overset{\text{iid}}{\sim} \mathcal{N}(0,1)$

Figure 3: Assumed true additive Gaussian noise model (cf. Section 1.2).



true        enc($R$)        dec($T,C,M,D$)        dec($T,C,D$)

$-$                        0                        $+$

Figure 4: Toy brain maps of the true causes of $P$ in the model described in Figure 3 as well as the causes as inferred by three different models (cf. Section 1.2).

6

the analysis [Waldorp et al., 2011]. In general, encoding and decoding models are descriptive models that are fitted to explain observed data and are not designed to answer aforementioned questions. The cases in which they provide reliable information about effective relationships are somewhat coincidental and scarce [Weichwald et al., 2015].

We argue that, if the aim is to answer these questions, we need to advance beyond the encoding versus decoding dichotomy and consider and develop methods specifically tailored to investigate cause-effect relationships (for examples refer to [Chén et al., 2016; Grosse-Wentrup et al., 2016; Weichwald et al., 2016a,b]). We shall not let ourselves be put off with the slogan "correlation is not causation" and instead tackle and openly discuss the subtle problems in answering the core neuroimaing questions. In particular, there has been remarkable progress in very carefully and rigorously researching the required assumptions and theoretical underpinnings of causal inference from experimental data (e. g. [Hoyer et al., 2009; Meinshausen et al., 2016; Mooij et al., 2016; Peters et al., 2016; Schölkopf et al., 2016]). Future research should focus on the right tools for the right questions.

# References

D. R. Bach, M. Symmonds, G. Barnes, and R. J. Dolan. Whole-brain neural dynamics of probabilistic reward prediction. *Journal of Neuroscience*, 2017. ISSN 0270-6474. doi: 10.1523/JNEUROSCI. 2943-16.2017.

O. Y. Chén, C. M. Crainiceanu, E. L. Ogburn, B. S. Caffo, T. D. Wager, and M. A. Lindquist. High-dimensional Multivariate Mediation: with Application to Neuroimaging Data. *arXiv preprint arXiv:1511.09354v2*, 2016.

T. Davis, K. F. LaRocque, J. A. Mumford, K. A. Norman, A. D. Wagner, and R. A. Poldrack. What do differences between multi-voxel and univariate analysis mean? How subject-, voxel-, and trial-level variance impact fMRI analysis. *NeuroImage*, 97:271–283, 2014.

K. J. Friston, A. P. Holmes, K. J. Worsley, J.-B. Poline, C. D. Frith, and R. S. Frackowiak. Statistical parametric maps in functional imaging: A general linear approach. *Human Brain Mapping*, 2(4): 189–210, 1994.

M. Grosse-Wentrup, D. Janzing, M. Siegel, and B. Schölkopf. Identification of causal relations in neuroimaging data with latent confounders: An instrumental variable approach. *NeuroImage*, 125:825–833, 2016.

S. Haufe, F. Meinecke, K. Görgen, S. Dähne, J.-D. Haynes, B. Blankertz, and F. Bießmann. On the interpretation of weight vectors of linear models in multivariate neuroimaging. *NeuroImage*, 87: 96–110, 2014.

P. O. Hoyer, D. Janzing, J. M. Mooij, J. Peters, and B. Schölkopf. Nonlinear causal discovery with additive noise models. In *Advances in neural information processing systems*, pages 689–696, 2009.

A. G. Huth, T. Lee, S. Nishimoto, N. Y. Bilenko, A. Vu, and J. L. Gallant. Decoding the semantic content of natural movies from human brain activity. *Frontiers in Systems Neuroscience*, 10: 81, 2016.

N. Meinshausen, A. Hauser, J. M. Mooij, J. Peters, P. Versteeg, and P. Bühlmann. Methods for causal inference from gene perturbation experiments and validation. *Proceedings of the National Academy of Sciences*, 113(27):7361–7368, 2016.

T. M. Mitchell, R. Hutchinson, R. S. Niculescu, F. Pereira, X. Wang, M. Just, and S. Newman. Learning to decode cognitive states from brain images. *Machine Learning*, 57(1-2):145–175, 2004.

J. M. Mooij, J. Peters, D. Janzing, J. Zscheischler, and B. Schölkopf. Distinguishing cause from effect using observational data: methods and benchmarks. *Journal of Machine Learning Research*, 17(32):1–102, 2016.

T. Naselaris, K. N. Kay, S. Nishimoto, and J. L. Gallant. Encoding and decoding in fMRI. *NeuroImage*, 56(2):400–410, 2011.

J. Pearl. *Causality*. Cambridge University Press, 2009.

F. Pereira, T. Mitchell, and M. Botvinick. Machine learning classifiers and fMRI: a tutorial overview. *NeuroImage*, 45(1): S199–S209, 2009.

J. Peters, P. Bühlmann, and N. Meinshausen. Causal inference by using invariant prediction: identification and confidence intervals. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 78(5):947–1012, 2016.

B. Schölkopf, D. W. Hogg, D. Wang, D. Foreman-Mackey, D. Janzing, C.-J. Simon-Gabriel, and J. Peters. Modeling confounding by half-sibling regression. *Proceedings of the National Academy of Sciences*, 113(27):7391–7398, 2016.

P. Spirtes, C. N. Glymour, and R. Scheines. *Causation, Prediction, and Search*. MIT Press, 2000.

J. Taylor and R. J. Tibshirani. Statistical learning and selective inference. *Proceedings of the National Academy of Sciences*, 112(25):7629–7634, 2015.

M. T. Todd, L. E. Nystrom, and J. D. Cohen. Confounds in multivariate pattern analysis: Theory and rule representation case study. *NeuroImage*, 77:157–165, 2013.

L. Waldorp, I. Christoffels, and V. van de Ven. Effective connectivity of fMRI data using ancestral graph theory: Dealing with missing regions. *NeuroImage*, 54(4): 2695–2705, 2011.

S. Weichwald, B. Schölkopf, T. Ball, and M. Grosse-Wentrup. Causal and anticausal learning in pattern recognition for neuroimaging. In *Pattern Recognition in Neuroimaging, 2014 International Workshop on*, pages 1–4. IEEE, 2014.

S. Weichwald, T. Meyer, O. Özdenizci, B. Schölkopf, T. Ball, and M. Grosse-Wentrup. Causal interpretation rules for encoding and decoding models in neuroimaging. *NeuroImage*, 110:48–59, 2015.

S. Weichwald, A. Gretton, B. Schölkopf, and M. Grosse-Wentrup. Recovery of non-linear cause-effect relationships from linearly mixed neuroimaging data. In *Pattern Recognition in Neuroimaging (PRNI), 2016 International Workshop on*, pages 1–4. IEEE, 2016a.

S. Weichwald, M. Grosse-Wentrup, and A. Gretton. MERLiN: Mixture Effect

Recovery in Linear Networks. *IEEE Journal of Selected Topics in Signal Processing*, 10(7):1254–1266, 2016b.

A. Woolgar, P. Golland, and S. Bode. Coping with confounds in multivoxel pattern analysis: what should we do about reaction time differences? A comment on Todd, Nystrom & Cohen 2013. *NeuroImage*, 98:506–512, 2014.