

Machine Learning 2: advanced

Andre F. Marquand

a.marquand@donders.ru.nl

Outline



- 1 Alternative data-driven approaches
- 2 Going Nonlinear
- 3 Understanding model predictions
- 4 Recommendations

Outline

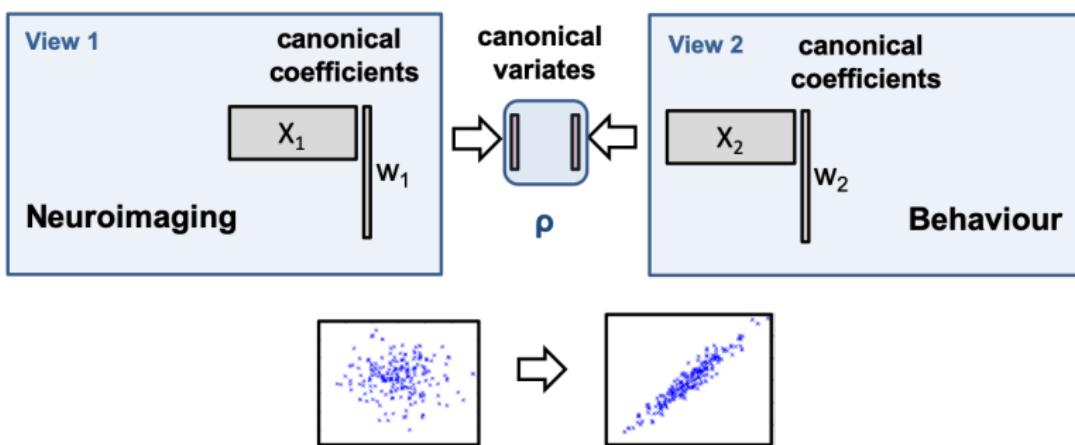


- 1 Alternative data-driven approaches
- 2 Going Nonlinear
- 3 Understanding model predictions
- 4 Recommendations

Finding mappings between brain and behaviour



- **Canonical Correlation Analysis** is a standard statistical tool for finding multivariate relationships between datasets
- Generalises Pearson correlation to multiple variables
- Finds projections of the data that maximise the correlation between “views” of the data



Canonical Correlation Analysis



- CCA is related to techniques such as partial least squares
- Formally, CCA solves the following objective function:

$$\max_{w_1, w_2} \text{corr}(X_1 w_1, X_2 w_2)$$

subject to $\|w_1^T X_1^T X_1 w_1\| \leq 1$ and $\|w_2^T X_2^T X_2 w_2\| \leq 1$

where the constraint is sometimes amended to:

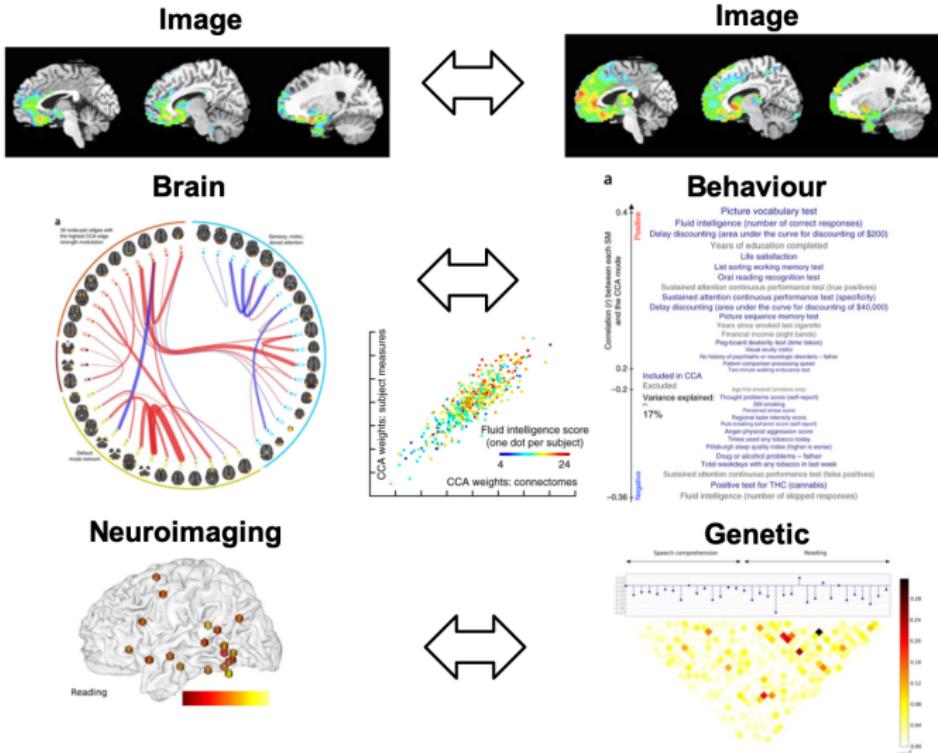
subject to $\|w_1\|^2 \leq 1$ and $\|w_2\|^2 \leq 1$

...and other constraints can be added (e.g. to promote sparsity)

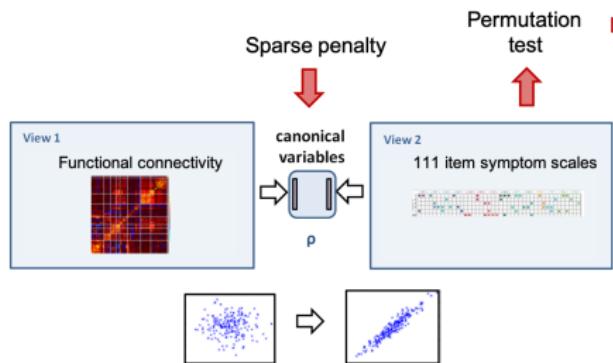
$$P(w_1) < c_1 \text{ and } P(w_2) < c_2$$

- if $n > p_1$ and p_2 , an analytical solution is available
- There are many variants (kernel CCA, Bayesian CCA, deep CCA...)

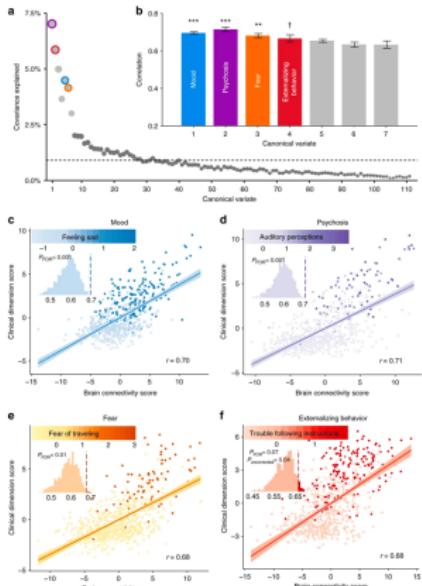
Applications of CCA



Applications of CCA

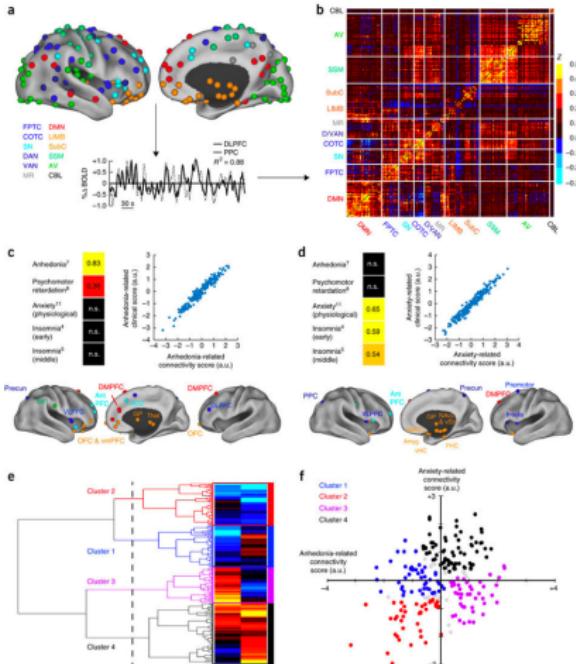


4 Linked dimensions

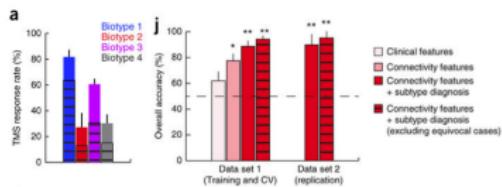


Xia et al. (2018)

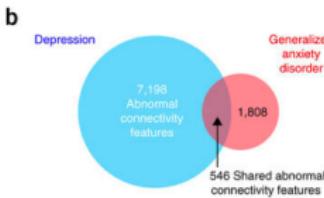
Stratification of major depression



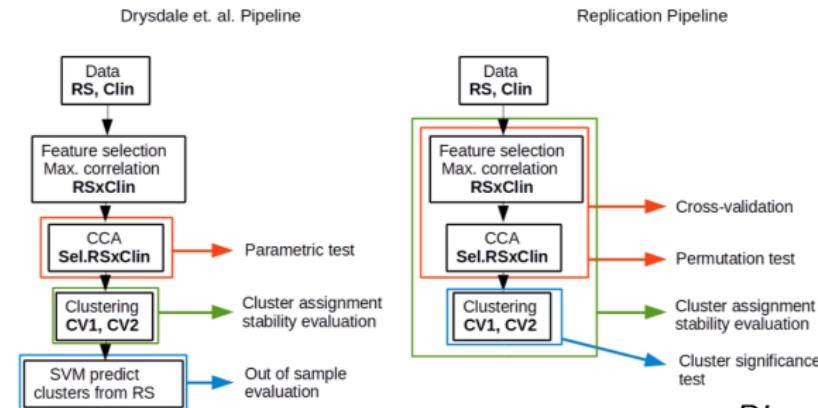
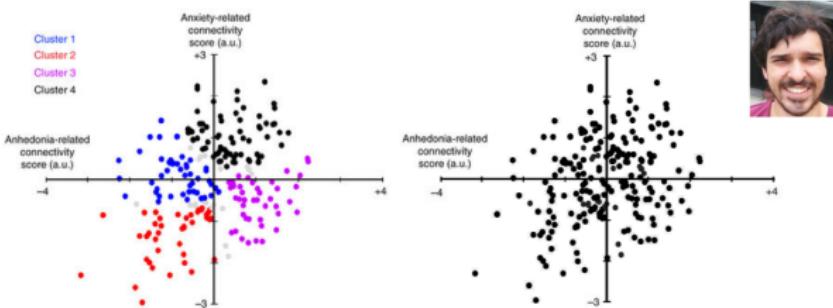
- Extensive validation
 - Predict treatment response (TMS)



- Cut across diagnoses

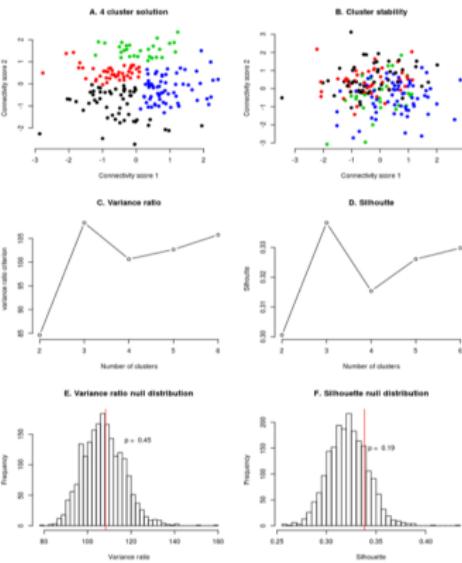
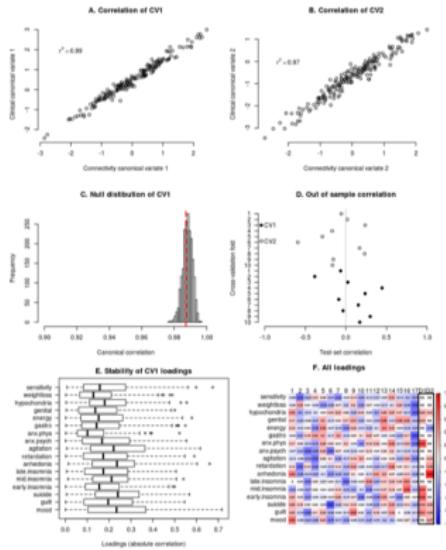


Stratification of major depression



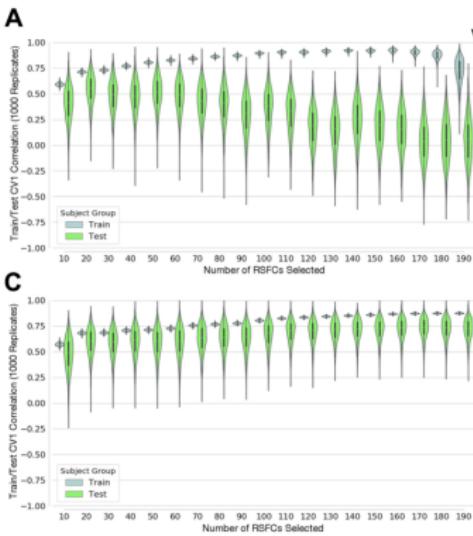
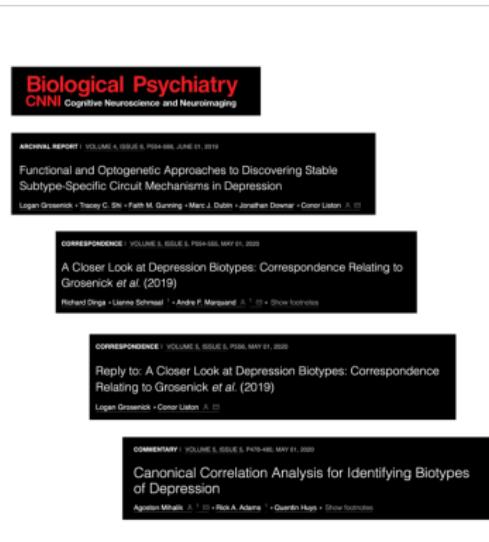
Dinga et al. (2019)

Stratification of major depression



Dinga et al. (2019)

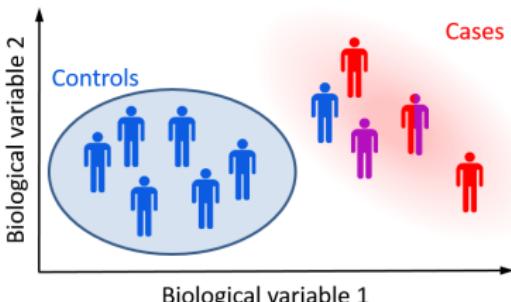
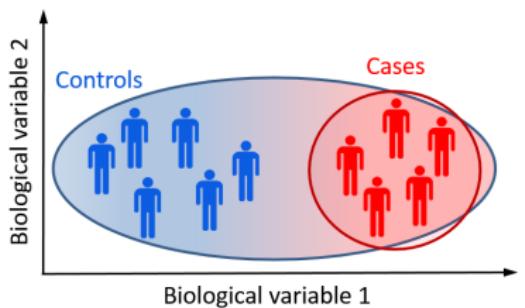
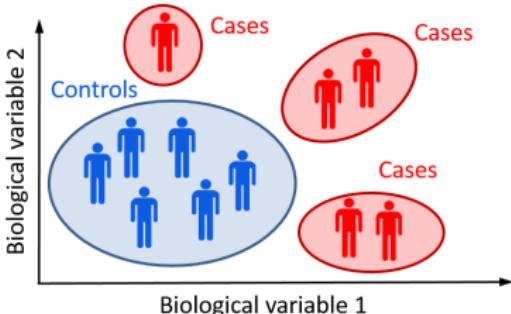
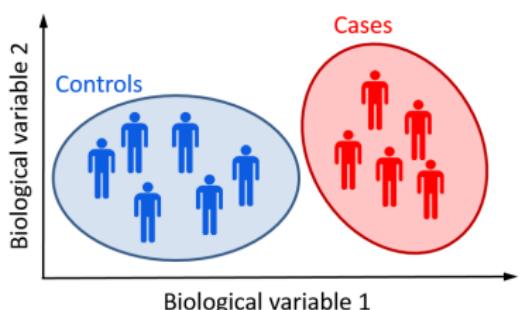
Applications of CCA



- CCA easily overfits, even when $N > p$
- Magnitude of in-sample canonical correlation is meaningless!
- Regularization and/or feature selection is very important
- Statistical evaluation should include the whole pipeline

Dinga et al. (2020); Grosenick and Liston (2020); Mihalik et al. (2020)

Many types of heterogeneity



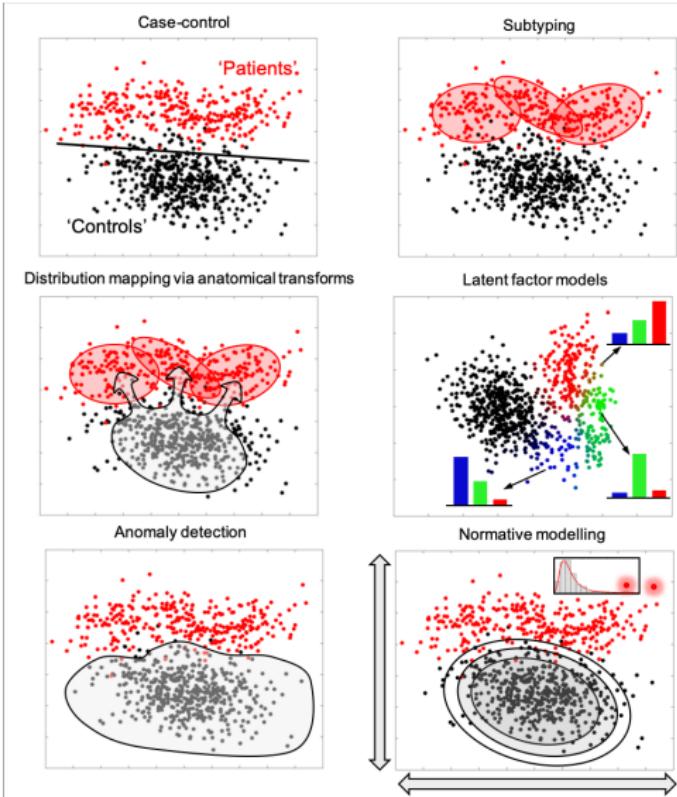
Many types of heterogeneity



Nature Reviews | Genetics

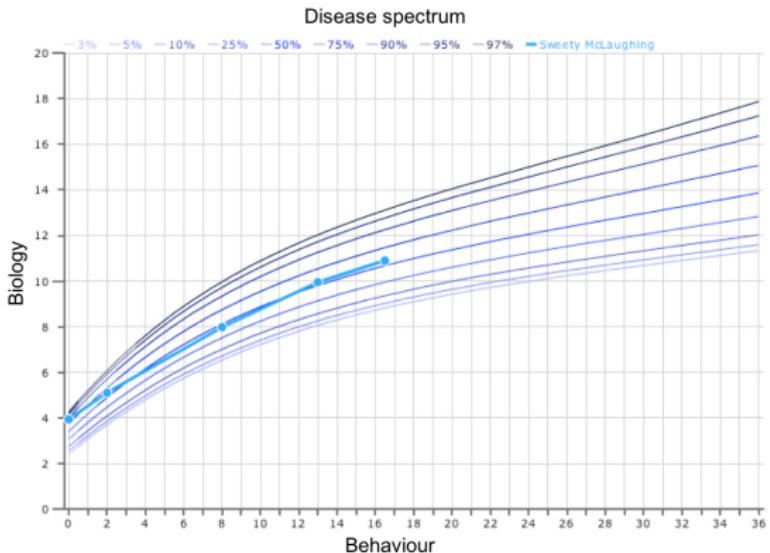
Burmeister et al. (2008)

Methods for addressing heterogeneity



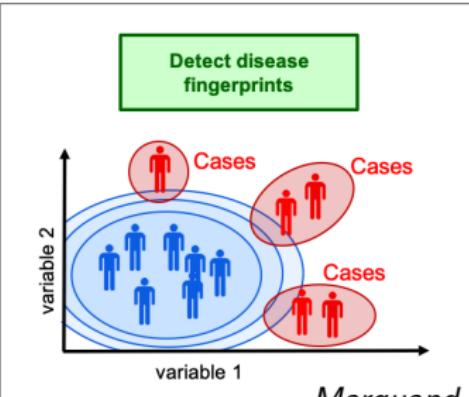
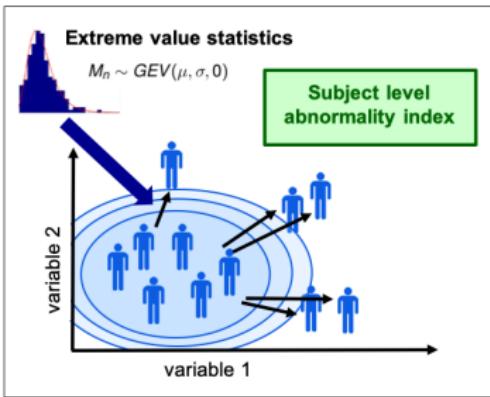
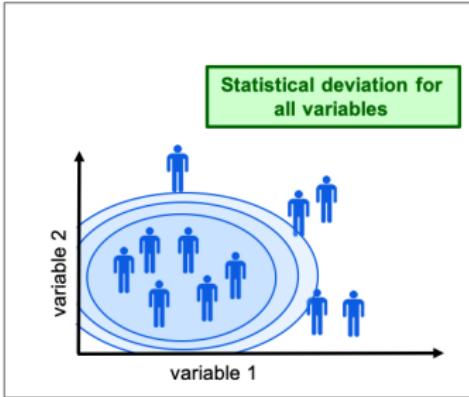
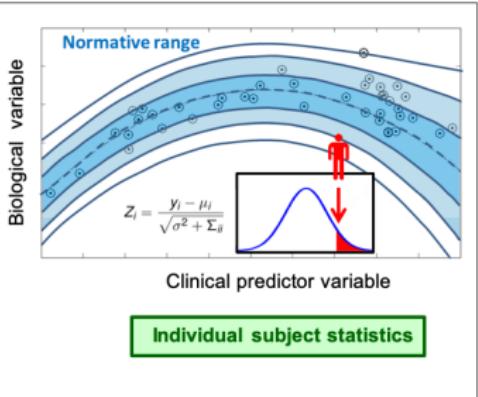
Dong et al. (2016); Zhang et al. (2016); Mourao-Miranda et al. (2011)

Normative modelling



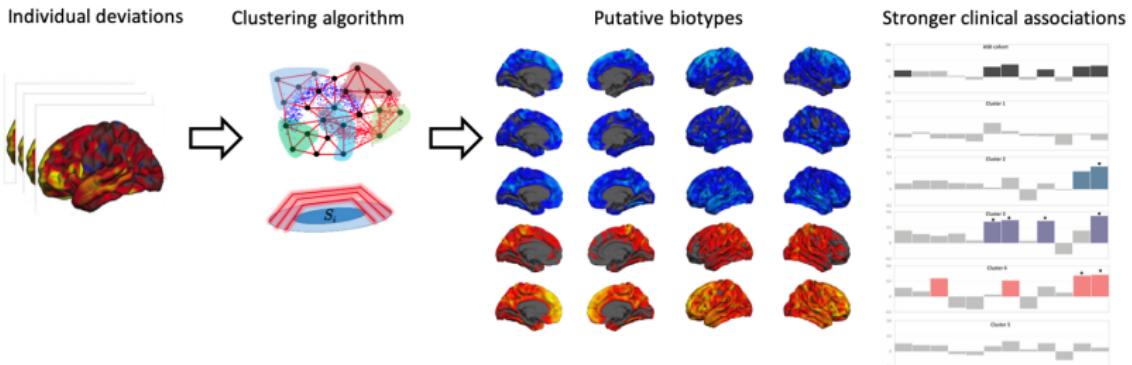
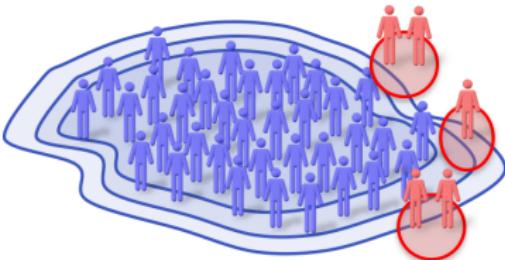
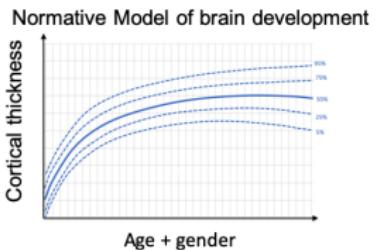
Marquand et al. (2016)

Normative modelling



Marquand et al. (2016)

Normative modelling of autism



Tutorial: https://github.com/saigerutherford/CPC_2020

Zabihí et al. (2019)

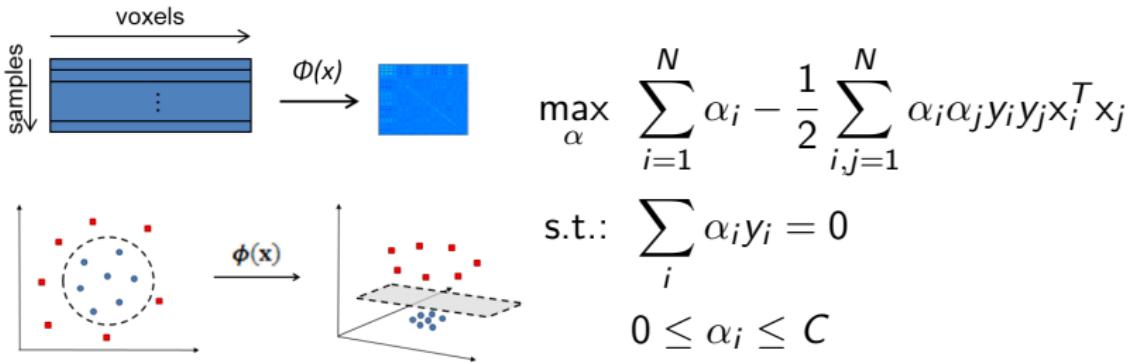
Outline



- 1 Alternative data-driven approaches
- 2 Going Nonlinear
- 3 Understanding model predictions
- 4 Recommendations



- Kernel methods (e.g. SVM, GPs) use the “kernel trick” to turn a linear model into a non-linear one



- In the dual form, the data appear as an inner product, which can be substituted with a kernel function

$$x_i^T x_j \Rightarrow k(x_i, x_j) \Rightarrow \phi(x_i)^T \phi(x_j)$$



- Any method that can be written in terms of inner products can be kernelized
- Kernel function must give rise to a positive definite matrix
- Many different functions are admissible

$$\text{linear : } k(x, x') = x^T x'$$

$$\text{RBF : } k(x, x') = \exp(-\gamma ||x - x'||^2)$$

$$\text{polynomial : } k(x, x') = (\gamma x^T x' + c)^d$$

$$\text{sigmoid : } k(x, x') = \tanh(\gamma x^T x' + c)$$

- linear operations on kernels also yield valid kernels, e.g.

$$k(x, x') = k_1(x, x') + 2k_2(x, x')k_3(x, x') + \dots$$

- This is the basis for *multi-kernel learning*

Multi-Kernel Learning



- Kernels can represent different modalities, different views or different regions

$$k(\mathbf{x}, \mathbf{x}') = \sum_{m=1}^M d_m k_m(\mathbf{x}, \mathbf{x}') \text{ with } d_m \geq 0 \text{ and } \sum_m d_m = 1$$

- Optimisation problem is a convex combination of kernels

$$\min_w \frac{1}{2} \sum_{m=1}^M ||\mathbf{w}_m||^2 + C \sum_{i=1}^n \xi_i$$

$$\text{subject to: } y_i \left(\sum_m (\mathbf{x}_i^T \mathbf{w}_m + b) \right) \geq 1 - \xi_i, \xi_i \geq 0 \quad \forall i$$

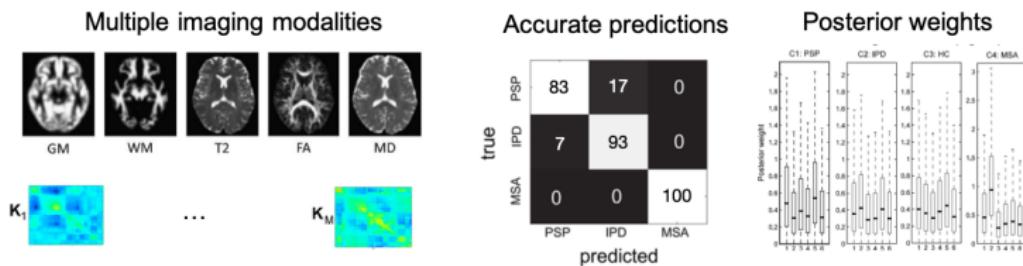
$$\sum_m d_m = 1 \text{ and } d_m \geq 0$$

Applications of MKL in neuroscience



Multimodal data fusion for predicting brain disorders

$$p(y_{ic}|x_1 \dots x_M) = \frac{\exp f_{ic}}{\sum_d \exp f_{id}} \quad p(y|x_1 \dots x_M) = \mathcal{GP}(0, K(\theta))$$



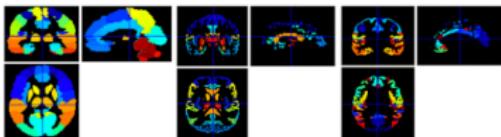
Filippone et al. (2012)

Applications of MKL in neuroscience

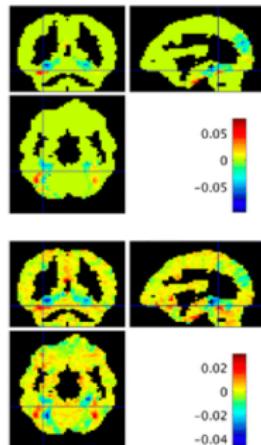


Better accommodating spatial variation

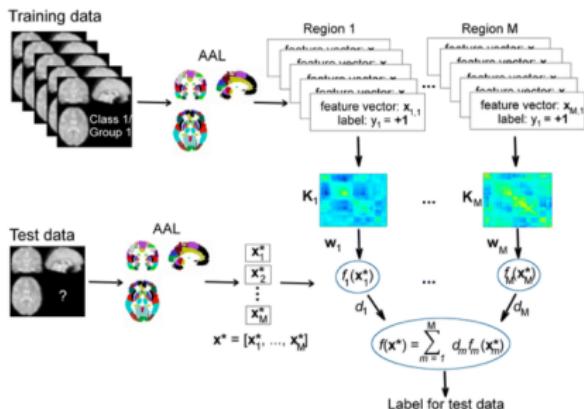
Atlas parcellation



Improve SNR



Estimate one kernel per region

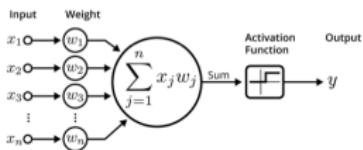


Deep Learning

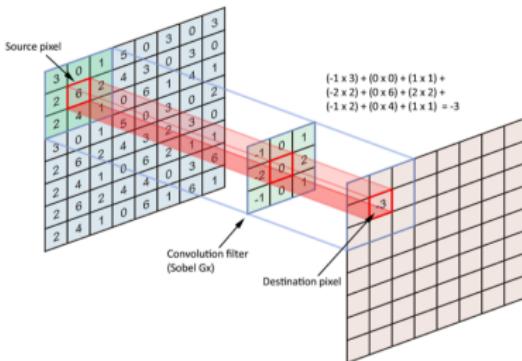
It is helpful to think of deep learning as combining matrix products with point-wise linearity, e.g.:

$$f(x) = \sigma(xW_1 + b_1)W_2$$

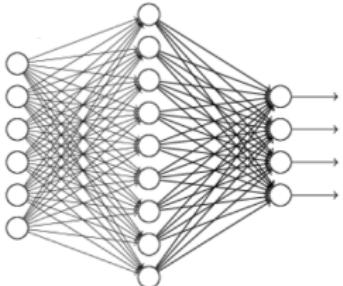
Artificial neuron
(pointwise non-linearity)



Convolution



Fully connected



Max pooling

1	1	2	4
5	6	7	8
3	2	1	0
1	2	3	4

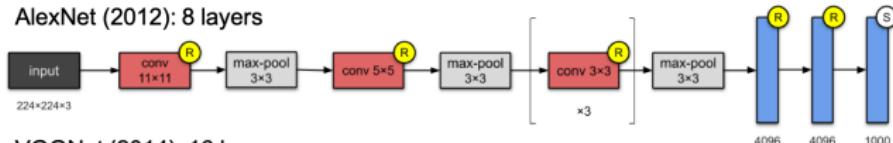
max pool with 2x2 filters and stride 2

6	8
3	4

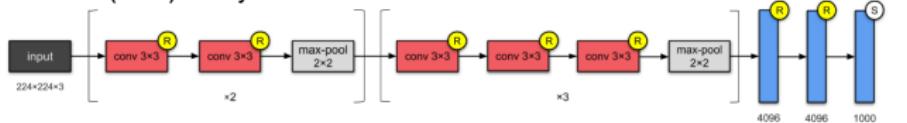
Convolutional Neural Networks



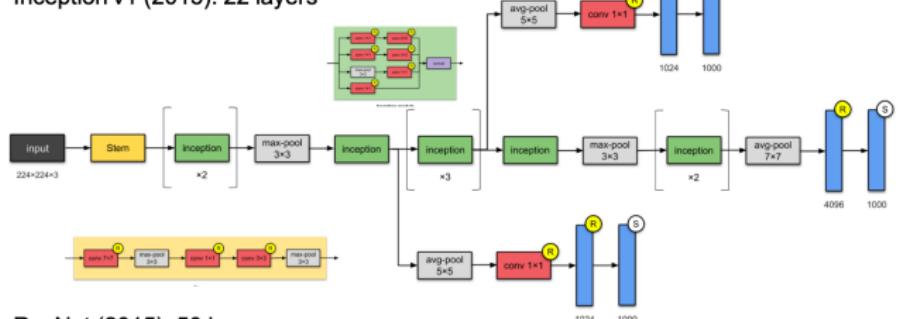
AlexNet (2012): 8 layers



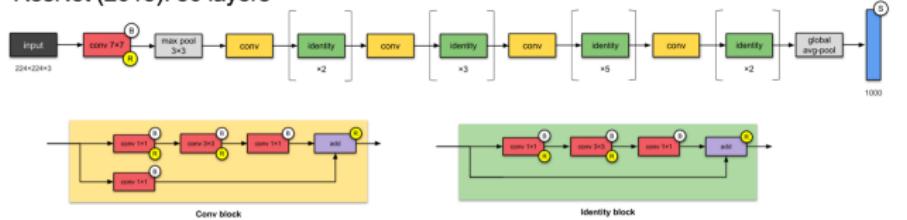
VGGNet (2014): 16 layers



Inception v1 (2015): 22 layers



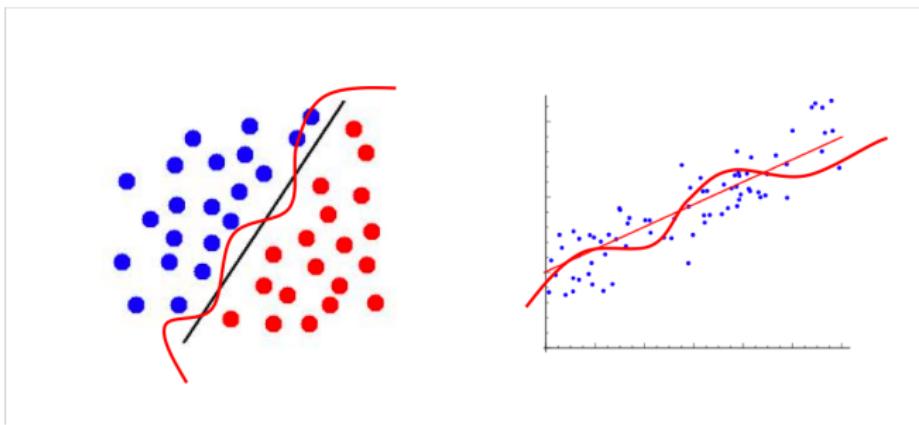
ResNet (2015): 50 layers



Overfitting (again)



- But if your problem is linear, your fancy nonlinear algorithm will just overfit

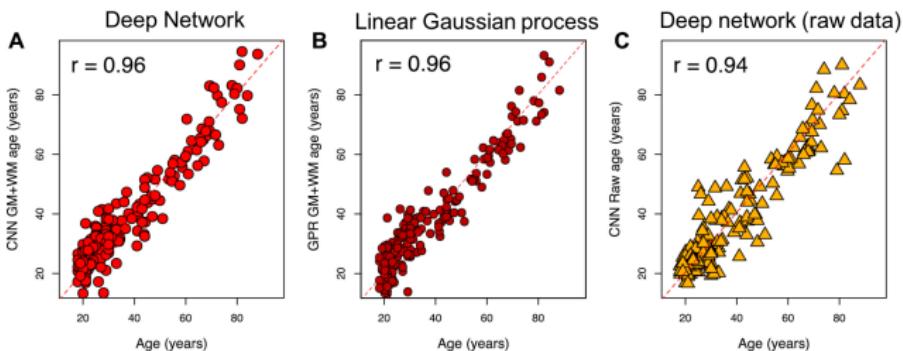
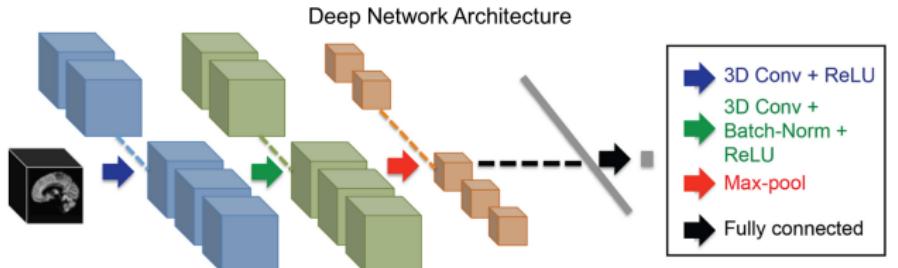


- The more complex the model, the easier it is to overfit
- In complex (deep) models it is often not possible to properly optimise all parameters
- This makes validation extremely important!

Deep Learning in Neuroscience?



- Predict age from $N = 2001$ structural MRI images

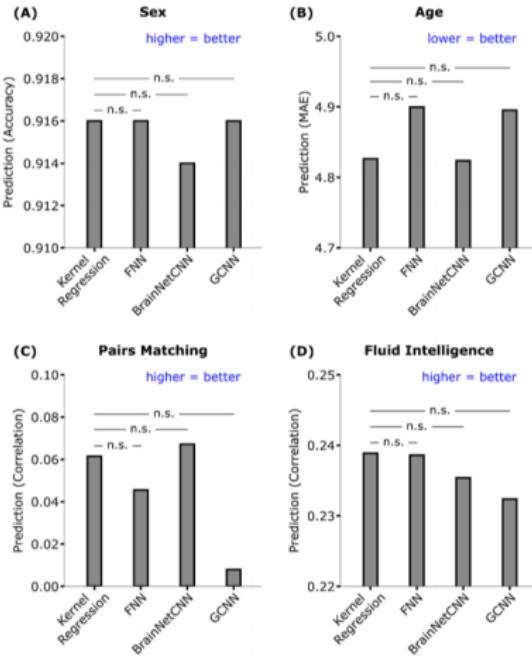


- Similar performance to a linear model on preprocessed data
- Better performance on minimally processed data *Cole et al. (2017)*

Deep Learning in Neuroscience?



biobank^{uk}



Deep Learning in Neuroscience?



New Results

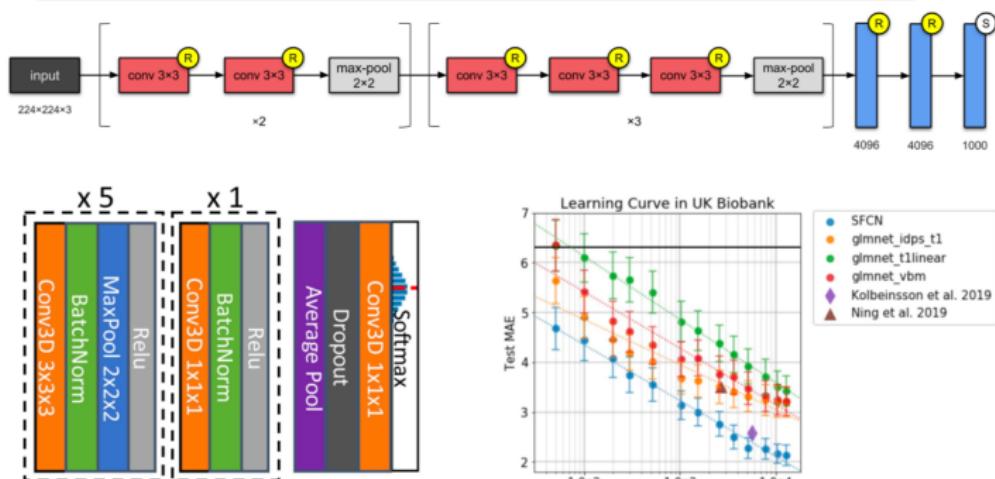
[Comment on this paper](#)

Accurate brain age prediction with lightweight deep neural networks

Han Peng, Weikang Gong, Christian F. Beckmann, Andrea Vedaldi, Stephen M. Smith

doi: <https://doi.org/10.1101/2019.12.17.879346>

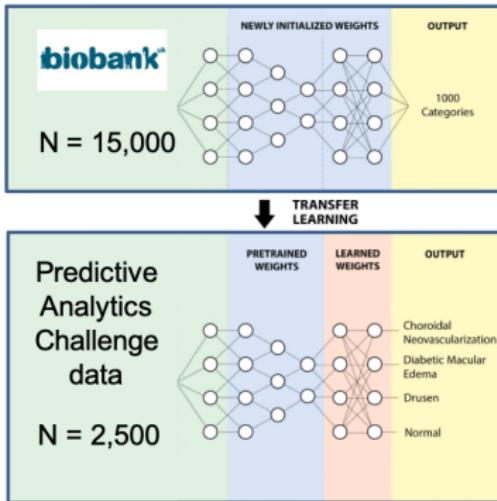
This article is a preprint and has not been certified by peer review [what does this mean?].



Won the 2019 PAC 'brainage' challenge

Peng et al. (2019)

Transfer Learning



- Basic idea: transfer knowledge (i.e. weights) from a large dataset to a small one (where it is harder to learn)
- Different variants depending on whether the targets are the same, similar or different

Deep Learning in Medicine?



- Predict mortality from electronic health records

ARTICLE

OPEN

Scalable and accurate deep learning with electronic health records

Alvin Rajkomar^{1,2}, Eyal Oren¹, Kai Chen¹, Andrew M. Dai¹, Nissan Hajaj³, Michaela Hardt¹, Peter J. Liu¹, Xiaobing Liu¹, Jake Marcus¹, Mimi Sun¹, Patrik Sundberg¹, Hector Yee¹, Kun Zhang¹, Yi Zhang¹, Gerardo Flores¹, Gavin E. Duggan¹, Jamie Irvine¹, Quoc Le¹, Kurt Litsch¹, Alexander Mossin¹, Justin Tansuwan¹, De Wang¹, James Wexler¹, Jimbo Wilson¹, Dana Ludwig², Samuel L. Volchenboum³, Katherine Chou¹, Michael Pearson¹, Srinivasan Madabushi¹, Nigam H. Shah⁴, Atul J. Butte², Michael D. Howell¹, Claire Cui¹, Greg S. Corrado¹ and Jeffrey Dean¹

	Hospital A	Hospital B
Inpatient Mortality, AUROC (95% CI)		
Deep learning 24 hours after admission	0.95 (0.94-0.96)	0.93 (0.92-0.94)
Full feature enhanced baseline at 24 hours after admission	0.93(0.92-0.95)	0.91(0.89-0.92)
Full feature simple baseline at 24 hours after admission	0.93(0.91-0.94)	0.90(0.88-0.92)
Baseline (aEWS ²) at 24 hours after admission	0.85(0.81-0.89)	0.86(0.83-0.88)
30-day Readmission, AUROC (95% CI)		
Deep learning at discharge	0.77 (0.75-0.78)	0.76 (0.75-0.77)
Full feature enhanced baseline at discharge	0.75(0.73-0.76)	0.75(0.74-0.76)
Full feature simple baseline at discharge	0.74(0.73-0.76)	0.73(0.72-0.74)
Baseline (mHOSPITAL ³) at discharge	0.70(0.68-0.72)	0.68(0.67-0.69)
Length of Stay at least 7 days AUROC (95% CI)		
Deep learning 24 hours after admission	0.86 (0.86-0.87)	0.85 (0.85-0.86)
Full feature enhanced baseline at 24 hours after admission	0.85(0.84-0.85)	0.83(0.83-0.84)
Full feature simple baseline at 24 hours after admission	0.83(0.82-0.84)	0.81(0.80-0.82)
Baseline (mLiu ⁴) at 24 hours after admission	0.76(0.75-0.77)	0.74(0.73-0.75)

Outline



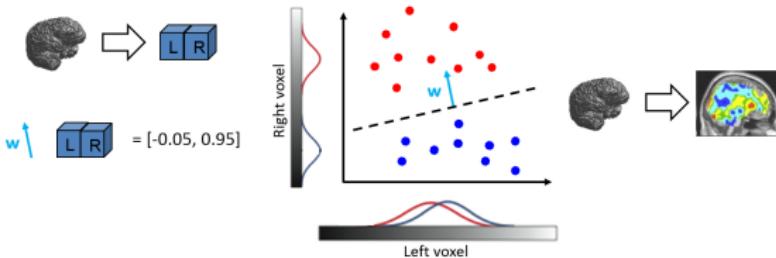
- 1 Alternative data-driven approaches
- 2 Going Nonlinear
- 3 Understanding model predictions
- 4 Recommendations

Mapping the discriminative pattern

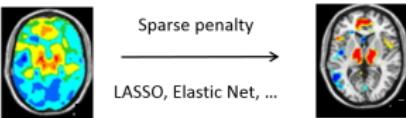


For clinical applications it is crucial to infer which variables drive the predictions. There are multiple options:

- Regional classification accuracy as a proxy (searchlight)
- For linear models, the weights can be directly visualised (“discriminative mapping”)



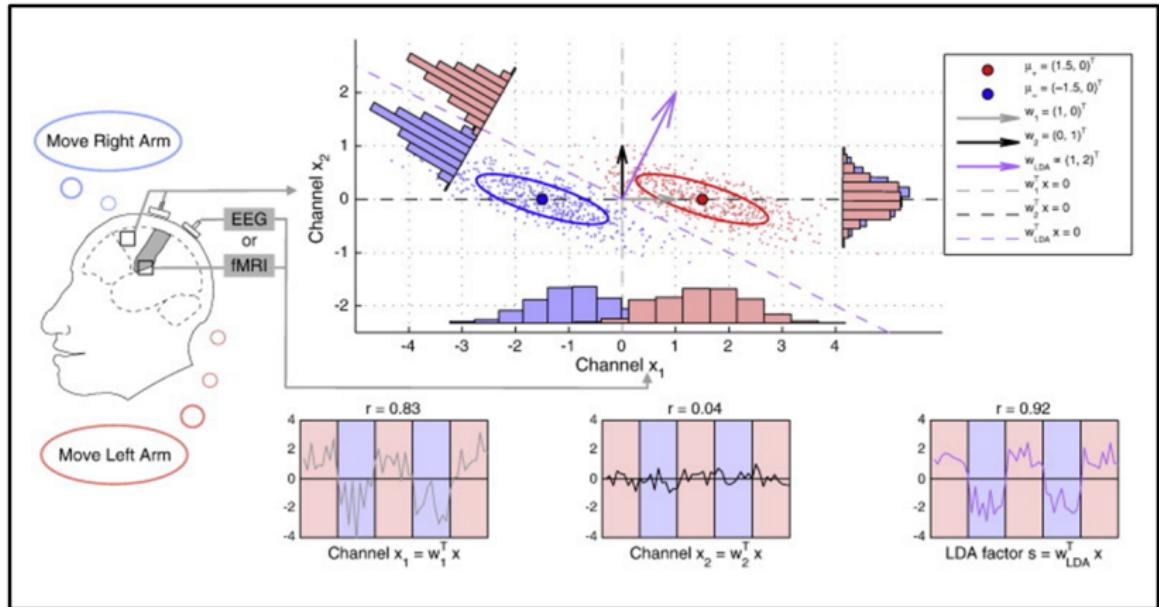
- Can also use regularization to enforce sparse weights



- Weights are often regarded as being difficult to interpret, but this is not always true

Weights do not reflect univariate differences

One proposal is to consider the weights from a forward model



Construct 'forward maps' by premultiplying by the data covariance

$$a = \frac{1}{\sigma_y^2} \sum_x w$$

Haufe et al. (2014)

Understanding weights of discriminative models

- The correct interpretation of the weights is the contribution of each feature to the predictions. This is the same as in a GLM
- Difficulty arises only due to multicollinearity between predictor variables which inflates the variance of the weights
- A variable can have a high weight because:
 - ① It is associated with the response variable
 - ② It acts as a 'suppressor' variable that helps to cancel out noise or mismatch in other covariates
- To distinguish between these possibilities, we can do the following (assumes standardized data)

$$a \propto \sum_x w = \text{cov}[X, \hat{y}] = \text{corr}[X, \hat{y}]$$

- These are **structure coefficients** from multivariate statistics
- The univariate association between covariate p and the predictions is given simply by:

$$\rho(x_p, \hat{y})$$

Kraha et al. (2012)

But what about the penalty?



- Collinearity is well-known in classical GLM settings, where advice is often given to avoid collinearity
- It is true that collinearity impacts on efficiency, but models with collinear predictors are still interpretable
- Collinearity also impacts penalised regression. Recall that:

$$f(x_i, w) = x_i^T w \quad \Rightarrow \hat{w} = \min_w \sum_{i=1}^n \ell(y_i, f_i) + \lambda J(w)$$

- Considering ridge regression, where the objective function is:

$$\hat{w} = \min_w \sum_{i=1}^n \mathcal{N}(w^T x_i, \sigma^2) + \frac{\lambda}{2} \|w\|_2^2$$

- This is equivalent to maximising the following

$$\min_w -\frac{1}{2\sigma^2} \left(\sum_{i=1}^n (y_i - x_i^T w)^2 + \lambda w^T w \right)$$

But what about the penalty?



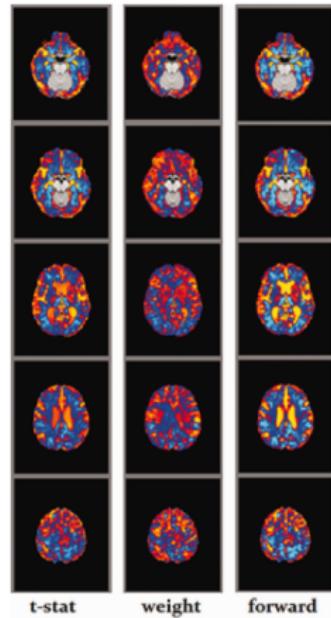
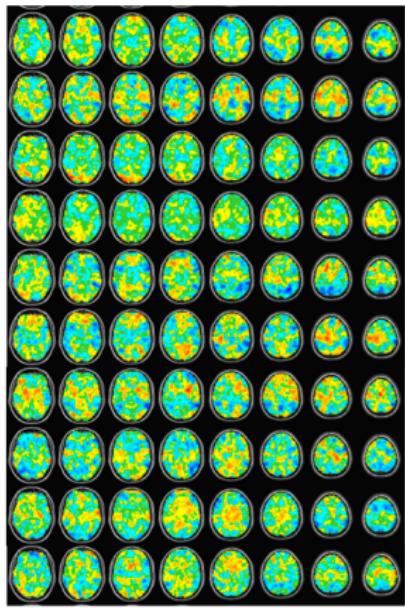
- This is exactly equivalent to finding the MAP estimate of a posterior distribution over w , with prior:

$$\mathcal{N}(0, \sigma^2 / \lambda I)$$

So what does this mean?

- The magnitude and sign of weights are influenced by collinearity
- Including and excluding variables can change the magnitude and sign of variables in the model
- It is obvious that when $p > n$, the problem is ill-posed in that there are many ways the same prediction can be achieved.
- Regularisation helps to stabilise coefficients, but this does not eliminate the problem

Examples of weights

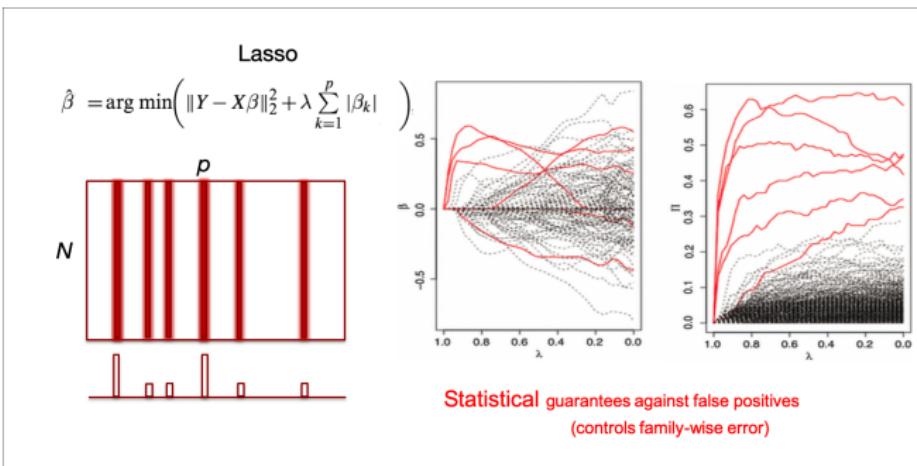


Aksman et al. (2016)

Stability selection



- It is important to recognise that coefficient stability is largely independent from accuracy, especially if features are correlated
- **stability selection** is a method that aims to identify truly important variables by examining stability under perturbation to the data

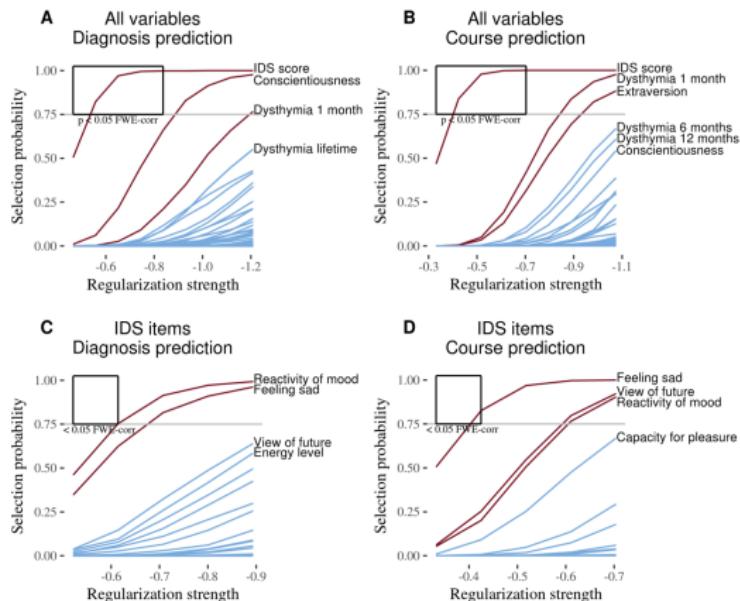
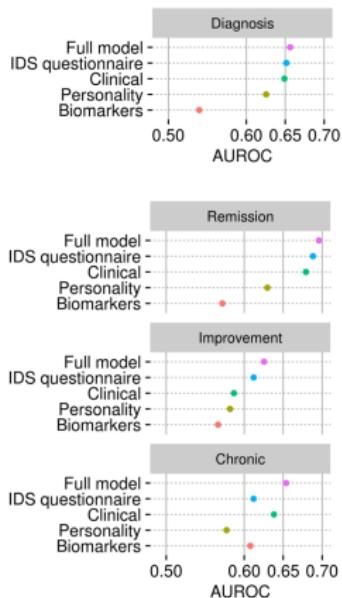


Meinshausen and Bühlmann (2010)

Stability selection in major depression



Aim: to identify prognostic markers for depression from an extensive panel (cognition, symptoms, biomarkers...)



Dinga et al. (2018)

Outline



- 1 Alternative data-driven approaches
- 2 Going Nonlinear
- 3 Understanding model predictions
- 4 Recommendations



- Machine learning provides powerful tools for single subject inference and detect spatially distributed effects
- Many different approaches beyond simple notions such as 'classification' or 'clustering'

Recommendations

- Linear models are often sufficient: they are fast, interpretable and often perform as well as non-linear methods
- Careful validation is extremely important for all methods to guard against overfitting
- Machine learning can be easily integrated with neurocognitive models (e.g. to assess candidate models)

Acknowledgments



Mariam
Zabihi



Richard
Dinga



Thomas
Wolfers



Saige
Rutherford



Predictive
Clinical
Neuroscience Lab

We are hiring!

Tutorial on Normative modelling:

https://github.com/saigerutherford/CPC_2020

References I



- Leon Aksman, David J. Lythgoe, Steven C.R. Williams, Martha Jokisch, Christoph Moenninghoff, Johannes Streffer, Karl Heinz Joeckel, Christian Weimar, and Andre F. Marquand. Making use of longitudinal information in pattern recognition. *Human Brain Mapping*, 37(12):4385–4404, 2016. doi: 10.1002/hbm.23317. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/hbm.23317>.
- James Cole, Rudra Poudel, Dimosthenis Tsagkasoullis, Matthan Caan, Claire Steves, Tim Spector, and Giovanni Montana. Predicting brain age with deep learning from raw imaging data results in a reliable and heritable biomarker. *Neuroimage*, (In Press), 2017.
- Richard Dinga, Andre F. Marquand, Dick J. Veltman, Aartjan T. F. Beekman, Robert A. Schoevers, Albert M. van Hemert, Brenda W. J. H. Penninx, and Lianne Schmaal. Predicting the naturalistic course of depression from a wide range of clinical, psychological, and biological data: a machine learning approach. 2018.
- Richard Dinga, Lianne Schmaal, Brenda W.J.H. Penninx, Marie Jose van Tol, Dick J. Veltman, Laura van Velzen, Maarten Mennes, Nic J.A. van der Wee, and Andre F. Marquand. Evaluating the evidence for biotypes of depression: Methodological replication and extension of drysdale et al. (2017). *NeuroImage: Clinical*, 22: 101796, 2019. ISSN 2213-1582. doi: <https://doi.org/10.1016/j.nicl.2019.101796>. URL <http://www.sciencedirect.com/science/article/pii/S2213158219301469>.
- Richard Dinga, Lianne Schmaal, and Andre F. Marquand. A closer look at depression biotypes: Correspondence relating to grosenick et al. (2019). *Biological Psychiatry: Cognitive Neuroscience and Neuroimaging*, 5(5):554 – 555, 2020. ISSN 2451-9022. doi: <https://doi.org/10.1016/j.bpsc.2019.09.011>. URL <http://www.sciencedirect.com/science/article/pii/S245190221930299X>.
- A. Dong, N. Honnorat, B. Gaonkar, and C. Davatzikos. Chimera: Clustering of heterogeneous disease effects via distribution matching of imaging patterns. *IEEE Transactions on Medical Imaging*, 35(2):612–621, Feb 2016. ISSN 0278-0062. doi: 10.1109/TMI.2015.2487423.
- Andrew T Drysdale, Logan Grosenick, Jonathan Downar, Katharine Dunlop, Farrokh Mansouri, Yue Meng, Robert N Fetcho, Benjamin Zebley, Desmond J Oathes, Amit Etkin, Alan F Schatzberg, Keith Sudheimer, Jennifer Keller, Helen S Mayberg, Faith M Gunning, George S Alexopoulos, Michael D Fox, Alvaro Pascual-Leone, Henning U Voss, BJ Casey, Marc J Dubin, and Conor Liston. Resting-state connectivity biomarkers define neurophysiological subtypes of depression. *Nature Medicine*, (23):28–38, 2017.
- Maurizio Filippone, Andre Marquand, Camilla Blain, Steven Williams, Janaina Mourao-Miranda, and Mark Girolami. Probabilistic prediction of neurological disorders with a statistical assessment of neuroimaging data modalities. *Annals of Applied Statistics*, 6:1883–1905, 2012.

References II



- Edith Le Floch, Vincent Guillemot, Vincent Frouin, Philippe Pinel, Christophe Lalanne, Laura Trinchera, Arthur Tenenhaus, Antonio Moreno, Monica Zilbovicius, Thomas Bourgeron, Stanislas Dehaene, Bertrand Thirion, Jean-Baptiste Poline, and Edouard Duchesnay. Significant correlation between a set of genetic polymorphisms and a functional brain network revealed by feature selection and sparse partial least squares. *NeuroImage*, 63 (1):11 – 24, 2012. ISSN 1053-8119.
- Logan Grosenick and Conor Liston. Reply to: A closer look at depression biotypes: Correspondence relating to grosenick et al (2019). *Biological Psychiatry: Cognitive Neuroscience and Neuroimaging*, 5(5):556, 2020. ISSN 2451-9022. doi: <https://doi.org/10.1016/j.bpsc.2019.11.002>. URL <http://www.sciencedirect.com/science/article/pii/S2451902219302988>.
- Stefan Haufe, Frank Meinecke, Kai Goergen, Sven Daehne, John-Dylan Haynes, Benjamin Blankertz, and Felix Biessmann. On the interpretation of weight vectors of linear models in multivariate neuroimaging. *NeuroImage*, 87(0):96 – 110, 2014.
- Tong He, Ru Kong, Avram Holmes, Minh Nguyen, Mert Sabuncu, Simon Eickhoff, Danilo Bzdok, Jiashi Feng, and B. T. Thomas Yeo. Deep neural networks and kernel regression achieve comparable accuracies for functional connectivity prediction of behavior and demographics. *BioRxiv*, 2019.
- Amanda Kraha, Heather Turner, Kim Nimon, Linda Zientek, and Robin Henson. Tools to support interpreting multiple regression in the face of multicollinearity. *Frontiers in Psychology*, 3:44, 2012. ISSN 1664-1078. doi: 10.3389/fpsyg.2012.00044. URL <https://www.frontiersin.org/article/10.3389/fpsyg.2012.00044>.
- Andre F. Marquand, lead Rezek, Jan Buitelaar, and Christian F. Beckmann. Understanding heterogeneity in clinical cohorts using normative models: beyond case control studies. In Press, 2016.
- Agoston Mihalik, Rick A. Adams, and Quentin Huys. Canonical correlation analysis for identifying biotypes of depression. *Biological Psychiatry: Cognitive Neuroscience and Neuroimaging*, 5(5):478 – 480, 2020. ISSN 2451-9022. doi: <https://doi.org/10.1016/j.bpsc.2020.02.002>. URL <http://www.sciencedirect.com/science/article/pii/S245190222030032X>.
- Janaina Mourao-Miranda, David R. Hardoon, Tim Hahn, Andre F. Marquand, Steve C.R. Williams, John Shawe-Taylor, and Michael Brammer. Patient classification as an outlier detection problem: An application of the one-class support vector machine. *NeuroImage*, 58(3):793 – 804, 2011. ISSN 1053-8119. doi: <https://doi.org/10.1016/j.neuroimage.2011.06.042>. URL <http://www.sciencedirect.com/science/article/pii/S1053811911006872>.

References III



- Han Peng, Weikang Gong, Christian F. Beckmann, Andrea Vedaldi, and Stephen M. Smith. Accurate brain age prediction with lightweight deep neural networks. *bioRxiv*, 2019. doi: 10.1101/2019.12.17.879346. URL <https://www.biorxiv.org/content/early/2019/12/18/2019.12.17.879346>.
- Stephen M Smith, Thomas E Nichols, Diego Vidaurre, Anderson M Winkler, Timothy E J Behrens, Matthew F Glasser, Kamil Ugurbil, Deanna M Barch, David C Van Essen, and Karla L Miller. A positive-negative mode of population covariation links brain connectivity, demographics and behavior. 18:1565–1567, 2015.
- Cedric Xia, Zongming Ma, Rastko Ceric, Shi Gu, Richard F. Betzel, Antonia N. Kaczkurkin, Monica E. Calkins, Philip A. Cook, Angel Garcia de la Garza, Simon N. Vandekar, Zaixu Cui, Tyler M. Moore, David R. Roalf, Kosha Ruparel, Daniel H. Wolf, Christos Davatzikos, Ruben C. Gur, Raquel E. Gur, Russell T. Shinohara, Danielle S. Bassett, and Theodore D. Satterthwaite. Linked dimensions of psychopathology and connectivity in functional brain networks. *Nature Communications*, 9, 2018.
- Mariam Zabihi, Marianne Oldehinkel, Thomas Wolfers, Vincent Frouin, David Goyard, Eva Loth, Tony Charman, Julian Tillmann, Tobias Banaschewski, Guillaume Dumas, Rosemary Holt, Simon Baron-Cohen, Sarah Durston, Sven Boelte, Declan Murphy, Christine Ecker, Jan K. Buitelaar, Christian F. Beckmann, and Andre F. Marquand. Dissecting the heterogeneous cortical anatomy of autism spectrum disorder using normative models. 4(6):567 – 578, 2019. ISSN 2451-9022. doi: <https://doi.org/10.1016/j.bpsc.2018.11.013>. URL <http://www.sciencedirect.com/science/article/pii/S245190221830329X>. The Bridging of Scales: Techniques for Translational Neuroscience.
- Xiuming Zhang, Elizabeth C. Mormino, Nanbo Sun, Reisa A. Sperling, Mert R. Sabuncu, and B. T. Thomas Yeo. Bayesian model reveals latent atrophy factors with dissociable cognitive trajectories in alzheimer's disease. *Proceedings of the National Academy of Sciences*, 113(42):E6535–E6544, 2016. ISSN 0027-8424. doi: 10.1073/pnas.1611073113. URL <https://www.pnas.org/content/113/42/E6535>.