

Decoding analyses

Neuroimaging: Pattern Analysis 2017



- Conceptual explanation of decoding/machine learning (ML) and related concepts;
- Not about mathematical basis of ML algorithms!



- Understand the difference (and surprising overlap) between statistics and machine learning;
- Recognize the major different "flavors" of decoding analyses;
- Understand the concepts of overfitting and cross-validation
- Know the most important steps in a decoding analysis pipeline;



- **PART 1**: what is machine learning?
 - ... and how does it differ from statistics?
 - What "flavors" of decoding analyses are there?
 - Important ML concepts: cross-validation, overfitting
- **PART 2**: how to set up a ML pipeline?
 - From partitioning to significance testing

PART 1: WHAT IS MACHINE LEARNING?



- Decoding ≈ machine learning
 - Decoding is the neuroimaging-specific application of machine learning algorithms and techniques;



- Defining machine learning is (in a way) trivial:
 - Algorithms that learn how to model the data without an explicit instruction how to do so;
- But how it that different from traditional statistics?
 - Like the familiar GLM (linear regression, t-tests)?
 - Similar, but different ...

- They have different **origins**:
 - Statistics is a subfield from mathematics
 - Machine learning is a subfield from computer science
 - "Science vs. engineering"
- They have a different goal:
 - Statistical models aim for **inference** about the population based on a limited sample
 - Machine learning models aim for accurate **prediction** of new samples

- Psychologists (you included) are taught statistics: how to make inferences about general psychological "laws" based on a limited sample:
 - "I've tested the reaction time of 40 people (sample) before and after drinking coffee ...
 - ... and based on the significant results of a paired sample t-test, t(39), p < 0.05 (statistical test) ...
 - ... I conclude that caffeine improves reaction times (statement about population)



- Crucially: we are quite certain that our findings in the sample will generalize to the population, if and only if assumptions of the model hold (and sample = truly random)
- Uses concepts like standard errors, confidence intervals, and p-values to generalize the model to the population



- ML models do not aim for inference, but aim for prediction
 - Instead of assuming findings will generalize to the population, ML analyses in fact *literally check* whether it generalizes;
 - It's like they're saying: "I don't give a shit about assumptions – if it works, it works."



- Instead of assuming that the findings from the model will generalize beyond the sample, ML tests this explicitly by applying this to ("predicting") a new sample
- New sample is concretely part of your dataset! (not like "the population")



- While having a different goal, in the end both types of models simply try to model the data;
- Take for example linear regression:
 - It has a **statistics** 'version' (as defined in the GLM) ...
 - ... and an ML 'version' (using a mathematical technique called gradient descent to find the optimal βs)

So ...?

- As said, statistics and ML are the **same**, yet different;
 - Both aim to model the data (with different techniques) ...
 - ... but have a different way to generalize findings
- "But why do we have to learn a whole new paradigm (ML), then?", you might ask ...



- Traditional statistical models do not fare well with high-dimensional problems
 - In decoding: dimensionality = amount of voxels
- Neuroimaging data likely violates many assumptions of traditional statistical models ...
- Sometimes, decoding analyses actually need prediction: e.g. predict clinical treatment outcome



"Features in the **world**"

"Features in the **brain**"



What happens in that arrow?



- Machine learning algorithms try to model the features (X) such that they approximate the target (y)
- In fMRI (decoding): can the voxel activities (X) be weighted such that they approximate the feature-of-interest (y)?





 $\hat{\mathbf{y}} = \mathbf{f}(\mathbf{X})$

β denotes the weighting parameters (or just "parameters" or "coefficients") for our features (in X)

βX denotes the (matrix) product
of the weighting parameters and
X - which means that y is
approximated as a weighted
linear combination of features



- Disclaimer: the "βX" term implies that it is a linear model (i.e. a linear weighting of features)
- Decoding analyses almost always use linear models, because most world-brain relations are probably linear (cf. Naselaris & Kay)
- Non-linear models exist, but are rarely used
 - Also because they often perform worse than linear models



$\hat{\mathbf{y}} = f(\mathbf{X})$









- ML models (f) find parameters (β) that weigh features (X) such that they optimally approximate the target (y)
- Applied to fMRI: we decode from the brain (X) to the world (y) by weighing voxels!

Questions so far?



ML algorithms – f() – exist in different 'flavors'



– Regression vs. classification

Regression

- Target (y) is continuous
- "Predict someone's weight based on someone's height"
- Example models: linear regression, ridge regression, LASSO

Classification

- Target (y) is categorical
- "Predict someone's gender based on someone's height"
- Example models: support vector machine (SVM), logistic regression, decision trees,



Regression models a continuous variable





• Classification models a categorical variable





Subjects perform a memory task in which they have to give responses. Their responses can be either correct or incorrect.

I want to analyze whether the patterns in parietal cortex are predictive of whether someone is going to respond (in)correctly.





During fMRI acquisition, subject see a set of images of varying emotional valence (from 0, very negative, to 100, very positive).

I want to decode stimulus valence from the bilateral insula.





Subjects perform an attentional blink task in the scanner (during which we measure fMRI).

I want to predict whether someone has a relatively high IQ (>100) or low IQ (<100) based upon the patterns in dorsolateral PFC during the attentional blink task.





- Both types of model try to approximate the target by 'weighting' the features (X)
 - Additionally, classification algorithms need a "squash" function to convert the outputs of βX to a categorical value
- The examples were simplistic (K = 1); usually,
 ML models operate on high-dimensional data!





K = >3? Difficult to visualize, but process is the same: weighting features such that a multidimensional plane is able to separate classes as well as possible in K-dimensional space


- We know what ML models do (find weighting parameters β to approximate y), but how do we evaluate the model?
- In other words, what is the model performance?









Accuracy = 18 / 20 = 90% [percent correct]



- Performance is often evaluated not only on the original sample, but also on a "new sample"
- This process is called **"cross-validation"**





 Model performance is often evaluated on a new ("unseen") sample: cross-validation





 Model performance is often evaluated on a new ("unseen") sample: cross-validation



Why do we want (need) to do cross-validation?



Model fit ≠ good prediction





- When your fit on your train-set is better than on your test-set, you're overfitting
- Overfitting means you're modelling noise instead of signal





- Overfitting = **modeling noise**
- Noise = random (uncorrelated from sample to sample)
- Therefore, a model based on noise will not generalize



- A small sample/feature-ratio often causes overfitting
- When there are few samples or many features, models may fit on random/accidental relationships



 When there are few samples, models may fit on random/accidental relationships













 When there are few samples, models may fit on random/accidental relationships





- Two options:
 - Gather **more data** (not always feasible)
 - **Reduce** the amount of **features** (more about this later)
 - [Regularization beyond the scope of this course!]
- Feature selection/extraction is an often-used technique in decoding analyses



- Sometimes, decoding analyses use a specific form of cross-validation to perform cross-decoding
- In cross-decoding, you aim to show "informational overlap" between two type of representations

- Cross-decoding!

• For example: suppose you have the hypothesis attractiveness drives the perception of friendliness





- ML models find weights to approximate a continuous (regression) or categorical (classification) dependent variable (y)
- Good fit ≠ good generalization ...
- ... therefore, cross-validate the model!
- Optimize the sample/feature ratio to reduce overfitting (spurious feature-DV correlations)
- Cross-decoding uses cross-validation to uncover shared representations ('informational overlap')

PART 2: BUILDING A DECODING PIPELINE



- O. Pattern extraction & preparation
- 1. Partitioning train/test
- 2. Feature selection/extraction
- 3. Model fitting (TRAIN)
- 4. Model generalization (TEST)
- 5. Statistical test of performance
- 6. Optional: plot weights



A typical decoding pipeline

O. Pattern extraction & preparation

Week

- 1. Partitioning train/test
- 2. Feature selection/extraction
- 3. Model fitting (TRAIN)
- 4. Model generalization (TEST)
- 5. Statistical test of performance
- 6. Optional: plot weights



- O. Pattern extraction & preparation
- 1. Partitioning train/test
- 2. Feature selection/extraction
- 3. Model fitting (TRAIN)
- 4. Model generalization (TEST)
- 5. Statistical test of performance
- 6. Optional: plot weights





















- 0. Pattern extraction & preparation
- 1. Partitioning train/test
- 2. Feature selection/extraction
- 3. Model fitting (TRAIN)
- 4. Model generalization (TEST)
- 5. Statistical test of performance
- 6. Optional: plot weights



- Goal: high sample/feature ratio!
- MRI: often many features (voxels), few samples (trials/instances/subjects)
- What to do???



- Reducing the dimensionality of patterns:
 - Select a subset of features (feature **selection**)
 - Transform features in lower-dimensional components (feature extraction)



- Ideas for feature selection?
 - ROI-based (e.g. only hippocampus);
 - (Independent!) functional mapper;
 - Data-driven selection: univariate feature selection



- Ideas for feature selection?
 - ROI-based (e.g. only hippocampus);
 - (Independent!) functional mapper;
 - Data-driven selection: univariate feature selection



- Use univariate difference scores (e.g. t-value/F-value) to select features
- Only select a subset of the voxels with the highest scores

Univariate feature selection

Steps:

- 1. Calculate test-statistic (t-test for difference happy/angry)
- 2. Select only the "best" 100 voxels (or a percentage)





- Importantly, data-driven feature selection needs to be cross-validated!
 - Performed on train-set only!


ToThink:

Why do you need to cross-validate your feature selection?

Isn't cross-validating the model-fit enough?

ToThink:

Why do you need to cross-validate your feature selection?





- Ideas for feature extraction?
 - PCA
 - Averaging within regions ("downsampling")





- Ideas for feature extraction?
 - PCA
 - Averaging within regions ("downsampling")
- Feature extraction (often) does not need to be cross-validated



- O. Pattern extraction & preparation
- 1. Partitioning train/test
- 2. Feature selection/extraction
- 3. Model fitting (TRAIN)
- 4. Model generalization (TEST)
- 5. Statistical test of performance
- 6. Optional: plot weights



- O. Pattern extraction & preparation
- 1. Partitioning train/test
- 2. Feature selection/extraction
- 3. Model fitting (TRAIN)
- 4. Model generalization (TEST)
- 5. Statistical test of performance
- 6. Optional: plot weights



- Remember I said ML does not use any statistical tests for inference?
 - I lied.
- Decoding analyses do use significance tests to infer whether decoding performance (R² or % accuracy) would generalize to the population



- Remember, if you want to statistically test something, you need to have a null- and alternative hypothesis:
 - H_0 : performance = chance level
 - H_a: performance > chance level



- What is chance level?
- R²?
 - Cross-validated R^2_{null} : O
- Accuracy?
 - 1 / number of classes
 - Decode negative/positive/neutral? Chance = 33%





- Observed performance is measured as the average cross-validated performance (R² / accuracy)
- Slightly different for within/between subject analyses:
 - Within: average performance **across subjects**
 - Between: average performance across folds



• Take e.g. a 3-fold cross-validation setup





- H_a: 58.5 > 50
- Subject-wise average performance estimates represent the data points
- Assuming independence between subjects, we can use simple parametric statistics
 - t-test(n-subjects 1) of observed performance (i.e. 58.5%) against the null performance (i.e. 50%)





- H_a: 54.8 > 50
- Fold-wise performance estimates represent data points
- Problem: different folds contain the same subjects







- Problem: different folds contain the same subjects
- Consequence: dependence between data points
 - Violates assumptions of many parametric statistical tests
- Solution: non-parametric (permutation) test



- Permutation tests do not assume (the shape of) a null-distribution, but "simulate" them
- To simulate the null-distribution (results expected when H_o is true), permutation tests literally simulate "performance at chance"





































```
"Non-parametric" p-value =
```

```
∑(null-scores > 52 served-scores)
```

Number of Que Omutations





- O. Pattern extraction & preparation
- 1. Partitioning train/test
- 2. Feature selection/extraction
- 3. Model fitting (TRAIN)
- 4. Model generalization (TEST)
- 5. Statistical test of performance
- 6. Optional: plot weights

If there is time left!



• Often, researchers (read: reviewers) ask:

"Which features are important for which class?"

- What they want:
- Intrinsic need for blobs?



Weight mapping

- Technically, weights (β) in *linear* models are interpretable: higher = more important
 - Negative weights: evidence for class 0
 - Positive weights: evidence for class 1



Weight mapping

- Technically, weights (β) in *linear* models are interpretable: higher = more important
 - Negative weights: evidence for class 0
 - Positive weights: evidence for class 1





Haufe et al. (2014, *NeuroImage*) showed that high weights ≠ class importance

... A widespread miscon-

ception about multivariate classifier weight vectors is that (the brain regions corresponding to) measurement channels^{*} with large weights are strongly related to the experimental condition. In fact, such conclusions can be unjustified. *voxels



- Haufe et al. (2014, *NeuroImage*) showed that high weights ≠ class importance
- Features (voxels) may function like (class-independent) "filters"
 - E.g. reflect physiological noise
- Inherent problem for decoding analyses
 (brain → world)



- Solution: mathematical trick to convert a decoding model to an encoding model
 - Weights from encoding models *are* interpretable
- Activation patterns = $(X'X)\beta$
- But: very similar to traditional activation-based analysis ...
- Conclusion: (like always) choose the analysis best suited for your question!



- **O.** Pattern extraction & preparation
- 1. Partitioning train/test
- 2. Feature selection/extraction
- 3. Model fitting (TRAIN)
- 4. Model prediction (TEST)
- 5. Statistical test of performance
- 6. Optional: plot weights

Estimate and extract patterns such that X = N-samples by N-features (voxels)



- O. Pattern extraction & preparation
- 1. Partitioning train/test
- 2. Feature selection/extraction
- 3. Model fitting (TRAIN)
- 4. Model prediction (TEST)
- 5. Statistical test of performance
- 6. Optional: plot weights

Use hold-out or K-fold partitioning



- O. Pattern extraction & preparation
- 1. Partitioning train/test
- 2. Feature selection/extraction
- 3. Model fitting (TRAIN)
- 4. Model prediction (TEST)
- 5. Statistical test of performance
- 6. Optional: plot weights




- O. Pattern extraction & preparation
- 1. Partitioning train/test
- 2. Feature selection/extraction
- 3. Model fitting (TRAIN)
- 4. Model prediction (TEST)
- 5. Statistical test of performance
- 6. Optional: plot weights





- O. Pattern extraction & preparation
- 1. Partitioning train/test
- 2. Feature selection/extraction
- 3. Model fitting (TRAIN)
- 4. Model prediction (TEST)
- 5. Statistical test of performance
- 6. Optional: plot weights





- O. Pattern extraction & preparation
- 1. Partitioning train/test
- 2. Feature selection/extraction
- 3. Model fitting (TRAIN)
- 4. Model prediction (TEST)
- 5. Statistical test of performance
- 6. Optional: plot weights





Permutation test against simulated null (between- subject)



- O. Pattern extraction & preparation
- 1. Partitioning train/test
- 2. Feature selection/extraction
- 3. Model fitting (TRAIN)
- 4. Model prediction (TEST)
- 5. Statistical test of performance
- 6. Optional: plot weights

Do not plot weights! A (complementary) univariate analysis would suffice

If you insist, plot the corresponding encoding model – $(X'X)\beta$



- Abraham et al: about the scikit-learn package for decoding analyses (read before tomorrow if possible!)
- Pereira et al: tutorial-style paper about decoding analyses