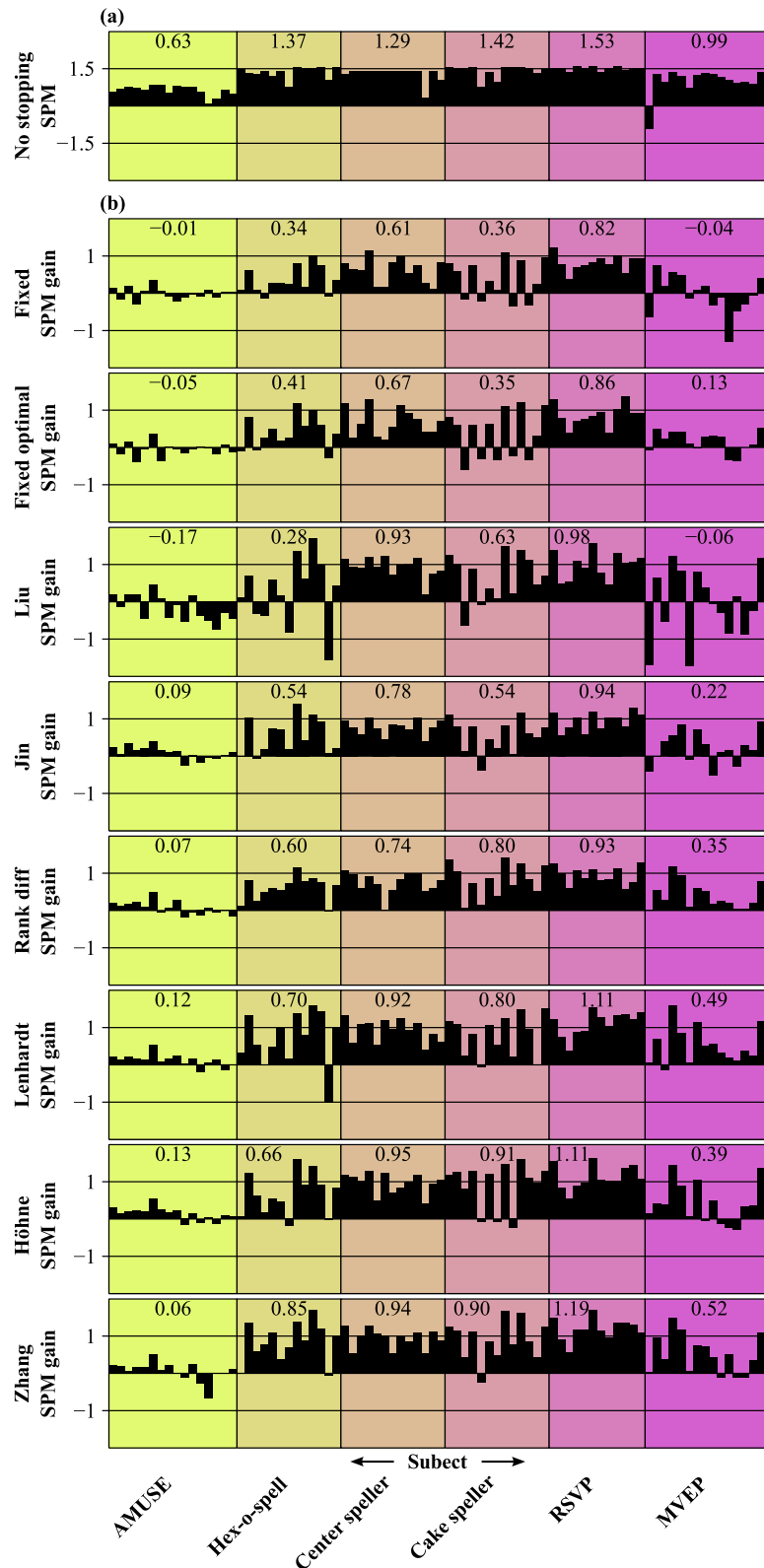
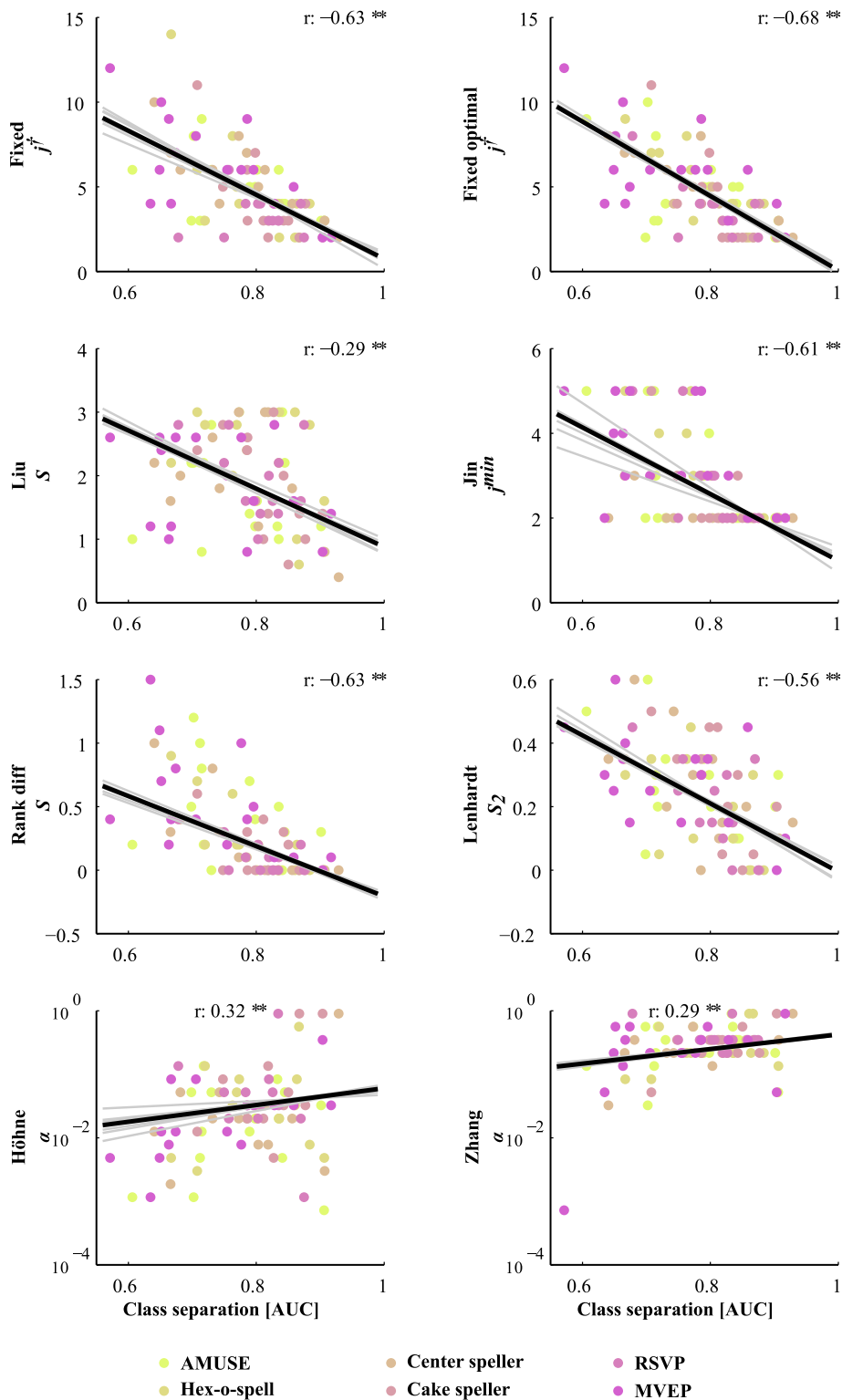


Supplementary Figure 1: Artificial datasets. A) To indicate the effect size, the baseline performance of *no stopping* is shown. B) Shown is the performance gain over the *no stopping* condition for all methods, on several artificially created datasets. For any method to be robust, it should not have negative gains, meaning it should not reduce performance with respect to *no stopping*. With many outliers in the online data, all methods do reduce the performance. All methods show a similar dependency on data separation, and apart from *Liu* and, to a much lesser extend, *Jin* and *Rank diff*, all methods reduce to baseline for inseparable data. *Liu* as the only method is susceptible to all sorts of data distortion, including drifts and scaling.



Supplementary Figure 2: Individual subject profiles. A) To indicate the effect size, the baseline performance of *no stopping* is shown. B) Shown are for each subject and each stopping methods the achieved gains over the *no stopping* condition. Numbers indicate the average gain per method per paradigm (see main manuscript, Figure 6). Negative gains mean that a method reduced the performance for a subject (negative bar) or on average over an entire paradigm (negative number).



Supplementary Figure 3: Estimating S^\dagger . The individually optimized scaling parameter S is plotted against the AUC of the training data. High linear correlations are found between these, which leads to a set of coefficients for each methods that allow the direct estimate of S^\dagger from the training data. The black line indicates the linear fit of all data, using an iteratively reweighted least squares fit. The gray lines indicate the fit of all but one dataset, showing high resemblance with the overall fit. Note that the α values for *Höhne* and *Zhang* are plotted on a log scale, given the nature of the parameter.

Appendix A. Höhne method in depth

The method is based on the Welch's t-test, which can be used to assess whether two normally distributed samples with individual variance have equal mean. The test statistic can be calculated as

$$t(X_1, X_2) = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{s_1^2}{N_1} + \frac{s_2^2}{N_2}}}. \quad (\text{A.1})$$

where \bar{X}_i , s_i^2 , and N_i are sample mean, sample variance and sample size, respectively.

In the BCI scenario, the two normally distributed samples are classifier outputs of class c vs. the joint set of classifier outputs of all other classes \tilde{c} , thus $t(\mathbf{D}_{c,1\dots j,t}, \mathbf{D}_{\tilde{c},1\dots j,t})$. The test statistic is transformed into a p-value in order to obtain a measure which is independent of sample size.

In the early stopping problem, we are finally testing, if there is a class c that fulfills the stopping criterion,

$$M(\mathbf{D}^{test}, t, j) < \alpha \quad (\text{A.2})$$

Sub-function $M(\mathbf{D}^{test}, t, j)$ is defined by

$$M(\mathbf{D}, t, j) = \min_c p(t(\mathbf{D}_{c,1\dots j,t}, \mathbf{D}_{\tilde{c},1\dots j,t}), \nu) \quad (\text{A.3})$$

As one can see in Equation A.6, the transformation of the t-value into the p-value requires an additional parameter ν , which can be directly computed from the data as well. The transformation from the t-value into the corresponding p-value is implemented in numerous standard statistical toolboxes. Nevertheless, a detailed analytic derivation is given below.

$$p(t, \nu) = \int_{-\infty}^t f(u) du \quad (\text{A.4})$$

$$\nu = \frac{(\frac{s_1^2}{N_1} + \frac{s_2^2}{N_2})^2}{\frac{s_1^4}{N_1^2(N_1-1)} + \frac{s_2^4}{N_2^2(N_2-1)}} \quad (\text{A.5})$$

$$f(t) = \frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\nu\pi} \Gamma(\frac{\nu}{2})} \left(1 + \frac{t^2}{\nu}\right)^{-\frac{\nu+1}{2}} \quad (\text{A.6})$$

Here, Γ denotes the Gamma-function, in contrast to its previous usage as the training function.

Appendix B. Classifier output distributions in Zhang

In practice, the classifier output distributions f_a and f_u are modeled as univariate Gaussians

$$f_{a/u}(d) = \frac{1}{\sqrt{2\pi\sigma_{a/u}^2}} \exp\left(-\frac{(d - \mu_{a/u})^2}{2\sigma_{a/u}^2}\right), \quad (\text{B.1})$$

where the means and variances are estimated from the training data by

$$\hat{\mu}_a = \frac{1}{JT} \sum_{t=1}^T \sum_{j=1}^J \mathbf{D}_{t,j,t}^{train} \quad (\text{B.2})$$

$$\hat{\mu}_u = \frac{1}{JT(C-1)} \sum_{t=1}^T \sum_{j=1}^J \sum_{c \neq l_t} \mathbf{D}_{c,j,t}^{train} \quad (\text{B.3})$$

and

$$\hat{\sigma}_a^2 = \frac{1}{JT-1} \sum_{t=1}^T \sum_{j=1}^J (\mathbf{D}_{t,j,t}^{train} - \hat{\mu}_a)^2 \quad (\text{B.4})$$

$$\hat{\sigma}_u^2 = \frac{1}{JT(C-1)-1} \sum_{t=1}^T \sum_{j=1}^J \sum_{c \neq l_t} (\mathbf{D}_{c,j,t}^{train} - \hat{\mu}_u)^2. \quad (\text{B.5})$$

Such modeling is in line with the LDA classifier assuming the features of the two classes to be normally distributed as well. In fact, $\hat{\mu}_{a/u}$ and $\hat{\sigma}_{a/u}^2$ can be obtained from the feature means and covariance matrix $\mathbf{m}_{a/u}$ and Σ estimated by LDA via

$$\hat{\mu}_{a/u} = \mathbf{w}^\top \mathbf{m}_{a/u} \quad \text{and} \quad \hat{\sigma}_a^2 = \hat{\sigma}_u^2 = \mathbf{w}^\top \Sigma \mathbf{w}, \quad (\text{B.6})$$

where \mathbf{w} is the LDA projection vector.