# Classical (frequentist) inference

## Methods & models for fMRI data analysis
## 23 October 2018

**Klaas Enno Stephan**

With many thanks for slides & images to:

FIL Methods group
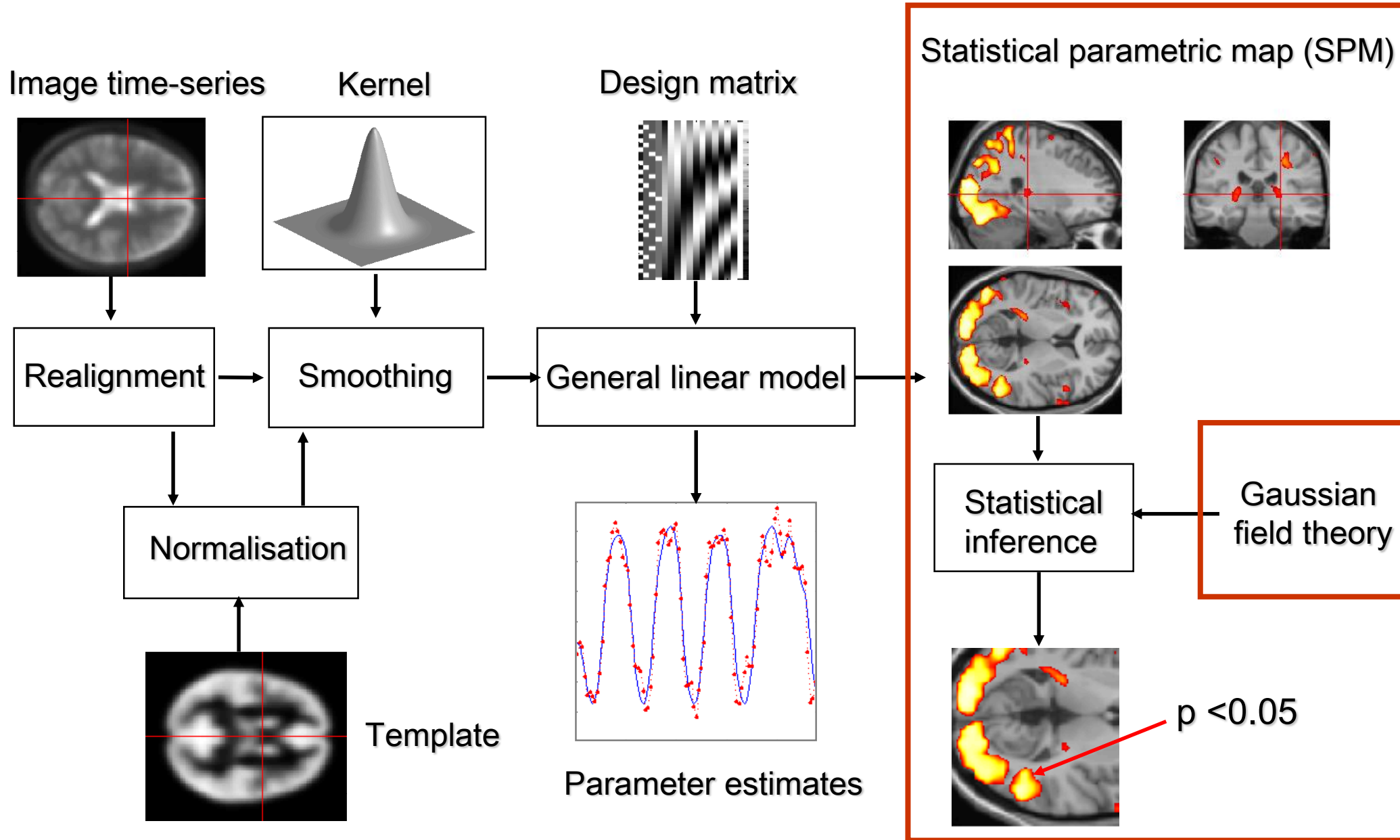
Translational Neuromodeling Unit

Universität Zürich UZH
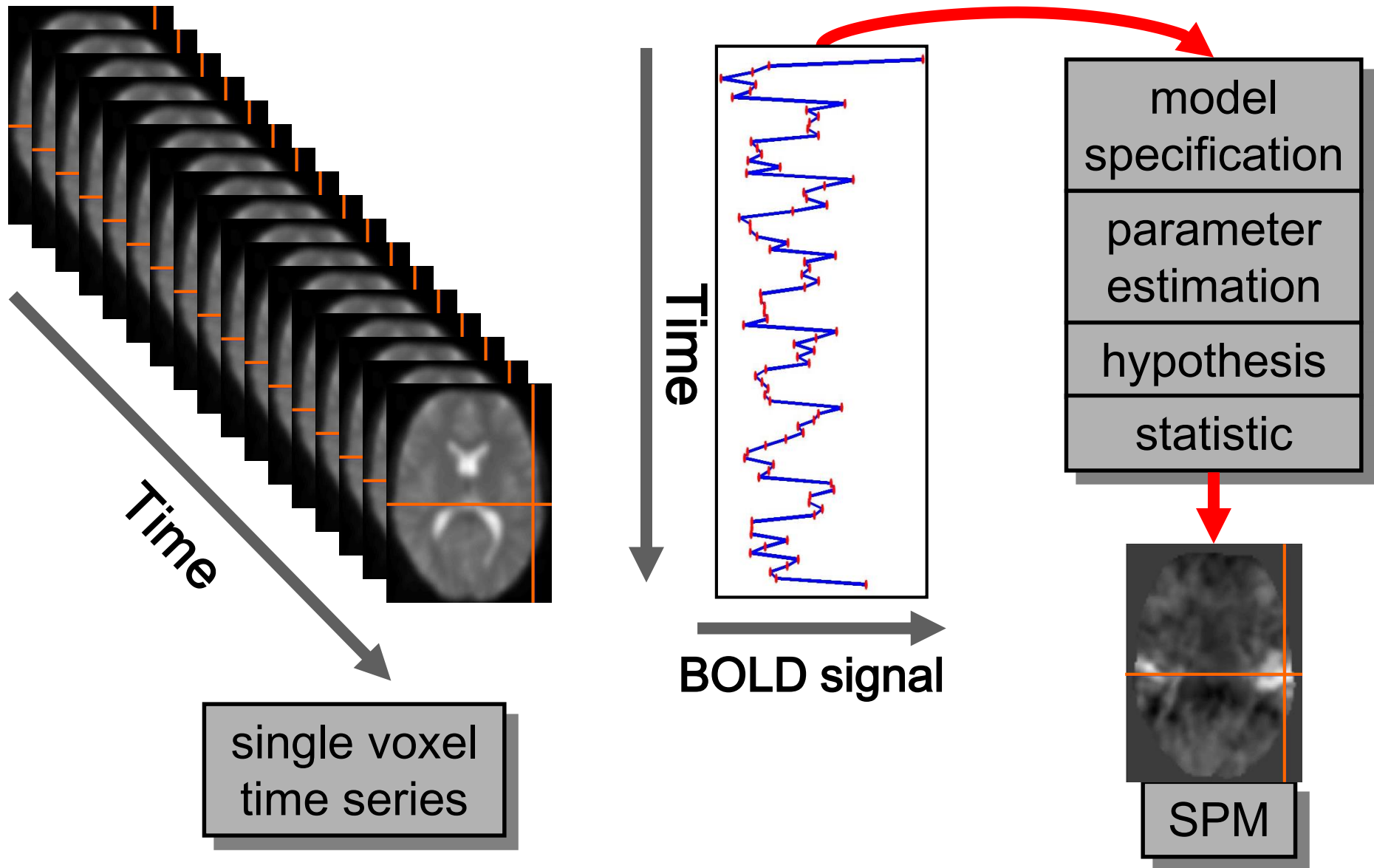
ETH
Eidgenössische Technische Hochschule Zürich
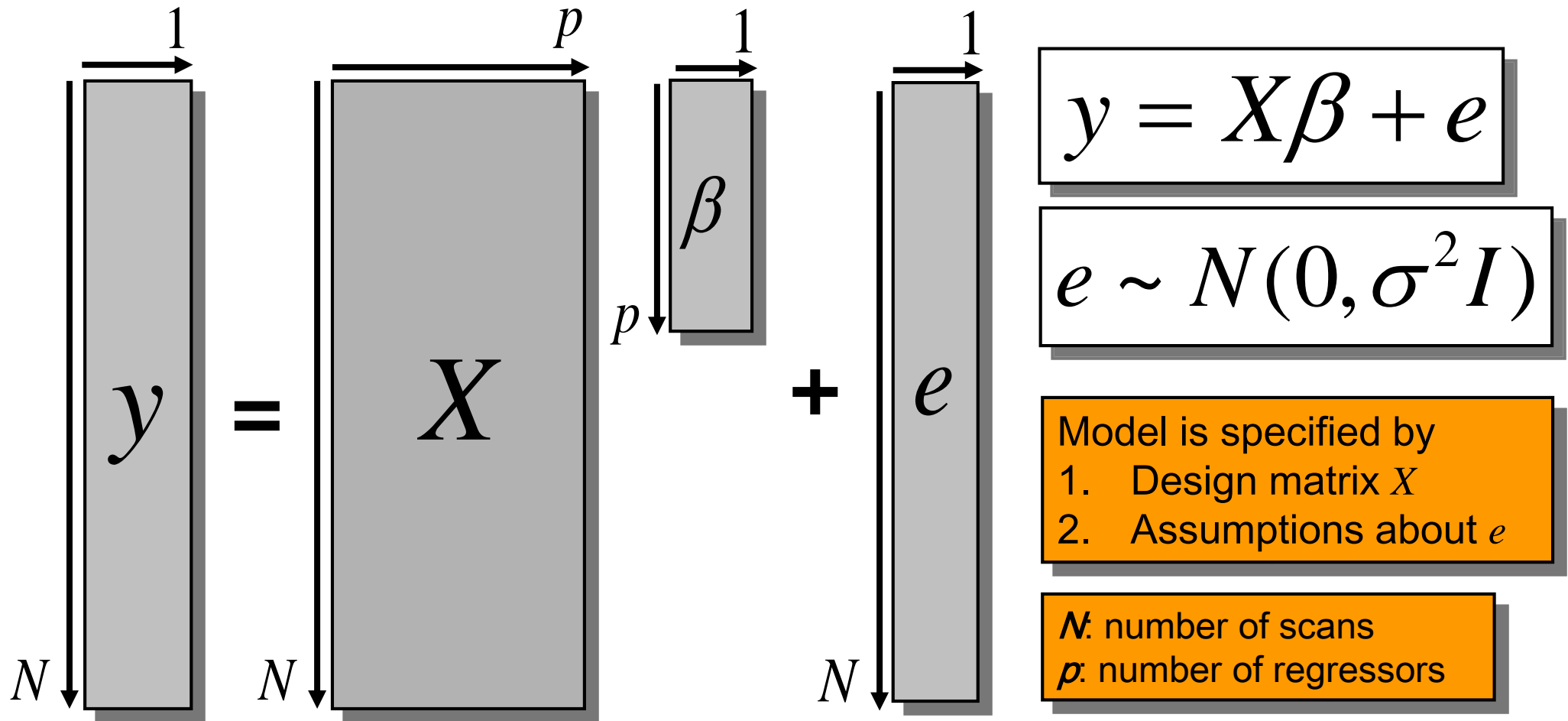Swiss Federal Institute of Technology Zurich

# Overview of SPM



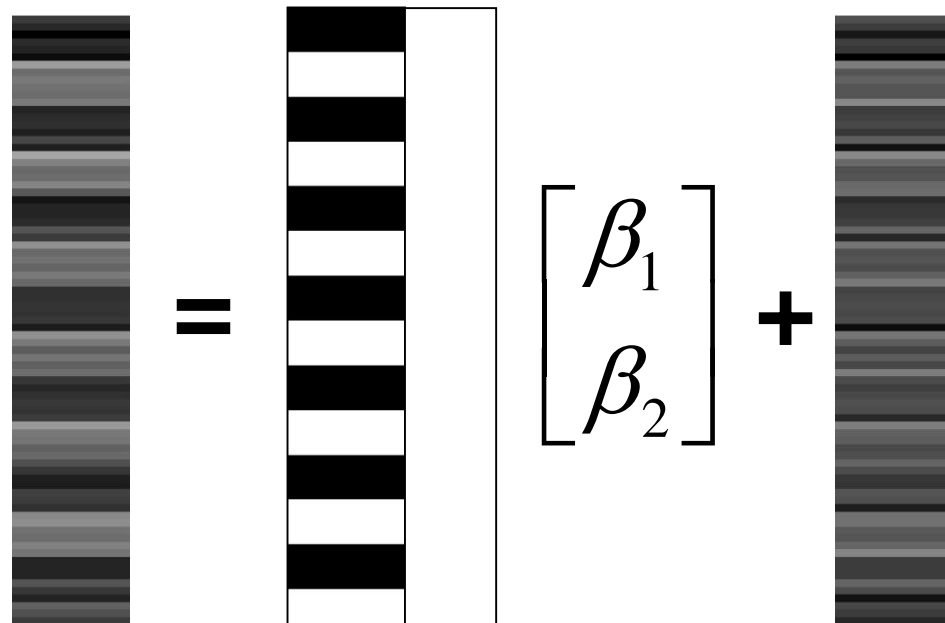Image time-series — Kernel — Design matrix — Statistical parametric map (SPM)

Realignment → Smoothing → General linear model → Statistical inference ← Gaussian field theory

Normalisation

Template

Parameter estimates

p <0.05

# Voxel-wise time series analysis



single voxel
time series

BOLD signal

model
specification

parameter
estimation

hypothesis

statistic

SPM

# Mass-univariate analysis: voxel-wise GLM



$$y = X\beta + e$$

$$e \sim N(0, \sigma^2 I)$$

Model is specified by
1. Design matrix $X$
2. Assumptions about $e$

$N$: number of scans
$p$: number of regressors

The design matrix embodies all available knowledge about experimentally controlled factors and potential confounds.

# Ordinary least squares (OLS) parameter estimation



$$y = \begin{bmatrix} X \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} + e$$

$y$    $X$    $e$

$$y = X\beta + e$$

**Objective:** estimate parameters to minimize $\sum_{t=1}^{N} e_t^2$

**Ordinary least squares estimation (OLS) (assuming i.i.d. error):**

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

# OLS parameter estimation

The Ordinary Least Squares (OLS) estimators are: $\hat{\beta} = (X^T X)^{-1} X^T y$

These estimators minimise $\sum e_t^2 = e^T e$. They are found solving either

$$\frac{\partial \left( \sum e_t^2 \right)}{\partial \hat{\beta}_t} = 0 \quad \text{or} \quad X^T e = 0$$

Under i.i.d. assumptions, the OLS estimates correspond to ML estimates:

$$e \sim N(0, \sigma^2 I) \longrightarrow Y \sim N(X\beta, \sigma^2 I)$$

$$\hat{\sigma}^2 = \frac{\hat{e}^T \hat{e}}{N - p}$$

$$\hat{\beta} \sim N(\beta, \boxed{\sigma^2 (X^T X)^{-1}})$$

NB: precision of our estimates depends on design matrix!

# Maximum likelihood (ML) estimation

probability density function ($\theta$ fixed!)

$$y \mapsto p(y \mid \theta)$$

likelihood function ($y$ fixed!)

$$\theta \mapsto p(y \mid \theta)$$
$$L(\theta \mid y) = p(y \mid \theta)$$

ML estimator

$$\hat{\theta} = \arg \max_{\theta} L(\theta \mid y)$$

For $cov(e) = \sigma^2 I$, the ML estimator is equivalent to the OLS estimator:

$$\hat{\beta} = (X^T X)^{-1} X^T y \quad \textbf{OLS}$$

For $cov(e) = \sigma^2 V$, the ML estimator is equivalent to a weighted least squares (WLS) estimate (with $W = V^{-1/2}$):

$$\hat{\beta} = (X^T W X)^{-1} X^T W y \quad \textbf{WLS}$$

# Recap: Dealing with non-sphericity by defining a filter matrix $W$

- Enhanced noise model

$$e \sim N(0, \sigma^2 V)$$

- Remember linear transform for Gaussians

$$x \sim N(\mu, \sigma^2), y = ax$$
$$\Rightarrow y \sim N(a\mu, a^2\sigma^2)$$

- Choose $W$ such that error covariance becomes spherical

$$We \sim N(0, \sigma^2 W^2 V)$$
$$\Rightarrow W^2 V = I$$

- **Conclusion: $W$ is a simple function of $V$ $\Rightarrow$ so how do we estimate $V$?**

$$\Rightarrow W = V^{-1/2}$$

$$Wy = WX\beta + We$$

# Recap: Estimating $V$: Multiple covariance components
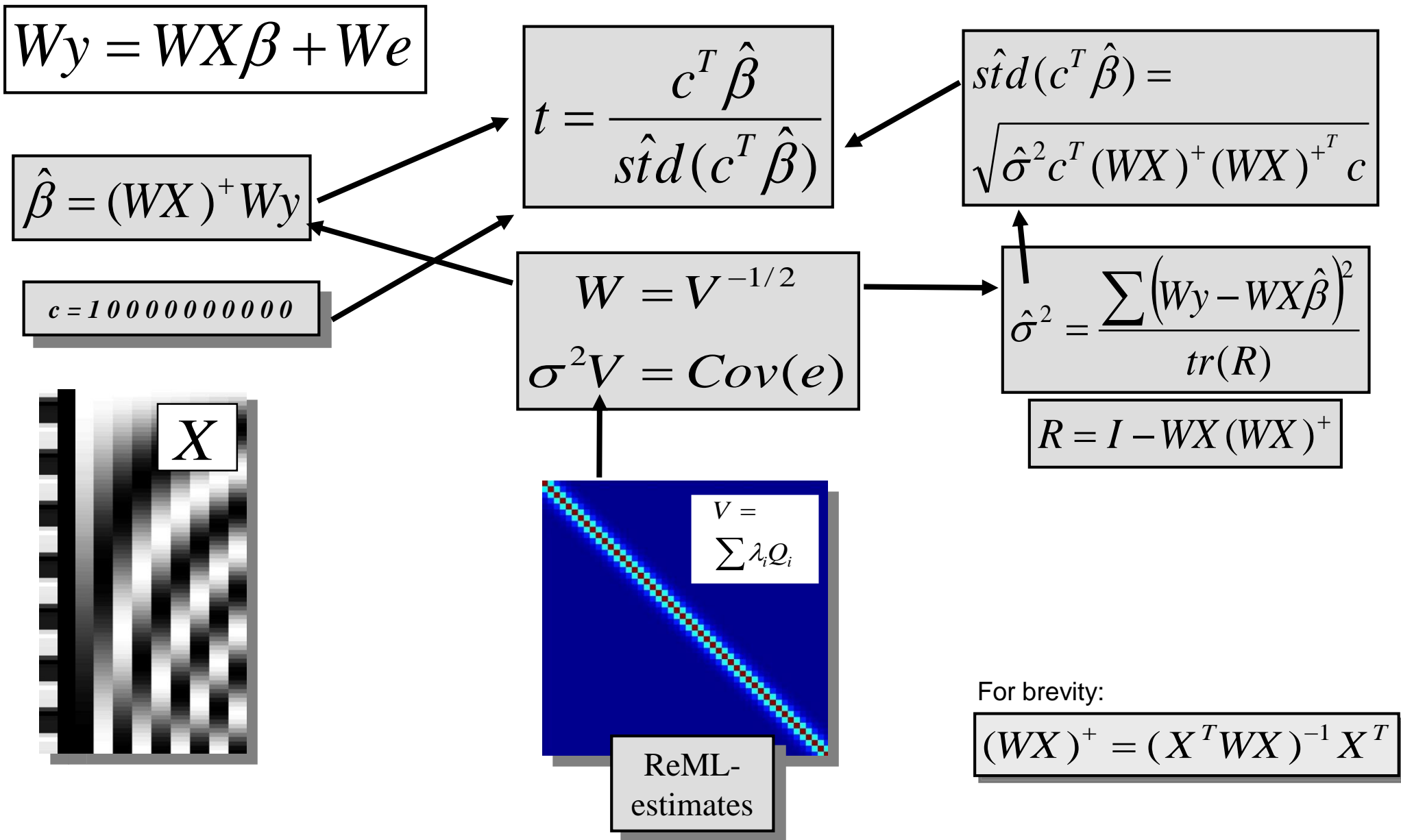
$$e \sim N(0, \sigma^2 V)$$

enhanced noise model

$$V \propto Cov(e)$$
$$V = \sum \lambda_i Q_i$$

error covariance components $Q$
and hyperparameters $\lambda$



$$V = \lambda_1 Q_1 + \lambda_2 Q_2$$

Estimation of hyperparameters $\lambda$ with ReML (restricted maximum likelihood).

# Bonus material: t-statistic based on ML estimates in SPM

$$Wy = WX\beta + We$$

$$\hat{\beta} = (WX)^{+}Wy$$

$$c = 1\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0$$

$$X$$

$$t = \frac{c^{T}\hat{\beta}}{s\hat{t}d(c^{T}\hat{\beta})}$$

$$W = V^{-1/2}$$
$$\sigma^{2}V = Cov(e)$$

$$V = \sum \lambda_{i}Q_{i}$$

ReML-estimates

$$s\hat{t}d(c^{T}\hat{\beta}) = \sqrt{\hat{\sigma}^{2}c^{T}(WX)^{+}(WX)^{+T}c}$$

$$\hat{\sigma}^{2} = \frac{\sum(Wy - WX\hat{\beta})^{2}}{tr(R)}$$

$$R = I - WX(WX)^{+}$$

For brevity:

$$(WX)^{+} = (X^{T}WX)^{-1}X^{T}$$

# Terminology

- A **statistic** is the result of applying a mathematical function to a **sample** (set of data).

- (More formally, a **statistic** is a function of a sample where the function itself is independent of the sample's distribution. The term is used both for the function and for the value of the function on a given sample.)

- A statistic is distinct from an unknown statistical **parameter**, which is a population property and can only be estimated approximately from a sample.

- A statistic used to estimate a parameter is called an **estimator**.
  For example, the sample mean is a statistic and an estimator for the population mean, which is a parameter.

# Hypothesis testing

To test an hypothesis, we construct a "test statistic".

- **"Null hypothesis" H$_0$ = "there is no effect"** $\Rightarrow c^T \beta = 0$

  This is what we want to disprove.

  $\Rightarrow$ The "alternative hypothesis" H$_1$ represents the outcome of interest.

- **The test statistic T**

  The test statistic summarises the evidence for H$_0$.

  $\Rightarrow$ We need to know the distribution of T under the null hypothesis.
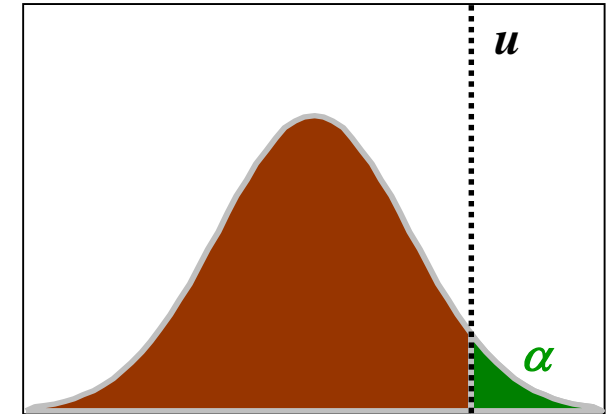


Null Distribution of T

# Hypothesis testing

- **Type I Error $\alpha$:**

  Acceptable *false positive rate* $\alpha$.

  Threshold $u$ controls the false positive rate

  $$\alpha = p(T > u \mid H_0)$$

- **Observation of test statistic t, a realisation of T:**

  A *p*-value summarises evidence against $H_0$.

  This is the probability of observing t, or a more extreme value, under the null hypothesis:

  $$p(T \geq t \mid H_0)$$

- **The conclusion about the hypothesis:**

  We reject $H_0$ in favour of $H_1$ if $t > u$



Null Distribution of T



Null Distribution of T

# Types of error

Actual condition

Test result

|  | H$_0$ true | H$_0$ false |
|---|---|---|
| **Reject H$_0$** | **False positive (FP)**<br><br>**Type I error** $\alpha$ | **True positive (TP)** |
| **Failure to reject H$_0$** | **True negative (TN)** | **False negative (FN)**<br><br>**Type II error** $\beta$ |

**specificity: 1-$\alpha$**
= TN / (TN + FP)
= proportion of actual
negatives which are
correctly identified

**sensitivity (power): 1-$\beta$**
= TP / (TP + FN)
= proportion of actual
positives which are
correctly identified

# One cannot accept the null hypothesis
# (one can only fail to reject it)



**Absence of evidence is not evidence of absence!**

If we do not reject $H_0$, then all can we say is that there is not enough evidence in the data to reject $H_0$. This does not mean that we can accept $H_0$.

**What does this mean for neuroimaging results based on classical statistics?**

A failure to find an "activation" in a particular area does not mean we can conclude that this area is not involved in the process of interest.

# Contrasts

- We are usually not interested in the whole $\beta$ vector.

- A contrast $c^T\beta$ selects a specific effect of interest:
  - $\Rightarrow$ a contrast vector $c$ is a vector of length $p$
  - $\Rightarrow$ $c^T\beta$ is a linear combination of regression coefficients $\beta$

$c^T = [1\ 0\ 0\ 0\ 0\ ...]$

$c^T\beta = \mathbf{1}\beta_1 + \mathbf{0}\beta_2 + \mathbf{0}\beta_3 + \mathbf{0}\beta_4 + \mathbf{0}\beta_5 + ...$

$c^T = [0\ \text{-}1\ 1\ 0\ 0\ ...]$

$c^T\beta = \mathbf{0}\beta_1 + \mathbf{-1}\beta_2 + \mathbf{1}\beta_3 + \mathbf{0}\beta_4 + \mathbf{0}\beta_5 + ...$

- Under i.i.d assumptions:

$$c^T \hat{\beta} \sim N(c^T \beta, \sigma^2 \boxed{c^T (X^T X)^{-1} c})$$

NB: the precision of our estimates depends on design matrix and the chosen contrast !

# Bonus material: Estimability of parameters

- If $X$ is not of **full rank** then different parameters can give identical predictions, i.e. $X\beta_1 = X\beta_2$ with $\beta_1 \neq \beta_2$.

- The parameters are therefore 'non-unique', 'non-identifiable' or '**non-estimable**'.

- For such models, $X^T X$ is not invertible so we must resort to generalised inverses (SPM uses the Moore-Penrose **pseudo-inverse**).

- This gives a parameter vector that has the smallest norm of all possible solutions.

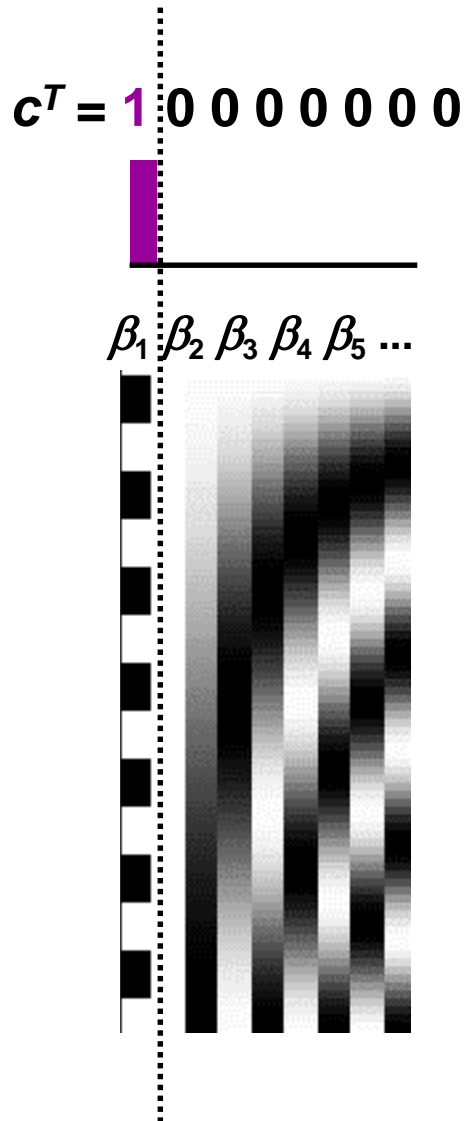- However, even when parameters are non-estimable, certain contrasts may well be!

One-way ANOVA
(unpaired two-sample *t*-test)

1 0 1
1 0 1
1 0 1
1 0 1
0 1 1
0 1 1
0 1 1
0 1 1



Rank(X)=2

# Bonus material: Estimability of parameters

- If $X$ is not of **full rank** then different parameters can give identical predictions, i.e. $X\beta_1 = X\beta_2$ with $\beta_1 \neq \beta_2$.

- The parameters are therefore 'non-unique', 'non-identifiable' or '**non-estimable**'.

- For such models, $X^TX$ is not invertible so we must resort to generalised inverses (SPM uses the Moore-Penrose **pseudo-inverse**).

- This gives a parameter vector that has the smallest norm of all possible solutions.

- However, even when parameters are non-estimable, certain contrasts may well be!



parameters

parameter estimability

(gray $\vec{\beta}$ not uniquely specified)

# Bonus material: Estimability of contrasts

- Linear dependency: there is one contrast vector $q$ for which $Xq = 0$.

- Thus: $y = X\beta + Xq + e = X(\beta + q) + e$

- So if we test $c^T\beta$ for a design matrix with linear dependencies, we implicitly also test $c^T(\beta + q)$, thus an estimable contrast has to satisfy $c^Tq = 0$.

- In the above ANOVA example (unpaired t-test), any contrast vector that is orthogonal to q=[1 1 -1] is estimable:
  [1  0  0], [0  1  0], [0  0  1] are not estimable.
  [1  0  1], [0  1  1], [1  -1  0], [0.5  0.5  1] are estimable.

# Student's t-distribution

- first described by William Sealy Gosset, a statistician at the Guinness brewery at Dublin

- t-statistic is a signal-to-noise measure:  t = effect / standard deviation

- t-distribution is an approximation to the normal distribution for small samples

- t-contrasts are simply linear combinations of the betas
  - $\Rightarrow$ the t-statistic does not depend on the scaling of the regressors or on the scaling of the contrast

- Unilateral test in SPM:
$$H_0 : c^T \beta = 0 \ \text{vs.} \ H_1 : c^T \beta > 0$$



Probability density function of Student's t distribution

# t-contrasts – SPM{t}

$c^T = $ **1** 0 0 0 0 0 0 0

$\beta_1\ \beta_2\ \beta_3\ \beta_4\ \beta_5\ ...$

**Question:** box-car amplitude > 0 ?

$$= $$

$$H_1 = c^T\beta > 0 \ ?$$

**Null hypothesis:** $H_0:\ c^T\beta = 0$

**Test statistic:**

$$t = \frac{\textit{contrast of parameter estimates}}{\text{variance estimate}}$$

$p(y\,|\,c^T\hat{\beta} = 0)$

$$t = \frac{c^T\hat{\beta}}{\hat{std}(c^T\hat{\beta})} = \frac{c^T\hat{\beta}}{\sqrt{\hat{\sigma}^2 c^T\left(X^T X\right)^{-1} c}} \sim t_{N-p}$$

# t-contrasts in SPM

For a given contrast $c$:



beta_???? images

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

ResMS image

$$\hat{\sigma}^2 = \frac{\hat{e}^T \hat{e}}{N - p}$$

con_???? image

$$c^T \hat{\beta}$$

spmT_???? image

SPM{$t$}

# t-contrast: a simple example

Passive word listening versus rest



$$c^T = [\; 1 \qquad 0\; ]$$



X

Design matrix

Q: activation during listening ?

Null hypothesis: $\beta_1 = 0$
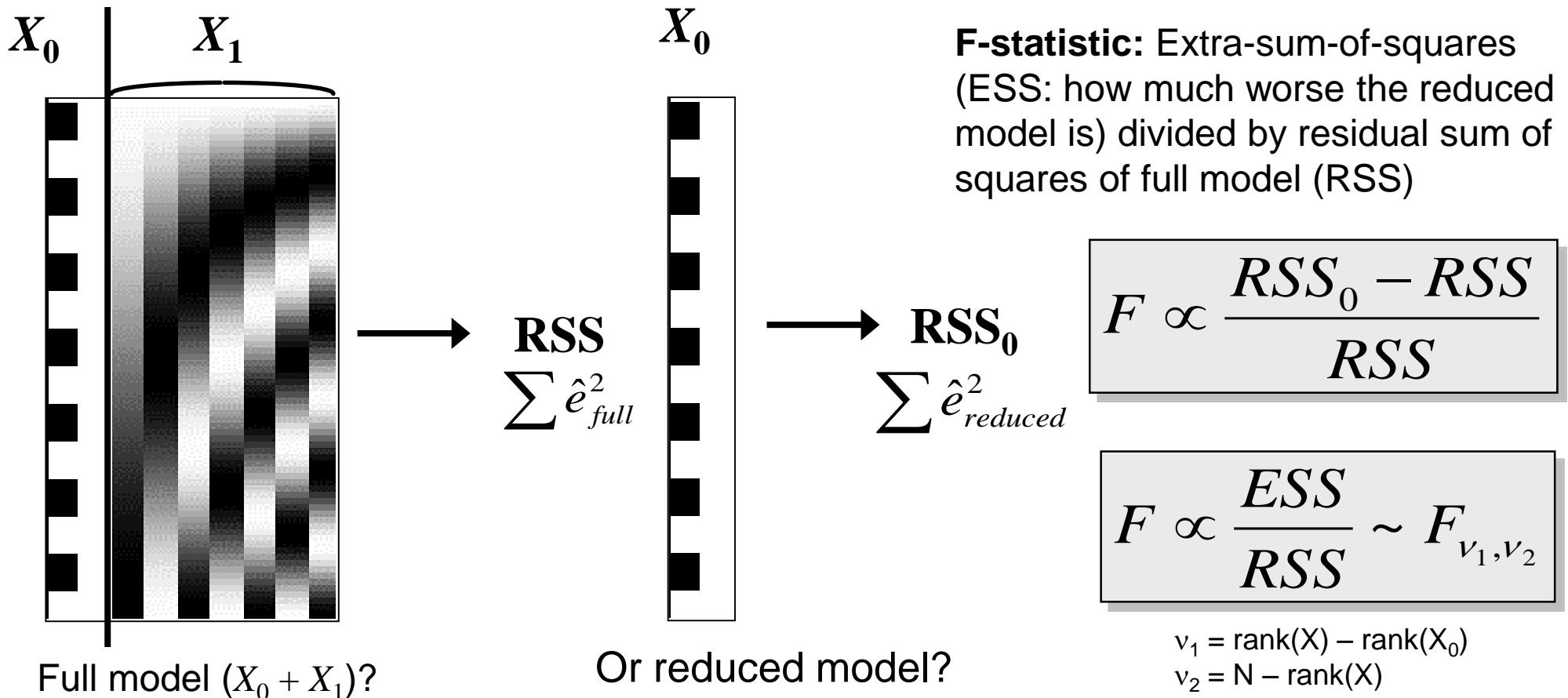
$$t = \frac{c^T \hat{\beta}}{Std(c^T \hat{\beta})}$$

$$p(y \mid c^T \hat{\beta} = 0)$$

**SPMresults:**
Height threshold T = 3.2057 {p<0.001}

**Statistics:** *p-values adjusted for search volume*

| set-level | | cluster-level | | | voxel-level | | | | | mm mm mm |
|---|---|---|---|---|---|---|---|---|---|---|
| p | c | p corrected | $k_E$ | p uncorrected | p FWE-corr | p FDR-corr | T | ($Z_=$) | p uncorrected | |
| 0.000  10 | | 0.000 | 520 | 0.000 | 0.000 | 0.000 | 13.94 | Inf | 0.000 | −63 −27  15 |
| | | | | | 0.000 | 0.000 | 12.04 | Inf | 0.000 | −48 −33  12 |
| | | | | | 0.000 | 0.000 | 11.82 | Inf | 0.000 | −66 −21   6 |
| | | 0.000 | 426 | 0.000 | 0.000 | 0.000 | 13.72 | Inf | 0.000 | 57 −21  12 |
| | | | | | 0.000 | 0.000 | 12.29 | Inf | 0.000 | 63 −12  −3 |
| | | | | | 0.000 | 0.000 | 9.89 | 7.83 | 0.000 | 57 −39   6 |
| | | 0.000 | 35 | 0.000 | 0.000 | 0.000 | 7.39 | 6.36 | 0.000 | 36 −30 −15 |
| | | 0.000 | 9 | 0.000 | 0.000 | 0.000 | 6.84 | 5.99 | 0.000 | 51   0  48 |
| | | 0.002 | 3 | 0.024 | 0.001 | 0.000 | 6.36 | 5.65 | 0.000 | −63 −54  −3 |
| | | 0.000 | 8 | 0.001 | 0.001 | 0.000 | 6.19 | 5.53 | 0.000 | −30 −33 −18 |
| | | 0.000 | 9 | 0.000 | 0.003 | 0.000 | 5.96 | 5.36 | 0.000 | 36 −27   9 |
| | | 0.005 | 2 | 0.058 | 0.004 | 0.000 | 5.84 | 5.27 | 0.000 | −45 42   9 |
| | | 0.015 | 1 | 0.166 | 0.022 | 0.000 | 5.44 | 4.97 | 0.000 | 48 27 24 |
| | | 0.015 | 1 | 0.166 | 0.036 | 0.000 | 5.32 | 4.87 | 0.000 | 36 −27 42 |

# F-test: the extra-sum-of-squares principle

Model comparison: Full vs. reduced model

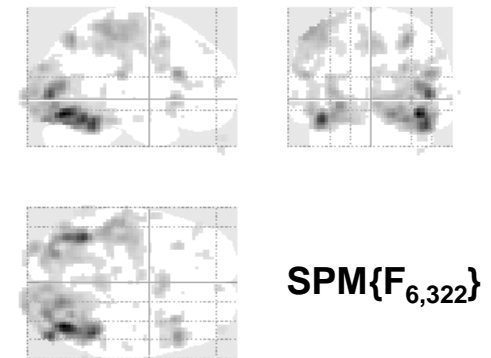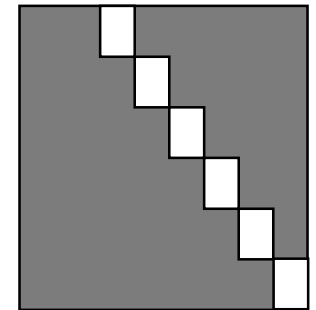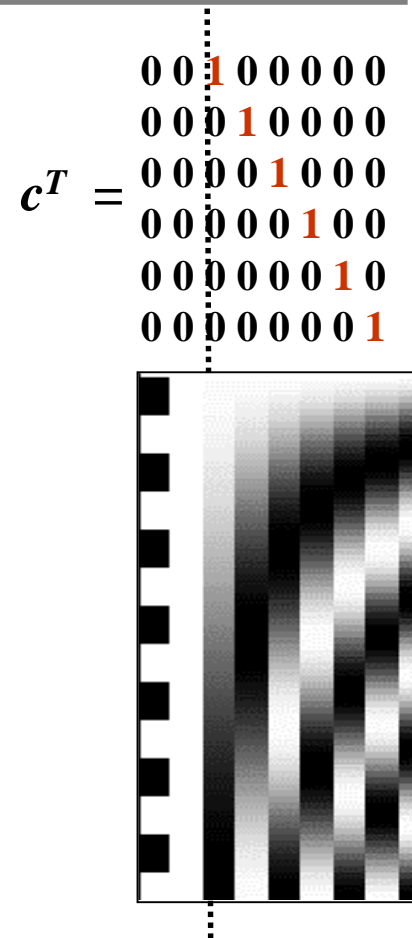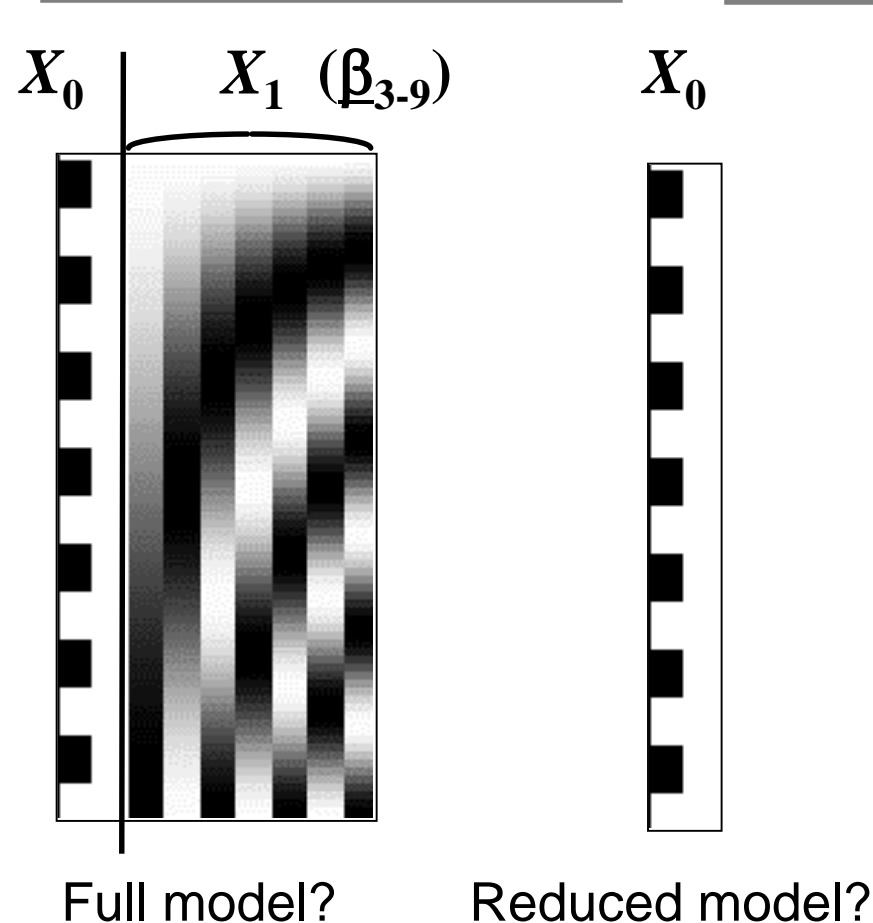**Null Hypothesis H$_0$:** True model is $X_0$ (reduced model)

$X_0$   $X_1$                    $X_0$

**F-statistic:** Extra-sum-of-squares (ESS: how much worse the reduced model is) divided by residual sum of squares of full model (RSS)

**RSS**
$$\sum \hat{e}^2_{full}$$

**RSS$_0$**
$$\sum \hat{e}^2_{reduced}$$

$$F \propto \frac{RSS_0 - RSS}{RSS}$$

$$F \propto \frac{ESS}{RSS} \sim F_{v_1, v_2}$$

Full model $(X_0 + X_1)$?        Or reduced model?

$v_1 = \text{rank}(X) - \text{rank}(X_0)$
$v_2 = N - \text{rank}(X)$

# F-test: multidimensional contrasts – SPM{F}

Tests multiple linear hypotheses:

$H_0$: True model is $X_0$

$H_0$: $\beta_3 = \beta_4 = ... = \beta_9 = 0$

test $H_0$ : $c^T \beta = 0$ ?

$X_0$   $X_1$ ($\beta_{3-9}$)   $X_0$

$$c^T = \begin{matrix} 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{matrix}$$

Full model?    Reduced model?

SPM{$F_{6,322}$}

# F-test: a few remarks

- F-tests can be viewed as testing for the additional variance explained by a larger model wrt. a simpler (nested) model ⇨ model comparison

- Hypotheses:

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

**Null hypothesis H$_0$:** $\qquad$ $\beta_1 = \beta_2 = ... = \beta_p = 0$

**Alternative hypothesis H$_1$:** $\qquad$ At least one $\beta_k \neq 0$

- F-tests are not directional:
  When testing a uni-dimensional contrast with an *F*-test, for example $\beta_1 - \beta_2$, the result will be the same as testing $\beta_2 - \beta_1$.

# Bonus material: Differential F-contrasts



contrast(s)

Design matrix

- equivalent to testing for effects that can be explained as a linear combination of the 3 differences

- useful when using informed basis functions and testing for overall shape differences in the HRF between two conditions

# F-contrast in SPM



beta_???? images

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

ResMS image

$$\hat{\sigma}^2 = \frac{\hat{e}^T \hat{e}}{N - p}$$

ess_???? images

$$( RSS_0 - RSS )$$

spmF_???? images

SPM{F}

# F-test example: movement related effects



**To assess movement-related activation:**

There is a lot of residual movement-related artifact in the data (despite spatial realignment), which tends to be concentrated near the boundaries of tissue types.

By including the realignment parameters in our design matrix, we can "regress out" linear components of subject movement, reducing the residual error, and hence improve our statistics for the effects of interest.

# Example: a suboptimal model



True signal (--) and observed signal

Fitting ($\beta_1 = 0.2$, $\beta_2$ (const.) = 0.11);
(here: blue solid line = total fit)

Residuals (still contain some signal)

⇨ Test for the green regressor not significant

# Example: a suboptimal model



$\beta_1 = 0.22$
$\beta_2 = 0.11$

**Y** **=** **X** $\beta$ **+** **e**

*Residual Var.* = 0.3

$p(Y/ \boldsymbol{b}_1 = 0) \Rightarrow$
$p$-value = 0.1
($t$-test)

$p(Y/ \boldsymbol{b}_1 = 0) \Rightarrow$
$p$-value = 0.2
($F$-test)

# A better model



True signal + observed signal

Model (green and red)
and true signal (blue ---)
Red regressor: temporal derivative of
the green regressor

Total fit (blue)
and partial fit (green & red)
Adjusted and fitted signal

Residuals (less variance & structure)

⇨ *t*-test of the green regressor almost significant
⇨ *F*-test very significant
⇨ *t*-test of the red regressor very significant

# A better model

$\beta_1 = 0.22$
$\beta_2 = 2.15$
$\beta_3 = 0.11$



**Y**          **X** $\beta$          **e**

*Residual Var.* $= 0.2$

$p(Y/ \boldsymbol{b}_1 = 0) \Rightarrow$
$p$-value $= 0.07$
($t$-test)

$p(Y/ \boldsymbol{b}_1 = 0, \boldsymbol{b}_2 = 0) \Rightarrow$
$p$-value $= 0.000001$
($F$-test)

# Recap from previous lecture:
# Correlation among regressors



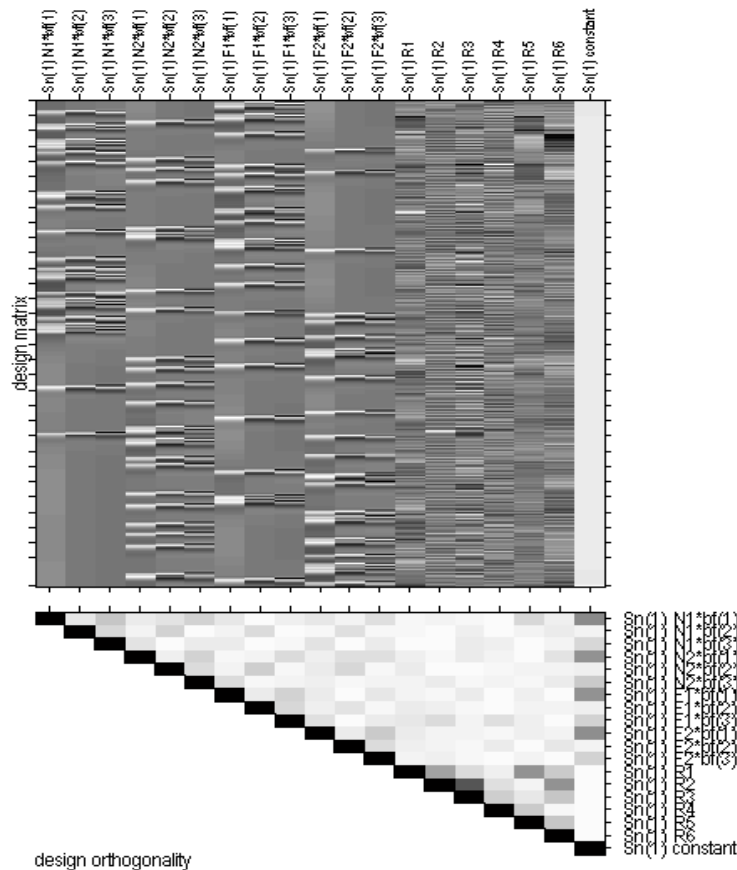$$y = x_1\beta_1 + x_2\beta_2 + e$$
$$\beta_1 = \beta_2 = 1$$

$$y = x_1\beta_1 + x_2^*\beta_2^* + e$$
$$\beta_1 > 1; \beta_2^* = 1$$

Correlated regressors = explained variance is shared between regressors

When $x_2$ is orthogonalized with regard to $x_1$, only the parameter estimate for $x_1$ changes, not that for $x_2$!

# Design orthogonality



Statistical analysis: Design orthogonality

design orthogonality

Measure : abs. value of cosine of angle between columns of design matrix
Scale : black - colinear (cos=+1/-1)
white - orthogonal (cos=0)
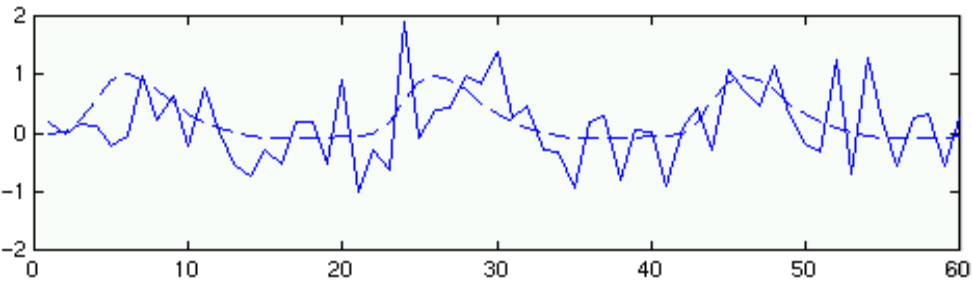gray - not orthogonal or colinear

- For each pair of columns of the design matrix, the orthogonality matrix depicts the magnitude of the **cosine of the angle** between them, with the range 0 to 1 mapped from white to black.

- The cosine of the angle between two vectors *a* and *b* is obtained by:

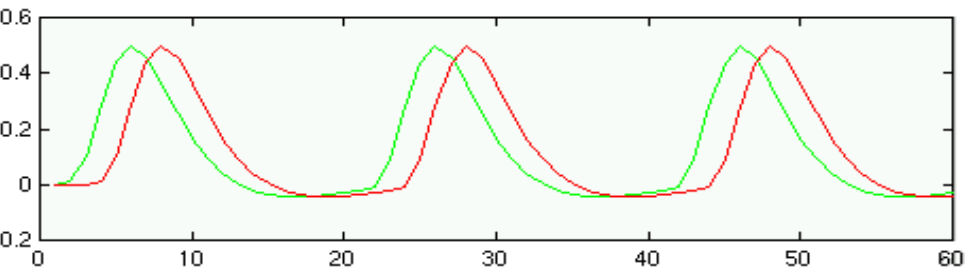$$\cos \alpha = \frac{ab}{|a||b|}$$

- For **zero-mean vectors**, the cosine of the angle between the vectors is the same as the **correlation** between the two variates:
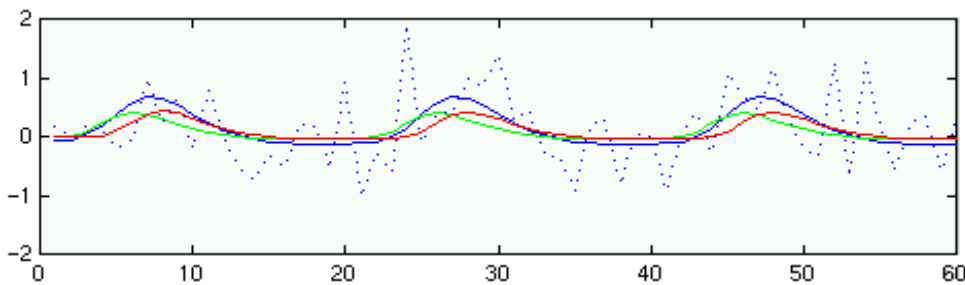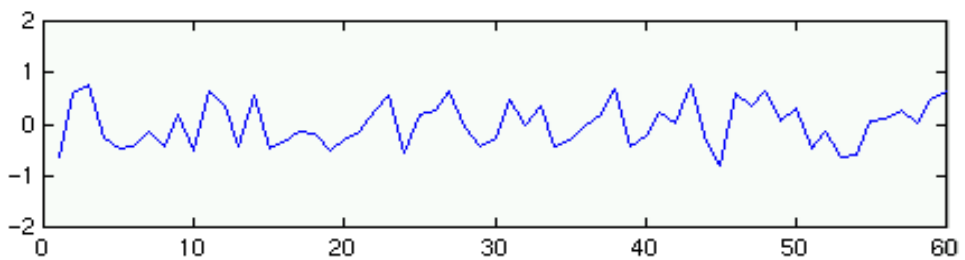
$$\cos \alpha = corr_{a,b}$$

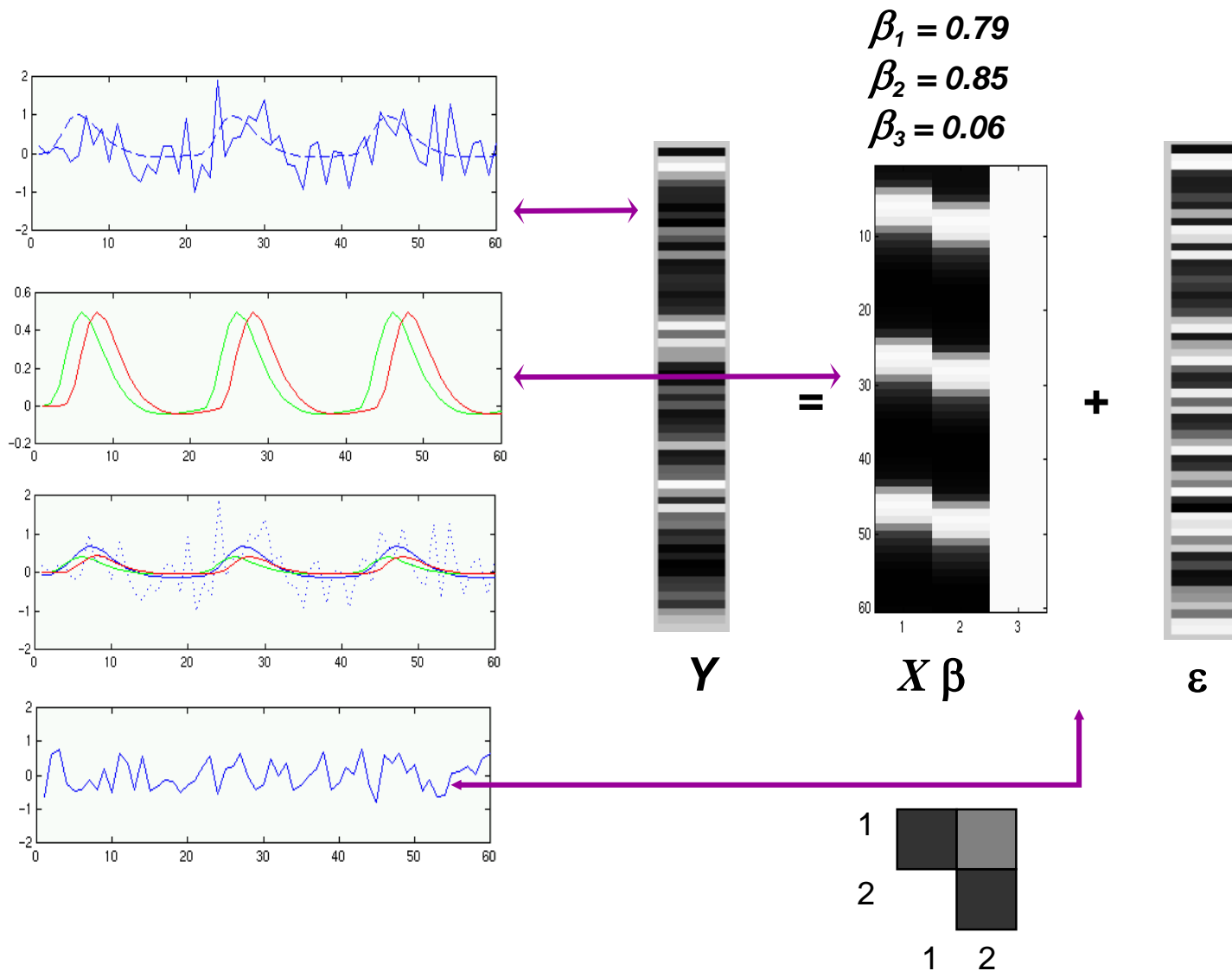# Correlated regressors



True signal

Model (green and red)

Fit (blue: total fit)

Residual

# Correlated regressors



$\beta_1 = 0.79$
$\beta_2 = 0.85$
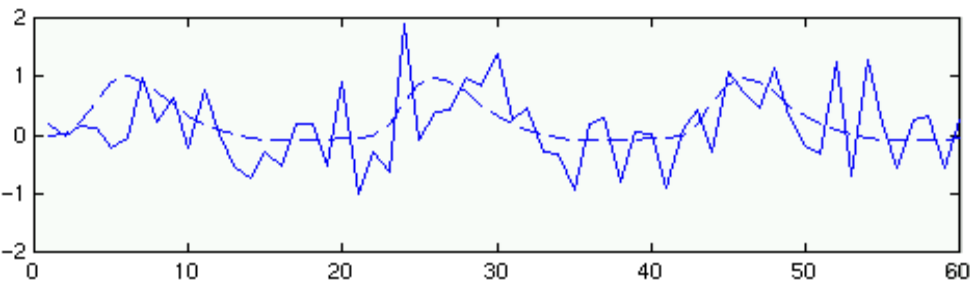$\beta_3 = 0.06$

$Y$ $=$ $X$ $\beta$ $+$ $\varepsilon$

*Residual var.* $= 0.3$

$p(Y/\boldsymbol{b}_1 = 0) \Rightarrow$
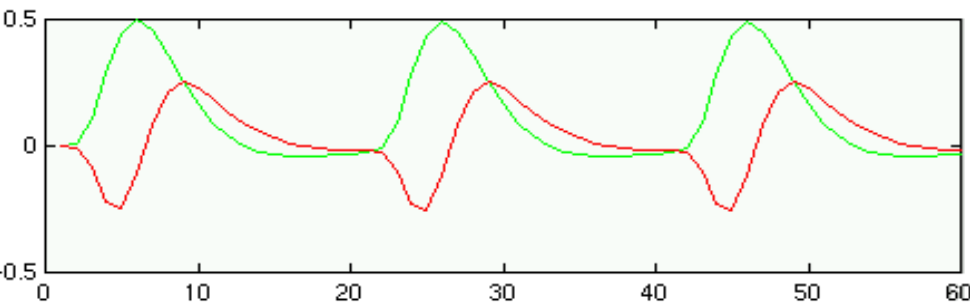$p\text{-}value = 0.08$
($t$-test)

$P(Y/\boldsymbol{b}_2 = 0) \Rightarrow$
$p\text{-}value = 0.07$
($t$-test)

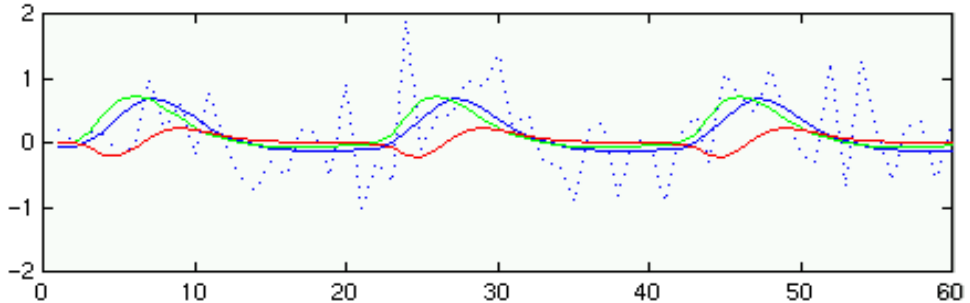$p(Y/\boldsymbol{b}_1 = 0, \boldsymbol{b}_2 = 0) \Rightarrow$
$p\text{-}value = 0.002$
($F$-test)

# After orthogonalisation
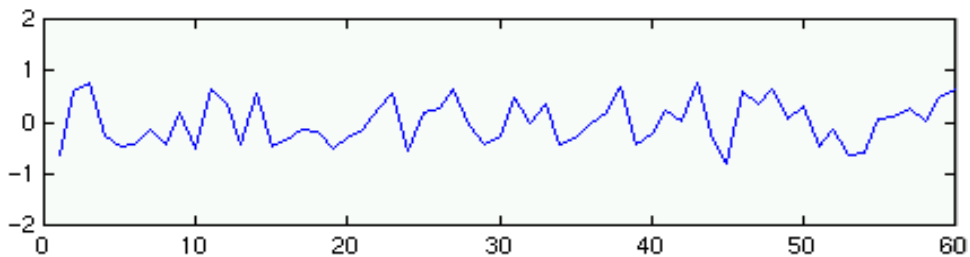


True signal

Model (green and red)
red regressor has been
orthogonalised with respect to the green
one
⇔ remove everything that correlates with
the green regressor

Fit (does not change)
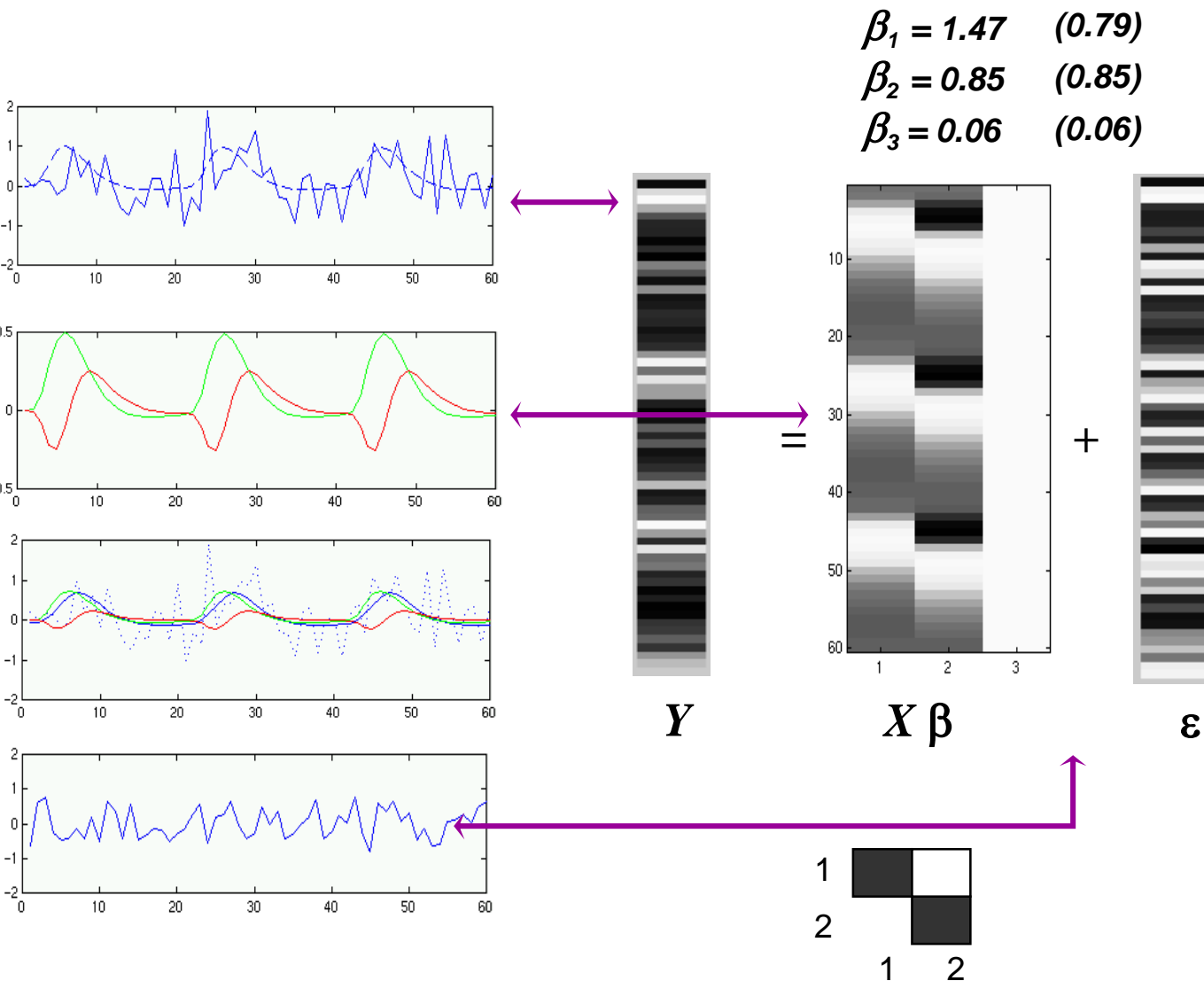
Residuals (do not change)

# After orthogonalisation



$\beta_1 = 1.47$   (0.79)
$\beta_2 = 0.85$   (0.85)
$\beta_3 = 0.06$   (0.06)

$Y$     $X \; \beta$     $\varepsilon$

*Residual var.* $= 0.3$

$p(Y/ \mathbf{b}_1 = 0)$
$p\text{-}value = 0.0003$   **does change**
($t$-test)

$p(Y/ \mathbf{b}_2 = 0)$
$p\text{-}value = 0.07$   **does not change**
($t$-test)

$p(Y/ \mathbf{b}_1 = 0, \mathbf{b}_2 = 0)$
$p\text{-}value = 0.002$   **does not change**
($F$-test)

# Bonus material: Design efficiency

- The aim is to minimize the standard error of a *t*-contrast (i.e. the denominator of a t-statistic).

$$\text{var}(c^T \hat{\beta}) = \hat{\sigma}^2 c^T (X^T X)^{-1} c$$

$$T = \frac{c^T \hat{\beta}}{\sqrt{\text{var}(c^T \hat{\beta})}}$$

- This is equivalent to maximizing the efficiency $\varepsilon$:

$$\varepsilon(\hat{\sigma}^2, c, X) = (\hat{\sigma}^2 c^T (X^T X)^{-1} c)^{-1}$$

Noise variance

Design variance

- If we assume that the noise variance is independent of the specific design:

$$\varepsilon(c, X) = (c^T (X^T X)^{-1} c)^{-1}$$

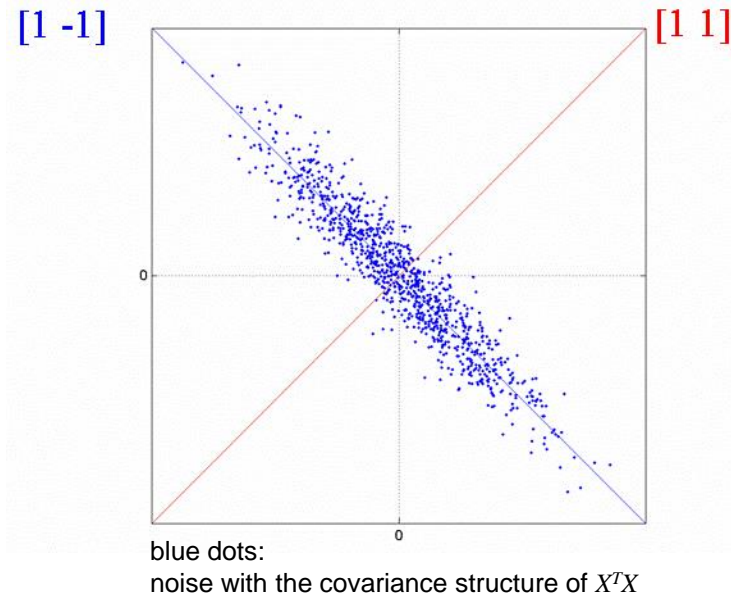NB: efficiency depends on design matrix and the chosen contrast !

- This is a relative measure: all we can say is that one design is more efficient than another (for a given contrast).

# Bonus material: Design efficiency

$$\varepsilon(c, X) = (c^T (X^T X)^{-1} c)^{-1}$$

- $X^T X$ is the covariance matrix of the regressors in the design matrix

- efficiency decreases with increasing covariance

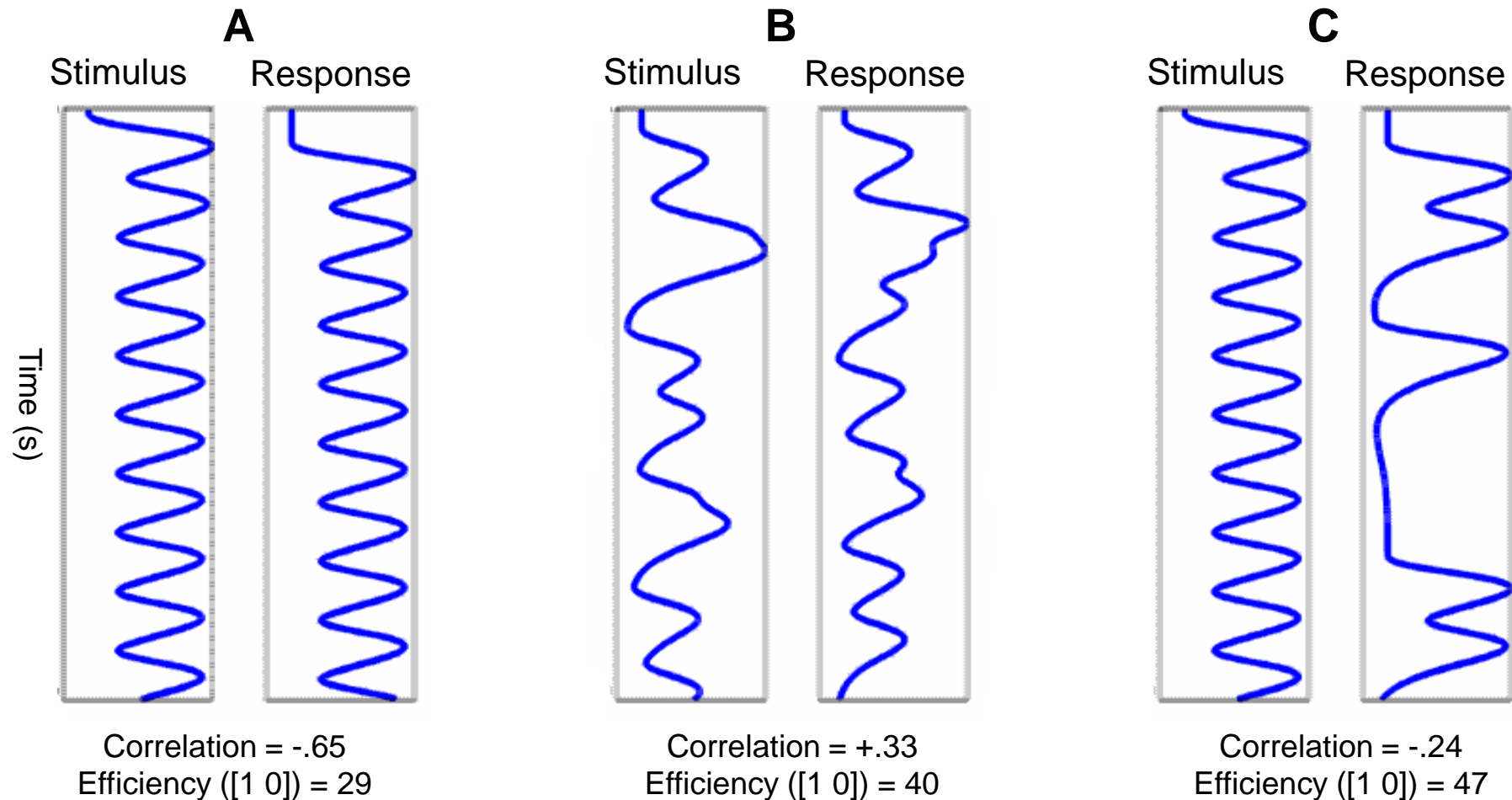- but note that efficiency differs across contrasts

[1 -1]        [1 1]



blue dots:
noise with the covariance structure of $X^T X$

$$X^T X \propto \begin{bmatrix} 1 & -0.9 \\ -0.9 & 1 \end{bmatrix}$$

$c^T = [1\ 0]$      $\rightarrow\ \varepsilon = 0.19$

$c^T = [1\ 1]$      $\rightarrow\ \varepsilon = 0.05$

$c^T = [1\ \text{-}1]$      $\rightarrow\ \varepsilon = 0.95$

# Bonus material: Example: working memory

| A | B | C |
|---|---|---|
| Stimulus   Response | Stimulus   Response | Stimulus   Response |

Time (s)

Correlation = -.65
Efficiency ([1 0]) = 29

Correlation = +.33
Efficiency ([1 0]) = 40

Correlation = -.24
Efficiency ([1 0]) = 47

- A: Response follows each stimulus with (short) fixed delay.
- B: Jittering the delay between stimuli and responses.
- C: Requiring a response only for half of all trials (randomly chosen).

# Bibliography

- Friston KJ et al. (2007) *Statistical Parametric Mapping: The Analysis of Functional Brain Images*. Elsevier.

- Christensen R (1996) *Plane Answers to Complex Questions: The Theory of Linear Models*. Springer.

- Friston KJ et al. (1995) Statistical parametric maps in functional imaging: a general linear approach. *Human Brain Mapping* 2: 189-210.

# Thank you