



### PROJECT: WHY & WHAT FOR?

Genome-Wide Association Studies (GWAS) = test the statistical association between the GV and the phenotype in a regression model

#### Family-based GWAS

$$y_{ij} = b_0 + b_1 * g_{ij} + \epsilon_{ij}$$

where i is indicator of family and j is subjects within families.  
y, b, g and  $\epsilon$  are vectors (n = number of phenotypes within family)

$$X = \begin{pmatrix} 1 & g_1 \\ 1 & g_2 \\ \vdots & \vdots \\ 1 & g_N \end{pmatrix} \quad b = \begin{pmatrix} b_0 \\ b_1 \end{pmatrix} \quad y = \begin{pmatrix} ph_1 \\ ph_2 \\ \vdots \\ ph_N \end{pmatrix}$$

Statistical Power - paramount in GWAS for:

- small effect genes: < 1% explained variance
- up to 6 million tests  $\rightarrow$  adapted  $\alpha = 10^{-8}$

**Aim:** Increase power by refinement of statistical methodologies and meta-analyses  
Retain computational speed

### SANDWICH CORRECTED SE

**Background:** Relatives resemble each other because they share genes (A) and environment (C). Resemblance is expressed in:

THE FAMILIAL COVARIANCE MATRIX  $V$

$$\epsilon | X \sim N(0, V)$$

$$V(\Theta)$$

$$\Theta = [\sigma^2_A, \sigma^2_C, \sigma^2_E]$$

What model for  $V$  is most powerful *and* fast?

**Methods:** Use simulations to compare the standard and sandwich corrected Unweighted Least Squares (ULS) and Maximum Likelihood (ML).

#### SIMULATION ACE trait 4-sib family

	ML <sub>standard</sub> ACE <sub>model</sub> (true)	Sandwich corrected ML (false: AE Model)	Sandwich corrected ML (false: CE model)	Sandwich corrected ULS (false: E model)
mean(b1)	-0.142	-0.142	-0.142	-0.142
mean (st.err.)	0.023	0.024	0.024	0.031
mean (t-value)	-6.03	-5.98	-5.98	-4.65
<b>power</b>	<b>75.7</b>	<b>74.2</b>	<b>74.2</b>	<b>25.1</b>

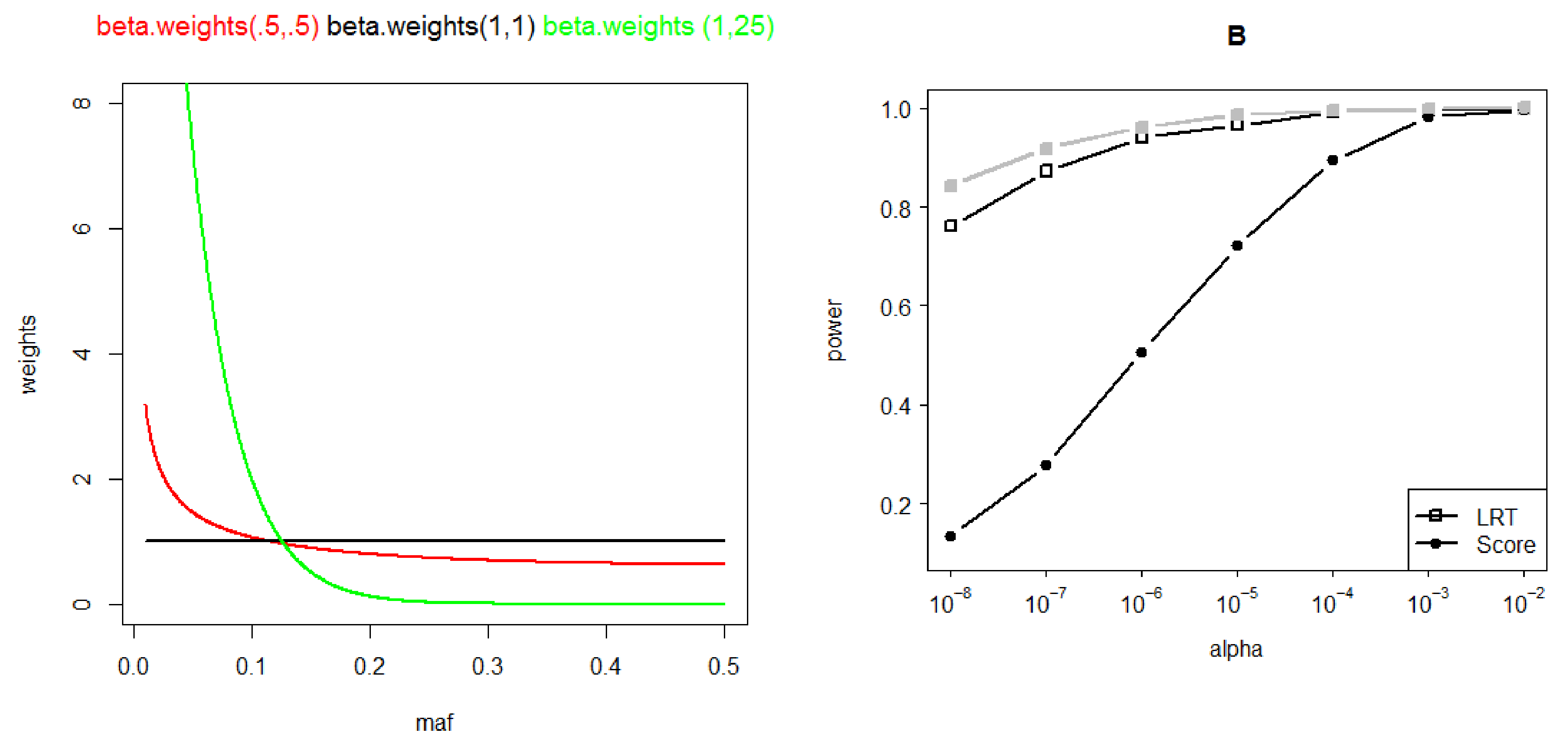
**Conclusion:** Model  $V$  as an AE or a CE & use ML with a SANDWICH!

### THE WEIGHTING IS THE HARDEST PART

**Background:** SKAT - important rare variants (RV) test based on a random effects model.  
Weights assigned to capture the likelihood of a RV being functional.  
Correct weighting increases power and yet correct weights are not known.  
What is the effect of weight misspecification in SKAT?

**Methods:** Compare LRT and score test under weight misspecification using simulations.

**Figure:** LEFT: Weights assigned based on frequency (maf)  
RIGHT: Simulated weights: beta.weights(1,1), Fitted weights: beta.weights(.5,.5).



**Conclusion:** LRT is more robust and powerful than score under weight misspecification. This is a paramount result, as misspecified models are likely to be the rule rather than the exception.

### MZ TWINS OR MZ SINGLETONS?

**Background:** Occasionally in family-based GWAS, including monozygotic (MZ) twins, the data from one MZ twin is dropped, thus reducing the MZ pairs to singletons.

Is this practice optimal?

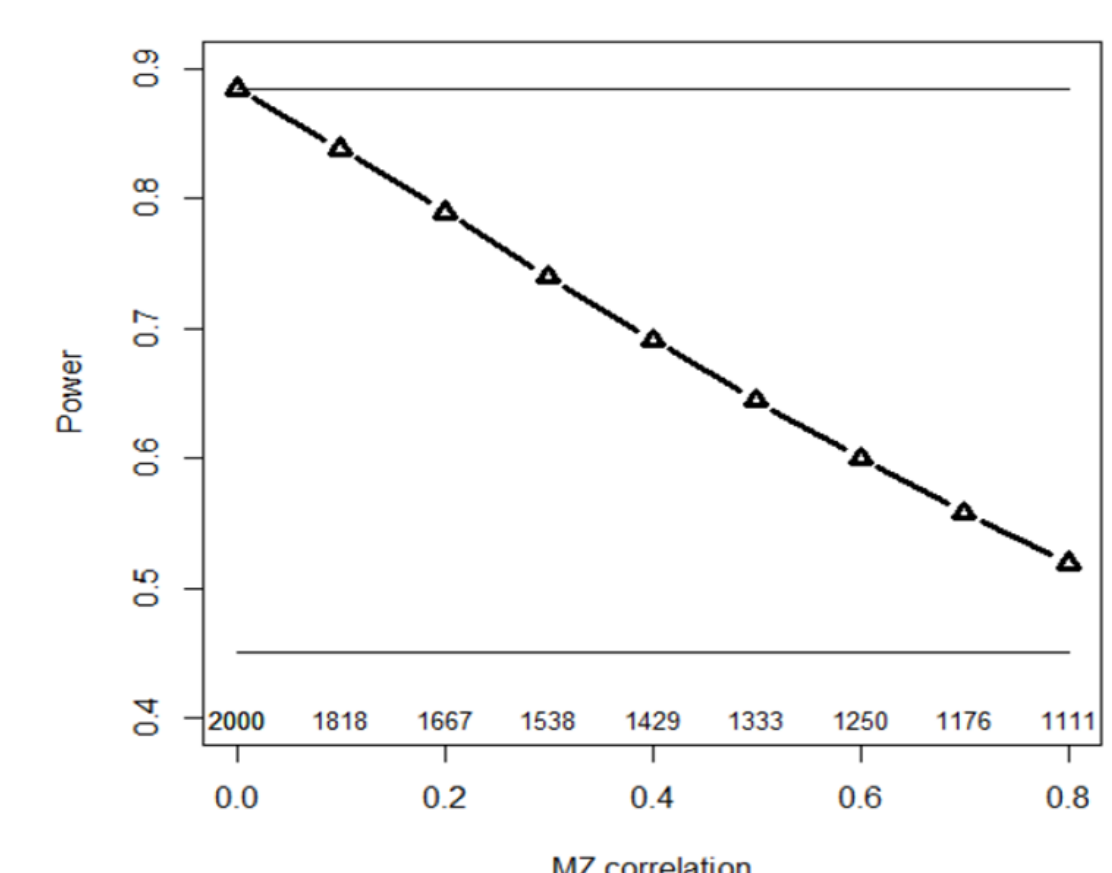
Compute effective sample size:

$$N_E = (2*N) / (1+r)$$

intraclass correlation

ranges from N (r=1) to 2\*N (r=0)

~~MZ pairs or MZ singletons?~~

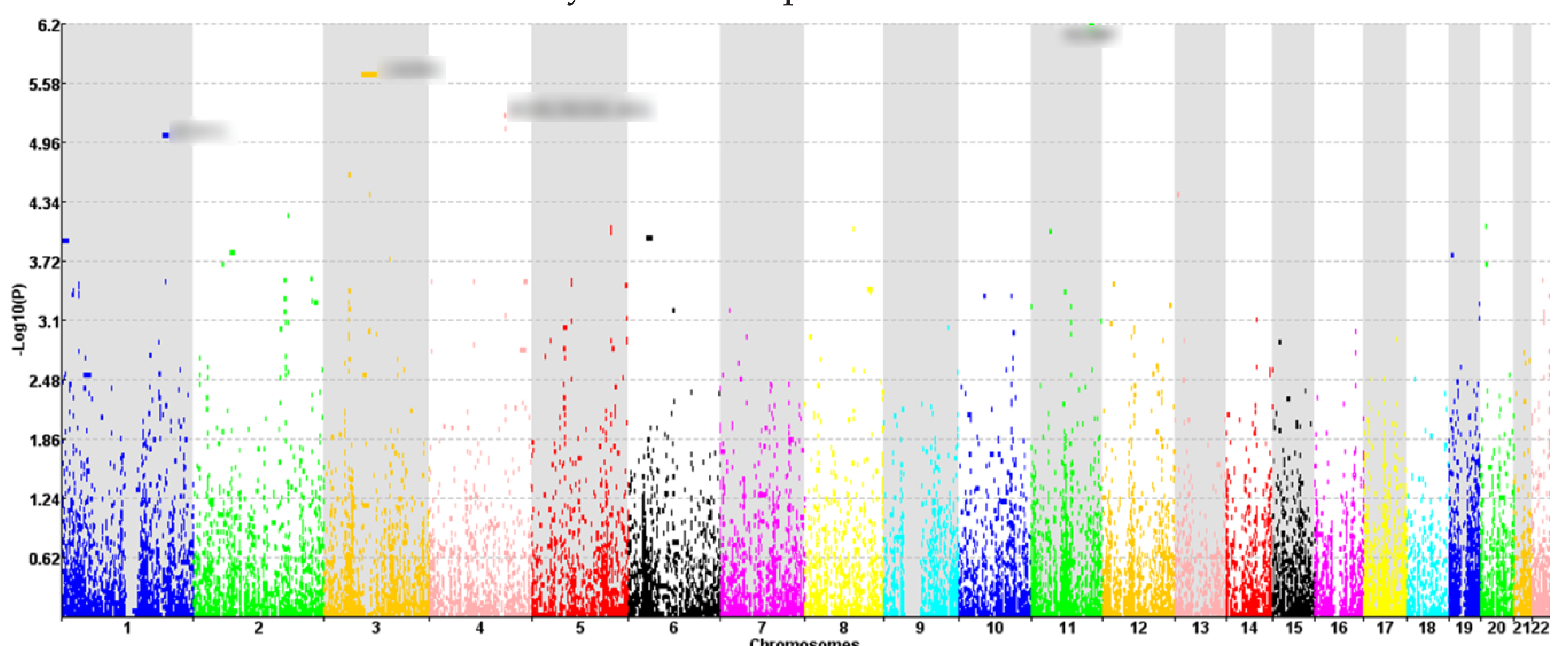


**Conclusion:** the presence of MZ twin pairs does not affect the type I error rate, and reducing MZ pairs to singletons reduces power.

### 5 GENES IMPLICATED IN CANNABIS USE: A META-ANALYSIS

**Background:** Regular cannabis use has been associated with health problems (mood and anxiety disorders) and predicts diminished educational and professional attainment.

**Methods:** Fixed effects meta-analysis in a sample >32.000 individuals.



**Implications:** One can start building a road map for developing drugs to treat cannabis dependence and abuse.