

# Differential expression analysis of time-series RNA-seq experiments

## Part II

Joyce Hsiao and Lauren Blake  
Gilad and Lynch lab meeting  
August 17, 2016

# Motivating study: “iPSC differentiation into endoderm”



4 chimpanzees, all replicated once



6 Humans, 2 are replicated once

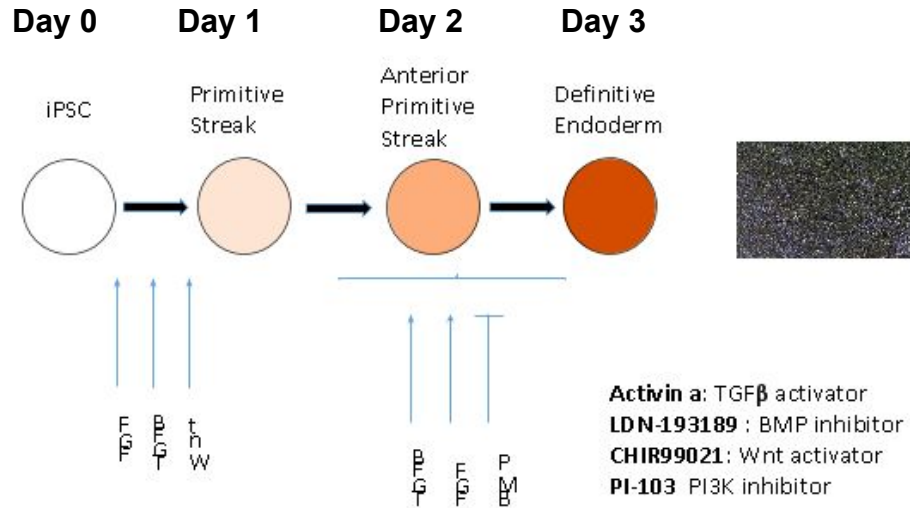


Figure Credit: Sammy Thomas

## Our big questions:

1. For each species, identify significant gene expression change since Day 0.
2. Between species, identify significant species difference in gene expression pattern during the experiment.
3. Among genes with significant species difference, characteristics and groupings of expression curves.

# Modeling time series gene expression data

**Functional data analysis (FDA)** models individual trajectories as a sample of random functions, which are not parametrically specified (Muller, 2006). Examples of parametric trend are linear or quadratic pattern, and an example of non-parametric trend is piecewise linear trend. For example, for gene  $i$  and individual  $j$  at time  $k$ ,

$$y_{ijk} = \mu_i(t_{jk}) + \epsilon_{ijk}$$

Two methods to represent these functions

1. Basis expansion methods: represent the mean function as a linear combination of basis function. An example is polynomial function.

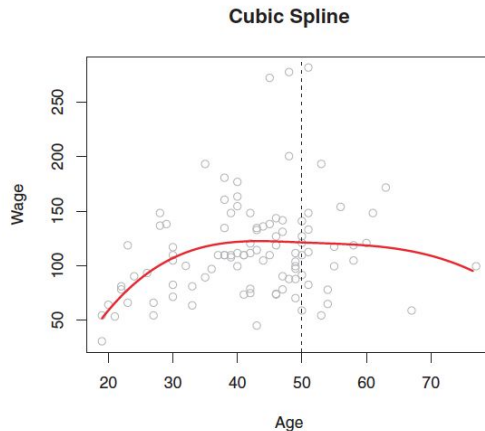
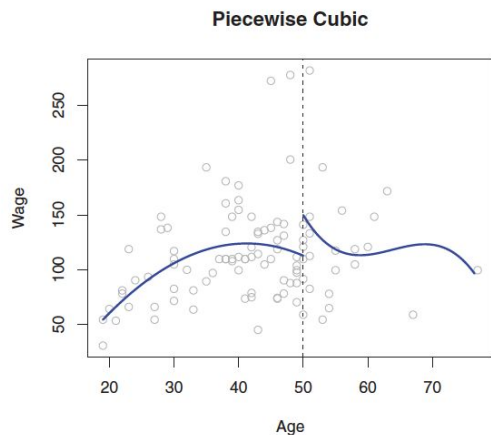
$$s(t) = \sum_{k=1}^P \beta_k s_k(t)$$

2. Smoothing splines: represent the mean function as a non-parametric function with smoothness constraint to control the shape of the curve.

$$\sum_{j=1}^N (y_{ij} - g_i(t_j))^2 + \lambda \int g_i''(t)^2 dt$$

# EDGE (Extracting Differential Gene Expression)

Storey et al, 2005: For every gene  $i$  and individual  $j$ , model each time series using a natural cubic spline basis, specifically fitting a piecewise third-order polynomial with continuous second-order derivative at the knots.

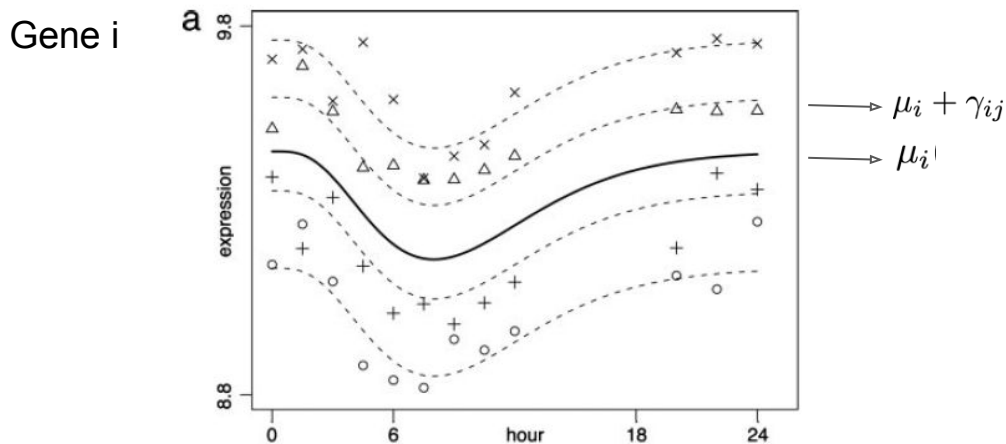


# Introducing the EDGE model

Storey et al, 2005: For every gene  $i$  and individual  $j$  at time point  $k$ ,

$$y_{ijk} = \mu_i(t_{jk}) + \gamma_{ij} + \epsilon_{ijk}$$

Individual random variation is model by  $\gamma_{ij}$ , with zero mean and gene-dependent variance.



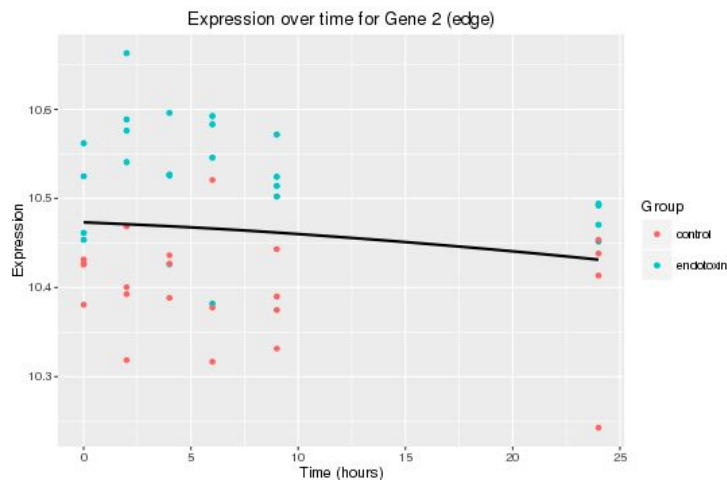
# Key analytical concerns

- Choosing shape of the fitted curve
  - How many degrees of polynomial?
  - Same or different for every gene?
- Stating the null hypothesis
  - Consider changes within a biological condition and changes between biological conditions
  - Within a biological condition
    - Compare to baseline or the average
    - Is there a control condition?
  - Between biological conditions
    - Condition-specific variation?
- Testing the null hypothesis
  - Bootstrapping or not

# Hypothesis testing: comparing expression trajectories

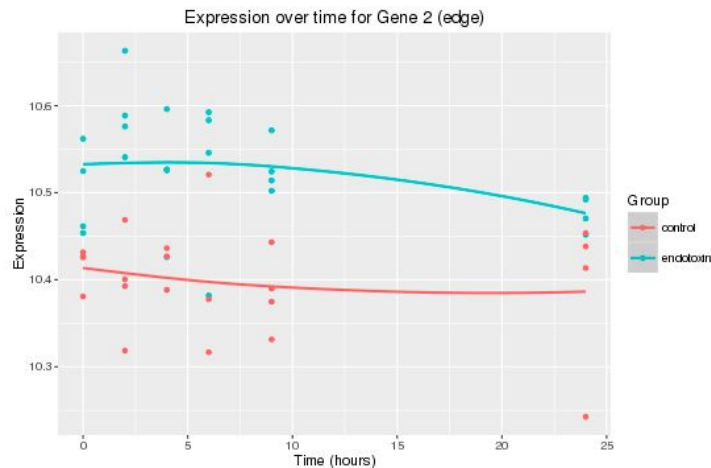
Under null hypothesis:  
Single trajectory across time  
(no group difference)

$$\hat{y}_{ij}^0 = \sum_{k=1}^p \hat{\beta}_{ik}^0 s(t_j)$$



Under alternative hypothesis:  
Group-specific trajectories

$$\hat{y}_{ij}^H = \sum_{k=1}^p \hat{\beta}_{ik}^H s(t_j) \quad \hat{y}_{ij}^C = \sum_{k=1}^p \hat{\beta}_{ik}^C s(t_j)$$



# Computing test statistic

1. Compute sum of squares of the residuals under and under alternative.

$$SS_i^{null} = \sum_{j=1}^N [y_{ij} - \hat{y}_{ij}^0]^2 \quad SS_i^{alt} = \sum_{j=1}^N [y_{ij} - \hat{y}_{ij}^H]^2 + \sum_{j=1}^N [y_{ij} - \hat{y}_{ij}^C]^2$$

2. Construct a statistic that is proportional to the typical F statistic.

$$F_i = \frac{SS_i^{null} - SS_i^{alt}}{SS_i^{alt}}$$



# Obtaining null distribution of the test statistic

## Bootstrap method

1. Compute residuals from the alternative model fit, so for human samples, we get within-individual residuals  $\epsilon_{ijk}^{\hat{}} = y_{ijk} - \hat{y}_{ijk}^H$  and between-individual residuals  $\hat{\gamma}_{ij}^H$
2. Make bootstrapped data  $y_{ijk}^* = (\hat{y}_{ijk}^0 - \hat{\gamma}_{ij}^0) + \epsilon_i^* + \gamma_i^*$   
where  $\epsilon_i^*$  and  $\gamma_i^*$  are randomly sampled with replacement from fitted residuals.
3. For each of the B iterations, compute F null statistic.

## Choosing p

$$\mu_i(t) = \sum_{k=1}^p \beta_{ik} s_k(t)$$
$$\beta_{i1} s_1(t) + \beta_{i2} s_2(t) + \cdots + \beta_{ip} s_p(t)$$

p is determined by the number of regions, degree of polynomial and the number of constraints at the endpoints. For natural cubic splines over 3 regions,

$$p = 3 \text{ regions} \times 3 \times 2 \text{ constraints} = 12$$

1. If p is too large, then we lose power because degrees of freedom are wasted.
2. If p is too small, then power is lost because expression is not properly modeled.
3. Varying p by genes resulted in over-fitting the data and artificially inflating the significance.

# EDGE versus maSigPro

	EDGE	maSigPro
Input measurement	Continuous	Continuous
Curve function	Natural cubic splines	Polynomial
Null hypothesis	Same trajectory across groups	Same trajectory across groups
Allow individual replicates	No	No
Test statistic	Proportion to F	F
Bootstrapped null distribution	Yes	No

Note. Next maSigPro is a RNA-seq version of maSigPro, which assumes negative binomial distribution of the data.

# Edge and maSigPro pipeline

<https://jhsiao999.github.io/diffTimeExpression/analysis/>

diffTimeExpression

Home

About

License

GitHub

## Home

- [Documents](#)
- [Practical examples](#)

Last updated: 2016-08-15

Code version: 09e9a3d1877cbebb88bbfb58c7cdb7a436b99ea

## Documents

- Notes
  - [Selected papers on statistical and analytical issues](#)
  - [Selected RNA-seq studies of time course gene expression](#)
- Presentations
  - Introduction to differential expression analysis of RNA-seq time-series experiments. [download](#)
  - Differential expression analysis of RNA-seq time-series experiments Part II: functional data analysis. [upcoming](#)

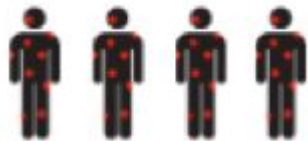
## Practical examples

- [EDGE](#)
- [maSigPro](#)
- [Comparison of EDGE and maSigPro](#)

# Running edge and maSigPro on the same dataset

Interested in differences bet. 2 groups; Study design (from Calvano et al. 2005)

Cases w/ endotoxin



Controls w/ placebo



Blood draws  
at time 0  
(before  
infusion), 2,  
4, 6, 9, and  
24 hours

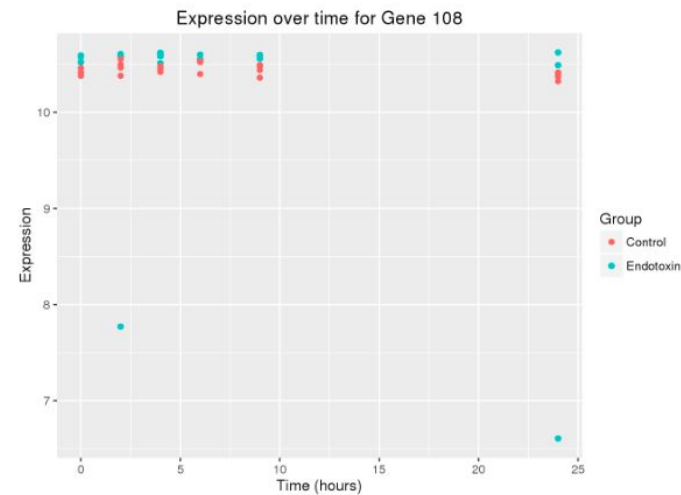
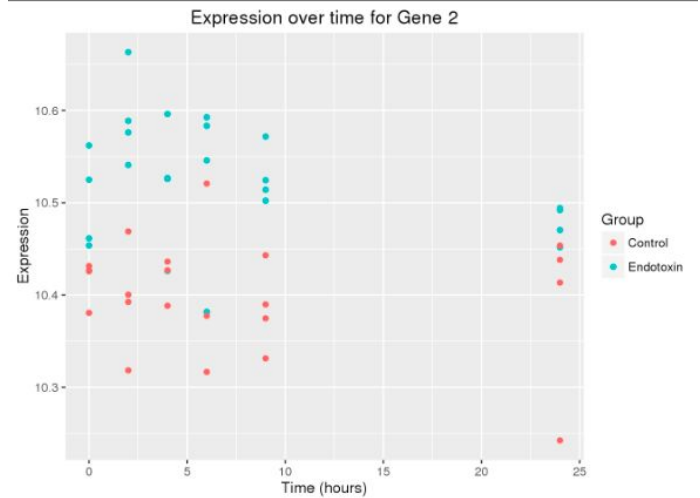
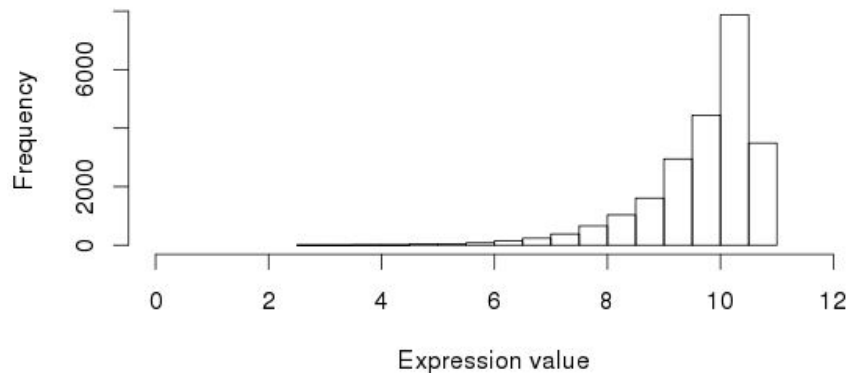


Leukocyte  
isolation /  
RNA-Seq



Normalized  
expression  
data for 500  
genes  
Cov: ind.,  
time, group.

### Expression values in the endotoxin data set (n = 23,000)



# Model comparison

**EDGE:** For every gene, person  $j$ , time  $k$ ,

$$H_o : \exp_{jk} = \alpha_j \times \text{group} + \beta_j * \text{spline basis} + \epsilon_{jk}$$

$$H_A : \exp_{jk} = \alpha_j \times \text{group} + \beta_j * \text{spline basis} + \gamma_j * \beta_j * \text{spline basis (interaction)} + \epsilon_{jk}$$

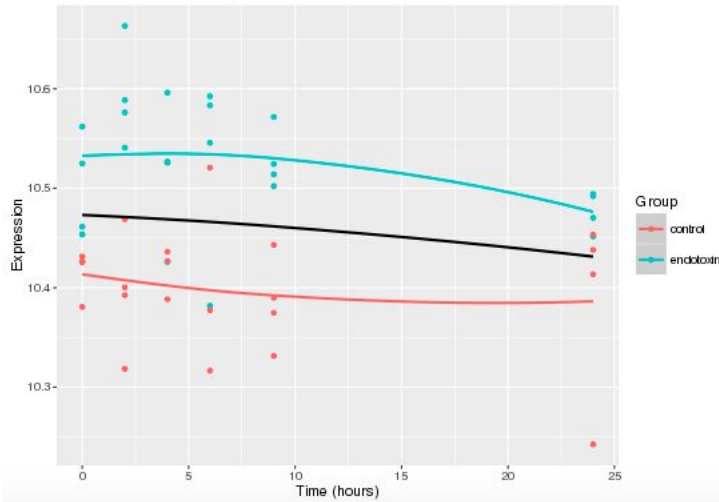
**maSigPro:** For every gene person  $j$ , time  $k$ ,

$$H_o : \mathbf{u}_o + \mathbf{u}_{1j}t + \mathbf{u}_{2j}t^2 + \epsilon_{jk}$$

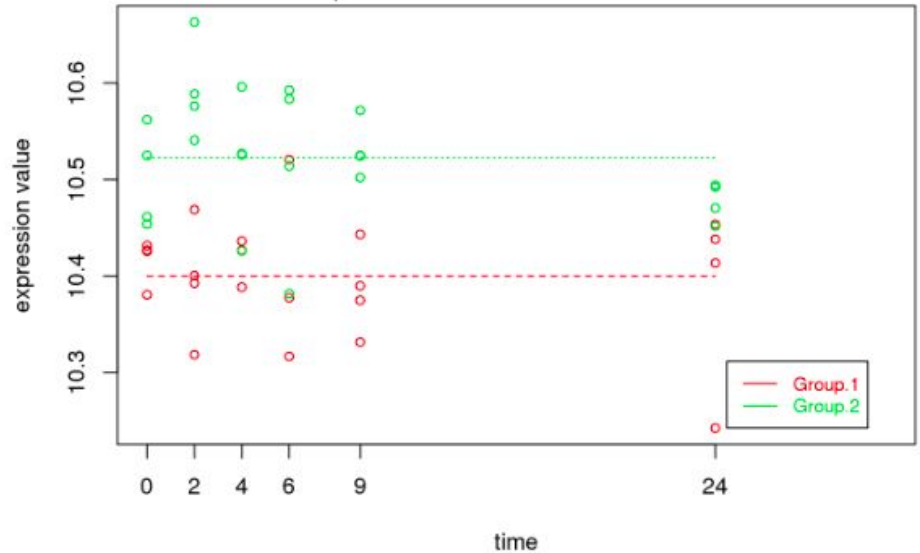
$$H_A : \exp_{jk} = \beta_o + \beta_{1j}t + \beta_{2j}t^2 + \mathbf{d}_o + \mathbf{d}_{1j}t + \mathbf{d}_{2j}t^2 + \epsilon_{jk}$$

# Comparison of results

Expression over time for Gene 2



Expression over time for Gene 2



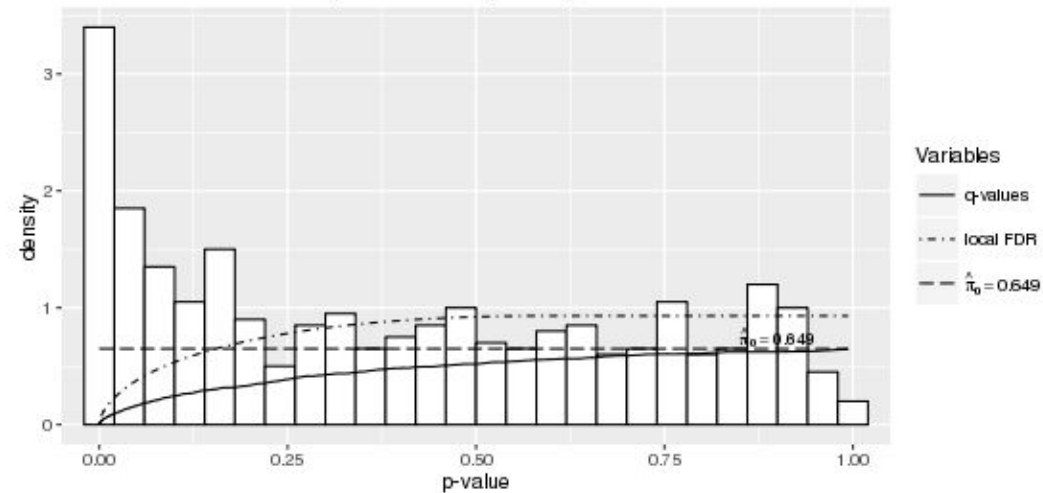


# Edge

Cumulative number of significant calls:

	<1e-04	<0.001	<0.01	<0.025	<0.05	<0.1	<1
p-value	14	24	51	75	100	132	500
q-value	0	0	22	25	42	72	500
local FDR	0	0	14	21	24	35	500

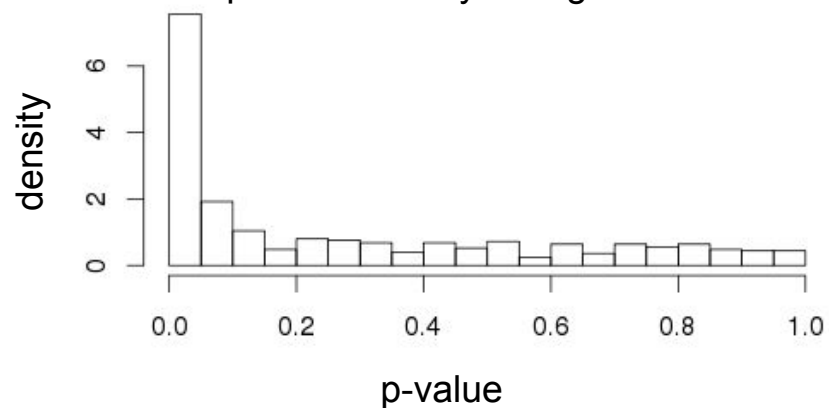
p-value density histogram



# maSigPro

	<1e-4	<1e3	<0.01	<0.025	<0.05	<0.1	<1
p-val	61	96	142	162	189	237	500

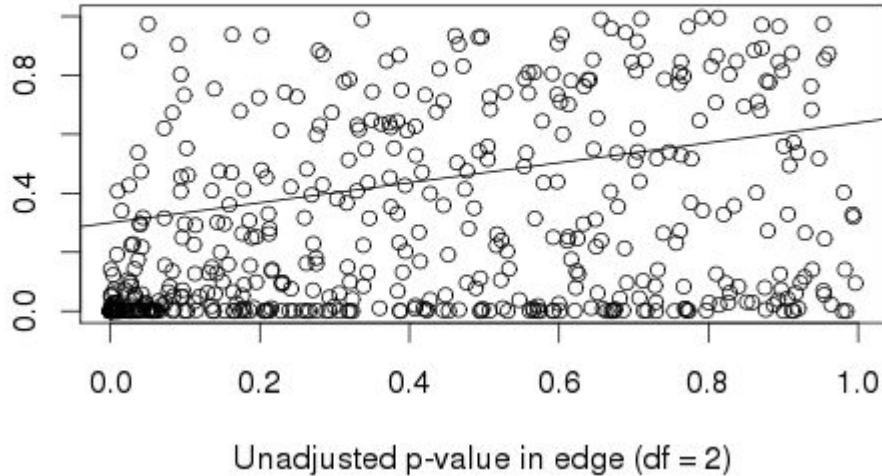
p-value density histogram



# Comparison of unadjusted p-values in 2 programs

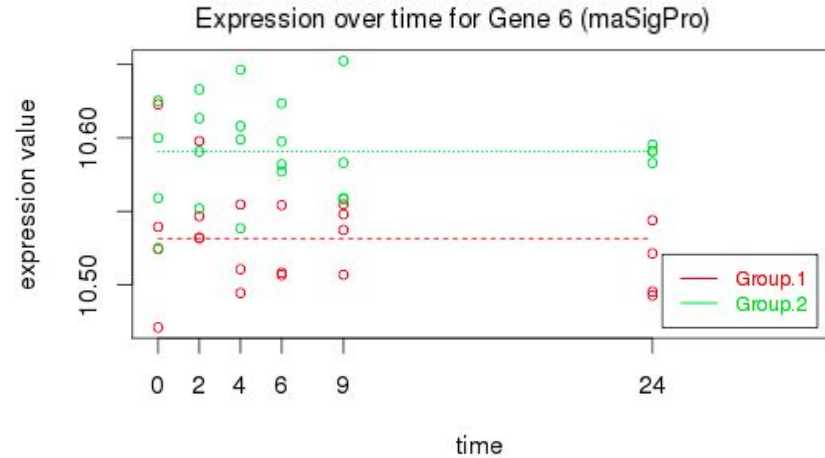
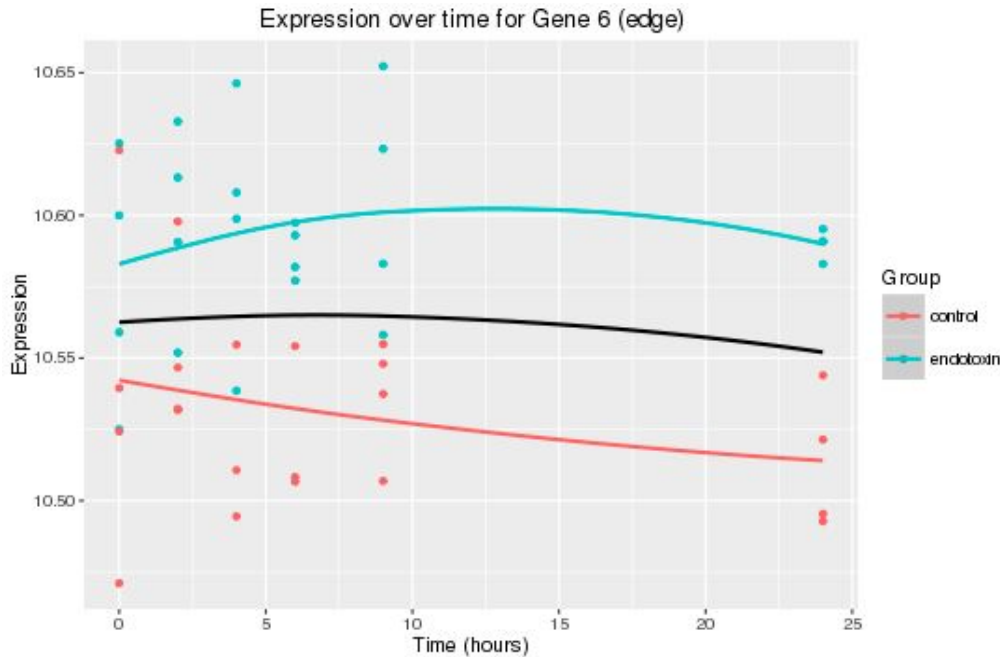
Unadjusted p-value in maSigPro (df = 2)

**Unadjusted p-values (Pearson's corr. = 0.341)**

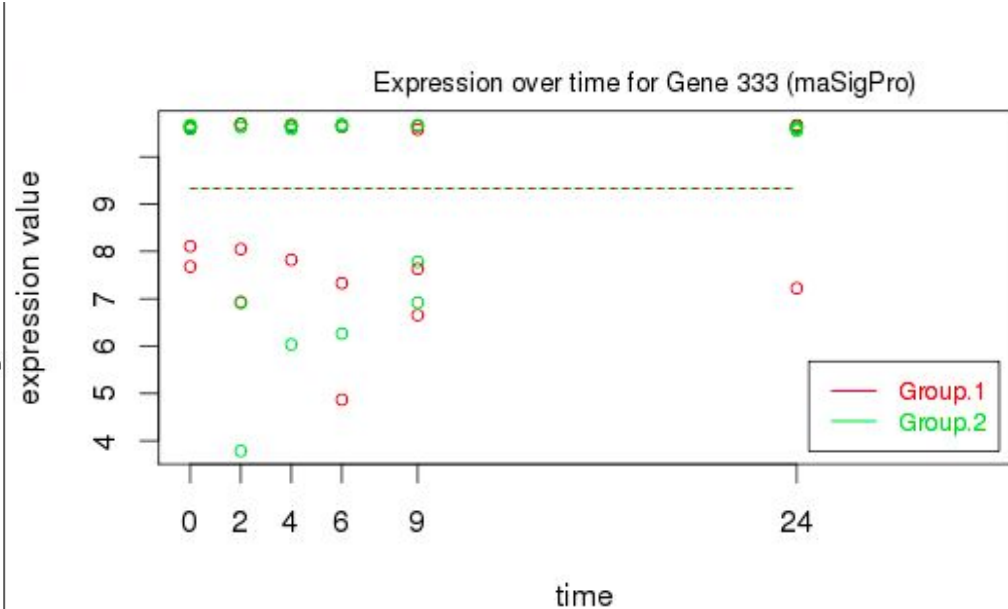
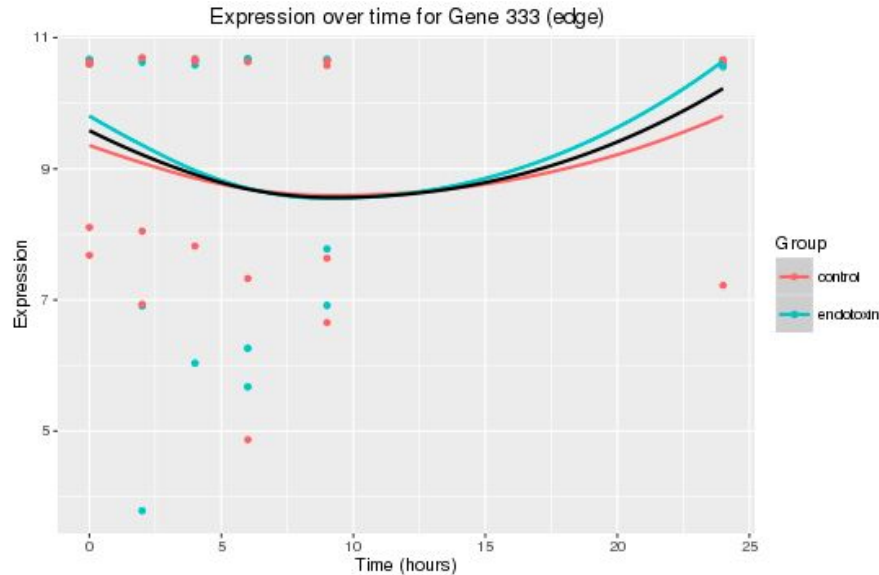


Program	unadj_pvalues_0.01	unadj_pvalues_0.1
edge	53	131
maSigPro	142	237
Genes in common	38	100

When maSigPro calls a gene significant (adj. p-val. =  $6.2 \cdot 10^{-4}$ ) and edge does not (adj. p-val = 0.90)



When edge calls a gene significant (adj. p-val. = 0.093)  
and maSigPro does not (adj. p-val. = 0.589)



# Edge Pipeline Part I

## 1) Create splines

```
# Null model: make the basis matrix for natural cubic splines
null_model <- ~grp + ns(tme, df = 2, intercept = FALSE)

# Full model: make the basis matrix for natural cubic splines
full_model <- ~grp + ns(tme, df = 2, intercept = FALSE) + (grp):ns(tme, df = 2, intercept = FALSE)

de_obj <- build_models(data = endoexpr, cov = cov, full.model = full_model, null.model = null_model)
```

## 2) Fit the null and full models using least sqs.

```
ef_obj <- fit_models(de_obj, stat.type = "odp")
```

# Edge Pipeline Part II

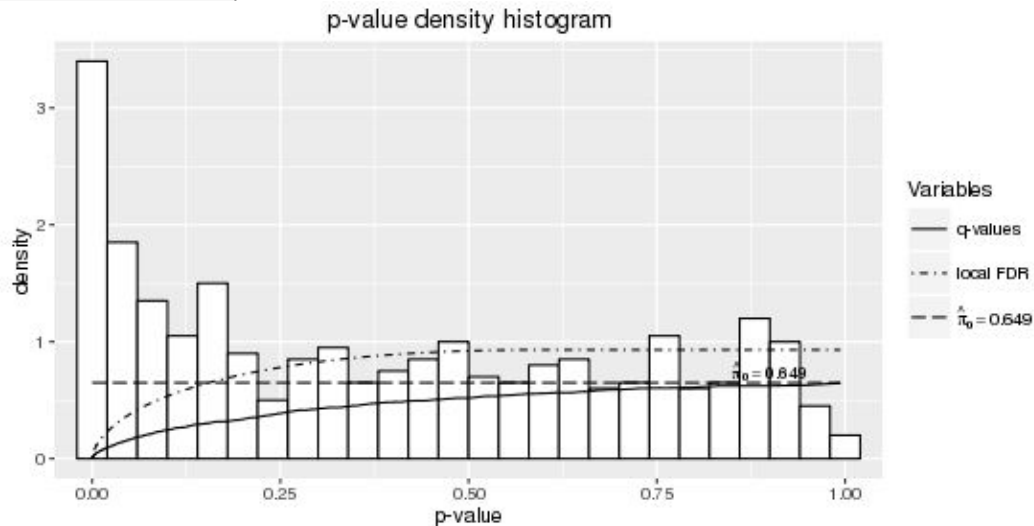
## 3) Significance testing

```
# Run odp
de_odp <- odp(de_obj, bs.its = 50, verbose = FALSE, n.mods = 50)

# See the results
summary(de_odp)
```

Cumulative number of significant calls:

	<1e-04	<0.001	<0.01	<0.025	<0.05	<0.1	<1
p-value	14	24	51	75	100	132	500
q-value	0	0	22	25	42	72	500
local FDR	0	0	14	21	24	35	500



# maSigPro Pipeline Part I

1) Make the experimental design matrix

	Time	Replicate	Group.1	Group.2
G1.T1.1	0	1	1	0
G1.T2.1	0	1	1	0
G1.T6.1	0	1	1	0
G1.T3.1	0	1	1	0
G1.T4.1	2	2	1	0
G1.T5.1	2	2	1	0

2) Make a regression matrix for the full regression model. **Discussion: degrees!**

```
matrix_endo_design <- make.design.matrix(final_endo_design, degree = 5)
```

	Group.2vsGroup.1	Time	TimeXGroup.2	Time2	Time2XGroup.2	Time3	Time3XGroup.2	Time4
G1.T1.1	0	0	0	0	0	0	0	0
G1.T2.1	0	0	0	0	0	0	0	0
G1.T6.1	0	0	0	0	0	0	0	0
G1.T3.1	0	0	0	0	0	0	0	0
G1.T4.1	0	2	0	4	0	8	0	16
G1.T5.1	0	2	0	4	0	8	0	16

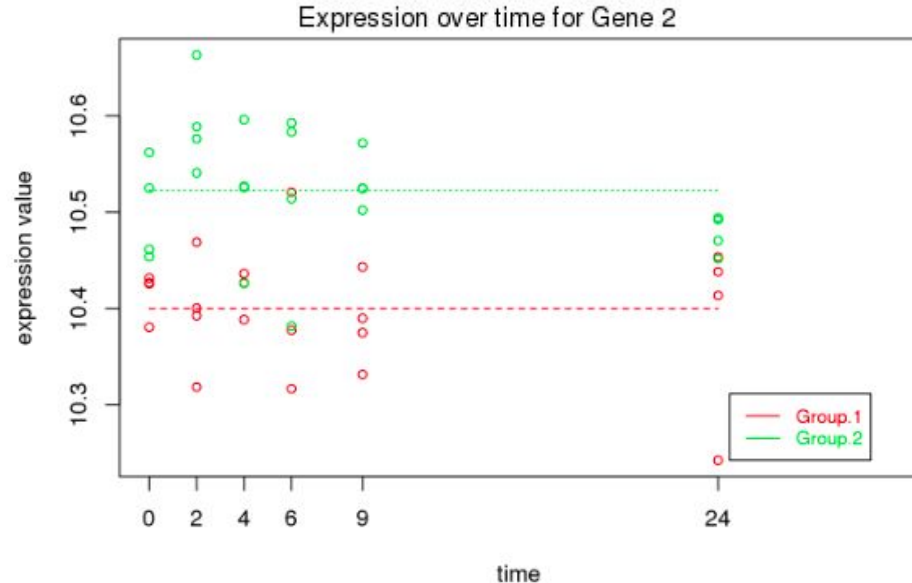
	Time4XGroup.2	Time5	Time5XGroup.2
G1.T1.1	0	0	0
G1.T2.1	0	0	0
G1.T6.1	0	0	0
G1.T3.1	0	0	0

# maSigPro Part II

3) Compute a regression fit for each gene. It will compute an unadjusted p-value and a p-value associated with the F-statistic of the model (here, FDR = 10%)

```
endo_p_vector <- p.vector(t_endo_data, matrix_endo_design, counts = TRUE, theta = 10, Q = 0.10, MT.adjust = "BH")
```

- Unadjusted and adjusted p-values for each gene
- Genes significant at FDR 10%







# maSigPro Part IV

4) For the significant genes in step 3, use forward stepwise regression to find which coefficients are statistically significant → different genes have different coefficients in the “best” model

```
endo_t_stat <- T.fit(endo_p_vector)
```

```
endo_sig_genes <- get.siggenes(endo_t_stat, rsq = 0.2, var = "groups")
```

\$group.coeffs

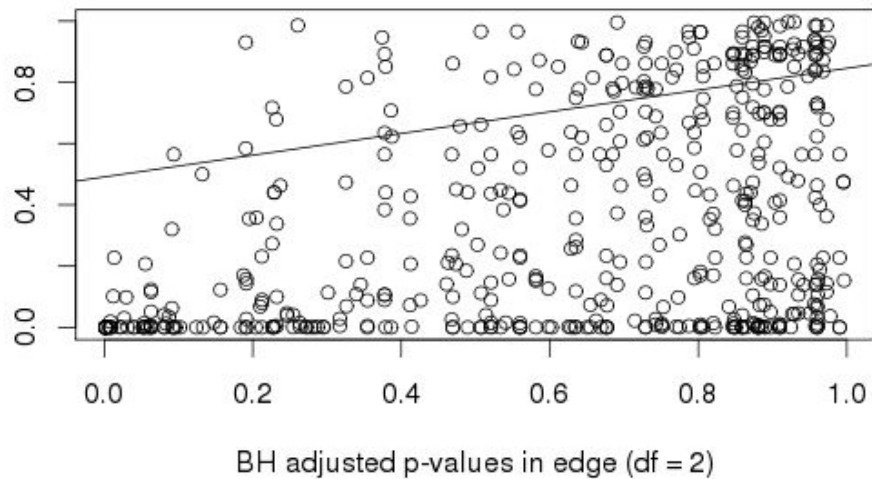
	Group.1_beta0	Group.1_beta1	Group.1_beta2	Group.2_beta0	Group.2_beta1	Group.2_beta2
2	2.341797	0.0000000000	0.000000e+00	2.353548	0.0000000000	0.000000e+00
4	2.332467	0.0000000000	0.000000e+00	2.332467	0.0046917019	-1.970172e-04
6	2.354382	0.0000000000	0.000000e+00	2.359995	0.0000000000	0.000000e+00
7	2.334190	0.0000000000	0.000000e+00	2.345320	0.0000000000	0.000000e+00
8	2.287567	0.0000000000	0.000000e+00	2.258974	0.0000000000	5.367678e-05

# References

1. Storey et al 2015
2. Muller 2006
3. Conesa 2006

BH adjusted p-values in maSigPro (df = 2)

**BH adjusted p-values (Pearson's corr. = 0.405)**



Program	unadj_pvalues_0.01	unadj_pvalues_0.1	adj_pvalues_FDR_0.01	adj_pvalues_FDR_0.1
edge	53	131	19	54
maSigPro	142	237	119	177
Genes in common	38	100	18	46