# Bioinformatics Analysis in R

## Advanced Gene Expression: Analysis of Cancer Genome Atlas

Ivan G. Costa, Martin Grasshoff

Institute for Computational Genomics
RWTH University Hospital
www.costalab.org

Institute for
Computational Genomics

RWTH AACHEN UNIVERSITY

# Summary

1. Obtain data from cancer patients from TCGA

2. Pre-process and analysis of RNA-seq data

3. Use machine learning to build a classifier for personalised medicine

4. Use interesting markers for survival analysis

Institute for
Computational Genomics
01011011010
10100100101

RWTH AACHEN UNIVERSITY

# The Cancer Genome Atlas

- TCGA  is a NCI (US) funded project to generate cohorts of cancers:
  - Currently 33 cancers with 80-780 patients
- Comprehensive data from tissues:
  - Histology, clinical, gene expression profiling, copy number variation, DNA methylation using arrays or sequencing
- Data is publicly available upon generation and deposited in a portal (portal.gdc.cancer.gov)

# The Cancer Genome Atlas - Portal

# The Cancer Genome Atlas - Portal



Check a gene or cancer type!
I will try liver ….

# LIHC - Liver Hepatocellular Carcinoma

Explore Project Data   ⬇ Biospecimen   ⬇ Clinical   ⬇ Manifest

## ⊞ Summary

| | |
|---|---|
| **Project ID** | TCGA-LIHC |
| **Project Name** | Liver Hepatocellular Carcinoma |
| **Disease Type** | Adenomas and Adenocarcinomas |
| **Primary Site** | Liver and intrahepatic bile ducts |
| **Program** | TCGA |

**CASES**
377

**FILES**
10,814

**ANNOTATIONS**
28

### Cases and File Counts by Data Category

| Data Category | Cases (n=377) | Files (n=10,814) |
|---|---|---|
| ▪ Raw Sequencing Data | 377 | 1,637 |
| ▪ Transcriptome Profiling | 376 | 2,122 |
| ▪ Simple Nucleotide Variation | 375 | 3,032 |
| ▪ Copy Number Variation | 376 | 1,536 |
| ▪ DNA Methylation | 377 | 430 |
| ▪ Clinical | 377 | 423 |
| ▪ Biospecimen | 377 | 1,634 |

### Cases and File Counts by Experimental Strategy

| Experimental Strategy | Cases (n=377) | Files (n=10,814) |
|---|---|---|
| ▪ Diagnostic Slide | 365 | 379 |
| ▪ Tissue Slide | 377 | 491 |
| ▪ WXS | 376 | 3,820 |
| ▪ RNA-Seq | 371 | 1,696 |
| ▪ miRNA-Seq | 373 | 1,275 |
| ▪ Genotyping Array | 376 | 1,536 |
| ▪ Methylation Array | 377 | 430 |

Institute for
Computational Genomics
01011011010
10100100101

RWTH AACHEN UNIVERSITY

# LIHC - Liver Hepatocellular Carcinoma



Gene expression data!

# Bioinformatics Pipeline / RNA-seq

# Bioinformatics Pipeline / RNA-seq



Practical part not covered!

# Bioinformatics Pipeline / RNA-seq

# Next Generation Sequencing

▸ **NGS take advantage of parallelization**

　▸ **reads millions/billions of reads per run**

　▸ **short reads (50-100 bps)**

　▸ **error rates (0.1-1%)**

# Read Types



Fragment DNA:

Single end

Paired end
Ins: 200-800 bp

Institute for
Computational Genomics
01011011010
10100100101

RWTH AACHEN
UNIVERSITY

# Bioinformatics Pipeline / RNA-seq

# Alignment

- a large reference sequence is given (genome)
  - up to billions of base pairs
- short reads (<200bps)
- find most probable position of the read in the genome (by inexact string matching)



100 bp

short reads

long genome

3.4 bn bp (human)

# Alignment - Split Read Mapping (RNA-Seq)



Processed mRNA

# Alignment - Split Read Mapping (RNA-Seq)



- **reads are split between exons when mapped to genome**
- **aligners use transcript information or try to find splice events (STAR & TOPHAT)**

# Reference based aligners - Overview

| | Time | Precision | Pairs | GAPs | Phred | Memory | Application (Comments) |
|---|---|---|---|---|---|---|---|
| **BOWTIE** | + | | + | - | - | 5GB | General (max. 3 missmatches) |
| **BWA** | + | | + | + | + | 8GB | General (max of 200bps reads) |
| **NOVOALIGN** | | + | + | + | + | 8GB | General (commercial license) |
| **STAR** | + | | + | - | + | 32GB | RNA-Seq (allow split-maps) |
| **BISMARK** | + | | + | + | + | 10GB | Bisulfite/reduced sequencing |

**Computers need large memory and a few hours of computation per experiment!**

RWTH AACHEN
UNIVERSITY

# Quantification (Count Matrix)



## Simple Counting Approaches

**Gene Level** - 17 reads
**Exon level** - exon 1 (8 reads), exon 2 (3 reads), exon 3 (6 reads)
**Transcript Level** - Exons 1,2 & 3 (10 reads) and exon 1 & 3 (7 reads) *
* complex computational methods required (RSe, or TopHAT needed for this)

## Fragments per Kilobase (FPKM)
- normalize counts by  read size (kb) and RNA-seq library size (mb)

# RNA-seq and Differential Analysis

## Arrays and RNA-seq have distinct distributions



**VOOM analysis is necessary to make variance similar to arrays.**

# Bioinformatics Pipeline / RNA-seq



We will see this today!

Data Size Computational Effort

Sequencing — Pre-processing — Alignment — Count Matrix → Clustering, PCA Differential Expression Survival Analysis

**Provided by TGCA or your Core Facility!**

Institute for Computational Genomics

RWTH AACHEN UNIVERSITY

# Personalized Medicine

Diagnosis and treatment choices is mostly carried on macromolecular features:

- morphology of tumours (image), symptoms, blood levels

Challenges: use molecular markers (expression or genetics) for diagnosis or treatment selection.

# Machine Learning - Classifier



Data

Expression matrix X
(genes vs samples)

classification vector *Y*
(diagnosis)

Find a function:

f(*x*) → *y*

# Machine Learning - Classifier



Gene 1

Gene 2

cancer type

- DLCL
- FL
- CLL

?

Data

Expression matrix X (genes vs samples)

classification vector $Y$ (diagnosis)

Find a function:

$f(x) \rightarrow y$

For new patients X':

$f(x') \rightarrow y'$

RWTH AACHEN UNIVERSITY

# Linear Classifier



cancer type

Gene 1

Gene 2

Linear Function:

$$f(x, A) = a_0 + a_1 x_1 + ... + a_L x_L$$

$$f(x, A) > 0 \Rightarrow \text{class A}$$

$$f(x, A) \leq 0 \Rightarrow \text{class B}$$

- Works for 2 classes only
  - Train a function for each cancer type
- Find coefficients *A*
  - estimated with neural networks or support vector machines

# Linear Classifier - Problems



- Most real world problems are not linearly separable!

- There will be always some error!

- Solution: non-linear functions

Institute for
Computational Genomics
01011011010
10100100101

RWTH AACHEN
UNIVERSITY

# Nonlinear Classifier - Problems



- Polinomial Function
- $f(x, A) = a_0 + a_{11}x^3_1 + \ldots + a_{L1}x^3_L$

$$a_{12}x^2_1 + \ldots + a_{L2}x^2_L$$

$$a_{12}x_1 + \ldots + a_{L2}x_L$$

- Third order polynomial
- Problem: overfitting

# Nonlinear Classifier - Problems



- Polinomial Function
- $f(x, A) = a_0 + a_{11}x^3_1 + \ldots + a_{L1}x^3_L$

$$a_{12}x^2_1 + \ldots + a_{L2}x^2_L$$

$$a_{12}x_1 + \ldots + a_{L2}x_L$$

- Third order polynomial
- Problem: overfitting

# Nonlinear Classifier - Problems



- Polinomial Function
- $f(x, A) = a_0 + a_{11}x_1^3 + \ldots + a_{L1}x_L^3$

$$a_{12}x_1^2 + \ldots + a_{L2}x_L^2$$

$$a_{12}x_1 + \ldots + a_{L2}x_L$$

- Third order polynomial
- Problem: overfitting

# Curse of Dimensionality

Size of a Euclidean space grows with dimension (number of genes)

Dots (patients) are sparsely distributed in space

# Curse of Dimensionality : Example

Sparse data

- three genes

- 2 patients with known cancer (red vs yellow)

- 1 unknown  (green)

# Curse of Dimensionality : Example



- Sparse data

  - three genes

  - 2 patients with known cancer (red vs yellow)

  - 1 unknown  (green)


Perfect classifier (on training)

# Curse of Dimensionality : Example



- Sparse data

  - three genes

  - 2 patients with known cancer (red vs yellow)

  - 1 unknown  (green)

Both are perfect classifiers (on training)

Hard to generalise!

# Curse of Dimensionality : Example



- There are millions of perfect linear classifiers

- And even more non-linear classifiers!

# Dealing with Curse of Dimensionality
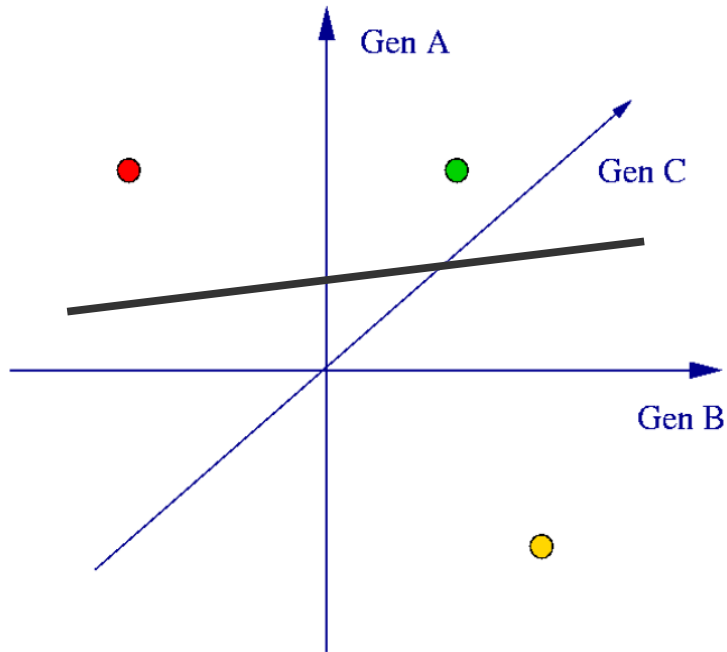
- Have a proper training / test evaluation procedure

- Use classifiers which are as simple as possible

- Reduce the dimension of your data (feature selection or PCA)

# Classifier Evaluation

**Measures for a two class problem (cancer + vs. non-cancer - )**



Source: Lever et al., Nat. Methods (2016)

# Classifier Evaluation

**Measures for a two class problem (cancer + vs. non-cancer - )**



Source: Lever et al., Nat. Methods (2016)

# Classifier Evaluation

- The performance of your classifier needs to be evaluated on your test data:

  - an independent "validation cohort"

  - or retain a set of samples (1/3) that has similar distribution of classes of your total data

$X$ → $X$ train / $X$ test

# Classifier Evaluation

- The performance of your classifier needs to be evaluated on your test data:

  - an independent "validation cohort"

  - or retain a set of samples (1/3) that has similar distribution of classes of your total data

| $X$ | → | $X$ train |
| | | $X$ test |

- Never use test data to improve classification (choose a better classifier or marker gene)

  - For this you need to establish validation data (or cross validation)

Institute for
Computational Genomics
01011011010
10100100101

RWTH AACHEN
UNIVERSITY

# Cross-validation

# Elastic Net

Is based on a linear function:

$f(x, A) = a_0 + a_1 x_1 + ... + a_L x_L$

$f(x, A) > 0 \Rightarrow$ classe A

$f(x, A) \leq 0 \Rightarrow$ classe B

- Find coefficients *A, while most of then have* 0.
  - A shrinkage factor ($\lambda$) controls the number of genes selected.
  - Shrinkage factor can be automatically identified with cross-validation.

RWTH AACHEN UNIVERSITY

# Hands on!

# Exercise (after the handout)

You should perform clustering of tissues with liver cancer. Tip: use code similar to the one seen in gene expression data (day 3). Since, we are interested in grouping patients, you can transpose the matrix with the function **t**.

1. Can you see nice clusters in the dendrogram?

2. What about genes associated to each group? Are they associated to some particular biological function? Use differential expression analysis and GO enrichment analysis to solve this task.

www.costalab.org

# Survival Analysis

Can be used to evaluate if characteristics of a patients
indicates an increase/decrease risk of survival
- clinical: tumour type, gender
- Molecular: expression of a gene, mutation

Common Survival Tests:
- Cox proportional hazards regression (not seen here)
- Compares survival with a numeric variable
- Kaplan-Meier graph / Log-rank test
- compares the survival of groups of individuals
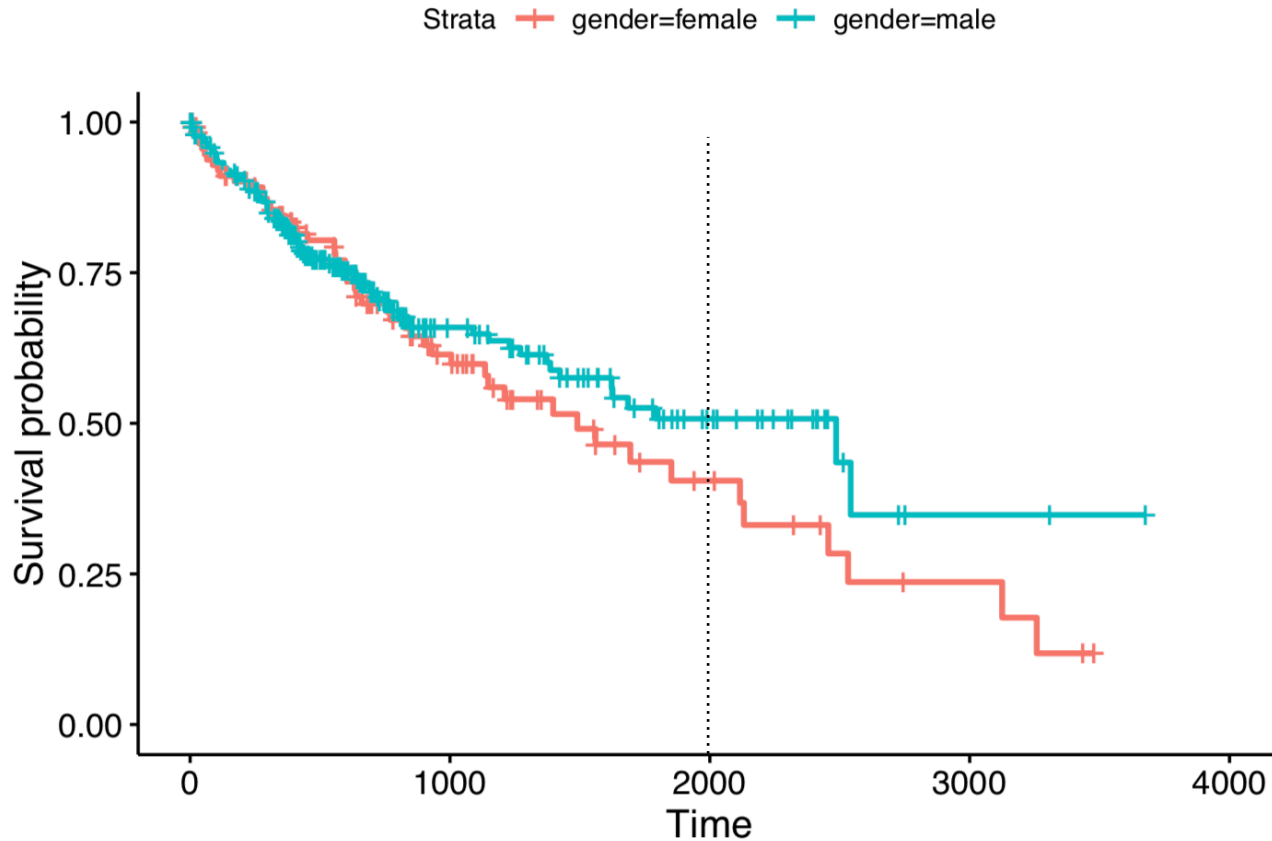
# Kaplan-Meier graph / Log-rank test

## Data:

- **Event**: death / alive
- **Time**: period between first and last observation.
- **Characteristics**: sex, tumor grade

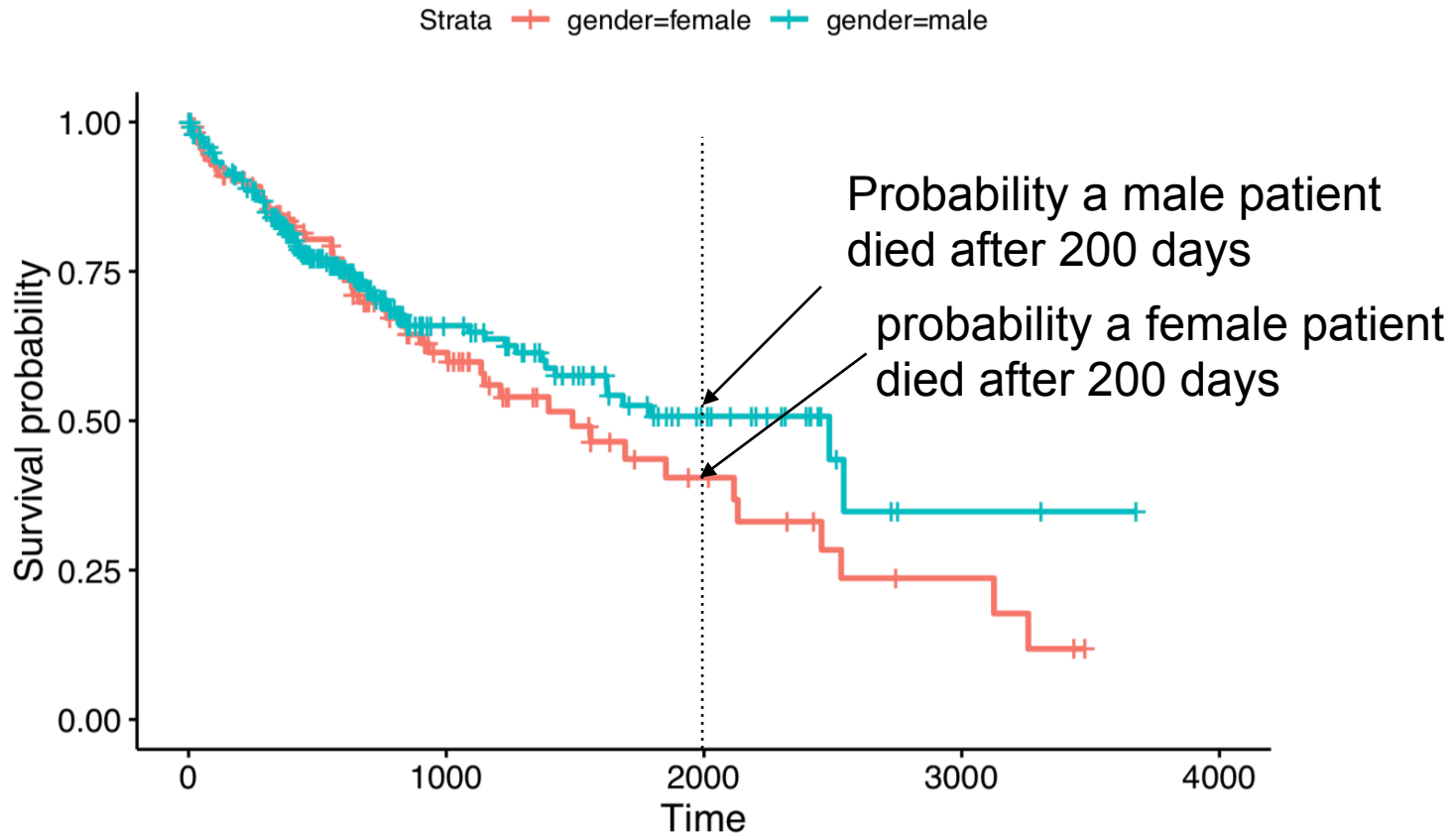| Patient | Status | Time | Sex |
|---------|--------|------|--------|
| 1 | Dead | 343 | Male |
| 2 | Alive | 20 | Male |
| 3 | Alive | 300 | Female |
| 4 | Dead | 200 | Male |

# Kaplan-Meier plot

## Survival of LIHC patients - male vs. Female
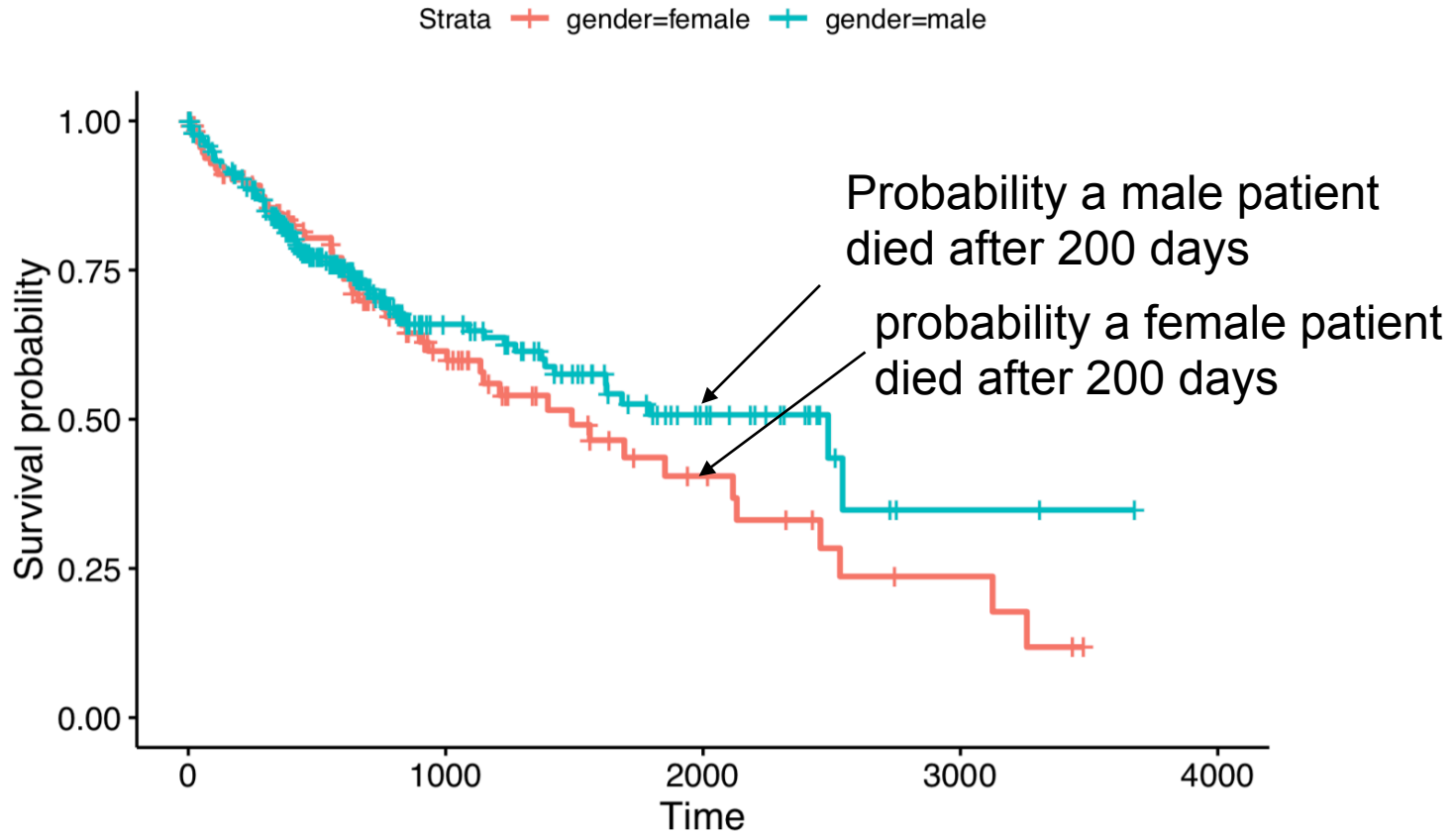
# Kaplan-Meier plot



Survival of LIHC patients - male vs. Female

Probability ( $X$ days) = $\dfrac{\text{\# cases alive after } X \text{ days}}{\text{\# cases measured after } X \text{ days}}$
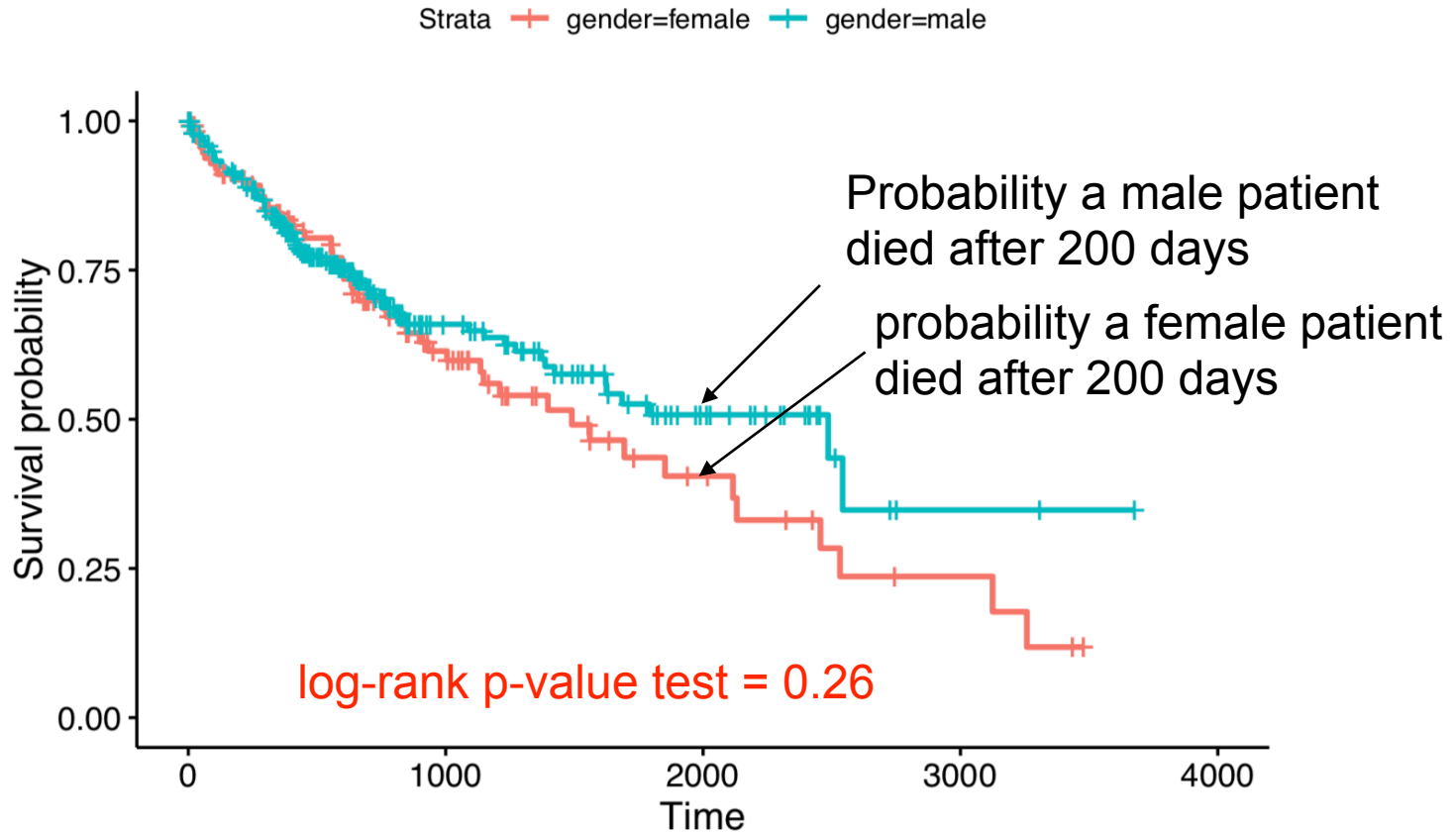
# Log-rank test



Is the survival difference significant?
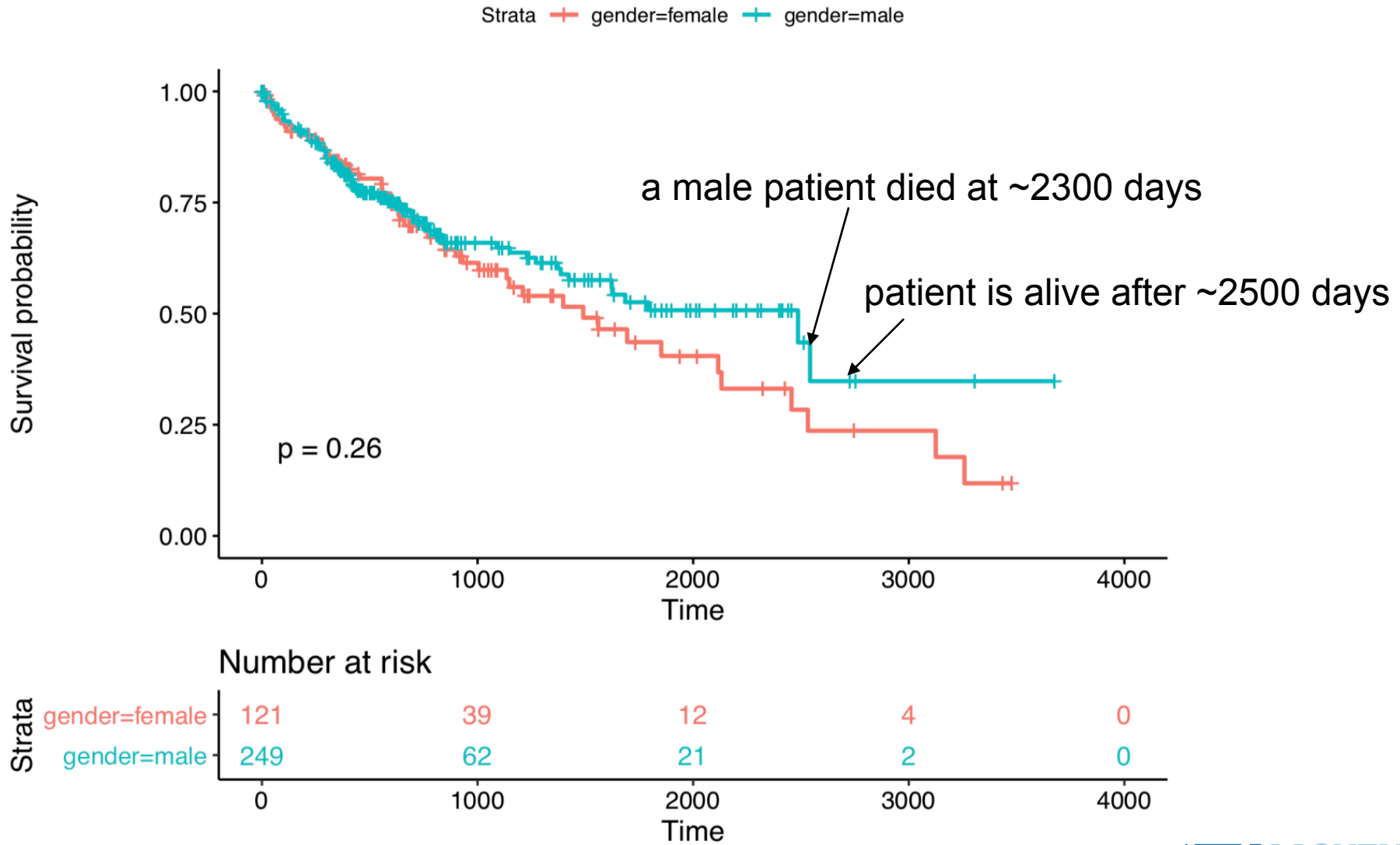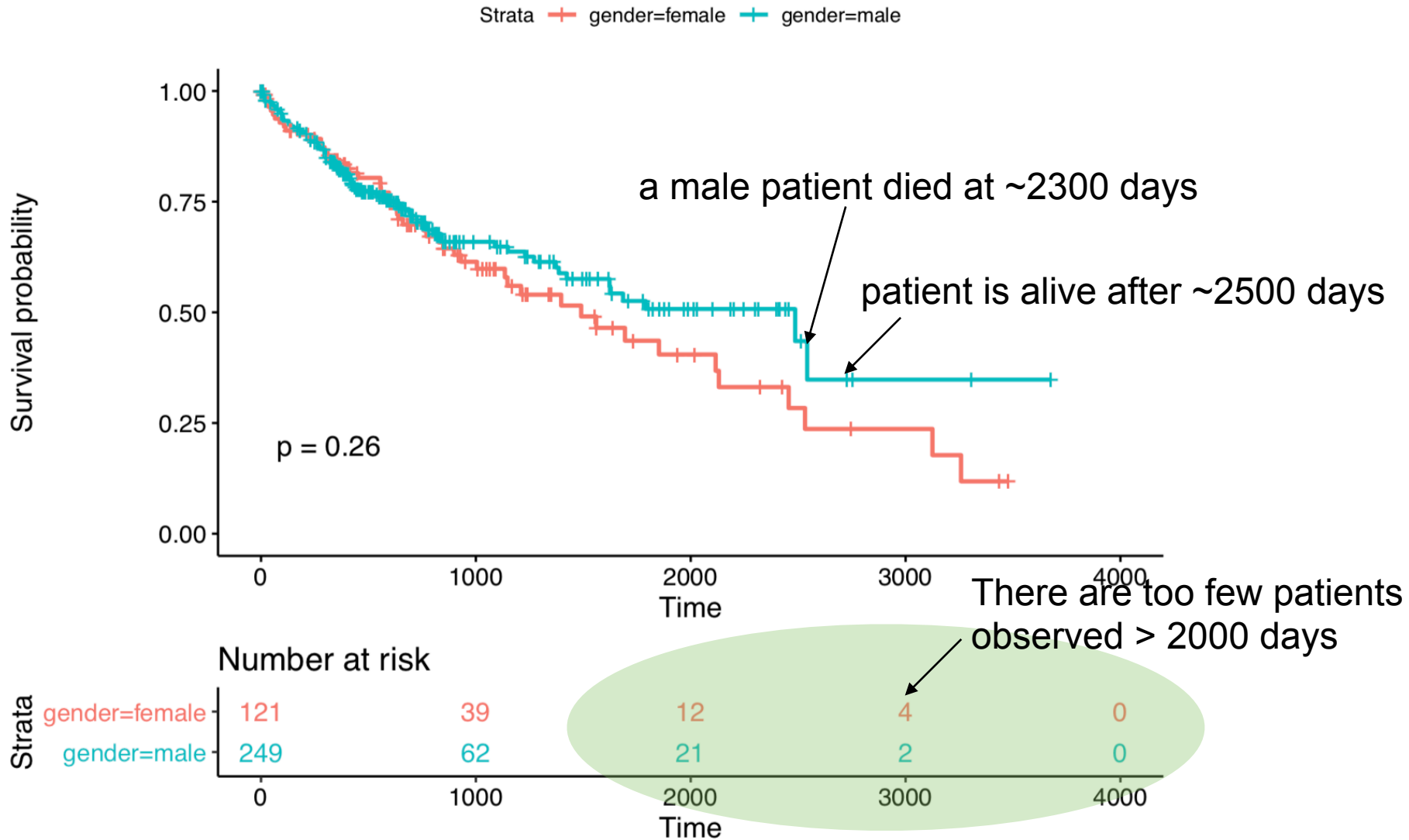
# Log-rank test



Is the survival difference significant?

Probability a male patient died after 200 days

probability a female patient died after 200 days

log-rank p-value test = 0.26

# Kaplan-Meier plot

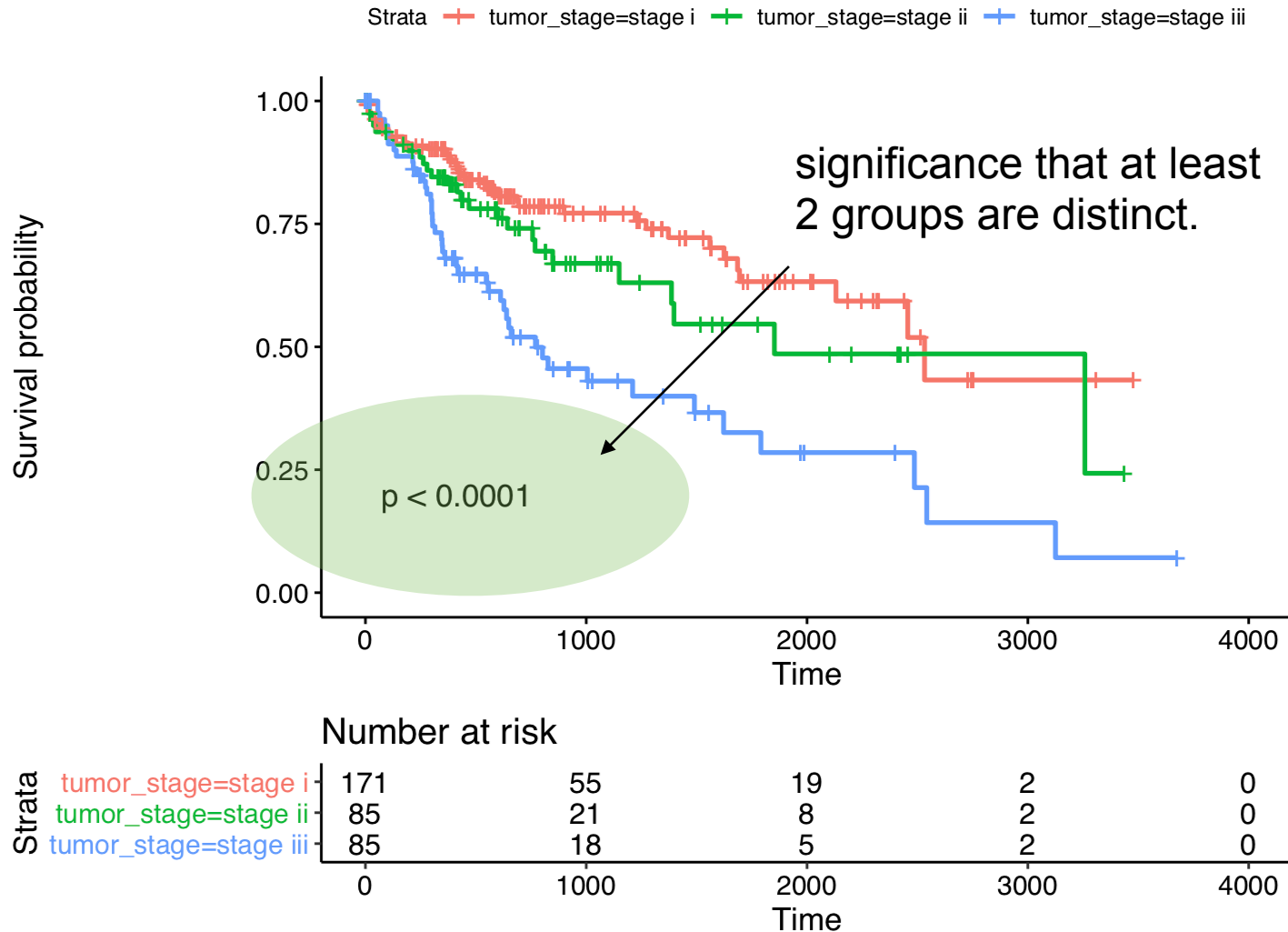# Kaplan-Meier plot

# Kaplan-Meier / Log-Rank Test

KM and LRT can compare several groups at a time.



Survival vs Tumour stage at diagnosis

significance that at least 2 groups are distinct.

p < 0.0001

Number at risk

| Strata | | | | |
|---|---|---|---|---|
| tumor_stage=stage i | 171 | 55 | 19 | 2 | 0 |
| tumor_stage=stage ii | 85 | 21 | 8 | 2 | 0 |
| tumor_stage=stage iii | 85 | 18 | 5 | 2 | 0 |

# Survival Analysis and Biological Markers

How to perform survival analysis on biological markers?

1. Given their continuous nature of gene expression, Cox hazards test is recommended.

2. An alternative is to group patients by expression of a gene (low/high expression) and use Kaplan-Meyer plots (seen in practical).

**Important: if you test several markers you need to correct for multiple testing!!!**

# Next week - Final Project

Ideas:
- Perform an analysis of a real gene expression data set
- Project can be developed in groups of 2-3 students
- Groups need to create an R code and a 10 minutes presentation showing the analysis

Schedule:

9:30 - Problem explanation
15:00 - Delivery of code and presentation slides
15:00 to 17:00 - Presentations

# Hands on!

# Hands on!