

Biostatistics

Applications in Medicine

Nuno Sepúlveda, 26.10.2023

Syllabus

1. General review

- a. What is Biostatistics?
- b. Population/Sample/Sample size
- c. Type of Data – quantitative and qualitative variables
- d. Common probability distributions
- e. Work example – Malaria in Tanzania

2. Applications in Medicine

- a. Construction and analysis of diagnostic tools – Binomial distribution, sensitivity, specificity, ROC curve, Rogal-Gladen estimator
- b. Estimation of treatment effects - generalized linear models
- c. Survival analysis - Kaplan-Meier curve, log-rank test, Cox's proportional hazards model

3. Applications in Genetics, Genomics, and other 'omics data

- a. Genetic association studies – Hardy-Weinberg test, homozygosity, minor allele frequencies, additive model, multiple testing correction
- b. Methylation association studies – M versus beta values, estimation of biological age
- c. Gene expression studies based on RNA-seq experiments – Tests based on Poisson and Negative-Binomial

4. Other Topics

- a. Estimation of Species diversity – Diversity indexes, Poisson mixture models
- b. Serological analysis – Gaussian (skew-normal) mixture models
- c. Advanced sample size and power calculations

Prevent

Diagnose

Medicine

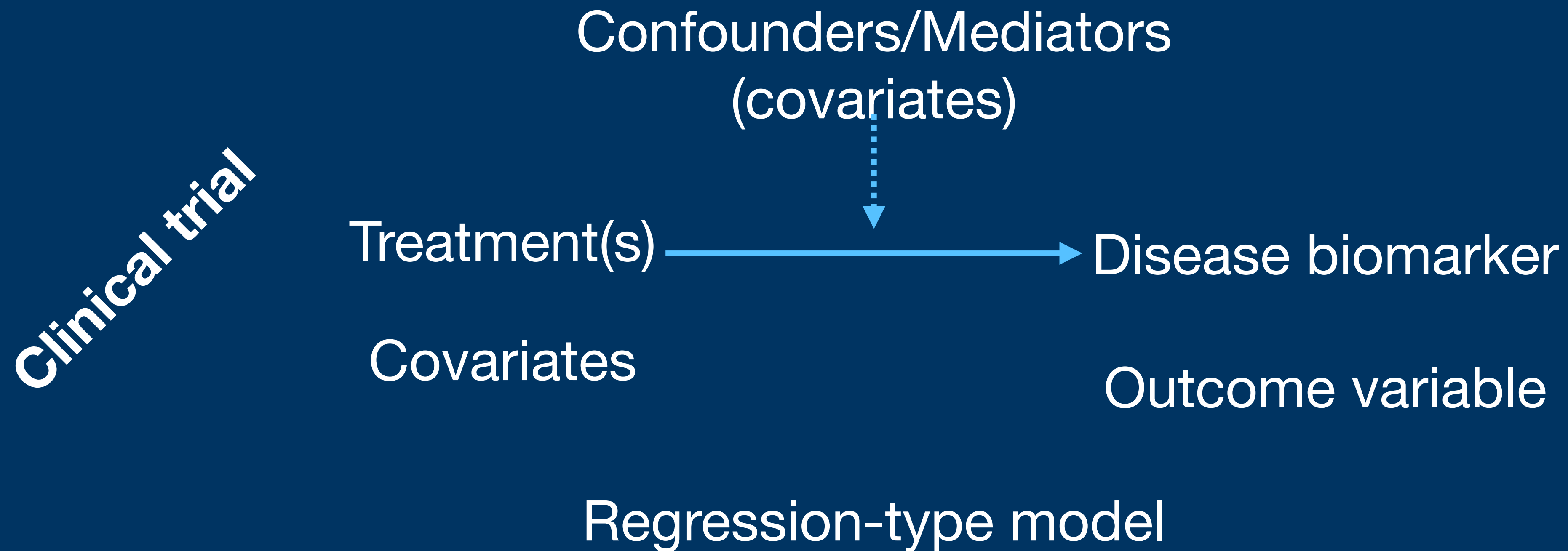
Treat

Improve

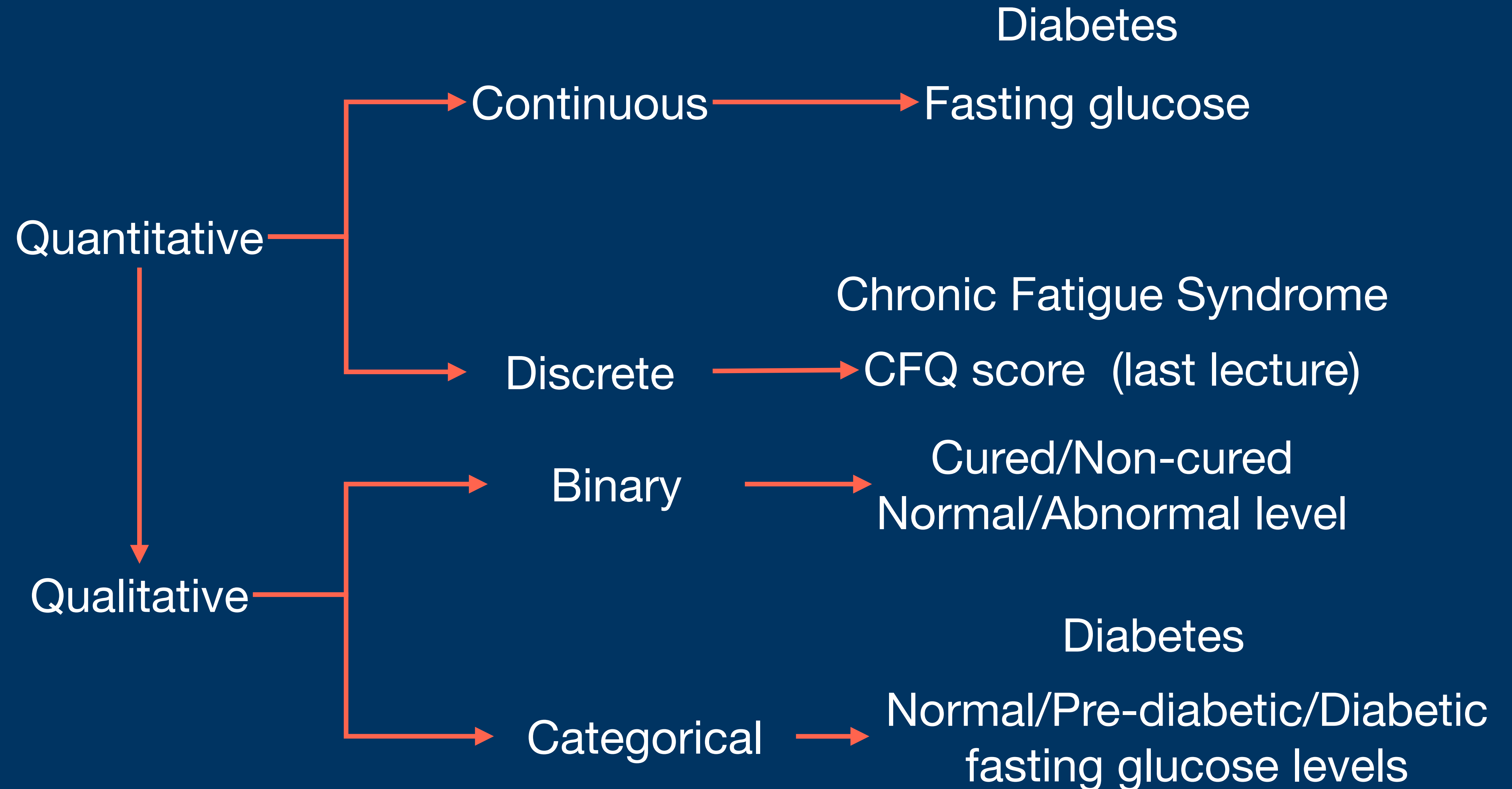
Develop

Basic question

What are the treatment effects on a disease biomarker?



Disease biomarker



Generalised linear models



Generalised linear models

$$Y | \theta \rightsquigarrow F(\theta)$$

Random component

Generalised linear models

Y_1, \dots, Y_n Outcomes

Y_i = random variable representing the biomarker value of individual i

x_{11}, \dots, x_{1p}

$\vdots \quad \vdots$ Covariates

x_{n1}, \dots, x_{np}

x_{ij} = value of covariate j of individual i

Generalised linear models

$$Y_i | \theta_i \rightsquigarrow F(\theta_i)$$

Random component

$$g(\theta_i) = \alpha + \sum_{j=1}^p \beta_j x_{ij}$$

Systematic component

$g(\cdot)$ = link function

Generalised linear models

$$Y_i | \theta_i \rightsquigarrow F(\theta_i)$$

Random component

$$g(\theta_i) = \alpha + \sum_{j=1}^p \beta_j x_{ij}$$

Systematic component

$g(\cdot)$ = link function

$F(\theta)$ should belong to the exponential family of distributions

Exponential family of distributions

$$f_{X_i}(x | \theta_i) = h(x) e^{\eta(\theta_i)T(x) - A(\theta_i)}$$

The support of the distribution does not depend on the parameter

$\eta(\cdot)$ = canonical link function

Exercise: Is Bernoulli distribution a member of exponential family?

$$f_{X_i}(x | \pi_i) = \pi_i^x (1 - \pi_i)^{1-x}$$

Generalised linear models

What are the main advantages of using these models?

Popular GLMs: linear regression

$$Y_i | \mu_i, \sigma \rightsquigarrow \text{Normal}(\mu_i, \sigma)$$

Random component

+

$$\mu_i = \alpha + \sum_{j=1}^p \beta_j x_{ij}$$

Systematic component

$$g(\mu_i) = \mu_i$$

Canonical link function

Popular GLMs: **logistic regression**

$$Y_i | \pi_i \rightsquigarrow \text{Bernoulli}(\pi_i)$$

Random component

+

$$g(\pi_i) = \alpha + \sum_{j=1}^p \beta_j x_{ij}$$

Systematic component

$$g(\pi_i) = \log \frac{\pi_i}{1 - \pi_i}$$

canonical link function

logit

Popular GLMs: **probit regression**

$$Y_i | \pi_i \rightsquigarrow \text{Bernoulli}(\pi_i)$$

Random component

+

$$g(\pi_i) = \alpha + \sum_{i=1}^p \beta_i x_{ij}$$

Systematic component

$$g(\pi_i) = \Phi^{-1}(\pi_i)$$

Probit link function

where $\Phi^{-1}(\cdot)$ is the quantile function of a standard Normal distribution

Popular GLMs: complementary log-log

$$Y | \pi \rightsquigarrow \text{Bernoulli}(\pi)$$

Random component

+

$$g(\pi) = \alpha + \sum_{i=1}^p \beta_i x_i$$

Systematic component

$$g(\pi) = \log(-\log(1 - \pi))$$

Complementary log-log link function

A theoretical note on the link functions for the Bernoulli model

$$g(\pi_i) = \log \frac{\pi_i}{1 - \pi_i}$$

$$g(\pi_i) = \Phi^{-1}(\pi_i)$$

$$g(\pi) = \log(-\log(1 - \pi))$$

A theoretical note on the link functions for the Bernoulli model

$$\eta_i = \log \frac{\pi_i}{1 - \pi_i} \Leftrightarrow \pi_i = \frac{e^{\eta_i}}{1 + e^{\eta_i}}$$

$$\eta_i = \Phi^{-1}(\pi_i) \Leftrightarrow \pi_i = \Phi(\eta_i)$$

$$\eta_i = \log(-\log(1 - \pi_i)) \Leftrightarrow \pi_i = 1 - e^{-e^{\eta_i}}$$

A theoretical note on the link functions for the Bernoulli model

$$\eta_i = \log \frac{\pi_i}{1 - \pi_i} \Leftrightarrow \pi_i = \frac{e^{\eta_i}}{1 + e^{\eta_i}}$$

Cumulative distribution of a standard
Logistic distribution

$$\eta_i = \Phi^{-1}(\pi_i) \Leftrightarrow \pi_i = \Phi(\eta_i)$$

Cumulative distribution of a Standard
Normal distribution

$$\eta_i = \log(-\log(1 - \pi_i)) \Leftrightarrow \pi_i = 1 - e^{-e^{\eta_i}}$$

1- Cumulative distribution of an Extreme
Value distribution

A theoretical note on the link functions for the Bernoulli model

$$\eta_i = \log \frac{\pi_i}{1 - \pi_i} \Leftrightarrow \pi_i = \frac{1}{1 + e^{-\eta_i}}$$

Cumulative distribution of a standard
Logistic distribution

$$\eta_i = \Phi^{-1}(\pi_i) \Leftrightarrow \pi_i = \Phi(\eta_i)$$

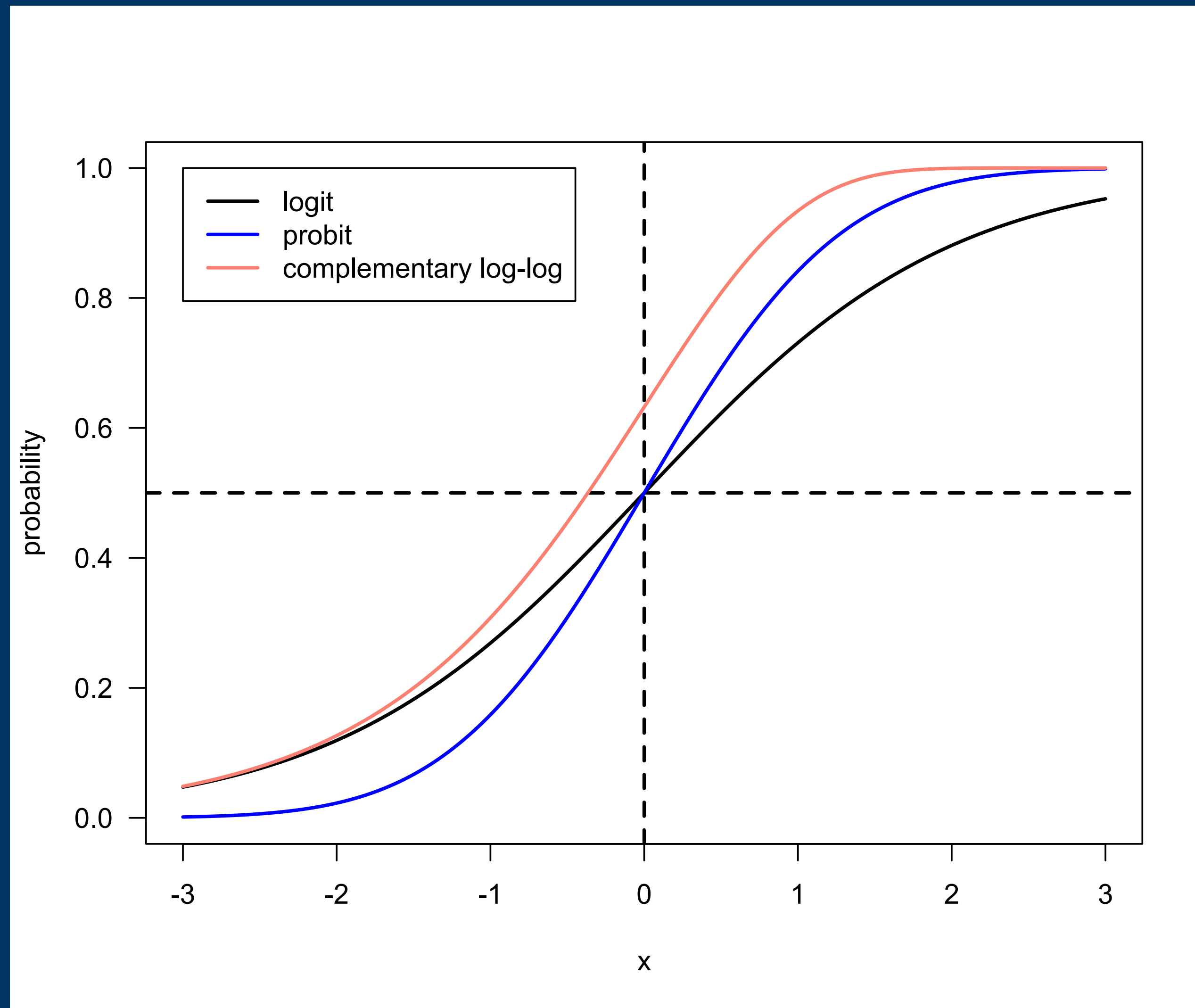
Cumulative distribution of a Standard
Normal distribution

$$\eta_i = \log(-\log(1 - \pi_i)) \Leftrightarrow \pi_i = 1 - e^{-e^{\eta_i}}$$

1- Cumulative distribution of an Extreme
Value distribution

Practical Implication: The inverse of any cumulative probability distribution can be used as a link function

A practical note on the link functions for the Bernoulli model



Exercise:

RESEARCH ARTICLE

B-Lymphocyte Depletion in Myalgic Encephalopathy/ Chronic Fatigue Syndrome. An Open-Label Phase II Study with Rituximab Maintenance Treatment

Øystein Fluge^{1*}, Kristin Risa¹, Sigrid Lunde¹, Kine Alme¹, Ingrid Gurvin Rekeland¹, Dipak Sapkota^{1,2}, Einar Kleboe Kristoffersen^{3,4}, Kari Sørland¹, Ove Bruland^{1,5}, Olav Dahl^{1,4}, Olav Mella^{1,4*}

- 1 Department of Oncology and Medical Physics, Haukeland University Hospital, Bergen, Norway,
- 2 Department of Clinical Medicine, University of Bergen, Haukeland University Hospital, Bergen, Norway,
- 3 Department of Immunology and Transfusion Medicine, Haukeland University Hospital, Bergen, Norway,
- 4 Department of Clinical Science, University of Bergen, Haukeland University Hospital, Bergen, Norway,
- 5 Department of Medical Genetics and Molecular Medicine, Haukeland University Hospital, Bergen, Norway



Exercise:

Abstract

Background

Myalgic Encephalopathy/Chronic Fatigue Syndrome (ME/CFS) is a disease of unknown etiology. We previously reported a pilot case series followed by a small, randomized, placebo-controlled phase II study, suggesting that B-cell depletion using the monoclonal anti-CD20 antibody rituximab can yield clinical benefit in ME/CFS.

Methods

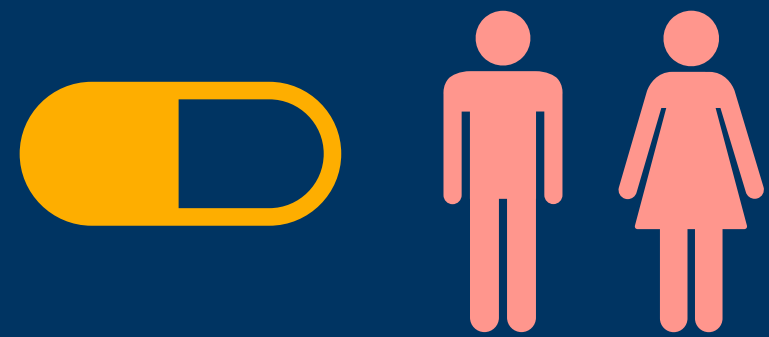
In this single-center, open-label, one-armed phase II study (NCT01156909), 29 patients were included for treatment with rituximab (500 mg/m²) two infusions two weeks apart, followed by maintenance rituximab infusions after 3, 6, 10 and 15 months, and with follow-up for 36 months.

Findings

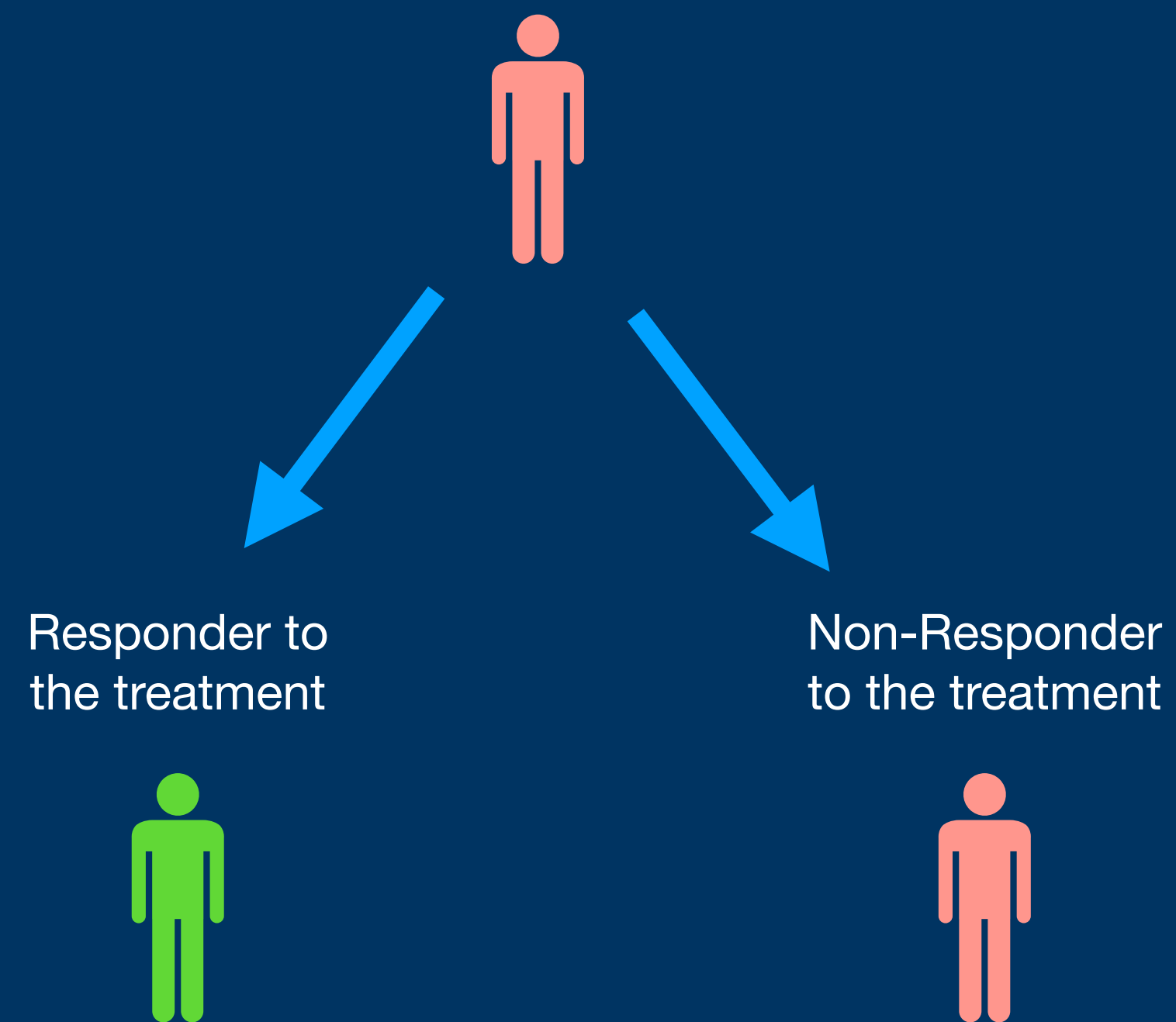
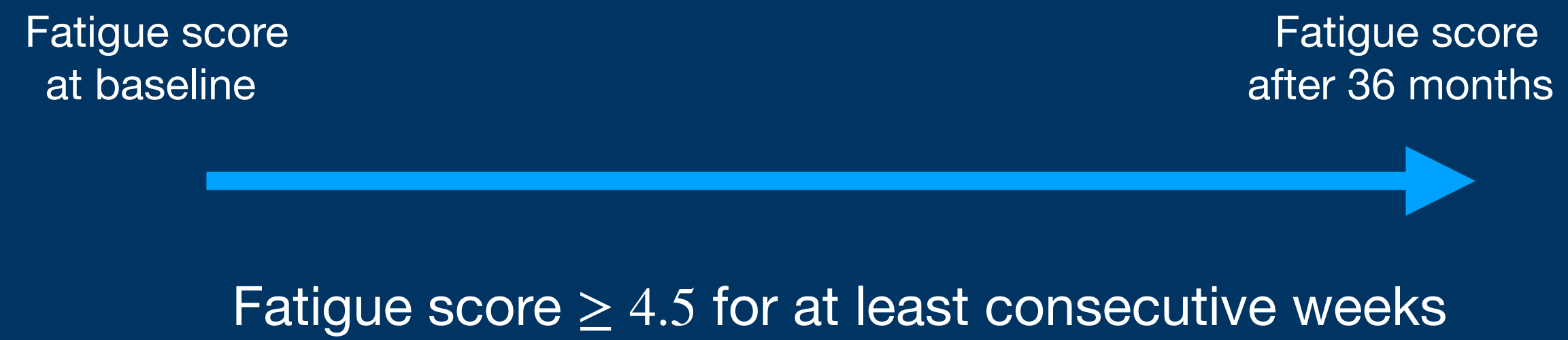
Major or moderate responses, predefined as lasting improvements in self-reported *Fatigue score*, were detected in 18 out of 29 patients (intention to treat). Clinically significant responses were seen in 18 out of 28 patients (64%) receiving rituximab maintenance treatment. For these 18 patients, the mean response durations within the 156 weeks study period were 105 weeks in 14 major responders, and 69 weeks in four moderate responders. At end of follow-up (36 months), 11 out of 18 responding patients were still in ongoing clinical remission. For major responders, the mean lag time from first rituximab infusion until start of clinical response was 23 weeks (range 8–66). Among the nine patients from the placebo group in the previous randomized study with no significant improvement during 12

Exercise:

Rituximab (n=29)



Biomarker
Fatigue score



Let's analyse the data

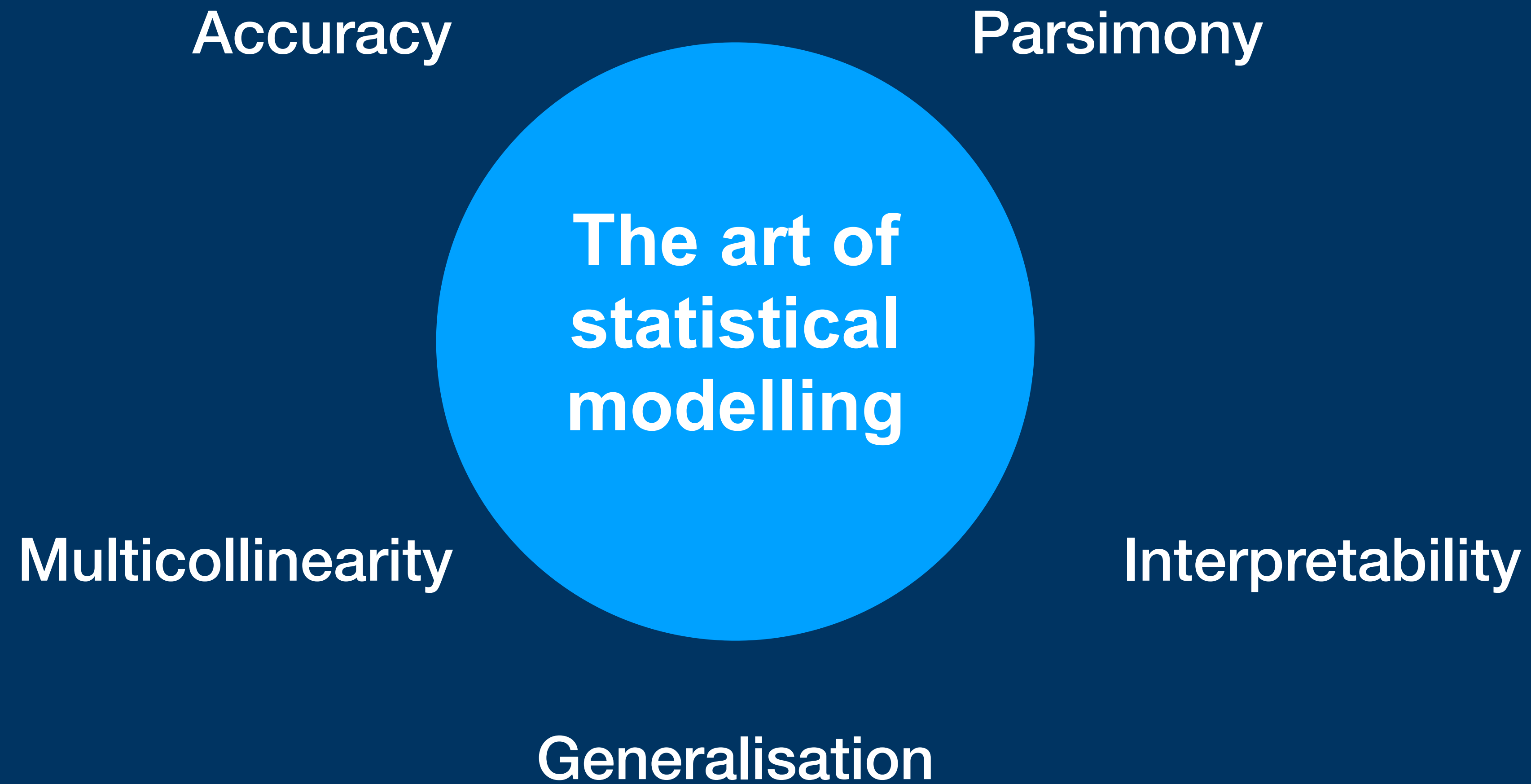
dataset: data_mecfs_rituximab.csv:

Estimate the probability of treatment response using statistical inference methods for the binomial distribution?

Use `binom.test` or `prop.test` functions

Construct a logistic regression model to understand whether age, gender, disease duration affect the treatment success?

Use `glm` function



The art of constructing a model

Select the best link function

- Fit models with different link functions and compare them

- Fit models with flexible link functions (e.g., Aranda-Ordaz link function for Bernoulli models)

Select the best subset of covariates (feature selection)

- Forward/Backward/Stepwise Regression

- Penalised regression (LASSO or Elastic-Net)

The art of constructing a model

Select the best link function

- Fit models with different link functions and compare them

- Fit models with flexible link functions (e.g., Aranda-Ordaz link function for Bernoulli models)

Select the best subset of covariates (feature selection)

- Forward/Backward/Stepwise Regression

- Penalised regression (LASSO or Elastic-Net)

Model comparison and selection

AIC - Akaike's Information Criterion

$$\text{AIC}(M) = (-2)\log\text{-L}(\hat{\theta} | M, \mathbf{x}) + 2p$$

BIC - Bayesian Information Criterion

$$\text{BIC}(M) = (-2)\log\text{-L}(\hat{\theta} | M, \mathbf{x}) + p \log(n)$$

$\log\text{-L}(\hat{\theta} | M, \mathbf{x})$ is the log-likelihood function evaluated on the parameter estimates

p is the number of parameters of model M

n is the sample size

Choose the model with the lowest values of one of these measures

Forward selection

“Empty” Model

Add covariate

Add covariate

Add covariate

⋮

Stop procedure

Increased accuracy **compensates**
increased model complexity

Increased accuracy **does not compensate**
increased model complexity

Backward elimination

“All covariates” Model

Remove covariate

Remove covariate

Remove covariate

⋮

Stop procedure

Decreased model complexity **does not have** an impact on model accuracy

Decreased model complexity **has an impact** on model accuracy

Stepwise regression

“Empty” Model

Add covariate 1

Add covariate 2

Remove covariate 1

Add covariate 3

Remove covariates 1, 2

⋮

Stop procedure



Increased accuracy **compensates**
increased model complexity

Increased accuracy **does not compensate**
increased model complexity

Stepwise regression

Advantages

Remove multicollinearity

Easy automation

Speed

Disadvantages

Overestimation of the number of predictors

Inflated type I errors

Unstable to slight changes in the data

Model comparison and selection

AIC - Akaike's Information Criterion

$$\text{AIC}(M) = (-2)\log\text{-L}(\hat{\theta} | M, \mathbf{x}) + 2p$$

BIC - Bayesian Information Criterion

$$\text{BIC}(M) = (-2)\log\text{-L}(\hat{\theta} | M, \mathbf{x}) + p \log(n)$$

$\log\text{-L}(\hat{\theta} | M, \mathbf{x})$ is the log-likelihood function evaluated on the parameter estimates

p is the number of parameters of model M

n is the sample size

Choose the model with the lowest values of one of these measures

Model validation

AIC - Akaike's Information Criterion

$$\text{AIC}(M) = (-2)\log\text{-L}(\hat{\theta} | M, \mathbf{x}) + 2p$$

BIC - Bayesian Information Criterion

$$\text{BIC}(M) = (-2)\log\text{-L}(\hat{\theta} | M, \mathbf{x}) + p \log(n)$$

$\log\text{-L}(\hat{\theta} | M, \mathbf{x})$ is the log-likelihood function evaluated on the parameter estimates

p is the number of parameters of model M

n is the sample size

Choose the model with the lowest values of one of these measures

Penalised regression

Estimation



Model selection

Accuracy



Bias

Penalised regression

$$\hat{\mathbf{b}} = \underset{\mathbf{b}}{\operatorname{argmin}} \left\{ \sum_{i=1}^n \left(y_i - b_0 - \sum_{j=1}^p b_j x_i \right)^2 \right\} .$$

subject to a constraint

$$pen \leq \lambda$$

pen = penalty function

λ = tuning parameter

Ridge Regression

$$\hat{\mathbf{b}} = \operatorname{argmin}_{\mathbf{b}} \left\{ \sum_{i=1}^n \left(y_i - b_0 - \sum_{j=1}^p b_j x_i \right)^2 \right\},$$

subject to $\sum_{j=1}^p b_j^2 \leq \lambda_2$

$$\lambda_2 \in \left[0, \sum_{j=1}^p (\hat{b}_j^*)^2 \right]$$

↑
OLS estimates

Geometrical interpretation (2D)

$$\sum_{j=1}^2 b_j^2 \leq \lambda_2$$

$$r^2(\cos^2 \theta + \sin^2 \theta) \leq \lambda_2$$

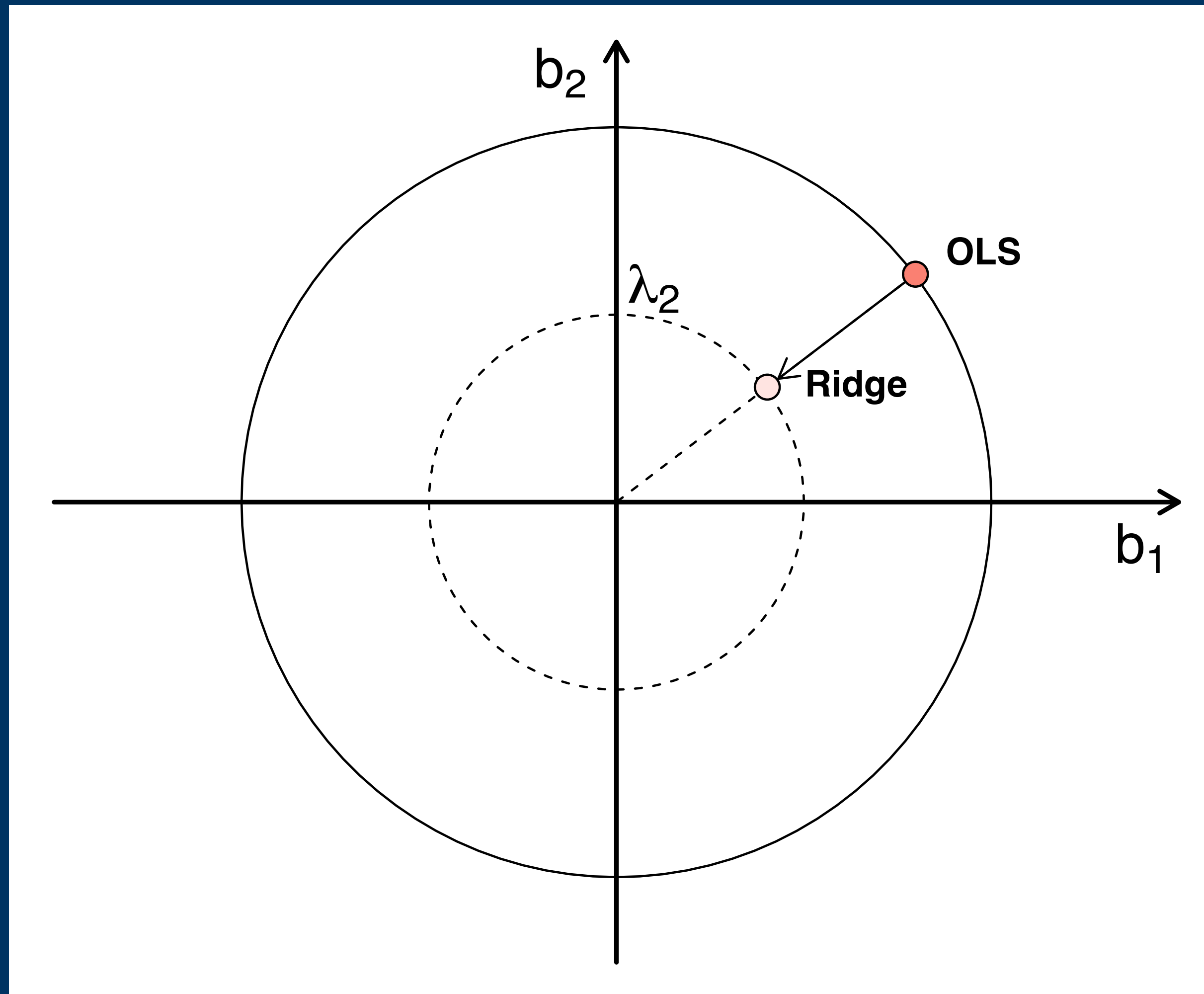
$$b_1 = r \cos \theta$$

$$r^2 \leq \lambda_2$$

$$b_2 = r \sin \theta$$

Ridge estimator is only dependent on the radius and not on the angle

Geometrical interpretation (2D)



Ordinary least squares estimator

$$\hat{\mathbf{b}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

Ridge estimator

$$\hat{\mathbf{b}} = (\mathbf{X}^T \mathbf{X} + \lambda_2 \mathbf{I})^{-1} \mathbf{X}^T \mathbf{Y}$$

Ridge Regression

$$\hat{\mathbf{b}} = \operatorname{argmin}_{\mathbf{b}} \left\{ \sum_{i=1}^n \left(y_i - b_0 - \sum_{j=1}^p b_j x_i \right)^2 \right\},$$

0% shrinkage

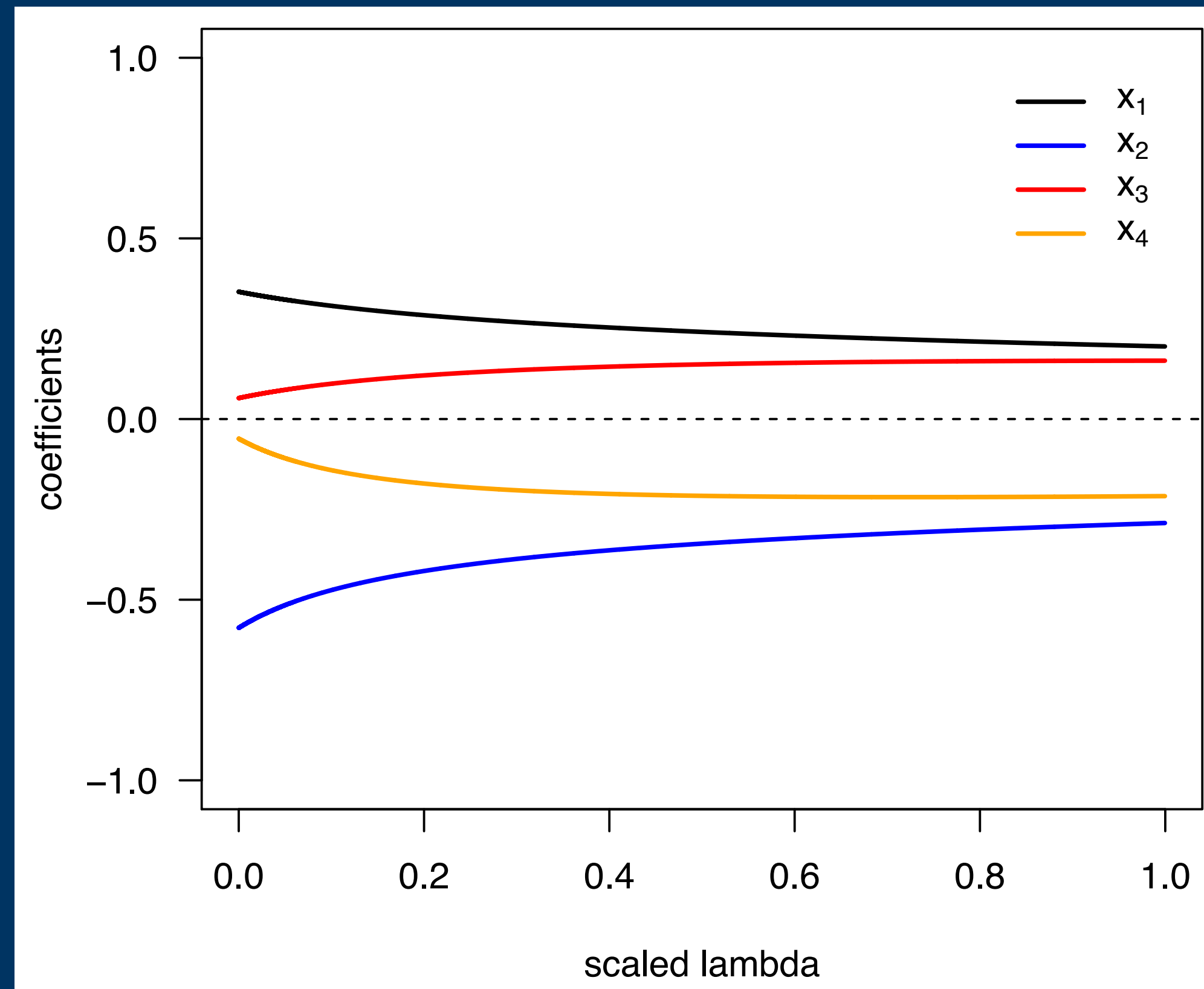
subject to

$$\frac{\sum_{j=1}^p b_j^2}{\sum_{j=1}^p (\hat{b}_j^*)^2} \leq 1 - \lambda^*$$

$$\lambda^* \in [0, 1]$$

“100%” shrinkage

Ridge trace plot



Ridge regression

Advantages

Remove multicollinearity

Estimator with a closed form

Shrinkage

Disadvantages

Biased estimators

No shrinkage to zero

(No model selection)

LASSO Regression

$$\hat{\mathbf{b}} = \underset{\mathbf{b}}{\operatorname{argmin}} \left\{ \sum_{i=1}^n \left(y_i - b_0 - \sum_{j=1}^p b_j x_i \right)^2 \right\},$$

subject to $\sum_{j=1}^p |b_j| \leq \lambda_1$

$$\lambda_1 \in \left[0, \sum_{j=1}^p |\hat{b}_j^*| \right]$$

OLS estimates

Geometrical interpretation (2D)

$$\sum_{j=1}^2 |b_j| \leq \lambda_1$$

$$b_1 = r \cos \theta$$

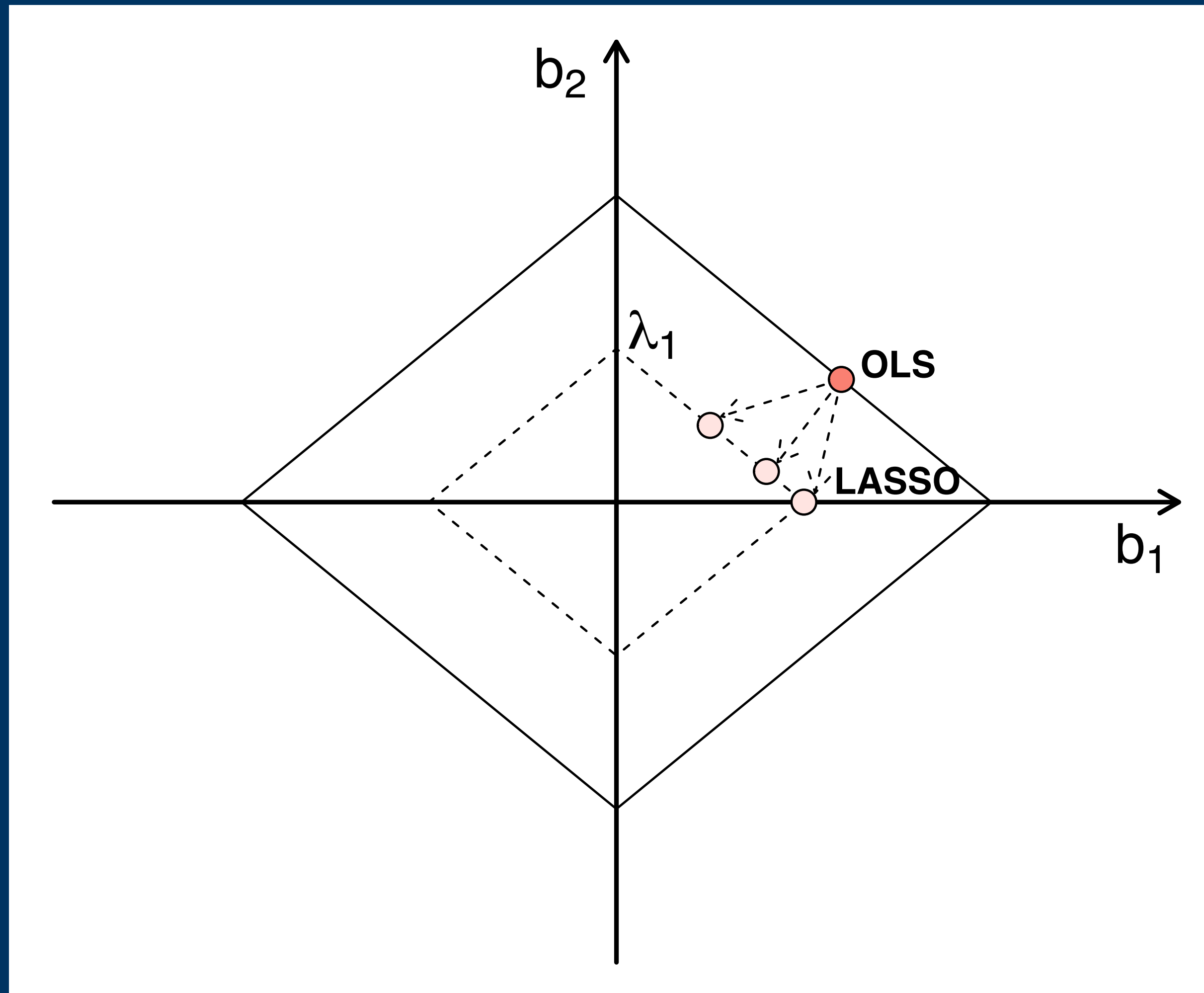
$$b_2 = r \sin \theta$$

$$r(\cos \theta + \sin \theta) \leq \lambda_2$$

$$r^2 \leq \lambda_2$$

LASSO estimator is dependent on
both radius and angle

Geometrical interpretation (2D)



LASSO Regression

$$\hat{\mathbf{b}} = \underset{\mathbf{b}}{\operatorname{argmin}} \left\{ \sum_{i=1}^n \left(y_i - b_0 - \sum_{j=1}^p b_j x_i \right)^2 \right\},$$

subject to

$$\frac{\sum_{j=1}^p |b_j|}{\sum_{j=1}^p |b_j^*|} \leq 1 - \lambda^*$$

0% shrinkage (OLS)

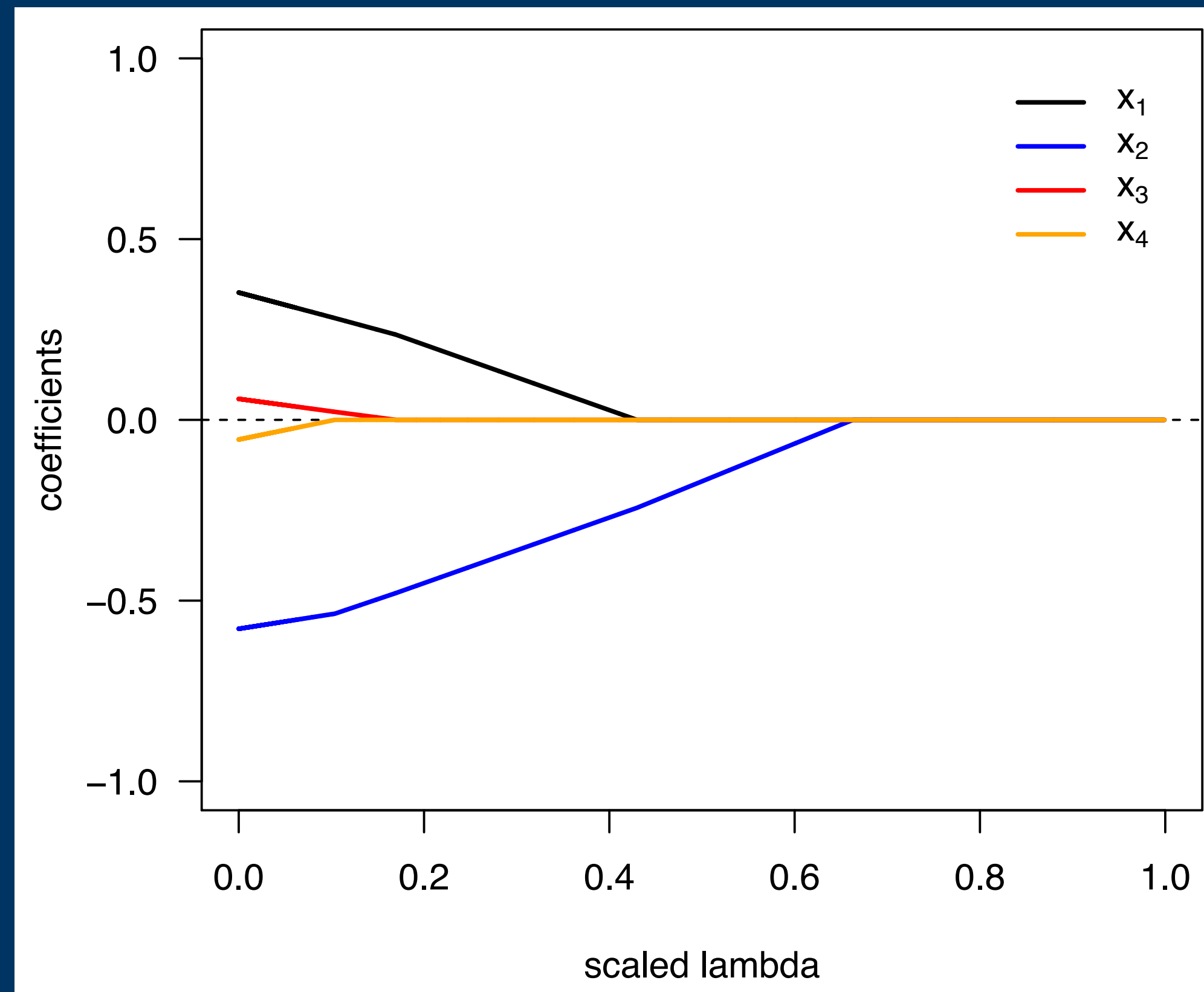


$$\lambda^* \in [0, 1]$$



100% shrinkage

LASSO trace plot



LASSO regression

Advantages

Remove multicollinearity

Shrinkage to zero

(Model selection)

Disadvantages

Random choice of highly correlated covariates

No closed-form expression

Problems with standard errors

Elastic Net Regression

$$\hat{\mathbf{b}} = \operatorname{argmin}_{\mathbf{b}} \left\{ \sum_{i=1}^n \left(y_i - b_0 - \sum_{j=1}^p b_j x_i \right)^2 \right\},$$

subject to $\alpha \|\mathbf{b}\|_1 + (1 - \alpha) \|\mathbf{b}\|^2 \leq \lambda$ for some λ and $\alpha \in [0,1]$.

$\alpha = 0 \Rightarrow$ Ridge regression

$\alpha = 1 \Rightarrow$ LASSO regression

Estimation of the tuning parameter(s)

Evaluate a grid of
possible values



Highest
accuracy

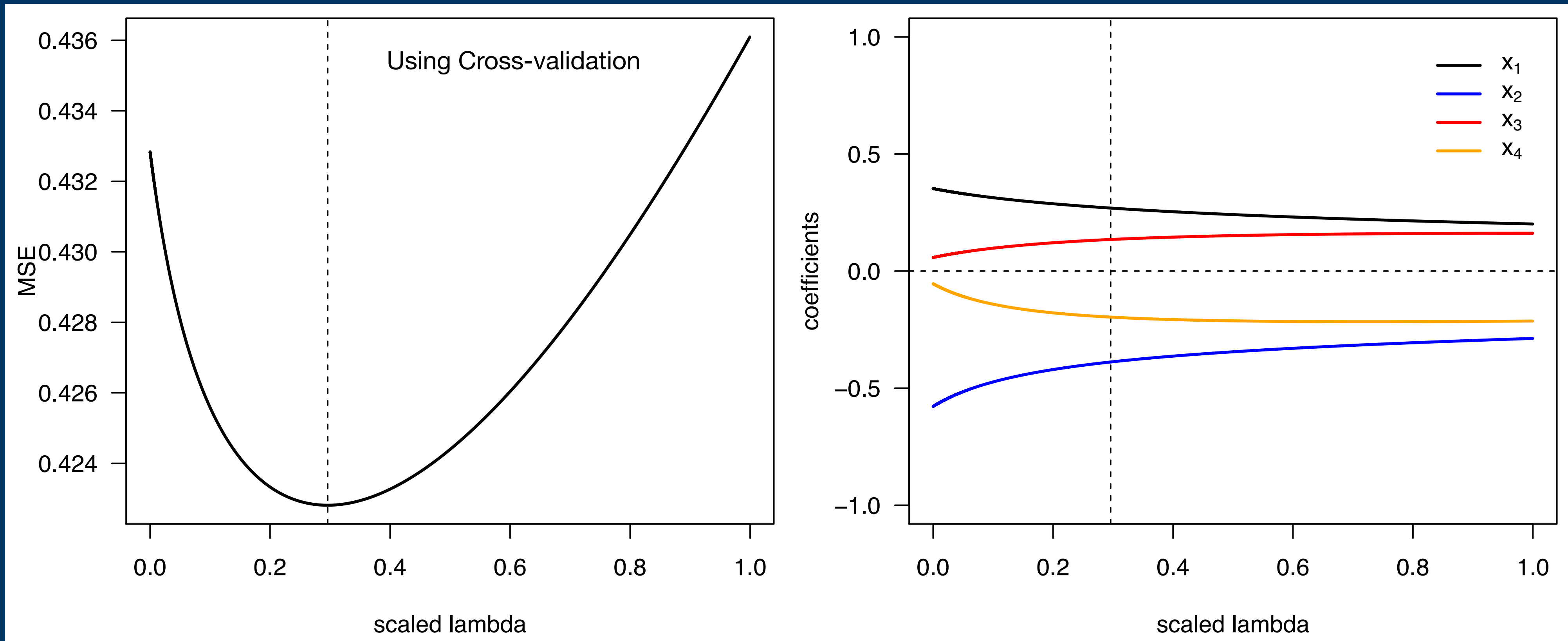


Cross-
validation

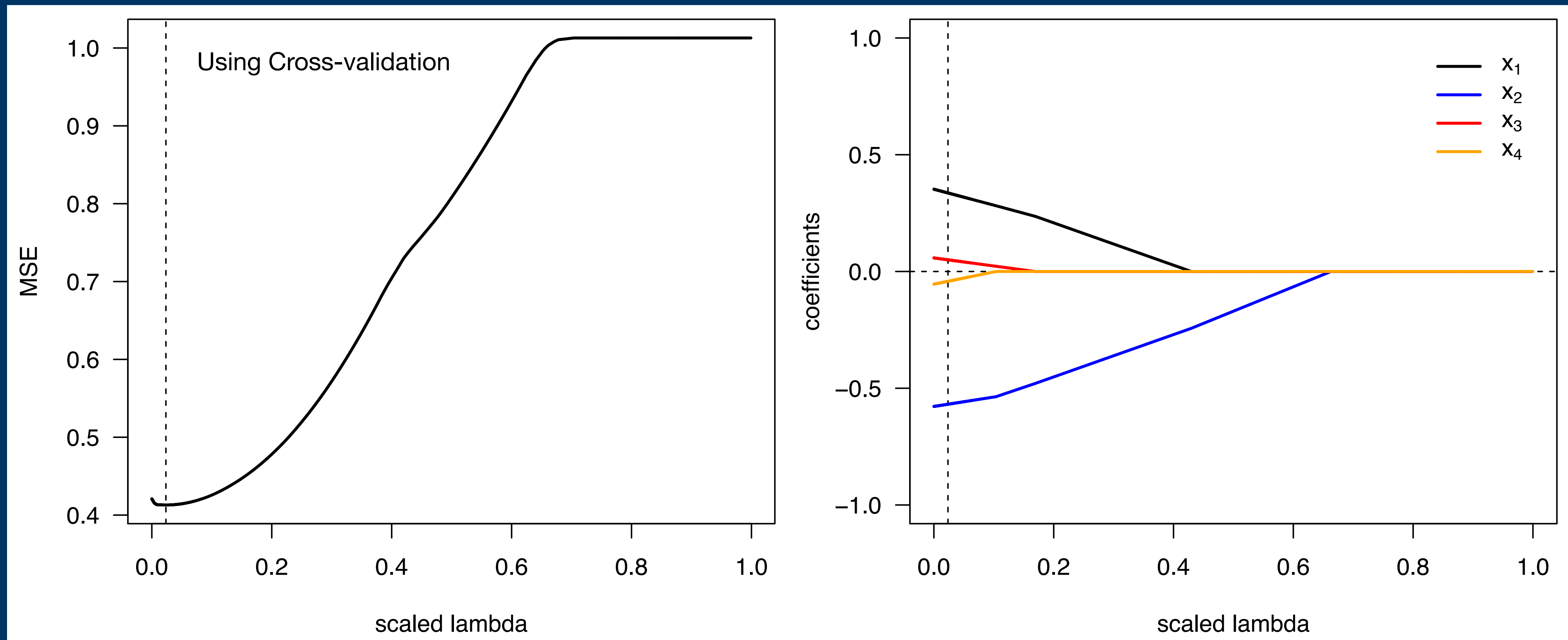


Lowest mean
squared error

Example: Ridge Regression



Example: LASSO Regression



Example: Elastic Net Regression

