

Lesson 6: validation

the BDC group

September 24, 2019

Outline

- scalable learning
- connecting the pieces together
- peer instructions
- a hands-on example

Scalable learning

Connecting the pieces together

The pieces of the puzzle, up to now:

- estimating the model parameters
 - validating the model structure
 - estimating the uncertainties
-
- validating the data

Connecting the pieces together - piece 1: estimating the model parameters

Most important strategies, up to now:

- Ordinary Least Squares
- Principal Component Regression
- Partial Least Squares

Ordinary Least Squares

$$\text{Assumptions: } \left\{ \begin{array}{l} \text{data generation model: } y_t = f(u_t; \theta) + v_t \\ \text{dataset: } \mathcal{D} = \{(u_t, y_t)\}_{t=1, \dots, N} \\ \text{hypothesis space: } \theta \in \Theta \end{array} \right.$$

$$\text{Formulation: } \hat{\theta}_{\text{LS}} = \arg \min_{\theta \in \Theta} \left\| \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} - \begin{bmatrix} f(u_1; \theta) \\ \vdots \\ f(u_N; \theta) \end{bmatrix} \right\|^2 = \arg \min_{\theta \in \Theta} \sum_{t=1}^N \left(y_t - f(u_t; \theta) \right)^2$$

Principal Components Regression

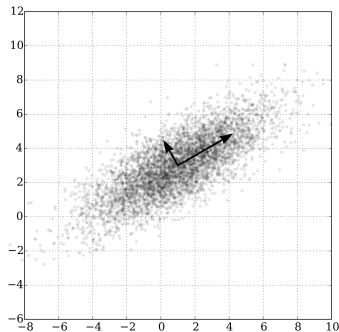
step 1: $X = U\Sigma V^T \rightarrow$ select U^{PC} (i.e., the first n components)

step 2: $\hat{\theta}_{\text{PCR}} = \arg \min_{\theta \in \Theta} \sum_{t=1}^N \left(y_t - U^{\text{PC}} \theta \right)^2$

Principal Components Regression

step 1: $X = U\Sigma V^T \rightarrow$ select U^{PC} (i.e., the first n components)

step 2: $\hat{\theta}_{\text{PCR}} = \arg \min_{\theta \in \Theta} \sum_{t=1}^N \left(y_t - U^{\text{PC}} \theta \right)^2$

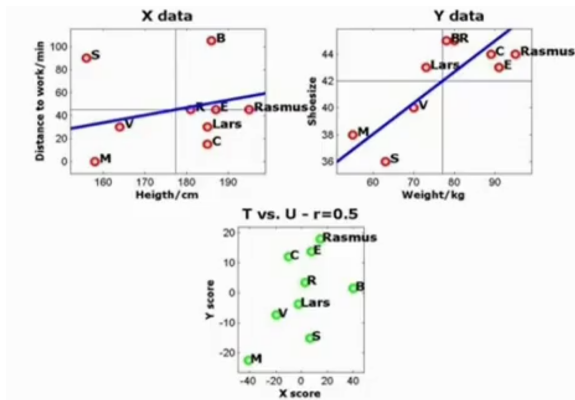


Partial Least Squares

$$\begin{cases} X = TP^T + E \\ Y = UQ^T + F \end{cases} \quad (1)$$

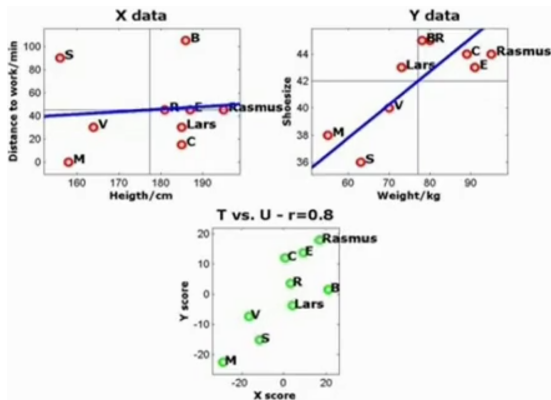
Partial Least Squares

$$\begin{cases} X = TP^T + E \\ Y = UQ^T + F \end{cases} \quad (1)$$



Partial Least Squares

$$\begin{cases} X = TP^T + E \\ Y = UQ^T + F \end{cases} \quad (1)$$



?

Connecting the pieces together

The pieces of the puzzle, up to now:

- estimating the model parameters
 - validating the model structure
 - estimating the uncertainties
-
- validating the data

Connecting the pieces together - piece 2: validating the model structure

First question, even before seeing the data: are the *internal* and *external* validity satisfied?

Connecting the pieces together - piece 2: validating the model structure

First question, even before seeing the data: are the *internal* and *external* validity satisfied? For us, in a rough way:

- drawing conclusions (i.e., make models) ignoring some unknown input \implies losing internal validity

Connecting the pieces together - piece 2: validating the model structure

First question, even before seeing the data: are the *internal* and *external* validity satisfied? For us, in a rough way:

- drawing conclusions (i.e., make models) ignoring some unknown input \implies losing internal validity
- drawing conclusions (i.e., make models) ignoring some part of the input space (AND being in a situation where this leads to poor generalization capabilities) \implies losing external validity

Connecting the pieces together - piece 2: validating the model structure

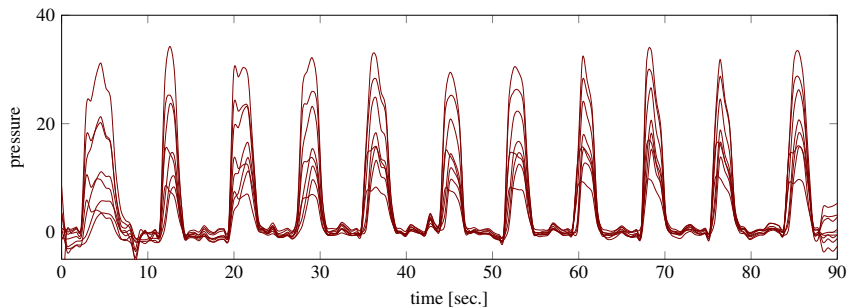
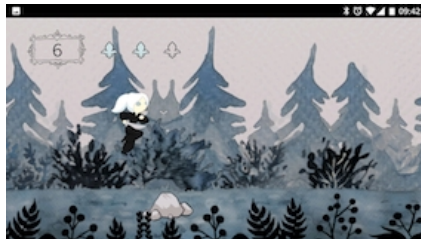
First question, even before seeing the data: are the *internal* and *external* validity satisfied? For us, in a rough way:

- drawing conclusions (i.e., make models) ignoring some unknown input \implies losing internal validity
- drawing conclusions (i.e., make models) ignoring some part of the input space (AND being in a situation where this leads to poor generalization capabilities) \implies losing external validity

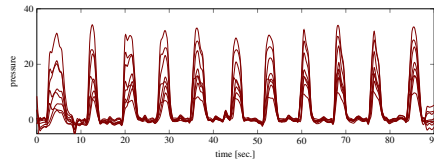
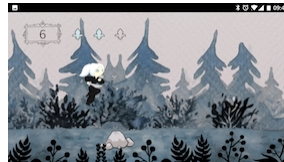
Question 1

why "... space AND being ... "?

Example - introduction

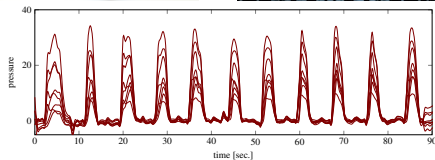
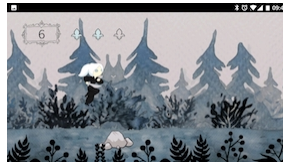


Example - population validity



Definition: population validity = how well the sample used can be extrapolated to a population as a whole (*type of external validity*)

Example - ecological validity



Definition: ecological validity = how well the findings can be extrapolated to real life settings (*type of external validity*)

?

Bias vs. variance

our approach: θ = unknown and deterministic; $\hat{\theta}$ = estimator of θ

Bias vs. variance

our approach: θ = unknown and deterministic; $\hat{\theta}$ = estimator of θ

Question 2

is $\hat{\theta}$ always a random variable?

- yes
- no

Bias vs. variance

our approach: θ = unknown and deterministic; $\hat{\theta}$ = estimator of θ

Question 2

is $\hat{\theta}$ always a random variable?

- yes
- no

usual case: $\hat{\theta}$ = random variable

Mean Squared Error

how do we weight the difference between θ and $\hat{\theta}$?

Mean Squared Error

how do we weight the difference between θ and $\widehat{\theta}$?

squared error committed by a specific realization of $\widehat{\theta}$: $\|\theta - \widehat{\theta}\|^2$

Mean Squared Error

how do we weight the difference between θ and $\widehat{\theta}$?

squared error committed by a specific realization of $\widehat{\theta}$: $\|\theta - \widehat{\theta}\|^2$

mean squared error committed by $\widehat{\theta}$: $\mathbb{E}[\|\theta - \widehat{\theta}\|^2]$

Remark 1: the MSE is a function of θ !

$$\text{MSE}(\theta) = \mathbb{E} \left[\|\theta - \widehat{\theta}\|^2 \right]$$

Remark 1: the MSE is a function of θ !

$$\text{MSE}(\theta) = \mathbb{E} \left[\|\theta - \widehat{\theta}\|^2 \right]$$

Example: $y_t \sim \mathcal{N}(\mu, 1)$ $\widehat{\mu} = 3$

Remark 1: the MSE is a function of θ !

$$\text{MSE}(\theta) = \mathbb{E} \left[\|\theta - \widehat{\theta}\|^2 \right]$$

Example: $y_t \sim \mathcal{N}(\mu, 1)$ $\widehat{\mu} = 3$

Question 3

MSE(3) = ?

- 0
- 1
- 10

Remark 1: the MSE is a function of θ !

$$\text{MSE}(\theta) = \mathbb{E} \left[\|\theta - \widehat{\theta}\|^2 \right]$$

Example: $y_t \sim \mathcal{N}(\mu, 1)$ $\widehat{\mu} = 3$

fundamental message: given $\widehat{\theta}$, that estimator may have excellent performance for certain specific θ s and awful performance for other ones!

Remark 2: the MSE is defined over the measure of y_1, \dots, y_N !

$$\mathbb{E}[\|\theta - \widehat{\theta}\|^2] = \int_{\mathcal{Y}^N} \|\theta - \widehat{\theta}(y_1, \dots, y_N)\|^2 dp(y_1, \dots, y_N; \theta)$$

Remark 2: the MSE is defined over the measure of y_1, \dots, y_N !

$$\mathbb{E}[\|\theta - \widehat{\theta}\|^2] = \int_{\mathcal{Y}^N} \|\theta - \widehat{\theta}(y_1, \dots, y_N)\|^2 dp(y_1, \dots, y_N; \theta)$$

important implication: the MSE cannot be computed!

Remark 2: the MSE is defined over the measure of y_1, \dots, y_N !

$$\mathbb{E} \left[\|\theta - \widehat{\theta}\|^2 \right] = \int_{\mathcal{Y}^N} \|\theta - \widehat{\theta}(y_1, \dots, y_N)\|^2 dp(y_1, \dots, y_N; \theta)$$

important implication: the MSE cannot be computed!

Strategy: estimate some alternative quantity from the data:

$$\begin{cases} y_t = f(u_t; \theta) + v_t \\ \widehat{y}_t = \widehat{f}(u_t; \widehat{\theta}) \end{cases} \quad \mapsto \quad \frac{1}{N} \sum_{t=1}^N (y_t - \widehat{y}_t)^2$$

Remark 2: the MSE is defined over the measure of y_1, \dots, y_N !

$$\mathbb{E}[\|\theta - \widehat{\theta}\|^2] = \int_{\mathcal{Y}^N} \|\theta - \widehat{\theta}(y_1, \dots, y_N)\|^2 dp(y_1, \dots, y_N; \theta)$$

important implication: the MSE cannot be computed!

Strategy: estimate some alternative quantity from the data:

$$\begin{cases} y_t = f(u_t; \theta) + v_t \\ \widehat{y}_t = \widehat{f}(u_t; \widehat{\theta}) \end{cases} \mapsto \frac{1}{N} \sum_{t=1}^N (y_t - \widehat{y}_t)^2$$

'training-vs-test' and cross-validation are examples

the bias - variance tradeoff

Decomposing the MSE in two interesting terms

$$\mathbb{E}[\|\widehat{\theta} - \theta\|^2]$$

Decomposing the MSE in two interesting terms

$$\mathbb{E} \left[\left\| \widehat{\theta} - \theta \right\|^2 \right] = \mathbb{E} \left[\left\| \widehat{\theta} - \mathbb{E} \left[\widehat{\theta} \right] + \mathbb{E} \left[\widehat{\theta} \right] - \theta \right\|^2 \right]$$

Decomposing the MSE in two interesting terms

$$\begin{aligned}\mathbb{E}\left[\|\widehat{\theta} - \theta\|^2\right] &= \mathbb{E}\left[\|\widehat{\theta} - \mathbb{E}[\widehat{\theta}] + \mathbb{E}[\widehat{\theta}] - \theta\|^2\right] \\ &= \mathbb{E}\left[\|\mathcal{V} + \mathcal{B}\|^2\right]\end{aligned}\quad \left\{ \begin{array}{l} \mathcal{V} := \widehat{\theta} - \mathbb{E}[\widehat{\theta}] \\ \mathcal{B} := \mathbb{E}[\widehat{\theta}] - \theta \end{array} \right.$$

Decomposing the MSE in two interesting terms

$$\begin{aligned}\mathbb{E}\left[\|\widehat{\theta} - \theta\|^2\right] &= \mathbb{E}\left[\|\widehat{\theta} - \mathbb{E}[\widehat{\theta}] + \mathbb{E}[\widehat{\theta}] - \theta\|^2\right] \\ &= \mathbb{E}\left[\|\mathcal{V} + \mathcal{B}\|^2\right] \\ &= \mathbb{E}\left[(\mathcal{V} + \mathcal{B})^T (\mathcal{V} + \mathcal{B})\right]\end{aligned}\quad \left\{ \begin{array}{l} \mathcal{V} := \widehat{\theta} - \mathbb{E}[\widehat{\theta}] \\ \mathcal{B} := \mathbb{E}[\widehat{\theta}] - \theta \end{array} \right.$$

Decomposing the MSE in two interesting terms

$$\begin{aligned}\mathbb{E}\left[\|\widehat{\theta} - \theta\|^2\right] &= \mathbb{E}\left[\|\widehat{\theta} - \mathbb{E}[\widehat{\theta}] + \mathbb{E}[\widehat{\theta}] - \theta\|^2\right] \\ &= \mathbb{E}\left[\|\mathcal{V} + \mathcal{B}\|^2\right] \\ &= \mathbb{E}\left[(\mathcal{V} + \mathcal{B})^T (\mathcal{V} + \mathcal{B})\right] \\ &= \mathbb{E}\left[\|\mathcal{V}\|^2 + \|\mathcal{B}\|^2 + 2\mathcal{V}^T \mathcal{B}\right]\end{aligned}\quad \left\{ \begin{array}{l} \mathcal{V} := \widehat{\theta} - \mathbb{E}[\widehat{\theta}] \\ \mathcal{B} := \mathbb{E}[\widehat{\theta}] - \theta \end{array} \right.$$

Decomposing the MSE in two interesting terms

$$\begin{aligned}\mathbb{E}\left[\|\widehat{\theta} - \theta\|^2\right] &= \mathbb{E}\left[\|\widehat{\theta} - \mathbb{E}[\widehat{\theta}] + \mathbb{E}[\widehat{\theta}] - \theta\|^2\right] \\ &= \mathbb{E}\left[\|\mathcal{V} + \mathcal{B}\|^2\right] \\ &= \mathbb{E}\left[(\mathcal{V} + \mathcal{B})^T (\mathcal{V} + \mathcal{B})\right] \\ &= \mathbb{E}\left[\|\mathcal{V}\|^2 + \|\mathcal{B}\|^2 + 2\mathcal{V}^T \mathcal{B}\right]\end{aligned}\quad \left\{ \begin{array}{l} \mathcal{V} := \widehat{\theta} - \mathbb{E}[\widehat{\theta}] \\ \mathcal{B} := \mathbb{E}[\widehat{\theta}] - \theta \end{array} \right.$$
$$\mathbb{E}[\mathcal{V}^T \mathcal{B}] = \mathbf{0}$$

Decomposing the MSE in two interesting terms

$$\begin{aligned}\mathbb{E} \left[\left\| \widehat{\theta} - \theta \right\|^2 \right] &= \mathbb{E} \left[\left\| \widehat{\theta} - \mathbb{E} \left[\widehat{\theta} \right] + \mathbb{E} \left[\widehat{\theta} \right] - \theta \right\|^2 \right] \\ &= \mathbb{E} \left[\left\| \mathcal{V} + \mathcal{B} \right\|^2 \right] \\ &= \mathbb{E} \left[(\mathcal{V} + \mathcal{B})^T (\mathcal{V} + \mathcal{B}) \right] \\ &= \mathbb{E} \left[\left\| \mathcal{V} \right\|^2 + \left\| \mathcal{B} \right\|^2 + 2\mathcal{V}^T \mathcal{B} \right] \\ &= \mathbb{E} \left[\left\| \mathcal{V} \right\|^2 \right] + \left\| \mathcal{B} \right\|^2\end{aligned}\quad \left\{ \begin{array}{l} \mathcal{V} := \widehat{\theta} - \mathbb{E} \left[\widehat{\theta} \right] \\ \mathcal{B} := \mathbb{E} \left[\widehat{\theta} \right] - \theta \end{array} \right.$$
$$\mathbb{E} \left[\mathcal{V}^T \mathcal{B} \right] = \mathbf{0}$$

Decomposing the MSE in two interesting terms

$$\begin{aligned}\mathbb{E} \left[\left\| \widehat{\theta} - \theta \right\|^2 \right] &= \mathbb{E} \left[\left\| \widehat{\theta} - \mathbb{E} \left[\widehat{\theta} \right] + \mathbb{E} \left[\widehat{\theta} \right] - \theta \right\|^2 \right] \\ &= \mathbb{E} \left[\left\| \mathcal{V} + \mathcal{B} \right\|^2 \right] \\ &= \mathbb{E} \left[(\mathcal{V} + \mathcal{B})^T (\mathcal{V} + \mathcal{B}) \right] \\ &= \mathbb{E} \left[\left\| \mathcal{V} \right\|^2 + \left\| \mathcal{B} \right\|^2 + 2\mathcal{V}^T \mathcal{B} \right] \\ &= \mathbb{E} \left[\left\| \mathcal{V} \right\|^2 \right] + \left\| \mathcal{B} \right\|^2 \\ &= \text{"variance"} + \text{"bias"}^2\end{aligned}$$
$$\begin{cases} \mathcal{V} := \widehat{\theta} - \mathbb{E} \left[\widehat{\theta} \right] \\ \mathcal{B} := \mathbb{E} \left[\widehat{\theta} \right] - \theta \end{cases}$$
$$\mathbb{E} \left[\mathcal{V}^T \mathcal{B} \right] = \mathbf{0}$$

Definitions

ideal model: $y_t = f(u_t) + v_t$

our model: $y_t = \widehat{f}(u_t, \widehat{\theta})$

underfitting = a $\widehat{f}(\cdot, \widehat{\theta})$ that misses the fundamental features of f_0

overfitting = a $\widehat{f}(\cdot, \widehat{\theta})$ that follows v_t instead of $f(u_t)$.

Remarking fact: the model complexity affects the bias - variance tradeoff

$$y_t = \prod_{k=1}^n \theta_k u_t^k + v_t$$

Remarking fact: the model complexity affects the bias - variance tradeoff

$$y_t = \prod_{k=1}^n \theta_k u_t^k + v_t$$

we will see how this strongly connects to the Ockham's razor

Quiz time!

Question 3

Underfitting is associated to

- *high bias and low variance*
- *low bias and high variance*

Quiz time!

Question 4

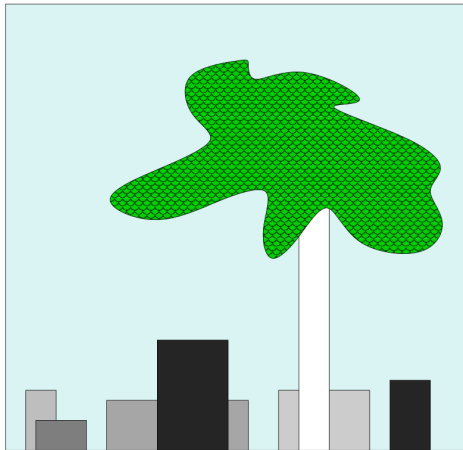
Overfitting is associated to:

- *high bias and low variance*
- *low bias and high variance*

?

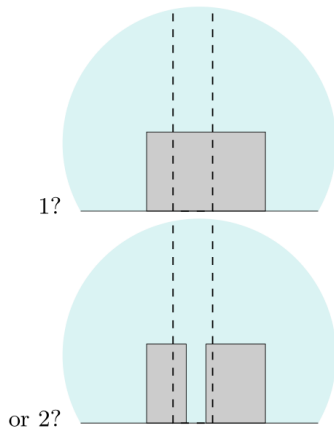
Ockham's razor

from David J.C. MacKay, *Information Theory, Inference, and Learning Algorithms*



Ockham's razor

from David J.C. MacKay, *Information Theory, Inference, and Learning Algorithms*



Connecting the pieces together - piece 2: validating the model structure

Strategy: estimate the MSE (or other performance indexes) from the data:

$$\begin{cases} y_t = f(u_t; \theta) + v_t \\ \hat{y}_t = \hat{f}(u_t; \hat{\theta}) \end{cases} \mapsto \frac{1}{N} \sum_{t=1}^N (y_t - \hat{y}_t)^2$$

Connecting the pieces together - piece 2: validating the model structure

Strategy: estimate the MSE (or other performance indexes) from the data:

$$\begin{cases} y_t = f(u_t; \theta) + v_t \\ \hat{y}_t = \hat{f}(u_t; \hat{\theta}) \end{cases} \mapsto \frac{1}{N} \sum_{t=1}^N (y_t - \hat{y}_t)^2$$

Most important strategies, up to now:

- dividing the dataset into training / test / validation
- cross-validation

Why is estimating the MSE on the training data a very bad idea?

Why is estimating the MSE on the training data a very bad idea?

$$y \sim \mathcal{N}(\mu, \sigma^2 I) \quad \widehat{\mu} : \mathbb{R}^n \mapsto \mathbb{R}^n \quad \widehat{\mu}(\cdot) = \text{estimator of } \mu \quad (2)$$

Why is estimating the MSE on the training data a very bad idea?

$$y \sim \mathcal{N}(\mu, \sigma^2 I) \quad \widehat{\mu} : \mathbb{R}^n \mapsto \mathbb{R}^n \quad \widehat{\mu}(\cdot) = \text{estimator of } \mu \quad (2)$$

$$R = \mathbb{E}[\|\mu - \widehat{\mu}\|_2^2]$$

Why is estimating the MSE on the training data a very bad idea?

$$y \sim \mathcal{N}(\mu, \sigma^2 I) \quad \widehat{\mu} : \mathbb{R}^n \mapsto \mathbb{R}^n \quad \widehat{\mu}(\cdot) = \text{estimator of } \mu \quad (2)$$

$$\begin{aligned} R &= \mathbb{E} \left[\|\mu - \widehat{\mu}\|_2^2 \right] \\ &= \mathbb{E} \left[\|\mu - y + y - \widehat{\mu}\|_2^2 \right] \end{aligned}$$

Why is estimating the MSE on the training data a very bad idea?

$$y \sim \mathcal{N}(\mu, \sigma^2 I) \quad \widehat{\mu} : \mathbb{R}^n \mapsto \mathbb{R}^n \quad \widehat{\mu}(\cdot) = \text{estimator of } \mu \quad (2)$$

$$\begin{aligned} R &= \mathbb{E} \left[\|\mu - \widehat{\mu}\|_2^2 \right] \\ &= \mathbb{E} \left[\|\mu - y + y - \widehat{\mu}\|_2^2 \right] \\ &= \mathbb{E} \left[\|\mu - y\|_2^2 \right] + \mathbb{E} \left[\|y - \widehat{\mu}\|_2^2 \right] + 2\mathbb{E} \left[(\mu - y)^T (y - \widehat{\mu}) \right] \end{aligned}$$

Why is estimating the MSE on the training data a very bad idea?

$$y \sim \mathcal{N}(\mu, \sigma^2 I) \quad \widehat{\mu} : \mathbb{R}^n \mapsto \mathbb{R}^n \quad \widehat{\mu}(\cdot) = \text{estimator of } \mu \quad (2)$$

$$\begin{aligned} R &= \mathbb{E} \left[\|\mu - \widehat{\mu}\|_2^2 \right] \\ &= \mathbb{E} \left[\|\mu - y + y - \widehat{\mu}\|_2^2 \right] \\ &= \mathbb{E} \left[\|\mu - y\|_2^2 \right] + \mathbb{E} \left[\|y - \widehat{\mu}\|_2^2 \right] + 2\mathbb{E} \left[(\mu - y)^T (y - \widehat{\mu}) \right] \\ &= n\sigma^2 + \mathbb{E} \left[\|y - \widehat{\mu}\|_2^2 \right] + 2\mathbb{E} \left[(y - \mu)^T (\widehat{\mu} - y) \right] \end{aligned}$$

Why is estimating the MSE on the training data a very bad idea?

$$y \sim \mathcal{N}(\mu, \sigma^2 I) \quad \widehat{\mu} : \mathbb{R}^n \mapsto \mathbb{R}^n \quad \widehat{\mu}(\cdot) = \text{estimator of } \mu \quad (2)$$

$$\begin{aligned} R &= \mathbb{E} \left[\|\mu - \widehat{\mu}\|_2^2 \right] \\ &= \mathbb{E} \left[\|\mu - y + y - \widehat{\mu}\|_2^2 \right] \\ &= \mathbb{E} \left[\|\mu - y\|_2^2 \right] + \mathbb{E} \left[\|y - \widehat{\mu}\|_2^2 \right] + 2\mathbb{E} \left[(\mu - y)^T (y - \widehat{\mu}) \right] \\ &= n\sigma^2 + \mathbb{E} \left[\|y - \widehat{\mu}\|_2^2 \right] + 2\mathbb{E} \left[(y - \mu)^T (\widehat{\mu} - y) \right] \\ &= -n\sigma^2 + \mathbb{E} \left[\|y - \widehat{\mu}\|_2^2 \right] + 2 \sum_{i=1}^n \text{cov}(y_i, \widehat{\mu}_i) \end{aligned}$$

Why is estimating the MSE on the training data a very bad idea?

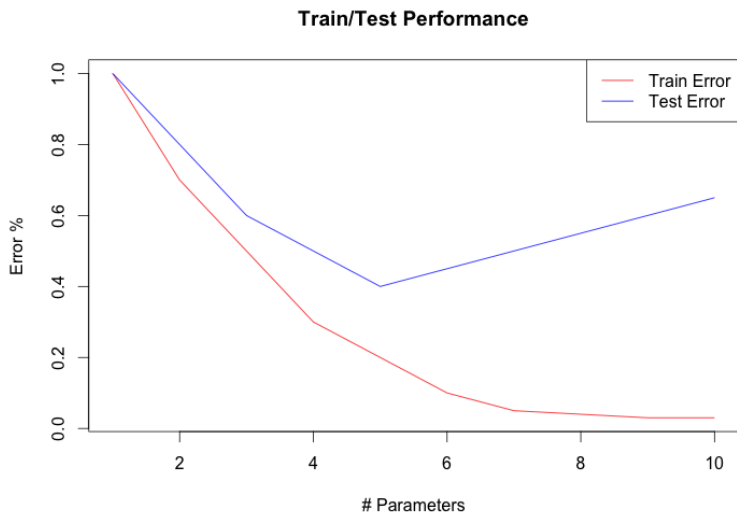
$$y \sim \mathcal{N}(\mu, \sigma^2 I) \quad \widehat{\mu} : \mathbb{R}^n \mapsto \mathbb{R}^n \quad \widehat{\mu}(\cdot) = \text{estimator of } \mu \quad (2)$$

$$\begin{aligned} R &= \mathbb{E} \left[\|\mu - \widehat{\mu}\|_2^2 \right] \\ &= \mathbb{E} \left[\|\mu - y + y - \widehat{\mu}\|_2^2 \right] \\ &= \mathbb{E} \left[\|\mu - y\|_2^2 \right] + \mathbb{E} \left[\|y - \widehat{\mu}\|_2^2 \right] + 2\mathbb{E} \left[(\mu - y)^T (y - \widehat{\mu}) \right] \\ &= n\sigma^2 + \mathbb{E} \left[\|y - \widehat{\mu}\|_2^2 \right] + 2\mathbb{E} \left[(y - \mu)^T (\widehat{\mu} - y) \right] \\ &= -n\sigma^2 + \mathbb{E} \left[\|y - \widehat{\mu}\|_2^2 \right] + 2 \sum_{i=1}^n \text{cov}(y_i, \widehat{\mu}_i) \\ \implies \quad \widehat{R} &= -n\sigma^2 + \|y - \widehat{\mu}\|_2^2 + 2 \sum_{i=1}^n \text{cov}(y_i, \widehat{\mu}_i) \end{aligned} \quad (3)$$

Potential practical difficulties when doing model validation

- lack of data
- lack of control of the input variables
- uncertainty about the underlying probability distributions and correlations

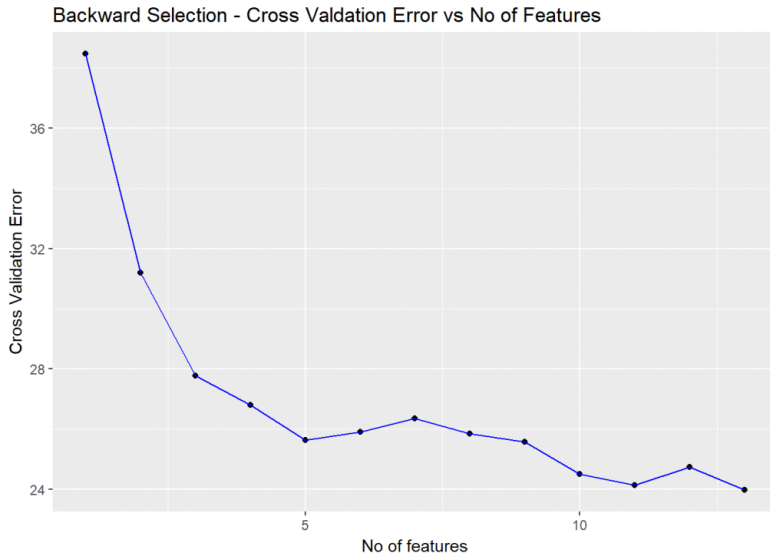
Practical examples: which model structure would you select?



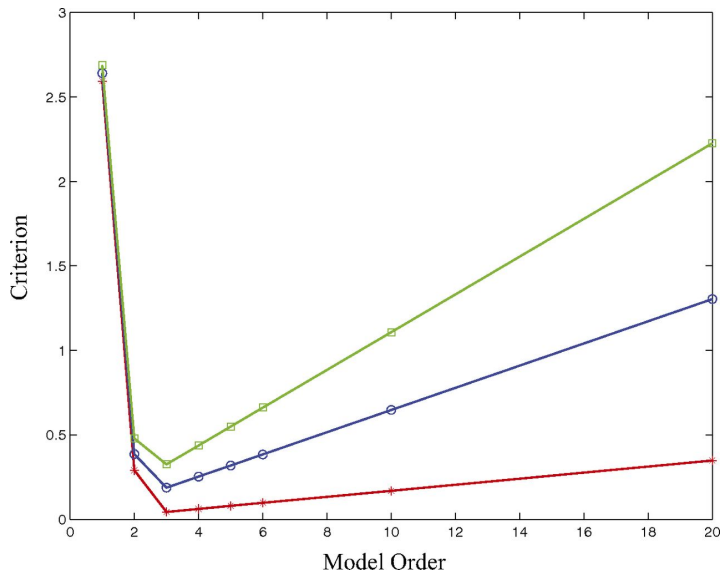
Practical examples: which model structure would you select?



Practical examples: which model structure would you select?



Practical examples: which model structure would you select?



?

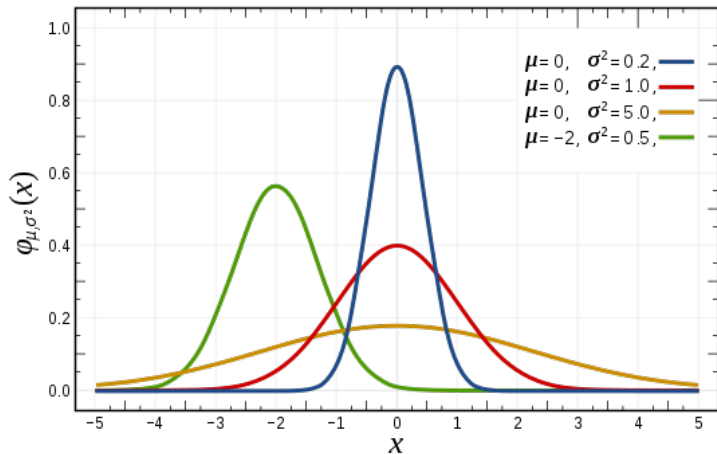
Connecting the pieces together

The pieces of the puzzle, up to now:

- estimating the model parameters
 - validating the model structure
 - estimating the uncertainties
-
- validating the data

Connecting the pieces together - piece 3: estimating the uncertainty on the model parameters

Typical aim: estimating mean and variance, assuming a Gaussian distribution:



Simplest strategy: jackknifing

Purpose: estimate *bias* and *variance* of the estimator;

Simplest strategy: jackknifing

Purpose: estimate *bias* and *variance* of the estimator; steps:

- 1 define the reduced datasets $X_{[i]} := \{X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n\}$

Simplest strategy: jackknifing

Purpose: estimate *bias* and *variance* of the estimator; steps:

- 1 define the reduced datasets $X_{[i]} := \{X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n\}$
- 2 compute the corresponding estimates $\hat{\theta}_{[i]} = \hat{\theta}(X_{[i]})$

Simplest strategy: jackknifing

Purpose: estimate *bias* and *variance* of the estimator; steps:

- 1 define the reduced datasets $X_{[i]} := \{X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n\}$
- 2 compute the corresponding estimates $\hat{\theta}_{[i]} = \hat{\theta}(X_{[i]})$
- 3 compute the average “reduced” estimate $\hat{\theta}_{\text{ave}} := \frac{1}{n} \sum_{i=1}^n \hat{\theta}_{[i]}$

Simplest strategy: jackknifing

Purpose: estimate *bias* and *variance* of the estimator; steps:

- 1 define the reduced datasets $X_{[i]} := \{X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n\}$
- 2 compute the corresponding estimates $\widehat{\theta}_{[i]} = \widehat{\theta}(X_{[i]})$
- 3 compute the average “reduced” estimate $\widehat{\theta}_{\text{ave}} := \frac{1}{n} \sum_{i=1}^n \widehat{\theta}_{[i]}$
- 4 compute the estimated bias and variance

$$\widehat{\text{bias}}(\widehat{\theta})_{\text{jk}} := \frac{n-1}{n} \sum_{i=1}^n (\widehat{\theta}_{[i]} - \widehat{\theta}) \quad \widehat{\text{var}}(\widehat{\theta})_{\text{jk}} := \frac{n-1}{n} \sum_{i=1}^n (\widehat{\theta}_{[i]} - \widehat{\theta}_{\text{ave}})^2$$

Simplest strategy: jackknifing

Purpose: estimate *bias* and *variance* of the estimator; steps:

- 1 define the reduced datasets $X_{[i]} := \{X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n\}$
- 2 compute the corresponding estimates $\widehat{\theta}_{[i]} = \widehat{\theta}(X_{[i]})$
- 3 compute the average “reduced” estimate $\widehat{\theta}_{\text{ave}} := \frac{1}{n} \sum_{i=1}^n \widehat{\theta}_{[i]}$
- 4 compute the estimated bias and variance

$$\widehat{\text{bias}}(\widehat{\theta})_{\text{jk}} := \frac{n-1}{n} \sum_{i=1}^n (\widehat{\theta}_{[i]} - \widehat{\theta}) \quad \widehat{\text{var}}(\widehat{\theta})_{\text{jk}} := \frac{n-1}{n} \sum_{i=1}^n (\widehat{\theta}_{[i]} - \widehat{\theta}_{\text{ave}})^2$$

→ can be generalized to k deletions

More sophisticated strategy: bootstrapping

More sophisticated strategy: bootstrapping

- 1 generate B new datasets with the same dimension as the original one sampling with replacement
- 2 compute the B estimates
- 3 compute some statistics on the B estimates

More sophisticated strategy: bootstrapping

- 1 generate B new datasets with the same dimension as the original one sampling with replacement
- 2 compute the B estimates
- 3 compute some statistics on the B estimates
- 4 if one wants to have a direct estimate of the predictive performance, average the performance of the various B estimates on the corresponding “out-of-bag” samples

More sophisticated strategy: bootstrapping

- 1 generate B new datasets with the same dimension as the original one sampling with replacement
- 2 compute the B estimates
- 3 compute some statistics on the B estimates
- 4 if one wants to have a direct estimate of the predictive performance, average the performance of the various B estimates on the corresponding “out-of-bag” samples

? how big shall B be ?

Flashback: bias and variance tradeoff

Alternative strategies:

- jackknifing
- bootstrapping
- CV (*i.e., go directly towards estimating the MSE / whatever performance index is desired*)

Comparisons

- the bootstrap handles skewed distributions better
- the jackknife is suitable for smaller original data samples

how to use the expressions for the estimated variance

Using \widehat{P} for finding confidence intervals on $\widehat{\theta}$

$$(\widehat{\theta} - \theta) \sim \mathcal{N}(0, \widehat{P})$$

Using \widehat{P} for finding confidence intervals on $\widehat{\theta}$

$$(\widehat{\theta} - \theta) \sim \mathcal{N}(0, \widehat{P}) \quad \Longrightarrow \quad (\widehat{\theta}^{(k)} - \theta^{(k)}) \sim \mathcal{N}(0, \widehat{P}_{(kk)})$$

Using \widehat{P} for finding confidence intervals on $\widehat{\theta}$

$$(\widehat{\theta} - \theta) \sim \mathcal{N}(0, \widehat{P}) \quad \implies \quad (\widehat{\theta}^{(k)} - \theta^{(k)}) \sim \mathcal{N}(0, \widehat{P}_{(kk)})$$

$$\implies \quad \mathbb{P} \left[\left| \widehat{\theta}^{(k)} - \theta^{(k)} \right| \geq \alpha \right] \approx \sqrt{\frac{1}{2\pi \widehat{P}_{(kk)}}} \int_{|x| \geq \alpha} \exp \left(-\frac{1}{2} \frac{x^2}{\widehat{P}_{(kk)}} \right) dx$$

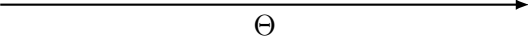
Using \widehat{P} for finding confidence intervals on $\widehat{\theta}$

$$(\widehat{\theta} - \theta) \sim \mathcal{N}(0, \widehat{P}) \quad \Longrightarrow \quad (\widehat{\theta}^{(k)} - \theta^{(k)}) \sim \mathcal{N}(0, \widehat{P}_{(kk)})$$

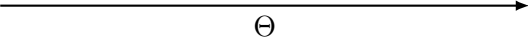
$$\Longrightarrow \quad \mathbb{P} \left[\left| \widehat{\theta}^{(k)} - \theta^{(k)} \right| \geq \alpha \right] \approx \sqrt{\frac{1}{2\pi \widehat{P}_{(kk)}}} \int_{|x| \geq \alpha} \exp \left(-\frac{1}{2} \frac{x^2}{\widehat{P}_{(kk)}} \right) dx$$

$$\Longrightarrow \quad (\widehat{\theta} - \theta)^T \widehat{P}^{-1} (\widehat{\theta} - \theta) \sim \chi^2(K)$$

Confidence intervals: these mysterious objects...

$$\begin{cases} \hat{\theta} \in \Theta \\ C \subseteq \Theta \end{cases}$$


Confidence intervals: these mysterious objects...

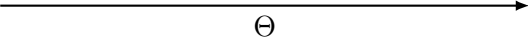
$$\begin{cases} \hat{\theta} \in \Theta \\ C \subseteq \Theta \end{cases}$$


Definition 1 (Confidence interval)

$C : \mathcal{D} \mapsto 2^{\Theta}$ is a C.I. with confidence level α if

$$\inf_{\theta \in \Theta} \mathbb{P}[\mathcal{D} : C(\mathcal{D}) \ni \theta] \geq \alpha$$

Confidence intervals: these mysterious objects...

$$\left\{ \begin{array}{l} \hat{\theta} \in \Theta \\ C \subseteq \Theta \end{array} \right.$$


Definition 1 (Confidence interval)

$C : \mathcal{D} \mapsto 2^{\Theta}$ is a C.I. with confidence level α if

$$\inf_{\theta \in \Theta} \mathbb{P}[\mathcal{D} : C(\mathcal{D}) \ni \theta] \geq \alpha$$

a C.I. is not a statement about θ

?

Connecting the pieces together

The pieces of the puzzle, up to now:

- estimating the model parameters
- validating the model structure
- estimating the uncertainties

- validating the data

Connecting the pieces together - piece 4: validating the data

Can happen at every step

- validating the data
- estimating the model parameters
- validating the data
- validating the model structure
- validating the data
- estimating the uncertainties
- validating the data

Connecting the pieces together - piece 4: validating the data

Most important strategies, up to now:

- Hotelling's T^2
- F- and Q-residuals (*to be seen now*)

Hotelling's T^2

Fundamental question: are these

$$x_1, \dots, x_{n_x} \quad y_1, \dots, y_{n_y} \quad (4)$$

identically distributed?

Hotelling's T^2

Fundamental question: are these

$$x_1, \dots, x_{n_x} \quad y_1, \dots, y_{n_y} \quad (4)$$

identically distributed? Algorithm:

$$\bar{x} := \frac{1}{n_x} \sum_i x_i \quad \bar{y} := \frac{1}{n_y} \sum_i y_i \quad (5)$$

$$\Sigma_x := \frac{1}{n_x - 1} \sum_i (x_i - \bar{x}) (x_i - \bar{x})^T \quad \Sigma_y := \frac{1}{n_y - 1} \sum_i (y_i - \bar{y}) (y_i - \bar{y})^T \quad (6)$$

$$\Sigma := \frac{(n_x - 1) \Sigma_x + (n_y - 1) \Sigma_y}{n_x + n_y - 2} \quad (7)$$

Hotelling's T^2

Fundamental question: are these

$$x_1, \dots, x_{n_x} \quad y_1, \dots, y_{n_y} \quad (4)$$

identically distributed? Algorithm:

$$\bar{x} := \frac{1}{n_x} \sum_i x_i \quad \bar{y} := \frac{1}{n_y} \sum_i y_i \quad (5)$$

$$\Sigma_x := \frac{1}{n_x - 1} \sum_i (x_i - \bar{x}) (x_i - \bar{x})^T \quad \Sigma_y := \frac{1}{n_y - 1} \sum_i (y_i - \bar{y}) (y_i - \bar{y})^T \quad (6)$$

$$\Sigma := \frac{(n_x - 1) \Sigma_x + (n_y - 1) \Sigma_y}{n_x + n_y - 2} \quad (7)$$

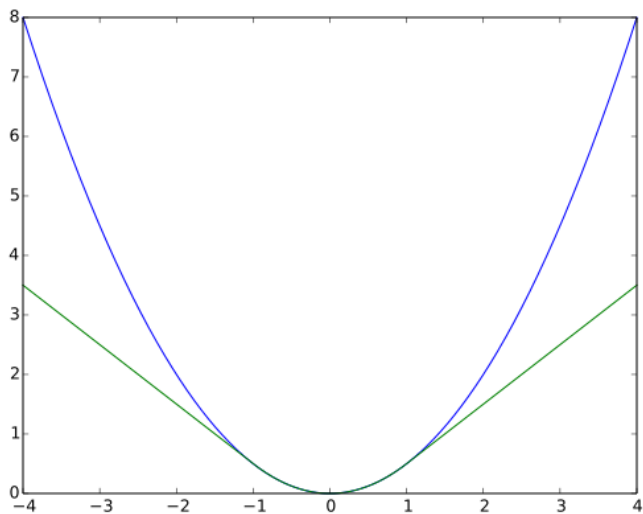
$$t^2 := \frac{n_x n_y}{n_x + n_y} (\bar{x} - \bar{y}) \Sigma (\bar{x} - \bar{y})^T \quad (8)$$

Causes of outliers

- measurement error
- heavy tail distributions
- mixture models

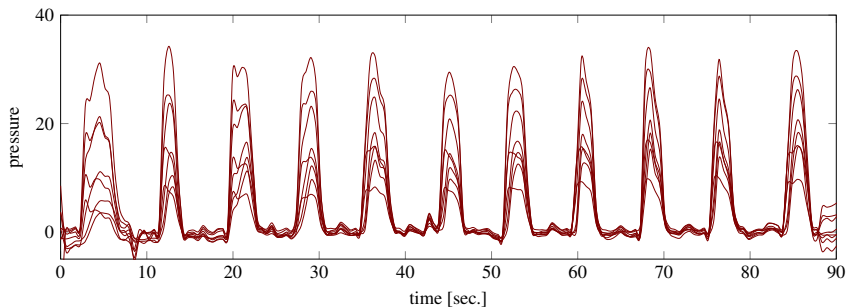
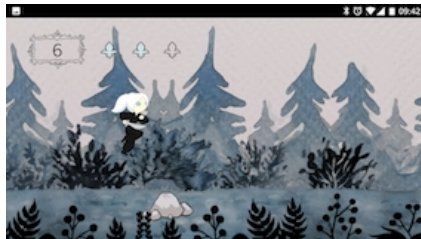
Norms to treat outliers

Norms to treat outliers



?

Hands-on example - vaginal pressure data



Vaginal pressure data: state variables

- $m_a(t) :=$ number of motor units that are in an *active state* at time t , and that are activated by a voluntary drive
- $m_f(t) :=$ number of motor units that are in a *fatigued state* at time t
- $m_r(t) :=$ number of motor units that are in a *resting state* at time t
- $u(t) :=$ muscular activation signal, sometimes referred to as the “brain stimulus” or “brain force”
- $M :=$ total number of motor units present in the muscles, assumed to be constant over time, i.e., $m_a(t) + m_f(t) + m_r(t) = M$ for all t

Vaginal pressure data: state dynamics

$$\dot{m}_a(t) = -\theta_{a \rightarrow f} m_a(t) + \theta_{f \rightarrow a} m_f(t) + u(t) \theta_{r \rightarrow a} m_r(t) - (1 - u(t)) \theta_{a \rightarrow r} m_a(t) \quad (9)$$

Vaginal pressure data: state dynamics

$$\dot{m}_a(t) = -\theta_{a \mapsto f} m_a(t) + \theta_{f \mapsto a} m_f(t) + u(t) \theta_{r \mapsto a} m_r(t) - (1 - u(t)) \theta_{a \mapsto r} m_a(t) \quad (9)$$

$$\left\{ \begin{array}{l} m_f(k+1) = \phi_{f \mapsto a} m_f(k) + (1 - \phi_{a \mapsto f}) m_a(k) \\ m_a(k+1) = \phi_{a \mapsto f} m_a(k) + (1 - \phi_{f \mapsto a}) m_f(k) \\ \quad + u(k) \phi_{r \mapsto a} m_r(k) - (1 - u(k)) \phi_{a \mapsto r} m_a(k) \\ m_r(k) = M - m_a(k) - m_f(k) \end{array} \right. \quad (10)$$

Vaginal pressure data: state dynamics

$$\dot{m}_a(t) = -\theta_{a \mapsto f} m_a(t) + \theta_{f \mapsto a} m_f(t) + u(t) \theta_{r \mapsto a} m_r(t) - (1 - u(t)) \theta_{a \mapsto r} m_a(t) \quad (9)$$

$$\left\{ \begin{array}{l} m_f(k+1) = \phi_{f \mapsto a} m_f(k) + (1 - \phi_{a \mapsto f}) m_a(k) \\ m_a(k+1) = \phi_{a \mapsto f} m_a(k) + (1 - \phi_{f \mapsto a}) m_f(k) \\ \quad + u(k) \phi_{r \mapsto a} m_r(k) - (1 - u(k)) \phi_{a \mapsto r} m_a(k) \\ m_r(k) = M - m_a(k) - m_f(k) \end{array} \right. \quad (10)$$

$$\left\{ \begin{array}{l} \phi_{f \mapsto a} := 1 - \theta_{f \mapsto a} T \\ \phi_{a \mapsto f} := 1 - \theta_{a \mapsto f} T \\ \phi_{a \mapsto r} := \theta_{a \mapsto r} T \\ \phi_{r \mapsto a} := \theta_{r \mapsto a} T \end{array} \right. \quad (11)$$

Vaginal pressure data: estimation

$$\begin{aligned} m_a(k+1) = & \left(\phi_{a \rightarrow f} - \phi_{a \rightarrow r} - (\phi_{r \rightarrow a} - \phi_{a \rightarrow r}) u(k) \right) m_a(k) \\ & + \left(1 - \phi_{f \rightarrow a} - \phi_{r \rightarrow a} u(k) \right) (1 - \phi_{a \rightarrow f}) \left(\sum_{\tau=0}^{k-1} \phi_{f \rightarrow a}^{k-1-\tau} m_a(\tau) \right) + \phi_{r \rightarrow a} M u(k) \end{aligned} \quad (12)$$

Vaginal pressure data: estimation

$$\begin{aligned} m_a(k+1) = & \left(\phi_{a \rightarrow f} - \phi_{a \rightarrow r} - (\phi_{r \rightarrow a} - \phi_{a \rightarrow r}) u(k) \right) m_a(k) \\ & + \left(1 - \phi_{f \rightarrow a} - \phi_{r \rightarrow a} u(k) \right) (1 - \phi_{a \rightarrow f}) \left(\sum_{\tau=0}^{k-1} \phi_{f \rightarrow a}^{k-1-\tau} m_a(\tau) \right) + \phi_{r \rightarrow a} M u(k) \end{aligned} \quad (12)$$

or, in a more compact way,

$$m_a(k+1) = f(m_a(k), u(k); \theta) \quad (13)$$

Vaginal pressure data: estimation

$$\begin{aligned} m_a(k+1) = & \left(\phi_{a \rightarrow f} - \phi_{a \rightarrow r} - (\phi_{r \rightarrow a} - \phi_{a \rightarrow r}) u(k) \right) m_a(k) \\ & + \left(1 - \phi_{f \rightarrow a} - \phi_{r \rightarrow a} u(k) \right) (1 - \phi_{a \rightarrow f}) \left(\sum_{\tau=0}^{k-1} \phi_{f \rightarrow a}^{k-1-\tau} m_a(\tau) \right) + \phi_{r \rightarrow a} M u(k) \end{aligned} \quad (12)$$

or, in a more compact way,

$$m_a(k+1) = f(m_a(k), u(k); \theta) \quad (13)$$

\implies *naturally leads to a nonlinear LS formulation*

What do we want to do?

- estimate the parameters of person A and the associated uncertainty on this estimate (*i.e., her physiological status*)
- check if the model structure is meaningful
- detect outliers in the measurements stream (*i.e., detect if the sensor is breaking*)

Vaginal pressure data: first step

how to compute $m_a(k)$ and $u(k)$ from the data?

