

# Paraméterek hibájának becslése II. Az illeszkedés jósága

Kormányos Andor

Komplex Rendszerek Fizikája Tanszék

2020. február 20.

# A paraméterek hibájának becslése más módon

A paraméterek hibájának becslésére már láttuk, hogy a Hesse-mátrix használható (ha a modell analitikusan ismer).

Egy jellegében más módszer: nézzük az illesztett paraméterek stabilitását úgy, hogy új “mérési eredményeket” szimulálunk: Monte Carlo módszerek (bonyolult)

# A paraméterek hibájának becslése más módon

A paraméterek hibájának becslésére már láttuk, hogy a Hesse-mátrix használható (ha a modell analitikusan ismer).

Egy jellegében más módszer: nézzük az illesztett paraméterek stabilitását úgy, hogy új “mérési eredményeket” szimulálunk: Monte Carlo módszerek (bonyolult)

Két egyszerűbb módszer:

- Jackknife módszer
- Bootstrapping

# Jackknife<sup>1</sup> módszer

Tekintsük a mérési pontokat, de minden lépésben hagyjunk ki egyet az illesztésből

- hagyjuk ki az  $i$ . pontot
- illesszük a modellt  $N - 1$  pontra
- legyen az illesztett paraméterek vektora  $\mathbf{a}_i$

Minden egyes mérési pontra megismételve összesen  $N$  különböző paramétervektort kapunk

- ezek átlaga lesz a becsült paramétervektor

$$\mathbf{a}_{\text{jack}} = \frac{1}{N} \sum_i \mathbf{a}_i$$

- ezek szórása az illesztett paraméterek standard hibája

$$\sigma_{\text{jack}}^2 = \frac{N-1}{N} \sum_i (\mathbf{a}_i - \mathbf{a}_{\text{jack}})^2$$

---

<sup>1</sup>jackknife = bicska

# Bootstrapping

Jelöljük az eredeti,  $N$  db adatot tartalmazó adathalmazt  $\mathcal{D}_{(0)}$ -val. Ennek segítségével új, szintetikus adathalmazokat generálunk:  $\mathcal{D}^{(S)}, \mathcal{D}_2^{(S)}, \dots$ , melyek szintén  $N$  db adatot tartalmaznak ("S": synthetic).

Hogyan?

# Bootstrapping

Jelöljük az eredeti,  $N$  db adatot tartalmazó adathalmazt  $\mathcal{D}_{(0)}$ -val. Ennek segítségével új, szintetikus adathalmazokat generálunk:  $\mathcal{D}^{(S)}, \mathcal{D}_2^{(S)}, \dots$ , melyek szintén  $N$  db adatot tartalmaznak ("S": synthetic).

Hogyan?

- véletlenszerűen kiválasztunk  $N$  adatot  $\mathcal{D}_{(0)}$ -ból, úgy, hogy minden választás után "visszahelyezzük" a kiválasztott adatot  $\mathcal{D}_{(0)}$ -ba
- $\Rightarrow \mathcal{D}_1^{(S)}, \mathcal{D}_2^{(S)}, \dots$ -ben egyes adatok *kétszer* (esetleg többször is) szerepelhetnek
- minden  $\mathcal{D}_i^{(S)}$  segítségével kiszámoljuk az  $\mathbf{a}_i^{(S)}$  illesztett paramétervektort
- $\mathbf{a}_i^{(S)}$  segítségével kapunk egy eloszlást az illesztett paraméterekre
- ennek az eloszlásnak a segítségével számolhatunk pl átlagot és szórást a paraméterekre

# Konfidencia tartomány

A bootstrapping-gal kapunk egy  $M$ -dimenziós eloszlást az  $\mathbf{a}$  paramétervektorra

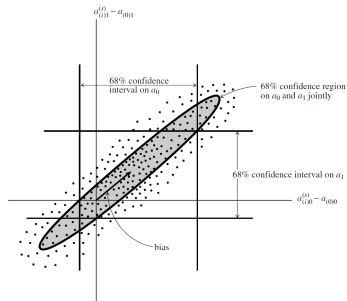
**konfidencia tartomány:** az a tartomány, amely adott valószínűséggel (pl 95%) tartalmazza a valódi paramétervektort

# Konfidencia tartomány

A bootstrapping-gal kapunk egy  $M$ -dimenziós eloszlást az  $\mathbf{a}$  paramétervektorra

**konfidencia tartomány:** az a tartomány, amely adott valószínűséggel (pl 95%) tartalmazza a valódi paramétervektort

Kétdimenziós példa:  $\mathbf{a}_i^{(s)} = (a_{(i)0}^{(s)}, a_{(i)1}^{(s)})$



**Figure:** A szimulált paraméterértékek 68% van a függőleges illetve vízszintes vonalakkal jelzett intervallumban



# További fontos kérdések

A paraméterek meghatározásával még nem ért véget a feladat:

- mekkora a meghatározott paraméterek hibája? ✓
- egyáltalán mennyire jó a modell? Hiába kicsi a meghatározott paraméterek hibája, ha rossz a modell, amit használunk

# Mennyire illeszkedik jól a modell?

A  $\chi^2$  definíció szerint:

$$\chi^2 = \sum_i \frac{[y_i - y(x_i|\mathbf{a})]^2}{\sigma_i^2}$$

Ha a mérést többször megismételnénk és mindig kiszámolnánk a  $\chi^2$ -t, akkor a  $\chi^2$ -ekre találnánk egy eloszlást

# Mennyire illeszkedik jól a modell?

A  $\chi^2$  definíció szerint:

$$\chi^2 = \sum_i \frac{[y_i - y(x_i|\mathbf{a})]^2}{\sigma_i^2}$$

Ha a mérést többször megismételnénk és mindig kiszámolnánk a  $\chi^2$ -t, akkor a  $\chi^2$ -ekre találnánk egy eloszlást

**Bizonyítás nélkül:** ha a modell lineárisan függ  $a_0, a_1, \dots, a_M$ -től, akkor a különböző  $\chi^2$  értékek eloszlása az összeg minimuma körül  $\nu = N - M$  szabadsági fokú  $\chi^2$  eloszlást követ. Itt  $N$  a mérési pontok, és  $M$  az illesztett paraméterek száma.

# Mennyire illeszkedik jól a modell?

A  $\chi^2$  definíció szerint:

$$\chi^2 = \sum_i \frac{[y_i - y(x_i|\mathbf{a})]^2}{\sigma_i^2}$$

Ha a mérést többször megismételnénk és mindig kiszámolnánk a  $\chi^2$ -t, akkor a  $\chi^2$ -ekre találnánk egy eloszlást

**Bizonyítás nélkül:** ha a modell lineárisan függ  $a_0, a_1, \dots, a_M$ -től, akkor a különböző  $\chi^2$  értékek eloszlása az összeg minimuma körül  $\nu = N - M$  szabadsági fokú  $\chi^2$  eloszlást követ. Itt  $N$  a mérési pontok, és  $M$  az illesztett paraméterek száma.

**Figyelem:** Ne keverjük össze a mérések és a modell segítségével számolt  $\chi^2$  összeget és a  $\chi^2$  eloszlást!

# $\chi^2$ eloszlás

**Definíció:** ha  $\xi_1, \xi_2, \dots, \xi_\nu$  valószínűségi változók standard normál eloszlással, akkor  $\xi_1^2 + \xi_2^2 + \dots + \xi_\nu^2$  eloszlása  $\nu$  szabadsági fokú  $\chi^2$  eloszlás

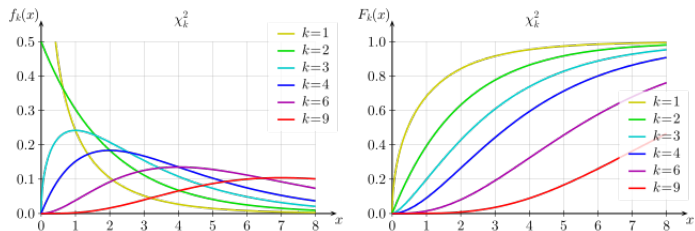


Figure: Különböző szabadsági fokú  $\chi^2$  sűrűség és eloszlás függvények

# Mennyire illeszkedik jól a modell?

**Megjegyzés:** mivel az  $a_0, a_1, \dots, a_M$  értékeket úgy határoztuk meg, hogy  $\chi^2$  összeg minimális legyen, ezért a fenti összegben nem minden tag független  $\Rightarrow$  ezért van  $\nu = N - M$  szabadsági fok

# Mennyire illeszkedik jól a modell?

**Megjegyzés:** mivel az  $a_0, a_1, \dots, a_M$  értékeket úgy határoztuk meg, hogy  $\chi^2$  összeg minimális legyen, ezért a fenti összegben nem minden tag független  $\Rightarrow$  ezért van  $\nu = N - M$  szabadsági fok

## Mikor fogadhatjuk el a modellt?

Egyszerű szabály: normáljuk  $\chi^2$ -et  $\nu = N - M$ -mel. A modell elfogadhatóan illeszkedik, ha  $\frac{\chi^2}{\nu} \approx 1$

# Mennyire illeszkedik jól a modell?

**Megjegyzés:** mivel az  $a_0, a_1, \dots, a_M$  értékeket úgy határoztuk meg, hogy  $\chi^2$  összeg minimális legyen, ezért a fenti összegben nem minden tag független  $\Rightarrow$  ezért van  $\nu = N - M$  szabadsági fok

## Mikor fogadhatjuk el a modellt?

Egyszerű szabály: normáljuk  $\chi^2$ -et  $\nu = N - M$ -mel. A modell elfogadhatóan illeszkedik, ha  $\frac{\chi^2}{\nu} \approx 1$

Pontosabban:  $\chi^2$  statisztikus átlaga  $\nu$ , ekörül a szórása  $\sqrt{2\nu}$ .



## Példa: parabola illesztése 5 pontra

Emlékeztető: feltételezve, hogy  $\sigma_i = 0.1$ , azt kaptuk, hogy

$$\mathbf{a} = \begin{bmatrix} 1.049 \\ -0.020 \\ 0.986 \end{bmatrix}$$

Az illesztés jósága:

$$\begin{aligned} \chi^2 &= (\mathbf{X}\mathbf{a} - \mathbf{y}/\sigma_i)^2 \\ &= 4.11 \end{aligned}$$

$$\frac{\chi^2}{\nu} = \frac{4.11}{5 - 3} = 2.06$$

Ez még a szóráson belül van.

