

Zbrani zapiski za 3. letnik

Patrik Žnidaršič

Prevedeno dne 2. junij 2024

Zahvala Jakobu Schraderju in Matiji Fajfarju za raznovrstne popravke v zapiskih.

Kazalo

1	Analiza	3	7
1.1	Splošno		8
1.2	Linearna NDE prvega reda		9
1.3	Prvi integral enačbe		11
1.4	Parametrično reševanje		12
1.4.1	Lagrangeova in Clairontova enačba		13
1.4.2	Ovojnice družin krivulj		14
1.5	Enačbe drugega reda		15
1.6	Eksistenčni izrek		16
1.7	Sistemi linearnih NDE		20
1.8	Linearne NDE višjega reda		26
1.8.1	Enačbe s konstantnimi koeficienti		28
1.8.2	Linearizacija		30
1.9	Variacijski račun		30
1.9.1	Vezani ekstremini		34
2	Mehanika		37
2.1	Osnove Newtonove mehanike		38
2.2	Premočrtno gibanje		44
2.3	Gibanje po krivulji		46
2.4	Gibanje v polju centralne sile		48
2.5	Relativno gibanje		52
2.6	Sistem materialnih točk		57
2.7	Togo telo		58
2.7.1	Prosta vrtavka		63
2.7.2	Eulerjevi koti		64
3	Uvod v numerične metode		67
3.1	Računske napake		68
3.2	Nelinearne enačbe		69
3.2.1	Bisekcija		70
3.2.2	Navadna iteracija		70
3.2.3	Tangentna metoda		71
3.2.4	Sekantna metoda		72
3.2.5	Ostale metode		73
3.2.6	Ničle polinomov		73

3.2.7	Durand-Kernerjeva metoda	75
3.3	Sistemi linearnih enačb	75
3.3.1	Matrične norme	75
3.3.2	Občutljivost sistema linearnih enačb	79
3.3.3	LU razcep	80
3.3.4	Razcep Choleskega	84
3.4	Sistemi nelinearnih enačb	85
3.5	Linearni problemi najmanjših kvadratov	87
3.5.1	Normalni sistem	87
3.5.2	QR razcep	88
3.5.3	Gram-Schmidtova ortogonalizacija	88
3.5.4	Givensove rotacije	89
3.5.5	Householderjeva zrcaljenja	90
3.6	Lastne vrednosti	91
3.6.1	Potenčna metoda	92
3.6.2	Inverzna iteracija	93
3.6.3	Ortogonalna iteracija	94
3.6.4	QR iteracija	95
3.7	Polinomska interpolacija	98
3.7.1	Lagrangeova oblika	98
3.7.2	Deljene difference	99
3.8	Numerično integriranje	100
3.8.1	Newton-Cotesove formule	101
3.8.2	Napake pri numeričnem integriranju	102
3.8.3	Gaussove kvadraturene formule	103
3.9	Diferencialne enačbe	104
3.9.1	Runge-Kutta metode	105
4	Verjetnost	107
4.1	Izidi, dogodki, verjetnosti	108
4.1.1	Pogojna verjetnost in neodvisnost	110
4.1.2	Neodvisnost dogodkov	111
4.2	Slučajne spremenljivke in porazdelitve	112
4.2.1	Slučajni vektorji	117
4.2.2	Neodvisnost slučajnih spremenljivk	117
4.2.3	Pričakovana vrednost diskretnih spremenljivk	118
4.2.4	Večrazsežne zvezne porazdelitve	119
4.2.5	Pogojne pričakovane vrednosti	122
4.3	Rodovne funkcije	123
4.3.1	Procesi razvejanja	125
4.3.2	Panjerjeva rekurzija	125
4.4	Tabele	126

5	Algebra 3	127
5.1	Reševanje polinomskih enačb	128
5.1.1	Rešljive grupe	132
5.1.2	Rešljivost polinomskih enačb z radikali	133
5.2	Moduli	134
5.2.1	Projektivni moduli	139
5.2.2	Tenzorski produkt modulov	141
5.2.3	Skrčitev in razširitev skalarjev	143
5.2.4	Eksaktna zaporedja modulov	143
5.3	Teorija kategorij	145
5.3.1	Univerzalne konstrukcije	147
6	Analiza 4	149
6.1	Osnovni tipi PDE	150
6.2	Kvazilinearne enačbe prvega reda v dveh spremenljivkah	151
6.2.1	Linearna PDE	154
6.2.2	Ovojnica družine ravnin	154
6.3	Nelinearne enačbe prvega reda	155
6.3.1	Cauchyjeva naloga za PDE prvega reda	157
6.4	Lagrangeova metoda	160
6.4.1	Odvisnost funkcij	162
6.5	Enačbe drugega reda	164
6.5.1	Cauchyjev problem	166
6.6	Valovna enačba na realni osi	168
6.7	Toplotna enačba	170
6.8	Sturm-Liouvilleova teorija	173
6.9	Harmonične funkcije	176
6.9.1	Dirichletov problem	177
7	Izbrane teme iz analize podatkov	179
7.1	Linearna regresija	181
7.2	Logistična regresija	182
7.3	Najbližji sosedi	183
7.4	Vrednotenje napovednih modelov	183
7.5	Odločitvena drevesa	185
7.6	Metoda podpornih vektorjev	186
7.7	Ansambli napovednih modelov	188
7.8	Nevronske mreže	190
7.9	Nenadzorovano učenje	192
7.10	Krčenje razsežnosti	193
7.11	Manjkajoče vrednosti	194
7.12	Neenakomerne porazdelitve	195

8	Numerična linearna algebra	197
8.1	Singularni razcep	198
8.1.1	Aproksimacija z matrikami nižjega ranga	200
8.1.2	Regularizacija	200
8.2	Nesimetrični problem lastnih vrednosti	202
8.2.1	Implicitna QR iteracija	204
8.3	Simetrični problem lastnih vrednosti	205
8.3.1	Rayleighova iteracija	207
8.3.2	QR iteracija	208
8.3.3	Bisekcija	209
8.3.4	Jacobijeva metoda	210
8.3.5	Deli in vladaj	212
8.4	Računanje singularnega razcepa	213
8.4.1	Enostranska Jacobijeva metoda	213
8.4.2	Enostranska QR iteracija	214
8.4.3	Weylov izrek	214
8.5	Posplošitve problema lastnih vrednosti	216
8.5.1	Posplošen problem lastnih vrednosti	216
8.5.2	Kvadratni problem lastnih vrednosti	218
9	Statistika	219
9.1	Centralni limitni izrek	220
9.2	Konvergenca porazdelitev	224
9.3	Uvod v statistiko	227
9.3.1	Sklepna statistika	227
9.3.2	Opisna statistika	229
9.3.3	Ocenjevanje in napovedovanje	230
9.3.4	Vrednotenje cenilk in prediktorjev	231
9.4	Pričakovana vrednost in varianca slučajnih vektorjev	232
9.5	Pridobivanje cenilk	233
9.5.1	Metoda empirične porazdelitve	233
9.5.2	Metoda največjega verjetja	235
9.6	Večrazsežna normalna porazdelitev	236
9.7	Statistično sklepanje z nadzorovanim tveganjem	239

1 Analiza 3

1.1 Splošno

Definicija. Naj bo $F : I \times \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ zvezna funkcija in I interval v \mathbb{R} . NAVADNA DIFERENCIALNA ENAČBA PRVEGA REDA je enačba oblike $F(x, y(x), y'(x)) = 0$, kjer je $y(x)$ neka funkcija. Rešitev enačbe je vsaka funkcija $y_r(x) : I \rightarrow \mathbb{R}$, za katero velja enačba.

Opomba. NDE n -tega reda definiramo podobno kot enačbo oblike

$$F(x, y, y', y'', \dots, y^{(n)}) = 0.$$

Opomba. Smiselno je opazovati tudi enačbe, kjer je $F = (F_1, \dots, F_m)$ vektorska funkcija. Temu pravimo SISTEM NDE.

Opomba. Naj bo $y^{(n)} = F(x, y, y', \dots, y^{(n-1)})$ enačba reda n . Ta enačba je ekvivalentna primernemu sistemu $n \times n$ prvega reda; definirajmo $y_1 = y$, $y_2 = y'$, ..., $y_{n-1} = y^{(n-2)}$. Tedaj dobimo enačbo $y'_n = F(x, y_1, \dots, y_n)$.

Definicija. INTEGRALSKA KRIVULJA γ vektorskega polja $F : \Omega \subseteq \mathbb{R}^n \rightarrow \mathbb{R}^n$ skozi točko $x_0 \in \Omega$ je krivulja $\gamma : [0, b) \rightarrow \Omega$, za katero velja

- v vsaki točki t je $\dot{\gamma}(t) = F(\gamma(t))$,
- $\gamma(0) = x_0$.

Vprašanje 1. Definiraj integralske krivulje.

Če prvi pogoj iz definicije zapišemo v koordinatah,

$$\begin{bmatrix} \dot{x}_1 \\ \vdots \\ \dot{x}_n \end{bmatrix} (t) = \begin{bmatrix} F_1(x_1, \dots, x_n) \\ \vdots \\ F_n(x_1, \dots, x_n) \end{bmatrix},$$

dobimo sistem n NDE prvega reda z n neznankami. Ta sistem ni eksplicitno odvisen od t ; takim sistemom pravimo AVTONOMNI SISTEMI. Pokazali bomo, da za vsako izbiro x_0 obstajata interval $[0, a)$ in krivulja γ , za katero veljata pogoja v definiciji.

Vsak neavtonomen sistem lahko prepišemo v avtonomnega, z uvedbo nove odvisne spremenljivke $v(t) = t$. Dobimo nov sistem

$$\begin{aligned} \dot{v} &= 1, \\ \dot{x}_1 &= F_1(v, x_1, \dots, x_n), \\ &\vdots \\ \dot{x}_n &= F_n(v, x_1, \dots, x_n). \end{aligned}$$

Partikularna rešitev tega sistema je tedaj integralska krivulja vektorskega polja $\vec{F}(v, \vec{x})$ v RAZŠIRJENEM FAZNEM PROSTORU $\mathbb{R} \times \Omega$ (Ω je običajen fazni prostor), ki ustreza primernemu začetnemu pogoju.

Vprašanje 2. Kako spremenimo neavtonomni sistem v avtonomnega?

V nekaterih primerih poznamo rešitev NDE. Če imamo enačbo z ločljivima spremenljivkama

$$\dot{x} = f(t)g(x),$$

lahko enačbo delimo z $g(x)$, in definiramo $h(x) = 1/g(x)$. Dobimo

$$h(x)\dot{x} = f(t).$$

Sedaj definiramo $H(x)$ kot primitivno funkcijo $h(x)$, in $F(t)$ kot primitivno funkcijo $f(t)$. Velja $\dot{H}(x) = \dot{F}(t)$, torej je $x(t) = H^{-1}(F(t) + C)$.

Vprašanje 3. Kako rešiš enačbo z ločljivima spremenljivkama?

Če imamo enačbo s homogeno desno stranjo, torej $\dot{x} = f(t, x)$, kjer velja $f(t, x) = f(\lambda t, \lambda x)$ za $\lambda \in \mathbb{R} \setminus \{0\}$, potem velja $f(t, x) = f(1, x/t)$. Vpeljemo novo spremenljivko $v = x/t$, in s kratkim računom pridemo do $\dot{x} = t\dot{v} + v$. Po drugi strani velja $\dot{x} = f(1, v)$, torej

$$\dot{v} = \frac{1}{t} (f(1, v) - v).$$

To je enačba z ločljivima spremenljivkama, ki jo znamo rešiti.

Vprašanje 4. Kako rešiš enačbo s homogeno desno stranjo?

1.2 Linearna NDE prvega reda

LINEARNA NDE PRVEGA REDA je enačba oblike

$$y' = f(x)y + g(x),$$

kjer sta $f(x)$ in $g(x)$ znani funkciji. Ta enačba je NEHOMOGENA z NEHOMOGENOSTJO $g(x)$. Njena HOMOGENIZACIJA je enačba

$$y' = f(x)y.$$

Predpostavimo, da je $y(x) \in \mathcal{C}^1([a, b])$. Oglejmo si operator $A : \mathcal{C}^1([a, b]) \rightarrow \mathcal{C}([a, b])$, definiran kot

$$Ay(x) = y'(x) - f(x)y.$$

Trditev. Preslikava A je linearen operator.

Dokaz je trivialen, in zato izpuščen. Vidimo, da je y rešitev homogene enačbe natanko tedaj, ko je $A(y) = 0$. Rešitev homogene enačbe je torej jedro preslikave A . Homogena enačba je enačba z ločljivimi spremenljivkami, torej jo znamo rešiti. Rešitve so oblike

$$y(x) = C \exp \left(\int_a^x f(\xi) d\xi \right)$$

za $C \in \mathbb{R}$. Množico teh rešitev označimo z R_h .

Trditev. Naj bosta y_1, y_2 rešitvi nehomogene enačbe. Tedaj je $y(x) = y_1(x) - y_2(x)$ rešitev homogene enačbe.

Dokaz. Izračun odvoda nam da

$$(y_1(x) - y_2(x))' = f(x)y_1(x) + g(x) - (f(x)y_2(x) + g(x)) = f(x)(y_1(x) - y_2(x)).$$

□

Definicija. Naj bo V nek vektorski prostor in $W \subseteq V$. Če obstaja tak vektorski podprostor $H \subseteq V$, da za poljubna $w_1, w_2 \in W$ velja $w_1 - w_2 \in H$, je W AFIN PODPROSTOR v V , modeliran z vektorskim podprostorom H .

Rešitve nehomogene enačbe so torej afin prostor, modeliran s prostorom R_h rešitev homogene enačbe. Če želimo poiskati splošno rešitev, poiščemo rešitev homogenega sistema, in neko partikularno rešitev. Partikularno rešitev dobimo z nastavkom

$$y_p(x) = C(x) \exp \left(\int_a^x f(\xi) d\xi \right),$$

temu postopku pravimo VARIACIJA KONSTANTE.

Vprašanje 5. Kako rešiš linearno NDE prvega reda? Utemelji postopek.

S tem znanjem lahko rešimo še dve posebni NDE. Prva je Bernoulijeva enačba

$$p(x)y' + q(x)y = r(x)y^\alpha(x)$$

za $\alpha \in \mathbb{R}$. V primeru $\alpha = 0$ ali $\alpha = 1$, je to nehomogena linearna enačba prvega reda. Sicer vpeljemo $z(x) = (y(x))^{1-\alpha}$ in računamo

$$p(x)z'(x) \frac{1}{1-\alpha} + q(x)z(x) = r(x)$$

oziroma

$$z'(x) + \frac{q(x)}{p(x)}(1-\alpha)z(x) = \frac{r(x)}{p(x)}(1-\alpha).$$

To je nehomogena linearna NDE prvega reda, torej jo znamo rešiti.

Vprašanje 6. Kako rešiš Bernoulijevo enačbo?

Druga taka enačba je Riccatyjeva enačba

$$y'(x) = a(x)y^2(x) + b(x)y(x) + c(x),$$

ki je v splošnem ne znamo rešiti. Poznamo pa dva načina obravnave, ki nas lahko včasih pripeljeta do rešitve. Denimo, da uganemo neko partikularno rešitev $y_p(x)$. Enačbo tedaj rešujemo z nastavkom $y(x) = y_p(x) + z(x)$ za neko neznano funkcijo z . Če to vstavimo v enačbo, dobimo

$$y_p' + z' = ay_p^2 + 2ay_pz + az^2 + by_p + bz + c,$$

členi y'_p , ay_p^2 , by_p in c odpadejo, ker tvorijo rešitev enačbe. Ostane torej

$$z' = (2ay_p + b)z + az^2,$$

kar je Bernoulijeva enačba, ki jo znamo rešiti.

Drug način za reševanje Riccatijeve enačbe je s pretvorbo na linearni sistem prvega reda. Vpeljemo $y = u/v$, s čimer dobimo

$$u'v - uv' = au^2 + buv + cv^2.$$

Ker imamo dve neznanki, potrebujemo še eno enačbo. Izberemo $u'v = buv + cv^2$. Iz tega izpeljemo $v' = -au$ in $u' = bu + cv$. Zapisano matrično

$$\begin{bmatrix} v' \\ u' \end{bmatrix} = \begin{bmatrix} 0 & -a \\ c & b \end{bmatrix} \begin{bmatrix} v \\ u \end{bmatrix}$$

Sistema v splošnem ne znamo rešiti, ker funkcije a, b, c niso konstantne. Lahko pa rešitev zapišemo v obliki neskončne vrste.

Vprašanje 7. Kako rešiš Riccatijevo enačbo?

1.3 Prvi integral enačbe

Splošna rešitev enačbe $y' = f(x, y)$ je enoparametrična družina funkcij $y = \phi(x, C)$. Denimo, da obstaja taka funkcija $u(x, y) : [a, b] \times M \rightarrow \mathbb{R}$ na razširjenem faznem prostoru, da zanjo velja $u(x, y(x, C)) = \text{konst.}$ za vsak C (konstanta je lahko drugačna za različne C). Taki funkciji pravimo PRVI INTEGRAL ENAČBE. Recimo, da velja $\partial_y u \neq 0$. Potem lahko iz enakosti izračunamo funkcijo $y(x, D)$, da velja $u(x, y(x, D)) = D$.

Trditev. Vsaka krivulja $y(x)$, ki je implicitno podana z enačbo $u(x, y) = D$, kjer je u prvi integral enačbe $y' = f(x, y)$, je rešitev te enačbe.

Dokaz. Naj bo $y_0(x)$ dana krivulja. Tedaj velja $u(x, y_0(x)) = D$. Če odvajamo po x , dobimo

$$\partial_x u + y'_0 \partial_y u = 0,$$

torej

$$y'_0 = -\frac{\partial_x u}{\partial_y u}.$$

Po drugi strani za vsako rešitev velja $y' = f(x, y)$, iz česar izpeljemo

$$f(x, y) = -\frac{\partial_x u}{\partial_y u}.$$

Sledi, da je y_0 res rešitev enačbe. □

Vsaka diferencialna enačba ima neskončno mnogo prvih integralov, vsakega lahko še transformiramo s poljubno $\psi : \mathbb{R} \rightarrow \mathbb{R}$.

Vprašanje 8. Kaj je prvi integral enačbe? Kako iz njega dobiš rešitev enačbe?

Imejmo dano vektorsko polje $F : \Omega \subseteq \mathbb{R}^2 \rightarrow \mathbb{R}^2$, podano z

$$F(x, y) = \begin{bmatrix} P(x, y) \\ Q(x, y) \end{bmatrix}$$

Poiščimo družino krivulj, ortogonalnih na polje $F(x, y)$, in jih parametrizirajmo z $\gamma(x) = (x, y(x))$. Izpeljemo lahko pogoj

$$y'(x) = -\frac{P(x, y)}{Q(x, y)},$$

kar je diferencialna enačba prvega reda. Recimo, da je polje potencialno. Tedaj obstaja taka funkcija $u : [a, b] \times \mathbb{R} \rightarrow \mathbb{R}$, da velja $\partial_x u = P$ in $\partial_y u = Q$. Krivulje, ki so ortogonalne na $\vec{\nabla} u$, so natanko izohipse ploskve $(x, y, u(x, y))$. Za izohipso velja $u(x, y(x)) = C$, iz česar z odvajanjem izpeljemo

$$y' = -\frac{\partial_x u}{\partial_y u}.$$

Potencial u je torej prvi integral zgornje enačbe. Ker lahko potencial poiščemo z integralom, enačbo v tem primeru znamo rešiti.

Če polje ni potencialno, imamo še vedno ortogonalne krivulje, torej še vedno velja $y' = -P/Q$. Denimo, da je $y(x, C)$ splošna rešitev. Če lahko poiščemo prvi integral u , mora obstajati funkcija $\lambda(x, y)$, za katero je polje $(\lambda P, \lambda Q)$ potencialno, oziroma $\partial_x u = \lambda P$ in $\partial_y u = \lambda Q$. Taki funkciji pravimo INTEGRIRUJOČI MNOŽITELJ. Če ga lahko najdemo, lahko rešimo enačbo; tega pa v splošnem ne znamo.

Vprašanje 9. Kako poiščeš družino krivulj, pravokotnih na dano vektorsko polje $F = (P, Q)$? Kaj je integrirujoči množitelj?

1.4 Parametrično reševanje

Naj bo NDE prvega reda podana implicitno,

$$F(x, y, y') = 0.$$

Denimo, da y' ne moramo eksplicitno izraziti z x, y , ali pa je ekspliciten izraz nepripraven. Na F pogledajmo nekoliko drugače; vsaka dovolj lepa funkcija treh spremenljivk podaja družino ploskev, $F(\xi, \eta, \zeta) = C$ je implicitna enačba ploskve v \mathbb{R}^3 za vsak $C \in \mathbb{R}$. Tako podano ploskev lahko parametriziramo. Naj bo $(u, v) \mapsto (\varphi(u, v), \psi(u, v), \chi(u, v))$ neka parametrizacija. Imamo tri pristope za reševanje, v odvisnosti od F .

Če y ne nastopa eksplicitno, torej $F(x, y') = 0$, nam enačba definira krivuljo. Parametriziramo jo z $\xi = \varphi(t)$ in $\eta = \psi(t)$, da velja $F(\varphi, \psi) = 0$. Za poljubno rešitev $t \mapsto (x(t), y(t))$ velja $\dot{y} = y'\dot{x}$, torej za $\varphi(t) = x(t)$ in $\psi(t) = y'(t)$ dobimo $\dot{y} = \chi\dot{\varphi}$, oziroma

$$y(t) = \int_0^t \chi(\tau) \dot{\varphi}(\tau) d\tau.$$

Dobimo parametrično izraženo rešitev $t \mapsto (\varphi(t), y(t))$.

Vprašanje 10. Kako parametrično rešiš enačbo $F(x, y') = 0$?

Če x ne nastopa eksplicitno, torej $F(y, y') = 0$, dobimo enačbo krivulje $F(\xi, \eta) = 0$, ki jo parametriziramo s $t \mapsto (\chi(t), \psi(t))$. Če označimo $\psi = y$ in $\chi = y'$, velja $\dot{\psi} = \chi\dot{x}$ oziroma

$$x(t) = \int_0^t \frac{\dot{\psi}(\tau)}{\chi(\tau)} d\tau,$$

torej je $t \mapsto (x(t), \psi(t))$ parametrično podana rešitev.

Vprašanje 11. Kako parametrično rešiš enačbo $F(y, y') = 0$?

V splošnem nam $F(x, y, y') = 0$ definira ploskev. Parametriziramo jo kot zgoraj z

$$x = \varphi(u, v) \qquad y = \psi(u, v) \qquad y' = \chi(u, v)$$

Naj bo $t \mapsto (x(u(t), v(t)), y(u(t), v(t)))$ neka rešitev naše enačbe. Potem je $t \mapsto (\varphi, \psi, \chi)$ krivulja na ploskvi. V nadaljevanju predpostavimo, da je preslikava $(u, v) \mapsto (x, y)$ obrnljiva, in izračunajmo

$$\begin{aligned} \dot{y} &= y_u \dot{u} + y_v \dot{v} = \psi_u \dot{u} + \psi_v \dot{v}, \\ \dot{x} &= \varphi_u \dot{u} + \varphi_v \dot{v}. \end{aligned}$$

Ker tudi tu velja $\dot{y} = y'\dot{x}$, izrazimo

$$u' = \frac{\dot{u}}{\dot{v}} = -\frac{\psi_v - \chi\varphi_v}{\psi_u - \chi\varphi_u}.$$

Dobili smo eksplicitno enačbo prvega reda v spremenljivki $u = u(v)$.

Vprašanje 12. Kako parametrično rešiš $F(x, y, y') = 0$?

1.4.1 Lagrangeova in Clairontova enačba

Lagrangeova enačba je enačba oblike

$$y = x\varphi(y') + \psi(y').$$

Rešujemo jo parametrično; $x = u$, $y' = v$ in $y = u\varphi(v) + \psi(v)$. Velja $dy = y'dx$, iz česar izpeljemo

$$(\varphi(v) - v)du + (u\varphi'(v) + \psi'(v))dv = 0.$$

Če je $\varphi(v) \neq v$, dobimo

$$(\phi(v) - v) \frac{du}{dv} + u\varphi'(v) + \psi'(v) = 0,$$

kar je linearna diferencialna enačba prvega reda, če pa je $\varphi(v) = v$, pa imamo Clairontovo enačbo

$$y = xy' + \psi(y').$$

To predelamo v

$$(u + \psi'(v))dv = 0,$$

in obravnavamo dva primera. Če je $dv = 0$, je y' konstanta, torej dobimo družino rešitev $y = Cx + \psi(C)$ (to vstavimo v enačbo; ni nujno vsak C dober). Če pa je $u + \psi'(v) = 0$, pa dobimo še eno rešitev.

Vprašanje 13. Kaj sta Lagrangeova in Clairontova enačba? Kako ju rešimo?

1.4.2 Ovojnice družin krivulj

Imejmo družino krivulj, podano implicitno z enačbo $F(x, y, C) = 0$. Denimo, da je družina taka, da obstaja krivulja, ki se v vsaki svoji točki dotika natanko enega člana družine. Taki krivulji pravimo OVOJNICA družine. Smiselno jo je parametrizirati z $C \mapsto (x(C), y(C))$, pri čemer se ovojnica v točki $(x(C), y(C))$ dotika člana družine s tem C . Denimo, da je parametrizacija regularna, torej za vsak C

$$(\partial_C x)^2 + (\partial_C y)^2 \neq 0.$$

Definirajmo

$$\phi(C) = F(x(C), y(C), C).$$

Ker funkcija izračuna F v točki na krivulji, je $\phi = 0$. Torej

$$\phi'(C) = \partial_x F \partial_C x + \partial_y F \partial_C y + \partial_C F = 0.$$

Če je $t \mapsto (x(t), y(t))$ parametrizacija C_0 -tega člana družine, velja

$$\partial_t F(x(t), y(t), C_0) = \partial_x F \dot{x} + \partial_y F \dot{y} = 0$$

v točki dotika z ovojnico. Vektor $[\dot{x}, \dot{y}]^T$ je vzporeden z $[\partial_C x, \partial_C y]^T$ v tej točki, torej je

$$\begin{bmatrix} \partial_x F \\ \partial_y F \end{bmatrix} \perp \begin{bmatrix} \dot{x} \\ \dot{y} \end{bmatrix} \parallel \begin{bmatrix} \partial_C x \\ \partial_C y \end{bmatrix}$$

in zato

$$\partial_x F \partial_C x + \partial_y F \partial_C y = 0.$$

Torej v točki dotika velja

$$\partial_C F = 0.$$

Iz para enačb $F(x, y, C) = 0$ in $\partial_C F = 0$ dobimo vse točke na ovojnici.

Vprašanje 14. Kaj je ovojnica družine krivulj? Kako jo izračunaš? Izpelji.

1.5 Enačbe drugega reda

Najpomembnejša enačba drugega reda je drugi Newtonov zakon. Malce posplošeno ima obliko

$$\ddot{x}_i = F_i(x_1, \dots, x_n)$$

za $i = 1, \dots, n$. Če vpeljemo $p = \dot{x}$ in $q = x$, dobimo sistem prvega reda

$$\begin{aligned}\dot{q}_i &= p_i \\ \dot{p}_i &= F_i(q_1, \dots, q_n)\end{aligned}$$

Sistem lahko še posplošimo. Naj bosta $F, G : M \subseteq \mathbb{R}^{2n} \rightarrow \mathbb{R}^n$ preslikavi iz faznega prostora M . Zanima nas časovni razvoj sistema, ki je podan z enačbami

$$\begin{aligned}\dot{q} &= G(q, p), \\ \dot{p} &= F(q, p).\end{aligned}$$

Posebej pomembni so sistemi, za katere obstaja funkcija $H : M \rightarrow \mathbb{R}$, za katero velja

$$\begin{aligned}\frac{\partial H}{\partial q_i} &= -F_i, \\ \frac{\partial H}{\partial p_i} &= G_i.\end{aligned}$$

Taki funkciji pravimo HAMILTONIAN.

Izrek. Če Hamiltonian obstaja, potem je prvi integral sistema.

Dokaz. Naj bo $t \mapsto (q(t), p(t))$ neka rešitev sistema. Potem imamo

$$\partial_t H(q, p) = \partial_q H \dot{q} + \partial_p H \dot{p} = \begin{bmatrix} -F & F \end{bmatrix} \cdot \begin{bmatrix} G \\ G \end{bmatrix} = 0$$

□

Vprašanje 15. Kaj je Hamiltonian? Dokaži, da je prvi integral.

Definicija. Naj bo podana funkcija $H : M \subseteq \mathbb{R}^{2n} \rightarrow \mathbb{R}$. Sistem enačb

$$\begin{aligned}\dot{q} &= \partial_p H \\ \dot{p} &= -\partial_q H\end{aligned}$$

je HAMILTONSKI SISTEM S HAMILTONSKO FUNKCIJO H .

Da bo sistem $\dot{q} = G(q, p), \dot{p} = F(q, p)$ Hamiltonski, mora obstajati funkcija H , za katero velja $\partial_p H = G$ in $\partial_q H = -F$. Zapisano v drugačni obliki

$$\begin{bmatrix} \dot{q} \\ \dot{p} \end{bmatrix} = \begin{bmatrix} 0 & I \\ -I & 0 \end{bmatrix} \cdot \begin{bmatrix} \partial_q H \\ \partial_p H \end{bmatrix} = \begin{bmatrix} 0 & I \\ -I & 0 \end{bmatrix} \vec{\nabla} H,$$

torej

$$\vec{\nabla} \cdot H = \begin{bmatrix} -F \\ G \end{bmatrix}.$$

Da bo sistem hamiltonski, mora biti torej polje $[-F, G]^T$ potencialno.

Vprašanje 16. Pod katerim pogojem je sistem hamiltonski? Dokaži.

Definicija. ELIPTIČNI INTEGRAL PRVE VRSTE z modulom m je funkcija, podana s predpisom

$$F(x; m) = \int_0^x (1 - m \sin^2 \xi)^{-1/2} d\xi.$$

Inverzna funkcija te funkcije se imenuje JACOBIJEVA AMPLITUDA, velja

$$y = F(x; m) \Leftrightarrow x = \operatorname{am}(y; m).$$

Vprašanje 17. Obravnavaj gravitacijsko nihalo.

Odgovor: Gravitacijsko nihalo je oblike

$$\ddot{q} = -\sin q.$$

Temu sistemu pripada hamiltonska funkcija

$$H(q, p) = \frac{1}{2}p^2 - (\cos q - 1).$$

Vzdolž neke rešitve $t \mapsto (q(t), p(t))$ je to konstanta, in velja

$$\frac{1}{2}\dot{q}^2 - \cos q + 1 = E.$$

Enačbo lahko prevedemo v

$$\frac{dq}{\sqrt{1 - \frac{2}{E} \sin^2 q/2}} = \sqrt{2E} dt.$$

Rešitev je

$$q = 2 \operatorname{am} \left(\sqrt{\frac{E}{2}} t + C; \frac{2}{E} \right).$$

☒

1.6 Eksistenčni izrek

Izrek (Eksistenčni). Naj bo vektorsko polje $F(t, x)$ podano na valju

$$\mathcal{C}_{a,b} = \{(t, x) \mid |t - t_0| \leq a, \|x - x_0\| \leq b\}$$

za neki par t_0, x_0 in $a, b \in \mathbb{R}^+$. Naj bo $F(t, x)$ na $\mathcal{C}_{a,b}$ zvezno in naj bo $F(t, x) : \mathcal{C}_{a,b} \rightarrow \mathbb{R}^n$ Lipschitzova glede na x pri vsakem t . Alternativno je lahko preslikava $F(t, x)$ odvedljiva po x pri vsakem t in $\|D_x F\|$ na $\mathcal{C}_{a,b}$ omejeno število. Potem obstaja natanko ena rešitev začetnega problema

$$\dot{x} = F(t, x) \quad x(t_0) = x_0$$

za vsak x_0 . Rešitev $\varphi(t)$ obstaja na $[t_0 - a', t_0 + a']$ za nek $a' \leq a$. Še več: za družino začetnih problemov $\dot{x} = F(t, x), x(t_0) = \hat{x}$, kjer je $\|x_0 - \hat{x}\|$ dovolj majhno, obstajata $0 < a' \leq a$ in funkcija $g(t, x) : \mathcal{C}_{a',b'} \rightarrow \mathbb{R}^n$, za katero velja

- je zvezna na obe spremenljivki,
- $\partial_t g(t, x) = F(t, x)$,
- $g(t_0, \hat{x}) = \hat{x}$.

Vprašanje 18. Formuliraj eksistenčni izrek.

Začetni problem $\dot{x} = F(t, x), x(t_0) = x_0$ je ekvivalenten integralni enačbi

$$x(t) = x_0 + \int_{t_0}^t F(\tau, x(\tau)) d\tau.$$

Dokaz bo uporabil Picardov operator

$$Af(x) = x_0 + \int_{t_0}^t F(\tau, f(\tau)) d\tau$$

na posebnem funkcijskem prostoru.

Naj bo

$$\mathcal{C}_{a,b} = \{(t, x) \mid |t - t_0| \leq a, \|x - x_0\| \leq b\}.$$

Predpostavimo, da F ustreza predpostavkam izreka. Naj bo L Lipschitzova konstanta za F na $\mathcal{C}_{a,b}$ in

$$C = \max_{(t,x) \in \mathcal{C}_{a,b}} \|F(t, x)\|.$$

Velikokrat lahko za L vzamemo kar maksimum norme Jacobijeve matrike na $\mathcal{C}_{a,b}$. Označimo s K_0 stožec

$$K_0 = \{(t, x) \mid |t - t_0| \leq a', \|x - x_0\| \leq C |t - t_0|\},$$

kjer je a' dovolj majhen, da velja $K_0 \subseteq \mathcal{C}_{a,b}$. Sedaj sprostimo začetno vrednost x_0 . Naj bo $\|\hat{x} - x_0\| < b'$, $K_{\hat{x}} = K_0 + (\hat{x} - x_0)$ premik prostora in

$$K = \bigcup_{\|\hat{x} - x_0\| < b'} K_{\hat{x}},$$

kjer je b' dovolj majhen, da je $K \subseteq \mathcal{C}_{a,b}$. Za nas bo pomemben nov prostor $\mathcal{C}_{a',b'}$.

Začetni problem zapišimo nekoliko drugače. Spomnimo se: Iščemo $g(t, \hat{x})$, da bo $\partial_t g = F(t, g)$ in $g(t_0, \hat{x}) = \hat{x}$. Vpeljimo novo funkcijo $h(t, x) : \mathcal{C}_{a,b} \rightarrow \mathbb{R}^n$,

$$g(t, x) = x + h(t, x).$$

Velja

$$\begin{aligned}\partial_t h(t, x) &= F(t, g(t, x)) \\ h(t_0, x) &= g(t_0, x) - x = 0.\end{aligned}$$

Torej je h pri vsakem x rešitev začetnega problema $\dot{h} = F(t, x + h(t, x)), h(t_0, x) = 0$.

Definiramo

$$M = \{h(t, x) : \mathcal{C}_{a',b'} \rightarrow \mathbb{R}^n \mid h \text{ zvezna}, \|h(t, x)\| \leq C |t - t_0|\}.$$

Te preslikave zavzemajo vrednosti v stožcu K_0 .

Vprašanje 19. Povej postopek konstrukcije funkcijskega prostora v dokazu eksistenčnega izreka.

Opremimo M z maksimum normo (in s tem z metriko in topologijo)

$$\|h\| = \max_{(t,x) \in \mathcal{C}_{a',b'}} \|h(t, x)\|.$$

Trditev. *Prostor M je poln metrični prostor.*

Dokaz. Naj bo $\{h_n(t, x)\}_n$ Cauchyjevo zaporedje v M . Ker je \mathbb{R}^n poln, obstaja limita $\lim_{n \rightarrow \infty} h_n(t, x)$ za poljubna t, x . Ker je norma definirana z maksimumom, je konvergenca glede na to normo enakomerna, torej je

$$h(t, x) = \lim_{n \rightarrow \infty} h_n(t, x)$$

zvezna. Če je $h_n \in M$, velja $\|h_n(t, x)\| \leq C |t - t_0|$. To očitno velja tudi v limiti. \square

Vprašanje 20. Dokaži, da je ta funkcijski prostor poln.

Našo rešitev poiščemo kot limito iteracij Picardove preslikave. Označimo

$$h_n(t, x) = A^n(h_0(t, x))$$

za $h_0 = 0$. Dokazati moramo, da za vsak $n \in \mathbb{N}$ velja $\|h_n(t, x)\| \leq C |t - t_0|$. To naredimo z indukcijo na n . Pri $n = 0$ to očitno velja, indukcijski korak pa pokažemo z računom

$$\|h_{n+1}(t, x)\| = \left\| \int_{t_0}^t F(\tau, x + h_n(\tau, x)) d\tau \right\| \leq \int_{t_0}^t \|F(\tau, x + h_n(\tau, x))\| d\tau.$$

Po indukcijski predpostavki vemo, da točka $h_n(t, x)$ leži v K_0 za vsak (t, x) , zato $x + h_n(t, x)$ leži v $K_x \subseteq K \subseteq \mathcal{C}_{a,b}$. Sledi

$$\|F(\tau, x + h_n(\tau, x))\| \leq C,$$

zato

$$\|h_{n+1}(t, x)\| \leq \left| \int_{t_0}^t C d\tau \right| = C |t - t_0|.$$

Pokazali smo, da je $h_n \in M$ za vsak n . Ker je M poln, je tudi limita v M , če obstaja.

Pokazati moramo še, da je A na M skrčitev. Naj bosta $h_1, h_2 \in M$ poljubni. Oglejmo si

$$\begin{aligned} \|Ah_1(t, x) - Ah_2(t, x)\| &= \left\| \int_{t_0}^t F(\tau, x + h_1(\tau, x)) - F(\tau, x + h_2(\tau, x)) d\tau \right\| \\ &\leq \int_{t_0}^t \|F(\tau, x + h_1(\tau, x)) - F(\tau, x + h_2(\tau, x))\| d\tau. \end{aligned}$$

Ker je F Lipschitzova glede na x , velja

$$\begin{aligned} \|Ah_1(t, x) - Ah_2(t, x)\| &\leq \int_{t_0}^t L \|h_1(\tau, x) - h_2(\tau, x)\| d\tau \\ &\leq \int_{t_0}^t L \|h_1 - h_2\| d\tau \\ &= L \|h_1 - h_2\| |t - t_0| \\ &\leq L \|h_1 - h_2\| a' \end{aligned}$$

Po potrebi še zmanjšamo a' , da bo $La' < 1$.

Ker je A skrčitev, limita zaporedja h_n obstaja in je fiksna točka preslikave A . Torej je rešitev začetnega problema

$$\partial_t h = F(t, x + h(t, x)) \quad h(t_0, x) = 0,$$

iz katere dobimo preslikavo $g(t, x) = x + h(t, x)$. Naša limita je po konstrukciji zvezna (ker leži v M), torej je tudi g zvezna glede na oba argumenta.

To je konec dokaza eksistenčnega izreka.

Vprašanje 21. Dokaži eksistenčni izrek.

Trditev. Naj bo $f : U \subseteq \mathbb{R}^n \rightarrow \mathbb{R}^m$ zvezno odvedljiva na konveksnem kompaktu $M \subseteq U$. Potem je na M Lipschitzova.

Dokaz. Naj bosta $x, y \in M$ poljubni točki. Definiramo $z(t) = x + t(y - x)$ kot daljico med x in y . Velja

$$f(y) - f(x) = \int_0^1 \partial_\tau f(z(\tau)) d\tau = \int_0^1 Df(z(\tau))(y - x) d\tau.$$

Ker je f zvezno odvedljiva na kompaktu, norma Jacobijeve matrike doseže maksimum, torej velja

$$\|f(x) - f(y)\| = \left\| \int_0^1 Df(z(\tau))(y - x) d\tau \right\| \leq \int_0^1 \|Df\| \|y - x\| d\tau = \|Df\| \|y - x\|.$$

□

Vprašanje 22. Dokaži, da je zvezno odvedljiva funkcija na konveksnem kompaktu Lipschitzova.

Imejmo NDE $\dot{x} = F(t, x)$. Tok te enačbe je preslikava

$$\phi : (a, b) \times (\alpha, \beta) \times U \rightarrow \mathbb{R}^n,$$

definirana s predpisom

$$\phi(t, t_0, x) = \gamma(t),$$

kjer je γ rešitev začetnega problema $\dot{\gamma}(t) = F(t, \gamma(t)), \gamma(t_0) = x$.

Trditev. Za tok enačbe velja

- Za vsak t_0 , za katerega rešitve začetnih problemov $\gamma(t_0) = x$ obstajajo, in za vsak t dovolj blizu t_0 , je $x \mapsto \phi(t, t_0, x)$ difeomorfizem U na svojo sliko.
- Za t_1, t_2 dovolj blizu t_0 velja

$$\phi(t_2, t_1, \phi(t_1, t_0, x)) = \phi(t_2, t_0, x).$$

Dokaz. Prva točka: Uporabimo izrek o inverzni preslikavi na $\phi(t, t_0, \cdot)$. Ker je $\phi(t_0, t_0, x) = x$, obstaja okolica t_0 , v kateri je $\det D_x(t, t_0, x) \neq 0$, in dobimo difeomorfizem.

Druga točka sledi iz edinosti, ki nam jo da eksistenčni izrek. □

Vprašanje 23. Kaj je tok enačbe? Kakšne lastnosti ima?

1.7 Sistemi linearnih NDE

Naj bodo podane funkcije $a_{ij} : [a, b] \rightarrow \mathbb{R}$ in $b_k : [a, b] \rightarrow \mathbb{R}$, ki so na $[a, b]$ omejene. Sistem NDE prvega reda s koeficienti $a_{ij}(t)$ in desno stranjo $b_k(t)$ je sistem

$$\begin{aligned} \dot{x}_1 &= a_{11}x_1 + \dots + a_{1n}x_n \\ &\vdots \\ \dot{x}_n &= a_{n1}x_1 + \dots + a_{nn}x_n \end{aligned}$$

Naj bo matrika $A(t)$ podana s koeficienti a_{ij} in b vektor podan s komponentami b_k . Za $x = [x_1 \dots x_n]^T$ sistem zapišemo kot $\dot{x} = Ax + b$. Če je $b = 0$ pravimo, da je sistem HOMOGEN.

Izrek. Če je $A : [a, b] \rightarrow \mathbb{R}^{n \times n}$ zvezna in omejena, je množica rešitev homogenega sistema $\dot{x} = Ax$ n -dimenzionalen vektorski prostor v prostoru $C^1([a, b])$.

Dokaz. Naj bo R prostor rešitev. Najprej moramo pokazati, da je R vektorski podprostor. Za vsak par rešitev x_1, x_2 in $\alpha, \beta \in \mathbb{R}$ velja

$$A(\alpha x_1 + \beta x_2) = \alpha Ax_1 + \beta Ax_2 = \partial_t(\alpha x_1 + \beta x_2).$$

Naj bo $t_0 \in [a, b]$. Po eksistenčnem izreku obstaja rešitev vsakega začetnega problema $\dot{x} = Ax, x(t_0) = c \in \mathbb{R}^n$. Naj bo e_1, \dots, e_n kanonična baza v \mathbb{R}^n . Obstajajo torej rešitve x_i začetnih problemov $\dot{x} = Ax, x(t_0) = e_i$. Dokazati moramo še, da so x_i tudi globalne rešitve; ta del pustimo za kasneje, preostanek dokaza je lokalni.

Trdimo, da je $\{x_i\}_i$ baza R . Linearna neodvisnost je trivialna. Naj bo x rešitev sistema in $c = x(t_0)$. Vektor c razvijemo po bazi e_i v $c = \alpha_1 e_1 + \dots + \alpha_n e_n$ in definiramo

$$\tilde{x} = \alpha_1 x_1 + \dots + \alpha_n x_n.$$

Ker sta tako x kot \tilde{x} rešitvi začetnega problema $\dot{x} = Ax, x(t_0) = c$, sta po eksistenčnem izreku enaki. \square

Vprašanje 24. Kaj je množica rešitev homogenega sistema linearnih NDE? Dokaži.

Definicija. FUNDAMENTALNA MATRIKA sistema $\dot{x} = Ax$ je matrika

$$\phi(t, t_0) = \begin{bmatrix} x_{11}(t) & \dots & x_{1n}(t) \\ \vdots & \ddots & \vdots \\ x_{n1}(t) & \dots & x_{nn}(t) \end{bmatrix},$$

v kateri je i -ti stolpec enak i -ti rešitvi iz dokaza.

Za matriko $\phi(t, t_0)$ velja $\phi(t_0, t_0) = I$. Naj bo x rešitev začetnega problema $\dot{x} = Ax, x(t_0) = c$. Potem velja $x = \phi(t, t_0)c$.

Trditev. Za vsak t veja $\det \phi(t, t_0) \neq 0$.

Dokaz. Recimo, da obstaja t_1 , da je $\det \phi(t_1, t_0) = 0$. Potem je matrika $\phi(t_1, t_0)$ singularna, zato ima netrivialno jedro, torej obstaja vsaj en neničeln vektor $d \in \mathbb{R}^n$, da $\phi(t_1, t_0)d = 0$. Torej

$$\phi(t_1, t_0)d = \sum_{i=1}^n x_i(t_1)d_i = 0.$$

Definirajmo

$$z(t) = \sum_{i=1}^n d_i x_i(t).$$

Ta funkcija je rešitev sistema, zanjo velja $z(t_1) = 0$. Tudi funkcija $w(t) = 0$ je rešitev začetnega problema $\dot{x} = Ax, x(t_1) = 0$, torej po eksistenčnem izreku $z = 0$. Ker so v točki t_0 vektorji x_i linearno neodvisni, velja $d_i = 0$. \square

Vprašanje 25. Kaj je fundamentalna matrika homogenega sistema linearnih NDE? Dokaži, da je nesingularna.

Oglejmo si preslikavo $\mathcal{F} : (t_0 - \varepsilon, t_0 + \varepsilon) \times \mathbb{R}^n \rightarrow \mathbb{R}^n$, definirano z

$$\mathcal{F}(t, y) = \phi(t, t_0)y.$$

To je tok enačbe $\dot{x} = Ax$. Za vsak dovolj majhen t je preslikava $y \mapsto \mathcal{F}(t, y)$ difeomorfizem, saj je obrnljiva linearna preslikava $\mathbb{R}^n \rightarrow \mathbb{R}^n$. Vzemimo $t_0 = 0$ in označimo $\phi(t, 0) = \phi(t)$.

Trditev. Velja $\phi(t_1 + t_2) = \phi(t_1)\phi(t_2)$.

Dokaz. Po eni strani imamo za vsak $x \in \mathbb{R}^n$

$$\begin{aligned}\phi(t_1)x &= \mathcal{F}(t_1, x), \\ \phi(t_2)\phi(t_1)x &= \mathcal{F}(t_2, \phi(t_1)x) = \mathcal{F}(t_2 + t_1, x),\end{aligned}$$

ker je \mathcal{F} tok, po drugi strani pa

$$\phi(t_1 + t_2)x = \mathcal{F}(t_1 + t_2, x).$$

□

Vprašanje 26. Pokaži, da velja $\phi(t_1 + t_2) = \phi(t_1)\phi(t_2)$.

Trditev. Splošna rešitev sistema $\dot{x} = Ax + b$ je afin podprostor v $\mathcal{C}^1([a, b])$, modeliran nad prostorom R rešitev homogenega sistema.

To pomeni, da obstajajo vektorji $x_p \in \mathcal{C}^1([a, b])$, da je množica rešitev $\dot{x} = Ax + b$ enaka

$$W = \{x_h + x_p \mid x_h \in R\}.$$

Če imamo R in želimo poiskati W , potrebujemo eno partikularno rešitev nehomogenega sistema. To dobimo z variacijo konstante. Za vsak konstanten vektor $c \in \mathbb{R}^n$ je $\phi(t)c$ rešitev homogenega sistema. Poskusimo poiskati kakšno rešitev $\dot{x} = Ax + b$ z nastavkom $x_p = \phi(t)c(t)$, kjer je $c(t) : \mathbb{R} \rightarrow \mathbb{R}^n$ neznana funkcija. Začnemo z

$$\dot{x}_p = \dot{\phi}c + \phi\dot{c} = Ax_p + b,$$

iz česar dobimo $\phi\dot{c} = b$, oziroma $\dot{c} = \phi(-t)b(t)$. Sledi

$$c(t) = \int_0^t \phi(-\tau)b(\tau)d\tau$$

in

$$x_p(t) = \phi(t)c(t) = \int_0^t \phi(t - \tau)b(\tau)d\tau.$$

Dokazali smo

Trditev. Splošna rešitev nehomogenega problema $\dot{x} = Ax + b$ je

$$x(t, c) = \phi(t)c + \int_0^t \phi(t - \tau)b(\tau)d\tau,$$

kjer je c začetni pogoj pri $t_0 = 0$.

Vprašanje 27. Kako poiščeš množico rešitev $\dot{x} = Ax + b$?

Kako pa izračunamo ϕ ? V zaključeni obliki za splošen A izračunati ne moremo, lahko pa dobimo izrazitev z Dysonovo vrsto. Oglejmo si začetni problem

$$\dot{\phi} = A\phi, \phi(0) = I.$$

Ta je ekvivalenten integralni enačbi

$$\phi(t) = I + \int_0^t A(\tau)\phi(\tau)d\tau.$$

To lahko razvijemo naprej v

$$\phi(t) = I + \int_0^t A(\tau_1) \left(I + \int_0^{\tau_1} A(\tau_2)\phi(\tau_2) \right) d\tau_1,$$

in nadaljujemo. Na koncu dobimo

$$\phi(t) = I + \sum_{n=1}^{\infty} \int_0^t A(\tau_1) \int_0^{\tau_1} A(\tau_2) \dots \int_0^{\tau_{n-1}} A(\tau_n) d\tau_n \dots d\tau_1$$

Trditev. Naj bo matrična funkcija $A(t) : [0, T] \rightarrow \mathbb{R}^{n \times n}$ omejena po normi $\|A(t)\| \leq M$, Potem Dysonova vrsta konvergira.

Dokaz. Ocenimo lahko

$$\|\phi\| \leq 1 + \sum_{n=1}^{\infty} \int_{\Delta_n(t)} \|A(\tau_1) \dots A(\tau_n)\| d\tau \leq 1 + \sum_{n=1}^{\infty} \int_{\Delta_n(t)} M^n d\tau,$$

kjer je $\Delta_n(t)$ urejeni n -simpleks. Nadalje velja

$$\|\phi\| \leq 1 + \sum_{n=1}^{\infty} V(\Delta_n(t)) M^n = 1 + \sum_{n=1}^{\infty} \frac{t^n}{n!} M^n = e^{Mt} \leq e^{MT}.$$

□

Vprašanje 28. Pod katerim pogojem Dysonova vrsta konvergira? Dokaži.

Recimo, da je A konstanta matrika. V tem primeru se Dysonova matrika glasi

$$\phi(t) = e^{At}.$$

Vprašanje 29. Kakšna je Dysonova vrsta, če je matrika koeficientov konstanta?

Trditev. Naj bo $A : [0, T] \rightarrow \mathbb{R}^{n \times n}$ matrična funkcija, za katero velja $A(t_1)A(t_2) = A(t_2)A(t_1)$ za vsaka $t_1, t_2 \in [0, T]$. Potem velja

$$\phi(t) = \exp \left(\int_0^t A(\tau) d\tau \right).$$

Dokaz. Označimo $\square_n(t) = [0, t]^n$. Oglejmo si

$$\int_{\square_n(t)} A(\tau_1) \dots A(\tau_n) d\tau.$$

Za skoraj vsak $\tau \in \square_n(t)$ obstaja natanko ena permutacija $\sigma \in S_n$, da velja $\sigma \cdot \tau \in \Delta_n(t)$. Označimo

$$\Delta^\sigma(t) = \{\tau \in \square_n(t) \mid \sigma \cdot \tau \in \Delta_n(t)\}.$$

Razen na množici z mero 0 velja

$$\square_n(t) = \bigcup_{\sigma \in S_n} \Delta^\sigma(t),$$

zato za vsako funkcijo $\mathcal{A} : \square_n(t) \rightarrow \mathbb{R}^{n \times n}$ velja

$$\int_{\square_n(t)} \mathcal{A}(\tau) d\tau = \sum_{\sigma \in S_n} \int_{\Delta^\sigma(t)} \mathcal{A}(\tau) d\tau = \sum_{\sigma \in S_n} \int_{\Delta_n(t)} \mathcal{A}(\tau_{\sigma(1)}, \dots, \tau_{\sigma(n)}) \underbrace{|\det \sigma|}_{=1} d\tau.$$

Če matrike komutirajo, torej velja

$$\int_{\square_n(t)} A(\tau_1) \dots A(\tau_n) d\tau = n! \int_{\Delta_n(t)} A(\tau_1) \dots A(\tau_n) d\tau.$$

Torej je

$$\int_{\Delta_n(t)} A(\tau_1) \dots A(\tau_n) d\tau = \frac{1}{n!} \int_0^t A(\tau_1) d\tau_1 \dots \int_0^t A(\tau_n) d\tau_n = \frac{1}{n!} \left(\int_0^t A(\tau) d\tau \right)^n.$$

□

Vprašanje 30. Kakšna je fundamentalna matrika, če $A(t_1)$ komutira z $A(t_2)$ za vsaka t_1, t_2 ? Dokaži.

Trditev (Liouvilleova formula). Za fundamentalno matriko $\phi(t)$ sistema $\dot{x} = Ax + b$ velja

$$\det \phi(t) = \exp \left(\int_0^t \text{sl}(A(\tau)) d\tau \right)$$

Dokaz. Naj bo

$$\phi(t) = \begin{bmatrix} x_{11} & \cdots & x_{1n} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{nn} \end{bmatrix}.$$

Če definicijo determinante odvajamo, dobimo

$$\partial_t \det \phi(t) = \sum_{\pi \in S_n} (-1)^{s(\pi)} \sum_{i=1}^n x_{1,\pi(1)} \cdots \dot{x}_{i,\pi(i)} \cdots x_{n,\pi(n)}.$$

Velja $\dot{\phi} = A\phi$, torej

$$\begin{bmatrix} \dot{x}_{11} & \cdots & \dot{x}_{1n} \\ \vdots & \ddots & \vdots \\ \dot{x}_{n1} & \cdots & \dot{x}_{nn} \end{bmatrix} = \begin{bmatrix} \sum_i a_{1i} x_{i1} & \cdots & \sum_i a_{1i} x_{in} \\ \vdots & \ddots & \vdots \\ \sum_i a_{ni} x_{i1} & \cdots & \sum_i a_{ni} x_{in} \end{bmatrix},$$

iz česar dobimo $\dot{x}_{ij} = \sum_k a_{ik} x_{kj}$. To vstavimo v prejšnji zapis

$$\partial_t \det \phi(t) = \sum_{\pi \in S_n} (-1)^{s(\pi)} \sum_{i=1}^n x_{1,\pi(1)} \cdots \left(\sum_{j=1}^n a_{ij} x_{j,\pi(j)} \right) \cdots x_{n,\pi(n)},$$

ki ga prvo seštejemo po π . Pri vsakem i dobimo determinanto

$$\begin{vmatrix} x_{11} & \cdots & x_{1n} \\ \vdots & \ddots & \vdots \\ \sum_j a_{ij} x_{j1} & \cdots & \sum_j a_{ij} x_{jn} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{nn} \end{vmatrix} = \begin{vmatrix} x_{11} & \cdots & x_{1n} \\ \vdots & \ddots & \vdots \\ a_{ii} x_{i1} & \cdots & a_{ii} x_{in} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{nn} \end{vmatrix} = a_{ii} \det \phi$$

Torej

$$\partial_t \det \phi(t) = a_{11} \det \phi + a_{22} \det \phi + \cdots + a_{nn} \det \phi = \text{sl } A \det \phi.$$

To je diferencialna enačba, katere rešitev je

$$\det \phi = \exp \left(\int_0^t \text{sl } A d\tau \right).$$

□

Vprašanje 31. Povej in dokaži Liouviloevo formulo.

1.8 Linearne NDE višjega reda

Obravnavamo enačbe oblike

$$a_n(t)x^{(n)} + \dots + a_1(t)\dot{x} + a_0(t)x = b(t).$$

Definicija. Linearni diferencialni operator s koeficienti $a_i(t)$ je preslikava

$$L : \mathcal{C}^1([a, b]) \rightarrow \mathcal{C}([a, b]),$$

podana s predpisom

$$Lx(t) = a_n(t)x^{(n)} + \dots + a_0(t)x.$$

Splošna rešitev homogene enačbe $Lx = 0$ je $\ker L$. Vemo, da je splošna rešitev n -dimenzionalni vektorski prostor. Enačbo s substitucijo

$$\begin{aligned} x_1 &= x \\ x_2 &= \dot{x} \\ &\vdots \\ x_n &= x^{(n-1)} \end{aligned}$$

prepišemo v sistem

$$\begin{aligned} \dot{x}_0 &= x_1 \\ \dot{x}_1 &= x_2 \\ &\vdots \\ \dot{x}_n &= \frac{1}{a_n}(b - a_0x_1 - a_2x_1 - \dots - a_{n-1}x_{n-2}) \end{aligned}$$

Označimo

$$p_i(t) = \frac{a_i(t)}{a_n(t)},$$

s čimer izrazimo matriko koeficientov zgornjega sistema

$$A(t) = \begin{bmatrix} 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ -p_0 & -p_1 & -p_2 & \dots & -p_{n-1} \end{bmatrix}.$$

Izrek. Naj bo $L : \mathcal{C}^n([a, b]) \rightarrow \mathcal{C}^0([a, b])$ regularen diferencialni operator. Za vsak začetni pogoj $x(t_0) = c_0, \dot{x}(t_0) = c_1, \dots, x^{(n-1)}(t_0) = c_{n-1}$ ima enačba $Lx = b$ natanko eno rešitev na vsem intervalu $[a, b]$.

Opomba. Diferencialni operator: $Lx = a_n x^{(n)} + \dots + a_1 \dot{x} + a_0 x = b$ je regularen, če je $a_n(t) \neq 0$ za vsak t in če so $a_i(t)$ omejene.

Množica rešitev homogene linearne enačbe $Lx = 0$ je n -dimenzionalen vektorski prostor v $\mathcal{C}^n([a, b])$. Vsaka rešitev $x(t)$ namreč na enoličen način določa rešitev sistema $\dot{\vec{x}} = A\vec{x}$ za

$$\dot{\vec{x}} = \begin{bmatrix} x(t) \\ \dot{x}(t) \\ \vdots \\ x^{(n-1)}(t) \end{bmatrix}.$$

Tudi obratno je res: vsak vektor \vec{x} na enoličen način določa svojo prvo komponento.

Oglejmo si fundamentalno matriko

$$\phi(t) = \begin{bmatrix} x_1 & x_2 & \cdots & x_n \\ \dot{x}_1 & \dot{x}_2 & \cdots & \dot{x}_n \\ \vdots & \vdots & \ddots & \vdots \\ x_1^{(n-1)} & x_2^{(n-1)} & \cdots & x_n^{(n-1)} \end{bmatrix}.$$

Denimo, da so funkcije $x_1(t), x_2(t), \dots, x_n(t)$ baza rešitev enačbe $Lx = 0$. Za bazo velikokrat vzamemo take vektorje $\vec{x}_1, \dots, \vec{x}_n$, da je $\phi(t=0) = I$. Če je ta baza dobljena iz baze rešitev enačbe $Lx = 0$, potem za te rešitve velja $x_i(0) = 0, \dots, x_i^{(i-1)}(t) = 0, x_i^{(i)}(0) = 1, x_i^{(i+1)}(0) = 0, \dots, x_i^{(n-1)}(0) = 0$. Determinanta

$$W(t) = \det \phi(t)$$

se imenuje DETERMINANTA WRONSKEGA. V tem primeru se Liouvilleova formula glasi

$$W(t) = W(t_0) \exp \left(- \int_{t_0}^t p_{n-1}(\tau) d\tau \right).$$

Rešitev nehomogene enačbe dobimo s pomočjo variacije konstante;

$$\vec{x} = \int_{t_0}^t \phi(t) \phi^{-1}(\tau) \vec{b}(\tau) d\tau,$$

oziroma, ker ima \vec{b} v tem primeru le eno neničelno komponento $b(t)$, bo prva komponenta \vec{x} enaka

$$x(t) = \int_{t_0}^t \sum_{i=1}^n \phi_{1i}(t) \phi_{in}^{-1}(\tau) b(\tau) d\tau = \sum_{i=1}^n x_i(t) \int_{t_0}^t \phi_{in}^{-1}(\tau) b(\tau) d\tau,$$

kjer je $(x_i(t))_i$ baza rešitev homogene enačbe $Lx = 0$.

Vprašanje 32. Kako izračunaš rešitev linearne NDE višjega reda?

1.8.1 Enačbe s konstantnimi koeficienti

Naj bo sedaj linearen diferencialni operator L podan z

$$Lx = x^{(n)} + a_1 x^{(n-1)} + \dots + a_n x.$$

Oglejmo si enačbo $Lx = 0$. Če vstavimo nastavek $x(t) = e^{\lambda t}$:

$$\lambda^n e^{\lambda t} + a_1 \lambda^{n-1} e^{\lambda t} + \dots + a_{n-1} \lambda e^{\lambda t} + a_n e^{\lambda t} = 0$$

oziroma (ker $e^{\lambda t} \neq 0$ tudi za $\lambda \in \mathbb{C}$)

$$\lambda^n + a_1 \lambda^{n-1} + \dots + a_{n-1} \lambda + a_n = 0.$$

Ta polinom imenujemo KARAKTERISTIČNI POLINOM ENAČBE $Lx = 0$ in označimo s $P(\lambda)$.

Trditev. Če so $\lambda_1, \dots, \lambda_n$ različne ničle karakterističnega polinoma $P(\lambda)$, potem so funkcije $x_i(t) = e^{\lambda_i t}$ baza rešitev homogene enačbe $Lx = 0$.

Dokaz. Vemo, da je rešitev sistema vektorski prostor, dokazati moramo samo, da so te rešitve linearno neodvisne. Priredimo našim rešitvam pripadajoče rešitve sistema, ki je prirejen $Lx = 0$,

$$x_i(t) \mapsto \vec{x}_i(t) = \begin{bmatrix} x_i(t) \\ \dot{x}_i(t) \\ \vdots \\ x_i^{n-1}(t) \end{bmatrix}.$$

Te stolpce zložimo v matriko in dobimo kandidatko za fundamentalno matriko $\phi(t)$

$$\phi(t) = \begin{bmatrix} e^{\lambda_1 t} & e^{\lambda_2 t} & \dots & e^{\lambda_n t} \\ \lambda_1 e^{\lambda_1 t} & \lambda_2 e^{\lambda_2 t} & \dots & \lambda_n e^{\lambda_n t} \\ \lambda_1^2 e^{\lambda_1 t} & \lambda_2^2 e^{\lambda_2 t} & \dots & \lambda_n^2 e^{\lambda_n t} \\ \vdots & \vdots & \ddots & \vdots \\ \lambda_1^{n-1} e^{\lambda_1 t} & \lambda_2^{n-1} e^{\lambda_2 t} & \dots & \lambda_n^{n-1} e^{\lambda_n t} \end{bmatrix}$$

Matrika $\phi(0)$ je vandermondova, torej

$$W(t) = \det \phi(t) = W(0) \exp \left(\int_0^t \text{sl}(A) d\tau \right) = \prod_{i>j} (\lambda_i - \lambda_j) \cdot e^{-a_1 t}.$$

Ker so vsi λ_i različni, velja $W(0) \neq 0$, torej je $W(t) \neq 0$ za vsak t , in so funkcije x_i res linearno neodvisne in so baza prostora rešitev enačbe $Lx = 0$. \square

Vprašanje 33. Kaj je karakteristični polinom homogene linearne NDE višjega reda? Kako z njim poiščemo rešitve enačbe, če so vse ničle različne? Dokaži.

Z razvojem po prvem stolpcu lahko izračunamo

$$\det(A - \lambda I) = \begin{vmatrix} -\lambda & 1 & 0 & \cdots & 0 \\ 0 & -\lambda & 1 & \cdots & 0 \\ 0 & 0 & -\lambda & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ -a_n & -a_{n-1} & -a_{n-2} & \cdots & -a_1 - \lambda \end{vmatrix} = \lambda^n + a_1\lambda^{n-1} + \dots + a_{n-1}\lambda + a_n.$$

Lastne vrednosti matrike A so torej res ničle karakterističnega polinoma $P(\lambda)$ enačbe $Lx = 0$.

Kaj pa če ima A večkratne lastne vrednosti?

Trditev. Naj bo λ k -kratna ničla polinoma $P(\lambda)$. Potem so funkcije $x_0(t) = e^{\lambda t}$, $x_1(t) = te^{\lambda t}$, \dots , $x_{k-1}(t) = t^{k-1}e^{\lambda t}$ linearne neodvisne rešitve $Lx = 0$.

Dokaz. Opazimo, da velja $x_i(t) = \frac{\partial^i}{\partial \lambda^i} e^{\lambda t}$ za $i = 0, \dots, k-1$. Kot že vemo, velja $Lx_0(t) = P(\lambda)e^{\lambda t}$. Odvedemo to enačbo i -krat po λ . Na levi dobimo

$$\frac{\partial^i}{\partial \lambda^i} Lx_0(t) = L\left(\frac{\partial^i}{\partial \lambda^i} e^{\lambda t}\right) = L(x_i(t)).$$

To je res, ker so vsi koeficienti a_i konstantni na t in λ . Imamo torej

$$\begin{aligned} Lx_i(t) &= \frac{\partial^i}{\partial \lambda^i} Lx_0(t) = \frac{\partial^i}{\partial \lambda^i} (P(\lambda)e^{\lambda t}) \\ &= \sum_{l=0}^i \binom{i}{l} P^{(l)}(\lambda) \frac{\partial^{i-l}}{\partial \lambda^{i-l}} e^{\lambda t} \\ &= \sum_{l=0}^i \binom{i}{l} P^{(l)}(\lambda) t^{i-l} e^{\lambda t} \\ &= \sum_{l=0}^i \binom{i}{l} P^{(l)}(\lambda) x_{i-l}(t). \end{aligned}$$

Naj bo sedaj λ ničla k -te stopnje in $i \leq k$. Potem velja $P^{(l)}(\lambda) = 0$, torej $Lx_i(t) = 0$, funkcija $x_i(t) = t^i e^{\lambda t}$ je torej res rešitev enačbe $Lx = 0$ za $i = 0, \dots, k-1$. \square

Vprašanje 34. Kako s karakterističnim polinomom poiščemo rešitve linearne NDE višjega reda s konstantnimi koeficienti, če niso vse ničle različne? Dokaži.

Naj bo sedaj λ kompleksna ničla $\lambda = a + ib$. Če so koeficienti operatorja L realni, potem je tudi $\bar{\lambda}$ ničla P . Če je $\lambda = a + ib$ ničla k -tega reda, je tudi $\bar{\lambda}$ ničla k -tega reda. Ti dve lastni vrednosti dasta $2k$ baznih rešitev. Če jih želimo na najpreprostejši način izraziti

z realnimi funkcijami, dobimo bazo

$$\begin{aligned}x_0(t) &= e^{ta} \cos(bt), x_1(t) = e^{ta} \sin(bt), \\x_2(t) &= te^{ta} \cos(bt), x_3(t) = te^{ta} \sin(bt), \\&\vdots \\x_{2k-2}(t) &= t^{k-1} e^{ta} \cos(bt), x_{2k-1}(t) = t^{k-1} e^{ta} \sin(bt).\end{aligned}$$

1.8.2 Linearizacija

Imejmo nelinearen sistem NDE $\dot{\vec{x}} = F(t, x)$. Naj bo $\vec{x}_0(t)$ neka rešitev tega sistema. Definiramo nelinearen operator $\mathcal{F}(\vec{x}(t)) = \dot{\vec{x}}(t) - F(t, \vec{x})$. Funkcija \vec{x}_0 je rešitev sistema natanko tedaj, ko je $\mathcal{F}(\vec{x}_0) = 0$. Splošna rešitev \mathcal{S} , ki je n -parametrični nelinearen prostor, je nivojska ploskev $\mathcal{S} = \mathcal{F}^{-1}(0)$. Naj bo sedaj $s \mapsto \vec{x}(t, s) \in \mathcal{S}$ pot v prostoru rešitev, za katero velja $\vec{x}(t, 0) = \vec{x}_0(t)$. Oglejmo si odvod

$$\left. \frac{d}{ds} \right|_{s=0} \mathcal{F}(\vec{x}(t, s)) = D_{\vec{x}_0} \mathcal{F} \left(\left. \frac{d}{ds} \right|_{s=0} \vec{x}(t, s) \right) =: D_{\vec{x}_0} \mathcal{F}(\vec{u}(t)).$$

Če je $\vec{x}(t, s)$ rešitev za vsak s , potem velja $\mathcal{F}(\vec{x}(t, s)) = 0$, zato

$$\left. \frac{d}{ds} \right|_{s=0} \mathcal{F}(\vec{x}(t, s)) = D_{\vec{x}_0} \mathcal{F}(\vec{u}(t)) = 0.$$

Operator $D_{\vec{x}_0} \mathcal{F}$ je linearen. Če v predpis odvoda vstavimo definicijo \mathcal{F} , dobimo

$$\left. \frac{d}{ds} \right|_{s=0} \mathcal{F}(\vec{x}(t, s)) = \left. \frac{d}{ds} \right|_{s=0} \left(\dot{\vec{x}}(t, s) - F(t, \vec{x}(t, s)) \right) = \dot{\vec{u}} - \left. \frac{d}{ds} \right|_{s=0} F(t, \vec{x}(t, s)).$$

Če desni člen posredno odvajamo, dobimo sistem $\dot{\vec{u}} - A(t)\vec{u} = 0$ za

$$A(t) = \begin{bmatrix} \partial_{x_1} F_1(t, \vec{x}_0(t)) & \cdots & \partial_{x_n} F_1(t, \vec{x}_0(t)) \\ \vdots & \ddots & \vdots \\ \partial_{x_1} F_n(t, \vec{x}_0(t)) & \cdots & \partial_{x_n} F_n(t, \vec{x}_0(t)) \end{bmatrix}.$$

Dobili smo linearizacijo začetnega sistema okoli rešitve \vec{x}_0 . Splošna rešitev te linearizacije je jedro operatorja $D_{\vec{x}_0} \mathcal{F}$, oziroma tangentni prostor na \mathcal{S} v točki \vec{x}_0 .

Vprašanje 35. Izpelj linearizacijo nelinearnega sistema NDE. Kakšne so rešitve linearizacije?

1.9 Variacijski račun

Definicija. Naj bosta U in V Banachova prostora in $P : U \rightarrow V$ operator. Naj bo $u \in U$ točka v prostoru. GATEAUXOV ODVOD P v u in v smeri $v \in U$ je podan s predpisom

$$D_u^G P(v) = \left. \frac{d}{dt} \right|_{t=0} P(u + tv) = \lim_{t \rightarrow 0} \frac{P(u + tv) - P(u)}{t}.$$

Ta odvod se imenuje tudi SMERNI ali ŠIBKI odvod.

Definicija. Naj bosta U in V Banachova prostora in $P : U \rightarrow V$ operator. FRECHETOV ali KREPKI odvod P v točki u je omejen linearen operator $\mathcal{A} : U \rightarrow V$, za katerega velja $P(u + tv) = P(u) + \mathcal{A}(v) + o(u, v)$ in

$$\lim_{\|v\| \rightarrow 0} \frac{\|o(u, v)\|}{\|v\|} = 0.$$

Pišemo $\mathcal{A} = D_u^F P$.

Vprašanje 36. Definiraj šibki in krepki odvod. Kako se še imenujeta?

Trditev. Če je P v $u \in U$ Frechetovo odvedljiv, je tam tudi Gateauxovo odvedljiv. Tedaj sta odvoda enaka.

Dokaz. Za vsak $v \in U$ velja

$$P(u + tv) = Pu + D_u^F P(tv) + o(u, tv) = Pu + tD_u^F P(v) + o(u, tv),$$

iz česar pride

$$\frac{P(u + tv) - Pu}{t} = D_u^F P(v) + \frac{o(u, tv)}{t}.$$

Brez škode za splošnost predpostavimo $\|v\| = 1$. Ker Frechetov odvod obstaja, velja

$$\lim_{t \rightarrow 0} \frac{P(u + tv) - Pu}{t} = D_u^F P(v) + \lim_{\|tv\| \rightarrow 0} \frac{o(u, tv)}{\|tv\|} = D_u^F P(v).$$

□

Trditev. Če obstaja $D_u^G P$ in je limita

$$\lim_{t \rightarrow 0} \frac{P(u + tv) - Pu}{t} = D_u^G P(v)$$

enakomerna glede na v na enotski sferi $S \subseteq U$, potem obstaja tudi $D_u^F P$.

Dokaz. Enakomernost pomeni, da za vsak $\varepsilon > 0$ obstaja $\delta > 0$, da $|t| < \delta$ implicira

$$\left\| \frac{P(u + tv) - Pu}{t} - D_u^G P(v) \right\| < \varepsilon$$

ne glede na $v \in S \subseteq U$, oziroma

$$\|P(u + tv) - Pu - D_u^G P(tv)\| < \varepsilon t.$$

Označimo $h = tv$ in dobimo

$$\|P(u + h) - Pu - D_u^G P(h)\| < \varepsilon \|h\|.$$

Po definiciji limite potem velja

$$\lim_{\|h\| \rightarrow 0} \frac{\|P(u+h) - Pu - D_u^G P(h)\|}{\|h\|} = 0.$$

□

Vprašanje 37. Povej zadostni pogoj za obstoj Frechetovega odvoda.

Trditev. Naj bo $u \in U$ lokalni minimum funkcionala $\mathcal{L} : U \rightarrow \mathbb{R}$, in naj bo \mathcal{L} krepko odvedljiv. Potem velja $D_u \mathcal{L}(v) = 0$ za vsak v .

Dokaz. Ker je v u dosežen minimum, velja za vsak v in vsak dovolj majhen t $\mathcal{L}(u+tv) > \mathcal{L}(u)$. Ker je $\mathcal{L}(u+tv) = \mathcal{L}(u) + D_u \mathcal{L}(tv) + o(tv)$, velja $D_u \mathcal{L}(tv) + o(tv) > 0$. Če je $t > 0$, je $tD_u \mathcal{L}(v) + o(tv) > 0$ in

$$D_u \mathcal{L}(v) + \frac{o(tv)}{t} > 0.$$

Ker drug člen limitira k 0, je za dovolj majhen t predznak izraza odvisen le od predznaka $D_u \mathcal{L}(v)$, ki mora torej biti pozitiven. Če enako naredimo za $t < 0$, dobimo, da mora biti predznak $D_u \mathcal{L}(v)$ negativen; torej je enak nič. □

Vprašanje 38. Dokaži, da je odvod funkcionala v lokalnem minimumu enak 0.

Če je

$$\mathcal{L} = \int_a^b L(x, u, u') dx,$$

je odvod funkcionala enak

$$D_u \mathcal{L}(v) = \int_a^b \left(\frac{\partial L}{\partial u}(x) v(x) + \frac{\partial L}{\partial u'}(x) v'(x) \right) dx.$$

Definicija. Naj bo U Banachov prostor funkcij $y : [a, b] \rightarrow \mathbb{R}$ z metriko, porojeno iz maksimum norme. Naj bosta $A, B \in \mathbb{R}$ konstanti. Prostor

$$V = \{y \in U \mid y(a) = A, y(b) = B\}.$$

imenujemo PROSTOR DOPUSTNIH FUNKCIJ.

Nekoliko splošneje: če so $l_i : U \rightarrow \mathbb{R}$ omejeni linearni funkcionali, je prostor dopustnih funkcij podan z

$$V = \{y \in U \mid \forall i. l_i(y) = A_i\}.$$

Definicija. DOPUSTNA VARIACIJA je vsaka funkcija $v \in U$, za katero velja $u + tv \in V$ za vsak t .

Za dopustne variacije velja $l_i(u+tv) = l_i(u) + tl_i(v)$. Če je $u+tv \in V$, je $l_i(u+tv) = A_i$, in torej $l_i(v) = 0$ za vsak i . Prostor dopustnih variacij je torej podan s predpisom

$$\text{Var} = \{v \in U \mid l_i(v) = 0 \forall i\} = \bigcap_i \ker l_i.$$

Prostor Var je torej linearen podprostor v U , prostor V pa je afin podprostor v U , modeliran s prostorom Var . Če se vrnemo k osnovnemu variacijskemu problemu, velja $l_1(u) = u(a)$ in $l_2(u) = u(b)$.

Definicija. TESTNE FUNKCIJE na intervalu $[a, b]$ so funkcije $\varphi : [a, b] \rightarrow \mathbb{R}$, za katere velja

- $\varphi \in \mathcal{C}^\infty$,
- $\overline{\text{supp } \varphi} \subsetneq [a, b]$.

Vprašanje 39. Kaj so dopustne funkcije, dopustne variacije in testne funkcije? Kaj mora veljati za dopustne variacije?

Izrek (Osnovni izrek variacijskega računa). Naj bo $f : [a, b] \rightarrow \mathbb{R}$ zvezna funkcija. Če za vsako testno funkcijo φ velja

$$\int_a^b f(x)\varphi(x)dx = 0,$$

potem je $f(x) = 0$ na $[a, b]$.

Dokaz. Denimo, da ni tako, da obstaja $x_0 \in (a, b)$, za katerega je $f(x_0) = c > 0$. Zaradi zveznosti f obstaja $\delta > 0$, da velja $f(x) > c/2$ za $x \in (x_0 - \delta, x_0 + \delta)$. Naj bo φ testna funkcija, za katero je $\text{supp } \varphi \subset (x_0 - \delta, x_0 + \delta)$ in $\varphi(x) > 0$ na $\text{supp } \varphi$. Potem imamo

$$\int_a^b f(x)\varphi(x)dx = \int_{x_0-\delta}^{x_0+\delta} f(x)\varphi(x)dx > \frac{c}{2} \int_{x_0-\delta}^{x_0+\delta} \varphi(x)dx > 0,$$

kar je protislovno. □

Vprašanje 40. Povej in dokaži osnovni izrek variacijskega računa.

Predpostavimo, da je funkcija $\partial_{u'}L(x, u, u')$ odvedljiva po x , in z integracijo po delih izračunamo

$$\int_a^b \partial_{u'}Lv'dx = - \int_a^b \frac{d}{dx} (\partial_{u'}L) v dx + \partial_{u'}Lv|_a^b.$$

Odvod je torej enak

$$D_u\mathcal{L}(v) = \int_a^b \left(\partial_u L - \frac{d}{dx} \partial_{u'}L \right) v dx + \partial_{u'}Lv|_a^b = 0.$$

Ker so testne funkcije tudi dopustne variacije, dobimo, da za vse testne funkcije v velja

$$D_u \mathcal{L}(v) = \int_a^b \left(\partial_u L - \frac{d}{dx} \partial_{u'} L \right) v dx = 0,$$

oziroma, po osnovnem izreku variacijskega računa,

$$\partial_u L(x, \hat{u}(x), \hat{u}'(x)) - \frac{d}{dx} \partial_{u'} L(x, \hat{u}(x), \hat{u}'(x)) = 0.$$

Temu pravimo Euler-Lagrangeova enačba.

Vprašanje 41. Izpelji Euler-Lagrangeovo enačbo.

Če imamo samo en robni pogoj $u(a) = A$, bo za testne funkcije še vedno veljala Euler-Lagrangeova enačba, za ostale dopustne variacije pa dobimo novo enačbo

$$\partial_{u'} L(b) = 0.$$

Temu pravimo DINAMIČNI POGOJ.

1.9.1 Vezani ekstremini

Naj bodo $\phi_i : V \rightarrow \mathbb{R}$ nelinearni funkcionali. Množica $W \subseteq V$ naj bo podana s predpisom

$$W = \{u \in V \mid \forall i. \phi_i(u) = l_i\}.$$

Označimo $\vec{\phi} : V \rightarrow \mathbb{R}^n$ s komponentami $\vec{\phi} = (\phi_1, \dots, \phi_n)$ in $\vec{l} = (l_1, \dots, l_n)$. Velja $W = \vec{\phi}^{-1}(\vec{l})$. Predpostavljamo, da so v vseh točkah $u \in W$, ki nas bodo zanimala, funkcionali $\{D_u \phi_i\}_i$ linearno neodvisni v dualnem prostoru Var^* . Ker je za pot $v \in V$ z $v(0) = u$ odvod $\dot{v}(0)$ v tangentnem prostoru $T_u V = \text{Var}$, in ker velja

$$\left. \frac{d}{dt} \right|_{t=0} \phi_i(v(t)) = D_u \phi_i(\dot{v}(0)),$$

so $D_u \phi_i$ res v Var^* .

Naj bo $\mathcal{L} : V \rightarrow \mathbb{R}$ funkcional. Iščemo njegove ekstreme na podmnožici $W \subseteq V$. Videli smo, da je $\hat{u} \in W$ stacionarna točka $\mathcal{L} : W \rightarrow \mathbb{R}$, če za vsako krivuljo $v : (-\varepsilon, \varepsilon) \rightarrow W$, za katero je $v(0) = \hat{u}$, velja

$$D_u \mathcal{L}(\dot{v}(0)) = 0.$$

En način iskanja bi bil, da za vsak $u \in W$ poiščemo $T_u W$ in najdemo tisti \hat{u} , za katerega je $D_{\hat{u}} \mathcal{L} = 0$ na prostoru $T_{\hat{u}} W$. Tangentni prostor je enak

$$T_u W = \{\dot{v}(0) \mid \forall i. D_u \phi_i(\dot{v}(0)) = 0\} = \bigcap_{i=1}^n \ker D_u \phi_i.$$

Ker sta tako odvod kot tangentni prostor v vsaki točki različna, je ta način iskanja nepraktičen.

Vprašanje 42. Povej in razloži primitiven način iskanja vezanih ekstremov.

Boljši način je, da \mathcal{L} modificiramo tako, da bo za novi $\tilde{\mathcal{L}}$ veljal sklep, da če je $D_{\hat{u}}\tilde{\mathcal{L}}(v) = 0$ za vse dopustne variacije v , bo \hat{u} stacionarna točka \mathcal{L} . Vzemimo poljubno dopustno variacijo v in jo dekomponirajmo v obliko $v = v_u + v_u^\perp$, kjer je $v_u \in T_u W$, v_u^\perp pa pravokoten nanj. Naj bo $\{\varphi_1(u), \dots, \varphi_n(u)\}$ dualna baza baze $\{D_u\phi_i\}_i$, torej $D_u\phi_i \cdot \varphi_j(u) = \delta_{ij}$. Definiramo

$$v_u^\perp = \sum_{i=1}^n D_u\phi_i(v) \cdot \varphi_i(u)$$

in trdimo, da velja $v - v_u^\perp \in T_u W$. To je res, ker je

$$\begin{aligned} D_u\phi_j(v - v_u^\perp) &= D_u\phi_j(v) - D_u\phi_j\left(\sum_{i=1}^n D_u\phi_i(v) \cdot \varphi_i(u)\right) \\ &= D_u\phi_j(v) - \left(\sum_{i=1}^n D_u\phi_i(v) \cdot D_u\phi_j(\varphi_i(u))\right) \\ &= D_u\phi_j(v) - D_u\phi_j(v) \\ &= 0. \end{aligned}$$

Če definiramo

$$\tilde{\mathcal{L}}(u) = \mathcal{L}(u) - \sum_{i=1}^n D_u\mathcal{L}(\varphi_i(u))\phi_i(u),$$

bo veljalo $D_u\tilde{\mathcal{L}}(v_u^\perp) = 0$, kar lahko preverimo s podobnim računom. Označimo $\lambda_i(u) = D_u\mathcal{L}(\varphi_i(u))$. Odvajanje nam da

$$D_u\tilde{\mathcal{L}} = D_u\mathcal{L} - \sum_{i=1}^n \lambda_i(u) D_u\phi_i - \sum_{i=1}^n D_u\lambda_i \cdot \phi_i(u).$$

Če ustrezno modificiramo λ_i , bo res veljalo $D_{\hat{u}}\lambda_i = 0$ v ekstremnih točkah, vendar to vodi v zelo kompliciran predpis. Alternativno lahko rečemo, da so λ_i konstantne.

Vprašanje 43. Izpelj drugi način reševanja problemov vezanih ekstremov.

Strategija reševanja variacijskih problemov z vezmi je tedaj takšna: Prvo rešimo Euler-Lagrangeovo enačbo za

$$\tilde{\mathcal{L}} = \mathcal{L} - \sum_{i=1}^n \lambda_i \phi_i,$$

kjer dobimo rešitev \hat{u} , odvisno od λ_i . Parametre poiščemo s pomočjo pogojev $\phi_i(\hat{u}) = l_i$.

2 Mehanika

2.1 Osnove Newtonove mehanike

Definicija. AFIN PROSTOR \mathcal{A} nad vektorskim prostorom V je množica z binarno operacijo $+: \mathcal{A} \times V \rightarrow \mathcal{A}$, za katero velja:

- Za poljuben $\mathbf{A} \in \mathcal{A}$ ter $a, b \in V$ velja $(\mathbf{A} + a) + b = \mathbf{A} + (a + b)$
- Za poljubna $\mathbf{A}, \mathbf{B} \in \mathcal{A}$ obstaja natanko določen $a \in V$, da je $\mathbf{B} = \mathbf{A} + a$.

DIMENZIJA afinega prostora je enaka dimenziji vektorskega prostora V .

Definicija. Naj bo \mathcal{A} afin prostor nad vektorskim prostorom V . Definiramo operacijo odštevanja $\mathcal{A} \times \mathcal{A} \rightarrow V$ s predpisom

$$\mathbf{B} - \mathbf{A} = a \Leftrightarrow \mathbf{B} = \mathbf{A} + a.$$

Trditev. V afinem prostoru veljajo naslednje zveze:

- $\mathbf{A} - \mathbf{A} = 0$.
- $(\mathbf{A} - \mathbf{B}) + (\mathbf{B} - \mathbf{A}) = 0$.
- $(\mathbf{A} - \mathbf{B}) + (\mathbf{B} - \mathbf{C}) + (\mathbf{C} - \mathbf{A}) = 0$.
- $(\mathbf{A} - \mathbf{B}) + a = (\mathbf{A} + a) - \mathbf{B}$.
- $(\mathbf{A} - \mathbf{B}) + \mathbf{C} = (\mathbf{C} - \mathbf{B}) + \mathbf{A}$.

Definicija. Preslikava $g: \mathcal{A} \rightarrow \mathcal{A}'$ med afinima prostoroma je AFINA, če obstaja $dg \in L(V, V')$, da za vsaka $\mathbf{A}, \mathbf{B} \in \mathcal{A}$ velja $g(\mathbf{A}) - g(\mathbf{B}) = dg(\mathbf{A} - \mathbf{B})$.

Za afino preslikavo g si lahko izberemo POL \mathbf{O} , ter izpeljemo

$$g(\mathbf{A}) = g(\mathbf{O}) + dg(\mathbf{A} - \mathbf{O}).$$

Vrednosti funkcije seveda niso odvisne od izbire pola.

Vprašanje 1. Definiraj afin prostor in afino preslikavo.

Definicija. GALILEJEVA STRUKTURA je trojica $\mathcal{G} = (\mathcal{A}, \mathfrak{t}, \rho)$, kjer je \mathcal{A} štirirazsežni afin prostor nad V , $\mathfrak{t} \in L(V, \mathbb{R})$ in ρ euklidska metrika na $\ker \mathfrak{t}$, porojena z normo $\|\cdot\|$. Funkciji \mathfrak{t} pravimo ČASOVNOST, elementom \mathcal{A} pa pravimo DOGODKI. Pretečeni čas med dogodkoma \mathbf{A} in \mathbf{B} označimo s $\mathfrak{t}(\mathbf{A}, \mathbf{B})$. Dogodka sta ISTOČASNA, če je $\mathfrak{t}(\mathbf{A}, \mathbf{B}) = 0$. Za istočasne dogodke lahko definiramo razdaljo $\rho(\mathbf{A}, \mathbf{B}) = \|\mathbf{B} - \mathbf{A}\|$ (uporabimo isto oznako kot za metriko v $\ker \mathfrak{t}$).

Definicija. Galilejevi strukturi $\mathcal{G} = (\mathcal{A}, \mathfrak{t}, \rho)$ in $\mathcal{G}' = (\mathcal{A}', \mathfrak{t}', \rho')$ sta EKVIVALENTNI, če obstaja afina bijekcija $g: \mathcal{A} \rightarrow \mathcal{A}'$, ki ohranja časovnost in razdaljo med istočasnimi dogodki;

$$\mathfrak{t}'(g(\mathbf{A}) - g(\mathbf{B})) = \mathfrak{t}(\mathbf{A}, \mathbf{B}), \quad \rho'(g(\mathbf{A}), g(\mathbf{B})) = \rho(\mathbf{A}, \mathbf{B}).$$

Taki transformaciji pravimo GALILEJEVA TRANSFORMACIJA.

Vprašanje 2. Definiraj Galilejevo strukturo in Galilejeve transformacije.

Modelni primer je naravna Galilejeva struktura na $\mathcal{A} = \mathbb{R} \times \mathbb{E}$, kjer je \mathbb{E} trirazsežni Evklidski prostor. Za elemente $A_i = (t_i, \mathbf{P}_i) \in \mathcal{A}$ naravne strukture velja

- $t(A_1 - A_2) = t_1 - t_2$,
- $\rho(A_1, A_2) = \|\mathbf{P}_1 - \mathbf{P}_2\|$.

Definicija. KOORDINATNI SISTEM na \mathcal{A} je bijekcija $\phi : \mathcal{A} \rightarrow \mathbb{R} \times \mathbb{E}$ s komponentami $\phi(A) = (\tau\phi(A), \pi\phi(A))$, in pri kateri je $\tau \circ \phi$ linearna preslikava.

Opomba. Če sta ϕ in ϕ' koordinatna sistema, je preslikava $\phi' \circ \phi^{-1} : \mathbb{R} \times \mathbb{E} \rightarrow \mathbb{R} \times \mathbb{E}$ bijekcija.

Vprašanje 3. Kaj je koordinatni sistem?

Izrek. Galilejeva transformacija $g : \mathbb{R} \times \mathbb{E} \rightarrow \mathbb{R} \times \mathbb{E}$ je oblike

$$g(t, \mathbf{P}) = (t'_0 + t, \mathbf{P}'_0 + \vec{c}t + Q(\mathbf{P} - \mathbf{P}_0)),$$

kjer je $Q \in O(3)$ ortogonalna transformacija.

Dokaz. Ker je g afina preslikava, jo lahko zapišemo kot

$$g(t, \mathbf{P}) = g(t_0, \mathbf{P}_0) + dg(t - t_0, \mathbf{P} - \mathbf{P}_0),$$

kjer je $dg \in L(\mathbb{R}^4, \mathbb{R}^4)$. Če označimo $g(t_0, \mathbf{P}_0) = (t'_0, \mathbf{P}'_0)$, in zapišemo dg kot bločno matriko, dobimo

$$g(t, \mathbf{P}) = (t'_0, \mathbf{P}'_0) + \begin{bmatrix} \alpha & \vec{a}^T \\ \vec{c} & Q \end{bmatrix} \begin{bmatrix} t - t_0 \\ \mathbf{P} - \mathbf{P}_0 \end{bmatrix} = (t'_0, \mathbf{P}'_0) + \begin{bmatrix} \alpha(t - t_0) + \vec{a} \cdot (\mathbf{P} - \mathbf{P}_0) \\ (t - t_0) \cdot \vec{c} + Q(\mathbf{P} - \mathbf{P}_0) \end{bmatrix}.$$

Za dogodka (t_1, \mathbf{P}_1) in (t_2, \mathbf{P}_2) zahtevamo

$$t_2 - t_1 = \tau(g(t_2, \mathbf{P}_2) - g(t_1, \mathbf{P}_1)).$$

Če razvijemo desno stran zahteve po izpeljani formuli, dobimo pogoj

$$t_2 - t_1 = \alpha(t_2 - t_1) + \vec{a} \cdot (\mathbf{P}_2 - \mathbf{P}_1).$$

Iz tega sledi $\alpha = 1$ in $\vec{a} = \vec{0}$. Drug pogoj je, da se mora razdalja med istočasnimi dogodki ohranjati. Iz spodnjega dela bločne matrike dobimo pogoj

$$\|\mathbf{P}_2 - \mathbf{P}_1\| = \|Q(\mathbf{P}_2 - \mathbf{P}_1)\|,$$

torej mora biti Q ortogonalna. □

Vprašanje 4. Kakšno obliko imajo Galilejeve transformacije $\mathbb{R} \times \mathbb{E} \rightarrow \mathbb{R} \times \mathbb{E}$? Dokaži.

Če definiramo $\vec{v} = \dot{\mathbf{P}}$ in $\vec{a} = \dot{\vec{v}}$, lahko opazujemo, kako se ti količini obnašata pri Galilejevi transformaciji. V koordinatnem sistemu $\phi'(t', \mathbf{P}')$ velja $\vec{v}' = \partial_{t'} \mathbf{P}' = \dot{\mathbf{P}}'$ in $\vec{a}' = \dot{\vec{v}}'$. Izpeljemo $\vec{v}' = \vec{c} + Q\dot{\mathbf{P}}(t' - t'_0) = \vec{c} + Q\dot{\mathbf{P}}(t)$ in $\vec{a}' = Q\ddot{\mathbf{P}}(t)$.

Za sistem materialnih točk $\mathcal{P} = \{\mathbf{P}_1, \dots, \mathbf{P}_n\}$ lahko definiramo

$$\begin{aligned}\underline{\mathbf{P}} &= (\mathbf{P}_1, \dots, \mathbf{P}_n) \\ \underline{\mathbf{P}}'_0 &= (\mathbf{P}'_0, \dots, \mathbf{P}'_0) \\ \underline{\vec{c}} &= (\vec{c}, \dots, \vec{c}) \\ \underline{\mathbf{P}}' &= (\mathbf{P}'_1, \dots, \mathbf{P}'_n) = \underline{\mathbf{P}}'_0 + \underline{\vec{c}}t + Q(\underline{\mathbf{P}} - \underline{\mathbf{P}}_0)\end{aligned}$$

Gibanje lahko tedaj zapišemo s tremi principi.

- *Princip determiniranosti:* Trajektorija sistema materialnih točk \mathcal{P} je v danem koordinatnem sistemu natanko določena z začetnim položajem in hitrostjo. To pomeni, da obstaja funkcija interakcije \vec{f} , da velja

$$\ddot{\underline{\mathbf{P}}} = \vec{f}(t, \underline{\mathbf{P}}, \dot{\underline{\mathbf{P}}}).$$

- *Princip relativnosti:* Obstaja tak razred koordinatnih sistemov, v katerem je funkcija interakcije invariantna na Galilejeve transformacije. Temu razredu pravimo RAZRED INERCIJALNIH KOORDINATNIH SISTEMOV. To pomeni, da je funkcija interakcije invariantna v tem razredu,

$$\ddot{\underline{\mathbf{P}}}' = \vec{f}(t', \underline{\mathbf{P}}', \dot{\underline{\mathbf{P}}}').$$

- *Princip o sorazmernosti:* Obstajajo pozitivne konstante α_{ij} , da za vsako interakcijo med materialnimi točkami sistema $\mathcal{P} = (\mathbf{P}_1, \dots, \mathbf{P}_n)$ velja

$$\vec{f}_i = - \sum_{j \neq i} \alpha_{ji} \vec{f}_j.$$

Te konstante so enake za vse možne interakcije v sistemu.

Vprašanje 5. Kateri so principi gibanja?

Z ozirom na princip relativnosti izpeljemo $Q\ddot{\underline{\mathbf{P}}} = Q\vec{f}(t, \underline{\mathbf{P}}, \dot{\underline{\mathbf{P}}})$. Če v to enakost vstavimo vrednosti $t' = t'_0 + t$, $\vec{c} = \vec{0}$, $Q = I$ ter $\mathbf{P}'_0 = \mathbf{P}_0$, dobimo $\vec{f}(t'_0 + t, \underline{\mathbf{P}}, \dot{\underline{\mathbf{P}}}) = \vec{f}(t, \underline{\mathbf{P}}, \dot{\underline{\mathbf{P}}}_0)$, kar mora veljati za vsak t'_0 . Sledi, da funkcija \vec{f} ne more biti eksplicitno odvisna od časa. Tej ugotovitvi pravimo HOMOGENOST ČASA.

Če sedaj vstavimo $\vec{c} = \vec{0}$, $Q = I$ in $\mathbf{P}'_0 = \mathbf{P}_0 + \vec{a}$, kjer je \vec{a} poljuben vektor (in ne pospešek), izpeljemo $\mathbf{P}' = \mathbf{P} + \vec{a}$, in sledi $\vec{f}(\underline{\mathbf{P}} + \underline{\vec{a}}, \dot{\underline{\mathbf{P}}}) = \vec{f}(\underline{\mathbf{P}}, \dot{\underline{\mathbf{P}}})$, torej \vec{f} ne more biti odvisna od absolutnih položajev. Seveda je še vedno lahko odvisna od relativnih položajev (v tem primeru se \vec{a} odšteje). Tej lastnosti pravimo HOMOGENOST PROSTORA.

S poljubno izbiro vektorja \vec{c} in $Q = I$ lahko podobno izpeljemo, da je \vec{f} lahko odvisna le od relativnih hitrosti, čemur pravimo HOMOGENOST PROSTORA HITROSTI.

Če nenazadnje relaksiramo še pogoj na Q , dobimo

$$\vec{f}(Q(\mathbf{P}_i - \mathbf{P}_j), Q(\dot{\mathbf{P}}_i, \dot{\mathbf{P}}_j)) = \underline{Q}\vec{f}(\mathbf{P}_i - \mathbf{P}_j, \dot{\mathbf{P}}_i - \dot{\mathbf{P}}_j).$$

Funkcijam, ki zadoščajo temu pogoju, pravimo IZOTROPIČNE FUNKCIJE.

V posebnem primeru za $n = 1$ je \vec{f} konstantna funkcija (ker ne more biti odvisna od ničesar). Ker za vsak $Q \in O(3)$ velja $\vec{f} = Q\vec{f}$, mora biti $\vec{f} = \vec{0}$. Torej se prosta materialna točka v inercialnem koordinantem sistemu premika premočrtno s konstantno hitrostjo. To je ena od implikacij v prvem Newtonovem zakonu.

Vprašanje 6. Izpelji homogenost časa in faznega prostora iz principov gibanja.

Definicija. Interakcija \vec{f} je PARSKA, če lahko zapišemo

$$\vec{f}_i = \sum_{j \neq i} \vec{f}_{ji}(\mathbf{P}_i - \mathbf{P}_j, \dot{\mathbf{P}}_j - \dot{\mathbf{P}}_i)$$

za vse indekse i .

Definicija. Interakcija \vec{f} je LOKALNA, če je parska in če velja

$$\lim_{\mathbf{P}_i - \mathbf{P}_j \rightarrow \infty} \vec{f}_{ji} = \vec{0}.$$

Vprašanje 7. Definiraj parske in lokalne interakcije.

Lema. Za števila α_{ij} iz principa sorazmernosti velja

- $\alpha_{ij}\alpha_{ji} = 1$,
- $\alpha_{ij}\alpha_{jk}\alpha_{kj} = 1$.

Dokaz. Prva točka: Izberemo si take interakcije \vec{f}_k , ki so parske in lokalne in ki so neodvisne od relativnih hitrosti. Vse točke razen i in j pošljemo v neskončnost, da je njihov vpliv ničeln. Tedaj velja $\vec{f}_i = -\alpha_{ji}\vec{f}_j$ in $\vec{f}_j = -\alpha_{ij}\vec{f}_i$, torej $\vec{f}_i = \alpha_{ij}\alpha_{ji}\vec{f}_i$.

Druga točka: Izberemo si indekse i, j, k in podobno kot prej pošljemo druge točke v neskončnost. Ob predpostavki parske in lokalne interakcije tako dobimo

$$\begin{aligned} \vec{f}_i &= -\alpha_{ji}\vec{f}_j - \alpha_{ki}\vec{f}_k, \\ \vec{f}_j &= -\alpha_{ij}\vec{f}_i - \alpha_{kj}\vec{f}_k. \end{aligned}$$

Če vstavimo drugo enačbo v prvo,

$$\vec{f}_i = \alpha_{ji}\alpha_{ij}\vec{f}_i + \alpha_{ji}\alpha_{kj}\vec{f}_k - \alpha_{ki}\vec{f}_k,$$

2 Mehanika

nam člen na levi in prvi člen na desni po prvi točki odpadeta. Dobljeno enačbo še pomnožimo z α_{ik} in nam ostane

$$\vec{f}_k = \alpha_{ji}\alpha_{kj}\alpha_{ik}\vec{f}_k.$$

□

Lema. Naj za pozitivna števila α_{ij} velja ugotovitev prejšnje leme. Potem obstajajo števila m_i , da je $\alpha_{ji} = m_j/m_i$.

Dokaz. Števila α_{ij} so definirana le za $i \neq j$. Definicijo lahko razširimo, da je $\alpha_{ii} = 1$. Definiramo $l_{ij} = \log \alpha_{ij}$. Velja $l_{ii} = 0$ in $l_{ij} = -l_{ji}$, poleg tega pa tudi $l_{ij} + l_{jk} + l_{ki} = 0$.

Izberemo si indeks i_0 , ki nam bo definiral enoto mase. Velja $l_{i_0j} + l_{jk} + l_{ki_0} = 0$, kar odštejemo od prejšnje vsote treh členov in dobimo

$$l_{ij} - l_{i_0j} + l_{ki} - l_{ki_0} = 0.$$

Od tu izpeljemo, da za poljubna j in k velja

$$l_{ij} - l_{i_0j} = l_{ik} - l_{i_0k},$$

torej je $n_{ii_0} = l_{ij} - l_{i_0j}$ dobro definirana količina. Opazimo, da za $i = j$ velja $n_{ii_0} = l_{ii_0}$.

Definiramo $m_i = \exp n_{ii_0}$. Sledi

$$\log \alpha_{ij} = l_{ij} = l_{i_0j} + n_{ii_0} = -l_{ji} + n_{ii_0} = -n_{ji_0} + n_{ii_0} = \log m_i - \log m_j = \log \frac{m_i}{m_j}.$$

□

Opomba. Številom m_i pravimo INERCIJSKE MASE.

Vprašanje 8. Kaj so inercialne mase? Dokaži, da res obstajajo.

Produktu $m\vec{f} = \vec{F}$ pravimo SILA. Iz parskosti sledi

$$\vec{F}_i = \sum_{j \neq i} \vec{F}_{ji}(\mathbf{P}_i - \mathbf{P}_j, \dot{\mathbf{P}}_i - \dot{\mathbf{P}}_j).$$

Naj velja $\sum_{i \neq j} \vec{F}_{ji} = \vec{0}$. Predpostavimo, da so sile lokalne, in fiksiramo indeksa $k \neq l$. Če vsa ostala telesa pošljemo v neskončnost, ostane

$$\vec{F}_{kl} + \vec{F}_{lk} = \vec{0}.$$

S tem smo dokazali tretji Newtonov zakon.

Trditev (tretji Newtonov zakon). Če so vse sile parske in lokalne, velja $\vec{F}_{kl} = -\vec{F}_{lk}$.

Za nadaljevanje potrebujemo še dodaten princip gibanja, ki ga imenujemo *princip o masi*. Pravi, da je inercialna masa enaka v vseh koordinatnih sistemih.

Vprašanje 9. Kaj je princip o masi?

Najpreprostejši primer sile je gravitacija. Med točkama (m_1, \mathbf{P}_1) in (m_2, \mathbf{P}_2) deluje sila

$$\vec{F}_{21} = \frac{\kappa M_1 M_2}{|\mathbf{P}_1 - \mathbf{P}_2|^2} \frac{\mathbf{P}_2 - \mathbf{P}_1}{|\mathbf{P}_2 - \mathbf{P}_1|}.$$

Številoma M_1 in M_2 pravimo GRAVITACIJSKI MASI. Z eksperimentiranjem je Newton ugotovil, da so pravzaprav enake inercialnim masam.

Definicija. Zunanja sila $\vec{F} = \vec{F}(t, \mathbf{P}, \dot{\mathbf{P}})$ je POTENCIALNA, če obstaja potencial U , da je $\vec{F} = -\vec{\nabla}_{\mathbf{P}} U$.

Definicija. DELO sile \vec{F} pri gibanju materialne točke od \mathbf{P}_1 do \mathbf{P}_2 je krivuljni integral

$$A = \int_{\mathbf{P}_1}^{\mathbf{P}_2} \vec{F} \cdot d\mathbf{P} = \int_{t_1}^{t_2} \vec{F} \cdot \dot{\mathbf{P}} dt,$$

kjer smo pot parametrizirali s $\mathbf{P}(t)$. Produktu $\vec{F} \cdot \dot{\mathbf{P}}$ pravimo MOČ.

Definicija. KINETIČNA ENERGIJA T je enaka $\frac{1}{2}m |\dot{\mathbf{P}}|^2$.

Za rezultanto vseh sil \vec{F} na telo m lahko izpeljemo

$$A = \int_{t_1}^{t_2} \vec{F} \cdot \dot{\mathbf{P}} dt = \int_{t_1}^{t_2} m \ddot{\mathbf{P}} \cdot \dot{\mathbf{P}} dt = m \int_{t_1}^{t_2} \partial_t \left(\frac{1}{2} \dot{\mathbf{P}} \cdot \dot{\mathbf{P}} \right) dt = T_2 - T_1,$$

kar lahko zapišemo v izrek.

Izrek (izrek o delu). Delo rezultante vseh sil je enako razliki kinetične energije telesa.

Vprašanje 10. Povej in dokaži izrek o delu.

Definicija. Sila je KONZERVATIVNA v danem razredu inercialnih koordinatnih sistemov, če obstaja inercialni koordinatni sistem, v katerem je \vec{F} potencialna in odvisna samo od položaja.

Tedaj je \vec{F} potencialna, torej velja $\vec{F} = -\vec{\nabla} U$ za nek potencial U , ki mu pravimo POTENCIALNA ENERGIJA. Velja

$$A = \int_{t_1}^{t_2} \vec{F} \cdot \dot{\mathbf{P}} dt = - \int_{t_1}^{t_2} \vec{\nabla} U \cdot \dot{\mathbf{P}} dt = U(\mathbf{P}_1) - U(\mathbf{P}_2).$$

Vidimo, da je delo odvisno le od začetnega in končnega položaja. Sledi $T_2 - T_1 = U_1 - U_2$, torej je $T_1 + U_1 = T_2 + U_2 = E_0$ konstantna vrednost.

Izrek (izrek o energiji). Če je rezultanta vseh sil konzervativna, je vsota kinetične in potencialne energije konstanta gibanja.

Vprašanje 11. Povej in dokaži izrek o energiji.

2.2 Premočrtno gibanje

Definicija. Gibanje je PREMOČRTNO, če ima pospešek konstantno smer.

Primer takega gibanja je poševni met. Opazimo, da lahko vedno izberemo koordinatni sistem, v katerem tir poti leži na premici: Če je $\vec{a} = a\vec{e}$, kjer je \vec{e} konstanten vektor, velja

$$\vec{v} = \vec{e} \int_{t_0}^t a dt + \vec{v}_0.$$

Izberemo lahko sistem, kjer je \vec{v}_0 enak $\vec{0}$, in bo torej \vec{v} vzporeden \vec{e} .

Če gibanje poteka pod vplivom konzervativne sile, lahko zapišemo potencial U , in velja izrek o energiji

$$\frac{1}{2}m\dot{x}^2 + U(x) = E_0.$$

Od tod izpeljemo

$$\dot{x} = \pm \sqrt{\frac{2}{m}(E_0 - U(x))}.$$

Enačbo z ločljivimi spremenljivkami tedaj integriramo in dobimo

$$\pm \int_{x_0}^x \frac{dx}{\sqrt{\frac{2}{m}(E_0 - U(x))}} = \int_{t_0}^t dt = t - t_0.$$

Dobimo funkcijo $t = t(x)$. Če se na poti ne ustavimo, po izreku o inverzni preslikavi obstaja funkcija $x = x(t)$. Pravimo, da je premočrtno gibanje INTEGRABILNO.

Če v kvalitativni analizi ugotovimo, da je neko gibanje periodično med točkama a in b , lahko periodo izračunamo kot

$$\begin{aligned} T &= \int_{x_0}^b \frac{dx}{\sqrt{\frac{2}{m}(E_0 - U(x))}} - \int_a \frac{dx}{\sqrt{\frac{2}{m}(E_0 - U(x))}} + \int_a^{x_0} \frac{dx}{\sqrt{\frac{2}{m}(E_0 - U(x))}} \\ &= \sqrt{2m} \int_a^b \frac{dx}{\sqrt{E_0 - U(x)}}. \end{aligned}$$

Ker je $E = U(x)$ v krajiščih, je to posplošen integral. Situacija v obeh krajiščih je simetrična, torej preverimo le za levo krajišče, da integral res konvergira. V prvem koraku razvijemo preslikavo U v Taylorjev polinom prve stopnje v točki a , kjer se pojavi vrednost odvoda U v neki točki ξ blizu a . Ker je odvod zvezen, obstaja tak $\delta > 0$, da za $\xi \in [a, a + \delta]$ velja $2\partial_x U(a) < \partial_x U(\xi) < \frac{1}{2}\partial_x U(a)$, torej

$$\int_a^{a+\delta} \frac{dx}{\sqrt{E_0 - U(x)}} = \int_a^{a+\delta} \frac{dx}{\sqrt{-\partial_x U(\xi)(x-a)}} \leq \int_a^{a+\delta} \frac{1}{\sqrt{-\frac{1}{2}\partial_x U(a)}} \frac{1}{\sqrt{x-a}} dx < \infty.$$

Vprašanje 12. Izpelji izraz za periodo premočrtnega potencialnega gibanja.

Lema. Za $a < b$ velja

$$\int_a^b \frac{dx}{\sqrt{(b-x)(x-a)}} = \pi.$$

Dokaz. Uvedemo novo spremenljivko $x = \frac{1}{2}(a+b) + \frac{1}{2}(b-a)z$, s čimer se integral spremeni v

$$\int_{-1}^1 \frac{dz}{\sqrt{(1-z)(1+z)}} = \pi.$$

□

Primer. Oglejmo si harmonični oscilator, ki deluje pod potencialom $U = \frac{1}{2}kx^2$. Tedaj velja $F = -\partial_x U = -kx$, torej $m\ddot{x} = -kx$. Rešitev tega sistema je $x = A \cos \omega t + B \sin \omega t$ za $\omega = \sqrt{k/m}$. Iz tega lahko kar direktno preberemo $T = 2\pi/\omega$. Posebnost harmoničnega oscilatorja je, da je T neodvisen od E_0 . Takemu gibanju pravimo **IZOHRONIČNO**, harmonični potencial je edini primer izohroničnega potenciala, ki je simetričen glede na svoj minimum.

Če ima potencial lokalni minimum v x_0 , lahko za določanje potenciala uporabimo harmonično aproksimacijo. Zapišemo

$$\hat{U}(x) = U(x_0) + \partial_x U(x_0)(x - x_0) + \frac{1}{2}\partial_{x^2} U(x_0)(x - x_0)^2,$$

kar je harmonični potencial s periodo

$$T = 2\pi \sqrt{\frac{m}{\partial_{x^2} U(x_0)}}.$$

Ta aproksimacija je dobra, če velja $E_0 - U(x) \ll 1$.

Vprašanje 13. Izpelj harmonično aproksimacijo.

Druga vrsta aproksimacije, ki jo lahko uporabimo, je **LIBRACIJSKA**. Računamo

$$t = \text{sgn } \dot{x} \int_{x_0}^x \frac{dx}{\sqrt{\frac{2}{m}(E_0 - U(x))}} = \sqrt{\frac{m}{2}} \text{sgn } \dot{x} \int_{\theta_0}^{\theta} \frac{\frac{1}{2}(b-a)(-\sin \theta)d\theta}{\frac{1}{2}(b-a)\sqrt{\chi(\theta)}\sqrt{1 - \cos^2 \theta}}$$

za substitucijo $x = \frac{1}{2}(a+b) + \frac{1}{2}(b-a)\cos \theta$ in $E_0 - U(x) = (x-a)(b-x)\chi(x)$. Pri tem smo si x predstavljali kot kosinus kota v krožnici, ki poteka skozi točki a in b in ima središče na njuni zveznici. Če računamo dalje, dobimo

$$t = \sqrt{\frac{m}{2}} \int_{\theta_0}^{\theta} \frac{1}{\sqrt{\chi(\theta)}} d\theta.$$

Če želimo dobiti periodo gibanja, bo θ tekel od 0 do π .

$$T = 2\sqrt{\frac{m}{2}} \int_0^{\pi} \frac{d\theta}{\sqrt{\chi(\theta)}}.$$

Ta integral aproksimiramo s trapezno formulo, ki je natančna v primeru, da je funkcija v integralu afina. Rešitev je tedaj

$$T \doteq \pi \sqrt{\frac{m}{2}} \left(\frac{1}{\sqrt{\chi(a)}} + \frac{1}{\sqrt{\chi(b)}} \right).$$

Vprašanje 14. Izpelj libracijsko aproksimacijo.

Trditev. Za premočrtno potencialno periodično gibanje velja $\frac{dS}{dE_0} = \frac{T}{m}$ za ploščino S faznega diagrama.

Dokaz. Velja $S = 2\sqrt{\frac{m}{2}} \int_a^b \sqrt{E_0 - U} dx$, torej

$$\frac{dS}{dE_0} = 2\sqrt{\frac{m}{2}} \left(b' \sqrt{E_0 - U(b)} - a' \sqrt{E_0 - U(a)} + \int_a^b \frac{dx}{2\sqrt{E_0 - U}} \right).$$

Ker je $U(a) = U(b) = E_0$, sta prva dva člena v oklepaju enaka 0, torej

$$\frac{dS}{dE_0} = \sqrt{\frac{m}{2}} \int_a^b \frac{dx}{\sqrt{E_0 - U}} = \frac{T}{m}.$$

□

Vprašanje 15. Kako se pri nihanju ploščina faznega portreta spreminja z energijskim nivojem E_0 ?

2.3 Gibanje po krivulji

Dana je krivulja $\vec{r} = \vec{r}(s(t))$, kjer je s naravni parameter. Če s \mathbf{P} označimo trenutno lokacijo, velja

$$\vec{v} = \frac{d\mathbf{P}}{dt} = \frac{d\mathbf{P}}{ds} \frac{ds}{dt} = \vec{e}_t \dot{s},$$

kjer je \vec{e}_t enotski vektor, tangenten na krivuljo, in

$$\vec{a} = \ddot{s}\vec{e}_t + \dot{s} \frac{d\vec{e}_t}{dt} = \ddot{s}\vec{e}_t + \dot{s}^2 \kappa \vec{e}_n.$$

V enačbi κ predstavlja ukrivljenost, \vec{e}_n pa normalo na krivuljo. Drugi Newtonov zakon poleg rezultante vseh sil vsebuje tudi silo vezi \vec{S} . Razpisan v smereh krivuljnega koordinatnega sistema ima obliko

$$\begin{aligned} m\ddot{s} &= \vec{F} \cdot \vec{e}_t + \vec{S} \cdot \vec{e}_t \\ m\kappa\dot{s}^2 &= \vec{F} \cdot \vec{e}_n + \vec{S} \cdot \vec{e}_n \\ 0 &= \vec{F} \cdot \vec{e}_b + \vec{S} \cdot \vec{e}_b \end{aligned}$$

Tu imamo štiri neznanke (s, \vec{S}) ter tri enačbe, torej potrebujemo še dodatno konstitutivno relacijo za silo vezi. Če se omejimo na gladke krivulje (take, kjer ni trenja), dobimo dodatno enačbo

$$\vec{S} \cdot \vec{e}_t = 0.$$

Delo take sile vezi je enako 0. Če je \vec{F} konzervativna sila, $\vec{F} = -\vec{\nabla}U$, dobimo

$$m\ddot{s} = -\frac{dU}{ds},$$

iz česar lahko izpeljemo energijsko enačbo

$$\frac{1}{2}m\dot{s}^2 + U(s) = E_0.$$

Torej je gibanje po gladki krivulji pod vplivom konzervativne sile reducibilno na premočrtno gibanje v ločni dolžini.

Vprašanje 16. Na kaj se reducira gibanje po gladki krivulji pod vplivom konzervativne sile? Izpelji.

Vprašanje 17. Obravnavaj matematično nihalo kot gibanje po gladki krožnici.

Odgovor: Če je l polmer krožnice, velja $s = l\theta$. Na točko poleg sile vezi deluje tudi teža, ki ima potencial

$$U = -m\vec{g}\vec{r} = -mgl \cos \frac{s}{l}.$$

Za dovolj majhen E_0 je gibanje periodično, in velja

$$T = \sqrt{2m} \int_{-s_0}^{s_0} \frac{ds}{\sqrt{E_0 + mgl \cos \frac{s}{l}}}.$$

Če uporabimo substitucijo $s = l\theta$ in začetno energijo zapišemo z začetnim odklonom, dobimo

$$T = \sqrt{2ml} \int_{-\theta_0}^{\theta_0} \frac{d\theta}{mgl(\cos \theta - \cos \theta_0)}.$$

Na tej točki upoštevamo, da je funkcija soda, in uporabimo $\cos \theta = 1 - 2 \sin^2 \frac{\theta}{2}$;

$$T = 2\sqrt{\frac{l}{g}} \int_0^{\theta_0} \frac{d\theta}{\sqrt{\sin^2 \frac{\theta_0}{2} - \sin^2 \frac{\theta}{2}}},$$

kar se s substitucijo $\sin \frac{\theta}{2} = u \sin \frac{\theta_0}{2}$ končno predela na

$$T = 4\sqrt{\frac{l}{g}} \int_0^1 \frac{du}{\sqrt{(1-u^2)(1-u^2 \sin^2 \frac{\theta_0}{2})}}.$$

To je eliptični integral, odgovor je

$$T = 4\sqrt{\frac{l}{g}} K\left(\sin^2 \frac{\theta_0}{2}\right)$$

☒

2.4 Gibanje v polju centralne sile

Definicija. Sila $\vec{F} = \vec{F}(\mathbf{P})$ je CENTRALNA, če obstaja točka \mathbf{O} (pol sile), da \vec{F} deluje v smeri zveznice med \mathbf{P} in \mathbf{O} , in da je njena velikost odvisna le od razdalje.

Trditev. Konzervativna sila \vec{F} je centralna natanko tedaj, ko obstaja pol, okoli katerega je vrtilna količina konstantna.

Dokaz. Recimo, da je \vec{F} centralna. Tedaj za vrtilno količino okoli \mathbf{O} , velja

$$\begin{aligned}\vec{l}(\mathbf{O}, \mathbf{P}) &= (\mathbf{P} - \mathbf{O}) \times m\dot{\mathbf{P}} \\ \partial_t \vec{l} &= \dot{\mathbf{P}} \times m\dot{\mathbf{P}} + (\mathbf{P} - \mathbf{O}) \times m\ddot{\mathbf{P}} = (\mathbf{P} - \mathbf{O}) \times \vec{F}\end{aligned}$$

Če je \vec{F} centralna, je vzporedna $\mathbf{P} - \mathbf{O}$, in se izniči tudi drugi člen.

Recimo, da je vrtilna količina konstantna. Po zgornjem izračunu $\dot{\vec{l}} = (\mathbf{P} - \mathbf{O}) \times \vec{F} = 0$, torej je \vec{F} vzporedna zveznici. Pokazati moramo še, da je velikost sile odvisna le od razdalje. Po predpostavki je sila konzervativna, torej $\vec{F} = -\vec{\nabla}U$ in

$$\vec{F} = -\left(\partial_r U \vec{e}_r + \frac{1}{r} \partial_\theta U \vec{e}_\theta + \frac{1}{r \sin \theta} \partial_\varphi U \vec{e}_\varphi\right)$$

v sferičnih koordinatah. Ker je sila vzporedna zveznici \vec{e}_r , mora veljati $U = U(r)$. \square

Vprašanje 18. Definiraj centralno silo in jo karakteriziraj ob predpostavki, da je konzervativna. Dokaži karakterizacijo.

Trditev. Zvezna centralna sila je konzervativna.

Dokaz. Definiramo

$$U = \int_{|\mathbf{P}_0 - \mathbf{O}|}^{|\mathbf{P} - \mathbf{O}|} F(r) dr + U(\mathbf{P}_0).$$

\square

Vprašanje 19. Dokaži: zvezna centralna sila je konzervativna.

Trditev. Gibanje v polju centralne sile je ravninsko. Dogaja se na ravnini, ki vsebuje center sile in ima normalo v smeri \vec{l} .

Dokaz. Računamo

$$(\mathbf{P} - \mathbf{P}_0) \cdot \vec{l} = (\mathbf{P} - \mathbf{O}) \cdot \vec{l} + (\mathbf{O} - \mathbf{P}_0) \cdot \vec{l} = 0 + 0 = 0,$$

upoštevaje definicijo \vec{l} . \square

Vprašanje 20. Dokaži, da je gibanje v polju centralne sile ravninsko.

Izpeljimo kinematiko v polarnem koordinatnem sistemu. Definiramo

$$\begin{aligned}\vec{e}_r &= \cos \theta \vec{i} + \sin \theta \vec{j}, \\ \vec{e}_\theta &= \partial_\theta \vec{e}_r = -\sin \theta \vec{i} + \cos \theta \vec{j},\end{aligned}$$

in izračunamo

$$\begin{aligned}\vec{r} &= r \vec{e}_r, \\ \vec{v} &= \dot{r} \vec{e}_r + r \dot{\theta} \vec{e}_\theta, \\ \vec{a} &= (\ddot{r} - r \dot{\theta}^2) \vec{e}_r + (r \ddot{\theta} + 2 \dot{r} \dot{\theta}) \vec{e}_\theta.\end{aligned}$$

Prvi komponenti hitrosti pravimo RADIALNA HITROST, drugi OBODNA HITROST. Podobno prvi komponenti pospeška pravimo RADIALNI POSPEŠEK, drugi pa OBODNI POSPEŠEK.

Za vrtilno količino velja

$$\vec{l} = \vec{r} \times m \vec{v} = m r^2 \dot{\theta} \vec{k}.$$

Pogledamo lahko tudi ploščinsko hitrost

$$\dot{\vec{A}} = \frac{1}{2} \vec{r} \times \vec{v},$$

iz česar dobimo $\vec{l} = 2m \dot{\vec{A}}$ oziroma $\dot{\vec{A}} = \frac{1}{2} r^2 \dot{\theta} \vec{k}$. Za gibanje v polju centralne sile je \vec{l} konstantna, torej sta konstantni tudi ploščinska hitrost in DVOJNA PLOŠČINSKA HITROST

$$c_0 = r^2 \dot{\theta}.$$

Vprašanje 21. Izpelji kinematiko v polarnem koordinatnem sistemu. Kaj je dvojna ploščinska hitrost?

Računamo lahko

$$\dot{r} = \frac{dr}{d\theta} \dot{\theta} = \frac{dr}{d\theta} \frac{c_0}{r^2} = -c_0 \partial_\theta \left(\frac{1}{r} \right),$$

iz česar s spremenljivko $u = 1/r$, $u' = \partial_\theta u$ izpeljemo

$$\ddot{r} = -c_0 u'' \dot{\theta} = -c_0^2 u^2 u''.$$

To uporabimo v BINETOVİ FORMULI

$$a_r = \ddot{r} - r \dot{\theta}^2 = -c_0^2 u^2 (u + u'').$$

Vprašanje 22. Izpelji Binetovo formulo.

Prvi Keplerjev zakon pravi, da se planeti gibljejo okoli Sonca v elipsah. Elipso lahko parametriziramo kot

$$r = \frac{p}{1 + \varepsilon \cos \theta}.$$

2 Mehanika

Z uporabo Binetove formule dobimo

$$a_r = -c_0^2 \frac{1}{pr^2}.$$

Za silo gravitacije $\vec{F} = -\kappa m M u^2 \vec{e}_r$ bo potem veljalo

$$\frac{c_0^2}{p} = \kappa M = \text{konst.}$$

Za tako parametrizacijo elipse velja

$$a = \frac{p}{1 - \varepsilon^2}$$
$$b = \frac{p}{\sqrt{1 - \varepsilon^2}}$$

Če je T perioda gibanja, je

$$A = \frac{1}{2} c_0 T,$$

torej

$$T = \frac{2\pi ab}{c_0}.$$

Drugi Keplerjev zakon pravi, da je kvadrat periode gibanja sorazmeren kubu večje polosi elipse, $T^2 = k a^3$, torej

$$\frac{p}{c_0^2} = \frac{k}{4\pi^2} = \frac{1}{\kappa M}.$$

Dobimo, da je k konstanten za vse planete.

Vprašanje 23. Povej drugi Keplerjev zakon. Pokaži, da je koeficient enak za vse planete.

Centralna sila je potencialna, $\vec{F} = -\vec{\nabla} V$. Velja energijska enačba

$$\frac{1}{2} m v^2 + V(r) = E_0,$$

ki jo lahko predelamo v

$$\frac{1}{2} m \dot{r}^2 + \frac{1}{2} m r^2 \dot{\theta}^2 + V(r) = E_0.$$

Upošteva je $r^2 \dot{\theta} = c_0$ dobimo

$$\frac{1}{2} m \dot{r}^2 + \frac{l^2}{2mr^2} + V(r) = E_0.$$

Če zadnja dva člena na levi strani pospravimo v EFEKTIVNI POTENCIAL $U(r)$, smo reducirali gibanje na premočrtno s potencialom. Če to razrešimo na $r = r(t)$, lahko zapišemo

$$\theta(t) = \int_{t_0}^t \frac{c_0}{r^2(t)} dt.$$

Pravimo, da je gibanje v polju centralne sile INTEGRABILNO.

Vprašanje 24. Pokaži, da je gibanje v polju centralne sile integrabilno.

V primeru gravitacije integral žal ni zaprte oblike. Dobimo pa lahko enačbo trajektorije:

$$\frac{dr}{d\theta} = \frac{dr}{dt} \frac{dt}{d\theta} = \dot{r} \frac{1}{\dot{\theta}} = \pm \frac{1}{c_0} r^2 \sqrt{\frac{2}{m}(E_0 - U(r))}.$$

Za $c_0 = l/m$ po integraciji dobimo

$$\theta - \theta_0 = \pm \frac{l}{\sqrt{2m}} \int_{r_0}^r \frac{dr}{r^2 \sqrt{E_0 - U(r)}}.$$

Za gravitacijsko silo in $\gamma = \kappa m M$ velja

$$U(r) = \frac{l^2}{2mr^2} - \frac{\gamma}{r}.$$

Po integriranju dobimo

$$p = \frac{l^2}{m\gamma}, \quad \varepsilon = \sqrt{\frac{2l^2}{m\gamma^2} E_0 + 1}.$$

Za $E_0 < 0$ je $\varepsilon < 1$, in dobimo elipso, pri $E_0 = 0$ dobimo parabolo $\varepsilon = 1$ in pri $E_0 > 0$ imamo hiperbolo za $\varepsilon > 1$.

Vprašanje 25. Kakšna je oblika tira planeta glede na energijo? Izpelji.

Če je gibanje periodično glede na efektivni potencial, se točka giblje med krožnicama med dvema APSIDNIMA RADIJEMA. Tir točke se dotika teh krožnic. Manjšemu od radijev pravimo PERICENTER, večjemu pa APOCENTER.

Trditev. Tir je simetričen glede na apsidni radij.

Dokaz. Velja

$$\begin{aligned} \theta^+ - \theta_0 &= \frac{l}{\sqrt{2m}} \int_{r_a}^r \frac{dr}{r^2 \sqrt{E_0 - U(r)}} \\ \theta^- - \theta_0 &= -\frac{l}{\sqrt{2m}} \int_{r_a}^r \frac{dr}{r^2 \sqrt{E_0 - U(r)}} \end{aligned}$$

Sledi $\theta^+ - \theta_0 = \theta_0 - \theta^-$. □

Izračunamo lahko OVOJNO ŠTEVILO

$$\Delta\theta = \frac{l}{\sqrt{2m}} \int_{r_a}^{r_b} \frac{dr}{r^2 \sqrt{E_0 - U(r)}}.$$

Tir gibanja bo zaprt takrat, ko je $\frac{\Delta\theta}{\pi} \in \mathbb{Q}$.

Trditev. *Tir je zaprt ali pa je gosta množica v kolobarju $K(0, r_a, r_b)$.*

Izrek (Bertrand). *Vsi tiri v okolici krožnega tira so zaprti natanko tedaj, ko je V gravitacijski ali Hookov potencial.*

Vprašanje 26. Kakšen je tir gibanja točke v polju centralne sile?

2.5 Relativno gibanje

Definicija. Koordinatni sistem $\varphi(t, \mathbf{P})$ se GIBLJE glede na koordinatni sistem $\varphi'(t', \mathbf{P}')$, če obstaja trojica $(\mathbf{P}_0, \mathbf{P}'_0, Q)$, kjer je $\mathbf{P}_0 \in \mathbb{E}$, $\mathbf{P}'_0 : \mathbb{R} \rightarrow \mathbb{E}$, in $Q : \mathbb{R} \rightarrow SO(3)$, da velja $t = t'$ in

$$\mathbf{P}'(t) = \mathbf{P}'_0(t) + Q(t)(\mathbf{P} - \mathbf{P}_0).$$

Predpostavimo, da je φ' inercialen, in ga imenujmo ABSOLUTNI KOORDINATNI SISTEM, φ pa je RELATIVNI KOORDINATNI SISTEM.

Trditev. *Rotacijski del gibanja je neodvisen od izbire trojice.*

Dokaz. Recimo

$$\mathbf{P}' = \mathbf{P}'_0 + Q(\mathbf{P} - \mathbf{P}_0) = \tilde{\mathbf{P}}'_0 + \tilde{Q}(\mathbf{P} - \tilde{\mathbf{P}}_0).$$

Za $\mathbf{P} = \mathbf{P}_0$ dobimo

$$\mathbf{P}'_0 = \tilde{\mathbf{P}}'_0 + \tilde{Q}(\mathbf{P}_0 - \tilde{\mathbf{P}}_0).$$

Če to vstavimo nazaj gor, pridemo do

$$Q(\mathbf{P} - \mathbf{P}_0) = \tilde{Q}(\mathbf{P} - \mathbf{P}_0),$$

torej $Q = \tilde{Q}$. □

Vprašanje 27. Kdaj se koordinatni sistem giblje glede na nek drugi koordinatni sistem? Dokaži, da je rotacijski del neodvisen od izbire trojice.

Za odvod velja

$$\vec{v}' = \dot{\mathbf{P}}'_0 + \dot{Q}(\mathbf{P} - \mathbf{P}_0) + Q\dot{\mathbf{P}} = Q(Q^T \vec{v}'_0 + Q^T \dot{Q}(\mathbf{P} - \mathbf{P}_0) + \vec{v}_{\text{rel}}).$$

Prvemu členu pravimo TRANSLATORNA HITROST, drugemu ROTACIJSKA HITROST, tretjemu pa RELATIVNA HITROST.

Trditev. $Q^T \dot{Q}$ je poševno simetrični tenzor.

Dokaz. Velja $Q^T Q = I$. Če to odvajamo, dobimo

$$\partial_t(Q^T)Q + Q^T \dot{Q} = 0.$$

□

Izrek. Naj bo W poševno simetričen na trirazsežnem evklidskem prostoru. Potem obstaja vektor $\vec{\omega}$ tako, da je $W\vec{a} = \vec{\omega} \times \vec{a}$ za vsak \vec{a} .

Dokaz. Vemo, da obstaja lastna vrednost $W\vec{p} = \lambda\vec{p}$. Ker je

$$\lambda |\vec{p}|^2 = \vec{p} \cdot W\vec{p} = W^T \vec{p} \cdot \vec{p} = -\lambda |\vec{p}|^2,$$

torej $\lambda = 0$.

BŠS je $|\vec{p}| = 1$. Dopolnimo ga lahko do ortonormirane baze prostora $\{\vec{p}, \vec{q}, \vec{r}\}$, kjer je $\vec{r} = \vec{p} \times \vec{q}$. Ker je W poševno simetričen, je $\vec{a} \cdot W\vec{a} = 0$ za poljuben \vec{a} , in dobimo

$$W\vec{q} = W\vec{q} \cdot \vec{r} \cdot \vec{r}, \quad W\vec{r} = W\vec{r} \cdot \vec{q} \cdot \vec{q}.$$

Za poljuben \vec{a} je

$$W\vec{a} = \vec{a} \cdot \vec{q} \cdot W\vec{q} \cdot \vec{r} \cdot \vec{r} + \vec{a} \cdot \vec{r} \cdot W\vec{r} \cdot \vec{q} \cdot \vec{q}$$

Velja $\vec{q} \times \vec{r} = \vec{p}$ in $\vec{r} \times \vec{p} = \vec{q}$, torej

$$W\vec{a} = W\vec{q} \cdot \vec{r} \cdot (\vec{a} \cdot \vec{q} \cdot \vec{p} \times \vec{q} + \vec{a} \cdot \vec{r} \cdot \vec{p} \times \vec{r}) = W\vec{q} \cdot \vec{r} \cdot \vec{p} \times \vec{a}.$$

□

Vprašanje 28. Dokaži, da poševno simetričen tenzor deluje kot vektorski produkt.

Definicija. $[W_1, W_2] = W_1 W_2 - W_2 W_1$

Prostor poševno simetričnih tenzorjev z operacijama $+$ in $[\cdot]$ je algebra.

Trditev. Preslikava $W \mapsto \vec{\omega}(W)$ je homomorfizem.

Dokaz. Pokazati želimo $\vec{\omega}([W_1, W_2]) = \vec{\omega}(W_1) \times \vec{\omega}(W_2)$. Velja

$$W_1 W_2 \vec{a} = \vec{\omega}_1 \times (\vec{\omega}_2 \times \vec{a}) = (\vec{\omega}_1 \cdot \vec{a}) \cdot \vec{\omega}_2 - (\vec{\omega}_1 \cdot \vec{\omega}_2) \cdot \vec{a},$$

torej dobimo

$$[W_1, W_2] \vec{a} = W_1 W_2 \vec{a} - W_2 W_1 \vec{a} = (\vec{\omega}_1 \times \vec{\omega}_2) \times \vec{a}.$$

□

Trditev. $(A\vec{a}) \times (A\vec{b}) = A^*(\vec{a} \times \vec{b})$

Dokaz. Če so $\vec{a}, \vec{b}, \vec{c}$ linearno neodvisni, velja

$$\det A = \frac{[A\vec{a}, A\vec{b}, A\vec{c}]}{[\vec{a}, \vec{b}, \vec{c}]}.$$

Tudi če niso neodvisni, velja

$$(A\vec{a} \times A\vec{b}) \cdot \vec{c} = [A\vec{a}, A\vec{b}, \vec{c}].$$

Za začetek predpostavimo, da je A obrnljiva. Tedaj

$$\begin{aligned}
 (A\vec{a} \times A\vec{b}) \cdot \vec{c} &= [A\vec{a}, A\vec{b}, AA^{-1}\vec{c}] \\
 &= \det A \cdot [\vec{a}, \vec{b}, A^{-1}\vec{c}] \\
 &= \det A \cdot (\vec{a} \times \vec{b}) \cdot A^{-1}\vec{c} \\
 &= (\vec{a} \times \vec{b}) \cdot (A^*)^T \vec{c} \\
 &= A^*(\vec{a} \times \vec{b}) \cdot \vec{c},
 \end{aligned}$$

kjer smo upoštevali $A^{-1} = (A^*)^T / \det A$. Če A ni obrnljiva, pa obstaja A_ε , da je $|A - A_\varepsilon| < \varepsilon$, torej v limiti velja za vse matrike. \square

Posledica. Če je $Q \in SO(3)$, je $Q(\vec{a} \times \vec{b}) = Q\vec{a} \times Q\vec{b}$.

Vprašanje 29. Dokaži: če je $Q \in SO(3)$, je $Q(\vec{a} \times \vec{b}) = Q\vec{a} \times Q\vec{b}$.

Definicija. Vektor kotne hitrosti rotacije Q je vektor $\vec{\omega}' = Q\vec{\omega}$, kjer je $\vec{\omega}$ osni vektor poševno simetričnega tenzorja $W = Q^T \dot{Q}$.

Trditev. Osni vektor poševno simetričnega tenzorja $\dot{Q}Q^T$ je vektor kotne hitrosti rotacije Q .

Dokaz. Računamo

$$\begin{aligned}
 \dot{Q}Q^T \vec{a} &= QQ^T \dot{Q}Q^T \vec{a} \\
 &= Q\vec{\omega} \times (Q^T \vec{a}) \\
 &= Q\vec{\omega} \times (QQ^T \vec{a}) \\
 &= (Q\vec{\omega}) \times \vec{a}.
 \end{aligned}$$

\square

Trditev. Vektor kotne hitrosti rotacije $R(\vec{e}, \varphi)$ okoli stalne osi \vec{e} za kot φ je $\vec{\omega}' = \dot{\varphi}\vec{e}$.

Dokaz. Imenujmo to rotacijo Q . Velja $Q\vec{e} = \vec{e}$; če to odvajamo po času, dobimo $\dot{Q}\vec{e} = 0$. Če to množimo z leve s Q^T , pridemo do $Q^T \dot{Q}\vec{e} = \vec{\omega} \times \vec{e} = 0$, torej je $\vec{\omega}$ vzporeden \vec{e} .

Dokazati moramo še, da je $\omega = \dot{\varphi}$. Naj bo \vec{f} poljuben vektor dolžine 1. Velja $\vec{f} \cdot Q\vec{f} = \cos \varphi$. Če to odvajamo po času, dobimo

$$\vec{f} \cdot \dot{Q}\vec{f} = -\sin \varphi \dot{\varphi}$$

oziroma

$$Q^T \vec{f} \cdot Q^T \dot{Q}\vec{f} = -\dot{\varphi} \sin \varphi.$$

Levo stran enakosti lahko razpišemo v

$$Q^T \vec{f} \cdot Q^T \dot{Q}\vec{f} = Q^T \vec{f} \cdot (\vec{\omega} \times \vec{f}) = \vec{\omega} \cdot (\vec{f} \times Q^T \vec{f}) = \vec{\omega} \cdot (-\sin \varphi \vec{e}).$$

Sledi $\omega = \dot{\varphi}$. \square

Vprašanje 30. Kaj je vektor kotne hitrosti rotacije? Kako ga izrazimo v primeru stalne osi?

Trditev. Naj bo W poševno simetričen tenzor z enotskim osnim vektorjem \vec{e} . Potem je $e^{\theta W}$ rotacija okoli \vec{e} za kot θ .

Dokaz. Naj bo $A = e^{\theta W}$. Prvo dokažimo, da je $A \in SO(3)$. Če označimo $B = A^T A$ in to enakost odvajamo po θ , dobimo $B' = 0$, torej je B konstanta. Za $\theta = 0$ dobimo $B = I$. Pokazati moramo še, da je determinanta A enaka 1. Za to prvo pokažimo $\det e^{\theta W} = e^{\text{sl} W}$. Če je $f(\theta) = \det e^{\theta W}$, je

$$f' = \frac{\partial \det A}{\partial \theta} * WA = A^* * WA = I * WA(A^*)^T = I * WAA^{-1} \det A = \det A \text{sl} W,$$

kjer $*$ predstavlja skalarni produkt, če matrike vektoriziramo. Sledi $f' = \text{sl} W \cdot f$, rešitev diferencialne enačbe je $f(\theta) = e^{\theta \text{sl} W}$. Ker je $W^T = -W$, je $\text{sl} W = 0$ in $\det A = 1$.

Razpis A v Taylorjevo vrsto pokaže $A\vec{e} = \vec{e}$. Če definiramo $Q(t) = e^{t\theta W}$, lahko izračunamo $Q^T \dot{Q} = \theta W$, iz česar dobimo osni vektor $\vec{\omega} = \theta \vec{e}$. To je enak osni vektor kot za rotacijo $R(\vec{e}, \theta)$. \square

Izrek. $R(\vec{e}, \varphi) = \cos \varphi I + (1 - \cos \varphi) \vec{e} \otimes \vec{e} + \sin \varphi W(\vec{e})$

Dokaz. V izrazu nastopa tenzorski produkt $(\vec{a} \otimes \vec{b})\vec{c} = (\vec{b} \cdot \vec{c})\vec{a}$.

Velja $R(\vec{e}, \varphi) = e^{\varphi W(\vec{e})}$. Kratek račun pokaže

$$W^k = \begin{cases} (-1)^{n+1}(\vec{e} \otimes \vec{e} - I) & k = 2n, n \geq 1 \\ (-1)^n W & k = 2n + 1 \end{cases}$$

Če rotacijsko matriko razvijemo v Taylorjevo vrsto in zberemo lihe in sode člene, dobimo natanko zeleno obliko. \square

Vprašanje 31. Kako izrazimo rotacijo?

Iz vse te izražave lahko končno izračunamo

$$\begin{aligned} \dot{\mathbf{P}}' &= \dot{\mathbf{P}}'_0 + \dot{Q}(\mathbf{P} - \mathbf{P}_0) + Q\dot{\mathbf{P}} \\ &= \dot{\mathbf{P}}'_0 + QQ^T \dot{Q}(\mathbf{P} - \mathbf{P}_0) + Q\dot{\mathbf{P}} \\ &= \dot{\mathbf{P}}'_0 + Q(\vec{\omega} \times (\mathbf{P} - \mathbf{P}_0)) + Q\dot{\mathbf{P}} \\ &= \vec{v}'_0 + \vec{\omega}' \times Q(\mathbf{P} - \mathbf{P}_0) + Q\vec{v}_{\text{rel}} \\ &= \vec{v}'_0 + \vec{\omega}' \times \vec{\zeta}' + \vec{v}'_{\text{rel}} \end{aligned}$$

Za $\zeta = \mathbf{P} - \mathbf{P}_0$. Če je \vec{u} poljuben vektor, velja

$$\partial_t \vec{u}' = \partial_t (Q\vec{u}) = \dot{Q}\vec{u} + Q\dot{\vec{u}} = Q(\vec{\omega} \times \vec{u}) + Q\dot{\vec{u}},$$

torej transformacija in odvod po času komutirata le v primeru, da je $\vec{\omega}' \times \vec{u}' = 0$. Primer takega vektorja je $\vec{\omega}$, torej velja $\partial_t \vec{\omega}' = (\dot{\vec{\omega}})'$. Sedaj lahko izračunamo tudi pospešek

$$\begin{aligned}\vec{a}' &= \partial_t \vec{v}' \\ &= \partial_t \vec{v}'_0 + \dot{\vec{\omega}}' \times \vec{\zeta}' + \vec{\omega}' \times \partial_t \vec{\zeta}' + \partial_t \vec{v}'_{\text{rel}} \\ &= \vec{a}'_0 + \dot{\vec{\omega}}' \times \vec{\zeta}' + \vec{\omega}' \times (\vec{\omega}' \times \vec{\zeta}') + 2\vec{\omega}' \times \vec{v}'_{\text{rel}} + Q\vec{v}'_{\text{rel}} \\ &= \vec{a}'_0 + \dot{\vec{\omega}}' \times \vec{\zeta}' + \vec{\omega}' \times (\vec{\omega}' \times \vec{\zeta}') + 2\vec{\omega}' \times \vec{v}'_{\text{rel}} + \vec{a}'_{\text{rel}}.\end{aligned}$$

Prvi člen tega predpisa imenujemo POSPEŠEK IZHODIŠČA RKS, drugi člen EULERJEV POSPEŠEK, tretji CENTRIFUGALNI POSPEŠEK, četrti CORIOLISOV POSPEŠEK in peti RELATIVNI POSPEŠEK.

Vprašanje 32. Izpelji predpis za pospešek relativnega koordinatnega sistema in poi-menuj člene.

Izrek. Koordinatna sistema $\varphi(t, \mathbf{P})$ in $\varphi'(t, \mathbf{P}')$ sta v istem razredu Galilejevih koordinatnih sistemov natanko tedaj, ko koordinatna transformacija $\mathbf{P} \mapsto \mathbf{P}'$ preslika premočrtno gibanje s konstantno brzino v premočrtno gibanje s konstantno brzino.

Dokaz. V desno je očitno iz zapisa pospeška. V levo računamo

$$\begin{aligned}\vec{0} &= \vec{a}_0 + \dot{\vec{\omega}}' \times \vec{\zeta}' + \vec{\omega}' \times (\vec{\omega}' \times \vec{\zeta}') + 2\vec{\omega}' \times \vec{v}'_{\text{rel}}, \\ \mathbf{P}' &= \mathbf{P}'_0 + tv_0 \vec{e}, \\ \vec{0} &= \vec{a}'_0 + tv_0(\dot{\vec{\omega}}' \times \vec{e} + \vec{\omega}' \times (\vec{\omega}' \times \vec{e})) + 2v_0 \vec{\omega}' \times \vec{e}.\end{aligned}$$

To velja za vsak t , torej tudi za $t = 0$,

$$\vec{0} = \vec{a}'_0 + 2v_0 \vec{\omega}' \times \vec{e}.$$

Ta del je neodvisen od t , torej se pri poljubnem t pokrajša in dobimo

$$\vec{0} = \dot{\vec{\omega}}' \times \vec{e} + \vec{\omega}' \times (\vec{\omega}' \times \vec{e}).$$

Enačbo skalarno množimo z \vec{e} ,

$$0 = \vec{\omega}' \times (\vec{\omega}' \times \vec{e}) \cdot \vec{e}.$$

To je mešani produkt, ki ga lahko ciklično zamenjamo in dobimo

$$0 = (\vec{\omega}' \times \vec{e}) \cdot (\vec{e} \times \vec{\omega}'),$$

iz česar sledi $\vec{\omega}' \parallel \vec{e}$. Sedaj upoštevamo prejšnjo enačbo in vidimo $\vec{a}'_0 = \vec{0}$. Če vzamemo drug primer premočrtnega gibanja $\mathbf{P}' = \mathbf{P}'_0 + tv_0 \vec{f}$, kjer \vec{f} ni vzporeden \vec{e} , spet dobimo $\vec{\omega}' \parallel \vec{f}$, torej mora veljati $\vec{\omega}' = 0$. \square

Vprašanje 33. Kaj velja za premočrtno gibanje pri koordinatni transformaciji? Doka-ži.

2.6 Sistem materialnih točk

Imamo točke \mathbf{P}_i, m_i za $i = 1, \dots, N$. Definiramo masno središče

$$\mathbf{P}_* = \mathbf{O} + \frac{1}{m} \sum_{i=1}^N m_i (\mathbf{P}_i - \mathbf{O}).$$

Izračunamo lahko, da je neodvisno od izbire točke \mathbf{O} , poleg tega pa velja tudi

$$m\ddot{\mathbf{P}}_* = \sum_{i=1}^N \vec{F}_i$$

za zunanje sile \vec{F}_i . Za vsako točko velja

$$m_i \ddot{\mathbf{P}}_i = \vec{F}_i + \sum_{j \neq i} \vec{F}_{ji}$$

Če to z leve vektorsko pomnožimo s $\mathbf{P}_i - \mathbf{O}$ in seštejemo po i , dobimo

$$\sum_{i=1}^N (\mathbf{P}_i - \mathbf{O}) \times m_i \ddot{\mathbf{P}}_i = \sum_{i=1}^N (\mathbf{P}_i - \mathbf{O}) \times \vec{F}_i + \sum_{i=1}^N \sum_{j \neq i} (\mathbf{P}_i - \mathbf{O}) \times \vec{F}_{ji}.$$

Prvi člen predstavlja rezultanto navora zunanjih sil, drugi pa rezultanto navora notranjih sil. Če definiramo vrtilno količino $\vec{l}(\mathbf{O}, \mathbf{P}_i) = (\mathbf{P}_i - \mathbf{O}) \times m_i \dot{\mathbf{P}}_i$ in $\vec{l} = \sum_i \vec{l}(\mathbf{O}, \mathbf{P}_i)$, lahko torej zapišemo

$$\dot{\vec{l}} = \vec{N}(\mathbf{O}) + \vec{N}_*(\mathbf{O}).$$

Definicija. Sila \vec{F}_{ji} je CENTRALNA, če je oblike

$$\vec{F}_{ji} = F_{ji}(|\mathbf{P}_i - \mathbf{P}_j|) \frac{\mathbf{P}_i - \mathbf{P}_j}{\|\mathbf{P}_i - \mathbf{P}_j\|}.$$

Trditev. Če so vse notranje sile centralne, je navor notranjih sil enak 0.

Dokaz. Zapišemo

$$\sum_i \sum_j (\mathbf{P}_i - \mathbf{O}) \times \vec{F}_{ji} = \frac{1}{2} \left(\sum_i \sum_j (\mathbf{P}_i - \mathbf{O}) \times \vec{F}_{ji} + \sum_j \sum_i (\mathbf{P}_j - \mathbf{O}) \times \vec{F}_{ij} \right).$$

To razpišemo in dobimo 0. □

Izrek (Izrek o vrtilni količini). Če so vse notranje sile centralne, je odvod vrtilne količine enak rezultanti navora zunanjih sil, $\dot{\vec{l}}(\mathbf{O}) = \vec{N}(\mathbf{O})$. Pri tem predpostavimo, da je pol \mathbf{O} fiksna točka.

Vprašanje 34. Povej in dokaži izrek o vrtilni količini.

Če je \mathbf{P}_0 poljubna točka, lahko izračunamo

$$\vec{l}(\mathbf{P}_0) = \vec{l}(\mathbf{O}) + m(\mathbf{O} - \mathbf{P}_0) \times \dot{\mathbf{P}}_*.$$

V nadaljevanju predpostavimo, da so vse notranje sile centralne. Če zgornjo enačbo odvajamo, dobimo

$$\dot{\vec{l}}(\mathbf{P}_0) = \vec{l}(\mathbf{P}_0) - \dot{\mathbf{P}}_0 \times m\dot{\mathbf{P}}_*.$$

Izrek o vrtilni količini torej ohrani obliko za premikajoč pol, če velja:

- $\dot{\mathbf{P}}_0 = \vec{0}$,
- $\dot{\mathbf{P}}_* = \vec{0}$,
- $\mathbf{P}_0 = \mathbf{P}_*$.

Vprašanje 35. Za katere pole izrek o vrtilni količini ohranja obliko?

Če je sistem zaprt, velja $\vec{F} = 0$, iz česar dobimo $\vec{N} = 0$. Sklepamo, da se masno središče giblje premočrtno s konstantno hitrostjo, in da je vrtilna količina konstanta. Dodatno je konstantna $m\dot{\mathbf{P}}_*$, torej velja izrek o ohranitvi gibalne količine.

2.7 Togo telo

Definicija. Gibanje sistema materialnih točk je TOGO, če gibanje ohranja razdalje med točkami.

Definicija. Telo je TOGO, če nima koles. Alternativno je telo togo, če je edini njegov način gibanja togo gibanje.

Izrek. Izometrija je afina preslikava.

Dokaz. Izberimo \mathbf{P}_0 in točke $\mathbf{A}, \mathbf{B}, \mathbf{C}$, ki niso kolinearne. Naj bo \mathbf{P} poljubna točka. Za $\vec{r} = \mathbf{P} - \mathbf{P}_0$, $\vec{a} = \mathbf{A} - \mathbf{P}_0$, $\vec{b} = \mathbf{B} - \mathbf{P}_0$ in $\vec{c} = \mathbf{C} - \mathbf{P}_0$ velja

$$\vec{r} = \alpha\vec{a} + \beta\vec{b} + \gamma\vec{c}.$$

Izometrija naj slika \mathbf{T} v \mathbf{T}' , poleg tega naj je

$$\vec{r}' = \alpha'\vec{a}' + \beta'\vec{b}' + \gamma'\vec{c}'.$$

Če enačbo za \vec{r} skalarno množimo z \vec{a} , \vec{b} oziroma \vec{c} , dobimo enačbe

$$\vec{r} \cdot \vec{a} = \alpha\vec{a} \cdot \vec{a} + \beta\vec{b} \cdot \vec{a} + \gamma\vec{c} \cdot \vec{a}$$

$$\vec{r} \cdot \vec{b} = \alpha\vec{a} \cdot \vec{b} + \beta\vec{b} \cdot \vec{b} + \gamma\vec{c} \cdot \vec{b}$$

$$\vec{r} \cdot \vec{c} = \alpha\vec{a} \cdot \vec{c} + \beta\vec{b} \cdot \vec{c} + \gamma\vec{c} \cdot \vec{c}$$

To zložimo v matriko A , da dobimo

$$\begin{bmatrix} \vec{r} \cdot \vec{a} \\ \vec{r} \cdot \vec{b} \\ \vec{r} \cdot \vec{c} \end{bmatrix} = A \begin{bmatrix} \alpha \\ \beta \\ \gamma \end{bmatrix}$$

in podobno za količine s črto. Izometrija ohranja tako razdalje kot kote, torej sta levi strani sistemov in matriki A, A' enaki. Sledi torej $\alpha = \alpha', \beta = \beta'$ in $\gamma = \gamma'$. Dobimo

$$\mathbf{P}' = \mathbf{P}'_0 + \vec{r}' = \mathbf{P}'_0 + A\vec{r} = \mathbf{P}'_0 + A(\mathbf{P} - \mathbf{P}_0)$$

in $A \in O(3)$. □

Togo gibanje lahko zapišemo v obliki

$$\mathbf{P}(t) = \mathbf{P}_0(t) + Q(t)(\mathbf{P}(t=0) - \mathbf{P}_0(t=0)).$$

Predstavimo ga kot relativno gibanje, kjer je RKS položaj telesa, kot ga vidimo ob času $t = 0$, AKS pa položaj telesa ob času t . RKS se giblje skupaj s telesom, zato mu pravimo TELESNI KS, absolutnemu pa pravimo PROSTORSKI KS.

Vprašanje 36. Kako opišeš gibanje togega telesa?

Za telo B definiramo volumensko in masno mero in predpostavimo, da za $B' = B(t)$ velja $m(B') = m(B)$. Sledi

$$m(B') = \iiint_{B'} \rho' dV' = \iiint_B \rho'(\mathbf{P}'(\mathbf{P})) \left| \det \frac{\partial \mathbf{P}'}{\partial \mathbf{P}} \right| dV.$$

Ker pa je $\mathbf{P}' = \mathbf{P}'_0 + Q(\mathbf{P} - \mathbf{P}_0)$, je torej odvod enak Q in njegova determinanta enaka 1. Torej

$$\iiint_B \rho'(\mathbf{P}'(\mathbf{P})) dV = m(B') = m(B) = \iiint_B \rho(\mathbf{P}) dV,$$

torej (ker to velja tudi na delih telesa) $\rho'(\mathbf{P}'(\mathbf{P})) = \rho(\mathbf{P})$. Podobno izpeljemo, da je

$$\mathbf{P}'_* = \frac{1}{m(B)} \iiint_{B'} \mathbf{P}' - \mathbf{P}'_0 dm + \mathbf{P}'_0$$

enak $\mathbf{P}'_* = \mathbf{P}'_0 + Q(\mathbf{P}_* - \mathbf{P}_0)$ za \mathbf{P}_* , definiran s podobnim integralom. Podobno kot v diskretnem primeru lahko tudi tu izpeljemo

$$\dot{\mathbf{P}}'_* = \frac{1}{m(B')} \iiint_{B'} \dot{\mathbf{P}}' dm.$$

Vrtilno količino izračunamo kot

$$\begin{aligned}
 \vec{l}(\mathbf{P}'_0) &= \iiint_{B'} (\mathbf{P}' - \mathbf{P}'_0) \times \partial_t(\mathbf{P}' - \mathbf{P}'_0) dm' \\
 &= \iiint_{B'} \vec{\zeta}' \times (\vec{\omega}' \times \vec{\zeta}') dm' \\
 &= \iiint_{B'} |\vec{\zeta}'|^2 \vec{\omega}' - (\vec{\zeta}' \cdot \vec{\omega}') \vec{\zeta}' dm' \\
 &= \iiint_{B'} \left(|\vec{\zeta}'|^2 I - \vec{\zeta}' \otimes \vec{\zeta}' \right) dm' \cdot \vec{\omega}'.
 \end{aligned}$$

Iz tega dobimo definicijo prostorskega vztrajnostnega tenzorja

$$J'(\mathbf{P}'_0) = \iiint_{B'} \left(|\vec{\zeta}'|^2 I - \vec{\zeta}' \otimes \vec{\zeta}' \right) dm' = \iiint_{B'} |\mathbf{P}' - \mathbf{P}'_0|^2 I - (\mathbf{P}' - \mathbf{P}'_0) \otimes (\mathbf{P}' - \mathbf{P}'_0) dm'.$$

Podobno definiramo telesni vztrajnostni tenzor, izračunamo lahko $J' = QJQ^T$. Vztrajnostni tenzor je simetričen, poleg tega pa ga lahko diagonaliziramo

$$J(\mathbf{P}_0) = \sum_{i=1}^3 J_i \vec{w}_i \otimes \vec{w}_i,$$

kjer so J_i GLAVNE (lastne) vrednosti.

Trditev. Vztrajnostni tenzor je nenegativen. Če B ni tanka palica, je pozitivno definiten.

Dokaz. Trdimo, da je $\vec{e} \cdot J\vec{e} \geq 0$ za enotski \vec{e} . Po Cauchy-Schwarzu velja

$$\iiint_B |\vec{\zeta}|^2 - (\vec{e} \cdot \vec{\zeta})^2 dm \geq 0,$$

enakost dobimo natanko tedaj, ko je \vec{e} vzporeden $\vec{\zeta}$, kar se zgodi le, če je B tanka palica. \square

Vprašanje 37. Definiraj vztrajnostni tenzor in dokaži, da je pozitivno definiten.

Trditev. Normala na ravnino zrcalne simetrije telesa je glavna smer vztrajnostnega tenzorja s polom na tej ravnini.

Dokaz. Naj bo \vec{e} ta normala. Trdimo $J(\mathbf{P}_0)\vec{e} = J\vec{e}$ za neko lastno vrednost J . To je ekvivalentno temu, da je deviator

$$D\vec{e} = \vec{\zeta} \otimes \vec{\zeta}\vec{e}$$

vzporeden \vec{e} . Naj bo $\vec{f} \perp \vec{e}$ in izračunajmo

$$\vec{f} \cdot D\vec{e} = \iiint_B \vec{f} \cdot \vec{\zeta} \cdot \vec{e} \cdot \vec{\zeta} dm = \iiint_B \vec{f} \cdot \vec{\zeta}' \cdot \vec{e} \cdot \vec{\zeta}' dm$$

za zrcalno sliko $\vec{\zeta}'$. Izračunamo lahko, da velja $\vec{e} \cdot \vec{\zeta}' = -\vec{e} \cdot \vec{\zeta}$ in $\vec{f} \cdot \vec{\zeta}' = \vec{f} \cdot \vec{\zeta}$, torej velja $\vec{f} \cdot D\vec{e} = 0$. \square

Trditev. Naj ima B dve ravnini simetrije z normalama \vec{e}_1 in \vec{e}_2 . Potem ima $J(\mathbf{P}_0)$, kjer je \mathbf{P}_0 v preseku obeh ravnin, glavne smeri $\vec{e}_1, \vec{e}_2, \vec{e}_1 \times \vec{e}_2$.

Vprašanje 38. Kaj lahko poveš o vztrajnostnemu tenzorju, če ima eno ali dve ravnini zrcalne simetrije? Dokaži.

Če razpišemo predpis za tenzorski produkt z uporabo $\mathbf{P} - \mathbf{P}_0 = (\mathbf{P} - \mathbf{P}_*) + (\mathbf{P}_* - \mathbf{P}_0)$, dobimo

$$J(\mathbf{P}_0) = J(\mathbf{P}_*) + m |\mathbf{P}_* - \mathbf{P}_0|^2 I - m(\mathbf{P}_* - \mathbf{P}_0) \otimes (\mathbf{P}_* - \mathbf{P}_0).$$

Tej formuli pravimo STEINERJEV IZREK.

Vprašanje 39. Kaj pravi Steinerjev izrek?

Za kinetično energijo razpišemo

$$T = \frac{1}{2} \iiint_{B'} |\dot{\mathbf{P}}'|^2 dm',$$

pri čemer upoštevamo $\dot{\mathbf{P}}' = \dot{\mathbf{P}}_* + \vec{\omega}' \times (\mathbf{P}' - \mathbf{P}_*)$, ker se telo ne giblje v svojem sistemu in je zato $Q\dot{\mathbf{P}} = 0$. Na koncu dobimo

$$T = \frac{1}{2} m |\dot{\mathbf{P}}_*|^2 + \frac{1}{2} \vec{\omega} \cdot J(\mathbf{P}_*) \vec{\omega}.$$

Vprašanje 40. Izpelji predpis za kinetično energijo togega telesa.

Pri dinamiki togega telesa imamo tri principe:

- $m\ddot{\mathbf{P}}_* = \vec{F}'_{\text{zun}}$ (enačba gibanja masnega središča)
- $\partial_t \vec{l}'(\mathbf{P}'_*) = \vec{N}'(\mathbf{P}'_*)$ (princip o vrtilni količini)
- $\vec{N}'(\mathbf{P}'_0) = \vec{N}'(\mathbf{P}'_*) + (\mathbf{P}'_0 - \mathbf{P}'_*) \times \vec{F}'_{\text{zun}}$

Izrek. Naj bo \mathbf{P}'_0 točka trenutnega ali stalnega mirovanja togega telesa. Potem velja $\partial_t \vec{l}'(\mathbf{P}'_0) = \vec{N}'(\mathbf{P}'_0)$.

Dokaz. Računamo

$$\begin{aligned} \vec{l}'(\mathbf{P}'_0) &= Q\vec{l}(\mathbf{P}_0) \\ &= QJ(\mathbf{P}_0)\vec{\omega} \\ &= Q(J(\mathbf{P}_*) + m |\mathbf{P}_0 - \mathbf{P}_*|^2 I - m(\mathbf{P}_0 - \mathbf{P}_*) \otimes (\mathbf{P}_0 - \mathbf{P}_*))\vec{\omega} \\ &= Q\vec{l}(\mathbf{P}_*) + m |\mathbf{P}_0 - \mathbf{P}_*|^2 \vec{\omega} - m((\mathbf{P}_0 - \mathbf{P}_*) \cdot \vec{\omega})(\mathbf{P}'_0 - \mathbf{P}'_*) \\ &= \vec{l}'(\mathbf{P}'_*) + m((\mathbf{P}'_0 - \mathbf{P}'_*) \times \vec{\omega}') \times (\mathbf{P}'_0 - \mathbf{P}'_*) \\ &= \vec{l}'(\mathbf{P}'_*) + m\dot{\mathbf{P}}_* \times (\mathbf{P}_0 - \mathbf{P}_*). \end{aligned}$$

Zadnji sklep uporabi dejstvo $\dot{\mathbf{P}}'_* = \dot{\mathbf{P}}'_0 + \vec{\omega}' \times (\mathbf{P}'_* - \mathbf{P}'_0)$, ki ga dobimo z odvajanjem $\mathbf{P}'_* = \mathbf{P}'_0 + Q(\mathbf{P}_* - \mathbf{P}_0)$. Upoštevamo še dejstvo, da je $\dot{\mathbf{P}}'_0 = 0$, ker je to točka trenutnega mirovanja. Zgornjo enakost nato še odvajamo do

$$\begin{aligned}\partial_t \vec{l}'(\mathbf{P}'_0) &= \vec{N}'(\mathbf{P}'_*) + m\ddot{\mathbf{P}}'_* \times (\mathbf{P}_0 - \mathbf{P}_*) + m\dot{\mathbf{P}}'_* \times (-\dot{\mathbf{P}}'_*) \\ &= \vec{N}'(\mathbf{P}'_*) + (\mathbf{P}'_* - \mathbf{P}'_0) \times \vec{F}'_{\text{zun}} \\ &= \vec{N}'(\mathbf{P}'_0).\end{aligned}$$

□

Navor okoli masnega središča ali navor okoli točke mirovanja lahko razpišemo v

$$\begin{aligned}\vec{N}'(\mathbf{P}'_*) &= \partial_t \vec{l}'(\mathbf{P}'_*) \\ &= \partial_t (QJ(\mathbf{P}_*)\vec{\omega}) \\ &= \dot{Q}J(\mathbf{P}_*)\vec{\omega} + QJ(\mathbf{P}_*)\dot{\vec{\omega}} \\ &= \vec{\omega}' \times QJ(\mathbf{P}_*)\vec{\omega} + QJ(\mathbf{P}_*)\dot{\vec{\omega}} \\ &= Q(\vec{\omega} \times J(\mathbf{P}_*)\vec{\omega} + J(\mathbf{P}_*)\dot{\vec{\omega}}).\end{aligned}$$

Iz tega dobimo Eulerjeve dinamične enačbe

$$\vec{N}(\mathbf{P}_*) = J(\mathbf{P}_*)\dot{\vec{\omega}} + \vec{\omega} \times J(\mathbf{P}_*)\vec{\omega}.$$

Enako dobimo za točko trenutnega mirovanja.

Vprašanje 41. Izpelj Eulerjeve dinamične enačbe. Dokaži, da delujejo tudi v točki mirovanja.

V lastnem KS vztrajnostnega momenta je J diagonalen. Dobimo

$$\begin{aligned}J_1\dot{\omega}_1 - \omega_2\omega_3(J_2 - J_3) &= N_1 \\ J_2\dot{\omega}_2 - \omega_1\omega_3(J_3 - J_1) &= N_2 \\ J_3\dot{\omega}_3 - \omega_1\omega_2(J_1 - J_2) &= N_3\end{aligned}$$

Podobno lahko naredimo za rotacijo okoli stalne osi \vec{e}

$$\begin{aligned}J_{13}\ddot{\varphi} - J_{23}\dot{\varphi}^2 &= N_1 \\ J_{23}\ddot{\varphi} + J_{13}\dot{\varphi}^2 &= N_2 \\ J_{33}\ddot{\varphi} &= N_3.\end{aligned}$$

Vprašanje 42. Izpelj Eulerjeve dinamične enačbe v lastnem koordinatnem sistemu in za rotacijo okoli stalne osi.

Poznamo volumenske sile (B_1, \vec{f}_1) , za katere velja

$$\vec{F}_1 = \iiint_{B_1} \vec{f}_1 dV.$$

Trditev. Homogena volumenska sila je ekvivalentna točkovni sili s prijemališčem v masnem središču.

Dokaz. Za silo je očitno, za navor pa velja

$$\vec{N}(\mathbf{O}) = \iiint_B (\mathbf{P} - \mathbf{O}) \times \vec{f} dm = \left(\iiint_B \mathbf{P} - \mathbf{O} dm \right) \times \vec{f} = (\mathbf{P}_* - \mathbf{O}) \times \vec{F}.$$

□

Za odvod kinetične energije velja

$$\partial_t T = m \dot{\mathbf{P}}_*' \cdot \ddot{\mathbf{P}}_*' + \vec{\omega} \cdot J_* \dot{\vec{\omega}} = \dot{\mathbf{P}}_*' \cdot \vec{F}' + \vec{\omega}' \cdot \vec{N}'_*,$$

kjer smo upoštevali Eulerjeve dinamične enačbe. Recimo, da je \vec{F} volumenska sila, in da velja $\vec{f}' = -\vec{\nabla} \cdot u$. Definiramo

$$U = \iiint_{B'} u(\mathbf{P}') dm'$$

in računamo

$$\partial_t U = \iiint_{B'} \partial_t u(\mathbf{P}') dm' = \iiint_{B'} \vec{\nabla} \cdot u \cdot \dot{\mathbf{P}}' dm = - \iiint_{B'} \vec{f}' \cdot (\dot{\mathbf{P}}_*' + \vec{\omega} \times (\mathbf{P}' - \mathbf{P}'_*)) dm'$$

Sedaj ciklično zamenjamo člene mešanega produkta v drugem členu, do

$$\partial_t U = -\vec{F}' \cdot \dot{\mathbf{P}}_*' - \iiint_{B'} \vec{\omega} \cdot ((\mathbf{P}' - \mathbf{P}'_*) \times \vec{f}') dm = -\vec{F}' \cdot \dot{\mathbf{P}}_*' - \vec{\omega}' \cdot \vec{N}'_*.$$

Torej je $T + U$ konstantna. To je izrek o energiji za togo telo.

Vprašanje 43. Izpelji izrek o energiji za togo telo.

Definicija. Telo B_1 se kotili po B_2 , če imata telesi v dotikališču enako hitrost.

2.7.1 Prosta vrtavka

Vrtavka je prosta, če je $\vec{F}' = 0$ in $\vec{N}' = 0$. Eulerjeve dinamične enačbe imajo tedaj obliko

$$\begin{aligned} J_1 \dot{\omega}_1 - \omega_2 \omega_3 (J_2 - J_3) &= 0, \\ J_2 \dot{\omega}_2 - \omega_1 \omega_3 (J_3 - J_1) &= 0, \\ J_3 \dot{\omega}_3 - \omega_1 \omega_2 (J_1 - J_2) &= 0. \end{aligned}$$

Enačba ima več rešitev:

- Trivialna: $\omega_1 = \omega_2 = \omega_3 = 0$

2 Mehanika

- Enakomerno vrtenje okoli ene osi: $\omega_1 = \omega_2 = 0$, $\omega_3 = \text{konst.}$.
- Netrivialne rešitve

Če je vrtavka simetrična, $J_1 = J_2$, je ω_3 konstantna, in velja

$$\begin{aligned}\dot{\omega}_1 - \frac{\omega_3(J_1 - J_3)}{J_1}\omega_2 &= 0 \\ \dot{\omega}_2 - \frac{\omega_3(J_3 - J_1)}{J_2}\omega_1 &= 0\end{aligned}$$

Za $\Omega = \omega_3(J_1 - J_3)/J_1$ je rešitev sistema

$$\begin{aligned}\omega_1 &= A \cos(\Omega t - \delta) \\ \omega_2 &= -A \sin(\Omega t - \delta)\end{aligned}$$

Vektor $\vec{\omega}$ torej precesira s kotno hitrostjo Ω okoli tretje osi.

Vprašanje 44. Reši prosto simetrično vrtavko.

Če se preselimo v koordinatni sistem, v katerem je J diagonalen, velja $l_i = J_i \omega_i$. Velja $l^2 = l_1^2 + l_2^2 + l_3^2 = \text{konst.}$ in

$$T = \frac{1}{2}J_1\omega_1^2 + \frac{1}{2}J_2\omega_2^2 + \frac{1}{2}J_3\omega_3^2 = \frac{l_1^2}{2J_1} + \frac{l_2^2}{2J_2} + \frac{l_3^2}{2J_3}.$$

Iz tega dobimo

$$1 = \frac{l_1^2}{2TJ_1} + \frac{l_2^2}{2TJ_2} + \frac{l_3^2}{2TJ_3},$$

čemur pravimo BINETOV ELIPSOID. Vektor \vec{l} leži na preseku sfere $|\vec{l}| = l$ in Binetovega elipsoida. Če uredimo $J_1 \leq J_2 \leq J_3$, dobimo $1 \leq l^2/(2TJ_1)$, torej za polos velja $a_1 \leq l$. Podobno dobimo $a_3 \geq l$. Če je rotacija le okoli največje ali najmanjše osi, se sfera in elipsoid dotikata v točki; za malo zmoten elipsoid imamo v preseku dve krožnici. Za srednjo polos pa majhna perturbacija elipsoida razklene točko v dve krivulji, ki prideta daleč od točke. Enakomerna rotacija okoli srednje osi je torej nestabilna.

Vprašanje 45. Analiziraj stabilnost enakomerne rotacije proste vrtavke.

2.7.2 Eulerjevi koti

Rotacijo predstavimo kot kompozitum rotacij

$$R = R(\vec{k}'', \psi)R(\vec{i}', \theta)R(\vec{k}, \varphi).$$

Lema. Za $\vec{f}' = R(\vec{e}, \varphi)\vec{f}$ velja $R(\vec{f}', \psi)R(\vec{e}, \varphi) = R(\vec{e}, \varphi)R(\vec{f}, \psi)$.

Trditev. Velja $R(\vec{k}'', \psi)R(\vec{i}', \theta)R(\vec{k}, \varphi) = R(\vec{k}, \varphi)R(\vec{i}, \theta)R(\vec{k}, \psi)$.

Dokaz za lemo in trditev je preprost račun.

Trditev. Vsako rotacijo, ki ni rotacija okoli osi \vec{k} , lahko enolično napišemo z Eulerjevo rotacijo s koti $\theta, \varphi, \psi \in [0, 2\pi)$.

Dokaz. Naj bo $\vec{\varepsilon}_1, \vec{\varepsilon}_2, \vec{\varepsilon}_3$ rotiran koordinatni sistem. Za vozliščnico $\vec{k} \times \vec{\varepsilon}_3 = \vec{e}$ definiramo kote

- φ je kot med \vec{i} in \vec{e} ,
- ψ je kot med \vec{e} in $\vec{\varepsilon}_1$,
- θ je kot med \vec{k} in $\vec{\varepsilon}_3$.

Za dokaz enoličnosti pogledamo delovanje preslikave

$$R(\vec{i}, \theta_1) = R(\vec{k}, \varphi_2 - \varphi_1) R(\vec{i}, \theta_2) R(\vec{k}, \psi_2 - \psi_1)$$

na vektorjih \vec{i} in \vec{k} . □

Vprašanje 46. Dokaži, da se lahko vsako rotacijo zapiše na enoličen način z Eulerjevimi rotacijami.

Trditev. Vektor kotne hitrosti Eulerjeve rotacije je $\vec{\omega}' = \dot{\varphi}\vec{k} + \dot{\theta}\vec{i}' + \dot{\psi}\vec{k}''$.

Dokaz. Označimo $R = R(\vec{k}, \varphi) R(\vec{i}, \theta) R(\vec{k}, \psi) = R_1 R_2 R_3$. Po definiciji kotne hitrosti je

$$\begin{aligned} \dot{R}\vec{a} &= \vec{\omega}' \times R\vec{a} \\ &= (\dot{R}_1 R_2 R_3 + R_1 \dot{R}_2 R_3 + R_1 R_2 \dot{R}_3) \vec{a} \\ &= \vec{\omega}'_1 \times R_1 R_2 R_3 \vec{a} + R_1 (\vec{\omega}'_2 \times R_2 R_3 \vec{a}) + R_1 R_2 (\vec{\omega}'_3 \times R_3 \vec{a}) \\ &= (\vec{\omega}'_1 + R_1 \vec{\omega}'_2 + R_1 R_2 \vec{\omega}'_3) \times R\vec{a}. \end{aligned}$$

Sledi

$$\begin{aligned} \vec{\omega}' &= \vec{\omega}'_1 + R_1 \vec{\omega}'_2 + R_1 R_2 \vec{\omega}'_3 \\ &= \dot{\varphi}\vec{k} + R_1 \dot{\theta}\vec{i} + R_1 R_2 \dot{\psi}\vec{k} \\ &= \dot{\varphi}\vec{k} + \dot{\theta}\vec{i}' + \dot{\psi}\vec{k}'' . \end{aligned}$$

□

Vprašanje 47. Kako izraziš vektor kotne hitrost Eulerjeve rotacije? Dokaži.

Če razpišemo \vec{k}'', \vec{i}' in \vec{j}' , dobimo

$$\vec{\omega}' = (\dot{\theta} \cos \varphi + \dot{\psi} \sin \varphi \sin \theta) \vec{i} + (\dot{\theta} \sin \varphi - \dot{\psi} \cos \varphi \sin \theta) \vec{j} + (\dot{\varphi} + \dot{\psi} \cos \theta) \vec{k},$$

po drugi strani pa je $\vec{\omega}' = Q\vec{\omega}$. Za rotacijo Q velja

$$\vec{\omega}(Q) \times \vec{a} = Q^T \dot{Q} \vec{a} = -\dot{Q}^T Q \vec{a} = -\vec{\omega}'(Q^T) \times Q^T Q \vec{a} = -\vec{\omega}'(Q^T) \times \vec{a},$$

2 Mehanika

torej $\vec{\omega}(Q) = -\vec{\omega}'(Q^T)$. Sledi

$$\vec{\omega} = (\dot{\theta} \cos \psi + \dot{\varphi} \sin \theta \sin \psi) \vec{i} + (-\dot{\theta} \sin \psi + \dot{\varphi} \sin \theta \cos \psi) \vec{j} + (\dot{\psi} + \dot{\varphi} \cos \theta) \vec{k}$$

oziroma

$$\vec{\omega}' = (\dot{\theta} \cos \psi + \dot{\varphi} \sin \theta \sin \psi) \vec{\varepsilon}_1 + (-\dot{\theta} \sin \psi + \dot{\varphi} \sin \theta \cos \psi) \vec{\varepsilon}_2 + (\dot{\psi} + \dot{\varphi} \cos \theta) \vec{\varepsilon}_3,$$

kjer je $\vec{\varepsilon}_1 = \vec{i}'''$ in podobno za drugi komponenti.

Vprašanje 48. Izpelji predpis za $\vec{\omega}'$ v rotirani bazi.

Trditev. Ne obstaja parametrizacija $SO(3)$, da je $\vec{\omega}' = \dot{\vec{x}}$.

Dokaz. Če bi obstajala taka parametrizacija $\vec{x} = \vec{x}(\vec{y}) = \vec{x}(\varphi, \theta, \psi)$, potem velja

$$\vec{\omega}' = \dot{\vec{x}} = \frac{\partial \vec{x}}{\partial \vec{y}} \dot{\vec{y}} = \begin{bmatrix} \sin \theta \sin \psi & \cos \psi & 0 \\ \sin \theta \cos \psi & -\sin \psi & 0 \\ \cos \theta & 0 & 1 \end{bmatrix} \dot{\vec{y}}.$$

Računamo lahko

$$\frac{\partial^2 x_1}{\partial y_1 \partial y_2} = \cos \theta \sin \psi \qquad \frac{\partial^2 x_1}{\partial y_2 \partial y_1} = 0.$$

To ne mora biti res. □

3 Uvod v numerične metode

3.1 Računske napake

Kadar z numerično metodo nekaj izračunamo, ne dobimo točne vrednosti, vendar nek približek. ABSOLUTNO NAPAKO definiramo kot razliko med približkom in točno vrednostjo:

$$d_a = \hat{x} - x.$$

Po drugi strani je RELATIVNA NAPAKA kvocient

$$d_r = \frac{\hat{x} - x}{x}.$$

Približek lahko izrazimo kot $\hat{x} = x(1 + d_r)$.

Vprašanje 1. Definiraj absolutno in relativno napako.

Števila predstavljamo s plavajočo vejico, ki je pravzaprav eksponentni zapis

$$x = \pm m \cdot b^e,$$

kjer je m MANTISA, zapisana kot $m = 0.c_1c_2 \dots c_t$ za $c_i \in \{0, \dots, b-1\}$, število b je BAZA zapisa, e pa EKSPONENT v mejah $L \leq e \leq U$. Števila običajno zapišemo NORMALIZIRANA, torej s $c_1 \neq 0$. V primeru najnižje možne potence dovoljujemo tudi SUBNORMALIZIRANA števila, kjer je $c_1 = 0$. Predstavljiva števila v takšnem zapisu označujemo s $P(b, t, L, U)$.

V standardu IEEE imamo dve števili:

- Enojni zapis: $P(2, 24, -125, 128)$,
- Dvojni zapis: $P(2, 53, -1021, 1023)$.

Vprašanje 2. Kaj je $P(b, t, L, U)$? Kakšne vrednosti imata `float` in `double`?

Pri zaokroževanju številu odrežemo decimalke za neko vrednostjo, in po potrebi prištejemo b^{-t} . Boljšega od teh približkov označimo s $\text{fl}(x)$.

Izrek. Če za x velja, da $|x|$ leži na intervalu med najmanjšim in največjim pozitivnim predstavljivim normaliziranim številom, potem velja

$$\frac{|\text{fl}(x) - x|}{|x|} \leq u,$$

za OSNOVNO ZAOKROŽITVENO NAPAKO $u = \frac{1}{2}b^{1-t}$.

Vprašanje 3. Kaj je osnovna zaokrožitvena napaka? Povej izrek.

Standard IEEE zagotavlja omejeno napako tudi pri osnovnih operacijah:

- $\text{fl}(x \oplus y) = (x \oplus y)(1 + \delta)$ za $|\delta| \leq u$ za osnovne operacije $+$, $-$, \cdot , $/$,
- $\text{fl}(\sqrt{x}) = \sqrt{x}(1 + \delta)$ za $|\delta| \leq u$.

Drug vir napak je občutljivost problema, ki ni povezana z numeriko. Obravnavamo vprašanje, kako se pri majhni spremembi v vhodnih podatkih spremeni pravilni odgovor. Za zvezno odvedljivo f lahko absolutno občutljivost merimo z odvodom

$$|f(x + \delta x) - f(x)| \approx |f'(x)| |\delta x|.$$

Poznamo tri vrste napak, ki skupaj sestavljajo celotno napako:

- Neodstranljiva napaka: napaka zaradi zaokroževanja podatkov
- Napaka metode: nenatančnost metode
- Zaokrožitvena napaka: napaka zaradi zaokroževanja znotraj metode

Vprašanje 4. Katere vrste napak poznamo?

3.2 Nelinearne enačbe

Iščemo rešitve enačbe $f(x) = 0$ za nek $f : \mathbb{R} \rightarrow \mathbb{R}$. Pri tem lahko pridemo do več različnih situacij glede obstoja in enoličnosti rešitve; možnost je, da rešitev obstaja in je ena sama, da je rešitev več, ampak končno mnogo, da jih je neskončno mnogo, ali pa da rešitve sploh ni.

Naj bo α ničla za zvezno odvedljivo f . Ničla je enostavna natanko tedaj, ko je $f'(\alpha) \neq 0$. Če ni enostavna, je m -kratna natanko tedaj, ko je prvi neničelni odvod v α reda m . V primeru enostavne ničle lokalno obstaja inverzna funkcija, da je $\alpha = f^{-1}(0)$. Absolutna občutljivost problema je tedaj enaka $|(f^{-1})'(0)| = \frac{1}{|f'(\alpha)|}$. Če je α dvojna ničla, uporabimo Taylorjev približek druge stopnje

$$f(x) \approx \underbrace{f(\alpha)}_{=0} + \underbrace{f'(\alpha)}_{=0}(x - \alpha) + \frac{1}{2}f''(\alpha)(x - \alpha)^2,$$

torej za $|f(x)| \leq \varepsilon$ velja

$$|x - \alpha| \leq \sqrt{\frac{2\varepsilon}{|f''(\alpha)|}}.$$

Višje kot so ničle, bolj občutljiv je problem iskanja. Za večkratno ničlo v splošnem velja

$$|x - \alpha| \leq \left(\frac{m!\varepsilon}{|f^{(m)}(\alpha)|} \right)^{1/m}$$

Vprašanje 5. Analiziraj problem iskanja ničel.

3.2.1 Bisekcija

Pri implementaciji bisekcije si hranimo zaporedja $(a_n)_n$ levih mej, $(b_n)_n$ desnih mej, $(c_n)_n$ sredinskih približkov, in $(e_n)_n$ polovičnih velikosti intervala. Nove člene izračunamo po predpisih

$$\begin{aligned}e_{n+1} &= e_n/2, \\a_n &= a_{n-1} \text{ ali } c_{n-1}, \\b_n &= b_{n-1} \text{ ali } c_{n-1}, \\c_n &= a_n + e_n.\end{aligned}$$

Pri tem zmanjšamo število računskih operacij, se izognemo problemom glede možnih nepredvidenih zaokrožitev, skokov izven območja ali računskih napak. Bisekcija nam lahko poišče liho ničlo, ne pa sode, prav tako lahko poišče lih pol (ne pa sodega).

Vprašanje 6. Opiši delovanje bisekcije.

3.2.2 Navadna iteracija

Pri navadni iteraciji iskanje rešitve enačbe $f(x) = 0$ prevedemo na iskanje rešitve enačbe $x = g(x)$ za ustrezno izbrano funkcijo g . Splošna primerna izbira je recimo $g(x) = x - f(x)$, ali pa $g(x) = x - h(x)f(x)$ za neničelno funkcijo h . Da postopek $x_{r+1} = g(x_r)$ konvergira, mora biti g v okolici α skrčitev.

Izrek. Naj bo $\alpha = g(\alpha)$ in naj g na intervalu $I = [\alpha - \delta, \alpha + \delta]$ za nek $\delta > 0$ zadošča Lipschitzovem pogoju $|g(x) - g(y)| \leq m|x - y|$ za nek $m \in [0, 1)$, in poljubna $x, y \in I$. Potem za vsak $x_0 \in I$ zaporedje $x_{r+1} = g(x_r)$ konvergira k α , in velja

- $|x_r - \alpha| \leq m^r |x_0 - \alpha|$,
- $|x_{r+1} - \alpha| \leq \frac{m}{1-m} |x_{r+1} - x_r|$.

Dokaz. Dokažimo prvo, da zaporedje ne zapusti intervala I ; če velja $x_r \in I$, je $|x_r - \alpha| \leq \delta$, torej

$$|x_{r+1} - \alpha| = |g(x_r) - g(\alpha)| \leq m|x_r - \alpha| < |x_r - \alpha| \leq \delta,$$

torej je tudi $x_{r+1} \in I$. S ponavljanjem tega postopka tudi dokažemo prvo točko, za drugo točko pa ocenimo

$$\begin{aligned}|x_{r+1} - \alpha| &= |x_{r+1} - x_{r+2} + x_{r+2} - x_{r+3} + x_{r+3} - \dots - \alpha| \\&\leq |x_{r+1} - x_{r+2}| + |x_{r+2} - x_{r+3}| + \dots \\&\leq (m + m^2 + m^3 + \dots) |x_{r+1} - x_r| \\&= \frac{m}{1-m} |x_{r+1} - x_r|.\end{aligned}$$

□

Posledica. Če je $\alpha = g(\alpha)$, če je g zvezno odvedljiva in če velja $|g'(\alpha)| < 1$, potem obstaja nek $\delta > 0$, da za vsak x_0 , $|x_0 - \alpha| < \delta$, zaporedje $x_{r+1} = g(x_r)$ konvergira k α .

Vprašanje 7. Povej in dokaži izrek o navadni iteraciji.

Definicija. Naj bo $\lim x_r = \alpha$. Pravimo, da $(x_r)_r$ KONVERGIRA K α Z REDOM p , če velja

$$\lim_{r \rightarrow \infty} \frac{|x_{r+1} - \alpha|}{|x_r - \alpha|^p} = C$$

za neko konstanto $C > 0$.

Izrek. Naj bo $\alpha = g(\alpha)$ za p -krat zvezno odvedljivo funkcijo g in naj velja $g'(\alpha) = \dots = g^{(p-1)}(\alpha) = 0$ ter $g^{(p)}(\alpha) \neq 0$. Tedaj v bližini α zaporedje $x_{r+1} = g(x_r)$ konvergira k α z redom p .

Dokaz. Izraz $x_{r+1} = g(x_r)$ v okolici α razvijemo v Taylorjevo vrsto:

$$x_{r+1} = g(\alpha) + g'(\alpha)(x_r - \alpha) + \dots + \frac{g^{(p-1)}(\alpha)}{(p-1)!}(x_r - \alpha)^{p-1} + \frac{g^{(p)}(\xi)}{p!}(x_r - \alpha)^p$$

za nek ξ v bližini α . Sledi

$$\frac{x_{r+1} - \alpha}{(x_r - \alpha)^p} = \frac{g^{(p)}(\xi)}{p!}.$$

□

Vprašanje 8. Kaj je red konvergence zaporedja? Kako ga poiščeš z odvodom?

3.2.3 Tangentna metoda

Če vzamemo $\alpha = x_r + \Delta x_r$, in razvijemo dva člena Taylorjeve vrste

$$0 = f(x_r + \Delta x_r) = f(x_r) + f'(x_r)\Delta x_r + \frac{1}{2}f''(\xi_r)\Delta x_r^2,$$

ter nato zanemarimo zadnji člen, dobimo

$$\Delta x_r = \frac{-f(x_r)}{f'(x_r)},$$

s čimer smo izpeljali tangentno metodo, ki je pravzaprav poseben primer naravne iteracije. Če je α enostavna ničla, lahko izračunamo, da ima $g(x) = x - f(x)/f'(x)$ ničelni odvod v α (ob predpostavki $f \in \mathcal{C}^2$), in je torej konvergenca vsaj kvadratična. V primeru m -kratne ničle za $m \geq 2$ pa po dolgotrajni izpeljavi dobimo, da je konvergenca zagotovljena, a linearna. Za dvakrat zvezno odvedljive funkcije je vsaka ničla f torej privlačna negibna točka.

Vprašanje 9. Izpelji tangentno metodo in pokaži, kakšen red konvergence ima.

3 Uvod v numerične metode

Izrek. Naj bo f na $I = [0, \infty)$ dvakrat zvezno odvedljiva, strogo naraščajoča, konveksna in naj ima na I ničlo α . Potem za vsak $x_0 \in I$ tangenta metoda konvergira k α .

Dokaz. Velja $f'(x) > 0$ in $f''(x) > 0$. S Taylorjevim razvojem $f(\alpha)$ okoli točke x_0 dobimo oceno

$$x_1 - \alpha = \frac{f''(\xi)}{2f'(x_0)}(x_0 - \alpha)^2 \geq 0$$

pokažemo, da je ne glede na x_0 točka x_1 vedno desno od α . Dokažimo še, da je x_{r+1} nujno med α in x_r :

$$x_{r+1} = x_r - \frac{f(x_r)}{f'(x_r)} < x_r,$$

vedno pa velja $x_r > \alpha$. To je torej strogo padajoče omejeno zaporedje, ki mora nekam konvergirati; to bo seveda α . \square

Vprašanje 10. Za kakšne funkcije lahko globalno zagotoviš konvergenco tangentne metode? Dokaži.

3.2.4 Sekantna metoda

Če je izračun odvoda zahteven, ga lahko aproksimiramo z diferenčnim kvocientom. Namesto tangente na f v točki x_r tako uporabimo sekanto skozi točki $(x_r, f(x_r))$ in $(x_{r-1}, f(x_{r-1}))$. Dobljena metoda tehnično ni navadna iteracija, ker uporablja zadnja dva približka, vendar se obnaša sorodno. Naslednji približek izračunamo s predpisom

$$x_{r+1} = x_r - \frac{f(x_r)(x_r - x_{r-1})}{f(x_r) - f(x_{r-1})}.$$

Analiza sekantne metode je težja kot analiza metod navadne iteracije. Izkaže se, da velja

$$|e_{r+1}| \approx c |e_r| |e_{r-1}|$$

za neko konstanto c , ki je pravzaprav enaka

$$c = \frac{|f''(\alpha)|}{2|f'(\alpha)|}.$$

Označimo red sekantne metode s p . Obstaja konstanta $D > 0$, da je $|e_{r+1}| \approx D |e_r|^p$, torej

$$|e_{r+1}| \approx CD |e_{r-1}|^{p+1} = D^{p+1} |e_{r-1}|^{p^2},$$

iz česar sledi $p^2 = p + 1$ oziroma $p = \phi$, torej je konvergenca superlinearna.

Vprašanje 11. Razloži sekantno metodo in izpelji njen red konvergence.

3.2.5 Ostale metode

Pri *Mullerjevi metodi* uporabimo tri približke x_r , x_{r-1} in x_{r-2} , ter skozi točke $(x_r, f(x_r))$, $(x_{r-1}, f(x_{r-1}))$, $(x_{r-2}, f(x_{r-2}))$ potegnemo polinom stopnje 2. Za naslednji približek vzamemo tisto izmed dveh ničel polinoma, ki je bližnja x_r . Ena od prednosti te metode je, da lahko išče kompleksne ničle tudi z realnimi začetnimi približki. Izkaže se, da je red konvergence približno 1.84.

Vprašanje 12. Razloži Mullerjevo metodo.

Če zamenjamo vlogi x in y , in najdemo polinom $p(y)$, ki poteka skozi točke $(f(x_r), x_r)$, $(f(x_{r-1}), x_{r-1})$, $(f(x_{r-2}), x_{r-2})$, dobimo približek za f^{-1} , in lahko za naslednji približek vzamemo $x_{r+1} = p(0)$. Metodo imenujemo *inverzna interpolacija*, ima pa isti red konvergence kot Mullerjeva metoda.

Vprašanje 13. Razloži metodo inverzne interpolacije.

3.2.6 Ničle polinomov

Ničle polinomov lahko iščemo na več načinov:

- Poiščeš eno ničlo in reduciraš polinom.
- Računaš vse ničle hkrati.
- Prevedeš problem na problem iskanja lastnih vrednosti.

Za redukcijo na problem lastnih vrednosti uporabljamo *spremljevalno matriko* polinoma $p(x) = a_0x^n + \dots + a_n$

$$C_p = \begin{bmatrix} 0 & 1 & & & \\ & 0 & 1 & & \\ & & \ddots & \ddots & \\ & & & 0 & 1 \\ -\frac{a_n}{a_0} & -\frac{a_{n-1}}{a_0} & \dots & -\frac{a_2}{a_0} & -\frac{a_1}{a_0} \end{bmatrix}$$

Vprašanje 14. Kako izgleda spremljevalna matrika polinoma?

Ena od metod, ki računa vse ničle hkrati, je *Laguerrova metoda*. Za polinom $p(x) =$

3 Uvod v numerične metode

$a_0x^n + \dots + a_n$ z ničlami $\alpha_1, \dots, \alpha_n$ definiramo

$$\begin{aligned} S_1(x) &= \sum_{i=1}^n \frac{1}{x - \alpha_i} = \frac{p'(x)}{p(x)}, \\ S_2(x) &= \sum_{i=1}^n \frac{1}{(x - \alpha_i)^2} = -S_1'(x) = \frac{(p'(x))^2 - p(x)p''(x)}{p^2(x)}, \\ a(x) &= \frac{1}{x - \alpha_n}, \\ b(x) &= \frac{1}{n-1} \sum_{i=1}^{n-1} \frac{1}{x - \alpha_i}, \end{aligned}$$

da velja $S_1(x) = a(x) + (n-1)b(x)$. Tedaž za

$$\begin{aligned} d_i(x) &= \frac{1}{x - \alpha_i} - b(x), \\ d(x) &= \sum_{i=1}^{n-1} d_i^2(x) \end{aligned}$$

dobimo $S_2 = a^2 + (n-1)b^2 + d$, ker je $\sum_i d_i = 0$. Dobili smo sistem enačb v spremenljivkah a, b , ki ga lahko rešimo in dobimo

$$a_{1,2} = \frac{1}{n} \left(S_1 \pm \sqrt{(n-1)(nS_2 - S_1^2 - nd)} \right).$$

Če x obravnavamo kot približek za ničlo α_n , bo člen nd v bližini α_n majhen, zato ga zanemarimo. Iz tega izrazimo

$$\alpha_n = x - \frac{n}{S_1 \pm \sqrt{(n-1)(nS_2 - S_1^2)}}.$$

Laguerrova metoda nam torej da postopek za izračun približka ničle

$$x_{r+1} = x_r - \frac{np(x_r)}{p'(x_r) \pm \sqrt{(n-1)[(n-1)(p'(x_r))^2 - np(x_r)p''(x_r)]}}.$$

Za odločitev, kateri predznak pripišemo korenu v imenovalcu, imamo tri možnosti:

- Vedno izberemo plus,
- Vedno izberemo minus,
- *Stabilna varianta*: izbereš tistega, ki ti v imenovalcu da večjo absolutno vrednost.

V prvih dveh primerih fiksno iščemo v eni smeri od začetnega približka, v tretjem pa to ni zagotovljeno.

Izrek. Če ima polinom p same realne ničle, potem za vsak začetni približek x_0 stabilna verzija Laguerrove metode konvergira proti najbližji desni oz. levi ničli, pri čemer si mislimo, da sta kraka realne osi pri $+\infty$ in $-\infty$ združena. Konvergenca v bližini enostavne ničle je kubična.

Če ima polinom kompleksne ničle, metoda konvergira za skoraj vse začetne približke.

Vprašanje 15. Izpelji Laguerrovo metodo in razloži, kako delujejo vse možnosti za izbiro naslednjega približka.

3.2.7 Durand-Kernerjeva metoda

Izberimo približke x_1, \dots, x_n za ničle polinoma $p(x)$ z vodilnih koeficientom 1. Iščemo popravke $\Delta x_1, \dots, \Delta x_n$, da bodo $x_i + \Delta x_i$ točne ničle. Velja

$$\begin{aligned} p(x) &= (x - (x_1 + \Delta x_1))(x - (x_2 + \Delta x_2)) \dots (x - (x_n + \Delta x_n)) \\ &= \prod_{i=1}^n (x - x_i) - \sum_{j=1}^n \Delta x_j \prod_{i \neq j} (x - x_i) + \dots, \end{aligned}$$

člene drugega in večjega reda pa zanemarimo (torej vse, kar je v tropičju). Če s $q(x)$ označimo nezanemaren del, velja

$$q(x_l) = -\Delta x_l \prod_{i \neq l} (x_l - x_i),$$

iz česar lahko izračunamo Δx_l .

Vprašanje 16. Razloži Durand-Kernerjevo metodo.

3.3 Sistemi linearnih enačb

3.3.1 Matrične norme

Definicija. Preslikava $\|\cdot\| : \mathbb{C}^n \rightarrow \mathbb{R}$ je VEKTORSKA NORMA, če velja

- $\|x\| \geq 0$ za vse x , in $\|x\| = 0 \Leftrightarrow x = 0$,
- $\|\alpha x\| = |\alpha| \|x\|$,
- $\|x + y\| \leq \|x\| + \|y\|$.

Vse vektorske norme so ekvivalentne, za poljubni normi $\|x\|_A$ in $\|x\|_B$ obstajata konstanti C_1, C_2 , da velja

$$C_1 \|x\|_A \leq \|x\|_B \leq C_2 \|x\|_A.$$

3 Uvod v numerične metode

Konkretno za 1-normo, 2-normo in supremum normo veljajo ocene

$$\begin{aligned}\|x\|_2 &\leq \|x\|_1 \leq \sqrt{n} \|x\|_2 \\ \|x\|_\infty &\leq \|x\|_1 \leq n \|x\|_\infty \\ \|x\|_\infty &\leq \|x\|_2 \leq \sqrt{n} \|x\|_\infty\end{aligned}$$

Vprašanje 17. Definiraj vektorsko normo. V kakšnem razmerju so znane vektorske norme?

Definicija. Preslikava $\|\cdot\| : \mathbb{C}^{n \times n} \rightarrow \mathbb{R}$ je MATRIČNA NORMA, če velja

- $\|A\| \geq 0$ in $\|A\| = 0 \Leftrightarrow A = 0$,
- $\|\alpha A\| = |\alpha| \|A\|$,
- $\|A + B\| \leq \|A\| + \|B\|$,
- $\|AB\| \leq \|A\| \|B\|$ (submultiplikativnost).

Matrika je tudi vektor, na njej so tudi definirane običajne vektorske norme. Definiramo funkcije

$$\begin{aligned}N_1(A) &= \sum_{i,j} |a_{ij}|, \\ N_2(A) &= \sqrt{\sum_{i,j} |a_{ij}|^2}, \\ N_\infty(A) &= \max_{i,j} |a_{ij}|.\end{aligned}$$

Te funkcije ustrezajo prvim trem točkam definicije matrične norme, to pa še ni dovolj. Za matriki

$$A = B = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

velja $N_\infty(A) = N_\infty(B) = 1$, ampak $N_\infty(AB) = 2$, torej N_∞ ni matrična norma. V nasprotju N_1 in N_2 dejansko sta matrični normi. Funkcijo N_2 imenujemo FROBENIUSOVA NORMA in označimo z $N_2(A) = \|A\|_F$.

Vprašanje 18. Dokaži, da N_∞ ni matrična norma. Kaj je Frobeniusova norma?

Izrek. Naj bo $\|\cdot\|_v$ vektorska norma na \mathbb{C}^n . Potem je

$$\|A\|_m = \max_{x \neq 0} \frac{\|Ax\|_v}{\|x\|_v}$$

matrična norma. Taki normi pravimo OPERATORSKA NORMA.

Dokaz. Prve tri točke so očitne, preverimo samo submultiplikativnost. Za vsak $x \neq 0$ velja $\|Ax\|_v \leq \|A\|_m \|x\|_v$ po definiciji, torej

$$\|AB\| = \max_{x \neq 0} \frac{\|ABx\|_v}{\|x\|_v} \leq \max_{x \neq 0} \frac{\|A\|_m \|Bx\|_v}{\|x\|_v} = \|A\|_m \|B\|_m.$$

□

Vprašanje 19. Dokaži, da so operatorske norme res matrične norme.

Za matrično normo $\|\cdot\|_m$ in vektorsko normo $\|\cdot\|_v$ pravimo, da sta USKLAJENI, če za vsako matriko A in vektor x velja

$$\|Ax\|_v \leq \|A\|_m \|x\|_v.$$

Lema. Za vsako matrično normo obstaja usklajena vektorska norma.

Dokaz. Za vektor x definiramo

$$\|x\|_v = \left\| \begin{bmatrix} x & 0 & \dots & 0 \end{bmatrix} \right\|_m,$$

kjer smo vektor dopolnili do kvadratne matrike. To je očitno vektorska norma, normi sta očitno usklajeni. □

Vprašanje 20. Dokaži, da ima vsaka matrična norma usklajeno vektorsko normo.

Posledica. Za vsako matrično normo in poljubno lastno vrednost λ matrike A velja $|\lambda| \leq \|A\|$.

Dokaz. Naj bo $Ax = \lambda x$ za nek $x \neq 0$. Velja

$$|\lambda| \|x\|_v = \|\lambda x\|_v = \|Ax\|_v \leq \|A\| \|x\|_v.$$

□

Vprašanje 21. Kakšna je povezava med lastnimi vrednostmi in matrično normo? Dokaži.

Lema. Norma $\|A\|_1$ je enaka največji 1-normi stolpca matrike A .

Dokaz. Naj bo x vektor. Velja

$$\|Ax\|_1 = \left\| \sum_i x_i a_i \right\|_1 \leq \sum_i |x_i| \|a_i\|_1 \leq \max_{j=1,\dots,n} \sum_i |x_i| \|a_j\|_1 = \max_j \|x\|_1 \|a_j\|_1.$$

Sledi

$$\frac{\|Ax\|_1}{\|x\|_1} \leq \max_j \|a_j\|_1.$$

Enakost dobimo, če za x vzamemo e_k , kjer je k stolpec z največjo 1-normo. □

3 Uvod v numerične metode

Podobno pokažemo, da je $\|A\|_\infty$ enaka največji 1-normi vrstice.

Vprašanje 22. Čemu sta enaki normi $\|A\|_1$ in $\|A\|_\infty$? Dokaži za eno.

Lema. Velja $\|A\|_2 = \max_j \sqrt{\lambda_j}$, kjer so λ_j lastne vrednosti matrike $A^H A$.

Dokaz. Matrika $A^H A$ je hermitska, torej so vse njene lastne vrednosti realne. Velja

$$x^H A^H A x = (Ax)^H Ax = \|Ax\|_2^2 \geq 0,$$

torej je $A^H A$ pozitivno semidefinitna, in ima nenegativne lastne vrednosti. Izraz je torej res dobro definiran.

Naj bodo $\sigma_1^2 \leq \sigma_2^2 \leq \dots \leq \sigma_n^2$ singularne vrednosti matrike A (lastne vrednosti $A^H A$). Ker je $A^H A$ hermitska, lahko lastne vektorje izberemo tako, da so ortonormirani. Naj za v_i torej velja $A^H A v_i = \sigma_i^2 v_i$. Za poljuben x velja

$$x = \sum_i \alpha_i v_i,$$

torej

$$\|Ax\|_2^2 = (Ax)^H (Ax) = x^H A^H A x = \sum_i |\alpha_i|^2 \sigma_i^2 \leq \sigma_n^2 \sum_i |\alpha_i|^2 = \sigma_n^2 \|x\|_2^2.$$

Sledi neenakost

$$\frac{\|Ax\|_2}{\|x\|_2} \leq \sigma_n,$$

kjer dobimo enačaj, če vzamemo $x = v_n$. □

Vprašanje 23. Čemu je enaka 2-norma matrike? Dokaži.

Frobeniusova norma ni operatorska, ker za vse operatorske norme velja $\|I\| = 1$, ampak $\|I\|_F = \sqrt{n}$. Kljub temu pa so vse matrične norme ekvivalentne. Za konkretne primere velja

$$\begin{aligned} \frac{1}{\sqrt{n}} \|A\|_F &\leq \|A\|_2 \leq \|A\|_F, \\ \frac{1}{\sqrt{n}} \|A\|_1 &\leq \|A\|_2 \leq \sqrt{n} \|A\|_1, \\ \frac{1}{\sqrt{n}} \|A\|_\infty &\leq \|A\|_2 \leq \sqrt{n} \|A\|_\infty. \end{aligned}$$

Velja tudi $\|A\|_2 \leq \sqrt{\|A\|_1 \|A\|_\infty}$ in $N_\infty(A) \leq \|A\|_2 \leq n N_\infty(A)$.

Vprašanje 24. Kako oceniš 2-normo matrike?

Lema. Normi $\|\cdot\|_2$ in $\|\cdot\|_F$ sta invariantni na množenje z unitarno matriko.

Dokaz. Za $x \in \mathbb{C}$ velja $\|Ux\|_2 = \|x\|_2$. Pri 2-normi torej

$$\|UA\|_2 = \max_x \frac{\|UAx\|_2}{\|x\|_2} = \max_x \frac{\|Ax\|_2}{\|x\|_2} = \|A\|_2,$$

za Frobeniusovo normo po

$$\|UA\|_F^2 = \|U[a_1 \dots a_n]\|_F^2 = \sum_i \|Ua_i\|_2^2 = \sum_i \|a_i\|_2^2 = \|A\|_F^2.$$

V drugo smer uporabimo dejstvo $\|A^H\|_2 = \|A\|_2$ in $\|A^H\|_F = \|A\|_F$. \square

Vprašanje 25. Dokaži, da sta 2-norma in Frobeniusova norma invariantni na množenje z unitarno matriko.

Lema. Naj bo $\|X\| < 1$. Potem je $I - X$ nesingularna, inverz je enak

$$(I - X)^{-1} = \sum_{k=0}^{\infty} X^k,$$

in če je $\|I\| = 1$, velja

$$\|(I - X)^{-1}\| \leq \frac{1}{1 - \|X\|}.$$

Dokaz. Recimo, da je $I - X$ singularna. Tedaj obstaja vektor $w \neq 0$, da je $(I - X)w = 0$, torej $Xw = w$ in je w lastni vektor za lastno vrednost 1. To je protislovje, ker mora biti norma večja od lastne vrednosti.

Računamo

$$(I - X) \sum_{k=0}^m X^k = I - X^{m+1} \xrightarrow{m \rightarrow \infty} I,$$

ker je $\|X\| < 1$ in $\|X\|^{m+1} \geq \|X^{m+1}\|$. Dodatno velja

$$\left\| \sum_{k=0}^{\infty} X^k \right\| \leq \|I\| + \|X\| + \|X\|^2 + \dots = \frac{1}{1 - \|X\|}.$$

\square

Vprašanje 26. Povej predpostavke in dokaži formulo za inverz matrike $I - X$.

3.3.2 Občutljivost sistema linearnih enačb

Imejmo sistem $Ax = b$, kjer je A nesingularna matrika. Denimo, da A in b zmotimo v $A + \Delta A$ in $b + \Delta b$, kjer je $A + \Delta A$ še vedno nesingularna. Nov sistem ima potem obliko

$$(A + \Delta A)(x + \Delta x) = b + \Delta b.$$

3 Uvod v numerične metode

Lema. Če je A nesingularna in $\|\Delta A\| < \frac{1}{\|A^{-1}\|}$, je $A + \Delta A$ nesingularna.

Dokaz. Računamo

$$A + \Delta A = A(I + A^{-1}\Delta A),$$

in ocenimo normo

$$\|A^{-1}\Delta A\| \leq \|A^{-1}\| \|\Delta A\| < 1,$$

torej velja po prejšnji lemi. □

Naj torej velja ta pogoj za eno izmed znanih operatorskih norm. Računamo

$$\begin{aligned} (A + \Delta A)\Delta x &= \Delta b - \Delta Ax, \\ \Delta x &= (I + A^{-1}\Delta A)^{-1}(\Delta b - \Delta Ax), \\ \frac{\|\Delta x\|}{\|x\|} &\leq \frac{\|A\| \|A^{-1}\|}{1 - \|A^{-1}\| \|\Delta A\|} \left(\frac{\|\Delta b\|}{\|b\|} + \frac{\|\Delta A\|}{\|A\|} \right), \end{aligned}$$

kjer smo v zadnjem delu uporabili sklep $\|b\| \leq \|A\| \|x\|$. Kvaliteta ocene je odvisna od vrednosti

$$\kappa(A) = \|A^{-1}\| \|A\|,$$

ki ji pravimo OBČUTLJIVOST MATRIKE ali POGOJENOSTNO ŠTEVILO.

Za 2-normo velja

$$\|A^{-1}\|_2 = \max_{x \neq 0} \frac{\|A^{-1}x\|_2}{\|x\|_2} = \max_{y \neq 0} \frac{\|y\|_2}{\|Ay\|_2} = \left(\min_{y \neq 0} \frac{\|Ay\|_2}{\|y\|_2} \right)^{-1} = \frac{1}{\sigma_n(A)},$$

kjer je $\sigma_n(A)$ najmanjša singularna vrednost. Dobimo torej

$$\kappa_2(A) = \frac{\sigma_1(A)}{\sigma_n(A)}.$$

Vprašanje 27. Kaj je občutljivost matrike? Kako jo uporabiš za oceno občutljivosti sistema linearnih enačb?

3.3.3 LU razcep

Naj bo dan vektor $w \in \mathbb{R}^n$ in naj velja $w_k \neq 0$. Definiramo ELIMINACIJSKO MATRIKO

$$L_k = \begin{bmatrix} 1 & & & & & \\ & 1 & & & & \\ & & \ddots & & & \\ & & & 1 & & \\ & & & -l_{k+1,k} & 1 & \\ & & & \vdots & & \ddots \\ & & & -l_{n,k} & & 1 \end{bmatrix}$$

za $l_{i,k} = \frac{w_i}{w_k}$. Velja

$$L_k w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_k \\ 0 \\ \vdots \\ 0 \end{bmatrix}.$$

Če je $L_k = I - l_k e_k^T$, lahko izračunamo $L_k^{-1} = I + l_k e_k^T$. Če množimo matriko A z leve z L_1 , uničimo prvi stolpec, razen prvega elementa. Če to matriko množimo z L_2 z leve, uničimo poddiagonalne elemente v drugem stolpcu (L_2 gradimo iz elementom te druge matrike). Tako lahko nadaljujemo in pridemo do

$$L_{n-1} L_{n-2} \dots L_2 L_1 A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ 0 & a_{22}^{(1)} & \dots & a_{2n}^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & a_{nn}^{(n-1)} \end{bmatrix} = U$$

Za matriko $L = L_1^{-1} L_2^{-1} \dots L_{n-1}^{-1}$ tedaj velja $A = LU$. Izračunamo lahko

$$L = I + \sum_k l_k^T e_k,$$

iz česar vidimo, da je L spodnje trikotna z enicami na diagonalni. Da LU razcep lahko naredimo, morajo biti elementi $a_{11}, a_{22}^{(1)}, \dots, a_{nn}^{(n-1)}$ na diagonalni neničelni. Pravimo jim PIVOTI.

Vprašanje 28. Izpelji in razloži osnovni LU razcep. Kaj so pivoti?

Algorithm 1 LU razcep brez pivotiranja

```

for  $j = 1, \dots, n-1$  do
  for  $i = j+1, \dots, n$  do
     $l_{ij} = a_{ij}/a_{jj}$ 
    for  $k = j+1, \dots, n$  do
       $a_{ik} = a_{ik} - l_{ij}a_{jk}$ 
    end for
  end for
end for

```

Osnovni postopek za izračun LU razcepa je prikazan v algoritmu 1. S tem dobimo elemente matrike L , razen enic na diagonalni, ter elemente matrike U nad diagonalno. Algoritem deluje v $O(n^3)$ s prefaktorjem $2/3$.

3 Uvod v numerične metode

Sistem $Ax = b$ rešimo tako, da prvo razcepimo $A = LU$, nato s premo substitucijo rešimo trikotni sistem $Ly = b$, in nazadnje z obratno substitucijo rešimo še $Ux = y$.

Vprašanje 29. Zapiši osnovni algoritem za LU razcep. Kakšno časovno zahtevnost ima?

Izrek. Za matriko A je ekvivalentno

- Obstaja enoličen LU razcep $A = LU$, kjer je L spodnje trikotna z enicami na diagonalah, in U zgornje trikotna nesingularna.
- Vse vodilne podmatrike A so nesingularne.

Dokaz. V desno: A_k je produkt ustreznih vodilnih podmatrik L_k in U_k . V levo: Indukcija na n . Za matrike velikosti 1 ni nič za dokazati. Naj sedaj velja za n , dokažimo da velja tudi za $n + 1$. Definiramo

$$\tilde{A} = \begin{bmatrix} A & b \\ c^T & \delta \end{bmatrix}.$$

Po indukcijski predpostavki lahko matriko A razcepimo v $A = LU$. Določimo lahko torej $u = L^{-1}b$, $l = U^{-T}c$ in $\xi = \delta - l^T u$, da dobimo

$$\tilde{A} = \begin{bmatrix} L & 0 \\ l^T & 1 \end{bmatrix} \cdot \begin{bmatrix} U & u \\ 0^T & \xi \end{bmatrix}.$$

Po predpostavki sta ξ in $\det U$ različni od 0, torej je $\det \tilde{U} = \xi \det U \neq 0$. □

Vprašanje 30. Kdaj je LU razcep enoličen? Dokaži.

Ničle in majhni elementi na diagonalah so problem za LU razcep. Rešimo ga tako, da uvedemo delno oz. kompletno pivotiranje. Pri delnem pivotiranju na vsakem koraku namesto elementa v diagonalah vzamemo element pod diagonalo, ki je največji po absolutni vrednosti, ter ga z menjavo vrstic postavimo na pivotno mesto. Rezultat tega je razcep $PA = LU$, kjer je P permutacijska matrika, ki ustreza menjavi vrstic.

Lema. Če je A nesingularna, obstaja taka permutacijska matrika P , da za PA obstaja LU razcep brez pivotiranja.

Dokaz. Vmesne matrike so oblike $L_{j-1}P_{j-1} \dots L_2P_2L_1P_1A$. Ker so vse našteje matrike nesingularne, mora biti tudi produkt nesingularen, torej obstaja neničelni element v ostanku stolpca. □

Vprašanje 31. Kako deluje LU razcep z delnim pivotiranjem? Pod katerim pogojem ga lahko naredimo?

Pri kompletnem pivotiranju poiščemo največji element v neobdelani podmatriki, in ga postavimo na pivotno mesto z zamenjavo stolpcev in vrstic. Dobimo razcep $PAQ = LU$, kjer P permutira vrstice in Q stolpce. V algoritmu dobimo $O(n^3)$ dodatnih primerjanj.

Vprašanje 32. Opiši LU razcep s kompletnim pivotiranjem.

Sistem $Ax = b$ rešimo numerično z eno izmed variant LU razcepa, in dobimo rešitev \hat{x} . Zanima nas ocena obratne stabilnosti $\|\Delta A\|$, kjer je $(A + \Delta A)\hat{x} = b$. Za analizo si mislimo, da smo pivotiranje že naredili, tako da lahko analiziramo le osnovni LU razcep.

Lema. Naj bo $L \in \mathbb{R}^{n \times n}$ nesingularna spodnje trikotna matrika. Če sistem $Ly = b$ numerično rešimo s premo substitucijo, potem za izračunani \hat{y} velja $(L + \Delta L)\hat{y} = b$, kjer je $|\Delta L| \leq nu|L|$ po elementih.

Dokaz. V i -tem koraku nastavimo

$$\hat{y}_i = \frac{1}{l_{ii}(1 + \alpha_i)(1 + \beta_i)} \left(b_i - \sum_{k=1}^{i-1} l_{ik}\hat{y}_k(1 + \gamma_{ik}) \right),$$

kjer so $\alpha_i, \beta_i, \gamma_{ik}$ manjše od osnovne računske napake u . Za $\gamma_{ii} = (1 + \alpha_i)(1 + \beta_i)$ velja $|\gamma_{ii}| \leq 2u$, in $|\gamma_{ik}| \leq (i - 1)u$, torej za poljubna j, l velja $|\gamma_{jl}| \leq nu$. \square

Lema. Če je U nesingularna zgornje trikotna matrika, in če sistem $Ux = y$ numerično rešimo z obratno substitucijo, izračunani \hat{x} zadošča enačbi $(U + \Delta U)\hat{x} = y$, kjer je $|\Delta U| \leq nu|U|$.

Lema. Naj bo A taka matrika, da se izvede LU razcep brez pivotiranja. Za izračunani matriki \hat{L} in \hat{U} tedaj velja $\hat{L}\hat{U} = A + E$, kjer je $|E| \leq nu|\hat{L}||\hat{U}|$.

Izrek. Če sistem $Ax = b$ rešimo z LU razcepom, potem za izračunani \hat{x} velja $(A + \Delta A)\hat{x} = b$, kjer je $|\Delta A| \leq 3nu|L||U| + O(u^2)$.

Posledica tega je, da velja $\|\Delta A\|_\infty \leq 3nu\|L\|_\infty\|U\|_\infty$, to pa nam ne pomaga zares, ker je lahko produkt norm L in U poljubno velik v primerjavi z normo A .

Če ne pivotiramo, postopek ni obratno stabilen, če pa pivotiramo, pa je $|l_{ij}| \leq 1$, torej je $\|L\|_\infty \leq n$. Če vpeljemo PIVOTNO RAST

$$g = \frac{\max_{ij} |u_{ij}|}{\max_{ij} |a_{ij}|},$$

velja $\|U\|_\infty \leq ng\|A\|_\infty$, in $\|\Delta A\|_\infty \leq 3n^3gu\|A\|_\infty$. Pri delnem pivotiranju velja $g \leq 2^{n-1}$, kar se v splošnem tudi lahko kdaj zgodi, recimo za matriko

$$\begin{bmatrix} 1 & & & & & 1 \\ -1 & 1 & & & & 1 \\ -1 & -1 & 1 & & & 1 \\ \vdots & \vdots & \vdots & \ddots & & \vdots \\ -1 & -1 & -1 & \cdots & 1 & 1 \\ -1 & -1 & -1 & \cdots & -1 & 1 \end{bmatrix},$$

torej LU razcep z delnim pivotiranjem tudi ni obratno stabilen, v praksi pa se to redko zgodi.

Točne ocene za LU razcep s kompletnim pivotiranjem ne poznamo, smatramo pa, da je obratno stabilno.

Vprašanje 33. Analiziraj obratno stabilnost LU razcepa. Kaj je pivotna rast?

3.3.4 Razcep Choleskega

Izrek. Veljajo naslednje točke:

- Če je A simetrična pozitivno definitna matrika, je vsaka vodilna podmatrika simetrično pozitivno definitna.
- Če je A simetrična pozitivno definitna, obstaja enoličen razcep $A = LU$, kjer je L spodnje trikotna matrika z enicami na diagonalni in U zgornje trikotna matrika s pozitivnimi diagonalnimi elementi.
- A je simetrična pozitivno definitna natanko tedaj, ko je $A = VV^T$ za neko spodnje trikotno matriko V , ki ima pozitivne diagonalne elemente.

Dokaz. Prva točka: Velja $x^T A_k x = \tilde{x}^T A \tilde{x}$, kjer je \tilde{x} vektor x , dopolnjen z ničlami.

Druga točka: Vse vodilne podmatrice so simetrične pozitivno definitne, torej nesingularne in obstaja enoličen LU razcep. Če je $A = LU$, je $\det A_k = u_{11}u_{22} \dots u_{kk}$, torej so $u_{ii} > 0$.

Tretja točka: Razcepimo $A = LU$ in to dodatno razcepimo v $A = LDW$, kjer je D diagonalna matrika z elementi u_{11}, \dots, u_{nn} , in $W = D^{-1}U$ zgornje trikotna matrika z enicami na diagonalni.

Ker je $A = A^T$, je $W^T D L^T$ LU razcep matrike A (če združimo desni dve matriki) in velja $W^T = L$ ter $D L^T = U$. Torej je $A = LDL^T$, in lahko definiramo $V = L\sqrt{D}$. \square

Vprašanje 34. Karakteriziraj pozitivno definitnost matrike z razcepom Choleskega in dokaži karakterizacijo.

Če imamo izračunan razcep Choleskega, za $j < k$ velja

$$a_{jk} = \sum_{i=1}^{k-1} v_{ji}v_{ki} + v_{jk}v_{kk},$$

pri $j = k$ pa dobimo

$$a_{kk} = \sum_{i=1}^{k-1} v_{ki}^2 + v_{kk}^2.$$

Vidimo, da če poznamo vse elemente V pred v_{jk} , ga lahko direktno izračunamo. Iz tega dobimo algoritem 2. Ta porabi $\frac{1}{3}n^3$ operacij za račun razcepa, pri čemer pa računamo n korenov.

Algorithm 2 Razcep Choleskega

for $k = 1, \dots, n$ **do**

 Nastavi

$$v_{kk} = \sqrt{a_{kk} - \sum_{i=1}^{k-1} v_{ki}^2}$$

for $j = k + 1, \dots, n$ **do**

 Nastavi

$$v_{jk} = \frac{1}{v_{kk}} \left(a_{jk} - \sum_{i=1}^{k-1} v_{ji} v_{ki} \right)$$

end for

end for

Vprašanje 35. Zapiši algoritem za izračun razcepa Choleskega.

Če rešujemo sistem z razcepom Choleskega, izračunamo rešitev \hat{x} . Tedaj vemo $(A + \Delta A)\hat{x} = b$, kjer velja $|\Delta A| \leq 3nu |V| |V^T|$. Ocenimo lahko

$$[|V| |V^T|]_{jk} = \sum_{i=1}^{\min(j,k)} |v_{ji}| |v_{ki}| \leq \sqrt{\sum_{i=1}^j |v_{ji}|^2} \sqrt{\sum_{i=1}^k |v_{ki}|^2}$$

po Cauchy-Schwarzu. To je nadalje enako

$$\leq \sqrt{a_{jj}} \sqrt{a_{kk}} \leq \|A\|_{\infty}.$$

Reševanje sistema z razcepom Choleskega je torej obratno stabilno, in velja

$$\|\Delta A\|_{\infty} \leq 3n^2 u \|A\|_{\infty},$$

kjer dodaten n na desni pride od tega, da smo na levi vzeli normo namesto absolutne vrednosti.

Vprašanje 36. Analiziraj stabilnost razcepa Choleskega.

3.4 Sistemi nelinearnih enačb

Rešujemo sistem

$$\begin{aligned} f_1(x_1, \dots, x_n) &= 0 \\ &\vdots \\ f_n(x_1, \dots, x_n) &= 0 \end{aligned}$$

3 Uvod v numerične metode

za $f_i : \mathbb{R} \rightarrow \mathbb{R}$ ali $\mathbb{C} \rightarrow \mathbb{C}$. Ekvivalentno $F(x) = 0$ za $F : \mathbb{R}^n \rightarrow \mathbb{R}^n$ (ali $\mathbb{C}^n \rightarrow \mathbb{C}^n$).

Prvi možni pristop reševanja je navadna iteracija. Sistem $F(x) = 0$ zapišemo v ekvivalentni obliki $x = G(x)$, izberemo $x^{(0)}$ ter iteriramo.

Izrek. Naj bo $G : \mathbb{R}^n \rightarrow \mathbb{R}^n$ zvezno odvedljiva na zaprti množici $\Omega \subseteq \mathbb{R}^n$. Če za $x \in \Omega$ velja

- $G(x) \in \Omega$,
- $\rho(JG(x)) \leq m < 1$,

kjer je JG Jacobijeva matrika, ρ pa spektralni radij (po absolutni vrednosti največja lastna vrednost), potem ima G na Ω natanko eno negibno točko α , in za vsak $x^{(0)} \in \Omega$ zaporedje $x^{(r+1)} = G(x^{(r)})$ konvergira k α .

Zadosten pogoj za konvergenco je že, da je $\|JG(\alpha)\| < 1$ v neki matrični normi. Za kvadratično konvergenco mora biti $JG(\alpha) = 0$ po komponentah.

Vprašanje 37. Kako poiščeš rešitev sistema nelinearnih enačb z navadno iteracijo? Povej izrek.

Podobno kot v enodimenzionalnem primeru lahko uporabimo razvoj v Taylorjevo vrsto in zanemarimo višje člene. Dobimo izraz za popravek

$$x^{(r+1)} = x^{(r)} - (JF(x^{(r)}))^{-1} F(x^{(r)}).$$

V praksi raje uporabimo algoritem 3.

Algorithm 3 Newtonova metoda

```

Izberi  $x^{(0)}$ .
for  $r = 0, 1, 2, \dots$  do
    Reši sistem  $JF(x^{(r)})\Delta x^{(r)} = -F(x^{(r)})$ .
     $x^{(r+1)} = x^{(r)} + \Delta x^{(r)}$ .
end for

```

Vprašanje 38. Razloži Newtonovo metodo za rešitev sistema nelinearnih enačb.

Ker je računanje Jacobijeve matrike zahtevno, se lahko poslužimo kakšne kvazi-Newtonove metode. Pri taki metodi na različne načine aproksimiramo Jacobijevo matriko in zmanjšamo zahtevnost enega koraka. S tem običajno pade red konvergence na superlinearno. Najbolj znana kvazi-Newtonova metoda je BROYDNOVA METODA, kjer približek Jacobijeve matrike B_{r+1} določimo kot najbližjo matriko, ki zadošča t.i. *sekantnemu pogoju*

$$B_{r+1}(x^{(r+1)} - x^{(r)}) = F(x^{(r+1)}) - F(x^{(r)}).$$

Ker je $B_r \Delta x^{(r)} = -F(x^{(r)})$, mora torej veljati

$$\Delta B_r \Delta x^{(r)} = F(x^{(r+1)}),$$

matrika ΔB_r pa je taka, da je $\|\Delta B_r\|_2$ minimalna.

Lema. Dana sta neničelna vektorja x, y . Matrika A z minimalno normo, ki preslika x v y , je

$$A = \frac{yx^T}{\|x\|_2^2}.$$

Dokaz. Očitno je $Ax = y$. Če za matriko B velja $Bx = y$, je $\|y\|_2 = \|Bx\|_2 \leq \|B\|_2 \|x\|_2$, torej

$$\|B\|_2 \geq \frac{\|y\|_2}{\|x\|_2}.$$

Po drugi strani se da preveriti, da za matrike ranga 1 velja $\|yx^T\|_2 = \|y\|_2 \|x\|_2$. \square

Algorithm 4 Broydnova metoda

Določi $x^{(0)}$ in B_0 .

for $r = 0, 1, \dots$ **do**

 Reši $B_r \Delta x^{(r)} = -F(x^{(r)})$.

 Izračunaj

$$B_{r+1} = B_r + \frac{F(x^{(r+1)})(\Delta x^{(r)})^T}{\|\Delta x^{(r)}\|_2^2}.$$

end for

Vprašanje 39. Izpelji Broydnovo metodo.

3.5 Linearni problemi najmanjših kvadratov

3.5.1 Normalni sistem

Dana je matrika $A \in \mathbb{R}^{m \times n}$ za $m > n$ in vektor $b \in \mathbb{R}^m$. Iščemo $x \in \mathbb{R}^n$, ki minimizira napako $\|Ax - b\|_2$. Ta napaka bo minimalna, ko bo Ax pravokotna projekcija b na sliko $\text{im } A$. Velja $Ax - b \perp Az$ za vse $z \in \mathbb{R}^n$ natanko tedaj, ko je za vsak z

$$z^T A^T (Ax - b) = 0.$$

Sledi $A^T(Ax - b) = 0$ oziroma $A^T Ax = A^T b$, čemur pravimo NORMALNI SISTEM. Pri tem smo tiho predpostavili, da je $\text{rang } A = n$, sicer sistem nebi bil enolično rešljiv.

Velja $w^T A^T A w = \|Aw\|_2^2 > 0$ za $w \neq 0$, torej je GRAMOVA MATRIKA $A^T A$ simetrična pozitivno definitna, in zanjo obstaja razcep Choleskega. Pri reševanju sistema seveda uporabimo ta razcep.

Vprašanje 40. Izpelji normalni sistem. Kaj je Gramova matrika?

3.5.2 QR razcep

Izrek. Naj bo $A \in \mathbb{R}^{m \times n}$ za $m \geq n$ polnega ranga. Potem obstaja enoličen razcep $A = QR$, kjer je $Q \in \mathbb{R}^{m \times n}$ z ortonormiranimi stolpci ($Q^T Q = I_n$) in $R \in \mathbb{R}^{n \times n}$ zgornje trikotna s pozitivnimi diagonalnimi elementi.

Dokaz. Če bi veljalo $A = QR$, je $A^T A = R^T Q^T Q R = R^T R$. Matrika $A^T A$ je simetrična pozitivno definitna, torej je $R^T R$ njen razcep Choleskega, in velja $R = V^T$. Iz $A = QR$ sledi $Q = AR^{-1}$. \square

Vprašanje 41. Povej in dokaži izrek o obstoju in enoličnosti QR razcepa.

Če poznamo $A = QR$, je $\text{im } A = \text{im } Q$. V drugem primeru bo vsakršno delo stabilnejše, ker so stolpci ortonormirani. Normalni sistem se tedaj prevede na $Rx = Q^T b$, ki ga lahko rešimo s premo substitucijo.

3.5.3 Gram-Schmidtova ortogonalizacija

Poznamo tri načine za izračun QR razcepa. Najenostavnejši pristop je Gram-Schmidtova ortogonalizacija. Velja

$$a_k = \sum_{i=1}^{k-1} r_{ik} q_i + r_{kk} q_k.$$

Če to enačbo množimo z leve s q_j^T , nam ostane

$$r_{jk} = q_j^T a_k.$$

Celoten postopek je prikazan v algoritmu 5, ki ima računsko zahtevnost $2mn^2$.

Algorithm 5 QR razcep s klasičnim Gram-Schmidtovim postopkom

```

for  $k = 1, \dots, n$  do
   $q_k = a_k$ 
  for  $i = 1, \dots, k - 1$  do
     $r_{ik} = q_i^T a_k$ 
     $q_k = q_k - r_{ik} q_i$ 
  end for
   $r_{kk} = \|q_k\|_2$ 
   $q_k = \frac{1}{r_{kk}} q_k$ 
end for
```

Vprašanje 42. Izpelji in zapiši klasični Gram-Schmidtov postopek za izračun QR razcepa.

V algoritmu naredimo še popravek, ki bo povečal stabilnost; pri računanju r_{ik} namesto formule $r_{ik} = q_i^T a_k$ uporabimo $r_{ik} = q_i^T q_k$. Novemu postopku pravimo MODIFICIRAN

GRAM-SCHMIDT, in je v teoriji ekvivalenten klasičnemu. Modificiran postopek moramo tudi bolj pametno uporabiti; izračunamo QR razcep razširjene matrike

$$\begin{bmatrix} Q & q_{n+1} \end{bmatrix} \cdot \begin{bmatrix} R & z \\ 0 & \rho \end{bmatrix} = \begin{bmatrix} A & b \end{bmatrix},$$

in dobimo

$$Ax - b = \begin{bmatrix} Q & q_{n+1} \end{bmatrix} \cdot \begin{bmatrix} R & z \\ 0 & \rho \end{bmatrix} \cdot \begin{bmatrix} x \\ -1 \end{bmatrix} = \begin{bmatrix} Q & q_{n+1} \end{bmatrix} \cdot \begin{bmatrix} Rx - z \\ -\rho \end{bmatrix}.$$

Najboljšo rešitev dobimo, ko je $Rx = z$. Dejansko je $z = Q^T b$, le da smo ga v tem postopku izračunali z modificiranim Gram-Schmidtovim postopkom, kar je numerično bolje.

Vprašanje 43. Kaj je modificiran Gram-Schmidtov postopek? Kako ga pravilno uporabiš za reševanje sistema najmanjših kvadratov?

3.5.4 Givensove rotacije

Če je $c = \cos \varphi$ in $s = \sin \varphi$, matrika

$$R_{ik}^T(\varphi) = \begin{bmatrix} 1 & & & & & & & & \\ & \ddots & & & & & & & \\ & & 1 & & & & & & \\ & & & c & & & s & & \\ & & & & 1 & & & & \\ & & & & & \ddots & & & \\ & & & & & & 1 & & \\ & & & -s & & & c & & \\ & & & & & & & 1 & \\ & & & & & & & & \ddots & \\ & & & & & & & & & 1 \end{bmatrix},$$

ki ima elemente na diagonali in v stolpcih in vrsticah i, k , predstavlja rotacijo za φ v ravnini, ki jo razpenjata e_i in e_k v \mathbb{R}^m . Z ustrezno izbiro c in s lahko slikamo (x_i, x_k) v $(y_i, 0)$. Če je $r = \sqrt{x_i^2 + x_k^2}$, je taka izbira $c = x_i/r$ in $s = x_k/r$. Če te rotacije ustrezno kombiniramo, dobimo QR razcep, kakor je prikazano v algoritmu 6, ki ima zahtevnost $3mn^2 - n^3$ če ne računamo Q , za računanje Q pa porabimo še $6m^2n - 3mn^2$ operacij.

Vprašanje 44. Izpelji QR razcep z Givensovimi rotacijami. Kakšna je njegova časovna zahtevnost?

Algorithm 6 QR razcep z Givenssonovimi rotacijami

```

 $Q = I_m$ 
for  $i = 1, \dots, n$  do
  for  $k = i + 1, \dots, m$  do
     $r = \sqrt{a_{ii}^2 + a_{ki}^2}$ 
     $c = a_{ii}/r$ 
     $s = a_{ki}/r$ 
    Izračunaj
  
```

$$A([i, k], i : n) = \begin{bmatrix} c & s \\ -s & c \end{bmatrix} \cdot A([i, k], i : n)$$

$$b([i, k]) = \begin{bmatrix} c & s \\ -s & c \end{bmatrix} \cdot b([i, k])$$

$$Q([i, k], :) = \begin{bmatrix} c & s \\ -s & c \end{bmatrix} \cdot Q([i, k], :)$$

```

  end for
end for
 $Q = Q^T$ 

```

3.5.5 Householderjeva zrcaljenja

Vzemimo $w \in \mathbb{R}^m$, ki je različen od 0, in definirajmo

$$P = I - \frac{2}{w^T w} w w^T.$$

Velja $P = P^T$ in $P^2 = I$, poleg tega pa je w lastni vektor za P z lastno vrednostjo -1 . Če je $u \perp w$, je $Pu = u$. Preslikavo lahko torej obravnavamo kot zrcaljenje čez ravnino, katere normala je w .

Če imamo dana dva enako dolga vektorja x, y , lahko z izbiro $w = x - y$ dobimo $y = Px$. Z izbiro $w = x \mp \|x\|_2 e_1$ se x preslika v

$$P \cdot \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} \pm \|x\|_2 \\ 0 \\ \vdots \\ 0 \end{bmatrix}.$$

Za numerično stabilnost si izberemo, da prištevamo, če je x_1 pozitiven, sicer odštevamo; $w = x + \operatorname{sgn} x_1 \|x\|_2 e_1$, kjer je $\operatorname{sgn} 0 \neq 0$. Z zrcaljenjem na enem koraku uničimo celoten stolpec matrike. Postopek izračuna QR razcepa je prikazan v algoritmu 7. Algoritem ima časovno zahtevnost $2mn^2 - \frac{2}{3}n^3$, če nas ne zanima Q .

Vprašanje 45. Izpelji QR razcep s Householderjevimi zrcaljenji.

Algorithm 7 QR razcep s Householderjevimi zrcaljenji

```

 $Q = I_m$ 
for  $i = 1, \dots, n$  do
    Določi  $w_i \in \mathbb{R}^{m-i+1}$  iz  $A(i : m, i)$ 
     $A(i : m, i : n) = P_i A(i : m, i : n)$ 
     $b(i : m) = P_i b(i : m)$ 
     $Q(i : m, :) = P_i Q(i : m, i)$ 
end for
 $Q = Q^T$ 

```

Občutljivost predoločenega sistema $Ax = b$ je odvisna od $\kappa_2(A) + \|r\|_2 \kappa_2^2(A)$ za $r = b - Ax$. Če uporabimo normalni sistem $A^T Ax = A^T b$, rešujemo z občutljivostjo $\kappa_2(A^T A) = \frac{\sigma_1^2(A)}{\sigma_n^2(A)} = \kappa_2^2(A)$, če pa uporabimo QR razcep, pa velja $\kappa_2(R) = \kappa_2(A)$.

Vprašanje 46. Kakšna je občutljivost reševanja predoločenega sistema?

3.6 Lastne vrednosti

Dana je matrika $A \in \mathbb{R}^{n \times n}$. Iščemo njene lastne vrednosti, in morda lastne vektorje. Iščemo lahko tudi leve lastne vektorje $y^H A = \lambda y^H$.

Lema. Če je x desni lastni vektor A za lastno vrednost λ in je y levi lastni vektor za lastno vrednost $\mu \neq \lambda$, potem je $y^H x = 0$.

Dokaz. Velja $\lambda y^H x = y^H A x = \mu y^H x$, torej $y^H x = 0$. □

Posledica. Če je A simetrična in sta x, y lastna vektorja za različni lastni vrednosti, je $y^T x = 0$.

Vprašanje 47. Kakšni so lastni vektorji simetrične matrike? Dokaži.

Vprašanje 48. Zakaj Jordanova forma ni primerna za numerično računanje?

Odgovor: Poglejmo si primer

$$A(\varepsilon) = \begin{bmatrix} 0 & 1 \\ \varepsilon & 0 \end{bmatrix}.$$

Matrika $A(0)$ je že svoja Jordanova kletka, ki pa ni blizu Jordanove forme za $\varepsilon \neq 0$, ki je

$$J = \begin{bmatrix} \sqrt{\varepsilon} & \\ & -\sqrt{\varepsilon} \end{bmatrix}.$$

☒

Izrek. Za vsako matriko $A \in \mathbb{R}^{n \times n}$ obstaja Schurova forma $A = U S U^H$, kjer je U unitarna in S zgornje trikotna.

3 Uvod v numerične metode

Dokaz. Dokažemo z indukcijo na n . Za $n = 1$ je primeren razcep $A = [1]A[1]$. Za splošen n pa: Vsaka matrika ima vsaj en lastni vektor, torej velja $Ax = \lambda x$ za nek x velikosti 1, in nek λ . Za U_1 izberemo tako unitarno matriko, da je $U_1 e_1 = x$, in definiramo $B = U_1^H A U_1$. Velja $B e_1 = \lambda e_1$, torej je B oblike

$$B = \begin{bmatrix} \lambda & y^T \\ 0 & C \end{bmatrix}.$$

Za matriko C velja indukcijska predpostavka: obstajata U_2 in S_2 , da je $C = U_2 S_2 U_2^H$, in je U_2 unitarna, S_2 pa zgornje trikotna. Definiramo

$$U = U_1 \begin{bmatrix} 1 & \\ & U_2 \end{bmatrix}, \\ S = U^H A U.$$

□

Vprašanje 49. Dokaži, da za vsako matriko obstaja Schurova forma.

Izrek. Če je $A \in \mathbb{R}^{n \times n}$, obstajata ortogonalna Q in kvazi zgornje trikotna T , obe realni, da je $A = QTQ^T$.

Kvazi zgornje trikotna matrika je zgornje trikotna, razen da na diagonali dopuščamo bloke 2x2. V vsakem takem bloku se nahajajo konjugirani pari kompleksnih lastnih vrednosti. Razcepu $A = QTQ^T$ pravimo REALNA SCHUROVA FORMA.

Vprašanje 50. Kaj je realna Schurova forma?

3.6.1 Potenčna metoda

Algorithm 8 Potenčna metoda

Izberemo nek $z_0 \neq 0$.

for $k = 1, 2, \dots$ **do**

$y_k = A z_{k-1}$

$z_k = y_k / \|y_k\|$

end for

Izrek. Naj bo λ_1 dominantna lastna vrednost matrike A (največja po absolutni vrednosti in strogo večja od druge največje). Za naključno izbrani začetni vektor z_0 vektorji z_k iz potenčne metode po smeri konvergirajo k lastnemu vektorju za λ_1 .

Dokaz. Dokažemo le za primer, ko je A diagonalizabilna. Naj bo $A = XDX^{-1}$ za $X = [x_1, \dots, x_n]$ in $D = \text{diag}(\lambda_1, \dots, \lambda_n)$. Začetni vektor z_0 lahko razvijemo v bazi lastnih vektorjev

$$z_0 = \sum_{i=1}^n \alpha_i x_i.$$

Za vsak k velja

$$z_k = \frac{A^k z_0}{\|A^k z_0\|} = \frac{\sum_i \alpha_i \lambda_i^k x_i}{\left\| \sum_i \alpha_i \lambda_i^k x_i \right\|}.$$

Ulomek na obeh straneh delimo z λ_1^k :

$$z_k = \frac{\sum_i \alpha_i \lambda_i^k \lambda_1^{-k} x_i}{\left\| \sum_i \alpha_i \lambda_i^k \lambda_1^{-k} x_i \right\|}$$

Če je $\alpha_1 \neq 0$, bo to po smeri konvergiralo k x_1 . □

Vprašanje 51. Zapiši algoritem potenčne metode in dokaži pravilnost za digonalizabilne matrike.

Konvergenca metode je linearna in odvisna od razmerja $|\lambda_2/\lambda_1|$. Če je z približek za lastni vektor, lahko približek za pripadajočo lastno vrednost dobimo z Rayleighovim kvocientom

$$\lambda = \frac{z^H A z}{z^H z} = \rho(z, A).$$

Pri tem dejansko rešimo predoločen sistem $Az = \lambda z$, kjer je z matrika, λ spremenljivka in Az desna stran. Za kvocient veljata naslednji lastnosti:

- $\rho(\alpha z, A) = \rho(z, A)$ za $\alpha \neq 0$
- $\rho(x_i, A) = \lambda_i$

To nam poda ustavitveni kriterij za algoritem; izračunamo $\rho(z_k, A)$ in primerjamo $\|Az_k - \rho(z_k, A)z_k\| < \varepsilon$.

Vprašanje 52. Kaj je Rayleighov kvocient?

Denimo, da smo našli λ_1 in x_1 . Kako sedaj poiščemo λ_2 in x_2 ? Podobno kot pri dokazu obstoja Schurove forme poiščemo unitarno matriko U_1 , da je $U_1 e_1 = x$ (npr. s Householderjevim zrcaljenjem). Vzamemo

$$B = U^H A U = \begin{bmatrix} \lambda_1 & \alpha^T \\ 0 & C \end{bmatrix},$$

in nadaljujemo s potenčno metodo na C , ki pa je seveda ne izrazimo eksplicitno. Množenje izvedemo tako, da množimo z B , in iz produkta vzamemo spodnjih $n - 1$ elementov.

Vprašanje 53. Kako s potenčno metodo poiščeš vse lastne vrednosti?

3.6.2 Inverzna iteracija

Naj bo σ zelo dober približek za eno izmed lastnih vrednosti matrike A . Matrika $(A - \sigma I)^H$ ima lastne vrednosti $\frac{1}{\lambda_i - \sigma}$, in velja

$$\frac{1}{|\lambda_j - \sigma|} \gg \frac{1}{|\lambda_i - \sigma|},$$

kjer je λ_j tista lastna vrednost, za katero je σ dober približek. To dejstvo izrabimo v algoritmu 9.

Algorithm 9 Inverzna iteracija

```

Izberi naključen  $z_0 \neq 0$ .
for  $k = 1, 2, \dots$  do
    Reši  $z_{k-1} = (A - \sigma I)y_k$ 
     $z_k = y_k / \|y_k\|$ 
end for
  
```

Vprašanje 54. Opiši delovanje inverzne iteracije.

3.6.3 Ortogonalna iteracija

Definicija. Prostor N je INVARIANTNI PODPROSTOR za matriko A , če je za vsak $x \in N$ tudi $Ax \in N$.

Izrek. Naj bo $S = [S_1 S_2]$ nesingularna matrika. Naj bo $B = S^{-1}AS$ oblike

$$B = \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix}.$$

Potem stolpci S_1 razpenjajo invariantni podprostor za A natanko tedaj, ko je $B_{21} = 0$.

Dokaz. Ker je $SB = AS$, velja $S_1 B_{11} + S_2 B_{21} = AS_1$. □

Recimo, da je res $B_{21} = 0$. Naj za lastne vrednosti matrike A velja $|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_p| > |\lambda_{p+1}| \geq \dots$, kjer je p število stolpcev v S_1 . V tem primeru je invariantni podprostor velikosti p , ki mu pripadajo največje lastne vrednosti, dominanten.

Algorithm 10 Ortogonalna iteracija

```

Izberi naključno matriko  $Z_0 \in \mathbb{R}^{n \times p}$ .
for  $k = 1, 2, \dots$  do
     $Y_k = AZ_{k-1}$ 
    Za  $Z_k$  vzemi  $Q_k$  iz QR razcepa  $Y_k = Q_k S_k$ .
end for
  
```

Izrek. Če za lastne vrednosti A velja $|\lambda_1| \geq \dots \geq |\lambda_p| > |\lambda_{p+1}|$, potem za naključno izbrano matriko $Z_0 \in \mathbb{R}^{n \times p}$ matrika Z_k konvergira proti ortonormirani bazi za dominantni invariantni podprostor dimenzije p .

Dokaz. Dokažemo samo za primer, kjer lahko matriko A diagonaliziramo v $A = XDX^{-1}$. Označimo $D = \text{diag}(D_1, D_2)$ in $X = [X_1, X_2]$, kjer sta D_1 in X_1 dimenzije p . Ker je $|\lambda_p| > 0$, velja $\det D_1 \neq 0$. Če izrazimo Z_0 v bazi lastnih vektorjev kot

$$Z_0 = X \begin{bmatrix} W_1 \\ W_2 \end{bmatrix},$$

potem lahko za naključen Z_0 predpostavimo, da je W_1 nesingularna. Pokazati moramo, da $\text{im } Z_k \rightarrow \text{im } X_1$ za $k \rightarrow \infty$. Velja

$$\text{im } Z_k = \text{im } AZ_{k-1} = \dots = \text{im } A^k Z_0 = \text{im } X D^k X^{-1} X \begin{bmatrix} W_1 \\ W_2 \end{bmatrix} = \text{im } X \begin{bmatrix} D_1^k W_1 \\ D_2^k W_2 \end{bmatrix}$$

Ker je W_1 obrnljiva, je to enako

$$\text{im } Z_k = \text{im } X \begin{bmatrix} I \\ D_2^k W_2 W_1^{-1} D_1^{-k} \end{bmatrix} D_1^k W_1 = \text{im } X \begin{bmatrix} I \\ D_2^k W_2 W_1^{-1} D_1^{-k} \end{bmatrix}.$$

Elemente spodnje matrike lahko razpišemo kot

$$\left[D_2^k W_2 W_1^{-1} D_1^{-k} \right]_{ij} = [W_2 W_1^{-1}]_{ij} \frac{\lambda_{p+i}^k}{\lambda_j^k}$$

Ker so vsi elementi D_1 večji od vseh v D_2 , ulomek na desni konvergira k 0 za $k \rightarrow \infty$, torej $\text{im } Z_k \rightarrow \text{im } X_1$. \square

Vprašanje 55. Pod katerim pogojem deluje ortogonalna iteracija? Dokaži za diagonalizabilne matrike.

Posledica. Če za matriko A velja $|\lambda_1| > \dots > |\lambda_n|$, potem za naključno izbrano matriko Z_0 matrika $Z_k^T A Z_k$ konvergira proti Schurovi formi.

Dokaz. Vzemimo $p \in \{1, \dots, n\}$ in razcepimo $Z_k = [Z_{k1} Z_{k2}]$. Po izreku Z_{k1} konvergira proti ONB za dominantni invariantni podprostor dimenzije p . Matrika $Z_{k2}^T A Z_{k1}$ torej konvergira k 0 (po izreku od prej), to pa velja za vsak p , torej smo dobili zgornje trikotno matriko, kjer so elementi na diagonali lastne vrednosti, urejene po absolutni vrednosti. \square

Vprašanje 56. Kako z ortogonalno iteracijo poiščeš Schurovo formo? Dokaži.

3.6.4 QR iteracija

Algorithm 11 QR iteracija

```

 $A_0 = A$ 
for  $k = 0, 1, \dots$  do
    Izračunaj QR razcep  $A_k = Q_k R_k$ 
     $A_{k+1} = R_k Q_k$ 
end for

```

Če za lastne vrednosti velja $|\lambda_1| > |\lambda_2| > \dots > |\lambda_n|$, potem A_k konvergira k Schurovi formi. Velja $A_{k+1} = Q_k^T A_k Q_k$, torej so si vse A_k ortogonalno podobne.

3 Uvod v numerične metode

Izrek. Za matriko A_k iz QR iteracije velja $A_k = Z_k^T A Z_k$, kjer je Z_k matrika iz ortogonalne iteracije z začetkom $Z_0 = I$.

Dokaz. Dokažemo z indukcijo na k , kjer je baza indukcije trivialna. Če velja za k , računamo $A Z_k = Z_{k+1} S_{k+1}$, zato $Z_k^T A Z_k = Z_k^T Z_{k+1} S_{k+1}$. Matrika $Z_k^T Z_{k+1}$ je ortogonalna, S_{k+1} pa zgornje trikotna, torej je to QR razcep matrike $Z_k^T A Z_k = A_k$ po induksijski predpostavki. Velja torej

$$A_{k+1} = R_k Q_k = S_{k+1} Z_k^T Z_{k+1} = Z_{k+1}^T A Z_k Z_k^T Z_{k+1} = Z_{k+1}^T A Z_{k+1}.$$

□

Vprašanje 57. Zapiši algoritem QR iteracije. Kako deluje? Dokaži.

Definicija. Matrika H je ZGORNJE HESSENBERGOVA, če je $h_{ij} = 0$ za $i > j + 1$.

Lema. Če je matrika A zgornje Hessenbergova, se oblika med QR iteracijo ohranja.

Dokaz. Če je A zgornja Hessenbergova, je Q zgornja Hessenbergova.

□

Algorithm 12 Redukcija matrike v zgornje Hessenbergovo obliko

$Q = I$

for $k = 1, \dots, n - 2$ **do**

 Določi $w_k \in \mathbb{R}^{n-k}$ za Householderjevo zrcaljenje, ki slika $A(k + 1 : n, k)$ v $\pm x e_1$

$A(k + 1 : n, k : n) = P_k A(k + 1 : n, k : n)$

$A(:, k + 1 : n) = A(:, k + 1 : n) P_k$

$Q(:, k + 1 : n) = Q(:, k + 1 : n) P_k$

end for

Lema nam pove, da si lahko prihranimo veliko računanja med iteracijo. Prvo lahko s Householderjevimi zrcaljenji pretvorimo A v zgornjo Hessenbergovo obliko, kot v algoritmu 12 Izračun QR razcepa take matrike lahko naredimo z $n - 1$ Givenssonovimi rotacijami, torej v $O(n^2)$. Če je R_{ij} rotacija, velja

$$R_{n-1,n}^T \dots R_{23}^T R_{12}^T A_k = R_k$$

oziroma

$$A_{k+1} = R_k Q_k = R_k R_{12} R_{23} \dots R_{n-1,n}.$$

Vprašanje 58. Kako optimiziraš korak QR iteracije na $O(n^2)$?

Definicija. Zgornje Hessenbergova matrika H je NERAZCEPNA, če so vsi elementi v spodnji diagonalni neničelni. V nasprotnem primeru je RAZCEPNA.

Algorithm 13 QR iteracija z enojnim premikom

```

Naredi redukcijo na Hessenbergovo obliko  $A_0 = Q^T A Q$ 
for  $k = 0, 1, \dots$  do
    Izberi premik  $\sigma_k$ 
    Naredi QR razcep  $A_k - \sigma_k I = Q_k R_k$ 
     $A_{k+1} = R_k Q_k + \sigma_k I$ 
end for

```

Če je matrika razcepna, lahko problem razdelimo na dva manjša. Predpostavimo lahko torej, da je začetna Hessenbergova matrika nerazcepna. Pri numeričnem reševanju element $h_{i+1,i}$ proglasimo za 0, če je $|h_{i+1,i}| \leq \varepsilon(|h_{i,i}| + |h_{i+1,i+1}|)$.

Število iteracij lahko zmanjšamo z uporabo premikov, kot je prikazano v algoritmu 13.

Matriki A_k in A_{k+1} sta še vedno ortogonalno podobni, ker velja

$$A_{k+1} = R_k Q_k + \sigma_k I = Q_k^T (A_k - \sigma_k I) Q_k + \sigma_k I = Q_k^T A_k Q_k.$$

Lema. Če je A nerazcepna zgornje Hessenbergova matrika in za premik σ izberemo lastno vrednost A , potem za matriko $B = RQ + \sigma I$, kjer je $A - \sigma I = QR$, velja $b_{n,n-1} = 0$ in $b_{n,n} = \sigma$.

Dokaz. Matrika $A - \sigma I$ je singularna, zaradi nerazcepnosti pa je prvih $n - 1$ stolpcev linearno nedovisnih, torej mora veljati $r_{n,n} = 0$. Torej je zadnja vrstica R enaka 0, in je B predpisane oblike. \square

Vprašanje 59. Kako deluje QR iteracija s premikom? Kaj je njena prednost? Dokaži.

Idealno je za premik torej izbrati čim boljši približek za lastno vrednost. Poznamo dve pogosti izbiri premika;

- Enojni premik: izberemo $\sigma_k = (A_k)_{nn}$. Deluje dobro za matrike s samimi realnimi lastnimi vrednostmi.
- Dvojni premik: za σ_{k1} in σ_{k2} izberemo lastni vrednosti matrike $A(n-1 : n, n-1 : n)$, in v okviru ene iteracije naredimo dva premika

$$\begin{aligned} A_k - \sigma_{k1} I &= Q_k R_k & A_{k+1/2} - \sigma_{k2} I &= \tilde{Q}_k \tilde{R}_k \\ A_{k+1/2} &= R_k Q_k + \sigma_{k1} I & A_{k+1} &= \tilde{R}_k \tilde{Q}_k + \sigma_{k2} I \end{aligned}$$

Pokažemo lahko, da je za realno A_k tudi A_{k+1} realna.

Vprašanje 60. Kako izberemo premik za QR iteracijo?

3.7 Polinomska interpolacija

3.7.1 Lagrangeova oblika

Problem z interpolacijo v standardni bazi je, da je Vandermondova matrika zelo občutljiva.

Izrek. Za paroma različne točke x_0, \dots, x_n in vrednosti y_0, \dots, y_n obstaja natanko en polinom stopnje $\leq n$, da je $p(x_i) = y_i$.

Dokaz. Obstoj dokažemo s konstrukcijo, kjer uporabimo Lagrangeove bazne polinome

$$l_{n,i}(x_j) = \begin{cases} 1 & i = j \\ 0 & i \neq j \end{cases}.$$

Če take polinome lahko konstruiramo (kjer je $l_{n,i}$ stopnje $\leq n$), bo veljalo

$$p = \sum_i y_i l_{n,i}.$$

Če definiramo

$$l_{n,i}(x) = \frac{(x - x_0) \dots (x - x_{i-1})(x - x_{i+1}) \dots (x - x_n)}{(x_i - x_0) \dots (x_i - x_{i-1})(x_i - x_{i+1}) \dots (x_i - x_n)},$$

potem $l_{0,n}, \dots, l_{n,n}$ sestavljajo bazo za polinome stopnje $\leq n$.

Dokazati moramo še enoličnost: denimo, da je \tilde{p} tudi polinom stopnje $\leq n$, ki zadošča $\tilde{p}(x_i) = y_i$. Tedaj je $q = p - \tilde{p}$ polinom stopnje $\leq n$, za katerega velja $q(x_i) = 0$. Polinom stopnje $\leq n$, ki ima vsaj $n + 1$ ničel, je enak $q = 0$. \square

Vprašanje 61. Kako interpoliraš polinom z Lagrangeovo bazo? Dokaži, da je interpolacija enolična.

Izrek. Če je f $(n + 1)$ -krat zvezno odvedljiva na $[a, b]$, ki vsebuje paroma različne vozle x_0, \dots, x_n ter x , in je p interpolacijski polinom za f , potem velja

$$f(x) - p(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!} \omega(x)$$

za $\omega(x) = (x - x_0) \dots (x - x_n)$ in $\min(x_0, \dots, x_n, x) \leq \xi \leq \max(x_0, \dots, x_n, x)$.

Dokaz. Predpostavimo lahko, da je x različen od x_0, \dots, x_n . Definiramo $g(z) = f(z) - p(z) - C\omega(z)$, kjer nastavimo C tako, da ima g pri x ničlo. Poleg tega ima g tudi ničle v vozlih x_0, \dots, x_n . Funkcija g je $(n + 1)$ -krat zvezno odvedljiva in ima $n + 2$ različnih ničel. Po Rollovem izreku ima g' vsaj $n + 1$ ničel. Če postopek nadaljujemo, dobimo, da ima $g^{(n+1)}$ ničlo ξ . Če upoštevamo, da sta p in ω polinoma, dobimo

$$C = \frac{f^{(n+1)}(\xi)}{(n+1)!}.$$

\square

Vprašanje 62. Kako izraziš napako interpolacije? Dokaži.

3.7.2 Deljene difference

Definicija. Za paroma različne točke x_0, \dots, x_k je DELJENA DIFERENCA $[x_0, \dots, x_k]_f$ vodilni koeficient (pri x^k) interpolacijskega polinoma za f v točkah x_0, \dots, x_k .

Izrek. Za paroma različne točke x_0, \dots, x_n lahko interpolacijski polinom za f na x_0, \dots, x_n zapišemo kot

$$p(x) = [x_0]_f + [x_0, x_1]_f(x - x_0) + \dots + [x_0, \dots, x_n]_f(x - x_0) \dots (x - x_{n-1}).$$

Dokaz. Dokažemo z indukcijo na n . Baza indukcije je očitna, dokažimo korak. Naj bo $p_n(x)$ polinom stopnje $\leq n$, ki se z f ujema na x_0, \dots, x_n . Dodamo še točko x_{n+1} in iščemo p_{n+1} . Za

$$p_{n+1}(x) = p_n(x) + C(x - x_0) \dots (x - x_n)$$

in ustrezen C velja $p_{n+1}(x_i) = f(x_i)$. Izračunamo torej

$$C = \frac{f(x_{n+1}) - p_n(x_{n+1})}{(x_{n+1} - x_0) \dots (x_{n+1} - x_n)}.$$

□

Vprašanje 63. Kakšna je Newtonova oblika interpolacije? Dokaži.

Izrek. Za deljene difference velja

- $[x_0, \dots, x_n]_f$ je simetrična funkcija argumentov
- $[x_0, \dots, x_n]_f$ je linearen funkcional v f
- Velja rekurzivna formula

$$[x_0, \dots, x_k]_f = \frac{[x_1, \dots, x_k]_f - [x_0, \dots, x_{k-1}]_f}{x_k - x_0}.$$

Dokaz. Prvi točki sta očitni. Za tretjo: Naj p_a interpolira f v točkah x_0, \dots, x_{k-1} , in p_b v točkah x_1, \dots, x_k . Hitro lahko preverimo, da je ustrezen interpolacijski polinom

$$p(x) = \frac{x - x_k}{x_0 - x_k} p_a(x) + \frac{x - x_0}{x_k - x_0} p_b(x),$$

ker je stopnja p manjša ali enaka k in $p(x_0) = p_a(x_0) = f(x_0)$ ter $p(x_k) = p_b(x_k) = f(x_k)$. Za ostale točke tudi velja $p(x_i) = f(x_i)$. □

Vprašanje 64. Povej in dokaži rekurzivno formulo za deljene difference.

3 Uvod v numerične metode

Formula nam pove, da je smiselna definicija, če se točke x_i ponavljajo, naslednja:

$$[x_0, \dots, x_k]_f = \begin{cases} \frac{f^{(k)}(x_0)}{k!} & x_0 = x_1 = \dots = x_k \\ \frac{[x_1, \dots, x_k]_f - [x_0, \dots, x_{k-1}]_f}{x_k - x_0} & \text{sicer} \end{cases}$$

kjer pri drugi definiciji poskrbimo, da $x_0 \neq x_k$ (vrstni red točk je nepomemben).

Izrek. Za k -krat zvezno odvedljivo f velja

$$[x_0, \dots, x_k]_f = \int_0^1 dt_1 \int_0^{t_1} dt_2 \dots \int_0^{t_{k-1}} f^{(k)}(\xi_k) dt_k,$$

kjer je $\xi_k = t_k(x_k - x_{k-1}) + \dots + t_1(x_1 - x_0) + x_0$.

Posledica. Za k -krat zvezno odvedljivo f velja

$$[x_0, \dots, x_k]_f = \frac{f^{(k)}(\xi)}{k!},$$

kjer je $\min(x_0, \dots, x_k) \leq \xi \leq \max(x_0, \dots, x_k)$.

Dokaz. Uporabimo zadnji izrek in izrek o povprečni vrednosti. Upoštevamo, da je volumen simpleksa enak $k!^{-1}$. \square

Izrek. Če je p interpolacijski polinom za f v točkah x_0, \dots, x_n , potem velja

$$f(x) - p(x) = [x_0, \dots, x_n, x]_f (x - x_0) \dots (x - x_n).$$

Dokaz. Če desni strani prištejemo $p(x)$, dobimo interpolacijski polinom za f točkah x_0, \dots, x_n in x . \square

Posledica tega je, da je ocena za napako enaka kot prej, tudi če uporabljamo ponovljene točke.

Vprašanje 65. Kako izraziš oceno za napako interpolacije z deljenimi diferencami?

3.8 Numerično integriranje

Želimo izračunati integral

$$I(f) = \int_a^b f(x) dx.$$

Ideja je, da funkcijo aproksimiramo z interpolacijskim polinomom

$$\int_a^b f(x) dx = \int_a^b \sum_{i=0}^n f(x_i) l_{n,i}(x) dx + R(f) = \sum_{i=0}^n f(x_i) \int_a^b l_{n,i}(x) dx + R(f),$$

kjer integrale interpolacijskih polinomov imenujemo UTEŽI ali KOEFICIENTI α_i , napaka $R(f)$ pa je oblike

$$R(f) = \int_a^b \frac{f^{(n+1)}(\xi)}{(n+1)!} \omega(x) dx.$$

Tako dobljene kvadrature formule so določene z izbiro vozlov. Napaka pri računanju se razdeli na dve komponenti; napaka metode $R(f)$ ter neodstranljiva napaka, ki jo dobimo, ker ne poznamo točnih vrednosti f v vozlih.

Vprašanje 66. Izpelji obliko kvadrature formul.

3.8.1 Newton-Cotesove formule

Pri NC formulah vozle izberemo enakomerno, $x_i = a + ih$. Ločimo zaprti tip formul, kjer uporabimo vse vozle, in odprti tip, kjer izpustimo vozla v krajiščih.

Najenostavnejša kvadratura formula je TRAPEZNA FORMULA, t.j. formula zaprtega tipa za $n = 1$

$$\int_{x_0}^{x_1} f(x) dx = \frac{h}{2}(f(x_0) + f(x_1)) + \int_{x_0}^{x_1} \frac{1}{2} f''(\xi_x)(x - x_0)(x - x_1) dx,$$

kjer za napako po izreku o povprečni vrednosti velja

$$R(f) = \frac{1}{2} f''(\xi) \int_{x_0}^{x_1} (x - x_0)(x - x_1) dx = -\frac{h^3}{12} f''(\xi).$$

Vprašanje 67. Opiši trapezno formulo.

Če namesto tega vzamemo $n = 2$, dobimo SIMPSONOVO FORMULO

$$\int_{x_0}^{x_2} f(x) dx = \frac{h}{3}(f(x_0) + 4f(x_1) + f(x_2)) + R(f)$$

za napako

$$R(f) = \int_{x_0}^{x_2} \frac{1}{6} f'''(\xi_x)(x - x_0)(x - x_1)(x - x_2) dx.$$

Če je f polinom stopnje 3, je f''' konstantna, in velja $R(f) = 0$. Podobno se zgodi tudi pri vseh NC formulah za sode n . Oblika napake za vse formule je vedno $R(f) = C f^{(m)}(\xi)$, kjer je m stopnja najnižjega polinoma, za katerega formula ni točna. Za Simpsonovo formulo je to $m = 4$, če vstavimo $f(x) = (x - x_0)^4$ dobimo $C = -h^5/90$.

Vprašanje 68. Povej Simpsonovo formulo. Za katero stopnjo polinomov je točna? Izpelji predpis za napako.

Druga vrsta so Newton-Cotesove formule zaprtega tipa, ki so smiselne za $n \geq 2$. Pri $n = 2$ dobimo SREDINSKO FORMULO

$$\int_{x_0}^{x_2} f(x) dx = 2hf(x_1) + \frac{h^3}{3} f''(\xi),$$

3 Uvod v numerične metode

za $n = 4$ pa MILNEOVO FORMULO

$$\int_{x_0}^{x_4} f(x)dx = \frac{4h}{3}(2f(x_1) - f(x_2) + 2f(x_3)) + \frac{28h^5}{90}f^{(4)}(\xi),$$

ki NC formula z najmanj vozli in negativno utežjo.

Vprašanje 69. Kaj sta sredinska in Milneova formula? Zakaj je druga omembe vredna?

3.8.2 Napake pri numeričnem integriranju

Velja

$$\int_a^b f(x)dx = \sum_{i=0}^n \alpha_i f(x_i) + D_n,$$

kjer za napako metode D_n ne bo veljalo nujno $D_n \rightarrow 0$ za $n \rightarrow \infty$. Če npr. integriramo $f(x) = (1+x^2)^{-1}$ na $[-5, 5]$ z enakomerno razporejenimi točkami, bo veljalo

$$\max_{x \in [-5, 5]} |f(x) - p_n(x)| \rightarrow \infty,$$

kjer je p_n interpoliran polinom za n točk. Pri izračunu vsote se pojavita še neodstranljiva D_n in zaokrožitvena napaka. Recimo, da velja $|f(x_i) - f_i| \leq \varepsilon$, kjer je f_i izračunan približek in $f(x_i)$ točna vrednost. Potem lahko ocenimo

$$|D_n| = \left| \sum_{i=0}^n \alpha_i f(x_i) - \sum_{i=0}^n \alpha_i f_i \right| \leq \sum_{i=0}^n |\alpha_i| \varepsilon.$$

Če je $\alpha_i \geq 0$ za vse i , lahko to nadalje ocenimo z $|D_n| \leq (b-a)\varepsilon$, ker velja

$$\sum_{i=0}^n \alpha_i = \sum_{i=0}^n \int_a^b l_{i,n}(x)dx = \int_a^b 1dx = b-a.$$

Če to ne velja, pa imamo težave. Pri NC formulah vsota absolutnih vrednosti α_i hitro divergira za $n \rightarrow \infty$.

Vprašanje 70. Analiziraj neodstranljivo napako pri numeričnem integriranju.

Temu problemu se lahko ognemo tako, da razdelimo interval na manjše kose in integriramo na vsakem posebej. Temu pravimo SESTAVLJENO PRAVILO, sedaj pa imamo težavo, da moramo računati veliko vrednosti, tudi če tega ne potrebujemo na celotni domeni. Rešitev so adaptivne metode.

Oglejmo si adaptivno Simpsonovo metodo. Velja $I(f) = S_h(f) + R_h(f) = S_{h/2}(f) + R_{h/2}(f)$, s konkretnima napakama

$$R_h(f) = \frac{h^4(b-a)}{180}f^{(4)}(\xi_1), \quad R_{h/2}(f) = \frac{h^4(b-a)}{16 \cdot 180}f^{(4)}(\xi_2).$$

Pri predpostavki, da je $f^{(4)}(\xi_1) \approx f^{(4)}(\xi_2)$, dobimo

$$R_h(f) \approx 16R_{h/2}(f),$$

iz česar izpeljemo oceno za napako

$$R_{h/2}(f) \approx \frac{S_{h/2}(f) - S_h(f)}{15}.$$

Potem lahko dobimo ekstrapoliran približek

$$I(f) = S_{h/2}(f) + R_{h/2}(f) \approx \frac{16S_{h/2}(f) - S_h(f)}{15}.$$

Če želimo izračunati integral funkcije f med a in b , po (sestavljani) Simpsonovi formuli izračunamo S_h in $S_{h/2}$ ter preverimo, če je ocena za napako manjša od ε . Če je, končamo, sicer pa interval razdelimo na dva, in rekurzivno izračunamo integrala f na teh intervalih, pri čemer zahtevamo, da je napaka manjša od $\varepsilon/2$. Postopek se bo končal, ker se mera za napako zmanjšuje s faktorjem 16, zahtevana natančnost pa s faktorjem 2.

Vprašanje 71. Razloži adaptivno Simpsonovo metodo.

3.8.3 Gaussove kvadrature formule

Imejmo

$$I(f) = \int_a^b f(x)\rho(x)dx,$$

kjer je ρ nenegativna utež. Kvadratura formula ima obliko

$$I(f) = \sum_{i=0}^n \alpha_i f(x_i) + R(f),$$

kjer je

$$\alpha_i = \int_a^b l_{i,n}(x)\rho(x)dx.$$

Formula je vedno točna za polinome stopnje $\leq n$ za poljubno izbiro vozlov x_0, \dots, x_n . Če definiramo skalarni produkt

$$\langle f, g \rangle = \int_a^b f(x)g(x)\rho(x)dx$$

in uporabimo Gram-Schmittovo ortogonalizacijo na standardni polinomski bazi, dobimo polinome $\varphi_0, \varphi_1, \dots$, za katere velja $\langle \varphi_i, \varphi_j \rangle = \delta_{ij}$.

Izrek. Če so $\varphi_0, \varphi_1, \dots$ ortogonalni polinomi na $[a, b]$ z utežjo ρ , ima φ_k same realne enostavne ničle, ki vse ležijo v (a, b) .

3 Uvod v numerične metode

Dokaz. Naj ima φ_k v (a, b) l različnih ničel, kjer je $l < k$. Označimo jih z z_1, \dots, z_l , in definiramo

$$g(x) = (x - z_1)^{j_1} \cdots (x - z_l)^{j_l},$$

kjer je j_i enak 1, če je z_i liha ničla, in 0 sicer. Ker je integrand pozitiven, velja

$$\int_a^b g(x) \varphi_k(x) \rho(x) dx > 0.$$

Polinom g je manjše stopnje kot k , ker pa je φ_k pravokoten na vse polinome stopnje manjše od k , bi moral biti integral enak 0. To je protislovje. \square

Če za vozle vzamemo ničle polinoma φ_{n+1} , bo veljalo $\omega(x) = c\varphi_{n+1}(x)$, torej je ω pravokoten na vse polinome stopnje $\leq n$. Naj bo f polinom stopnje $\leq 2n+1$. Zapišemo ga lahko kot $f(x) = h(x)\omega(x) + g(x)$, kjer sta stopnji h in g manjši od n . Računamo

$$\int_a^b f(x) \rho(x) dx = \underbrace{\int_a^b h(x) \omega(x) \rho(x) dx}_{\langle h, \omega \rangle = 0} + \int_a^b g(x) \rho(x) dx = \sum_{i=0}^n \alpha_i g(x_i) = \sum_{i=0}^n \alpha_i f(x_i).$$

Formula je natančna za vse polinome stopnje $\leq n$, torej tudi za g ; sledi, da je natančna tudi za f . Izkaže se tudi, da so vse uteži pri Gaussovih formulah pozitivne, tako da ni problemov s stabilnostjo.

Vprašanje 72. Povej idejo za Gaussovimi kvadrturnimi formulami in dokaži, da so natančne za polinome stopnje $\leq 2n+1$.

3.9 Diferencialne enačbe

Numerična rešitev diferencialne enačbe je sestavljena iz zaporedja x_0, x_1, \dots in pripadajočih približkov y_0, y_1, \dots . Metode delimo na enokoračne, kjer y_{n+1} izračunamo iz y_n , ter večkoračne, kjer y_{n+1} izračunamo iz prejšnjih nekaj približkov. Poleg tega ločimo metode na eksplcitne, kjer imamo formulo za y_{n+1} , in implicitne, kjer y_{n+1} izračunamo z reševanjem nelinearne enačbe.

Če je $f = f(x, y)$ Lipschitzova v y s konstanto L , in je $y(x)$ točna rešitev začetnega problema $y' = f(x, y), y(x_0) = y_0$, ter $\tilde{y}(x)$ rešitev začetnega problema z zmotenim začetnim pogojem $y(x_0) = \tilde{y}(x_0)$, potem za poljuben $x \geq x_0$ velja

$$|\tilde{y}(x) - y(x)| \leq e^{L(x-x_0)} |\tilde{y}_0 - y_0|.$$

Podobno slabo mejo dobimo tudi, če zmotimo še f .

Najenostavnejša eksplcitna metoda je EKSPPLICITNA EULERJEVA METODA

$$\begin{aligned} y_{n+1} &= y_n + hf(x_n, y_n), \\ x_{n+1} &= x_n + h. \end{aligned}$$

Poznamo tudi implicitno obliko

$$y_{n+1} = y_n + hf(x_{n+1}, y_{n+1}).$$

3.9.1 Runge-Kutta metode

Pri Runge-Kutta metodah najprej izračunamo koeficiente

$$k_i = hf \left(x_n + \alpha_i h, y_n + \sum_{j=1}^m \beta_{ij} k_j \right)$$

za $i = 1, \dots, m$, nato pa

$$y_{n+1} = y_n + \sum_{i=1}^m \gamma_i k_i.$$

Pri tem je m stopnja metode, konstante $\alpha_i, \beta_{ij}, \gamma_i$ pa določimo tako, da se y_{n+1} čim bolj ujema z razvojem $y(x_n + h)$ v Taylorjevo vrsto. Veljati mora

$$\alpha_i = \sum_{j=1}^m \beta_{ij}.$$

Metoda je eksplisitna, če za $i \leq j$ velja $\beta_{ij} = 0$, sicer pa je implicitna. Stopnja metode m je različna od REDA metode k , ki je enak eksponentu v zadnjem členu Taylorjeve vrste, s katerim se metoda natančno ujema.

Vprašanje 73. Razloži Runge-Kutta metode za reševanje diferencialnih enačb.

4 Verjetnost

Komentar za učenje: poglej si tudi vserazne primere v zvezku, in jih poračunaj za vajo.

4.1 Izidi, dogodki, verjetnosti

Vprašanje 1. Kaj je množica Ω vseh možnih izidov? Povej nekaj primerov.

Odgovor: To je množica, ki hrani vse možne rezultate nekega poskusa. Pri mešanju kupa n kart velja $\Omega = S_n$, pri n -kratnem metu kovanca je to $\Omega = \{G, S\}^n$, itd. \square

Definicija. Družina \mathcal{F} podmnožic množice Ω je σ -ALGEBRA, če velja:

- $\Omega \in \mathcal{F}$,
- $A \in \mathcal{F} \implies A^c \in \mathcal{F}$,
- $A_1, A_2, \dots \in \mathcal{F} \implies \bigcup_i A_i \in \mathcal{F}$.

Definicija. Naj bo Ω množica možnih izidov, in \mathcal{F} σ -algebra nad Ω . VERJETNOST je preslikava $P : \mathcal{F} \rightarrow [0, 1]$, za katero velja $P(\Omega) = 1$, in kjer za disjunktne dogodke $A_1, A_2, \dots \in \mathcal{F}$ velja $P(\bigcup_i A_i) = \sum_i P(A_i)$.

Opomba. To sta aksioma Kolmogorova.

Vprašanje 2. Kaj je verjetnost?

Izrek (Formula za vključitve in izključitve). Naj bodo A_1, \dots, A_n dogodki. Potem velja

$$P\left(\bigcup_{i=1}^n A_i\right) = \sum_{i=1}^n P(A_i) - \sum_{1 \leq i < j \leq n} P(A_i \cap A_j) + \dots + (-1)^{n-1} P\left(\bigcap_{i=1}^n A_i\right).$$

Dokaz. Definirajmo dogodke

$$B_r = \{\omega \in \Omega \mid \omega \text{ je vsebovan v natanko } r \text{ množicah } A_i\}.$$

To so disjunktne dogodke, za katere velja $\bigcup_i A_i = \bigcup_r B_r$. Sledi

$$P\left(\bigcup_{i=1}^n A_i\right) = \sum_{r=1}^n P(B_r).$$

Poglejmo si, kolikokrat smo v formuli v izreku šteli vsako izmed množic B_r . Ta množica je vsebovana v preseku do r dogodkov, torej se v prvem členu pojavi r -krat, v drugem $\binom{r}{2}$, v tretjem $\binom{r}{3}$, itd. Vsota je tedaj

$$\binom{r}{1} - \binom{r}{2} + \binom{r}{3} - \dots + (-1)^{r-1} \binom{r}{r} = 1,$$

kar lahko izpeljemo iz razvoja izraza $0 = (1 - 1)^r$. \square

Vprašanje 3. Povej formulo za izključitve in izključitve. Kaj je ideja dokaza?

Lema. Naj bodo A_1, A_2, \dots dogodki. Če je $A_1 \subseteq A_2 \subseteq A_3 \subseteq \dots$, je verjetnost unije

$$P\left(\bigcup_{i=1}^{\infty} A_i\right) = \lim_{n \rightarrow \infty} P(A_n).$$

Če namesto tega velja $A_1 \supseteq A_2 \supseteq A_3 \supseteq \dots$, je

$$P\left(\bigcap_{i=1}^{\infty} A_i\right) = \lim_{n \rightarrow \infty} P(A_n).$$

Dokaz. Druga formula sledi iz De Morganovih pravil, dokažemo samo prvo. Zapišemo

$$\bigcup_{i=1}^{\infty} A_i = A_1 \cup (A_2 \setminus A_1) \cup (A_3 \setminus (A_1 \cup A_2)) \cup \dots$$

To so disjunktni dogodki, torej zanje velja

$$\begin{aligned} P\left(\bigcup_{i=1}^{\infty} A_i\right) &= P(A_1) + \sum_{k=2}^{\infty} P(A_k \setminus (A_1 \cup \dots \cup A_{k-1})) \\ &= \lim_{n \rightarrow \infty} \left(P(A_1) + \sum_{k=2}^n P(A_k \setminus (A_1 \cup \dots \cup A_{k-1})) \right) \\ &= \lim_{n \rightarrow \infty} P\left(\bigcup_{k=1}^n A_k \setminus (A_1 \cup \dots \cup A_{k-1})\right) \\ &= \lim_{n \rightarrow \infty} P(A_n). \end{aligned}$$

□

Lema (Prva Borel-Cantorjeva lema). Naj bodo A_1, A_2, \dots dogodki, za katere velja $\sum_i P(A_i) < \infty$. Definiramo $\bar{A} = \{\omega \in \Omega \mid \omega \text{ je vsebovan v neskončno mnogo } A_k\}$. Tedaj velja $P(\bar{A}) = 0$.

Dokaz. Prepričamo se lahko, da velja $\bar{A} = \bigcap_{n=1}^{\infty} \bigcup_{m=n}^{\infty} A_m$. Te unije so padajoče za $n \rightarrow \infty$, zato je po prejšnji lemi velja

$$P(\bar{A}) = \lim_{n \rightarrow \infty} P\left(\bigcup_{m=n}^{\infty} A_m\right).$$

Iz dokaza prešnje leme vidimo, da velja sklep

$$P\left(\bigcup_{k=1}^n A_k\right) \leq \sum_{k=1}^n P(A_k) \implies P\left(\bigcup_{k=1}^{\infty} A_k\right) \leq \sum_{k=1}^{\infty} P(A_k).$$

4 Verjetnost

Torej velja

$$P(\overline{A}) \leq \lim_{n \rightarrow \infty} \sum_{k=n}^{\infty} P(A_k).$$

Izraz na desni pa je rep konvergenčne vrste, torej je limita enaka 0. \square

Vprašanje 4. Povej in dokaži prvo Borel-Cantorjevo lemo.

4.1.1 Pogojna verjetnost in neodvisnost

Definicija. Naj bo B dogodek s $P(B) > 0$. **POGOJNA VERJETNOST** dogodka A glede na B je

$$P(A|B) = \frac{P(A \cap B)}{P(B)}.$$

Vprašanje 5. Kaj je pogojna verjetnost?

Primer (Bertrandov paradoks). Imamo tri škatle. V prvi sta dva zlatnika, v drugi zlatnik in srebrnik, in v zadnji dva srebrnika. Izberemo eno škatlo tako, da ima vsaka verjetnost $1/3$. Iz izbrane škatle tedaj naključno izberemo kovanec. Definiramo dogodka A , drugi kovanec v škatli je zlatnik, in B , izbrani kovanec je zlatnik. Z izpisom izidov izračunamo $P(A|B) = 2/3$.

Definicija. Družina dogodkov $\{H_1, \dots, H_n, \dots\}$ je **PARTICIJA** Ω , če je njihova unija enaka Ω in če so paroma disjunktni.

Vprašanje 6. Kaj je particija? Izpelji formulo za popolno verjetnost.

Odgovor: Za definicijo glej zgoraj. Naj bo A dogodek. Računamo

$$\begin{aligned} P(A) &= P(A \cap \Omega) \\ &= P(A \cap \bigcup_i H_i) \\ &= P(\bigcup_i A \cap H_i) \\ &= \sum_i P(A \cap H_i) \\ &= \sum_i \frac{P(A \cap H_i)}{P(H_i)} P(H_i) \\ &= \sum_i P(A|H_i) P(H_i). \end{aligned}$$

Če je $P(H_i) = 0$, lahko člen izpustimo. \boxtimes

4.1.2 Neodvisnost dogodkov

Definicija. Dogodki $\{A_i\}_{i \in I}$ so NEODVISNI, če za vsako končno poddružino A_1, A_2, \dots, A_n velja

$$P(A_1 \cap \dots \cap A_n) = P(A_1) \dots P(A_n).$$

Vprašanje 7. Kdaj so dogodki neodvisni?

Definicija. Družina dogodkov $\mathcal{P} = \{A_1, \dots, A_n\}$ je π -SISTEM, če za vsaka $A_i, A_j \in \mathcal{P}$ velja $A_i \cap A_j \in \mathcal{P}$.

Opomba. Če π -sistemu dodamo \emptyset in Ω , spet dobimo π -sistem.

Izrek. Če je $\mathcal{P} = \{B_1, \dots, B_n\}$ π -sistem in je A neodvisen od vseh B_k , je A neodvisen od vseh dogodkov, ki jih lahko sestavimo iz dogodkov v \mathcal{P} s komplementiranjem, preseki in unijami.

Dokaz. S preprostim izračunom lahko pokažemo, da če je A neodvisen od dogodkov C_1, \dots, C_m , ki so vsi disjunktni od A , je A neodvisen tudi od njihove unije. Poleg tega opazimo, da so vsi dogodki, ki jih sestavimo v izreku, končne unije dogodkov $B_1^* \cap \dots \cap B_m^*$, kjer je B_i^* bodisi enak B_i bodisi B_i^c .

V luči teh ugotovitev je dovolj dokazati, da je A neodvisen od vsakega dogodka $B_1^* \cap \dots \cap B_m^*$. Če izberemo vse dogodke, kjer ni komplementa, je presek v \mathcal{P} , zato jih lahko nadomestimo z enim samim. Brez škode za splošnost se torej omejimo na dogodke oblike $B_1^c \cap \dots \cap B_m^c \cap B_{m+1}$. Velja

$$P\left(A \cap \left(\bigcup_i B_i\right)^c \cap B_{m+1}\right) = P(A \cap B_{m+1}) - P\left(\left(\bigcup_i B_i\right) \cap A \cap B_{m+1}\right),$$

kjer smo uporabili pomožni sklep $P(A \cap B^c) = P(A) - P(A \cap B)$, ki ga izpeljemo iz dejstva $P(A) = P(A \cap B) + P(A \cap B^c)$. Zgornji izraz je nadalje enak

$$P(A)P(B_{m+1}) - P\left(\bigcup_i A \cap B_i \cap B_{m+1}\right),$$

ker sta A in B_{m+1} neodvisna. Drugi člen razvijemo po formuli za vključitve in izključitve in dobimo

$$\begin{aligned} P(A)P(B_{m+1}) - \sum_i P(A \cap B_i \cap B_{m+1}) + \sum_{i,j} P(A \cap B_i \cap B_j \cap B_{m+1}) \\ - \dots + (-1)^m P(A \cap B_1 \cap \dots \cap B_{m+1}). \end{aligned}$$

V vseh členih dobimo presek A z dogodkom v \mathcal{P} , torej lahko izpostavimo $P(A)$;

$$P(A) \left(P(B_{m+1}) - \sum_i P(B_i \cap B_{m+1}) + \dots \right).$$

4 Verjetnost

V drugem členu produkta smo dobili razvoj dogodka po formuli za vključitve in izključitve, ki ga lahko skrčimo v

$$P(A) \left(P(B_{m+1}) - P\left(\bigcup_i B_i \cap B_{m+1}\right) \right).$$

Nazadnje še uporabimo zgornji sklep v drugo smer in dobimo

$$P(A)P\left(B_{m+1} \left(\bigcup_i B_i\right)^c\right),$$

kar zaključimo dokaz. □

4.2 Slučajne spremenljivke in porazdelitve

Definicija. SLUČAJNA SPREMENLJIVKA X je funkcija $\Omega \rightarrow \mathbb{R}$, da je za $a < b$ množica $X^{-1}((a, b])$ dogodek v σ -algebri dogodkov \mathcal{F} .

Opomba. Ekvivalentno definicijo dobimo, če namesto polodprtih intervalov vzamemo odprte ali zaprte. Izbiro je predpisal ISO standard.

Opomba. Funkcija sama po sebi je popolnoma deterministična, naključna je izbira argumenta.

Definicija. Slučajna spremenljivka je DISKRETNA, če je njena zaloga vrednosti števna ali končna množica.

Definicija. PORAZDELITEV diskretne slučajne spremenljivke X z vrednostmi $(x_i)_i$ je dana z verjetnostmi $P(X^{-1}(x_i))$.

Vprašanje 8. Definiraj slučajne spremenljivke. Kdaj je slučajna spremenljivka diskretna? Kaj je porazdelitev?

Obstaja nekaj standardnih diskretnih porazdelitev.

Primer (Hipergeometrijska porazdelitev). Imamo posodo z B belimi in R rdečimi kroglicami. Označimo $N = B + R$ in naključno izberemo $n \leq N$ kroglic tako, da so vse podmnožice enako verjetne. Če z X označimo število izbranih belih kroglic, dobimo slučajno spremenljivko. Za $\max\{0, n - R\} \leq k \leq \min\{n, B\}$ je

$$P(X = k) = \frac{\binom{B}{k} \binom{R}{n-k}}{\binom{N}{n}}.$$

Na kratko označimo $X \sim \text{HiperGeom}(n, B, N)$.

Vprašanje 9. Opiši hipergeometrijsko porazdelitev.

Primer (Binomska porazdelitev). Kovanec z maso m vržemo n -krat zaporedoma, pri čemer so vsi meti medsebojno neodvisni, verjetnost grba pa je $p \in (0, 1)$. Naj bo X število grbov v teh n metih. Tedaj za $k = 0, \dots, n$ velja

$$P(X = k) = \binom{n}{k} p^k (1-p)^{n-k}.$$

Označimo $X \sim \text{Bin}(n, p)$.

Vprašanje 10. Opiši binomsko porazdelitev.

Primer (Geometrijska porazdelitev). Naj bo X število metov kovanca, potrebnih, da pade prvi grb. Pri tem so meti neodvisni, kovanec pade na grb z verjetnostjo p . Možne vrednosti za X so vsi $k \in \mathbb{N}$, velja

$$P(X = k) = (1-p)^{k-1} p.$$

Na kratko označimo $X \sim \text{Geom}(p)$.

Vprašanje 11. Opiši geometrijsko porazdelitev.

Primer (Negativna binomska porazdelitev). Mečemo kovanec in čakamo na m grbov; naj bo X število potrebnih metov. Možne vrednosti X so tedaj $k = m, m+1, \dots$, pri čemer velja

$$P(X = k) = \binom{k-1}{m-1} p^m (1-p)^{k-m}.$$

Oznaka je $X \sim \text{NegBin}(m, p)$.

Vprašanje 12. Opiši negativno binomsko porazdelitev.

Definicija. POCHHAMMERJEV SIMBOL $(a)_n$ je definiran kot

$$(a)_n = a(a+1) \dots (a+n-1).$$

Opomba. Izračunamo ga lahko tudi kot

$$(a)_n = \frac{\Gamma(a+n)}{\Gamma(a)}.$$

Primer (Poissonova porazdelitev). Oglejmo si dogajanje binomske porazdelitve, ko velja $np = \lambda$ konstanta, in ko $n \rightarrow \infty$. Tedaj

$$\begin{aligned} \lim_{n \rightarrow \infty} \binom{n}{k} \left(\frac{\lambda}{n}\right)^k \left(1 - \frac{\lambda}{n}\right)^{n-k} &= \lim_{n \rightarrow \infty} \frac{\lambda^k}{k!} \frac{n(n-1) \dots (n-k+1)}{n^k} \left(1 - \frac{\lambda}{n}\right)^n \left(1 - \frac{\lambda}{n}\right)^{-k} \\ &= \frac{\lambda^k}{k!} e^{-\lambda}. \end{aligned}$$

4 Verjetnost

Če je za $k = 0, 1, \dots$ verjetnost

$$P(X = k) = e^{-\lambda} \frac{\lambda^k}{k!},$$

pravimo, da ima X Poissonovo porazdelitev, in označimo $X \sim \text{Po}(\lambda)$.

Vprašanje 13. Opiši Poissonovo porazdelitev.

Definicija. Porazdelitev zvezne slučajne spremenljivke je podana z verjetnostmi $P(X \in (a, b])$ za $a < b$.

Opomba. Pogosto želimo izračunati $P(X \in A)$, kjer $A \subseteq \mathbb{R}$ ni interval. V tem primeru lahko verjetnost izračunamo, če je A sestavljena iz števnih unij, števnih presekov in komplementov polodprtih intervalov. Taki družini pravimo BORELOVE MNOŽICE, in jo označimo z $B(\mathbb{R})$. Tehnično so to najmanjša σ -algebra na \mathbb{R} , ki vsebuje vse polodprte intervale.

Definicija. Slučajna spremenljivka X ima ZVEZNO PORAZDELITEV, če obstaja nenegativna funkcija $f_X : \mathbb{R} \rightarrow [0, \infty)$, da je

$$P(X \in (a, b]) = \int_a^b f_X(x) dx.$$

Funkciji f_X pravimo GOSTOTA PORAZDELITVE.

Vprašanje 14. Definiraj zvezno porazdelitev in gostoto porazdelitve.

Primer (Normalna porazdelitev). Pravimo, da ima X normalno porazdelitev s parametroma μ in σ^2 , če je gostota enaka

$$f_X(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right).$$

Pri tem je σ razdalja od μ do prevoja, μ pa središče porazdelitve.

Vprašanje 15. Opiši normalno porazdelitev.

Primer (Eksponentna porazdelitev). Pravimo, da ima X eksponentno porazdelitev s parametrom λ , če velja

$$f_X(x) = \begin{cases} \lambda e^{-\lambda x} & x \geq 0, \\ 0 & x < 0. \end{cases}$$

Označimo z $X \sim \exp(\lambda)$.

Primer (Gama porazdelitev). Slučajna spremenljivka X ima gama porazdelitev s parametroma $a, \lambda > 0$, če je gostota enaka

$$f_X(x) = \begin{cases} \frac{\lambda^a}{\Gamma(a)} x^{a-1} e^{-\lambda x} & x \geq 0, \\ 0 & x < 0. \end{cases}$$

Oznaka: $X \sim \Gamma(a, \lambda)$.

Primer (Enakomerna porazdelitev). Predvidljivo je

$$f_X(x) = \begin{cases} \frac{1}{b-a} & a \leq x \leq b, \\ 0 & \text{sicer.} \end{cases}$$

Oznaka: $X \sim U(a, b)$.

Primer (Beta porazdelitev). Spremenljivka X ima beta porazdelitev, če je

$$f_X(x) = \begin{cases} \frac{1}{B(a,b)} x^{a-1} (1-x)^{b-1} & 0 < x < 1, \\ 0 & \text{sicer.} \end{cases}$$

Oznaka: $X \sim \text{Beta}(a, b)$.

Vprašanje 16. Opiši eksponentno, gama, enakomerno in beta porazdelitev.

Če je X slučajna spremenljivka, kakšna mora biti funkcija $f : \mathbb{R} \rightarrow \mathbb{R}$, da bo $Y = f(X)$ tudi slučajna spremenljivka? Za poljubna $a < b$ mora biti $X^{-1}(f^{-1}((a, b)))$ dogodek. Če je funkcija (odsekoma) zvezna, že zadošča; potreben in zadosten pogoj pa je, da je f merljiva, torej da je $f^{-1}(A)$ dogodek za vse $A \in \mathcal{B}(\mathbb{R})$.

Definicija. PORAZDELITVENA FUNKCIJA slučajne spremenljivke X je funkcija $F_X : \mathbb{R} \rightarrow \mathbb{R}$, podana z $F_X(x) = P(X \leq x)$.

Če ima X gostoto f_X , je

$$F_X(x) = \int_{-\infty}^x f_X(u) du.$$

Izrek. Naj bo F_X porazdelitvena funkcija slučajne spremenljivke X . Tedaj velja

- F_X je nepadajoča,
- $\lim_{x \rightarrow \infty} F_X(x) = 1$ in $\lim_{x \rightarrow -\infty} F_X(x) = 0$,
- F_X je desno zvezna.

Dokaz. Prva točka: za $x < y$ velja $F_X(y) - F_X(x) = P(X \in (x, y]) \geq 0$.

Druga točka: Definiramo $A_n = \{X \leq n\}$. Tedaj velja

$$\bigcup_{n=1}^{\infty} A_n = \Omega,$$

in ti dogodki so naraščajoči. Potem je

$$1 = P(\Omega) = P\left(\bigcup_n A_n\right) = \lim_{n \rightarrow \infty} P(A_n) = \lim_{n \rightarrow \infty} F_X(n).$$

Ker je F_X nepadajoča, velja tudi $\lim_{x \rightarrow \infty} F_X(x) = 1$ zvezno. Za drugo formulo podobno definiramo $B_n = \{X \leq -n\}$, kar je padajoče zaporedje dogodkov s praznim presekom. Za limito velja podoben sklep kot prej.

Tretja točka: Naj bo $x_n \downarrow x$. Definiramo $C_n = \{X \leq x_n\}$, velja

$$\bigcap_{n=1}^{\infty} C_n = \{X \leq x\}.$$

Ker so C_n padajoči, je

$$F_X(x) = P\left(\bigcap_{n=1}^{\infty} C_n\right) = \lim_{n \rightarrow \infty} P(X \leq x_n) = \lim_{n \rightarrow \infty} F_X(x_n),$$

torej je F_X res desno zvezna. \square

Vprašanje 17. Definiraj porazdelitveno funkcijo slučajne spremenljivke. Kakšne lastnosti ima?

Vprašanje 18. Naj velja $X \sim N(\mu, \sigma^2)$ in $Y = aX + b$. Kakšna je gostota Y ?

Odgovor: Računamo

$$F_Y(y) = P(Y \leq y) = P(X \leq \frac{y-b}{a}) = \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{(y-b)/a} \exp\left(-\frac{(u-\mu)^2}{2\sigma^2}\right) du.$$

Ker je F_X zvezno odvedljiva, je

$$f_Y(y) = F'_Y(y) = \frac{1}{\sqrt{2\pi}\sigma a} \exp\left(-\frac{(y-b-a\mu)^2}{2a^2\sigma^2}\right),$$

torej $Y \sim N(a\mu + b, a^2\sigma^2)$. \boxtimes

Definicija. Če je $Z \sim N(0, 1)$, rečemo, da ima Z STANDARDIZIRANO NORMALNO PORAZDELITEV. Porazdelitveno funkcijo Z označimo s ϕ , torej

$$\phi(z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z \exp\left(-\frac{1}{2}u^2\right) du.$$

Če je $X \sim N(\mu, \sigma^2)$, je torej

$$F_X(x) = \phi\left(\frac{x-\mu}{\sigma}\right).$$

Vprašanje 19. Kaj je standardizirana normalna porazdelitev?

Vprašanje 20. Kaj je verjetnostna transformacija?

Odgovor: Naj bo X zvezno porazdeljena s porazdelitveno funkcijo F_X , za katero predpostavimo, da je zvezna. Definiramo $Y = F_X(X)$ in računamo za $y \in (0, 1)$

$$P(Y \leq y) = P(X \in F^{-1}((-\infty, y])) = F_X(\sup\{x \mid F_X(x) \leq y\}) = y,$$

torej $Y \sim U(0, 1)$. \boxtimes

Definicija. Naj bo $p \in (0, 1)$. Vsakemu številu x_p , za katerega je $P(X \leq x_p) = p$, rečemo p -TI KVANTIL porazdelitve slučajne spremenljivke X .

Opomba. p -ti kvantil ni enolično določen.

Vprašanje 21. Kaj je p -ti kvantil porazdelitve slučajne spremenljivke?

4.2.1 Slučajni vektorji

Primer (Multinomska porazdelitev). Imamo r škatel, vanje mečemo n kroglic. Meti so neodvisni, škatlo k zadenemo z verjetnostjo p_k . Velja $\sum_k p_k = 1$. V vsaki škatli je pristalo slučajno število kroglic $X_k \sim \text{Bin}(n, p_k)$. Te spremenljivke zložimo v vektor $\underline{X} = (X_1, \dots, X_r)$, ki ima porazdelitev

$$P(X_1 = k_1, \dots, X_r = k_r) = p_1^{k_1} \dots p_r^{k_r} \binom{n}{k_1, \dots, k_r}.$$

Pravimo, da ima \underline{X} multinomsko porazdelitev s parametroma n in \underline{p} , ter označimo $\underline{X} \sim \text{Multinom}(n, \underline{p})$.

Vprašanje 22. Kaj je multinomska porazdelitev?

Definicija. SLUČAJNI VEKTOR \underline{X} s komponentami X_1, \dots, X_r je funkcija $\underline{X} : \Omega \rightarrow \mathbb{R}^r$, da je

$$\underline{X}^{-1} \left(\prod_{k=1}^r (a_k, b_k] \right)$$

dogodek za vse $a_k < b_k$.

Definicija. Slučajni vektor je DISKRETEN, če ima vrednosti v končni ali števni množici.

Vprašanje 23. Definiraj slučajne vektorje.

Izrek. Naj bosta \underline{X} in \underline{Y} diskretna slučajna vektorja. Za vse možne vrednosti \underline{x} vektorja \underline{X} velja

$$P(\underline{X} = \underline{x}) = \sum_{\underline{y}} P(\underline{X} = \underline{x}, \underline{Y} = \underline{y}).$$

Formuli pravimo FORMULA ZA ROBNO PORAZDELITEV.

Vprašanje 24. Povej formulo za robno porazdelitev.

4.2.2 Neodvisnost slučajnih spremenljivk

Definicija. Slučajni spremenljivki X in Y sta NEODVISNI, če za vsaki Borelovi A in B velja $P(X \in A, Y \in B) = P(X \in A)P(Y \in B)$.

Definicija. Slučajne spremenljivke X_1, \dots, X_n so NEODVISNE, če za vsak nabor Borelovih množic A_1, \dots, A_n velja

$$P(\forall i. X_i \in A_i) = \prod_j P(X_j \in A_j).$$

Definicija. Slučajne spremenljivke $\{X_i\}_{i \in I}$ so NEODVISNE, če so neodvisne vse končne poddružine.

Vprašanje 25. Definiraj neodvisnost slučajnih spremenljivk.

Izrek. Naj za diskretni slučajni spremenljivki X in Y velja $P(X = x, Y = y) = f(x)g(y)$ za funkciji $f : R(X) \rightarrow \mathbb{R}$ in $g : R(Y) \rightarrow \mathbb{R}$ (R je zaloga vrednosti). Potem sta X in Y neodvisni.

Dokaz. Po formuli za robne porazdelitve je

$$P(X = x) = f(x) \sum_y g(y) = f(x)C_1.$$

Podobno $P(Y = y) = C_2g(y)$. Predpišemo

$$P(X = x, Y = y) = P(X = x)P(Y = y)C_1^{-1}C_2^{-1}.$$

Dokazati moramo še, da velja $C_1C_2 = 1$. Seštejmo

$$\begin{aligned} 1 &= \sum_{x,y} P(X = x, Y = y) \\ &= \frac{1}{C_1C_2} \sum_{x,y} P(X = x)P(Y = y) \\ &= \frac{1}{C_1C_2} \left(\sum_x P(X = x) \right) \left(\sum_y P(Y = y) \right) \\ &= \frac{1}{C_1C_2}. \end{aligned}$$

□

Vprašanje 26. Kako še lahko določiš, da sta slučajni spremenljivki neodvisni? Dokaži.

4.2.3 Pričakovana vrednost diskretnih spremenljivk

Definicija. Naj bo X diskretna slučajna spremenljivka z vrednostmi x_1, x_2, \dots . PRIČAKOVANA VREDNOST $E(X)$ je število, dano z

$$E(X) = \sum_i x_i P(X = x_i).$$

Če slučajno spremenljivko X vstavimo v funkcijo, spet dobimo slučajno spremenljivko $Y = f(X)$. Če je X diskretna, je taka tudi Y , torej

$$E(Y) = \sum_y yP(Y = y) = \sum_x f(x)P(X = x).$$

Vprašanje 27. Definiraj pričakovano vrednost diskretne spremenljivke. Kako se preslika s funkcijo?

Izrek. Naj bodo X_1, \dots, X_n slučajne spremenljivke in $\alpha_1, \dots, \alpha_n$ konstante. Če obstaja $E(X_i)$ za $i = 1, \dots, n$, obstaja tudi

$$E\left(\sum_i \alpha_i X_i\right) = \sum_i \alpha_i E(X_i).$$

Definicija. Slučajna spremenljivka I ima BERNOULLIJEVO PORAZDELITEV, če je njena zaloga vrednosti enaka $\{0, 1\}$. Če označimo $p = P(I = 1)$, pišemo $I \sim \text{Bernoulli}(p)$.

Vprašanje 28. Kaj je Bernoullijeva porazdelitev?

4.2.4 Večrazsežne zvezne porazdelitve

Definicija. Slučajni vektor \underline{X} ima ZVEZNO PORAZDELITEV, če obstaja nenegativna funkcija $f_{\underline{X}}(\underline{x})$, da za $A \subseteq \mathbb{R}^n$ velja

$$P(\underline{X} \in A) = \int_A f_{\underline{X}}(\underline{x}) d\underline{x}.$$

Funkciji $f_{\underline{X}}$ pravimo GOSTOTA.

Vprašanje 29. Definiraj gostoto porazdelitve slučajnega vektorja.

Izrek. Naj bo $f_{\underline{X}}(\underline{x})$ gostota vektorja \underline{X} in $m < n$. Privzemimo, da je funkcija

$$(x_1, \dots, x_n) \mapsto \int_{\mathbb{R}^{n-m}} f_{\underline{X}}(x_1, \dots, x_n) dx_{m+1} \dots dx_n$$

Riemannovo integrabilna (lahko tudi v izlimitiranem smislu) po vseh Jordanovo izmerljivih množicah. Potem je to funkcija gostote vektorja $\underline{X}' = (x_1, \dots, x_m)$.

Izrek. Slučajna vektorja $\underline{X}, \underline{Y}$ sta neodvisna natanko tedaj, ko je

$$f_{\underline{X}, \underline{Y}}(\underline{x}, \underline{y}) = f_{\underline{X}}(\underline{x}) f_{\underline{Y}}(\underline{y})$$

skoraj povsod.

4 Verjetnost

Dokaz. V desno: Velja

$$P(X \in A, Y \in B) = \int_{A \times B} f_{\underline{X}, \underline{Y}}(\underline{x}, \underline{y}) d\underline{x} d\underline{y},$$

$$P(X \in A)P(Y \in B) = \int_A f_{\underline{X}}(\underline{x}) d\underline{x} \int_B f_{\underline{Y}}(\underline{y}) d\underline{y} = \int_{A \times B} f_{\underline{X}} f_{\underline{Y}}.$$

Ker sta vektorja neodvisna, sta ti količini enaki, torej sta integrirani funkciji enaki skoraj povsod.

V levo: Velja

$$P(\underline{X} \in A, \underline{Y} \in B) = \int_{A \times B} f_{\underline{X}}(\underline{x}) f_{\underline{Y}}(\underline{y}) d\underline{x} d\underline{y}$$

$$= \int_A f_{\underline{X}}(\underline{x}) d\underline{x} \int_B f_{\underline{Y}}(\underline{y}) d\underline{y}$$

$$= P(\underline{X} \in A)P(\underline{Y} \in B).$$

□

Vprašanje 30. Karakteriziraj neodvisnost zvezno porazdeljenih slučajnih vektorjev in dokaži karakterizacijo.

Izrek. Naj bo $f_{\underline{X}, \underline{Y}}(\underline{x}, \underline{y}) = f(\underline{x})g(\underline{y})$ za nenegativni f, g . Potem sta \underline{X} in \underline{Y} neodvisni.

Dokaz je praktično enak kot pri podobnem izreku za diskretne spremenljivke, le da pišemo integrale namesto vsot.

Izrek. Naj bo \underline{X} slučajni vektor z gostoto $f_{\underline{X}}(\underline{x})$. Predpostavimo $P(\underline{X} \in U) = 1$ za neko odprto množico $U \subseteq \mathbb{R}^n$. Naj bo $\phi : U \rightarrow V$ difeomorfizem in $\underline{Y} = \phi(\underline{X})$. Velja $P(\underline{Y} \in V) = 1$ in

$$f_{\underline{Y}}(\underline{y}) = \begin{cases} f_{\underline{X}}(\phi^{-1}(\underline{y})) |\det D\phi^{-1}(\underline{y})| & \underline{y} \in V \\ 0 & \underline{y} \notin V \end{cases}$$

Vprašanje 31. Povej transformacijsko formulo.

Definicija. Naj bo X zvezno porazdeljena slučajna spremenljivka z gostoto $f_X(x)$. Definiramo

$$E(X) = \int_{-\infty}^{\infty} x f_X(x) dx,$$

$$E(g(x)) = \int_{-\infty}^{\infty} g(x) f_X(x) dx.$$

Vprašanje 32. Kako je definirana pričakovana vrednost zvezne slučajne spremenljivke?

Definicija. Za slučajno spremenljivko X imenujemo količino $E(X^m)$ m -TI MOMENT.

Definicija. Za slučajno spremenljivko X imenujemo količino $E((X - E(x))^m)$ m -TI CENTRALNI MOMENT.

Vprašanje 33. Kaj je moment in kaj centralni moment slučajne spremenljivke?

Definicija. Če imamo množico števil x_1, \dots, x_n , njihov RAZTROS definiramo kot

$$\sigma^2 = \frac{1}{n} \sum_{k=1}^n (x_k - \bar{x})^2,$$

kjer je \bar{x} povprečje.

Definicija. VARIANCA slučajne spremenljivke je količina

$$\text{var}(X) = E(X^2) - (E(X))^2.$$

Pravimo, da varianca obstaja, če obstajata obe pričakovani vrednosti.

Definicija. Naj bo X slučajna spremenljivka, za katero $\text{var}(X)$ obstaja. Količini $\sqrt{\text{var}(X)}$ rečemo STANDARDNI ODKLON slučajne spremenljivke X in jo označimo s $\text{SD}(X)$.

Definicija. KOVARIANCA dveh slučajnih spremenljivk je $\text{cov}(X, Y) = E(XY) - E(X)E(Y)$. Rečemo, da obstaja, če obstajajo vse pričakovane vrednosti.

Vprašanje 34. Definiraj raztros, varianco, standardni odklon in kovarianco.

Izrek. Naj bodo X_1, \dots, X_m in Y_1, \dots, Y_n slučajne spremenljivke in $\alpha_1, \dots, \alpha_m, \beta_1, \dots, \beta_n$ konstante. Velja

$$\text{cov}\left(\sum_{k=1}^m \alpha_k X_k, \sum_{l=1}^n \beta_l Y_l\right) = \sum_{k,l} \alpha_k \beta_l \text{cov}(X_k, Y_l).$$

Vprašanje 35. Povej nekaj lastnosti kovariance.

Odgovor:

- bilinearnost,
- $\text{cov}(X, X) = \text{var}(X)$,
- $\text{cov}(X, Y) = \text{cov}(Y, X)$,
- če sta X in Y neodvisni, je $\text{cov}(X, Y) = 0$.

☒

Definicija. Količina

$$\rho = \frac{\text{cov}(X, Y)}{\sqrt{\text{var}(X) \text{var}(Y)}}$$

se imenuje KORELACIJSKI KOEFICIENT.

Vprašanje 36. Definiraj korelacijski koeficient.

4.2.5 Pogojne pričakovane vrednosti

Definicija. POGOJNA PORAZDELITEV diskretne slučajne spremenljivke X glede na dogodek B je dana z verjetnostjo $P(X = x | B)$. POGOJNA PRIČAKOVANA VREDNOST slučajne spremenljivke X glede na dogodek B je dana z

$$E(X | B) = \sum_x xP(X = x | B).$$

Alternativno lahko izračunamo

$$E(X \cdot I_B) = \sum_x xP(X \cdot I_B = x) = P(B) \sum_x xP(X = x | B) = P(B)E(X | B),$$

torej $E(X | B) = E(X \cdot I_B)/P(B)$. Iz te formule sledi linearnost pogojne pričakovane vrednosti.

Vprašanje 37. Definiraj pogojno pričakovano vrednost in izpelji alternativno izražavo.

Izrek. Naj bo $\{H_1, H_2, \dots\}$ particija. Naj bo X diskretna slučajna spremenljivka. Velja

$$E(X) = \sum_k E(X | H_k)P(H_k).$$

Dokaz. Računamo

$$\sum_k E(X | H_k)P(H_k) = \sum_k \sum_x xP(X = x | H_k)P(H_k) = \sum_x P(X = x) = E(X).$$

□

Vprašanje 38. Dokaži formulo za popolno pričakovano vrednost.

Definicija. Naj bosta $\underline{X}, \underline{Y}$ slučajna vektorja z gostoto $f_{\underline{X}, \underline{Y}}(\underline{x}, \underline{y})$. Za \underline{x} , za katere je $f_{\underline{X}}(\underline{x}) > 0$, definiramo POGOJNO GOSTOTO \underline{Y} glede na $\{\underline{X} = \underline{x}\}$ kot

$$f_{\underline{Y} | \underline{X} = \underline{x}}(\underline{y}) = \frac{f_{\underline{X}, \underline{Y}}(\underline{x}, \underline{y})}{f_{\underline{X}}(\underline{x})}.$$

Vprašanje 39. Definiraj pogojno gostoto.

Izrek. Naj bosta $\underline{X}, \underline{Y}$ slučajna vektorja. Za $p = \dim \underline{X}$ velja

$$E(g(\underline{Y})) = \int_{\mathbb{R}^p} E(g(\underline{Y}) | \underline{X} = \underline{x}) f_{\underline{X}}(\underline{x}) d\underline{x}.$$

Dokaz. Računamo

$$\begin{aligned}
 \int_{\mathbb{R}^p} E(g(\underline{Y}) \mid \underline{X} = \underline{x}) f_{\underline{X}}(\underline{x}) d\underline{x} &= \int_{\mathbb{R}^p} f_{\underline{X}}(\underline{x}) d\underline{x} \int_{\mathbb{R}^q} g(\underline{y}) f_{\underline{Y} \mid \underline{X}=\underline{x}}(\underline{y}) d\underline{y} \\
 &= \int_{\mathbb{R}^{p+q}} g(\underline{y}) f_{\underline{X}}(\underline{x}) f_{\underline{Y} \mid \underline{X}=\underline{x}}(\underline{y}) d\underline{x} d\underline{y} \\
 &= \int_{\mathbb{R}^{p+q}} g(\underline{y}) f_{\underline{X}, \underline{Y}}(\underline{x}, \underline{y}) d\underline{x} d\underline{y} \\
 &= E(g(\underline{Y})).
 \end{aligned}$$

□

Vprašanje 40. Dokaži zvezno verzijo formule za popolno pričakovano vrednost.

4.3 Rodovne funkcije

Definicija. RODOVNA FUNKCIJA nenegativne celoštevilске slučajne spremenljivke X je dana z $G_X(s) = \sum_{k=0}^{\infty} s^k P(X = k)$.

Potenčna vrsta konvergira enakomerno na $[-1, 1]$, ker je tam dominirana z vrsto $\sum_{k=0}^{\infty} P(X = k) = 1$. Torej je $G_X(s)$ zvezna na $[-1, 1]$. Na intervalu $(-1, 1)$ je vrsta neskončnokrat odvedljiva. Iz formule $E(f(x)) = \sum_k f(k)P(X = k)$ za $f(x) = s^x$ dobimo

$$E(s^X) = \sum_{k=0}^{\infty} s^k P(X = k) = G_X(s).$$

Izrek. Naj bosta X, Y neodvisni spremenljivki. Velja $G_{X+Y}(s) = G_X(s)G_Y(s)$.

Dokaz. Sledi iz dejstva $E(s^{X+Y}) = E(s^X s^Y)$.

□

Vprašanje 41. Povej nekaj lastnosti rodovnih funkcij.

Izrek. Naj bo X slučajna spremenljivka. Potem velja

- $E(X) = \lim_{s \uparrow 1} G'_X(s)$
- $E(X(X-1)(X-2)\dots(X-k+1)) = \lim_{s \uparrow 1} G_X^{(k)}(s)$

Dokaz. Samo za prvo točko. Opazimo, da ima G_X pozitivne koeficiente, torej imajo vsi odvodi pozitivne koeficiente. Sledi, da so odvodi na $(0, 1)$ povsod naraščajoči.

Za začetek predpostavimo $E(X) < \infty$. Za $s \in (0, 1)$ velja

$$\sum_{k=1}^N kP(X = k)s^{k-1} \leq G'_X(s) \leq E(X),$$

kjer prva neenakost velja, ker je vrsta na levi glava vrste za G'_X . Funkcija G'_X je na $(0, 1)$ naraščajoča in omejena, torej limita $\lim_{s \uparrow 1} G'_X(s)$ obstaja. Levi člen v zgornji neenakosti konvergira k $E(X)$ za $N \rightarrow \infty$, torej trditev velja po izreku o sendviču.

Če je $E(X) = \infty$, potem levi člen v zgornji neenačbi divergira za $N \rightarrow \infty$, in je posledično tudi

$$\lim_{s \uparrow 1} G'_X(s) = \infty.$$

□

Vprašanje 42. Kako z odvodom rodovne funkcije izračunaš pričakovano vrednost? Dokaži.

Izrek. Naj bodo N, X_1, \dots neodvisne celoštevilске slučajne spremenljivke in naj bodo X_1, X_2, \dots enako porazdeljene. Naj bo $X = X_1 + \dots + X_N$. Potem velja $G_X(s) = G_N(G_{X_i}(s))$.

Dokaz. Računamo

$$\begin{aligned} G_X(s) &= E(s^X) \\ &= \sum_{n=0}^{\infty} E(s^X \mid N = n) P(N = n) \\ &= \sum_{n=0}^{\infty} E(s^{X_1 + \dots + X_n} \mid N = n) P(N = n). \end{aligned}$$

Na tej točki upoštevamo neodvisnost spremenljivk in dobimo

$$\begin{aligned} G_X(s) &= \sum_{n=0}^{\infty} E(s^{X_1 + \dots + X_n}) P(N = n) \\ &= \sum_{n=0}^{\infty} E(s^{X_1}) \dots E(s^{X_n}) P(N = n) \\ &= \sum_{n=0}^{\infty} (G_{X_1}(s))^n P(N = n) \\ &= G_N(G_{X_1}(s)). \end{aligned}$$

□

Vprašanje 43. Kako izračunaš porazdelitev vsote slučajno mnogo slučajnih spremenljivk? Dokaži.

4.3.1 Procesi razvejanja

Predpostavimo, da so porazdelitve števila sinov za vsakega posameznika enake, da so generacije sočasne in da so števila sinov medsebojno neodvisne. Kakšna je verjetnost, da kraljeva rodbina izumre?

Označimo z Z_n število posameznikov v n -ti generaciji. Naj bodo $\xi_{n,k}$ vse neodvisne z rodovno funkcijo G . Definiramo rekurzivno $Z_0 = 1$ in

$$Z_{n+1} = \xi_{n+1,1} + \dots + \xi_{n+1,Z_n}.$$

Zaradi predpostavke o neodvisnosti so izpolnjene vse predpostavke zadnjega izreka in lahko izračunamo za $G_n = G_{Z_n}$

$$G_{n+1}(s) = G_n(G(s)).$$

Ker je $G_1 = G$, dobimo $G_k = G \circ G \circ \dots \circ G$. Velja

$$\{\text{rodbina izumre}\} = \bigcup_{n=1}^{\infty} \{Z_n = 0\}$$

in ti dogodki so naraščajoči, torej je

$$\eta = P(\text{rodbina izumre}) = \lim_{n \rightarrow \infty} P(Z_n = 0).$$

Ker je $P(Z_n = 0) = G_n(0)$, dobimo

$$\eta = \lim_{n \rightarrow \infty} G_n(0) = \lim_{n \rightarrow \infty} G_{n+1}(0) = G(\lim_{n \rightarrow \infty} G_n(s)) = G(\eta).$$

Ker je $s = 1$ vedno fiksna točka G na $[0, 1]$, ni pa nujno edina. Naj bo $\bar{\eta}$ poljubna fiksna točka na $[0, 1]$. Funkcija G je nepadajoča na $[0, 1]$, torej $G(0) \leq G(\bar{\eta}) = \bar{\eta}$. Z večkratno uporabo G na tej neenakosti dobimo $G_n(0) \leq \bar{\eta}$, iz česar sledi $\eta = \lim_{n \rightarrow \infty} G_n(0) \leq \bar{\eta}$ in je η najmanjša fiksna točka.

Vprašanje 44. Kako dobiš verjetnost izumrtja procesa razvejanja? Dokaži.

4.3.2 Panjerjeva rekurzija

Definicija. Celoštevilska slučajna spremenljivka N ima porazdelitev PANJERJEVEGA RAZREDA, če obstajata konstanti a, b , da je

$$P(N = n) = \left(a + \frac{b}{n}\right) P(N = n - 1).$$

Binomska, Poissonova in negativna binomska so vse v tem razredu. Zanima nas vsota $X = X_1 + \dots + X_N$. Po dolgem računanju pridemo do formule

$$P(X = n + 1) = \frac{1}{1 - aP(X_1 = 0)} \sum_{k=1}^{n+1} \left(a + \frac{bk}{n+1}\right) P(X_1 = k) P(X = n - k + 1).$$

Vprašanje 45. Povej formulo za Panjerjevo rekurzijo.

4.4 Tabele

Porazdelitev X	X^2
$N(0, 1)$	$\Gamma\left(\frac{1}{2}, \frac{1}{2}\right)$

Tabela 4.1: Porazdelitve

Porazdelitev X	$E(X)$	$E(X^2)$	Varianca
$\text{Bin}(n, p)$	np	$n^2p^2 + np(1-p)$	npq
$\text{NegBin}(n, p)$	$\frac{n}{p}$	$\frac{m(m+1)}{p^2} - \frac{m}{p}$	$\frac{mq}{p^2}$
$N(\mu, \sigma^2)$	μ	$\sigma^2 + \mu^2$	σ^2
$\Gamma(a, \lambda)$	$\frac{a}{\lambda}$	$\frac{a(a+1)}{\lambda^2}$	$\frac{a}{\lambda^2}$

Tabela 4.2: Pričakovane vrednosti

Porazdelitev X	Porazdelitev Y	Porazdelitev $X + Y$
$\text{Po}(\lambda)$	$\text{Po}(\mu)$	$\text{Po}(\lambda + \mu)$
$\text{Bin}(m, p)$	$\text{Bin}(n, p)$	$\text{Bin}(m + n, p)$
$\text{Polya}(a, \beta)$	$\text{Polya}(b, \beta)$	$\text{Polya}(a + b, \beta)$
$N(\mu, \sigma^2)$	$N(\nu, \tau^2)$	$N(\mu + \nu, \sigma^2 + \tau^2)$
$\Gamma(a, \lambda)$	$\Gamma(b, \lambda)$	$\Gamma(a + b, \lambda)$

Tabela 4.3: Vsote

Porazdelitev $\text{Polya}(a, \beta)$ je dana z

$$P(X = k) = \frac{\beta^a (a)_k}{k! (1 + \beta)^{a+k}}.$$

5 Algebra 3

5.1 Reševanje polinomskih enačb

Vemo, da za polinomsko enačbo $x^n + a_{n-1}x^{n-1} + \dots + a_0 = 0$ obstaja razširitev polja \mathbb{F} , v kateri je enačba razcepna. Rešitev polinomske enačbe se izraža z radikali, če se da zapisati rešitve enačbe s pomočjo računskih operacij in korenov.

Definicija. Naj bo K polje. Razširitvi oblike $K(\sqrt[n]{a})/K$ za $a \in K$ pravimo **RADIKALSKA RAZŠIRITEV** polja K .

Polinomska enačba $p(X) = 0$ je rešljiva z radikali natanko tedaj, ko rešitve živijo v neki razširitvi E/F , pri čemer obstaja končna veriga razširitev $F \subseteq E_0 \subseteq \dots \subseteq E_k = E$, kjer so zaporedne razširitve radikalske. Problem se prevede na vprašanje iz teorije grup. Od tu naprej predpostavimo, da so vse razširitve končne.

Definicija. Naj bo E/F končna razširitev. To je **NORMALNA RAZŠIRITEV**, če za vsak nerazcepen polinom $p(X) \in F[X]$ velja ena od naslednjih možnosti:

- p nima ničle v E
- p ima vse ničle v E

Ekvivalentno: vsak nerazcepen polinom $p(X) \in F[X]$ z vsaj eno ničlo v E razpade v linearne faktorje v $E[X]$.

Primer. $\mathbb{Q}(\sqrt[3]{2})$ nad \mathbb{Q} ni normalna. Polinom $p(X) = X^3 - 2$ ne razpade.

Primer. $\mathbb{Q}(\sqrt{2})$ je normalna.

Vprašanje 1. Definiraj radikalske in normalne razširitve. Povej primer normalne razširitve in razširitve, ki ni normalna.

Izrek. Naj bo E/F končna razširitev. Potem je ta razširitev normalna natanko tedaj, ko je E razpadno polje nekega polinoma s koeficienti iz F .

Dokaz. V desno: Recimo, da je E normalna. Naj bo $p_i(X)$ minimalni polinom a_i . Definiramo $p(X) = p_1(X) \dots p_k(X)$. Zaradi normalnosti so vse ničle polinoma $p_i(X)$ v E , torej so vse ničle polinoma $p(X)$ v E . Velja $F(p) \subseteq E$, ampak $a_1, \dots, a_k \in F(p)$, torej $F(p) \supseteq F(a_1, \dots, a_k) = E$.

V levo: Recimo, da je $E = F(p)$. Vzemimo poljuben nerazcepen $q(X) \in F[X]$, ki ima neko ničlo $a \in E$. Naj bo b poljubna ničla q . Ničli a in b imata isti minimalni polinom, q , torej je $F(a) \cong F(b)$. Izomorfizem δ lahko razširimo na izomorfizem razpadnih polj τ . Predpostavimo lahko, da je $\delta(a) = b$, torej tudi $\tau(a) = b$. Ničla a je racionalen izraz, odvisen od a_1, \dots, a_k , torej je tudi b racionalen izraz, odvisen od a_1, \dots, a_k , in posledično $b \in E$. \square

Vprašanje 2. Karakteriziraj normalnost za končne razširitve. Dokaži karakterizacijo.

Definicija. Polinom $p(X) \in F[X]$ je SEPARABILEN, če ima same enostavne ničle. Končna razširitev E/F je SEPARABILNA, če je minimalni polinom vsakega elementa $a \in E$ separabilen.

Opomba. V karakteristiki 0 je vsaka končna razširitev separabilna.

Izrek (Primitivni element). *Vsaka separabilna razširitev je enostavna.*

Dokaz. Če je F končno polje, je tudi E končno in je (E^*, \cdot) ciklična grupa, generirana z a . Velja $E = F(a_1, \dots, a_n)$, brez škode za splošnost $a_i \neq 0$, torej obstajajo n_i , da je $a_i = a^{n_i}$. Torej je $E = F(a)$.

Če pa je F neskončno polje, zapišimo $E = F(a_1, \dots, a_n)$, ker je razširitev končna. Poleg tega so a_i algebraični. Ker velja $F(a_1, \dots, a_n) = F(a_1, \dots, a_{n-2})(a_{n-1}, a_n)$, izrek zadošča dokazati za $n = 2$.

Naj bo torej $E = F(b, c)$ ter $p(X)$ in $q(X)$ minimalna polinoma b in c . Označimo z E_1 razpadno polje polinoma $p(X)q(X)$ in z $b = b_1, \dots, b_r$, $c = c_1, \dots, c_s$ njune ničle v tem polju. Izberimo $\lambda \in F$ tako, da $\lambda \neq (b_j - b)(c - c_k)^{-1}$ za vse j in k . Tak λ res obstaja, ker je polje neskončno. Sedaj definiramo $a = b + \lambda c$. Očitno je $F(a) \subseteq F(b, c)$; trdimo, da velja tudi vsebovanost v drugi smeri. Vpeljimo $f(X) = p(a - \lambda X) \in F(a)[X]$. Velja $f(c) = p(a - \lambda c) = p(b) = 0$, torej je c hkrati ničla $f(X)$ in $q(X) \in F[X] \subseteq F(a)[X]$. Naj bo $\tilde{q}(X) \in F(a)[X]$ minimalni polinom elementa c nad poljem $F(a)$. Velja $\tilde{q}|f$ in $\tilde{q}|q$.

Možne ničle polinoma \tilde{q} so skupne ničle polinomov q in f . Vemo, da imata skupno ničlo c ; katere od c_i pa so še ničle f ? Če je $f(c_j) = p(a - \lambda c_j) = 0$, potem je $a - \lambda c_j = b_k$ za nek k , torej $b + \lambda c - \lambda c_j = b_k$ in posledično $\lambda = (b_j - b)(c - c_k)^{-1}$, kar je po predpostavki nemogoče. Torej je c edina ničla \tilde{q} . Ker je to minimalni polinom in polje F separabilno, je ničla enostavna in zato $\tilde{q}(X) = X - c \in F(a)[X]$, torej $c \in F(a)$. \square

Primer. Razširitev $E = F_2(\sqrt{t})$ je algebraična razširitev, ki ni separabilna. Minimalni polinom za \sqrt{t} je $x^2 - t = (x + \sqrt{t})^2$.

Vprašanje 3. Kdaj je razširitev separabilna? Podaj primer končne razširitve, ki ni algebraična. Povej izrek o primitivnem elementu in ga dokaži za končna polja.

Definicija. GALISOVA GRUPE RAZŠIRITVE E/F je grupa

$$\text{Gal}(E/F) = \{\sigma \in \text{Aut}(E) \mid \sigma|_F = \text{id}_F\}.$$

Lema. Naj bo E/F razširitev, $a \in E$ ničla $p(X) \in F[X]$ ter $\sigma \in \text{Gal}(E/F)$. Potem je $\sigma(a)$ tudi ničla $p(X)$.

Dokaz. Uporabimo σ na zapisu $p(X) = a_n X^n + \dots + a_0$. \square

Definicija. Naj bo $p(X) \in F[X]$ polinom. GALISOVA GRUPE polinoma p je $\text{Gal}(p) = \text{Gal}(F(p)/F)$.

Vprašanje 4. Definiraj Galoisovo grupo razširitve in polinoma.

Definicija. Normalnim separabilnim razširitvam pravimo GALISOVE RAZŠIRITVE.

Trditev. Naj bo E/F končna Galoisova razširitev. Potem je $|\text{Gal}(E/F)| = [E : F]$.

Dokaz. Po izreku o primitivnem elementu obstaja $a \in E$, da je $E = F(a)$. Naj bo $p(X) \in F[X]$ minimalni polinom a . Zaradi normalnosti so vse ničle p v E , torej je za $\sigma \in \text{Gal}(E/F)$ slika $\sigma(a)$ lahko katerakoli ničla p . Ničel je ravno toliko kot je stopnja razširitve, ker ima p paroma različne ničle. Torej je za σ natanko toliko možnosti. \square

Opomba. Če je E/F končna separabilna razširitev, je $|\text{Gal}(E/F)| \leq [E : F]$.

Vprašanje 5. Kolikšen je red Galoisove grupe končne Galoisove razširitve? Dokaži.

Lema. Naj bo $F \subseteq L \subseteq E$.

- Če je E/F končna, je tudi E/L končna.
- Če je E/F normalna, je tudi E/L normalna.
- Če je E/F separabilna, je tudi E/L separabilna.

Dokaz. Prva točka je enostavna, dokaz tretje pa je podoben dokazu druge, torej dokažemo le to. Naj bo $p(X) \in L[X]$ nerazcepen z ničlo $a \in E$ in naj bo $q(X) \in F[X]$ minimalen polinom za a nad F . Trdimo, da $p(X)$ deli $q(X)$. Vzemimo največji skupni delitelj $d(X) \in L[X]$ polinomov $p(X)$ in $q(X)$. Ker imata skupno ničlo a , d ni konstanten, $d(a) = 0$ in d deli p . Ker je p nerazcepen, mora biti $d(X) = p(X)$, torej p res deli q . Vsaka ničla b polinoma p je torej tudi ničla q . Ker je E normalna, je $b \in E$. \square

Izrek (Fundamentalni izrek Galoisove teorije). Naj bo E/F končna Galoisova razširitev.

- Predpisa $L \mapsto \text{Gal}(E/L)$ in $G \mapsto E^G$ sta paroma inverzni preslikavi med vmesnimi polji razširitve in podgrupami Galoisove grupe.
- Za poljubni vmesni polji $F \subseteq L \subseteq M \subseteq E$ velja $[M : L] = |\text{Gal}(E/F) : \text{Gal}(E/M)|$.
- Za vmesno polje $F \subseteq L \subseteq E$ je L/F normalna natanko tedaj, ko je $\text{Gal}(E/L) \triangleleft \text{Gal}(E/F)$. Velja še $\text{Gal}(L/F) \cong \text{Gal}(E/F) / \text{Gal}(E/L)$.

Dokaz. Za prvo točko: Dokažimo, da za vsak L velja

$$E^{\text{Gal}(E/L)} = L.$$

Očitno je $L \subseteq E^{\text{Gal}(E/L)}$. Trdimo, da za vsak $a \in E \setminus L$ obstaja $\sigma \in \text{Gal}(E/L)$, ki ga ne fiksira. Naj bo $q(X) \in L[X]$ minimalni polinom a . Ker $a \notin L$, je $\deg q > 1$, torej ima q od a različno ničlo $b \in E$. Polji $L(a)$ in $L(b)$ sta izomorfni, izomorfizem lahko razširimo do izomorfizma razpadnih polj polinoma q ; ta izomorfizem fiksira elemente L in slika a v b .

Sedaj pokažimo $\text{Gal}(E/E^G) = G$. Očitno je $G \subseteq \text{Gal}(E/E^G)$, za drugo inkluzijo pa je dovolj dokazati (ker so vse grupe končne), da je $[E : E^G] \leq |G|$. Zaradi separabilnosti je $E = E^G(a)$. Naj bo $q(X) \in E^G[X]$ minimalni polinom a . Velja $\text{st } q = [E : E^G]$. Definiramo

$$p(X) = \prod_{\tau \in G} (X - \tau(a)) \in E[X].$$

Trdimo, da je $p(X) \in E^G[X]$, za kar moramo pokazati, da vsak $\sigma \in G$ fiksira koeficiente polinoma. Označimo polinom slik koeficientov z

$$\sigma(p(X)) = \prod_{\tau \in G} (X - \sigma(\tau(a))).$$

Ko τ preteče cel G , tudi $\sigma\tau$ preteče cel G , torej $\sigma(p(X)) = p(X) \in E^G[X]$. Velja $p(a) = 0$, torej q deli p in zato $[E : E^G] = \text{st } q \leq \text{st } p = |G|$.

Druga točka: Preprost račun

$$[M : L] = \frac{[E : L]}{[E : M]} = \frac{|\text{Gal}(E/L)|}{|\text{Gal}(E/M)|} = |\text{Gal}(E/L) : \text{Gal}(E/M)|.$$

Tretja točka: Dokaz v več korakih. Prvo pokažimo, da je L/F normalna natanko tedaj, ko za vsak $\sigma \in \text{Gal}(E/F)$ velja $\sigma(L) = L$. V desno naj bo $\sigma \in \text{Gal}(E/F)$. Velja $L = F(a_1, \dots, a_k)$, označimo s p_i minimalni polinom a_i . Ker je razširitev normalna, so vse ničle p_i v L , torej se a_i s σ slika v neko ničlo p_i . Torej je $\sigma(a_i) \in L$ za vse i , drugo inkluzijo pa dobimo z inverzom σ^{-1} . Podobno v drugo smer.

Sedaj pokažimo, da za $F \subseteq L \subseteq E$ in poljuben $\sigma \in \text{Gal}(E/F)$ velja $\text{Gal}(E/\sigma(L)) = \sigma \text{Gal}(E/F) \sigma^{-1}$. To velja zato, ker je $\tau \in \text{Gal}(E/\sigma(L))$ natanko tedaj, ko je $\tau|_{\sigma(L)} = \text{id}$, oziroma $\tau(\sigma(x)) = \sigma(x)$ za vsak $x \in L$. Če množimo z leve s σ^{-1} , dobimo $\sigma^{-1}\tau\sigma(x) = x$ za $x \in L$, kar je ekvivalentno zahtevi $\sigma^{-1}\tau\sigma \in \text{Gal}(E/L)$.

V tretjem koraku opazimo, da je desni del ekvivalence v prvem koraku enaka zahtevi, da za vsak $\sigma \in \text{Gal}(E/F)$ velja $\text{Gal}(E/F) = \text{Gal}(E/\sigma(L))$, oziroma $\sigma \text{Gal}(E/L) \sigma^{-1} = \text{Gal}(E/L)$, kar pa je po definiciji enako $\text{Gal}(E/L) \triangleleft \text{Gal}(E/F)$.

V zadnjem koraku recimo, da je L/F normalna, in definirajmo $\phi : \text{Gal}(E/F) \rightarrow \text{Gal}(L/F)$ kot

$$\phi(\sigma) = \sigma|_L.$$

Po prvem koraku dokaza prejšnje točke zožena preslikava res slika L v L , iz izreka o izomorfizmu torej dobimo $\text{Gal}(E/F)/\text{Gal}(E/L) \cong \text{Gal}(L/F)$. \square

Vprašanje 6. Povej in dokaži fundamentalni izrek Galoisove teorije.

5.1.1 Rešljive grupe

Definicija. Grupa G je REŠLJIVA, če obstaja končno zaporedje podgrup

$$\{e\} = G_0 \subseteq G_1 \subseteq \cdots \subseteq G_k = G,$$

da je za vsak i $G_i \triangleleft G_{i+1}$ ter G_{i+1}/G_i Abelova grupa.

Vprašanje 7. Kaj je rešljiva grupa? Povej nekaj pozitivnih in negativnih primerov.

Odgovor: Abelove grupe so rešljive, enostavne nekomutativne grupe pa niso. Za $H = \langle (123) \rangle$ je $\{id\} \triangleleft H \triangleleft S_3$, za $K = \langle (12)(34), (14)(32) \rangle$ pa $\{id\} \triangleleft K \triangleleft A_4 \triangleleft S_4$, preverimo lahko, da sta tako S_3 kot S_4 rešljivi. \square

Trditev. Vsaka podgrupa rešljive grupe je rešljiva.

Dokaz. Naj bo G rešljiva, $\{e\} = G_0 \triangleleft \cdots \triangleleft G_k = G$, in naj bo $H \leq G$. Potem je $G_i \cap H \triangleleft G_{i+1} \cap H$, za kvocient pa velja

$$\frac{G_{i+1} \cap H}{G_i \cap H} = \frac{G_{i+1} \cap H}{G_i \cap G_{i+1} \cap H} \cong \frac{(G_{i+1} \cap H) \cdot G_i}{G_i} \leq \frac{G_{i+1}}{G_i},$$

kjer smo uporabili drugi izrek o izomorfizmu. Grupa G_{i+1}/G_i je Abelova, torej je

$$\{e\} = G_0 \cap H \triangleleft \cdots \triangleleft G_k \cap H = H.$$

□

Vprašanje 8. Pokaži, da je podgrupa rešljive grupe rešljiva.

Trditev. Kvocient rešljive grupe je rešljiv.

Dokaz. Naj bo G rešljiva in $N \triangleleft G$. Preverimo lahko, da velja

$$\{e\} \triangleleft \frac{G_0 N}{N} \triangleleft \cdots \triangleleft \frac{G_k N}{N} = \frac{G}{N}.$$

Z uporabo tretjega in drugega izreka o izomorfizmu izračunamo

$$\frac{G_{i+1} N / N}{G_i N / N} \cong \frac{G_{i+1} / N}{G_i / N} = \frac{G_{i+1} G_i / N}{G_i / N} \cong \frac{G_{i+1}}{G_{i+1} \cap G_i N} \cong \frac{G_{i+1} / G_i}{(G_{i+1} \cap G_i N) / G_i}.$$

Kvocient Abelove grupe je Abelova grupa. □

Vprašanje 9. Pokaži, da je kvocient rešljive grupe rešljiv.

Trditev. Naj bo $N \triangleleft G$ in naj bosta N in G/N rešljivi. Potem je G rešljiva grupa.

Dokaz. Naj bosta $N_0 \triangleleft \cdots \triangleleft N_k$ in $G_0/N \triangleleft \cdots \triangleleft G_l/N$ zaporedji edink. Iz $G_i/N \triangleleft G_{i+1}/N$ sledi $G_i \triangleleft G_{i+1}$. Po tretjem izreku o izomorfizmu je G_{i+1}/G_i Abelova (ker je $\frac{G_{i+1}/N}{G_i/N}$ Abelova). Sedaj imamo zaporedje

$$N_0 \triangleleft N_i \triangleleft \cdots \triangleleft N_k = N = G_0 \triangleleft G_1 \triangleleft \cdots \triangleleft G_l = G.$$

□

Vprašanje 10. Kako sestaviš rešljivo grupo iz rešljive edinke in kvocienta? Dokaži.

5.1.2 Rešljivost polinomskih enačb z radikali

Dan je polinom $p(X) \in F[X]$. Rešljivost polinomske enačbe $p(X) = 0$ z radikali je ekvivalentna obstoju razširitve E/F , ki jo dobimo z zaporedjem radikalskih razširitev $F = E_0 \subseteq E_1 \subseteq \cdots \subseteq E_k = E$.

Lema. Naj bo F polje, ki vsebuje primitivni n -ti koren enote. Naj bo $a \in F$. Potem je Galoisova grupa polinoma $X^n - a$ ciklična.

Dokaz. Naj bo ε primitivni n -ti koren enote. Če je b neka ničla $X^n - a$, so vse ničle oblike $b, \varepsilon b, \dots, \varepsilon^{n-1}b$. Potem je $F(X^n - a) = F(b, \varepsilon b, \dots, \varepsilon^{n-1}b) = F(b)$, ker je $\varepsilon \in F$. Preslikava $\sigma \in \text{Gal}(X^n - a)$ je določena s sliko b ;

$$b \mapsto \varepsilon^i b.$$

S tem je določena injektivna preslikava $\text{Gal}(X^n - a) \rightarrow \mathbb{Z}_n$, ki je hkrati homomorfizem grup. □

Izrek. Enačba $p(X) = 0$ je rešljiva z radikali natanko tedaj, ko je Galoisova grupa polinoma p rešljiva grupa.

Dokaz. Dokažemo le implikacijo v desno; v drugo smer je podobno. Enačba je rešljiva, torej obstajajo razširitve $F = E_0 \subseteq E_1 \subseteq \cdots \subseteq E_k = E$, da je $E_{i+1} = E_i(a_i)$ in $a_i^{n_i} \in F$. Definiramo $n = n_0 \dots n_{k-1}$. Naj bo ε primitivni n -ti koren enote. Za ta dokaz se ne bomo ukvarjali s separabilnostjo; privzamemo npr. da je karakteristika polja 0.

Po izreku o primitivnem elementu obstaja a , da je $E = F(a)$. Naj bo $g(X) \in F[X]$ minimalni polinom a in $\Omega = F(g(X)(X^n - 1))$. To je normalna razširitev F , ki vsebuje E in ε . Vzemimo še \tilde{E} kot normalno zaprtje polja $E(\varepsilon)$ v Ω tj. najmanjša normalna razširitev $E(\varepsilon)$, ki je vsebovana v Ω .

Polje \tilde{E} vsebuje vse $\sigma(E(\varepsilon))$ za $\sigma \in \text{Gal}(\Omega/F)$, torej je enako

$$\tilde{E} = F(\varepsilon, a_0, a_1, \dots, a_{k-1}, \sigma_1(a_0), \dots, \sigma_2(a_0), \dots).$$

Razširitev \tilde{E}/F lahko naredimo po korakih, po vrstnem redu kot zgoraj. V vsaki vmesni razširitvi dodamo ničlo polinoma $X^{n_i} - b_i$ za nek b_i , torej so razširitve radikalske. Preverimo lahko, da so tudi normalne.

Po Galoisovi korespondenci dobimo verigo podgrup

$$\text{Gal}(\tilde{E}/F) \supseteq \text{Gal}(\tilde{E}/F(\varepsilon)) \supseteq \cdots \supseteq \text{Gal}(\tilde{E}/\tilde{E}) = \{e\}.$$

Zaradi normalnosti je to zaporedje edink. Kvocienti so po lemi ciklične grupe, torej je $\text{Gal}(\tilde{E}/F)$ rešljiva. Grupa

$$\frac{\text{Gal}(\tilde{E}/F)}{\text{Gal}(\tilde{E}/F(p))} \cong \text{Gal}(F(p)/F) = \text{Gal}(p)$$

je kvocient rešljive grupe, torej je rešljiva. \square

Vprašanje 11. Dokaži, da je rešljivost polinomske enačbe $p(X) = 0$ z radikali implicira rešljivost grupe $\text{Gal}(p)$.

5.2 Moduli

Definicija. Naj bo R kolobar in M neprazna množica. LEVI R -MODUL je $(M, +, \cdot)$, kjer sta $+: M \times M \rightarrow M$ in $\cdot: R \times M \rightarrow M$. Pri tem zahtevamo

- $(M, +)$ je Abelova grupa
- $r(m_1 + m_2) = rm_1 + rm_2$
- $(r_1 + r_2)m = r_1m + r_2m$
- $(r_1r_2)m = r_1(r_2m)$
- $1 \cdot m = m$

Analogno lahko definiramo desne module. Vsak levi R -modul je pravzaprav tudi desni R^{opp} modul, kjer je R^{opp} kolobar z enakim seštevanjem kot v R in množenjem v obratnem vrstnem redu.

Vprašanje 12. Definiraj modul.

Primer. Vektorski prostor nad poljem F je levi (in desni) F -modul.

Primer. Abelova grupa (tu pisana aditivno) je \mathbb{Z} -modul za množenje

$$n \cdot g = \underbrace{g + g + \cdots + g}_n.$$

Primer. Naj bo I levi ideal kolobarja R . Potem je I levi R -modul.

Primer. Če je V vektorski prostor nad poljem F in R množica endomorfizmov $V \rightarrow V$, je V levi R -modul za množenje

$$A \cdot v = A(v).$$

Vprašanje 13. Naštej nekaj enostavnih primerov modulov.

Definicija. Naj bo M R -modul. Množica $N \subseteq M$ je **PODMODUL** v M , če je N za podedovano seštevanje in množenje s skalarjem tudi modul. Oznaka $N \leq M$.

Definicija. Naj bo M R -modul in $X \subseteq M$ neprazna množica. Najmanjšemu podmodulu, ki vsebuje X , pravimo **PODMODUL, GENERIRAN Z MNOŽICO** X . Oznaka $\langle X \rangle$. Modul je **KONČNO GENERIRAN**, če obstaja končna X , da je $M = \langle X \rangle$. Če je modul generiran z enim samim elementom, je **CIKLIČNI**.

Definicija. Naj bosta N, M R -modula. Preslikava $\varphi : M \rightarrow N$ je **HOMOMORFIZEM** R -modulov, če velja $\varphi(m_1 + m_2) = \varphi(m_1) + \varphi(m_2)$ in $\varphi(rm) = r\varphi(m)$. **JEDRO** homomorfizma je množica

$$\ker \varphi = \{m \in M \mid \varphi(m) = 0\},$$

SLIKA pa

$$\operatorname{im} \varphi = \{\varphi(m) \mid m \in M\}.$$

Definicija. Naj bo M R -modul in $N \leq M$. Za relacijo

$$m_1 \sim m_2 \Leftrightarrow m_1 - m_2 \in N$$

definiramo **KVOCIENTNI MODUL** $M/\sim = M/N$.

Izrek. Za kvocientne module veljajo naslednje lastnosti:

- če je $\varphi : M \rightarrow N$ homomorfizem, potem je $M/\ker \varphi \cong \operatorname{im} \varphi$,
- če sta N, K podmodula v M , potem je $(N + K)/K \cong N/N \cap K$,
- velja

$$\frac{M/N}{L/N} \cong M/L.$$

Definicija. Naj bodo N_1, \dots, N_k R -moduli. **ZUNANJI DIREKTNI PRODUKT** je R -modul $N_1 \times N_2 \times \dots \times N_k$ s seštevanjem in množenjem s skalarjem po komponentah. Oznaka $N_1 \oplus N_2 \oplus \dots \oplus N_k$.

Opomba. Če imamo neskončno družino modulov, lahko naredimo podobno konstrukcijo, ki ji pravimo **KARTEZIČNI PRODUKT**.

Definicija. Naj bodo $N_1, \dots, N_k \leq M$. Vsoti $N_1 + \dots + N_k$ pravimo **NOTRANJA DIREKTNA VSOTA**, če za vsak i velja $N_i \cap (N_1 + \dots + N_{i-1} + N_{i+1} + \dots + N_k) = \{0\}$.

Opomba. Zunanja in notranja definicija sta ekvivalentni.

Definicija. Naj bo M R -modul in $X \subseteq M$. Pravimo, da je X BAZA modula M , če je $M = \langle X \rangle$ in če za vsak $k \in \mathbb{N}$ in vse $x_1, \dots, x_k \in X$ iz enakosti $r_1x_1 + \dots + r_kx_k = 0$ sledi $r_i = 0$.

Primer. Obstajajo moduli, ki nimajo baz. Modul $M = \mathbb{Z}_n = \mathbb{Z}/n\mathbb{Z}$ je \mathbb{Z} -modul. Recimo, da je nek $x + n\mathbb{Z}$ v bazi. Ampak $n(x + n\mathbb{Z}) = 0$, n pa v \mathbb{Z} ni enak 0.

Definicija. Modul je PROST, če ima bazo.

Vprašanje 14. Kaj je baza modula? Podaj primer modula, ki ni prost.

Izrek. Naj bo M R -modul. Naslednje trditve so ekvivalentne:

- M je prost R -modul.
- Obstaja indeksna množica Λ , da je M kot R -modul izomorfen zunanji direktni vsoti kopij R -modula R :

$$M \cong \bigoplus_{\lambda \in \Lambda} R.$$

- Obstaja indeksna množica Λ in podmoduli $M_\lambda \leq M$, $M_\lambda \cong R$, da je

$$M = \bigoplus_{\lambda \in \Lambda} M_\lambda.$$

Dokaz. Druga in tretja točka sta očitno ekvivalentni. Naj bo X baza modula M . Trdimo $M = \bigoplus_{x \in X} \langle x \rangle$. Ker je $M = \langle X \rangle$, je M vsota podmodulov $\langle x \rangle$. Vzemimo $z \in \langle x \rangle \cap \sum_{y \neq x} \langle y \rangle$. Velja $z = rx = r_1y_1 + \dots + r_ky_k$, iz definicije baze pa sledi $r = r_i = 0$. Preveriti moramo še, da je $\langle x \rangle \cong R$. Ustrezen izomorfizem $R \rightarrow \langle x \rangle$ bo $\varphi(r) = rx$.

Sedaj dokažimo, da iz druge točke sledi prva. Imamo izomorfizem

$$\varphi : \bigoplus_{\lambda \in \Lambda} R \rightarrow M.$$

Izberimo e_λ kot Λ -terico, ki ima na mestu λ enico, drugje pa 0, in definirajmo $x_\lambda = \varphi(e_\lambda)$. Preverimo lahko, da je $X = \{x_\lambda \mid \lambda \in \Lambda\}$ baza. \square

Vprašanje 15. Karakteriziraj prostost modula in dokaži karakterizacijo.

Posledica. Vsak R -modul je kvocient nekega prostega R -modula.

Dokaz. Naj bo M R -modul. Definiramo

$$\varphi : \bigoplus_{m \in M} R \rightarrow M$$

tako, da za element e_m direktne vsote, ki ima na m -tem mestu enico, drugje pa nič, predpišemo $\varphi(e_m) = m$. Ti elementi tvorijo bazo, torej je s tem natanko določen epimorfizem R -modulov. \square

Vprašanje 16. Pokaži, da je vsak modul kvocient prostega modula.

Primer. Podmodul prostega modula ni nujno prost. \mathbb{Z}_8 je prost \mathbb{Z}_8 -modul, njegov podmodul $2\mathbb{Z}_8 = \{0, 2, 4, 6\}$ pa ni prost.

Vprašanje 17. Podaj primer podmodula prostega modula, ki ni sam prost.

Izrek (univerzalna lastnost). *R -modul M je prost natanko tedaj, ko obstaja množica X in preslikava $\iota : X \rightarrow M$, da velja: za vsak R -modul N in vsako preslikavo $\kappa : X \rightarrow N$ obstaja natanko en homomorfizem R -modulov $f : M \rightarrow N$, da velja $f \circ \iota = \kappa$.*

Dokaz. V desno: X naj bo baza M , ι pa vložitev, torej

$$M = \bigoplus_{x \in X} R, \\ \iota(x) = e_x.$$

Naj bo N poljubna R -modul in $\kappa : X \rightarrow N$ poljubna preslikava. Recimo, da obstaja f , kot zahtevamo. Potem mora veljati $f \circ \iota = \kappa$, torej $f(e_x) = \kappa(x)$. Ker je $\{e_x\}_x$ baza, to enolično določi f , če le obstaja. Definiramo ga kot

$$f\left(\sum_{x \in X} a_x e_x\right) = \sum_{x \in X} a_x \kappa(x),$$

kar ustreza zahtevi.

V levo: Dokazujemo, da je

$$M \cong \bigoplus_{x \in X} R.$$

Definiramo κ , kar nam po univerzalni lastnosti poda homomorfizem f , da diagram

$$\begin{array}{ccc} X & \xrightarrow{\kappa} & \bigoplus_x R \\ \downarrow \iota & \searrow g & \nearrow f \\ M & & \end{array}$$

komutira. Pokazali smo, da prosti moduli imajo univerzalno lastnost, torej obstaja natanko en homomorfizem R -modulov g kot zgoraj, torej $g \circ \kappa = \iota$. Pokazati želimo, da sta f in g inverzni. Velja $f \circ g \circ \kappa = f \circ \iota = \kappa = \text{id} \circ \kappa$. Oglejmo si torej diagram

$$\begin{array}{ccc} X & \xrightarrow{\kappa} & \bigoplus_x R \\ \downarrow \kappa & \searrow f \circ g & \nearrow \text{id} \\ \bigoplus_x R & & \end{array}$$

Za oba homomorfizma diagram komutira, torej velja $f \circ g = \text{id}$. Podobno v drugo smer. \square

Vprašanje 18. Povej in dokaži izrek o univerzalni lastnosti prostih modulov.

Primer. Naj bo F polje in R množica endomorfizmov

$$R = \text{End}(\oplus_{n \in \mathbb{N}} F),$$

torej neskončnih kvadratnih matrik elementov iz F . Kot R -modul je R prost z bazo $\{\text{id}\}$. Imamo pa tudi bazo iz dveh elementov:

$$B_1 = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & \cdots \\ 0 & 0 & 1 & 0 & 0 & \cdots \\ 0 & 0 & 0 & 0 & 0 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix} \quad B_2 = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 & \cdots \\ 0 & 0 & 0 & 1 & 0 & \cdots \\ 0 & 0 & 0 & 0 & 0 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}$$

Vsaka matrika A je R -linearna kombinacija

$$A = [a_1, a_3, a_5, \dots]B_1 + [a_2, a_4, a_6, \dots]B_2,$$

torej je $\{B_1, B_2\}$ res baza (preverimo lahko, da sta res linearno neodvisni).

Definicija. Kolobar R ima LASTNOST ENOLIČNEGA RANGA, če imajo vse baze poljubnega prostega R -modula isto kardinalnost.

Vprašanje 19. Povej primer kolobarja, ki nima lastnosti enoličnega ranga.

Primer. Polja in obsegi imajo lastnost enoličnega ranga.

Definicija. Naj bo R kolobar, $I \triangleleft R$ in M R -modul. Potem je

$$IM = \left\{ \sum_{i=1}^k u_i m_i \mid u_i \in I, m_i \in M, k \in \mathbb{N} \right\}$$

podmodul v M .

Lema. Če je M prost R -modul z bazo X in $I \triangleleft R$ različen od R , potem je M/IM prost R/I -modul z bazo $X + IM$.

Dokaz. Očitno je $\langle X + IM \rangle = M/IM$, ker je $\langle X \rangle = M$. Linearna neodvisnost:

$$\begin{aligned} \sum_i (r_i + I)(x_i + IM) &= 0 + IM \\ \Leftrightarrow \sum_i (r_i x_i + IM) &= 0 + IM \\ \Leftrightarrow \sum_i r_i x_i &\in IM. \end{aligned}$$

Velja

$$\sum_i r_i x_i = \sum_j u_j m_j = \sum_j u_j \left(\sum_k s_k x_k \right) = \sum_{j,k} \underbrace{u_j s_k}_{\in I} x_k.$$

To sta dva razvoja po bazi X , torej $r_i = \sum_j u_j s_i \in I$, in $r_i + I = 0 + I$. □

Vprašanje 20. Dokaži: če je M prost R -modul z idealom I , je M/IM prost R/I -modul.

Lema. Ob oznakah prejšnje leme je $|X| = |X + IM|$.

Dokaz. Recimo, da je $x + IM = x' + IM$. Potem je $x - x' \in IM$, torej

$$x - x' = \sum_j u_j x_j.$$

Če je kakšen $u_j = 1$, je $I = R$, kar po predpostavki ne velja. Torej so vsi različni od 1, in zaradi enoličnosti razvoja po bazi velja $x = x'$. \square

Izrek. Veljata naslednji točki o lastnosti enoličnega ranga:

- Naj bo R kolobar. Recimo, da obstaja pravi ideal $I \triangleleft R$, da ima R/I lastnost enoličnega ranga. Potem ima R lastnost enoličnega ranga.
- Vsak komutativen kolobar ima lastnost enoličnega ranga.

Dokaz. Prva točka: Naj bo M poljuben prost R -modul z bazama X in Y . Vemo, da je M/IM prost R/I -modul z bazama $X + IM$ in $Y + IM$. Ker ima R/I lastnost enoličnega ranga, je $|X + IM| = |Y + IM|$, torej je po lemi $|X| = |Y|$.

Druga točka: Naj bo R komutativen kolobar. Če je R polje, ima lastnost enoličnega ranga. Če pa ni polje, ima nek pravi ideal, ki je vsebovan v nekem maksimalnem idealu I . Potem je R/I polje, torej ima lastnost enoličnega ranga. Po prvi točki ima tako R lastnost enoličnega ranga. \square

Vprašanje 21. Dokaži, da ima vsak komutativen kolobar lastnost enoličnega ranga.

5.2.1 Projekтивni moduli

Definicija. R -modul P je PROJEKTIVEN, če za vsak homomorfizem R -modulov $f : P \rightarrow M$ in vsak epimorfizem R -modulov $g : M' \rightarrow M$ obstaja homomorfizem R -modulov $h : P \rightarrow M'$, da naslednji diagram komutira:

$$\begin{array}{ccc} P & & \\ f \downarrow & \searrow h & \\ M & \xleftarrow{g} & M' \end{array}$$

Vprašanje 22. Kaj je projektiven modul?

Trditev. Vsak prost R -modul je projektiven.

Dokaz. Naj bo F prost modul z bazo X . Naj bo $\iota : X \rightarrow M$ vložitev, f homomorfizem $F \rightarrow M$ in $g : M' \rightarrow M$ epimorfizem. Definiramo $\varphi : X \rightarrow M'$ tako, da $\varphi(x)$ izberemo nek element M' , ki se z g slika v $f(\iota(x))$. Ker je F prost, obstaja natanko en homomorfizem $h : F \rightarrow M'$, da je $\varphi(x) = h(\iota(x))$ za $x \in X$. Ta ustreza zahtevi projektivnosti. \square

Vprašanje 23. Dokaži: vsak prost R -modul je projektiven.

Izrek. Za R -modul P so ekvivalentne naslednje trditve:

- P je projektiven
- za vsak epimorfizem $\varphi : M \rightarrow P$ je $M = P \oplus \ker \varphi$
- obstaja R -modul M , da je $P \oplus M$ prost R -modul

Dokaz. 1 v 2: Ker je P projektiven, obstaja homomorfizem R -modulov $\psi : P \rightarrow M$, da je $\varphi \circ \psi = \text{id}$. Iz tega sledi injektivnost ψ , torej je $\text{im } \psi \cong P$. Dokazali bomo $M = \ker \varphi \oplus \text{im } \psi$. Naj bo $x \in \ker \varphi \cap \text{im } \psi$, torej $x = \psi(y)$. Sledi $\varphi(\psi(y)) = 0$, torej je $x = 0$. Naj bo $m \in M$. Zapišemo

$$m = \underbrace{\psi(\varphi(m))}_{\in \text{im } \psi} + \underbrace{(m - \psi(\varphi(m)))}_{\in \ker \varphi}.$$

2 v 3: Vsak modul je kvocient prostega modula, torej obstaja nek prost F , in epimorfizem $f : F \rightarrow P$. Po predpostavki je $F = \ker \varphi \oplus P$.

3 v 1: Obstaja R -modul N , da je $P \oplus N$ prost. Ta je projektiven, torej za projekcijo na prvo komponento $\pi_1 : P \oplus N \rightarrow P$ obstaja $\tilde{h} : P \oplus N \rightarrow M'$, da naslednji diagram komutira:

$$\begin{array}{ccc} & P \oplus N & \\ \tilde{h} \swarrow & \downarrow f \circ \pi_1 & \\ M' & \longrightarrow & M \end{array}$$

Definiramo $h : P \rightarrow M'$ kot $h(p) = \tilde{h}(p, 0)$. \square

Vprašanje 24. Karakteriziraj projektivnost in dokaži karakterizacijo.

Vprašanje 25. Podaj primer projektivnega modula, ki ni prost.

Odgovor: \mathbb{Z}_2 in \mathbb{Z}_3 sta \mathbb{Z}_6 -modula, $\mathbb{Z}_6 = \mathbb{Z}_2 \oplus \mathbb{Z}_3$ pa je prost \mathbb{Z}_6 -modul, torej je \mathbb{Z}_2 projektiven, ni pa prost. \boxtimes

5.2.2 Tenzorski produkt modulov

V razdelku gledamo le module nad komutativnimi kolobarji z enico.

Definicija. Naj bosta M, N R -modula. Naj bo F prost R -modul nad množico $M \times N$ in T podmodul v F , generiran z naslednjimi elementi:

- $(m_1 + m_2, n) - (m_1, n) - (m_2, n)$
- $(m, n_1 + n_2) - (m, n_1) - (m, n_2)$
- $(rm, n) - r(m, n)$
- $(m, rn) - r(m, n)$

Potem je TENZORSKI PRODUKT R -modulov M, N R -modul

$$M \otimes_R N = F/T.$$

Vprašanje 26. Definiraj tenzorski produkt.

Za elementarne tenzorje veljajo naslednje lastnosti:

- $(m_1 + m_2) \otimes n = m_1 \otimes n + m_2 \otimes n$
- $m \otimes (n_1 + n_2) = m \otimes n_1 + m \otimes n_2$
- $(rm) \otimes n = r(m \otimes n)$
- $m \otimes (rn) = r(m \otimes n)$

Opomba. Preslikava $\beta : M \times N \rightarrow M \otimes_R N$, definirana z $\beta(m, n) = m \otimes n$, je bilinearna.

Izrek (univerzalna lastnost). *Naj bodo M, N, K R -moduli. Za vsako bilinearne preslikavo $\gamma : M \times N \rightarrow K$ obstaja natanko en homomorfizem R -modulov $f : M \otimes_R N \rightarrow K$, da naslednji diagram komutira:*

$$\begin{array}{ccc} M \times N & \xrightarrow{\gamma} & K \\ \beta \downarrow & \nearrow f & \\ M \otimes_R N & & \end{array}$$

Dokaz. Naj bosta F, T kot v definiciji tenzorskega produkta. Ker je F prost nad $M \times N$, obstaja natanko en homomorfizem $\tilde{f} : F \rightarrow K$, da diagram

$$\begin{array}{ccc} M \times N & \xrightarrow{\gamma} & K \\ \iota \downarrow & \nearrow \tilde{f} & \\ F & & \end{array}$$

komutira. Velja $\tilde{f}(m, n) = \gamma(m, n)$. Trdimo, da je $T \subseteq \ker \tilde{f}$; velja

$$\tilde{f}((m_1 + m_2, n) - (m_1, n) - (m_2, n)) = \gamma(m_1 + m_2, n) - \gamma(m_1, n) - \gamma(m_2, n) = 0,$$

ker je γ bilinearna. Podobno pokažemo, da je \tilde{f} -slika ostalih generatorjev 0, torej \tilde{f} inducira natanko en homomorfizem $f : F/T \rightarrow K$. \square

Vprašanje 27. Dokaži, da za tenzorski produkt velja univerzalna lastnost.

Izrek. Naj bodo M, N, P R -moduli. Naj bo $\beta : M \times N \rightarrow P$ poljubna bilinearna preslikava. Recimo, da za vsak R -modul K in vsako bilinearno preslikavo $\gamma : M \times N \rightarrow K$ obstaja natanko en homomorfizem R -modulov $f : P \rightarrow K$, da naslednji diagram komutira:

$$\begin{array}{ccc} M \times N & \xrightarrow{\gamma} & K \\ \beta \downarrow & \nearrow f & \\ P & & \end{array}$$

Potem sta P in $M \otimes_R N$ izomorfna.

Dokaz. Za $\alpha(m, n) = m \otimes n$ po prejšnjem izreku obstaja natanko en $g : M \otimes_R N \rightarrow R$, da naslednji diagram komutira:

$$\begin{array}{ccc} M \times N & \xrightarrow{\beta} & P \\ \alpha \downarrow & \nearrow g & \\ M \otimes_R N & & \end{array}$$

Po predpostavki obstaja natanko en $f : P \rightarrow M \otimes_R N$, da je $f \circ \beta = \alpha$. Preverimo lahko, da sta f in g inverzna. \square

Trditev. Naj bo M R -modul. Potem je $R \otimes_R M \cong M$.

Dokaz. Definiramo $\beta : R \times M \rightarrow M$ s predpisom $\beta(r, m) = rm$. To je bilinearna preslikava, torej po univerzalni lastnosti inducira homomorfizem R -modulov $R \otimes_R M \rightarrow M$, ki slika $r \otimes m \mapsto rm$. V drugo smer podamo homomorfizem $M \rightarrow R \otimes_R M$ s predpisom $m \mapsto 1 \otimes m$. Preverimo lahko, da sta preslikavi inverzni. \square

Vprašanje 28. Kaj je $R \otimes_R M$? Dokaži.

Izrek. Naj bo M prost R -modul nad $X = \{m_i \mid i \in I\}$ in N prost R -modul nad $Y = \{n_j \mid j \in J\}$. Potem je $M \otimes_R N$ tudi prost R -modul nad $Z = \{m_i \otimes n_j \mid i \in I, j \in J\}$.

Dokaz. Množica Z očitno generira $M \otimes_R N$. Naj bo

$$\sum_{ij} r_{ij}(m_i \otimes n_j) = 0.$$

Naredimo homomorfizem $f_k : N \rightarrow R$ za vsak $k \in J$, s predpisom $f_k(n_j) = \delta_{kj}$. Preslikava je definirana na bazi, lahko jo razširimo do homomorfizma $N \rightarrow R$. Oglejmo

si preslikavo $\text{id}_M \otimes_R f_k : M \otimes_R N \rightarrow M \otimes_R R$. Slika zgornje linearne kombinacije s to preslikavo je

$$0 = \sum_{ij} r_{ij}(m_i \otimes \delta_{kj}) = \sum_i r_{ik} m_i \otimes 1.$$

Ker je $M \otimes_R R \cong M$, velja $\sum_i r_{ik} m_i = 0$ (slikamo z izomorfizmom), torej je $r_{ik} = 0$ za vsak i . To lahko naredimo za poljuben k . \square

Vprašanje 29. Pokaži, da je tenzorski produkt prostih modulov prost.

Trditev. Naj bodo A, B, C R -moduli. Potem velja:

- $A \otimes_R B \cong B \otimes_R A$
- $A \otimes_R (B \oplus C) \cong (A \otimes_R B) \oplus (A \otimes_R C)$
- $A \otimes_R (B \otimes_R C) \cong (A \otimes_R B) \otimes_R C$

5.2.3 Skrčitev in razširitev skalarjev

Naj bo $\varphi : S \rightarrow R$ homomorfizem komutativnih kolobarjev. Če je M R -modul, postane tudi S -modul preko predpisa

$$s \circ m = \varphi(s) \cdot m.$$

Predpis je dobro definiran, temu pravimo SKRČITEV SKALARJEV. Podobno lahko naredimo tudi v obratni smeri: če je M S -modul, je po ravno ugotovljenem R tudi S -modul, in je $R \oplus_R M$ R -modul z množenjem

$$r(r' \otimes m) = rr' \otimes m.$$

Temu pravimo RAZŠIRITEV SKALARJEV.

Vprašanje 30. Kaj je razširitev in kaj skrčitev skalarjev?

Trditev. Naj bo $\varphi : R \rightarrow S$ homomorfizem komutativnih kolobarjev. Naj bo M R -modul in N S -modul. Z uporabo skrčitve oz. razširitve skalarjev lahko definiramo Abelovi grupi $\text{Hom}_R(M, N)$ in $\text{Hom}_S(S \otimes_R M, N)$, ki sta izomorfni.

Izrek (Dualnost med Hom in \otimes). Naj bodo M, N, P R -moduli za komutativen kolobar R . Potem sta $\text{Hom}_R(M \otimes_R N, P)$ in $\text{Hom}_R(M, \text{Hom}_R(N, P))$ izomorfna kot R -modula.

Vprašanje 31. V kakšnem smislu sta množica homomorfizmov modulov in tenzorski produkt dualna?

5.2.4 Eksaktna zaporedja modulov

Definicija. Naj bodo A, B, C R -moduli. Imejmo homomorfizme R -modulov

$$A \xrightarrow{f} B \xrightarrow{g} C$$

To zaporedje je EKSAKTNO pri B , če velja $\text{im } f = \ker g$.

Vprašanje 32. Kaj je eksaktno zaporedje? Karakteriziraj injektivnost in surjektivnost z eksaktnostjo.

Odgovor: Za definicijo glej zgoraj. Oglejmo si zaporedje

$$0 \longrightarrow A \xrightarrow{f} B \xrightarrow{g} C \longrightarrow 0$$

To zaporedje je eksaktno pri A natanko tedaj, ko je f injektivna, in eksaktno pri C natanko tedaj, ko je g surjektivna. \square

Za zaporedje

$$0 \longrightarrow A \xrightarrow{f} B \xrightarrow{g} C \longrightarrow 0$$

pravimo, da je **KRATKO EKSAKTNO**, če je eksaktno pri A , B in C . Izkaže se, da so vsa kratka eksaktna zaporedja oblike

$$0 \longrightarrow B \hookrightarrow A \longrightarrow A/B \longrightarrow 0$$

Vprašanje 33. Kaj je kratko eksaktno zaporedje? Kakšno obliko ima?

Trditev. Naj bo $0 \rightarrow A \xrightarrow{f} B \xrightarrow{g} C \rightarrow 0$ kratko eksaktno zaporedje R -modulov. Naslednje trditve so ekvivalentne:

- Obstaja homomorfizem $p : B \rightarrow A$, da je $p \circ f = \text{id}_A$
- Obstaja homomorfizem $s : C \rightarrow B$, da velja $g \circ s = \text{id}_C$
- Obstajata homomorfizma $p : B \rightarrow A$ in $s : C \rightarrow B$, da velja $p \circ f = \text{id}_A$, $g \circ s = \text{id}_C$ in $f \circ p + s \circ g = \text{id}_B$.

V vseh treh primerih je $B \cong A \oplus C$.

Dokaz. Dokažimo le implikacijo iz prve točke v drugo; ostalo gre podobno. Najprej dokažimo, da je $B = \ker p \oplus \text{im } f$. Naj bo x v preseku. Potem je $p(x) = 0$ in $x = f(a)$, torej $a = p(f(a)) = 0$ in $x = 0$. Vzemimo $b \in B$. Potem je

$$b = f(p(b)) + (b - f(p(b))),$$

velja $p(b - f(p(b))) = p(b) - p(f(p(b))) = 0$.

Velja $\text{im } f \cong A$. Trdimo, da je $\ker p \cong C$. Dokažimo, da je g , skrčen na $\ker p$, izomorfizem. Naj bo $c \in C$. Ker je g surjektiv, obstaja $b \in B$, da je $g(b) = c$. Po prejšnjem je $b = x + y$ za $x \in \ker p$ in $y \in \text{im } f$, torej je $c = g(b) = g(x)$, ker je $\text{im } f = \ker g$. Za injektivnost vzemimo $z \in \ker p$, za katerega je $g(z) = 0$. Ker je $\ker g = \text{im } f$, je $z = f(y)$ za nek $y \in A$. Potem je $0 = p(z) = p(f(y)) = y$, torej $z = 0$. Za s lahko torej vzamemo inverz skrčitve g na $\ker p$. \square

Vprašanje 34. Dokaži: če za kratko eksaktno zaporedje obstaja homomorfizem $p : B \rightarrow A$, za katerega je $p \circ f = \text{id}_A$, obstaja tudi homomorfizem $s : C \rightarrow B$, da je $g \circ s = \text{id}_C$.

Izrek (kratka lema o petih). *Imejmo diagram R -modulov in homomorfizmov*

$$\begin{array}{ccccccccc} 0 & \longrightarrow & A & \xrightarrow{f} & B & \xrightarrow{g} & C & \longrightarrow & 0 \\ & & \alpha \downarrow & & \beta \downarrow & & \gamma \downarrow & & \\ 0 & \longrightarrow & A & \xrightarrow{f'} & B & \xrightarrow{g'} & C & \longrightarrow & 0 \end{array}$$

pri čemer sta obe vrstici eksaktni in oba kvadrata komutativna diagrama. Če sta α in γ injektivna, je tudi β injektiven. Če sta α in γ surjektivna, je tudi β surjektiven.

Dokaz. Dokažemo samo za injektivnost. Vzemimo $b \in \ker \beta$. Velja $\gamma(g(b)) = g'(\beta(b)) = g'(0) = 0$, torej je $g(b) = 0$ in $b \in \ker g = \operatorname{im} f$. Torej je $b = f(a)$. Velja $0 = \beta(f(a)) = f'(\alpha(a))$, torej $\alpha(a) \in \ker f' = \{0\}$ in $a = 0$. Torej $b = 0$. \square

Vprašanje 35. Povej in dokaži kratko lemo o petih.

5.3 Teorija kategorij

Definicija. KATEGORIJA je struktura $\underline{\mathcal{C}}$, sestavljena iz

- razreda objektov $\operatorname{Ob} \underline{\mathcal{C}}$,
- množice morfizmov $\underline{\mathcal{C}}(A, B)$ med vsakima objektoma A, B ,
- preslikav $\circ : \underline{\mathcal{C}}(A, B) \times \underline{\mathcal{C}}(B, C) \rightarrow \underline{\mathcal{C}}(A, C)$, kjer so A, B, C poljubni objekti kategorije $\underline{\mathcal{C}}$.

Zahtevamo, da so preslikave \circ asociativne, in da za vsak objekt A obstaja morfizem $1_A \in \underline{\mathcal{C}}(A, A)$, da velja $f \circ 1_A = f$ za vsak $f \in \underline{\mathcal{C}}(A, B)$ ter $1_A \circ g = g$ za vsak $g \in \underline{\mathcal{C}}(B, A)$. Temu pogoju pravimo UNITALNOST.

Vprašanje 36. Definiraj kategorijo.

Definicija. Naj bo $\underline{\mathcal{C}}$ kategorija in $f \in \underline{\mathcal{C}}(A, B)$. Potem je f IZOMORFIZEM, če obstaja morfizem $g \in \underline{\mathcal{C}}(B, A)$, da je $f \circ g = 1_B$ in $g \circ f = 1_A$. Pravimo, da sta A in B IZOMORFNA, če med njima obstaja izomorfizem.

Definicija. Naj bo $f \in \underline{\mathcal{C}}(A, B)$ morfizem. Potem je

- PREREZ, če obstaja $g \in \underline{\mathcal{C}}(B, A)$, da je $g \circ f = 1_A$,
- RETRAKT, če obstaja $g \in \underline{\mathcal{C}}(B, A)$, da je $f \circ g = 1_B$,
- MONOMORFIZEM, če za vsaka morfizma $g, h \in \underline{\mathcal{C}}(C, A)$ iz enakosti $f \circ g = f \circ h$ sledi $g = h$.
- EPIMORFIZEM, če za vsaka morfizma $g, h \in \underline{\mathcal{C}}(B, C)$ iz enakosti $g \circ f = h \circ f$ sledi $g = h$.

Vprašanje 37. Definiraj prerez, retrakt, monomorfizem in epimorfizem.

Vprašanje 38. Povej primer epimorfizma, ki ni surjektiven.

Odgovor: Vložitev $\iota : \mathbb{Z} \rightarrow \mathbb{Q}$ je homomorfizem kolobarjev. Če je $g \circ \iota = h \circ \iota$, velja $g(n) = h(n)$ za vsako celo število n . Potem velja tudi

$$g\left(\frac{a}{b}\right) = g(a)g\left(\frac{1}{b}\right) = h(a)g\left(\frac{1}{b}\right) = h\left(b\frac{a}{b}\right)g\left(\frac{1}{b}\right) = h\left(\frac{a}{b}\right).$$

☒

Definicija. Naj bo $\underline{\mathcal{C}}$ kategorija. Objekt Z je ZAČETNI OBJEKT, če za vsak objekt A velja $|\underline{\mathcal{C}}(Z, A)| = 1$. Objekt K je KONČNI OBJEKT, če za vsak objekt A velja $|\underline{\mathcal{C}}(A, K)| = 1$.

Trditev. Poljubna začetna objekta dane kategorije sta izomorfna.

Dokaz. Naj bosta Z_1 in Z_2 začetna objekta. Potem obstaja $f \in \underline{\mathcal{C}}(Z_1, Z_2)$ in $g \in \underline{\mathcal{C}}(Z_2, Z_1)$. Velja $f \circ g \in \underline{\mathcal{C}}(Z_2, Z_2)$; ker pa je $|\underline{\mathcal{C}}(Z_2, Z_2)| = 1$, je $f \circ g = 1$. Podobno v drugo smer. \square

Opomba. Enako velja za končne objekte.

Vprašanje 39. Dokaži: poljubna začetna objekta v dani kategoriji sta izomorfna.

Definicija. Naj bo $\underline{\mathcal{C}}$ kategorija. NASPROTNA KATEGORIJA $\underline{\mathcal{C}}^{\text{op}}$ je definirana na naslednji način:

- objekti so enaki objektom $\underline{\mathcal{C}}$
- če sta A, B objekta, je $\underline{\mathcal{C}}^{\text{op}}(A, B) = \underline{\mathcal{C}}(B, A)$
- za $f \in \underline{\mathcal{C}}^{\text{op}}(A, B)$ in $g \in \underline{\mathcal{C}}^{\text{op}}(C, A)$ je kompozitum $f * g = g \circ f$.

Vprašanje 40. Definiraj nasprotno kategorijo.

Definicija. Naj bosta $\underline{\mathcal{C}}$ in $\underline{\mathcal{D}}$ kategoriji. FUNKTOR med kategorijama je predpis F , za katerega velja

- za vsak objekt A v $\underline{\mathcal{C}}$ imamo natanko določen objekt $F(A)$ v $\underline{\mathcal{D}}$,
- za vsak $f \in \underline{\mathcal{C}}(A, B)$ imamo natanko določen morfizem $F(f) \in \underline{\mathcal{D}}(F(A), F(B))$,
- $F(1_A) = 1_{F(A)}$ za vsak objekt A v $\underline{\mathcal{C}}$,
- $F(f \circ g) = F(f) \circ F(g)$

KOFUNKTOR med $\underline{\mathcal{C}}$ in $\underline{\mathcal{D}}$ je funktor v nasprotni kategoriji.

Vprašanje 41. Definiraj funktorje in kofunktorje.

Definicija. Naj bosta F, G funktorja iz $\underline{\mathcal{C}}$ v $\underline{\mathcal{D}}$. NARAVNA TRANSFORMACIJA med F in G je nabor morfizmov $\mu_C : F(C) \rightarrow G(C)$ za objekt C kategorije $\underline{\mathcal{C}}$, da za vsak morfizem $f : A \rightarrow B$ v kategoriji $\underline{\mathcal{C}}$ diagram

$$\begin{array}{ccc} F(A) & \xrightarrow{F(f)} & F(B) \\ \mu_A \downarrow & & \downarrow \mu_B \\ G(A) & \xrightarrow{G(f)} & G(B) \end{array}$$

komutira. Pravimo, da sta funktorja F in G NARAVNO IZOMORFNA, če obstaja naravna transformacija $\mu : F \rightarrow G$, za katero so vsi morfizmi μ_C izomorfizmi.

Vprašanje 42. Kaj je naravna transformacija?

5.3.1 Univerzalne konstrukcije

Definicija. Naj bosta A, B objekta kategorije $\underline{\mathcal{C}}$. PRODUKT objektov A in B je objekt P , skupaj z morfizmoma $p : P \rightarrow A$ in $q : P \rightarrow B$, da za poljuben objekt X in poljubna morfizma $f : X \rightarrow A$ ter $g : X \rightarrow B$ obstaja natanko en morfizem $h : X \rightarrow P$, da naslednji diagram komutira:

$$\begin{array}{ccccc} A & \xleftarrow{p} & P & \xrightarrow{q} & B \\ & \nwarrow f & \uparrow h & \nearrow g & \\ & & X & & \end{array}$$

Definicija. KOPRODUKT objektov A in B je produkt v nasprotni kategoriji.

Vprašanje 43. Definiraj produkt in koprodukt.

Definicija. Kategorija je KONKRETNA, če so objekti množice, morfizmi pa preslikave.

Definicija. Naj bo $\underline{\mathcal{C}}$ konkretna kategorija in X neprazna množica. PROSTI OBJEKT nad množico X v kategoriji $\underline{\mathcal{C}}$ je objekt F , skupaj s preslikavo $\iota : X \rightarrow F$, da za vsak objekt C in vsako preslikavo $f : X \rightarrow C$ obstaja natanko en morfizem $\tilde{f} : F \rightarrow C$, da naslednji diagram komutira:

$$\begin{array}{ccc} X & \xrightarrow{f} & C \\ \iota \uparrow & \nearrow \tilde{f} & \\ F & & \end{array}$$

Opomba. Prost objekt, če obstaja, je do izomorfizma natanko enolično določen. Defini-

ramo kategorijo $\hat{\mathcal{C}}$ z objekti $X \xrightarrow{f} C$ in morfizmi

$$\begin{array}{ccc} X & \xrightarrow{f} & C \\ \text{id} \downarrow & & \vdots \\ X & \xrightarrow{g} & C \end{array}$$

kjer je diagram komutativen. Prosti objekt v C je ravno začetni objekt kategorije $\hat{\mathcal{C}}$.

Vprašanje 44. Definiraj prosti objekt in pokaži, da je do izomorfizma natanko enolično določen.

6 Analiza 4

6.1 Osnovni tipi PDE

Uporabljamo notacijo s predavanj:

- $\alpha = (\alpha_1, \dots, \alpha_n) \in \mathbb{N}_0^n$ je multiindeks,
- $|\alpha| = \alpha_1 + \dots + \alpha_n$,
- odvod funkcije u z multiindeksom α je

$$D^\alpha u = \frac{\partial^{|\alpha|}}{\partial x_1^{\alpha_1} \dots \partial x_n^{\alpha_n}},$$

- za $k \in \mathbb{N}$ označimo skupek odvodov k -tega reda kot $D^k u = \{D^\alpha u \mid |\alpha| = k\}$. To včasih obravnavamo kot množico, včasih pa kot vektor ali matriko.

Definicija. PARCIALNA DIFERENCIALNA ENAČBA je enačba, ki vsebuje neznano funkcijo u vsaj dveh spremenljivk ter nekatere njene parcialne odvode.

Definicija. PARCIALNA ENAČBA k -TEGA REDA je enačba oblike

$$F(x, u(x), Du(x), \dots, D^k u(x)) = 0,$$

kjer je $x \in U^{\text{odp}} \subseteq \mathbb{R}^n$ in $F : U \times \mathbb{R} \times \mathbb{R}^n \times \dots \times \mathbb{R}^{n^k} \rightarrow \mathbb{R}$ dana.

Rešitev PDE je vsaka funkcija, ki enačbi zadošča. Lahko je klasična ali posplošena rešitev; za klasično obstajajo vsi odvodi $Du, \dots, D^k u$, ki jih vstavimo v F , posplošena rešitev pa je rešitev v smislu integracije po delih.

Primer. Dana je enačba $\Delta u = f$. Če je $u \in \mathcal{C}^2$ in označimo $\Delta u = g$, za vsak $\phi \in \mathcal{C}^\infty(U)$ s kompaktnim nosilcem velja

$$\langle \Delta u, \phi \rangle = \langle g, \phi \rangle = \int_U g(x) \phi(x) dx.$$

Iz integracije po delih sledi $\langle \Delta u, \phi \rangle = \langle u, \Delta \phi \rangle$ oziroma $\langle u, \Delta \phi \rangle = \langle g, \phi \rangle$. Tu zahtevamo samo, da je u dovolj regularen, da lahko izračunamo integral v skalarnem produktu. Šibka rešitev je vsaka $u \in \mathcal{C}(U)$, ki zadošča $\langle u, \Delta \phi \rangle = \langle f, \phi \rangle$ za vsak ϕ .

Poznamo štiri osnovne tipe parcialnih diferencialnih enačb.

- Linearna PDE je linearna kombinacija parcialnih odvodov

$$\sum_{|\alpha| \leq k} a_\alpha(x) D^\alpha u(x) = f(x).$$

- Semilinearna PDE je linearna samo v odvodu najvišjega reda

$$\sum_{|\alpha|=k} a_\alpha(x) D^\alpha u + b(x, u, Du, D^{k-1}u) = 0.$$

- Kvazilinearna PDE je enačba oblike

$$\sum_{|\alpha|=k} a_\alpha(x, u, Du, \dots, D^{k-1}u) \cdot D^\alpha u + b(x, u, \dots, D^{k-1}u) = 0.$$

- Povsem nelinearna PDE je enačba, ki je nelinearno odvisna od odvodov najvišjega reda.

Vprašanje 1. Definiraj linearne, semilinearne in kvazilinearne PDE.

Definicija. Pravimo, da je problem PDE z začetnim ali robnim pogojem DOBRO POSTAVLJEN, če ima naslednje lastnosti:

- rešitev obstaja,
- rešitev je enolična,
- rešitev je zvezno odvisna od začetnih pogojev.

6.2 Kvazilinearne enačbe prvega reda v dveh spremenljivkah

Obravnavamo enačbe oblike

$$a(x, y, u)u_x + b(x, y, u)u_y = c(x, y, u). \quad (6.1)$$

Reševanja se bomo lotili z metodo karakteristik. Enačbo prvo prepišemo v obliko

$$\langle (a, b, c), (u_x, u_y, -1) \rangle = 0.$$

Ploskvi $z = u(x, y)$, ki reši kvazilinearno enačbo, pravimo INTEGRALNA PLOSKEV enačbe (??). Če je $u : D \subseteq \mathbb{R}^2 \rightarrow \mathbb{R}$, lahko njen graf parametriziramo z $\vec{r}(x, y) = (x, y, u(x, y))$. Ploskev opišemo z normalo $\vec{r}_x \times \vec{r}_y = -(u_x, u_y, -1)$. Če u že imamo, tedaj vektor (a, b, c) leži v tangentni ravnini na graf.

Definicija. Vektor (a, b, c) se imenuje KARAKTERISTIČNA SMER za PDE.

Pišimo $V = (a, b, c)$, čemur pravimo pripadajoče (karakteristično) vektorsko polje. Iz eksistenčnega izreka sledi, da za vsak $p_0 = (x_0, y_0, z_0)$ obstaja natanko ena tokovnica polja V , ki poteka skozi točko p_0 . Če označimo $\varphi_t(p) = \gamma_p(t)$, je φ tok, in velja $\varphi_{t+s} = \varphi_t \circ \varphi_s$.

Trditev. Naj bodo $a, b, c \in C^1(\mathbb{R}^3)$, $S = \{z = u(x, y)\}$ poljubna integralna ploskev za (6.1) in $p_0 \in S$. Tedaj vsaka karakteristična krivulja za (6.1), ki gre skozi točko p_0 , cela leži v S .

Dokaz. Naj bo $\gamma(t) = (x, y, z)$ karakteristična krivulja z $\gamma(0) = p_0$. Dokazati želimo, da velja $z = u(x, y)$, za kar je dovolj pokazati $\dot{z} = \partial_t u(x, y)$. Označimo $w(t) = z(t) -$

$u(x(t), y(t))$ in računamo

$$\begin{aligned}\dot{w} &= \dot{z} - u_x \dot{x} - u_y \dot{y} \\ &= c(\gamma(t)) - u_x(x, y)a(\gamma(t)) - u_y(x, y)b(\gamma(t)) \\ &= c(x, y, w + u(x, y)) - u_x(x, y)a(x, y, w + u(x, y)) - u_y(x, y)b(x, y, w + u(x, y)).\end{aligned}$$

Velja $w(0) = 0$. Dobili smo Cauchyjevo nalogo oblike $\dot{w}(t) = f(t, w(t))$, $w(0) = 0$. Preverimo lahko, da za f lokalno velja Lipschitzov pogoj, torej po eksistenčnem izreku obstaja enolična rešitev. Ker je u po predpostavki rešitev (6.1), je $w = 0$ rešitev Cauchyjeve naloge, ki je enolična. \square

Vprašanje 2. Pokaži, da tokovnice ne zapustijo integralne ploskve.

Posledica. Ob predpostavkah trditve je vsaka integralna ploskev unija karakterističnih krivulj.

Posledica. Ob predpostavkah trditve:

- Če se dve integralski ploskvi S_1, S_2 sekata v točki $p \in S_1 \cap S_2$, cela tokovnica skozi p leži v $S_1 \cap S_2$.
- Če je $C = S_1 \cap S_2$ krivulja, ki je presek dveh integralskih ploskev S_1, S_2 , ki se sekata ne-tangentno, je karakteristična krivulja.

Dokaz. Tokovnica ne zapusti ne S_1 ne S_2 , torej mora biti vsebovana v preseku. Za drug del naj bo $p \in S_1 \cap S_2$. Tangentni ravnini $T_p S_1$ in $T_p S_2$ se sekata ne-tangentno, torej sta različni, in je njun presek enodimenzijski. Ker sta S_1 in S_2 integralski ploskvi, njuni tangentni ravnini vsebujeta smer karakteristike in je zato

$$T_p S_1 \cap T_p S_2 = \{\lambda(a(p), b(p), c(p)) \mid \lambda \in \mathbb{R}\} = T_p(S_1 \cap S_2) = T_p C.$$

Torej je C v vsaki točki tangentna na smer karakteristike, torej je karakteristična krivulja. \square

Vprašanje 3. Dokaži, da je presek integralskih ploskev karakteristična krivulja.

Vsaka integralska ploskev $z = u(x, y)$ za (6.1) je unija karakterističnih krivulj polja $V = (a, b, c)$. Naj bo Γ gladka krivulja, parametrizirana z $\gamma(s) = (f(s), g(s), h(s))$. Skozi $\gamma(s)$ napeljemo karakteristično krivuljo, ki vsebuje točko $\gamma(s)$. Privzamemo lahko, da gre ta krivulja skozi točko $\gamma(s)$ ob času $t = 0$. Iščemo vektorsko funkcijo $R(s, t) = (x(s, t), y(s, t), z(s, t))$, ki zadošča

- za vsak s funkcija $R(s, \cdot)$ reši karakteristični sistem $\dot{x} = a(x, y, z)$, $\dot{y} = b(x, y, z)$, $\dot{z} = c(x, y, z)$,
- začetni pogoji: $x(s, 0) = f(s)$, $y(s, 0) = g(s)$, $z(s, 0) = h(s)$.

Naj bo J kompakten interval, $\gamma : J \rightarrow \mathbb{R}^3$, $s \in J$ in $\vec{F}_s = (f(s), g(s), h(s))$. Rešujemo sistem

$$\frac{d}{dt}\vec{x}_s = V(\vec{x}_s), \quad \vec{x}_s(0) = \vec{F}_s.$$

Rešitev bo podana z $R(s, t) = \vec{x}_s(t)$. Po eksistenčnem izreku za sisteme NDE obstaja $\varepsilon > 0$, da na intervalu $(-\varepsilon, \varepsilon)$ obstaja natanko ena rešitev $\vec{x}_s : (-\varepsilon, \varepsilon) \rightarrow \mathbb{R}^3$. Želeli bi, da je R parametrizacija integralske ploskve. Problem je, če je Γ karakteristična krivulja; tedaj je slika R enodimenzionalna, torej potrebujemo dodatne pogoje.

Poskusimo sistem enačb $x = x(s, t)$, $y = y(s, t)$ prevesti na sistem $s = s(x, y)$, $t = t(x, y)$ in vstaviti v tretjo komponento R . To nam da $z = z(s(x, y), t(x, y)) = u(x, y)$. Kdaj pa to smemo narediti? Spomnimo se: če za neki $s_0 \in J$ velja

$$\frac{\partial(x, y)}{\partial(s, t)}(s_0, 0) = \begin{vmatrix} \frac{\partial x}{\partial s} & \frac{\partial x}{\partial t} \\ \frac{\partial y}{\partial s} & \frac{\partial y}{\partial t} \end{vmatrix} \neq 0,$$

inverz lahko poiščemo lokalno na neki okolici $(s_0, 0)$. Po privzetkih na \vec{F}_s je determinanta enaka

$$\begin{vmatrix} f'(s_0) & a(p_0) \\ g'(s_0) & b(p_0) \end{vmatrix} = b(p_0)f'(s_0) - a(p_0)g'(s_0)$$

za $p_0 = (f(s_0), g(s_0), h(s_0))$. Geometrijsko želimo, da je projekcija na prvi dve komponenti tangente $\dot{\gamma}(s_0)$ linearno neodvisna od projekcije vektorja $V(p_0)$ na prvi dve komponenti.

Trditev. Naj bo $\gamma = (f, g, h)$ pot v \mathbb{R}^3 razreda \mathcal{C}^1 , ki v dani točki $p_0 = \gamma(s_0)$ zadošča pogoju transverzalnosti

$$\begin{vmatrix} f'(s_0) & a(p_0) \\ g'(s_0) & b(p_0) \end{vmatrix} \neq 0.$$

Tedaj v neki okolici točke $(x_0, y_0) = (f(s_0), g(s_0))$ obstaja natanko ena rešitev $u = u(x, y)$ kvazilinearne PDE (6.1), ki zadošča pogoju $h(s) = u(f(s), g(s))$ za vsak s blizu s_0 .

Dokaz. Vse razen enoličnosti smo že pokazali. Denimo, da neka integralna ploskev vsebuje točko p_0 . Tedaj po že dokazanem vsebuje vse karakteristične krivulje skozi točko p_0 . Posledično (vsaj lokalno) celo ploskev parametriziramo z $(x(s, t), y(s, t), z(s, t))$. Če privzamemo, da v neki okolici p_0 obstaja še ena rešitev $w = w(x, y)$, za katero je $h = w(f, g)$ na okolici s_0 , tedaj je $W = \{z = z(x, y)\}$ unija tokovnic. Če gre neka tokovnica iz W skozi točko p_0 , je cela vsebovana v U . Torej $W \subseteq U$. \square

Vprašanje 4. Kako poiščeš rešitev kvazilinearne PDE prvega reda dveh spremenljivk? Katere predpostavke potrebuješ? Pokaži, da je rešitev enolična.

6.2.1 Linearna PDE

To je enačba oblike

$$a(x, y)u_x + b(x, y)u_y = c(x, y)u + d(x, y).$$

V pripadajočem polju $V = (a, b, cu + d)$ funkcije a, b, c, d niso odvisne od u . Enačbe karakteristik imajo obliko

$$\begin{aligned}\dot{x} &= a(x, y), \\ \dot{y} &= b(x, y), \\ \dot{z} &= c(x, y)z + d(x, y).\end{aligned}$$

Vprašanje 5. Kakšne so enačbe karakteristik za linearno PDE?

6.2.2 Ovojnica družine ravnin

Naj bo

$$\Pi_\lambda = \{(x, y, z) \in \mathbb{R}^3 \mid z = \psi(x, y, \lambda) = a(\lambda)x + b(\lambda)y\}$$

družina ravnin, kjer za interval $I \subseteq \mathbb{R}$ funkciji $a, b \in \mathcal{C}^2(I)$ taki, da je $\vec{n}(\lambda) = (a(\lambda), b(\lambda), -1)$ injektivna in

$$W(a', b') = \begin{vmatrix} a' & b' \\ a'' & b'' \end{vmatrix} \neq 0$$

za vsak λ . OVOJNICA družine Π_λ je ploskev

$$\mathcal{O} = \{(x, y, z) \in \mathbb{R}^3 \mid z = \varphi(x, y)\},$$

za katero velja:

- $\mathcal{O} \subseteq \bigcup_\lambda \Pi_\lambda$,
- če je $p \in \mathcal{O} \cap \Pi_\lambda$, potem je normala na \mathcal{O} v točki p vzporedna $\vec{n}(\lambda)$.

Trivialna izbira je $\mathcal{O} = \Pi_\lambda$ za nek λ , te pa v nadaljnje ne bomo upoštevali.

Vprašanje 6. Definiraj ovojnico družine ravnin v \mathbb{R}^3 .

Vzemimo projekcijo $\pi : \mathbb{R}^3 \rightarrow \mathbb{R}^2$ na prvi dve komponenti. Ker zahtevamo, da je $\vec{n}(\lambda)$ injektivna, za vsak par $(x, y) \neq (0, 0)$ obstaja natanko določena točka $p \in \mathcal{O}$, da je $(x, y) = \pi(p)$. To nam definira preslikavo $\Lambda : \pi(\mathcal{O}) \rightarrow \mathbb{R}$, ki slika par (x, y) v pripadajoči λ .

Naš cilj je poiskati funkcijo φ , ki definira \mathcal{O} . Ker je $\mathcal{O} \subseteq \bigcup_\lambda \Pi_\lambda$, velja

$$\varphi(x, y) = a(\Lambda(x, y))x + b(\Lambda(x, y))y.$$

Zahtevali smo, da je normala vzporedna z $\vec{n}(\lambda)$, torej lahko parametriziramo \mathcal{O} z

$$u(x, y) = (x, y, \varphi(x, y)).$$

Upošteva je definicijo Π_λ lahko zapišemo normalo kot

$$u_x \times u_y = (1, 0, \varphi_x) \times (0, 1, \varphi_y) = (-\psi_x - \psi_\lambda \Lambda_x, -\phi_y - \psi_\lambda \Lambda_y, 1)$$

oziroma nasprotno vrednot kot

$$-u_x \times u_y = (\psi_x, \psi_y, -1) + \psi_\lambda(\Lambda_x, \Lambda_y, 0).$$

Velja $\psi_x = a$ in $\psi_y = b$, torej $-u_x \times u_y = \vec{n}(\Lambda(x, y)) + \psi_\lambda(\Lambda_x, \Lambda_y, 0)$. Zahtevamo, da je to vzporedno z $\vec{n}(\Lambda(x, y))$, torej mora biti drugi člen ničeln. Če je $\Lambda_x = \Lambda_y = 0$, je Λ konstantna preslikava in je $\mathcal{O} = \Pi_\Lambda$; to je izključen trivialni primer. Velja torej $\psi_\lambda = 0$, kar nam da dodatno enačbo za ovojnico.

Vprašanje 7. Kako poiščeš ovojnico družine ploskev v \mathbb{R}^3 ? Izpeljži potrební pogoji.

Ta enačba je ekvivalentna pogoju

$$D(\lambda) := \frac{a'(\lambda)}{b'(\lambda)} = \frac{-y}{x}.$$

Po predpostavki o determinanti Wronskega je odvod D neničeln, torej jo lahko na množici $\{x \neq 0\}$ obrnemo in dobimo $\Lambda_1(x, y) = D^{-1}(-y/x)$. Podobno lahko naredimo za $y \neq 0$, iz česar izpeljemo $\Lambda_2(x, y) = E^{-1}(-x/y)$ za $E(\lambda) = b'(\lambda)/a'(\lambda)$. Če pokažemo, da je $\Lambda_1 = \Lambda_2$, kjer sta obe definirani, bomo lahko zapisali

$$\Lambda(x, y) = \begin{cases} \Lambda_1(x, y), & x \neq 0 \\ \Lambda_2(x, y), & y \neq 0 \end{cases}$$

Vemo, da za $(x, y) \in \pi(\mathcal{O})$ obstaja natanko določen λ , pri katerem je (x, y) slika neke točke na \mathcal{O} s π . Ker je tretja koordinata $z = a(\Lambda_1)x + b(\Lambda_1)y = a(\Lambda_2)x + b(\Lambda_2)y$ pri tem natančno določena, iz pogoja injektivnosti \vec{n} lahko izpeljemo $\Lambda_1 = \Lambda_2$.

Iz zgornje izpeljave hitro vidimo, da velja $\Lambda(\sigma x, \sigma y) = \Lambda(x, y)$ za $\sigma \neq 0$, iz česar sledi naslednja trditev.

Trditev. Če je $p \in \mathcal{O}$, potem je $\sigma p \in \mathcal{O}$ za poljuben $\sigma \neq 0$.

6.3 Nelinearne enačbe prvega reda

Rešujemo enačbo oblike

$$F(x, y, u, u_x, u_y) = 0. \quad (6.2)$$

Tu je $F \in \mathcal{C}^1$ funkcija petih realnih spremenljivk in $u = u(x, y)$ iskana funkcija. Vpeljemo oznake $p = u_x$, $q = u_y$. Normala na ploskev $\{z = u(x, y)\}$ je $(u_x, u_y, -1) = (p, q, -1)$. Enačba (6.2) je zveza med točko (x, y, z) na Γ_u in normalo. Družino potencialnih normal parametriziramo z $\vec{n}(\lambda) = (p(\lambda), q(\lambda), -1)$. Za vsak izbrani λ dobimo potencialno tangentno ravnino na ploskev. Enačba te ravnine, ki jo označimo s $\Pi_{\vec{r}_0, \lambda} = \Pi_\lambda$, je

$$\langle \vec{r} - \vec{r}_0, \vec{n}(\lambda) \rangle = 0.$$

Definicija. Ogrinjača ravnin Π_λ se imenuje MONGEOV STOŽEC v točki (x_0, y_0, z_0) .

Ogrinjačo dobimo tako, da enačbo ravnin odvajamo po λ in enačimo z 0. Imamo torej $\langle \vec{r} - \vec{r}_0, \vec{n}(\lambda) \rangle = 0$ in $\langle \vec{r} - \vec{r}_0, \vec{n}'(\lambda) \rangle = 0$. Vektor $\vec{r} - \vec{r}_0$ je torej vzporeden z vektorskih produktom $\vec{n}(\lambda) \times \vec{n}'(\lambda)$. Če je $F(x_0, y_0, z_0, p(\lambda), q(\lambda)) = 0$ in odvajamo po λ , dobimo

$$F_p p' + F_q q' = 0.$$

Sledi $(q', -p') \parallel (F_p, F_q)$, torej (v kombinaciji s prejšnjo enačbo) $\vec{r} - \vec{r}_0 \parallel (F_p, F_q, pF_p + qF_q)$. Dobili smo parametrizacijo Mongeovega stožca

$$\begin{aligned}\vec{r} &= \vec{r}_0 + \mu(P(\lambda), Q(\lambda), R(\lambda)), \\ P(\lambda) &= F_p(x_0, y_0, z_0, p(\lambda), q(\lambda)), \\ Q(\lambda) &= F_q(x_0, y_0, z_0, p(\lambda), q(\lambda)), \\ R(\lambda) &= p(\lambda)P(\lambda) + q(\lambda)Q(\lambda).\end{aligned}$$

Privzemimo, da imamo integralno ploskev $S = \{z = u(x, y)\}$. Naj bo $\gamma(t) = (x(t), y(t), z(t))$ neka krivulja na tej ploskvi. Tedaj $\dot{\gamma}(t)$ leži na tangentni ravnini, zato $\vec{r}_0 + \dot{\gamma}(t)$ leži v Mongeovem stožcu. Posledično mora veljati (oz. lahko zahtevamo)

$$\begin{aligned}\dot{x} &= F_p, \\ \dot{y} &= F_q, \\ \dot{z} &= pF_p + qF_q,\end{aligned}$$

kjer je $p = p(\lambda)$ in $q = q(\lambda)$ za neki (ne vemo kateri) λ . Ker p in q v splošnem ne poznamo, ju obravnavamo kot novi neznanki. Vsaki točki na γ bomo dodali še en košček ravnine, določene z normalo $(p, q, -1)$, s čimer dobimo KARAKTERISTIČNI TRAK.

Trditev. Če krivulja $\gamma(t) = (x, y, z)$ zadošča pogoju $\dot{z} = p\dot{x} + q\dot{y}$, leži v karakterističnem traku, določenem s $p = u_x$ in $q = u_y$.

Dokaz. Računamo $\langle (\dot{x}, \dot{y}, \dot{z}), (p, q, -1) \rangle = \langle (\dot{x}, \dot{y}, p\dot{x} + q\dot{y}), (p, q, -1) \rangle = 0$. □

Če odvajamo zgornjo enačbo, dobimo

$$\begin{aligned}F_x + F_u u_x + F_p p_x + F_q q_x &= 0, \\ F_y + F_u u_y + F_p p_y + F_q q_y &= 0.\end{aligned}$$

Za $x = x(t)$ in $y = y(t)$ imamo torej

$$\begin{aligned}\dot{p} &= p_x \dot{x} + p_y \dot{y} = p_x F_p + p_y F_q, \\ \dot{q} &= q_x F_p + q_y F_q,\end{aligned}$$

iz česar sledi

$$\dot{p} = -F_x - F_u u_x - F_q q_x + p_y F_q = F_q(p_y - q_x) - F_x - F_u u_x$$

Pričakujemo $p = u_x$ in $q = u_y$, torej $p_y - q_x = u_{xy} - u_{yx} = 0$. Posledično je smiselno zahtevati $\dot{p} = -F_x - F_u p$ in $\dot{q} = -F_y - F_u q$. S tem smo dobili karakteristični sistem za x, y, z, q, p :

$$\begin{aligned}\dot{x} &= F_p \\ \dot{y} &= F_q \\ \dot{z} &= pF_p + qF_q \\ \dot{p} &= -F_x - F_u p \\ \dot{q} &= -F_y - F_u q\end{aligned}$$

Rešitve imenujemo KARAKTERISTIKE ENAČBE (6.2). Funkcija F je konstantna vzdolž rešitev tega sistema.

Vprašanje 8. Izpelji karakteristični sistem za nelinearno PDE prvega reda.

6.3.1 Cauchyjeva naloga za PDE prvega reda

Vzemimo karakteristiko $\lambda = \lambda(t) = (x, y, z, p, q) = (\gamma, p, q)$. Imamo še začetni pogoj $(x_0, y_0, z_0, p_0, q_0) \in \mathbb{R}^5$.

Definicija. Krivulja $\lambda = \lambda(s)$, ki zadošča $F(\lambda) = 0$, se imenuje INTEGRALNA KRIVULJA, če velja $z_s = px_s + qy_s$.

Rešitev karakterističnega sistema je integralna krivulja. Cauchyjeva naloga je reševanje (6.2), pri čemer graf rešitve vsebuje neko vnaprej dano krivuljo $\gamma(s) = (f(s), g(s), h(s))$. Želimo poiskati smiselna začetna pogoja za p in q , torej iščemo funkciji φ in ψ , da bo $(f, g, h, \varphi, \psi)(s)$ začetni pogoj za sistem. Rešitev sistema bo tedaj funkcija $R = R(s, t)$, za katero bo veljalo $R(s, 0) = (f, g, h, \varphi, \psi)(s)$, in da $(x(s, t), y(s, t), z(s, t))$ parametrizira $\{z = u(x, y)\}$, kjer je u rešitev enačbe (6.2).

Da bo R parametrizirala rešitev enačbe, mora veljati $F(R(s, t)) = 0$. Ker je F vzdolž tokovnic konstantna, zadošča zahtevati, da je $F(R(s, t)) = 0$, torej $F(f, g, h, \varphi, \psi) = 0$. Veljati mora

$$\begin{aligned}\dot{z} &= u_x \dot{x} + u_y \dot{y}, \\ z_s &= u_x x_s + u_y y_s.\end{aligned}$$

Ker želimo $p = u_x$ in $q = u_y$, pa dodatno

$$\begin{aligned}\dot{z} &= p\dot{x} + q\dot{y}, \\ z_s &= px_s + qy_s.\end{aligned}$$

Iz pogojev pri $t = 0$ torej dobimo

$$h'(s) = \varphi(s)f'(s) + \psi(s)g'(s).$$

Zgornjo enačbo lahko zapišemo tudi v matrični obliki

$$\begin{bmatrix} z_s \\ \dot{z} \end{bmatrix} = \begin{bmatrix} x_s & y_s \\ \dot{x} & \dot{y} \end{bmatrix} \cdot \begin{bmatrix} p \\ q \end{bmatrix}. \quad (6.3)$$

Če je determinanta te matrike neničelna, lahko lokalno izrazimo $s(x, y)$ in $t(x, y)$. Posledično dobimo

$$\begin{bmatrix} p \\ q \end{bmatrix} = \begin{bmatrix} s_x & t_x \\ s_y & t_y \end{bmatrix} \cdot \begin{bmatrix} z_s \\ \dot{z} \end{bmatrix},$$

iz česar sklepamo $u_x = p$ in $u_y = q$.

Na Γ je $x = x(s, 0) = f(s)$, zato je $x_s(\cdot, 0) = f'$ in podobno $y_s(\cdot, 0) = g'$. Za vsak s naj $x(\cdot, t)$ reši karakteristični sistem. Od tod sledi $\dot{x}(s, 0) = F_p(\Gamma(s))$ in $\dot{y}(s, 0) = F_q(\Gamma(s))$. To pomeni, da lahko pogoj neničelne determinante izrazimo s podatki kot

$$\det \begin{bmatrix} f' & g' \\ F_p \circ \Gamma & F_q \circ \Gamma \end{bmatrix} \neq 0.$$

Geometrijska interpretacija tega je, da projekcija Γ in tokovnice nista vzporedni v presečišču $(s, 0)$. Zato temu pravimo **POGOJ TRANSVERZALNOSTI**.

Če povzamemo, začetne pogoje za Cauchyjevo nalogo dobimo tako, da ob dani krivulji $\gamma(s) = (f, g, h)(s)$ poiščemo še $p = \phi(s)$ in $q = \psi(s)$, da za $\Gamma = (f, g, h, \phi, \psi)$ velja

- $F(\Gamma(s)) = 0$,
- $h' = \phi f' + \psi g'$,
- $f'(s)F_q(\Gamma(s)) - g'(s)F_p(\Gamma(s)) \neq 0$.

Za vsak s rešimo karakteristični sistem z začetnimi pogoji zgoraj. Na ta način dobimo $R(s, t) = (x, y, z, p, q)(s, t)$, za katero velja

- $R(s, 0) = (f, g, h, \phi, \psi)(s)$,
- projekcija R na prve tri komponente je lokalno graf funkcije.

Druga točka sledi iz dejstva, da je Jacobijeva determinanta

$$\det \frac{\partial(x, y)}{\partial(s, t)} \neq 0.$$

Če označimo $u(x, y) = z(s(x, y), t(x, y))$, bo to rešitev (6.2), dokaz sledi.

Preverili bomo, da funkcija u reši sistem (6.3). Ker na ploskvi $\Sigma = \{R(s, t)\}$ velja $F = 0$, bo sledilo $p = u_x$ in $q = u_y$. Verižno pravilo namreč pove

$$\begin{bmatrix} z_s \\ \dot{z} \end{bmatrix} = A \begin{bmatrix} u_x \\ u_y \end{bmatrix},$$

kjer smo z A označili matriko iz (6.3). Ker je A lokalno obrnljiva, bo torej veljalo $p = u_x$ in $q = u_y$. Drugi pogoj iz (6.3) (tj. $\dot{z} = p\dot{x} + q\dot{y}$) sledi iz karakterističnega sistema, za prvi pogoj pa označimo $B(s, t) = z_s - px_s - qy_s$. Velja

$$B(s, 0) = h'(s) - \phi(s)f'(s) - \psi(s)g'(s) = 0$$

in

$$\begin{aligned}\partial_t B(s, t) &= z_{st} - \dot{p}x_s - px_{st} - \dot{q}y_s - qy_{st} \\ &= \partial_s(\dot{z} - p\dot{x} - q\dot{y}) + p_s\dot{x} + q_s\dot{y} - \dot{p}x_s - \dot{q}y_s \\ &= F_p p_s + F_q q_s + (F_x + pF_z)x_s + (F_y + qF_z)y_s \\ &= (F \circ R)_s - F_z z_s + pF_z x_s + qF_z y_s \\ &= F_z(px_s + qy_s - z_s) \\ &= -F_z B.\end{aligned}$$

Torej B reši diferencialno enačbo

$$\dot{B} = -F_z B,$$

iz česar sledi

$$B(s, t) = B(s, 0) \exp\left(-\int_0^t F_z(\tau) d\tau\right) = 0.$$

S tem smo dokazali eksistenčni izrek za nelinearne PDE prvega reda.

Izrek. Naj bo $F(x, y, z, p, q)$ dana \mathcal{C}^2 funkcija, $\gamma = (f, g, h)$ dana \mathcal{C}^1 krivulja, $\phi, \psi \in \mathcal{C}^1$, $V_0 = \Gamma(s_0) = (f, g, h, \phi, \psi)(s_0)$ za neki s_0 in naj na neki okolici točke s_0 veljajo pogoji

- $F(\Gamma(s)) = 0$,
- $h' = \phi f' + \psi g'$,
- $f'(s)F_q(\Gamma(s)) - g'(s)F_p(\Gamma(s)) \neq 0$.

Tedaj obstaja $\varepsilon > 0$, da na krogli $K = K((s_0, 0), \varepsilon)$ obstaja natanko ena funkcija $R(s, t) = (x, y, z, p, q)(s, t)$, da velja

- $R(\cdot, 0) = \Gamma$ za $|s - s_0| < \varepsilon$,
- $F \circ R = 0$ za K ,
- za vsak $s \in K(s_0, \varepsilon) \subseteq \mathbb{R}$ funkcija $R(s, \cdot)$ reši karakteristični sistem za (6.2),
- projekcija R na prve tri komponente določa parametrično podano integralno ploskev $\{z = u(x, y)\}$ za (6.2).

Vprašanje 9. Izpelji rešitev Cauchyjeve naloge za nelinearno PDE prvega reda.

6.4 Lagrangeova metoda

Definicija. Realna funkcija $\phi(x, y, z)$ se imenuje PRVI INTEGRAL \mathcal{C}^1 polja $V = (a, b, c)$, če je konstantna vzdolž vsake krivulje $\gamma(t)$, ki reši $\dot{\gamma} = V(\gamma)$.

Vprašanje 10. Kaj je prvi integral?

Velja

$$\partial_t \phi(\gamma(t)) = \left\langle \vec{\nabla} \cdot \phi(\gamma(t)), \dot{\gamma}(t) \right\rangle,$$

torej je ϕ prvi integral za V natanko tedaj, ko je za vsako tokovnico γ za V ta skalarni produkt enak 0.

Trditev. Če u reši kvazilinearno enačbo $au_x + bu_y = c$, je $\psi(x, y, z) = z - u(x, y)$ prvi integral za prirejeno polje $V = (a, b, c)$.

Dokaz. Izračunamo $\left\langle \vec{\nabla} \cdot \psi, V \right\rangle = -au_x - bu_y + c = 0$. □

Trditev. Vsaka nivojna ploskev $\Sigma = \{\phi = C\}$ prvega integrala ϕ je unija karakteristik polja V .

Dokaz. Vzemimo točko $p_0 \in \Sigma$ in rešimo Cauchyjevo nalogo $\dot{\gamma} = V(\gamma)$, $\gamma(0) = p_0$. Tako dobimo tokovnico, ki vsebuje p_0 . Obratno, tokovnica, ki seka Σ v točki t_0 , cela leži v Σ , saj je ϕ prvi integral. □

Integralne ploskve za kvazilinearno enačbo so torej nivojnice prvih integralov polja $V = (a, b, c)$. Nivojnica prvega integrala pa je integralna ploskev samo, če na njej lahko izrazimo z kot funkcijo spremenljivk x, y .

Trditev. Naj bo $\phi \in \mathcal{C}^1$ prvi integral za $V = (a, b, c)$ v okolici točke $p_0 \in \mathbb{R}^3$. Označimo $d = \phi(p_0)$. Če je $\partial_z \phi(p_0) \neq 0$, je nivojnica $S = \{\phi = d\}$ v okolici integralna ploskev za kvazilinearno enačbo.

Dokaz. Ker je odvod neničeln, je S lokalno graf neke \mathcal{C}^1 funkcije u . Dokazati želimo, da u reši kvazilinearno enačbo. Lokalno velja $\phi(x, y, u(x, y)) = d$; če to odvajamo po x in y , dobimo

$$\left\langle \vec{\nabla} \cdot \phi, (1, 0, u_x) \right\rangle = 0,$$

$$\left\langle \vec{\nabla} \cdot \phi, (0, 1, u_y) \right\rangle = 0.$$

Torej je $\vec{\nabla} \cdot \phi$ vzporeden vektorskemu produktu $(1, 0, u_x) \times (0, 1, u_y) = (-u_x, -u_y, 1)$. Ker pa je ϕ prvi integral polja $V = (a, b, c)$, je nanj pravokoten, torej sta vektorja (a, b, c) in $(-u_x, -u_y, 1)$ pravokotna, kar je ravno kvazilinearna enačba. □

Vprašanje 11. Dokaži, da so nivojnice prvega integrala polja (a, b, c) integralne ploskve za kvazilinearno enačbo $au_x + bu_y = c$.

Definicija. Če je ϕ prvi integral za V in $d \in \mathbb{R}$, se družina nivojnic $\{\phi = d\}$ imenuje SPLOŠNA REŠITEV za kvazilinearno enačbo $au_x + bu_y = c$.

Izrek. Naj bosta ϕ_1 in ϕ_2 neodvisna prva integrala polja (a, b, c) v okolici točke $p \in \mathbb{R}^3$. Splošna rešitev $au_x + bu_y = c$ je v okolici točke p podana z $\{F(\phi_1, \phi_2) = 0\}$, kjer je $F = F(x, y)$ poljubna funkcija razreda \mathcal{C}^1 .

Dokaz. Preverimo, da je $F \circ (\phi_1, \phi_2) = 0$ res posplošena rešitev. Če je γ tokovnica in je $\phi_1 = d_1, \phi_2 = d_2$ na γ , je $F(\phi_1, \phi_2) \circ \gamma$ konstanta.

V drugo smer preverimo, da je rešitev res take oblike. Vzemimo $p_0 \in \mathbb{R}^3$, ki leži v definicijskem območju V . Vemo, da je $\langle \vec{\nabla} \cdot \phi(p_0), V(p_0) \rangle = 0$ za vsak prvi integral ϕ . Naj bo ϕ_3 še en prvi integral. Tedaj je $V(p_0) \in \ker D(\phi_1, \phi_2, \phi_3)$, ker je pravokoten na vse gradiente. Hkrati je $V \neq 0$ povsod, ker bi sicer nekje veljalo $\dot{\gamma} = 0$. Torej je Jacobijeva matrika izrojena in so funkcije ϕ_1, ϕ_2 in ϕ_3 odvisne, zato obstaja neničelna funkcija $G = G(x, y, z)$, da lokalno velja $G(\phi_1, \phi_2, \phi_3) = 0$. Nivojnica $\{\phi_3 = d\}$ je lokalno vsebovana v $\{G(\phi_1, \phi_2, \phi_3) = 0\}$, kar je natanko $\{F(\phi_1, \phi_2) = 0\}$, če za F vzamemo $F(x, y) = G(x, y, d)$. \square

Vprašanje 12. Utemelji pravilnost metode prvih integralov.

Izrek. Naj bo V vektorsko polje razreda \mathcal{C}^1 v okolici točke $p \in \mathbb{R}^3$ in naj bo $V(p) \neq 0$. Tedaj obstaja \mathcal{C}^1 vektorsko polje ψ , definirano na okolici točke p , da je $\langle \vec{\nabla} \cdot \psi_j, V \rangle_{\mathbb{R}^3} = \delta_{1j}$ v okolici p .

Dokaz. Privzemimo prvo $p = 0$ in $V(p) = (1, 0, 0)$. Naj bo ϕ_t tok polja V , tj. družina \mathcal{C}^1 vektorskih polj, za katera je $\phi_t(x) = V(\phi_t(x))$ in $\phi_0(x) = x$. V okolici točke 0 definiramo še $g(x_1, x_2, x_3) = \phi_{x_1}(0, x_2, x_3)$. Dobimo

$$\partial_{x_1} g(x) = \lim_{h \rightarrow 0} \frac{\phi_{x_1+h}(0, x_2, x_3) - \phi_{x_1}(0, x_2, x_3)}{h} = \partial_{x_1} \phi((0, x_2, x_3), x_1) = V(g(x))$$

za $\phi(x, t) = \phi_t(x)$. Velja še $g(0, x_2, x_3) = (0, x_2, x_3)$, zato je $\partial_{x_j} g = e_j$ na ravnini $x_1 = 0$. Iz $V(g(0)) = v(0) = (1, 0, 0)$ sledi $Dg(0) = I_3$, zato je g difeomorfizem na neki okolici točke 0. Označimo $\psi = g^{-1}$. Če pišemo $x = (x_1, x_2, x_3)$ in $\psi = (\psi_1, \psi_2, \psi_3)$, tedaj iz zveze $\psi \circ g = \text{id}$ sledi $\psi_j(g(x)) = x_j$. Odvajamo po x_1 in dobimo

$$\delta_{1j} = \sum_{k=1}^3 \partial_{y_k} \psi_j(g(x)) \partial_{x_1} g_k(x) = \sum_{k=1}^3 \partial_{y_k} \psi_j(g(x)) V_k(g(x)) = \langle \vec{\nabla} \cdot \psi_j(g(x)), V(g(x)) \rangle.$$

Ker je g difeomorfizem, velja $\langle \vec{\nabla} \cdot \psi_j, V \rangle = 0$ na neki okolici 0.

Sedaj obravnavajmo splošni primer. Če sta $p, V(p) \in \mathbb{R}^3$ poljubna, vzamemo ortogonalno $A \in \mathbb{R}^{3 \times 3}$, s katero se $\frac{V(p)}{\|V(p)\|}$ slika v $(1, 0, 0)$. Označimo $\lambda = \|V(p)\|^{-1}$ in $w(x) = \lambda A(V(x + p))$. Ta w slika iz okolice točke 0, velja $w(0) = (1, 0, 0)$. Torej obstaja

vektorsko polje η , definirano na okolici 0, da je $\langle \vec{\nabla} \cdot \eta_j, w \rangle = \delta_{1j}$. Sedaj definiramo še ψ tako, da velja $\eta(x) = \lambda^{-1} A\psi(x+p)$. Izračunamo lahko, da velja

$$\langle \vec{\nabla} \cdot \eta_j(x), w(x) \rangle = \langle \vec{\nabla} \cdot \psi_j(x+p), V(x+p) \rangle.$$

□

Vprašanje 13. Dokaži, da za vsako vektorsko polje V in točko p , za katero $V(p) \neq 0$, obstaja polje ψ , da je v okolici p skalarni produkt $\langle \vec{\nabla} \cdot \psi_j, V \rangle = \delta_{1j}$.

Posledica. Naj bo V vektorsko polje razreda \mathcal{C}^1 v okolici točke $p \in \mathbb{R}^3$ in naj bo $V(p) \neq 0$. Tedaj v okolici točke p obstajata dva neodvisna prva integrala.

Dokaz. To sta ψ_1 in ψ_2 iz dokaza izreka. □

Vprašanje 14. Utemelji, da obstajata dva neodvisna prva integrala kvazilinearne enačbe.

6.4.1 Odvisnost funkcij

Definicija. Naj bo $U^{\text{odp}} \subseteq \mathbb{R}^n$, $S \subseteq U$ in $f_1, \dots, f_n \in \mathcal{C}^1(U)$. Funkcije f_1, \dots, f_n so FUNKCIJSKO ODVISNE na S , če obstaja $V^{\text{odp}} \subseteq \mathbb{R}^n$, za katero je $(f_1, \dots, f_n)_*(S) \subseteq V$ in če obstaja \mathcal{C}^1 funkcija $F : V \rightarrow \mathbb{R}$, za katero je $F(f_1, \dots, f_n) = 0$ na S , hkrati pa ni prasluka $\{F = 0\}$ nikjer gosta.

Vprašanje 15. Definiraj funkcijsko odvisnost.

Definicija. KRITIČNA TOČKA preslikave $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$ je vsaka točka $p \in \mathbb{R}^m$, kjer je $\text{rang } Df(p) < n$.

Izrek (izrek o rangju). Naj bo $m \geq n$. Naj bosta $U^{\text{odp}} \subseteq \mathbb{R}^m$, $V^{\text{odp}} \subseteq \mathbb{R}^n$ odprti množici, $f \in \mathcal{C}^1(U, V)$ ter $p \in U$ točka maksimalnega ranga za f . Potem obstajata odprti množici $U_1, U_2 \subseteq \mathbb{R}^m$, za kateri je $p \in U_1 \subseteq U$, in difeomorfizem $X : U_1 \rightarrow U_2$, da velja

$$(f \circ X^{-1})(u_1, \dots, u_m) = (u_{m-n+1}, \dots, u_m)$$

za vse $u \in U_2$.

Dokaz. Za $m = n$ je to natanko izrek o inverzni preslikavi. Naj bo torej $m > n$ in $k = m - n$. Ker je $Df(p)$ maksimalnega ranga, ima n linearno neodvisnih stolpcev; brez škode za splošnost je torej $Df(p) = [C, D]$, kjer je D obrnljiva $n \times n$ matrika. Pišimo $x = (y, z)$ in definirajmo $g(y, z) = (y, f(y, z))$. Potem je

$$Dg(p) = \begin{bmatrix} I_k & 0 \\ C & D \end{bmatrix}$$

in $\det Dg(p) \neq 0$. Po izreku o inverzni preslikavi obstaja okolica U_1 točke p , za katero je $X = g|_{U_1}$ difeomorfizem. Preverimo lahko, da res ustreza želeni lastnosti. □

Vprašanje 16. Povej in dokaži izrek o rangju.

Posledica. Pod enakimi predpostavkami obstaja okolica U točke p , za katero je $f_*(U)$ odprta.

Dokaz. Vzemimo $U = U_1$ iz izreka. Potem je

$$f_*(U_1) = f \circ X^{-1} \circ X_*(U_1) = f \circ X_*^{-1}(U_2),$$

ta preslikava pa je projekcija in zato odprta. \square

Izrek (Sard). Naj bo $\Omega^{\text{odp}} \subseteq \mathbb{R}^m$ in $f : \Omega \rightarrow \mathbb{R}^m$ preslikava ranga C^r za $r = 1 + \max\{0, m - n\}$. Označimo množico kritičnih točk preslikave f z A . Potem ima $f_*(A)$ mero 0.

Vprašanje 17. Povej Sardov izrek.

Lema. Zaprta podmnožica v \mathbb{R}^n je množica ničel gladke funkcije $\mathbb{R}^n \rightarrow \mathbb{R}$.

Dokaz. Naj bo X zaprta podmnožica \mathbb{R}^n . Definiramo $U = \mathbb{R}^n \setminus X$ in

$$K_m = \{x \in U \mid d(x, X) \geq 2^{-m}\} \cap \overline{K(0, m)}$$

za $m \in \mathbb{N}$. Dodatno definiramo $L_m = K_m \setminus \text{Int } K_{m-1}$. Ta množica je kompaktna, torej jo lahko pokrijemo s končno mnogo krogli $B_1^m, \dots, B_{k_m}^m$ z radijem 2^{-m-1} . Velja $B_j^m \cap X = \emptyset$.

Množica $B = \{B_1^m, \dots, B_{k_m}^m\}_m$ pokrije unijo

$$\bigcup_m L_m = \mathbb{R}^n \setminus X$$

in je lokalno končna. Sedaj lahko definiramo gladko funkcijo φ_j^m tako, da je neničelna natanko na B_j^m , in

$$\varphi = \sum_{m,j} \varphi_j^m.$$

\square

Vprašanje 18. Dokaži: zaprte množice v \mathbb{R}^n so natanko ničle gladih funkcij.

Izrek. Naj bo $\Omega^{\text{odp}} \subseteq \mathbb{R}^m$ in $f_1, \dots, f_n \in C^r(\Omega)$ za $r = 1 + \max\{0, m - n\}$. Potem je množica $\{f_1, \dots, f_n\}$ funkcijsko odvisna na vsaki kompaktni podmnožici $S \subseteq \Omega$ natanko tedaj, ko je $\text{rang } D(f_1, \dots, f_n) < n$ za vsak $x \in \Omega$.

Dokaz. V desno: Recimo, da obstaja točka $p \in \mathbb{R}^n$, kjer je Df ranga n . Po izreku o rangju obstaja odprta okolica $U \ni p$, katere f -slika je odprta. Potem $f_*(\overline{U})$ ni nikjer gosta, torej ne obstaja F iz definicije funkcijske odvisnosti.

V levo: Po Sardovem izreku ima $f_*(\Omega)$ mero nič. Če je $K \subseteq \Omega$ kompakt, je kompaktna tudi njegova f -slika, ki ima mero nič. Torej je $f_*(K)$ zaprta in obstaja $\varphi \in \mathcal{C}^\infty(\Omega)$, da je $f_*(K) = \varphi^{-1}(0)$. Ta preslikava ustreza pogojem za F . \square

Vprašanje 19. Karakteriziraj funkcijsko odvisnost in dokaži karakterizacijo.

6.5 Enačbe drugega reda

Zanimajo nas enačbe v dveh spremenljivkah oblike

$$Lu = au_{xx} + 2bu_{xy} + cu_{yy} + du_x + eu_y + fu = g,$$

kjer so a, b, \dots, g dane funkcije. Vpeljimo kvadratno formo, ki je prirejena glavnemu delu operatorja L :

$$\sigma_{(x,y)}(\xi, \eta) = a(x, y)\xi^2 + 2b(x, y)\xi\eta + c(x, y)\eta^2.$$

Ob izbrani točki $(x, y) \in \mathbb{R}^2$ je to kvadratna forma na \mathbb{R}^2 . Imenuje se **SIMBOL** operatorja L . Sedaj izračunamo

$$\sigma_{(x,y)}(\xi, \eta) = \left\langle \begin{bmatrix} a(x, y) & b(x, y) \\ b(x, y) & c(x, y) \end{bmatrix} \cdot \begin{bmatrix} \xi \\ \eta \end{bmatrix}, \begin{bmatrix} \xi \\ \eta \end{bmatrix} \right\rangle_{\mathbb{R}^2},$$

in označimo zgornjo matriko z $A(x, y)$. Poleg tega označimo $\delta = b^2 - ac$.

Definicija. Operator L je v točki (x, y)

- ELIPTIČEN, če je $\delta < 0$,
- HIPERBOLIČEN, če je $\delta > 0$,
- PARABOLIČEN, če je $\delta = 0$.

Vprašanje 20. Kakšne vrste diferencialnih operatorjev drugega reda poznamo?

Imejmo $L_0 u = au_{xx} + 2bu_{xy} + cu_{yy}$ in funkciji $\xi, \eta : \mathbb{R}^2 \rightarrow \mathbb{R}$. Naj bo $w(\xi, \eta) = w(\xi(x, y), \eta(x, y)) = u(x, y)$ pretvorba v nove koordinate. Z verižnim pravilom pride-mo do

$$Lu = \tilde{L}w = Aw_{\xi\xi} + 2Bw_{\xi\eta} + Cw_{\eta\eta} + Dw_\xi + Ew_\eta + Fw$$

za

$$\begin{aligned} A &= a\xi_x^2 + 2b\xi_x\xi_y + c\xi_y^2 \\ B &= a\xi_x\eta_x + b(\xi_x\eta_y + \xi_y\eta_x) + c\xi_y\eta_y \\ C &= a\eta_x^2 + 2b\eta_x\eta_y + c\eta_y^2. \end{aligned}$$

Enačbe lahko prepisemo v obliko

$$\begin{bmatrix} A & B \\ B & C \end{bmatrix} = \begin{bmatrix} \xi_x & \xi_y \\ \eta_x & \eta_y \end{bmatrix} \cdot \begin{bmatrix} a & b \\ b & c \end{bmatrix} \cdot \begin{bmatrix} \xi_x & \xi_y \\ \eta_x & \eta_y \end{bmatrix},$$

iz česar sledi

$$\begin{vmatrix} A & B \\ B & C \end{vmatrix} = \begin{vmatrix} a & b \\ b & c \end{vmatrix} \cdot \begin{vmatrix} \xi_x & \xi_y \\ \eta_x & \eta_y \end{vmatrix}^2,$$

torej sta determinanti istega predznaka in se tip parcialnega diferencialnega operatorja ne spremeni s spremembo koordinat.

Izrek. Naj bo L hiperboličen na domeni $D \subseteq \mathbb{R}^2$. Teda v okolici točke $(x_0, y_0) \in D$ obstaja koordinatni sistem (ξ, η) , v katerem ima L obliko $\tilde{L}(w) = 2Bw_{\xi\eta} + l_1(w)$, kjer je l_1 linearen PDO prvega reda.

Dokaz izreka je izpeljava. Želimo, da bo $A = C = 0$, kar razpišemo v enačbo

$$a\xi_x^2 + 2b\xi_x\xi_y + c\xi_y^2 = 0$$

oz. podobno za η . Enačbi razrešimo na ξ_x/ξ_y oz. η_x/η_y , in dobimo

$$a\xi_x + (b \pm \sqrt{\delta})\xi_y = 0,$$

$$a\eta_x + (b \pm \sqrt{\delta})\eta_y = 0.$$

Hiperboličnost zagotavlja $\delta > 0$, torej imata enačbi realni rešitvi. Vzemimo $+$ pri ξ in $-$ pri η . To sta linearni PDE prvega reda, za kateri dobimo sistem $\dot{x} = a$ in $\dot{y} = b \pm \sqrt{\delta}$. Če je (lokalno) $a \neq 0$, lahko zapišemo sistem

$$\frac{dy}{dx} = \frac{b \pm \sqrt{\delta}}{a}.$$

Rešitvi $\xi(x, y)$ in $\eta(x, y)$ predstavljata novi koordinati, v katerih ima L zeleno obliko.

Vprašanje 21. Izpelji spremembo koordinat za pretvorbo hiperbolične enačbe v standardno obliko.

Izrek. Naj bo L paraboličen na $D \subseteq \mathbb{R}^2$. Teda obstajajo koordinate (ξ, η) , da ima operator L obliko $Aw_{\xi\xi} + l_1(w) = 0$, kjer je l_1 linearen diferencialni operator prvega reda.

Brez škode za splošnost privzamemo $|a| > 0$ (sicer to velja za c in zamenjamo vlogo x, y). Tokrat želimo $B = C = 0$; ker pa paraboličnost zagotavlja $B^2 = AC$, je dovolj zahtevati $C = 0$. Upošteva je $b^2 = ac$ lahko zapišemo

$$C = \frac{1}{a}(a\eta_x + b\eta_y)^2,$$

torej iščemo η , ki bo rešitev $a\eta_x + b\eta_y = 0$. Karakteristični sistem nam da enačbo

$$\frac{dy}{dx} = \frac{b}{a}.$$

Konstanta v rešitvi nam da $\eta(x, y)$, za ξ pa lahko izberemo poljubno funkcijo, neodvisno od η .

Vprašanje 22. Izpelji spremembo koordinat za pretvorbo parabolične enačbe v standardno obliko.

Izrek. Naj bo $L : u \mapsto au_{xx} + 2bu_{xy} + cu_{yy}$ eliptičen parcialni diferencialni operator z realno analitičnimi koeficienti a, b, c na $D \subseteq \mathbb{R}^2$. Tedaj lokalno obstajata koordinati (ξ, η) , v katerih ima L kanonično obliko $\tilde{L}w = D(w_{\xi\xi} + w_{\eta\eta}) + l_1(w)$, kjer je l_1 linearen diferencialen operator prvega reda.

Tokrat zahtevamo $A = C$ in $B = 0$, iz česar dobimo

$$\begin{aligned} a\xi_x^2 + 2b\xi_x\xi_y + c\xi_y^2 &= a\eta_x^2 + 2b\eta_x\eta_y + c\eta_y^2, \\ a\xi_x\eta_x + b(\xi_x\eta_y + \xi_y\eta_x) + c\xi_y\eta_y &= 0. \end{aligned}$$

Sedaj vpeljimo $\phi = \xi + i\eta$, s čimer pridemo do

$$a\phi_x^2 + 2b\phi_x\phi_y + c\phi_y^2 = 0.$$

Enako smo dobili v hiperboličnem primeru, le da je bil tam $\delta > 0$. Sistem je torej ekvivalenten

$$a\phi_x + (b \pm i\sqrt{ac - b^2})\phi_y = 0.$$

Tudi to enačbo rešujemo s pomočjo karakteristik. Funkcije, ki nastopajo, so po predpostavki realno analitične, torej jih lahko razširimo do kompleksnih holomorfnih funkcij dveh spremenljivk. Potem je ϕ rešitev

$$\frac{dy}{dx} = \frac{b \pm i\sqrt{ac - b^2}}{a}$$

v smislu, da je konstantna vzdolž karakteristik. Na koncu dobimo $\xi = \operatorname{Re} \phi$ in $\eta = \operatorname{Im} \phi$.

Vprašanje 23. Izpelji spremembo koordinat za pretvorbo eliptične enačbe v standardno obliko. Katero dodatno predpostavko potrebujemo?

6.5.1 Cauchyjev problem

Imamo enačbo oblike $au_{xx} + 2bu_{xy} + cu_{yy} = d$ in krivuljo $\gamma(s) = (f(s), g(s))$, na kateri želimo predpisati vrednosti rešitve u in njenega gradienta, $u = h(s)$, $u_x = \varphi(s)$ ter $u_y = \psi(s)$. Za poljubno \mathcal{C}^1 funkcijo $v(x, y)$ mora tako veljati

$$\partial_s v|_{\gamma} = v_x(x, y)f'(s) + v_y(x, y)g'(s),$$

iz česar dobimo zvezo, ki mora veljati za podatke, če naj bo problem dobro postavljen:

$$h'(s) = \varphi(s)f'(s) + \psi(s)g'(s).$$

Vprašanje 24. Opiši Cauchyjev problem za PDE drugega reda. Kateri pogoj mora veljati, da je dobro postavljen?

Pogoj lahko nadaljujemo z odvodi višjega reda, do

$$\begin{aligned} \partial_s u_x &= u_{xx}(\gamma)f' + u_{xy}(\gamma)g', \\ \partial_s u_y &= u_{yx}(\gamma)f' + u_{yy}(\gamma)g', \end{aligned}$$

iz česar na γ sledi

$$\begin{aligned}f'u_{xx} + g'u_{xy} &= \varphi', \\f'u_{xy} + g'u_{yy} &= \psi'.\end{aligned}$$

To je sistem treh linearnih enačb za u_{xx} , u_{xy} in u_{yy} :

$$\begin{aligned}au_{xx} + 2bu_{xy} + cu_{yy} &= d, \\f'u_{xx} + g'u_{xy} &= \varphi', \\f'u_{xy} + g'u_{yy} &= \psi'.\end{aligned}$$

Sistem je enolično rešljiv natanko tedaj, ko je njegova determinanta neničelna. Torej

$$\Delta = \begin{vmatrix} a & 2b & c \\ f' & g' & 0 \\ 0 & g' & g' \end{vmatrix} = ag'^2 - 2bf'g' + cf'^2 \neq 0$$

Definicija. Krivulja $\gamma = (f, g)$ je KARAKTERISTIČNA, če je $\Delta = 0$ vzdolž γ in je NEKARAKTERISTIČNA, če je $\Delta \neq 0$ na γ .

Vprašanje 25. Kdaj je krivulja karakteristična za Cauchyjev problem za PDE drugega reda? Iz česa to izhaja?

Enačba $\Delta = 0$ je ekvivalentna

$$a \left(\frac{g'}{f'} \right)^2 - 2b \frac{g'}{f'} + c = 0.$$

Če pišemo $\frac{dy}{dx} = \frac{g'}{f'}$, dobimo

$$\frac{dy}{dx} = \frac{b \pm \sqrt{b^2 - ac}}{a},$$

kar je navadna diferencialna enačba (funkcije a, b, c so dane).

Za nekarakteristično krivuljo γ odvajamo originalno enačbo po x :

$$a_x u_{xx} + a u_{xxx} + 2b_x u_{xy} + 2b u_{xxy} + c_x u_{yy} + c u_{yyy} = d_x,$$

kar prepišemo v

$$a u_{xxx} + 2b u_{xxy} + c u_{xyy} = -a_x u_{xx} - 2b_x u_{xy} - c_x u_{yy} + d_x.$$

Na krivulji γ potem velja

$$\begin{aligned}\partial_s u_{xx} &= u_{xxx} f' + u_{xxy} g', \\ \partial_s u_{xy} &= u_{xxy} f' + u_{xyy} g',\end{aligned}$$

kar je nov sistem za u_{xxx} , u_{xy} in u_{yy} . Zaradi predpostavke o nekarakterističnosti ima sistem natanko eno rešitev vzdolž γ , torej te vrednosti lahko izračunamo. Podobno dobimo, če odvajamo po y , in če postopek ponavljamo; rezultat so odvodi poljubno visokega reda. Dobimo torej vrsto

$$\sum_{k=0}^{\infty} \sum_{|\alpha|=k} \partial^{\alpha} u \frac{1}{\alpha_1! \alpha_2!} (x - x_0)^{\alpha_1} (y - y_0)^{\alpha_2}.$$

Izrek (Cauchy-Kowalewski). Če je Cauchyjev problem podan z realno analitičnimi podatki in je krivulja γ nekarakteristična za problem, potem vrsta zgoraj konvergira v neki okolici točke $(x_0, y_0) \in \gamma$ in predstavlja rešitev Cauchyjevega problema.

Vprašanje 26. Kako rešuješ Cauchyjevo nalogo na karakteristični in kako na nekarakteristični krivulji? Povej Cauchy-Kowalewskijev izrek.

6.6 Valovna enačba na realni osi

Imamo podan koeficient $c > 0$. Gledamo enačbo

$$u_{tt} - c^2 u_{xx} = 0,$$

kjer je $x \in (a, b)$ in $t > 0$. To je hiperbolična enačba, ki ima kanonično obliko $w_{\xi\eta} = 0$, pri čemer sta novi koordinati oblike $\xi = x - ct$ in $\eta = x + ct$. Če integriramo $w_{\xi\eta} = 0$, dobimo $w = D(\xi) + E(\eta)$ oziroma $u(x, t) = F(x - ct) + G(x + ct)$ za poljubni $F, G \in \mathcal{C}^2$. Predpis deluje tudi, če funkciji nista odvedljivi. Če sta le odsekoma zvezni, obstajata zaporedji gladkih funkcij F_n in G_n , ki konvergirata k F in G po točkah in enakomerno na kompaktnih, kjer sta F in G zvezni. Potem bo zaporedje $u_n = F_n(x - ct) + G_n(x + ct)$ na enak način konvergiralo k u .

Definicija. Če sta F, G odsekoma zvezni na \mathbb{R} , se $u(t, x) = F(x - ct) + G(x + ct)$ imenuje POSPLOŠENA REŠITEV za valovno enačbo.

Vprašanje 27. Kakšne oblike so rešitve valovne enačbe? Kaj pa splošene rešitve?

Očitno je $F(x - ct)$ konstantna na množici $\{x - ct = \text{konst}\}$ in $G(x + ct)$ konstantna na množici $\{x + ct = \text{konst}\}$. Fiksirajmo t_0 in privzemimo, da je u dvakrat zvezno odvedljiva povsod razen v točki (x_0, t_0) . Skozi to točko potujeta dve karakteristiki $x - ct = x_0 - ct_0$ in $x + ct = x_0 + ct_0$. Za $t_1 \neq t_0$ je u gladka povsod, razen na (x_-, t_1) in (x_+, t_1) . Singularnosti u torej potujejo vzdolž karakteristik.

Poglejmo si Cauchyjev problem za valovno enačbo. V homogenem primeru za dan c rešujemo $u_{tt} - c^2 u_{xx} = 0$ za $x \in \mathbb{R}$ ter $t > 0$, pri pogojih $u(x, 0) = f(x)$ in $u_t(x, 0) = g(x)$. Zahtevamo, da je rešitev zvezna na $\mathbb{R} \times [0, \infty)$ in \mathcal{C}^2 na $\mathbb{R} \times (0, \infty)$. Splošna rešitev je oblike $u(x, t) = F(x - ct) + G(x + ct)$. Veljati mora

$$u_t(x, 0) = -cF'(x) + cG'(x) = g(x),$$

kar integriramo v

$$-F(x) + G(x) = \frac{1}{c} \int_0^x g(\xi) d\xi + d.$$

Drug pogoj je

$$u(x, 0) = F(x) + G(x) = f(x),$$

iz česar lahko izpeljemo

$$\begin{aligned} G(x) &= \frac{1}{2}f(x) + \frac{1}{2c} \int_0^x g(\xi) d\xi + d, \\ F(x) &= \frac{1}{2}f(x) - \frac{1}{2c} \int_0^x g(\xi) d\xi - d. \end{aligned}$$

Če to sestavimo, dobimo D'ALEMBERTOVO FORMULO

$$u(x, t) = \frac{1}{2}(f(x - ct) + f(x + ct)) + \frac{1}{2c} \int_{x-ct}^{x+ct} g(\xi) d\xi.$$

Vprašanje 28. Izpelji d'Alembertovo formulo za homogeno valovno enačbo.

Za vrednost $u(x, t)$ moramo poznati f v krajiščih in g na notranjosti intervala vpliva $[x - ct, x + ct]$, ki mu pravimo INTERVAL VPLIVA. Trikotnik

$$\Delta(x, t) = \{(\chi, \tau) \mid 0 \leq \tau \leq t, x - c(t - \tau) \leq \chi \leq x + c(t - \tau)\}$$

pa imenujemo KARAKTERISTIČNI TRIKOTNIK.

Za nehomogen primer rešujemo $u_{tt} - c^2 u_{xx} = \varphi(x, t)$. Če je u rešitev Cauchyjevega problema, je

$$\iint_{\Delta(x, t)} \varphi dS = \iint_{\Delta(x, y)} (u_{tt} - c^2 u_{xx}) dS.$$

Za $P = -u_t$ in $Q = -c^2 u_x$ po Greenovi formuli dobimo

$$\iint_{\Delta(x, t)} (u_{tt} - c^2 u_{xx}) dS = \int_{\partial\Delta(x, t)} (-u_t, -c^2 u_x) \cdot d\vec{r}.$$

Integral izračunamo na vsaki stranici trikotnika posebej. Na koncu dobimo

$$u(x, t) = \frac{1}{2c} \iint_{\Delta(x, y)} \varphi dS + \frac{1}{2}(f(x + ct) + f(x - ct)) + \frac{1}{2c} \int_{x-ct}^{x+ct} g(\xi) d\xi.$$

Vprašanje 29. Izpelji d'Alembertovo formulo za nehomogeno valovno enačbo.

Izrek. Naj bodo $f \in \mathcal{C}^2(\mathbb{R})$, $g \in \mathcal{C}(\mathbb{R})$, $\varphi \in \mathcal{C}(\mathbb{R}^2)$ in $T > 0$. Cauchyjeva naloga $u_{tt} - c^2 u_{xx} = \varphi$, $u(x, 0) = f(x)$, $u_t(x, 0) = g(x)$ je za vsak $x \in \mathbb{R}$ ter $t \in [0, T]$ dobro postavljen, torej

- rešitev obstaja,
- rešitev je enolična,
- rešitev je zvezno odvisna od začetnih pogojev.

Dokaz. Obstoj in enoličnost sta dokazana zgoraj (začeli smo z rešitvijo, in izpeljali formulo). Zveznost preverjamo za supremum normo, kjer za f in g dovolimo neskončne vrednosti v normi. Dokazujemo, da za vsak $\varepsilon > 0$ obstaja $\delta > 0$, da pogoji $\|\varphi_1 - \varphi_2\|_\infty < \delta$, $\|f_1 - f_2\|_\infty < \delta$ in $\|g_1 - g_2\|_\infty < \delta$ implicirajo $\|u_1 - u_2\|_\infty < \varepsilon$.

Za razliko $u_1 - u_2$ ocenimo vsak kos posebej:

$$\begin{aligned} \left| \frac{f_1(x+ct) - f_2(x+ct) + f_1(x-ct) - f_2(x-ct)}{2} \right| &< \delta, \\ \left| \frac{1}{2c} \int_{x-ct}^{x+ct} (g_1(\xi) - g_2(\xi)) d\xi \right| &< \frac{1}{2c} 2ct\delta, \\ \left| \frac{1}{2c} \iint_{\Delta(x,t)} (\varphi_1 - \varphi_2) dS \right| &< \frac{1}{2c} S(\Delta(x,t))\delta, \end{aligned}$$

torej bo primerna izbira $\delta = \frac{1}{2} \frac{\varepsilon}{1+T+\frac{1}{2}T^2}$. □

Vprašanje 30. Dokaži, da je rešitev valovne enačbe zvezno odvisna od začetnih pogojev.

6.7 Toplotna enačba

Naj bo $U \subseteq \mathbb{R}^n$ omejena odprta množica z gladkim robom. Obravnavamo Cauchyjev problem za toplotno enačbo

$$u_t - \Delta u = 0$$

na $U \times (0, \infty)$. Poleg tega zahtevamo $u = 0$ na $\partial U \times (0, \infty)$ in $u(x, 0) = g(x)$ na $U \times \{0\}$. Iskali bomo rešitve oblike $u(x, t) = v(t)w(x)$. Izračunamo $u_t - \Delta u = w(x)v'(t) - v(t)\Delta w(x) = 0$, iz česar dobimo

$$\frac{\Delta w(x)}{w(x)} = \frac{v'(t)}{v(t)}.$$

Leva stran te enačbe je neodvisna od t , desna pa neodvisna od x , torej morata biti konstantni. Označimo to konstanto z $-\lambda$. Iskali bomo funkcije w , ki rešijo Cauchyjevo nalogo $-\Delta w = \lambda w$ in $w = 0$ na ∂U . Potem bo funkcija $u(x, t) = de^{-\lambda t}w(x)$ rešila enačbo $u_t - \Delta u = 0$ na $U \times (0, \infty)$ s pogojem $u = 0$ na ∂U . Pri $t = 0$ velja $u(x, 0) = dw(x) = g(x)$, torej je rešitev oblike $u(x, t) = g(x)e^{-\lambda t}$, a le v primeru, da je g oblike $dw(x)$, kjer je

$d \in \mathbb{R}$ in w reši omenjen problem. Ker je enačba linearna, je vsota rešitev take oblike tudi rešitev. Poiskali bomo števno družino rešitev w_j, λ_j, d_j , da bo

$$u(x, t) = \sum_{j=1}^{\infty} d_j e^{-\lambda_j t} w_j(x).$$

Vprašanje 31. Izpelji splošno rešitev za toplotno enačbo.

Obravnavajmo primer $n = 1$. Če je $\lambda \leq 0$, iz robnih pogojev izpeljemo, da mora biti w trivialen, kar za nas ni zanimivo. V primeru $\lambda > 0$ dobimo

$$w(x) = A \cos(\sqrt{\lambda}x) + B \sin(\sqrt{\lambda}x).$$

Iz robnih pogojev dobimo

$$\begin{bmatrix} \cos(\sqrt{\lambda}a) & \sin(\sqrt{\lambda}a) \\ \cos(\sqrt{\lambda}b) & \sin(\sqrt{\lambda}b) \end{bmatrix} \cdot \begin{bmatrix} A \\ B \end{bmatrix} = 0.$$

Dobimo dodatno zahtevno, da je ta matrika nesingularna. Sistem se prevede na $\sin(\sqrt{\lambda}(b-a)) = 0$, torej mora veljati

$$\lambda = \lambda_k = \left(\frac{k\pi}{b-a} \right)^2.$$

Upoštevamo $a = 0$ in $b = L$, da dobimo $\lambda_k = (k\pi/L)^2$. Iz $w(a) = w(b) = 0$ potem sledi $A = 0$, in so rešitve oblike

$$w_k(x) = B_k \sin\left(\frac{k\pi x}{L}\right).$$

Nastavek sedaj deluje za

$$g = \sum_{k=1}^{\infty} d_k w_k,$$

in iščemo koeficiente B_k , da bo ta enakost veljala.

Vprašanje 32. Izpelji rešitev Cauchyjevega problema za enodimenzionalno toplotno enačbo.

Oglejmo si enodimenzionalni nehomogen problem

$$u_t = k u_{xx} + F(x, t)$$

z robnimi pogoji $u(0, t) = a(t)$, $u(\pi, t) = b(t)$ in začetnimi pogoji $u(x, 0) = f(x)$. Rešitev sestavimo iz rešitev bolj enostavnih problemov:

$u_t - k u_{xx} = F(x, t)$	$u_t = k u_{xx}$	$u_t = k u_{xx}$
$u(0, t) = 0$	$u(0, t) = a(t)$	$u(0, t) = 0$
$u(\pi, t) = 0$	$u(\pi, t) = b(t)$	$u(\pi, t) = 0$
$u(x, 0) = 0$	$u(x, 0) = 0$	$u(x, 0) = f(x)$

Rešitve tretjega problema iščemo s separacijo spremenljivk. Iz tega dobimo nastavek

$$u(x, t) = \sum_{n=1}^{\infty} C_n e^{-kn^2 t} \sin(nx),$$

koeficiente pa dobimo iz Fourierovega razvoja f .

Za prvi problem razvijemo

$$F(x, t) = \sum_{n=1}^{\infty} F_n(t) \sin(nx)$$

in

$$u(x, t) = \sum_{n=1}^{\infty} T_n(t) \sin(nx).$$

Iz $u(x, 0) = 0$ sledi $T_n(0) = 0$ za vsak n . To vstavimo v enačbo $u_t = ku_{xx} + F(x, t)$ in dobimo

$$T'_n + kn^2 T_n = F_n,$$

kar je nehomogena NDE prvega reda z rešitvijo

$$T_n(t) = T_{n,p}(t) + B_n e^{-kn^2 t},$$

ki mora ustrezati pogoju $T_{n,p}(0) = -B_n$.

Za drugi problem enak pristop ne deluje zaradi robnih pogojev. Enačbo homogeniziramo z

$$w(x, t) = \frac{b-a}{\pi} x + a.$$

Velja $w(0, t) = w(\pi, t) = 0$. Sedaj razcepimo $u = v + w$ ter iščemo pogoje za v ;

$$v_t + w_t = k(v_{xx} + w_{xx}).$$

Iz tega dobimo

$$v_t = kv_{xx} - \frac{b'-a'}{\pi} x - a'.$$

To je problem prvega tipa za $\tilde{F}(x, t) = -\frac{b'-a'}{\pi} x - a'$.

Vprašanje 33. Kako rešiš nehomogeno toplotno enačbo v enodimenzionalnem primeru?

Rešitve take enačbe so enolične. Naj bosta u_1 in u_2 rešitvi; potem $u = u_1 - u_2$ reši enačbo $u_t = ku_{xx}$ s pogoji $u(0, t) = u(\pi, t) = u(x, 0) = 0$. Definiramo energijo

$$E(t) = \int_0^\pi u^2(x, t) dx.$$

Velja $E(0) = 0$ in $E(t) \geq 0$. Če predpis odvajamo, dobimo

$$E'(t) = 2k \int_0^\pi uu_{xx} dx = 2kuu_{xx}|_0^\pi - 2k \int_0^\pi u_x^2 dx = -2k \int_0^\pi u_x^2 dx \leq 0$$

za $k \geq 0$. Torej je E konstanta in $u = 0$.

Vprašanje 34. Dokaži, da je rešitev nehomogene toplotne enačbe v eni dimenziji enolična.

Izrek (princip maksima za toplotno enačbo). *Naj bo D omejena odprta množica v \mathbb{R}^n , $0 < T < \infty$ in $u \in \mathcal{C}(\overline{D} \times [0, \infty)) \cap \mathcal{C}^\infty(D \times (0, \infty))$ rešitev enačbe $u_t = k\Delta u$. Tedaj u zavzame maksimum na $D \times \{0\} \cup \partial D \times [0, T]$.*

Dokaz. Naj bo $\varepsilon > 0$. Definiramo $v = u + \varepsilon|x|^2$. Potem je $v_t = u_t$ in $\Delta v = \Delta u + 2n\varepsilon$. Velja $v_t - \Delta v = u_t - \Delta u - 2n\varepsilon < 0$. Če ima v maksimum v $D \times (0, T)$, potem v tej točki velja $v_t = 0$ in $\Delta v \leq 0$ (to je ravno sled Hessejeve matrike), torej $v_t - \Delta v \geq 0$, kar je protislovje. Podobno ugotovimo, da v ne mora imeti maksimuma v $D \times \{T\}$. Za dovolj majhen ε je v le majhna sprememba u , torej tudi u nima maksimuma na teh množicah. \square

Vprašanje 35. Dokaži princip maksima za toplotno enačbo.

6.8 Sturm-Liouvilleova teorija

Od tu naprej SL teorija. Naj bo $I = [a, b]$ in $p, q, r : I \rightarrow \mathbb{R}$ funkcije, ki zadoščajo

- $p \in \mathcal{C}^1$
- $q, r \in \mathcal{C}$
- p, r nenegativni

Na $\mathcal{C}^2(I)$ definiramo pridruženi STURM-LIOUVILLEOV operator

$$Lu = (pu')' + qu.$$

Naj bodo $\alpha_0, \alpha_1, \beta_0, \beta_1 \in \mathbb{R}$ take, da je $|\alpha_0| + |\alpha_1| > 0$ in $|\beta_0| + |\beta_1| > 0$ (torej paroma ne obe 0). Označimo $B_a(u) = \alpha_0 u(a) + \alpha_1 u'(a)$ in $B_b(u) = \beta_0 u(b) + \beta_1 u'(b)$. SL problem sprašuje po netrivialnih rešitvah enačbe

$$Lu = -\lambda ru \tag{6.4}$$

ob pogojih $B_a(u) = 0$ in $B_b(u) = 0$. Funkcijo r imenujemo UTEŽ ali GOSTOTA. Problem je REGULAREN, če sta $p, r > 0$, in SINGULAREN, če sta p ali q enaka 0 v točki a ali b , ali če imata v kakšni od teh točk skok, ali pa če je I neomejen. Če par (λ, u) reši SL problem, je λ LASTNA VREDNOST, u pa LASTNA FUNKCIJA.

Vprašanje 36. Definiraj Sturm-Liouvilleov problem. Kdaj je regularen in kdaj singularen?

Trditev. Vsako homogeno regularno linearno NDE drugega reda lahko zapišemo v obliki SL problema.

Dokaz. Imejmo enačbo $Av'' + Bv' + Cv = 0$. Velja $Lv = pv'' + p'v' + (\dots)$, kjer smo člene ničtega reda skrili. Veljati bi moralo $B = A'$. Če to ni res, lahko zaradi homogenosti enačbo množimo z neničelno funkcijo m in dobimo $m(Av'' + Bv') + mCv = 0$, kjer m določimo tako, da bo veljalo $(mA)' = mB$; izračun pokaže

$$m = \frac{1}{A} \exp\left(\int \frac{B}{A}\right).$$

Na koncu dobimo

$$mAv'' + mBv' + mCv = (pv')' + (mCv) = 0$$

za $p = \exp(B/A)$ in $q + \lambda r = \frac{C}{A}p$. □

Vprašanje 37. Pokaži, da se vsako enačbo drugega reda da zapisati kot SL problem.

Linearen operator na Banachovem prostoru je zvezen natanko tedaj, ko je omejen. Ker diferencialni operatorji običajno niso omejeni, je to sila neprijetno, zato se omejimo na podprostore. SL problem gledamo na definicijskem območju L , ki bo v prihodnje bodisi

$$\mathcal{U} = \{u \in \mathcal{C}^2(I) \mid B_a(u) = B_b(u) = 0\}$$

bodisi

$$\mathcal{V} = \{u \in \mathcal{C}^2(I) \mid u(a) = u(b), u'(a) = u'(b)\},$$

odvisno če gledamo običajen ali periodičen problem (druga množica so ravno pogoji periodičnega problema). Operatorja $L_{\mathcal{U}}$ in $L_{\mathcal{V}}$ (torej zožitvi L na ta prostora) sta definirana na gosti podmnožici $l^2(I)$. Za zvezno funkcijo $r : I \rightarrow \mathbb{R}$ označimo

$$l^2(r) = \{f : I \rightarrow \mathbb{C} \mid f \text{ merljiva}, \|f\|_{l^2(r)}^2 = \int_I |f|^2 r < \infty\}.$$

Ker je r zvezna na kompaktu, obstajata c_1 in c_2 , da velja $c_1 \leq r(x) \leq c_2$, kar označimo z $r \sim 1$.

Opomba. Oznako lahko razširimo: $r \sim g$ natanko tedaj, ko obstajata konstanti c_1 in c_2 , da velja $c_1 g(x) \leq r(x) \leq c_2 r(x)$.

Vidimo, da je $\|f\|_{l^2(r)} \sim \|f\|_{l^2}$.

Trditev. Naj bo $Lu = (pu')' + qu$ standardni SL operator. Če je $v \in \mathcal{C}^2(I)$, velja $uLv - vLu = [p \cdot (uv' - u'v)]'$ (LAGRANGEOVA IDENTITETA) in

$$uLv - vLu = p \cdot (uv' - u'v) \Big|_{\partial I}.$$

Dokaz je račun.

Posledica. Operator $L_{\mathcal{U}}$ je simetričen. Enako velja za \mathcal{V} , če je $p(a) = p(b)$.

Dokaz. Naj bosta u in v realni funkciji. Velja

$$\langle Lu, v \rangle - \langle u, Lv \rangle = \int_I (vLu - uLv) = p(u'v - uv')|_{\partial I}.$$

Če je $p(a) = p(b)$ in $L = L_{\mathcal{V}}$, je to enako 0. V drugem primeru prepišemo robni pogoj v

$$\begin{bmatrix} u(a) & u'(a) \\ v(a) & v'(a) \end{bmatrix} \cdot \begin{bmatrix} \alpha_0 \\ \alpha_1 \end{bmatrix} = 0.$$

Privzeli smo, da nista oba α_0, α_1 enaka 0, torej je matrika singularna. Enako velja za b , iz česar sledi $u'v - uv'|_a = 0$ in $u'v - uv'|_b = 0$.

Če sta u in v kompleksni funkciji, torej $u = u_1 + iu_2$ in $v = v_1 + iv_2$, pa velja

$$\langle Lu, v \rangle = \langle Lu_1, v_1 \rangle + \langle Lu_2, v_2 \rangle - i(\langle Lu_1, v_2 \rangle + \langle Lu_2, v_1 \rangle).$$

Sedaj le upoštevamo realni primer, dobimo $\langle Lu, v \rangle = \langle u, Lv \rangle$. □

Vprašanje 38. Dokaži: če je \mathcal{U} množica $\mathcal{C}^2(I)$ funkcij, ki ustrezajo robnemu pogoju SL problema, je zožitev operatorja $L_{\mathcal{U}}$ simetrična.

Trditev. Vse lastne vrednosti operatorjev L glede na utež r so realne.

Dokaz. Če je u lastna funkcija, velja

$$0 = \langle Lu, u \rangle - \langle u, Lu \rangle = \langle -\lambda ru, u \rangle - \langle u, -\lambda ru \rangle = (\bar{\lambda} - \lambda) \langle u, ru \rangle.$$

□

Vprašanje 39. Dokaži: vse lastne vrednosti SL problema so realne.

Lema. Če so funkcije f_1, \dots, f_n linearno odvisne, je determinanta Wronskega identično enaka 0. Obratno, če je Λ diferencialen operator z realnimi koeficienti in y_1, \dots, y_n rešijo enačbo $\Lambda y = 0$, ter je $W(y_1, \dots, y_n) = 0$, tedaj so linearno odvisne.

Trditev. Naj bo λ lastna vrednost za SL problem $L_{\mathcal{U}} = -\lambda ru$. Tedaj je pripadajoči lastni podprostor enodimenzionalen.

Dokaz. Naj za $u, v \in \mathcal{U}$ velja $Lu = -\lambda ru$ in $Lv = -\lambda rv$. Sledi $uLv - vLu = 0$, a hkrati $uLv - vLu = [p(u'v - uv')]'$, torej je $f = p(u'v - uv')$ konstantna. Ker sta $u, v \in \mathcal{U}$, je $f(a) = f(b) = 0$ in zato $f = 0$. Ker pa je $p \neq 0$, sledi še $u'v - vu' = 0$, kar je ravno determinanta Wronskega. Iz leme sledi, da sta u in v linearno odvisni, saj obe ležita v jedru diferencialnega operatorja $\Lambda y = Ly + \lambda ry$. □

Vprašanje 40. Dokaži: lastni podprostor za SL problem je enodimenzionalen.

Izrek (spektralni izrek za SL problem). *Za vsak regularen SL problem obstaja zaporedje lastnih vrednosti $\lambda_1 < \lambda_2 < \dots < \infty$, da velja*

- $\lim \lambda_n = \infty$
- *Pripadajoč normaliziran sistem lastnih funkcij $\{v_n\}_n$ sestavlja kompleten ortonormiran sistem v $l^2(r)$.*

Enako velja za periodičen SL problem, pri čemer imamo lahko dvodimenzionalne lastne podprostore.

Vprašanje 41. Povej spektralni izrek za SL problem.

6.9 Harmonične funkcije

Trditev (veržno pravilo). *Za $f, g \in \mathcal{C}^1$ velja*

$$\begin{aligned}\partial_{\bar{z}}g \circ f &= \partial_w g(f) \cdot \partial_{\bar{z}}f + \partial_{\bar{w}}g(f) \cdot \partial_{\bar{z}}\bar{f}, \\ \partial_zg \circ f &= \partial_w g(f) \cdot \partial_zg + \partial_{\bar{w}}g(f) \cdot \partial_z\bar{f}.\end{aligned}$$

Definicija. Naj bo $\Omega^{\text{odp}} \subseteq \mathbb{R}^2$ in $u \in \mathcal{C}^2(\Omega)$. Funkcija u je **HARMONIČNA** na Ω , če je $\Delta u = 0$. Funkcija je **SUBHARMONIČNA**, če je $\Delta u \geq 0$.

Definicija. Funkcija u zadošča **LASTNOSTI POVPREČNE VREDNOSTI (LPV)**, če za vsak $z_0 \in \Omega$ in $r > 0$, pri katerih je $\bar{K}(z_0, r) \subseteq \Omega$, velja

$$u(z_0) = \int_0^1 u(z_0 + re^{2\pi it}) dt.$$

Izrek. *Funkcija $u \in \mathcal{C}^2(\Omega)$ je harmonična na Ω natanko tedaj, ko tam zadošča LPV.*

Dokaz. Naj bo $z_0 \in \Omega$ ter $r > 0$ dovolj majhen. Definirajmo $\varphi : [0, r] \rightarrow \mathbb{C}$ s predpisom

$$\varphi(s) = \int_0^1 u(z_0 + se^{2\pi it}) dt.$$

Funkcija φ je zvezna na $[0, r]$, diferenciable na $(0, r)$ in velja $\varphi(0) = u(z_0)$. Izračunamo lahko

$$\varphi'(s) = \int_0^1 \partial_s u(x_0 + s \cos 2\pi t, y_0 + s \sin 2\pi t) dt = \int_0^1 \left\langle \vec{\nabla} \cdot u(\vec{r}(t)), (\cos 2\pi t, \sin 2\pi t) \right\rangle dt$$

za $\vec{r}(t) = z_0 + se^{2\pi it}$. Drugi vektor v skalarnem produktu je enak normalni $\vec{n}(\vec{r}(t))$, torej integriramo smerni odvod

$$\varphi'(s) = \int_0^1 \frac{\partial u}{\partial \vec{n}}(\vec{r}(t)) dt = \frac{1}{2\pi s} \iint_{K(z_0, s)} \Delta u dS$$

po Greenovi formuli, kjer konstanta pred integralom pride iz spremembe spremenljivke. Obe smeri implikacije enostavno sledita. \square

Vprašanje 42. Dokaži: funkcija je harmonična natanko tedaj, ko zadošča lastnosti povprečne vrednosti.

Trditev. Naj bosta $u, v \in \mathcal{C}^2(\Omega)$ realni. Če je $u + iv$ holomorfná, sta u in v harmonični.

Trditev. Naj bo $\Omega \subseteq \mathbb{R}^2$ enostavno povezana odprta množica. Če je u realna harmonična funkcija, obstaja holomorfná funkcija f na Ω , da je $u = \operatorname{Re} f$.

Posledica. Če je u harmonična in F holomorfná, je $u \circ F$ harmonična.

Dokaz. Vemo, da je lokalno $u = \operatorname{Re} f$. Sledi $u \circ F = \operatorname{Re} f \circ F = \operatorname{Re}(f \circ F)$. \square

Vprašanje 43. Dokaži: če je u harmonična in F holomorfná, je $u \circ F$ harmonična.

6.9.1 Dirichletov problem

Označimo enotski disk z D in njegov rob s $\mathbb{T} = \partial D$. Funkcije na \mathbb{T} lahko identificiramo z 1-periodičnimi funkcijami na \mathbb{R} na očiten način. Dirichletov problem je naslednji: Če imamo dano funkcijo $f \in \mathcal{C}(\mathbb{T})$, ali lahko najdemo $u \in \mathcal{C}(\overline{D})$, da bo u harmonična na D in da se bo ujemala z f na \mathbb{T} ?

Primer. Funkcijo z^{-k} razširimo iz ∂D na D kot \bar{z}^k . Realni del je enak realnemu delu z^k , torej je funkcija res harmonična.

Vprašanje 44. Kaj je Dirichletov problem na enotskem disku? Kako ga rešiš za z^{-k} ?

Označimo $e_m(t) = \exp(2\pi i m t)$ za $m \in \mathbb{Z}$ in $t \in [0, 1]$. Pripadajoča razširitev funkcije e_m na D je funkcija $E_m : D \rightarrow \mathbb{C}$, definirana kot

$$E_m(re^{2\pi i t}) = \begin{cases} r^m e^{2\pi i m t} & m \in \mathbb{N} \\ r^{-m} e^{2\pi i m t} & m \in -\mathbb{N} \end{cases}$$

Za $f = \sum_k c_k e_k$ dobimo razširitev $\sum_k c_k E_k$. Ta prehod lahko zapišemo na eleganten način. Funkcije e_k tvorijo ortonormiran sistem: če je $f = \sum_k c_k e_k$, je

$$c_k = \int_0^1 f(\tau) \overline{e_k(\tau)} d\tau.$$

Domnevamo, da je

$$F(re^{2\pi i t}) = \sum_k \int_0^1 f(\tau) \overline{e_k(\tau)} r^{|k|} e^{2\pi i k t} d\tau = \int_0^1 f(\tau) \sum_k r^{|k|} e^{2\pi i k(t-\tau)} d\tau.$$

Če vsoto v zadnjem integralu označimo z $G_r(t - \tau)$, je to enako $(f * G_r)(t)$. Z nekaj algebrske manipulacije lahko ugotovimo, da je

$$G_r(s) = P_r(s) := \frac{1 - r^2}{1 - 2r \cos(2\pi s) + r^2}.$$

Tej funkciji pravimo POISSONOVO JEDRO.

Vprašanje 45. Izpelji predpis za rešitev Dirichletovega problema na enotskem disku. Kaj je Poissonovo jedro?

7 Izbrane teme iz analize podatkov

Pri strojnem učenju obravnavamo MNOŽICO OPAZOVANIH OBJEKTOV O , in opazujemo lastnosti teh objektov. LASTNOST modeliramo kot preslikavo $V : O \rightarrow D_V$, kjer je zaloga vrednosti D_V lahko bodisi poljubna končna množica (v primeru diskretne spremenljivke), ali pa podmnožica \mathbb{R} , v primeru numerične spremenljivke. Za diskretne spremenljivke dodatno zahtevamo, da so vrednosti D_V neurejene. Lastnosti ustrezajo SPREMENLJIVKAM, ki jih ločimo na dva tipa. Prve so NAPOVEDNE SPREMENLJIVKE $\mathbf{X} = (X_1, \dots, X_p)$, ki predstavljajo podatke, druge pa so CILJNE SPREMENLJIVKE, ki jih želimo napovedati. Vsako ciljno spremenljivko lahko obravnavamo posebej, torej si mislimo, da imamo le eno, Y . Pri nenadzorovanem učenju ciljnih spremenljivk nimamo, tam nas namesto napovedi zanimajo drugačna vprašanja.

Vprašanje 1. Opiši podatkovno množico pri strojnem učenju.

Napovedni model predstavimo s funkcijo $m : D_{X_1} \times D_{X_2} \times \dots \times D_{X_p} \rightarrow D_Y$. Funkcija za podane vrednosti \mathbf{x} izračuna napoved ciljne spremenljivke $m(\mathbf{x}) = \hat{y}$. Napovedne modele ločimo glede na to, kakšno spremenljivko napovedujejo. REGRESIJSKI MODEL napoveduje numerično spremenljivko, KLASIFIKACIJSKI pa diskretno.

Za mero, kateri napovedni model je boljši od drugega, definiramo FUNKCIJO IZGUBE $L : D_Y \times D_y \rightarrow [0, \infty)$, ki izračuna napako ene napovedi \hat{y} glede na točno vrednost y . Pri regresijskih modelih običajno uporabljamo kvadratno napako

$$L_{SE}(y, \hat{y}) = (y - \hat{y})^2,$$

pri klasifikacijskih pa

$$L_{01}(y, \hat{y}) = \begin{cases} 0 & y = \hat{y}, \\ 1 & \text{sicer.} \end{cases}$$

Napovedna napaka modela je potem povprečna napaka na vseh primerih iz S ,

$$\text{Error}(m, S) = \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} L(y, m(\mathbf{x})).$$

Pri nadzorovanem učenju razdelimo množico S na učno in testno množico, S_{train} in S_{test} , ki sta disjunktni in pokrijeta vse primere. Pravimo, da je model TOČEN, če ima majhno napako na učni množici, in SPLOŠEN, če ima majhno napako na testni množici.

Vprašanje 2. Kako definiramo napako napovednega modela? Kaj sta točnost in splošnost?

Optimalen in nepristranski napovedni model bo takšen, ki bo minimiziral kvadratno napako $E((Y - m(\mathbf{X}))^2)$. V idealnem svetu bo tak model imel obliko $m^*(\mathbf{x}) = E(Y | \mathbf{X} = \mathbf{x})$. Če predpostavimo, da so vsa opažanja v S enako verjetna, je približek tega modela funkcija

$$m^*(\mathbf{x}_0) = \frac{1}{|S_0|} \sum_{(\mathbf{x}, y) \in S_0} y,$$

kjer je $S_0 = \{(\mathbf{x}, y) \in S_{\text{train}} \mid \mathbf{x} = \mathbf{x}_0\}$ množica primerov, kjer je $\mathbf{x} = \mathbf{x}_0$. Problem je v tem, da je v praksi ta množica najverjetneje prazna, ali pa vsebuje zelo malo primerov. Namesto tega lahko za S_0 vzamemo k najbližjih primerov iz S_{train} , čemur pravimo SOSEŠČINA točke \mathbf{x}_0 .

Vprašanje 3. Opiši delovanje metode najbližjih sosedov.

7.1 Linearna regresija

Pri numeričnih napovednih spremenljivkah in numerični ciljni spremenljivki lahko uporabimo linearno regresijo,

$$\hat{y} = \beta_0 + \sum_{i=1}^p \beta_i X_i,$$

kjer so β_i neznani koeficienti. Če jih zložimo v vektor $\boldsymbol{\beta} = [\beta_0, \dots, \beta_p]^T$, lahko model prepišemo v $\hat{y} = \mathbf{X}\boldsymbol{\beta}$ za $\mathbf{X} = [1, X_1, \dots, X_p]$. Optimalna izbira za $\boldsymbol{\beta}$ bo tedaj

$$\arg \min_{\boldsymbol{\beta}} (Y - \mathbf{X}\boldsymbol{\beta})(Y - \mathbf{X}\boldsymbol{\beta})^T$$

za stolpec ciljnih vrednosti Y . Funkciji, ki jo zgoraj minimiziramo, pravimo tudi RSS. Optimalen $\boldsymbol{\beta}$ dobimo z reševanjem predločenega sistema.

Vprašanje 4. Pojasni linearno regresijo.

Napako linearnega modela lahko merimo na več načinov. Ena možnost je RESIDUALNA STANDARDNA NAPAKA

$$\text{RSE} = \sqrt{\frac{\text{RSS}}{|S| - p - 1}},$$

kjer je podobno kot prej

$$\text{RSS} = \sum_{(\mathbf{x}, y) \in S} (y - \mathbf{x}^T \boldsymbol{\beta})^2.$$

Dober model imamo, ko je RSE blizu 0.

Drug način merjenja napake je DELEŽ POJASNJENE VARIANCE

$$R^2 = 1 - \frac{\text{RSS}}{\text{TSS}}$$

za

$$\text{TSS} = \sum_{(\mathbf{x}, y) \in S} (y - \bar{y})^2,$$

kjer z \bar{y} označimo povprečje. Mera R^2 ima vrednosti v intervalu $[0, 1]$, dober model pa dobimo za R^2 blizu 1.

Nazadnje imamo še POPOLNO NAPAKO MODELA oziroma RMSE

$$\text{RMSE} = \sqrt{\frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} (y - \mathbf{x}^T \boldsymbol{\beta})^2}.$$

Tudi tu ima dober model vrednost blizu 0.

Vprašanje 5. Kako meriš napako linearnega modela?

Če vemo (ali predpostavljamo), da ima neka spremenljivka nelinearen vpliv, lahko poskusimo dodati nove spremenljivke, npr. $X_1^2, X_1 X_2$, ipd.

7.2 Logistična regresija

Če imamo diskretno spremenljivko V z domeno $\{v_1, \dots, v_k\}$, jo spremenimo v k numeričnih spremenljivk, ki so indikatorji dogodka $V = v_i$. Dovolj je torej znati napovedati vrednost binarne diskretne spremenljivke. V nadaljevanju napovedujemo Y z vrednostmi $D_Y = \{\oplus, \ominus\}$. Namesto vrednosti \oplus in \ominus lahko poskusimo napovedati verjetnost $P(Y = \oplus | X = x)$, problem pa je, da nam linearna regresija hitro zbeži izven intervala $[0, 1]$. Rešitev je, da napovedujemo logaritem obetov $\log \frac{p}{1-p}$, ki leži na intervalu $(-\infty, \infty)$. Odločitvena meja pri $p = 0.5$ sovпада z vrednostjo $z = 0$.

Težava je v tem, da verjetnosti $P(Y = \oplus | X = x)$ ne poznamo, torej ne moremo formulirati najmanjših kvadratov. Deluje pa metoda največjega verjetja

$$L(\mathbf{X}, Y, \boldsymbol{\beta}) = \prod_{(\mathbf{x}, y) \in S} P(Y = y | \mathbf{X} = \mathbf{x}, \boldsymbol{\beta}) = \prod_{y=1} \frac{e^{\mathbf{x}^T \boldsymbol{\beta}}}{1 + e^{\mathbf{x}^T \boldsymbol{\beta}}} \prod_{y=0} \frac{1}{1 + e^{\mathbf{x}^T \boldsymbol{\beta}}}.$$

Najlažje maksimiziramo logaritem verjetja, torej

$$\log L = \sum_{y=1} (\mathbf{x}^T \boldsymbol{\beta} - \log(1 + e^{\mathbf{x}^T \boldsymbol{\beta}})) - \sum_{y=0} \log(1 + e^{\mathbf{x}^T \boldsymbol{\beta}}).$$

Če prvo vsoto pomnožimo z y , drugo pa z $1 - y$, nič ne spremenimo, torej je to enako

$$\log L = \sum_{(\mathbf{x}, y) \in S} (\mathbf{x}^T \boldsymbol{\beta} y - \log(1 + e^{\mathbf{x}^T \boldsymbol{\beta}})).$$

Maksimum funkcije poiščemo z odvodom.

Vprašanje 6. Razloži logistično regresijo.

Pri logistični regresiji lahko merimo KLASIFIKACIJSKO NAPAKO

$$\text{CE} = \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} \mathbb{1}(y \neq \mathbf{x}^T \boldsymbol{\beta}),$$

kjer imamo dober model, če je ta napaka blizu 0.

Vprašanje 7. Kaj je klasifikacijska napaka?

7.3 Najbližji sosedi

V metodi najbližjih sosedov definiramo SOSEŠČINO točke \mathbf{x}_0 kot množico S_0 najbližjih k meritev tej točki. Razdalja je običajno psevdometrika, kjer ne zahtevamo, da iz $d(a, b) = 0$ sledi $a = b$. Pri regresijskem modelu vzamemo (uteženo) povprečje množice S_0 , v klasifikaciji pa večinsko glasovanje, tj. najpogostejšo vrednost v soseščini.

Vprašanje 8. Razloži metodo najbližjih sosedov.

7.4 Vrednotenje napovednih modelov

Recimo, da obstaja idealen napovedni model $Y = m(\mathbf{X}) + \varepsilon$, kjer je za napako $E(\varepsilon) = 0$. Zanima nas pričakovana napaka napovednega modela \hat{m} , naučenega na neki množici S , pri $\mathbf{X} = \mathbf{x}_0$. Model m je determinističen, torej velja

$$E(Y | \mathbf{X} = \mathbf{x}_0) = E(m(\mathbf{x}_0) + \varepsilon) = E(m(\mathbf{x}_0)) = m(\mathbf{x}_0)$$

in

$$\text{var}(Y | \mathbf{X} = \mathbf{x}_0) = E(\varepsilon^2) = \sigma_\varepsilon^2.$$

S podobno enostavnim računom pridemo do

$$E(\text{Err}(\mathbf{x}_0)) = E((Y - \hat{m}(\mathbf{x}_0))^2) = \sigma_\varepsilon^2 + (E(\hat{m}(\mathbf{x}_0)) - m(\mathbf{x}_0))^2 + \text{var}(\hat{m}(\mathbf{x}_0)).$$

To nas vodi do definicije dveh novih količin: PRISTRANSKOST

$$E(\hat{m}(\mathbf{x}_0)) - m(\mathbf{x}_0)$$

in VARIANCA

$$\text{var}(\hat{m}(\mathbf{x}_0)).$$

Vprašanje 9. Izpelj predpis za pričakovano napako napovednega modela. Kaj sta pristranskost in varianca?

Pri učenju modelov ločimo UČNO in TESTNO napako. Učna napaka je oblike

$$\text{Err}_{\text{train}}(m, S) = E(L(y, m(\mathbf{x})) | (\mathbf{x}, y) \in S) = \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} L(y, m(\mathbf{x})),$$

kjer je L funkcija izgube. Ta napaka je seveda optimistična, zato merimo tudi testno napako

$$\text{Err}_{\text{test}}(m, S) = E(L(y, m(\mathbf{x})) | (\mathbf{x}, y) \notin S).$$

Razliki med tema napakama pravimo OPTIMIZE

$$o = \text{Err}_{\text{test}} - \text{Err}_{\text{train}}.$$

Izkaže se, da v splošnem velja

$$E(o) = \frac{2}{|S|} \sum_{(\mathbf{x}, y) \in S} \text{cov}(m(\mathbf{x}), y),$$

torej je optimizem obratno sorazmeren številu učnih primerov $|S|$. V skrajnem primeru za $|S| \rightarrow \infty$ sploh ne potrebujemo testnih primerov. Pri linearnih modelih se izkaže $\text{cov}(m(\mathbf{x}), y) = p\sigma_\varepsilon^2$, torej je optimizem premo sorazmeren s številom napovednih spremenljivk.

Vprašanje 10. Kaj je optimizem? Kako se izraža v primeru linearnega modela?

Poznamo več načinov za izbiro učnih in testnih podatkov. Najbolj enostaven je naključno vzorčenje, kjer vzamemo nekaj naključnih primerov za učne, ostale pa pustimo za testne. Problem s tem je, da je ocena napake občutljiva na izbiro primerov.

Boljši način je prečno preverjanje, kjer razdelimo S na k enako velikih množic, ki imajo približno enako porazdelitev ciljne spremenljivke. Potem naredimo k modelov, vsakič vzamemo en kos za testno in ostale za učno množico. Končna napaka je potem povprečje napak na posamičnih testnih.

Vprašanje 11. Razloži prečno preverjanje.

Druga možnost je zankanje, kjer vzamemo B naključnih vzorcev množice S s ponavljanjem. Te množice so vse enako velike kot S , imajo vlogo učnih, njihov komplement pa vlogo testnih množic. Verjetnost, da nek primer ni v učni množici, je potem $(1 - \frac{1}{|S|})^{|S|} \approx e^{-1}$. Napako izmerimo za vsak model (na komplementu njegovih učnih podatkov), končna ocena napake je povprečje teh.

Vprašanje 12. Razloži zankanje.

Pri dvojiški klasifikaciji imamo dva pristopa: ali napovedujemo razred, ali napovedujemo verjetnost. Če napovedujemo razred, je funkcija izgube L ravno indikator, če smo se zmotili. Problem tu je, da ta funkcija izgube ne loči med lažnimi pozitivnimi in lažnimi negativnimi primeri, kar bomo poskusili popraviti kasneje. Pred tem si oglejmo

$$\text{Err} = \frac{\text{FP} + \text{FN}}{n}$$

in NAPOVEDNO TOČNOST

$$\text{Acc} = \frac{\text{TP} + \text{TN}}{n} = 1 - \text{Err}.$$

Vprašanje 13. Kaj sta napovedna napaka in točnost pri dvojiški klasifikaciji?

Boljši meri sta DELEŽ PRAVILNO RAZVRŠČENIH PRIMEROV ali OBČUTLJIVOST

$$\text{TPR} = \frac{\text{TP}}{\text{P}},$$

ki ima za idealne modele vrednost 1, in DELEŽ NAPAČNO RAZVRŠČENIH PRIMEROV

$$\text{FPR} = \frac{\text{FP}}{N} = 1 - \frac{\text{TN}}{N},$$

ki ima v idealnem primeru vrednost 0. Ti meri lahko uporabimo kot ordinatno in abscisno os v t.i. prostoru ROC – RECEIVER OPERATING CHARACTERISTIC. Modeli z enako napako so v tem prostoru na isti (pozitivni) diagonali, z idealnim modelom v zgornjem levem kotu in naključnim modelom na simetrali lihih kvadrantov. Če je ena napaka tipa FP ali FN slabša kot druga, lahko primerno naklonimo premice enakovrednih modelov, (ki zdaj ne bodo $y = x + c$, temveč $y = kx + c$), in izbiramo med modeli na primeren način.

Vprašanje 14. Kaj je prostor ROC? Kako v njem poiščeš idealni model, če so napake FP α -krat slabše od napake FN?

Pri klasifikaciji s pomočjo verjetnosti si zadamo prag θ , da je napoved enaka $\mathbb{1}(m(\mathbf{x}) \geq \theta)$. Če spreminjamo θ od 0 do 1, se pomikamo v ROC diagramu od levega spodnjega kota do desnega zgornjega. Idealen model potem izberemo na podoben način kot prej, s sweepom premice določenega naklona.

Ploščina pod krivuljo ROC nam poda dodatno mero AUC, ki je neodvisna od izbire θ (oz. se nanaša na vse možne izbire θ). Poda nam oceno razdalje med verjetnosti za pozitivne in negativne primere.

Vprašanje 15. Kako izbereš optimalen odločitveni prag pri klasifikaciji z verjetnostjo?

7.5 Odločitvena drevesa

Odločitveno drevo razbije prostor na dele, na katerih je napoved konstantna. Pri tem uporabljamo le teste oblike $X_i < a_j$; ena veja predstavlja resnično, ena pa neresnično vrednost tega izraza. Končna vozlišča podajo napovedi. Za konstrukcijo uporabimo algoritem TDIDT (TOP-DOWN INDUCTION OF DECISION TREES). Definiramo mero nečistoče (IMPURITY) za neko množico; potem je nečistoča drevesa enaka

$$\text{Impurity}(T) = \sum_{l \in T} \frac{|S_l|}{|S|} \text{Impurity}(S_l),$$

kjer so l listi drevesa. Želimo poiskati drevo z najmanjšo nečistočo, to pa je žal NP-poln problem, zato uporabimo požrešni algoritem; v vsakem koraku vzamemo poddrevesi, ki imata skupaj najmanjšo nečistočo. Taki poddrevesi poiščemo s polnim iskanjem možnih testov. Ko napovedujemo vrednosti, je napoved lista enaka povprečju vrednosti ciljne spremenljivke v tem listu (ali večinska vrednost v diskretnem primeru).

Vprašanje 16. Kako deluje algoritem TDIDT?

Pri izbiri mere nečistoče ločimo numerične in diskretne primere. V numeričnem merimo varianco

$$\text{Impurity}(S) = \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} (y - \bar{y})^2,$$

pri diskretnem primeru pa je $\text{Impurity}(S) = \phi(p_1, \dots, p_c)$, kjer so

$$p_i = P(Y = y_i | S) = \frac{|\{y = y_i\}|}{|S|},$$

ϕ pa je neka funkcija, za katero želimo, da čim boljše modelira varianco. Zahtevamo, da ima največjo vrednost pri enakomerni porazdelitvi, minimalno pri determinističnih, ter da je neodvisna od permutacije argumentov. Imamo dve pogosti možnosti:

- Entropija: $\phi = \sum_i p_i \log_2 p_i$
- Indeks Gini: $\phi = 1 - \sum_i p_i^2$

Vprašanje 17. Kakšne mere nečistoče poznamo pri gradnji odločitvenih dreves?

Če so drevesa prevelika, lahko v listih dovoliš $\text{Impurity} > 0$, čemur pravimo TREE PRUNING. To lahko naredimo na dva načina. Če si nastavimo minimalno število primerov, ki jih ne bomo več ločevali, temu pravimo SPROTNO REZANJE, če pa prvo naredimo celo drevo, in nato neka notranja vozlišča spremenimo v končna, pa govorimo o NAKNADNEM REZANJU. Pri tem gledamo napako poddrevesa

$$\text{Err}_\alpha(T) = \text{Err}(T) + \alpha |T|,$$

kjer je T število listov v drevesu, α pa izbran parameter. V algoritmu primerjamo napako v staršu in v otrocih, in vzamemo manjšo.

Vprašanje 18. Kako upravljaš velikost odločitvenega drevesa?

7.6 Metoda podpornih vektorjev

Pri metodi podpornih vektorjev se ukvarjamo z dvojiško klasifikacijo. Iščemo najširši rob, podan z afino preslikavo, da bo za pozitivne primere veljalo $\beta^T x + \beta_0 \geq 1$, za negativne pa $\beta^T x + \beta_0 \leq -1$. V obeh primerih postavimo ciljno spremenljivko na $y_i = \pm 1$, da je $y_i(\beta^T x_i + \beta_0) \geq 1$. Vsak učni primer nam da eno tako neenakost. Če vzamemo dve točki na nasprotnem robu in upoštevamo, da je β normala na hiperravnini, v kateri se nahajata en pozitiven primer x_+ oz. en negativen primer x_- , potem je širina pasu, ki ga hiperravnini definirata, enaka

$$\frac{\beta^T(x_+ - x_-)}{\|\beta\|} = \frac{1 - \beta_0 - (-1 - \beta_0)}{\|\beta\|} = \frac{2}{\|\beta\|}.$$

Želimo imeti čim širši pas, torej iščemo maksimum tega izraza; ekvivalentno iščemo minimum $\frac{1}{2} \|\beta\|^2$ pod pogojem, da veljajo zgoraj zapisane neenakosti.

Minimum poiščemo z Lagrangeovimi multiplikatorji

$$\begin{aligned} L &= \frac{1}{2} \|\beta\|^2 - \sum_i \alpha_i (y_i(\beta^T \mathbf{x}_i + \beta_0) - 1), \\ \partial_\beta L &= \beta - \sum_i \alpha_i y_i \mathbf{x}_i, \\ \partial_{\beta_0} L &= - \sum_i \alpha_i y_i \end{aligned}$$

Če enačimo zadnji enačbi z 0 in vstavimo v L , dobimo

$$L = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j.$$

Karush-Kuhn-Tuckerjevi pogoji nam dajo dodatne omejitve $\alpha_i \geq 0$ ter

$$\sum_i \alpha_i y_i = 0.$$

Zahtevali smo $y_i(\beta^T \mathbf{x}_i + \beta_0) - 1 \geq 0$, iz KKT robnih pogojev pa dobimo $\alpha_i(y_i(\beta^T \mathbf{x}_i + \beta_0) - 1) = 0$. V primerih, ko je $\alpha_i \neq 0$, torej velja $y_i(\beta^T \mathbf{x}_i + \beta_0) = 1$; to so ravno primeri na robu. Pravimo jim **PODPORNI VEKTORJI**. Napoved modela bo enaka

$$\hat{y}_0 = g\left(\sum_i \alpha_i y_i \mathbf{x}_i^T \mathbf{x}_0 + \beta_0\right),$$

kjer je g indikatorska funkcija

$$g(x) = \begin{cases} 1, & x \geq 0 \\ 0, & x < 0 \end{cases}$$

Vprašanje 19. Izpeljite metodo podpornih vektorjev za linearno ločljive podatke.

Če podatki niso linearno ločljivi, lahko v optimizacijsko funkcijo dodamo člen, ki meri oddaljenost od pravičnega dela prostora, iščemo

$$\min_{\beta, \beta_0} \frac{1}{2} \|\beta\|^2 + C \sum_i \xi_i,$$

kjer je ξ_i razdalja i -tega primera od roba, ki temu primeru pripada (torej od pozitivnega roba, če je pozitiven, in od negativnega, če je negativen), oziroma 0, če je primer na pravični strani roba. Omejitve so sedaj oblike

$$\begin{aligned} y_i(\beta^T \mathbf{x}_i + \beta_0) &\geq 1 - \xi_i, \\ \xi_i &\geq 0. \end{aligned}$$

Tudi za ta problem lahko najdemo optimalno rešitev s podobnim postopkom. Parameter C določa, kako pomembno je manjšanje števila podpornih vektorjev, ki so sedaj lahko tudi znotraj roba (ali na drugi strani). Večja vrednost C pomeni, da bo model bolj ločeval razrede, s čimer se zniža predsodek, a poveča varianca.

Vprašanje 20. Kaj narediš, če podatki v metodi podpornih vektorjev niso linearno ločljivi? Kaj je vloga parametra cene?

Če imamo več kot dva razreda, lahko klasifikacijo naredimo na dva načina:

- Zgradimo model za vsak par možnih vrednosti v klasifikaciji
- Zgradimo model za vsako možnost, ki ločuje podatke tega tipa od vseh ostalih

Če podatki niso linearni, jih transformiramo z neko funkcijo ϕ , da nelinearna odločitvena meja postane linearna. Raje pa definiramo

$$K(u, v) = \langle \phi(u), \phi(v) \rangle,$$

ki nam dovoli, da se ne ubadamo s prehodom v transformiran prostor. Jedrno funkcijo uporabimo namesto skalarnega produkta v predpisu za optimizacijo. Na splošno lahko ϕ slika v poljuben Hilbertov prostor, zato so lahko jedra konkretno komplicirana; če želimo, lahko ϕ recimo slika v višjedimenzionalni prostor kot naš prostor podatkov.

Vprašanje 21. Kako modificiraš metodo podpornih vektorjev z jedrnimi funkcijami?

7.7 Ansambli napovednih modelov

V ansamblu \hat{M} imamo več modelov \hat{m}_i . Ideja je, da vrednosti ciljnih spremenljivk kombiniramo iz napovedi posameznih modelov v ansamblu. Običajno to naredimo s povprečjem, za razvrščanje pa vzamemo gostišnico. Za razvrščanje lahko pridemo do pričakovane napake: če vsakemu od T modelov pripišemo naključno spremenljivko $Z_i = \mathbb{1}(y \neq \hat{y}_i)$, in predpostavimo, da so te spremenljivke neodvisne z $E(Z_i) = \varepsilon_i$, potem je

$$\begin{aligned} \text{Err}(\hat{M}) &= P\left(\sum_i Z_i > \frac{T}{2}\right) = P\left(\sum_i (Z_i - E(Z_i)) + T\varepsilon > \frac{T}{2}\right) \\ &= P\left(\frac{1}{T} \sum_i (Z_i - E(Z_i)) > \frac{1}{2} - \varepsilon\right) \leq \exp\left(-2T\left(\frac{1}{2} - \varepsilon\right)^2\right) \end{aligned}$$

po Hoeffdingovi neenakosti. Napaka je torej reda e^{-T} , če je le $\varepsilon < \frac{1}{2}$.

Vprašanje 22. Ocení napako klasifikacijskega ansambla.

Obravnavajmo homogen ansambel regresijskih modelov. Razlike potem pridejo iz različnih podatkovnih množic S , ki jih obravnavajmo kot slučajno spremenljivko.

$$E_S(\hat{M}(\mathbf{x}_0)) = E_S\left(\frac{1}{T} \sum_i \hat{m}_i(\mathbf{x}_0)\right) = \frac{1}{T} \sum_i E_S(\hat{m}_i(\mathbf{x}_0)) = E_S(\hat{m}(x_0))$$

Torej je pristranskost $E_S(\hat{M}(\mathbf{x}_0)) - m(\mathbf{x}_0)$ enaka kot pristranskost osnovnega modela \hat{m} .

Vprašanje 23. Pokaži, da ansambel ne spremeni pristranskosti.

Izračunamo lahko, da za korelacijski faktor

$$\rho_S(\mathbf{x}_0) = \frac{E_S((\hat{m}_i(\mathbf{x}_0) - E(\hat{m}_i(\mathbf{x}_0)))(\hat{m}_j(\mathbf{x}_0) - E(\hat{m}_j(\mathbf{x}_0))))}{\sqrt{\text{var}_S(\hat{m}_i(\mathbf{x}_0)) \text{var}_S(\hat{m}_j(\mathbf{x}_0))}}$$

velja

$$\text{var}_S(\hat{M}(\mathbf{x}_0)) = \left(\rho + \frac{1 - \rho}{T} \right) \text{var}_S(\hat{m}(\mathbf{x}_0)).$$

Levemu členu produkta pravimo FAKTOR ZMANJŠEVANJA VARIANCE.

Vprašanje 24. Kaj je faktor zmanjševanja variance? Izpelji predpis.

Za pridobivanje različnih modelov v homogenem ansamblu imamo več možnosti. Ena od njih je vrečenje, kjer spreminjamo učne podatke za vsak model z vzorčenjem skupne učne množice s ponavljanjem. Napako ansambla lahko potem ocenimo kot povprečje napak na učni množici, kjer napovedi za posamičen primer iz učne množice generiramo le z modeli, ki tega testa niso imeli med učnimi podatki. Temu pravimo OCENA OUT-OF-BAG.

Vprašanje 25. Kaj je ocena out-of-bag?

Drug način za pridobivanje različnih modelov so naključni podprostori, kjer vsak model naučimo na $q < p$ naključno izbranih spremenljivkah z vsemi podatki. Tu ni vzorčenja, tako da ne moramo uporabljati ocen out-of-bag.

Če je naš osnovni model odločitveno drevo, dobimo naključen gozd. Tukaj vsako drevo učimo na naključni podmnožici $q \leq p$ spremenljivk, kar pomaga dekorelirati drevesa, če imamo neko močno spremenljivko. Naključnih dreves običajno ne režemo.

Vprašanje 26. Kako narediš naključni gozd?

Ansambli sami po sebi ne podajajo jasne interpretacije. Edino, kar lahko naredimo, je da pogledamo napovedno moč spremenljivke, kar lahko naredimo na več načinov. V primeru naključnega gozda imamo POVPREČNO ZMANJŠEVANJE NEČISTOČE, ki je enaka povprečni vrednosti zmanjševanja nečistoče v vseh vozliščih drevesa, ki testirajo vrednost določene spremenljivke. Napovedna moč v ansamblu je potem povprečje vseh teh povprečij po modelih, ki ansambel sestavljajo. Običajno jo še normaliziramo, da je največja vrednost enaka 1.

Vprašanje 27. Opiši povprečno zmanjševanje nečistoče.

Druga strategija je povprečno zmanjševanje točnosti. Tukaj opazujemo razliko v napaki pravega modela in modela, treniranega na množici, kjer smo opazovano spremenljivko po stolpcih naključno permutirali, torej izničili njen vpliv. Če je razlika v napakah majhna, spremenljivka nima velikega vpliva.

Vprašanje 28. Opiši povprečno zmanjševanje točnosti.

7.8 Nevronske mreže

Nevronske mreže so sestavljene iz nevronov in sinaps. Vsak nevron hrani stanje $v \in \mathbb{R}$ iz izhod y , sinapse pa so utežene povezave med nevroni. Uteži označimo z w_i . Poleg tega ima vsak nevron aktivacijsko funkcijo ϕ , s pomočjo katere izračuna izhod: $y = \phi(v)$. Pri evalvaciji stanja nevrona izračunamo iz stanja njegovih vhodnih sinaps kot

$$v = w_0 + \sum_i w_i x_i.$$

Mrežo ustavimo, ko se neha spreminjati, zato običajno nima ciklov.

Pri usmerjeni nevronske mreži imamo nevrone razporejene v $L+2$ plasti. Plasti 0 pravimo VHODNA PLAST, kjer imamo en nevron za vsako vhodno spremenljivko, plasti $L+1$ pa IZHODNA PLAST. Tu je en nevron v primeru regresije in en nevron za vsak razred v primeru klasifikacije. Preostane L SKRITI PLASTI, ki nam dovoljujejo kompleksnejše računanje. Sinapse so postavljene tako, da je vsak nevron v neki plasti povezan z vsemi v naslednji plasti, ni pa povezan s prejšnjimi plastmi ali svojo plastjo.

Vprašanje 29. Opiši postavitev usmerjene nevronske mreže.

Najenostavnejša mreža je ENOSTAVNI PERCEPTRON. To je pravzaprav (skoraj) linearni model z $L = 0$ in stopničasto aktivacijsko funkcijo $\phi(v) = \mathbb{1}(v > 0)$. Enostavneje ga zapišemo kot

$$\hat{y} = \phi(\mathbf{w}^T \mathbf{x}).$$

Učenje poteka z gradientnim spustom, torej imamo predpis

$$\Delta \mathbf{w} = \eta(y - \phi(\mathbf{w}^T \mathbf{x}))\mathbf{x},$$

kjer v vsaki iteraciji obravnavamo en primer, parameter η pa se imenuje LEARNING RATE.

Vprašanje 30. Opiši enostavni perceptron.

Za splošnejšo klasifikacijsko nevronske mrežo imamo v izhodnem sloju toliko nevronov, kolikor je možnosti v klasifikaciji. Kot odgovor mreže dobimo neka števila v \mathbb{R} , ki jih pretvorimo v verjetnosti s funkcijo SOFTMAX

$$\phi(v_i) = \frac{e^{v_i}}{\sum_j e^{v_j}}.$$

Vprašanje 31. Kaj je softmax?

Če imamo več plasti, za gradientni spust potrebujemo odvode napake po w_{ji}^l , torej po utežeh med l -to in $(l+1)$ -to plastjo. Pri tem odvajamo posredno:

$$\frac{\partial E}{\partial w_{ji}^l} = \frac{\partial E}{\partial y_i^l} \frac{\partial y_i^l}{\partial v_i^l} \frac{\partial v_i^l}{\partial w_{ji}^l} = \frac{\partial E}{\partial y_i^l} \phi'(v_i^l) y_j^{l-1},$$

odvod napake po y_i^l pa izračunamo kot

$$\frac{\partial E}{\partial y_i^l} = \sum_k \frac{\partial E}{\partial v_k^{l+1}} \frac{\partial v_k^{l+1}}{\partial y_i^l} = \sum_k \frac{\partial E}{\partial y_k^{l+1}} \frac{\partial y_k^{l+1}}{\partial v_k^{l+1}} \frac{\partial v_k^{l+1}}{\partial y_i^l} = \sum_k \frac{\partial E}{\partial y_k^{l+1}} \phi'(v_k^{l+1}) \frac{\partial v_k^{l+1}}{\partial y_i^l},$$

kjer seštevamo po nevronih v naslednji plasti. Računamo lahko torej od desne proti levi, da dobimo vse odvode. V plasti $L + 1$ namesto drugega predpisa za $\partial E / \partial y_i^{L+1}$ uporabimo $-(y_i - y_i^{L+1})$, ki izhaja iz predpisa napake $E = \frac{1}{2} \sum_i (y_i - \hat{y}_i)^2$ in dejstva, da dejansko odvajamo po \hat{y}_i .

Vprašanje 32. Izpelji predpis za računanje gradienta v večplastni nevronske mreži.

V primeru klasifikacije uporabljamo drugačno funkcijo izgube, prečno entropijo

$$H = - \sum_i P_i \log_2 Q_i,$$

kjer je P prava (diskretna) porazdelitev cilje porazdelitve Y , torej $P_k = 1$ natanko tedaj, ko je trenutno opazovana vrednost enaka k , Q pa je porazdelitev, ki jo napove model, $Q_i = y_i^{L+1}$. Potem je $E = -\log_2 y_k^{L+1}$, odvod pa

$$\frac{\partial E}{\partial \hat{y}_i} = \begin{cases} 0 & i \neq k \\ \frac{1}{\hat{y}_k \ln 2} & i = k \end{cases}$$

z enakim nadaljnjim delom.

Vprašanje 33. Kako izračunaš odvode za primer klasifikacijske nevronske mreže?

Poznamo tudi konvolucijske nevronske mreže, kjer so nevroni v posamičnih plasteh organizirani v matrike. Namesto linearne obtežene vsote po vseh nevronih uporabimo konvolucijski filter, tj. matriko $c \times c$ uteži w_{ij} , s katero kombiniramo nevrone. Če je plast $l - 1$ dimenzije $x \times y$, je plast l potem dimenzije $(x - c + 1) \times (y - c + 1)$. Pogosto želimo tudi zmanjšati dimenzije teh matrik, kar običajno naredimo z max-akumulacijo, ki vzame največjo vrednost v neki izbrani podmatriki; ponavadi velikosti 2×2 .

Vprašanje 34. Opiši delovanje konvolucijskih nevronske mreže.

Poznamo tudi rekurenčne nevronske mreže, kjer dovolimo pojavitev zanke. Za možnost računanja gradientov diskretiziramo čas in pri zanki za vhod uporabimo izhod prejšnje iteracije; tako si lahko zapomnimo prejšnje podatke. S tem smo pravzaprav dobili zaporedje enakih nevronske mreže, ki ustrezajo časovnim točkam. Problem se pojavi pri učenju; napaka je sedaj funkcija časa

$$E = \frac{1}{T} \sum_{t=1}^T L(y_t, \hat{y}_t),$$

kjer je L funkcija izgube. Na podoben način kot prej lahko izračunamo odvod te napake, pri čemer dobimo rekurzivno zvezo v t .

Vprašanje 35. Opiši rekurenčne nevronske mreže.

7.9 Nenadzorovano učenje

Pri nenadzorovanem učenju nimamo ciljne spremenljivke Y . Namesto napovedi iščemo povezane skupine podatkov („clustering“). Kot vhodni podatek dobimo množico podatkov S z metriko d , in želimo število skupin k . Prva možnost tu je algoritem HAC. Začnemo s $|S|$ enojci, dokler imamo več kot eno skupino, najbližji dve združimo. Rezultat je dendrogram, kjer y koordinata povezave med skupinama prikaže medsebojno razdaljo. Razdaljo med skupinama C_1 in C_2 lahko definiramo na več načinov, ali kot najmanjšo/največjo/povprečno razdaljo med primeri, ali pa bolj ekonomično; za tretjo skupino C je

$$d(C_1 \cup C_2, C) = \alpha_1 D(C_1, C) + \alpha_2 D(C_2, C) + \beta D(C_1, C_2)$$

za

$$\alpha_1 = \frac{|C_1| + |C|}{|C_1| + |C_2| + |C|},$$

$$\beta = -\frac{|C|}{|C_1| + |C_2| + |C|}.$$

Temu pravimo WARDOVA FORMULA.

Vprašanje 36. Opiši delovanje algoritma HAC in povej Wardovo formulo.

Druga metoda je k -means clustering. Tukaj iščemo

$$\min_{|C|=k} \sum_{C \in \mathcal{C}} |C| \text{var}(C),$$

kar lahko naredimo za evklidsko razdaljo. Izračunamo CENTROID (težišče)

$$c = \arg \min_u \sum_{v \in S} d(u, v) = \frac{1}{|S|} \sum_{v \in S} \vec{v},$$

kar dobimo z računanjem odvodov. Varianca skupine je vsota kvadratov razdalj med primeri in centroidom.

Razbitje izračunamo iterativno. Začnemo z naključno razporeditvijo, in potem iterativno izračunamo centroide vseh skupin ter vse primere povežemo s tistim centroidom, ki jim je najbližje. Iteracijo ustavimo, ko ni več spremembe.

Če imamo neevklidsko mero razdalje, se lahko zgodi, da centroida ne znamo poiskati. V tem primeru ga nadomestimo z MEDOIDOM

$$m = \arg \min_{u \in S} \sum_{v \in S} d(u, v).$$

Vprašanje 37. Opiši delovanje algoritma k -means clustering. Kaj naredimo, če nimamo evklidske metrike?

Kvaliteto razvrščanja merimo s povezanostjo in ločenostjo primerov od ostalih. POVEZANOST je definirana kot

$$a(u) = \frac{1}{|C_u| - 1} \sum_{v \in C_u} d(u, v),$$

kjer je C_u skupina, kateri pripada u . Delimo s $|C_u| - 1$, ker je razdalja $d(u, u) = 0$. LOČENOST pa je definirana kot

$$b(u) = \min_{C \in \mathcal{C}, C \neq C_u} \frac{1}{|C|} \sum_{v \in C} d(u, v).$$

Ti vrednosti kombiniramo v OBRIS

$$s(u) = \frac{b(u) - a(u)}{\max(a(u), b(u))},$$

ki ima vrednosti med -1 in 1 . Primer je razvrščen dobro za s blizu 1 in slabo za s blizu -1 ; mera kvalitete celotnega razvrščanja bo povprečje $s(u)$ čez vse u .

Vprašanje 38. Kako vrednotimo razvrščanje?

7.10 Krčenje razsežnosti

Prvo vprašanje je, zakaj bi razsežnost sploh krčili. Odgovor je prekletstvo večrazsežnosti: Naj bodo X_1, \dots, X_n točke v enotski krogli v \mathbb{R}^p , porazdeljene enakomerno, in naj bo M najmanjša razdalja teh točk do središča. Izračunamo lahko

$$P(M \leq x) = 1 - (1 - x^p)^n.$$

Potem je mediana (tj. razdalja, pri kateri je najbližja točka z verjetnostjo $\frac{1}{2}$ v krogli s tem radijem) enaka

$$x_m = F_M^{-1}\left(\frac{1}{2}\right) = \left(1 - \frac{1}{2^{1/n}}\right)^{1/p}.$$

Z veliko točkami se oddaljujemo od 1 , z veliko dimenzijami pa se bližamo 1 . Pozorni moramo biti na razmerje n/p .

Vprašanje 39. Demonstriraj problem večrazsežnosti.

Želeli bi si RANGIRANJE napovednih spremenljivk X_i po pomembnosti. Za to definiramo RELEVANTNOST spremenljivke kot funkcijo $r_S : \mathbf{X} \rightarrow \mathbb{R}$, ki za podatkovno množico S pove, kako pomembna je spremenljivka $X_i \in \mathbf{X}$ (tu je \mathbf{X} množica spremenljivk).

Mere relevantnosti ločimo na nadzorovane in nenadzorovane. Za nenadzorovane opazimo, da je relevantnost je pozitivno korelirana z varianco spremenljivke, saj večja varianca pomeni, da imamo več prostora za odločitve. Ta mera pa se izgubi, če podatke normaliziramo. V primeru diskretne spremenljivke lahko namesto variance vzamemo entropijo

$$H(X_i) = - \sum_v p_v \log_2 p_v,$$

kjer je $p_v = |\{x_i \in S \mid x_i = v\}| / |S|$.

Vprašanje 40. Opiši nenadzorovano mero relevantnosti.

Pri nadzorovanih merah relevantnosti začnemo z uni-variantnimi metodami. Pri teh gledamo samo povezanost med X_i in Y za nek i . Ker uporabljajo Y , so nadzorovane; ker ignorirajo vse ostale X_j , so uni-variantne. Če sta X_i in Y numerični, je ena možnost korelacijski faktor, imamo pa še cel kup drugih možnosti. V primeru numerične X_i in diskretne Y lahko opazujemo RELIEF, kjer za naključno izbran primer $(\mathbf{x}, y) \in S$ poiščemo najbližjega sosedo v projekciji na X_i , \mathbf{x}_H , ki pripada istemu razredu kot \mathbf{x} , in najbližjega sosedo \mathbf{x}_M , ki pripada različnemu razredu. Če je razdalja do \mathbf{x}_M veliko večja od razdalje do \mathbf{x}_H , je to dober znak za relevantnost spremenljivke, na katero smo projicirali. V algoritmu je relevantnost X_i potem enaka

$$\frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} (|\mathbf{x}_{Mi} - \mathbf{x}_i| - |\mathbf{x}_{Hi} - \mathbf{x}_i|).$$

Vprašanje 41. Opiši metodo reliefa.

Zadnja možnost je ocenjevanje relevantnosti na podlagi modela. Če imamo razumljiv model, npr. linearni, lahko relevantnost direktno preberemo iz njega. V primeru manj razumljivega modela lahko permutiramo vrednosti v stolpcu i v podatkih in natreniramo nov model; spremenljivka je tako uničena, relevantnost dobimo kot razliko v napovedni napaki. Ta metoda je najpočasnejša, a je zelo kvalitetna.

Vprašanje 42. Kako oceniš relevantnost spremenljivke z danim modelom?

Pri zmanjševanju razsežnosti pogosto želimo tvoriti nove spremenljivke, ki bodo bolj relevantne za nas. Običajno ima to obliko linearne transformacije $Z = XW$, kjer je Z dimenzije $q \ll p$. Za to imamo metodo glavnih komponent, kjer dodatno zahtevamo, da je W ortogonalna, in da dimenzije Z pojasnijo čim več variance v osnovnih podatkih. Predpostavimo, da imajo vse spremenljivke v X povprečje 0, torej je kovariančna matrika enaka $C = X^T X$. Ta je simetrična nenegativno definitna, torej jo lahko s pomočjo singularnega razcepa zapišemo kot $C = WDW^T$, kjer je D diagonalna matrika kvadratov singularnih vrednosti.

Vprašanje 43. Razloži metodo glavnih komponent.

7.11 Manjkajoče vrednosti

Včasih v testnih ali učnih primerih manjka kakšen podatek. Če manjka y , primer samo izbrišemo, ker nam ne pomaga. To lahko naredimo tudi za manjkajoče X_i , čemur pravimo COMPLETE-CASE ANALYSIS. Tu lahko pride do problemov; eno je predsodek v podatkih, drugo pa, da ima tako učna množica lahko premalo primerov. Podobna ideja je AVAILABLE-CASE ANALYSIS, kjer namesto primerov brišemo spremenljivke (stolpce), kjer podatki manjkajo.

Vprašanje 44. Kakšne so težave z brisanjem neznanih podatkov?

Če podatkov ne brišemo, delamo IMPUTACIJO. Numerično spremenljivko lahko nadomestimo s povprečjem vseh primerov v podatkovni množici, diskretno pa z gostišnico, ali pa neznane podatke postavimo v razred zase. Dobra ideja je, da dodamo še novo indikatorsko spremenljivko, ki pove, ali je bila originalna spremenljivka prisotna, ali pa je imputirana. Druga možnost, ki je baje „relativno korektna“, je, da neznano vrednost nadomestimo z naključno vrednostjo iz tega stolpca. To je dejanje iz obupa; v praksi je bolje, da uporabiš znanje o problemu. Korektno pa je, ker tako ne prinašaš predsodka, ker se porazdelitev ne spremeni. Še zadnja osnovna možnost je nadomeščanje s povprečno vrednostjo najbližjih sosedov (po znanih spremenljivkah).

Vprašanje 45. Opiši osnovne tehnike imputacije.

Bolj napredno je nadomeščanje z napovednimi modeli. To je iterativni postopek, kjer v prvem koraku napovemo manjkajoče vrednosti z eno od enostavnih metod, v vseh nadaljnjih korakih pa izberemo eno od spremenljivk (v nekem vrstnem redu), naučimo napovedni model na vseh X_i brez te, in napovemo vrednost te spremenljivke. Tu ne uporabljamo podatkov iz Y , ker bi potem bile generirane vrednosti pristranske, kar nam efektivno uniči testne primere.

Vprašanje 46. Opiši imputacijo z napovednimi modeli.

Izjema tu so modeli, ki sprejemajo tudi neznane vrednosti, recimo odločitvena drevesa. V algoritmu TDIDT primere z neznano vrednostjo X_i enakomerno porazdelimo med vse naslednike; te primere v računih obravnavamo, kot da bi bili na pol v eni in na pol v drugi množici.

Vprašanje 47. Kako prilagodiš TDIDT tako, da lahko dela tudi z neznanimi vrednostmi?

7.12 Neenakomerne porazdelitve

Težava pri neenakomernih porazdelitvah Y je, da bomo za pogoste vrednosti dobili manjšo napako. Za porazdelitev

$$Y \sim \begin{pmatrix} 0.01 & 0.99 \\ \oplus & \ominus \end{pmatrix}$$

je model $Y = \ominus$ perfekten, a ga ne želimo uporabiti iz očitnih razlogov. Za diskreten Y vzorčimo, torej upoštevamo manj velikih razredov (PODVZORČENJE) ali pa primere majhnih razredov večkrat kopiramo (PREVZORČENJE).

Vprašanje 48. Kaj je podvzorčenje in kaj prevzorčenje?

Noben od teh pristopov ni idealen. Malo boljši je algoritem SMOTE, kjer tvorimo nove sintetične primere. Algoritem za izbran primer (x, y) poišče k najbližjih sosedov in za

naključnega (x_n, y) z enakim y tvori sintetični primer $(x + g \cdot (x_n - x), y)$, kjer je $g \in [0, 1]$ naključno izbran. Dodatno algoritem sprejme tudi parametra o stopnji podvzorčenja in prevzorčenja; to sta faktorja v številu novih razredov manjšinskega oz. večinskega razreda glede na originalno število.

Vprašanje 49. Opiši algoritem SMOTE.

Pri regresiji neenakomerna porazdelitev izgleda tako, da imamo osamelce z zelo različnimi vrednostmi Y . Definiramo FUNKCIJO POMEMBOSTI $\phi : Y \rightarrow [0, 1]$, ki da visoko pomembnost redkim primerom in zvezno pada k pogostejšim. Primere z visoko pomembnostjo prevzorčimo, primere z nizko pomembnostjo pa podvzorčimo. V regresijski obliki algoritma SMOTE v sintetičnih primerih uporabimo konvolucijo tudi v spremenljivki Y .

Vprašanje 50. Kako obravnavaš vprašanje o neenakomerni porazdelitvi v primeru regresije?

8 Numerična linearna algebra

8.1 Singularni razcep

Izrek. Za vsako matriko $A \in \mathbb{R}^{m \times n}$, kjer je $m \geq n$, obstaja razcep $A = U\Sigma V^T$, kjer je $U \in \mathbb{R}^{m \times m}$ ortogonalna, $V \in \mathbb{R}^{n \times n}$ ortogonalna in $\Sigma \in \mathbb{R}^{m \times n}$ oblike

$$\Sigma = \begin{bmatrix} \sigma_1 & & & \\ & \ddots & & \\ & & \sigma_n & \\ & & 0 & \end{bmatrix}$$

za singularne vrednosti $\sigma_1 \geq \dots \geq \sigma_n \geq 0$ matrike A .

Dokaz. Če je $A = U\Sigma V^T$, potem je $A^T A = V\Sigma^T \Sigma V^T$. Ta matrika je simetrična nenegativno definitna, torej so njene lastne vrednosti realne in nenegativne ter jih lahko uredimo padajoče. Lastne vektorje lahko izberemo ortonormirane in jih zložimo v stolpce V .

Iz $AV = U\Sigma$ sledi $Av_i = \sigma_i u_i$ za stolpce v_i in u_i . Če je $\sigma_i \neq 0$, lahko tak u_i poiščemo. Če je $\sigma_r > \sigma_{r+1} = \sigma_{r+2} = \dots = \sigma_n = 0$, tako določimo prvih r stolpcev U . Ker je v_j lastni vektor za $A^T A$, velja

$$u_i^T u_j = \frac{1}{\sigma_i \sigma_j} v_i^T A^T A v_j = \frac{1}{\sigma_i \sigma_j} \sigma_j^2 v_i^T v_j = \delta_{ij}.$$

Sedaj imamo

$$\begin{aligned} V &= [V_1 \quad V_2], \\ \Sigma &= \begin{bmatrix} S & 0 \\ 0 & 0 \end{bmatrix}, \\ U &= [U_1 \quad U_2], \end{aligned}$$

kjer smo U_2 določili tako, da dopolnimo stolpce U_1 do ortonormirane baze prostora.

Preverimo, da te matrike res tvorijo singularni razcep:

$$U^T A V = \begin{bmatrix} U_1^T A V_1 & U_1^T A V_2 \\ U_2^T A V_1 & U_2^T A V_2 \end{bmatrix}.$$

Velja $(AV_2)^T AV_2 = V_2^T A^T AV_2 = V_2^T \cdot 0 = 0$, ker so stolpci V_2 lastni vektorji za lastno vrednost 0. Iz definicije u_i pa vidimo, da je $U_1^T AV_1 = S$. Velja $AV_1 = [\sigma_1 u_1, \dots, \sigma_r u_r]$, ti stolpci so pravokotni vrednostim u_{r+j} po definiciji. \square

Če je $m < n$, singularni razcep dobimo tako, da transponiramo razcep za A^T . V dokazu je r rang matrike A . Numerično je singularni razcep najboljše orodje za računanje ranga; dobimo tudi bazi za jedro in sliko matrike A (stolpci V_2 in U_1).

Vprašanje 1. Pokaži, da za vsako matriko obstaja singularni razcep.

Če je A ranga $r < n$, potem za vektor x , ki reši predoločen sistem $Ax = b$, izberemo tistega, ki minimizira $\|Ax - b\|_2$ (ta ni več enolično določen), in ki ima med takimi minimalno normo $\|x\|_2$.

Izrek. Naj bo $A \in \mathbb{R}^{m \times n}$, $A = U\Sigma V^T$ ranga r . Potem je rešitev predločenega sistema $Ax = b$ enaka

$$x = \sum_{i=1}^r \frac{u_i^T b}{\sigma_i} v_i.$$

Dokaz. Zapišimo $U = [U_1 U_2]$, $V = [V_1 V_2]$ in $\Sigma = \text{diag}(S, 0)$, kjer so prve komponente širine r , druge pa širine $n - r$. Potem je

$$\|Ax - b\|_2 = \|U\Sigma V^T x - b\|_2 = \|\Sigma V^T x - U^T b\|_2 = \left\| \begin{bmatrix} Sy_1 - U_1^T b \\ -U_2^T b \end{bmatrix} \right\|$$

za $V^T x = y = (y_1, y_2)$. Minimum bo dosežen, če je $Sy_1 = U_1^T b$, pri čemer je y_2 lahko poljuben. Velja $\|x\|_2 = \|y\|_2$, torej izberemo $y_2 = 0$. \square

Vprašanje 2. Kakšna je rešitev problema najmanjših kvadratov $Ax = b$ za matriko $A \in \mathbb{R}^{m \times n}$? Dokaži.

Definicija. Za matriko $A \in \mathbb{R}^{m \times n}$ je PSEVDONVERZ matrika $X \in \mathbb{R}^{n \times m}$, ki zadošča naslednjim točkam:

- $AXA = A$
- $XAX = X$
- $(AX)^T = AX$
- $(XA)^T = XA$

Označimo $X = A^+$.

Če je rang matrike A enak $r < n$, potem je

$$A^+ = (A^T A)^{-1} A^T.$$

Vprašanje 3. Definiraj psevdoinverz. Čemu je enak?

Izrek. Naj bo $A = U\Sigma V^T \in \mathbb{R}^{m \times n}$ in $\text{rang } A = r$. Potem je psevdoinverz enak

$$A^+ = \sum_{i=1}^r \frac{1}{\sigma_i} v_i u_i^T$$

oziroma $A^+ = V\Sigma^+ U^T$ za $\Sigma^+ = \text{diag}(S^{-1}, 0)$.

Dokaz. Iščemo $X \in \mathbb{R}^{n \times m}$, ki zadošča Moore-Penroseovim pogojem. Velja $X = VYU^T$ za nek Y . Veljati mora $(AX)^T = AX$, iz česar z zapisom po blokkih dobimo, da je zgornji desni blok Y enak 0; podobno iz drugih pogojev izpeljemo enakosti v ostalih blokkih. \square

Vprašanje 4. Kako izračunaš psevdoinverz s pomočjo singularnega razcepa?

8.1.1 Aproksimacija z matrikami nižjega ranga

Izrek (Eckart-Young-Mirskog). Če je $A = U\Sigma V^T$ in $\text{rang } A = r$, je za $k < r$ matrika

$$A_k = \sum_{i=1}^k \sigma_i u_i v_i^T$$

tista, ki minimizira

- $\min_{\text{rang } B=k} \|A - B\|_2,$
- $\min_{\text{rang } B=k} \|A - B\|_F.$

Dodatno velja $\|A - A_k\|_2 = \sigma_{k+1}$ in $\|A - A_k\|_F = \sqrt{\sigma_{k+1}^2 + \dots + \sigma_n^2}.$

Dokaz. Dokažemo samo prvo točko. Naj bo B matrika ranga k . Definiramo $V_{k+1} = [v_1, \dots, v_{k+1}]$. Velja $\dim(\ker B) = n - k$ in $\dim(\text{im } V_{k+1}) = k + 1$, torej obstaja vektor $w \neq 0$ v preseku. Brez škode za splošnost je $\|w\|_2 = 1$. Velja

$$\|A - B\|_2 = \max_{\|x\|_2=1} \|(A - B)x\|_2 \geq \|(A - B)w\|_2 = \|Aw\|_2 = \sqrt{w^T A^T A w} \geq \sigma_{k+1}.$$

Po drugi strani za A_k očitno velja $\|A - A_k\|_2 = \sigma_{k+1}$. □

Vprašanje 5. Kaj je najboljša aproksimacija matrike z matriko nižjega ranga? Dokaži.

Če je $\sigma_{k+1} < \sigma_k$, je iz dokaza razvidno, da je najboljša aproksimacija enolična. Velikost σ_{k+1} nam poda mero za razdaljo od matrik ranga k .

Aproksimacijo lahko uporabimo za delo z velikimi matrikami; če vemo, da je matrika nizkega ranga, moramo shraniti za red velikosti manj podatkov.

8.1.2 Regularizacija

Pri Fredholmovi integralni enačbi prve vrste

$$\int_0^1 K(s, t) f(t) dt = g(s)$$

sta podana jedro K in funkcija g , izračunati pa želimo f . Vsako jedro lahko razvijemo v

$$K(s, t) = \sum_{i=1}^{\infty} \sigma_i u_i(s) v_i(t),$$

kjer so u_i in v_i leve in desne singularne funkcije, $\sigma_1 \geq \sigma_2 \geq \sigma_3 \geq \dots$ pa singularne vrednosti, velja $\lim \sigma_n = 0$. Funkcije u_i, v_i so ortonormirane za skalarni produkt

$$\langle g, h \rangle = \int_0^1 g(t)h(t)dt,$$

poleg tega velja

$$\int_0^1 K(s, t)v_i(t)dt = \sigma_i u_i(s).$$

Podobno kot za singularni razcep lahko rešitev poiščemo z

$$f(t) = \sum_{i=1}^{\infty} \frac{1}{\sigma_i} \langle u_i, g \rangle v_i(t).$$

Pričakujemo, da bodo prispevki poznih členov majhni, torej da $\langle u_i, g \rangle$ limita k 0 hitreje kot σ_i . Če namesto g poznamo $g + \Delta g$, lahko izračunamo

$$\tilde{f} = \sum_{i=1}^{\infty} \frac{\langle u_i, g \rangle}{\sigma_i} v_i + \sum_{i=1}^{\infty} \frac{\langle u_i, \Delta g \rangle}{\sigma_i} v_i.$$

Problem je v tem, da pričakujemo, da bo šum razporejen po vseh komponentah funkcijskega prostora, torej bodo $\langle u_i, \Delta g \rangle$ majhni, a ne bodo konvergirali k 0. Majhna sprememba g lahko torej povzroči poljubno veliko spremembo v rezultatu.

Pri numeričnem reševanju bomo skalarni produkt izračunali z neko kvadraturno formulo

$$\int_0^1 f(t)dt = \sum_{i=0}^n \alpha_i f(t_i),$$

torej kot rešitev sistema

$$\begin{bmatrix} \alpha_0 K(s_0, t_0) & \alpha_1 K(s_0, t_1) & \cdots & \alpha_n K(s_0, t_n) \\ \alpha_0 K(s_1, t_0) & \alpha_1 K(s_1, t_1) & \cdots & \alpha_n K(s_1, t_n) \\ \vdots & \vdots & \ddots & \vdots \\ \alpha_0 K(s_n, t_0) & \alpha_1 K(s_n, t_1) & \cdots & \alpha_n K(s_n, t_n) \end{bmatrix} \cdot \begin{bmatrix} f(t_1) \\ \vdots \\ f(t_n) \end{bmatrix} = \begin{bmatrix} g(s_0) \\ \vdots \\ g(s_n) \end{bmatrix}.$$

Singularne vrednosti matrike so približki največjih singularnih vrednosti funkcije. Za večje matrike bo aproksimacija boljša, ampak občutljivost visoka.

Rešitev teh težav je regularizacija. Rešujemo sistem $Ax = b$ z $A = U\Sigma V^T$ in rešitvijo

$$x = \sum_{i=1}^n \frac{u_i^T b}{\sigma_i} v_i.$$

Če je prisoten še šum, dobimo rešitev $x + \Delta x$,

$$\Delta x = \sum_{i=1}^n \frac{u_i^T \Delta b}{\sigma_i} v_i.$$

Če $\frac{1}{\sigma_i} u_i^T b$ padajo k 0, $|u_i^T \Delta b|$ pa je enakega velikostnega reda za vse i , lahko uporabimo odrezan singularni razcep in psevdoinverz najboljše aproksimacije A ranga k ;

$$x_k = \sum_{i=1}^k \frac{u_i^T b}{\sigma_i} v_i = A_k^+ b.$$

Druge možnost je regularizacija Tihonova

$$x_{\text{reg}} = \sum_{i=1}^n \phi_i \frac{u_i^T b}{\sigma_i} v_i,$$

kjer so FAKTORJI FILTRA definirani kot

$$\phi_i = \frac{\sigma_i^2}{\sigma_i^2 + \alpha^2}$$

za PARAMETER REGULARIZACIJE α .

Izrek. Regularizacija Tihonova s parametrom $\alpha > 0$ vrne tisti $x \in \mathbb{R}^n$, ki minimizira $\|Ax - b\|_2^2 + \alpha^2 \|x\|_2^2$.

Dokaz. Zapisan x je natanko rešitev normalnega sistema

$$\begin{bmatrix} A \\ \alpha I \end{bmatrix} x = \begin{bmatrix} b \\ 0 \end{bmatrix}.$$

□

Vprašanje 6. Kaj je problem, ki ga rešujemo z regularizacijo? Povej izrek o regularizaciji Tihonova in ga dokaži.

8.2 Nesimetrični problem lastnih vrednosti

Naj bosta $A, E \in \mathbb{R}^{n \times n}$ matriki in $\varepsilon \in \mathbb{R}$. Označimo lastne vrednosti matrike $A + \varepsilon E$ z $\lambda_i(\varepsilon)$.

Izrek (Bauer-Fike). Naj bo $A = XDX^{-1}$ diagonalizabilna z $D = \text{diag}(\lambda_1, \dots, \lambda_n)$ in $X = [x_1, \dots, x_n]$. Za vsak $\varepsilon > 0$ lastne vrednosti $A + \varepsilon E$ ležijo v uniji krogov

$$K_i = \{z \mid |z - \lambda_i| \leq \varepsilon \|E\| \|X\| \|X^{-1}\|\},$$

kjer je $\|\cdot\|$ 1-, ∞ - ali 2-norma. Če unija razpade na več povezanih komponent, vsaka vsebuje toliko lastnih vrednosti, kolikor krogov jo sestavlja.

Dokaz. Naj bo $\lambda(\varepsilon)$ lastna vrednost $A + \varepsilon E$. Matrika $A + \varepsilon E - \lambda(\varepsilon)I$ je singularna, velja

$$A + \varepsilon E - \lambda(\varepsilon)I = X(D + \varepsilon X^{-1}XEX - \lambda(\varepsilon)I)X^{-1}.$$

Če je $\lambda(\varepsilon) = \lambda_i$ za nek i , očitno leži v uniji. Sicer je $D - \lambda(\varepsilon)I$ nesingularna, ker pa je srednji člen produkta zgoraj singularen, je singularna tudi

$$I + \varepsilon(D - \lambda(\varepsilon)I)^{-1}X^{-1}EX.$$

Torej je

$$1 \leq \|\varepsilon(D - \lambda(\varepsilon)I)^{-1}X^{-1}EX\| \leq \varepsilon \|(D - \lambda(\varepsilon)I)^{-1}\| \|X^{-1}\| \|X\| \|E\|.$$

Elementi matrike $(D - \lambda(\varepsilon)I)^{-1}$ so oblike $(\lambda_i - \lambda(\varepsilon))^{-1}$, torej je res

$$\min_i |\lambda_i - \lambda(\varepsilon)| \leq \varepsilon \|X^{-1}\| \|E\| \|X\|.$$

□

Vprašanje 7. Povej in dokaži Bauer-Fikeov izrek o občutljivosti problema lastnih vrednosti.

Posledica. Če je A simetrična, potem za vsako lastno vrednost $\tilde{\lambda}$ matrike $A + \varepsilon E$ obstaja lastna vrednost λ_i matrike A , da je $|\tilde{\lambda} - \lambda_i| \leq \varepsilon \|E\|_2$.

Dokaz. Matriko lahko diagonaliziramo z ortogonalno bazo, torej sta normi $\|X\|_2 = \|X^{-1}\|_2 = 1$. □

Vprašanje 8. Kaj pravi Bauer-Fikeov izrek o simetričnih matrikah?

Naj bo $Ax = \lambda x$ in $(A + \Delta A)(x + \Delta x) = (\lambda + \Delta \lambda)(x + \Delta x)$. Recimo, da je λ enostavna lastna vrednost in y levi lastni vektor, $y^H A = \lambda y^H$. Potem je

$$Ax + \Delta Ax + A\Delta x + \Delta A\Delta x = \lambda x + \Delta \lambda x + \lambda \Delta x + \Delta \lambda \Delta x.$$

Zadnji člen na vsaki strani zanemarimo in množimo z leve z y^H , da dobimo

$$y^H \Delta Ax + y^H A\Delta x = \Delta \lambda y^H x + \lambda y^H \Delta x,$$

oziroma

$$\Delta \lambda = \frac{y^H \Delta Ax}{y^H x}.$$

Vprašanje 9. Izpelji predpis za motnjo v lastni vrednosti ob dani motnji matrike ΔA .

Izrek. Naj bo λ_i enostavna lastna vrednost A z normiranimi desnim in levim lastnim vektorjem x_i in y_i . Če je $\lambda_i(\varepsilon)$ ustrezna lastna vrednost $A + \varepsilon E$, potem velja

$$\lambda_i(\varepsilon) = \lambda_i + \varepsilon \frac{y_i^H E x_i}{y_i^H x_i} + O(\varepsilon^2).$$

Definiramo

$$s_i = \frac{y_i^H x_i}{\|y_i\|_2 \|x_i\|_2}.$$

Potem je občutljivost lastne vrednosti enaka $1/|s_i|$. Ocenimo lahko $|\lambda_i(\varepsilon) - \lambda_i| \leq \varepsilon \frac{\|E\|_2}{|s_i|} + O(\varepsilon^2)$. Do maksimalne spremembe pride pri $E = y_i x_i^H$; tedaj je $|y_i^H E x_i| = 1$.

Vprašanje 10. Kaj je občutljivost lastne vrednosti?

Če je λ p -kratna lastna vrednost in ima v Jordanovi formi kletke velikosti m_1, \dots, m_k za $m_1 + \dots + m_k = p$, lahko pričakujemo $|\lambda_i(\varepsilon) - \lambda_i| = O(\varepsilon^{1/m_1})$ za ureditev $m_1 \geq m_2 \geq \dots \geq m_k$.

Izrek. Naj bo $A = XDX^{-1}$ diagonalizabilna ter $Y^H A = DY^H$, kjer je $X = [x_1, \dots, x_n]$, $Y = [y_1, \dots, y_n]$ in $\|x_i\|_2 = \|y_i\|_2 = 1$. Če je λ_i enostavna lastna vrednost, potem za dovolj majhen ε za lastni vektor $x_i(\varepsilon)$ za $\lambda_i(\varepsilon)$ velja

$$x_i(\varepsilon) = x_i + \sum_{j \neq i} \varepsilon \frac{y_j^H E x_i}{(\lambda_i - \lambda_j) s_j} x_j + O(\varepsilon^2).$$

Izrek. Če je λ_i enostavna lastna vrednost A , je $s_i \neq 0$.

Dokaz. Recimo, da je $y_i^H x_i = 0$. Naj bo U taka unitarna matrika, da je $Ue_1 = x_i$. Oglejmo si

$$B = U^H A U = \begin{bmatrix} \lambda_i & \cdots \\ 0 & C \end{bmatrix}.$$

Velja $Be_1 = \lambda_i e_1$, levi lastni vektor je oblike $z^H = y_i^H U$. Vemo $y_i^H x_i = 0$, torej je $y_i^H U U^H x_i = 0$, ker pa je $U^H x_i = e_1$, velja $y_i^H U = [0 w^H]$. Torej $[0 w^H] B = \lambda_i [0 w^H]$, oziroma $w^H C = \lambda_i w^H$. Torej je λ_i vsaj dvojna lastna vrednost A . \square

Vprašanje 11. Pokaži, da je občutljivost enostavne lastne vrednosti dobro definirana.

8.2.1 Implicitna QR iteracija

Izrek (Izrek o implicitnem Q). Naj bo $Q = [q_1, \dots, q_n]$ taka ortogonalna matrika, da je $H = Q^T A Q$ nerazcepna zgornje Hessenbergova matrika. Potem so stolpci q_2, \dots, q_n do predznaka natančno določeni s q_1 .

Dokaz. Naj bo tudi $V^T A V = G$, kjer je G nerazcepna zgornje Hessenbergova, V ortogonalna in $v_1 = q_1$. Velja $AV = VG$ in $Q^T A = H Q^T$. Če prvo enačbo množimo z leve s Q^T in drugo z desne z V , dobimo $Q^T V G = Q^T A V = H Q^T V$, za $W = Q^T V$ torej velja $WG = HW$. Matrika W je ortogonalna z $w_1 = e_1$, ker je $q_1 = v_1$. Poglejmo sedaj i -ti stolpec matrike $WG = HW$. Velja

$$\sum_{j=1}^{i+1} g_{ji} w_j = H w_i,$$

torej

$$g_{i+1,i}w_{i+1} = Hw_i - \sum_{j=1}^i g_{ji}w_j.$$

Pokažimo, da ima w_i neničelnih le prvih i elementov. Indukcija na i ;

$$w_{i+1} = \frac{1}{g_{i+1,i}}Hw_i - \frac{1}{g_{i+1,i}}\sum_{j=1}^i g_{ji}w_j.$$

Matrika H je zgornje Hessenbergova, drug člen razlike pa po indukcijski predpostavki vsebuje le elemente na mestih 1 do i . Torej je W zgornje trikotna; ker je tudi ortogonalna, mora biti diagonalna z ± 1 na diagonalni. \square

Vprašanje 12. Povej in dokaži izrek o implicitnem Q.

Obravnavajmo QR iteracijo z enojnim premikom. Če začnemo z zgornje Hessenbergovo matriko A_0 , v vsakem koraku izračunamo QR razcep $A_k - \sigma_k I = Q_k R_k$, in nato določimo $A_{k+1} = R_k Q_k + \sigma_k I$. Pišemo lahko tudi $A_{k+1} = Q_k^T A_k Q_k$. Ker smo začeli z zgornje Hessenbergovo matriko, naredimo QR razcep z Givensovimi rotacijami. Prvo z leve množimo z rotacijo R_{12}^T , ki nam spremeni prvi vrstici matrike; pri množenju z desne z R_{12} pa dobimo nov element na mestu $(3, 1)$. V nadaljevanju bomo množili z novimi rotacijami tako, da ne bomo spremenili prvega stolpca v prehodni matriki na desni, nov element pa bomo počasi premikali navzdol izven matrike.

Vprašanje 13. Kako uporabiš izrek o implicitnem Q pri QR iteraciji z enojnim premikom?

Za dvojni premik velja

$$N_k = A_k^2 - (\sigma_{k1} + \sigma_{k2})A_k + \sigma_{k1}\sigma_{k2}I = Q_k R_k$$

in $A_{k+1} = Q_k^T A_k Q_k$. Prvi stolpec matrike N_k lahko izračunamo v konstantno mnogo operacijah. Iz tega določimo Householderjevo zrcaljenje P_1 , ki stolpec zrcali v e_1 . Tukaj dobimo nove elemente tako pri množenju z leve kot pri množenju z desne, tako da premikamo grbo velikosti 2×2 .

Vprašanje 14. Kako uporabiš izrek o implicitnem Q pri QR iteraciji z dvojnimi premiki?

8.3 Simetrični problem lastnih vrednosti

Matrika $A = A^T$ ima same realne lastne vrednosti, torej jih lahko uredimo kot $\lambda_n \leq \dots \leq \lambda_1$. Lastne vektorje lahko izberemo tako, da so ortonormirani, $Ax_i = \lambda_i x_i$, $x_i^T x_j = \delta_{ij}$.

Lema. Za Rayleighov kvocient velja $\lambda_n \leq \rho(x, A) \leq \lambda_1$.

Dokaz. Zapišemo x v bazi lastnih vektorjev, in dobimo

$$\rho(x, A) = \frac{\sum_i \alpha_i^2 \lambda_i}{\sum_i \alpha_i^2}.$$

□

Izrek (Courant-Fisher). *Za simetrično matriko A velja*

$$\lambda_i = \min_{\dim S = n-i+1} \max_{x \in S} \rho(x, A) = \max_{\dim R = i} \min_{x \in R} \rho(x, A).$$

Dokaz. Označimo prvi izraz z L in drugega z D . Za poljuben par podprostorov S in R je $\dim S + \dim R = n + 1 > n$, torej obstaja neničelni vektor x_{RS} v preseku. Naj bo na levi strani minimum dosežen pri \hat{S} , na desni pa pri \hat{R} . Potem je

$$\max_{x \in \hat{S}} \rho(x, A) \geq \rho(x_{\hat{R}\hat{S}}, A) \geq \min_{x \in \hat{R}} \rho(x, A),$$

iz česar sledi $L \geq D$. Če pa izberemo $\tilde{S} = \text{Lin}(x_i, x_{i+1}, \dots, x_n)$ in $\tilde{R} = \text{Lin}(x_1, \dots, x_i)$, pa je očitno

$$L \leq \max_{x \in \tilde{S}} \rho(x, A) = \lambda_i = \min_{x \in \tilde{R}} \rho(x, A) \leq D,$$

torej $L = D = \lambda_i$. □

Vprašanje 15. Povej in dokaži Courant-Fisherjev izrek.

Posledica. *Če sta A in E simetrični matriki, potem za lastne vrednosti $\tilde{\lambda}_n \leq \dots \leq \tilde{\lambda}_1$ matrike $A + E$ velja*

$$\lambda_i + \lambda_n(E) \leq \tilde{\lambda}_i \leq \lambda_i + \lambda_1(E).$$

Dokaz. Uporabimo Courant-Fisherjev izrek za $\tilde{\lambda}_i$. Velja

$$\tilde{\lambda}_i = \min_{\dim S = n-i+1} \max_{x \in S} \rho(x, A + E),$$

upoštevamo $\rho(x, A + E) = \rho(x, A) + \rho(x, E)$ in ocenimo $\lambda_n(E) \leq \rho(x, E) \leq \lambda_1(E)$. □

Posledica (Weylov izrek). *Če so $\lambda_n \leq \dots \leq \lambda_1$ lastne vrednosti $A = A^T$ in $\tilde{\lambda}_n \leq \dots \leq \tilde{\lambda}_1$ lastne vrednosti $A + E$ za simetrično E , je $|\tilde{\lambda}_i - \lambda_i| \leq \|E\|_2$.*

Dokaz. Velja $\|E\|_2 = \max_i |\lambda_i(E)|$. □

Vprašanje 16. Povej in dokaži Weylov izrek.

Algorithm 14 Rayleighova iteracija

```

Izberi  $z_0 \neq 0$ 
for  $k = 0, 1, 2, \dots$  do
     $\sigma_k = \rho(z_k, A)$ 
    Reši  $(A - \sigma_k I)y_{k+1} = z_k$ 
     $z_{k+1} = y_{k+1} / \|y_{k+1}\|_2$ 
end for

```

8.3.1 Rayleighova iteracija

Lema. Naj bo $\|z_0\|_2 = 1$ približek za lastni vektor x simetrične matrike A z lastnimi vrednostmi $|\lambda_n| \leq \dots \leq |\lambda_2| < |\lambda_1|$. Če izvedemo en korak potenčne metode za A z začetnim vektorjem z_0 , potem za z_1 velja

$$\|z_i \pm x_1\| \leq \frac{|\lambda_2|}{|\lambda_1|} \|z_0 - x_1\|_2 + O(\|z_0 - x_1\|_2^2),$$

kjer vzamemo tisto možnost v $z_1 \pm x_1$, ki da manjšo normo.

Dokaz. Predpostavimo, da je z_0 dober približek za x_1 , torej

$$z_0 = x_1 + \sum_{i=2}^n \alpha_i x_i$$

za $|\alpha_i| \ll 1$. Potem je

$$z_1 \approx x_1 + \sum_{i=2}^n \alpha_i \frac{\lambda_i}{\lambda_1} x_i$$

in

$$\|z_1 - x_1\|_2^2 = \sum_{i=2}^n \alpha_i^2 \frac{\lambda_i^2}{\lambda_1^2} \leq \frac{\lambda_2^2}{\lambda_1^2} \sum_{i=2}^n \alpha_i^2 = \frac{\lambda_2^2}{\lambda_1^2} \|z_0 - x_1\|_2^2.$$

□

Vprašanje 17. Ocenite napako po enem koraku potenčne metode.

Posledica. Naj bo A simetrična in $|\sigma - \lambda_i| \ll |\sigma - \lambda_j|$ za $j \neq i$. Če je z_0 dober približek za lastni vektor x_i in naredimo en korak inverzne iteracije, je

$$\|z_1 \pm x_i\| = O(|\sigma - \lambda_i|) \|z_0 - x_i\|_2$$

Lema. Naj bo A simetrična in z približek za lastni vektor x_i za λ_i . Potem je $|\rho(z, A) - \lambda_i| \leq 2 \|A\|_2 \cdot \|z - x_i\|_2^2$.

Dokaz. Brez škode za splošnost je $i = 1$. Naj bo $z = \sum_i \alpha_i x_i$ in $\sum_i \alpha_i^2 = 1$. Potem je

$$\|z - x_1\|_2^2 = (1 - \alpha_1)^2 + \alpha_2^2 + \cdots + \alpha_n^2 = 2 - 2\alpha_1.$$

Velja $\rho(z, A) = \sum_i \lambda_i \alpha_i^2$, torej

$$\rho(z, A) - \lambda_1 = \sum_{i=1}^n \lambda_i \alpha_i^2 - \lambda_1 \sum_{i=1}^n \alpha_i^2 = \sum_{i=2}^n (\lambda_i - \lambda_1) \alpha_i^2.$$

Sledi

$$|\rho(z, A) - \lambda_1| \leq \sum_{i=2}^n |\lambda_i - \lambda_1| \alpha_i^2 \leq \sum_i (|\lambda_i| + |\lambda_1|) \alpha_i^2 \leq \sum_i 2 \|A\|_2 \alpha_i^2 \leq 2 \|A\|_2 (1 - \alpha_1^2)$$

Predpostavimo, da je $\alpha_i \approx 1$, torej lahko to ocenimo še z

$$2 \|A\|_2 \cdot 2(1 - \alpha_1) = 2 \|A\|_2 \|z - x_1\|_2^2.$$

□

Vprašanje 18. Ocení napako, ki jo dobimo z Rayleighovim kvocientom.

Izrek. *Rayleighova iteracija ima v bližini lastne vrednosti simetrične matrike kubično konvergenco.*

Dokaz. Naj bo z_k blizu x_i . Vemo, da je

$$\|z_{k+1} \pm x_i\| = O(|\sigma_k - \lambda_i| \|z_k - x_i\|_2) = O(\|z_k - x_i\|_2^2 \cdot \|z_k - x_i\|_2).$$

□

Vprašanje 19. Kakšen red konvergence ima Rayleighova iteracija za simetrične matrike? Dokaži.

8.3.2 QR iteracija

Če je A simetrična, začetna redukcija na zgornje Hessenbergovo obliko vrne tridiagonalno simetrično matriko. Tako A reduciramo na

$$T = \begin{bmatrix} a_1 & b_1 & & & \\ b_1 & a_2 & b_2 & & \\ & b_2 & \ddots & \ddots & \\ & & \ddots & a_{n-1} & b_{n-1} \\ & & & b_{n-1} & a_n \end{bmatrix}.$$

Algoritem bo še vedno porabil $O(n^3)$ operacij, ampak ga lahko zaradi simetrije malo pohitrimo. Predpostavimo lahko, da je T nerazcepna. En korak QR iteracije lahko

izvedemo v $O(n)$ operacijah, če pa želimo posodabljati še produkt matrik Q_k , je to dodatnih $O(n^2)$ operacij.

Imamo dva načina za izvajanje premika. Prvi je Rayleighov premik $\sigma_k = \rho(e_n, T_k) = a_n^{(k)}$.

Izrek. Če izvajamo QR iteracijo za tridiagonalno simetrično matriko T , potem se Rayleighovi premiki σ_k ujemamo z Rayleighovimi kvocienti, ki bi jih dobili pri Rayleighovi iteraciji z $z_0 = e_n$.

Vprašanje 20. Kako deluje QR iteracija z Rayleighovimi premiki?

Druga možnost so Wilkinsonovi premiki, kjer pogledamo

$$W_k = \begin{bmatrix} a_{n-1}^{(k)} & b_{n-1}^{(k)} \\ b_{n-1}^{(k)} & a_n^{(k)} \end{bmatrix}$$

in izberemo tisto lastno vrednost, ki je bližja $a_n^{(k)}$. Tukaj imamo celo globalno konvergenco, ki je v bližini lastne vrednosti kubična.

Vprašanje 21. Kaj so Wilkinsonovi premiki?

8.3.3 Bisekcija

Izrek. Če je T nerazcepna simetrična tridiagonalna matrika, ima vse lastne vrednosti enostavne.

Dokaz. Naj bo λ lastna vrednost T . Ker je matrika simetrična, sta geometrijska in algebraična večkratnost enaki. Ker so $b_i \neq 0$, je prvih $n-1$ stolpcev $T - \lambda I$ neodvisnih, torej je rang te matrike enak $n-1$ in je večkratnost λ enaka 1. \square

Vprašanje 22. Dokaži: nerazcepna simetrična tridiagonalna matrika ima vse lastne vrednosti enostavne.

Definicija. Matriki A in B sta KONGRUENTNI, če obstaja nesingularna X , da je $B = X^T A X$.

Če je A simetrična, je $\varphi(x) = x^T A x$ kvadratna forma. Potem za $x = Zy$ dobimo $x^T A x = y^T Z^T A Z y$ v drugi bazi.

Definicija. Za vsako $n \times n$ simetrično matriko A definiramo INERCIJO A kot trojico (ν, z, p) , kjer je ν število negativnih, z število ničelnih in p število pozitivnih lastnih vrednosti.

Izrek (Sylvester). Kongruentni simetrični matriki imata enako inercijo.

Dokaz. Naj bo (ν, z, p) inercija A in (ν', z', p') inercija $X^T A X$. Recimo $\nu' < \nu$. Vsi lastni vektorji matrike A za negativne lastne vrednosti razpenjajo nek invarianten podprostor

\mathcal{N} dimenzije ν . Podobno vsi lastni vektorji nenegativnih lastnih vrednosti X^TAX razpenjajo invarianten podprostor \mathcal{P} dimenzije $n - \nu'$ za X^TAX . Prostor $X\mathcal{P}$ je tudi podprostor dimenzije $n - \nu'$, ki se mora sekati z \mathcal{N} po predpostavki $\nu' < \nu$. Naj bo $v \in \mathcal{N} \cap X\mathcal{P}$ neničeln. Velja

$$\rho(v, A) = \frac{v^T Av}{v^T v} < 0,$$

po drugi strani pa obstaja $w \in \mathcal{P}$, da je $v = Xw$. Velja

$$\rho(v, A) = \frac{w^T X^T AX w}{w^T X^T X w} \geq 0,$$

kar je protislovje. Podobno v primeru $p' < p$. □

Vprašanje 23. Povej in dokaži Sylvestrov izrek.

Kako izračunamo inercijo za $T - \lambda I$? Uporabimo LDL^T razcep, z

$$L = \begin{bmatrix} 1 & & & & \\ l_1 & 1 & & & \\ & l_2 & 1 & & \\ & & \ddots & \ddots & \\ & & & l_{n-1} & 1 \end{bmatrix}$$

in $D = \text{diag}(d_1, \dots, d_n)$. Razcep poiščemo podobno kot LU razcep. Vemo, da je inercija $T - \lambda I$ enaka inerciji D . Lastne vrednosti poiščemo z bisekcijo, kjer upoštevamo informacije iz inercije.

V postopku za račun LDL^T razcepa smo nekje delili z diagonalnim elementom, kar lahko načeloma da vmesni rezultat $d_i = -\infty$. To nas ne moti, tako vrednost štejemo kot negativno, v naslednjem koraku pa bomo tako in tako dobili 0.

Vprašanje 24. Opiši metodo bisekcije za iskanje lastnih vrednosti.

8.3.4 Jacobijeva metoda

V Jacobijevi metodi matrike ne reduciramo na tridiagonalno obliko. Vzamemo Jacobijevo (Givensovo) rotacijo

$$R_{pq}^T = \begin{bmatrix} 1 & & & & & \\ & \ddots & & & & \\ & & 1 & & & \\ & & & c & & s \\ & & & & \ddots & \\ & & & -s & & c \\ & & & & & \ddots \\ & & & & & & 1 \end{bmatrix},$$

kjer izberemo c in s tako, da bo $\tilde{a}_{pq} = \tilde{a}_{qp} = 0$. S tem si lahko pokvarimo prejšnja postavljanja.

Lema. Če \tilde{A} dobimo iz A kot $\tilde{A} = R_{pq}^T A R_{pq}$, kjer R_{pq} določimo tako, da je $\tilde{a}_{pq} = \tilde{a}_{qp} = 0$, potem $\text{off}(\tilde{A})^2 = \text{off}(A)^2 - 2a_{pq}^2$, kjer je off norma izvendiagonalnega dela.

Dokaz. Ker je R_{pq} ortogonalna, velja

$$\text{off}(\tilde{A})^2 + \sum_i \tilde{a}_{ii}^2 = \|\tilde{A}\|_F^2 = \|A\|_F^2 = \text{off}(A)^2 + \sum_i a_{ii}^2.$$

Vsi diagonalni elementi razen pp in qq se ujemajo, torej se to poenostavi v

$$\text{off}(\tilde{A})^2 = \text{off}(A)^2 + a_{pp}^2 + a_{qq}^2 - \tilde{a}_{pp}^2 - \tilde{a}_{qq}^2.$$

Tudi 2×2 podmatriki mest (p, q) imata enako Frobeniusovo normo, torej je ostanek enak $-2a_{pq}^2$. \square

Vprašanje 25. Dokaži, da se pri Jacobijevi iteraciji norma izvendiagonalnih elementov zmanjšuje z iteracijami.

S preprostim izračunom lahko pridemo do naslednjega predpisa za s in c :

$$\begin{aligned} \tau &= \frac{a_{pp} - a_{qq}}{2a_{pq}}, \\ t &= \frac{\text{sgn } \tau}{|\tau| + \sqrt{1 + \tau^2}}, \\ c &= \frac{1}{\sqrt{1 + t^2}}, \\ s &= ct. \end{aligned}$$

V enem koraku metode potem posodobimo $A = R_{pq}^T A R_{pq}$ in $Q = Q R_{pq}$.

Vprašanje 26. Opiši en korak Jacobijeve iteracije.

Jacobijeva metoda je običajno počasnejša od ostalih, lahko pa z njo natančneje izračunamo majhne lastne vrednosti spd matrik. Metoda ima več variant:

- **KLASIČNA VARIANTA:** v vsakem koraku uničimo po absolutni vrednosti največji izvendiagonalni element. Preverjanje, kateri element je največji, lahko optimiziramo iz $O(n^2)$ v $O(n)$, če si shranjujemo največji element po vrsticah in stolpcih.
- **CIKLIČNA VARIANTA:** elemente uničujemo po vrsti, na koncu se vrnemo na začetek.
- **PRAGOVNA VARIANTA:** podobno kot ciklična metoda, a uničimo le dovolj velike elemente.

Vprašanje 27. Opiši tri variante Jacobijeve metode.

8.3.5 Deli in vladaj

Naj bo T nerazcepna simetrična tridiagonalna $n \times n$ matrika in $m \approx \frac{n}{2}$. Razdelimo T na dve matriki;

$$T = \begin{bmatrix} T_1 & \\ & T_2 \end{bmatrix} + b_m v v^T,$$

kjer je $v = e_m + e_{m+1}$, matriki T_1 in T_2 pa sta nerazcepni tridiagonalni. Rekurzivno lahko poiščemo $T_1 = Q_1 D_1 Q_1^T$ in $T_2 = Q_2 D_2 Q_2^T$. Potem velja

$$T = \begin{bmatrix} Q_1 & \\ & Q_2 \end{bmatrix} \left(\begin{bmatrix} D_1 & \\ & D_2 \end{bmatrix} + b_m u u^T \right) \begin{bmatrix} Q_1^T & \\ & Q_2^T \end{bmatrix}$$

za

$$u = \begin{bmatrix} Q_1^T & \\ & Q_2^T \end{bmatrix} v.$$

Recimo, da znamo učinkovito rešiti problem lastnih vrednosti za matriko $D + \rho u u^T = \tilde{Q} \Lambda \tilde{Q}^T$. Potem je

$$T = \begin{bmatrix} Q_1 & \\ & Q_2 \end{bmatrix} \tilde{Q} \Lambda \tilde{Q}^T \begin{bmatrix} Q_1^T & \\ & Q_2^T \end{bmatrix}.$$

Izrek. Naj bo $A = D + \rho u u^T$ za $D = \text{diag}(d_1, \dots, d_n)$ in $d_1 \geq d_2 \geq \dots \geq d_n$.

- Če je $u_i = 0$, je d_i lastna vrednost A z lastnim vektorjem e_i .
- Če je $d_i = d_{i+1}$, potem je d_i lastna vrednost A za lastni vektor $-u_{i+1}e_i + u_i e_{i+1}$.

Vprašanje 28. Opiši osnovno delovanje metode deli in vladaj. Kakšni so posebni primeri za lastne vrednosti?

Preostane problem lastnih vrednosti za $A = D + \rho u u^T$, kjer je $d_1 > \dots > d_n$ in $u_i \neq 0$ za vsak i . Naj bo λ lastna vrednost A . Predpostavimo, da je $\lambda \neq d_i$. Vemo, da je $A - \lambda I$ singularna. Ker je

$$A - \lambda I = D - \lambda I + \rho u u^T = (D - \lambda I)(I + \rho(D - \lambda I)^{-1} u u^T)$$

in ker je prva matrika v produktu nesingularna, je druga singularna. Za $x = (D - \lambda I)^{-1} u$ in $y^T = u^T$ velja naslednja lema:

Lema. $\det(I + x y^T) = 1 + y^T x$.

Dokaz. Če je $z \perp y$, je $(I + x y^T)z = z$, takih je $n - 1$ linearno neodvisnih z , torej je 1 $(n - 1)$ -kratna lastna vrednost. Velja $(I + x y^T)x = (1 + y^T x)x$, torej je $1 + y^T x$ zadnja lastna vrednost. \square

Računamo

$$\det(I + \rho(D - \lambda I)^{-1} u u^T) = 1 + \rho u^T (D - \lambda I)^{-1} u = 1 + \rho \sum_{i=1}^n \frac{u_i^2}{d_i - \lambda} = f(\lambda),$$

čemu pravimo SEKULARNA FUNKCIJA.

Vprašanje 29. Izpelji sekularno funkcijo.

Lastne vrednosti se rešitve sekularne enačbe $f(\lambda) = 0$. Funkcija ima n enostavnih polov, za $\rho > 0$ je naraščajoča in za $\rho < 0$ padajoča z asimptoto $y = 1$. Torej ima natanko n ničel na znanih intervalih (d_i, d_{i-1}) oziroma na enem od neskončnih intervalov. Če je α lastna vrednost, lahko izračunamo lastni vektor $(D - \alpha I)^{-1}u$.

Krivulja je strma okoli ničle in položna na velikem delu intervala, kar je težava, ker tako tangentna metoda kot bisekcija ne delujeta dobro. Namesto tega f aproksimiramo z

$$h(x) = \frac{c_1}{d_{i+1} - \lambda} + \frac{c_2}{d_i - \lambda} + c_3,$$

kjer konstante izračunamo tako, da se vrednosti h in f ter njunih odvodov ujemajo v približku x_0 .

Vprašanje 30. Kako poiščeš lastne vrednosti in lastne vektorje matrike $D + \rho uu^T$?

Metoda je hitrejša od QR iteracije za $n > 30$, če potrebujemo tudi lastne vektorje. Ne deluje pa brez računanja Q .

8.4 Računanje singularnega razcepa

Dano imamo matriko $A \in \mathbb{R}^{m \times n}$, iščemo $A = U\Sigma V^T$. Vemo, da je $A^T A = VDV^T$, kjer je $V = \sigma_1^2, \dots, \sigma_n^2$. To lahko uporabimo za numerični postopek:

- izračunaj $B = A^T A$
- poišči $B = VDV^T$ z rešitvijo simetričnega problema lastnih vrednosti
- izračunaj Σ iz D
- izračunaj vektorje $u_i = \frac{1}{\sigma_i} A v_i$ za $i = 1, \dots, \text{rang}(A)$
- dopolni U do ortogonalne matrike

Problem je, da nočemo eksplicitno računati B .

8.4.1 Enostranska Jacobijeva metoda

Implicitno izvajamo Jacobijevo metodo na $B = A^T A$. V iteraciji predpišemo $\tilde{B} = R_{pq}^T B R_{pq}$, namesto tega pa množimo $\tilde{A} = A R_{pq}$. Za določitev R_{pq} potrebujemo le elemente $[A^T A]_{pp}$, $[A^T A]_{pq}$ in $[A^T A]_{qq}$. Potem \tilde{A} konvergira proti matriki z ortogonalnimi stolpci. Norme stolpcev so singularne vrednosti, produkt rotacij R_{pq} pa nam da matriko V . Najbolje je uporabljati pragovno varianto Jacobijeve iteracije; element B_{pq} uničimo, če je dovolj velik v primerjavi z B_{pp} in B_{qq} . Jacobijeva metoda lahko za določene tipe matrik majhne singularne vrednosti izračunati natančneje od ostalih.

Vprašanje 31. Opiši enostransko Jacobijevo metodo.

8.4.2 Enostranska QR iteracija

Vsako matriko $A \in \mathbb{R}^{m \times n}$ lahko reduciramo na $A = U_1 B V_1^T$, kjer je B bidiagonalna in U_1, V_1 ortogonalni. Če potem izračunamo $B = U_2 \Sigma V_2^T$, bo $A = U_1 U_2 \Sigma (V_1 V_2)^T$. Redukcijo izvedemo z uporabo Householderjevih zrcaljenj z leve in z desne. Zahtevnost te redukcije je $8mn^2 - \frac{8}{3}n^3$.

Če je B bidiagonalna, lahko izračunamo $B^T B$ in delamo QR iteracijo z Wilkinsonovimi premiki. Matrika $B^T B$ je simetrična nenegativno definitna. Prav tako je spodnja desna 2×2 podmatrika simetrična nenegativno definitna, torej so premiki oblike σ^2 . Matrika $C = B^T B$ je simetrična tridiagonalna, en korak QR iteracije s premikom σ^2 nam da $\tilde{C} = Q^T C Q$, kjer je Q iz QR razcepa $C - \sigma^2 I$. Če nastavimo $\tilde{B} = P B \hat{Q}$, kjer je \hat{Q} ortogonalna s prvim stolpcem, enakim normiranemu prvemu stolpcu $C - \sigma^2 I$, P pa je taka ortogonalna, da je \tilde{B} spet bidiagonalna. Potem je $\tilde{C} = \tilde{B}^T \tilde{B} = \hat{Q}^T C \hat{Q}$ tridiagonalna, torej zgornja Hessenbergova. Po izreku o implicitnem Q je $\hat{Q} = Q \text{diag}(1, \pm 1, \dots, \pm 1)$. To je en korak enostranske QR iteracije s premikom σ^2 .

Matriko izjemoma množimo z rotacijami z leve in z desne, s čimer premikamo grbo skozi diagonalo. Matriki P in \hat{Q} dobimo kot produkt teh rotacij. En cikel skozi vse elemente na diagonalni vzame $30n$ iteracij.

Vprašanje 32. Opiši enostransko QR iteracijo.

8.4.3 Weylov izrek

Lema. Naj bo $B \in \mathbb{R}^{m \times n}$ za $m \geq n$. Potem ima matrika

$$C = \begin{bmatrix} 0 & B^T \\ B & 0 \end{bmatrix}$$

lastne vrednosti $\pm \sigma_i$ z lastnimi vektorji $(v_i, \pm u_i)$, in $m - n$ lastnih vrednosti 0 z lastnimi vektorji $(0, u_j)$ za $j = n + 1, \dots, m$. Pri tem je $B = U \Sigma V^T$.

Dokaz je račun. Če je

$$B = \begin{bmatrix} a_1 & b_1 & & & \\ & a_2 & b_2 & & \\ & & \ddots & \ddots & \\ & & & a_{n-1} & b_{n-1} \\ & & & & a_n \end{bmatrix}$$

in $P = [e_1, e_{n+1}, e_2, e_{n+2}, \dots, e_n, e_{2n}]$, je za zgornji C potem

$$P^T C P = \begin{bmatrix} 0 & a_1 & & & \\ a_1 & 0 & b_1 & & \\ & b_1 & 0 & \ddots & \\ & & \ddots & \ddots & a_n \\ & & & a_n & 0 \end{bmatrix}$$

Privzamemo, da je ta matrika nerazcepna (torej $a_i, b_i \neq 0$), sicer lahko problem razdelimo na dva manjša.

Izrek (Weylov izrek za singularne vrednosti). *Naj bosta $A, E \in \mathbb{R}^{m \times n}$ za $m \geq n$, $\sigma_1 \geq \dots \geq \sigma_n$ singularne vrednosti A in $\tilde{\sigma}_1 \geq \dots \geq \tilde{\sigma}_n$ singularne vrednosti $A + E$. Potem velja $|\tilde{\sigma}_i - \sigma_i| \leq \|E\|_2$.*

Dokaz. Naj bo

$$C = \begin{bmatrix} 0 & A^T \\ A & 0 \end{bmatrix} \quad F = \begin{bmatrix} 0 & E^T \\ E & 0 \end{bmatrix}$$

Vemo, da je $|\lambda_i(C) - \lambda_i(C + F)| \leq \|F\|_2 = \|E\|_2$. □

Vprašanje 33. Povej in dokaži Weylov izrek.

Singularne vrednosti A lahko izračunamo na dva načina:

- Izračunamo $B = A^T A$ in na B uporabimo obratno stabilno metodo za simetrični problem lastnih vrednosti. Te potem korenimo.
- A damo v obratno stabilno metodo za računanje singularnega razcepa.

V drugem načinu so numerično izračunane singularne vrednosti točne singularne vrednosti matrike $A + \Delta A$, kjer je $\|\Delta A\|_2 = O(\|A\|_2 u)$. Če gledamo prvi način in ignoriramo korenjenje, so numerično izračunane lastne vrednosti točne lastne vrednosti $B + \Delta B$ za $\|\Delta B\|_2 = O(\|B\|_2 u)$. Sledi $|\hat{\sigma}_i - \sigma_i| = O(\sigma_1 u)$. Po Weylovem izreku za simetrične matrike je $|\hat{\sigma}_i^2 - \sigma_i^2| = O(\|B\|_2 u) = O(\sigma_1^2 u)$. Ocenimo lahko $\hat{\sigma}_i \approx \sigma_i$, iz česar potem sledi

$$|\hat{\sigma}_i - \sigma_i| = O\left(\frac{\sigma_1^2 u}{\sigma_i}\right).$$

Oba načina dobro izračunata velike singularne vrednosti, za majhne pa je prvi način boljši.

Vprašanje 34. Kako stabilno izračunaš majhne singularne vrednosti?

8.5 Posplošitve problema lastnih vrednosti

8.5.1 Posplošen problem lastnih vrednosti

Imamo matriki $A, B \in \mathbb{C}^{n \times n}$. Množico vseh matrik $A - \lambda B$ za $\lambda \in \mathbb{C}$ imenujemo MATRIČNI ŠOP. Včasih ga označimo z (A, B) . Definiramo lahko karakteristični polinom $p(\lambda) = \det(A - \lambda B)$. Če je $p = 0$, pravimo, da je šop SINGULAREN, sicer pa je REGULAREN.

Naj bo (A, B) regularen šop. Če je $Ax = \lambda Bx$ za $x \neq 0$, je λ končna lastna vrednost in x desni lastni vektor. Podobno za $y^H A = \lambda y^H B$. Če pa je $Bx = 0$ za nek $x \neq 0$, je ∞ lastna vrednost in x desni lastni vektor; podobno za $y^H B = 0$. Za regularen šop ima problem n lastnih vrednosti, tj. ničle karakterističnega polinoma, ki jim dodamo še toliko neskončnih lastnih vrednosti, da pridemo do n .

Vprašanje 35. Opiši posplošen problem lastnih vrednosti.

Izrek. Naj bo (A, B) regularen šop.

- Če je B obrnljiva matrika, so lastne vrednosti tega šopa enake lastnim vrednostim $B^{-1}A$ oziroma AB^{-1} .
- ∞ je lastna vrednost šopa natanko tedaj, ko je $\det B = 0$.
- Če je A nesingularna, potem so lastne vrednosti (A, B) recipročne lastnim vrednostim matrike $A^{-1}B$ oz. BA^{-1} . Tu je $\infty = \frac{1}{0}$.

Vprašanje 36. Kakšne so lastne vrednosti regularnega šopa (A, B) , če je A obrnljiva?

Definicija. Če sta U, V nesingularni matriki, je šop (A, B) EKVIVALENTEN šopu (UAV, UBV) .

Izrek. Ekvivalentna šopa imata enake lastne vrednosti. Če je x desni lastni vektor za (A, B) , je $V^{-1}x$ desni lastni vektor za (UAV, UBV) . Če je y levi lastni vektor za (A, B) , je $U^{-H}y$ levi lastni vektor za (UAV, UBV) .

Vprašanje 37. Kaj lahko poveš o lastnih vektorjih ekvivalentnih šopov?

Izrek. Za regularen matrični šop (A, B) obstajata taki unitarni matriki Q, Z , da je $(Q^H A Z, Q^H B Z) = (S, T)$, kjer sta S in T zgornje trikotni matriki. Lastne vrednosti (A, B) so potem

$$\lambda_i = \frac{s_{ii}}{t_{ii}}$$

za $t_{ii} \neq 0$ in $\lambda_i = \infty$ za $t_{ii} = 0$.

Dokaz. Indukcija na n . Pri $n = 1$ je trivialno, pogledjmo primer $n - 1 \rightsquigarrow n$. Ker je (A, B) regularen šop, ima vsaj eno lastno vrednost λ in lastni vektor x . Če je λ končna, je $Ax = \lambda Bx$, sicer pa je $Bx = 0$. V vsakem primeru sta vektorja Ax in Bx kolinearna, torej obstaja y z $\|y\|_2 = 1$, da je $Ax = \alpha y$ in $Bx = \beta y$, ter $(\alpha, \beta) \neq (0, 0)$.

Dopolnimo x do unitarne matrike $X = [xX_1]$ in y do unitarne $Y = [yY_1]$. Potem je

$$Y^H AX = \begin{bmatrix} \alpha & \cdots \\ 0 & A_1 \end{bmatrix}$$

in

$$Y^H BX = \begin{bmatrix} \beta & \cdots \\ 0 & B_1 \end{bmatrix}.$$

Obstajata Q_1 in Z_1 , da je $Q_1^H(A_1 - \lambda B_1)Z_1 = S_1 - \lambda T_1$; velja

$$\begin{bmatrix} 1 & \\ & Q_1^H \end{bmatrix} Y^H AX \begin{bmatrix} 1 & \\ & Z_1 \end{bmatrix} = \begin{bmatrix} \alpha & \cdots \\ & S_1 \end{bmatrix} = S$$

$$\begin{bmatrix} 1 & \\ & Q_1^H \end{bmatrix} Y^H BX \begin{bmatrix} 1 & \\ & Z_1 \end{bmatrix} = \begin{bmatrix} \beta & \cdots \\ & T_1 \end{bmatrix} = T.$$

□

Vprašanje 38. Dokaži, da za regularen matrični šop obstaja posplošena Schurova forma.

Za računanje posplošene Schurove forme imamo na voljo QZ algoritem. Predpriprava je, da z ortogonalno ekvivalenčno transformacijo pretvorimo (A, B) v (A_1, B_1) , kjer je A_1 zgornje Hessenbergova in B_1 zgornje trikotna. Iščemo torej ortogonalni Q_1, Z_1 , da velja $A = Q_1 A_1 Z_1$, $B = Q_1 B_1 Z_1$. V prvem koraku uporabimo algoritem za QR razcep in poiščemo ortogonalno U , da je $U^T B$ zgornje trikotna. Potem z rotacijami iz leve uničujemo elemente v A (oz. $U^T A$). Z vsako rotacijo si pokvarimo en element v B , ki ga popravimo z rotacijo z desne (ta pa nam ne pokvari ustvarjenih ničel v A).

Vprašanje 39. Opiši postopek predpriprave za QZ algoritem.

Če ima B na diagonali ničlo, jo lahko prestavimo na spodnje desno mesto in nadaljujemo na podmatriki $(n-1) \times (n-1)$. Predpostavimo torej, da je B nesingularna. Matrika $C_0 = A_0 B_0^{-1}$ je zgornja Hessenbergova. Na njej naredimo korak QR iteracije z dvojnimi premiki, da dobimo $C_1 = Q_0^T C_0 Q_0$, kjer je Q_0 matrika iz QR razcepa $N_0 = C_0^2 - (\sigma_1 + \sigma_2)C_0 + \sigma_1\sigma_2 I$.

QZ iteracija nam iz matrik A_0, B_0 da matriki $A_1 = Q_0^T A_0 Z_0$ in $B_1 = Q_0^T B_0 Z_0$, kjer je Q_0 ortogonalna in taka, da je prvi stolpec enak normiranemu prvemu stolpcu N_0 , ter Z_0 ortogonalna in taka, da je A_1 zgornje Hessenbergova in B_1 zgornje trikotna. Po izreku o implicitnem Q je $A_1 B_1^{-1} = Q_0^T A_0 B_0^{-1} Q_0 = C_1$. Za določitev premika potrebujemo spodnjo desno 2×2 podmatriko C_0 , torej moramo izračunati le inverz spodnje desne 3×3 podmatrike B_0 . Za prvi stolpec N_0 potrebujemo tudi inverz vodilne 2×2 podmatrike B_0 .

Prvi stolpec N_0 ima tri elemente, torej za redukcijo potrebujemo 3×3 Householderjevo zrcaljenje. Ko s tem množimo A in B , dobimo grbi; v A je grba velikosti 1×1 , v B pa

velikosti 2×2 . Potem uničimo grbo v B z dvema zrcaljenjema iz desne, kar nam poveča grbo v A na velikost 2×2 . To grbo spravimo nazaj v 1×1 z zrcaljenjem iz leve, kar nam zopet prinese grbo v B velikosti 2×2 . Ko grbo premaknemo do konca, celoten postopek ponavljamo. QZ iteracija je ekvivalentna inverzni iteraciji na matriki AB^{-1} . Je bolj numerično stabilna kot množenje z inverzom in računanje lastnih vrednosti, a počasnejša.

Vprašanje 40. Opiši QZ algoritem. Čemu je ekvivalenten?

8.5.2 Kvadratni problem lastnih vrednosti

Naj bo $Q(\lambda) = \lambda^2 M + \lambda C + K$, kjer so M , C in K matrike dimenzij $n \times n$. Potem je $\det Q(\lambda)$ polinom stopnje manjše ali enake $2n$. Če ta polinom ni konstantno enak 0, pravimo, da je problem REGULAREN.

Če je za regularen problem $Q(\lambda)x = 0$, je λ lastna vrednost in x desni lastni vektor. Podobno je za $y^H Q(\lambda) = 0$ y levi lastni vektor. V primeru $Mx = 0$ za neničeln x pa pravimo, da je ∞ lastna vrednost za x . Lastne vrednosti so ničle karakterističnega polinoma, ki jih dopolnimo do $2n$ z neskončnimi.

Klasičen način reševanja je linearizacija. Če je $\det M \neq 0$, imamo $2n$ končnih lastnih vrednosti, in poiščemo lastne vrednosti matrike

$$\begin{bmatrix} 0 & I \\ -M^{-1}K & -M^{-1}C \end{bmatrix}.$$

V primeru $\det M = 0$, pa poiščemo rešitve $(A - \lambda B)u = 0$ za

$$A = \begin{bmatrix} C & K \\ K & 0 \end{bmatrix}, \quad B = \begin{bmatrix} -M & 0 \\ 0 & K \end{bmatrix}.$$

Vprašanje 41. Opiši kvadratni problem lastnih vrednosti. Kako ga rešiš?

9 Statistika

9.1 Centralni limitni izrek

Naj bo $S_n \sim \text{Bin}(n, \frac{1}{2})$, da velja $P(S_n = k) = \binom{n}{k} 2^{-n}$. Poglejmo si razmerje

$$\frac{P(S_n = k+1)}{P(S_n = k)} = \frac{\binom{n}{k+1}}{\binom{n}{k}} = \frac{n-k}{k+1}.$$

Za sedaj se omejimo na sode $n = 2m$. Potem je $P(S_n = k)$ največja za $k = m$, in za $d > 0$ dobimo

$$\begin{aligned} \frac{P(S_{2m} = m+d)}{P(S_{2m} = m)} &= \frac{P(S_{2m} = m+d)}{P(S_{2m} = m+d-1)} \frac{P(S_{2m} = m+d-1)}{P(S_{2m} = m+d-2)} \cdots \frac{P(S_{2m} = m+1)}{P(S_{2m} = m)} \\ &= \frac{2m-m-d+1}{m+d} \cdots \frac{2m-m}{m+1} \\ &= \frac{m-d+1}{m-d} \cdots \frac{m}{m+1} \\ &= \frac{1 + \frac{1-d}{m}}{1 + \frac{d}{m}} \cdots \frac{1}{1 + \frac{1}{m}} \\ &= \frac{1 - \frac{d-1}{m}}{1 + \frac{1}{m}} \cdots \frac{1}{1 + \frac{d}{m}}. \end{aligned}$$

Če je d dovolj majhen, je to približno enako

$$\begin{aligned} \frac{P(S_{2m} = m+d)}{P(S_{2m} = m)} &\approx \frac{e^{-\frac{d-1}{m}} e^{-\frac{d-2}{m}} \cdots e^0}{e^{\frac{1}{m}} e^{\frac{2}{m}} \cdots e^{\frac{d}{m}}} \\ &= \exp\left(-\frac{1}{m} - \frac{2}{m} - \cdots - \frac{d}{m} - \frac{d-1}{m} - \cdots - \frac{1}{m}\right) \\ &= \exp\left(-\frac{d^2}{m}\right). \end{aligned}$$

Pri $d < 0$ je porazdelitev simetrična, torej dobimo enak rezultat. Izkaže se

$$\lim_{m \rightarrow \infty} \sum_{d=-\infty}^{\infty} \left| \frac{P(S_{2m} = m+d)}{P(S_{2m} = m)} - e^{-d^2/m} \right| = 0.$$

Vrsta

$$\sum_{d=-\infty}^{\infty} \frac{P(S_{2m} = m+d)}{P(S_{2m} = m)} = \frac{1}{P(S_{2m} = m)}$$

je približno enaka

$$\sum_{d=-\infty}^{\infty} e^{-d^2/m} \approx \int_{-\infty}^{\infty} e^{-x^2/m} dx = \sqrt{m\pi},$$

torej je verjetnost $P(S_{2m} = m) \approx \frac{1}{\sqrt{m\pi}}$. Upoštevaje zgornjo formulo potem dobimo

$$P(S_{2m} = m+d) \approx \frac{1}{\sqrt{m\pi}} e^{-\frac{d^2}{m}}$$

oziroma

$$P(S_n = k) \approx \sqrt{\frac{2}{n\pi}} \exp\left(-\frac{2(k - n/2)^2}{n}\right).$$

Vse to je res tudi za lihe n . Rezultatu pravimo DE MOIVREOVA LOKALNA FORMULA.

Vprašanje 1. Izpelj de Moivreovo lokalno formulo.

Za $a, b \in \mathbb{Z}$ lahko izpeljemo

$$\begin{aligned} P(a \leq S_n \leq b) &= \sum_{k=a}^b P(S_n = k) \\ &\approx \sqrt{\frac{2}{n\pi}} \sum_{k=a}^b \exp\left(-\frac{2(k - \frac{m}{2})^2}{n}\right) \\ &\approx \sqrt{\frac{2}{n\pi}} \sum_{k=a}^b \int_{k-\frac{1}{2}}^{k+\frac{1}{2}} \exp\left(-\frac{2(x - \frac{m}{2})^2}{n}\right) dx \\ &= \sqrt{\frac{2}{n\pi}} \int_a^b \exp\left(-\frac{2(x - \frac{m}{2})^2}{n}\right) dx, \end{aligned}$$

kar se da razširiti za poljubna $a < b$. Dobljeni približek za verjetnost $P(S_n = k)$ je gostota normalne porazdelitve $N(\frac{n}{2}, \frac{n}{4})$ v točki k . Velja $E(S_n) = \frac{n}{2}$ in $\text{var}(S_n) = \frac{n}{4}$. Binomska porazdelitev, s katero smo začeli, je enaka vsoti n neodvisnih Bernoullijevih slučajnih spremenljivk, ki je torej porazdeljena približno normalno.

Izrek (centralni limitni izrek). *Naj bodo X_1, X_2, \dots neodvisne in enako porazdeljene slučajne spremenljivke s končnim drugim momentom. Označimo $\mu_1 = E(X_i)$ in $\sigma_1^2 = \text{var}(X_i)$. Za $S_n = X_1 + \dots + X_n$ veljata naslednji točki:*

- Če definiramo $W_n = \frac{S_n - n\mu_1}{\sigma_1\sqrt{n}}$, za vsak $w \in \mathbb{R}$ velja

$$\lim_{n \rightarrow \infty} P(W_n \leq w) = \phi(w).$$

- *Limita*

$$\lim_{n \rightarrow \infty} \sup_{w \in \mathbb{R}} |P(W_n < w) - \phi(w)| = \lim_{n \rightarrow \infty} \sup_{x \in \mathbb{R}} \left| P(S_n < x) - \phi\left(\frac{x - n\mu_1}{\sigma_1\sqrt{n}}\right) \right| = 0.$$

Vprašanje 2. Formuliraj centralni limitni izrek.

Dokažemo lahko tudi drugačno formulacijo: za vsako zvezno in omejeno funkcijo $h : \mathbb{R} \rightarrow \mathbb{R}$ z največ kvadratno rastjo (tj. obstaja M , da je $|h(x)| \leq M(1 + x^2)$) velja

$$\lim_{n \rightarrow \infty} E(h(W_n)) = E(h(Z)),$$

kjer je $Z \sim N(0, 1)$. Funkcijam h pravimo POIZKUSNE ali TESTNE funkcije.

Trditev (Jensenova neenakost). Naj bo X slučajna spremenljivka z vrednostmi na intervalu $I \subseteq \mathbb{R}$ in $\varphi : I \rightarrow \mathbb{R}$ konveksna funkcija. Tedaj je $\varphi(E(X)) \leq E(\varphi(X))$.

Ideja dokaza je, da poiščemo linearno funkcijo, ki se v iskani točki dotika grafa φ , in vrednost aproksimiramo s pomočjo nje. Trditev je pomembna zaradi posledice.

Posledica. Če je $p \geq q$, je $E(|X|^q) \leq (E(|X|^p))^{q/p}$.

Posledica. Če obstaja p -ti moment in je $q \leq p$, obstaja tudi q -ti moment.

Za nadaljevanje naj bodo Y_1, \dots, Y_n neodvisne z vsoto S in Z_1, \dots, Z_n neodvisne z vsoto T . Privzemimo še, da tretji momenti Y_k in Z_k obstajajo, da je $E(Y_k) = E(Z_k) = 0$ in da je $\text{var}(Y_k) = \text{var}(Z_k) = \sigma_k^2$. Tedaj je $E(T) = E(S) = 0$ in

$$\text{var}(S) = \text{var}(T) = \sum_k \sigma_k^2.$$

Za primerno $h : \mathbb{R} \rightarrow \mathbb{R}$ bomo ocenili razliko $E(h(S)) - E(h(T))$. Privzeli bomo, da je $h \in \mathcal{C}^3$ in $M_3 = \sup_{x \in \mathbb{R}} |h'''(x)| < \infty$. Definirajmo kombinirane vsote $V_k = Y_1 + \dots + Y_k + Z_{k+1} + \dots + Z_n$.

Dodatno privzemimo, da so Y_i in Z_i vsi medsebojno neodvisni. Velja

$$E(h(S)) - E(h(T)) = \sum_{k=1}^n E(h(V_k) - h(V_{k-1})) = \sum_{k=1}^n E(h(U_k + Y_k) - h(U_k + Z_k))$$

za $U_k = Y_1 + \dots + Y_{k-1} + Z_{k+1} + \dots + Z_n$. Velja

$$h(V_k) = h(U_k) + h'(U_k)Y_k + \frac{1}{2}h''(U_k)Y_k^2 + R_k,$$

kjer je R_k omejen z $\frac{1}{6}M_3|Y_k|^3$. Poleg tega je

$$h(V_{k-1}) = h(U_k) + h'(U_k)Z_k + \frac{1}{2}h''(U_k)Z_k^2 + \tilde{R}_k,$$

kjer je \tilde{R}_k podobno omejena. Zaradi neodvisnosti U_k , Y_k in Z_k potem dobimo

$$E(h(V_k)) = E(h(U_k)) + E(h'(U_k))E(Y_k) + \frac{1}{2}E(h''(U_k))E(Y_k^2) + E(R_k),$$

$$E(h(V_{k-1})) = E(h(U_k)) + E(h'(U_k))E(Z_k) + \frac{1}{2}E(h''(U_k))E(Z_k^2) + E(\tilde{R}_k).$$

Torej

$$|E(h(V_k) - h(V_{k-1}))| \leq E(R_k) + E(\tilde{R}_k) \leq \frac{1}{6}M_3(E(|Y_k|^3) + E(|Z_k|^3)),$$

iz česar naposled dobimo

$$|E(h(S)) - E(h(T))| \leq \frac{1}{6}M_3 \sum_{k=1}^n (E(|Y_k|^3) + E(|Z_k|^3)).$$

Vprašanje 3. Izpelji oceno za razliko $|E(h(S)) - E(h(T))|$, kjer je $S = Y_1 + \dots + Y_n$ in $T = Z_1 + \dots + Z_n$, in so Y_i enako porazdeljeni neodvisni, ter Z_i enako porazdeljeni in neodvisni z $E(Y_i) = E(Z_i) = 0$, $\text{var}(Y_i) = \text{var}(Z_i) = \sigma_i^2$.

Če vzamemo $Z_k \sim N(0, \sigma_k^2)$, je $T \sim N(0, \sigma^2)$. Tedaj lahko z integralom izračunamo

$$E(|Z_k|^3) = \sigma_k^3 \frac{4}{\sqrt{2\pi}} = \frac{4}{\sqrt{2\pi}} (E(Y_k^2))^{3/2} \leq \frac{4}{\sqrt{2\pi}} E(|Y_k|^3),$$

kjer smo v zadnjem koraku uporabili posledico Jensenove neenakosti. V tem primeru torej lahko ocenimo

$$|E(h(S)) - E(h(T))| \leq \left(\frac{1}{6} + \frac{2}{3\sqrt{2\pi}} \right) M_3 \sum_{k=1}^n E(|Y_k|^3).$$

Tu žal ne moramo vzeti $h(w) = \mathbb{1}(w \leq a)$. Velja pa naslednji izrek.

Izrek. *Obstaja taka univerzalna konstanta C , da za poljubne neodvisne slučajne spremenljivke Y_1, \dots, Y_n , za katere je $E(Y_k) = 0$, velja*

$$\sup_{x \in \mathbb{R}} \left| P(S \leq x) - \phi\left(\frac{x}{\sigma}\right) \right| = \sup_{w \in \mathbb{R}} \left| P\left(\frac{S}{\sigma} \leq w\right) - \phi(w) \right| \leq \frac{C}{\sigma^3} \sum_{k=1}^n E(|Y_k|^3),$$

pri čemer je $S = Y_1 + \dots + Y_n$ in $\sigma^2 = \text{var}(S) = \sum_k \text{var}(Y_k)$.

Naj bodo X_1, X_2, \dots neodvisne in enako porazdeljene z $E(X_1) = \mu_1$ ter $\text{var}(X_1) = \sigma_1^2$ in $E(|X_1 - \mu_1|^3) = \gamma^3$. Definiramo $S_n = X_1 + \dots + X_n$ in

$$W_n = \frac{S_n - n\mu_1}{\sigma_1 \sqrt{n}}.$$

Za fiksno n označimo

$$Y_k = \frac{X_k - \mu_1}{\sigma_1 \sqrt{n}},$$

da je $E(Y_k) = 0$ in $W_n = \sum_k Y_k$. Velja $\sigma^2 = \text{var}(W_n) = 1$ in

$$\sum_{k=1}^n E(|Y_k|^3) = n E(|Y_1|^3) = \frac{\gamma_1^3}{\sigma_1^3 \sqrt{n}} \xrightarrow{n \rightarrow \infty} 0.$$

Posledica. *S temi oznakami*

$$\sup_{x \in \mathbb{R}} \left| P(S_n < x) - \phi\left(\frac{x - n\mu_1}{\sigma_1 \sqrt{n}}\right) \right| = \sup_{w \in \mathbb{R}} |P(W_n < w) - \phi(w)| \leq \frac{C}{\sqrt{n}} \frac{\gamma_1^3}{\sigma_1^3},$$

za vsako funkcijo $h \in \mathcal{C}^3(\mathbb{R})$ pa velja še

$$|E(h(W_n)) - E(h(Z))| \leq \left(\frac{1}{6} + \frac{2}{3\sqrt{2\pi}} \right) \frac{\gamma_1^3}{\sigma_1^3} \frac{M_3}{\sqrt{n}}.$$

Vprašanje 4. Kaj lahko poveš o oceni razlike porazdelitve vsote od normalne porazdelitve?

9.2 Konvergenca porazdelitev

Definicija. Zaporedje realnih slučajnih spremenljivk X_1, X_2, \dots v PORAZDELITVI KONVERGIRA proti slučajni spremenljivki X , če za vsak $x \in \mathbb{R}$, za katerega je $P(X = x) = 0$, velja

$$\lim_{n \rightarrow \infty} P(X_n \leq x) = P(X \leq x).$$

Pišemo $X_n \xrightarrow[n \rightarrow \infty]{d} X$.

Opomba. Temu pravimo tudi ŠIBKA KONVERGENCA.

Definicija. Zaporedje slučajnih spremenljivk X_1, X_2, \dots z vrednostmi v topološkem prostoru S v PORAZDELITVI KONVERGIRA PROTI X , če za vsako zvezno in omejeno funkcijo $h : S \rightarrow \mathbb{R}$ velja

$$\lim_{n \rightarrow \infty} E(h(X_n)) = E(h(X)).$$

Funkcijam h rečemo tudi PREIZKUSNE ali TESTNE funkcije.

Izrek (Helly-Bray). *Za porazdelitve na realni osi z običajno topologijo definiciji sovpadata.*

Izrek (Centralni limitni izrek). *Če so X_1, X_2, \dots neodvisne in enako porazdeljene z $E(X_1) = \mu_1$ in $\text{var}(X_1) = \sigma_1^2$, velja*

$$\frac{X_1 + \dots + X_n - n\mu_1}{\sqrt{n}} \xrightarrow[n \rightarrow \infty]{d} N(0, \sigma_1^2).$$

Vprašanje 5. Definiraj konvergenco v porazdelitvi in reformuliraj centralni limitni izrek.

Trditev. *Naj bodo X_1, X_2, \dots, X slučajne spremenljivke z vrednostmi v topološkem prostoru S in T še en topološki prostor. Naj bo $g : S \rightarrow T$ zvezna. Če velja $X_n \xrightarrow[n \rightarrow \infty]{d} X$, velja tudi $g(X_n) \xrightarrow[n \rightarrow \infty]{d} g(X)$.*

Dokaz. Naj bo $h : T \rightarrow \mathbb{R}$ zvezna in omejena. Tedaj je tudi $h \circ g : S \rightarrow \mathbb{R}$ zvezna in omejena, zato je

$$E(h(g(X))) = \lim_{n \rightarrow \infty} E(h(g(X_n))) = \lim_{n \rightarrow \infty} E(h \circ g(X_n)) = E(h \circ g(X)).$$

□

Trditev. *Naj bo $S_0 \subseteq S$ in naj imajo X_1, X_2, \dots, X vrednosti v S_0 . Če gre $X_n \xrightarrow[n \rightarrow \infty]{d} X$ v okviru S_0 , to velja tudi v okviru prostora S .*

Dokaz. Naj bo $h : S \rightarrow \mathbb{R}$ zvezna in omejena, ter naj bo h_0 zožitev h na S_0 . Tudi ta funkcija je zvezna in omejena, zato je $\lim_{n \rightarrow \infty} E(h_0(X_n)) = E(h_0(X))$. □

Trditev. Naj bo S metrizabilen prostor, $S_0 \subseteq S$ pa zaprt podprostor. Če imajo X_1, X_2, \dots, X vrednosti v S_0 in $X_n \xrightarrow[n \rightarrow \infty]{d} X$ v okviru S , to velja tudi v okviru S_0 .

Dokaz. Vsaka zvezna omejena funkcija $h_0 : S_0 \rightarrow \mathbb{R}$ se da po Tietzejevem izreku razširiti do zvezne in omejene funkcije $h : S \rightarrow \mathbb{R}$. \square

Vprašanje 6. V katerem primeru se konvergenca v porazdelitvi razširi iz večjega prostora v podprostor? Kaj je ideja dokaza?

Izrek. Naj bo S metrizabilen in $S_0 \subseteq S$ odprt podprostor. Naj bodo X_1, X_2, \dots slučajne spremenljivke z vrednostmi v S in X z vrednostmi v S_0 . Naj gre $X_n \xrightarrow[n \rightarrow \infty]{d} X$ v okviru prostora S . Nadalje naj bodo X_1^*, X_2^*, \dots slučajne spremenljivke z vrednostmi v S_0 ter naj bo $X_n^*(\omega) = X_n(\omega)$ brž ko je $X_n(\omega) \in S_0$. Tedaj gre $X_n \xrightarrow[n \rightarrow \infty]{d} X$ tudi v okviru S_0 .

Trditev. Naj bodo X_1, X_2, \dots slučajne spremenljivke z vrednostmi v metričnem prostoru (S, d) . Naj bo $c \in S$ konstantna. Tedaj gre $X_n \xrightarrow[n \rightarrow \infty]{d} c$ natanko tedaj, ko za vsak $r > 0$ velja

$$\lim_{n \rightarrow \infty} P(d(X_n, c) \geq r) = 0.$$

Dokaz. V desno: Obstaja taka zvezna funkcija $h : S \rightarrow [0, 1]$, da je $h(c) = 0$ in $h(x) = 1$ za $d(x, c) \geq r$. Velja $\mathbb{1}(d(x, c) \geq r) \leq h(x)$, na tej neenakosti pa lahko uporabimo pričakovano vrednost in dobimo

$$P(d(X_n, c) \geq r) \leq E(h(X_n)) \xrightarrow[n \rightarrow \infty]{} E(h(c)) = 0.$$

V levo: Naj bo $h : S \rightarrow [0, 1]$ zvezna in $\varepsilon > 0$. Obstaja tak $r > 0$, da za vsak $x \in S$ z $d(x, c) \geq r$ velja $|h(x) - h(c)| < \varepsilon/2$. Sledi

$$\begin{aligned} & |E(h(X_n)) - h(c)| \\ & \leq E(|h(X_n) - h(c)|) \\ & = E(|h(X_n) - h(c)| \mathbb{1}(d(X_n, c) < r)) + E(|h(X_n) - h(c)| \mathbb{1}(d(X_n, c) \geq r)) \\ & \leq \frac{\varepsilon}{2} + P(d(X_n, c) \geq r). \end{aligned}$$

Za dovolj pozne n bo to manj od ε . \square

Vprašanje 7. Kdaj zaporedje slučajnih spremenljivk konvergira k konstanti? Dokaži.

Trditev (Neenačba Markova). Za nenegativno slučajno spremenljivko W in poljuben $a > 0$ velja

$$P(W \geq a) \leq \frac{E(W)}{a}.$$

Dokaz. Velja $\mathbb{1}(w \geq a) \leq w/a$. Na tem uporabimo pričakovano vrednost. \square

Vprašanje 8. Povej in dokaži neenačbo Markova.

Trditev. Če je $p > 0$ in $\lim_{n \rightarrow \infty} E((d(X_n, c))^p) = 0$, gre $X_n \xrightarrow[n \rightarrow \infty]{d} c$.

Dokaz. Velja $P(d(X_n, c) \geq r) = P((d(X_n, c))^p \geq r^p)$. Na drugem delu uporabimo neenačbo Markova. \square

Trditev (šibki zakon velikih števil). Če so X_1, X_2, \dots neodvisne in enako porazdeljene z $E(X_1^2) < \infty$ in $E(X_1) = \mu$, gre

$$\frac{X_1 + \dots + X_n}{n} \xrightarrow[n \rightarrow \infty]{d} \mu.$$

Dokaz. Pričakovana vrednost izraza je očitno μ . Potem je

$$E\left(\left(\frac{X_1 + \dots + X_n}{n} - \mu\right)^2\right) = \text{var}\left(\frac{X_1 + \dots + X_n}{n}\right) = \frac{\text{var}(X_1)}{n},$$

to pa konvergira k 0 za $n \rightarrow \infty$. \square

Vprašanje 9. Povej in dokaži šibki zakon velikih števil.

Opomba. Velja celo

$$P\left(\lim_{n \rightarrow \infty} \frac{X_1 + \dots + X_n}{n} = \mu\right) = 1.$$

Temu pravimo KREPKI ZAKON VELIKIH ŠTEVIL.

Izrek. Naj bo S metrizabilen prostor, ki je števna unija svojih kompaktnih podprostorov. Naj bodo X_1, X_2, \dots slučajne spremenljivke z vrednostmi v S . Naj bo T še en metrizabilen prostor in naj imajo Y_1, Y_2, \dots vrednosti v T . Če gre $X_n \xrightarrow[n \rightarrow \infty]{d} X$ in $Y_n \xrightarrow[n \rightarrow \infty]{d} c$, kjer je $c \in T$ konstanta, gre tudi $(X_n, Y_n) \xrightarrow[n \rightarrow \infty]{d} (X, c)$.

Vprašanje 10. Kaj velja za konvergenco parov slučajnih spremenljivk?

Naslednje tri trditve se imenujejo IZREKI SLUCKEGA.

Posledica. Naj bodo $X_1, X_2, \dots, X, Y_1, Y_2, \dots$ slučajni vektorji z vrednostmi v \mathbb{R}^m in $c \in \mathbb{R}^m$ konstanta. Če gre $X_n \xrightarrow[n \rightarrow \infty]{d} X$ in $Y_n \xrightarrow[n \rightarrow \infty]{d} c$, gre tudi $X_n + Y_n \xrightarrow[n \rightarrow \infty]{d} X + c$ in $X_n - Y_n \xrightarrow[n \rightarrow \infty]{d} X - c$.

Posledica. Naj bodo X_1, X_2, \dots, X slučajni vektorji z vrednostmi v U_1, U_2, \dots realne slučajne spremenljivke in c konstanta. Če gre $X_n \xrightarrow[n \rightarrow \infty]{d} X$ in $U_n \xrightarrow[n \rightarrow \infty]{d} c$, gre tudi $U_n X_n \xrightarrow[n \rightarrow \infty]{d} cX$.

Trditev. Naj bodo $X_1, X_2, \dots, X, U_1, U_2, \dots$ in c kot prej. Privzamemo še $c \neq 0$.

- Če so Z_1, Z_2, \dots taki slučajni vektorji z vrednostmi v \mathbb{R}^m , da je $Z_n = X_n/U_n$ za $U_n \neq 0$, potem gre $Z_n \xrightarrow[n \rightarrow \infty]{d} X/c$.
- Naj bo $m = 1$. Tedaj za vsak $a \in \mathbb{R}$, za katerega je $P(X/c = a) = 0$, velja

$$\lim_{n \rightarrow \infty} P\left(U_n \neq 0, \frac{X_n}{U_n} \leq a\right) = P\left(\frac{X}{c} \leq a\right)$$

ter

$$\lim_{n \rightarrow \infty} P\left(U_n \neq 0, \frac{X_n}{U_n} \geq a\right) = P\left(\frac{X}{c} \geq a\right).$$

Dokaz. Prva točka: Po izreku gre

$$\begin{bmatrix} X_n \\ U_n \end{bmatrix} \xrightarrow[n \rightarrow \infty]{d} \begin{bmatrix} X \\ c \end{bmatrix}.$$

Slučajni vektor (X, c) zavzame vrednosti $S_0 = \mathbb{R}^m \times (\mathbb{R} \setminus \{0\})$. V tej množici zavzamejo vrednosti tudi slučajni vektorji

$$\begin{bmatrix} X_n^* \\ U_n^* \end{bmatrix} = \begin{cases} (X_n, U_n), & U_n \neq 0 \\ (Z_n, 1), & U_n = 0 \end{cases}$$

Velja $Z_n = X_n^*/U_n^*$. Potem tudi $(X_n^*, U_n^*) \xrightarrow[n \rightarrow \infty]{d} (X, c)$ v okviru S_0 . Ker je preslikava $(x, u) \mapsto x/u$ zvezna $S_0 \rightarrow \mathbb{R}^m$, velja

$$Z_n = \frac{X_n^*}{U_n^*} \xrightarrow[n \rightarrow \infty]{d} \frac{X}{c}.$$

Druga točka: Definirajmo $Z_n = X_n/U_n$ za $U_n \neq 0$ ter $Z_n = a + 1$ za $U_n = 0$. Tedaj je $\{U_n \neq 0, X_n/U_n \leq a\} = \{Z_n \leq a\}$. Rezultat sledi po prvi točki, drugo limito pokažemo podobno. \square

Vprašanje 11. Povej izreke Sluckega. Dokaži izrek za deljenje.

9.3 Uvod v statistiko

9.3.1 Sklepna statistika

Naj bodo X_1, X_2, \dots neodvisne in enako porazdeljene slučajne spremenljivke z $E(X_1) = \mu$ in $\text{var}(X_1) = \sigma^2$. Recimo, da je σ fiksni, vrednosti μ pa ne poznamo. Lahko jo ocenimo na podlagi opaženih vrednosti X_1, \dots, X_n .

Zakon velikih števil pove, da gre $\overline{X_n} \xrightarrow[n \rightarrow \infty]{d} \mu$. Za koliko si upamo reči, da lahko $\overline{X_n}$ še odstopa od μ ? Želeli bi zatrditi, da je $\overline{X_n} - \delta < \mu < \overline{X_n} + \delta$. Postavimo torej INTERVAL

ZAUPANJA $(\overline{X}_n - \delta, \overline{X}_n + \delta)$ za μ . Maksimalni dopustni verjetnosti, da se ne zmotimo, pravimo STOPNJA TVEGANJA α . Tipična izbira je $\alpha = 0.05$. Zahtevamo torej

$$P(|\overline{X}_n - \mu| \geq \delta) \leq \alpha.$$

Če vemo, da je $X_1 \sim N(\mu, \sigma^2)$ in poznamo σ , potem vemo, da je $\overline{X}_n - \mu \sim N(0, \sigma^2/n)$. Za $Z \sim N(0, 1)$ torej želimo

$$\begin{aligned} P(|\overline{X}_n - \mu| \geq \delta) &= P\left(\left|\frac{\sigma}{\sqrt{n}}Z\right| \geq \delta\right) \\ &= P\left(|Z| \geq \frac{\delta\sqrt{n}}{\sigma}\right) \\ &= 2P\left(Z \geq \frac{\delta\sqrt{n}}{\sigma}\right) \\ &= 2\left(1 - \Phi\left(\frac{\delta\sqrt{n}}{\sigma}\right)\right) \\ &\leq \alpha. \end{aligned}$$

Ustrezna izbira je potem

$$\delta = \frac{\sigma}{\sqrt{n}}\Phi^{-1}\left(1 - \frac{\alpha}{2}\right).$$

Vprašanje 12. Izpelji širino intervala zaupanja za μ v primeru neodvisnih normalno porazdeljenih spremenljivk X_n .

Če nismo prepričani, da so X_i porazdeljene normalno, si lahko pomagamo s CLI. Velja

$$\frac{\overline{X}_n - \mu}{\sigma}\sqrt{n} \xrightarrow[n \rightarrow \infty]{d} N(0, 1).$$

Pri dovolj velikih n bo konstruirani interval zaupanja še vedno približno zadoščal izbrani stopnji tveganja, vendar je α le še NOMINALNA STOPNJA TVEGANJA.

Vprašanje 13. Kaj pa narediš, če X_n niso porazdeljene normalno?

Če ne poznamo σ , ga nadomestimo z EMPIRIČNIM STANDARDNIM ODKLONOM

$$\hat{\sigma}_n = \sqrt{\frac{1}{n}((X_1 - \overline{X}_n)^2 + \dots + (X_n - \overline{X}_n)^2)}.$$

To bo imelo smisel, če $\hat{\sigma}_n$ konvergira k σ . Po zakonu velikih števil vemo

$$\frac{(X_1 - \mu)^2 + \dots + (X_n - \mu)^2}{n} \xrightarrow[n \rightarrow \infty]{d} \sigma^2.$$

Za empirično porazdeljeno

$$W \sim \begin{pmatrix} X_1 & X_2 & \dots & X_n \\ 1/n & 1/n & \dots & 1/n \end{pmatrix}$$

vemo $\text{var}(W) = E(W^2) - (E(W))^2 = E((W - \mu)^2) - (E(W - \mu))^2$, torej

$$\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 - (\bar{X}_n - \mu)^2,$$

kjer prvi člen konvergira k σ^2 in drugi k 0. Po izreku Sluckega gre tudi

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 \xrightarrow[n \rightarrow \infty]{d} \sigma^2$$

načeloma v okviru \mathbb{R} , a po enem od izrekov prejšnjega poglavja tudi v $[0, \infty)$, torej konvergira tudi $\hat{\sigma} \rightarrow \sigma$.

Vprašanje 14. Kaj pa narediš v primeru enako porazdeljenih normalnih spremenljivk, če ne poznaš σ ?

9.3.2 Opisna statistika

Tukaj STATISTIKA pomeni povzetek opaženih vrednosti; če opazimo x_1, \dots, x_n , potem je statistika $f(x_1, \dots, x_n)$. Številu n pravimo NUMERUS. Za empirično porazdelitev definiramo ARITMETIČNO SREDINO kot povprečje vrednosti x_i , EMPIRIČNI STANDARDNI ODKLON pa kot

$$\sigma = \sqrt{\frac{1}{n} \sum_{k=1}^n (x_k - \bar{x})^2}.$$

Če uredimo vrednosti kot $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$, pravimo, da je $x_{(k)}$ k -TA VRSTILNA STATISTIKA. Za neko verjetnost $\alpha \in (0, 1)$ je KVANTIL slučajne spremenljivke X tako število x_α , da je $P(X < x_\alpha) \leq \alpha \leq P(X \leq x_\alpha)$. Kvantil ni enolično določen, razen v primeru, ko je gostota porazdelitve strogo pozitivna na nekem intervalu, drugje pa enaka 0. Nekaj kvantilov ima tudi svoja imena:

- $x_{1/2}$ je MEDIANA,
- $x_{1/3}$ in $x_{2/3}$ sta prvi in drugi TERCIL,
- $x_{1/4}, x_{1/2}, x_{3/4}$ so TERCILI
- desetine dajo DECILE, stotine CENTILE, itd.

Če kvantili niso enolično določeni, jim določamo arbitrarne vrednosti. Mediani običajno pripišemo vrednost

$$\frac{1}{2}(x_{1/2}^{\min} + x_{1/2}^{\max}).$$

Povezan pojem je INTERKVARTILNI RAZMIK (interquartile range)

$$\text{IQR} = x_{3/4}^{\max} - x_{1/4}^{\min}.$$

Vprašanje 15. Kaj so kvantili? Kaj je IQR?

Številu pojavitev posamične vrednosti pravimo FREKVENCA. Frekvence lahko sestavimo skupaj in prikažemo v histogramu. Želeli bi si, da le-ta čim manj odstopa od teoretične gostote porazdelitve. V primeru Gaussove gostote po kriteriju najmanjših kvadratov pridemo do Scottovega pravila

$$\text{širina} \approx \sigma \sqrt[3]{\frac{24\sqrt{2\pi}}{n}},$$

kar pa veselo ignoriramo, ko je to pripravno.

Vprašanje 16. Kako postaviš dober histogram?

Po dogovoru so OSAMELCI vrednosti izven intervala $[x_{1/4}^{\min} - \frac{3}{2}\text{IQR}, x_{3/4}^{\max} + \frac{3}{2}\text{IQR}]$.

9.3.3 Ocenjevanje in napovedovanje

Naše opažanje sestoji iz nekih podatkov X . Zanima nas druga količina Y , ki bi jo radi vsaj približno ocenili. Postavimo MATEMATIČNI MODEL, ki je opis mehanizma, za katerega domnevamo, da je generiral X in Y . V statistiki vzamemo verjetnostni model, kjer je vsaj ena od X oz. Y slučajna. Iščemo opazljivo količino \hat{Y} (tj. tako, ki jo lahko izračunamo z merljivo funkcijo $h(X)$), ki bo čim bližje Y .

V Bayesovem modelu sta X in Y slučajni spremenljivki na fiksnem verjetnostnem prostoru (Ω, \mathcal{F}, P) . Če poznamo X in Y živi na množici \mathbb{R} , je smiselna napoved $\hat{Y} = E(Y | X)$. Postopek je naslednji:

- Poznamo $P(X = x | Y = y)$
- Privzamemo apriorno porazdelitev Y
- Po Bayesovi formuli izračunamo $P(Y = y | X = x)$
- Izračunamo $E(Y | X = x)$

Bayesova statistika je pripravna, če podatke dobivamo postopoma.

Drug koncept je postavil Ronald Fisher. Tu imamo več verjetnostnih mer P_θ za $\theta \in I$. Modeliramo tako, da je X slučajna spremenljivka, ki jo poznamo, $y = g(\theta)$ pa je deterministična in je ne poznamo. Karakteristiko y ocenimo z njeno CENILKO \hat{y} . O cenilki govorimo, preden opazimo podatke; ko jih poznamo, iz njih izračunamo vrednost cenilke, in ji pravimo OCENA.

Vprašanje 17. Kakšna je razlika med Bayesovo in Fisherjevo statistiko?

9.3.4 Vrednotenje cenilk in prediktorjev

Naj bo $\hat{y} = h(X)$ cenilka za y . PRIČAKOVANA ali SREDNJA KVADRATIČNA NAPAKA cenilke je

$$\text{MSE}_\theta(\hat{y} | y) = E_\theta((\hat{y} - y)^2).$$

Pogosto jo še korenimo, da dobimo RMSE. Definiramo tudi PRISTRANSKOST

$$\text{Bias}_\theta(\hat{y} | y) = E_\theta(\hat{y}) - y.$$

Cenilka je NEPRISTRANSKA, če je pristranskost enaka 0 za vse θ . Opazimo, da je za nepristransko cenilko pričakovana kvadratna napaka enaka varianci. V tem primeru količini RMSE pravimo STANDARDNA NAPAKA. V splošnem je

$$\text{var}_\theta(\hat{y}) = \text{var}_\theta(\hat{y} - y) = E_\theta((\hat{y} - y)^2) - (E_\theta(\hat{y} - y))^2 = \text{MSE}_\theta(\hat{y} | y) - (\text{Bias}_\theta(\hat{y} | y))^2.$$

Vprašanje 18. Definiraj pričakovano napako in pristranskost.

Če podatke dobivamo postopoma, dobimo zaporedje cenilk \hat{y}_n . Zaporedje je ŠIBKO DOSLEDNO, če gre

$$\hat{y}_n \xrightarrow[n \rightarrow \infty]{d|\theta} y$$

za vse $\theta \in I$.

Trditev. Če je

$$\lim_{n \rightarrow \infty} \text{MSE}_\theta(\hat{y}_n | y) = 0$$

za vse $\theta \in I$, je zaporedje $(\hat{y}_n)_n$ šibko dosledno.

Dokaz. Sledi iz ene od trditev od prej; če je

$$\lim_{n \rightarrow \infty} E((d(X_n, c))^p) = 0,$$

potem $X_n \xrightarrow[n \rightarrow \infty]{d} c$. □

Vprašanje 19. Kdaj je zaporedje cenilk šibko dosledno? Povej zadostni pogoj.

Posledica. Če je zaporedje cenilk $(\hat{y}_n)_n$ asimptotično nepristransko, tj. da velja $\lim E_\theta(\hat{y}_n) = y$ za vsak $\theta \in I$, in če je $\lim \text{var}_\theta(\hat{y}_n) = 0$, je to zaporedje šibko dosledno.

Če so X_1, X_2, \dots enako porazdeljene z $E(X_1) = \mu$ in $E(X_1^2) < \infty$ ter nekorelirane, je \bar{X} dosledna cenilka za μ . Poleg tega je \bar{X} NAJBOLJŠA NEPRISTRANSKA LINEARNA CENILKA (NNLC) za μ . Najboljša tu pomeni, da ima najmanjšo srednjo kvadratno napako. Pri tem pazimo, da je \bar{X}^2 pristranska cenilka za μ^2 ; velja

$$E(\bar{X}^2) = \text{var}(\bar{X}) + (E(\bar{X}))^2 = \frac{\sigma^2}{n} + \mu^2,$$

je pa asimptotično nepristranska in šibko dosledna.

9.4 Pričakovana vrednost in varianca slučajnih vektorjev

Definicija. Za slučajni vektor $\underline{X} = (X_1, \dots, X_n)$ definiramo $E(\underline{X}) = (E(X_1), \dots, E(X_n))$.

Trditev. Za deterministično matriko A primerne velikosti in slučajni vektor \underline{X} velja $E(A\underline{X}) = AE(\underline{X})$.

Posledica. Če je \underline{u} determinističen, \underline{X} pa slučajen vektor, je $E(\underline{u} \cdot \underline{X}) = \underline{u} \cdot E(\underline{X})$.

Definicija. Za slučajno matriko M definiramo $E(M) = [E(M_{ij})]_{ij}$.

Trditev. Naj bo M slučajna matrika. Potem je $E(M^T) = E(M)^T$. Za poljubni deterministični matriki A in B primerne velikosti velja $E(AMB) = AE(M)B$.

Trditev. Za poljubna primerna slučajna vektorja \underline{X} in \underline{Y} ter za poljubni primerni slučajni matriki M in N velja $E(\underline{X} + \underline{Y}) = E(\underline{X}) + E(\underline{Y})$ in $E(M + N) = E(M) + E(N)$.

Definicija. KOVARIANČNA MATRIKA med slučajnimi vektorji \underline{X} in \underline{Y} je

$$\text{cov}(\underline{X}, \underline{Y}) = \begin{bmatrix} \text{cov}(X_1, Y_1) & \cdots & \text{cov}(X_1, Y_n) \\ \vdots & \ddots & \vdots \\ \text{cov}(X_n, Y_1) & \cdots & \text{cov}(X_n, Y_n) \end{bmatrix}.$$

Trditev. $\text{cov}(\underline{X}, \underline{Y}) = E((\underline{X} - E(\underline{X}))(\underline{Y} - E(\underline{Y}))^T) = E(\underline{XY}^T) - E(\underline{X})E(\underline{Y})^T$.

Vprašanje 20. Kako izraziš kovariančno matriko?

Izračunamo lahko

$$\text{cov}(A\underline{X}, B\underline{Y}) = A \text{cov}(\underline{X}, \underline{Y}) B^T,$$

iz česar za vektor \underline{u} dobimo

$$0 \leq \text{var}(\underline{u} \cdot \underline{X}) = \underline{u}^T \text{var}(\underline{X}) \underline{u} = \langle \text{var}(\underline{X}) \underline{u}, \underline{u} \rangle,$$

torej je kovariančna matrika slučajnega vektorja pozitivno semidefinitna.

Trditev. Pozitivno semidefinitna matrika je kovariančna matrika slučajnega vektorja.

Dokaz. Naj bo Σ pozitivno semidefinitna matrika. Vzemimo $\underline{Z} = (Z_1, \dots, Z_n)$, kjer so Z_i neodvisne in porazdeljene standardno normalno. Velja $\text{var}(\underline{Z}) = I_n$. Za Σ obstaja razcep Choleskega $\Sigma = VV^T$. Potem je $\text{cov}(A\underline{Z}, A\underline{Z}) = AA^T = \Sigma$. \square

Vprašanje 21. Pokaži, da je matrika enaka kovariančni matriki slučajnega vektorja natanko tedaj, ko je pozitivno semidefinitna.

Definicija. Slučajna vektorja \underline{X} in \underline{Y} sta NEKORELIRANA, če je $\text{cov}(\underline{X}, \underline{Y}) = 0$.

9.5 Pridobivanje cenilk

9.5.1 Metoda empirične porazdelitve

Opazimo vrednosti X_1, \dots, X_n , ki so enako porazdeljene. Te predstavljajo vzorec preučevane porazdelitve. Naj se y izraža kot karakteristika te porazdelitve, $y = g(\theta) = \text{Char}_\theta(X_1)$, kjer je Char_θ odvisna le od porazdelitve X_1 pri verjetnostni meri P_θ . Cenilka za y je potem $\hat{y} = \text{Char}(\text{Emp}(X_1, \dots, X_n))$.

Vprašanje 22. Opiši metodo empirične porazdelitve. Kakšne so predpostavke?

Kot primer vzemimo $\mu = E(X_1)$, ki ga ocenimo z $\hat{\mu} = E(\text{Emp}(X_1, \dots, X_n)) = \bar{X}_n$. To je nepristranska cenilka za μ . Če so X_1, \dots, X_n nekorelirane, vemo, da je

$$\text{SE} = \frac{\sigma}{\sqrt{n}}$$

kjer je $\sigma^2 = \text{var}(X_1)$. Drug pogled je enostavno slučajno vzorčenje. Na populaciji z N enotami naj bo definirana statistična spremenljivka, tj. funkcija $\{1, \dots, N\} \rightarrow \mathbb{R}$. Njene vrednosti na posameznih enotah izrazimo z x_1, \dots, x_N . Iz populacije vzamemo vzorec, ki ga sestavljajo enote K_1, \dots, K_n . Enostavno slučajno vzorčenje pomeni, da ima vektor (K_1, \dots, K_n) vrednosti v množici $\{(k_1, \dots, k_n) \mid k_i \neq k_j\}$. Označimo $X_i = x_{K_i}$, da velja $X_i \sim \text{Emp}(X_1, \dots, X_n)$. Potem je $\mu = E(X_1) = \mu$ in

$$\text{var}(X_1) = \frac{1}{N} \sum_{k=1}^N (x_k - \mu)^2.$$

Standardna napaka cenilke \bar{X} za \bar{x} je

$$\text{SE}^2 = \text{var}(\bar{X}) = \frac{1}{n^2} \text{var} \left(\sum_i X_i \right) = \frac{1}{n^2} \sum_{i,j} \text{cov}(X_i, X_j) = \frac{N-n}{N-1} \frac{\sigma^2}{n}.$$

Vprašanje 23. Opiši uporabo metode empirične porazdelitve za pričakovano vrednost.

Če ocenjujemo $\sigma^2 = \text{var}(X_1)$, je cenilka

$$\hat{\sigma}^2 = \frac{1}{n} \sum_k (X_k - \bar{X})^2$$

Za nepristranskost izračunajmo $E(\hat{\sigma}^2)$. Velja

$$\hat{\sigma}^2 = \frac{1}{n} \|(X_1 - \bar{X}, \dots, X_n - \bar{X})\|^2 = \frac{1}{n} \|\underline{X} - \bar{X} \cdot \underline{1}\|^2 = \frac{1}{n} \left\| \left(I - \frac{1}{n} \underline{1} \underline{1}^T \right) \underline{X} \right\|^2.$$

Označimo $H = I - \frac{1}{n} \underline{1} \underline{1}^T$. To je ortogonalni projektor, velja

$$\hat{\sigma}^2 = \frac{1}{n} (H \underline{X})^T (H \underline{X}) = \frac{1}{n} \underline{X}^T H \underline{X}.$$

Sled skalarja je enaka skalarju, velja

$$E(\hat{\sigma}^2) = \frac{1}{n} \text{sl}(HE(\underline{X}\underline{X}^T)).$$

Če zamenjamo spremenljivke $\underline{X}' = \underline{X} - \underline{\mu}$, lahko izrazimo

$$E(\hat{\sigma}^2) = \frac{1}{n} \text{sl}(HE(\underline{X}'\underline{X}'^T)) = \frac{1}{n} \text{sl}(H \text{var}(\underline{X})),$$

ker je $E(\underline{X}') = 0$.

Sedaj ločimo dva primera. Če so X_1, \dots, X_n nekorelirane, je $\text{var}(\underline{X}) = \sigma^2 I$ in posledično

$$E(\hat{\sigma}^2) = \frac{n-1}{n} \sigma^2,$$

torej je $\hat{\sigma}^2$ pristranska cenilka za σ^2 , cenilka $\hat{\sigma}_+^2 = \frac{n}{n-1} \hat{\sigma}^2$ pa je nepristranska.

Vprašanje 24. Pokaži, da je empirična varianca pristranska.

Če je SE standardna napaka cenilke \bar{X} za μ , potem je cenilka

$$\hat{\text{SE}}^2 = \frac{1}{n^2} \sum_{i=1}^n (X_i - \bar{X})^2$$

pristranska cenilka. Nepristranska cenilka bo

$$\hat{\text{SE}}_+^2 = \frac{1}{n(n-1)} \sum_{i=1}^n (X_i - \bar{X})^2.$$

Pri enostavnem slučajnem vzorčenju iz populacije z N enotami imamo

$$\text{var}(\underline{X}) = \sigma^2 \begin{bmatrix} 1 & \frac{-1}{N-1} & \cdots & \frac{-1}{N-1} \\ \frac{-1}{N-1} & 1 & \cdots & \frac{-1}{N-1} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{-1}{N-1} & \frac{-1}{N-1} & \cdots & 1 \end{bmatrix} = \sigma^2 \left(\frac{-1}{N-1} \underline{1}\underline{1}^T + \frac{N}{N-1} I \right).$$

Za matriko $H = I - \frac{1}{n} \underline{1}\underline{1}^T$ lahko izračunamo $H \text{var}(\underline{X}) = \frac{N\sigma^2}{N-1} H$, torej dobimo

$$E(\hat{\sigma}^2) = \frac{1}{n} \text{sl}(H \text{var}(\underline{X})) = \frac{1}{n} \text{sl} \left(\frac{N\sigma^2}{N-1} H \right) = \frac{N}{N-1} \frac{n-1}{n} \sigma^2$$

in je cenilka

$$\hat{\sigma}_+^2 = \frac{N-1}{N} \frac{n}{n-1} \hat{\sigma}^2$$

nepristranska za σ^2 .

Vprašanje 25. Kaj je nepristranska cenilka za σ^2 pri enostavnem slučajnem vzorčenju? Izpelji.

Poseben primer metode empirične porazdelitve je METODA MOMENTOV. Pri njej teoretični moment $m_k = E(X_1^k)$ ocenimo z empiričnim momentom, količino $y = g(m_1, \dots, m_k)$ pa ocenimo z $\hat{y} = g(\hat{m}_1, \dots, \hat{m}_k)$.

Vprašanje 26. Opiši metodo momentov.

9.5.2 Metoda največjega verjetja

VERJETJE je verjetnostna funkcija ali gostota, ki ga gledamo kot funkcijo parametra model θ . Za cenilko $\hat{\theta} = h(x)$ pravimo, da je CENILKA PO MNV, če verjetje $L(\theta | x)$ pri vsakem x doseže maksimum pri $\hat{\theta} = h(x)$. V praksi se izkaže, da je lažje maksimizirati $l = \ln L$. Če je $\underline{X} = (X_1, \dots, X_n)$, kjer so komponente neodvisne in enako porazdeljene, je potem l vsota verjetij za vsak X_i .

Vprašanje 27. Opiši metodo največjega verjetja.

Če je $I \subseteq \mathbb{R}^p$ in $\underline{\theta} = (\theta_1, \dots, \theta_p)$, je stacionarna točka verjetja ničla funkcije $\vec{\nabla} l$. Temu gradientu pravimo ZBIRNA FUNKCIJA. Recimo, da verjetje pride iz neke gostote $f(\underline{x} | \underline{\theta})$. Potem je

$$1 = \int_{\mathbb{R}^n} f(\underline{x} | \underline{\theta}) d\underline{x}.$$

Če to odvajamo po θ_j , in upoštevamo $\partial_{\theta_j} l = \partial_{\theta_j} L / L$, dobimo

$$E_{\theta} \left(\frac{\partial l}{\partial \theta_j}(\underline{\theta} | \underline{X}) \right) = \int_{\mathbb{R}^n} L(\underline{\theta} | \underline{x}) \frac{\partial l}{\partial \theta_j} d\underline{x} = 0.$$

Z računom lahko pokažemo

$$\frac{\partial^2 l}{\partial \theta_j \partial \theta_k} = \frac{\partial_{\theta_j} \partial_{\theta_k} L}{L} - \frac{\partial l}{\partial \theta_j} \frac{\partial l}{\partial \theta_k}.$$

Če odvajamo $\int L = 1$ po θ_j in θ_k , dobimo

$$\int_{\mathbb{R}^n} \frac{\partial_{\theta_j} \partial_{\theta_k} L}{L} f = 0,$$

iz česar sledi, da je pričakovana vrednost prvega člena dvojnega odvoda zgoraj enaka 0. Torej je

$$E_{\theta} \left(\frac{\partial^2 l}{\partial \theta_j \partial \theta_k} \right) = -E_{\theta} \left(\frac{\partial l}{\partial \theta_j} \frac{\partial l}{\partial \theta_k} \right) = -\text{cov}_{\theta} \left(\frac{\partial l}{\partial \theta_j}, \frac{\partial l}{\partial \theta_k} \right).$$

Pričakovana vrednost Hessejeve matrike Hl je torej nasprotna vrednost kovariančne matrike zbirne funkcije, in je negativno semidefinitna. Torej je stacionarna točka v povprečju maksimum. Definiramo FISHERJEVO INFORMACIJO

$$\text{FI}(\underline{\theta}) = \text{var}(\vec{\nabla} l(\underline{\theta} | \underline{x})) = -E(Hl(\underline{\theta} | \underline{x})).$$

Vprašanje 28. Kaj je Fisherjeva informacija? Izpelji povezavo med predstavitvama.

Če je $\underline{X} = (X_1, X_2, \dots, X_n)$, kjer so X_i neodvisni in enako porazdeljeni, potem lahko definiramo Fisherjevo informacijo ene komponente, da je $\text{FI} = n \text{FI}_1$.

Izrek. Naj bo $I^{\text{odp}} \subseteq \mathbb{R}^p$, FI_1 obrnljiva za vsak $\underline{\theta} \in I$ ter naj veljajo določeni dodatni tehnični pogoji. Če so X_i neodvisni in enako porazdeljeni, tedaj gre

- $\hat{\underline{\theta}} \xrightarrow[n \rightarrow \infty]{} \underline{\theta}$
- $\sqrt{n} \text{Bias}(\hat{\underline{\theta}} | \underline{\theta}) \xrightarrow[n \rightarrow \infty]{} 0$
- $n \text{var}(\hat{\underline{\theta}})(\text{FI}_1(\underline{\theta}))^{-1}$
- $n \text{MSE}(\hat{\underline{\theta}} | \underline{\theta}) \xrightarrow[n \rightarrow \infty]{d} (\text{FI}_1(\underline{\theta}))^{-1}$

Če je $g : I \rightarrow \mathbb{R}$ primerna funkcija ter $y = g(\underline{\theta})$, za $\hat{y} = g(\hat{\underline{\theta}})$ velja

- $\hat{y} \xrightarrow[n \rightarrow \infty]{d} y$
- $\sqrt{n} \text{Bias}(\hat{y} | y) \xrightarrow[n \rightarrow \infty]{} 0$
- $n \text{var}(\hat{y}) \xrightarrow[n \rightarrow \infty]{} (\vec{\nabla} \cdot g(\underline{\theta}))^T (\text{FI}_1(\underline{\theta}))^{-1} (\vec{\nabla} \cdot g(\underline{\theta}))$
- $n \text{MSE}(\hat{y} | y) \xrightarrow[n \rightarrow \infty]{} (\vec{\nabla} \cdot g(\underline{\theta}))^T (\text{FI}_1(\underline{\theta}))^{-1} (\vec{\nabla} \cdot g(\underline{\theta}))$

Vprašanje 29. Povej izrek o metodi največjega verjetja.

9.6 Večrazsežna normalna porazdelitev

Definicija. STANDARDNA VEČRAZSEŽNA NORMALNA PORAZDELITEV je porazdelitev slučajnega vektorja (Z_1, \dots, Z_n) , kjer so $Z_i \sim N(0, 1)$ neodvisne.

Opomba. To je n -razsežna porazdelitev z gostoto

$$\phi_n(z_1, \dots, z_n) = \frac{1}{(2\pi)^{n/2}} \exp\left(-\frac{z_1^2}{2} - \dots - \frac{z_n^2}{2}\right)$$

Definicija. VEČRAZSEŽNA NORMALNA PORAZDELITEV je poljubna porazdelitev slučajnega vektorja oblike $A\underline{Z} + \underline{\mu}$, kjer je \underline{Z} slučajni večrazsežni normalni vektor, A deterministična matrika, $\underline{\mu}$ pa deterministični vektor primerne velikosti.

Opomba. Za vektor $\underline{X} = A\underline{Z} + \underline{\mu}$ je $E(\underline{X}) = \underline{\mu}$ in $\text{var}(\underline{X}) = A A^T$.

Vprašanje 30. Definiraj večrazsežno normalno porazdelitev.

Če je \underline{Z} standardni n -razsežni normalni slučajni vektor, $A \in \mathbb{R}^{n \times n}$ deterministična matrika in $\underline{\mu} \in \mathbb{R}^n$. Teda j ima $\underline{X} = A\underline{Z} + \underline{\mu}$ gostoto

$$f_{\underline{X}}(\underline{x}) = \frac{1}{(2\pi)^{n/2} \sqrt{\det \Sigma}} \exp\left(-\frac{1}{2}(\underline{x} - \underline{\mu})^T \Sigma^{-1}(\underline{x} - \underline{\mu})\right)$$

za $\Sigma = AA^T$.

Trditev. Če je $m \leq n$, $\underline{\mu} \in \mathbb{R}^m$ in ima $A \in \mathbb{R}^{m \times n}$ poln rang, je $A\underline{Z} + \underline{\mu}$ porazdeljena $N(\underline{\mu}, AA^T)$.

Dokaz. Razcepimo lahko $A = BPQ$, kjer je $Q \in \mathbb{R}^{n \times n}$ ortogonalna, $P \in \mathbb{R}^{m \times n}$ koordinatna projekcija na prvih m komponent, in $B \in \mathbb{R}^m$ neizrojena. Ker je $QQ^T = I$, je $Q\underline{Z} \sim N(0, I)$, in so komponente tega vektorja neodvisne standardno normalne. Vektor $PQ\underline{Z}$ je prav tako standardni normalni z enakimi komponentami, le da jih ima m namesto n . Torej je $A\underline{Z} + \underline{\mu} = BPQ\underline{Z} + \underline{\mu}$ porazdeljen $N(\underline{\mu}, BB^T) = N(\underline{\mu}, AA^T)$. \square

Vprašanje 31. Kako je porazdeljen vektor $A\underline{Z} + \underline{\mu}$, kjer je A polnega ranga?

Trditev. Porazdelitev poljubnega večrazsežnega normalnega slučajnega vektorja je natančno določena z njegovo pričakovano vrednostjo in kovariančno matriko.

Dokaz. Če je \underline{X} m -razsežen normalen slučajni vektor, po definiciji obstajata $A \in \mathbb{R}^{m \times n}$ in $\underline{\mu} \in \mathbb{R}^m$, da je $\underline{X} = A\underline{Z} + \underline{\mu}$ za $\underline{Z} \sim N(0, I_n)$.

Naj bo $\tilde{\underline{Z}} \sim N(0, I_m)$ neodvisen od \underline{Z} . Definiramo $\underline{X}_k = A\underline{Z} + \frac{1}{k}\tilde{\underline{Z}} + \underline{\mu}$. Iz izrekov Sluckega sledi $\underline{X}_k \xrightarrow[k \rightarrow \infty]{d} \underline{X}$. Ker sta \underline{Z} in $\tilde{\underline{Z}}$ neodvisna, je vektor $(\underline{Z}, \tilde{\underline{Z}}) \sim N(0, I_{n+m})$. Pišimo

$$\underline{X}_k = \begin{bmatrix} A & \frac{1}{k}I_m \end{bmatrix} \cdot \begin{bmatrix} \underline{Z} \\ \tilde{\underline{Z}} \end{bmatrix} + \underline{\mu}.$$

Ta matrika je polnega ranga, zato je porazdelitev slučajnega vektorja \underline{X}_k odvisna le od pričakovane vrednosti $\underline{\mu}$ in kovariančne matrike $\text{var}(\underline{X}_k) = AA^T + \frac{1}{k^2}I_m$. Porazdelitev \underline{X} je limita verjetnostnih porazdelitev $N(\underline{\mu}, AA^T + \frac{1}{k^2}I_m)$; ta je natančno določena, torej je porazdelitev \underline{X} natančno določena z $\underline{\mu}$ in AA^T . \square

Vprašanje 32. Dokaži: porazdelitev normalnega slučajnega vektorja je natančno določena s pričakovano vrednostjo in kovariančno matriko.

Trditev. Če je bločni slučajni vektor $(\underline{X}_1, \underline{X}_2)$ porazdeljen večrazsežno normalno in sta \underline{X}_1 ter \underline{X}_2 nekorelirana, sta tudi neodvisna.

Dokaz. Velja

$$(\underline{X}_1, \underline{X}_2) \sim N\left(\begin{bmatrix} \underline{\mu}_1 \\ \underline{\mu}_2 \end{bmatrix}, \begin{bmatrix} \Sigma_1 & \\ & \Sigma_2 \end{bmatrix}\right),$$

to pa je tudi porazdelitev vektorja, kjer sta prvi in drugi blok nekorelirana. \square

Če je \underline{Z} porazdeljen standardno normalno z n prostorskimi stopnjami, je $\|\underline{Z}\|^2 \sim \chi^2(n) = \Gamma(\frac{n}{2}, \frac{1}{2})$.

Trditev. Če je H ortogonalni projektor ranga p in $\underline{X} \sim N(0, H)$, je $\|\underline{X}\|^2 \sim \chi^2(p)$.

Dokaz. Obstaja ortogonalna matrika Q , da je

$$QH Q^T = \begin{bmatrix} I_p & \\ & 0_{n-p} \end{bmatrix}.$$

Velja $Q\underline{X} \sim N(0, QH Q^T)$ in $\|Q\underline{X}\| = \|\underline{X}\|$. □

Vprašanje 33. Opiši porazdelitev χ^2 .

Opomba. Če so $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ neodvisne, je

$$\frac{1}{\sigma^2} \sum_i (X_i - \bar{X})^2$$

porazdeljena $\chi^2(n-1)$.

Definicija. STUDENTOVA PORAZDELITEV s p prostorskimi stopnjami je porazdelitev slučajne spremenljivke $Z/\sqrt{H/p}$, kjer sta $Z \sim N(0, 1)$ in $H \sim \chi^2(p)$ neodvisni.

Vprašanje 34. Definiraj Studentovo porazdelitev.

Če so $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ neodvisne slučajne spremenljivke, vemo, da je

$$\frac{\bar{X} - \mu}{\sigma} \sqrt{n} \sim N(0, 1).$$

Če ne poznamo σ , bi ga radi nadomestili s cenilko $\hat{\sigma}_+$. Pišimo

$$\frac{\bar{X} - \mu}{\hat{\sigma}_+} \sqrt{n} = \frac{\frac{\bar{X} - \mu}{\sigma} \sqrt{n}}{\hat{\sigma}_+/\sigma} = \frac{Z}{\sqrt{G/(n-1)}},$$

kjer je $Z \sim N(0, 1)$ in

$$G = \frac{1}{\sigma^2} \sum_i (X_i - \bar{X})^2.$$

Za $H = I_n - \frac{1}{n} \underline{1} \underline{1}^T$ velja $\bar{X} = \frac{1}{n} \underline{1}^T \underline{X}$ in $\sum (X_i - \bar{X})^2 = \|H\underline{X}\|^2$. Vektor

$$\begin{bmatrix} H\underline{X} \\ \underline{1}^T \underline{X} \end{bmatrix} = \begin{bmatrix} H \\ \underline{1}^T \end{bmatrix} \underline{X}$$

je porazdeljen večrazsežno normalno. Izračunamo lahko $\text{cov}(H\underline{X}, \underline{1}^T \underline{X}) = 0$, torej sta Z in G neodvisna, in je

$$\frac{\bar{X} - \mu}{\hat{\sigma}_+} \sqrt{n} \sim \text{Student}(n-1).$$

Vprašanje 35. Kako je porazdeljen $\frac{\bar{X} - \mu}{\hat{\sigma}_+} \sqrt{n}$, če so $X_i \sim N(\mu, \sigma^2)$ neodvisni?

9.7 Statistično sklepanje z nadzorovanim tveganjem

Opazimo X , zanima nas Y . Želimo zatrditi, da z visoko verjetnostjo velja $Y_{\min} < Y < Y_{\max}$, kjer sta Y_{\min} in Y_{\max} opazljivi. Intervalu med njima pravimo NAPOVEDNI INTERVAL, če pa je $Y = y = g(\theta)$ determinističen, pa INTERVAL ZAUPANJA.

Splošneje lahko izjavimo $Y \in C_X$, kjer je $\{C_X\}_X$ nek nabor množic. Naš cilj je nadzorovati verjetnost zmote $p_\theta = P_\theta(Y \notin C_X)$, za kar vnaprej izberemo STOPNJO TVEGANJA α , tipično 0.05. To lahko naredimo na več načinov:

- eksaktna stopnja tveganja: $p_\theta = \alpha$ za vsak θ ,
- nominalna stopnja tveganja: p_θ je čim bližje α ,
- asimptotično eksaktna stopnja tveganja: opazimo čim več vrednosti; zahtevamo

$$\lim_{n \rightarrow \infty} p_\theta = \alpha,$$

- konzervativna stopnja tveganja: $p_\theta \leq \alpha$ za vsak θ .

Nasprotni verjetnosti $1 - \alpha$ pravimo STOPNJA ZAUPANJA.

Vprašanje 36. Na katere načine lahko izberemo stopnjo zaupanja?

Za iskanje mej definiramo MERO RAZHAJANJA $\rho(x, y)$, ki meri neskladje med x in y . Če imamo srečo, je ρ PIVOTNA FUNKCIJA, kar pomeni, da je porazdelitev slučajne spremenljivke $\rho(X, Y)$ neodvisna od θ . Tedaj postavimo $C = \{(x, y) \mid \rho(x, y) < c_\alpha\}$, kjer c_α izračunamo tako, da je verjetnost zmote enaka α .

Osnovna ideja za iskanje $\rho(x, y)$ je, da nastavimo $\rho(x, y) = |h(x) - y|$, kjer je $h(x)$ cenilka za y . Če je $y = g(\theta)$ lastnost porazdelitve, lahko zanjo uporabimo metodo največjega verjetja, da dobimo \hat{y} . Potem velja naslednji izrek.

Izrek. Če definiramo

$$\widehat{RMSE} = \sqrt{\frac{1}{n} \vec{\nabla} \cdot g(\hat{\theta})^T \text{FI}_1^{-1}(\hat{\theta}) \vec{\nabla} \cdot g(\hat{\theta})},$$

kjer je $\hat{\theta}$ cenilka po MNV, pod podobnimi pogoji kot za prejšnji rezultat velja

$$\frac{\hat{y} - y}{\widehat{RMSE}} \xrightarrow[n \rightarrow \infty]{d} N(0, 1).$$

Dokaz. Skica. Pogledamo logaritem verjetja

$$l(\theta \mid x_1, \dots, x_n) = \sum_i l_1(\theta \mid x_i).$$

Normalna porazdelitev bo sledila iz CLI. Problem je, da to potrebujemo z cenilko; izpeljemo lahko, da je \hat{y} približno linearna funkcija vsote nekaj IID slučajnih spremenljivk. Iz večrazsežne variante CLI izpeljemo, da

$$\frac{\hat{y} - y}{\sqrt{\frac{1}{n} \vec{\nabla} \cdot g(\theta)^T \text{FI}_1^{-1}(\theta) \vec{\nabla} \cdot g(\theta)}} \xrightarrow[n \rightarrow \infty]{d} N(0, 1).$$

Manjka še, da smemo uporabiti $\hat{\theta}$ namesto θ , kar dobimo iz izrekov Slutkega. \square

Iz tega sledi asimptotični interval zaupanja

$$\hat{y} - \phi^{-1}\left(1 - \frac{\alpha}{2}\right) \text{RMSE} < y < \hat{y} + \phi^{-1}\left(1 - \frac{\alpha}{2}\right) \text{RMSE}.$$

Vprašanje 37. Kako dobiš asimptotični interval zaupanja z uporabo cenilke po MNV? Povej idejo dokaza izreka.