

Tereprendezés

1. Bevezetés

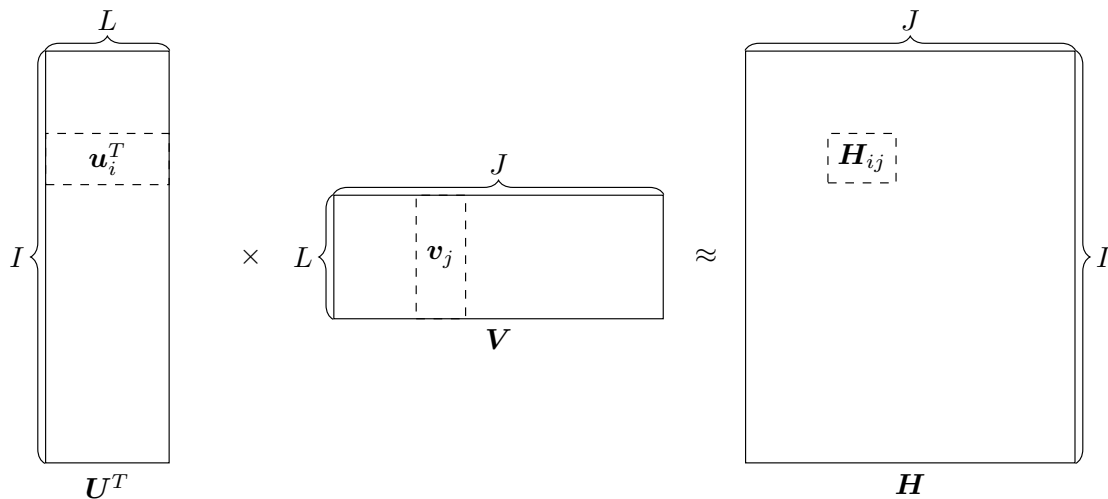
Csapatával stratégiai játékot fejlesztenek. A vezetés úgy dönt, hogy a fejlesztői eszköztárat a modderek számára ingyenesen elérhetővé kell tenni; ennek fontos része a pályagenerátor is. Kollégája el is készül egy terepgenerátorral, amely a jól bevált Diamond-Square algoritmust használja, ám a kapott terepek egytől egyig „tűskések”, zajosak, ráadásul sok helyet foglalnak a memóriában. Önre esett a választás, hogy orvosolja a problémát egy intelligens simító-tömörítő algoritmus implementációjával¹.

1.1. Alacsony rangú approximáció

Legyen adott egy $\mathbf{H} \in \mathbb{R}^{I \times J}$ magasságtérkép (*heightmap*), amely tehát egy $I \times J$ mátrix, amely elemeiben az egyes mezők magasságát tartalmazza. A feladatunk, hogy találjunk egy $\hat{\mathbf{H}}$ magasságtérképet, amely kelőképpen „sima”, mégis „hasonlít” az eredetihez. Nyilvánvalóan valamiféle közelítő megoldásra törekszünk, amely egyúttal a tárigény csökkentéséről is gondoskodik. Keressük hát a $\hat{\mathbf{H}}$ mátrixot a következő alakban:

$$\mathbf{H} \approx \hat{\mathbf{H}} = \mathbf{U}^T \mathbf{V},$$

ahol $\mathbf{U} \in \mathbb{R}^{L \times I}$ és $\mathbf{V} \in \mathbb{R}^{L \times J}$. Legyen $L \ll I, J$, így a \mathbf{U} és \mathbf{V} sokkal kisebb helyen tárolható, mint $\hat{\mathbf{H}}$ (ún. alacsony rangú approximáció). Az alábbi ábra ezt az elgondolást szemlélteti:

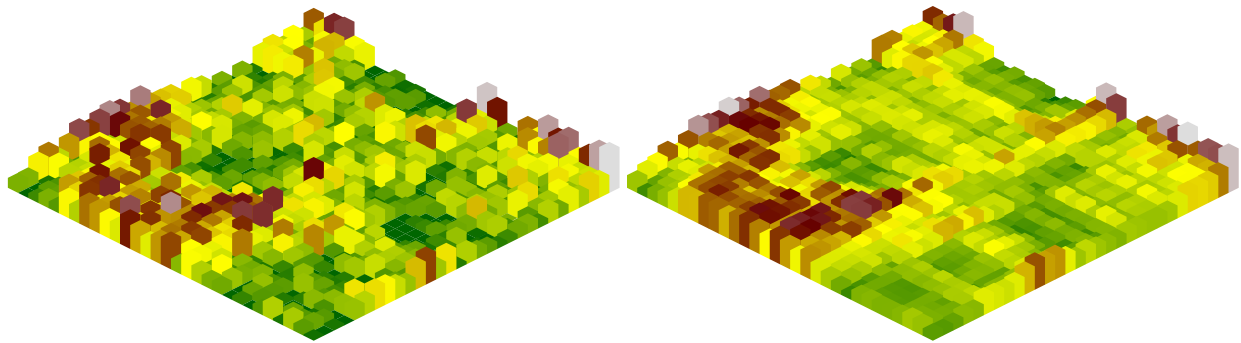


A \mathbf{H} mátrix i, j . elemét így is írhatjuk:

$$\mathbf{H}_{ij} \approx \mathbf{u}_i^T \mathbf{v}_j,$$

ahol \mathbf{u}_i és \mathbf{v}_j a megfelelő mátrix i . illetve j . oszlopvektora; ezek a vektorok L hosszúak.

¹Megjegyzendő, hogy a problémára ennél sokkal praktikusabb algoritmusok is léteznek. Ez a feladat inkább egy módszer bemutatását célozza, semmint egy gyakorlatban is jól használható algoritmus implementációját.



1. ábra. Bal oldalt az eredeti 30×30 magasságtérkép, jobb oldalt a simított változat $L = 3$ rang esetén. A tárigény $30^2 = 900$ értékről $2 \times 30 \times 3 = 180$ értékre csökkent.

1.2. A bayesi megközelítés

Éljünk azzal naiv feltevéssel, hogy a generált, egyenetlen terep valójában egy normál zajjal terhelt sima terep, ahol az egyes mezők függetlenek egymástól. Ez durva egyszerűsítése a problémának, a számításokat viszont jelentősen megkönnyíti, hiszen ekkor a \mathbf{H} mátrix eloszlása a következőképpen írható:

$$p(\mathbf{H}|\mathbf{U}, \mathbf{V}, \beta) = \prod_{i=1}^I \prod_{j=1}^J p(\mathbf{H}_{ij}|\mathbf{u}_i, \mathbf{v}_j, \beta) = \prod_{i=1}^I \prod_{j=1}^J \mathcal{N}(\mathbf{H}_{ij}|\mathbf{u}_i^T \mathbf{v}_j, \beta^{-1}), \quad (1)$$

ahol \mathcal{N} jelöli a normál eloszlást. Az első egyenlőség azt állítja, hogy a \mathbf{H}_{ij} értékek függetlenek egymástól; a második pedig azt, hogy \mathbf{H}_{ij} értéke az előbb látott $\mathbf{u}_i^T \mathbf{v}_j$ szám körül mozog β^{-1} szórással. A bayesi megközelítésben szükségünk van a $p(\mathbf{U})$ és $p(\mathbf{V})$ priorokra is; tudjuk, hogy ha ezeket szintén normál eloszlásnak választjuk, a poszterior is ilyen típusú lesz („konjugált prior”), azaz legyen

$$p(\mathbf{U}|\alpha_u) = \prod_{i=1}^I p(\mathbf{u}_i|\alpha_u) = \prod_{i=1}^I \mathcal{N}(\mathbf{u}_i|\mathbf{0}, \alpha_u^{-1} \mathbf{I}), \quad (2)$$

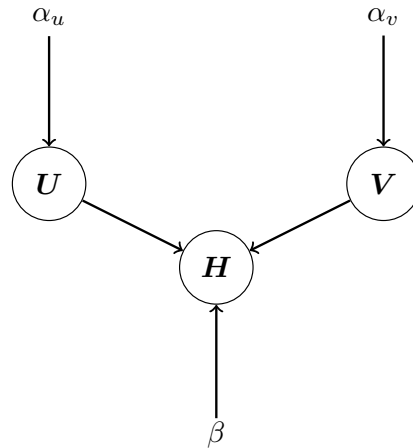
$$p(\mathbf{V}|\alpha_v) = \prod_{j=1}^J p(\mathbf{v}_j|\alpha_v) = \prod_{j=1}^J \mathcal{N}(\mathbf{v}_j|\mathbf{0}, \alpha_v^{-1} \mathbf{I}). \quad (3)$$

Másként fogalmazva, a függetlenség az \mathbf{U} , \mathbf{V} mátrixok oszlopvektoraira is teljesül, valamint ezek $\mathbf{0}$ várható értékű, α^{-1} szórássú körszimmetrikus kovarianciamátrixú normál eloszlást követnek. A fenti modell a 2. ábrán látható Bayes-hálóval írható le.

1.3. Következtetés Gibbs-mintavételezéssel

Feladatunk, hogy az ismert \mathbf{H} terep alapján megtaláljuk az \mathbf{U} és \mathbf{V} mátrixokat, ehhez a 2. ábrán látható Bayes-hálóban kell következtetnünk. A feladat valószínűségi természetéből fakadóan az eredmény nem is egy pontbecslés lesz \mathbf{U} -ra és \mathbf{V} -re, hanem egy egész eloszlás. Ezen eloszlás konstrukciójára az egyik legegyszerűbb, széles körben alkalmazott algoritmus a Gibbs-mintavételezés, amelynek általános leírása megtalálható a tankönyv 14. fejezetében. Az eljárás során az \mathbf{U} és \mathbf{V} változókat felváltva mintavételezzük, az aktuálisan kívül minden más változót fixálva; a két keresett eloszlás ezen két mintahalmazzal becsülhető. Így az is nyilvánvaló, hogy az algoritmus nem determinisztikus, különböző futtatásoknál más-más eredményeket kaphatunk; megmutatható azonban, hogy az iterációk számának növelésével azonban a megoldás konvergál a „valódi” eloszláshoz. Lépésenként:

1. Megadjuk az α_u , α_v hiperparamétereket. Minél magasabbak, annál nagyobb szórást engedünk meg az \mathbf{U} és \mathbf{V} mátrixoknál.



2. ábra. A modell gráfos ábrázolása. A csomópontok a változókat, a nyilak feltételes függőségi relációkat jelölnék.

2. Megadjuk a β precizitási paramétert. Minél magasabb, annál jobban szeretnénk közelíteni az értékeket az eredeti \mathbf{H} mátrixban, más szóval annál inkább „bízunk” az eredeti értékekben.
3. \mathbf{H} értékeit a zajos terep alapján állítjuk be, \mathbf{U} és \mathbf{V} értékeit pedig a prior alapján, véletlenszerűen.
4. Mintavételezzük (frissítjük) \mathbf{U} -t a $p(\mathbf{U}|\mathbf{V}, \mathbf{H}, \alpha_u, \alpha_v, \beta)$ eloszlás alapján.
5. Mintavételezzük (frissítjük) \mathbf{V} -t a $p(\mathbf{V}|\mathbf{U}, \mathbf{H}, \alpha_u, \alpha_v, \beta)$ eloszlás alapján.
6. Ismétljük a 4-6. lépéseket; az első meghatározott számú mintát eldobjuk („burn-in”), a továbbiakat rögzítjük, amíg el nem érjük az előre megszabott maximális iterációs számot. Végül, mivel valójában csak pontbecslésre van szükségünk, a rögzített minták átlagát vesszük.

Megfelelő paraméterezéssel a fenti sémát „elég ideig” iterálva olyan megoldáshoz jutunk, ahol az $\mathbf{U}^T \mathbf{V}$ szorzat valóban jól becsüli \mathbf{H} -t. A számításigény csökkentésére több lehetőség is adódik, de ezekre most nem térünk ki.

1.1. Állítás. Az (1), (2), (3) egyenletekkel definált modell esetén \mathbf{U} feltételes eloszlása szintén normál, a következő paraméterekkel:

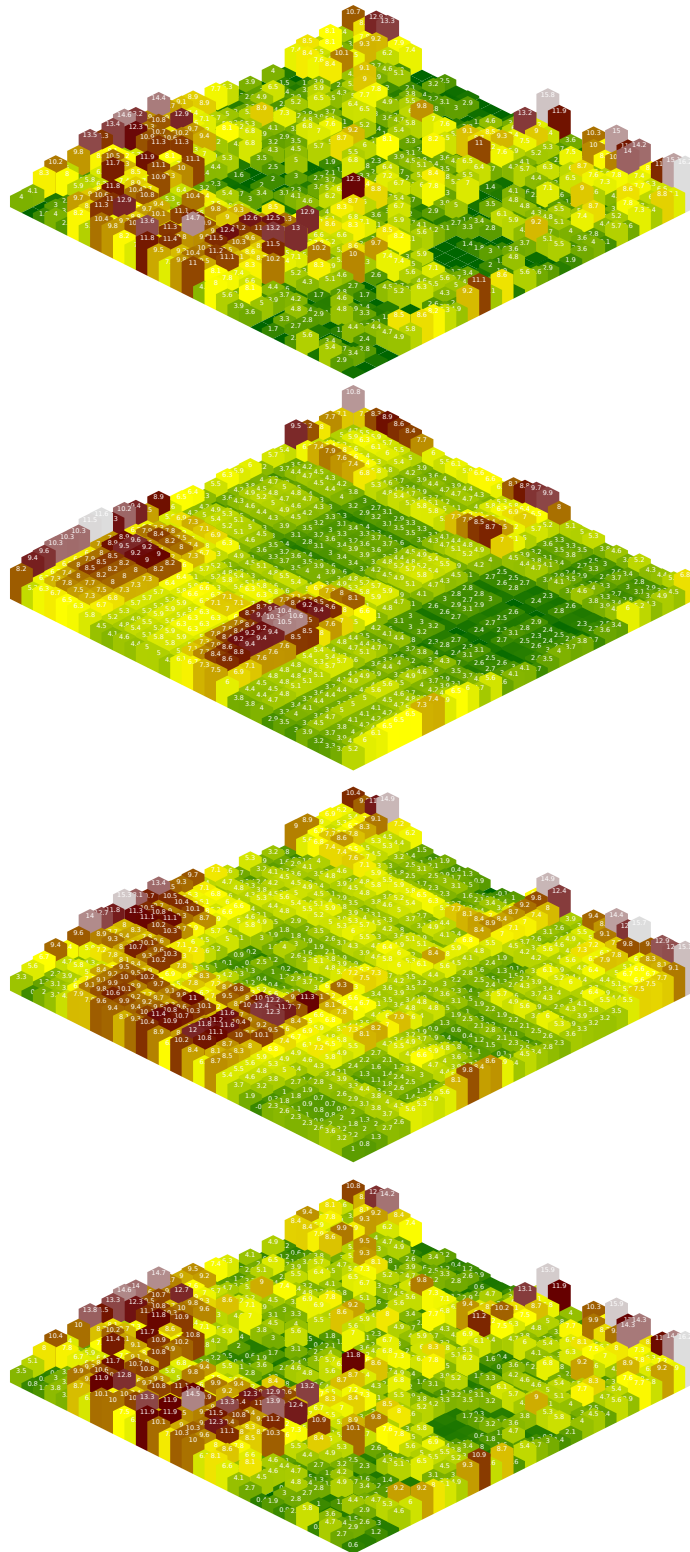
$$p(\mathbf{U}|\mathbf{V}, \mathbf{H}, \alpha_u, \alpha_v, \beta) = \prod_{i=1}^I \mathcal{N}(\mathbf{u}_i | \boldsymbol{\psi}_i, \boldsymbol{\Lambda}_i^{-1}),$$

$$\boldsymbol{\Lambda}_i = \beta \sum_{j=1}^J \mathbf{v}_j \mathbf{v}_j^T + \alpha_u \mathbf{I},$$

$$\boldsymbol{\psi}_i = \boldsymbol{\Lambda}_i^{-1} \beta \sum_{j=1}^J \mathbf{H}_{ij} \mathbf{v}_j.$$

Bizonyítás. Lásd a Függelékben. □

\mathbf{U} mintavételezéséhez tehát sorra ki kell számítani a $\boldsymbol{\Lambda}_i$ mátrixokat és $\boldsymbol{\psi}_i$ vektorokat a fenti képletek alapján, ezek ismeretében az egyes \mathbf{u}_i oszlopok egymástól függetlenül frissíthetők. \mathbf{V} mintavételezése hasonlóképpen alakul, csak \mathbf{u}_i és \mathbf{v}_j szerepet cserélnek.



3. ábra. Eredeti terep és simított verziói $L = 1$, $L = 3$ ($\alpha_u, \alpha_v = 1$, $\beta = 100$), valamint $L = 20$ ($\alpha_u, \alpha_v = 10$, $\beta = 1000$) paraméterekkel. Látható, hogy utóbbi már igen pontos közelítést ad. Minden esetben 20 burn-in után 100 iterációt használtunk.

2. Feladatok

Implementálja a Gibbs-mintavételezést Java nyelven. A lineáris algebrai számításokhoz és a normál eloszlásból való mintavételezéshez használja az Apache Commons könyvtárat (<http://commons.apache.org>), ezen kívül kódot ne emeljen át külső forrásból. A megoldás tartalmazzon egy Main osztályt, ezen belül egy main() függvényt. A bemeneti mátrixot a standard inputon várja, a kimenetet a standard outputra írja. A program forráskódját zip file-ba tömörítve töltsse fel.

2.1. Bemenet

A bemenet az I, J, L számok vesszővel elválaszva, utána a β paraméter, majd \mathbf{H} magasságtérkép csv formátumban. Példa:

```
5,6,2,100.0
0.20,0.50,2.51,1.12,6.98,4.43
1.08,0.20,1.87,1.31,1.44,3.96
0.67,4.29,4.48,5.68,6.03,3.49
3.93,1.64,1.84,3.32,2.08,4.63
4.89,7.30,3.20,2.38,3.04,0.73
```

(ennél nagyobb, kb. 50×50 mátrixokra számítson).

2.2. Kimenet

A kimenetre az \mathbf{U} és \mathbf{V} mátrixokat írja csv formátumban, egy üres sorral elválaszva. Példa:

```
0.25,2.65,1.61,7.63,2.88
3.79,2.31,5.40,3.88,6.00

1.68,5.91,3.38,3.24,0.57,2.61
2.37,7.10,5.43,0.05,0.01,3.14
```

2.3. Értékelés

Az értékelés során a \mathbf{H} és a $\mathbf{U}^T \mathbf{V}$ mátrixok elemenként vett négyzetes eltéréseinek átlagát vesszük alapul, több tesztesetre kiszámítva. Tanácsos ezt a megoldás részeként is kiszámítani és figyelembe venni a paraméterek finomhangolásánál. A sikerhez minden tesztnek meg kell felelni; sikertelenség esetén természetesen lehet újra próbálkozni. Tanácsos az algoritmust több, véletlenszerűen generált tesztesettel kipróbálni.

Függelék

Az 1.1. Állítás bizonyítása. A bizonyításhoz először felírjuk az együttes eloszlást. Ehhez kell a Bayes-tétel:

$$p(\mathbf{H}, \mathbf{U}, \mathbf{V} | \beta, \alpha_u, \alpha_v) \propto p(\mathbf{H} | \mathbf{U}, \mathbf{V}, \beta) p(\mathbf{U} | \alpha_u) p(\mathbf{V} | \alpha_v)$$

Ennek vesszük a logaritmusát, így a szorzatok helyett összegeket írjatunk. Ezenkívül az \mathbf{U} -t nem tartalmazó tagokat konstansnak vesszük. Így az (1) és (2) egyenletek alapján

$$\ln p(\mathbf{H}, \mathbf{U}, \mathbf{V} | \beta, \alpha_u, \alpha_v) = \sum_{i=1}^I \sum_{j=1}^J \ln \mathcal{N}(\mathbf{H}_{ij} | \mathbf{u}_i^T \mathbf{v}_j, \beta^{-1}) + \sum_{i=1}^I \ln \mathcal{N}(\mathbf{u}_i | \mathbf{0}, \alpha_u^{-1} \mathbf{I}) + \text{const.}$$

Most felhasználjuk a normál eloszlás képletét, így

$$\ln p(\mathbf{H}, \mathbf{U}, \mathbf{V} | \beta, \alpha_u, \alpha_v) = \sum_{i=1}^I \sum_{j=1}^J -\frac{1}{2} \beta (\mathbf{H}_{ij} - \mathbf{u}_i^T \mathbf{v}_j)^2 - \frac{1}{2} \sum_{i=1}^I \alpha_u \mathbf{u}_i^T \mathbf{I} \mathbf{u}_i + \text{const.}$$

ahol az U -t nem tartalmazó tényezőket ismét beépítettük a konstans tagba. Ebből kifejezhetjük u_i -t:

$$\sum_{i=1}^I -\frac{1}{2} \mathbf{u}_i^T \left(\beta \sum_{j=1}^J \mathbf{v}_j \mathbf{v}_j^T + \alpha_u \mathbf{I} \right) \mathbf{u}_i + \left(\beta \sum_{j=1}^J \mathbf{H}_{ij} \mathbf{v}_j \right)^T \mathbf{u}_i + \text{const},$$

ahonnan látható, hogy ezek u_i -ban kvadratikusan tagok. Így nincs más dolgunk, mint kiegészíteni teljes négyzetre, és azonnal észrevesszük, hogy a normál eloszlás kitevőjének is pontosan ilyen alakja van:

$$\sum_{i=1}^I -\frac{1}{2} (\mathbf{u}_i - \boldsymbol{\Lambda}_i^{-1} \boldsymbol{\varphi}_i)^T \boldsymbol{\Lambda}_i (\mathbf{u}_i - \boldsymbol{\Lambda}_i^{-1} \boldsymbol{\varphi}_i) + \text{const}.$$

A logaritmust visszalakítva tehát normál eloszlást kapunk:

$$p(\mathbf{U} | \mathbf{H}, \mathbf{V}, \alpha_u, \alpha_v, \beta) = \prod_i \mathcal{N}(\mathbf{u}_i | \boldsymbol{\psi}_i, \boldsymbol{\Lambda}_i^{-1}),$$

$$\boldsymbol{\Lambda}_i = \beta \sum_{j=1}^J \mathbf{v}_j \mathbf{v}_j^T + \alpha_u \mathbf{I},$$

$$\boldsymbol{\psi}_i = \boldsymbol{\Lambda}_i^{-1} \boldsymbol{\varphi}_i = \boldsymbol{\Lambda}_i^{-1} \beta \sum_{j=1}^J \mathbf{H}_{ij} \mathbf{v}_j,$$

ahol a konstans tagok automatikusan bekerültek a normalizációs tagba. □