

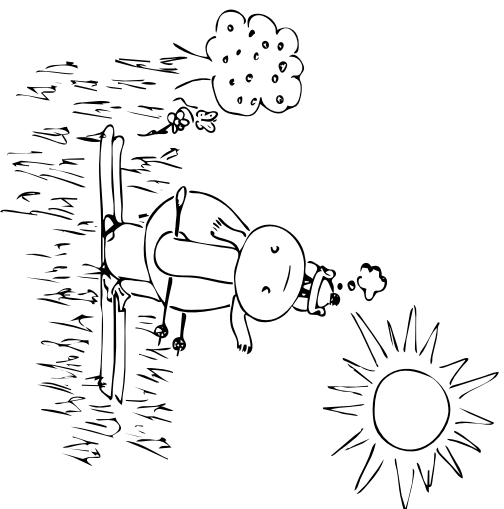
Korespondenční Seminář z Programování

28. ročník

KSP

Březen 2016

řešitel	škola	ročník	série	3-1	3-2	3-3	3-4	3-5	3-6	3-7	3-8	série	celkem
56. Ondřej Boryšek	GJarošBO	3	2									0,0	3,3
57. Adam Husník	GArabskáPH	2	1									0,0	3,0
58. Jitř Müller	G.Roundice	3	1									0,0	2,5
59. David Nápravník	GLHometPH	3	1									0,0	2,2
60.–61. Michael Bausano	GTrš	4	1									0,0	2,0
Martin Zoula	GNadKavaPH	4	4									0,0	2,0
Peter Marta	G KošiceS	4	1									0,0	1,0
62.–63. Michal Rickwood	G.CTřebová	2	1									1,0	1,0



Milí řešitelé a řešitelky!

Zitřiv! Sluníčko se směje na obloze a říká cosi o jarní. Organizátoři KSPřka kvíri pod péřnou a říkají cosi o zimním spánku. Strony opatřně vystřikují listy z pupenů. Organizátoři opatřně vytahují zpod polštáře papír s nápady na úlohy. Po nebi se prohánějí mráčky a mravíky přitáčkám. Po papíře se prohánějí tužky a sepišují zadání. Vstávat a řešit! A proč? No protože je jaro!

Při příležitosti jarní rovnodennosti (věčejší) vám pošláme zadání čtvrté série, jak stvořené ke čtení pod rozkvetlým stromem.

Připomínáme, že každému řešiteli, který v tomto ročníku z každé série dostane alespoň 5 bodů, darujeme propisku, blok, tužku, a možná i něco navíc.

Termín série: Pondělí 2. května 2016 v 8:00 SELČ (CoDEX má termín stejný)

Odevzdávání: Přes web na adrese <https://ksp.mff.cuni.cz/submit/>.

Odměna série: Čokoládu pošleme každému, kdo z níh této série získá alespoň 42 bodů.

Čtvrtá série dvacátého osmého ročníku KSP

Příběh pěti domů

Udávují poselství šroub, ještě pro jistotu naposled změní napětí, balím si nářadí a vyprázdním zptčky domů. Už se tu nedví zahrávat ani minutu. Je pátek po deváté večer a mě ještě doma čekají přípravy na víkend a taky jsem dětem sblíh pohádku.

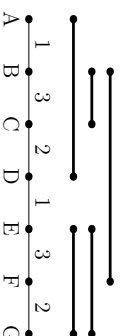
Pracuji jako technik pro jednu telefonní společnost. Práce to není špatná, docela mě i baví a khamé mě v běžném životě příliš nezabíje. Akorát když v mém okruhu nastane náhly výpadek, jako třeba teď, tak to musím hned jít opravit. Ale co nadělán, aspoň že se to tentokrát stalo přímo v našem bloku, tak jen stačí projít ulici a jsem doma.

28-4-1 Sledování telefonů 9 bodů

V ulici stojí v řadě N domů. Každý z nich má jeden telefon a je propojen telefonní linkou s dvěma sousedními domy (jedním v případě krajních). Grafové řečeno tvoří telefonní síť cestu: vrcholy představují domy a hrany spoje. Pokud zavoláme z domu a do b , musí hovor projít přes všechny spoje ležící na cestě mezi a a b .

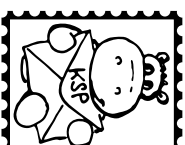
U každého spoje víme, kolik přes něj za posledního týdne prošlo hovorů. Na základě této informace bychom chtěli zjistit, kdo komu volal. To samozřejmě nelze určit jednoznačně, stejný provoz na spojích může vytvořit mnoho různých řešení. Vyberte to s nejmenším celkovým počtem hovorů (pokud je více takových, libovolně z nich).

Například pro $N = 7$ a počty hovorů na spojích postupně 1,3,2,1,3,2 muselo proběhnout nejméně pět hovorů a mohl vypadat například takto:



Tedy probably hovory $A \rightarrow D$, $B \rightarrow C$, $B \rightarrow F$ a $E \rightarrow G$ a znovu $E \rightarrow G$.

S menším počtem hovorů nelze vytvořit zadané vyřazení jinak.



① **Leďá varianta (za 5 bodů):** Řešte za předpokladu, že se počty zacyčených hovorů na libovolných dvou sousedních spojích liší maximálně o 10.

V **me ulici stojí pět domů** a při takové večerní procházce si aspoň zns jednou uvidím obrázek, jak se dří sousedům. Věciami zde žijeme už téměř deset let a k našemu nastěhování se užé společný příběh.

Bylo mi tehdy 25 let a o tomto roce můžu jednoznačně mluvit jako o roce smůly. Firma, ve které jsem byl zaměstnán, prodávavla, a tak musela částit. Jako nezkušený mladík jsem to odnesl já a další práci dlouho nemohl sehnat. A aby toho nebylo málo, tak mi navíc o pár týdnů později uvěřovali být a stopy zbyly tím, že ho prosit započít. Pochytil se dopadnout nepodatli, vedoucí kamery někoho nezaznamenaly a ani nebyly patrné jakkoliv známky umění. Prosit zohada a pojištění jsem nebyl. Takže jsem na jednou neměl ani práci, ani kde bydlet, a nápad já z toho ven, už vůbec ne.

Večer, když jsem svůj žal zapíjěl v místním lokle, neměl jsem tohž kam jít, ani jak se vyprázdnit, co mě trápí. Stručně jsem mu popsal svou situaci, soslil hlavu a zhrzele seděl dál. Opravdu jsem neměl náladu se s někým vykecávat.

Muž si na papír nacináral pár poznámek a pak povídal: „Pracuji pro společnost jménem Druhá šance a mám pro vás nabídku. S kolegou zrona zakládáme nový projekt a vy byste pro nás byli ideální kandidáti. Vaší účasť byste uvěřili řešení problémů s prací a bydlením, které vás momentálně sužují. Rozhodnutí je teď na vás. Pokud byste měl zájem se dozvědět více, zavolejte na toto číslo a domluvíme si schůzku. Noc ale nechtejte, naše prostředky jsou omezené a mohli bychom místo vás sehnat někoho jiného.“

Předal mi svou vizitku, rozloučil se a s úsměvem odešel. O nabídku jsem dlouze nepřemýšlel a hned ráno jsem volal, že přijímám. Schůzka byla další pondělí.

Přišel jsem na zadanou adresu a byl usazen do čekárny. Bylo nás tam celkem pět, tři muži a dvě ženy. Na pohled jsme všichni byli ve věku kolem pětáctácti let. Z kanceláře vyšel muž, kterého jsem již znal, v ruce držel nějaké papíry a povídal: „Výborně! Tak jste tu všichni. Vyprávějí jsme

laci. Vypočítej je ale stále pomalý a je jen malá šance, že tak hladově řešení překonáme. Hlavní důvodem si mnešl trochu děle počkat.

Co udeláme nyní, je, že znovu použijeme mrvavenčí algoritmus, ale omezení umožníví měst, že kterých budeme vyhledávat následně. Jelikož se pohybujeme v rovině, tak se nám takřka nikdy nemyplati jí moc daleko. Další město můžeme tedy vyhledat například z nejbližších K , kde K bude rozměrně malé. Treba 10.

Tím se algoritmus dost zrychlil. Už si můžeme dovolit použít více mrvavenčí. Použijeme jich aspoň 100 až 1000, aby v jedné iteraci zvládl označit rozumné možnosti hran feromony a další iterace tak měly větší variabilitu.

Pokud takové řešení s parametry $\alpha = 2$, $\beta = 2$, $\rho = 0.5$ a pro 1000 mrvavenčí nechtáme běžet 500 iterací (několik hodin), dostaneme se na řešení o hodnotě 20 550 678, což není o moc lepší než hladový algoritmus, ale aspoň něco. Tento algoritmus můžete vidět ve združeném kódu.

Já pak ještě experimentoval s tím, že feromony zanechává jen W nejlepší mrvavenčí z jedné iterace, ale dostal jsem se tak jen na víceméně stejně dobrý výsledek.

Program (C++):
<http://ksp.mf.f.cuni.cz/viz/28-3-8.cpp>

Další možné postupy

Evoluční a jin podobné algoritmy jsou většinou navrženy obecně, aby se jimi dala řešit širší škála problémů. Pro-

to neuvěřivaj charakteristiky žádného konkrétního z nich. Na jednu stranu je výhodou, že se dají prakticky vždy nějakým způsobem použít. Na druhou stranu ale mají nevýhodu v tom, že jsou málokdy nejlepší. Když se totiž zaměříme na konkrétní problém a navrhneme algoritmus přímo vhodný pro něj, tak obvykle dosáhneme vyššího úspěchu.

Kromě toho existují ještě heuristiky, které se snaží vylepšovat nějaké již existující řešení. Pro takové heuristiky se hodí jako počáteční řešení zvolit buď výsledek nějakého hladového algoritmu nebo právě výsledek některého z evolučních algoritmů.

Pro náš problém se dají použít například následující heuristiky:

- Prohodíme 2 hrany. Pokud se řešení zlepší, necháme je tak, jinak je vrátíme zpátky. (Díky této heuristice se postupně zbavujeme i křížení hran)
- Prohodíme 3 hrany. Pakliže už nefunguje prohazování dvou hran, můžeme začít prohazovat 3, což nám dává další možnosti.
- Přeházení K po sobě jdoucích měst. Také může pomoci.

U všech těchto heuristik nově řešení přijímame pouze pokud nám heuristika pomohla. Tyto heuristiky uvádíme jen pro příklad. Ty však nesoúvisí s evolučním algoritmy. Pomocí nich se řešení dalo vylepšit až na 18 000 000.

Karel Tesar

dětema kuchyněmi, třemi koupelnami a ložnice pro hosty je samozřejmostí. Interieri je navíc zkrášlen řadou obrazů historického i moderního umění...“ představenol Dalimír.

28-4-4 Podivuhodný obraz

12 bodů

V domě visí podivuhodný obraz. Je na něm nakreslených N bodů, které jsou spojeny dohromady M čarami. Celý obraz je černo-červený a má zajímavou vlastnost. Pokud z bodů vedou alespoň 2 čáry, některá z nich je černá a některá červená.

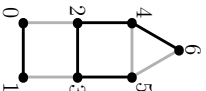
Co kdyby ale obraz vypadal jinak? Sel by pořádek takto nakreslit? Vmyslete algoritmus, který na vstupu dostane nameranou graf a obarví jeho hrany dvěma barvami tak, že každá vrchol se dotýká hran obou barev (případně zjistíte, že to nejde). Pozor: vstupní graf nemusí být souvislý.

🔗 **Lehká varianta (za 8 bodů):** Řešte za předpokladu, že všechny vrcholy mají sudé stupně.

Formální vstup: Na prvním řádku dostanete dvě celá čísla N a M udávající počet vrcholů a hran. Každý z následujících řádků popisuje jednu hranu pomocí dvojice čísel vrcholů (vrcholy číslujeme od nuly).

Formální výstup: Vypíše celkem M řádků popisujících barvy hran ve stejném pořadí, jako byly na vstupu. Barva každé hrany je buď 0 (černá), nebo 1 (červená), na obdržku šedá). Pokud graf nejde správně obarvit, vypíše jediný řádek obsahující číslo -1 .

Ukazkový vstup:	Ukazkový výstup:
7 9	0
0 1	1
0 2	1
1 3	0
2 3	0
2 4	0
3 5	1
4 5	0
4 6	1
5 6	1



Toto je praktická open-data úloha. V odevzrávacím systému si můžete vygenerovat vstupy a odevzdaté příslušné výstupy. Zaleží jen na vás, jak výstupy vytvoříte.

„To zní úplně jako sen!“ vyřičkla nadsazená dáma v klubovně.

„A to ještě není všechno,“ pokracoval Dalimír, „stanut se o takový dům není snadná. Proto také od nás máte k dispozici komorníka, uklízečku, zahradařku, osobního šoféra a hlídače k vašim službám. Součástí je také finanční dar, který by měl stačit na veškerou údržbu na alespoň dalších pět až deset let.“

„A tento dům se všemni jeho výhodami jsme připravili... Chcete napřít... Pro vás, slečno!“ ukázal na dámu v klubovně.

„To... to je naprosto úžasné,“ rozpučité dokonala dáma, „určitě se budu snažit vás neklamat a budu domu dělat dobré jméno. Určitě vám dokážu, že si jej zasloužím.“

„Ano, ano. Umytí na vás čeká komorník, který Vám vše ukáže,“ ukazuje na vchod Dalimír, „tak se běžte seznamit a my osluní budeme pokračovat dál.“

Tak to je teda husťj, říkám si pro sebe. Jestli všechny domy budou takovy, tak jsem toudné za vodou.

Nyní přicházíme ke druhému domu,“ znova mluví Dalimír, „tento dům se pšístí těmi nejnmodernějšími technologiemi, které doposud lidstvo vyvinulo. Nejen že budete mít

k dispozici ty nejlepší počítače a další hračky, ale ještě k tomu vám každý rok budeme dodávat nové. Společnost vám bude dělat domácí roboty, kteří za vás vyspí, automatiky ovládaní droni, kteří doletí ke dveřím pro poškozené záslabky, vyzvednou krabici s jídlem, a čeká na vás ještě řadu dalších technologických vychytávek. Navíc je všechno centrálně ovládané a synchronizované. Když na to přijde, tak celý den můžete pospat u počítače, a přitom se o všechno zokladnete postarat. Vysokorychlostní internet a pokrytí Wi-Fi v celém areálu je samozřejmostí.“

„Tak to už se nemůžu dočkat, až se nastěluju, to je přesně pro mě,“ usměla se ironicky hospořka. Všichni už jsme tušili, kdo bude novým vlastníkem domu. Ujahněnou oči-kou zouchly svitli oči a pozorně poslouchal každé slovo jako do té doby nikdy.

Vtom se zablasklo a přel domem se objevil umotlý muž s mečem v ruce. Na jeho rukoce a nohou publikovali první morderó světa. Pár vítěin tam stál, rozhlížel se, ale pak se zase zablasklo a byl jít.

„A teď jste mohli vidět... To byla asi... ukážka automatického zustršování pomocí holografické strážě,“ pokracoval Dalimír, „neml to ale jediný bezpečnostní prvek, který zde je. Celý areál je pokrytý kamerami, které zemnit můžete sledovat z libovolného přístroje a vstoupit můžete jen po identifikaci svou oční duhovkou.“

„Tak už asi tušíte, tak tento dům je přímo dělaný pro vás, pane,“ ukázal na očička a ten se zamouval. „Já... nemůžu se dočkat až si všechny tyhle super věci vykouším a pořádně nakouřujujít. Tenhle dům je to nejlepší, co jsem zatím ve svém životě viděl!“

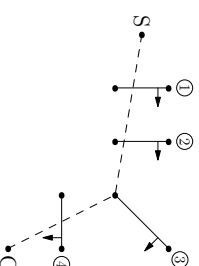
„Běžte ke vchodu. Tam vám kolega naskenuje duhovku a provede další inicializaci. My budeme pokračovat dál,“ povůdla Dalimír.

„Nyní přicházíme ke třetímu domu. V tom se určitě nikdy nudit nebudete! Umníť najdete bowlingové drátky, kalceřnky, saunu, minního s vřitkov, venkovní party bazén, minigolfové hřiště a kroleťové hřiště. Ať budete chtít podniknout cokoli, nikdy nebudete muset chodit moc daleko...“

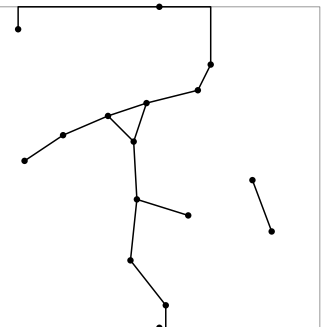
28-4-5 Hra kroket

11 bodů

Ve hře kroket přesouváme míček po hřišti pomocí několika iderů holi. Míček považujeme za bod nulové velikosti, který se po úderu pohybuje po úsečce. Celá trasa míčka tedy tvoří lomenou čáru. Na hřišti jsou také rozmištané branky, které si budeme představovat jako úsečky. Cílem hry je, aby míček projel všechny branky, a to v předepsaném pořadí. Navíc každá branka má určený směr, kterým je třeba jí projít.



Na vstupu dostanete startovní a cílovou pozici míčku a seznam branek (každá je určena svými krajními body). Najdte nejkratší trasu míčku, která začne na startu, skončí v cílovém bodě a projede všechny branky správným směrem a ve správném pořadí.



Štaci tedy zjistit, zda ℓ a p leží v téže komponente souvislosti. To můžeme otestovat třeba prohloubkáváním do šířky. Jistou nevyhnutelnou ovšem je, že náš graf má nejvýš $O(S^2)$ hran (rozmyslete si, jak by faktory les vypadaly). Bude tedy lepší změřit si celý graf v paměti a vyhodnotit ho podle potřeby za běhu.

Budeme si udržovat frontu stromů, které máme zpracovávat. Na počátku do ní vložíme stromy dosažitelné z ℓ , čím ty, jejichž x -ová souřadnice je menší než D . Pak postupně odibiráme stromy z fronty. Pro každý strom najdeme všechny jeho sousedy; tedy stromy ve vzdálenosti menší než D , a pokud jsme je ještě neviděli, přidáme je do fronty. Skončíme, pokud se fronte vyprázdí, nebo pokud objevíme strom s x -ovou souřadnicí větší než $L - D$.

To je všechno. Šmahdň, že? Ale pomalé... Navštívíme až S vrcholů a pro každý z nich prohledáme S potenciálních sousedů. Jehoňho souseda našetřit stiháme zpracovávat v konstantním čase (myslo vzdálenuost porovnávání jejich druhé mocniny; abychom nemuseli odmocňovat), takže celý výpočet trvá $O(S^2)$.

Program (C):
<http://ksp.mff.cuni.cz/viz/28-3-7-quad.c>

Pomohlo políčka?

Existuje řádka technik, kterými se dá hledat sousedních stromů urychlit. Můžeme například rozdělit les na políčka velikosti $D \times D$ a pro každé políčko si předpočítat, jaké stromy v něm leží. Souseďe daného stromu se pak nacházejí buď ve stejném políčku, nebo v některém z osmi okolních.

Tento trik může pro některé kombinace L a D výrazně pomoci. Limity v zadání úlohy bohužel připouštěly i případy, kdy je D řádově stejné jako L . Tehdy je políček konstantní počet, takže okolí políčka obsahující podstatnou část všech stromů. Škoda, tudíž cesta nereálná.

Znovu zametání

Není náhodou, že jsme k této sérii přihlíželi geometrickou kuchařičku, v níž se ukazuje princip zametání roviny. Hodi se i pro tuto úlohu.

Les budeme zametat zleva doprava svišlou přímkou. Budeme si při tom pamatovat, které stromy jsme už pokládali. Pokaždě, když zametací přímka narazí na nový strom, najdeme jeho sousedy mezi předchozími stromy a natáheme do nich hrany. Hledání si urychlíme dvěma triky:

- 12 Pomale log * říká, kolikrát musíme číslo x opakovaně zlogaritmovat, než se poprvé dostaneme pod 1. Roste tedy velice pomalu. Funkci $\alpha(x)$ zde deňnovat nelze, ale pozorujeme, že roste ještě mnohem pomaleji.
- 13 <http://ksp.mff.cuni.cz/viz/kuchacky/kostry>
- 14 <http://mj.ucw.cz/vyuka/ads/43-geom.pdf>

- Nemá svunyi zkonšet stromy, které se v y -ové souřadnici liší od nového stromu o více než D .

- Pokud se dva stromy vyskytnou v téže řádce, stačí uzavřít pravý z nich. Ten levý je buďto mo delčko, nebo už leží v téže komponente jako pravý, takže stačí hranu přivést do pravého. (Nezapomeňte, že souřadnice stromů jsou celočíselné, a proto je možných řádek málo.)

Počítáme si tedy pole indexované číslem řádku $(0 \dots L - 1)$. V něm si budeme pro každý řádek pamatovat, jaký nejpravější strom jsme tam viděli. Objevil-li se nový strom v y -tém řádku, stačí se podívat na nejpravější stromy v řádcích $y - D$ až $y + D$ a pro každý z nich zkonštrovat, jak je daleko.

Sousedě jednoho vrcholů tedy prohledáme v čase $O(D)$. Celý algoritmus nejprve stráví čas $O(S \log S)$ seřazením všech stromů a $O(D)$ inicializací pomocného pole. Poté prohledá S stromů a nad každým z nich stráví čas $O(D)$. Tím vytvoří graf o nejvýše $O(SD)$ hranách, který pak prohledá do šířky v čase $O(SD)$.

Čalková časová složitost je tedy $O(S \log S + L + SD)$.

Implementační (h)trky

Náš ukazkový program se drží předchozího nápadu, ale pro jednoduchost se v některých detaalech odchyluje.

Především si místo řízení stromů podle souřadnice vytvoří matrici $L \times L$, která o jednotlivých bodech česá říká, zda tam leží strom. Znamenáťi lesa pak prostě procházá tuto matici po sloupcích a hledá jedničky. Místo $O(S \log S)$ tedy potřebuje čas $O(L^2)$. To pro limity ze zadání je spíše lepší.

Také místo toho, abychom graf udržovali v paměti a pak ho celý prohledali, udržujeme komponenty souvislosti přibližně v datové struktuře Union-Find. Pro každou z $O(SD)$ hran tedy voláme Union, který trvá $O(\log S)$. Kdybychom byli pilnější, použili jsme rychlejší Union-Find se složitostí $O(\log^* S)^{12}$ nebo $O(\alpha(S))^{13}$, který je popsany v kuchařice o minimálních kostkách.¹³

I s jednoduchým Union-Findem má náš program složitost $O(L^2 + SD \log S)$, což je pro silbené velikosti vstupů naprosto postačující.

Program (C):

<http://ksp.mff.cuni.cz/viz/28-3-7-grid.c>

Dvoutvrný dígram

Zbyvá hustka pochlypností, zda by úlohu šlo řešit efektivněji, zejména v případě, kdy by nám nikdo nestihl celoušně souřadnice. Pochlypnosti samozřejmě nejsou na místě, známe řešení v čase $O(S \log S)$ s použitím jednoho miliónu kraňka z kompehndičného kloubného výpočítacího geometrii. Jan je to trochu pracnějši, takže vše pouze načrtáme. Detaily můžete najít v geometrické kapitole metvýčí knihy.¹⁴

Necht S je množina stromů v rovině. Pro každý strom bodu $s \in S$ definujeme jeho *oblast* R_s jako množinu všech bodů roviny, pro něž je s nejbližší strom (respektive jeden z nejbližších stromů, pokud jich je víc).

Podívejme se na nějaká dva stromy x a y . Všechny body roviny, které mají k x bližě než k y , tvoří polovornu (hranici přímku této polovorny je osa úsečky xy). Proto musí každá oblast R_s být přímkuam konečné mnoha polovornu.

může doprít pocít fungující robný a namísto matýského objekt jim prougnána chuou, aby se mohla napho věrnout smému útesu a nehltnu. Tuto dama dostala před deseti lety neuvěřitelné možnosti. Bohužel se ale snažla utvrdit více, než dokázala sama unest. Delšnu svou to bohužel umožňuje. Když se nad tím zamyslíte, tak už vlastně neexistuje mnoho věcí, které by se nedaly zřítit online. Jeho život byl naprosto oddan vrtulnímu reálné a moderním technologim. Bohužel vrtulnímu světu dával výraznou přednost před tím fyzickým a igitornou potřebu svého vlastního těla. A tak začal typě otehranoutm svou, zářný zážitek a nakonec se mu začala borit pádř.

Před necelým rokem ochrnul na dněi polovornu těla, když upadl na zem po cestě na záchod. Začtenáí jej až kamauadí z online hry, když se do ní dva dny po sobe nepřihlípl a ani o sobe neudá vědět. Dostal se pak k němu a poskyhnut mu lékařskou pomoc byl také problém. Projt přes všechna technologická zabezpečení, která v domě měl, by byl úkol nímáhné pro CIA. Ještě že měl jen obyčejná skleněná okna. Po tomto incidentu mu byl přidělen osobní pečovatel. Ten se následně stal jeho nejbližším přítелеm ve fyzickém světe za všechny ty roky. Vlastně byl zároveň jediným člověkem, který trávil v jeho přítomnosti více jak deset minut denně. Když jsem mýřel jeho dnu, zrovna se sanitka rozhoukala a začala odjíždět.

28-4-7 Jizda sanitkou

Sanitka má naloženého stabilizovaného pacienta a potřebné je její bezpečné odvězť z jeho domu do některé z nemocnic. Ve městě ale probíhají na různých místech opravy. V jejichž okolí to nebezpečně dmná a navíc poblíž nich hrozí úvazní v zácpě. Proto by se řídit řád dtžel od míst opravy co nejřál. Máte zadanou mapu města jako čtvercovou mřížku s vyznačenou počáteční pozicí, pozicemi nemocnic a pozicemi všech míst, kde probíhají opravy. Najdte největší K takové, že existuje cesta ze startu do nějaké nemocnice, která se po celou dobu drží ve vzdálenosti alespoň K od všech míst opravy. Zároveň také najdte libovolnou cestu splňující toto omezení.

Po mapě se polyhujeme pouze vodorovně a svisle. Vzdálenosti měříme v manhattanství metrice, tedy jako součet rozdů x -ových a y -ových souřadnic. Například políčka $(1, 2)$ a $(5, 3)$ mají vzdálenost $4 + 1 = 5$.

Pro následující mapu (S značí start, N nemocnici a X místo opravy) je největší $K = 2$ a jedna z možných optimálních cest je vyznačena šedě. Šmahno nahlédnete, že pro $K = 3$ žádná cesta neexistuje.

	N					
			X			N
				X		
					X	
	N		X		S	

Ve třetím domě byl hlíd. Neudá se to vůbec sromatnat s tím, co se tam dělo před pěti až deseti lety. Byvala tam dvoká party mamináhné každý druhý den, které se bez výjimky potahovala minimálné do pozdního rána dalšího dne. Já sám jsem se tam také dvakrát vydal a můžu říct, že to bylo super! Pak jsem ale další dva dny téměř jen ležel, spal a sřtřáhněl. Takovouhle obca jsem dopkzdal přežití maximálné dvoúdně měřící, tak je pro mě nastrojo nepochopitelné, že to ta zelenonlosu, dnes už černoňosá, hispěrka zvuáála táhnout přes čtři roky v kuse žhrba čtyřřátit týdně.

Celou dobu to vypadalo, že si vše maximálně užívá. Pak ale jednou z náčelo nuc, když se púřví zase dobře rozjela, proskočila zavřejm oknem a dopadla přímo do bozenu. Nic všázňeho se jí nestálo a nklou nepřišle neví, proč to vlastně udeřelo. Nikdo jí v to době ještě dost bláží, jeho dopkzdal něco šašit, nadoř pak přepřonášat. Z tohoto incidentu se ale nedokázala vzpamotovat a další skoro dva roky strvala na psychiatrickém oddělení. Asi to houka s tou volností, sponátosť a nevzácností přehnála. Zlá napho každý moment a jen pro ten moment. Už si ale nenašla dvířku, aby se sama ohrála odksu, kam, za čím, pro koho a proč jde. A tak se jí typle nevyřešená a odkládané pocity hraváthly tak dlouho, až to nageňnou vybuchou.

Dnška ale žije docela jiny život. Nášla si stáloho přítelce a společně přezímají jejího domu přestavěti v nočního klubu na moderní školbu. Jsme spolu přítele a občas se navštěvujeme. Myslím si, že navážou svému dřívějšímu mláďi je to celkem fajn a rozumná houka.

Ted už se blížím ke čtvertemu domu. Tam se loho za těch deset let moc neměnilo. Přítelkám tam žij podimný věčný student, jehož život byl řádně přisňjání pravdy. Druce tam žije podimný věčný student, jehož život je řádně ještě přisňjání pravdy. Akorát ngm na svou sňtku tihlou oblacovnou o mohe dluži obory. Každý den veškerá a chodí spát ve stejnou dobu, má pečlivě navrženy jákničky, třkákat týdně cvičí a vůbec v jeho životě není žádná nepravidelnost. Namísto kamarádů má pouze spoluhpracovníky a spoluhadatelce.

Ze začátku jsem jej púřřví klusit nekam pozvat, ale v jeho roznou se těžko dá hledat minuta, který by nebyla zabraňna v rámci jeho složitého týdenního režimu. Ted o něm absolutně nebolžíš říct, zda je v životě šťastný a zda má to, co hledal. Nerovnomám mu, nemáme si spolu co říct, emoce na sobě neudává znít, ale třeba je tak šťastný.

Ted už se konečně blížím ke svému domu. Tam na mě čeká manželka a mě dle brsné děti, pro které jsem na zhradě v posledním roce postatal malé dětské hřiště. Myslím, že ted žiju spokojený život, ve kterém ještě stáhu jdu za svým clem. Posledních deset let pro mě nebylo žijé jednoduchých. Masei jsem se hodně snažil v práci a dělat přesčasny, abych získal ujši kvalifikaci a zvedl mi plat. Domu nám zas nedarmo prasklo potrubí a děti jsme museli složité dovažet k bobčíkam a vůbec celkem často musíme řešit takové náhlé problémy. Ale zatím jsme to vždy všechno nějak zvládli.

Jsem opravdu rád, že mi kdysi byl přidělen právě tento dům a ne žalný jiny. Jsem opravdu spokojený. V tom se tedyi Dabmar třelí.

Všem mým souseďm se v ten den spínal jejich sen a začali si hlavě vnuovat a po nějaké době se toho přesřáli. Nikdy v životě se nemouřili, jak o věci v životě bojuvat a jak se o své celé snažit. Já oproti nám dostal jen pokrněny, které mě umožňovaly si za svůjí cti jít a snažit se je plnit. Samotný fakt, že se mi to daří a postupně za mnou jsou lidé vysledky, mě naphyuje opravdovym štěstím. Takovým

tím početem zadosťuchčentů, který mají sousedé dlanio nebo dokonce nikdy nemějí šanci poznat.

Oaži do nadržce dostali snůl cñl a o tom okamžiku se pře-stali zajímat o jakoukoliv cestu, která by je někým mohla dovést. Já jsem dostal možnost skončit si svou užasnú cestu a jít po ní. A to je přesně ono, protože právě ta cesta je můj cíl! Avo, je to tak: „Cesta je cíl!“

Příběh, pro vás napsal

Karel Tesoř

28-4-8 Strojoové učení 20 bodů

Co je strojoové učení? Lidské učení je schopnost adaptace na nové situace. Když budeme poprvé hrát šachy, asi nám to nepůjde moc dobře, ale po chvíli to půjde trochu lépe, a když budeme hrát dost dlouho, můžeme se stát experty. Můžeme se naučit hrát dobře šachy i když na začátku dostaneme jenom jejich pravidla, která nám umožňují dělat špatné i dobře tahy. Někdo nám nedá přesný postup, jak být expert – v naší hlavě hravě šachní vznikne schopnost hrát je dobře.

Strojoové učení se snaží replikovat tuhle schopnost v počítačích. Je velmi užitečné umět řešit problémy, na které neznáme žádné jasné algoritmy. To platí obzvlášť, když se snažíme automatizovat nějaký aspekt lidské inteligence. Jedna praktická aplikace je třeba *počítačové vidění*. Dnes počítačové umění čítá psaný text, rozpoznávání lidské tváře, nebo třeba hledat ve fotkách dopravní značky, a to všechno s nadsílkou přesnosti. Algoritmus strojoového učení například dostane 5000 příkladů 10 různých dopravních značek a jeho výsledkem bude postup, jak je od sebe rozpoznať.

Většina materiálů o strojoovém učení je buď v angličtině nebo používá anglickou terminologii. Aby pro vás bylo snadnější si to zájmu dohledat víc informací, budeme české terminology zavádět i s jejich anglickými ekvivalenty. Jestli byste chtěli strojoové učení včlenovat víc samostatná, mžžete si třeba najít materiály ke kurzu Machine Learning na Coursera.²

Druhy strojoového učení

Strojoové učení se dělí na tři široké kategorie: *učení s učitelem* (*supervised learning*), *bez učitele* (*unsupervised learning*) a *zpeřnavaovací učení* (*reinforcement learning*).

Zmíněný problém klasifikace dopravních značek patří pod *učení s učitelem*. Jakýsi učitel nám ukázal příklady a řekl nám, „tohle je stopka“, „tohle je zákaz vjezdu“, a tak dále. Naším úkolem je podle těchto *trénovačích dat* vytvořit algoritmus, který bude fungovat dobře nejen na příkladech, které jsme dostali od učitele, ale i na těch, které jsme ještě nikdy neviděli.

Učení bez učitele se snaží najít nějaké pravidelnosti, ale nemá zvenku zadáno, čím se takové pravidelnosti budou vysztavovat. Když máme nějakou sadu dat, o které nic nevíme, učení bez učitele nám třeba mžžže pomoci najít nějaké jeřhky, „významné vlastnosti“, na které se pak mžžžeme zaměřit zvlášť. Mnozí učení bez učitele patří třeba *detekce anomálií*, která hledá ve vstupních datech vzorky, které významně vybočují z „pravidelné struktury“. Když třeba máme datacentrum plně serverů, mžžžeme se snažit detekci anomálií najít servery, se kterými je něco nějakým způsobem špatně, i když předem nevíme, jakým způsobem se mohou porouchat a jak se to projeví v parametrech, které měříme.

² <http://coursera.org/learn/machine-learning>, nejblížší termín kurzu začíná 21. března.

Zpeřnavaovací učení je o něco speciálnější. Používá se na učení chování v prostředí, ve kterém nemusí být podle výstupu nautčeného systému hned jasné, jestli se dává čtyřte.

Představte si třeba, že se snažíme naučit hrát Pacmana. Algoritmu říkáme, jak vypadá labrynt, kde se Pacman nachází, kde jsou dutkové a kde je jídlo a on nám říká, kam dce, abychom šli. Chceme, aby se naučil, jak se má hřbat, aby sebral co nejvíce jídla a pokud možno nemřel. Když algoritmus řekne „jít vlevo“, tak nevíme hned, jestli to byl dlouhodobě dobrý krok, nebo špatný krok – to záleží na tom, jak se algoritmus zachová v budoucnosti (například jestli ho vítěm tohoto kroku za 5 tahů zabije duď). Pacman provádí posoupnost akcí v nějakém prostředí a snaží se, aby se nakonec naučil co nejlepší algoritmus pro hrání. Učení bez učitele na to použít zřejmě nejde (protože máme konkrétní cíl – snažíme se maximalizovat středně jídlo).

Učení s učitelem se taky nehodí. Intuitivně bychom se totiž neměli hrát dobře – učili bychom se jenom napodobovat učitele. Pokud by tedy učitel byl třeba mistr světa v Pacmanovi, tak by ho náš naučený algoritmus neměl porazit, protože se jenom naučil to, co umí učitel. Náš cíl je najít co nejlepší algoritmus, který je v rámci pravidel Pacmana možný.

Dnes si ukážeme pár základních algoritimů učení s učitelem.

Učení s učitelem

Zkusme se třeba naučit podle výšky a obvodu pasu předpovídat váhu lidí. Nejřřive najdeme nějaké dobrovolníky (necht je jich N) a posřřibáme jejich výšky, obvodů a váhy. Vstupním informacím, podle kterých se snažíme předpovídat hmotnost, řřikáme *přřiznaky (features)*. Počet přřiznaku, tedy u nás 2, označíme jako p ; první přřiznák je výška a druhý je obvod pasu. Uložíme si je do vektorů $x^{(1)}, x^{(2)}, \dots, x^{(N)}$. To, co se snažíme naučit předpovídat – tedy váhu – si uložíme do $y^{(1)}, \dots, y^{(N)}$. Všem vstupním i výstupním datům jeřř jako D .

Když třeba $x^{(8)} = (1.87, 93)$ a $y^{(8)} = 85$, tak osmý dobrovolník měl 1.87 m, měl obvod pasu 93 cm a vážil 85 kg. Informacím, které máme o jednom dobrovolníku, tedy dvojici $(x^{(8)}, y^{(8)})$, řřikáme společně *vzorek (sample)*. Používáme standardní vektorové značení, takže $x_1^{(8)} = 1.87$ a $x_2^{(8)} = 93$.

Většina aplikací učení s učitelem je *klasifikace* nebo *regrese*. Klasifikace je přřizzení vstupního vzorku do jedno z *kategorů*. Regrese je předpověď hodnoty nějaké funkce, která vede do reálných čísel. Předpovídaní hmotnosti je regrese. Učování, jaký typ dopravní značky jsme vřřifotili, je klasifikace.

Náš cíl je najít nějakou funkci f takovou, že pokud máme pro každý vzorek (x, y) je $y = f(x)$, nebo aspoň mezi nimi nebude velký rozdíl. Kromě toho ale chceme ještě jednu důležitou vlastnost: f by měla *dobře generalizovat*. Učení je hledání vhodné funkce f .

První požadavek, tedy aby f na známých vzorcích odpovídala správně, se dá splnit sponusem způsobů. Jedem z nich je třeba ten, že by si f při učení zapamatovala všechna trénovač data, a když by dostala vstup x , tak by se podívala, jestli tenhle vstup byl v trénovačích datech, a jestli byl, vrátila by jeho přřřisřřený výstup, a jinak by vrátila třeba 9999. Měle-li nepřřřijemný pooc, že na tomhle nápadu je něco špatně, máte ho zcela správně.

Máme tedy algoritmus, který nalazne plán pro daný strom začínající v kořeni a končící v prvním patře: Rekurzivním zvoláním téhož algoritmu nalazneme plán pro první podstrom kořene a obrátíme v něm pořadí vrcholů (jako první helmeči obřřáží nějaký vrchol v prvním patře tohoto podstromu, jako poslední jeho kořen). Za tento přřřevrácený plán přřijijeme rovněž přřřevrácený plán pro druhý, ... až poslední podstrom. Helmeči tedy skončí v kořeni posledního podstromu, tedy v prvním patře celého stromu, což přesně potřebujeme.

Pro zjednodušení implementace nemusíme plány ani dodatečně obracet. Stačí rekurzivním volání navíc předat jeden parametr, zda má generovat normální (z kořene do prvního patra), nebo obrácený (z prvního patra do kořene) plán. Díky tomu si nemusíme plány ukládat, nřřbrž mžžžeme vrcholů rovnomě vyřřisovat.

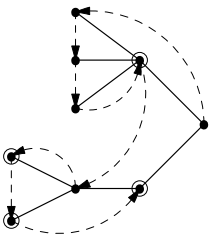
Snažmo si rozmyslete, že sousední vrcholů plány jsou vzájemně maximálně řřitřř hrany. Zřřývá ještě vyřřěšit okrajové podněmky rekurze: Zastavíme se na jednovrcholovém stromě (vrořřeném listem přřřívodního stromu), pro který je plánem prostě jednovrcholová posoupnost obsahující tento vrchol.

Celý algoritmus v pseudokódu:

`NAJDIPLAN(u, obrat):`

1. Pokud *obrat* = 0: vřřipš u .
2. Pro všechny v syny u :
 3. NAJDIPLAN(v , $1 - obrat$)
 4. Pokud *obrat* = 1: vřřipš u .

Rozmyslete si, že opravdu odřřovřřáá vyřře popsanému. Výsledný plán mžžže vypadat třeba takto:



Zakroužkovované jsou ty vrcholů, ve kterých hledáme obrácený plán. To jsou právě vrcholů v lichých patrech.

Obecně grařřy

Zobecnit řešení na všechny grařřy už je jednoduché. Z libovolného grařřu mžžžeme udělat strom tak, že najdeme jeho kořtm (napřř. prohlédnám do hloubky) a na ni sponusem vyřře popsaný algoritmus. Protože pro každý strom existuje řešení, hrany mimo kořtm jsou přřřebyřřeně a mžžžeme je ignorovat.

Navic hledání kořtm a konstrukci plánu nemusíme provádět odděleně, nřřbrž je mžžžeme spojit do jednoho DFS přřřichodu grařřem. Podrobneřřij v vzorovém programu.

Program (Python 3):

`http://exp.mff.cuni.cz/viz/28-3-6-dfs.py`

Filiřř Stědrnský

28-3-7 Legie v lese

Les je obědlmík o hraně L , v němž ležřř S bodových stromů. Horní stranou do obědlmíku vstupují kruhovřřá legie o přřřmě m a dce vyřřět řřpodní stranou. Nesmí přřřim narazít na

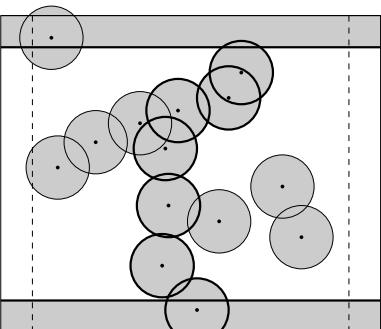
žádný strom, ani na levou řř pravou stranu obědlmíku. Šřřed legie tedy musí stále nudřřovat vzdálenost alespoň $D/2$ od všech stromů i od bočníh stran lesa.

Přřěkážky a ploty

Navigacní problémy tohoto druhu se řřešřř snažřř, pokud se lesem mřřsto kruhovřřá formace pohybůje jediný bod. Uložím proto trochu přřřevráceným: každý strom „nafunkcme“ na kruh o poloměru $D/2$ a boční strany lesa rozřřázáme smřřtem dorvnřř na obědlmřřky široké rovněž $D/2$. Aby se všechny přřřěkážky vřřšly do lesa, rozřřázáme les nahore i dole o $D/2$.

Legie napok sřřmskřřeme do jednéhodu. Ten mžžže stát právě na těch mřřstřřech, která neležřř v žádném kruhu ani obědlmřřku.

Dostaneme tedy nějaké přřěkážky ve tvaru kruhů a obědlmřřků a přřáme se, zda bod mžžže projít shora dolřř a vřřhnout se všem přřěkážkám.



Na přřědchozím obrázku to možné není, protože v něm existuje *plot* – posoupnost na sebe navazujících přřěkážků, která spojuje levý okraj s pravým. Každá trasa shora dolřř musí plot alespoň na jednom místě protnoout, takže trasa narazřř na alespoň jednu přřěkážku.

Dokonce platí, že kdykoliv nejde projít shora dolřř, existuje nějaký plot, který tomu brání. Vstrukturu množina bodů, do kterých se dá shora dojřř, tvořř nějakou uzavřřenou oblast. Hranice této oblasti se skláá z řřstřř horní, levřř a pravě strany lesa a nějakých řřstřř přřěkážků. Přřěkřřáme si, že hranici této oblasti obědlzíme od levého horního rohu lesa proti smřřtu hodinovřřých mřřček. Přřěkřřáme si, že dem z levě strany a následujícím přřřchodem na tu pravou jsme mřřseli projřř přřes na sebe navazujřřící řřstřř přřěkážků. To ovšem znamena, že tyto přřěkážky tvořří plot.

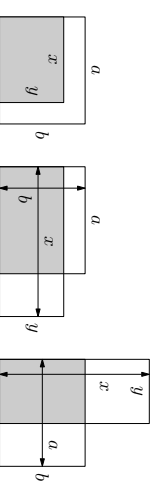
Gigantický grařř

Potřebujeme tedy vymyslet, jak hledat plot. Pomůžeme si grařřem: jeho vrcholů budou reprezentovat stromy, navíc přřřdáme vrchol ℓ pro levou stranu lesa a p pro pravou. Hranou spojíme každě dva vrcholů, jejichž vzdálenost v rovině je menší než D . To nastane právě tehdy, když se přřřisřře přřěkážky protínají. Plot pak odřřovřřáá cesty z vrcholů ℓ do vrcholů p v tomto grařřu. Pro příklad z přřědchozího obrázku vypadá grařř následovně:

Úloha má nepřijatelně mnoho stupňů volnosti: nejen, že si vybíráme, které desky použijeme, ale pro každou ještě určujeme natočení. Pojďme se otáčení zbavit. Dokážeme, že pokud každou desku předem natočíme „naležato“, tedy tak, aby její šířka byla větší nebo rovna výšce, o optimální řešení nepojedneme.

Nejprve uvažujme situaci se dvěma deskami rozměrů $a \times b$ a $x \times y$ (první rozměr je šířka, druhý výška). Jistě můžeme předpokládat, že $a \geq b$ a $x \geq y$ (jinak desku otočíme) a navíc $b \geq y$ (první deska je ta vyšší). Aby vznikl kvádr, musíme oběma deskami, chceme desky posunout a otočit tak, aby se překlápely v co největší ploše. Překryv jistě nezmůžeme, přiložeme-li na sebe levé dolní rohy obou desk.

Nyní rozlišme dva případy: V prvním je $x \leq a$, čili nižší deska je i užší (levý obrázek). Tehdy se desky překrývají celou plochou té menší z nich, takže se jistě nevyplácí kterkoliv desku otáčet.



Zbyvá nám případ $x > a$. Pokud obě desky ponecháme orientované stejně (prostorově obrázek), jejich společná část bude mít obsah $a \cdot y$. Pak lze je větší šobe otočme (pravý obrázek), vyjde obsah $b \cdot y$. První varianta je ovšem stejná nebo lepší, neboť $a \geq b$.

V obou případech platí, že překrývající se část má opět šířku větší nebo rovnou výšce, takže postup můžeme opakovat a o všech použitých deskách dokázat, že jsme je mohli nescházet naležato. Ve zbytku řešení tedy budeme bezelstně předpokládat, že deskami není možné otáčet.

Elementární řešení

Náš první pokus o algoritmus bude založený na jednoduchém pozorování: výsledky kvádr zdělí jak svou šířku, tak svou výšku od nlekeré desky každou možná od jiné. Slučí tedy výzkoušet n možných šířek, k nim n možných výšek a pokusíme probrat všechny desky a použít ty, které jsou dostatečně široké a vysoké. To zvládneme v čase $O(n^3)$ a získáme za to pár bodů.

Mimochodem, řešení tohoto druhu by fungovalo, i kdybychom brali v úvahu otáčení. Výšku i šířku kvádr bychom vybrali ze všech výšek i šířek (celkem $4n^2$ možností) a při testování, zda se deska vejde, bychom zkoušeli obě možné orientace.

Optimější řešení

V minulém řešení počítáme stále dokola podobné věci, takže zkusíme výpočet přeuspořádat, abychom se tomu vyhnutí.

Nadále budeme zkoušet všechny n možných šířek. Pro každou šířku w vybereme všechny dostatečně široké desky. Budeme je procházet od nejvyšší k nejnižší a průběžně přepočítávat aktuální kvádr. Pro nejvyšší desku samotnou má kvádr šířku w , výšku této desky a tloušťku l . Když přidáme další, nižší desku, šířka zůstane, výška klesne a tloušťka vzroste o l . Tak pokračujeme a pokusíme spočítáme objem aktuálního kvádra a započítáme ho do průběžného maxima.

Toho vše se dá stihnout v čase $O(n^3)$: na počátku výpočtu seřídíme všechny desky podle výšky. Pak zkusíme (v libovolném pořadí) všech n možných šířek, pro každou z nich probíráme desky v seříděném pořadí, přeskakujeme ty příliš těžké a pro ostatní v konstantním čase přepočítáváme objem kvádra.

Autor přiznává, že stále věří v existenci rychlejšího řešení, ale postupně si všechna taková vyvratil. Kdybyste o nějakém věděli, dejte nám prosím vědět.

Program (C):

`http://ksp.mff.cuni.cz/viz/28-3-5-5.c`

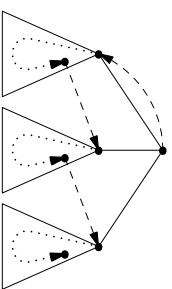
Martin „Matějka“ Mareš

28-3-6 Počítání přeživších

Aby řešení mohlo existovat, graf musí být souvislý. To budeme nadále předpokládat. Nejprve vytvášíme jednodušší variantu, kdy graf přátelství je stromem. Hledanému předpisu, kdo komu má helmicí předit, budeme říkat *plán* pro daný strom. Ten si lze představit prostě jako posloupnost vrcholů, kde dva sousední jsou vždy vzdálené ve stromě největšími hrany.

Strom si zakočíme a použijeme obvyklý trik na stromové úlohy. Nejprve rekurzivně najdeme plán pro každý podstrom kořene a pak tyto plány nějak napojíme, abychom z nich poskládali plán pro celý strom. Předpokládejme pro začátek (později se ukáže, že je to trochu složitější), že každým podstromem začne helmice putovat od kořene.

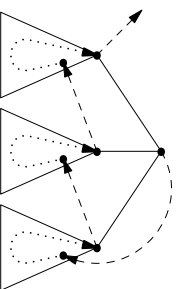
To znamená, že když helmice opouští první podstrom, musí-
me ji předat rovnou kořeni dříveho podstromu (přes kořen celého stromu to nejde, tam už byla):



Testované číry značí rekurzivně spočtené plány pro jednotlivé podstromy, čárkované jejich napojení. Aby předání mezi podstromy bylo korektní, musí plán pro každý podstrom končit v prvním patře tohoto podstromu (lince pod jeho kořenem). Kdyby končil hlouběji, předání bude přes více než tři hrany.

Abychom mohli naše řešení použít rekurzivně, musí tedy i výsledný plán pro celý strom končit v prvním patře. Ale plán nastříhový v obrázku výše zřejmě končí v patře druhém...

Leč není třeba propadat panice. Všimneme si, že kdybychom z kořene poslali helmici do vrcholu, který je aktuálně posledním a celý průběh příchoď príchod podstromy obrátili (díky neorientovanosti hran můžeme), skončíme v křízku prvním patře:



Taková funkce f by fungovala perfektně na trénovacích datech, ale jakmile bychom chtěli předpovědět hmotnost někoho, koho jsme ještě neviděli, byla by zcela k ničemu. *Generalizace* je právě tato schopnost fungovat dobře i na lidi, které jsme při učení ještě neviděli.

Přiznaky

Když chceme mít dobré předpovědi, samozřejmě velmi záleží na tom, abychom zvolili dobré příznaky. Mnsí obsahovat nějakou užitečnou informaci, na kterou závisí ta veličina, kterou předpovídáme. Když se snažíme naučit předpovídat hmotnost lidí, nejspíš nám ponauče vědět fyzikální vlastnosti, které s hmotností souvisí: třeba výška, obvod pasu, věk nebo procento tuků v těle. Naopak nám asi nepomůžou vědět barvu očí nebo oblíbenou kapku.

Než spustíme algoritmus strojeového učení, vyplácí se nejlépe zamyslet nad reálnou strukturou problému a získat pro něj významné co nejzhuštěnější příznaky. Zkusme třeba předpověď vidět podle výšce a hmotnosti pravděpodobnost srdečního příchoď v následujícím roce. Samotná hmotnost a výška sice jsou užitečné informace (když vážíme 250 kg, jsem rizikovější, než kdybych vážil 70 kg), ale lepší nejspíš bude si předat jako příznak BMI (hmotnost v kg)/(výška v m)². BMI zohledňuje, že různí lidé mají různou postavu – je asi zhruba stejně vzáží 100 kg a být vysoký 2 metry, než váží 80 kg a být vysoký 1,5 metry.

Příznaky jde široce dělit na kategoričké a numeričké. Numeričké příznaky se dají přirozeně vyjádřit jako číslo, třeba výška v centimetrech nebo barva pixelu v obrázku. Kategorie příznaky můžou být třeba krevní skupina nebo státní občanství. Existuje pro ně nějaký poměrně malý počet možných hodnot (třeba krevní skupiny jsou $\{0^+, 0^-, A^+, A^-, B^+, B^-, AB^-, AB^-\}$ a na těhle hodnotách nemusi třeba označit krevní skupiny místo názvu pořadovým čísly ($0^+ = 1, \dots, AB^- = 8$), ale operace nad takovými označením nejsou užitečné – sice nad našim označením můžeme tvrdit, že $(0^-) + (B^-) = (AB^-)$, ale to neodpovídá žádné mu vztahu v reálném světě. Stejně tak není B^- v žádném smyslu „větší než“ 0^- . Oproti tomu třeba může dávat smysl porovnávat výšky dvou lidí nebo počítat jejich rozdíl.

Hodně algoritmů potřebuje, aby jejich vstupy byla čísla. Tehdy potřebujeme kategoričké příznaky „zakódovat“ do číselných. Nejpopulárnější takové kódování s příznakem, který obsahuje jednu z K kategorií, vyrojí K číselných příznaků, jeden pro každou kategorii. Když vzorek patří do i -té kategorie, nastavíme i -tý příznak na 1 a ostatní na 0.

Předpokládáme, že existuje nějaký *skutečný* vztah f , který z x dává y (tedy z výšce a obvodu pasu dává hmotnost). Příznaky, které měříme, ale nemusi stačit na zcela přesnou odpověď: i když mají Jana a Katka stejnou výšku a obvod pasu, můžou mít jinou hmotnost, protože se Katka ráno nenamastala. Existuje nějaký vliv vnějších příznaků, které neměříme (nebo možná ani z principu měřit nelžou – třeba máme nepřesnou váhu). Dá se to neformálně napsat jako

$$y = f(x) + \epsilon; \text{ předpovídaná hodnota } y \text{ se skládá ze složky, která závisí na } x, \text{ a nějakého náhodného (a snad malého) } \epsilon, \text{ ve kterém jsou schované vlivy, které neumíme měřit.}$$

Funkci f , kterou se snažíme naučit, se říká *model* – snaží se co nejpřesněji *modelovat*, co by dělala f , kdybychom se

³ Kromě accuracy se pro klasifikátory měří i metricky *precision* a *recall*. České překlady tedy bohužel nemají tak dobře zavedený význam, jako anglické termíny.

ji mohli zeprat.

Měření chyby modelem

Po modelech chceme, aby byly co nejpřesnější a aby dobře generalizovaly. Chceme tedy, aby na *neznaměném* vzorku, který *nebyl k dispozici učícím algoritmu*, daly dobrou předpověď.

Existují různé metricky pro to, jak *dobře* předpověď je. Většina z nich jsou nějaká míra *chyby*.

Pro náš předpovědní hmotnosti se třeba hodí velmi obvykle používaná *kvadratická odchylka*. Kvadratická odchylka modelů na vzorku (x, y) je rovna $(y - f(x))^2$. Intuitivně jí velké odchylky vadí mnohem více než malé.

Pro klasifikaci udává se jako metrika hočí *accuracy*.³ Accuracy na jednom vzorku je 1 tehdy, pokud je f předpověď správně, a 0, když na něm udělá chybu. Když třeba předpovídáme, jakou dopravu značek obsahuje obrázek, zajímá nás jenom, jestli našleme tu správnou. Když si spleteme slopek se záznamen vzadu, je to pro accuracy stejně špatné, jako bychom si ji spleli s příkazným směrem jízdy.

Zatím jsme si ukázali definice dvou různých chyby modelem na jednom vzorku. Po modelem chceme, aby měl co nejmenší *střední hodnotu chyby*, neboli aby na náhodně vybratém *neznaměném* vzorku byl co nejpřesnější.

Když máme nějaký model f , jak zjistíme střední hodnotu chyby na *neznaměném* vzorku? Neznáme vzorky jsou pro nás nedostupné – nemůžeme jít změřit 7 miliard lidí a spočítat, jakou chybu průměrně děláme. Máme k dispozici jenom data od dobrovolníků. Dělá se to tak, že učícím algoritmem nedáme všechna data, která máme. Rozdělíme dataset \mathcal{D} na *trénovací množinu S* a *testovací množinu T*. Třeba v poměru 90%/10%. Učícím algoritmem dáme k dispozici jenom trénovací data. Testovací vzorky příd nám skryjeme. Až nám učící algoritmus dá model f , spočítáme jeho průměrnou chybu na testovací množině. S trochu statistiky se dá ukázat, že průměrná chyba na testovací množině rozumně odhaduje střední chybu na *vešech neznaměných vzorcích*.

Úkol 1 [1b]: Je potřeba, aby rozdelení \mathcal{D} obsahuje nějaký vzácí množin bylo *náhodné*. At dataset \mathcal{D} obsahuje nějakých 100 značek „slep“, pak 100 značek „del předstoup v jízde“ a nakonec 100 značek „slepá ulice“. Vymyslete, co a jak by se mohlo rozbit, kdybychom testovací množinu vybrali náhodně.

Když třeba po modelem chceme malé kvadratické odchylky, chceme malou střední kvadratickou odchylku na *neznaměných* datech. Tu odhadneme podle střední kvadratické odchylky na testovací množině T :

$$E \approx E_T = \text{MSE}_T = \frac{1}{|T|} \sum_{(x,y) \in T} (y - f(x))^2.$$

Střední kvadratické odchylce se říká *anglicky mean square error (MSE)*. Když nás zajímá accuracy na *neznaměných* datech, odhadneme ji podobně: pomocí průměrné accuracy na T .

Čím více dat dáme k dispozici algoritmu strojeového učení, tím více se jim bude moct přizpůsobit a tím bude naučený model dávat přesnější předpovědi. Na druhou stranu, čím větší máme testovací množinu, tím přesnější chyba na testovací množině E_T odhadne skutečnou chybu na *neznaměných* vzorcích E .

Vzpomenete si na f , která si uložila všechna trénovací data do tabulky a na všechno kromě nich vrátila 9 999. Takový model má na trénovacích datech nulovou chybu ($E_S = 0$). Na testovacích datech T , která nemá v tabulce, ale odpovídá stráně špatně, takže průměrná chyba na testovacích datech E_T nám správně řekne, jak strašně moc špatný tenhle model je.

Přučení a porovnávání modelů

Pokud se naučíme model, který je hodně dobrý na trénovacích datech, ale podstatně horší na testovacích datech, může to být kvůli *přeučení* (neboli *overfittingu*). K přučení dochází, pokud učícím algoritmem umožníme, aby se přizpůsobilá adaptováním na jednotlivé datové vektory, které ale obecně nepatří.

Hodně algoritmy strojového učení funguje tak, že postupně po *epochách* víc a víc adapty model na trénovací data. Vytvoří tedy posloupnost modelů f_1, f_2, \dots ve které jsou modely postupně čím dále tím adaptovány na trénovací data, ale po nějaké době se začnou přetřevovat, a tedy začínou být méně užitečné pro obecné použití. Podobná situace nastane, když zkonštruujeme na jednotlivých datech trénovacího učení a chceme z naučených modelů vybrat ten nejvýhodnější.

Ohavný přístup, jak vybrat z modelů f_1, f_2, \dots ten nejlepší, je zněit chybu všech modelů na testovacích množině a vrátit ten, který ji má nejméně, ale tady přístup je rozbitý.

Proč je rozbitý? Vzpomenete si, že testovací množina se používá k *oddělování skutečné chyby*. Když si z modelů vybereme ten, který na nejméně testovací chybu, bude jeho testovací chyba *přilíš optimistický* odhad skutečné chyby. Instrukujeme si tenhle problém malým myšlenkovým experimentem. Představte si, že naše modely jsou tři ferové mince. Pokud padne panna, model dá správný výsledek, a když padne orl, dá špatný výsledek. Každý model tedy ve skutečnosti dá správný výsledek v 50 % případi.

Testovací množinu vytvoříme tak, že každou minci desetkrát hodíme. Na první vyrobíme 3 orli, na druhé 7 orli a na třetí 5 orli. Vybereme si tedy drhlý model a budeme si o něm myslet, že je přesný v 70 % případi. To je víc než jeho skutečných 50% – druhý model jenom měl to štěstí, že při našem testu naházal nejvíce orli. Jsme moc optimističtí o jeho výkonu, a to o 20%.

Co kdybychom neházali třemi mincemi, ale 10 000 mincemi? Skoro určitě (s pravděpodobností asi 0,99994) se stane, že náhodou některá z nich naházá 10 orli. Být bychom *ceťremně optimističtí* – mysleli bychom si, že jsme našli minci, na které vždycky padají orli, i když je ferová. Přičítání dalších modelů způsobilo, že náš odhad je *horší* – teď už se mylíme o 50 % místo 20%.

Správné řešení je udělat „výběr nejlepší mince“ a „odhad pravděpodobnosti“ jako nezávislé experimenty: nejdrví 10× hodit a vybrat nejperspektivnější minci, a pak ji znovu 10× hodit a podle drhlých hodit odhadnout její „chílkost“. Když jsme v první fázi vybrali minci, co naházala 10 orli, ale ve skutečnosti je ferová a jenom měla štěstí, tak druhá fáze pořád bude 10 hodit ferovou minci a nejspíš nám řekne, že skutečná pravděpodobnost padnutí orla na vybrané

⁴ Pokud znáte skalární součin, všimnete si, že $f(\vec{x}) = \langle \vec{x}, \mathbf{1} \rangle \cdot \langle \vec{x}, \mathbf{1} \rangle$.

⁵ Pro lineární regresi existuje i vzorec, ze kterého jde nejlepší model přímo spočítat, ale gradientová metoda je jednodušší a pochoptí a obecněji – používá se ve spoustě dalších učicích algoritmy.

minci je 0,5.

Když vyberáme z více modelů ten nejlepší a pak chceme vědět, jaký výkon od něho můžeme očekávat na nových datech, jedno ze správných řešení je rozdělit data na tři množiny: *trénovací, validací* a *testovací* (těeba v poměru 80%/10%/10%). Trénovací množina se dá k dispozici učícím algoritmem (ono mád ní iteruje jeden algoritmus a leze z něj posloupnost modelů). Jako nejlepší model vybereme ten, co má nejméně chybu na *validacích* množině. Tím pádem „neznámejšíme testovací množinu“ a budeme ji moci dále používat k dobrému odhadování skutečné chyby.

Chyby na validacích a testovacích množinách odpovídají naměřeným „chílkostem“ v první a druhé fázi myšlenkového experimentu s mincemi.

Tedy už se trochu významně v základních termínech a souvislostech, tak se konečně vřhujeme na něco programovacího.

Lineární regrese

Algoritmy obecného učení jsou si obecně hodně podobné:

- Předepší, jaký obecný tvar budou mít modely f , které z nich budou padat. Tady předpis bude obsahovat vstupní vektor x a nějaký vektor *parametrů*, který se označuje β . Konkrétní model dostaneme, když do předpisu dosadíme parametry.
- Rikajit, jakou chybu se snaží minimalizovat.
- Popišit, jak elektrane nařít takové β , že předpis f s danými parametry β bude mít co nejméně chybu.

Lineární regrese konkrétně:

- Hledá *lineární model*. Lineární model je funkce f , která na vstup x dá jako výstup $f(x) = \theta + \sum_{i=1}^p \alpha_i x_i$. Konstanty θ a α_i jsou parametry určující konkrétní lineární model. Dohromady je těchto parametrů $(p+1)$, tedy vektor parametrů β je $(p+1)$ -rozměrný; nejdříve obsahuje p složek $\beta_i = \alpha_i, \beta_p = \alpha_p, \dots, \beta_p = \alpha_p$, a pak jednu složku $\beta_{(p+1)} = \theta$. Graf lineárního modelu je přímkou v p -rozměrném prostoru.⁴
- Chyba, kterou minimalizuje, je *sřední kvadratická odchylka* na trénovacích množině, proto se se jí taky říká *metoda nejménších čtverců*:

$$E = \text{MSE}_S(\beta) = \frac{1}{|S|} \sum_{(x,\theta) \in S} (y - f_\beta(x))^2$$

- Efektivně najde dobrou hodnotu vektoru β pomocí *gradientové metody*.⁵

Kdyby nám nezáleželo na času na naučení, stačil by nám první dva „modely“. Mohli bychom si jenom vygenerovat bilianru modelů (těeba pro všechny hodnoty všech složek β od -100 do 100 s krokem $0,1$), pro každý spočítat chybu a vybrat ten, co ji má nejméně. To ale rychle přestane být efektivní pro hodně parametrů.

Ukol 2 [10]: Jak zkrátit časová složitost tohoto hloupého přístupu „vygenerujeme si tabulku se spoustou modelů a vybereme ten nejlepší“ na volikosti trénovací množiny a počtu příznaků p ?

Efektivnější minimalizace chyby vyžaduje trochu matematické analýzy. O kvadratické chybě se dáji dokázat užitečné

$H[x] \leftarrow \mathbf{R}$

- Označ x_i až x_p všechny závislosti x .
- Pro $i = 1$ až k :
 - $h \leftarrow \text{OHODNOŤ}(x_i)$
 - Pokud $h = \mathbf{X}$:
 - $H[x] \leftarrow \mathbf{X}$, vrať \mathbf{X}
- $H[x] \leftarrow f(H[x_1], \dots, H[x_k])$, přičenž f je logická funkce předepsaná pro formuli x (AND, OR nebo XOR).
- Vrať $H[x]$.

Protože chceme zjistit ohodnocení všech formuláři, prostě funkci OHODNOŤ zavoláme postupně na všechny vchody. Celkem budeme potřebovat lineární čas, přesněji $O(N + M)$, kde N je počet formuláři a M je celkový počet závislostí. To nahlébáme následovně. Vnitřní část funkce OHODNOŤ provedeme pro každý vchod nejlíže jednou. Tím pádem nejvíše jednou provedeme vnitřní cyklus přes všechny hrany vyznačující z daného vchodu. Tudiž na každou hranu se podíváme maximálně jednou.

Jestli musíme uvěřit čas strávený na prvích třech řádcích funkce OHODNOŤ, které mohou být provedeny vidět. Ale funkce OHODNOŤ je volána jen ze dvou míst:

- Z hlavní smyčky, kde je volána jednou pro každý vchod (celkem $N \times$).
- Rekurzivně z jiného volání OHODNOŤ, ale to se děje právě v místě, kde procházíme výstupní hranu. A už víme, že každou hranu navštívíme maximálně jednou, tedy těchto volání funkce OHODNOŤ je dohromady maximálně M .

Z tohoto odhadu je taky vidět, že teď už se náš algoritmus nemůže začtyčit.

Program (Python 3):
<http://ksp.mf.cuni.cz/viz/28-3-3-dts.py>

Filip Stětnýský

28-3-4 Katapulty

Podíváme se první na leké varianty. Nabízí se umísťovat katapulty hladově, tedy dát každý katapult co nejlépeji to půjde.

Řekneme, že projdeme bitvoni lajnu zleva a při procházení si udržujeme množinu dosud neumísťených katapultů, které můžeme umístit na aktuální políček.

Na každém políčku pak přidáme do množiny katapulty, které můžeme nově umístit, a poté z této množiny vezmeme katapult, jehož úsek končí nejlíže. Zkontrolujeme, že právě hranice povoleného úseku je maximálně rovna aktuální pozici, a pokud ano, katapult umísťme na toto políčko. Pokud náhodou ne, zahlásíme, že řešení neexistuje. Budeme-li mít zrovna prázdňou množinu, jednoduše se posuneme na další políčko.

Standardně uhlédneme, že pokud náš algoritmus vydá řešení, bude toto řešení korektní – každý katapult bude ve svém úseku a zároveň bude na každém políčku jen jeden.

Musíme ale také ukázat, že pokud naopak ohlásíme neexistenci řešení, skutečně žádné řešení neexistuje. Když jsme vyhlásili neúspěch, máme nějaký katapult P , který jsme nedokázali umístit. Na každém políčku v úseku, kam smíme P umístit, už se ale vystrkuje jiný katapult, jehož úsek končí nejpodelji stejně, tedy žádný z nich nemůžeme posunout doprava.

Jestli musíme uvěřit, jestli jsme některý z těch katapultů nemohli umístit více vlevo. Podobným argumentem ale ukážeme, že ne – když vezmeme nejvíše levou hranici těchto katapultů a podíváme se na úsek mezi ní a začátkem úseku pro P , opět máme na každém políčku katapult, který nemůžeme přesunout doprava.

Takto dojdeme až na začátek lajny. Vlastně to znamená, že se snažíme umístit $L+1$ katapultů na L políček, a to evidentně jí nemůžeme. Náš algoritmus tedy ohlásí neúspěch, pouze když řešení neexistuje.

Připomeneme teď, že počet katapultů označujeme K a délkou lajny N .

Nejprve si všechmy katapulty seřadíme podle levé hranice veszstupně. Tím pádem se můžeme postupně posouvat v poli katapultů a nalezneme katapultů, které máme přidat do množiny; trvá lineárně v jejich počtu, za celý algoritmus tedy $O(K)$. Samotné řazení nám zabere $O(K \log K)$.

K udržování katapultů, resp. k jejich přidávání a nalezení toho s minimální pravou hranicí, nám dobře poslouží haldka. Tak bude každá operace trvat $O(\log K)$.

Nejprve tedy seřadíme katapulty a pak projdeme celou bojovou lajnu. Tento příchod nám trvá $O(N)$, na každém políčku se plháme na katapult s minimální souřadnicí a možná přidávame katapulty do množiny. Takový přidání je ale dohromady $O(K)$, takže celková časová složitost algoritmu je $O(N \log K)$. Přesněji tedy $(K \log K + N \log K)$, ale klidně předpokládáme $K \leq N$. Paralelová složitost je lineární v počtu katapultů, tedy $O(K)$.

Překně je si všimnout, že políčka, na kterých se nic neděje, jsou nezapínavá a vlastně nás jen zdržují. Išlo by se nám umět skákat jen po těch, kde se něco děje. To ale umíme zadržet – o každém katapultu víme, kde jeho úsek začíná a končí, takže se můžeme hrd posunout o jedno políčko (námělo-li pro umístění katapultu neprázdnou množinu), nebo rovnom skocit na příští udlost.

Události si můžeme udržovat v haldě, na každém políčku tedy ještě přidáme připadně zjištěný příští udlosti. Jelikož pro umístění K katapultů použijeme nejvíš K políček, bude celková složitost $O(K \log K)$.
 Tedy těží varianta, a pozor, přijde trik :) Všimneme si, že souřadnice katapultů jsou nezavíslé – bez ohledu na to, do kterého řádku katapult umísťme, můžeme ho umístit do stejných sloupců.
 Řešením úlohy tak není nic jiného než dvoji spuštění leké varianty, jednou pro řádky, jednou pro sloupce. A jelikož dvojka je konstanta, časová složitost zůstává $O(K \log K)$.

Program (C):
<http://ksp.mf.cuni.cz/viz/28-3-4.c>
 Program (C++):
<http://ksp.mf.cuni.cz/viz/28-3-4.cpp>

Karoly Bursňová

28-3-5 Závaží z fosen

Nejprve fosny ze zadání obohojíme až na struhou geometrickou realitu: Dostali jsme n deskek, každá má danou šířku a výšku, všechny mají jednotkovou tloušťku. Chceme vytvořit kvádr. Uděláme to tak, že si vybereme podmožnost deskek, položíme je na sebe a přídáme některé z nich o 90° otočime. Nakonec je otočime na šířku nejlíže desky a výška té nejnižší. Chceme přitom, aby vznikl kvádr o největším možném objemu.

že jsou položene na mížce o velikosti $|A|$, $|B|$, ve které souřadnice vrcholů určují indexy jednodílných třetozí s tím, že vlevo nahore máme vrchol $S = [0, 0]$ a pravou dole je vrchol $C = [a_{n-1}, b_{n-1}]$.

Pro každý vrchol dále bude platit, že z něj vede hrana doprava a dolů vždy, pokud není poslední ve své souřadnici, a dále si zavazujeme zkratkové hrany, které nám budou vést přes diagonální dolů doprava. Tyto hrany budou vést z těch vrcholů, které nám označují znak společný oběma třetozím.

Nyní, když máme takový graf postavený, nalozeme nejkratší cestu z vrcholu $[0, 0]$ do vrcholu $[A, B]$. Jakmile jsme jednu takovou cestu našli, vydáme se po ní a zastdíme se podle toho, kterým směrem se rozhodneme vydat z vrcholu $[a_i, b_j]$:

- Pokud je následující vrchol cesty napravo od aktuálního, vypíšeme na výstup znak a_r .
- Pokud je následující vrchol cesty směrem dolů od aktuálního, vypíšeme znak b_r .
- Pokud se následující vrchol cesty nachází směrem po diagonále, je celkem jedno, který znak necháme vypsat, protože aktuální vrchol označuje společný znak.

Tímto jsme vypsalí na výstup hledanou nejkratší spočtenou nadpostlopnost. Samotné sestavení grafu nám potrvá $O(|A| \cdot |B|)$, jelikož tento graf má počet vrcholů stejný, jako je součet délek obou řetězů $A \cdot B$. Nejkratší cestu potom umíme nalozit v lineárním čase jednodílným prohlédáním, takže nám to časovou složitost neporisí. Výsledný algoritmus má tímto stejnou asymptotickou složitost jako algoritmus popsaný výše.

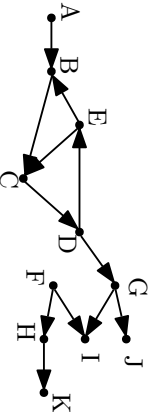
Program (C):
http://ksp.mff.cuni.cz/viz/28-3-2.c
http://ksp.mff.cuni.cz/viz/28-3-2.cpp

Václav Konečický a Karly Buršovná

28-3-3 Formulář na zbroj

Nejprve si uolnu převeďeme na grafovou. Pro každý formulář vytvořime jeden vrchol. Řekšme, že formulář x (přimo závisí na formuláři y , pokud k získání x musime předložít vyplněné y . Pro každý takovouto dvojici vytvořime v grafu orientovanou hranu $x \rightarrow y$. Některé formuláře jsou vydávány „zadarnu“ (nemají závislosti). Odpovídající vrcholy v grafu nemají žádné výstupní hrany – takovéto vrcholy obvykle nazýváme *stokly*.

Graf závislosti může vypadat třeba takto:



Stokly jsou v tomto případě vrcholy I, J, K. Dále řekneme, že nějaký formulář x *nepřimo závisí* na y , pokud z x vede orientovaná cesta do y . Takový x nemůžeme získat, aniž bychom někdy předtim získali y . Například F nepřimo závisí na K, a vskrátku: bez K nedostaneme H, a tedy ani F.

¹⁰ http://ksp.mff.cuni.cz/viz/kucharcky/grafy
¹¹ http://ksp.mff.cuni.cz/viz/kucharcky/dynamika

Pokud nějaký formulář leží na orientovaném cyklu, nepřimo závisí sám na sobě. Takový určitě nemůžeme získat: například abychom získali E, potřebujeme C, který vyžaduje D, a ten zase E. Jinými slovy abychom dostali E, museli bychom nejprve mít E. To určitě nejdě, protože třetozí nedávrají formuláře na dlnh.

Dalším formulárem, který určitě nemůžeme získat, je A, protože závisí na B, který leží na cyklu.

Řešení prohlédáváním do hloubky

Pro každý formulář budeme chtít určit jeho *ohodnocení*, které může nabývat jedné ze tří následujících hodnot:

- 0, pokud formulář lze získat vyplněním záporné
- 1, pokud formulář lze získat vyplněním kladné
- X, pokud formulář není možné získat

Ohodnocení formuláře x nalozeme jednodílnou rekurzivní funkcí OHODNOU(x). Ta nejprve rekurzivním zavoláním získá ohodnocení všech přímých závislostí x , na jejich základě určí ohodnocení x ; a to vrátí. To vlastně není nic jiného než prohlédávání do hloubky.¹⁰

Protože funkce OHODNOU může být na tenký formulář zavolána mnohokrát a nechtíme pýřvat časem opakovaním toho samého výpočtu, použijeme dynamické programování.¹¹ Již spočítaná ohodnocení si budeme ukládat do pole H indexovaného těly formulářů. Pokud je při zavolání OHODNOU(x) hodnota H[x] již známa, prostě jí rovnou vrátíme a nepočítáme znovu.

Kostra ohodnocovací logiky je jednodílná: pokud některou ze závislostí x nejdě získat (má ohodnocení X), ani x nemůžeme získat, tedy vrátíme rovněž X. Pokud napopak jde získat všechny, pak jde získat i x a jeho ohodnocení vytvořime tak, že ohodnocení závislosti zkombinujeme předepsanou logickou funkcí.

Ale to nestačí. Pokud graf obsaduje cykly, takovéto implementace by skončila nekonečnou rekurzí. Potřebujeme mít nějak cykly detekovat a vrcholy na nich označit jako nezávísitelné.

To můžeme udělat například tak, že na začátku funkce OHODNOU(x) si vrchol x označime jako *rozpracovaný* (a na konci toto označení zase zrušime). Pokud potom někdy při prohlédání závislosti narazíme na rozpracovaný vrchol, víme, že leží na cyklu. Proč? Pokud narazíme při prohlédání na rozpracovaný vrchol x , znamená to, že funkce OHODNOU(x) aktuálně běží (v nějakém vyšším patře rekurze). Všechny vrcholy, které navštívime, zatímco běží OHODNOU(x), patří mezi (nepřimo) závislosti x . Takže pokud narazíme na x , znamená to, že x závisí na x .

Zavedeme si dvě pomocné hodnoty, které se mohou vyskytovat v poli H: pokud ohodnocení daného vrcholu ještě neznáme:

- ? , pokud jsme dané pole ještě nenavštívili
- R, pokud je vrchol právě rozpracovaný

Upravovaný algoritmus bude vypadat následovně:

OHODNOU(x):

1. Pokud $H[x] = R$:
2. $H[x] \leftarrow X$, vrať X.
3. Pokud $H[x] \neq ?$: vrať $H[x]$.

vlastnosti. Zaprvé když si vyjdíme E jako funkci parametru $\beta = (\alpha_1, \dots, \alpha_n, \theta)$, zjistíme, že je *spojitá*. Zadržme má E jediné lokální minimum, které je i její jediné globální minimum. Záčetí nemá žádné inflexní body.

Gradientová metoda

Gradientová metoda (*gradient descent*) umí takové funkce minimalizovat. Obecně tato metoda funguje tak, že začneme v nějakém libovolném, třeba nulovém počítacím bodu β_0 . Potom se podíváme, kterým směrem bychom mohli trochu posunout β_0 tak, abychom tím co nejvíc snížili E .

Je vám ten nápad pověřomý? Je vám pověřomý zcela správně. Ve druhém dle seriálu jsme si ukazovali, jak hledat „horolezení“ extrémů reálných funkcí. Gradientová metoda je horolezení velmi podobná. Liší se tím, že „má kompas“ – umí najít zlepšující bod v okolí efektněji než nahodným zkoušením.

Kompassu se říká *gradient*. Gradient E v bodu β_0 se značí $\nabla E(\beta_0)$ a je to vektor, který říká, kterým směrem máme jít z β_0 , abychom šli co nejméně směrem zvyšující se E . V každém bodu může obecně gradient vést jiným směrem.

Když gradient otočime, dostaneme směr největšího poklesu E . Když je délka gradientu v bodě β velká, znamená to, že E v β rychle roste/klesá. Naopak malý gradient znamená, že se nacházime v placaté oblasti a nulový gradient znamená, že β je lokální minimum, maximum, nebo inflexní bod. Protože E inflexní body nemá, značí v ní nulový gradient lokální minimum nebo maximum.

Krok gradientové metody číslo t začne v souřadnicích β_{t-1} . Spočítá si v tomhle bodě gradient $\nabla E(\beta_{t-1})$, nějaký jeho γ -násobek odečte od souřadnic a tím spočítá β_t .

Jak velké má být γ ? Čím menší, tím dle bude gradientová metoda bžet, ale tím přesněji se zase treí do lokálního optima. Když bude γ moc velká, budou kroky gradientové metody lokální minimum „přeskokovat“, což se projeví tak, že se po pár třetozích dostaneme k nekonečné velikým hodnotám parametru a nic se nemaneme. Pokud se vám taková věc stane, zkuste γ zmenšit o pár řádů. Hodně učících algoritmů má podobný parametr ovlivňující velikost jednoho učícího kroku. Většinou se mu říká *learning rate*, rychlost učení.

Po odečetí γ -násobku gradientu přejdeme na krok $(t + 1)$: spočítáme gradient v β_t , odečteme jeho γ -násobek, a tak dále. Skončime, když je splněno nějaké kritérium, například když už jsme počítali dost dlouho nebo když je gradient hodně malý. Pokud gradientová metoda skončí v místě, kde je gradient skoro nulový, našla lokální minimum.

Gradient je rovný nule i v lokálním maximum. Když bychom měli tu neuvěřitelnou smůlu, že bychom začali gradientovou metodu dokonale přesně v lokálním maximum, tak bychom hned skončili a vrátili lokální maximum. To se ale skoro určitě nestane. Stane se možná, že začneme gradientovou metodu blízko lokálního maxima. Potom nás od nuly ale první krok trochu oddálí (protože jde směrem klesající E), druhý krok nás oddálí ještě víc a tak dále, takže v lokálním maximum neskončime.

Gradientová metoda tedy skončí poblíž lokálního minima, které je kvůli vlastnostem E i globální minimum, a tedy bod, ve kterém skončime, bude určovat parametry lineárního modelu s malou chybou.

⁶ http://ksp.mff.cuni.cz/viz/evolutione

Zbyvá už poslední kussek skladačky: jak gradient E spočítáme?

Pokud bychom o E nevěděli nic, mohli bychom místo počítání gradientu zkusiť jenom o trochu pohnout každou složkou β a vybrat to druhou polnání, které E nejvíc zmenší. To bude poměrně pomalé, ale jednodílné a bude to fungovat na všechna E , co jsou spojitá, mají jediné lokální minimum, které je zároveň i globální minimum a nemají inflexní body. Gradient se dá pro naši definici chybry spočítat explicitně. Je to vektor o $(t + 1)$ složkách a jeho t -tá složka se spočítá tímto průměrem přes všechny trénovací vzorky:

$$\nabla E(\beta)_t = 2 \cdot \frac{1}{|S|} \sum_{(x,y) \in S} x_t(f_{\beta}(x) - y).$$

Poslední složka (odpovídající členu θ) se spočítá jako:

$$\nabla E(\beta)_{(t+1)} = 2 \cdot \frac{1}{|S|} \sum_{(x,y) \in S} (f_{\beta}(x) - y).$$

V obou rovnicích je f_{β} lineární model určený parametry β . Bližně plán: Záčneme v libovolném (třeba nulovém) bodě $\beta_0 = (\tilde{\alpha}_0, \theta_0)$. Spočítáme gradient, odečteme jeho γ -násobek a opakujeme, dokud gradient není malý vektor nebo nás to nepřestane bavit. Uložíme výsledný model.

Je zaručeno, že gradientová metoda minimalizuje spojitě funkce, které mají jediné globální a lokální minimum a nemají inflexní body. Často se ale používá i na spojitě funkce, které tunle vlastnost nemají. Potom může skončit v lokálním, ale ne nutně globálním optimu, nebo v inflexním bodě. V praxi nám to ale často neradí – i tak dává poněkud dobré výsledky. Gradientové metody také můžeme pomoot *náhodnými vesměry* – pustime ji někdeklířte za sebou z různých náhodně zvolených β_0 . Dá se pak domluvit, že projdeme větší kus dělního oboru E , takže v nejlepší nalezeném bodě budeme blíž globálnímu minimum.

Ukol 3 [6b]: Stáhneme si ze stránky seriálu⁶ dataset o spotřebě benzínu v USA. Soubor má 5 sloupců dat oddělených čárkami. Každý řádek jsou data z jednoho státu. První sloupec je vyběrána data z benzínu v centech na galon, druhý je průměrná mzda v dolarce, třetí délka dálnic ve státních milcích, čtvrtý je počet obyvatel s řidičským průkazem a pátý je roční spotřeba benzínu ve státech v miliolech galonů. Naprogramujme lineární regresi a pomoci ml odvoďte vzorec na předpovídání spotřeby benzínu podle ostatních hodnot. Data nemusíte dělit na testovací a trénovací množiny. Pro nás fungovalo dobře začít v nulových parametrech a pak provést 100 000 třetací gradientové metody s $\gamma = 10^{-8}$. Pošlete svůj zdrojový kód a nalozeme střední kvadratickou odchylku na celém datasetu. Zkuste ji mít menší než 21 800.

Algoritmus K nejbližších sousedů

Algoritmus K nejbližších sousedů (anglicky K -nearest neighbors nebo KNN) je založený na jednodílné myšlence: když chceme poznat mají výšky a obvodu pasu předpovědět mou hmotnost, odhadneme ji podle lidí, kteří jsou mi podobní vyškou a obvodem pasu.

Model si bude pamatovat všechna trénovací data a jeho jediný parametr je přirozené číslo K . Když nám přijde nový vstup x , nalozíme ve trénovacích datech K vzorků, které jsou mu nejbližší. Podíváme se na známé výstupy pro nejbližší sousedy, nějak je zaveragejeme a výsledek vrátíme.

Konkrétní použití se liší podle toho, jak nadefinujeme, že je vstup x „blízký“ k nějakému trénovacímu vektoru, a podle toho, jak x vstoupí na nejbližších susedech výrobníe předpověď pro neznaný vstup.

Blízkost se může definovat třeba podle malé euklidovské vzdálenosti

$$d(x, x^{(j)}) = \sum_{j=1}^p (x_j - x_j^{(j)})^2.$$

Euklidovská vzdálenost se hodí na numerické příznaky. Pokud jí ale použijeme, je potřeba dát si pozor na to, aby rozdíly v příznacích byly stejně významné.

Co tím myslíme? Vzpomeneme si na příklad s předpovídáním hmotnosti a podívejme se na tři vzorky s těmito vstupními příznaky:

příznak	$i = 1$	$i = 2$	$i = 3$
$x_j^{(i)}$: výška v m	1.93	1.92	1.4
$x_2^{(i)}$: obvod pasu v cm	83	86	82

Dobrovolek 1 se liší od dobrovolek 2 jedním centimetrem výšky a třemi centimetry obvodu pasu. Dobrovolek 1 se od dobrovolek 3 liší 53 centimetry výšky a jedním centimetrem obvodu pasu. Nemusíme být experti na antropologii, abychom očekávali, že hmotnost dobrovolek 1 bude podobněji hmotnosti dobrovolek 2 než dobrovolek 3.

Rozhodli jsme se, že první složka (výška) bude číslo v metrech a druhá složka (obvod pasu) v centimetrech. Když si spočítáme euklidovskou vzdálenost tak, jak jsme si ji nadefinovali, dostaneme:

$$d(x^{(1)}, x^{(2)}) = (1.93 - 1.92)^2 + (83 - 86)^2 = 9.0001$$

$$d(x^{(1)}, x^{(3)}) = (1.93 - 1.4)^2 + (83 - 82)^2 = 1.2809$$

Dobrovolek 3 je tedy navzdory naší intuíci dobrovolekovi 1 mnohem „euklidovskyy blíž“ než dobrovolek 2. Je to tím, že jsme zvolili pro příznaků taková měřítka, že změna ve výšce je mnohem podstatnější než numericky stejně velká změna v obvodu pasu. Když bychom tedy hledali blízké vzorky podle euklidovské metricky, neodpovídalo by „blízké okolí“ našim představám.

Tohle se dá částečně řešit tím, že si vstupní data *normalizujeme*. Jeden ze způsobů normalizace je spočítat pro každý příznak jeho minimum a maximum na trénovací množině, pak všechny jeho hodnoty přeskálátovat do intervalu $[0; 1]$ a počítat euklidovské vzdálenosti na přeskálovanných příznacích. Stejně přeskálátování provedeme, když máme udělat předpověď pro nový vstup.

Když jsou všechny naše příznaky diskrétní, můžeme definovat *blízkost* podle takzvané Hammingovy vzdálenosti. Hammingova vzdálenost dvou vektorů je počet jejich slotů, které se liší.

Zhrvta agregace předpovědi z okolí. Pro regresi je nejčastější vztik výšky z nejbližších susedů a vrátit jejich průměr. Pro klasifikaci můžeme třeba vrátit tu kategorii, která má mezi K susedy nejvíce hlasů (a když jich není více hlasů dostalo nekoliik. vybírem z takových třeba tu první). Počítání průměru na kategoriích totiž nemusí být dobře definovaná operace.

⁷ Například signaly od receptorů bolesti nebo teploty se šíří rychlostí asi 0.5 až 2 m/s. Oproti tomu signaly od proprioceptorů – detektorů relativní pozice jednotlivých částí těla – se šíří rychlostí mezi 80 a 120 m/s. lidský mozek obecně funguje mnohem pomaleji než setřové počítače. Procesy v něm jsou však vysoce paralelizované.

Úkol 4 [6b]: Stráhněte si ze stránky seriálu dataset o kvalitě bíleho vína. Obsahuje 12 sloupečků oddělených čárkami. První řádek popisuje význam sloupečů. Prvních 11 sloupečů obsahuje různé chemické vlastnosti vína a posledních 11 je kvalita mezi 1 a 10, kterou se naučíte předpovídat.

Naprogramujte algoritmus K nejbližších susedů s normalizací příznaků a euklidovskou vzdáleností. Z kvality od K susedů vytrváte předpověď prostým aritmetickým průměrem.

Rozdělte dataset na trénovací, validační a testovací množiny v poměru 60%/20%/20%.

Mějte střední kvadratickou chybu na trénovací a validační množině v závislosti na K . Zkusťte všechna K od 1 do 40. Jak na K záleží chyba na trénovací množině? Proč? Jak na K záleží chyba na validační množině? Proč?

Vyběre nejbližší vypadající K a pomocí testovací množiny odhadněte, jaká bude skutečná accuracy vašeho modelu. Tabele úloha nejspíš pokržeí celkem dlouho. Jestli jí chcete pomoci, můžete zkusit využít více jader procesoru. K řešení přibalte svůj program a výsledky.

Opisování od evluce

Nejprošnější nám známý učební systém byl vytvořen přibližně čtyřmi miliardami let evoluce. Inspirování jeho architekturní výstavby mnoho koren *umělejší neuronových síti*. Otvírajší v posledních několika letech jsme díky několika novým trikem a rychlejším paralelním počítačům dosáhl úžasných výsledků, mimo jiné ve zpracování obrázu a zvuku, ale třeba i přirozeného jazyka.


Informace v lidském lidským mozku poskytli elektrické a chemické signaly. Dovedáme se, že nejtěžší větší „výpočetní jednotkou“ je neuron. Většina neuronů má nějaké „vstupní kanály“, kterým se říká *dendrity*, a jeden „výstup“ – *axon*. Dendrity tvoří jakousi „stromatitickou strukturu“, která sahá řádově stovky mikrometrů od těla neuronu. Zatímco dendritický strom některých neuronů má až tisíc větší, některé neuronů jich mají pouze jednotky. Délka některých lidských axonů dosahuje až 1 metru. Každý neuron má síce nejvíce jeden axon, ale ten se může mnohonásobně větvit a posílat jeho signály až stovkami jiných neuronů.

Elektrické signaly v nervovém systému jsou kodované pulzně. Neuron „sbírá“ signaly od svých vstupů a pokud se v neuronu nahromadí za určitou dobu dostatečně množství potenciálů, „vystřelí“ signál na výstup. Spojením množství neuronů se říká *synapse*. Některé synapse jsou *excitabilní* a některé jsou *inhibibilní*. Signály posílané po excitabilních synapsích „zvyšuje vnitřní potenciál“ clového neuronu, zatímco inhibiče mu „praní vstřelí“. Možná překvapivě se informace v mozku šíří relativně pomalu. V závislosti na typu signálu jsou to řádově jednotky až stovky metrů za sekundu.⁷

Systému neuronových spojení v mozku se souhrnně říká *konektivita*. Můžeme si její představu velmi zjednodušeně jako poněkud gigantický ortanovací graf. Jehož vrcholy tvoří jednotlivé neuronů. Hraný reprezentují spojení a vedou směrem, jakým se přenáší signály od neuronu, který je vypáli, do neuronu, který je má na výstupu. Kromě toho, že hrany mohou reprezentovat excitabilní nebo inhibiční spoje-

Máme před sebou dvě symetrické pyramidy (které už možná jako pyramidy vůbec nevyjadují) a chceme získat maximální součet, aby se stavělo právě k nějaké koprovací strategii. Jenže my začínáme a nemáme ani tu nejmenší jistotu, že součet bude naše řady koprovat. Co víc, nemáme ani jistotu, že součet bude hrát poslední.

Budeme tedy muset na dlíkat toho, že naše strategie funguje, ji jinak. Díže si klobouky, pojdeme z kopce.

 Nejprve si zavědeme další označení, tentokrát pro helmicce, na kterýž nelze žádná další (takže jejich slouzováním neshodíme nic dalšího). Tím budeme říkat *vohré*:

Všimněme si, že mezi všemi helmiciemi s momentálně nimaými hodnotou je alespoň jedna volná. To plyne z toho, že sněrem dohů hodnoty helmic neklesá. Také si snadno rozmyslíme, že my bereme vždy volnou helmici – alespoň nějaká z nich má nejmenší hodnotu, takže se nám určitě nvyplatí shazovat helmice víc.

Připrdíme teď jednotlivým helmiciem v jedné ze symetrických pyramid označení x_i taková, že $x_i \leq x_{i+1}$. Každé x_i se ve hře (resp. této hladové části hry) vyskytuje dvakrát – na začátku je jedinou v první pyramidě, jednou ve druhé.

Helmicce můžeme rozdělit na naše a na soupeřovy podle toho, kdo je shodil. K tomu si zavědeme další označení, a to *hlavní helmice*. To bude pro každý tah nejmenší z helmice, které hráč v daném tahu shodil. Pokud helmice shodil více, označme ty ostatní jako *vedlejší*.

Snadno si všimneme, že všechny naše helmice jsou hlavní (jelikož nikdy neshodíme víc než jednu helmici naráz).

Nyní bychom chtěli spátovat své helmice se soupeřovými, a to tak, aby vždy naše helmice měla maximálně takovou hodnotu, jakou má spátovaná soupeřova helmice. Některé soupeřovy helmice můzná znástanou nespátované, ale to nám nijak nevaří. Když by se nám takové párování podařilo najít, musí být náš počet kaminůk maximálně stejný jako soupeřovi.

Představme si, že vždy spátujeme naši hlavní helmici s hlavní helmici, kterou soupeř vezme v příštím tahu. Pro takové párování by nevostnosti určité platila. Jenže se může stát, že v posledním tahu hrájneme my, tedy zhrvta jedna nespátovaná helmice, kterou už není s čím spátovat. Ukážeme ale, že v takovém případě můžeme helmice přepárovat.

Označme nespátovanou helmici jako X a její hodnotu jako x_i . Nyní se podíváme na sufixy x_{i+1}, \dots, x_n všech hodnot, které byly ve hře. Existují dvě možnosti:

Zaprvé, existuje vedlejší soupeřova helmice s hodnotou z toho nebo sufixu. Jelikož jsme zatím párovávali vždy s hlavními helmiciemi, je tato vedlejší helmice nespátovaná a navíc má určitě aspoň stejně hodnotu jako X . Tedy můžeme X spátovat s touto vedlejší helmici. Tím máme vše spátováno a nevostnosti platí.

Zadruhé, žádná taková vedlejší helmice neexistuje. To znamená, že sufixu odpovídají pouze hodnoty hlavních helmice. Zároveň ořšen sufixu odpovídají hodnoty snadno počinu helmicce, a jedna z těchto helmice je X . To znamená, že alespoň jedna hlavní helmice je spátována s nějakou, jejíž hodnota je menší než x_i .

Označme helmice v tomto páru jako A a B , necht $A \leq B$ (tedy hodnota A v sufixu neleží, hodnota B ano). V dosavádím párování platí, že hodnota naší helmice je menší než soupeřovy, tedy A je naše. Přepárování nyní X s B ; A se stane nespátovanou.

Jelikož x_i bylo v sufixu nejvyšší, je určité hodnota X neostře menší než hodnota B , takže nevostnosti nám stále platí. Může se zdát, že jsme si nepomohli, opět máme jednu nespátovanou helmici. Její hodnota ale určité leží v řadě hodnot o jedna víc než vlevo, tedy opakovaně přepárování spátujeme všechny helmice v konečném čase.

Taký ještě zhrázáme, že pokud se dostaneme až do situace, kdy má nespátovaná helmice hodnotu x_i , musí mít soupeř nutně alespoň jednu vedlejší helmici. V opakem případě by muselo existovat párování vedoucí před sufixu, ale před x_i už nic není.

Vždy tedy umíme vytvořit párování takové, že naše helmice má maximálně stejnou hodnotu, jakou soupeřova, tedy i z hladové části hry budeme mít maximálně tolik, kolik získá soupeř.

Tím jsme dokázali, že popsaná strategie je skutečně vyhrávající. UHFH!

Kerry Burnsont

28-3-2 Liný Pisat

Označme si vždy jako řetězece A a B , bez útiny na obecnosti předpokládejme, že $|A| \geq |B|$. Podívejme se nejprve, jakým způsobem sestrojíme řetězece S takový, aby obsahoval po odstranění některých znaků každou větu a také byl nejkratší, neboli aby tento řetězece byl nejkratší spočetnou nadposloupeností A, B . Jak ale tento řetězece bude vypadat? Určíte se nám vůbec nvyplatí mít délku $|S| > |A| + |B|$, jelikož bychom museli přidat další zbyřené znaky. Řetězece dále můžeme krátit tím, že znaky společné oběma větám zapíšeme do našeho řetězece S jen jednou.

Tyto společné znaky však vůbec nemusí být na stejných pozicích u obou řetězců. Příkladem jsou řetězece

$$A = \text{xkxBCkxDXE}, \quad B = \text{klmABCDEkIm}.$$

kde společné znaky ABCDE tvoří nejlépeší spočetnou podobnost (zakrácené NSP). Znaky této NSP se ale v obou řetězcích vyskytují na různých indexech: v A na indexech $[1, 3, 5, 7, 9]$ a v B na indexech $[3, 4, 5, 6, 8]$. My potřebujeme zjistit jak samotnou NSP, tak indexy jejích znaků v obou řetězcích.

K vytvoření problému nalezení NSP lze využít dynamické ho programování, algoritmus je přímo popsán v kuchařce. Tento algoritmus nás stojí $\mathcal{O}(|A| \cdot |B|)$ času.

Jakmile máme indexy společných znaků z obou vět, můžeme sestavit náš požadovaný řetězece S . Máme i -ý společný znak, a_i a b_i indexy tohoto společného znaku v A a B . Potom postupně přidáme zvyšovat i od 1 po délku NSP a při každé iteraci vyjdeme všechny znaky z A, B mezi $(i - 1)$ -ním a i -tým společným znakem a nakoniec samový i -tý znak. Tento algoritmus nám zabere $\mathcal{O}(|A| + |B|)$ času.

Celkové algoritmus zabere $\mathcal{O}(|A| \cdot |B|)$ času. Dívod je takový, že nám právě nejvíce času potřvá vylíedání NSP. Samotné vypsaní řetězece trvá jen lineárně, což nám asymptotickou složitosť nezmenší.

liný pohled na útinu

Na tuto útinu existuje i řešený založené na odlišné myšljenje, které nakoniec bude stejně efektivní. Představme si orientovaný graf, jehož vrcholy jsou reprezentovány uspořádanou dvojicí $[a_i, b_i]$. Tyto vrcholy si můžeme graficky znázornit,

Vzorová řešení třetí série dracáčého osmého ročníku KSP

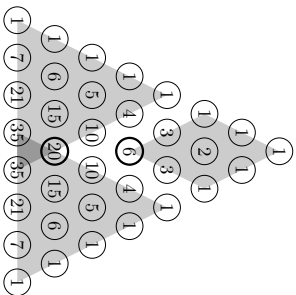
28-3-1 Pyramida z helmice

Úvodem dnešních řešení se vám všem musíme omluvit za neposrohodně poradující ohoďnocení pyramidové úlohy (a to ještě vypadalo krásně, že nějaká úloha vyšla jako první!). Pokud jste si tedy nad úlohou lámali hlavu a nemohli na je kloudného přijít, možná jste jen projevili víc rozumu než obvykle, kteří kdysi určili tomu, že úloha je přece jednoduchoučka.

Doufáme, že příště bude bodové ohoďnocení lépe vyhovovat o obřízosti úlohy, a hlavně že jsou někoho od řešení neodradili. Ale teď už hoňam na řešení.

Snadno si můžeme rozmyslet, že pro $N \leq 2$ vždy vyhrají druhý hráč. Naopak dost komplikované můžeme dojít k tomu, že pro $N \geq 3$ existuje vyhávající strategie pro prvního hráče.

Nejprve tuto vyhávající strategii popíšeme, a pak teprve dokážeme, že opravdu funguje a soupeř nemá šanci. Zaverdme si teď pář oznamení, ať se nám strategie popíše snáz. Předtím, *lehkým řádkem* budeme myslet řádek obsahující l dlý počet helmic. *Druhý stříd* pak bude prostřední helmice na nejspodnějším lichém řádku a *levý stříd* prostřední helmice na druhém lichém řádku odspodu. Také si dovolíme používat pojem pyramida, přestože tvar, který budon helmice v průběhu hry vytvářet, nebude vždy pyramidou připomínat.



1. Hned v prvním tahu shodíme levý stříd. Tím jakoby vznikly dvě shodné pyramidy, které se částečně překývají.
2. Potom budeme „kopírovat“ soupeřovy taly, dokud soupeř nestříd do dráhého středu. Jinými slovy, soupeř stříd do nějaké helmice v jedné z pyramid, my strčíme do stejné helmice ve druhé pyramidě.
3. Když soupeř strčí do dráhého středu, začneme hrát hladově. To znamená, že vždy shodíme tu helmici, která právě obsahuje nejméně kaminček.

Dokazování vezmeme trochu na přesáčku. Začneme tím, že opravdu můžeme hrát podle dráhého bodu. Pyramidy vzniklé po našem prvním tahu totiž mají prímek (překryv) a nemají být jasné, že pokud soupeř strčí do něčeho z tohoto prímku, můžeme kopírovat.

Ten prímek ale není velký: obsahuje draly stříd a připadne



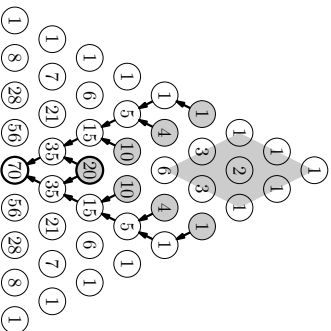
(Je-li N sudé) dvě helmice pod ním. Draly stříd máme ošetřeny zvlášť. Helmice nikdy neshodí helmici stojící vedle ní, takže shodí-li soupeř jednu, my shodíme druhou. Kopírovat tedy můžeme vždy.

To, že každý tah komé strčení do dráhého středu umíme zkopírovat, je důležité. Právě proto to musí být soupeř, kdo strčí do dráhého středu, připadne ho shodí spolu s helmici pod ním.

Jak se kopírovací část strategie projeví na počtu sebraných kaminček? Určité platí, že když okopírujeme tah, sebereme maximálně tolik, kolik sebral soupeř. Ve skutečnosti sebereme ostře méně právě tehdy, když soupeř shodí něco z prímku pyramid. V takovém případě ale určité shodí draly stříd. To dohornady znamená, že po konci kopírovací strategie se rozdíl naščí a soupeřových kaminček zvyší alespoň o hodnotu dráhého středu.

Nesmíme ale zapomenat, že před kopírovací částí jsme sami sebrali levý stříd. My ovšem ukážeme, že součet kaminček v helmicích, které shodíme shozem levného středu, je o jedna menší než hodnota dráhého středu. Když tohle platí, máme po kopírovací strategii alespoň o jeden kaminček méně než soupeř.

Indukci podle r budeme dokazovat obecnější tvrzení: Počet kaminček získaných shozem středu helmice v r -tém lichém řádku shora je o jedna menší než počet kaminček ve středu helmici na $(r+1)$ -ním lichém řádku.



Pro $r = 1$ naše tvrzení evidentně platí ($1 = 2 - 1$). Indukční krok sledujeme na obrázku. Pokud shodíme střed 3. lichého řádku, což je číslo 6, spadne celý „dřaman“ a z indukce už víme, že jeho součet je 19.

Shodíme-li stříd o řádek níž, tedy číslo 20, spadnou navíc dvě diagonály vedoucí od okraje ke středu 20. Šípky na obrázku ukazují, že každá z těchto diagonál se přeseče na 35, takže obě diagonály dohornady na 70, což je přesně stříd následujícího lichého řádku.

Ovšem číslo 20 leží v obou diagonálách, takže jsme ho započítali dvakrát. To je ale správné: jednou reprezentuje samo sebe, podruhé součet dřamanu nad diagonálami zvěřšený o jedničku. Analogická úvaha funguje pro libovolné r .

Hurá! Tím jsme ukázali, že po kopírovací části bude mít soupeř alespoň o jeden kaminček víc. Pokud ale soupeř shodí draly stříd přímou, pokrácujeme třetí částí (hladovou). Teď musíme ukázat, že 1 z této části získáme nejvýše tolik, kolik soupeř.

existuje samozřejmě velon mnoho dalších parametrů, které mají vliv na zpracování signálu. Učení probláh znamená vlastnosti konkatomu, například přidáváním nových symapsí nebo změnami parametrů těch symapsí, co už existují. Když symapsičky spojujeme neourony páři společně, jejich spojení se posiluje.

Umělé neuronové sítě modelují skutečnost s různu úrovni detailu a mají mnoho různých aplikací, podle kterých se odvíjí její architektura. Začneme od nejjednoduššího modelu, který dává něco užitečného.

Umělý neuron

Namísto posílání pulzních signálů budeme reprezentovat informací pomocí čísel. Umělý neuron je pro nás jednotka, která má n číselných vstupů x_1, \dots, x_n , a jeden výstup y . Neuron bude ze svých vstupů počítat *potenciál* ξ . Některé vstupy budou k potenciálu přispívat pozitivně, některé negativně. Výstup y bude záviset na potenciálu podle takzvané *aktivované funkce*, značené $\sigma: y = \sigma(\xi)$. Směr přispívání, tedy míra excitace nebo inhibice, bude pro každý vstup jiná a bude určena takzvanou *váhou* (*weight*). Váha i -tého vstupu se značí w_i .

Standardně se potenciál počítá jako součet vážených příspěvků všech vstupů. K tomuto součtu se také přidá takzvaný *bias* (český překlad by mohl být „předstudek“ nebo „sklon“) θ :

$$\xi = \theta + \sum_i w_i x_i.$$

Aktivací funkce se používají různé. Vždy jsou dehnované na celém \mathbb{R} , omezené (většinou na $[-1, 1]$ nebo $[0, 1]$) a rostoucí (respektive neklesající). První model, který si ukážeme, bude vracet 1, pokud je $\xi \geq 0$, a jindy -1 (tedy jeho aktivací funkce je funkce signumu).

Perceptron

Takovému umělému neuronu se říká *perceptron*. Jeho výstup je 1, pokud $\theta + \sum_i w_i x_i \geq 0$, a jindy -1 . Jsou-li vstupy perceptronu dva, je $\theta + \sum_i w_i x_i = 0$ rovnice přímky ve dvou rozměrech. Perceptron tohle přímku rozděljuje dvojnásobně prostor vstupů na dvě poloviny. V jedné vrací 1 a ve druhé vrací -1 . Přidáme-li další vstupy, analogicky perceptron dělí n -rozměrný vstupní prostor na dva poloprostory.

Perceptron lze použít jako jednoúčelý klasifikátor dvou kategorií, pro něž bude jeho očekávaný výstup 1 a -1 . Rozdělíme si vstupní příznaky trénovacích dat do *pozitivní* *kategorie* P a do *negativní* *kategorie* N .

Ukážeme si *algoritmus perceptronové učení*, o kterém víme, že pokud jsou kategorie P a N od sebe bez dých rozdílu (perceptronem (neboli nějakou nadrovinou), pak náš algoritmus nakonec najde váhy nějakého takového perceptronu. Nejdříve ke všem vzorkům v P a N přidáme jako novou složku na začátek jedničku. Tím pádem můžeme zapome-

nout na bias – stane se z něho nová, první váha. Vzorky „vyplešně“ o jedničky označme jako P' a N' .

Potom sledujeme vzorky do jedné množiny: N' mají být vzorky, pro které $\sum_i w_i x_i < 0$, a to platí právě tehdy, když $\sum_i w_i \cdot (-x_i) > 0$. Vyrobíme si množinu R , která se bude skládat z P' a minus jednotkou vynásobených vektory z N' . Perceptron, ve kterém jsou všechny vektory v R v pozitivní polovině, bude správně klasifikovat P' i N' .

Zinicializujeme perceptron libovolnými vahami, třeba nulovými. Učení probláh tak, že iterujeme přes množinu R a postupně perceptron učim jednotlivé vzorky x . Chceme po něm, aby na každém vzorku bylo $\sum_i w_i x_i$ větší než 0.

Pokud na zkontrovaném vzorku je $\sum_i w_i x_i > 0$, je vzorek perceptronem rozpoznáný správně, neděláme nic a jdeme na další vzorek. Pokud je sůma menší než 0, pak ke každé váze w_i přičteme $x_i \cdot \gamma$ a jdeme se učít další vzorek. γ je tedy opět rychlost učení, konstanta mezi 0 a 1. Pokud najdeme perceptron, který správně klasifikuje všechny vstupy, vrátíme ho. Pokud jenom chceme dobrý perceptron, který ne nutně klasifikuje všechno dobře, můžeme místo toho jenom běžet, dokud nám nedojde čas, a pak vrátit ten perceptron, který dokázal klasifikovat nejvíce vzorků po sobě správně.

Úkol 5 [ob]: Strámně si ze stránky seriálu datasetu o kosatcích. Tento slavný „ Iris dataset“ poprvé použil v 30. letech významný biolog a statistik Ronald Fischer. Každý řádek reprezentuje vlastnosti jednoho kosatce. První řádek popisuje význam sloupčků: první dva jsou délky a šířky kalíšních listků, další dva jsou délky a šířky okvětních listků. Máme za úkol předpovědět poslední sloupec – konkrétní druh kosatce: *setosa*, *versicolor*, nebo *virginica*.

Zkusíte naprogramovat rozpoznávání kosatců pomocí perceptronu. Jeden perceptron dokáže od sebe rozpoznat jenom 2 třídy; budete muset být kreativní. Nemí jedině správné řešení – překvapte nás! Jakou zvládnáte accuracy na testovacích datech? Naše autorství řešení umí z 60 trénovacích vzorků mít na zbytku datasetu 93,42%.

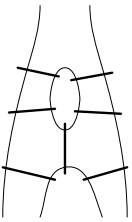
Perceptrony jsou užitečné, ale počítá docela slabě – pokud mají kategorie „složitou hranici“ a nejsou lineárně oddělit, potřebujeme na jejich naučení našim modelům povolit složitější strukturu.

Zábavné věci se začínou dít třeba, když začneme neourony *vrstvit* do takzvaných *dopředních usměrněných neuronových sítí* (*feedforward neural networks*). První vrstva takové sítě jsou vstupy data. Druhá vrsta konzumuje výstup první vrstvy jako svůj vstup. Sedí v ní například 10 neuronů a každý z nich má vlastní váhy a transformuje vstupy jiným způsobem. Vstupy druhé vrstvy jsou vstupy pro třetí vrstvu, a tak dále až do výstupní vrstvy. Informace se šíří jenom z nižších vrstev do vyšších. Vyšší vrstvy jsou schopné počítat složitější a složitější transformace vstupních dat.

Michal Pokorný

Historický problém

V roce 1736 se švýcarskému matematikovi Leonhardu Eulerovi na stůl dostal na první pohled jednodušeý problém, který mu předložil starosta města Královce (dnešní Kalinin-grad). Královcem teče řeka Pregola, na ní je několik ostrovů a ostrovy jsou spojeny se zbytkem města mosty. Dobaova ilustrace situaci vyzníhla takto (schematická kresba):



Pan starosta se pana matematika v dopise tázál, jestli je možné začít na některém z břehů (nebo ostrovu) a udělat si vycházku po městě tak, že se každým mostem projde právě jednou. Navíc chtěl procházkou skončit na knsu suché země, ze kterého vyšel.

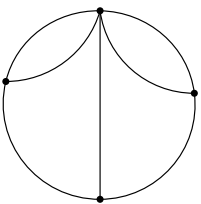
Profesor Euler jej nejprve chtěl poslat k šípku – problém jde snadno vyřešit rozložením případů, což by zvládli i tehdejší studenti střední školy (natož pak ti dnešní). Zachoval se ovšem jako první matematik – přišel na to, jak problém zohlednit, a mistrně vyřešil hádanku i pro všechna možná města, která kdy budou dříve potřádat podobné procházky.

Eulerovský tah

Pojdme si nyní problém popsat abstraktně a tím si připomenout grafovou terminologii. Vrcholy našeho grafu jsou kusy pevniny; ať už to budou části města nebo ostrovy. Mezi dvěma vrcholy povede hrana, pokud jsou spojeny mostem, a omen most odpovídá hraně.

V tomto zadání má smysl uvěřit, že mezi dvěma kusy pevniny povede mostů více – například v Praze jich vede tolik, že se na to přajší v leckteré zeměpisné olympiádě. Graf, kde mezi vrcholy vede více hran, nazýváme *multigraf*, a pokud dvě hrany vedou mezi stejnými vrcholy, mluvíme o nich jako o *paralelních* hranách.

Obecná procházka v grafu z vrcholu A do vrcholu B (posloupnost hran taková, že cílový vrchol předloží hrany je počátek vrcholu hrany následující) se nazývá *stezka z A do B*. Ve steh se mohou opakovat jak hrany, tak vrcholy; sled tedy není řešením našeho problému (ve sledu je možné se vrátit po hraně, ze které jsme právě přišli).



Po naší úlohu se hodí posloupnost hran taková, že vrcholy se opakují mnohokrát, ale hrany nikdy. Těto posloupnosti se říká *tah z A do B*. Kdyby se neopakovaly ani vrcholy, pak posloupnost označujeme jako *cesta*. Tah (respektive sled) je *uzavřený*, pokud začíná v A a končí také v A .

Podivněmi se tedy na mapu Královce jako na multigraf, přáme se, zdali existuje uzavřený tah takový, že každou hranu navštíví právě jednou. Takovému tahu pak říkáme *uzavřený eulerovský*.

Minulodnem, tahu se „tah“ někdy jen tak náhodou. Děti

se často ve škole přiklonávají v umění nakreslit obrázek

jedním tahem, aby se tužkou nemuselo vtáčet po už nakreslené čáře. Pokud si obrázek představíme jako graf (čárky jsou hrany, místa jejíh setkání vrcholy), pak eulerovský tah nalezneme jen v tom obrázku, který lze nakreslit jedním tahem. V uzavřeném eulerovském tahu se pak vrátíme i do mistra, kde jsme začali.

Podmínky tahu

Je na čase podívat řešení našeho problému s eulerovským tahem. Půjde nám na to jako matematici – nejprve ukážeme *nutnou* a hned nato *postčůující* podmínku. Nutná vlastnost grafu je taková, že bez ní eulerovský tah není možné najít; postačující vlastností je ta, se kterou vždy eulerovský tah najít umíme. Jsou-li obě podmínky stejné, pak se jedná o ekvivalenci, a tak tomu bude i nyní.

Představme si, že jsme kouzlem nějaký uzavřený eulerovský tah našli, ať už je jakýkoli. Vždy, když se dostaneme do jednoho vrcholu (a není důležité, jestli už jsme v něm byli, nebo ne), tak z něj musíme hned také odejít, abychom tah uzavřeli. A protože tah je eulerovský, každou hranou projdeme jen jednou, takže tyto dvě hrany (tu přichází a odchází) už nepoužijeme. U každého vrcholu mimo výchozí tedy platí, že hrany tvoří dvojici – jedna, co vedla dovnitř, a jedna, která z něj vedla ven.

Podobná věc platí i pro startovní vrchol. Sice do něj nevstoupíme poprvé pomocí hrany, takže počet navštívených hran u něj bude stále lichý – ale jen do chvíle, než se do něj naposledy vrátíme a skončíme, protože skončením jsme použili poslední hrany, která bude tvořit dvojici s hranou první.

Jakou vlastnost grafu jsme odhalili? Nepochá, že graf má sudý počet hran (protože trojhlavní jedním tahem nakreslíme a přesto má 3 hrany), ale platí, že do každého vrcholu vede sudý počet hran, tedy že graf má *všechny stupně sudé*. Nezapomíname také na to, že graf musí být souvislý – dva oddělené obrázky jedním tahem bez zvednutí tužky nemakreslime. Máme nutné podmínky!

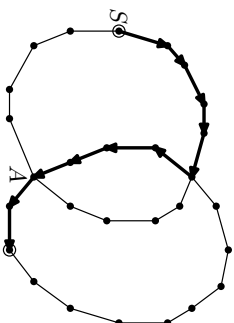
Nalezení tahu

Zbývá tedy ověřit, že podmínky jsou i postačující. Mějme souvislý graf, který má všechny stupně sudé. Umnme v něm vždy najít uzavřený eulerovský tah? Ověřme to, jak se na informatiky patří – algoritmem.

Předložený algoritmus je založený na vylepšeném prohlédávání do hloubky, tedy DFS, o kterém jste si mohli přečíst v první grafové kuchařce.⁸

Vyberme si vrchol, v něm začneme. Naš algoritmus musí umět označovat hrany jako „probrané“, jako to dělá DFS. Vyberme si tedy jednu hranu a pokračujme dále, zatím bez vypisování.

Po nějakém tom prohledání se jistě stane, že jsme se zastavili – vrchol už nemá žádné nepoužité hrany. Nutně to znamená, že to je ten vrchol, ve kterém jsme začínali. V prohledání do hloubky se vrátíme zpět, ale my k tomu přidáme vypisování cesty – postupně pozpátku vypisujeme hrany, kterými se vrátíme zpět v prohlédávání.



Na obrázku výše je příklad právě prohlédávho algoritmu. Začal ve zvýrazněném vrcholu vlevo, procházel po šípčkáč az do bodu A , kde vohl hrany tak, že hned skončil na začátku. Dále pokračoval vypisováním hran pozpátku, až došel zase do bodu A . Zde si vybral jednu ještě nepoužitou hranu a po ní přešel celou druhou kružnicí – zbytek hran – zpět do bodu A . Nyní vypisuje hrany pozpátku od bodu A .

Bud tímto vypíšem dojdeme až na začátek, nebo se dostaneme do vrcholu, který má ještě nějaké nepoužité hrany (stručne může vypadat třeba jako na obrázku). Potom vypisování zastavíme a pokračujme v prohlédávání DFS přes nepoužitou hranu. I tam se to může zastavit (a zastavit), i tam začneme vypisovat pozpátku. Nakonec dojdeme do původního místa rozbočení, a budeme opět pozpátku vypisovat hrany, které nás nakonec dostanou až na počátek, kde skončíme.

Najde tento algoritmus opravdu konkrétní uzavřený eulerovský tah? Graf byl souvislý a o algoritmu DFS se ví, že v takovém případě navštíví každou hranu právě jednou. Algoritmus opravd vypisuje cyklus – jen je u něj trochu zvláštní způsob, jak ho vypisuje. Když dojde na křžovátku s ještě nepoužitými hranami, tak vypis zastaví, těso po nich kráčí, označuje si je a vypisuje, až když se po nich vrátí. Ověřme si, že hrany opravd navazují.

V duchu argumentů z předcházející části víme, že jediný vrchol grafu s lichým počtem nepoužitých hran je právě ona křžovátka – a algoritmus DFS prochází graf podobně, jako jsme ho procházeli v minulé sekci, takže právě do tohoto vrcholu algoritmus dojde, až se přichod tomto části grafu zastaví.

Jakmile som program dojde (a nezhnou mu volné hrany), začne cestovat zpět a hrany vypisovat – a opravdu, pokračuje se tedy z mistra, kde naposledy přestal, a program vskuknu vypíše tah přes všechny hrany v grafu – uzavřený eulerovský tah.

Věta o eulerovském tahu v celé své kráse tedy zní: *(Mut- ti)gnof obsahuje uzavřený eulerovský tah právě tehdy, když má všechny stupně sudé a je souvislý.*

Je třeba podotknout, že složitost našeho algoritmu na bázi DFS je lineární vůči velikosti grafu (počtu vrcholů a hran). Existují i jiné algoritmy pro hledání eulerovského tahu, jed- na varianta například prochází grafem a vyběhá si na křžovkách takové hrany, které souvislé grafu pokud možno nepouždou. Tyto algoritmy už nemusí mít nutné lineární časovou složitost.

Jiné druhy procházek

Nějen kreslením obrázků ze stejného bodu živ je človek. Co kdybychom mohli začít a skončit v jiném místě, tedy ptali se po neuzavřených eulerovských tazích, znenúto by se něco? Není tomu tak, pouze mnme a postačující podmínky si vyžadají, aby všechny vrcholy měly sudý stupeň až na právě dva vrcholy, které mají lichý stupeň. Pokud nám to nevěříte, zkuste si to rozmyslet sami; opravdu to není těžké. Smysl také dává zkusit najít ne uzavřený tah, ale uzavřenou cestu – uzavřenou cestu přes všechny vrcholy, která navštíví každý vrchol právě jednou (říká se jí „Hamiltonovská cesta“). Bohužel, ačkoli jsou problémy přibuzné, musíme vás zklamané – není známe žádné efektivní (polynomiální) algoritmus na tento problém, a kdyžby jej někdo z vás našel, vyřešil by otázku „P vs. NP“, o níž se více dočtete v kuchařce o těžkých problémech.⁹

V matematice se také někdy zmnují „náhodné procházky“ po grafech – můžete si je představit tak, že se po mostech města Královce motá oplec, který si hází (oplnou nebo spravdivnou) minci a podle toho se rozhoduje, přes který most jít dál. Použití mají tyto modely hlavně v matematické teorii grafů a teorii pravděpodobnosti. O tom si můžeme povědět zase někdy jindy.

Martin Bohm

⁸ <http://ksp.mff.cuni.cz/viz/kucharka/grafoy>

⁹ <http://ksp.mff.cuni.cz/viz/kucharka/tezke-problemy>