

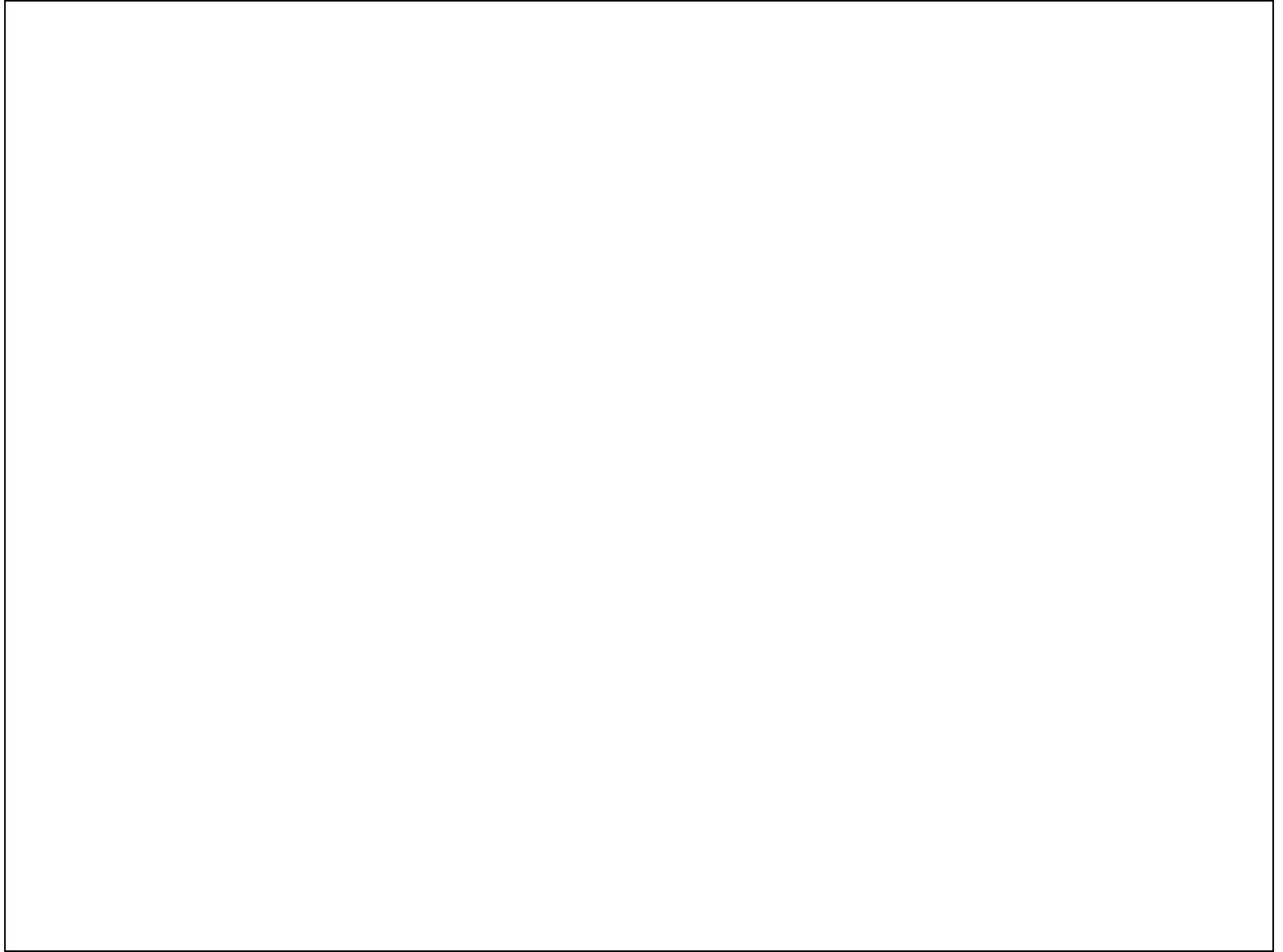
Hľadanie robustných a klinicky relevantných podtypov vo vysokopokryvných molekulárnych dátach

príkladová štúdia: kolorektálny karcinóm

Eva Budinská

Letná škola Matematickej Biológie

Mikulov, 14.9. 2012

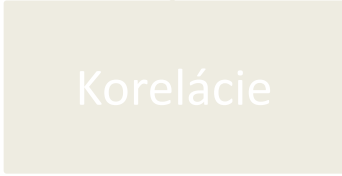
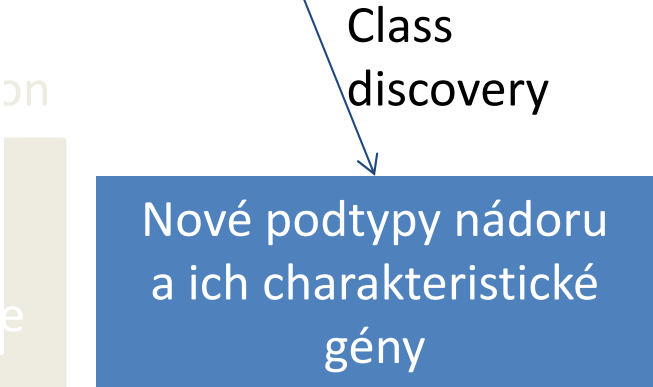


Genomické dáta
jedného druhu nádoru
Karcinóm kolorekta

Genomické dáta
Iného druhu nádoru
Karcinóm prsu



EVA BUDINSKA
Hľadanie robustných a
klinicky relevantných
podtypov vo
vysokopokryvných
molekulárnych dátach



Klinické dáta

Biologická a
klinická
interpretácia

Histopatológia

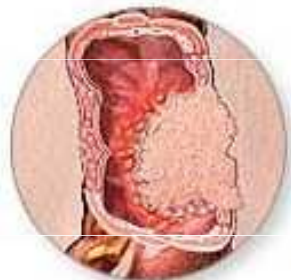
Prežitie

podtypmi génov

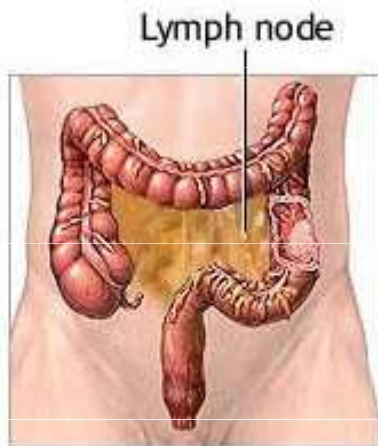
KOLOREKTÁLNY KARCINÓM



Stage I

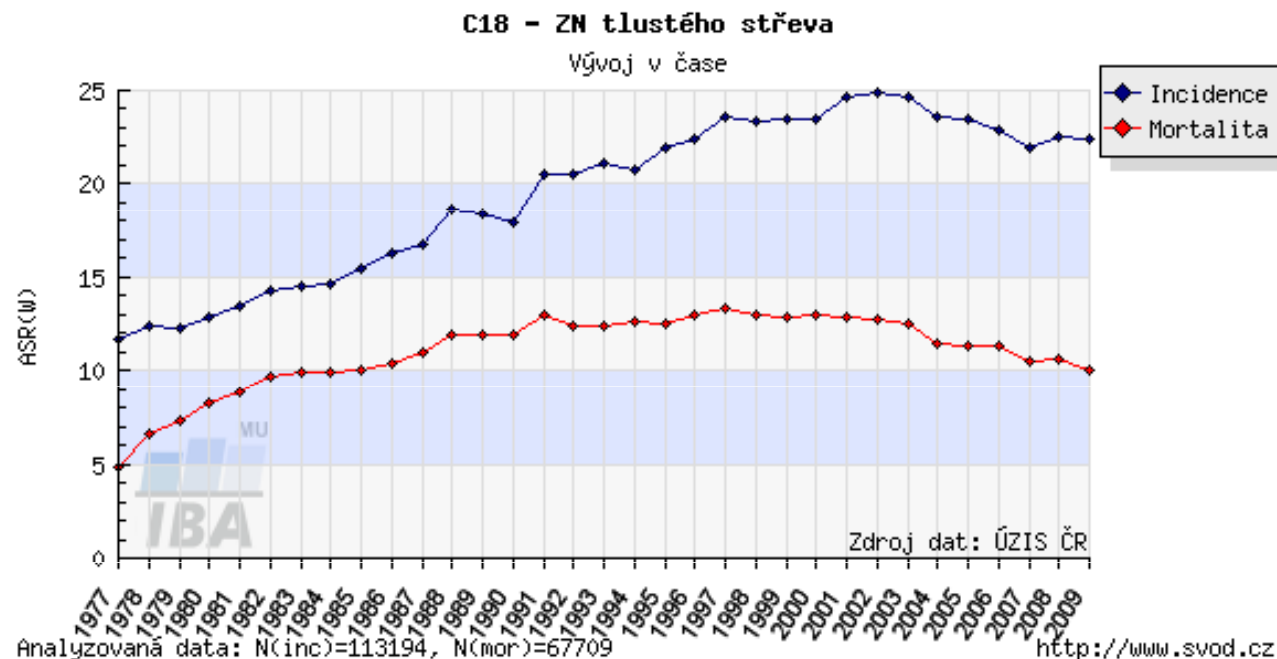


Stage II



Stage III

Colon Cancer



Heterogénne ochorenie s rozdielnou odpoveďou na terapiu.

Len niekoľko klinicky používaných markerov:

- *BRAF/KRAS mutácia – pre kvalifikáciu na antiEGFR terapiu (u štádia s metastázami)*
- *MSI – mikrosatelitová nestabilita – všeobecne považovaná za dobrý marker*

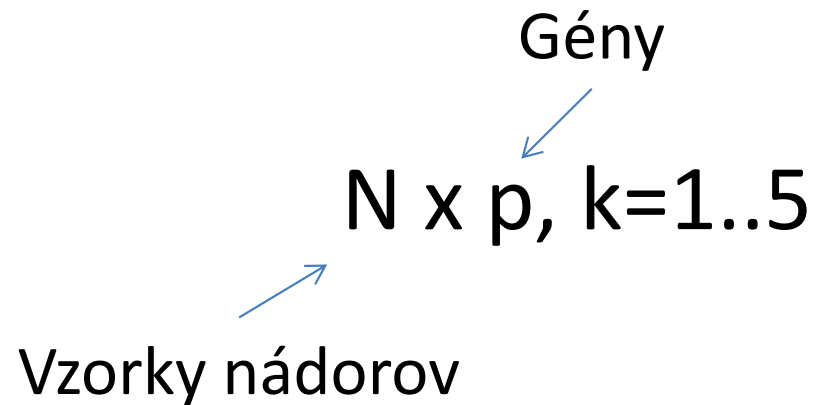
Cieľ:

*Nájsť skupiny nádorov kolorekta s podobnou expresiou génov (podobným génovým profilom)
~ podtypy*

Charakterizovať tieto podtypy pomocou klinických a známych molekulárnych parametrov.

Dátové súbory:

Matica obsahujúca kvantitatívnu expresiu génovej aktivity nádorov



Matica klinických a molekulárnych parametrov ku každej nádorovej vzorke, vrátane prežitia pacienta

Podtypy:

Neznáma pravda => Snaha o čo najobjektívnejšiu a najvšeobecnejšiu identifikáciu skutočnosti

Prístup:

Zhlukovanie – class discovery

Obmedzenia:

- (Ne) reprezentatívnosť populácie v dátach
- Technické možnosti
- Typ dát určuje uhol pohľadu – génová expresia, epidemiológia, histológia, DNA mutácie...)
- Odchýlky v dátach môžu podstatne ovplyvniť výsledok

Nevýhody zhlukovania:

- Rozdelenie dát do zhlukov i v prípade neexistencie skutočnej vnútornej štruktúry dát (štruktúra je náhodná)
- Potrebné metódy ad hoc stanovenia správneho počtu zhlukov
- Šum v dátach negatívne ovplyvňuje výsledok
- **Jedna veľká množina génov z biologického motívu ovplyvní zásadne zhluky**

VÝBER PREMENNÝCH / ZMENŠENIE DIMENZIE DÁT

Nezávislý výber – z 25 000 génov, výber podmnožiny 3025 génov s najvyššou variabilitou v súbore

Redukcia dimenzionality – práca s génovými modulmi

génový modul – sada génov s rovnakou génovou expresiou

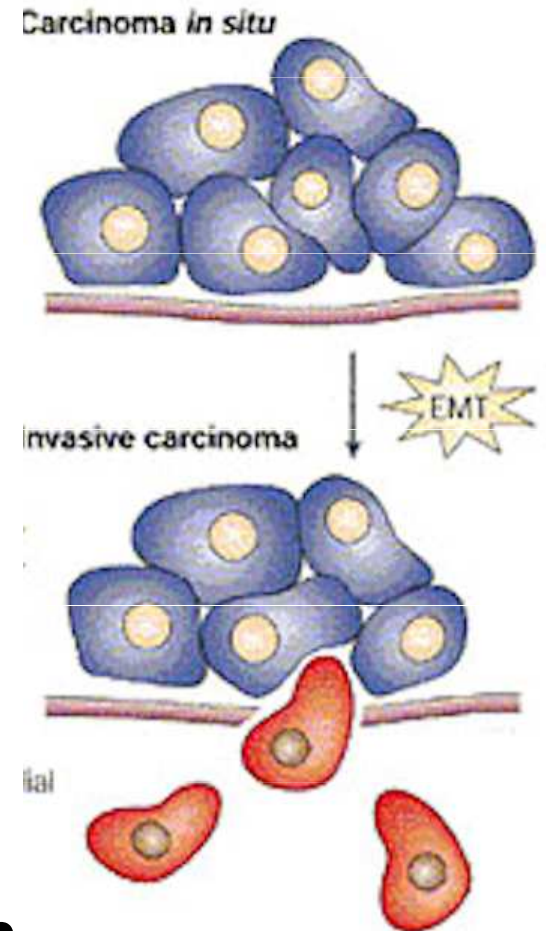
Jedná sa o istú formu váženia efektu biologických motívov

Predpoklad:

sada korelovaných génov \sim biologický motív

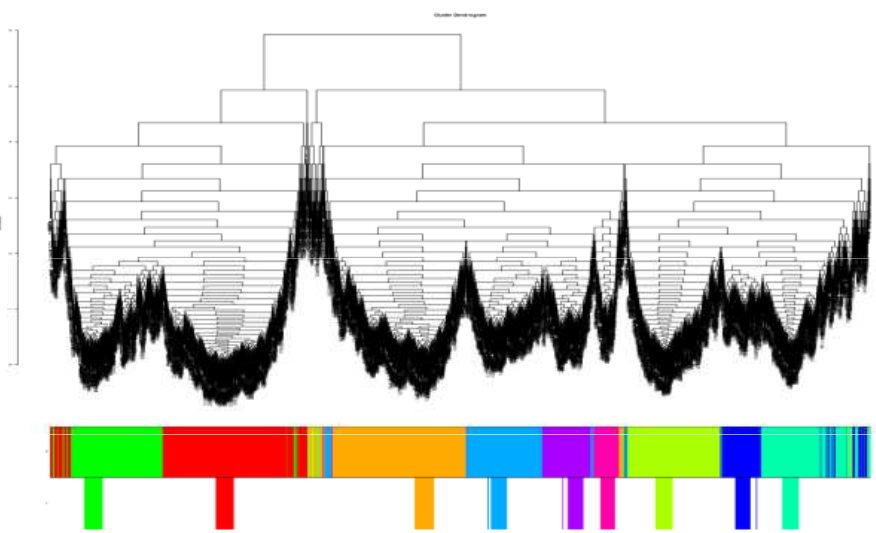
Príklad EMT

- EMT – epiteliálno mesenchymálny prechod
- Génová expresia podobná zdravému mezenchymálnemu tkanivu
- Obvykle reprezentovaný **zmenou v stovkách génov**
- Identifikácia modulu EMT a jeho reprezentácia jednou hodnotou (priemerom) **zmenší jeho efekt v zhlukovaní** a dá šancu **d ďalším dôležitým procesom** reprezentovaným menším množstvom génov



3025 génov
Dátový súbor 1

Pearsonova korelácia,
Hierarchické zhlukovanie
Complete linkage
Rezanie dendrogramu



150 génových modulov

Validácia ich korelácií v
4 ďalších dátových súboroch

Výber 54 modulov
(625 génov)

Medián génov v module

54 Meta-génov

Pearsonova korelácia,
Hierarchické zhlukovanie
Complete linkage
Rezanie dendrogramu

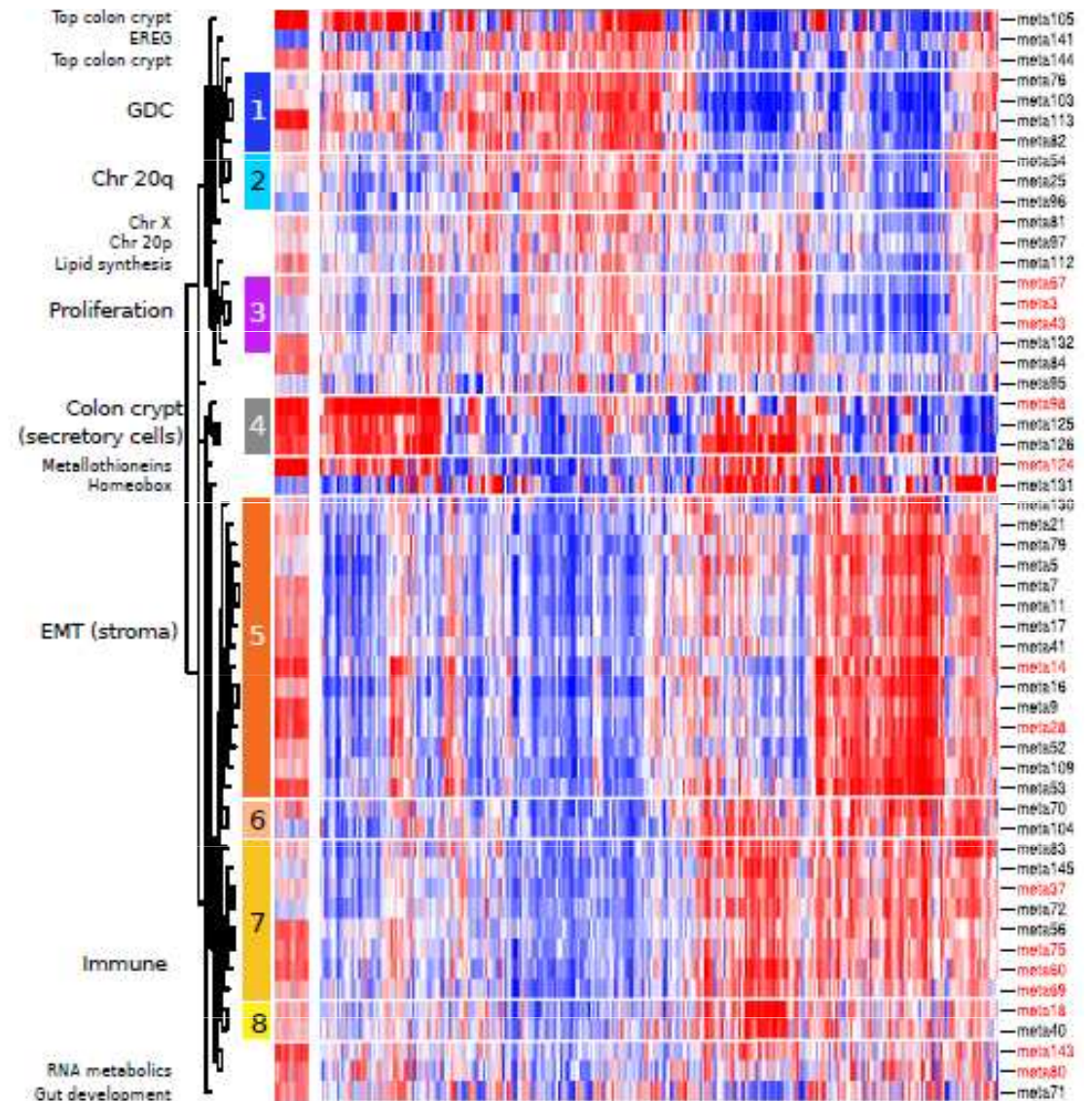
8 hlavných zhlukov
~ hlavné biologické
motívy

Identifikácia biologických motívov

- Analýza génových sád

| | V genóme | V biol motíve |
|----------|----------|---------------|
| V genóme | a | b |
| V module | c | d |

Biol motív:
EMT, proliferácia,
Chromozóm 20q...

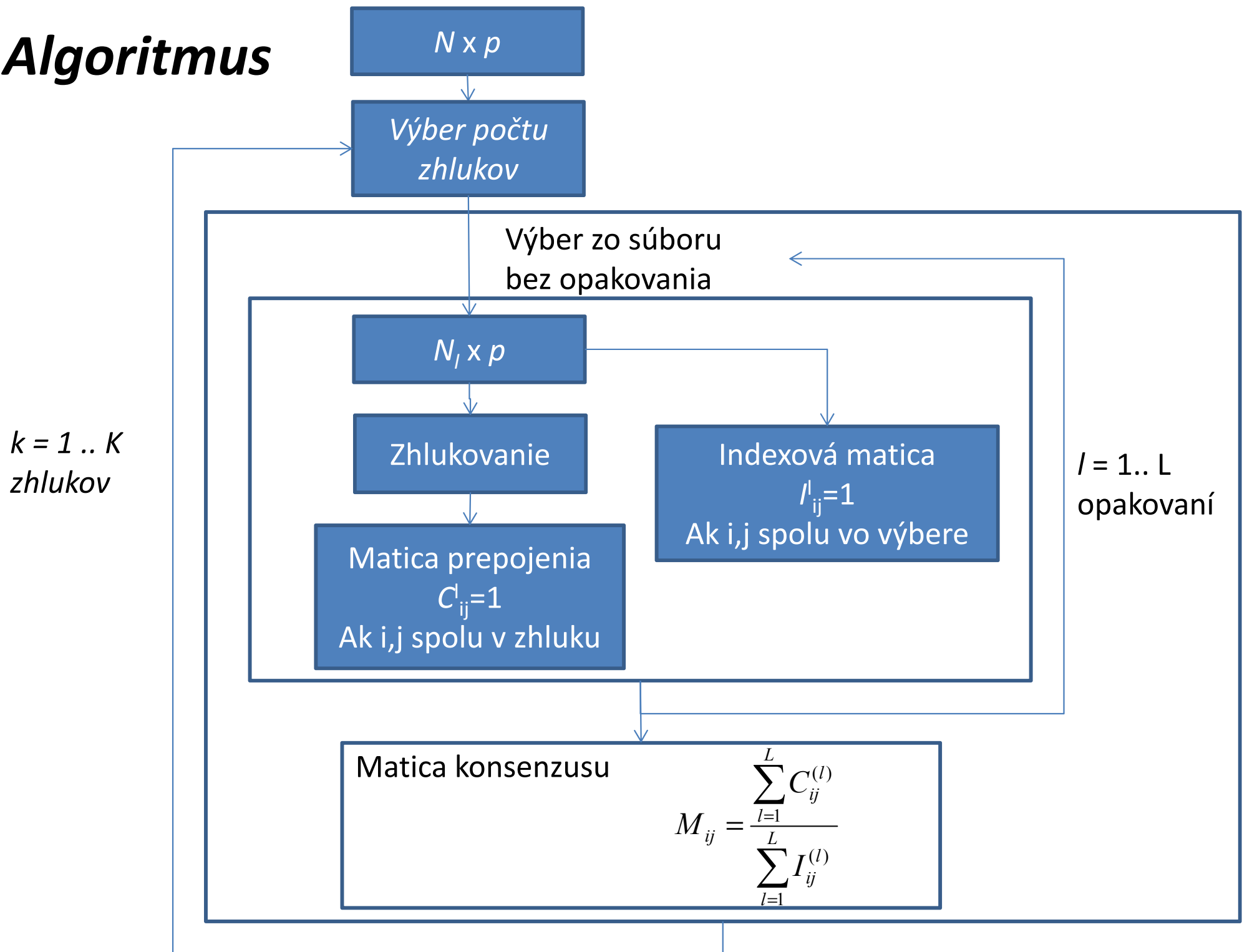


ROBUSTNOST V ZHLUKOVANÍ

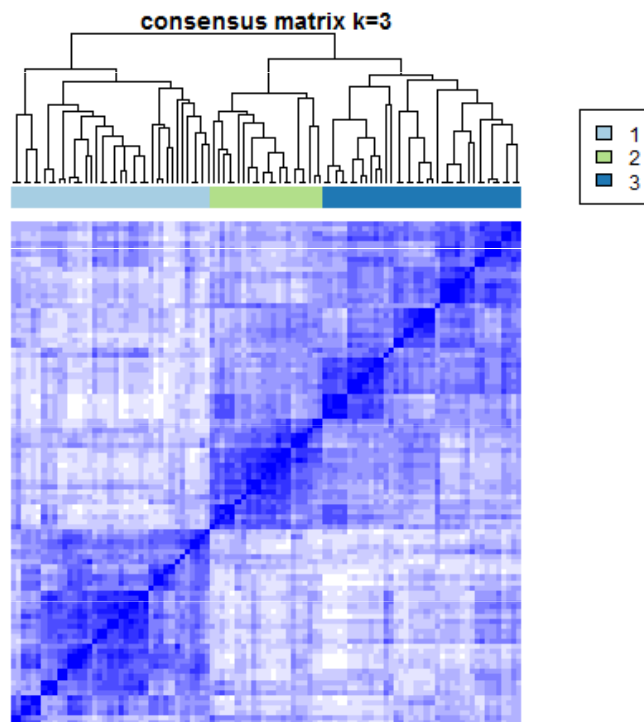
Consensus clustering

Metóda založená na opakovanom prevzorkovaní (<90% pôvodnej populácie a/alebo premenných) a zhlukovaní, za účelom vytvorenia konsenzusu medzi viacerými opakovaniami zhlukovania.

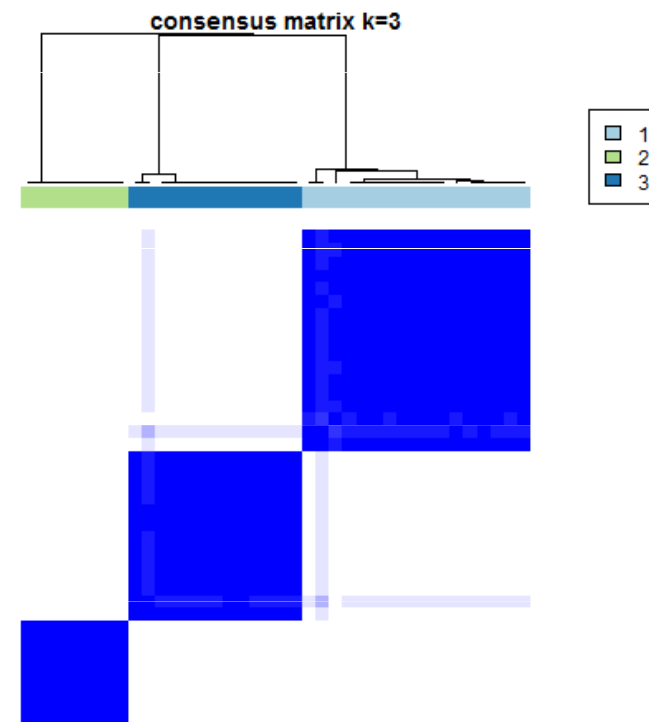
Algoritmus



- Ukazuje zhlukovanie stabilné aj pri drobnom narušení štruktúry dát – zhluky sú robustné
- Konsenzus – je nová metrika podobnosti (0,1)
- Pomáha určiť : počet a stabilitu zhlukov



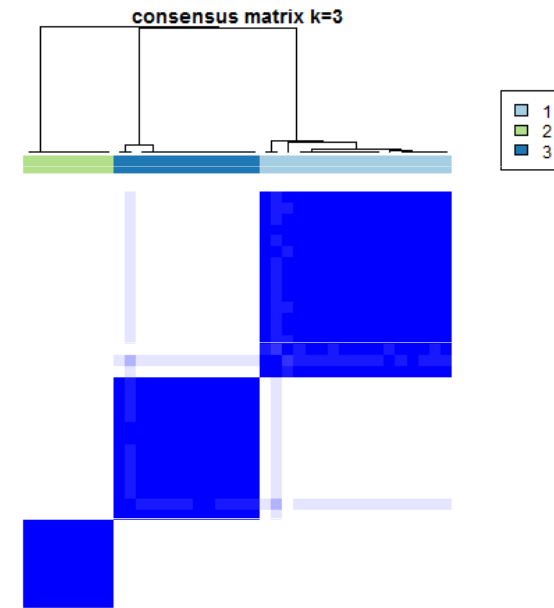
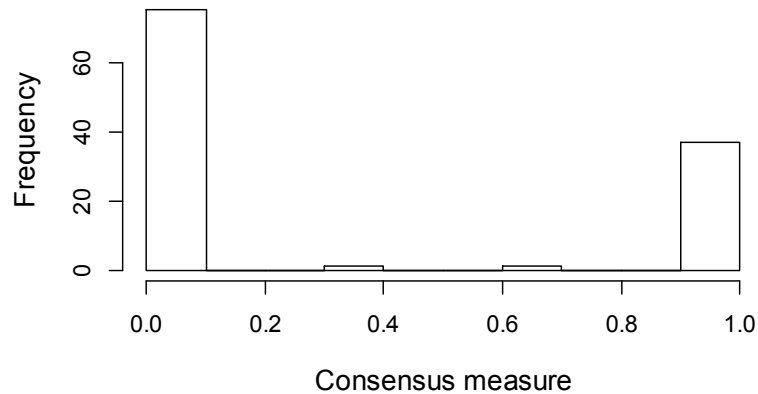
Náhodný výber z normálneho rozloženia



Leukémia:
ALL-T, ALL-B, AML

Odhad počtu zhlukov I

Histogram miery konsenzusu

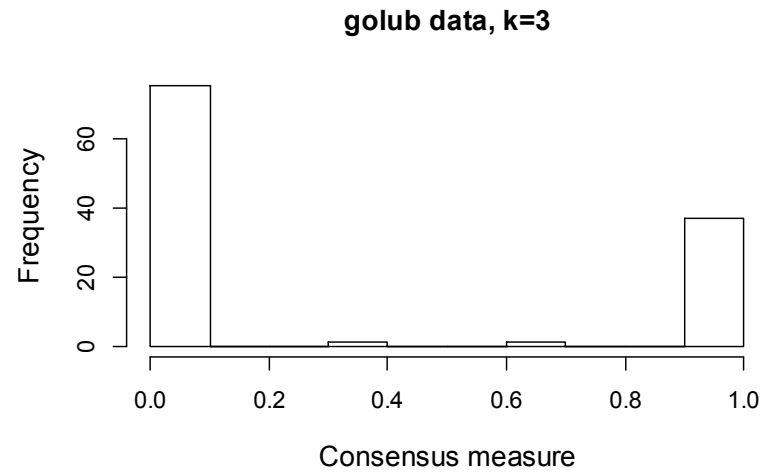


Konsenzus medzi dvomi vzorkami

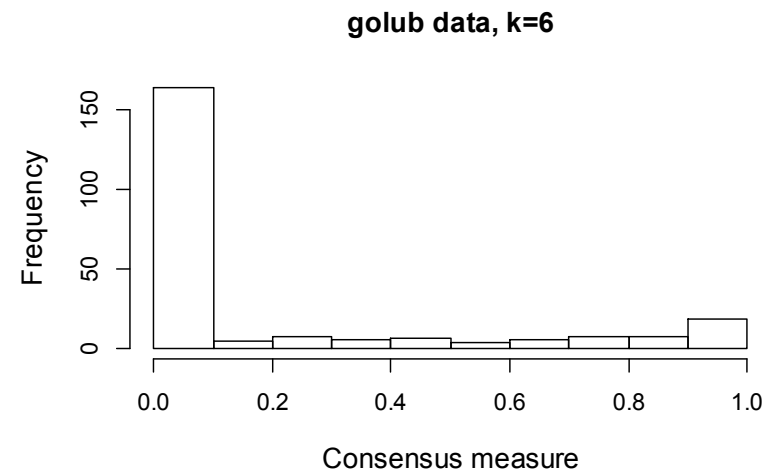
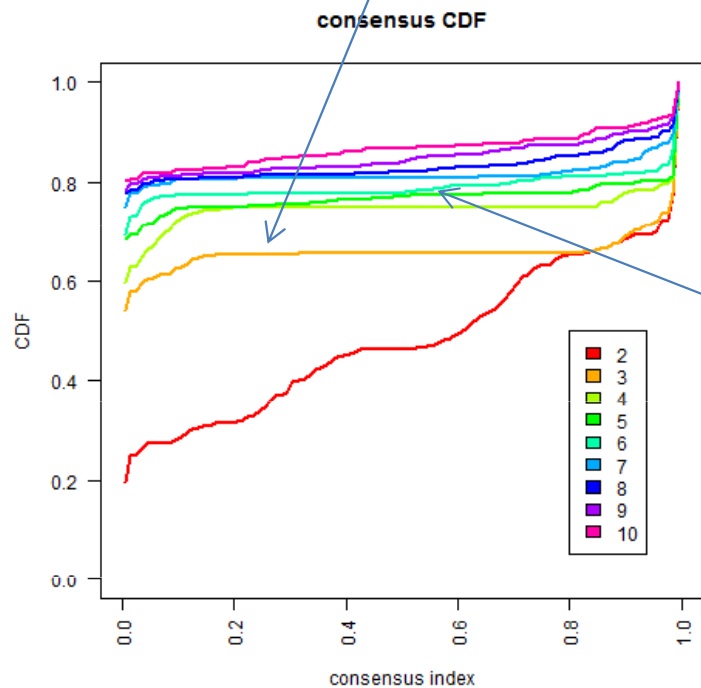
$$CDF^x = \frac{\sum_{i < j} 1\{M_{ij} \leq x\}}{N(N-1)/2}$$

Kumulatívna distribučná funkcia

Odhad počtu zhlukov II

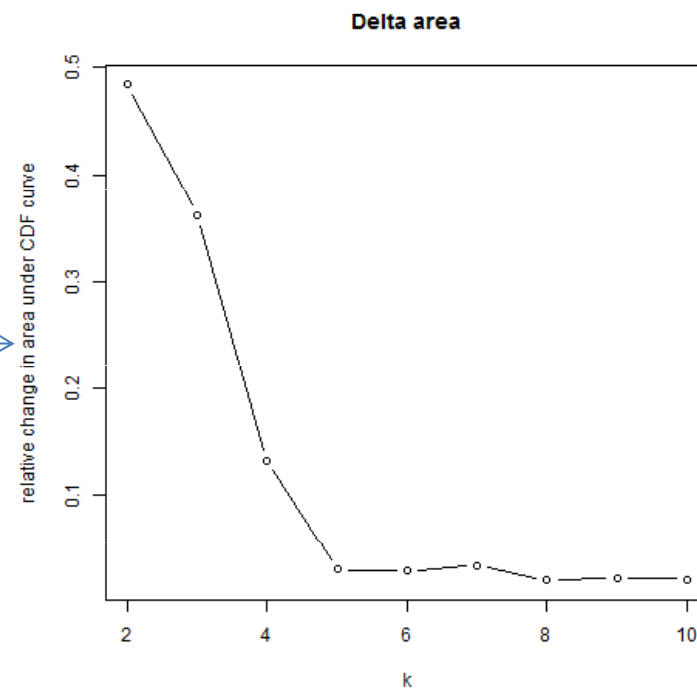
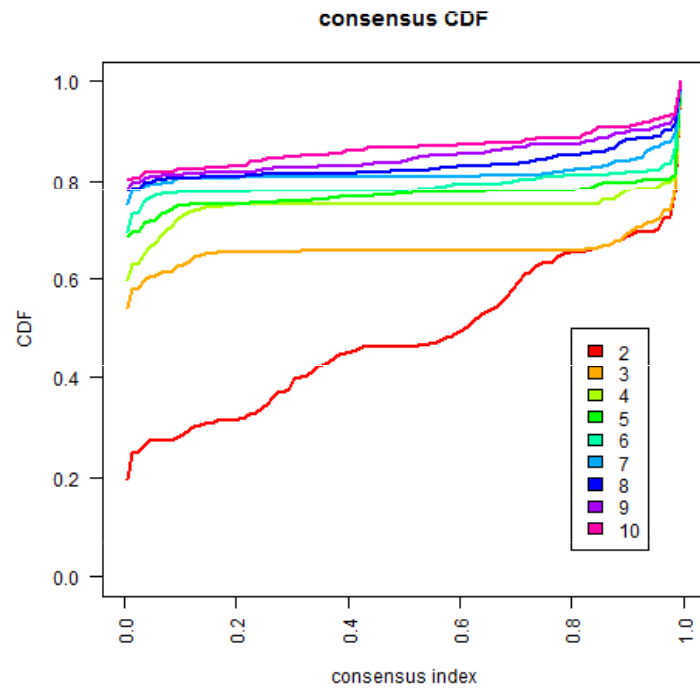


$$CDF^x = \frac{\sum_{i < j} 1\{M_{ij} \leq x\}}{N(N-1)/2}$$

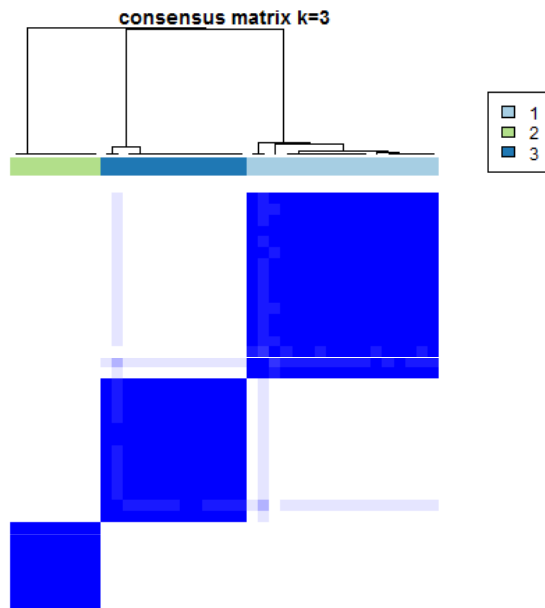


Odhad počtu zhlukov III

Delta = relatívna zmena plochy pod CDF krivkou medzi dvoma k



Odhad stability zhlukov



Kumulatívna distribučná funkcia

$$m^k = \frac{1}{N_l(N_l - 1) / 2} \sum_{\substack{i, j \in I_k \\ i < j}} M_{ij}$$

***Consensus clustering v hierarchickom
zhlukovaní - nevýhoda***

Statické rezanie dendrogramu

Všetky vzorky sú zaradené do zhluku

Dynamic hybrid cut

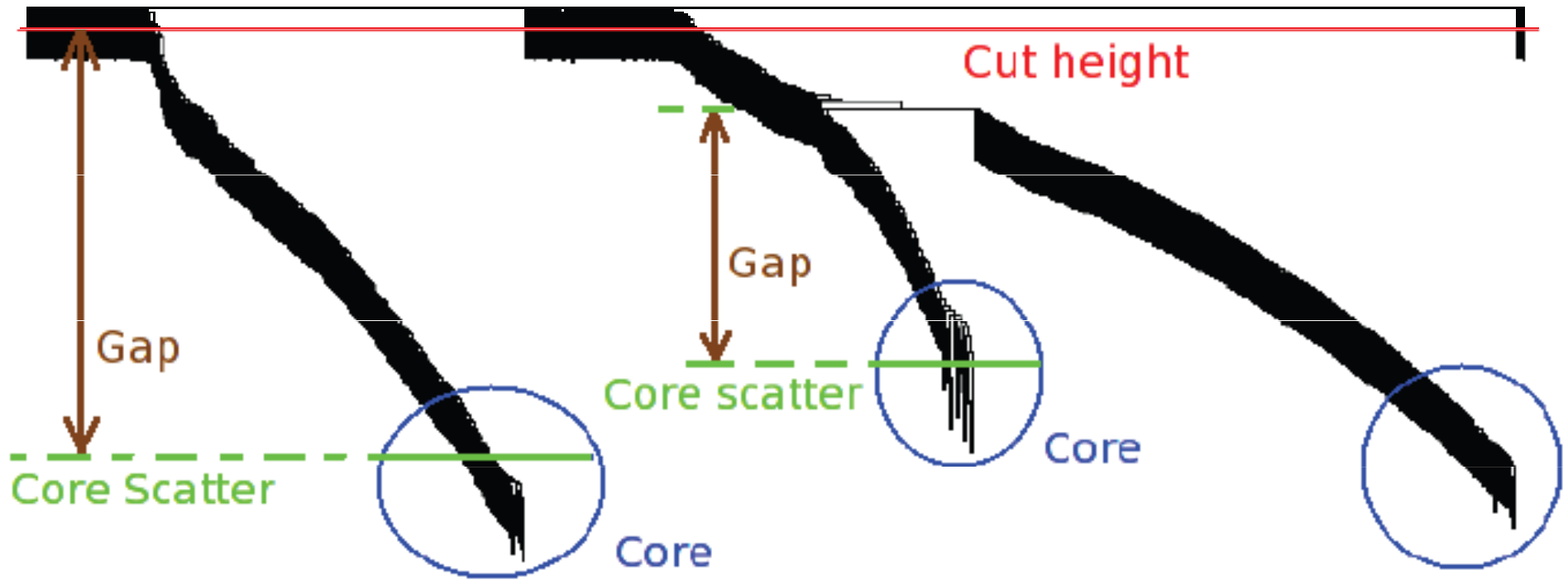
Poskytuje dynamické rezanie dendrogramu na základe minimálnej veľkosti zhlukov, maximálnej výšky rezu a ďalších parametrov.

Ak vzorka nespĺňa kritériá, nie je zaradená do zhluku.

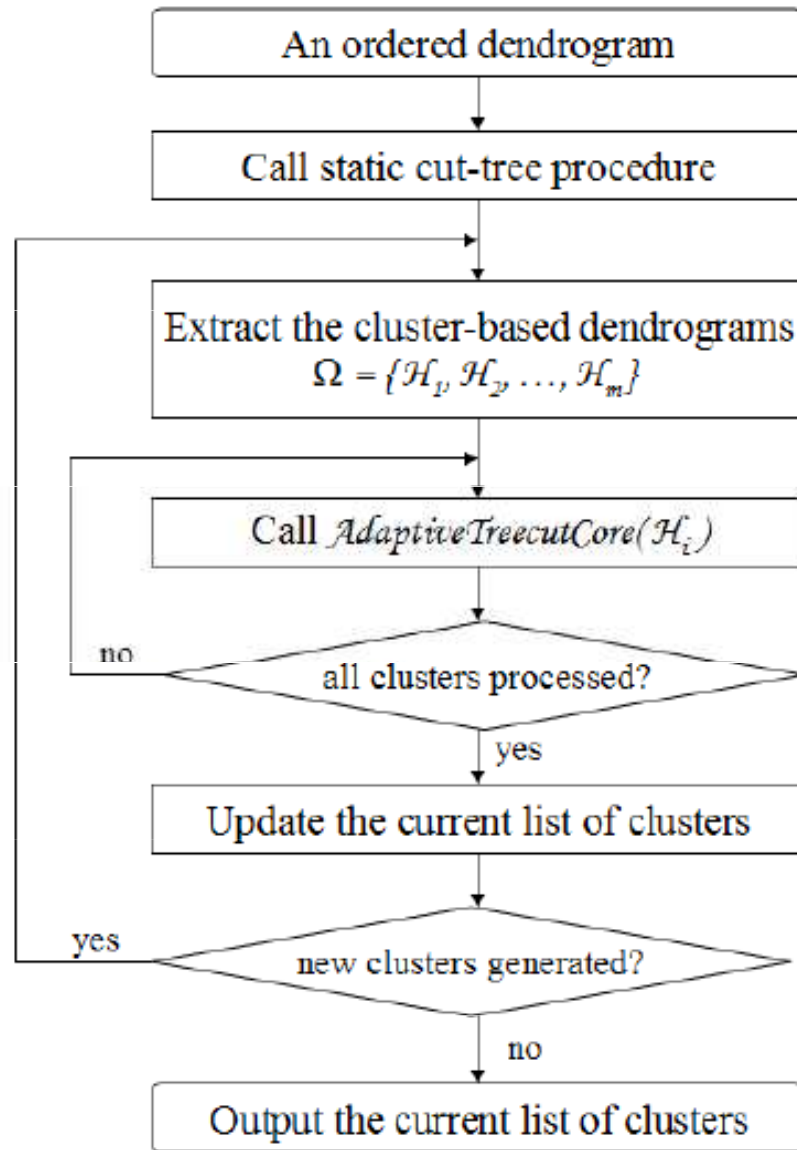
Dynamic hybrid cut - základy

- (1) Zhluk musí mať minimálny počet členov
- (2) Príliš vzdialené objekty sú zo zhuku vylúčené aj v prípade, že patria do rovnakého ramena dendrogramu
- (3) Každý zhluk by mal byť separovaný od okolia medzerou
- (4) Jadro (najnižšie prepojené objekty) každého zhuku musí byť silne prepojené

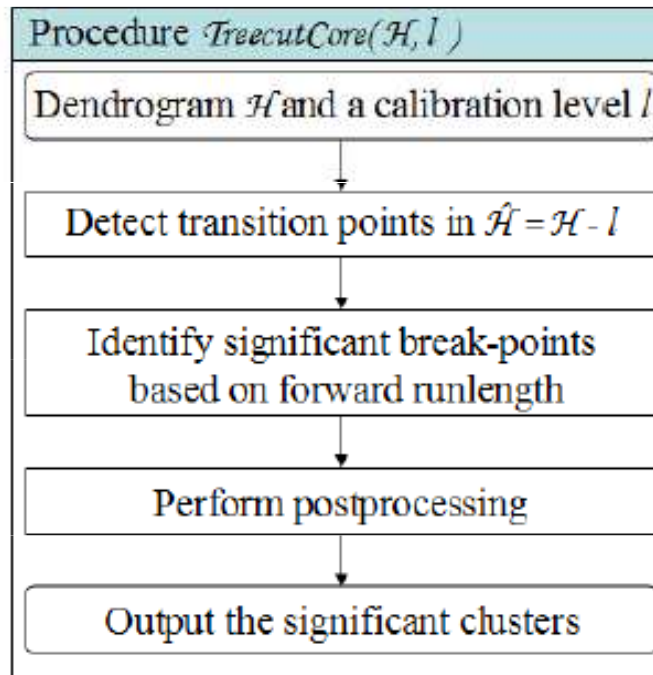
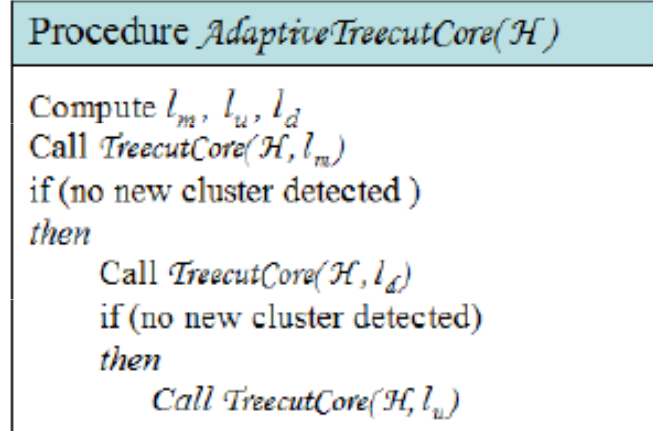
Parametre dynamic hybrid cut



Algorithmus dynamic hybrid cut



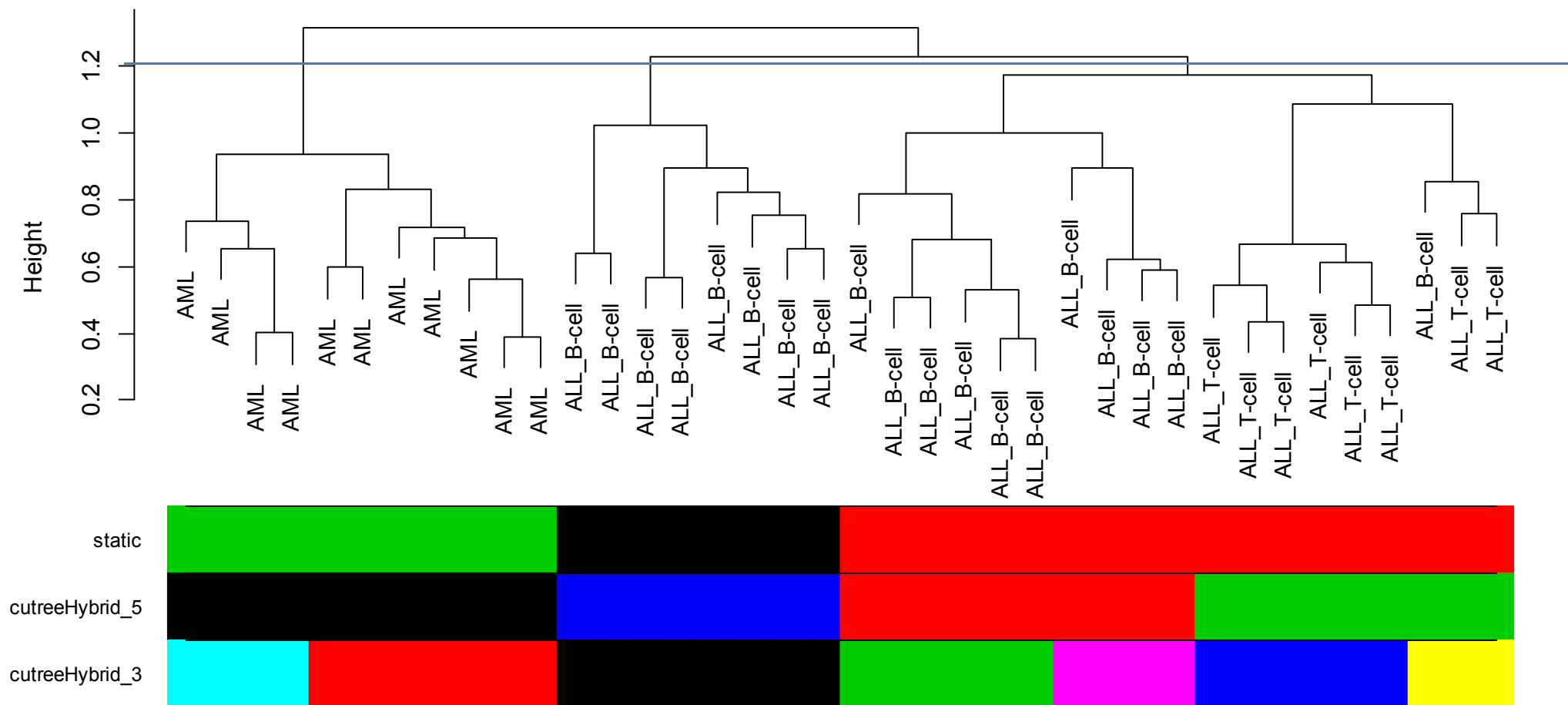
(a)



(b)

Porovnanie statického rezu s dynamickým rezom

Cluster Dendrogram



Cieľ:

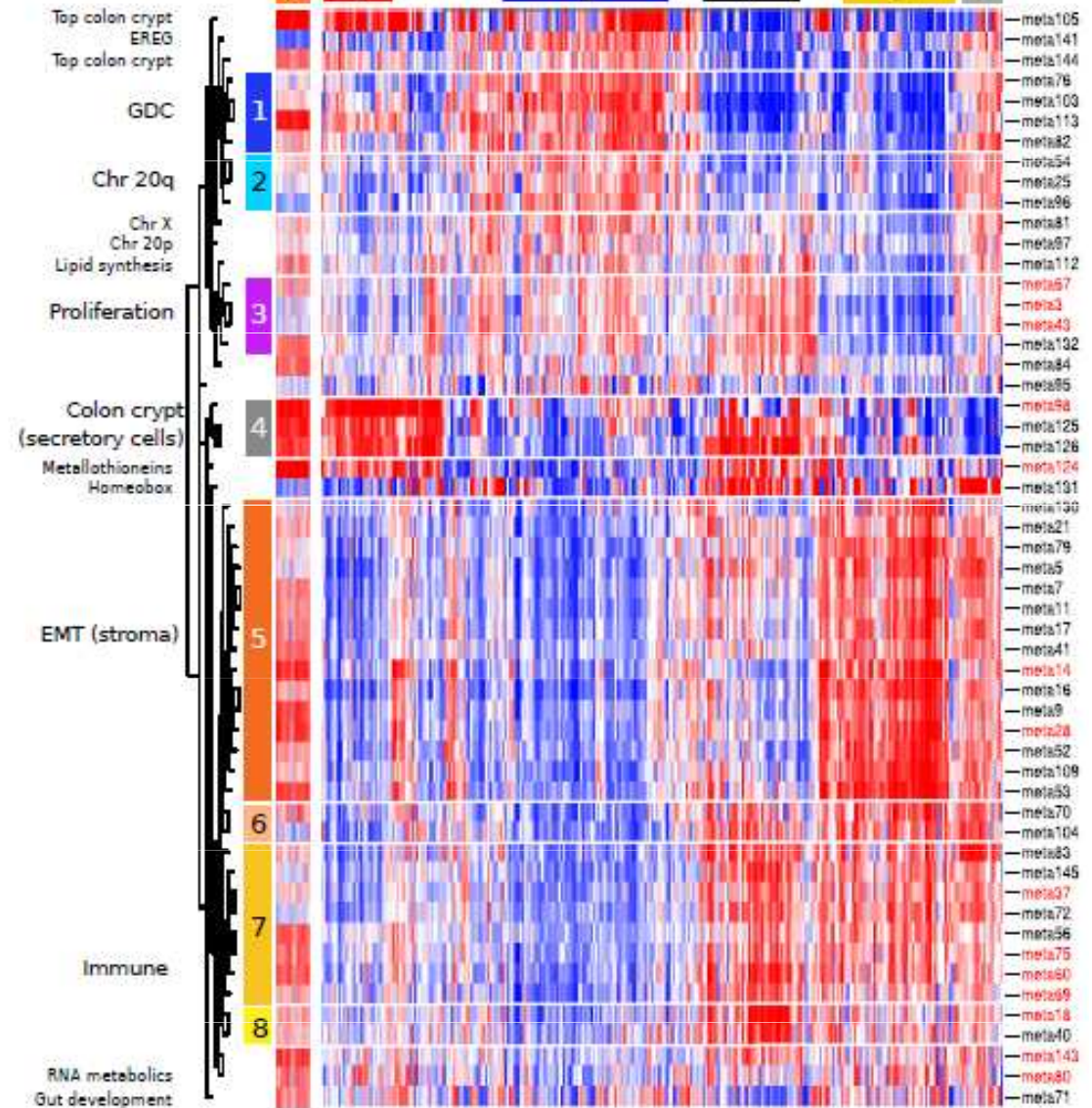
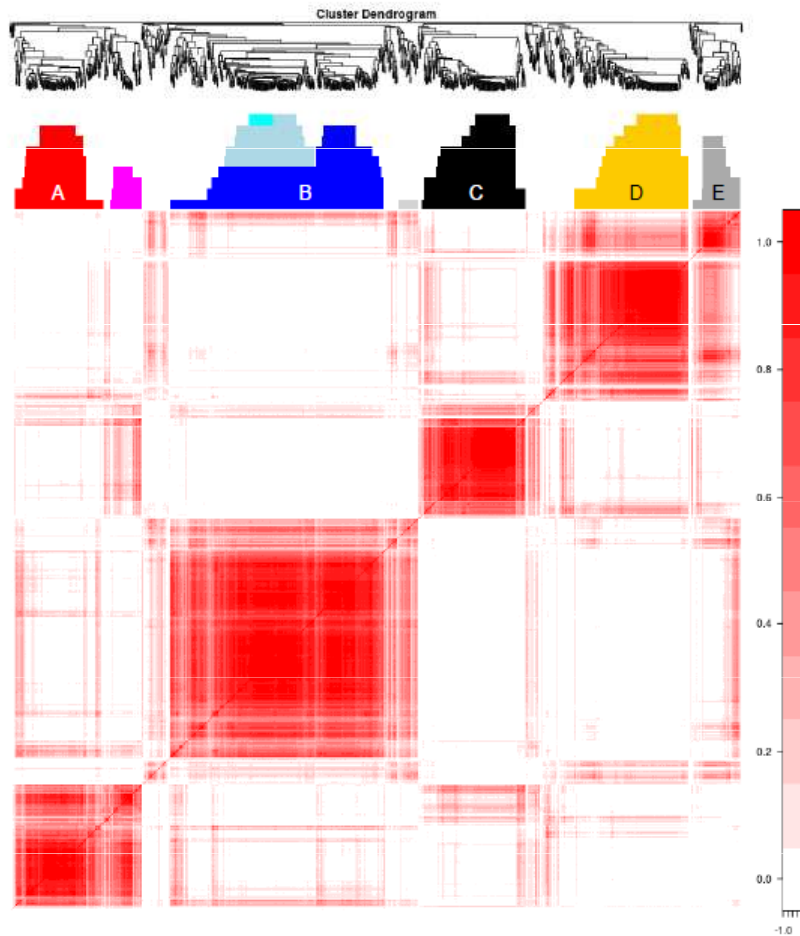
*Nájsť skupiny nádorov kolorekta s podobnou
expresiou génov (podobným génovým profilom)
~ podtypy*

*Charakterizovať tieto podtypy pomocou
klinických a známych molekulárnych
parametrov.*

Motívy génovej expresie v podtypoch

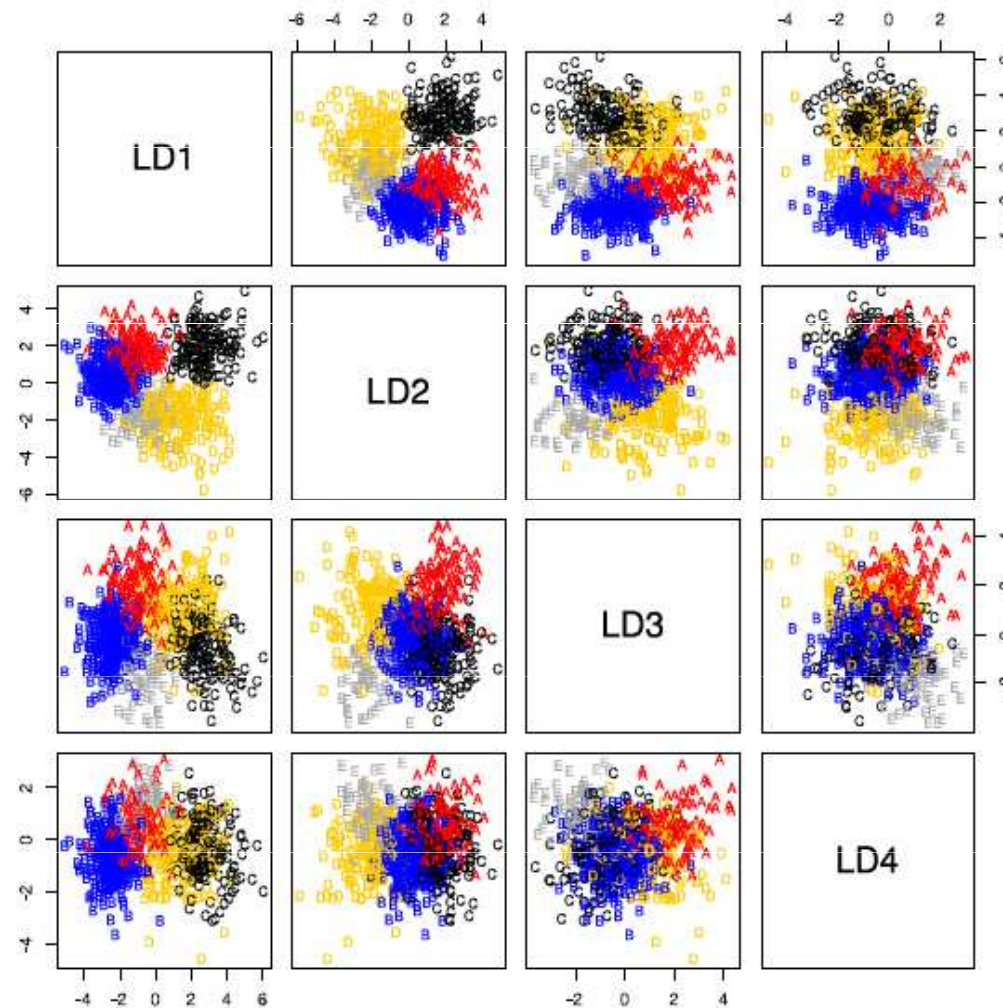
Analýza génových sád

Hierarchické zhlukovanie na matici konsenzusu



Expresné profily pre jednotlivé podtypy

- Klasifikátor LDA (linear discriminant analysis)



Minimálna génová sada

- Selekcia génov pomocou *elastic net* – počíta s korelovanými génmi (neodstraňuje len jeden gén, ale všetky korelované)
- Klasifikátor vytvorený na selektovaných génoch (shrunken centroids)

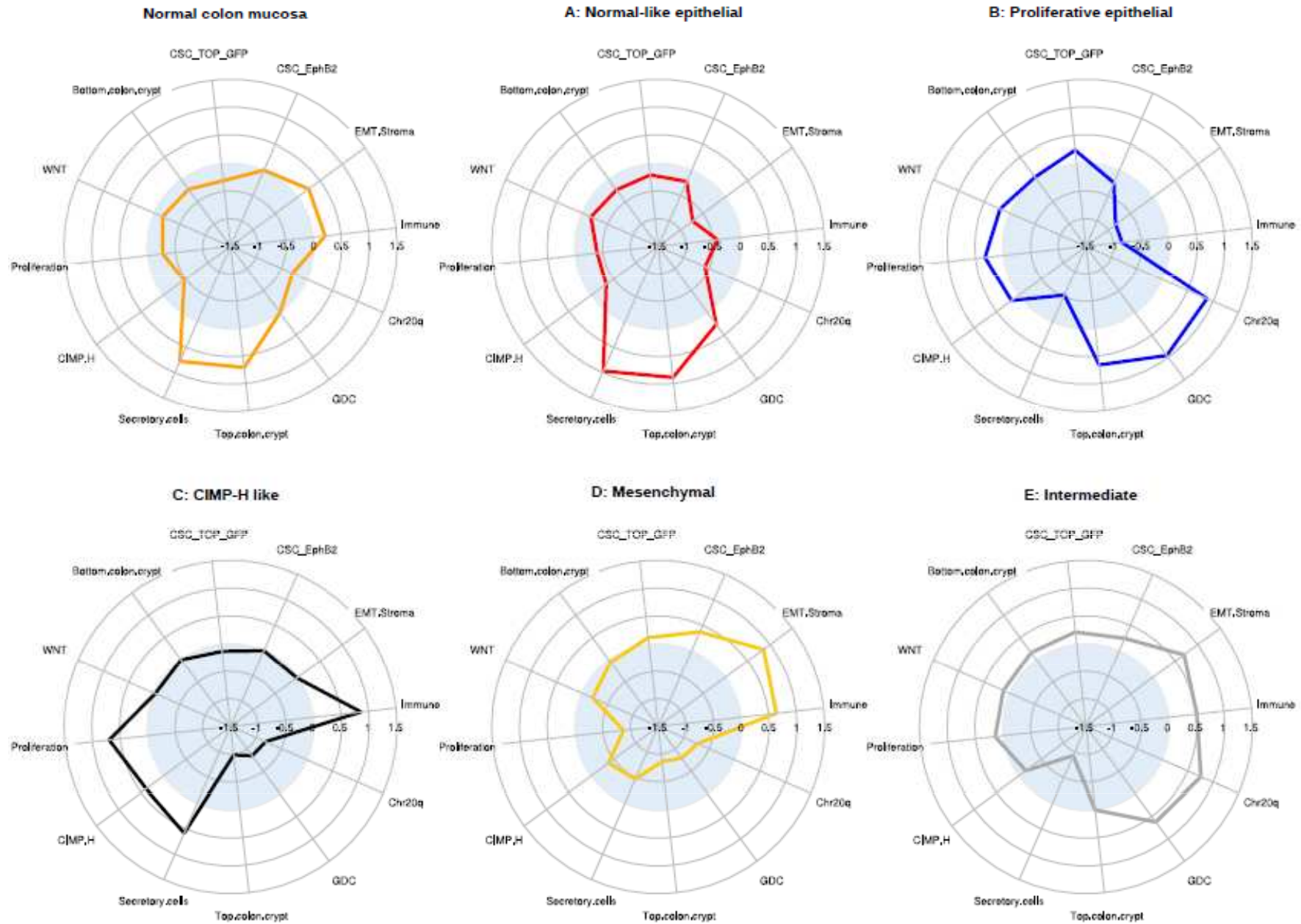
| Subtype | Minimálna génová sada |
|---------|---|
| A | CLCA1, PADI2, ADTRP, RETNLB, TIMP3, MUC2, FNDC1, NR3C2, SULF1, B3GNT7, STYK1, CHI3L1 |
| B | FARP1, ALOX5, FSCN1, HNF4A, RARRES3, MYRIP, GPSM2, TSPAN6, CCDC113, CDHR1, KCTD12, SGK1, BASP1, MT1E, GPX8, RPS6KA3, SOCS3, SLC5A6, PRR15, PLAGL2, IHH, CREB3L1, TP53RK, YAE1D1, EPB41L3, QPRT, KCNK5, RNF43, VAV3, CXCR4, ITPRIP, GRM8, GFPT2, KCNMA1, KIAA0226L, RNASE1 |
| C | TFAP2A, ATP9A, RAB27B, ANP32E, CXCL14, IDO1, RARRES3, EGLN3, KIAA0226L, C10orf99, RPL22L1, PLK2 |
| D | PRICKLE1, RBM47, TAGLN, BOC, HOOK1, C7, ANK2, DCHS1, DDR2, CRYAB, GEM |
| E | REG4, IL6, CXCL5, RAB27B, CEACAM6, PI15, MRPS31, RAP2A, UQCC, AGR3, HSD11B1, IL1B |

Cieľ:

*Nájsť skupiny nádorov kolorekta s podobnou expresiou génov (podobným génovým profilom)
~ podtypy*

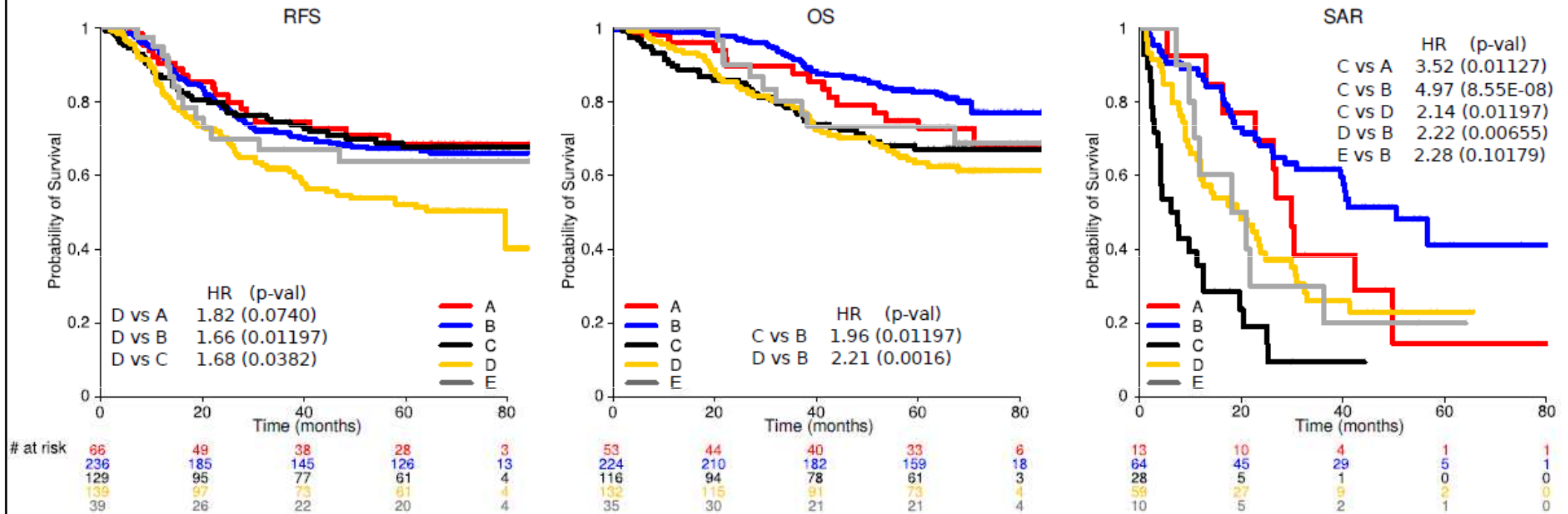
Charakterizovať tieto podtypy pomocou klinických a známych molekulárnych parametrov.

Vzor expresie biologických motívov



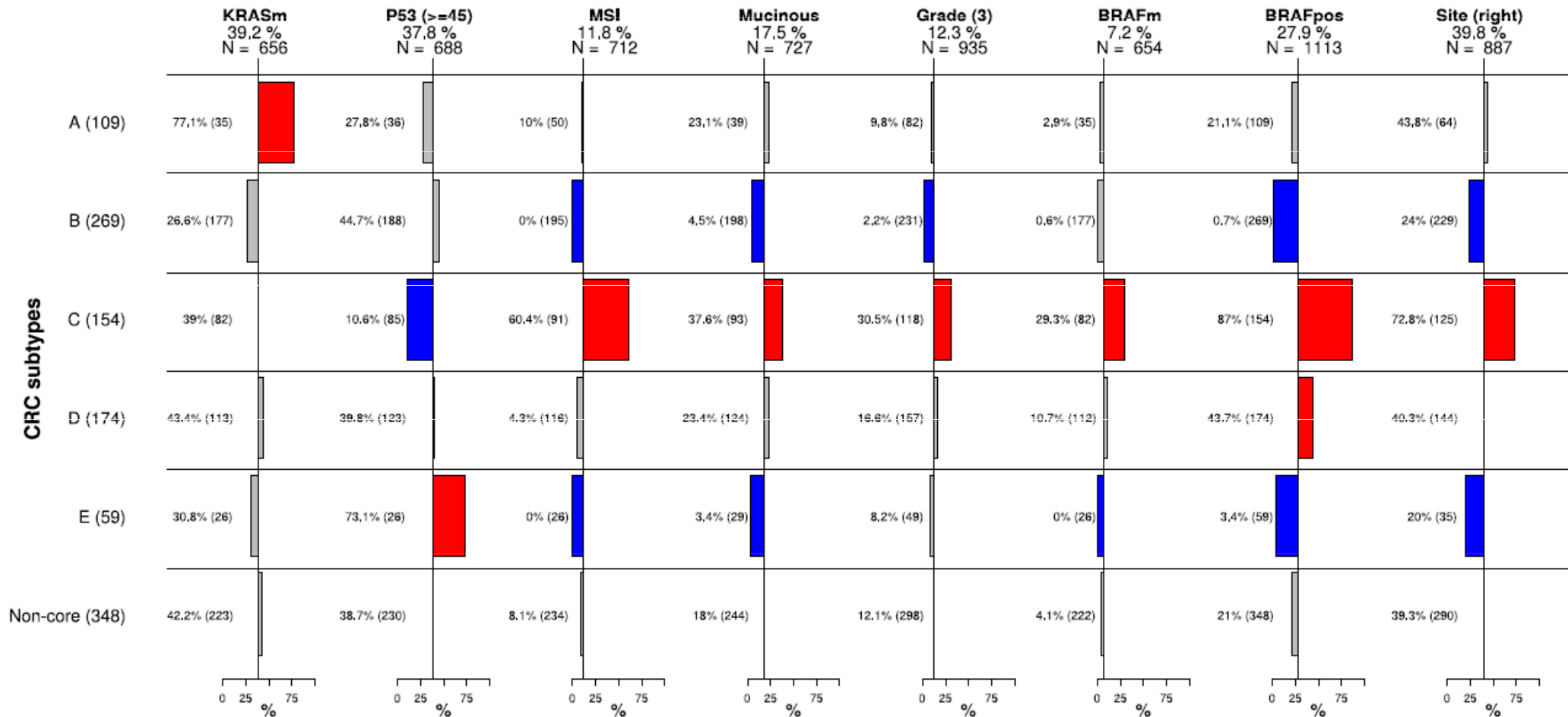
- Analýza génových sád pomocou KS testu

Rozdiely v prežití



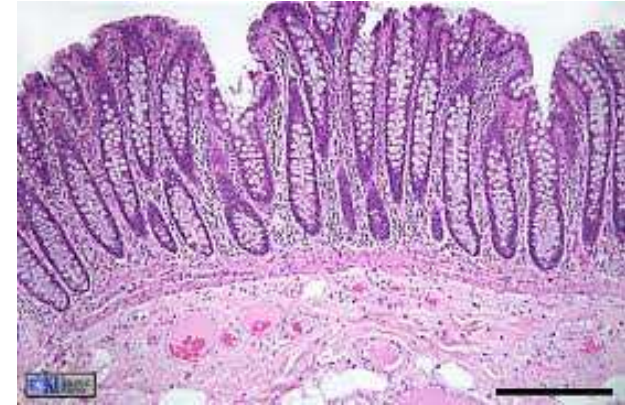
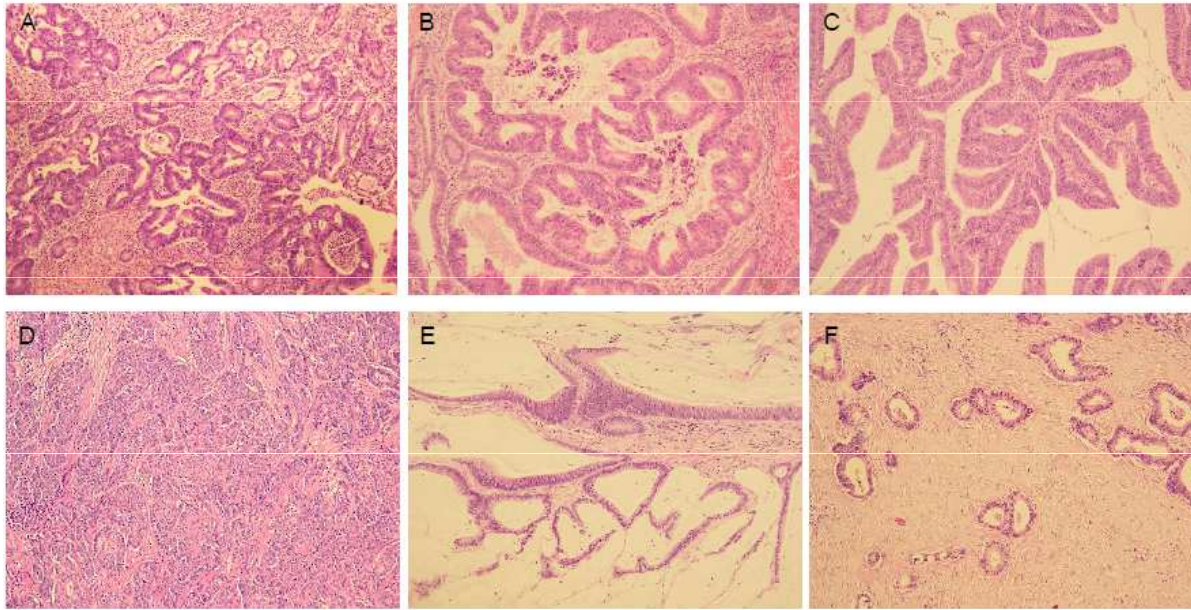
- Kaplan-Meierove krivky prežitia
- Coxov model proporcionálnych rizík (efekt štádia vs podtypov)

Charakterizácia podtypov klinickými a molekulárnymi premennými

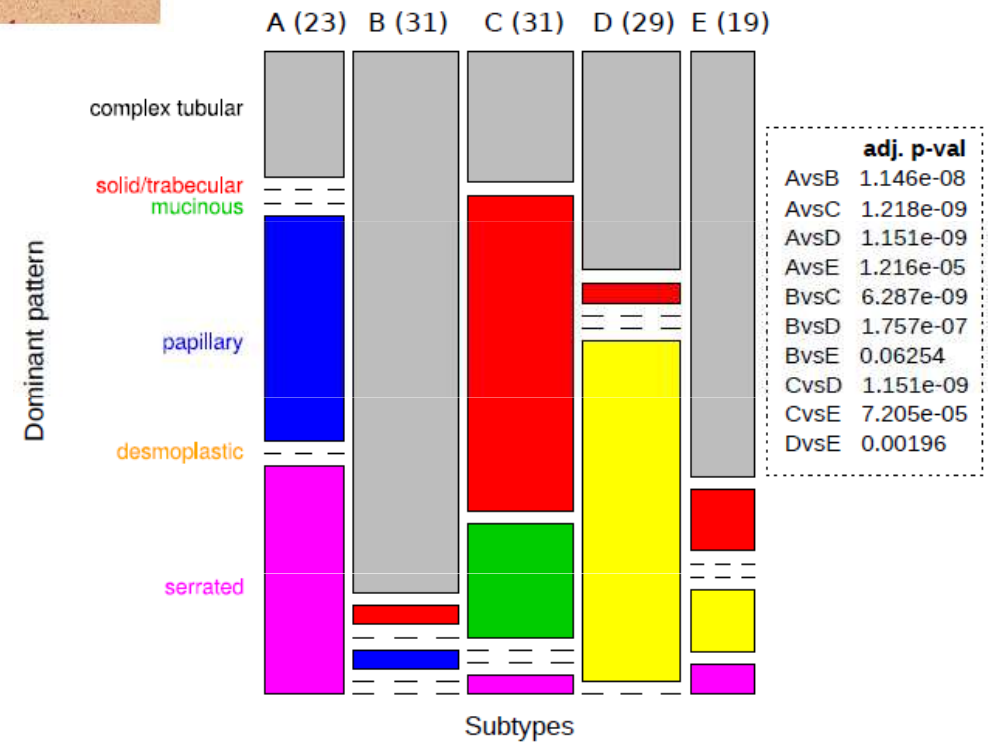


- Rozdiel od populačnej baseline u každého podtypu pomocou Fisherovho exaktného testu, FDR úprava p-hodnôt

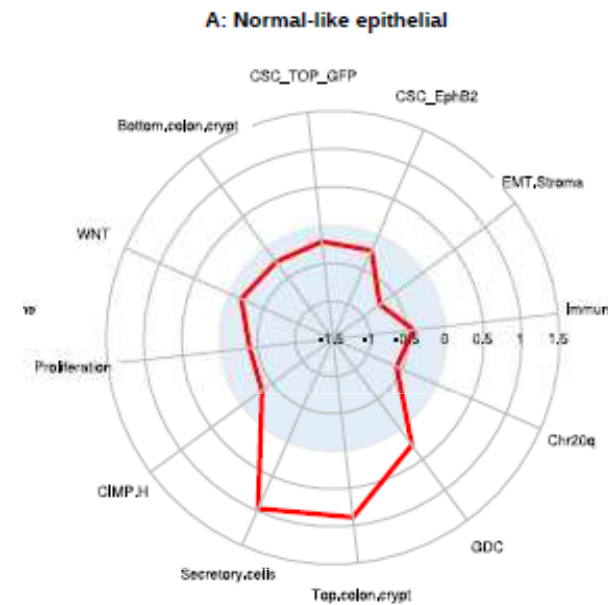
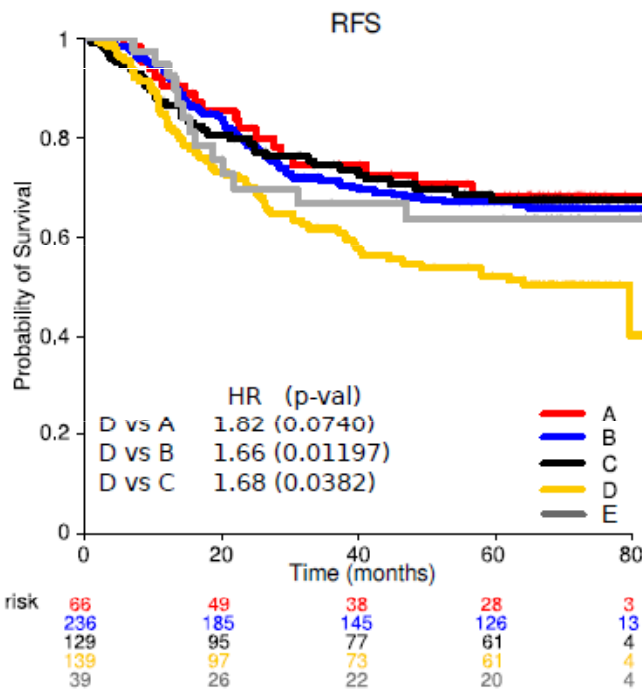
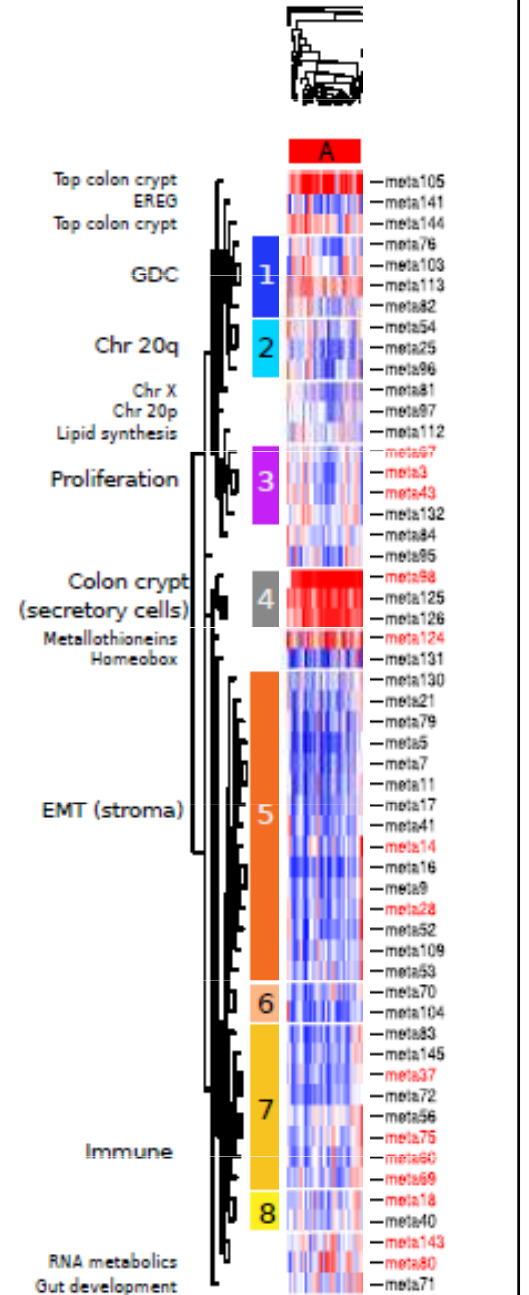
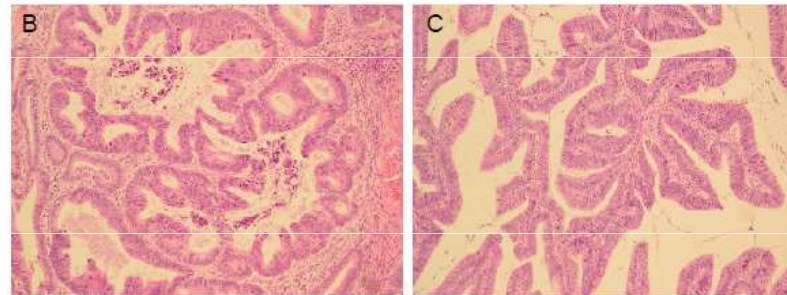
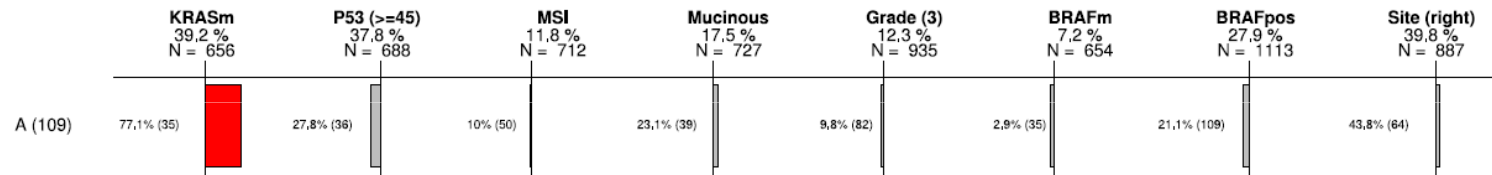
Histologické rozdiely



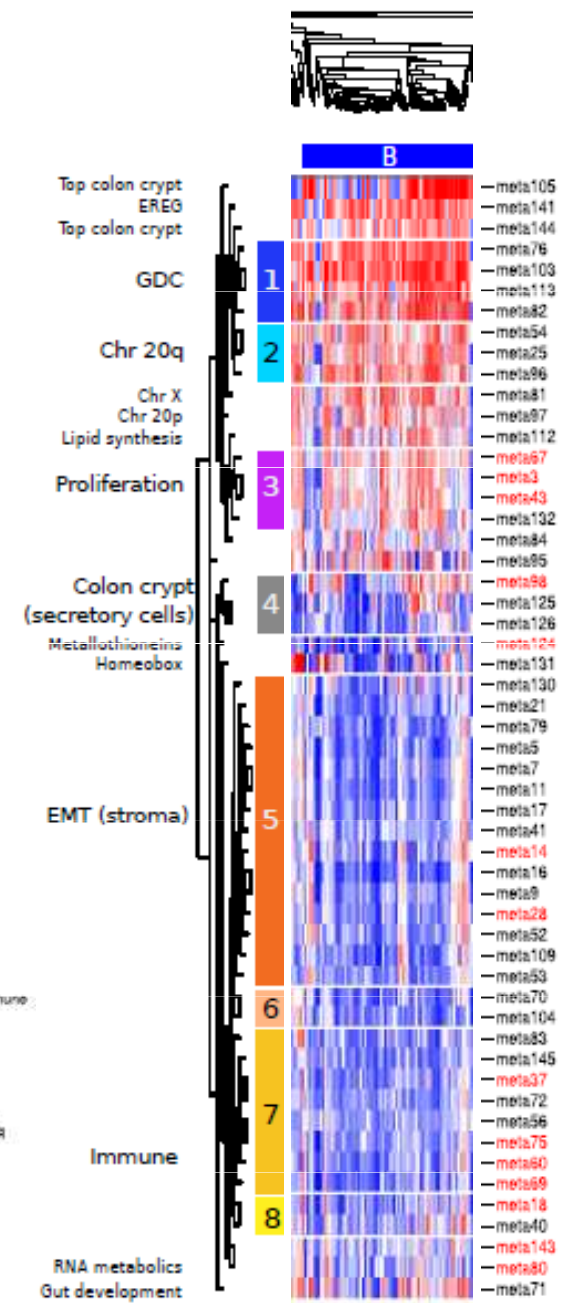
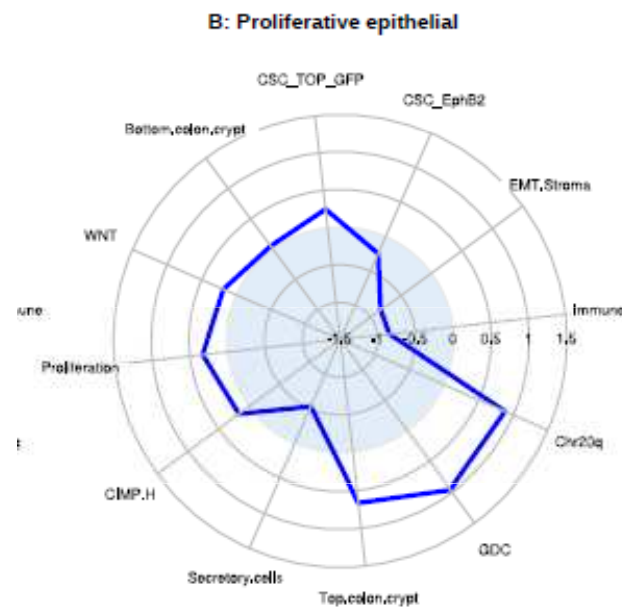
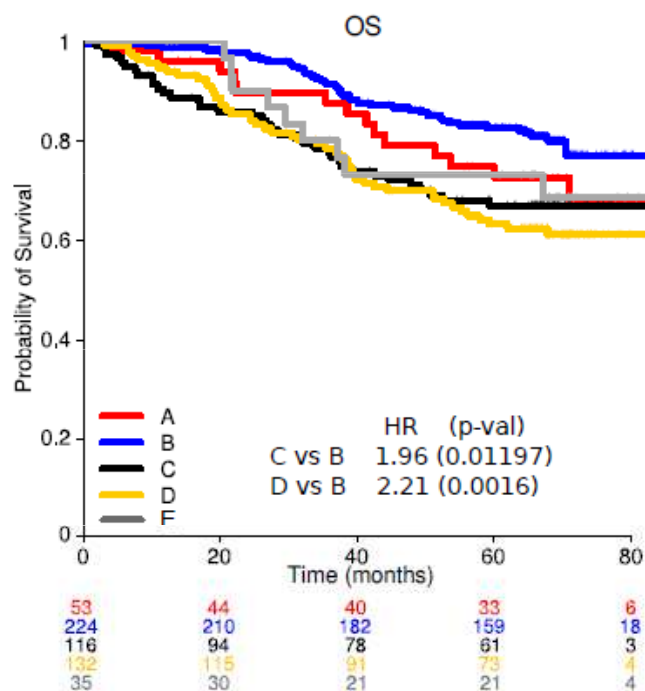
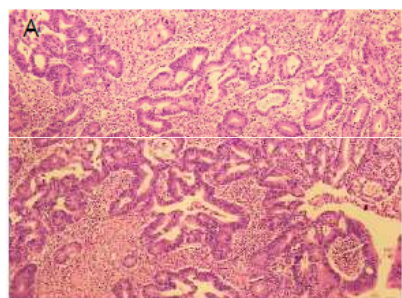
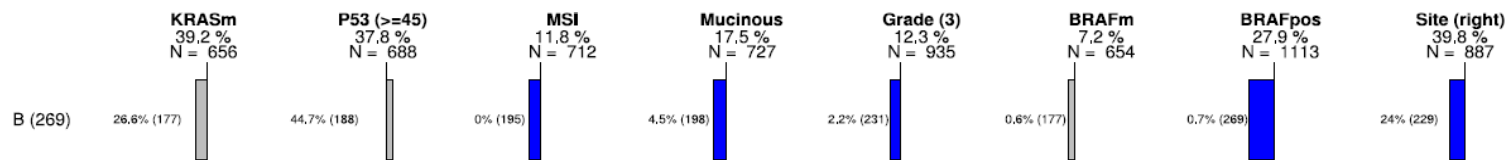
- Fisherov test, adjustácia na FDR



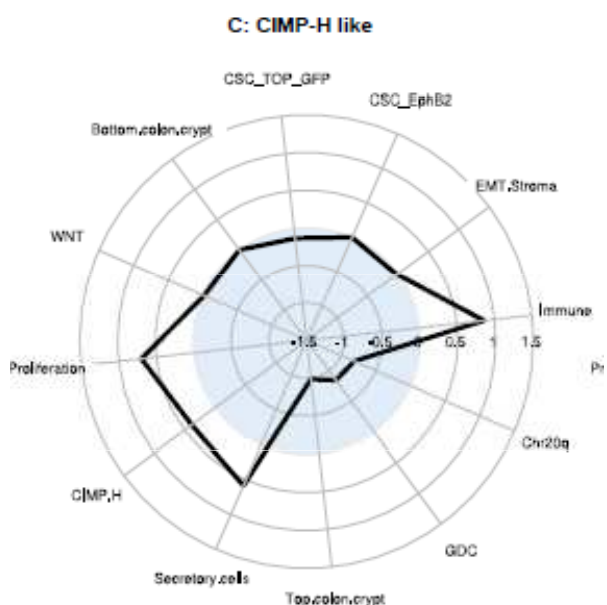
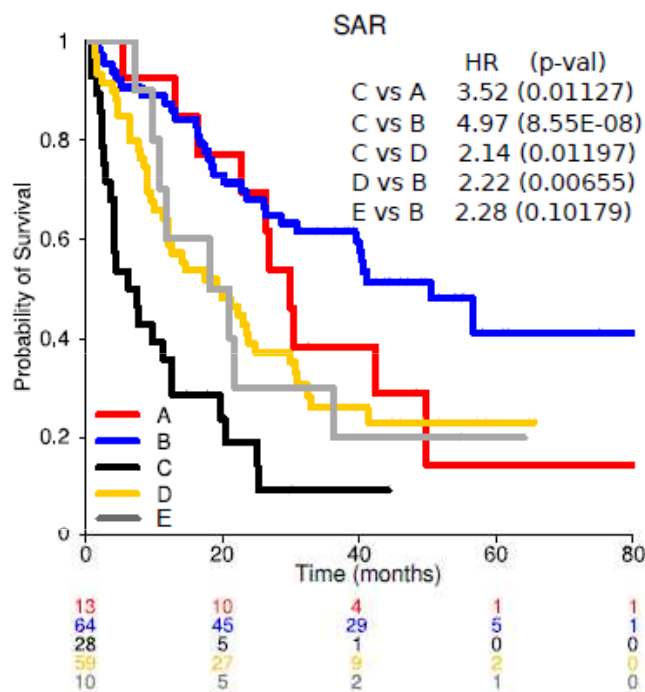
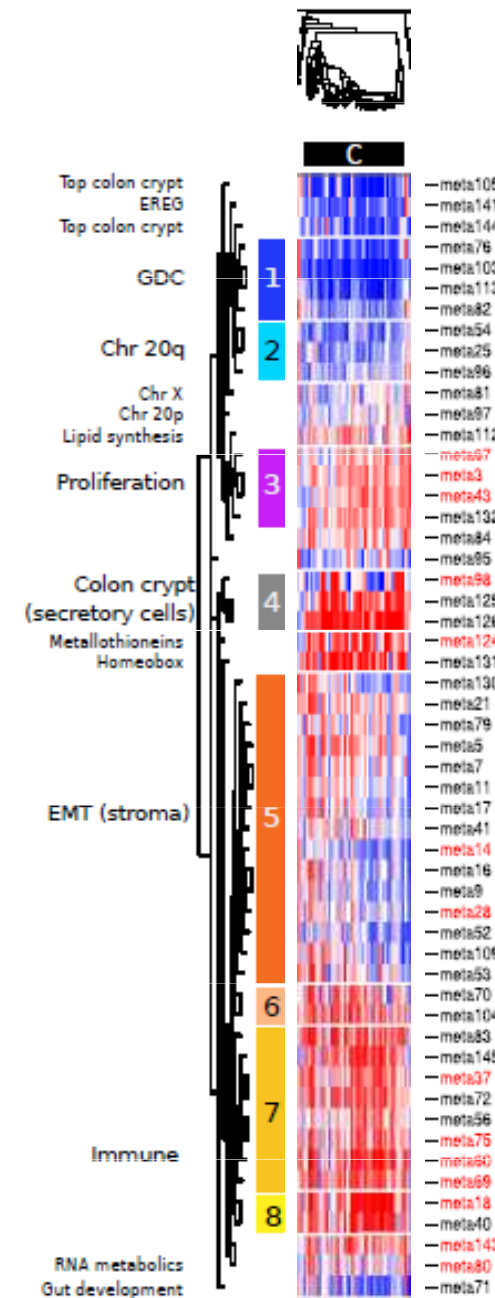
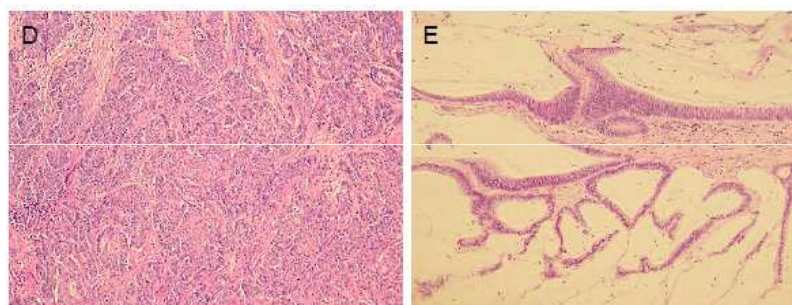
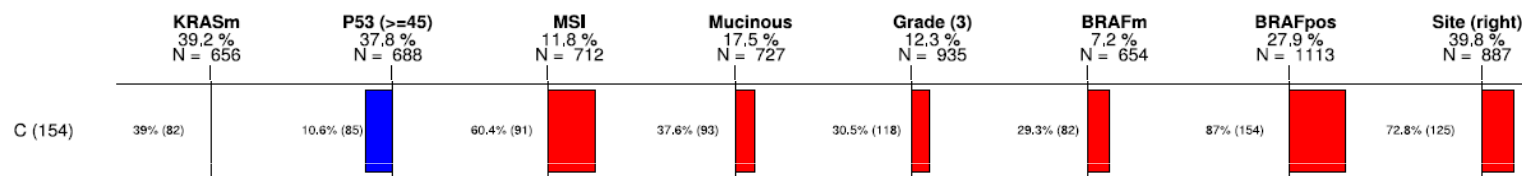
Podtyp A – Surface crypt like - KRAS mutanti, papillary a serrated morfotyp, najviac diferencovaný, bez aktívnej Wnt signálnej dráhy. Dobrý OS a RFS.



Podtyp B – Lower crypt like – diferencované ale bez sekrečných buniek, proliferujúce, a aktívnou Wnt signálnou dráhou. Komplexný tubulárny morfortyp. Časo MSS, BRAFwt, nižšieho grádu, dobré prežitie v OS, RFS i SAR.



Podtyp C – CIMP-H like - časo MSI, *BRAF*-mutantné, hypermutované, z pravej časti hrubého čreva. Histologicky - horšie diferencované, solídno-trabekulárne s mucínovým morfortypom. Aktívne proliferujú a majú silnú imunitnú reakciu. Dobrý RFS, ale zlý OS and SAR.



| |
|----|
| 13 |
| 64 |
| 28 |
| 59 |
| 10 |

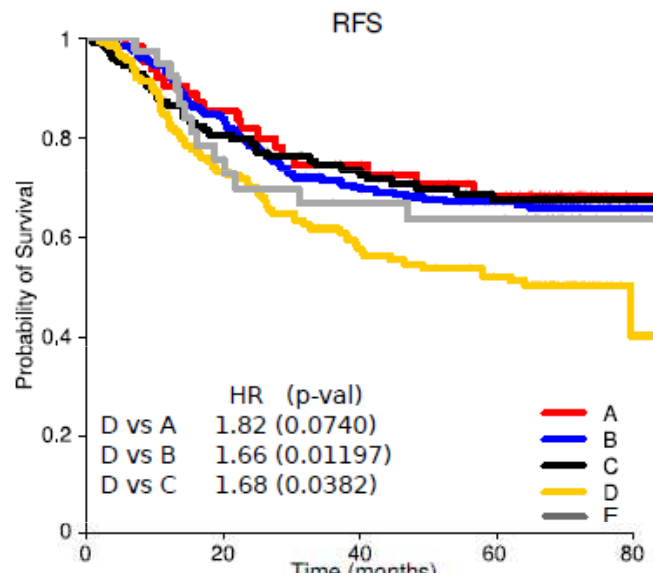
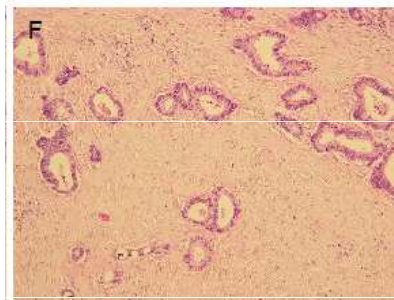
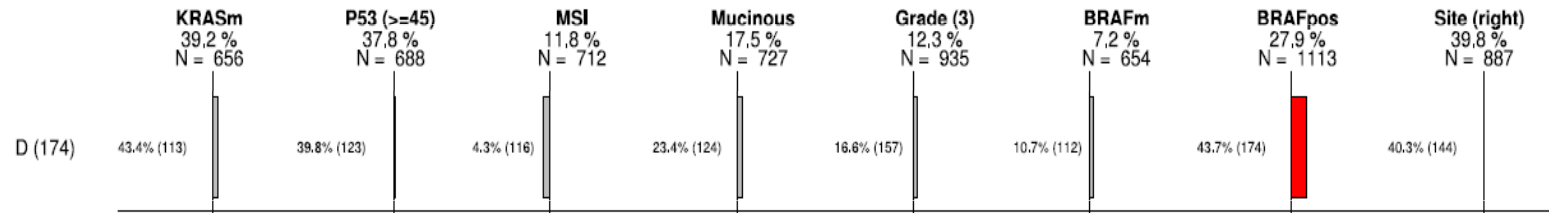
| |
|----|
| 10 |
| 45 |
| 5 |
| 27 |
| 5 |

| |
|----|
| 4 |
| 29 |
| 1 |
| 9 |
| 2 |

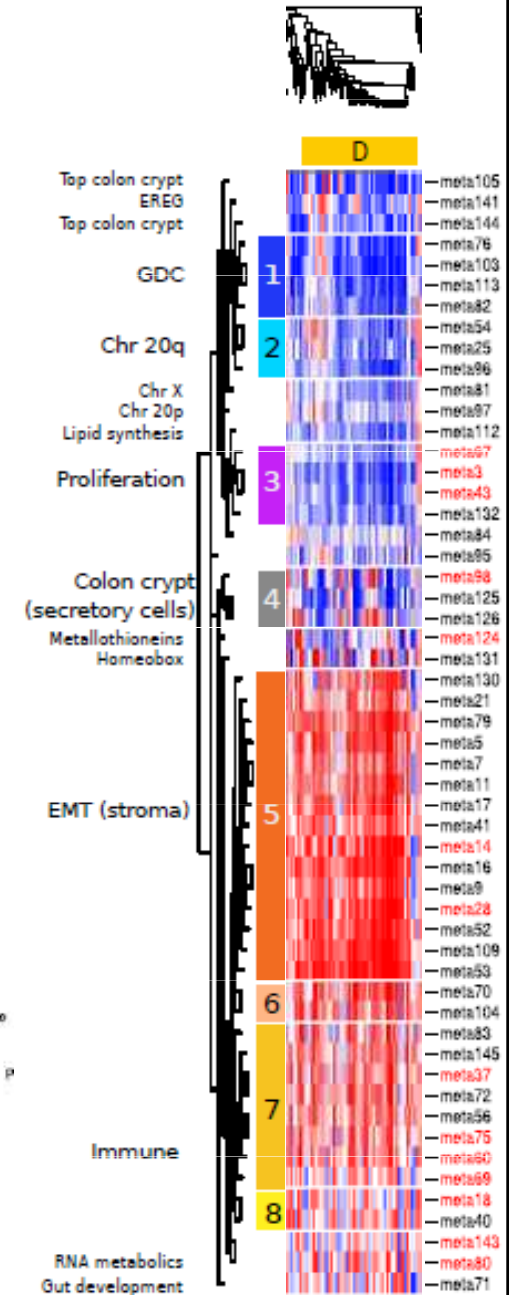
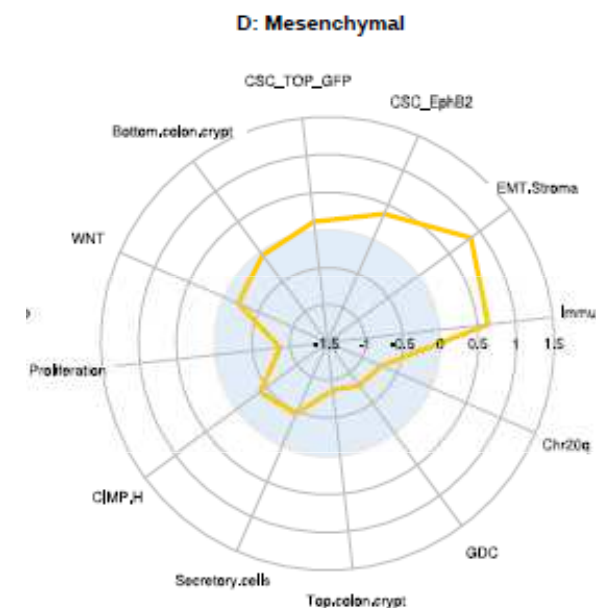
| |
|---|
| 1 |
| 5 |
| 0 |
| 2 |
| 1 |

| |
|---|
| 1 |
| 1 |
| 0 |
| 0 |
| 0 |

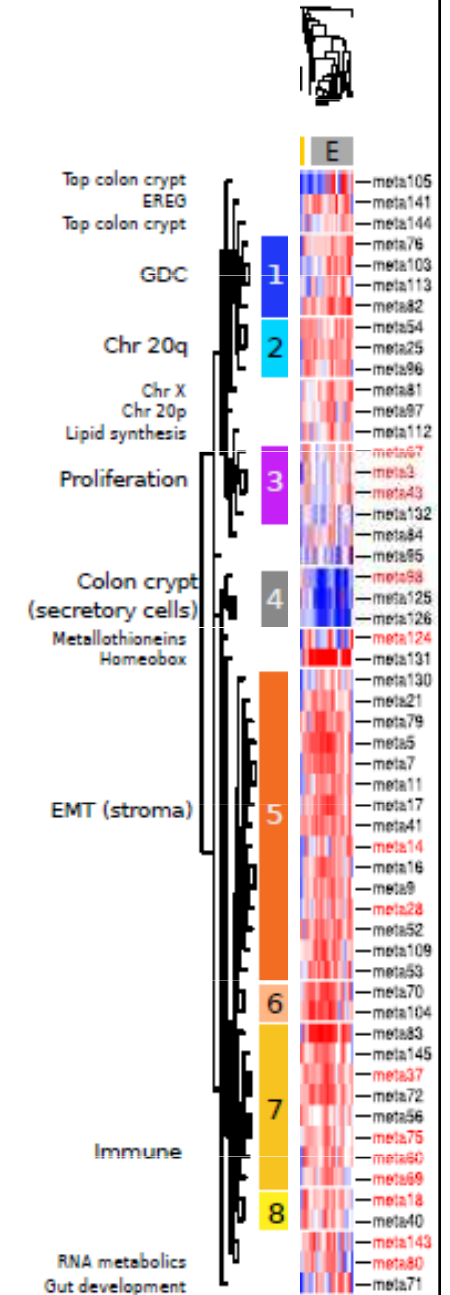
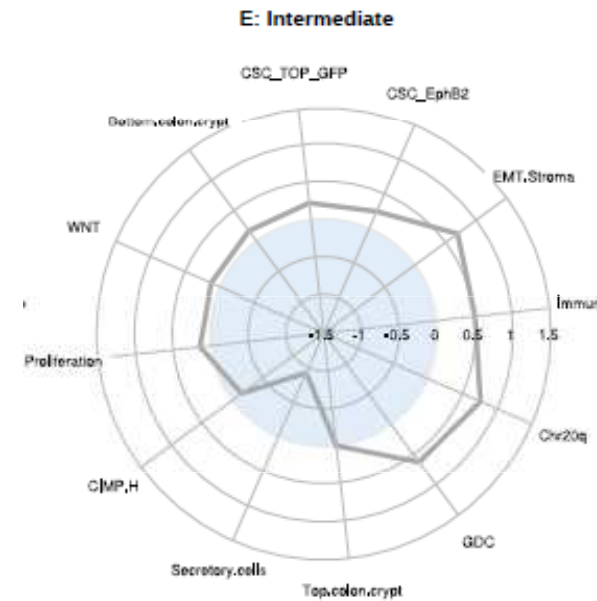
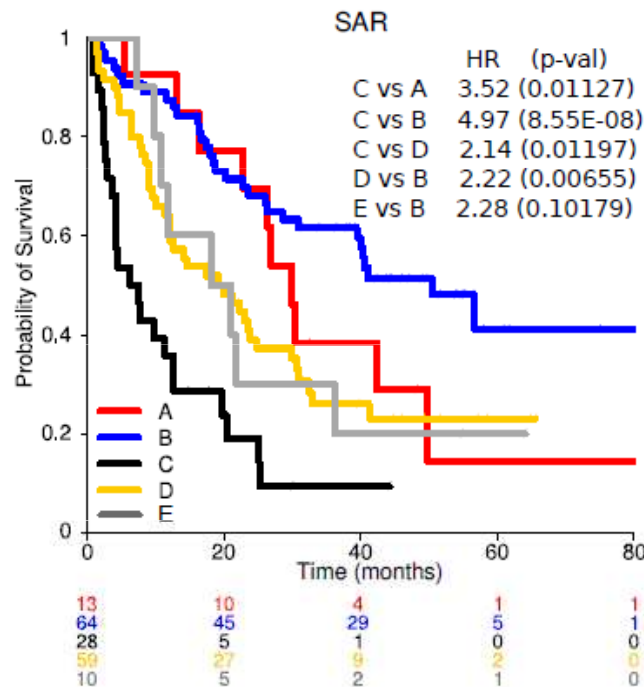
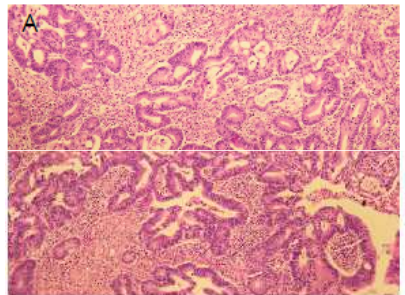
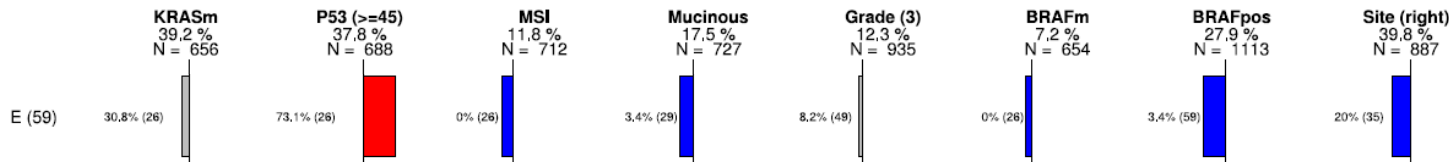
Podtyp D – Mesenchymal – markery kmeňových buniek, veľa mezenchymálnych buniek, ktoré sa prejavujú expresiou EMT génov. Wnt signálna dráha je neaktívna a proliferácia nízka. Klinické a mutačné charakteristiky sa nelíšia od populačnej baseline. Majú najkratšie prežitie do relapsu, zlý OS a SAR.



| # at risk | 0 | 20 | 40 | 60 | 80 |
|-----------|-----|-----|-----|-----|----|
| A | 66 | 49 | 38 | 28 | 3 |
| B | 236 | 185 | 145 | 126 | 13 |
| C | 129 | 95 | 77 | 61 | 4 |
| D | 139 | 97 | 73 | 61 | 4 |
| F | 39 | 26 | 22 | 20 | 4 |



Podtyp E – Mixed – často MSS, *BRAF*wt, z ľavej strany hrubého čreva. Podobne ako podtyp D exprimuje gény kmeňových buniek a EMT procesu, avšak podobne s B má vysokú aktiviu kanonickej Wnt dráhy a vyzerá viac diferenciovaný. Je podobný B – komplexný tubulárny, len častejšie obsahuje mutáciu *p53*.



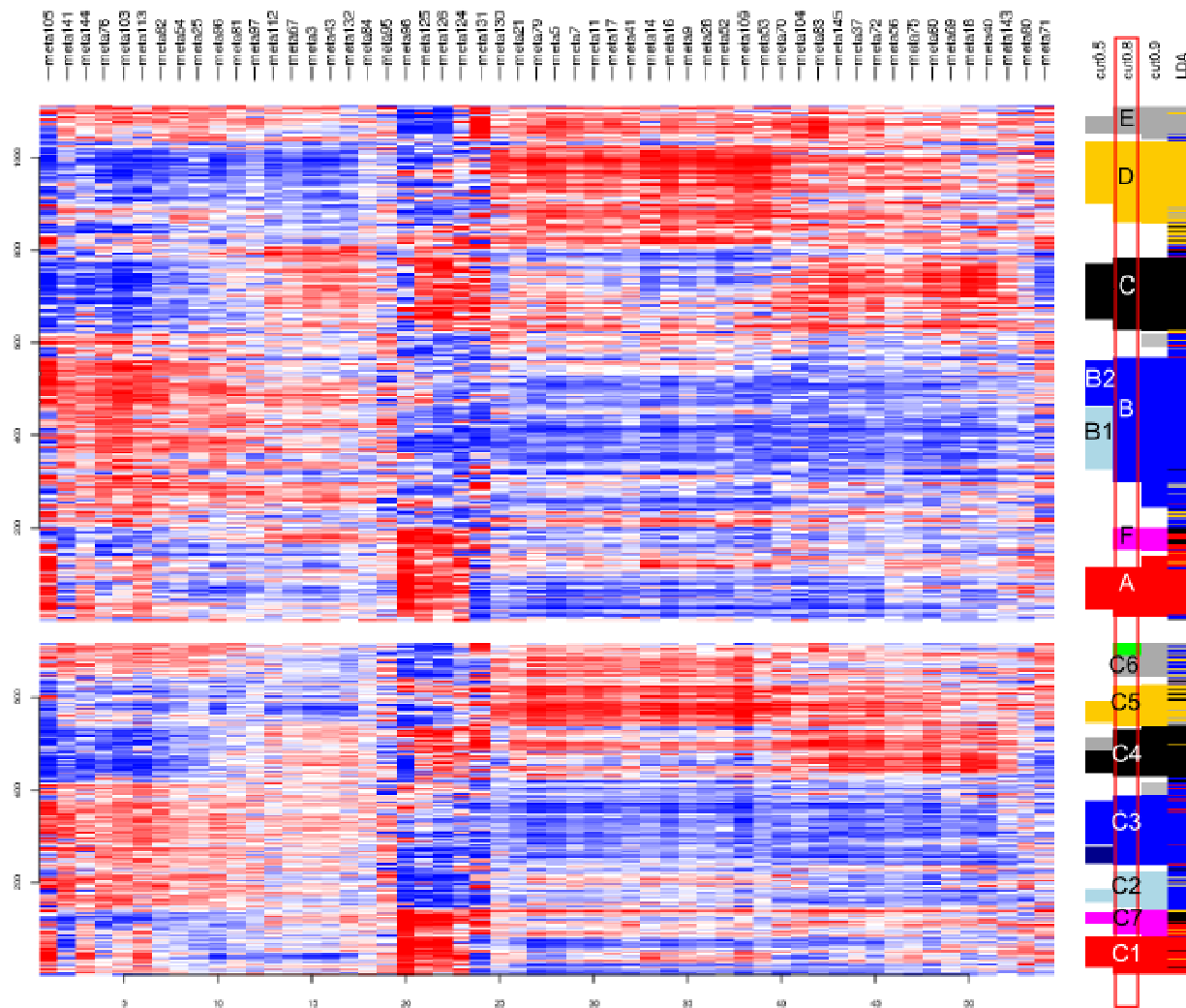
VALIDÁCIA V ZHLUKOVANÍ MOLEKULÁRNYCH DÁT

Validácia algoritmu a parametrov modelu na testovacom súbore

Keď zopakujem celú procedúru na inom súbore, dostanem podobné skupiny?

training set

validation set



| Cluster/subtype in validation set | LDA assignment | | | | | SUM |
|-----------------------------------|----------------|-----|-----|-----|----|-----|
| | A | B | C | D | E | |
| C1 / A | 74 | 4 | 3 | 3 | 0 | 84 |
| C2 / B1 | 1 | 58 | 0 | 2 | 13 | 74 |
| C3 / B2 | 12 | 134 | 1 | 0 | 1 | 148 |
| C4 / C | 1 | 2 | 99 | 4 | 0 | 106 |
| C5 / D | 0 | 3 | 12 | 64 | 7 | 86 |
| C6 / E | 1 | 17 | 0 | 17 | 13 | 48 |
| C7 / F | 23 | 1 | 22 | 9 | 1 | 56 |
| Non-core | 21 | 53 | 18 | 8 | 18 | 118 |
| SUM | 133 | 272 | 155 | 107 | 53 | 720 |

| Cluster/subtype in validation set | Subtypes from training set most correlated to validation subtypes | | | | | |
|-----------------------------------|---|------|-----------|----------------|------|-----------|
| | First subtype | | | Second subtype | | |
| | Subtype | Cor | P-val | Subtype | Cor | P-val |
| C1 / A | A | 0.85 | p<1.0E-15 | F | 0.41 | p<1.0E-15 |
| C2 / B1 | B | 0.71 | p<1.0E-15 | E | 0.47 | p<1.0E-15 |
| C3 / B2 | B | 0.91 | p<1.0E-15 | A | 0.36 | p<1.0E-15 |
| C4 / C | C | 0.89 | p<1.0E-15 | F | 0.29 | p<1.0E-15 |
| C5 / D | D | 0.93 | p<1.0E-15 | E | 0.37 | p<1.0E-15 |
| C6 / E | E | 0.63 | p<1.0E-15 | D | 0.58 | p<1.0E-15 |
| C7 / F | F | 0.61 | p<1.0E-15 | C | 0.55 | p<1.0E-15 |

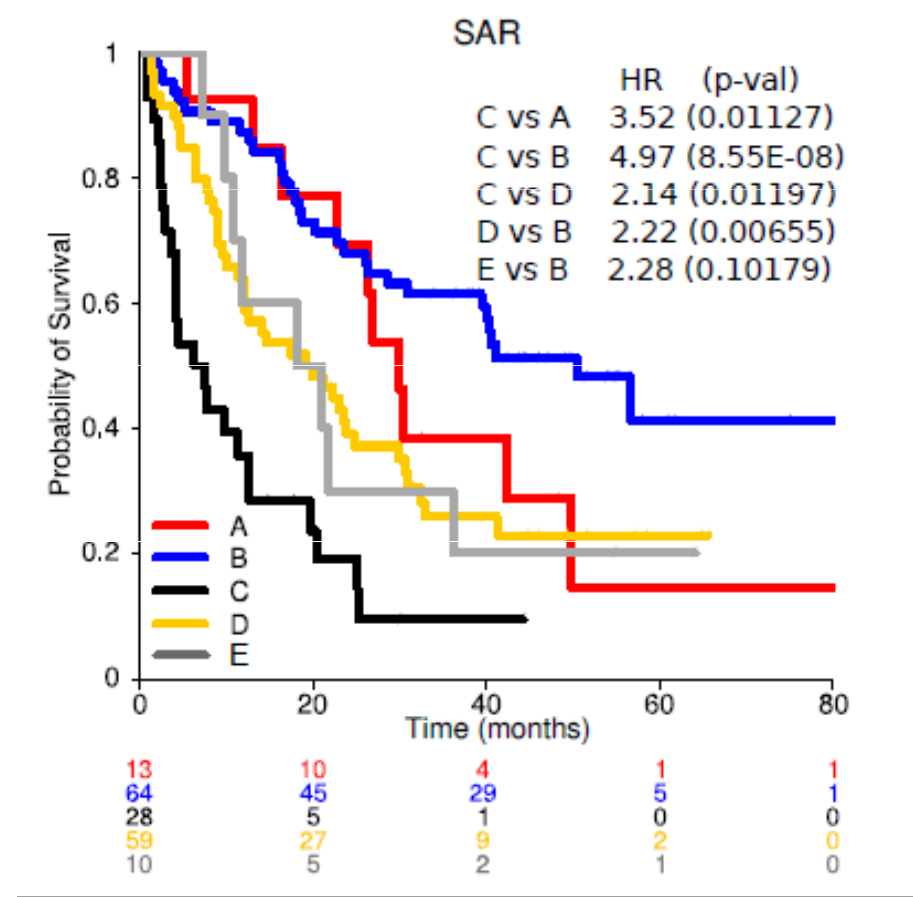
Validácia konceptu pomocou klinických,
molekulárnych a histologických charakteristík
objavených skupín

Majú objavené skupiny biologickú podstatu /
odrážajú známe vedecké poznatky?

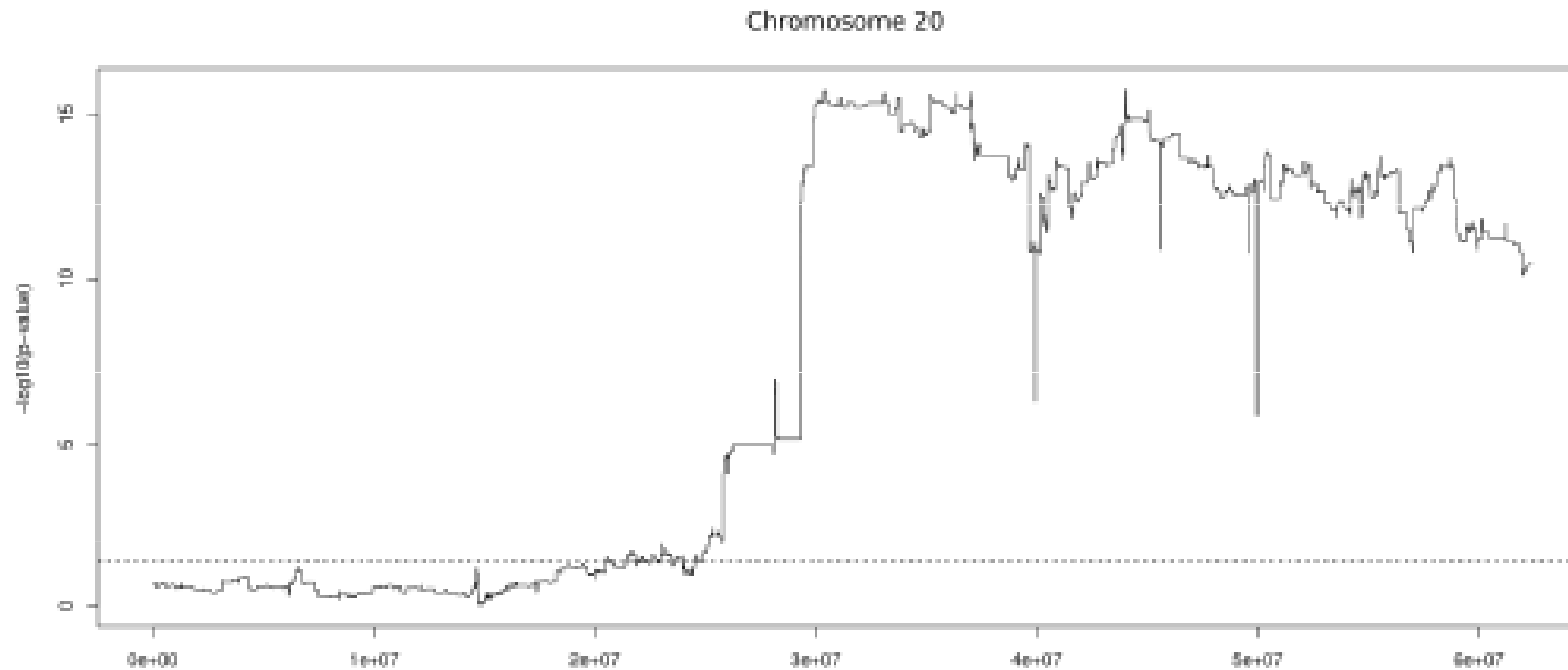
Je rozloženie týchto charakteristík medzi
podtypmi porovnateľné vo validačnom súbore?

***Majú objavené skupiny biologickú podstatu /
odrážajú známe vedecké poznatky?***

Podtyp C – pravostranné, BRAFm, MSI nádory, ktoré sú známe zlým prežitím po relapse



Zvýšená expresia génov chr. 20q v podtype B by mohla znamenať amplifikáciu chr20q regiónu.



Podtyp D – mezenchymálny – histologické vyhodnotenie: v nádore prítomná silná desmoplastická reakcia (mezenchymálne tkanivo)

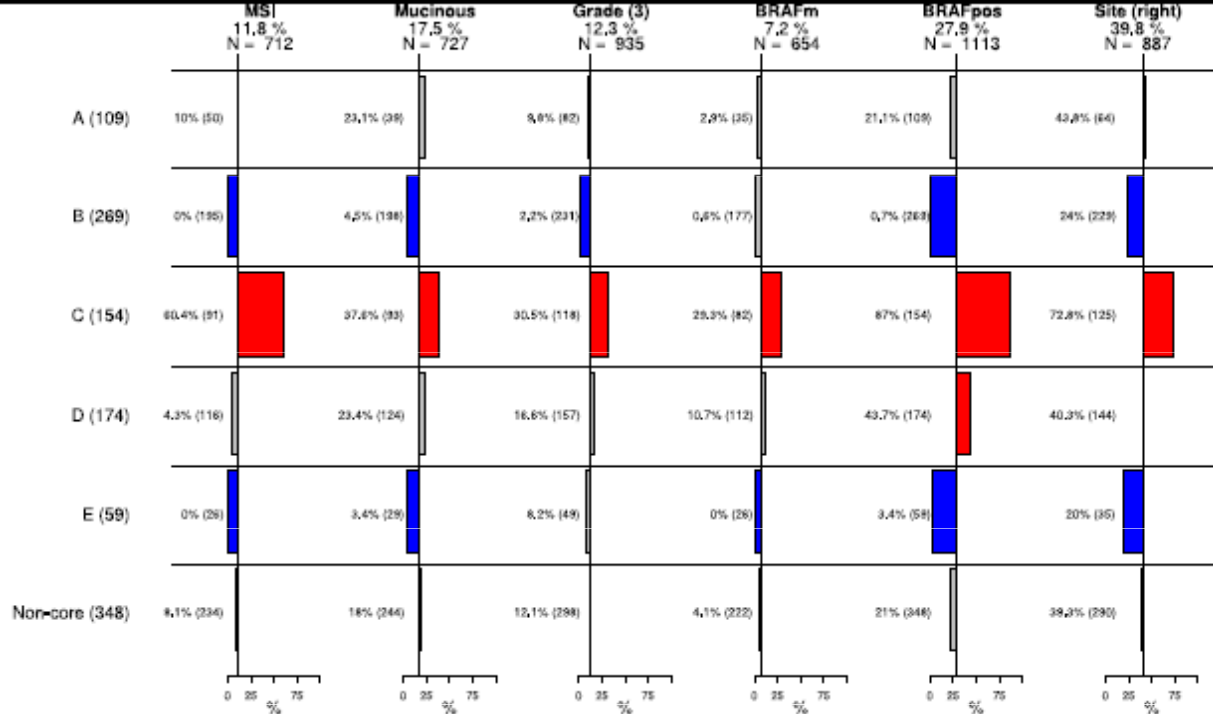
Nádor



***Je rozloženie klinických charakteristík medzi
podtypmi porovnateľné vo validačnom súbore?***

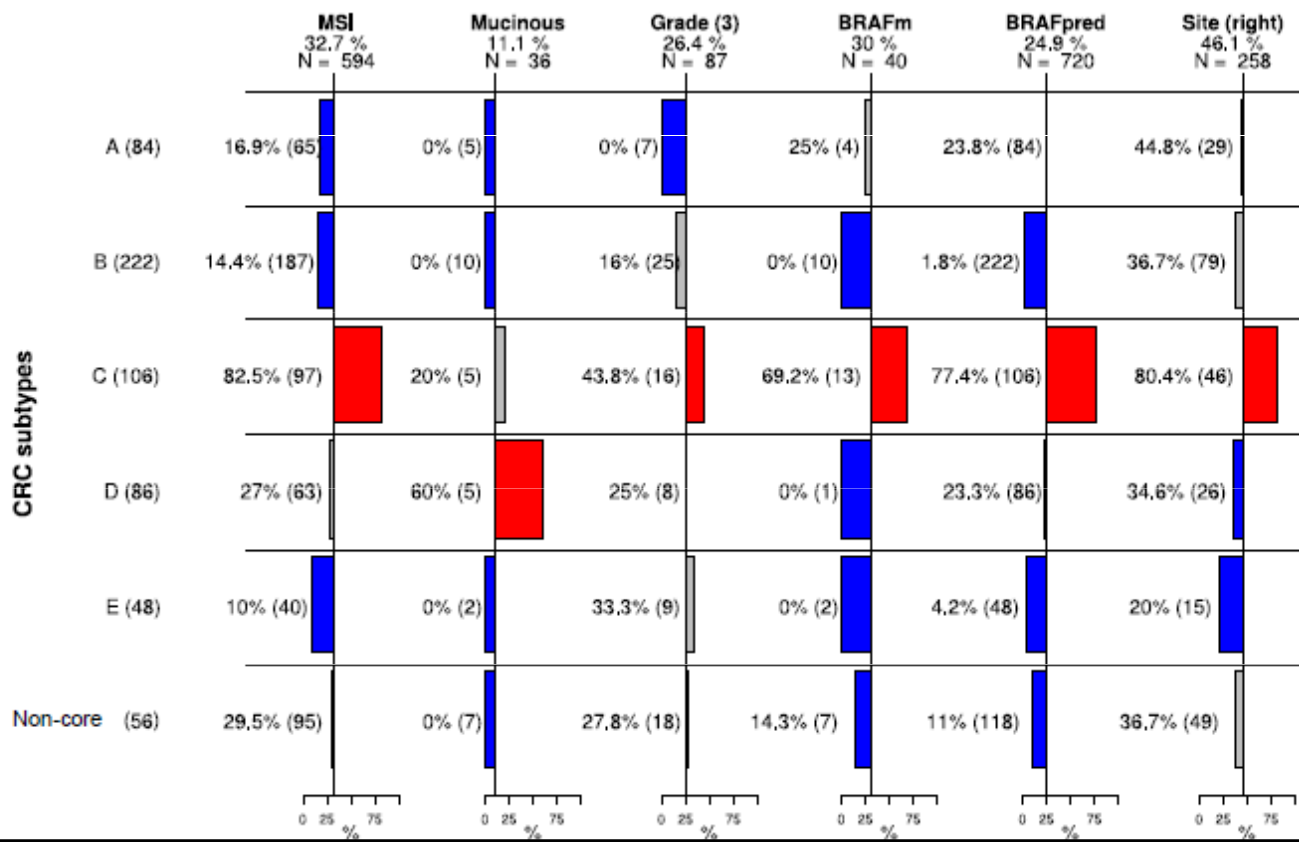
Discovery

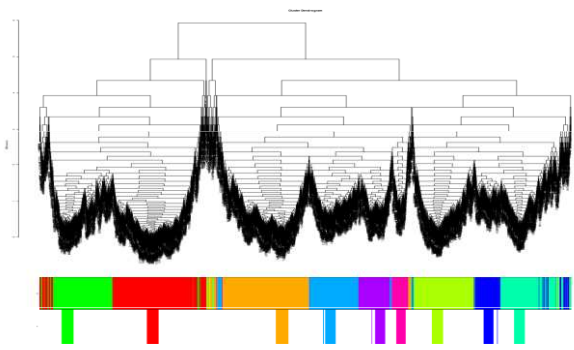
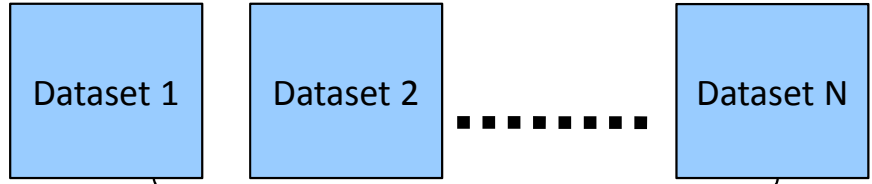
CRC subtypes



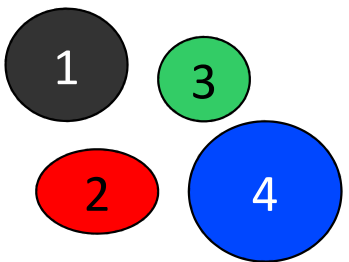
Validation

CRC subtypes





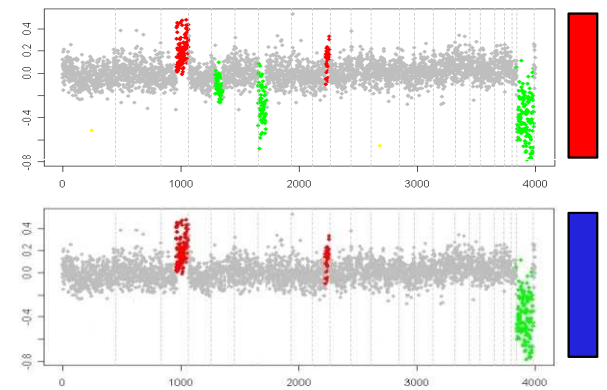
Klinické dáta



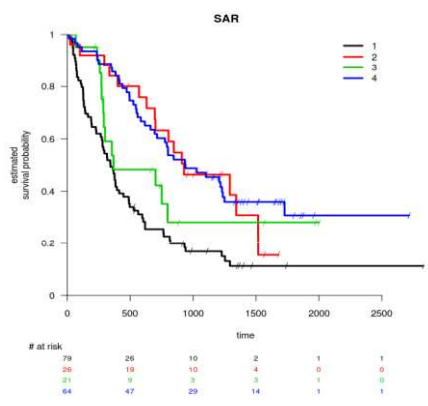
| | Right | MSI | KRAS | BRAF |
|-----------|-------|-----|------|------|
| ● (black) | 35% | 25% | 10% | 10% |
| ● (red) | 15% | ... | | |
| ● (green) | 0% | | ... | |
| ● (blue) | 50% | | | ... |

Definícia podtypov karcinómu

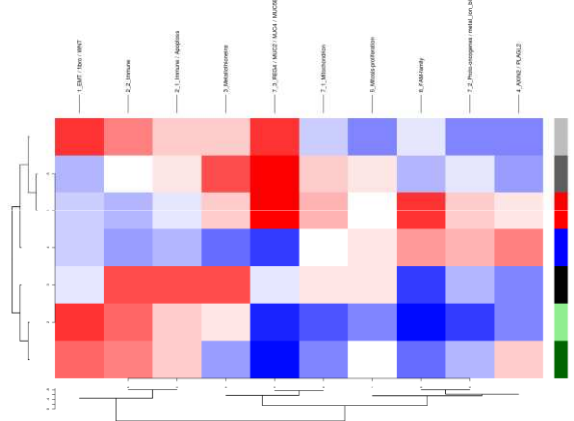
Iné zdroje dát (CNV, metylácia...)



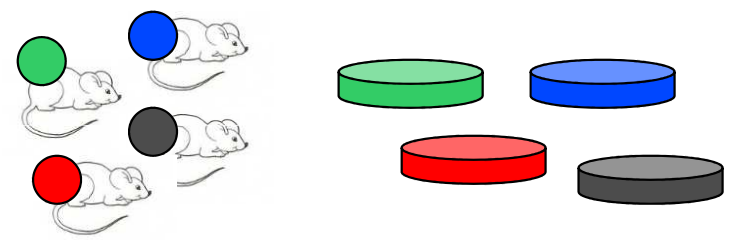
Prežitie



Molekulárny profil



Rezistencia na liečbu



Zhrnutie

- Konsenzusové zhukovanie
- Dynamické rezanie stromu
- Validácia neznámej pravdy – medzi dátovými súbormi
- Množstvo metód v rámci jednej štúdie
- Je to úžasné