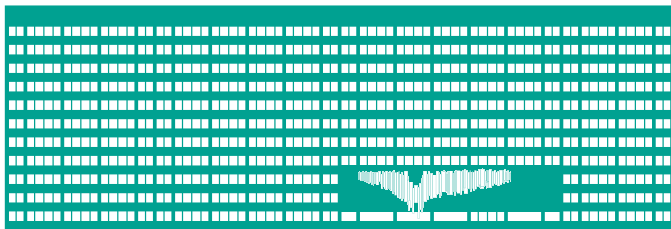


VŠB TECHNICKÁ
UNIVERZITA
OSTRAVA

VSB TECHNICAL
UNIVERSITY
OF OSTRAVA



www.vsb.cz

Algoritmy pro Bioinformatiku

Exprese genů

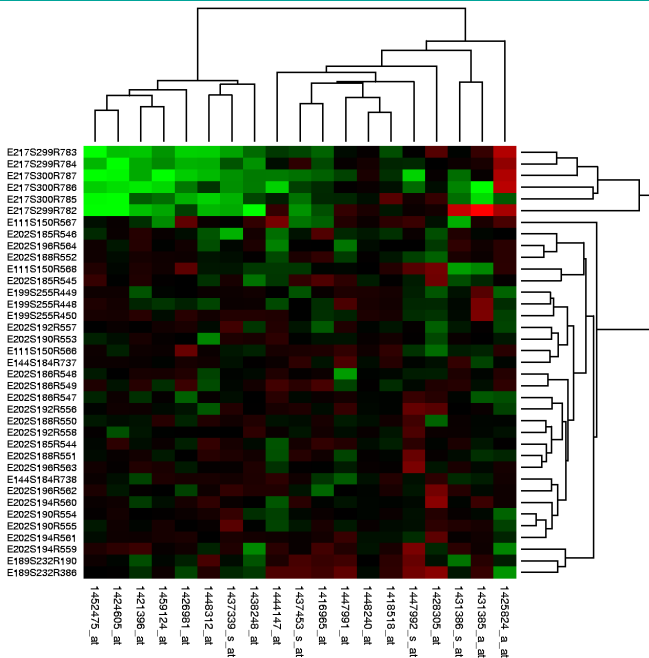
Michal Vašínek

VŠB – Technická univerzita Ostrava

FEI/EA404

michal.vasinek@vsb.cz

20. listopadu 2019





- Co to je exprese genů.
- Postup při analýze.
- Ukázka v R.



- Genom - množina genů jednoho organismu.
- Transcriptom - mRNA v buňce v jednom konkrétním okamžiku.
- Proteom - všechny polypeptidy v buňce v jednom konkrétním okamžiku.
- Interactom - množina všech protein-protein interakcí v buňce v jednom konkrétním okamžiku.



- Proces během kterého je informace uložená v DNA převedena na buněčnou funkci nebo strukturu formou přepisu DNA do proteinu.
- Exprese genu regulována na mnoha úrovních:
 - Navázání RNA polymerázy ovlivněno přítomností represorů a aktivátorů - váží se poblíž promotérů a ovlivňují schopnost polymerázy zahájit transkripci.
 - Regulace aktivátorů a represorů pomocí koregulátorů - koregulátory jsou přenašeči signálu, který způsobí navázání nebo odpojení regulačních proteinů.
 - Mutace v oblasti promotérů mohou ovlivnit schopnost polymerázy nasednout na DNA.
 - MicroRNA naváže se na mRNA zabrání translaci do proteinu.



- Které geny jsou aktivní při běžných buněčných procesech?
- Jaký je rozdíl aktivity genů mezi dvěma skupinami vzorků (např. pacienti vs. zdraví jedinci).



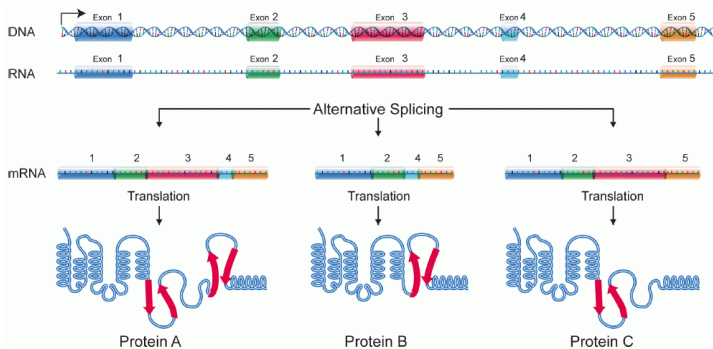
RNA sekvenování využívá standardních nástrojů sekvenování nové generace (Next-Generation Sequencing - NGS).

- Vzorek RNA je odebrán ze zkoumané tkáně, buňky, skupiny buňek.
- V daném okamžiku je množství RNA závislé na stavu buňky (je vystavena útoku, stresu, změně okolních podmínek, reakce na signály).
- Některá RNA se může vyskytovat více, tvorba jiné naopak potlačena.
- Chemickými reakcemi se RNA molekula reverzně přepíše na cDNA (chybí úseky DNA, které se nepřepisují na RNA).
- cDNA lze osekvenovat standardními nástroji NGS.



- RNA-seq technologie umožňuje detekovat transkripty s velmi nízkou i velmi vysokou úrovní exprese.
- Experimenty jsou reprodukovatelné - mají menší variabilitu způsobenou způsobem měření.
- Technologie RNA-seq není limitována na organismy s popsáním genomem.
- RNA-seq kombinuje znalost o mutacích v přepisovaných oblastech DNA s množstvím exprimovaných genů (prekurzorů proteinů), díky tomu technologie umožňuje nalezení dosud neobjevených genů, či úseků, které se přepisují do nějakého typu RNA.

- Rozdíl v expresi genů mezi různými vzorky.
- Studium alternativních spojení exonů při transkripci z DNA => vede k jiným proteinům.

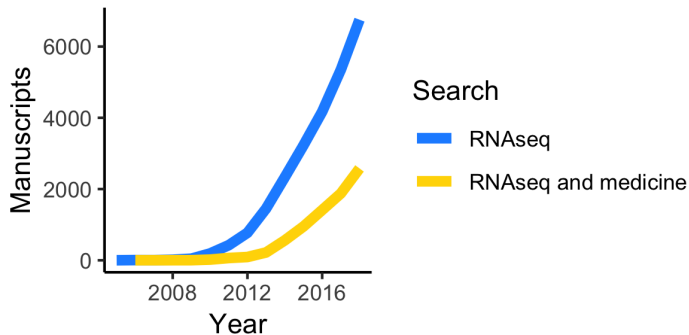




- Identifikace exprese specifické pro konkrétní alelu u bodových záměn, charakteristika rakovinných procesů.
- Další aplikací: RNA-seq jedné molekuly - nová oblast zkoumání komplexních biologických procesů zaměřených na chování jedné individuální buňky.
 - Výzkum kmenových buněk - kmenová buňka se diferencuje ve specifický typ buňky, jak je proces diferenciacce řízen, které geny (proteiny) jsou exprimovány, které potlačeny.



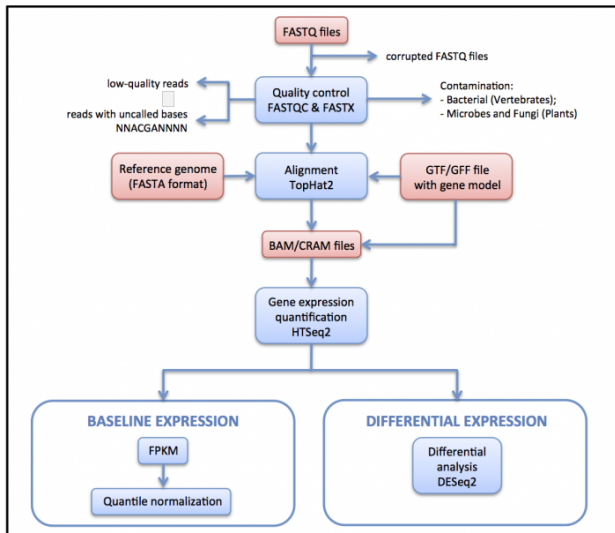
- Vyvinuto v 2005.
- Exponenciálně rostoucí počet vědeckých výstupů.





Ve chvíli, kdy je vzorek cDNA osekvenován, obdržíme k analýze standardní FASTQ soubory. Konkrétní metody se mohou lišit pro různé experimenty, ale v základu obsahují následující kroky:

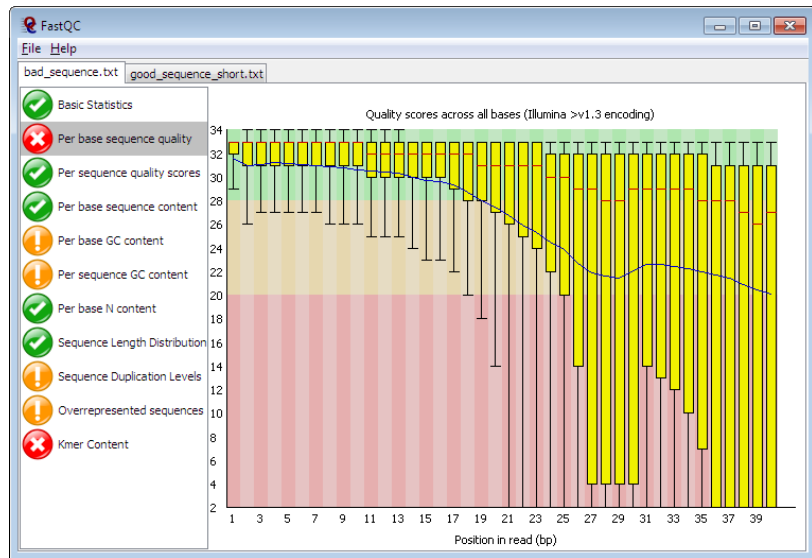
- Kontrola kvality sekvencí ve FASTQ souborech
- Přiložení sekvencí k referenční sekvenci, alternativně de-novo sestavení cDNA genomu.
- Kvantifikování exprese genů (výpočet konkrétních počtů sekvencí RNA překládaných na proteiny).
- Detekce diferenciálně exprimovaných genů (meziskupinové porovnání).



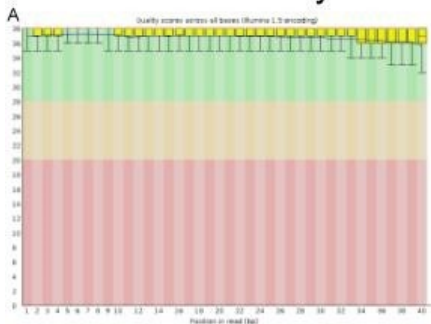


- Odstranění adaptérových sekvencí. Adaptérové sekvence obsahují specifické sekvence nukleotidů¹.
- Vyfiltrování sekvencí s nízkou kvalitou.
- Vyfiltrování nebo přepracování sekvencí s nepřečtenými nukleotidy.
- Odfiltrování kontaminovaných sekvencí (sekvence jiných organismů, bakterií, virů).
- Průmyslovým standardem ověření kvality je použití nástroje FastQC.
- Odstranění primerů, adaptérových sekvencí, případně odstranění nepřečtených bazí - standardem nástroj Trimmomatic.

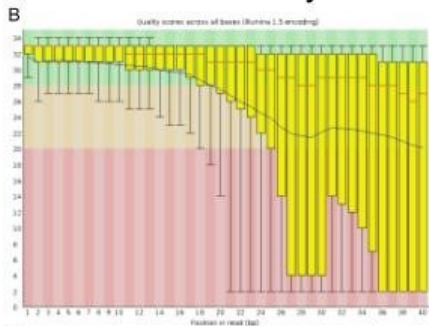
¹<https://www.youtube.com/watch?v=fCd6B5HRaZ8>



Good Quality



Bad Quality





- Na vstupu již předpokládáme sekvence o vysoké kvalitě. Dva možné postupy v závislosti na existenci referenční sekvence:
 - Namapování na existující referenční sekvenci
 - Namapování sekvencí nevyžaduje znalost pozic genů, nebo různých variant přepisu exonů.
 - Tento přístup zároveň umožňuje nalezení nových, neoanotovaných transkriptů.
 - Nutnost řešit sekvence překrývající více exonů.
 - De-novo sestavení transkriptů.
 - Sestavení sekvencí do skupiny překrývajících se sekvencí => tzv. contigy.
 - Na tyto contigy pak původní sekvence znovu namapujeme pro určení počtu sekvencí překrývajících určitý transkript.



Existuje celá řada nástrojů, které umožňují namapování krátkých sekvencí na referenční sekvenci, velmi populární TopHat, BWA, Bowtie, Hisat.

TopHat například pracuje ve dvou krocích:

- Nejdříve jsou sekvence nepřekrývající více exonů namapovány na referenční sekvenci.
- Ve druhém kroku se pokoušíme o alignment dosud nenamapovaných sekvencí. Například využijeme pigeonhole principle, část sekvence se nám podaří namapovat na jeden exon, druhou část na druhý exon. Očekáváme, že se exony vyskytnou relativně blízko sobě.
- Tento krok je následován kontrolou kvality namapování sekvencí (bereme v potaz především počet mutací a existence alternativních alignmentů). Používaný nástroj Picard.



Různé počty readů

Mějme dva vzorky C a P . První vzorek obsahuje 5 milionů readů a druhý 10 milionů readů.

Gen	C	P
gen1	4000	4000
gen2	4000	8000
...
Celkem	5 mil.	10 mil.

- Jsou geny 1 a 2 u vzorku P jinak exprimovány než u vzorku C .



Různé celkové počty readů u různých vzorků vedou na první úlohu o normalizaci:

Count per milion reads - CPM

Udává počet readů namapovaných na gen vztažený k jednomu milionu readů.

$$CPM(gene) = \frac{N_{gene} \times 10^6}{N} \quad (1)$$

- N_{gene} je počet readů namapovaných na konkrétní gen,
- N je celkový počet readů.



Mějme dva vzorky C a P . První vzorek obsahuje 5 milionů readů a druhý 10 milionů readů.

Gen	C	P	CPM(C)	CPM(P)
gen1	4000	4000	800	400
gen2	4000	8000	800	800
...
Celkem	5 mil.	10 mil.	1mil.	1 mil.



Různé délky genů

Mějme ve vzorku 5 milionů readů a např gen1 s počtem pokrývajících readů $N_{gen1} = 1000$ a délkou $|gen| = 2000\text{bp}$.

Gen	N_{gen}	$ gen $
gen1	1000	2000
gen2	1000	4000
gen3	2000	4000

Tabulka: Table caption

- Dva geny se stejným pokrytí ready.
- Určení, který gen je exprimován více.



Různé délky genů vedou k potřebě normalizace počtu sekvencí:

RPKM - reads per kilo base of transcript per million mapped reads

RPKM udává počet readů na kilobázi transkriptu (přepsaného genu) v miliónu namapovaných readů.

$$RPKM(gen) = \frac{CPM(gen) \times 10^3}{|gen|}$$

Normalizace není zapotřebí pro vyhodnocení exprese mezi různými skupinami vzorků, nicméně je nezbytná pro vzájemné porovnání exprese genů v jednom vzorku.



Mějme ve vzorku 5 milionů readů a např gen1 s počtem pokrývajících readů $N_{gen1} = 1000$ a délkou $|gen_1| = 2000\text{bp}$.

Gen	N_{gen}	$ gen $	RPKM
gen1	1000	2000	100
gen2	1000	4000	50
gen3	2000	4000	100



- Diferenciální analýza má na vstupu normalizované počty readů.
- Předpokládá alespoň dvě skupiny vzorků.
- Diferenciální analýza se snaží pomocí statistických metod prokázat změny v expresi genů oproti kontrolní skupině vzorků.
- Statistické testování slouží k ohodnocení významnosti pozorované změny, jinými slovy zda-li je rozdíl větší, než jaké hodnoty bychom očekávali na základě přirozené náhodné variability.



- Vytvořena celá řada nástrojů.
- limma - moderated t test.
- edgeR a DESeq založeny na modelování pomocí negativního binomického rozdělení.
- baySeq a EBSeq založeny na Bayesovské statistice.
- Samotný experiment obvykle předurčí jednu z metod, některé metody dokáží provést diferenciální analýzu pouze mezi dvojicí skupin. Nástroje jako edgeR, limma-voom, DeSeq či maSigPro dokáží testovat rozdílnost více skupin.




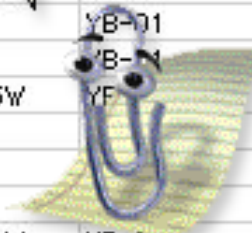
“Data don’t make any sense,
we will have to resort to statistics.”



UKR_Name	gene_name	plate_col	row_col
YBR124W	YBR124W	YB-01	d
YBR094W	YBR094W	YB-01	l
YBR091C	MRS5	YB-01	l
YBR078W	ECM33	YB-01	h
YBR075W	YBR075W	YB-01	h
YBR072W	HSP26	YB-01	h
YBR069C	VAP1	YB-01	h
YBR054W	YR02	YB-01	d
YBR051W	YBR051W	YB-01	d
YBR048W	RPS11B	YB-01	d

It looks like you're trying to do bioinformatics in Excel.

 Download R





- Náhodná proměnná z z normálního rozdělení má střední hodnotu μ a rozptyl σ^2 .
- Říkáme, že z má rozdělení $N(\mu, \sigma^2)$.
- Hustota pravděpodobnosti je dána:

$$g(z) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2}$$



Definition

Nechť z_1, z_2, \dots, z_n je náhodný výběr z normálního rozdělení $N(0, 1)$. Náhodný vektor $z = (z_1, z_2, \dots, z_n)$ má rozdělení $N_n(0, 1)$, pak dle definice:

$$\sum_{i=1}^n z_i^2 = z'z \text{ má rozdělení } \chi^2(n)$$

- Součet druhých mocnin n nezávislých normálně rozdělených náhodných proměnných má χ^2 rozdělení s n stupni volnosti.



Definition

Nechť z je $N(0, 1)$, u je $\chi^2(p)$ a z a u jsou nezávislé náhodné proměnné, pak je:

$$t = \frac{z}{\sqrt{u/p}} \text{ má rozdělení } t(p)$$

tzv. t rozdělení s p stupni volnosti.



- Mějme dvě skupiny vzorků, uvažujme, že tyto dvě skupiny se liší nějakým rysem, který nás zajímá.
- Pro obě skupiny vzorků jsme provedli sekvenování RNA.
- Chceme zjistit, které geny jsou charakteristické pro zkoumanou odlišnost,
- Předpokládejme, že rozdělení počtu sekvencí je normální náhodná proměnná a pro obě skupiny platí, že mají stejný rozptyl.
- Můžeme použít t test o shodě průměrů náhodných veličin se stejným rozptylem.

$$t = \frac{\bar{x} - \bar{y}}{s \sqrt{\frac{1}{n_x} + \frac{1}{n_y}}}$$



$$t = \frac{\bar{x} - \bar{y}}{s \sqrt{\frac{1}{n_x} + \frac{1}{n_y}}}$$

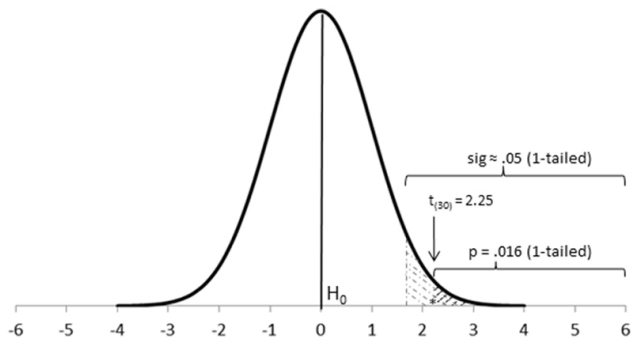
- n_x, n_y jsou počty vzorků v jednotlivých skupinách.
- Nulová hypotéza: oba průměry jsou shodné, tedy $\bar{x} - \bar{y} = 0$.
- Statistika t má rozdělení $T(n_x + n_y - 2)$.
- Odhad rozptylu

$$s^2 = \frac{(n_x - 1)s_x^2 + (n_y - 1)s_y^2}{(n_x + n_y - 2)}$$

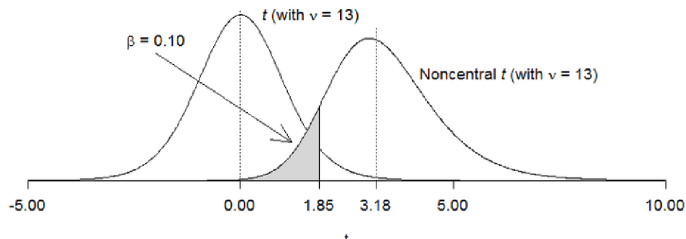


$$t = \frac{\bar{x} - \bar{y}}{s \sqrt{\frac{1}{n_x} + \frac{1}{n_y}}}$$

- Pro vypočtenou hodnotu t spočítáme pravděpodobnost, že hodnota je větší nebo rovna t tzv. p-hodnota (p-value).
- Podle předem zvolené hladiny významnosti α porovnáme zda je $p(X > t) < \alpha$, pokud ano, pak zamítáme nulovou hypotézu.
- Hladina významnosti nám udává chybu, že jsme nesprávně zamítli nulovou hypotézu.



- Pro vypočtenou hodnotu t spočítáme pravděpodobnost, že hodnota je větší nebo rovna t tzv. p-hodnota (p-value).
- Podle předem zvolené hladiny významnosti α porovnáme zda je $p(X > t) < \alpha$, pokud ano, pak zamítáme nulovou hypotézu.
- Hladina významnosti nám udává chybu, že jsme nesprávně zamítli nulovou hypotézu.



- Nulová hypotéza má centrální t rozdělení.
- Alternativní hypotéza má necentrální t rozdělení.
- Rozdělení alternativní hypotézy neznáme, máme pouze odhady.



- Normální rozdělení předpokládá spojité hodnoty počtu readů.
- Geny s nižší expresí mají větší roztyl a mají tendenci být zešikmené.
- Pro geny s nižší expresí neplatí předpoklad normálního rozdělení.
- Doporučuje se aplikovat \log_2 transformaci počtu readů před vyhodnocením významnosti.



Předpokládejme, že následující model popisuje nějakou modelovou lineární závislost pro n pozorování.:

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

- y a ϵ jsou náhodné proměnné a předpokládejme, že x je vždy nějaká konstanta.
- $E(\epsilon) = 0$ jinými slovy $E(y) = \beta_0 + \beta_1 x$.
- $var(\epsilon) = var(y) = \sigma^2$



Předpokládejme, že máme náhodný vzorek n pozorování y_1, y_2, \dots, y_n . Neznáme β_0 ani β_1 a potřebujeme získat nějaký jejich odhad $\hat{\beta}_0, \hat{\beta}_1$ a pomocí těchto odhadů vytvoříme predikci \hat{y} hodnoty y pro nějakou hodnotu x .

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

- $\hat{\beta}_0$ a $\hat{\beta}_1$ získáme metodou nejmenších čtverců - minimalizací rozdílu mezi skutečnou hodnotou y a predikovanou \hat{y}_i .

$$\hat{\epsilon}\hat{\epsilon} = \sum_i \hat{\epsilon}_i^2 = \sum_i (y_i - \hat{y}_i)^2 = \sum_i (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$$

- $y_i - \hat{y}_i$ nazýváme residuum y , nebo také chyba odhadu.
- $(y_i - \hat{y}_i)^2 = SSE$ nazýváme residuum součtu čtverců, nebo také chybu odhadu součtu čtverců.



Derivací podle $\hat{\beta}_0$ a $\hat{\beta}_1$ a položením derivací rovno 0, nalezneme odhady $\hat{\beta}_0$ a $\hat{\beta}_1$, které minimalizují $\hat{\epsilon}\hat{\epsilon} = SSE$:

$$\hat{\beta}_1 = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_i (x_i - \bar{x})^2}$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

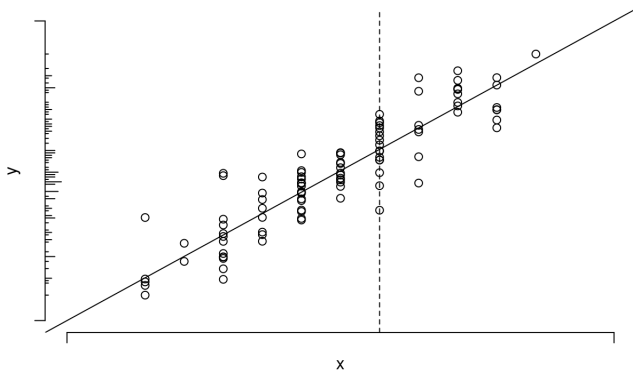


$$t = \frac{\hat{\beta}_1}{s / \sqrt{\sum_i (x_i - \bar{x})^2}}$$

- Statistické testování je obvykle založeno na předpokladu, že pokud mezi dvěma proměnnými neexistuje lineární vztah, pak je $\beta_1 = 0$.
- Hypotéza H_0 je tedy $H_0 : \beta_1 = 0$.
- Odhad rozptylu:

$$s^2 = \frac{\sum_i (y_i - \hat{y}_i)^2}{n - 2} = \frac{\sum_i (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2}{n - 2}$$

- Hodnoty x_i volíme např. $x_i = 0$ pokud vzorek je z kontrolní skupiny a $x_i = 1$ pokud je ze skupiny pacientů.



- Umožňuje testovat lineární závislost pouze mezi dvěma kategoriemi, např. zdraví vs. nemocní.



	A	B	C	D	E
1		baseMean	log2FoldChange	stat	pvalue
2	APOL1	2290.56530882017	3.10832963628912	8.61937720535332	6.73192495579598e-18
3	LYZ	77870.2427655811	-7.0510192436375	-7.74791509820844	9.34135604883871e-15
4	SLC6A8	678.327168089741	1.90544175497903	7.52749712316208	5.17221477079242e-14
5	SEMA4B	1741.57724244313	1.49340137828406	7.15320066774613	8.47775311396553e-13
6	STATH	3313.87839012742	-8.16963354067458	-7.06518543076518	1.60401654729597e-12
7	ZG16B	1487.34560622635	-5.7128835515497	-6.80027227078752	1.04421612857286e-11
8	KCNMA1	259.585606740072	-2.76284495857023	-6.79096191874476	1.11388327853096e-11
9	ATP2A3	730.994954105378	-3.69920010406926	-6.78428505072016	1.16662786413202e-11
10	EGLN3	706.580777036715	3.49475068397002	6.77969308994643	1.20431442080478e-11
11	SOX4	1503.45484443017	1.32943866871453	6.7756078105378	1.23884320257993e-11
12	PRH2	9053.73856673143	-7.84574482823239	-6.74997021728914	1.47875506823875e-11
13	MARK1	459.174938188555	0.785729633569354	6.62769799546138	3.40962140042419e-11
14	FERMT2	166.730996298386	-2.0893443692225	-6.4386924529594	1.20507094321743e-10
15	APOL2	943.627078393778	1.62708093956464	6.40488809297558	1.50479712245154e-10
16	NDUFA4L2	372.555153565341	3.77626062525379	6.35736307736515	2.05246477081785e-10
17	DMBT1	3703.80190168882	-8.62250414133728	-6.29771090292796	3.02073212034546e-10

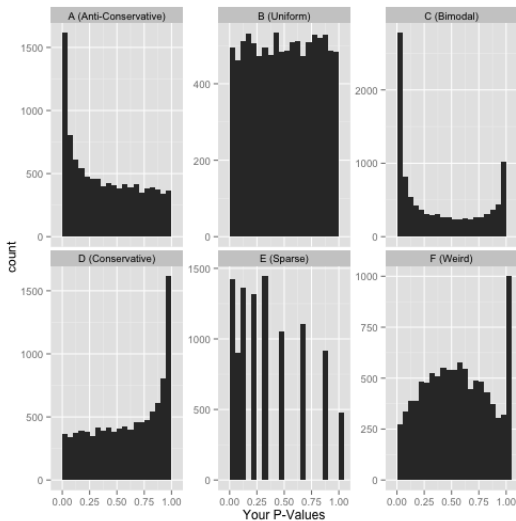


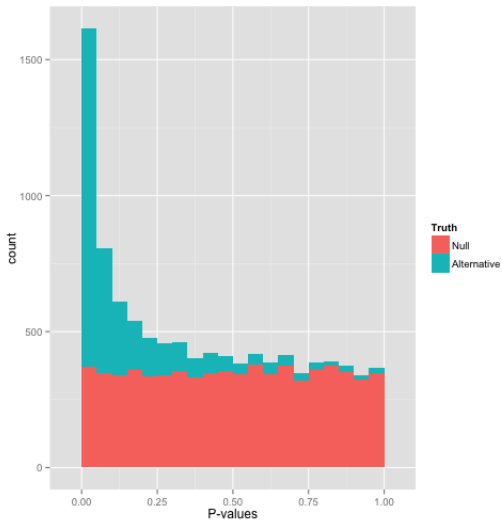
- Porovnání míry exprese mezi skupinami vzorků.
- Vztaženo ke kontrolní skupině.
- Průměrné exprese: kontrolní skupina $\overline{y_c}$, zájmová skupina (pacienti) $\overline{y_p}$
- Fold change:

$$fc = \frac{\overline{y_p}}{\overline{y_c}}$$

- V logaritmickém tvaru čteme změnu exprese stejně pouze s opačným znaménkem.

$$\log_2 fc = \log_2 \frac{\overline{y_p}}{\overline{y_c}}$$





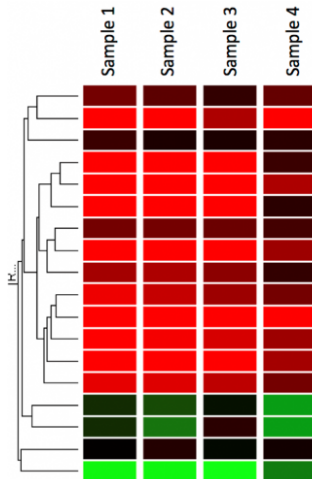


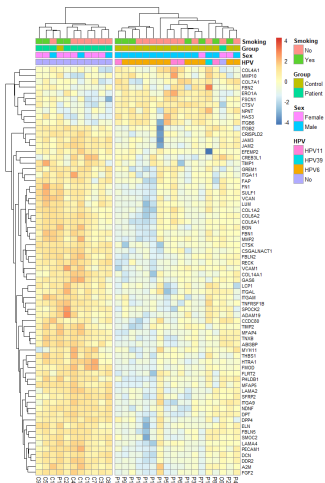
- Obecně počty readů spíše odpovídají negativnímu binomickému rozdělení.
- Sestavujeme obecný lineární regresní model.
- Waldův test.



Nejbvyklejší metodou vizualizace (vědecko-výzkumná úroveň) exprimovaných dat jsou tzv. heatmapy.

- Použito hierarchické shlukování.
- Geny shlukovány podle podobné exprese.
- Pacienti shlukováni podle podobné exprese genů.







- Obvykle se nám nepodaří rozdělit vzorky do skupin.
- Vzorky nejsou dokonalé.
- Vzorky nemusí v konkrétní časový okamžik odběru a sekvenování popisovat stejný stav buněk.
- Proces analýzy zahrnuje mnoho kroků, každý krok vnáší určitou chybu.
- Exprese genů je komplexní proces, v experimentu máme pouze RNA sekvence, důvodů exprese však může být více (např. fenotypy, nemoci).
- Vzorky měřeny na různých strojích => potřeba mezistrojového přizpůsobení, plánování experimentu - technické replikáty.
- Magické statistické metody pro zahrnutí variability měření a skrytých proměnných, RUV (remove unwanted variance), SVA (surrogate variable analysis)

DĚKUJI za pozornost

Michal Vašinek

VŠB – Technická univerzita Ostrava

FEI/EA404

michal.vasinek@vsb.cz

20. listopadu 2019