

Fikiina

Statnicová

Sněška

Ni-zi

1. Automatické plánování, plánovací graf, kompilace plánování do jiných formalismů jako je SAT nebo CSP, hierarchické plánování, plánování v prostoru plánů. Plánování pohybu a problém lokalizace v robotice.

NI-UMI

Automatické plánování

- Úloha hledání **posloupnosti akcí**, které vedou k dosažení určitého cíle
- **Klasické plánování** – plně pozorovatelné, deterministické, konečné, statické, diskrétní, offline (nejdříve naplánujeme, pak vykonáme)
- **Neklasické plánování** – stochastické, dynamické prostředí
- **Vyjadřovací prostředky**
 - o **Atom** – predikátový symbol s dosazenými konstantami
 - o **Jazyk** – predikátové symboly a konstanty (nulární funkce)
 - o **Stav** – konečné množiny základních atomů, uzavřený svět (co není zmíněno, to neplatí)
 - o **Cíl** – částečná specifikace stavu konečnou množinou literálů
- **Operátory**
 - o Trojice $o(\text{name}(o), \text{precond}(o), \text{effects}(o))$
 - **Name(o)** – $n(x_1 \dots x_n)$
 - **Precond(o)** – literály, které musí být splněny
 - **Effect(o)** – množina literálů, které budou platit po provedení operátoru
- STRIPS – dnes v podobě PDDL
- **Akce** – základní instance operátoru
 - o **Aplikovatelnost akce**:
 - Mějme akci a a stav s
 $\text{Precond}^+(a), \text{precond}^-(a)$
 $\text{Effect}^+(a), \text{effect}^-(a)$
 - Akce je aplikovatelná ve stavu s : $\text{Precond}^+(a) \subseteq s \wedge \text{precond}^-(a) \cap s = \emptyset$
 - o **Aplikace apl. Akce**:
 - Odebereme negativní efekty, přidáme pozitivní efekty:
$$\gamma(s, a) = (s - \text{effect}^-(a)) \cup (\text{effect}^+(a))$$
- **Plánovací doména** – jazyk (konstanty a predik. S.) L a operátory O
 - o Implicitně určuje množinu stavů a funkci následníka – $\gamma: S \times A \rightarrow S$
- **Plánovací problém** – předp. plánovací domény: L, O
 - o Zadání plánovacího problému: (O, S_0, g)
 - o Plánovací problém SP s přechody: (Σ, s_0, s_g)
 - $\Sigma = (s, A, \gamma)$
 - S_g = množina cílových stavů
- **Plány a řešení**
 - o **Plán** $P = (o, s_0, g)$ – **posloupnost akcí** $\pi = [a_1, \dots, a_n]$, kde a_i = instance operátoru z O
 - o Plán je řešením P, jestliže je proveditelný z s_0 a dosahuje cíle g
- **Dopředné plánování** – nedeterministický výběr akce do navštíveného stavu, tím se vytvoří nový stav
- **Zpětné plánování**
 - o **Inverzní aplikace** – inverzní funkce k funkci následníka, tvořím pod-cíl, ze kterého se dostat do cíle
 - Lze uvažovat jen relevantní akce – ty, které nás posunou směrem ke splnění cíle (splní alespoň 1 cílový literál), a zároveň nezruší žádný ze splněných literálů cíle
 - Cíl g a akce a – $\gamma^{-1}(g, a)$
 - Akce a je relevantní pro cíl g , když:
 - o $g \cap \text{effects}^+(a) \neq \emptyset$
 - o $g^+ \cap \text{effects}^-(a) = \emptyset$ a $g \cap \text{effects}^-(a) = \emptyset$
 - Pro relevantní a vzhledem k g definujeme: $\gamma^{-1}(g, a) = (g - \text{effect}^-(a)) \cup \text{precond}(a)$

- Jako u DP se vybírá nedeterministický nový stav z předchozího, ale tentokrát se jde od cíle k počátečnímu stavu
- **Liftování** – zavedení proměnných zmenší větvicí faktor
- **Prostor plánů**
 - Postupně opravujeme částečný plán – množina částečně instanciovaných operací/akcí, množina podmínek
 - Oprava – odstranění **kazu**
 - Uzly prohledávání jsou částečné plány
 - Kaz 1: otevřený cíl
 - Máme akci s předpokladem p, o kterém nevíme, jak jej nastavit/splnit
 - Léčba:
 1. Najít a, co by mohla produkovat p
 2. Instanciovat proměnné a/nebo nastavit vazbové podmínky
 3. Nastavit kauzální podmínku
 - Kaz 2: Hrozba
 - Existuje možnost, že si smažeme (c) předpoklad (p) produkováný pro nějakou akci (b)
 - Léčba:
 1. Přinutit b, aby předcházelo a
 2. Přinutit c, aby předcházelo a
 3. Svázat proměnné tak, aby ke smazání p nedošlo
- Plánování v prostoru plánů (PSP) – vrací částečně uspořádaný plán
- **Analýza dosažitelnosti**
 - Zaměření na větvicí faktor, zkoumání dosažitelnosti stavů (cílových)
 - Expanze stavového prostoru
 - Narazíme-li při expanzi na cílový stav, zrekonstruujeme plán
 - Expanze -> strom dosažitelných stavů ($O(k^d)$)
 - **Graf dosažitelných stavů** – sloučení opakujících se stavů ($O(k^d)$)
 - **Relaxace dosažitelnosti**
 - Původní: s_g dosažitelný \Leftrightarrow obražen v grafu d. s.
 - Relaxovaný: s_g dosažitelný \Rightarrow obsažen v plánovacím grafu d. s.
 - Sjednocení atomů úrovně stromu/grafu dos. Do 1 stavu/uzlu
- **Plánovací graf**
 - Střídání stavových (P_0, \dots) a akčních (A_1, \dots) vrstev
 - Hrany zpřesňují aproximaci oproti jen sjednocení atomů
 - Hrany z P_i do A_{i+1} – ukazují předpoklady akce
 - Hrany z A_i do P_i – ukazují efekty akce
 - Skrz negativní efekty nemažeme
 - $P_0 \subseteq P_1 \subseteq \dots$
 - Axiom rámce (frame) – atomy přenášíme do další vrstvy
- **Paralelní plány** – určeny plán. Grafem – posl. Množin akcí
 - Převod na **sekvenční plán** – posloupnost množiny uspořádání akcí
 - Nutno zajistit **nezávislost** – mezi paralelně proveditelnými akcemi
 - Dvojice akcí (a, b) je nezávislá \Leftrightarrow

$$\text{Effect}^+(a) \cap (\text{precond}(b) \cup \text{effect}^+(b)) = \emptyset$$

$$\wedge$$

$$\text{Effect}^+(b) \cap (\text{precond}(a) \cup \text{effect}^+(a)) = \emptyset$$
 - Množina akcí π je nezávislá, jestliže je v π nezávislá každá dvojice akcí
 - Nezávislou množinu akcí lze aplikovat paralelně – akce si vzájemně neškodí
- **Vzájemné vyloučení – mutex**
 - V plán. Grafu mezi atomy nebo akcemi
 - Zpřesnění aproximace – zákaz nežádoucích dvojic
 - Akce a a b z A_i jsou vzájemně vyloučené (mutex) = jsou závislé, nebo je a mutex s předpokl. b

- Atomy p a q z P_i jsou vzájemně vyloučené – každá akce a v A_i , že $p \in \text{effect}^+(a)$ je mutex s každou akcí b v A_i , že $q \in \text{effect}^+(b)$ a v A_i není akce c , že $p, q \in \text{effect}^+(c)$
- Algoritmus **Graphplan** – střídá 2 fáze:
 1. **Expanze** plánovacího grafu
 2. Pokus o **extrakci** paralelního plánu z plánovacího grafu
 - Neúspěch \rightarrow pokračuje expanze, + úroveň do p. g. – po čase stabilizace
- Extrakce plánu
 - GP-search – chce bez-mutexovou množinu akcí π splňující g
 - Extract – tabulka nogoodů ∇
 $\nabla(i) \dots$ zákazové cíle pro úroveň i
- Plánování jako SAT
 - Stav odpovídá ohodnocení některých výrokových proměnných
 - Nutno dávat pozor na realitu
 - Cíl – vynutíme platnost cíle klauzulí
 - Problém – vývoj v čase \rightarrow časový index
 - Časová expanze – naznačuje, jak se postupně mění stavy v čase
 - Zavedení proměnných pro akce (a časové kroky) – pro každou časovou jednotku a každý atom proměnná, znamenající, že atom platí v daném čase
 - Pro každou časovou jednotku a akci další proměnnou – akce se provedla v daném čase
- **SATplan**
 - Využívá **plánovací graf** na určení atomů a akcí, které reprezentovat v časových krocích pomocí proměnných
 - **Mutexy** – zakazují určité kombinace akcí, přirozeně vyplývají z plánovacího grafu, podporují jednotkovou propagaci
 - Splňování formule odpovídá extrakci plánu z plánovacího grafu
- Hierarchické plánování – zkrácení plánů, rozděluje problém na málo velkých akcí, kde každou velkou akci rozděluje na pár menších atd.

Plánování pohybu a problém lokalizace v robotice

- **Lokalizace**
 - **Počáteční znalost**
 - počáteční poloha známá – stopování
 - poč. poloha neznámá – globální lokalizace – změna na stopování hned po objev. polohy
 - **Přechodový model** – pohyb robota v rovině, známe přesnou mapu, $X_t = (x_t, y_t, \theta_t)$, posuvná rychlost v_t , rotační ω_t – deterministická stavová predikce
 - skuteční roboti nejsou determinističtí \rightarrow použití **Gaussovske distribuce**
 - senzorový model – vzhovávání k orientačnímu bodu, který senzory detekují (vzdálenost, natočení)
 - lokalizace **Monte Carlo (MCL)** – založeno na částicové filtraci
 - repr. Přesvědčení jako **mračno částic** (vzorků), které odpovídají stavům
 - přechodový a senzorový model na vstupu
 - vrací, jak odpovídá pozorování senzorem simulaci pozorování ve vzorku
- **Kalmanův filtr** – přesvědčení robota $P(X_t | z_{1:t}, a_{1:t-1})$ reprezentuje pomocí normálního rozd. $N(\mu_t, \Sigma_t)$
 - Funguje s lineárními (fah) přechodovými a senzorovými modely
 - Přechodový model $X_{t+1} = f(X_t, v_t, \omega_t)$ interpretujeme jako maticovou operaci, senzorový model $z_{t+1} = h(X_t)$ taky
 - Přesvědčení reprezentované jako normální rozdělení $N(\mu_t, \Sigma_t)$ je po lineárních transformacích opět popsáno normálním rozdělením $N(\mu_{t+1}, \Sigma_{t+1})$
- **Mapování** – otázka vytvoření reprezentace prostředí na základě údajů ze senzoru
 - Kombinace technik – lokalizace, dynam. Vytváření senzor. Modelu, někdy prohledávání prostředí
 - Technické aspekty – reprezentace pomocí mračna bodů, možnost spolupráce robotů, pohybující se objekty v prostředí
 - Aplikace – autonomní vozidla, Henry

- **Plánování pohybu**
 - o Typy pohybů – z bodu do bodu, souladný – v kontaktu s překážkou
 - o **Konfigurační prostor** – poloha, orientace, natočení kloubů
 - Hledání cesty – v robotice ve spojitém prostoru
 - o Spojitost → rozklad na buňky, skeletonizace
 - o Pro zjednodušení nehledíme na neurčitost, vše deterministické
- **Konfigurační prostor**
 - o Reprezentace pomocí pracovních souřadnic – pozice prvků robota
 - Plně popisuje pozici
 - Vhodná na detekci kolizí
 - Ne všechny souřadnice jsou možné
 - o Reprezentace pomocí konfigurací – natočení kloubů...
- **Rozklad na buňky**
 - o Konečný počet spojitých buněk – hledání cesty v rámci buňky ez
 - o Jednoduché na implementaci
 - o Komplikuje se s rostoucí dimenzí
 - o Rozdělování polo-obsazených buněk na 2^{dimenze}
- **Skeletonizace** – převod hledání cesty na úlohu v jedné dimenzi
 - o **Voroného diagram** – množina bodů, které mají stejnou vzdálenost ke 2 a více překážkám
 - Tvoří graf, v němž je problém hledání cesty 1D
 - Vstup i výstup po úsečce
 - Obtížná konstrukce
 - o **Pravděpodobnostní mapa** – generuje náhodné konfigurace, spojíme ty, mezi nimiž vede snadná cesta
- **Kinematika**
 - o **Dopředná kinematika** – známe otočení kloubů – konfiguraci, chceme určit polohu efektoru – lineární transformace
 - Konfigurace => pracovní souřadnice (z konfiguračního do pracovního prostoru)
 - o **Inverzní kinematika** – známe pozici efektoru reprezentovanou pracovními souřadnicemi, chceme určit konfiguraci – složitá transformace
 - Řešení rovnic, víc DOF -> víc řešení (nekonečno) – hledáme řešení minimalizující něco (energii, ...)
 - Pracovní souřadnice => konfigurace (z pracovního do konfiguračního prostoru)
- **Dynamika a řízení** – nutno uvažovat nejen kinematický stav (konfiguraci), ale i rychlost – dynamický stav
 - o Plánování s dynamickými stavy – hard, řešení diferenciálních rovnic
 - o Alternativa – **metoda kontroleru** – je-li hodnota nějaké stavové proměnné xt odchýlena od požadované $y(t)$, robot ji kompenzuje působením at (síla efektoru)
 - Kontroler P – proporcionální – může oscilovat
 - Kontroler PD – proporciálně derivační – derivace je tlumič, na náhlou velkou změnu pomalá reakce
 - **Kontroler PID** – proporciálně derivačně integrační – bere v potaz systematické externí působení (zahrnuto v integračním faktoru), nejčastější

2. Splňování omezení s konečnými doménami (CSP), pokročilé prohledávání (backjumping, dynamický backtracking), filtrace domén a lokální konzistenční techniky, globální omezení, rozhodovací heuristiky.

NI-UMI

Splňování omezení s konečnými doménami (CSP)

- **Stav** v CSP není atomický, má vnitřní strukturu
 - o Rozdělen do **proměnných**
 - o Cíl definován jako **splnění omezení**
- Faktorová reprezentace -> doménově nezávislé heuristiky
- **CSP** = problém splňování podmínek – trojice (X, D, C)
 - X ... konečná množina **proměnných**
 - D ... konečná **doména** (obor hodnot) proměnných
 - C ... množina **podmínek** nad X – omezení pro rozhodnutí
 - Binární, ternární, n-nární... (podle |X|)
- **Řešení** = přiřazení hodnot proměnným tak, že všechny podmínky jsou splněny (konjunkce)
- **Prohledávání** v CSP
 - o **Stav** = částečné ohodnocení proměnné
 - **Konzistentní** – přiřazené hodnoty neporušují podmínky
 - o **Počáteční stav** – prázdné přiřazení
 - o **Akce** – přiřazení hodnoty proměnné z její domény
 - o **Cílový stav** – všechny p. ohodnoceny + konzistentní
 - Komutativita ohodnocení
- **Chronologický backtracking** – rekurzivní prohledávání do hloubky (DFS)
 - o Dopředný chod – přiřazení
 - o Nelze přiřadit -> u poslední přiřazené zkoušíme jinou hodnotu
- Složitost: CSP je **NP-těžké**
- Constraint programming – ILOG, Minion, Gecode, MiniZinc, ...
- Neefektivita backtrackingu
 - o Pasivní využití podmínek
 - o Zapomínání již odvedené práce – hodněkrát nastaví pár proměnných na stejné ohodnocení, které je nesplnitelné

Pokročilé prohledávání

- **Backjumping** – pokud v průběhu ohodnocování x momentální proměnná nemůže být nastavena na žádnou hodnotu, najdou se všechny nesplněné omezení, z nich se vyberou konfliktní proměnné, následně se odnastaví všechny proměnné až do nejbližší konfliktní proměnné
 - o Rozšiřujeme ohodnocení S' do proměnné x, zbývá ale ohodnotit víc proměnných
 - Přiřazení $S'(x)=d$ nezpůsobí konflikt, ale přesto nejde rozšířit na plné ohodnocení (odhalení konfliktu až při dalším ohodnocování)
 - Nemožnost přiřazení $S'(x)=d$ odhalí prohledávací algoritmus (rekurzivní volání BJ) při pokusu ohodnotit zbylé proměnné
 - Uvažujeme **konfliktní množinu** Conf pro zbylé proměnné, které v pozdějším běhu nebylo možné ohodnotit – Conf je příčina nemožnosti rozšířit $S' \cup \{x=d\}$ na celé X
 - S' pro Conf \ {x} je pak příčina nemožnosti provést $S'(x)=d$ – $\text{Conf}_d = \text{Conf} \setminus \{x\}$ je konfliktní množina pro hodnotu $d \in D(x)$
 - Nelze-li provést přiřazení $S'(x)=d$ pro žádnou hodnotu d z domény x, BJ skočí zpět na poslední proměnnou z konfliktní množiny
 - o **Backjumping** v listu
 - Máme částečné ohodnocení s' , ohodnocujeme poslední proměnnou x. Předpokládejme, že s' do x rozš. Nelze -> skok zpět na poslední prom. z konfliktní množiny pro x

- **Dynamický backtracking** – při skoku zpět část ohodnocení, která je OK, ponechá
 - o Změna pořadí proměnných (-> dynamika) – aby se zachovalo ohodnocení nekonzistentních proměnných
- **Závislosti (nogoody)**
 - o S' ...částečné ohodnocení
 - o $Conf = \{y_1 \dots y_n\}$ množina konzistentních proměnných k S'
 - o **Nogood** = podmínka zakazující vytvořit nerozšiřitelné přiřazení
 - o BJ objevuje nogoody – lze ukládat – učení (lze i zapomínat – šetření paměti)
 - o logický důsledek v daném CSP
 - o Specifické lze skládat do obecnějších – rezoluční pravidlo

Filtrace domén, konzistence, globální omezení a heuristiky

- **Filtrace domén** – dopředná kontrola
 - o Zavádíme pracovní domény pro proměnné, ze kterých je možné hodnoty vyškrtávat (dynam. datová struktura)
 - o Při ohodnocení prom. kontrola souv. podmínek, případně vyškrtáváme hodnoty z pracovních domén proměnné zasažené podmínkou
- **Backtracking s filtrací** – přiřazení opatříme dopřednou kontrolou (inference)
- Důraznější filtrace – nebudeme čekat na úplné ohodnocení proměnné v podmínkách
 - o Propagace podmínek
- **Hranová konzistence** – obecnější druh filtrace pro CSP s binárními podmínkami
 - o Výpočet: opakovaně vynucujeme konzistenci na hranách (dokud se mění prac. domény)
 - o předpokládáme podmínku $c \in C$ nad proměnnými x_1 a x_2 , D = doména
 1. kontrolujeme orientovanou hranu (x_1, x_2)
 2. pro každou hodnotu $d_1 \in D(x_1)$ hledáme **podporu** $d_2 \in D(x_2)$ vzhledem k podm. c
 3. podpora – $d_2 \in D(x_2)$ taková, že $(d_1, d_2) \in c$
 4. nemá-li d_1 podporu v $D(x_2)$, odebereme d_1 z $D(x_1)$
 - o **Podpora** pro a z domény x : hodnota b je podpora v doméně y , jestliže (x, y) splňuje podmínku (a nemá podporu, jestliže je nemožné, aby x nabývala hodnoty a , aniž by porušila podmínku)
- **Heuristiky pro výběr proměnné:**
 - o Výběr nejvíce omezené proměnné
 - o Výběr „klíčové“ proměnné
 - o Výběr tak, abychom co nejdříve narazili na řešení
 - o Využití struktury – míra souvislosti grafu CSP je určující
 - Komponenty souvislosti lze řešit nezávisle
 - Nejjednodušší souvislý graf je strom
- **Zobecněná hranová konzistence (GAC)** – zobecněná podpora
 - o Podmínky jsou nebinární, pro jednu zvolenou proměnnou v podmínce je podpora tvořena všemi ostatními proměnnými
 - o Globální podmínky
 - **allDifferent**($x_1 \dots x_k$) – obvykle uvažována jako jedna podmínka nahrazující všechny dvojice nerovností
 - Modelujeme jako hledání párování – konzistentní ohodnocení odpovídají párování
 - o **Propagace** v allDifferent
 - Najdeme nějaké ohodnocení proměnných, které splňuje allDifferent, následně se odeberou zbytečné hrany a dle toho se upraví pracovní domény proměnných
 - Najdeme maximální párování (maximální tok) – odstraníme hrany, co nepatří do žádného max. párování
 - Vyřazení neoznačených promítneme zpět do pracovních domén

3. Systematické a lokální splňování v logice (DPLL, CDCL, WalkSAT, posílání zpráv). Automatické uvažování, rozhodování v teoriích prvního řádu, obecná rezoluce, princip SAT-modulovaných teorií (SMT). Zpracování přirozeného jazyka.

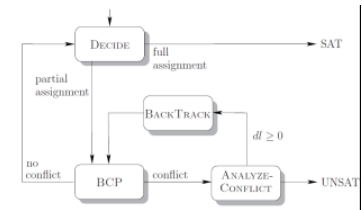
NI-UMI

Systematické a lokální splňování v logice

- ϕ je splnitelná, jestliže existuje α , že $\alpha \models \phi$
 - o NP-úplný problém
 - o $\alpha \models \phi$: „ ϕ je splněná pro ohodnocení α “
- **CNF** – konjunktivní normální tvar – podobá se CSP
- **Backbone** – podmnožina lit. Φ , které jsou pravdivé v každém splňujícím ohodnocení ϕ
- Backdoor – podmnožina proměnných ϕ , po jejichž ohodnocení je vyřešení splnitelnosti pro ϕ po dosazení v polynomiálním čase
- Autarky (autarkie) – částečné ohodnocení, které splní každou klauzuli, ve které nastavuje nějaký literál
- Unsatisfiable core – je-li ϕ nespíitelná, podmnožina jejích klauzulí, která je stále nespíitelná
- **Cejtinova transformace** – postupujeme podle derivačního stromu
 - o Pro každý vnitřní uzel derivačního stromu formule ϕ zavedeme **pomocnou proměnnou** a_i
 - Indikátorem splnění pod-formule odp. stromu
 - Listy = původní proměnné
- **Jednotková propagace** – pokud klauzule není splněná a všechny až na jednu proměnnou v ní jsou ohodnocené, poslední proměnná musí splnit danou klauzuli
- **Jednoduché splňování – DPLL**
 - o Backtracking s jednotkovou propagací čistých proměnných
 - **Čistá proměnná** – buď pouze neg., nebo pouze poz., lze ihned ohodnotit
 - o **Důvody konfliktu**
 - Předchůdcovská klauzule (Antecedent(x))
 - Jestliže byl literál l ohodnocen jednotkovou propagací kvůli klauzuli c – c je předchůdcovská pro l
 - **Implikační graf** pro ohodnocení α'
 - Zachycuje vznik α' , umožňuje vypátrat důvod nemožnosti α' dále rozšířit – příprava na skok zpět (backjumping)
 - Není třeba ho explicitně konstruovat (looking at you, CDCL)
 - Vrcholům je přiřazena informace o ohodnocení, K je speciální vrchol označující konflikt (nemožnost rozšířit dané částečné ohodnocení)
- **Učení klauzulí** – hranový řez v IG, který odděluje K a rozhodovací vrcholy, určuje konfliktní klauzuli c
 - o Lze využít k zpětnému skoku, lze si ji zapamatovat
 - o Spec. Příklad nogoodu
- **Splnitelnost s garancí** – DPLL, CDCL, ...
 - o Úplnost – vždy skončí a odpoví
 - o Mohou poskytnout formální zdůvodnění odpovědi (implikační graf) – Explainable AI
 - o Speciální třídy formulí – garantujeme rychlou odpověď
- **Trojice paradigmát** – CSP, SAT, IP (integer programming)
 - o CSP
 - + prohledávací algoritmy, silná formulace
 - Heterogenní omezení
 - o SAT
 - + homogenní, silné učení
 - absence aritmetiky
 - o IP
 - + dobrá aritmetika
 - slabé učení

- **CDCL** algoritmus (Conflict Driven Clause Learning) – využívá jednotkové propagace, staví implikační graf (NE EXPLICITNĚ)

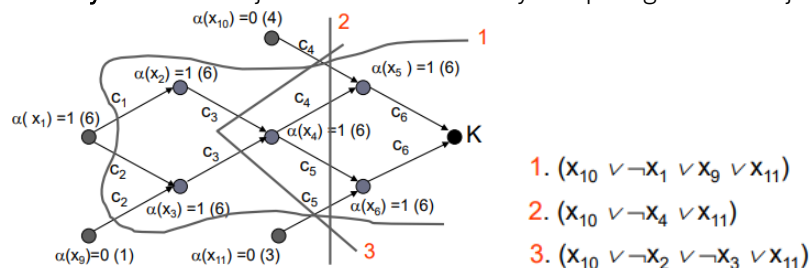
- o Základ systematických řešičů SATu – kombinuje BJ (skoky zpět), UP (jednot. propagaci) a učení
 - Analyzuje implikační graf pro nalezení vhodného konfliktu
- o Součástí:



- **DECIDE** – ohodnotí další neohodnocenou proměnnou
- **BCP** – jednotková propagace (boolean constraint propagation)
- **BackTrack(level)** – zruší všechna rozhodnutí na úrovních vyšších než level
- o Když dojde ke konfliktu, nalezne vhodný řez grafu, z něj vytvoří nogood, ten přidá jako klauzuli
 - **Vhodný řez** – z poslední vrstvy obsahuje pouze prom., u které se rozhodovalo ohodnocení
 - Přidaná klauzule následně vynutí změnu ohodnocení u poslední rozhodované proměnné pomocí UP
 - Dělá **backjumping** – protože když se dělá řez, skočí se do předposlední rozhodovací vrstvy, která obsahuje proměnnou v udělaném řezu – může přeskočit i nějakou vrstvu, která není na hraně řezu

- o **Analýza konfliktu**

- **Hranový řez** – odděluje K a rozhodovací uzly v implik. grafu – určuje konfl. klauzuli (nogood)



- Obecné požadavky – krátká konfliktní klauzule (kvůli UP, snadno uložitelná, obecnější)
- Budeme chtít **VYNUCUIJÍCÍ konfliktní klauzuli**
 - Obsahuje jediný literál z aktuální rozhodovací úrovně
 - Po návratu (BackTrack), je-li vynucující klauzule naučena (přidána), dojde skrz ni okamžitě k jednotkové propagaci
 - Nutnost, aby CDCL fungoval – v. k. k. zajistí větvení
- o **Unikátní implikační bod – UIP**
 - Def. Vzhledem k aktuální rozhodovací úrovni
 - (i) Vrchol implikačního grafu různý od K
 - (ii) Všechny cesty z rozhodovacího vrcholu do K jím procházejí
 - Vždy existuje, může jich být víc – zajímá nás 1. – nejbíže K

- o **Volba konfliktní klauzule**

- Vynucení a UIP – obsahuje první UIP jako jediný literál z aktuální rozhodovací úrovně
 - Klauzule bude vynucující a bude krátká
- Následuje návrat
 - Uložení KK (učení)
 - Návrat na 2. nejvyšší rozhodovací úroveň z KK
 - Singularity

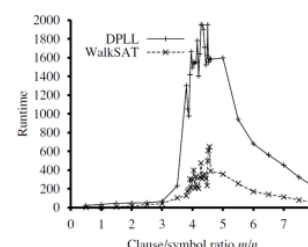
- o Vlastnosti CDCL: Vždy skončí, KK lze zapomínat, krátké klauzule, systematický

- o Implikační graf není třeba explicitně konstruovat

- o Heuristiky pro výběr

- VSIDS – každý literál má své skóre, literál s nejvyšším skóre ohodnocujeme True
- Berkmin – proměnné a literály mají VSIDS skóre
 - Konfliktní klauzule na zásobník
 - Nerozhodnutá klauzule se bere ze zásobníku
 - Prázdný zásobník

- Pomocí restartů se může učit z více stran zároveň
- Jednotková propagace efektivně – rychlé určení jednotkové klauzule
 - 2 sledované literály
 - V každé klauzuli označíme 2 neohodnocené literály (sledované)
 - U každé proměnné známe její sledované výskyty (ve kterých klauz. je sledován její literál)
- SAT řešen lokálně – **GSAT**
 - Hladové procházení úplných ohodnocení a restarty
 - Neúplný – nemůže-li najít splňující ohodnocení, nevíme, jestli je vstupní formule splnitelná, či ne
 - Najde takovou proměnnou, jejímž překlopením co nejvíce zvýším počet splněných klauzulí
 - Účelová funkce – počet splněných klauzulí (náchylný k uvážnutí v lokálním maximu)
 - **Restart** = solver vezme zpět každé přiřazení na cestě a udělá novou sérii rozhodnutí – zahodí podstrom, ale neodnaučí se naučené
- Náhodná procházka – **WalkSAT**
 - Procházení úplných ohodnocení s restarty
 - Z lokálních maxim se může dostat náhodnými kroky
 - Střídání hladového kroku (s pravděpodobností p) a náhodného kroku (pst 1-p)
- **Fázový přechod** – poměr mezi počtem klauzulí a počtem proměnných
 - Náhodné k-CNF formule ($k=3, 4, 5, \dots$)
 - Snadno splnitelné formule – málo klauzulí (málo omezení)
 - Snadno nesplnitelné formule – mnoho klauzulí (příliš omezené)
 - Ve fázovém přechodu jsou nejtěžší formule někde uprostřed
 - Limitní splnitelnost pro určité hodnoty: $c = \text{počet klauzulí} / \text{počet proměnných}$
- **Algoritmy posílání zpráv**
 - Na jednotlivé klauzule formule pohlíží jako na paralelní funkční jednotky – zkoumá, jak změna hodnoty proměnné ovlivní související klauzule
 - Faktorový graf (graf funkčních jednotek) – uzly jsou klauzule a proměnné, hrany jsou že proměnná je v klauzuli
 - **Propagace výstrah** – klauzule posílají proměnným zprávy (výstrahy), jak si přejí proměnnou nastavit, aby klauzule byla splněna
 - Aktualizační pravidlo pro výstrahy
 - **Výstrahami inspirovaná decimace** – postupně zjednodušuje formuli dosazováním
 - Preferovaná hodnota
 - Indikátor sporu
- Systematický CDCL se inspiroval u lokálních
 - **Restarty** v CDCL
 - **Náhodný krok** – s malou pst ignoruje heuristiku a vybere náhodnou hodnotu
- Lokální se inspirovaly u CDCL – 1 sledovaný literál
- **Inkrementální splnitelnost** – konflikt pro $\phi(t)$ je konfliktem i pro $\phi(t+1)$



Automatické uvažování

- Automatické dokazování vět, uvažování s neurčitostí, ...
- ATP – algoritmy ukazující, že zadané matematické tvrzení je logickým důsledkem daných axiomů
- APC – interaktivní verze ATP
- **Formalizace uvažování** – pomocí logiky prvního řádu (first order logic – FOL)
 - Silný vyjadřovací prostředek
 - Platíme nerozhodnutelností – nemůže existovat algoritmus, který by rozhodoval o platnosti zadaného tvrzení v logice 1. řádu
- **Rozhodování** – pro každý vstup skončit a dát správnou odpověď, zda tvrzení platí, či ne
 - Pro logiku prvního řádu máme jen semi-rozhodnutelnost
 - Algoritmus skončí, jestliže tvrzení platí a dokud neskončí, nic nevíme
 - Nebo rozhodnutelnost ve speciálních případech (fragmenty)

- **Nerozhodnutelnost**
 - o Příčiny nerozhodnutelnosti
 - Logika 1. řádu je tak silná, že v ní lze popisovat algoritmy
 - Programy a logika se v důsledku toho pak dokáží oklamat
 - o Algoritmy neumí rozhodovat tvrzení o algoritmech → neumí rozhodovat o tvrzeních obecně
- **Logika prvního řádu**
 - o Jazyk:
 - Proměnné pro individua – x, y, z, \dots
 - Spojky – unární, binární (negace, konjunkce, ...)
 - Kvantifikátory – všeobecný, existenční
 - Pomocné symboly – závorky, tečka
 - o Signatura (nelogické symboly):
 - Symboly pro funkce (transformace individuí) – $f, g, h, +, * \dots$ libovolné arity (nulární = konstanty)
 - Symboly pro predikáty (vlastnosti individuí) – $R, S, <, = \dots$ libovolné arity (nulární obvykle ne)
- **Formule 1. řádu**
 - o **Termy** – proměnná je term
 - Jestliže f je funkční symbol arity n a $t_1 \dots t_n$ jsou termy, pak $f(t_1 \dots t_n)$ je term
 - o **Atomy** – jestliže p je predikátový symbol arity n a $t_1 \dots t_n$ jsou termy, pak $p(t_1 \dots t_n)$ je atomická formule (atom)
 - o Formule je slovo (konečná posloupnost symbolů)
 - Atomická formule je formule
 - Jestliže p a q jsou formule, pak $(\neg p)$, $(p \wedge q)$, $(p \vee q)$, $(p \Rightarrow q)$, $(p \Leftrightarrow q)$ jsou formule
 - Jestliže p je formule a x proměnná, $\forall x(p(x))$ a $\exists x(p(x))$ jsou formule
 - Každá formule vznikne konečným počtem aplikací těchto pravidel
- **Pojmy v logice 1. řádu**
 - o **Teorie** – libovolná množina formulí
 - o **Důkaz formule** ϕ z teorie T – posloupnost formulí končících ϕ , každá f . patří do T nebo je z nějakých předchozích f . v posloupnosti odvozena odvozovacími pravidly
 - o **Pravdivost** (validita) – $T \models \phi$ (ϕ je pravdivá (validní, platná) v T), model pro T je modelem pro ϕ
 - Model – struktura, která splňuje všechny formule teorie
 - Sporná teorie nemá model
 - o **Struktura** (interpretace, realizace) signatury
 - Nosná množina M (prvky jsou individua)
 - Realizace funkcí nad M : $f: M \times M \times \dots \times M \rightarrow M$
 - Realizace predikátů nad M : $p \subseteq M \times M \times \dots \times M$
- **Semi-rozhodování: rezoluce**
 - o Rezoluční metoda – klauzální zamítací dokazování
 - Předp. Dokazování tvrzení C vzhledem k množině axiomů T
 - Uvážíme $T \cup \{\neg C\}$, snažíme se odvodit spor
 - o **Rezoluční pravidlo**
 - speciální výroková varianta (řez) – x, y výrokové literály, z výroková proměnná

$$\frac{(x_1 \vee x_2 \vee \dots \vee x_n \vee z), (y_1 \vee y_2 \vee \dots \vee y_m \vee \neg z)}{(x_1 \vee x_2 \vee \dots \vee x_n \vee y_1 \vee y_2 \vee \dots \vee y_m)}$$
 - Obecná varianta pro logiku 1. řádu

$$\frac{(p_1 \vee p_2 \vee \dots \vee p_n \vee r), (q_1 \vee q_2 \vee \dots \vee q_m \vee \neg s)}{(p_1 \vee p_2 \vee \dots \vee p_n \vee q_1 \vee q_2 \vee \dots \vee q_m) \theta}$$

 - p, q literály, r, s atomy
 - θ – nejobecnější unifikace r a s (tedy $r\theta = s\theta$ a každá jiná θ' , že $r\theta' = s\theta'$ lze vyjádřit pomocí θ , tedy $\theta' = \theta\rho$)

- Příprava na rozhodování
 1. převod T a C na CNF
 - standardizace proměnných – přejmenovávání
 - skolemizace – odstraňování existenčních kvantifikátorů
 - zahození všeobecných kvantifikátorů
 - distribuce skrz \wedge a \vee
 2. eliminace spojek jiných, než \neg , \wedge , \vee
 3. negace dovnitř až k atomům
- **Unifikační algoritmus** – hledá nejobecnější unifikaci (mgu) θ množiny termů x a y – chceme $x\theta = y\theta$
- **ANL loop** – pomocí DFS či BFS výběr klauzule c , vrací splnitelnost
- Holíčův paradox – „neexistuje holič, který by holil právě ty, kteří se neholí sami“
 - o ANL s rezolučním pravidlem nenajde důkaz
 - ➔ Pravidlo faktorizace – eliminace unifikovaného literálu \rightarrow zúplnění rezoluce pro logiku 1. řádu
- **Resoluce s faktorizací** – zdánlivě úplný algoritmus pro rozhodování o platnosti tvrzení v logice 1. řádu, problém teorie s nekonečně mnoho axiomů
- **Dokazování ve fragmentech** – omezení na podmnožinu jazyka logiky 1. řádu – **fragment**
 - o Často bez kvantifikátorů, speciální sada axiomů
 - o Speciální algoritmy na dokazování ve fragmentech
 - Poskytují úplnost rozhodovacího procesu
 - Převod otázky $T \models \phi$ na splnitelnost
- **Dokazování s rovností**
 - o Logika s rovností (fragment s rovností) – z predikátů pouze $=$, bez funkčních symbolů a kvantifikátorů
 - o Zjednodušení – jen spojky \neg , \wedge , \vee a převod na NNF (negační normální tvar)
 - o V NNF lze hovořit o literálech
- **Graf rovností** – def. Pro formuli ϕ s rovnostmi v NNF jako $G_{=} = (\text{Var}(\phi), E_{=}, E_{\neq})$
 - o Vrcholy – proměnné, hrany – $E_{=}$ rovnosti a E_{\neq} nerovnosti
 - o Sporná kružnice – cyklus v $G_{=}$, který obsahuje právě jednu nerovnost
- **Zjednodušovací algoritmus** – na vstupu formule ϕ_E s rovnostmi ve tvaru NNF, konstrukce grafu rovností
 - o Výstup je kratší ekvivalentně splnitelná formule
- **Výroková kostra** – zachycuje Booleovskou strukturu formule s rovnostmi
 - o Pro každou rovnost (atom) zavedeme výrokovou proměnnou, v pův. formuli rovnosti nahradíme zavedenými výrokovými proměnnými
 - o Nezachycuje tranzitivitu rovnosti
- **Tranzitivita rovnosti**
 - o Nepolární graf rovností $G_{=}^{\text{NP}} = (\text{Var}(\phi), E_{=} \cup E_{\neq})$ – zapomeneme na polaritu
 - Pro každý cyklus a každou hranu v cyklu přidáme tranzitivní podmínku – zakáže hraně přiřadit False, když jsou všechny ostatní hrany True
- **Líný přístup**
 - o Komponenty – decide_T (dokazovací procedura pro konjunktivní fragment teorie T), řešič SATu
 - o Spolupráce decide_T a řešiče SATu
 - Řešič vybírá, které literály splnit
 - decide_T kontroluje, jestli vybrané literály mohou současně platit (ne \rightarrow zakázané ohodnocení)
- SAT modulované teorie (**SMT**) – spolupráci SATu a DECIDE lze zabudovat do CDCL
 - o Propagace teorií přímo do řešiče SATu
 - o Na formuli t se kladou další podmínky, aby byla zaručena konečnost algoritmu
- Další metody – přirozená dedukce, metoda tablo, rozhodování s bitovými vektory, ukazateli, poli

Zpracování přirozeného jazyka

- **Formální jazyk** – množina řetězců, řetězec posloupnost terminálů
- **Přirozený jazyk** – nemá pevnou definici (přístup k nim jako by byly formální)
- **Gramatika** – sada pravidel specifikující jazyk
 - o $G = (N, T, S, P)$ – terminály, neterminály, počáteční neterminál, konečná množina pravidel
 - o **Chomského hierarchie** gramatik a jazyků – rekurzivně spočetné, kontextové, bezk., regulární
 - o **Bezkontextová gramatika** – pravidla $A \rightarrow \beta$, A neterminál a β řetězec t. a net.
- **Fáze komunikace**
 - o Záměr – sdělit skutečnost P
 - o Generování – přesvědčivého sdělení o skutečnosti $P \rightarrow$ věta
 - o Syntéza – sdělení, hlasový syntetizátor \rightarrow zvuk
 - o Vnímání – rozpoznání hlasu \rightarrow věta
 - o Analýza – věta \rightarrow derivační strom
 - o Odstranění dvojsmyslů – určení správné interpretace, pravděpodobnostně
 - o Začlenění – budeme věřit či ne
- **Analýza** (parsing)
 - o Pravděpodobnostní derivační strom – vnitřním uzlům přiřazena pravděpodobnost
 - o Algoritmus CYK – předpokládá bezkontextovou (pravděpodobnostní) gramatiku v Chomského normální formě (pravidla $X \rightarrow YZ$ nebo $X \rightarrow w$)
- **Učení PCFG**
 - o Učení pravděpodobností pro PCFG – podle banky stromů (treebank), korpus korektně vytvořených derivačních stromů vzhledem k dané bezkontextové gramatice
- **Rozšířené gramatiky**
 - o **Lexikalizované PCFG**
 - Pravděpodobnost pravidel závisí na vztahu slov v derivačním stromu, ne jen na jejich sousedství ve větě
 - Vztah všech dvojic je komplikované počítat, vybereme nejvýznamější slovo z daného spojení – hlava (head)
 - Stále mnoho dvojic, kro které potřebujeme pravděpodobnost
 - o **DCG** – definite-clause grammars (gramatiky s určenou klauzulí)
 - Používáme navíc **{podmínku}** – pravidlo platí, pokud je podmínka splněná
 - Problém – skloňování/pády, shoda přísudku s podmětem
- **Sémantika** – výrazy – aritmetické výrazy, lambda notace
 - o Kvantifikace – používá se přechodná forma mezi derivačním stromem a logickou interpretací
 - **Kvazi-logická forma** – akceptace ústupků od přísných pravidel ve prospěch zachycení kontextu, modelování významových relací mezi větami
 - o Další jevy
 - Časy – kalkulus událostí (event calculus)
 - Nejednocnačnost – lexikální, syntaktická (více derivací), doslovné x obrazné
- **Statistický překlad** – překládáme větu e z angličtiny do francouzštiny f
 - o $P(e|f)$ – **překladový model** – jak pravděpodobné je, že je e anglickým překladem dané francouzské věty f
 - o $P(f|e)$ – jak pravděpodobné je, že f je francouzským překladem dané anglické věty e
 - o $P(f)$ – **jazykový model** – jak pravděpodobné je, že f je věta ve francouzštině
 - o Hledáme řetězec f^*
$$f^* = \operatorname{argmax} P(f|e) = \operatorname{argmax} P(e|f)P(f)$$
 - o Lze počítat podle $P(e|f)P(f)$ – máme-li dobrý jazykový model pro $P(f)$, nebo přímo $P(f|e)$
 - o Průběh statistického překladu
 - Rozdělíme vstupní anglickou větu e na slovní spojení
 - Vyhledáme slovní spojení v anglické části korpusu
 - Najdeme odpovídající překlad slovních spojení ve francouzské části korpusu
 - Sestavíme výsledný překlad z nalezených francouzských překladů slovních spojení

4. Metody pro hodnocení a výběr příznaků (univariетní/multivariетní metody, filtrační/wrapper/embedded metody). Selektivní/adaptivní metody redukce počtu instancí: Condensed Nearest Neighbor (CNN), Delauney/Gabriel/RNG grafy, Wilsonova editace, Multi-edit metoda, Tomkovy spoje.

NI-PDD

Univariетní / multivariетní metody

- Univariетní – zvažuje každý příznak zvlášť

- o Jak relevantní je proměnná x_i pro predikci výsledku proměnné Y ?
- o Irelevantní příznak X_i : $P(X_i | Y = 1) = P(X_i | Y = -1)$
- o **T-test** – normálně rozdělené třídy, stejný rozptyl σ^2 neznámý, odhad z dat jako σ_{within}^2

- Nulová hypotéza: $\mu^+ = \mu^-$

- Pokud je nul. hyp. pravdivá: $t = (\mu^+ - \mu^-) / (\sigma_{within} \sqrt{1/m^+ + 1/m^-}) \sim \text{Student}(m^+ + m^- - 2 \text{ d.f.})$

- o **Korelace** – vzájemný lineární vztah mezi proměnnými

- +1 přímá závislost, -1 nepřímá závislost

- **Pearsonův korelační koeficient**

$$\rho_{X,Y} = \frac{\text{cov}(X,Y)}{\sigma_X \sigma_Y} = \frac{E\{(X - \mu_X)(Y - \mu_Y)\}}{\sigma_X \sigma_Y}$$

- **Korelační a kovariační matice**

$$\rho = \begin{bmatrix} 1 & \rho(X_1, X_2) & \dots & \rho(X_1, X_n) \\ \rho(X_2, X_1) & 1 & & \rho(X_2, X_n) \\ \vdots & & \ddots & \vdots \\ \rho(X_n, X_1) & \rho(X_n, X_2) & \dots & 1 \end{bmatrix} \quad \begin{matrix} \text{Variance } (\sigma^2) \\ \downarrow \\ C = \begin{bmatrix} D(X_1) & C(X_1, X_2) & \dots & C(X_1, X_n) \\ C(X_2, X_1) & D(X_2) & & C(X_2, X_n) \\ \vdots & & \ddots & \vdots \\ C(X_n, X_1) & C(X_n, X_2) & \dots & D(X_n) \end{bmatrix} \end{matrix}$$

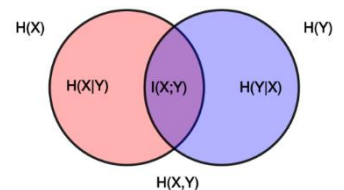
- **Spearmanův koeficient pořadové korelace**

- Máme 2 náhodné proměnné X a Y s neznámými rozděleními
 - Seřadíme n hodnot x_i a y_i podle velikosti a přiřadíme jim pořadová čísla p_i a q_i .
Hodnota koeficientu je:

$$\rho = 1 - \frac{6 \sum_i (p_i - q_i)^2}{n(n^2 - 1)}$$

- o **Entropie** – předpokládané množství informace v náhodné proměnné

$$H = - \sum_i P_i \log P_i$$



- o **Univariетní závislost**

- **Nezávislost**: $P(X,Y) = P(X)P(Y)$

- **Vzájemná informace** – hodnota vzájemné závislosti mezi dvěma náhodnými proměnnými

- Kvantifikuje množství informace získané o jedné proměnné zkoumáním jiné náhodné proměnné

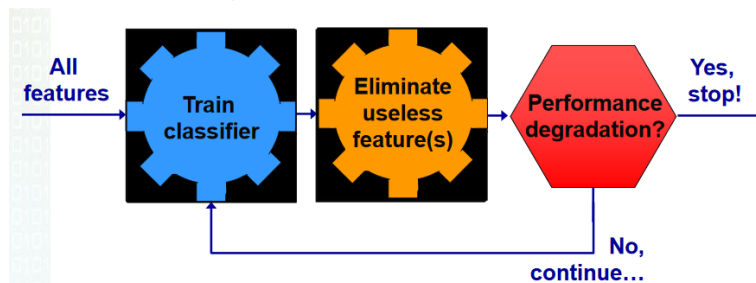
$$MI(X, Y) = \int P(X,Y) \log \frac{P(X,Y)}{P(X)P(Y)} dX dY$$

- Empirická vzájemná informace – z dat kontingenční tabulka, z ní empirická pravděpodobnost, ta je vstupem pro vzájemnou pravděpodobnost

- Multivariетní – zvažuje podmnožiny příznaků najednou

Filtrační / wrapper / embedded metody

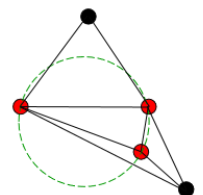
- **Filter metody** – hodnotí příznaky nebo jejich podmnožiny nezávisle na klasifikátoru
 - o *Kritéria* – hodnocení relevance příznaků
 - o *Hledání* – seřazení příznaků
 - o *Hodnocení* – statistické testy
 - o Robustní výsledky, ale nemusí vybrat úplně nejlepší příznaky
- **Wrapper metody** – používají klasifikátor k hodnocení příznaků / podmnožin
 - o Rozdělení dat na tréninkové, validační a testovací
 - o Problém přeučení, ale najde nejlepší příznaky
 - o *Kritéria* – měření užitečnosti podmnožiny
 - o *Hledání* – prohledávání prostoru podmnožin příznaků
 - o *Hodnocení* – **Křížová validace** – rozdělení do k skupin
 - Trénink na všech skupinách kromě 1
 - Výběr těch příznaků s nejlepším skóre podle poslední skupiny
 - Opakování, dokud nejsou všechny skupiny použity pro testování
 - Zprůměrovat
- **Embedded metody** – jako wrapper, hledání ovládáno algoritmem konstruujícím klasifikátor
 - o *Kritéria* – měření užitečnosti podmnožin



- o *Hledání* – naváděno učícím procesem
- o *Hodnocení* – křížová validace
- o Výsledky podobné wrapperům, ale menší riziko přeučení a míň výpočetně náročné

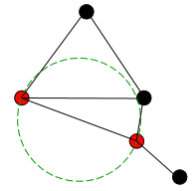
Selektivní/adaptivní metody redukce počtu instancí

- **Condensed Nearest Neighbour (CNN)** – inkrementální, závislá na pořadí, ani minimální ani decision boundary konzistentní
 - o $O(n^3)$ pro brute-force metodu
 - o Lze navázat na reduced NN
 - o Postup:
 - Inicializace podmnožiny jedním trénovacím případem
 - Klasifikace většiny vzorků pomocí té podmnožiny a přesun všech špatně klasifikovaných vzorků do podmnožiny
 - To se opakuje, dokud buď nedochází k přesunu nebo je podmnožina plná
- **Proximity grafy** – poskytují různé definice „sousedů“
 - o Grafická reprezentace dat – zobrazují, jak jsou prvky propojeny na základě vzdáleností/podobnosti
 - o **Delaunay**
 - **Delaunayho triangulace** – duální graf Voronoiho diagramu?
 - Voronoi diagram – rozděluje rovinu na oblasti na základě vzdálenosti k bodům (zdroje) – každá oblast obsahuje body, které jsou blíže k danému zdroji než k jakémukoliv jinému
 - Vrcholy Voronoi diagramu jsou rohy trojúhelníků v D.T.
 - Tři body jsou si sousedy, pokud kružnice jim opsaná neobsahuje žádné jiné body
 - Voronoi editing: ponechá body, jejichž sousedi jsou z jiné třídy
 - Rozhodovací hranice je stejná
 - Konzervativní podmnožina, nechává si body navíc, drahá na výpočet ve vysokých dimenzích



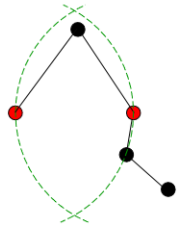
- **Gabriel**

- Podmnožina delaunayho triangulace
- Body jsou sousedi, pouze pokud je jejich sféra vlivu prázdná
- Nezachovává identickou rozhodovací hranici, ale změny se odehrávají jen vně konvexního obalu datových bodů
- Spoje mezi sousedy = Tomek links



- **RNG** – Relative Neighbourhood Graph

- Podmnožina Gabrielova gradu
- Dva body jsou sousedi, pokud „luna“ definovaná průsečíkem jejich radiálních sfér je prázdná
- Dále zmenšuje počet sousedů
- Rozhodovací hranice bývají drastické a nemusí být konzistentní s trénovací množinou



- **Editing** – vyjmutí šumových bodů a tvorba hladkých rozhodovacích hranic

- Často se zachovávají body daleko od hranice a vzniknou homogenní shluky bodů
- **Wilson editing** – vyjmutí bodů, které mají jinou třídu než většina jejich k nejbližších sousedů
- **Multi-edit** – opakovaně opakuje Wilsonův editing na náhodné části, klasifikuje pomocí 1-NN pravidla
 - **Difúze** – rozdělení dat do $N \geq 3$ náhodných podmnožin
 - **Klasifikace** – klasifikace S_i pomocí 1-NN za použití $S_{(i+1) \bmod N}$ jako trénovací set
 - **Editing** – zrušení všech špatně klasifikovaných vzorků
 - 1-NN pravidlo – klasifikace X podle jeho nejbližšího souseda z trénovacích bodů
 - **Konfúze** – přidání všech vzorků do nové množiny
 - **Terminace** – pokud posledních I iterací nemělo za následek žádný editing, konec, jinak od začátku
- Kombinace editingu a condensingu – první editing na odstranění šumu a hladkou hranici, potom kondenzace pro získání menší podmnožiny

- **Tomkovy spoje** – odstranění šumu a hraničních případů

- Hrany, které spojují vrcholy, které sdílejí společný okraj v diagramu
- Tomkův spoj
 - E_i, E_j patří do jiných tříd, $d(E_i, E_j)$ je vzdálenost mezi nimi
 - Pár (E_i, E_j) je Tomkův spoj, pokud neexistuje takové E_l , že $d(E_i, E_l) < d(E_i, E_j)$ nebo $d(E_j, E_l) < d(E_i, E_j)$

5. Algoritmy pro nahrazování chybějících hodnot. Detekce a ošetření odlehlých hodnot. Vyvažování skupin dat (undersampling/oversampling metody).

NI-PDD

Nahrazování chybějících hodnot

- Nelze rozlišit mezi chybějící hodnotou a prázdnou buňkou
- **Způsoby nahrazení:**
 - o **Nic nedělán** – použije se relevantní reprezentace tak, aby šel aplikovat model
 - Např. použití -1, když jsou všechny ostatní hodnoty kladné
 - o **Vymazání**
 - **Listwise** – vymaže všechny data pro případ, který á jednu či více chybějících hodnot
 - **Pairwise** – listwise ale až v pozdějším stádiu, kdy jsou určeny důležité příznaky – zaměření na opravdu potřebná data
 - **Příznaky s chybějícími daty** – odstraní celý příznak, pokud chybí většina dat
 - Užitečné jen pokud málo záznamů obsahuje chybějící data, nebo příznak obsahuje většinu chybějících dat
 - o **Imputace** – proces nahrazení chybějících dat substituovanými hodnotami
 - Neexistuje ultimátně nejlepší způsob
- **Mechanismy chybějících dat**
 - o **MCAR** = *missing completely at random* – pokud je pravděpodobnost, že příznak chybí nezávislá na hodnotě příznaku a hodnotách jiných příznaků
 - Chybějící hodnoty jsou nezávislé na datech
 - o **MAR** = *missing at random* – pokud je pravděpodobnost, že příznak chybí nezávislá na hodnotě příznaku, ale může záviset na hodnotách dalších příznaků
 - *Na otázku o platu častěji odpoví muži než ženy*
 - o **MNAR** = *missing not at random* – pokud je pravděpodobnost, že příznak chybí závislá na hodnotě příznaku
 - *Chybí údaj o platu, protože na otázku člověk neodpověděl kvůli jeho výši*
- **Imputace**
 - o **Průměr/medián** pomocí přítomných hodnot daného příznaku
 - o **k nejbližších sousedů** – vhodné pouze pro datasety nízkých dimenzí kvůli prokletí dimenzionality
 - Využití data encodingu – one-hot/dummy
 - o **Evaluace** – pomocí korelace

Odlehlé hodnoty

- **Odlehlé hodnoty** = anomálie = extrémní hodnoty, které se odchylují od ostatních pozorování
 - o Detekce anomálií – podvodné transakce, podezřelí cestující, zdravotní problémy
- **Detekce** – nesupervizované metody, hlavně clustering
 - o Díky absenci labelů – outlieři nejsou předem známi
 - o **Jeden příznak – rule of thumb** – outlieři jsou mimo interval ($Q1 - 1.5 * IQR$, $Q3 + 1.5 * IQR$)
 - $Q1$ = první kvartil
 - $Q3$ = třetí kvartil
 - IQR = inter quartile range = $Q3 - Q1$ (prostředních 50 %)
 - o **Více příznaků** – více-dimenzionální pohled
 - **Shlukovací přístupy** závislé na vzdálenostech (k-means, local outlier factor)
 - **Density-based** přístupy (SVM)
 - **Modelové** přístupy (neuronky – autoenkodéry, LSTM)
- Po identifikaci nutná analýza původu – Lidská chyba, Chyba měření, Experimentální chyba (extrakce dat), Samplingová chyba, Přírodní úkaz (novoty)
- **LOF = local outlier factor**
 - o Hledání anomálií měřením lokální odchylky daného datového bodu vzhledem k jeho sousedům

- Založeno na konceptu lokální hustoty, kde lokálnost je dána k nejbližšími sousedy, jejichž vzdálenost je použita na odhad hustoty
- Porovnáním lokální hustoty datového bodu s lokálními hustotami jeho sousedů lze identifikovat oblasti s podobnou hustotou
- Body, které mají podstatně nižší hustotu, než jejich sousedi jsou anomálie
- **LOF(A)** = průměrná lokální hustota sousedů A dělená lokální hustotou A (relativní)
 - Hodnota blízka 1 znamená, že je objekt srovnatelný se sousedy (není odchylka)
 - Hodnota pod 1 indikuje hustší oblast (inliner), hodnota hodně nad 1 indikují odchylky
- Výhody:
 - LOF umí identifikovat anomálie v části datasetu takové, které by v jiné části stejného datasetu nebyly považovány za anomálie
 - Dobrá geometrická intuice pro málo rozměrné vektorové prostory
 - Mnoho variant a odvětví použitelnosti
- Nevýhody:
 - Výsledné hodnoty se těžko interpretují
- Další metody detekce – **Isolation forest**, metody používající meta-learning (MetaOD)

Vyvažování skupin dat

- **Baseline metody**
 - Random over-sampling
 - Náhodná replikace příkladů z menšinových tříd
 - Zvyšuje pravděpodobnost overfittingu
 - Random under-sampling
 - Náhodná eliminace příkladů z většinových tříd
 - Může zničit potenciálně užitečná data, která by mohla být důležitá pro indukční proces
- **Under-sampling metody**
 - Tomek links
 - Odstranění šumu i hraničních případů
 - E_i, E_j patří do jiných tříd, $d(E_i, E_j)$ je vzdálenost mezi nimi
 - $A(E_i, E_j)$ je Tomkův spoj, pokud neexistuje takové E_l , že $d(E_i, E_l) < d(E_i, E_j)$ nebo $d(E_j, E_l) < d(E_i, E_j)$
 - **Condensed Nearest Neighbour Rule (CNN)** – výběr bodů blízko hranice mezi třídami
 - Algoritmus:
 - E je původní trénovací množina
 - E' obsahuje všechny pozitivní vzorky z S a jeden náhodně vybraný negativní vzorek
 - Klasifikace E s 1-NN pravidlem za použití vzorků z E'
 - Přesun všech špatně klasifikovaných vzorků z E do E'
 - Citlivý na šum – šumové vzorky se pravděpodobně špatně klasifikují a hodně jich bude přidáno do trénovací množiny
 - **One-Sided Selection (OSS)**
 - Tomek links + CNN (v tomto pořadí)
 - **CNN + Tomek links**
 - Nalezení Tomkových spojů je výpočetně náročné, bylo by levnější, kdyby se provádělo na redukovaném datasetu
 - **NCL = neighborhood cleaning rule**
 - Pro odstranění vzorků z většinové třídy
 - Více k čištění dat než redukcí
 - Algoritmus:
 - Nalezení 3 nejbližších sousedů pro každý vzorek E_i z trénovací množiny
 - Pokud E_i patří do většinové třídy a 3 nejbližší sousedi do menšinové, E_i se odstraní
 - Pokud E_i patří do menšinové třídy a 3 nejbližší sousedi do většinové, odstraní se ti sousedi

- **Over-sampling – SMOTE** = Synthetic Minority Over-Sampling Technique
 - Tvorba nových vzorků z menšinové třídy interpolací mezi několika vzorky z minoritní třídy které jsou si blízko
 - V prostoru příznaků spíš než v prostoru dat
 - Algoritmus:
 - Pro každý vzorek z minoritní třídy:
 - Vytvoř syntetické vzorky podél úseček spojujících nějaké/všechny z k nejbližších sousedů z menšinové třídy
 - V závislosti na množství over-samplingu se sousedi z k vybírají náhodně
 - **Generování syntetických vzorků:**
 - Vezme se rozdíl mezi vektorem příznaků, který se zvažuje a jeho nejbližší soused
 - Rozdíl se vynásobí náhodným číslem od 0 do 1
 - Toto se přičte k uvažovanému vektoru příznaků
- **Kombinace under-samplingu a over-samplingu**
 - **Smote + Tomek links**
 - Problém se Smote – může vyrobít umělou minoritní třídu moc hluboko v prostoru většinové třídy
 - Tomek links: čištění dat
 - Namísto odstranění pouze vzorků z většinové třídy které formují Tomek links, vzorky z obou tříd jsou odstraněny
 - **Smote + ENN**
 - ENN – extended nearest neighbour – odstraní vzorky, jejichž label se liší od třídy alespoň dvou z jeho tří nejbližších sousedů
 - ENN odstraní víc vzorků než Tomek links
 - ENN odstraňuje vzorky z obou tříd
 - Over-sampling bývá lepší než under-sampling, ale kombinace může být ještě výhodnější

6. Lineární projekce dat do prostoru o méně dimenzích: metoda hlavních komponent (PCA), lineární diskriminační analýza (LDA). Nelineární metody redukce dimensionalit (Sammonova projekce).

NI-PDD

- **Projekce dat** – nalezení transformace projektující data z n dimenzního prostoru do k dimenzního prostoru ($k < n$) při zachování nějaké formy informace (např. vzdálenost)
 - o Výhoda: vyloučení redundantní informace (korelace)
 - o Nevýhoda: fyzický význam nových atributů těžké interpretovat

Lineární projekce

- Jednoduché na výpočet:

$$Y = U \cdot X \quad (b_i = u_i^T a_i)$$

$\begin{matrix} \nearrow & \nwarrow & \nwarrow \\ k \times 1 & k \times n & n \times 1 \end{matrix} \quad (k \ll n)$

- **Náhodná projekce** (mapování) – původní d -dimenzní data se projektují do k -dimenzního podprostoru ($k \ll d$) za využití náhodné $k \times d$ matice R , jejíž sloupce mají jednotkové délky
 - o Využívá maticovou notaci kde $X_{d \times N}$ je původní množina N d -dimenzních pozorování,

$$X_{k \times N}^{RP} = R_{k \times d} X_{d \times N}$$

je projekce dat do nižšího k -dimenzního podprostoru

- o **Johnson-Lindenstrauss lemma**: pokud jsou body ve vektorovém prostoru projektovány na náhodně vybraný podprostor vhodné velké dimenze, pak jsou vzdálenosti mezi korespondujícími body (v originálním a novém prostoru) zhruba zachovány
- o Technicky vzato ta transformace není projekce, protože R není nutně ortogonální
- **Metoda hlavních komponent (PCA)** – cíl je snížit dimenzionalitu dat a zároveň **zachovat co nejvíce variance** v datasetu (koresponduje množství informace)

- o **Informační ztráta** – je implikována redukcí dimenzionalit

$$\min \| \underset{\text{original data}}{x} - \underset{\text{reconstructed data (after inverse transformation)}}{\hat{x}} \| \quad (\text{reconstruction error})$$

- PCA zachovává co nejvíce informace (variance) je možno:

- o „Nejlepší“ prostor nízké dimenze má **střed** v průměru vzorku a má **směry** rozhodnuté „nejlepšími“ vektory vlastních čísel kovarianční matice dat x
 - „Nejlepšími“ vektory vlastních čísel se myslí ty odpovídající **největším vlastním číslům** = **hlavní komponenty**
 - Protože je kovarianční matice reálná a symetrická, tyto vektory jsou ortogonální a formují množinu báze vektorů
- o Možný výpočet – prolnutí přímkou těžištěm, zkusit přímku rotovat a tím najít nejlepší úhel – PCA transformuje osy
- o Cíl je transformace, nalezení nového souřadného systému, který zachovává max. míru rozptýlení/informace
 - Hlavní komponenty = osy v cílovém prostoru
 - Pořadí hlavních komponent má význam
 - Hlavní komponenty jsou na sebe kolmé

- **Lineární diskriminační analýza (LDA)** = Fisherova projekce – snížení dimenzionalit za **zachování co nejvíce diskriminační informace** (chyba špatné klasifikace)

- o Nalezení směrů, ve kterých jsou třídy **nejlépe separovány**
- o Bere v úvahu rozptýlení uvnitř tříd, ale také rozptýlení mezi třídami (= **scatter**)
- o Notace:
 - C tříd
 - μ_i průměrový vektor třídy i , $i = 1, 2, \dots, C$

- M_i počet vzorků ve třídě i
- $M = \sum_{i=0}^C M_i$ celkový počet vzorků
- y_j matice velikostí $C \times M_i$ obsahující vzorky z jednotlivých tříd
- Scatter matice **uvnitř třídy**:

$$S_w = \sum_{i=1}^C \sum_{j=1}^{M_i} (y_j - \mu_i)(y_j - \mu_i)^T$$

- Scatter matice **mezi třídami**:

$$S_b = \sum_{i=1}^C (\mu_i - \mu)(\mu_i - \mu)^T$$

- S_b má hodnotu maximálně $C - 1$
- Každá pod-matice má hodnotu 1 nebo méně – vnější součin dvou vektorů
- $\mu = \frac{1}{C} \sum_{i=1}^C \mu_i$ průměr celého datasetu

- $y = U^T x$, U je projekční matice

- LDA počítá transformaci, která **maximalizuje scatter mezi třídami** a minimalizuje scatter uvnitř tříd:

$$\max \frac{|U^T S_b U|}{|U^T S_w U|} = \max \frac{|\tilde{S}_b|}{|\tilde{S}_w|} \rightarrow \text{products of eigenvalues}$$

- $\tilde{S}_b, \tilde{S}_w \dots$ scatter matice projektovaných dat y

- **LDA lineární transformace:**

- Řešení LDA je dáno vlastními vektory obecného problému vlastních vektorů:

$$S_B u_k = \lambda_k S_W u_k$$

- Lineární transformace je dána maticí U , jejíž sloupce jsou vlastní vektory problému výše

- Souřadnice v novém prostoru: $\begin{bmatrix} b_1 \\ b_2 \\ \dots \\ b_K \end{bmatrix} = \begin{bmatrix} u_1^T \\ u_2^T \\ \dots \\ u_K^T \end{bmatrix} (x - \mu) = U^T (x - \mu)$

- Jelikož S_b má hodnotu nejvýše $C - 1$, je maximální počet vlastních vektorů s nenulovými vlastními čísly $C - 1$ (takže maximální dimenze podprostoru je $C - 1$)

- Při malé trénovací množině je PCA výhodnější než LDA, LDA zas lepší na větších a reprezentativních datech

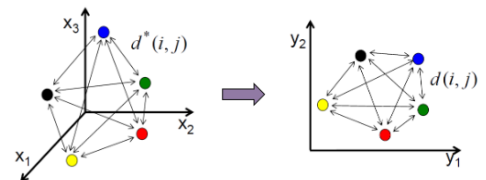
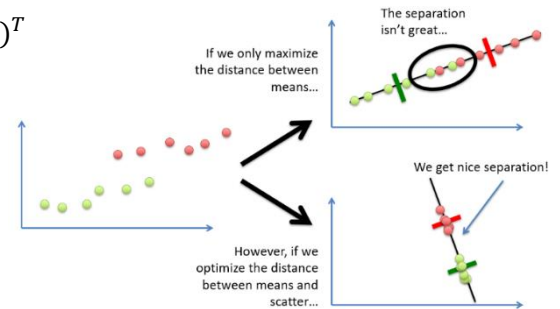
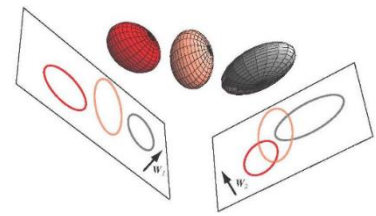
Nelineární redukce dimenzionality

- **Sammonova projekce** – mapuje vysoko-dimenzní prostor na prostor nižší dimenze při **zachování mezi-bodových vzdáleností** z vysoko-dimenzního prostoru v nízko-dimenzním prostoru

- Netransformuje souřadnice, místo toho reorganizuje pozice vzorů v novém prostoru
- Vzdálenost mezi i -tým a j -tým objektem v původním prostoru je d_{ij}^* , a vzdálenost jejich proj. je d_{ij}
- Sammonova projekce minimalizuje následující chybovou funkci – Sammon's stress neboli **Sammonova chyba**:

$$E = \frac{1}{\sum_{i < j} d_{ij}^*} \sum_{i < j} \frac{(d_{ij}^* - d_{ij})^2}{d_{ij}^*}$$

- Minimalizace se provádí např. gradientním sestupem



7. Učení dopředných neuronových sítí, konvoluční neuronové sítě a jejich regularizace.

NI-MVI

Dopředné neuronové sítě

- **Perceptron** – nejjednodušší neuronová síť, binární klasifikátor
 - o Mapuje vstupní vektor na výstup pomocí lineární kombinace
 - o 1 neuron, který vypočte lineární kombinaci vstupů a vah, přičte bias a tuto hodnotu vnitřního potenciálu prožene nelineární aktivační funkcí, jejíž výstup je zároveň výstupem celého perceptronu
 - o Aktivační funkce – často step function: 1 pro *vnitřní potenciál* ≥ 0 , jinak
 - o Perceptron nelze použít pro problémy, které nejsou lineárně separabilní (např. XOR).
- **Perceptronový algoritmus**
 1. Náhodně zvolený **vektor vah** kolmý na rovinu, která dělí 2 třídy (BINÁRNÍ klasifikátor)
 2. Opakovaně se vybírá náhodný prvek z trénovacích dat
 - a. Je z třídy P, ale hodnota wx je záporná – nový vektor vah w se vypočte přičtením vektoru x k vektoru vah w – posune se do pozitivní nadroviny určené vektorem w
 - b. Je z třídy N, ale wx je kladná – x se od w odečte, bod se posune do negativní poloroviny
 3. **Vektor vah je operacemi 2. posouván k optimálnímu rozdělení tříd**
 - o Prvotní nastavení vah – průměr všech kladných vektorů minus průměr všech záporných
- **Gradient learning** – konverguje i pro lineárně neseparabilní data a maximalizuje rozdělení tříd
 - o **Ztrátová funkce** (cost function) – např. **sum of squared errors**
 - o Funkce je minimalizována krokem proti směru gradientu – ovládán ještě parametrem learning rate
 - Potřeba parciální derivace ztrátové funkce vůči každé váze vektoru w
- **Cross-entropy loss (log loss)** – ztrátová funkce používaná pro klasifikaci c tříd
 - o Hodnota pro třídu Y a vektor p pravděpodobností, že bod na vstupu náleží do tříd:
$$L(Y, p) = -\log(p_y).$$
 - o Hodnota stoupá, čím jistější si je model klasifikací špatné třídy – nejvíc trestá predikce, kterými si je model jistý a jsou špatné
- **Multilayer perceptrony (MLP)** – více vrstev s více neurony, výstup každého neuronu vede do každého neuronu další vrstvy
 - o Nelineární aktivační funkce – **sigmoida, tanh, ReLU, ...** - musí být diferencovatelná
 - o I na lineárně neseparabilní data
 - o Vrstva MLP = funkce \rightarrow neuronová síť je složená funkce
 - o Pro použití gradientního sestupu pro učení vah musíme najít parciální derivaci Z.F. vůči každé váze – pomocí chain rule pro derivace složených funkcí
$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$
 - o **Backpropagation** – zpětné šíření chyby
 1. Pro všechna trén. Data spočítáme výstup sítě
 2. Vyčíslíme hodnotu cost funkce – např. průměr cross entropy
 3. **Postupně od poslední vrstvy ke vstupní** počítáme parc. Derivace vah vůči cost funkci \rightarrow gradient
 4. Gradientní sestup regulovaný learning rate
 5. Opakování do splnění terminálního kritéria
- **Momentum setrvačnosti** – Při updatu vah kromě aktuálního momentu přičítáme část gradientu z předchozího kroku
 - o Urychlení konvergence, únik lokálním optimům, zamezení přehnané reakce na vzorek
- **ADAM – Adaptive moment estimation**
 - o **Optimizátor** využívající exponentially decaying moving average gradientu akožto momentum
 - o Adaptivní learning rate parametry pro jednotlivé váhy – penalizace vah způsobujících velké oscilace – zamezení přeučení

- **Kohonen's self-organizing maps (SOM)** – typ neuronové sítě trénované nesupervizovaným učením
 - o Produkce méně dimenzionální (2D) diskretizované reprezentace vstupního prostoru = mapa -> pro **redukci dimenzionality**
 - o **Neighbourhood funkce** – snaha zachovat topologické vlastnosti prostoru
 1. Inicializace vah na malé hodnoty
 2. Vybrán vstup a neuron, který je mu nejbližší (např. euklidovsky, manhattansky)
 3. Neuron s jeho okolím jsou updatovány a posunuty směrem ke vstupu
 4. Opakování 2-3 a neurony a okolí updatovány a posouvány k datovým bodům, dokud nejsou rozprostřeny tak, že zachycují nějakou strukturu
 5. Váhy neuronů – pro popis dané oblasti



Konvoluční neuronové sítě

- Kromě fully connected vrstev konvoluční a pooling vrstvy
- **Fully connected vrstvy** – všechny neurony spojeny se všemi neurony v předchozí vrstvě
 - o Každý neuron má vlastní váhu
 - o Nepředpokládá nic o povaze dat, funguje zcela obecně, ale velmi výpočetně náročné
- **Konvoluční vrstvy** – neurony spojeny jen s hrstkou neuronů v předchozí vrstvě, které patří do nějakého okolí
 - o Stejná sada vah použita pro všechny neurony ve vrstvě
 - o Každý neuron si ze svého regionu (např. 3x3 grid) **extrahuje příznak** podle stejného předpisu (bo stejné váhy)
 - Dává smysl jen v případě, že data jsou prostorová a příznaky se vyskytují lokálně, a ne na jakékoliv pozici
 - o **Konvoluční neuron** – má **filtr** definovaný váhami, ten aplikuje konvoluci na svůj region
 - **Padding** – pro místa, kde filter zasahuje mimo rozsah konvoluce
 - **Stride** – pro délku kroku při přechodech konvoluce
 - Aktivace = vizuální stimulace
 - Extrakce nejdřív obecných příznaků – rohy, hrany, čáry, později textury atd
 - o **Pooling vrstvy** – Následují konvoluční vrstvy – **redukce dimenzionality**
 - Sjednocení extrahovaných příznaků (max pooling, average pooling, ...)
- **Regularizace**
 - o **Dropout** – během trénování **vypíná náhodné neurony** (= nastavuje váhy na 0) = regularizace
 - Sít se nenaučí spoléhat na malou sadu důležitých neuronů
 - o Další forma regularizace například L1 a L2, augmentace dat, batch normalizace (standardizace vstupů do vrstvy).
- Typy předurčených konvolučních sítí
 - o **AlexNet** – stackované konvoluční vrstvy
 - o **GoogLeNet** – později Inception sítě – inception moduly s méně parametry, místo poslední FC vrstvy average pooling
 - o **VGG** – velmi hluboká a velká síť
 - o **ResNet** – reziduální skip connections – pomoc s vanishing gradient problémem, zajištění information flow, batch normalization
 - o DenseNet
 - o MobileNet

8. Autoencodery a generativní neuronové sítě.

NI-MVI

Autoenkodéry

- Feedforward NN, **na výstupu se snaží napodobit vstup**, který je v průběhu dopředného chodu propuštěn malým bottleneckem
- Kódovací část do malé dimenze = **encoder**
- Část rekonstruuující vstup = **decoder**
- Vizualizace v menší dimenzi, učení abstraktních příznaků, komprese, rekonstrukce zašuměných či porušených dat
- **Stacked AE**
 - o Řada sparse autoencodérů jdoucích po sobě – každý bere na vstupu výstup předchozího
 - o U posledního encodéru odstranění decodéru – výstup jde do klasifikátoru/prediktoru
 - Poslední vrstva trénovaná supervizovaně na nějakou úlohu
 - o Každá vrstva naučena na greedy extrakci nejlepších příznaků – víme, že extrahujeme smysluplné příznaky
- Porovnání s PCA
 - o PCA je lineární transformace kdežto AE umí modelovat i komplexní nelineární funkce.
 - o Příznaky nalezené PCA jsou lineárně nekorelované, protože jsou ortogonální. Autoencoded příznaky mohou být korelované.
 - o PCA je rychlejší a výpočetně levnější.
 - o Jednovrstvý AE s lineární aktivační funkcí funguje velmi podobně jako PCA.
 - o AE je náchylnější na overfitting kvůli velkému množství parametrů.
- **Regularizace**
 - o **Dropout, L1 regularizace**
 - o Reprezentace nižší dimenze se musí naučit extrahovat důležité příznaky
 - o Přidání regularizačního termu – podpora hledání smysluplných příznaků
 - **Sparsity regularisation** – k loss funkci přičteme počet aktivních neuronů
 - o **Denoising AE** – vstupní data zašuměny, takže se AE musí naučit podstatné příznaky pro rekonstrukci vstupu

Generativní neuronové sítě

- Snaha **namodelovat distribuci** trénovacích dat a **generování dat** ze stejného rozdělení, aby nebylo poznat, že jsou syntetická
 - o Generování obrázků, textu, dogenerování dat pro jiné metody
- **Gaussian mixture models (GMM)**
 - o **Modelování distribuce** – vážený součet Gaussovských distribucí
 - o Trénink váh a parametrů rozdělení pomocí expectation maximization algoritmu
 - EM algoritmy – nejdřív náhodně zvolíme parametry modelu, pak upřesňujeme – k-means
 - Odhad parametrů jednotlivých rozdělení
 - o Využití autoenkodéru – zakódujeme jím vstupy do latentního prostoru nízké dimenze, tam použít GMM k naučení distribuce a pomocí naučeného rozdělení generujeme nová data -> dekodérem převod do původních dimenzí
- **Autoregresivní generativní modely**
 - o **Predikce budoucího chování** na základě známých historických dat
 - o **Forecasting časových řad**
 - o Trénovací data musí obsahovat **autokorelaci** – pravidlo určující korelaci signálu a jeho časového zpoždění (např. obrázek – přechod pozadí – pixely se postupně mění – pozvolný přechod = vysoká autokorelace) – můžeme pozorovat **postupné změny**
 - o **PixelRNN** – predikce **pixelů** jako časových řad
 - V obr. Jdeme od rohu a predikujeme, nebo jdeme proti sobě ze dvou směrů

- Díky paměti lze hledat i netriviální sekvence
- Pomalé, nelze paralelizovat
- **PixelCNN** – pomocí **konvolučních filtrů** výběr oblastí obrázků, podle kterých se dopočítávají pixely
 - Lze více paralelizovat
 - Postupné generování částí obrázku a konvolucí extrakce příznaků, podle ní pak dobarvení dalších pixelů
 - Super-resolution – nejdřív modelování obecných tvarů a postupně doplnění detailů

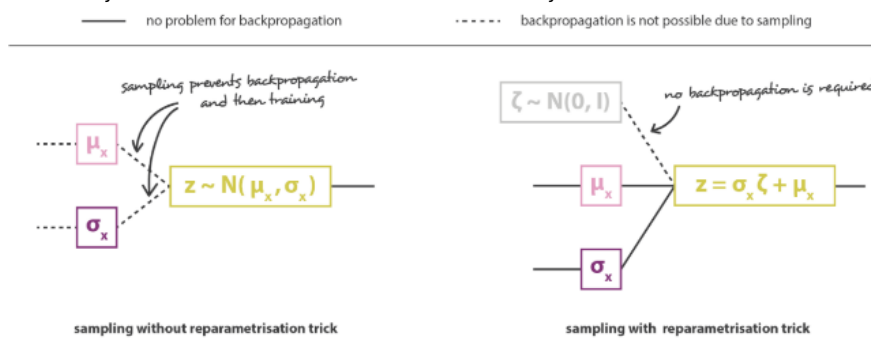
- Variational autoencoders (VAE)

- **Autoenkodér**, jehož **rozdělení** latentního prostoru (embeddingů) je **regularizováno** tak, aby mělo **dobré vlastnosti pro generování** nových dat
- Nedělají encoding vstupu na jeden bod latentního prostoru, ale **encodují ho jako distribuci** na



latentním prostoru, ze které je pak vzorkován bod – dekodován a porovnán se vstupem.

- VAE loss zakomponuje dva typy chyb:
 - Reconstruction error původního vstupu (least squares stejně jako AE)
 - Také jak moc se latentní distribuce liší od jednotkového normálního rozdělení



- Reparametrizace – pro propagaci chyby sítí – samplujeme jednotkové normální rozdělení

- Generative Adversarial Network (GAN)

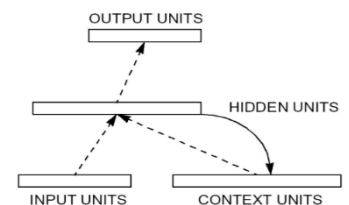
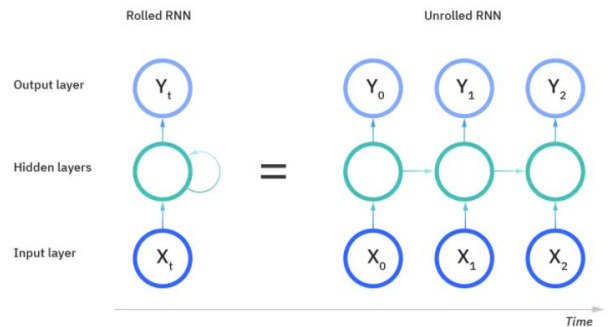
- Dvě komponenty – **Generátor a Diskriminátor** – hrají spolu **zero-sum hru**, kde mají společný loss a jeden se ho snaží minimalizovat a druhý maximalizovat.
- **Generátor** – z náhodného vektoru v latentním prostoru **generuje upsamplingem** (dekonvoluce) obrázek
- **Diskriminátor** se snaží **rozeznat skutečné obrázky** od těch uměle vygenerovaných
- Postupně se takhle navzájem vylepšují ve svých činnostech.
- Nestabilní trénování, často problém, že jeden model se moc rychle zlepší a druhý ho nemá šanci dohnat.
- Problém v počtu vygenerovaných objektů, problém s perspektivou,... Protože diskriminátor může hledat jen nějaké featury ale je mu jedno kolik jich je.
- Problémy:
 - Výstupy generují z náhodného šumu, takže pokud bychom chtěli generovat na základě nějakých konkrétních příznaků, tak nelze jednoduše určit zdrojový šum
 - Diskriminátor se učí jen rozpoznat skutečné a vygenerované vstupy, ale už neřeší, jestli ty vstupy vypadají tak jak mají

9. Rekurentní neuronové sítě a jejich učení, neuroevoluce.

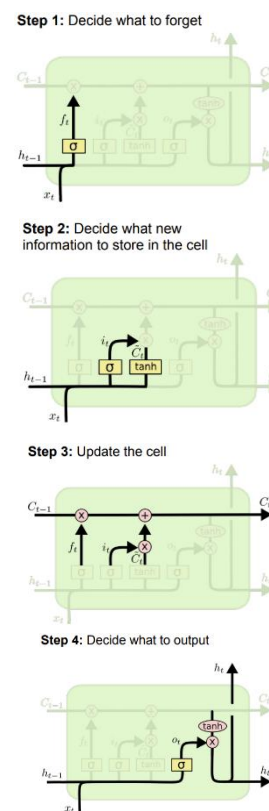
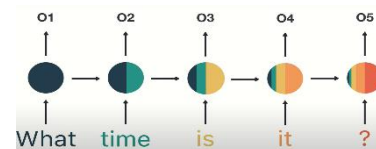
NI-MVI

Rekurentní neuronové sítě

- **NN s pamětí** – alespoň 1 cyklus, mezi časovými kroky předáván vstup – zachování dřívější informace, větší důraz na nedávné vstupy
 - o Aktivace neuronu i v případě, že v daném časovém kroku není žádný vstup
- Překlad, NLP, speech recognition, image captioning
- Vstupy (např. slova) různé délky, které se vždy zakódují na stejnou velikost pro zpracování stejnou architekturou
- Sdílení parametrů mezi časovými kroky – lze předpokládat, že určité příznaky se mohou vyskytovat v různých částech sekvence
- Trénovací data pro RNN mají podobu sekvence k input-output párů $\langle a_0, y_0 \rangle, \dots, \langle a_{k-1}, y_{k-1} \rangle$, přičemž výchozí hodnota vnitřního stavu x_0 musí být explicitně nastavena (typicky nulový vektor)
- **Backpropagation through time – BPTT**
 - o Způsob trénování RNN – podobný backpropagaci, ale chyby se napříč časovými kroky sčítají (protože sdílejí parametry)
 - Exploding/vanishing gradient problém
 - o Síť se rozbíjí do časových kroků (unroll), chyba kumulativně sečtena pro každý krok, poté složení sítě a aktualizace vah
 - Může vést k exploding/vanishing gradient problému
 - o Truncated Backpropagation in Time – kumulativní chyby pro podsekvence k kroků, prováděno postupně pro všechny podsekvence
- **Elmanova síť** – Vanilla RNN
 - o Feedforward síť s částečnou rekurencí
 - o Architektura – 4 vrstvy – vstup, skrytá vrstva, kontextová vrstva, výstup
 - o **Kontextová vrstva** – krátkodobě pamatuje výstupy skryté vrstvy – každý neuron má paměťovou buňku – detekce časově proměnlivých příznaků
 - o Využití BPTT – rekurence se rozbíjí a počítá se jako dopředná síť
- **Hopfieldova síť** – 1 vrstva n fully connected rekurentních neuronů (každý vstupem každého dalšího, vrstva je zároveň kontext)
 - o **Content-addressable memory** – asociativní paměť – do HS si lze ukládat data
 - Asociace – jsou referenční vzory
 - o **Hebbovské nesupervizované učení** – když se 2 neurony aktivují spolu, posilují se mezi nimi synapse
 - Pro dvojici neuronů součet násobku vstupů – váha synapse
 - o Lze modelovat jako energetickou funkci – pohyb v prostoru chyby sítě
 - Stabilní stavy – vyšší váhy
 - Snaha o minimalizaci energie
 - o Omezená kapacita
 - o Optimalizace, auto-asociace
- **Echo state networks** (reservoirs)
 - o Modelování složitých nelineárních dynamických systémů pomocí rekurentního modulu
 - o Signál vstupuje do RNN, tam se šíří a způsobuje oscilace
 - o Predikce časových řad
 - o Různé reakce neuronů na excitaci – nehomogenní, výstupní neurony tuto dynamiku převádí na cílový model dynamiky
 - o Minimalizace čtverců chyb



- Problém s krátkodobou pamětí (short term memory – způsobeno vanishing gradientem) – při zpracování vstupů problém udržovat informace z dřívějších (proto LSTM a GRU)
- Čistší RNN rychlejší, vhodné pro krátké sekvence a krátkodobé vztahy
- **Hradla** – tensor operace, které se učí, jaké informace přidávat/odebírat ze skrytého stavu
 - Umožňují dlouhodobé závislosti u GRU a LSTM
- **LSTM** – long short term memory
 - Informace z prvotních časových kroků se může projevit výrazně později
 - **Forget gate** – které informace ponechat a které zahodit
 - Vstup v daném kroce a vnitřní stav předchozího projdou sigmoidou, pokud je hodnota blízko 0, informace bude zapomenuta
 - **Input gate** – určuje, které informace ze současného vstupu budou uloženy do LSTM paměťové buňky
 - Vnitřní stav a vstup do sigmoidy a tanh, výstupy vynásobeny
 - **Cell state** – stav LSTM paměťové buňky
 - Podle výsledku Forget gate **stav buňky buď zapomenut nebo ponechán**
 - Výsledek sečten s výstupem Input gate – aktuální vstup ovlivněný předchozím stavem
 - Stav buňky přenášen jako kontext do příštího kroku + použit v aktuálním kroku pro výstup
 - **Output gate** – rozhodne hodnotu následujícího vnitřního stavu
 - Současný vstup a předchozí vnitřní stav do sigmoidy, nově pozměněný stav buňky do tanh, pak vynásobení – rozh., která část se přenesle do příštího vnitřního stavu
- **GRU** – Gated recurrent unit
 - Méně parametrů, než LSTM
 - Nemá stav buňky – používá vnitřní stav k přenosu informací
 - **Update gate** – rozhoduje, jak moc předchozí informace předat do budoucna
 - Co z minulých informací předávat dál, co ze současného vstupu
 - **Reset gate** – rozhoduje, jaké informace už můžeme zapomenout
- **Stack RNN** – realizace dlouhodobé paměti pomocí externí paměti – zásobník, fronta
 - Elmanova síť, ale místo kontextuálních buněk je zásobník
 - Může být i více paralelních zásobníků

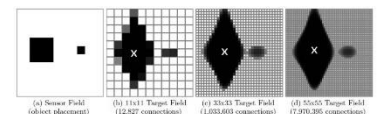


Neuroevoluce

- **Neuroevoluce** – užití evolučních algoritmů ke generování neuronových sítí
 - Kombinuje spojitý a diskrétní prostor – váhy a topologie
 - Síť může mít hodně ekvivalentních stavů, třeba prořezávat SP
 - Topology and Weight Evolving Artificial Neural Networks.
- GNARL rozděloval mutace na parametrickou (váhy mutující gaussovským šumem) a strukturální (přidávání a odebrání neuronů a spojů).
- **SANE = Symbiotic, Adaptive Neuro-Evolution**
 - Princip **koevoluce** – v jedné populaci se šlechtí váhy, v další neurony – **společná fitness funkce**
 - Vývoj neuronů jakožto vah spojů vstupujících do neuronu
 - Fitness = fitness 5 nejlepších sítí, ve kterých se daný neuron objevil
 - Blueprint se vyvíjí formou spojování neuronů do sítě – fitness se počítá jako fitness celé sítě
- **NEAT = Neuro-Evolution of Augmenting Topologies**
 - **Komplexifikace** – začíná se z malých topologií, evolucí se komplikují a propojují
 - Param. mutace (Gaussovské zašumění) a strukturální (přidávání neuronů a spojů, vypínání s., ...)
 - Chromozom = objekt, který má v sobě neurony a spoje

- Proměnlivá délka genomu
- **Mutate:**
 - **Přidáme spoj** mezi dva nespojené neurony
 - **Rozpojíme spoj a přidáme tam neuron** – propojuje původní dva sousedy
- Při tvorbě nových neuronů se změna označuje inovačním číslem – pořadí genetických změn
- **Křížení** – seřadíme genotypy rodičů podle inovačních čísel a postupně slučujeme do potomka
 - Potomek strukturálně komplexnější
- **Niching** – Populace je rozdělena na druhy – sítě **podobných vlastností** se **vyvíjí zvlášť** od jiných druhů
 - Protože v populaci je spousta sítí různých velikostí, menší se učí rychleji – přidáním genu dočasně snížíme fitness, ale kdybychom pořád chtěli jen vyšší fitness, nebudeme zvětšovat sítě a hledat komplexnější topologie
 - Výpočet podobnosti sítí přes vzdálenostní metriku
 - Dostatečně odlišný jedinec zakládá nový niche
 - **Fitness sharing** – fitness vydělen počtem jedinců v nichu – šance pro zajímavé slabší jedince se vyvinout – zabraňuje přemnožení zatím dobrých řešení
- **Přímé kódování** – všechny linky (spoje) jsou reprezentovány dedikovaným genem
- **Nepřímé kódování** – optimalizuje nějaké DNA (předpis), ze kterého se síť pak postaví
 - HyperNEAT, HyperGP
 - Stavění sítí, které mají nějaké symetrie
- **HyperNEAT**
 - NEAT není použit k vývinu topologie, ale k **vývinu jejího nepřímého kódování** stylem CPPN
 - **CPPN** – síť, která **přiřazuje neuronům váhy**
 - **Substrate** – prostor, ve kterém jsou rozmístěny neurony – např. rovina (neurony mají souřadnice) – CPPN dostane souřadnice, podle nějaké funkce určí váhu
 - Definuje možnosti propojování neuronů
 - Lze škálovat substrate density – hustota systému souřadnic
 - Jiný druh – HyperGP – místo neuronové sítě se nepřímo kóduje genetickým programováním
- **Novelty search** = místo abychom hledali něco konkrétního, co už známe, dáme šanci věcem zcela jiným
 - Prostor bývá plný lokálních minim, proto je vhodné se občas vydat zdánlivě odlišnou cestou, abychom dosáhli skutečného optima
 - Každé individuum při evoluci **odměněno**, pokud **objeví něco nového**
 - **Curiosity driven learning** – agent v prostředí získává **odměnu** za to, že se dostane do bodu, **kde předtím nebyl**
- **Covariance Matrix Adaptation Evolution Strategy (CMA-ES)**
 - **Strategie** pro složitou numerickou nelineární nekonvexní **black box optimalizaci**
 - Kombinuje vlastnosti evolučních algoritmů i gradient. technik – dobré překonávání lokálních minim
- **Ant Colony Optimization (ACO)**
 - Mravenci komunikují pomocí **feromonů** – algoritmus se toto snaží napodobit – agenti nanášejí v prostoru tolik „feromonu“, jak dobré je současné řešení
 - Čím více feromonu je někde nanášeno, tím více agentů tudy bude chodit.
 - Často uvázne v lokálním minimu, takže často kombinováno s tabu searchem.
 - Často používáno pro **hledání cest v grafech** – nejvíce feromonu se na hrany nanášejí, když je daná cesta dlouhá.
- **Particle Swarm Optimization (PSO)**
 - Částice mají **rychlostní vektor** udávající směr pohybu
 - Každý jedinec si pamatuje svoje lokální nejlepší řešení a snaží se k němu vrátit
 - Zároveň si celé **hejno** pamatuje globální nejlepší řešení a všichni míří k tomuto řešení.
 - Lokální a globální optima se balancují pomocí parametrů, postupně tak prohledávají prostor řešení
 - Záleží na počáteční rychlosti (možná oscilace)

1. Initialize population.
2. Compute fitness for all individuals.
3. Speciate by means of Explicit Fitness Sharing.
4. Adjust fitness $f' = f / \text{species_size}$.
5. Determine offspring count for all species proportionally to f' .
6. Eliminate inferior individuals of current generation.
7. Reproduce – replace current generation by its offspring.
8. While not satisfied, go to 2.

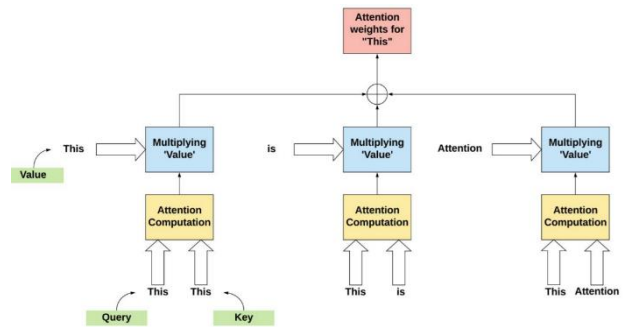


10. Transformer, pozornostní mechanismy, transfer a meta learning.

NI-MVI

Transformer

- Zpracovávají vstup ve formě sekvenčních dat, ale všechny najednou – mohou být paralelizovány
- Machine translation, text summarization, image description generation, ...
- Na rozdíl od RNN taky není tak problematický vanishing gradient.
- **Query** – jeden pro daný attention block
- Pro každou value se query vynásobí odpovídajícím klíčem key – tak získáme hodnotu parametru, kterým value násobíme -> získáme attention
- Poziční kódování – kvůli tomu, že jsou modelu předložena všechna data najednou – není info o pořadí slov
 - o Přičte ke každému vektorovému embeddingu **časový vektor**
 - Stejná slova budou zakódována jinak kvůli jiné pozici
- **Self-attention** – porovnání tokenů všech slov ve větě mezi sebou a převážení těchto embeddingů tak, že reprezentace zachycuje kontext
 - o **Keys = queries = values**
 - o Proces nezávislý na velikosti vstupu, nevyžaduje učení
 - o Skalární součin query a key, normalizace výsledku -> váhy pro daný query, 1 váha odpovídá 1 value – pronásobíme, sečteme -> kontextový vektor pro slovo použité jako query
 - Zopakování tak, že se jako query vystřídají všechna slova
 - o Parametrizované matice z keys, values a queries – trénovatelný attention
- **Multi head attention** – v rámci vstupů může být několik kontextů (jedno slovo má kontexty s více slovy) – je potřeba zachytit víc pozorností
 - o Několik vrstev lineárních vstupů k, v, q, každý má vliv na nezávislé trénování jednoho attentionu
 - o Natrénuje několik kontextových vektorů – spojení pomocí concatenation a napojení do dense vrstev
- **Masked attention**
 - o Při použití **maskingu** má model přístup jen ke vstupům, které byly před bodem, který se snaží predikovat
 - o Aby model **nemohl podvádět** a použít budoucí slova k predikci současného
 - o V dekodéru
- **Cross attention** – keys a values generovány z jiného výstupu než queries
 - o U dekodéru, keys a values se berou z enkodéru a queries z předchozího výstupu stackovaného dekodéru
- Transformer má enkodér a dekodér které oba používají multi head attention a mohou být stackovány
- V dekodéru se používá ještě masked attention a cross attention
- Z návrhu lze vidět hodně reziduálních spojení aby se potlačil gradient vanishing, takže před vstupem do attentionu se embeddingy ještě odpojí a znovu připojí po výstupu attentionu.
- Na výstupu decoder je pak Linera vrstva, která z výstupu dekoduje vektor velikosti slovníku a softmaxem určí, které slovo je nejpravděpodobnější.
- **BERT**
 - o **Masked language model** – self-supervizovaně se učí **doplňovat zamaskovaná (chybějící) slova** do textu
 - o **Next sentence prediction** – predikuje, **zda věta B následuje za větou A** v daném kontextu
 - o BERT trénován, aby minimalizoval loss obou úkolů
 - o 340M parametrů
 - o Seq2seq model využívající transformer architekturu se stackem enkodérů
 - o Použitelný jako jazykový základ, lze fine-tunovat na konkrétní úkoly



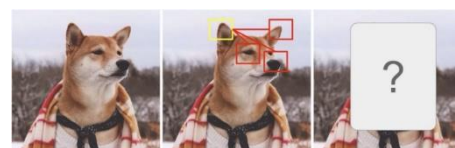
- Problém transformerů – vysoké paměťové nároky, rozšiřitelnost na dlouhé sekvence, složitá implementace, složitost attention vrstev (kvadratická vzhledem k délce sekvence)
- **Performer** – Obchází problém se škálovatelností attention vrstev – místo query a key matice dosazuje **aproximované matice**
 - o Lineární složitost vzhledem k délce sekvence
- **Synthesizer** – query a key nahrazeny a attention matice sestavena trénovatelnou sítí, až pak pronásobena values
 - o Random synthesizer – náhodná matice
- **Reformer** – lokální sémantické hashování – sousedé z podobných kontextů skončí s podobným hashem -> předpočítání podobnosti
 - o Rychlejší, ale náročnější na paměť

Pozornostní mechanismy

- **Attention** mechanismus **napodobuje kognitivní pozornost** – důležité části dat zvýrazněny, méně důležité potlačeny
- **Zpracování textu** – důraz na důležitá a více vypovídající slova
- **Zpracování obrazu** – zaměření na regiony v obraze, které obsahují důležité příznaky – nedůležité je např. pozadí
- Postup:
 - o Vstup = hodnoty (values v) – zakódováním slov nebo příznaků v obraze
 - o Výstup attentionu – lineární kombinace values s parametry alfa
 - Alfa normalizované na součet 1
 - o Pro získání alfa – key a query
 - Každá value má s sebou spojený klíč **key**
 - **Query** je jeden pro daný attention mechanismus
 - o Funkce, která zkombinuje i-tý key a query tak, že výsledkem je alfa pro i-tou value
 - Dot-product attention – skalární součin query a key, na výsledek aplikace nelineární transformace (třeba tanh), potom normalizace přes softmax
 - o **Alfa udává důležitost vstupu**
 - o **Attention** mechanismus je třeba trénovat – matice trénovatelných vah, které se lineárně kombinují s values, keys a queries
- **Attention block v dekodéru** – vnitřní stav enkodéru se pošle do attention bloku, projde FC vrstvou a softmax výstup určuje důležitost vstupů do dekodéru
- Pro seq2seq překlady – alignment model – tvorba **kontextového vektoru** – větší kontext = větší alfa
- Příklad, kdy values = keys
 - o Hledáme **relevanci různých pozic** jedné sekvence
 - o Pokud v rámci jedné věty hledáme slova, která odkazují na stejný objekt
- **Global/Soft attention** – Bere v potaz celý vstupní prostor, všechny attentiony se předají dál

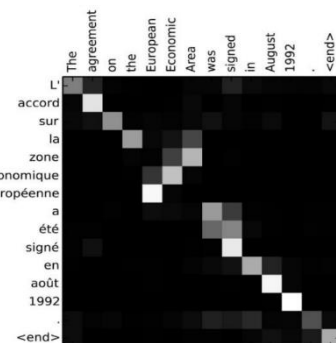
She is eating a green apple.

high attention
low attention



The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .

Fig. 6. The current word is in red and the size of the blue shade indicates the activation level. (Image source: Cheng et al., 2016)



- **Local/Hard attention** – Bere v potaz jen lokální region (např. patch v obrázku), pravděpodobnost že se zaměří na konkrétní část může být dána její relevancí (attention score)
- Výhoda: model se může soustředit na části vstupu, které nejlépe pomohou splnit úkol
- Nevýhoda: výpočetní složitost

Meta learning

- Stavba high-level **systému, který je přenositelný** na různé bottom-level AI problémy
- Ovlivňuje, jak vypadá daný podproblém, aplikovatelný na řadu úloh
- Modely se automatizovaně **učí se učit**
- Neuroevoluce je meta learning:
 - o meta-level = evoluční algoritmus modelující topologii sítě
 - o bottom-level = samotná neuronová síť
- Hyper Networks:
 - o Meta-level = malá síť
 - o Bottom-level = velká síť
- Použití LSTM k zapamatování sekvence updatů vah v neur. sítě – určuje nastavování vah bottom-level sítě
- **Optimalizace hyperparametrů (HPO)**
 - o Bottom-level AI = model, jehož hyperparametry optimalizujeme
 - o Meta-level = algoritmus, který provádí optimalizaci hyperparametrů
 - o **Grid search** – systematické prohledávání parameter grid (tabulka parametrů) – kombinace všech povolených hodnot
 - o **Random search** – náhodně prohledává parameter grid, překvapivě funkční
 - o **Bayesovská optimalizace** hyperparametrů – informované prohledávání, bere v potaz minulé ohodnocení pro vyzkoušené kombinace hyperparametrů
 - Funkce hodnotící výběr hyperparametrů podle přínosnosti
 - o **HyperBand = bandit-based** přístup k optimalizaci hyperparametrů
 - Kombinace **alokace zdrojů** a brzkého zastavení k prohledávání parameter gridu
 1. Volba většího množství n-tic hyperparametrů
 2. Přiřazení n-ticím budget několik iterací
 3. Po iteracích ohodnocení n-tic (trénink a ohodnocení modelu), zahazení horší poloviny
 4. Opět rozdělí budget na iterace mezi n-tice, po doběhnutí zase zahodí – tak iteruje a realokuje prostředky na trénování, dokud nepřežije jen jeden
- **Metadatabáze = databáze pro uchování metadat** o trénování různých problémů
 - o Informace, jaké metody se používají na jaké podproblémy a s jakými výsledky – může pak **doporučovat** vhodné algoritmy
 - o Podle extrahovaných atributů vypočteme podobné problémy, podle toho volba algoritmu
 - o Problém s kompatibilitou algoritmů a problémů – šum, neinformativita
- I pro rekomendační systémy – postupné upřesňování nastavení
- **Model Agnostic Meta Learning (MAML)**
 - o Meta learning přístup k **nastavování parametrů nezávisle na modelu a problému.**
 - o vybere vzorek ze sady problémů, problémy ohodnotí a provede update parametrů modelu pomocí gradientního sestupu
- **Neural architecture search (NAS)**
 - o Automatické učení a evoluce topologií neuronových sítí.
 - o **Prohledávání stavového prostoru** – přechody mezi stavy = různé operace jako konvoluční vrstvy, fully connected, pooling, ..., hyperparametry a propojování tak, aby tvořily validní neuronovou síť
 - o Stavový prostor bývá velmi velký a ohodnocení stavů (fitness) je výpočetně náročné
 - o Příklad – neuroevoluce
- Automated Machine Learning (AutoML)
 - o Předzpracování dat – extrakce a předzpracování příznaků
 - o Kritéria AutoML: Performance, speed, explainability, simplicity.
- **Few-shot learning**
 - o Učení na velmi **malých datasetech**
 - o Trénovací data = support set, podle kterého se model musí naučit se učit
 - o Využití předchozích znalostí o podobných případech, které už model zná
 - o Využití znalostí o samotném učení – pomocí omezení nutíme model generalizovat
 - o Využití znalostí o samotných datech a jejich distribuci a variabilitě

11. Ensemble metody: rozdíl mezi základními metodami (např. Bagging, Boosting, XGBoost).

NI-ADM

- **Ensemble metody** – více modelů použito k dosažení lepšího prediktivního výkonu v porovnání s použitím jednoho modelu
 - o Bagging snižuje rozptyl, boosting snižuje bias

Bagging

- **Bagging** = bootstrap aggregating – trénování **několika modelů nezávisle** na sobě a potom **kombinování** jejich predikcí
 - o Každý model natrénován na jiné podmnožině dat vytvořené **samplingem s nahrazením** (bootstrapping) z originálního datasetu
 - o Výhody:
 - Snižuje rozptyl – průměrováním predikcí z několika modelů
 - Snižuje overfitting – použitím několika modelů trénovaných na různých podmnožinách dat
 - o Nevýhody:
 - Zvyšuje složitost – využitím více algoritmů se zvyšuje složitost a výpočetní náročnost
 - Malý vliv na bias
- **Náhodné lesy** – algoritmus, který konstruuje množinu rozhodovacích stromů při tréninku a vyhodí třídu, která je **modus tříd** (klasifikace) nebo **průměr predikce** (regrese) jednotlivých stromů
 - o Každý **rozhodovací strom** stavěn na různé podmnožině originálních dat
 - Konstrukce např. pomocí ID3 algoritmu – pro každý příznak se spočítá kritérium (entropie/gini index), pak se zkusí rozdělit podle těch příznaků a počítá se kritérium, vybere se nejlepší rozdělení a opakuje se
 - o V každém uzlu stromu se k dělení používá pouze **náhodná podmnožina** příznaků
 - o Konečná predikce je **průměr (regrese) nebo většinový hlas (klasifikace)** napříč stromy
 - o Zvyšují prediktivní přesnost a kontrolují over-fitting zavedením náhody do ensmbu
 - Rozptyl se sníží bez zvýšení biasu
 - o **Hyperparametry**:
 - Počet stromů
 - Velikost bootstrapu
 - Počet příznaků

Boosting

- Kombinuje několik weak learnerů, které vytvoří strong learner
 - o **Weak learner** – algoritmus strojového učení, který klasifikuje s přesností o něco málo lepší než náhodné hádání
 - Decision stumps, naivní bayes, k-nejbližších sousedů
 - o **Strong learner** – má nízkou míru chyby
 - o Weak learneri se trénují v sekvenci, každá se snaží napravit předchůdce
- **AdaBoost** = Adaptive boosting
 - o Speciální případ Gradient boostingu
 - o Přiřazuje stejné váhy všem vzorkům a vybírá slabý klasifikátor, který minimalizuje chybu
 - o Pro evaluaci prvního learnera AdaBoost zvýší váhy špatně klasifikovaným vzorkům, aby tvořily větší část trénovací množiny pro následující klasifikátor
 - o Proces se opakuje, pokaždé přiřadí větší váhy špatně klasifikovaným vzorkům
 - o Závěrečná predikce je **vážený hlas** (v klasifikaci) nebo **vážená suma** (v regresi) predikcí tvořených jednotlivými learnery
 - o AdaBoost je **adaptivní** ve smyslu že následující weak learneri jsou upravováni ve prospěch instancí špatně klasifikovaných předchozími kroky

- Cíl je nastavit váhy a trénovat data tak, aby správně předpovídali neobvyklé pozorování
- Někdy se používají **Decision stumps** – jako stromy v rozhodovacím lese, ale mají jen jeden uzel a dva listy
- Algoritmus:
 - **Inicializace vah** (stejně hodnoty)
 - **Postaví se weak learner** – trénink základního modelu
 - **Výpočet chyby** – suma vah špatně klasifikovaných bodů
 - **Výpočet důležitosti learneru** – na základě chyby – menší chyba značí větší důležitost
 - **Update vah** – zvětší se váhy špatně klasifikovaných bodům
 - To se **opakuje** dokud se chyba nepřestane zlepšovat
 - **Formace finálního klasifikátoru** – vážený hlas mezi slabými klasifikátory
- Výhody:
 - **Všestrannost** – je robustní a všestranná metoda, pro mnoho praktických aplikací
 - **Komplexní klasifikace** – dobře zvládá špatně klasifikovatelné body z komplexních problémů
 - **Přesnost**
- Nevýhody
 - Riziko **přeučení** – hlavně na datasetech s vysokým šumem, ale dá se zlepšit regularizací
 - **Složitost** výběru a ladění – výběr slabého klasifikátorů a ladění hyperparametrů vyžaduje expertízu

- XGBoost

- **Spojité cílové proměnné** – může být pravděpodobnost, že Y patří do třídy c (podobně jako v logistické regresi)
- Závěrečná predikce dána **n weak learnery** – hyperparametr n je dané číslo a weak learneri rozhodovací stromy
- Během trénovací fáze se weak learneri f_k postupně konstruuji, aby se zlepšil ensemble model:

$$F^{(t)} = \sum_{k=1}^t f_k = F^{(t-1)} + f_t \text{ pro } n = 1, 2, \dots, n$$
- Závěrečný model je $F^{(n)} = \sum_{k=1}^n f_k = F^{(n-1)} + f_n$
- **Ztrátová funkce $l(y, \hat{y})$** – měří kvalitu predikce
 - Pro gradient boosting obecně musí mít všude první derivace, a v XGBoostu musí mít všude druhé derivace
 - Derivace existují pro obvyklou volbu čtverce reziduí: $l(y, \hat{y}) = (y - \hat{y})^2$
- Cíl trénovacího procesu je minimalizovat **objective funkci**

$$\sum_{i=1}^N \ell(y_i, \hat{y}_i^{(n)}), \quad \hat{y}_i^{(n)} = F^{(n)}(\mathbf{x}_i) = \sum_{k=1}^n f_k(\mathbf{x}_i)$$
 je predikce pro i -tý datový bod a N počet trénovacích datových bodl
- **Regularizace** – pro předejití overfittingu – s ní je objective funkce:

$$\sum_{i=1}^N \ell(y_i, \hat{y}_i^{(n)}) + \sum_{k=1}^n \Omega(f_k)$$
 - Ω je rostoucí **funkce složitosti stromu f_k**
 - Složitost se dá měřit hloubkou, počtem listů, součtem vah listů, ...
- Minimalizujeme

$$\sum_{i=1}^N \ell(y_i, \hat{y}_i^{(t)} + f_{t+1}(\mathbf{x}_i)) + \Omega(f_{t+1})$$
 - Podle volby ztrátové funkce to nemusí být možné
 - Zjednodušení – **aproximace pomocí Taylorova polynomu**
 - V XGBoostu polynom druhého stupně – aproximace v bodě $x_0 + d$

$$f(x_0 + d) \approx f(x_0) + f'(x_0)d + \frac{1}{2}f''(x_0)d^2.$$
 - Cílem je **minimalizace** ne y , ale „**délky kroku**“ **ve směru daného derivací** – gradientu, kvůli tomu jsou to gradientové metody duh
 - Dalšími úpravami můžeme dostat míru kvality stromu a předpis pro váhy
 - Stále NP-těžké konstruování stromů, lze použít greedy algoritmus jako ID3

12. Jádrové metody: jádrová regrese, báze funkce, Support Vector Machine (SVM): separabilní a neseperabilní případ.

NI-ADM

Připomenutí základů

- **Lineární regrese** – vysvětlovaná proměnná Y a vektory příznaků X : $Y = w_0 + w_1x_1 + \dots + w_px_p + \varepsilon$
 - o Vektor vah w a náhodná proměnná ε
 - o Odhadujeme vektor vah \hat{w} a pomocí něj predikce Y : $\hat{Y} = \hat{w}^T x$
- **OLS – ordinary least squares**
 - o Při tréninku minimalizace reziduální sumy čtverců $RSS(w) = \|Y - Xw\|^2$
 - o Minimum pomocí **normální rovnice**: $X^T Y - X^T X w = 0$ (odpovídá $\nabla RSS(w) = 0$)
 - o Pokud je $X^T X$ regulární, existuje právě jedno řešení: $\hat{w}_{OLS} = (X^T X)^{-1} X^T Y$
- **Hřebenová regrese**
 - o Minimalizujeme regularizovaný reziduální součet čtverců: $RSS_\lambda(w) = \|Y - Xw\|^2 + \lambda \sum_{i=1}^p w_i^2$
 - λ je prostě debilní parametr $\lambda \geq 0$
 - o Matice I' - na diagonále kromě 1. prvku jedničky, jinak všude nuly
 - $RSS_\lambda(w) = \|Y - Xw\|^2 + \lambda w^T I' w$
 - o To nám předělává normální rovnici a nové řešení (a zjistíme, že pro $\lambda = 0$ máme normální OLS)
 - $X^T Y - X^T X w - \lambda I' w = 0$ (to co odčítáme je regulární pro nezáporné λ)
 - $\hat{w}_\lambda = (X^T X + \lambda I')^{-1} X^T Y$

Lineární báze funkce

- Zobecnění, aby bylo možno použít pro nelineární funkce – nahrazení původních přízn. za transformované
 - o \mathcal{X} ... prostor obsahující všechny možné hodnoty vektoru příznaků $X = (X_1, \dots, X_p)^T$, typ. $\mathcal{X} = \mathbb{R}^p$
 - o $\varphi_1 \dots \varphi_M$... **báze funkce** = lineárně nezávislé funkce z \mathcal{X} do \mathbb{R}
 - Představují **transformace** pův. příznaků X_1, \dots, X_p do **nového M-dimenzního prost. příznaků**
 - Model pro vysvětlovanou proměnnou je teď v tom novém prostoru lineární
- **Model lineární báze funkce**
 - o Vysvětlovaná proměnná Y v hodnotě $x = (x_1, \dots, x_p)^T$ z \mathcal{X} je dána jako:
$$Y = w_1 \varphi_1(x) + \dots + w_M \varphi_M(x) + \varepsilon = w^T \varphi(x) + \varepsilon$$
 - $\varphi: \mathcal{X} \rightarrow \mathbb{R}^M$ je definována jako $\varphi(x) = (\varphi_1(x) \dots \varphi_M(x))^T$ pro všechna $x \in \mathcal{X}$
 - $w = (w_1, \dots, w_M)^T$, ... vektor neznámých parametrů
 - ε ... náhodná proměnná s $E\varepsilon = 0$

Odhad lineární báze funkce

- o **Model** pro tréninkový vzorek daný N dvojicemi $(Y_1, x_1) \dots (Y_N, x_N)$ v maticové formě:
$$Y = \Phi w + \varepsilon$$

kde

$$\Phi = \begin{pmatrix} \varphi(x_1)^T \\ \vdots \\ \varphi(x_N)^T \end{pmatrix} = \begin{pmatrix} \varphi_1(x_1) & \dots & \varphi_M(x_1) \\ \vdots & \ddots & \vdots \\ \varphi_1(x_N) & \dots & \varphi_M(x_N) \end{pmatrix}, \quad Y = \begin{pmatrix} Y_1 \\ \vdots \\ Y_N \end{pmatrix}, \quad \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_N \end{pmatrix}$$

- o **Minimalizace** (podobně jako v hřebenové regresi):

$$RSS_\lambda(w) = \|Y - \Phi w\|^2 + \lambda w^T w$$

- o **Normální rovnice**:

$$\Phi^T Y - \Phi^T \Phi w - \lambda w = 0$$

- o Pro existuje právě jedno **řešení**:

$$\hat{w}_\lambda = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T Y$$

- o **Predikce** Y v x :

$$\hat{Y} = \hat{w}_\lambda^T \varphi(x)$$

- Různé volby báze funkce zajišťují flexibilitu modelu
- Často skončíme s velkým počtem báze funkcí (větším než tréninkový vzorek), takže je nutná regularizace

Jádrový trik

- Duální reprezentace

- Převod problému do formy, kdy data vstupují do problému jako skalární součiny
- V duální reprezentaci jsou báze funkce dány implicitně skrze **jádrovou funkci**
- Máme hodnoty w ve tvaru $w = \Phi^T \alpha$, kde $\alpha \in \mathbb{R}^N$

- Φ ... matice báze funkcí
- **Omezení w** na podprostor \mathbb{R}^M generovaný vektory $\varphi(x_1) \dots \varphi(x_N)$
- Dosazení do $RSS_\lambda(w)$:

$$RSS_\lambda(\alpha) = \|Y - \Phi \Phi^T \alpha\|^2 + \lambda \alpha^T \Phi \Phi^T \alpha$$

- Protože jsme w omezili, dostaneme:

$$\min_{\alpha} RSS_\lambda(\alpha) = \min_{\alpha} RSS_\lambda(w(\alpha)) \geq \min_w RSS_\lambda(w)$$

- Pokud $\lambda > 0$, pak

$$\min_{\alpha} RSS_\lambda(\alpha) = \min_w RSS_\lambda(w)$$

Navíc, pokud w^* minimalizuje $RSS_\lambda(w)$, potom $\alpha^* = \frac{1}{\lambda}(Y - \Phi w^*)$ minimalizuje $RSS_\lambda(\alpha)$

Pokud α^* minimalizuje $RSS_\lambda(\alpha)$, potom $w^* = \Phi^T \alpha^*$ minimalizuje $RSS_\lambda(w)$

- Tím pádem jsou úlohy minimalizace $RSS_\lambda(\alpha)$ a $RSS_\lambda(w)$ **ekvivalentní**
- λ ... regularizační parametr – penalizuje hodnoty vah, aby nedošlo k přeučení

- Jádrová funkce $k: \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}$

$$k(x, y) = \varphi(x)^T \varphi(y)$$

- Gramova matice: $G = \Phi \Phi^T$

- $N \times N$ symetrická matice
- **Pozitivně semidefinitní:**

$$\alpha^T G \alpha = \alpha^T \Phi \Phi^T \alpha = (\Phi^T \alpha)^T (\Phi^T \alpha) = \|\Phi^T \alpha\|^2 \geq 0$$

- Pomocí ní se dá $RSS_\lambda(\alpha)$ vyjádřit jako:

$$RSS_\lambda(\alpha) = \|Y - G\alpha\|^2 + \lambda \alpha^T G \alpha$$

- (i, j) -tá komponenta Gramovy matice:

$$G_{i,j} = (\Phi \Phi^T)_{i,j} = \sum_{\ell=1}^M \Phi_{i,\ell} \Phi_{j,\ell} = \sum_{\ell=1}^M \varphi_\ell(x_i) \varphi_\ell(x_j) = k(x_i, x_j)$$

- Gramova matice je celá dána jádrovou funkcí

- Predikce v duální reprezentaci

- Minimalizujeme $\hat{w} = \Phi^T \hat{\alpha}$
- Predikce Y v x :

$$\hat{Y} = \hat{w}^T \varphi(x) = \hat{\alpha}^T \Phi \varphi(x) = \sum_{i=1}^N \sum_{j=1}^M \hat{\alpha}_i \Phi_{i,j} \varphi_j(x) = \sum_{i=1}^N \hat{\alpha}_i k(x_i, x) = \hat{\alpha}^T k(x)$$

- Kde $k(x) = (k(x_1, x), \dots, k(x_N, x))^T$
- Tím pádem nám na vyjádření účelové funkce $RSS_\lambda(\alpha)$ i predikce \hat{Y} stačí jádrová funkce (a stejně tak pro $\hat{\alpha}$)

- Odhad α

- Řešení uzavřené formy minimalizačního problému
 - Gradient $\nabla RSS_\lambda(\alpha) = -2G(Y - G\alpha) + 2G\lambda\alpha$
 - Položení rovno nule \rightarrow normální rovnice $G(Y - G\alpha - \lambda\alpha) = 0$
 - Kvůli pozitivní semidefinitnosti Gramovy matice je $(G + \lambda I)$ pozitivně definitní
- Dostáváme odhad $\hat{\alpha} = (G + \lambda I)^{-1} Y$
 - Tak se dá vyjádřit jenom pomocí jádrové funkce

- Jádrový trik

- Stvořili jsme ekvivalentní duální reprezentaci celého modelu spolu s účelovou funkcí, která je formulovaná tak, že vstupní vektor x vstupuje do výpočtu jako forma skalárních součinů
- Jádrový trik = náhrada skalárních součinů za jádrovou funkcí
- Rozšíření – začne se s jádrem bez explicitní specifikace báze funkce
 - Implicitně tak můžeme používat prostory příznaků vysokých, i nekonečných dimenzí

- Lineární jádro – pro všechny x platí $\varphi(x)$:

$$k(x, y) = x^T y$$

- Zobecnění – polynomiální jádro:

$$k(x, y) = (x^T y + 1)^n$$

Jádrové (kernelové) stroje

- Předpokládáme, že skutečný model vysvětlované proměnné Y v x je $Y = f(x) + \varepsilon$, kde f je neznámá funkce a ε je náhodná proměnná s nulovou střední hodnotou
- V lineární báze expanzi použijeme odhad:

$$f(x) = w_1 \varphi_1(x) + \dots w_M \varphi_M(x)$$

- **Jádrový stroj** – model, kdy začínáme z jádrové funkce k :

$$f(x) = \sum_{j=1}^K \alpha_j k(x, \mu_j)$$

- $\mu_1 \dots \mu_K \in \mathcal{X}$ jsou středy
- Jádrový stroj odpovídá **lineární báze expanzi** s $\varphi_j(\cdot) = k(\cdot, \mu_j)$
 - Proto se $\varphi(x) = (k(x, \mu_1), \dots, k(x, \mu_K))^T$ jmenuje kernelizovaný vektor příznaků

- **Vektorové stroje**

- Speciální případ jádrových strojů, kde **středy jsou dány body z trénovacího datasetu**:

$$f(x) = \sum_{j=1}^N \alpha_j k(x, x_j)$$

- Vektorový stroj můžeme dostat tím, že uděláme **jádrový trik na modelu lineární báze expanze**
- Zajímají nás **řídke vektorové stroje**, kde je pro hodně bodů platí $\alpha_j = 0$

Diskriminativní funkce

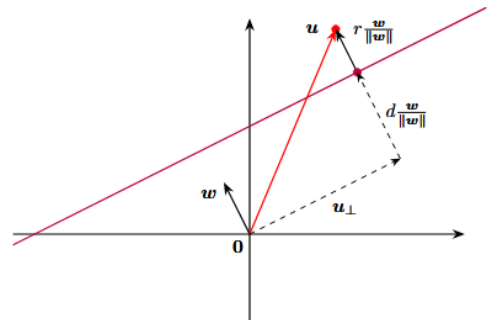
- Zaměřujeme se teď na binární klasifikaci, kde predikujeme $Y \in \{-1, 1\}$ v bodě $x \in \mathcal{X}$
- $f(x) = w^T \varphi(x) + w_0$

- Pokud $f(x) > 0$, pak $\hat{Y} = 1$, jinak -1 :

$$\hat{Y}(x) = \text{sgn}(f(x)) = \text{sgn}(w^T \varphi(x) + w_0)$$

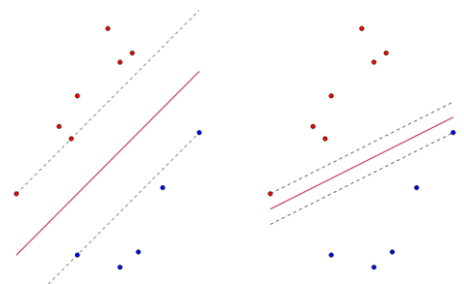
- **Rozhodovací hranice** v prostoru příznaků \mathbb{R}^M

- Definovaná **podmínkou $w^T u + w_0 = 0$**
- Řešení tvoří hyperrovinu v \mathbb{R}^M = rozdělovací hyperrovina (separating hyperplane)
 - Normála je dána vektorem w
 - Vzdálenost od 0 ve směru w : $d = -\frac{w_0}{\|w\|}$
 - Libovolný bod u lze rozložit na dvě části – u_{\perp} je kolmá na w a u_{\parallel} je s ním rovnoběžná
 - $r = u_{\parallel} - d$... vzdálenost u od rozhodovací hranice ve směru w
 - Můžeme dosadit do podmínky a získat $w^T u + w_0 = r \|w\|$, z čehož můžeme vyjádřit r a dosadit do diskriminativní funkce
 - Tím dostaneme, že **vzdálenost $\varphi(x)$ od rozhodovací hranice** je dána $r = \frac{f(x)}{\|w\|}$
 - Ta vzdálenost je **pozitivní, když $f(x) > 0$** a negativní, když $f(x) < 0$
 - **Znaménko $f(x)$ separuje prostor příznaků** na dva podprostory
 - Pozitivní je ve směru w , všechny predikce jsou 1
 - Negativní ve směru $-w$, predikce jsou -1



Princip velkého okolí (Large margin principle)

- Pokud máme diskriminativní funkci f , predikce v prostoru příznaků jde dobře geometricky interpretovat
 - Chceme nějak natrénovat neznámé parametry w a w_0
 - Trénovací vzorky N párů ve tvaru $(Y_1 x_1), \dots, (Y_N x_N)$
- **Lineárně separabilní případ** – máme trénovací body v prostoru příznaků $\varphi(x_1) \dots \varphi(x_N)$ a existuje alespoň jedna volba



parametrů w a w_0 taková, že $f(x_i) > 0$ pro body kde $Y_i = 1$ a $f(x_j) < 0$ pro body kde $Y_j = -1$ (binární klasifikace)

- Řešení může být nekonečno

- **Margin** = okolí = **nejmenší vzdálenost** mezi rozhodovací hranicí a kterýmkoliv ze vzorků

- Toto se snažíme maximalizovat

- Vzdálenost i -tého trénovacího bodu $\varphi(x_i)$ od rozhodovací hyperroviny je dána $r_i = \frac{f(x_i)}{\|w\|}$

- Takže **hledáme w a w_0** takové, aby minimální r_i bylo co největší, a zároveň aby každý bod byl na správné straně hranice: $Y_i f(x_i) > 0$

$$\max_{w, w_0} \min_i \frac{Y_i(w^T \varphi(x_i) + w_0)}{\|w\|}$$

- w a w_0 jdou libovolně přeškálovat, takže můžeme zvolit, že pro bod nejbližší k rozhodovací hranici bude platit $Y_j(w^T \varphi(x_j) + w_0) = 1$

- Tím pádem pro všechny body musí platit $Y_i(w^T \varphi(x_i) + w_0) \geq 1$

- Hledáme **minimum** $1/\|w\|$ – maximalizace $\|w\|^{-1}$, ekvivalentní minimalizaci $\|w\|^2$

- Řešíme následující **optimalizační problém**:

$$\min_{w, w_0} \frac{1}{2} \|w\|^2 \quad \text{subject to} \quad Y_i(w^T \varphi(x_i) + w_0) \geq 1 \quad \text{for all } i$$

- Omezení říká, že všechny body jsou na správné straně hranice s okolím minimálně 1 (vhodné škálování)
- Příklad kvadratického programování – minim. kvadratické funkce za podmínky lineární nerovnice
- **Řešení** ve formátu jádrové funkce je **support vector machine**

Support vector machine

- Pokud **nemáme lineárně separabilní data**, nemůžeme tu nerovnici dobře vyřešit pro všechna i

- Musíme **zrelaxovat podmínky** a povolit, ať jsou některé body na špatné straně hranice
- Zavedení **penalizace** se zvyšující se vzdáleností od hranice ve špatném směru

- **Volné proměnné**: $\xi_i \geq 0$, $i = 1, \dots, N$ takové, že $\xi_i = 0$, pokud je bod na nebo vně správného okolí na správné straně rozhodovací hranice, jinak $\xi_i = |Y_i - f(x_i)|$

- Pokud $\xi_i = 0$, pak $Y_i f(x_i) \geq 1$
- Pokud $0 < \xi_i \leq 1$, bod leží uvnitř okolí, ale na správné straně hranice a $Y_i f(x_i) > 1$
- Pokud $\xi_i > 1$, bod leží na špatné straně hranice a je špatně klasifikován – $Y_i f(x_i) < 0$

- Cíl SVM – díky volným proměnným můžeme **uvolnit hard podmínky** $Y_i f(x_i) \geq 1$ **soft margin podmínkami**:

$$Y_i f(x_i) \geq 1 - \xi_i \quad \forall i = 1 \dots N$$

- **Nový optimalizační problém**:

$$\min_{w, w_0, \xi} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^N \xi_i$$

S podmínkami:

$$\xi_i \geq 0 \quad \text{and} \quad Y_i(w^T \varphi(x_i) + w_0) + \xi_i - 1 \geq 0 \quad \text{for all } i.$$

- Protože $\xi_i > 1$ znamená že bod i je špatně klasifikovaný, můžeme $\sum_i \xi_i$ chápat jako **horní hranici** počtu **špatně klasifikovaných bodů**
- **C ... regularizační parametr** – kontroluje počet chyb, které můžeme tolerovat v trénovací množině
 - Větší C = větší důraz na minimalizaci chyb
 - Často $C = \frac{1}{\nu N}$, kde $0 < \nu \leq 1$ ovládá zlomek špatně klasifikovaných bodů během trénovací fáze - ν – **SVM klasifikátor**

- Pomocí Lagrangeových multiplikátorů a principu duality můžeme **maximalizovat ekvivalentní Lagrangeovský duální problém**:

$$\tilde{L}(a) = \sum_{i=1}^N a_i - \frac{1}{2} \sum_{i,j=1}^N a_i a_j Y_i Y_j k(x_i, x_j)$$

- Vzhledem k novým proměnným a s podmínkami:

$$0 \leq a_i \leq C \quad \text{and} \quad \sum_{i=1}^N a_i Y_i = 0$$

- $k(x_i, x_j)$... jádrová funkce, která má k původní báze funkci vztah $k(x, y) = \varphi(x)^T \varphi(y)$
- **Souvislost s KKT**¹
 - Abychom vyřešili (minimalizovali) duální problém, musíme maximalizovat Lagrangian s ohledem na w a w_0 – vytvoříme parciální derivace, položíme je rovny 0 a získáme dvě rovnice:
 - Parc. derivace w : $w = \sum a_i Y_i x_i$
 - Parc. derivace w_0 : $\sum a_i Y_i = 0$
 - Dosazení rovnic do optim. problému dostaneme Lagrangeovský duální problém (výše)
 - Jelikož $a \geq 0$, bude minimum v $a = 0$
 - Pro bod, ve kterém se nachází min. Lagrangeovy funkce L platí $\nabla L = 0$ a KKT podmínky
 - $a = 0$ pro všechny ne-support vektory (z KKT podmínek prý)
 - **KKT** podmínky jsou (snad) nutné a postačující podmínky pro optimální řešení problémů s omezením
 - Primal feasibility – všechny body musí splňovat margin podmínky: $Y_i f(x_i) \geq 1 - \xi_i$
 - Dual feasibility – všechny Lagrangeovy multiplikátory musí být nezáporné: $\alpha_i \geq 0$
 - Complementary slackness – buď $\alpha_i = 0$, nebo x_i leží přesně na rozhodovací hranici (marginu) nebo uvnitř ní: $\alpha_i (Y_i f(x_i) - 1 + \xi_i) = 0$
 - Stacionarita – gradient Lagrangianu podle w musí být nulový: $\nabla_w L = 0$
- Řešení = **support vector machine (SVM)** ve formě:

$$f(x) = \sum_{j=1}^N a_j Y_j k(x, x_j) + w_0$$

- **Predikce** v bodě x : $\hat{Y}(x) = \text{sgn}(f(x))$

- Vlastnosti řešení

- Řešení je **vektorový stroj**

$$f(x) = \sum_{j=1}^N \alpha_j k(x, x_j) + w_0$$

kde $\alpha_j = a_j Y_j$

- Je **řídke** - $a_j = 0$ pro všechny správně klasifikované body ležící vně okolí hranice
- **Support vector** = body, kde $a_j > 0$
 - Datové body, které leží nejbliž k rozhodovací rovině/hyperplane/hranici
 - Datové body, které je **nejtěžší klasifikovat**
 - Mají přímý vliv na polohu optimální rozhodovací hranice
 - Rozhodovací funkce je plně specifikovaná malou podmnožinou trénovacích vzorků = support vektorů
- Body s $a_j < C$ leží na hranici okolí
- Body s $a_j = C$ jsou buď správně klasifikované uvnitř okolí, nebo špatně klasifikované
- Čím větší je regularizační parametr C , tím řidší řešení dostaneme

- Vstup a výstup SVM

- Vstup – množina (vstup, výstup) párů trénovacích vzorků (vstup příznaky x , výstup y)
- Výstup – množina vah w , jedna pro každý příznak, jejichž lineární kombinace predikuje hodnotu y
- Používá se optimalizace na maximalizaci marginu (okolí, taková „šířka ulice“), abychom zmenšili počet vah, které jsou nenulové, na jenom pár, které odpovídají důležitým příznakům, na kterých záleží při stanovování rozhodovací hranice
 - Nenulové váhy odpovídají support vectorům („podporují rozhodovací hranici“)

- Support vector – **lineárně separabilní** případ

- Prvky trénovacího datasetu, které kdybychom odstranili, tak by se změnila pozice rozhod. hranice
- Jsou to kritické prvky trénovacího datasetu
- Problém nalezení optimální hyperroviny je optimalizační problém
 - Řešení pomocí Lagrangeových multiplikátorů

¹ Holeňova a Vašatova oblíbená otázka, MUSÍME umět

13. Algoritmy pro doporučování: základní přístupy a způsob vyhodnocení kvality, faktorizační metody pro doporučování.

NI-ADM

Základní principy

- **Interakce uživatele a položky** – systémy berou v potaz historické interakce mezi uživateli a položkami
 - o Explicitní – hodnocení filmů, implicitní – historie nákupů
- **Hodnocení podobnosti** – založeny na principu podobnosti – uživatelé se sdílenými zájmy pravděpodobně předvádějí stejné chování, a položky, které jsou často vybírány společně jsou nejspíš nějak spojeny
- **Personalizace** – systémy nabízejí personalizované návrhy spíše než univerzální množinu doporučení, na základně jednotlivých preferencí
- Typy úloh:
 - o **Top-N** – cílem je najít n věcí, které jsou pro uživatele nejvíc relevantní
 - Dobré doporučovat pouze položky, které uživatel nezná
 - o **Kolaborativní filtrování** – predikce o uživateli na základě sběru informací od hodně uživatelů
 - Pokud se dva uživatelé shodnou na jedné věci, nejspíš se shodnou i na dalších
 - o **Obsahové filtrování** – systém doporučuje položky podobné tomu, co si uživatel v minulosti oblíbil, v závislosti na vlastnostech položky
 - o **Hybridní metody** – kombinace kolaborativního a obsahového filtrování
 - o **Demografické doporučovací metody** – podle demografiky uživatele se doporučuje
 - o **Utilitní doporučovací úlohy** – tvorba utilitní funkce pro každého uživatele
- **k nejbližších sousedů** pro doporučovací systémy
 - o **User-based** kolaborativní filtrování – kNN identifikuje uživatele podobné danému uživateli
 - Hodnotí k nejpodobnějších uživatelů a jejich preference použije k doporučování
 - o **Item-based** kolaborativní filtrování – kNN identifikuje položky podobné těm, o které uživatel v minulosti projevil zájem
 - Systém doporučí položky nejpodobnější těm, které uživatel hodnotil kladně
 - Výběr k může prudce ovlivnit kvalitu doporučování, volí se empiricky vzhledem k výkonu systému
- **Scoring function** – predikce na základě hodnotící funkce – skóre $U \times I \rightarrow \mathbb{R}$ hodnotí jednotlivé itemy na základě kontextu (uživatelů)
- **Hodnocení**
 - o **Cumulative gain** – hodnotí, jak relevantních je prvních k položek
$$CG_k = \sum_{i=1}^k rel_i$$
 - **Discounted CG** – hodnota relevance logaritmicky proporcionálně zmenšena podle pozice výsledku
 - o **Precision** – kolik z doporučených ho skutečně zajímá
$$precision = \frac{|relevantní \cap doporučené|}{doporučené}$$
 - o **Recall** – kolik z položek, které ho zajímají, jsme mu doporučili
 - o F1 skóre, AUC (plocha pod křivkou ROC), RMSE pro spojitě hodnoty

Faktorizační metody pro doporučování

- **Maticová faktorizace** – rozkládá velkou, často řídkou user-item interaction matici na dvě menší matice
 - o **Odvozené matice** ukládají latentní atributy jak uživatelů, tak položek, a zjednodušují tak jejich komplexní vztah
 - o Např. doporučování filmů – latentní faktory mohou být preference žánru, režiséra nebo herce a korelace filmu s těmito faktory – predikce uživatelova hodnocení filmu může být zobrazena pomocí skalárního součinu vektoru latentních faktorů pro toho uživatele a položku

- **Optimalizace** v maticové faktorizaci

- Chceme **minimalizovat rozpor** mezi původní a rekonstruovanou maticí hodnocení
- Rozptyl měřen pomocí ztrátové funkce jako **MSE**:

$$\operatorname{argmin}_{U,V} \sum_{(i,j) \in \Omega} (r_{i,j} - \mathbf{u}_i^T \mathbf{v}_j)^2 + \lambda \left(\sum_x \|\mathbf{u}_x\|^2 + \sum_y \|\mathbf{v}_y\|^2 \right)$$

- Často pomocí **stochastického gradientního sestupu** a **alterujících nejmenších čtverců** (ALS) za použití regularizace pro předejití overfittingu
- **Řešení řídkosti** – uživatel obvykle interaguje jen s malou podmnožinou položek – maticová faktorizace to řeší zobrazením uživatelského hodnocení na všechny položky, včetně těch nepřítomných v trénovacích datech, a doplní tak nepozorované preference uživatele
- **Alterující nejmenší čtverce (ALS)** – technika zaměřená na minimalizaci mezery mezi pozorovanými a predikovanými hodnoceními
 - Alteruje mezi opravami uživatelské matice a položkové matice a zároveň optimalizuje tu druhou
 - **ALS algoritmus**:
 - Inicializace uživatelské (U) a položkové (V) matice náhodně
 - Opakování do konvergence:
 - Oprava V a minimalizace ztrátové funkce U
 - Oprava U a redukce ztrátové funkce V
 - Může zahrnovat **implicitní zpětná vazba** – kliknutí, historie nákupů, historie prohlížení
 - Nutné modifikace ztrátové funkce, aby byly zahrnuty konfidenční hladiny odvozené z množství interakce mezi uživatelem a položkami
 - Implicitní feedback = informace o preferencích, které nelze získat z explicitních akcí, jako jsou třeba hodnocení nebo recenze
- **Modelování implicitní zpětné vazby**
 - Potřeba když nejsou k dispozici explicitní informace (obvykle vzácné) – použití historie prohlížení apod.
 - V maticové faktorizaci je **přítomnost interakce** interpretovaná jako pozitivní signál, zatímco absence znamená nedostatek informace, ne přítomnost nezájmu
 - Např. matice P reprezentuje interakční matici a matice C indikuje konfidenční hladiny interakcí:

$$P = \begin{bmatrix} 1 & 0 & 0 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 & 1 & 1 \\ 1 & 1 & 1 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 & 1 \end{bmatrix} \quad C = \begin{bmatrix} 0.85 & 0 & 0 & 0.34 & 0 & 0.98 \\ 0 & 0.37 & 0.10 & 0 & 0.63 & 0.01 \\ 0.45 & 0.42 & 0.43 & 0 & 0 & 0.23 \\ 0 & 0 & 0.26 & 0 & 0 & 0.88 \end{bmatrix}$$

- Metody kolaborativního filtrování jako **Vážená Regularizovaná Maticová Faktorizace** (WRMF) řídí implicitní feedback přiřazením různých konfidenčních hladin pozorovaným a nepozorovaným interakcím
 - Optimalizační problém:

$$\min_{U,V} \sum_{i,j} C_{ij} (P_{ij} - \mathbf{u}_i^T \mathbf{v}_j)^2 + \lambda \|\mathbf{u}_i\|^2 + \lambda \|\mathbf{v}_j\|^2$$

- **Omezení a výzvy ALS**
 - **Cold Start** – ALS se musí potýkat s novými uživateli nebo položkami, které nemají historii dat o interakci
 - **Extrémní řídkost** – výkon ALS může klesnout, když je interakční matice příliš řídká
 - **Problém se škálováním** – příliš vysoký výkon u velkých systémů kvůli potřebě maticové inverze
 - I přesto je ALS důležitá metoda díky efektivnosti a adaptabilitě

14. Principy bayesovského modelování – pojmy model, apriorní a aposteriorní distribuce. Exponenciální třída distribucí, konjugovaná apriorní a jejich význam v bayesovském odhadu. Příklad konjugovaného apriorní.

NI-BML

- **Apriorní znalost** = počáteční znalost o studované proměnné před zohledněním dat (pozorování, měření)
 - o Vyjádřena **apriorní distribucí** – pravděpodobnost (u diskretních) nebo hustota (u spojitých)
 - o Zahrnutí nových informací → **aposteriorní znalost**
 - o Apriorní i aposteriorní kvantifikují **míru neurčitosti**
- **Vícerozměrné distribuce** – veličiny X, Y
 - o **Sdružené hustoty** $f(x, y) = f(y, x)$
 - o **Marginální hustoty** $f(x) = \int f(x, y) dy$, $f(y) = \int f(x, y) dx$
 - o **Podmíněné hustoty** $f(x|y) = \frac{f(x, y)}{f(y)}$ a naopak
 - o **Řetězové pravidlo** $f(x, y) = f(x|y)f(y) = f(y|x)f(x)$
- **Bayesova věta** – máme náhodné veličiny x, θ s hustotami $f(x|\theta)$ a $\pi(\theta)$:

$$\pi(\theta|x) = \frac{f(x|\theta)\pi(\theta)}{f(x)}, \quad f(x) > 0$$

- o $\pi(\theta|x)$... **aposteriorní podmíněná hustota** X
- o $\pi(\theta)$... **apriorní hustota**
- o $f(x|\theta)$... **model / věrohodnost** (likelihood) dat
- o $f(x)$... marginální hustota x / normální faktor / evidence
 - Jmenovatel je **normalizační faktor** nezávislý na odhadovaném θ a dostaneme ho jednoduše vysčítáním nebo integrací čitatele, $f(x) = \int f(x|\theta)\pi(\theta)d\theta$, zapisujeme často jen proporcionalitu:

$$\pi(\theta|x) \propto f(x|\theta)\pi(\theta)$$

- **Beta distribuce** - $p \in [0, 1] \rightarrow p \sim B(a, b)$
 - o Hustota $\pi(p|a, b) = \frac{1}{B(a, b)} p^{a-1} (1-p)^{b-1}$
 - o Střední hodnota $E[p] = E[p|a, b] = \frac{a}{a+b}$
 - o Beta funkce $B(a, b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$
 - Je v ní obsažená gamma funkce – rozšiřuje faktoriál na komplexní obor
 - o Beta distribuce modeluje pravděpodobnost úspěchu v binomickém rozdělení
 - o Aposteriorní distribuce je zároveň beta distribuce

Bayesovské odhadování

- o y_t ... pozorovaná veličina v okamžiku $t = 0, 1, 2$
- o Vektor hodnot $y_{0:t-1} [y_0 \dots y_{t-1}]$
- o y_t determinováno veličinou x_t a konstantním parametrem θ
- o y_t nezávislé, stejně rozdělené
- o **Apriorní hustota** $\pi(\theta|x_0, y_0)$ – vhodná znalost
 - x_0, y_0 ... pseudodata
- o $f(y_t|x_t, \theta)$ je hustota a $\pi(\theta|x_{0:t-1}, y_{0:t-1})$ apriorní hustota pro θ , pak aposteriorní hustota pro **jednokrokový update** je:

$$\begin{aligned} \pi(\theta|y_{0:t}, x_{0:t}) &= \frac{f(y_t|x_t, \theta) \pi(\theta|x_{0:t-1}, y_{0:t-1})}{\int f(y_t|x_t, \theta) \pi(\theta|x_{0:t-1}, y_{0:t-1}) d\theta} \\ &= \frac{f(y_t|x_t, \theta) \pi(\theta|x_{0:t-1}, y_{0:t-1})}{f(y_t|x_t)} \\ &\propto f(y_t|x_t, \theta) \pi(\theta|x_{0:t-1}, y_{0:t-1}). \end{aligned}$$

- Bez normalizující hustoty ve jmenovateli
- o Sekvenční update po jednom datu je stejný, jako update více daty najednou

Sekvenční odhad

- Ke stávající informaci přidáváme pouze nejnovější data = **update**

Sekvenční bayesovský update v krocích:

1. Použijeme apriorní distribuci
2. Přidáme nové pozorování
3. Dostaneme aposteriorní distribuci
4. Tu opět použijeme jako apriorní pro další pozorování

$$\pi(\theta|x_0, y_0) \xrightarrow[\text{Bayes}]{x_1, y_1} \pi(\theta|x_{0:1}, y_{0:1}) \xrightarrow[\text{Bayes}]{x_2, y_2} \pi(\theta|x_{0:2}, y_{0:2}) \rightarrow \dots \xrightarrow[\text{Bayes}]{x_t, y_t} \pi(\theta|x_{0:t}, y_{0:t}) \rightarrow \dots$$

Bodový odhad θ (neznámý důležitý parametr)

- o Střední hodnota $E[\theta] = E[\theta|x_{0:t}, y_{0:t}]$
- o Modus (MAP odhad), medián

Věrohodnost odhadu – rozptyl $\text{var}\theta$

- Problém bayesovského modelování = odvození aposteriorní distribuce a jejích vlastností

Exponenciální třída distribucí

$$f(y|x, \theta) = h(y, x)g(\theta) \exp[\eta^T T(y, x)]$$

- o $\eta \equiv \eta(\theta)$... přirozený parametr
- o $T(y, x)$... **suficientní statistika** fixního rozměru
 - Obsahuje všechny informace o parametrech modelu
- o $h(y, x)$... známá funkce
- o $g(\theta)$... **normalizační funkce**
- o Pokud $\eta(\theta) = \theta$, třída je **kanonická**

Konjugovaná apriorní distribuce

- o Apriorní distribuce, která má stejný tvar jako aposteriorní distribuce
- o $y|x, \theta$ má rozdělení z **exponenciální třídy distribucí** → apriorní distribuce θ s hyperparametry ξ a ν je k němu konjugovaná, pokud:

$$\pi(\theta) = q(\xi, \nu)g(\theta)^\nu \exp[\eta^T \xi]$$

- ξ má stejný rozměr jako $T(y, x)$
- $\nu \in \mathbb{R}^+$
- $q(\xi, \nu)$... známá funkce
- $g(\theta)$... normalizační funkce

- o Hyperparametry = parametry apriorní

Příklady konjugované distribuce:

Model	Použití	Konjugované apriorno
Normální, známý rozptyl	Všude možně	Normální
Normální, neznámý rozptyl	Všude možně	Normální inverzní – gamma
Bernoulli	Úspěch – neúspěch	Beta
Binomický	Úspěch – neúspěch	Beta
Poissonův	Řídké jevy (částice)	Gamma
Multinomický	Klasifikace do k tříd	Dirichletovo

Bayesovský odhad s konjugovaným apriornem

- Hyperparametry $\xi_{t-1}, \nu_{t-1} \rightarrow$ **Bayesův update**

$$\pi(\theta|y_{0:t}, x_{0:t}) \propto f(y_t|x_t, \theta) \pi(\theta|x_{0:t-1}, y_{0:t-1})$$

je jen triviální **součet**

$$\xi_t = \xi_{t-1} + T(y_t, x_t)$$

$$\nu_t = \nu_{t-1} + 1.$$

- Update pro více dat najednou – taky jen součet
- **Bayesova věta** se při použití konjugovaného apriorního **ztriviuje na přičtení sufficientní statistiky** k hyperparametru ξ_{t-1} a **inkrementaci** hyperparametru ν_{t-1}

15. Stavové modely: rovnice pro vývoj stavu a rovnice měření, rozdíly mezi nimi.
 Bayesovský sekvenční odhad stavových modelů a jejich vliv na apriorní distribuci (znalost).
 Možnosti odhadu stavů v případě nelinearity (pouze vyjmenovat).

NI-BML

Stavový model

- Charakterizován 3 veličinami:
 - o **Stav** x_t – nemůžeme pozorovat, ale můžeme odhadovat
 - Stav je parametr celého modelu, ale vyvíjí se
 - o **Vstup** = řídící veličina u_t – známá
 - o **Výstup** z_t – pozorovatelná veličina determinovaná x_t a u_t
- Model diskrétní (časové okamžiky) nebo spojitý (derivative) – tady jen diskrétní
- **Zápis modelu** – stochastické funkce stavů a vstupů:

$$\begin{aligned}x_t &= f_t(x_{t-1}, u_t) \\ z_t &= g_t(x_t)\end{aligned}$$

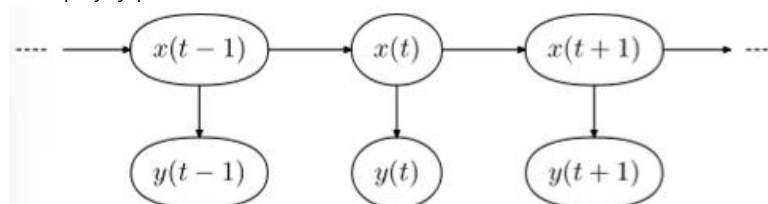
- Lineární systém:

$$\begin{aligned}x_t &= A_t x_{t-1} + B_t u_t + v_t \\ z_t &= H_t x_t + w_t\end{aligned}$$

- o x_t ... stav
- o u_t ... řídící veličina
- o w_t a v_t ... šum stavu a šum měření
- o A_t, B_t, H_t ... matice

Skrytý markovský model – HMM

- **Markovská vlastnost** – markovský proces (nebo markovský model 1. řádu) je model, v němž závisí pouze na stavu předchozím
 - o Popsán pravděpodobnostmi:
 - $p(x_t | x_{t-1})$... pravděpodobnost přechodu
 - $p(x_0)$... počáteční stav
- Skrytý markovský proces – není přímo pozorovatelný, ale lze na něj nahlížet prostřednictvím jiné pozorované veličiny y_t (nebo z_t)
 - o Předpokládá spojitý proces



Kalmanův filtr

- Uvažujme časově invariantní stavový model:

$$\begin{aligned}x_t &= A x_{t-1} + B u_t + v_t \\ z_t &= H x_t + w_t\end{aligned}$$

- o Obě šumové proměnné nezávislé a centrované v 0

$$\begin{aligned}v_t &\sim \mathcal{N}(0, Q) \\ w_t &\sim \mathcal{N}(0, R)\end{aligned}$$
- o Z normality:
 - $x_t \sim \mathcal{N}(A x_{t-1} + B u_t, Q)$ s hustotou $p(x_t | x_{t-1}, u_t)$
 - $z_t \sim \mathcal{N}(H x_t, R)$ s hustotou $f(z_t | x_t)$

- **Apriorní distribuce** pro x_t – model y_t je normální → konjugované apriorno bude tež normální se střední hodnotou x_{t-1}^+ a kovarianční maticí P_{t-1}^+

$$\pi(x_t | z_{0:t-1}, u_{0:t-1}) = \mathcal{N}(x_{t-1}^+, P_{t-1}^+)$$

- Kalmanův filtr má 2 kroky – **predikce stavu a update stavu**
 - o Predikce využívá 1. stavovou rovnici
 - o Update nutný kvůli šumu
 - o Predikovaný stav je apriorní informace do Bayesovy věty, kde je pomocí měření zkorrigován
- **Predikce**
 - o Užití modelu pro vývoj stavu a **předpověď na základě odhadu stavu v předchozím okamžiku**, jaký bude stav nyní
 - Řešíme časový vývoj stavu $x_{t-1} \rightarrow x_t$
 - o Použijeme **odhad reprezentovaný posledním aposteriornem**, nyní apriornem, pro další časový krok a ten proženeme modelem vývoje:
$$\pi(x_t | z_{0:t-1}, u_{0:t}) = \int p(x_t | x_{t-1}, u_t) \pi(x_{t-1} | z_{0:t-1}, u_{0:t-1}) dx_{t-1}$$
 - o **Násobíme 2 normální distribuce a marginalizujeme** \rightarrow máme zas normální distribuci $\mathcal{N}(x_{t-1}^-, P_{t-1}^-)$ s hyperparametry:
$$x_t^- = Ax_{t-1}^- + Bu_t$$

$$P_t^- = AP_{t-1}^- A^T + Q$$
 - Toto **vyčíslí rovnici** pro stav
 - Odhad stavu x_t^- - dosazení do rovnice
 - Kovariance odhadu P_t^- - **míra neurčitosti odhadu**
 - o Byla použita lineární transformace, platí
$$E[AX + B] = AE[X] + B$$

$$\text{var}(AX + B) = A^2 \text{var}(X) = A \text{var}(X) A$$
- **Update**
 - o **Korekce** – opraví predikovaný odhad novými pozorováními y_t
 - Pomocí Bayesovy věty:
$$\pi(x_t | z_{0:t}, u_{0:t}) \propto f(z_t | x_t) \pi(x_t | z_{0:t-1}, u_{0:t})$$
 - o Model se přepisuje do tvaru distribuce z **exponenciální třídy** + vhodné apriorno
 - Pak je bayesovský update jen **součet hyperparametrů a suficientní statistiky**
 - o **Kalmanovo zesílení** – zmenšuje střední kvadratickou chybu (optimální ji minimalizuje)
 - Čím větší je, tím větší je důraz na nová měření
$$K_t = P_t^- H^T (R + H P_t^- H^T)^{-1}$$
- **Rovnice KF se dají rozepsat:**

$$\hat{z}_t^- = H \hat{x}_t^-, \quad (\text{predikce měření})$$

$$\nu_t = z_t - \hat{z}_t^-, \quad (\text{innovace, chyba predikce } z_t)$$

$$S_t = H P_t^- H^T + R, \quad (\text{kovariance inovace } \nu_t)$$

$$K_t = P_t^- H^T S_t^{-1}, \quad (\text{kalmanovské zesílení - gain})$$

$$x_t^+ = \hat{x}_t^- + K_t \nu_t, \quad (\text{aposteriorní odhadu stavu } x_t)$$

$$P_t^+ = (I - K_t H) P_t^-. \quad (\text{aposteriorní kovariance})$$

Nelineární stavové modely

- Zase ve tvaru $x_t = f_t(x_{t-1}, u_t)$ a $z_t = g_t(x_t)$, ale jedna nebo obě **funkce nejsou lineární**
- Možná řešení v závislosti na míře linearity:
 - o **Slabá nelinearita** \rightarrow lokální **linearizace** derivacemi \rightarrow **rozšířený Kalmanův filtr** (EKF)
 - o **Velká nelinearita** \rightarrow **Unscented Kalmanův Filtr** (UKF)
 - o Vždy lze **brute force – Monte Carlo** vzorkování stavů ze souvisejícího stavového prostoru
 - Vygenerování hromady náhodných vektorů x_t
 - Vložení n. v. do stavového modelu
 - Zjištění, které jsou nejpravděpodobnější pro z_t
 - Nejpravděpodobnějším se přiřadí vysoké váhy a hledá se v jejich okolí
 - Např. částicový filtr (**particle filter**)

16. Rejection sampling (RS) a importance sampling (IS): důvody používání RS a IS, jejich základní principy a rozdíly, efektivita práce se vzorky. Stanovení vah v IS a možnosti jejich normování.

NI-BML

Monte Carlo integrace

- Integrál funkce $f(x)$ na $[a, b]$ je vlastně součet malých obdélníků s výškou $f(x_i)$ a délkou strany Δx
- Uvažujeme shodná $\Delta x \rightarrow$ můžeme je vytknout před sumu, rovnoměrně jimi rozdělit interval $[a, b]$ na N pod-intervalů a vlastně se přesunout ke střední hodnotě funkce

$$\int \dots = \frac{b-a}{N} \sum_{i=1}^N f(x_i) = (b-a)\langle f \rangle$$

- o $\langle f \rangle$... odhad střední hodnoty
- o Pro M-rozměrný případ analogicky, integrál je roven $V\langle f \rangle$, kde $V = b-a$ je „objem“ množiny V
- **MC integrace** – využívá princip výše s tím, že body x_i nebere jako deterministicky stanovené body v rovnoměrné síti, ale **vybírá tyto body náhodně z rovnoměrného rozdělení**, středování i objem zachovány
 - o Zákon velkých čísel zajišťuje **konvergenci ke skutečné hodnotě** při $N \rightarrow \infty$
 - Pro $N \rightarrow \infty$ konverguje chyba odhadu $\text{var}(V\langle f \rangle)$ k 0

Rejection sampling

- „accept-reject algoritmus“, využívá se, když nemůžeme použít konjugované apriorno
- Využívá **proposal distribuci** – vhodná distribuce, z níž umíme snadno vzorkovat (nemusí být normovaná)
- $f(x)$... hustota, z níž chceme vzorkovat
- Předpoklad algoritmu:

$$f(x) = \int_0^{f(x)} du = \int_0^{f(x)} \underset{=f(x,u)}{1_{0 < u < f(x)}} du$$

- o Na f je možno nahlížet jako na **marginální hustotu** sdružené distribuce $(x, u) \sim \mathcal{U}\{(x, u): 0 < u < f(x)\}$
- o **Základní teorém vzorkování** – vz. $x \sim f(x)$ je ekvivalentní k vz. $(x, u) \sim \mathcal{U}\{(x, u): 0 < u < f(x)\}$
 - Lze využít oklikou – navzorkovat (x, u) z větší množiny a vybrat takové dvojice, pro něž je podmínka $0 < u < f(x)$ splněná
- **Algoritmus**
 - o Předpokládáme $\int_a^b f(x)dx = 1, m > f(x) \forall x \in [a, b]$
 - o Vzorkujeme $(x', u) \sim \mathcal{U}(0 < u < m)$ tak, že:
 - Navzorkujeme $x' \sim \mathcal{U}(a, b)$
 - Navzorkujeme $u | x = x' \sim \mathcal{U}(0, m)$
 - Vzorek přijmeme, pokud $0 < u < f(x')$
 - o **Lze vylepšit**, bo lze vzorkovat rovnou přes množinu:
$$\mathcal{L} = \{(x', u): 0 < u < m(x'), m(x) \geq f(x')\}$$
 - $m(x)$ není hustota, můžeme ale využít vhodný **proposal** $g(x)$, který hustotou bude, a vhodné kladné číslo $M > 0$ a nastavit $m(x) = Mg(x)$
 - o Takhle dostaneme **lepší algoritmus**:
 - Navzorkujeme $x' \sim g$
 - Navzorkujeme $u \sim \mathcal{U}(0, 1)$
 - Vzorek x' přijmeme, když $u \leq \frac{f(x')}{Mg(x')}$
 - Kritérium přijetí je stejné, jako v základní metodě
 - Pravděpodobnost přijetí je $1/M$ – čím blíže je M k vzorkované hustotě, tím vyšší je četnost přijatých vzorků
- Rejection sampling je **efektivnější než MC**, ale **problémy s** efektivitou v oblastech, kde je **hustota koncentrovaná na malou podmnožinu**, nebo tam, kde nabývá moc nízkých hodnot a většina vz. není přijata (a taky nepotřebujeme integrovat)

Importance sampling

- Problém s efektivitou **rejection samplingu** – kompenzuje fakt, že k vzorkování používá jinou hustotu tím, že **přijímá** či **zahazuje** vzorky
- **Importance sampling** – kompenzuje jinou vzorkovací hustotu **přidělováním vah** všem vzorkům
- $f(x)$ cílová (komplikovaná) hustota a $g(x)$ proposal hustota, pak:

$$\int f(x)dx = \int g(x) \frac{f(x)}{g(x)} dx = \int g(x) w(x) dx$$

$= w(x)$

- o Podmínka $g(x) > 0$ tam, kde $f(x) > 0$
- o Výpočet vah – vložení hodnoty x do známých funkcí $f()$ a $g()$
- o Zobecnění – pro střední hodnotu $E[x]$ při $f(x)$:

$$\mathbb{E}_f[x] = \int x \cdot f(x) dx = \int x \cdot \frac{f(x)}{g(x)} g(x) dx = \int x \cdot w(x) g(x) dx = \mathbb{E}_g[x \cdot w(x)]$$

- Podmínka $g(x) > 0$ tam, kde $xf(x) \neq 0$

- **Odhad střední hodnoty** = vážený průměr vzorků:

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i w(x_i)$$

- **Normované váhy** – protože váhy z předchozí rovnice nejsou normované a nesčítají se do 1

- o **Normování** zde zastupuje faktor $1/N$ a mohli bychom tedy uvažovat:

$$w'(x_i) = \frac{w(x_i)}{N} \quad \text{kde} \quad \sum_{i=1}^N w'(x_i) \neq 1$$

- o Z praktických důvodů často použít **normované varianty vah**

$$W(x_i) = \frac{w(x_i)}{\sum_{i=1}^N w(x_i)}$$

- o **Odhad střední hodnoty** je pak vážený průměr v podobě:

$$\hat{\mu} = \sum_{i=1}^N x_i W(x_i)$$

- Odhad sice vychýlený, ale může mít nižší varianci a nezávisí na normalizační konstantě

- **Algoritmus IS** – je o **určování vah**, které pak používáme např. při výpočtu integrálů

- o Nageneryjeme N vzorků z proposal hustoty $x_i \sim g(x)$
- o Spočteme hodnotu hustoty $f(x_i)$
- o Spočteme hodnotu hustoty $g(x_i)$
- o Spočteme váhy $w(x_i) = \frac{f(x_i)}{g(x_i)}$
- o Váhy normalizujeme, buď $W(x_i) = w(x_i) / \sum w(x_i)$, nebo $w'(x_i) = w(x_i) / N$

17. QR rozklad: metody výpočtu, použití při výpočtu odhadu metodou nejmenších čtverců, QR algoritmus pro hledání vlastních čísel.

NI-PON

Ortogonalní matice

- **Ortogonalní vektory** – takové, jejichž skalární součin je 0: $(x, y) = x^T y = 0$
- **Ortogonalní matice** $Q \in \mathbb{R}^{m,n}$ – její matice jsou navzájem ortogonální a mají jednotkovou velikost – pro i -tý a j -tý sloupec musí platit:

$$q_i^T q_j = \begin{cases} 0 & \text{pro } i \neq j, \\ 1 & \text{pro } i = j. \end{cases}$$

- Důsledky:
 - o Pro ortogonalní Q platí, že $Q^T Q = I_n$
 - o Je-li Q (čtvercová) ortogonální, je její **transpozice i její inverzní**
 - o **Součin** ortogonálních matic je ortogonální matice:

$$(Q_1 Q_2 \cdots Q_n)^T Q_1 Q_2 \cdots Q_n = I_m$$

- o **Násobení** ortogonální maticí nemění **normu**:

$$\|Qx\| = \|x\|$$

Givensovy rotace

- **QR rozklad** = přepis matice A na součin ortogonální matice Q a horní trojúhelníkové matice R

$$A = QR$$

- Platí $Q^T A = R$ – matice R vznikne ze stejně velké matice A vynásobením ortogonální Q^T , která vynuluje v matici A všechny složky matice A pod diagonálou
 - o Složitě \rightarrow nacházíme ortogonální matice $Q_1 \dots Q_k$, které postupně zvětšují počet nul pod diagonálou A :
 - $Q_1 A$ má aspoň o 1 nulu pod diagonálou víc než A
 - $Q_2 Q_1 A$ má aspoň o 1 nulu pod diagonálou víc než $Q_1 A$
 - ...
 - $Q_k \dots Q_2 Q_1 A$ má pod diagonálou samé 0 a je rovna R
 - o Tzn. Hledáme ortogonální matice Q_i – součin o. matic je o. matice
- **Givensovy rotace** – matice Q_i jsou jednoduché – jen na 4 místech se liší od jednotkové matice – rotace 2 složek vektoru, roztažení rotace ve 2 dimenzích
 - o Hledáme ortogonální matici $S \in \mathbb{R}^{2,2}$, která otočí nenulový vektor $x = (a, b)^T$ tak, že výsledný vektor leží na ose $x \rightarrow 2.$ složka je 0
 - o 1. sloupec matice S : $(\alpha, \beta)^T$
 - o Celá matice S :

$$S = \begin{pmatrix} \alpha & -\beta \\ \beta & \alpha \end{pmatrix}$$

- o Platí $\alpha^2 + \beta^2 = 1 \rightarrow$ vlastnost rotace na osu x znamená, že musí pro vektor x platit:

$$\begin{pmatrix} \alpha & -\beta \\ \beta & \alpha \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} \pm \|x\| \\ 0 \end{pmatrix} = \begin{pmatrix} \pm \sqrt{a^2 + b^2} \\ 0 \end{pmatrix}$$

- 1. složka je $\pm \|x\|$, protože výsledný vektor musí mít stejnou normu jako původní
 - Otáčíme na kladnou nebo zápornou část osy x , podle toho znaménko
- o \rightarrow podmínka na parametry α, β : $a\beta + b\alpha = 0 \Leftrightarrow a\beta = -b\alpha$, což vede (i z definice) na soustavu:

$$a^2 \beta^2 = b^2 \alpha^2$$

$$\alpha^2 + \beta^2 = 1$$

$$\rightarrow \alpha = \frac{a}{\sqrt{a^2 + b^2}} = \cos \varphi, \beta = -\frac{b}{\sqrt{a^2 + b^2}} = -\sin \varphi$$

- QR rozklad pomocí rotací – hledáme QR rozklad $A \in \mathbb{R}^{5,3}$ a chceme matici, která vynuluje A_{21}

$$A = \begin{pmatrix} a_1 & \bullet & \bullet \\ b_1 & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{pmatrix}$$

- o Výpočet parametrů podle vzorců výše → získáme ortogonální Q_1 , která vynuluje b_1 :

$$\underbrace{\begin{pmatrix} \alpha_1 & -\beta_1 & 0 & 0 & 0 \\ \beta_1 & \alpha_1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix}}_{Q_1} \underbrace{\begin{pmatrix} a_1 & \bullet & \bullet \\ b_1 & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{pmatrix}}_{A} = \begin{pmatrix} \bullet & \bullet & \bullet \\ 0 & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{pmatrix} = \underbrace{\begin{pmatrix} a_2 & \bullet & \bullet \\ 0 & \bullet & \bullet \\ b_2 & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{pmatrix}}_{Q_1 A}$$

- Oranžová = prvek byl vynásobením změněn

- o V $Q_1 A$ provedeme zase stejný výpočet parametrů a podle toho zvolíme Q_2 :

$$\underbrace{\begin{pmatrix} \alpha_2 & 0 & -\beta_2 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ \beta_2 & 0 & \alpha_2 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix}}_{Q_2} \underbrace{\begin{pmatrix} a_2 & \bullet & \bullet \\ 0 & \bullet & \bullet \\ b_2 & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{pmatrix}}_{Q_1 A} = \begin{pmatrix} \bullet & \bullet & \bullet \\ 0 & \bullet & \bullet \\ 0 & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{pmatrix} = \underbrace{\begin{pmatrix} a_3 & \bullet & \bullet \\ 0 & \bullet & \bullet \\ 0 & \bullet & \bullet \\ b_3 & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{pmatrix}}_{Q_2 Q_1 A}$$

- o Analogicky Q_3, Q_4 – vynulování prvků b_3, b_4 :

- o 1. sloupec hotový – v 2. sloupci nulujeme b_2 :

- Q_5 mění jen 2. a 3. řádek – v 1. sloupci jsou 0, takže nevyrobíme ne0 tam, kde už byly

$$\underbrace{\begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & \alpha_5 & -\beta_5 & 0 & 0 \\ 0 & \beta_5 & \alpha_5 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix}}_{Q_5} \underbrace{\begin{pmatrix} \bullet & \bullet & \bullet \\ 0 & a_5 & \bullet \\ 0 & b_5 & \bullet \\ 0 & \bullet & \bullet \\ 0 & \bullet & \bullet \end{pmatrix}}_{Q_4 Q_3 Q_2 Q_1 A} = \begin{pmatrix} \bullet & \bullet & \bullet \\ 0 & \bullet & \bullet \\ 0 & 0 & \bullet \\ 0 & \bullet & \bullet \\ 0 & \bullet & \bullet \end{pmatrix} = \underbrace{\begin{pmatrix} \bullet & \bullet & \bullet \\ 0 & a_6 & \bullet \\ 0 & 0 & \bullet \\ 0 & b_6 & \bullet \\ 0 & \bullet & \bullet \end{pmatrix}}_{Q_5 Q_4 Q_3 Q_2 Q_1 A}$$

- o Postupně nulujeme další: $A = QR$, kde $Q = Q_1^T Q_2^T \dots Q_8^T Q_9^T$

$$Q_9 Q_8 \dots Q_2 Q_1 A = \begin{pmatrix} \bullet & \bullet & \bullet \\ 0 & \bullet & \bullet \\ 0 & 0 & \bullet \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} = R$$

- o To je už ten QR rozklad:

- Celkový odhad složitosti = $O(mn^2)$
- Je to vlastně jako GEM, ale s numerickou stabilitou díky ortogonálním operacím

Householderovy reflexe

- Zrcadlení vektoru v zrcadle, které reprezentuje nadrovina procházející počátkem souřadnice
 - o Např. \mathbb{R}^3 - hledáme 2 rozměrné zrcadlo, které nastavíme vektoru tak, aby odraz ležel na ose x
 - o Tím se vynulují všechny složky kromě 1. – rovna normě vektoru
 - o Hledáme **ortogonální matice P**, pro které:

$$Px = P \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} \pm \|x\| \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

- Zrcadlo v n -dimenzionálním prostoru = nadrovina procházející počátkem souřadnic
 - o Popsatelná pomocí normálového vektoru u , který je na ni kolmý (ortogonální), předp. $\|u\| = 1$
 - o Rovina = všechny vektory x ortogonální s u : $x^T u = 0$
 - Lineární rovnice pro n neznámých – množina řešení je podprostor dimenze $n - 1$

- Zrcadlení

- Vzdálenost x od $u = d \rightarrow$ bod na zrcadle nejbližší vektoru x = vektor $x - du$
- Zrcadlový obraz x získáme posunutím o stejnou vzdálenost za zrcadlo:

$$z.o.x = \tilde{x} = x - 2du$$

- Vzdálenost $d = u^T x \rightarrow \tilde{x} = x - 2u^T x u = x - 2uu^T x = (I - 2uu^T)x$
 - $P = I - 2uu^T$
 - P reprezentuje zrcadlení v zrcadle s normálovým jednotkovým vektorem u
 - P je rovná svojí transpozici \rightarrow je symetrická \rightarrow je ortogonální
 - Zrcadlení nemění velikost – když aplikujeme zrcadlo 2x, dostaneme se tam, kde jsme začali

- Volba vektoru u , aby zrcadlil x na osu odpovídající 1. souřadnici – aby platilo

$$Px = (I - 2uu^T) \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} \alpha \\ 0 \\ \vdots \\ 0 \end{pmatrix} = \alpha \underbrace{\begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}}_{e_1}$$

- $\alpha = \pm \|x\|$

$$x - 2uu^T x = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & & & & \\ 0 & & I_4 - 2u_2 u_2^T & & \\ 0 & & & & \\ 0 & & & & \end{pmatrix} u^T x = x - \alpha e_1.$$

- = vektor u je násobek vektoru $x - \alpha e$
- Velikost $u = 1 \rightarrow$ musí platit

$$u = \pm \frac{x - \alpha e}{\|x - \alpha e\|} = \pm \left(\frac{x_1 - \alpha}{\|x - \alpha e\|} \quad \frac{x_2}{\|x - \alpha e\|} \quad \cdots \quad \frac{x_n}{\|x - \alpha e\|} \right)^T$$

- Možnost vybrat si znaménko = zrcadlo má 2 opačné normálové vektory

$$\|x - \alpha e\| = \sqrt{(x_1 - \alpha)^2 + x_2^2 + \cdots + x_n^2} = \sqrt{\alpha^2 - 2\alpha x_1 + \underbrace{x_1^2 + x_2^2 + \cdots + x_n^2}_{=\|x\|^2 = \alpha^2}} = \sqrt{2\alpha(\alpha - x_1)}.$$

- \rightarrow složky u :

$$u_1 = \frac{x_1 - \alpha}{\sqrt{2\alpha(\alpha - x_1)}} = -\frac{\sqrt{\alpha - x_1}}{\sqrt{2\alpha}}$$

$$u_i = \frac{x_i}{\sqrt{2\alpha(\alpha - x_1)}} = \left\{ \text{platí } \sqrt{\alpha - x_1} = -\sqrt{2\alpha}u_1 \right\} = -\frac{x_i}{2\alpha u_1}.$$

- QR rozklad pomocí zrcadlení – analogický k ortogonálním rotacím

$$A = \begin{pmatrix} \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{pmatrix} = \begin{pmatrix} x_1^{(1)} & \bullet & \bullet \\ x_2^{(1)} & \bullet & \bullet \\ x_3^{(1)} & \bullet & \bullet \\ x_4^{(1)} & \bullet & \bullet \\ x_5^{(1)} & \bullet & \bullet \end{pmatrix} \rightarrow P_1 x^{(1)} = \begin{pmatrix} \alpha_1 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

- $\alpha_1 = \pm \|x^{(1)}\|$

$$\rightarrow P_1 A = P_1 \begin{pmatrix} \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet \end{pmatrix} = \underbrace{\begin{pmatrix} \alpha_1 & \bullet & \bullet \\ 0 & \bullet & \bullet \\ 0 & \bullet & \bullet \\ 0 & \bullet & \bullet \\ 0 & \bullet & \bullet \end{pmatrix}}_{P_1 A} = \begin{pmatrix} \bullet & \bullet & \bullet \\ 0 & x_1^{(2)} & \bullet \\ 0 & x_2^{(2)} & \bullet \\ 0 & x_3^{(2)} & \bullet \\ 0 & x_4^{(2)} & \bullet \end{pmatrix}$$

- 2. krok – vektor $x^{(2)}$ má jen 4 složky $\rightarrow P_2$ bude taková, aby se k 1. souřadnici chovala jako jednotková matice a zrcadlila pouze zbylé 4

- Vektor $u_2 \in \mathbb{R}^4$ tak, aby matice uvnitř zrcadlila $x^{(2)}$ na $(\alpha_2, 0, 0, 0)^T$
- 3. krok:

$$P_2 P_1 A = P_2 \underbrace{\begin{pmatrix} \bullet & \bullet & \bullet \\ 0 & x_1^{(2)} & \bullet \\ 0 & x_2^{(2)} & \bullet \\ 0 & x_3^{(2)} & \bullet \\ 0 & x_4^{(2)} & \bullet \end{pmatrix}}_{P_1 A} = \underbrace{\begin{pmatrix} \bullet & \bullet & \bullet \\ 0 & \alpha_2 & \bullet \\ 0 & 0 & x_1^{(3)} \\ 0 & 0 & x_2^{(3)} \\ 0 & 0 & x_3^{(3)} \end{pmatrix}}_{P_2 P_1 A}$$

- 4. krok - P_3 – vyrobí poslední 2 nuly analogicky
- Finální QR rozklad:

$$A = QR, \quad \text{kde } R = P_3 P_2 P_1 A \quad \text{a} \quad Q = P_1^T P_2^T P_3^T.$$

- Matic P_i potřeba tolik, kolik je sloupců matice $A = n$
- Složitost = $O(m^2 n)$ (# řádků = m , # sloupců = n)

QR rozklad a OLS

- **OLS = ordinary least squares** – obyčejné nejmenší čtverce
- Hledáme **minimum výrazu** $\|Xw - Y\|$, rozměry matice X jsou $N, p + 1$ a má lineárně nezávislé sloupce
- Najdeme **QR rozklad** matice X : $X = QR = (Q_L \quad Q_R) \begin{pmatrix} R_S \\ \theta \end{pmatrix}$
 - X a R mají stejnou hodnotu, protože násobení regulární maticí nemění hodnotu a Q je regulární
- **Výpočet OLS** – místo normy pracujeme s jejím kvadrátem, abychom se zbavili odmocniny

$$\|Xw - Y\|^2 = \|QRw - Y\|^2 = \|Q^T(QRw - Y)\|^2 = \|Rw - Q^T Y\|^2 = \dots$$

$$\dots = \|Rw - Q^T Y\|^2 = \left\| \begin{pmatrix} R_S w \\ \theta \end{pmatrix} - \begin{pmatrix} Q_L^T Y \\ Q_R^T Y \end{pmatrix} \right\|^2 = \left\| \begin{pmatrix} R_S w - Q_L^T Y \\ -Q_R^T Y \end{pmatrix} \right\|^2 = \dots$$

$$\dots = \left\| \begin{pmatrix} R_S w - Q_L^T Y \\ -Q_R^T Y \end{pmatrix} \right\|^2 = \|R_S w - Q_L^T Y\|^2 + \|Q_R^T Y\|^2$$

$$\|Xw - Y\|^2 = \|R_S w - Q_L^T Y\|^2 + \|Q_R^T Y\|^2$$

- Hledat výraz minimalizující to poslední vlevo je to stejné jako hledat w minimalizující výraz vpravo
 - Vpravo druhý sčítanec nezávisí na w , takže ho můžeme vynechat, protože neovlivní, kde bude maximum, takže řešíme úlohu:

$$\hat{w}_{OLS} = \underset{w \in \mathbb{R}^{p+1}}{\operatorname{argmin}} \|Xw - Y\| = \underset{w \in \mathbb{R}^{p+1}}{\operatorname{argmin}} \|R_S w - Q_L^T Y\|$$

- R_S je regulární, takže rovnice $R_S w = Q_L^T Y$ má právě jedno řešení
 - Protože je norma nezáporný výraz, je řešení nutně rovné hledanému minimu = vektoru \hat{w}_{OLS}
- Máme-li **QR rozklad** matice X , stačí pro nalezení odhadu metodou nejmenších čtverců vynásobit Y maticí Q_L^T a vyřešit **soustavu s horní trojúhelníkovou maticí**
 - Provedení zpětného chodu a postupně odzadu spočítání složek \hat{w}_{OLS}

QR algoritmus

- **QR algoritmus** = algoritmus pro hledání vlastních čísel
 - Využití toho, že podobné matice mají stejná vlastní čísla
 - Konstruujeme posloupnost matic, které jsou podobné matici, jejíž vlastní čísla hledáme
- **Algoritmus:**
 - Máme čtvercovou $A \in \mathbb{R}^{n,n}$ a hledáme její vlastní čísla
 - 1. Najdeme QR rozklad: $A = QR$
 - Matice A, R si nejsou podobné, ale podobné jsou si A a RQ
 - 2. $A = QR \Rightarrow R = Q^T A \Rightarrow RQ = Q^T A Q$
 - Díky ortogonalitě $Q^T = Q^{-1} \rightarrow A$ a RQ jsou podobné \rightarrow stejná vlastní čísla

3. Přes 2. konstruuji **posloupnost matic podobných matici A** :

$$A_0 = A = Q_0 R_0$$

$$A_i = R_{i-1} Q_{i-1} = Q_i R_i$$

→ konstruuujeme posloupnost matic $(A_i)_{i \geq 0}$, které jsou si navzájem podobné

- Často konvergence k horní trojúhelníkové matici
- **Zrychlení konvergence** pomocí Hessenbergovy formy

- **Hessenbergova forma matice** – matice $A \in \mathbb{R}^{n,n}$ je v H. formě, pokud pro $1 \leq i, j \leq n$ platí

$$i > j + 1 \Rightarrow A_{i,j} = 0$$

- Matice A je navíc **tridiagonální**, jestliže platí

$$(i > j + 1 \text{ nebo } i < j - 1) \Rightarrow A_{i,j} = 0$$

$$\begin{pmatrix} \bullet & \bullet & 0 & 0 \\ \bullet & \bullet & \bullet & 0 \\ 0 & \bullet & \bullet & \bullet \\ 0 & 0 & \bullet & \bullet \end{pmatrix}$$

- K libovolné matici $A \in \mathbb{R}^{n,n}$ existuje matice $B \in \mathbb{R}^{n,n}$ v Hessenbergově formě, která je A podobná. Matici podobnosti P lze navíc volit ortogonální a symetrickou:

$$B = PAP$$

- Pokud je matice A symetrická, můžeme navíc B volit tridiagonální.
- Matice P vznikne součinem $n - 2$ matic P_i provádějících Householderovy reflexe

- **Výpočet** singulárních hodnot – vlastních čísel symetrické matice $A^T A$, $A \in \mathbb{R}^{m,n}$

1. Matici $A^T A$ převedeme pomocí $n - 2$ matic P_i provádějících Householderovy reflexe na podobnou matici B , která je tridiagonální.
2. Na matici $B_0 = B$ aplikujeme QR algoritmus.
3. Najdeme QR rozklad $B_0 = Q_0 R_0$, ten můžeme najít pomocí $n - 1$ Givensových rotací, které nulují $n - 1$ prvků pod diagonálou. Násobení maticí provádějících Givensovu rotaci vždy ovlivňuje pouze 6 prvků matice a je tedy výpočetně nenáročný
4. Posloupnost matic $(B_i)_{i \geq 0}$:

$$B_i = R_{i-1} Q_{i-1}, \text{ její rozklad označíme } B_i = Q_i R_i$$

- Matice B_i zůstávají tridiagonální

18. Maticové faktorizace pomocí SVD, její výpočet, vlastnosti a použití ve strojovém učení: souvislost s metodou hlavních komponent (PCA)

NI-PON

SVD rozklad

- SVD = Singular value decomposition
- **Hlavní myšlenka** SVD rozkladu – máme matici $A \in \mathbb{R}^{m,n}$ – nehledáme její vlastní čísla, ale vl. č. matice $A^T A$
 - o $A^T A$ je čtvercová $n \times n$, **pozitivně semidefinitní, symetrická**
 - o Symetrická \rightarrow **diagonalizovatelná**, má reálná (nezáporná) vlastní čísla a z vlastních vektorů lze vytvořit ortonormální bázi prostoru \mathbb{R}^n
 - o $h(A) = r$, hodnota nesmí být větší než počet sloupců/řádků ($r \leq \min\{m, n\}$). $h(A^T A) = r$
 - o $A^T A$ má n vlastních čísel. Jelikož hodnota je r , je 0 vlastní číslo právě když $r < n$. V takovém případě má vlastní číslo 0 násobnost $n - r$.
 - o Matice má tedy **r kladných vlastních čísel**, které seřadíme od největšího po nejmenší a označíme σ_i^2 : $\sigma_1^2 \geq \sigma_2^2 \geq \dots \geq \sigma_r^2 > 0$.
 - o Ke každému vlastnímu číslu umíme najít lineárně nezávislých (+ ortonormálních) vlastních vektorů, kolik je násobnost daného vlastního čísla. Získáme soubor pro všechna $i = 1, 2, \dots, r$

$$(v_1, v_2, \dots, v_r), \quad \text{kde} \quad A^T A v_i = \sigma_i^2 v_i$$

- o Je-li $r = n$, tvoří tento soubor bázi \mathbb{R}^n . Je-li $r < n$, můžeme k němu připojit $n - r$ **vlastních vektorů** v příslušejících k vlastnímu číslu 0 a vytvořit tak ortonormální bázi i tak
- **Odvození SVD rozkladu**
 - o Pro každé $i = 1, 2, \dots, r$ definujeme vektor u_i :

$$u_i = \frac{1}{\sigma_i} A v_i$$

\rightarrow soubor vektorů, o kterém platí:

- Je ortonormální:

$$u_i^T u_j = \begin{cases} 0 & \text{pro } i \neq j \\ 1 & \text{pro } i = j \end{cases}$$

- Pro každé $i = 1, 2, \dots, r$:

$$A A^T u_i = \sigma_i^2 u_i$$

$\rightarrow u_i$ je **vlastní vektor** čtvercové, symetrické a pozitivně semidefinitní matice $A A^T$ příslušející jejímu kladnému vlastnímu číslu σ_i^2

$\rightarrow A^T A$ a $A A^T$ mají stejná kladná vlastní čísla

- o Ze souborů $(v_1 \dots v_n)$ a $(u_1 \dots u_m)$ vyrobíme matice $V \in \mathbb{R}^{n,n}$, $U \in \mathbb{R}^{m,m}$ tak, že tyto soubory napíšeme jako jejich sloupce.
- o Jelikož jsou oba soubory ortonormální, jsou obě matice ortogonální:

$$V^T V = I_n, U^T U = I_m$$

- o $\Sigma \in \mathbb{R}^{m,n}$ matice, která má na **diagonále postupně vlastní č. $\sigma_1 \dots \sigma_r$** , takto seřazená podle velikosti
 - Je-li $r < \min\{m, n\}$, **doplníme na diagonálu nuly** (doplnění nejmenšího vlastního čísla 0)
 - Např. Pokud $m > n = r + 1$ ($\rightarrow h(A) = n - 1$), bude matice Σ :

$$\begin{pmatrix} \sigma_1 & 0 & 0 & \dots & 0 & \dots & 0 \\ 0 & \sigma_2 & 0 & \dots & 0 & \dots & 0 \\ 0 & 0 & \sigma_3 & \dots & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & \sigma_r & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & & \vdots & \dots & \vdots \\ 0 & 0 & 0 & \dots & 0 & \dots & 0 \end{pmatrix} \in \mathbb{R}^{m,n}$$

\rightarrow s tímto značením a jelikož pro $i = 1, 2, \dots, r$: $A v_i = \sigma_i u_i$ a pro $i = r + 1, r + 2, \dots, n$: $A v_i = 0$ dostaneme:

$$A V = U \Sigma$$

- **Definice SVD rozkladu** (fucking finally) – Buď $A \in \mathbb{R}^{m,n}$ s hodnotí r . Potom existují ortogonální matice $V \in \mathbb{R}^{n,n}$, $U \in \mathbb{R}^{m,m}$ a diagonální matice $\Sigma \in \mathbb{R}^{m,n}$, která má na diagonále kladná čísla

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$$

doplňená příslušným počtem nul, takové, že

$$A = U\Sigma V^T$$

Tomuto rozkladu říkáme **SVD rozklad** a číslům σ_i pak **singulární hodnoty** matice A .

- o Singulární hodnoty σ_i jsou odmocniny vlastních čísel matic $A^T A$ a AA^T (proto ta nezápornost, AHA)
- o Sloupce matic V, U jsou tvořeny příslušnými vlastními vektory tvořícími ortonormální bázi
- SVD rozklad určen jednoznačně až na:
 - o Singulární hodnoty určeny jednoznačně, ale u vl. vektorů si můžeme vybrat znaménko (bo normované na jedničku)
 - o Pokud má singulární hodnota více vlastních vektorů (jako vlastní číslo má vyšší násobnost než 1), můžeme libovolně měnit pořadí těchto vektorů v matici V . Sloupce v U se pak dopočítají
- Supr trik – v maticích U, V, Σ lze uvažovat jen prvních r sloupců a řádků, ostatní jsou nulové
 $\rightarrow A = U_r \Sigma_r V_r^T$
 $\rightarrow A$ lze zrekonstruovat z těchto r sloupců
 - o Dokonce platí, že podprostor generovaný sloupci matice A = lineární obal souboru $(u_1 \dots u_r)$

Aproximace maticí s nižší hodnotí

- Když vynásobíme matice $m \times d$ a $d \times n$, dostaneme matice $m \times n$ – lze využít ke kompresi
- Mějme matice $A \in \mathbb{R}^{m,n}$, jaké je nejmenší d takové, že existují $B \in \mathbb{R}^{m,d}$ a $C \in \mathbb{R}^{d,n}$ pro které

$$A = BC$$

- Hodnota součinu nemůže být vyšší, než hodnota jednotlivých matic, $h(A) = r$ a $h(B) = h(C) = d$, musí být $d \geq r$ (kdyby ne, BC by mělo nižší hodnotu než A a nefungovalo by to). B a C lze najít pro $d = r$:

$$B = U_r, C = \Sigma_r V_r^T$$

\rightarrow nejmenší d je rovno $r = h(A)$

- Máme-li zadané d , jaké jsou matice $B \in \mathbb{R}^{m,d}$ a $C \in \mathbb{R}^{d,n}$ takové, že jejich součin BC je co nejbližší $A \in \mathbb{R}^{m,n}$?

\rightarrow SVD

- **Vzdálenost** 2 stejně velkých matic = norma jejich rozdílu: $\|A - \bar{A}\|$
 - o **Frobeniova norma** matice $A \in \mathbb{R}^{m,n}$ = suma čtverců odchylek v jednotlivých složkách

$$\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n a_{ij}^2}$$

- o Souvislost s SVD – je to odmocnina ze součtu kvadrátů singulárních hodnot

$$\|A\|_F = \sqrt{\sum_{i=1}^r \sigma_i^2}$$

- **Eckart-Young-Mirskyho věta** – Mějme matice $A \in \mathbb{R}^{m,n}$ hodnotí r a buď $U\Sigma V^T$ její SVD rozklad. Označme U_k, Σ_k, V_k sloupcové matice, které vznikly z U, Σ, V tak, že jsme vzali pouze jejich prvních k sloupců. Potom pro každé kladné přirozené číslo $d \leq r$ je řešením následující úlohy s neznámou $\bar{A} \in \mathbb{R}^{m,n}$

$$\text{minimalizujte } \|A - \bar{A}\|_F, \text{ kde } h(\bar{A}) \leq d$$

matice

$$A_d^* = U_d \Sigma_d V_d^T$$

- Neboli:

$$\underset{h(\bar{A}) \leq d}{\operatorname{argmin}} \|A - \bar{A}\|_F = U_d \Sigma_d V_d^T = \sum_{i=1}^d \sigma_i u_i v_i^T,$$

$\rightarrow d$ největších singulárních hodnot a k nim příslušné sloupce matic U, V definují nejlepší aproximaci maticí o hodnotí d

- Zbylých $r - d$ singulárních hodnot říká, jak dobrá je to aproximace:

$$\|\mathbf{A} - \mathbf{A}_d^*\|_F = \sqrt{\sum_{i=d+1}^r \sigma_i^2}.$$

- Místo Frobeniovy normy můžeme použít klasickou L_2 vektorovou normu:

$$\|\mathbf{A}\|_2 = \sup_{\mathbf{x} \in \mathbb{R}^n} \frac{\|\mathbf{A}\mathbf{x}\|_2}{\|\mathbf{x}\|_2}$$

= **spektrální norma**

- Rovná se největší singulární hodnotě matice \mathbf{A} $\|\mathbf{A}\|_2 = \sigma_1$

PCA a SVD

- SVD u PCA pro **redukci dimenzionality**
- PCA se počítá pro matici $\mathbf{X} \in \mathbb{R}^{N,p}$ (hodnost r), o které předpokládáme, že je **vycentrovaná** = průměry příznaků ve sloupcích jsou 0

$$0 = \frac{1}{N} \sum_{i=1}^N X_{ij} \quad \text{pro } j = 1, \dots, p$$

- Docílení **posunutím** dat o **průměry příznaků**
- Pro takovou matici lze spočítat matici kovariancí příznaků:

$$\frac{1}{p-1} \mathbf{X}^T \mathbf{X}$$

- **Hlavní komponenty** jsou dány **vlastními vektory** této matice a jejich **rozptyly** jsou příslušná **vlastní čísla**
- Souvislost se SVD rozkladem $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$:

$$\frac{1}{p-1} \mathbf{X}^T \mathbf{X} = \frac{1}{p-1} \mathbf{V} \mathbf{\Sigma}^T \mathbf{U} \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T = \frac{1}{p-1} \mathbf{V} \mathbf{\Sigma}^T \mathbf{\Sigma} \mathbf{V}^T$$

- Platí, že:

$$\mathbf{\Sigma}^T \mathbf{\Sigma} = \begin{pmatrix} (\sigma_1)^2 & 0 & 0 & \dots & 0 & \dots & 0 \\ 0 & (\sigma_2)^2 & 0 & \dots & 0 & \dots & 0 \\ 0 & 0 & (\sigma_3)^2 & \dots & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & (\sigma_r)^2 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 0 & \dots & 0 \end{pmatrix} \in \mathbb{R}^{p,p}$$

je čtvercová diagonální matice

→ dosazením jsme dostali **spektrální rozklad kovarianční matice** (vztah podobnosti matice s diagonální maticí)

- Vyplývá vztah mezi PCA a SVD:
 - **Hlavní komponenty** matice \mathbf{X} odpovídají **prvním r sloupcům** matice \mathbf{V}
 - Jejich **rozptyly**, navíc sestupně seřazené, najdeme jako **kvadráty singulárních hodnot** na diagonále $\mathbf{\Sigma}^T \mathbf{\Sigma}$, které stačí vydělit číslem $p - 1$

19. Hladká optimalizace (bez vazeb), spádové metody, volba směru a délky kroku.

NI-PON

- Totální derivace

- Mějme funkci $f: D \rightarrow \mathbb{R}^m$, kde $D \subset \mathbb{R}^n$ je otevřená množina. Řekneme, že f je **(totálně) diferencovatelná** v bodě $a \in D$, pokud existuje lineární zobrazení $\lambda: \mathbb{R}^n \rightarrow \mathbb{R}^m$ takové, že:

$$\lim_{h \rightarrow 0} \frac{\|f(a+h) - f(a) - \lambda(h)\|}{\|h\|} = 0.$$

- Pokud je f diferencovatelná v a , pak takové lineární zobrazení existuje právě jedno = **(totální) derivace** funkce f bodě $a = Df(a)$ nebo $f'(a)$.
- Jacobiho matice** = matice lineárního zobrazení $Df(a)$
- Mějme funkci $f: D \rightarrow \mathbb{R}^m$, kde $D \subset \mathbb{R}^n$ a $a \in D$. Pokud $\partial_j f_i(a)$ existuje a je spojitá na otevřeném okolí bodu a pro všechna $i, j: 1 \leq i \leq m, 1 \leq j \leq n$ potom $Df(a)$ existuje = **spojitě diferencovatelná** funkce

- Řetězové pravidlo

- Mějme funkci $g: D_g \rightarrow \mathbb{R}^m$, kde $D_g \subset \mathbb{R}^n$, a $f: D_f \rightarrow \mathbb{R}^q$, kde $D_f \subset \mathbb{R}^m$. Je-li funkce f diferencovatelná v $a \in D_g$ a funkce f diferencovatelná v $g(a) \in D_f$, potom je funkce $f \circ g: D_g \rightarrow \mathbb{R}^q$ diferencovatelná v a a platí: $D(f \circ g)(a) = Df(g(a)) \circ Dg(a)$

- Speciálně pro $n = q = 1$:

$$D(f \circ g)(a) = Df(g(a)) \circ Dg(a) = \nabla f(g(a))^T \cdot \begin{pmatrix} \partial g_1(a) \\ \vdots \\ \partial g_m(a) \end{pmatrix}$$

- Gradient** – směr největšího růstu – na derivaci ve směru s funkce $f: D_f \rightarrow \mathbb{R}$ v bodě $a \in D_f \subset \mathbb{R}^n$ můžeme nahlížet jako na derivaci jednorozměrné funkce $h(t) = f(a + ts)$ v bodě 0, $g(t) = a + ts$. Pak

$$Dh(t) = D(f \circ g)(t) = \nabla f(g(t))^T \cdot \begin{pmatrix} \partial g_1(t) \\ \vdots \\ \partial g_m(t) \end{pmatrix} = \nabla f(a + ts)^T \cdot \begin{pmatrix} s_1 \\ \vdots \\ s_n \end{pmatrix}$$

→ pro bod 0:

$$\nabla_s f(a) = Dh(0) = \nabla f(a)^T \cdot s.$$

- Optimalizace** → hledáme

$$\operatorname{argmin}_{x \in D \subset \mathbb{R}^n} f(x)$$

pro f (dvakrát) spojitě diferencovatelnou

- Řešení iteračními metodami → konstrukce posloupnosti aproximací x_0, x_1, x_2 konvergující k bodu, který je vhodným kandidátem na lokální minimum funkce f (typicky nulový gradient)
- Trust region** přístup – na okolí bodu x_k vytvoříme aproximaci funkce f , označené m_k , a hledáme minimum této aproximace m_k na okolí bodu x_k

$$p_k = \operatorname{argmin}_p \{m_k(x_k + p) \mid \|p\| \leq \Delta_k\}$$

- Další aproximace je $x_{k+1} = x_k + p_k$

- Line search** – aproximaci hledáme ve směru p_k :

$$x_{k+1} = x_k + \alpha_k p_k$$

kde $\alpha_k = \operatorname{argmin}_{\alpha \geq 0} f(x_k + \alpha_k p_k)$ nebo nějaká aproximace řešení této podúlohy

- Line search – směr spádu** $x_{k+1} = x_k + \alpha_k p_k$

kde p_k je zvolený směr a α_k je délka kroku

- Obecná volba směru = **směr spádu** (descent) = $\nabla f(x_k)^T \cdot p_k < 0$
- Časté volby směru: $p_k = -B_k^{-1} \nabla f(x_k)$:

- Metoda největšího spádu → $B_k = I$ (steepest descent)
- Newtonova metoda $B_k = \nabla^2 f(x_k)$
- Kvazi-Newtonova metoda $B_k \approx \nabla^2 f(x_k)$
- Je-li B_k pozitivně definitní, jedná se o spádový směr

- Ideální volba α_k do $x_{k+1} = x_k + \alpha_k p_k$ je minimem funkce

$$\phi(\alpha) = f(x_k + \alpha p_k) \quad \text{pro } \alpha \geq 0$$

= **exact** line search

- Výpočetně složité → hledání vhodné délky kroku, která zajistí dostatečný pokles funkční hodnoty f
= **inexact** line search

- **Armijova podmínka:**

$$\phi(\alpha_k) \leq \phi(0) + \alpha_k c \phi'(0),$$

kde $0 < c < 1$

→ funkční hodnota má ležet **pod zvolenou přímkou**, nalezený krok by ale neměl být příliš malý

- Klasický algoritmus:

Algorithm 3.1: ARMIJOSTEPBACKTRACK(ϕ, γ_0, c)

$i = 0$

while $\phi(\gamma_i) \geq \phi(0) + \gamma_i c \phi'(0)$

do $\begin{cases} \gamma_{i+1} = \rho \gamma_i \\ i = i + 1 \end{cases}$

return (γ_i)

- $\rho \in (0, 1)$

- **Goldsteinova podmínka:**

$$\phi(0) + \alpha_k(1 - c)\phi'(0) \leq \phi(\alpha_k) \leq \phi(0) + \alpha_k c \phi'(0)$$

kde $0 < c < 1/2$

→ funkční hodnota má ležet **mezi zvolenými přímkami**

→ omezuje se nevhodná možnost, že nalezený krok bude příliš malý

- Nevýhoda – můžeme minout optimální hodnotu pro α_k

- **Wolfeho podmínky:**

- Slabá podmínka:

$$\phi(\alpha_k) \leq \phi(0) + \alpha_k c_1 \phi'(0)$$

$$\phi'(\alpha_k) \geq c_2 \phi'(0)$$

kde $0 < c_1 < c_2 < 1$

- Silná podmínka:

$$\phi(\alpha_k) \leq \phi(0) + \alpha_k c_1 \phi'(0)$$

$$|\phi'(\alpha_k)| \leq c_2 |\phi'(0)|$$

kde $0 < c_1 < c_2 < 1$

- Přidání absolutní hodnoty → derivace $\phi'(\alpha_k)$ nebude mít moc velké hodnoty
- Pro spojitě diferencovatelnou a zdola omezenou f vždy existuje α_k splňující Wolfeho podmínky

- **Metoda největšího spádu**

- $x_{k+1} = x_k + \alpha_k p_k$ a $\alpha_k = \operatorname{argmin}_{\alpha \geq 0} \phi(\alpha)$
- **Největší spád:** $p_k = -\nabla f(x_k)$ (ajo my vlastně jdeme proti směru gradientu, co? PROTI SMĚRU NEJVĚTŠÍHO RŮSTU)
 - $(B_k = I)$
- Platí $p_{k+1}^T \cdot p_k = 0$

Věta 3.1 Necht f je dvakrát spojitě diferencovatelná a předpokládejme, že algoritmus největšího spádu konverguje k x^* , ve kterém je $\nabla^2 f(x^*)$ pozitivně definitní.

Necht r je číslo, pro které platí

$$r \in \left(\frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1}, 1 \right),$$

kde λ_n je největší a λ_1 nejmenší vlastní číslo matice $\nabla^2 f(x^*)$. Pro dostatečně velké k platí

$$f(x_{k+1}) - f(x^*) \leq r^2 (f(x_k) - f(x^*))$$

- $\frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1} = \frac{\frac{\lambda_n}{\lambda_1} - 1}{\frac{\lambda_n}{\lambda_1} + 1} = \frac{\kappa(\nabla^2 f(x^*)) - 1}{\kappa(\nabla^2 f(x^*)) + 1}$ poslední rovnost platí při použití mat. normy indukované E. normou
- Obecně vyžaduje vyšší počet kroků ke konvergenci, pro hůř podmíněnou je moc pomalá – při použití nepřesné volby délky kroku nebude rychlejší

- Newtonova metoda

- Máme aproximaci x_k a posunutí p . Pro $f(x_k + p)$ platí:

$$f(x_k + p) \approx f(x_k) + \nabla f(x_k)^T p + \frac{1}{2} p^T \nabla^2 f(x_k) p = m_k(p)$$

- Optimální $p \rightarrow$ **minimalizace** funkce m_k :

$$p_k = \underset{p}{\operatorname{argmin}} m_k(p)$$

- Po parciální derivaci 1. rovnice:

$$\nabla m_k(p) = 0 \iff \nabla f(x_k) + \nabla^2 f(x_k) p = 0.$$

$$p_k = -(\nabla^2 f(x_k))^{-1} \nabla f(x_k)$$

- Pokud není $\nabla^2 f(x_k)$ **pozitivně definitní**, pak p_k nemusí být spádový směr!

Věta 3.2 Necht f je dvakrát spojitě diferencovatelná, $\nabla^2 f(x)$ je Lipschitzovsky spojitá na okolí H bodu minima x^* ($\exists L, \forall x_0, x_1 \in H, \|\nabla^2 f(x_1) - \nabla^2 f(x_0)\| \leq L\|x_1 - x_0\|$), $\nabla f(x^*) = 0$ a $\nabla^2 f(x^*)$ je pozitivně definitní. Pro Newtonovu metodu s $\alpha_k = 1$ platí

1. je-li x_0 dostatečně blízko x^* , potom $(x_k)_{k=0}^{+\infty}$ konverguje k x^* ;
2. $(x_k)_{k=0}^{+\infty}$ konverguje kvadraticky;
3. posloupnost $(\|\nabla f(x_k)\|)_{k=0}^{+\infty}$ konverguje kvadraticky k nule.

-
- Poblíž řešení bude délka kroku $\alpha_k = 1$ splňovat Wolfeho podmínky \rightarrow s výchozí délkou kroku blízko 1 budou nalézat podobné kroky a kvůli výpočetní složitosti se délku nevyplatí optimalizovat

○ Modifikace Newtonovy metody

- Matice $\nabla^2 f(x_k)$ (= Hessián) nemusí být pozitivně definitní \rightarrow nemusí existovat inverze
- **Inverzi chceme**, abychom mohli vyřešit $B_k p_k = -\nabla f(x_k) \rightarrow$ modifikace
- Nastavení matice $B_k = \nabla^2 f(x_k) + E_k$ tak, aby B_k byla **pozitivně definitní**:
 - Změna záporných vlastních čísel
 - Přičtení kladné diagonální matice
 - Modifikace nějakého rozkladu matice $\nabla^2 f(x_k)$
- Pokud jsou čísla podmíněnosti matice B_k omezená, pak lze dokázat, že modifikace bude konvergovat

- Stochastický gradientní sestup

- Minimalizovaná funkce ve tvaru sumy:

$$f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x)$$

f_i ... ztrátová funkce pro i -tý vzorek

x ... parametry modelu

- **SGD** (Stochastic gradient descent) – začíná s aproximací x_0 a iteruje podle schématu:

$$x_{k+1} = x_k - \eta \nabla f_{i_k}(x_k).$$

i_k ... náhodné z uniformního rozdělení

η ... **míra učení** (learning rate)

- Výhoda = nízká výpočetní složitost, pro konvexní funkce konvergence skoro jistá
- **Modifikace SGD** – většinou modifikace míry učení nebo implementace setrvačnosti (vliv předchozích změn)

- **Adagrad:**

$$s_0 = 0,$$

$$s_k = s_{k-1} + (\nabla f_{i_k}(x_k))^2,$$

$$x_{k+1} = x_k - \eta \frac{\nabla f_{i_k}(x_k)}{\sqrt{s_k + \varepsilon}},$$

- Vektorové operace po složkách
- ε , abychom nedělili 0

- **Momentová metoda:**

$$x_{k+1} = x_k - (\eta \nabla f_{i_k}(x_k) + \gamma(x_k - x_{k-1})).$$

- γ ... faktor zapomínání menší než 1 (typicky 0.9)

20. Časové řady: aditivní a multiplikativní dekompozice, momenty (střední hodnota, rozptyl, autokovariance). Druhy stacionarity a rozdíl mezi nimi. Základní vlastnosti náhodné procházky a bílého šumu.

NI-SCR

Časové řady

- **Časová řada** = soubor pozorování x_t získaných v konkrétních časových okamžicích t
 - o Bud' (Ω, \mathcal{F}, P) pravděpodobnostní prostor a T množina indexů interpretovaných jako čas. Časovou řadou nazýváme množinu $\{x_t, t \in T\}$, kde x_t jsou náhodné veličiny z (Ω, \mathcal{F}, P)
 - o Pokud je t z celých čísel, je do řada s diskretním časem (tyhle tu teď řešíme), jinak se spojitým
- **Variabilita vývoje**
 - o **Trend** – dlouhodobý vývoj střední hodnoty
 - o **Sezónnost** – periodicky se opakující pravidelný vývoj časové řady
 - Zjištění periody – autokorelační funkce (ACF)
 - Pro analýzu sezónnosti nutno odstranit trend
 - o **Cyklické změny** – nepravidelné fluktuace – např. ekonomické cykly – ne pevná perioda
 - o Další **nepravidelné fluktuace**
- **Rozložení časové řady na složky**
 - o Sloučení cyklických změn a nepravidelných fluktuací do jednoho
 - Y_t ... pozorovaná veličina v čase t
 - T_t ... hodnota trendu
 - S_t ... sezónní složka
 - E_t ... nevysvětlitelná složka
 - o 2 modely:
 - **Aditivní:** $Y_t = T_t + S_t + E_t$
 - Amplituda sezónních složek je cca stejná
 - **Multiplikativní:** $Y_t = T_t \cdot S_t \cdot E_t$
 - S rostoucím trendem se zvyšuje i sezónní amplituda
 - o **Výběr modelu** – minimalizace součtu čtverců hodnot **autokorelační funkce** reziduí E_t – říká, jaká míra korelace v E_t zbyla
 - Od pozorované veličiny lze pro různé aplikace odečítat složku trendu nebo sezónnosti
 - Autokorelační koeficient – lineární korelace hodnot časové řady v různých časových okamžicích
 - Umožňuje odhalit opakující se vývoj řady
 - Značení $R(s, t)$ pro časy s, t
- **Náhodný proces** = posloupnost náhodných veličin x_t , kde t je z vhodné množiny indexů
- **Momenty** – popisují časové řady
 - o **Střední hodnota:** $\mu_t = E[x_t]$
 - o **Variance:** $\sigma_t^2 = \text{var } x_t$
 - o **Autokovariance:** $\gamma(t_1, t_2) = \text{cov}(x_{t_1}, x_{t_2}) = E[(x_{t_1} - \mu_{t_1}) \cdot (x_{t_2} - \mu_{t_2})]$

Stacionarita

- **Striktní stacionarita** (silná) – řada je striktně s., pokud sdružená distribuce x_{t_1}, \dots, x_{t_n} je stejná, jako sdružená distribuce $x_{t_1+\tau}, \dots, x_{t_n+\tau}$ pro všechna t_1, \dots, t_n, τ
 - o Libovolný posun o čas τ nemá vliv na sdruženou distribuci a ta tedy závisí jen na časech t_1, \dots, t_n pro libovolná n
- **Slabá stacionarita** – řada je slabě s., pokud je invariantní vůči posunům v čase pouze v rámci momentů rozdělení do druhého řádu:

$$\begin{aligned} E[x_t] &= \mu \\ \text{cov}(x_t, x_{t+\tau}) &= \gamma(\tau) \end{aligned}$$

- Trend → řada není stacionární
- Typy procesů podle stacionarity:
 - o **Stacionární** – bez trendu a jednotkových kořenů, silně či slabě
 - Pokud je $\{x_t, t \in T\}$ slabě stacionární proces, potom autokorelace závisí pouze na zpoždění mezi časy t a s : $R(s, t) = R(s - t) = R(\tau)$
 - o **Trend-stacionární** – stacionární po odstranění lineárního/nelineárního trendu
 - o **Stacionární po diferencování** (=procesy s jednotkovým kořenem) – stacionární po tolika diferencích, kolik mají jednotkových kořenů, ale jen za předpokladu, že nemají jiné kořeny uvnitř jednotkové kružnice
 - o **Nestacionární**
- Testy stacionarity
 - o **ADF** (Augmented Dickey-Fuller)
 - Uvažuje hypotézy:

$$H_0: \text{proces má jednotkový kořen}$$

$$H_A: \text{proces nemá jednotkový kořen (ale kořeny vně jednotkové kružnice)}$$
 - o **KPSS** – neparametrický test
 - Hypotézy:

$$H_0: \text{proces je trend – stacionární}$$

$$H_A: \text{proces není stacionární (má jednotkový kořen)}$$
 - Proces H_0 je **mean-reverting** = po šoku se vrací ke střední hodnotě
 - o Výsledky testů: (významný – V = zamítám)

ADF	KPSS	Pravděpodobná vlastnost
V	N	stacionární
N	V	Nestacionární, ex. Jednotkový kořen
N	N	Nedostatek evidence, možná trend-stacionární
V	V	Heteroskedasticita, strukturální změna, ...

Náhodná procházka a bílý šum

- **Bílý šum**:
 - o $E[x_t] = 0$
 - o $var(x_t) = \sigma^2 < \infty$
 - o $cov(x_t, x_{t+\tau}) = \gamma(\tau) = 0$
 - o Normální bílý šum:
 - o Zvuková syntéza, generátory náhodných čísel
- **Náhodná procházka** – uvažujeme diskrétní bílý šum $z_t \sim \mathcal{L}(0, \sigma^2)$. Proces $\{x_t\}$ je náhodný proces, pokud:

$$\begin{aligned}
 x_0 &= 0 \\
 x_t &= x_{t-1} + z_t \\
 \rightarrow x_t &= \sum_{i=1}^t z_i
 \end{aligned}$$

- o $E[x_t] = 0$
- o $var(x_t) = t \cdot \sigma^2$
- o Modelování cen akcií, odhad velikosti webu

21. Autoregresní modely (AR) a modely klouzavých průměrů (MA): základní vlastnosti modelů/procesů, jejich stacionarita. Zápis AR a MA, včetně zápisu pomocí operátoru zpoždění. Identifikace řádů AR a MA z autokorelačních funkcí a pomocí informačních kritérií.

NI-SCR

- Autokorelace

- **Lineární korelační koeficient** – určuje míru lineární nezávislosti mezi veličinami

$$\rho_{XY} = \frac{E[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y} = \frac{\sigma_{XY}}{\sigma_X \sigma_Y}$$

- $\rho \in [-1, 1]$
- $\rho_{XY} = 0 \Rightarrow X$ a Y jsou nezávislé (opačně neplatí)
- σ_{XY} = kovariance
- **Výběrový korelační koeficient** – pomocí výběrových variancí
- **Parciální korelační koeficient**
 - 2 náhodné veličiny X a Y , mezi kterými existuje závislost
 - X a Y ovlivněny třetí n -rozměrnou náhodnou veličinou Z
 - Pro měření korelace X a Y je potřeba je od vlivu Z očistit
 - Nalezneme regresní přímky:
$$\hat{X} = a + b^T Z$$

$$\hat{Y} = c + d^T Z$$
 - \hat{X}, \hat{Y} ... nejlepší lineární přiblížení k X a Y
 - $e_X = (X - \hat{X}), e_Y = (Y - \hat{Y})$... rezidua očištěná od vlivu Z
 - Jejich korelační koeficient = **parciální korelační koeficient** mezi X a Y při daném Z – $\rho_{XY \cdot Z}$
 - Projeví se v sezónnosti v časové řadě – vyšší hodnoty korelace v periodách sezón
- **Autokorelační koeficient** = lineární korelace hodnot časové řady v různých časových okamžicích
 - Umožňuje odhalit opakující se vývoj řady
 - Značení $R(s, t)$ pro časy s a t
 - Pokud je $\{X_t, t = 1, 2, \dots\}$ slabě stacionární proces (= existují časově invariantní první 2 momenty μ, σ^2), potom autokorelace závisí pouze na zpoždění mezi t a s :
$$R(s, t) = R(s - t) = R(\tau)$$
- **Parciální autokorelační funkce**
 - **Parciální autokorelace zpoždění $k \geq 1 - \alpha(k)$**
 - (auto)korelace mezi X_t a X_{t+k} s odstraněním lineárního vlivu mezilehlých hodnot $X_{t+1}, \dots, X_{t+k-1}$

- Informační kritéria

- **AIC – Akaikeho informační kritérium**

$$AIC = 2k - 2 \ln \mathcal{L}$$

- k ... počet odhadovaných parametrů
- \mathcal{L} ... maximální hodnota věrohodnosti při daném modelu
- Asymptoticky ekvivalentní ke křížové validaci

- **BIC – Bayesovské informační kritérium**

$$BIC = k \ln(n) - 2 \ln \mathcal{L}$$

- n ... počet pozorování

- Hodnotu AIC/BIC chceme minimalizovat

- Operátor zpoždění – lag operator B nebo L – pro zjednodušení zápisu

$$BX_t = X_{t-1}$$

$$B^{-1}X_t = X_{t+1}$$

$$B^k X_t = B \cdot B \cdot \dots \cdot BX_t = X_{t-k}$$

(B je tam k -krát)

Autoregresní modely (AR)

- Umožňují popis náhodného procesu na základě jeho předchozích realizací
- Autoregresní model řádu p :

$$x_t = c + \Phi_1 x_{t-1} + \Phi_p x_{t-p} + \varepsilon_t = c + \sum_{\tau=1}^p \Phi_\tau x_{t-\tau} + \varepsilon_t$$

- o ε_t ... bílý šum
- o $\Phi = [c, \dots, \Phi_p]^T$... vektor regresních koeficientů
- o Zápis pomocí operátoru zpoždění:

$$\varepsilon_t = \left(1 - \sum_{i=1}^p \Phi_i B^i\right) x_t$$

- p může být libovolné, ale vyšší řády nemusí dávat smysl – jak hlubokou minulost potřebujeme?
- AR procesy **nemusí být slabě stacionární** – kořeny charakteristické rovnice $1 - \Phi_1 z - \dots - \Phi_p z^p = 0$ musí ležet vně jednotkové kružnice (tzn. $\forall z_i \in \mathbb{C}: |z_i| > 1$)
- **Odhad parametrů** metodou nejmenších čtverců – odhad regresních koeficientů c a Φ_1
 - o X návrhová matice, Y vektor měření, ε šum
 - o $Y = X\Phi + \varepsilon \rightarrow \hat{\Phi} = (X^T X)^{-1} X^T Y$
- **Odhad řádu** AR modelu
 - o **ACF** – postupně klesá k nule, popř. klesá shora i zdola
 - Má mnoho významných lagů (stacionární AR procesy lze konvertovat do $MA(\infty)$)
 - o **PACF** – vrcholy do hodnoty řádu modelu, pak jdou strmě k nule
 - o Obecný příklad (při bílém šumu ε_t) $x_t = \Phi_1 x_{t-1} + \varepsilon_t$
 - Abs. hodnoty **ACF** budou zřejmě postupně klesat, protože x_t je přímo ovlivněno x_{t-1} atd.
 - Absolutní hodnoty **PACF** budou indikovat silnou korelaci mezi x_t a x_{t-1} , ale další hodnoty už by byly 0 kvůli očištění

Modely klouzavých průměrů (MA)

- Předchozí hodnota je **konstanta**, **šumová složka se propaguje** – nebere v potaz předchozí měření
- Model klouzavých průměrů řádu q :

$$x_t = c + \varepsilon_t + \theta_1 \varepsilon_{t-1} + \dots + \theta_q \varepsilon_{t-q}$$

- o $\varepsilon_t \sim iid N(0, \sigma^2)$
- o θ ... váha toho, jak se šum propaguje
- o Zápis pomocí operátoru zpoždění:

$$x_t = \left(1 + \sum_{i=1}^q \theta_i B^i\right) \varepsilon_t$$

- Prostřednictvím **bílého šumu** vystihuje **náhodné šoky** – nezávislé a stejně rozdělené
- $E[x_t] = c, var x_t = \sigma^2 \sum_{i=0}^q \theta_i^2$
- Veličiny x_t a x_{t+k} jsou pro k větší než řád modelu q nekorelované, korelace x_t a x_{t+k} pro $k < \text{řád}$ nenulová
- Prohodíme-li v kovarianci pořadí, změní se znaménko u k
- **Odhad** MA modelu je složitý \rightarrow numerická optimalizace
- **Invertibilita MA procesů** – u MA nás zajímá místo stacionarity
 - o MA proces 1. řádu je ekvivalentní AR procesu ∞ řádu
 - o **Charakteristický polynom** $MA(q)$:
$$\theta(x) = 1 - \theta_1 x - \theta_2 x^2 - \dots - \theta_q x^q$$
 - o Jsou-li kořeny char. polynomu vně jednotkové kružnice, potom je MA proces invertibilní - $AR(\infty)$
- **Odhad řádu** MA modelu
 - o ACF – vrcholy do hodnoty řádu modelu, pak jdou strmě k 0
 - o PACF – mnoho významných lagů (kvůli invertibilitě) – jako ACF u AR modelů
 - o Obecný příklad (při bílém šumu ε_t) $x_t = c + \varepsilon_t + \theta_1 \varepsilon_{t-1}$
 - x_t není ovlivněna x_{t-1}, x_{t-2}, \dots
 - Absolutní hodnoty **ACF** budou vysoké pro $k = 1$, ostatní nulové (šum iid)
 - Absolutní hodnoty **PACF** budou klesat k 0

22. Smíšené modely ARIMA: základní vlastnosti modelů/procesů, integrování a diferencování. Zápis ARIMA, včetně zápisu pomocí operátorů zpoždění a difference, speciální případy podle hodnot p, d, q. Problém redundance parametrů.

NI-SCR

- Smíšené modely ARMA (p, q)

$$x_t = c + \varepsilon_t + \sum_{i=1}^p \Phi_i x_{t-i} + \sum_{i=1}^q \theta_i \varepsilon_{t-i}$$

- o Neznámé $c, \Phi, \theta, \sigma^2$
- o Populární – flexibilní, ale mnoho neznámých
- o Většinou předpokládána normalita šumu
- o Ekonometrie – ceny akcií – MA šokové změny + AR vývoj na základě minulých cen
- o Rozšíření NARMA, ARIMA, SARIMA, VARMA, ...
- o Předpoklad = **slabá stacionarita** – co když je porušena?
 - Odstranění trendu
 - Metoda diferencí – pokud nepomůžou 1. difference, pomohou často 2.
 - ARMA model, který tohle dělá, je ARIMA
 - ARIMA (3,2,1) je ARMA (3,1) se dvěma diferencováními
- o Zápis pomocí operátoru zpoždění (bez c):

$$\Phi X_t = \theta \varepsilon_t$$

- Box-Jenkinsův přístup k ARMA modelům

1. **Identifikace modelu** – posouzení stacionarity a sezónnosti časové řady, odhalení přítomnosti AR a MA části a jejich řádů (ACF/PACF), popř. transformace pro zajištění stacionarity
2. **Odhad parametrů** – MLE
3. **Ověření modelu** – posouzení nekorelovanosti reziduí a jejich slabé stacionarity v čase (ACF/PACF)

Smíšené modely ARIMA (p, d, q)

- Parametry:

- o p – řád **autoregresní** části (AR)
- o q – řád modelu **klouzavých průměrů** (MA)
- o d – řád **diferencování**

- Zápis pomocí operátoru zpoždění (bez konstanty c):

$$\Phi(1 - B)^d X_t = \theta \varepsilon_t$$

$$\Phi = 1 - \sum_{i=1}^p \Phi_i B^i, \quad \theta = 1 + \sum_{i=1}^q \theta_i B^i$$

- Operátor difference ∇ :

$$\begin{aligned} \nabla X_t &= X_t - X_{t-1} = (1 - B)X_t \\ \nabla^k X_t &= (1 - B)^k X_t \end{aligned}$$

- Diferencování řady → odstranění trendu

- o Počítáme veličiny $Y'_t = Y_t - Y_{t-1}$
- o Diferencováním náhodné procházky dostaneme gaussovský bílý šum ($Y'_t = \varepsilon_t$)
 - Gaussovská náhodná procházka:

$$Y_0 = 0$$

$$Y_t = Y_{t-1} + \varepsilon_t$$

- Časové řady s **trendem nebo náhodnou procházkou** jsou vždy **silně pozitivně autokorelované**

- o V ACF vidíme velké korelace
- o V PACF typicky v 1. lagu a blízko 1

- Časové řady s ACF v 1. lagu s hodnotami -0.5 a méně mohou být **přediferencované**

- Modely s $d=1$ typicky předpokládají **konstantní průměrný trend** (náhodná procházka s driftem)
- AR charakteristika může značit **poddiferencovanost** časové řady, MA pak **přediferencovanost**

- Volba d :
 - o $d \leq 2$, 2. difference málokdy, 3. je výjimka
 - o Při $d = 2$ nepoužíváme v ARIMA modelu konstantu c
- Role konstanty v ARIMA modelu:
 - o $d = 0$ – konstanta zavádí nenulovou střední hodnotu a vyplatí se ji zkusit
 - o $d = 1$ – konstanta zavádí nenulový „průměrný“ trend, může se vyplatit
 - o $d = 2$ – konstanta by měla význam „trendu v trendu“, obvykle nechceme
- Běžné ARIMA modely:

ARIMA (0, 0, 0) + c	Konstantní model
ARIMA (0, 1, 0)	Model náhodné procházky
ARIMA (0, 1, 0) + c	Náhodná procházka s driftem
ARIMA (1, 0, 0) + c	AR (1)
ARIMA (2, 0, 0) + c	AR (2)
ARIMA (1, 1, 0) + c	AR (1) na 1x diferencovaných datech
ARIMA (2, 1, 0) + c	AR (2) na 1x diferencovaných datech
ARIMA (0, 1, 1)	Jednoduché exponenciální vyhlazování = MA (1) na 1x dif. Datech
ARIMA (0, 1, 1) + c	MA (1) na 1x diferencovaných datech s konstantním lineárním trendem
ARIMA (1, 1, 2)	Lineární exponenciální vyhlazování s tlumeným trendem
ARIMA (0, 2, 2)	Zobecněné lineární exponenciální vyhlazování
ARIMA (1, 0, 0) + c	Podle toho, jaké parametry má drift c : <ol style="list-style-type: none"> 1. $\Phi_1 = 0 \rightarrow$ bílý šum 2. $\Phi_1 = 0, c = 0 \rightarrow$ náhodná procházka 3. $\Phi_1 = 0, c \neq 0 \rightarrow$ náhodná procházka s driftem 4. $\Phi_1 < 0 \rightarrow$ řada oscilující mezi klad. a záp. hodnotami

- Obecná pravidla
 - o $p = 0$ nebo $q = 0$
 - o $p + q \leq 3$ – složitější modely jsou řídké
 - o Koeficienty Φ_i blízke θ_i mohou pracovat proti sobě
 - o Obecně volíme co **nejjednodušší** modely
 - o Modely lze porovnávat kritérii **AIC, BIC**
 - o Můžou existovat i složitější modely porušující pravidla
 - I tak by měly být p, q malé
 - Problém složitých modelů – **redundance parametrů**
 - AR a MA části modelu pracují proti sobě – zvýšíme-li řády u obou, mohou se navzájem vyrušit \rightarrow zvyšuje se složitost :(