

# A MEMÓRIA IDEGRENDSZERI MODELLJEI

MEMÓRIA: információ tárolás esetleges előhívás céljából  
↳ érzékelésből és annak feldolgozásából származik

TANULÁS, MEMÓRIA/EMLEKEZET, FELEJTÉS, ELŐHÍVÁS nem elválaszthatók

Az emlékezés folyamata információelméleti szempögből:

- KÓDOLÁS: információ befogadása, feldolgozása és kombinálása
- TÁROLÁS: kódolt információ tartós megörzése
- ELŐHÍVÁS: felidézés valamely inger hatására / egymásra / hatással vannak/

## IDEGRENDSZERI MODELL mit ér le?

Az a kérdés, hogy milyen tulajdonságokkal rendelkezik a memória

Előzmény: EMLEKEZETKUTATÁS: filozófia → kognitív pszichológia

KORAI EMLEKEZETKUTATÁS 3 FŐ ISKOLA'JA:

→ ÉSZAK-AMERIKA:

Hermann Ebbinghaus (1850-1909) német filozófus követői

↳ elsőként végzett kísérleteket (saját magán) → felejtési görbe (exponenciális)

verbális tanulási megközelítés szólistákkal

behaviorizmus: inger-válasz pszichológia (introspektív módszer elutasítása)

→ NEMETORSZÁG:

Gestalt-pszichológusok (alaklélektan)

belső reprezentációk hangsúlyozása

észleléskutatáshoz hasonló megközelítés → aktív folyamat

→ NAGY-BRITANNIA:

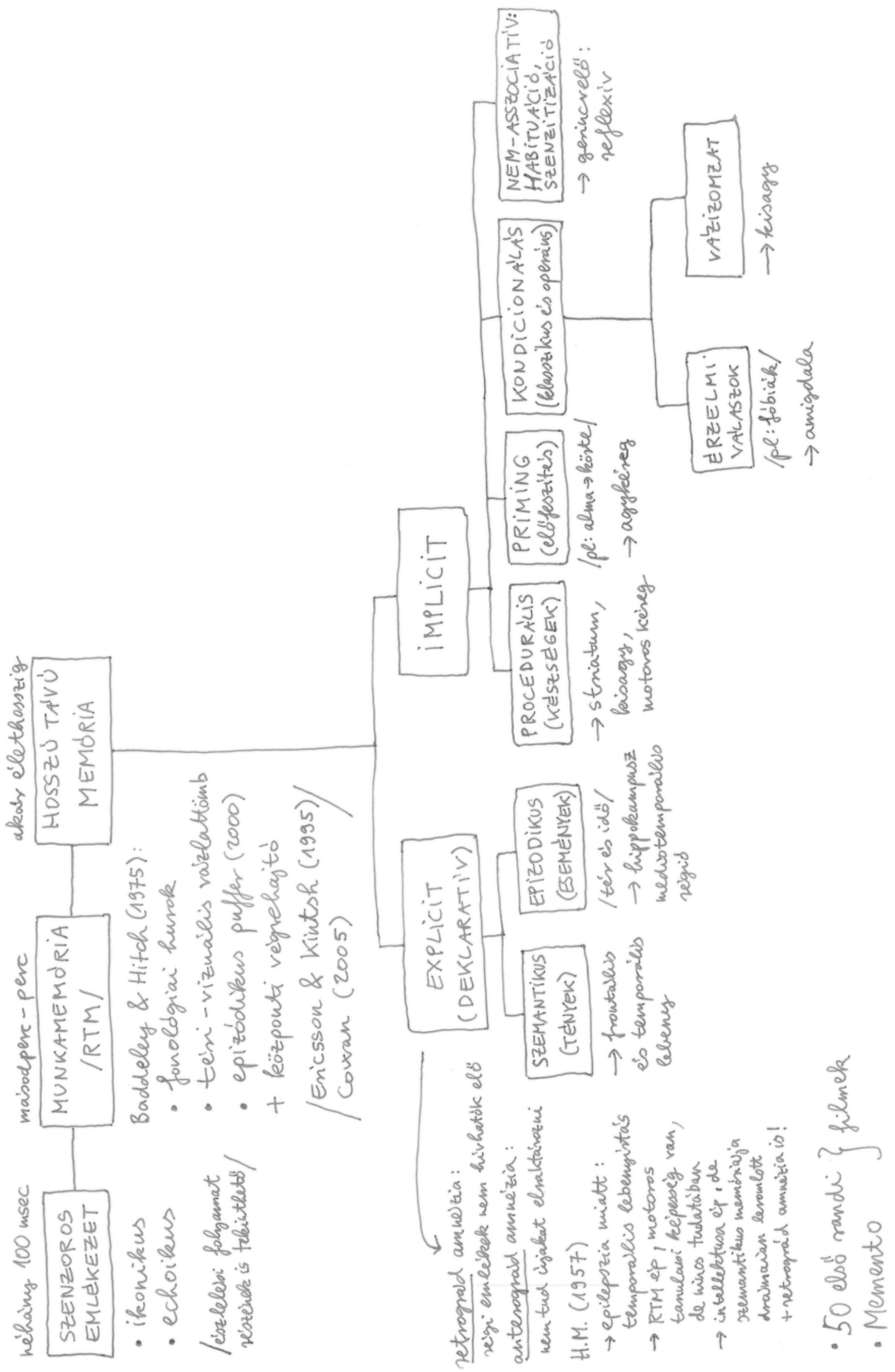
Frederic Charles Bartlett "Az emlékezés" (1932)

életszerű komplex situációk vizsgálata (idegen népmese)

hibázások vizsgálata

semmel: strukturális a tudást és befolyásolják az előhívást

# EMLÉKEZETI RENDSZEREK



Baddeley & Hitch (1975):

- fonológiai burok
- tér- és vizuális vázlatlomb
- epizodikus puffer (2000)
- + központi vegreható

/Ericsson & Kintosh (1995) /  
/Cowan (2005)

- ikonikus
  - echoikus
- /érzélelni folyamat részletek is felírhatók/

retrográd amnézia: régi emlékek nem hívhatók elő

anterográd amnézia: nem tud újakat elraktározni

H.M. (1957)

- epilepszia miatt: temporális lebenyintás
- RTM ép, motoros tanulási képesség van, de nincs tudatában
- intellektusa ép, de szemantikus memóriája drámaian leromlott + retrográd amnézia is!

- 50 elő randi } filmek
- Memento

SZENZOROS EMLÉKEZET: inger meghosszabbítása feldolgozás céljából

Vizsgálata: rövid ideig bemutatott ingerekkel

IKONIKUS EMLÉKEZET:

- 1740 Segner: iztó cigarettareg (fennmaradás  $\times 100$  msec)
- 1960 Sperling: 3x4 betű 50 msec-ig → 4-5 betű visszaidézés  
Ennyit láttak v. elfelejtették a beáramló alatt?  
Ha egy sorból kérdeznék vissza, akkor 3-4 betű  
⇒ halványuló emléksorom!

további effektusok:

- sötét mező előtte-utána ⇒ tovább fennmarad
- villanás töröl (rendelkezésre álló időben lineáris)
- maszkolási hatás retinális (nem mindegy melyik szem)
- de létezik mintázatmaszkolás, amelynek mindegy



AZ ÉSZLELÉS FOLYAMATA TÖBB SZINTEN IS TÁROLÁST IGÉNYEL!

ECHOIKUS EMLÉKEZET:

- 1970 Efron: szubjektív élmény szerint egy 30 v. 100 msec-es hang egyformán 130 msec-ig szól, de ez lerövidül, ha egy második hang követi
- Békeisy György: a jelenség fontos hangversenytermék akusztikájánál (visszhang)

megjegyzés: Az auditoros rendszerrel egyértelmű a szenzoros emlékezet szükségessége, hiszen a hang időbeli rezgés:

- tiszta zenei hang esetében is van hangnyomás változás:  
a változás felismeréséhez pedig összehasonlítás, azaz memória szükségeltetik

## RÖVID TAVÚ EMLÉKEZET:

modalitások integrációját teszi lehetővé!

⇒ előlött nem érdemes modalitás szerint osztályozni

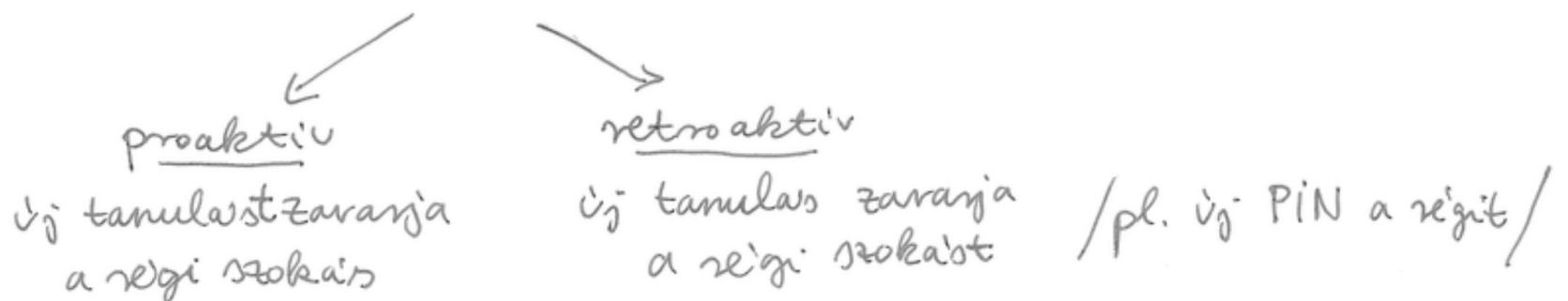
## RTM ES A TUDAT:

"Úgy tűnik, létezik elmémnek egy királyi fogadóterme, ahol a tudat ülészik, és ahova egyszerre két-három idea járul kihallgatásra." (Galton)

## RÖVIDTAVÚ FELEJTÉS

→ Például: megtarant bemutatkozás

Elméletek: → automatikus nyomelhalványulás } mindkét faktor fontos!  
→ interferencia (hasonló jobban zavar)

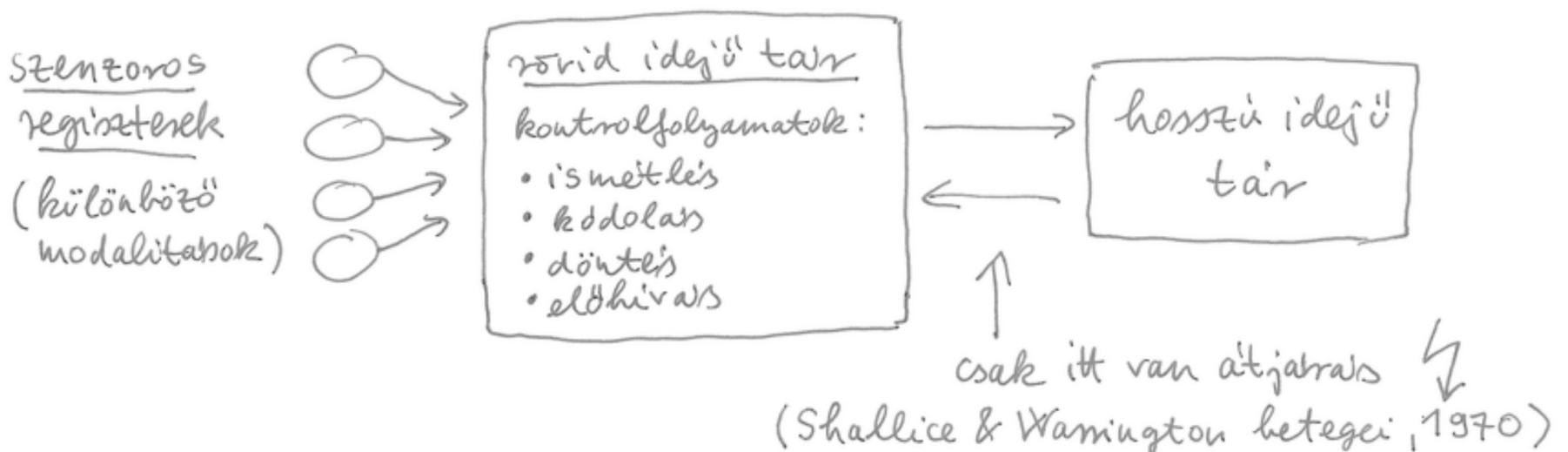


Különálló (rövid/hosszú tavú) rendszerek vagy sem?

→ Feldolgozási szintek elmélete:

Mivel mélyebben dolgozunk fel, annál jobban emlékezünk rá (átmenet).

→ Atkinson & Shiffrin modális modellje



INTERAKCIÓ A RENDSZEREK KÖZÖTT:

pl: deklaratív tudás leronthatja a procedurálisét, ha rásegítünk

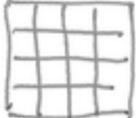
## MUNKAMEMÓRIA

Több komponens? → párhuzamosan végzett feladatokkal vizsgálható

FONOLÓGIAI HUROM: tár + artikulációs kontrollfolyamat

- nem a szótagszám, hanem a kiejtési idő számít
  - gyorsbeszédiek előnyben
  - bizonyos nyelvek előnyben-hátrányban (szónterjedelem)
- van olyan beteg, akinél csak a kontrollfolyamat sérült
- hasonló hangzásra a felidézés romlik
- olvasás tanuláshoz és nyelvi megértéshez is köze van!

Vizuális emlékezet: milyen mértékben hasonló a percepcióhoz?

- „lelki szemünk előtt” szubjektív értelemben nagy a variancia
- Shepard & Metzler 1971: mentális rotáció
  - lineáris és nem függ a komplexitástól (térben is)
  - de ha azt kell eldönteni, hogy része-e, akkor függ
- papírtábla lelki szemünkkel:
  - tovább tart válaszolni, ha a képelemek távolabb esnek korábbi képen
  - hasonló eredmény egyetemi város képzeletbeli térképén (mics mutatás)
- ha képzeletben közeledünk: jobban kitölti a látóteret az objektum
- de ha egy kockát testátló mentén ujjaink közé fogunk, akkor nem „látjuk” jól, hogy a köztér csúcsok nem egy síkba esnek!
- az abakusz mesterek tér- vizuális műemeltechnikát használják: szónterjedelmük oda-vissza vagy (de betűkre nem) és vizuális feladat zavarja őket
- Brookes - feladat:  mondat felidézések:

1-es feladat (vizuális): „A következő négyzetbe jobbra tegyen 2-est”  
2-es feladat (verbális): „A következő rossz négyzetbe tegyen 2-est”

(vizuális átlag: 8, verbális átlag: 6)

→ vizuális feladat (pontkövetés pálcával) rott az elsőt, de nem rott a verbálison

Tér- vagy vizuális?

→ bekötött szemmel csak hang alapján pontkövetés ⇒ tér- zavarás fontosabb!

Hol és Mi?

→ agykárosult beteg: tudja mi, de rossz irányba nyúl

# MUNKAMEMÓRIA IDEGRENDSZERI MODELL

Korábban ún. perisztens aktivitás modell:

a kódoló populáció a jel beírása után magasabb rátaival tüzel  
(az információ spike-okban tárolása metabolikusan költséges)

↓ helyette előfeszített állapot: facilitáció → kioldásakor sinkron tüzelés

Mongillo, Barak, Tsodyks (Science 2008):

## SYNAPTIC THEORY OF WORKING MEMORY

- prefrontális kortex: serkentő szinapszisok, facilitáció jellemző  
(ellentétben a szenzoros területekkel)
- $Ca^{2+}$  szint kódol a preszinaptikus oldalon (lecsengés  $\sim 15$  spike nélkül)

### SZIMULÁCIÓ:

- rekurrens integrate-and-fire neuronhálózat
- input serkentő populáció (random)
- Hebb prior: ugyanazon input serkentő neuronjai között erősebb szinapszisok
- Tsodyks-Markram modell dimenziótlan változó:  
forrás:  $x \in [0, 1]$  } spike után  $x \rightarrow x - ux$   
kihátrahúzás:  $u$
- alapszintek  $x=1, u=1$
- időállandók:  $\tau_D$  (depressió),  $\tau_F$  (facilitáció)  
 $x$  lecsengése  $\tau_D \ll \tau_F$
- előhívás: alacsony jel az egész hálózaton
- sinkronizált válasz: 20 ms-en belül tüzel az adott populáció
- egyúttal fűssül az előfeszítés ( $\tau_F$  lecsengésű)
- növelve a háttér: fixpont  $\rightarrow$  oszcilláció
- az ábrán 10.000 neuronból csak egy része van ábrázolva
- az ábrán  $x$  és  $u$  átlagok vannak
- több elem tárolása, lásd ábrán

## ASSZOCIATÍV MEMÓRIÁK

A rövid távú emlékezetre tekinthetünk úgy, mint egy átmeneti tárr.

A hosszú távú memóriába azonban tudást szeretnénk rögzíteni és később használni.

Minden zajos, ezért a tudással kapcsolatos következtetéseink asszociációk lesznek (oda-vissza).

FAJTAI : → autoasszociatív (töredékes jelből az eredeti)  
→ heteroasszociatív

↔ ellentéte közvetlen címzés: RAM, telefonkönyv

TÉNYEZŐK : kapacitás, stabilitás, vezési tartomány (medence)

MEGVÁLÓSÍTÁS ATTRAKTORHÁLÓZATOKKAL :

rekurrens neurális hálózat attraktorai :

→ stabil fixpont

→ periodikus oszcilláció (határciklus)

→ kaotikus

TANULÁS (SZINAPTIKUS SÚLYOK) :

→ offline (adatbázis a kezünkben)

→ online (folyamatosan jön az adat)

→ one-shot learning (adatbázis megvan, de egyszer nyúlhatunk hozzá)

ELŐHÍVÁSRA / ASSZOCIÁCIÓRA VONATKOZÓ MEGFIGYELÉSEK :

→ sokszor a régi ismerős meg kell szólaljon, hogy megismerjük

→ bűnügyek : azonosítás hang alapján

→ elérhetőség és rendelkezésre állás között különbség :

38% felidézés mellett 96% felismerés (100 szó)

→ (B)emutató, (F)elidézés → BFBFBF... } ugyanolyan hatékonyak  
→ BFFFBBBB... }

⇒ elem felidézése megnöveli az előhívási valószínűséget

→ dallam felismerés : dallamkontúr, hangmagasság intervallumok megőrződnek

Született tehetség szerepe? Nagymesteri szint?

pl: Mozart : Vatikáni kéms saját tulajdonú műve (tilos másolni)

Toscanini Karmester

→ kontextusfüggés : mélytengeri búvárok

→ gyakoriság paradoxon : gyakori szavakat jobban idézzük fel, de rosszabbul ismerjük fel

# HOPFIELD - HÁLÓZAT

(John Hopfield amerikai biológus 1982) → folytonos: '84

Mit várunk a hálózattól?

→ Tárolt információ zajos, torzított, esetleg hiányos változataiból (bemenet) az eredetire asszociáljon (stabil állapot)

Alapmodell: egyetlen teljesen összekötött réteg kétejtékű neuronokkal:

$$x_i \in \{-1, +1\} \rightarrow \{0, 1\}\text{-re való áttérés: } x_i = 2y_i - 1$$

DINAMIKA: 
$$x_i = \text{sgn} \left( \sum_j w_{ij} x_j - \theta_i \right)$$

→ szinkron frissítés: oszcillálhat (2-ciklus)

→ aszinkron frissítés tipikusabb (szekvenciális):  
minos távoli ugrás az állapotterben

Egy minta tárolása ( $\theta = 0$  eset)

$a$  minta stabil, ha  $\forall i$ : 
$$\text{sgn} \left( \sum_j w_{ij} a_j \right) = a_i$$

Milyen  $w_{ij}$  súlyokat válasszunk, hogy az  $a$  minta stabil legyen?

$$\boxed{w_{ij} = a_i a_j} \Rightarrow \forall i: \text{sgn} \left( \sum_j a_i a_j^2 \right) = \text{sgn} \left( a_i \sum_j a_j^2 \right) = \text{sgn}(a_i) = a_i$$

Van-e más stabil pont? Vegyük észre, hogy  $-a$  is megoldás:

$$\text{sgn} \left( \sum_j a_i a_j (-a_j) \right) = \text{sgn} \left( -a_i \sum_j a_j^2 \right) = -a_i$$

Nem csodálkozunk, hogy 2 megoldás  $\exists$ , hiszen  $\pm 1$ -re szimmetrikus a rendszer! De van-e más  $b \neq \pm a$  megoldás?

→ ha lenne ilyen stabil állapot, akkor  $\exists l, k: b_l = a_l, b_k \neq a_k$

miközben  $\forall i: \text{sgn} \left( \sum_j a_i a_j b_j \right) = \text{sgn} \left( a_i \sum_j a_j b_j \right) = a_i \text{sgn} \left( \sum_j a_j b_j \right) = b_i$

/ skalar szorzat előjele attól függ, hogy egyszeres vagy különbözőség a több, de mindkettő egyszerre nem lehet!

Mekkora a vonzási tartomány?

Egy adott  $i$  neuron akkor vált a kívánt  $a_i$  értékre, ha

$$\text{sgn} \left( \sum_j a_i a_j b_j \right) = a_i \text{sgn} \left( \sum_j a_j b_j \right) = a_i$$

+1 vagyis több az egyszeres

Több minta ( $\underline{a}^{(\alpha)}$  :  $\alpha = 1, 2, \dots, M$ ) tárolása:

$$W_{ij} := \sum_{\alpha=1}^M a_i^{(\alpha)} a_j^{(\alpha)} \quad / \text{tetszőleges konstans szorzó beledefiniálható:} \\ \text{irodalomban } 1/N \text{ normalizáció szokásos /}$$

$\underline{a}^{(k)}$  stabilitási feltétele:  $\forall i \quad a_i^{(k)} =$

$$= \operatorname{sgn} \left( \sum_{j=1}^N W_{ij} a_j^{(k)} \right) = \operatorname{sgn} \left( \sum_{j=1}^N \sum_{\alpha=1}^M a_i^{(\alpha)} a_j^{(\alpha)} a_j^{(k)} \right) = \\ = \operatorname{sgn} \left( \sum_{j=1}^N a_i^{(k)} \underbrace{a_j^{(k)} a_j^{(k)}}_{+1} + \sum_{j=1}^N \sum_{\alpha \neq k}^M a_i^{(\alpha)} a_j^{(\alpha)} a_j^{(k)} \right) = \\ = \operatorname{sgn} \left( N \cdot a_i^{(k)} + \underbrace{\sum_{\alpha \neq k}^M a_i^{(\alpha)} \sum_{j=1}^N a_j^{(\alpha)} a_j^{(k)}} \right)$$

Az a kérdés, hogy ennek a tagnak  $a_i^{(k)}$ -val megegyező az előjele, vagy  $N$ -nél kisebb az abszolút értéke?

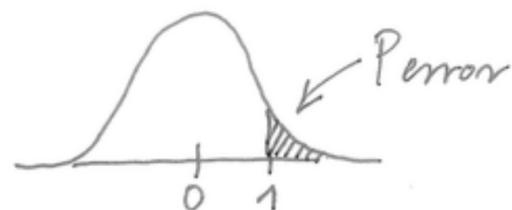
A stabilitás és a vezérsi tartomány is függ a minták számától és a relatív viszonyaitól, ezért általánosságban csak valószínűségeket számolhatunk:

→ feltesszük, hogy a fenti tag mindhárom komponense egymástól független és a minták komponensei egyforma (0.5) valószínűséggel  $\pm 1$ -ek

$N$ -el kevesebb: az a kérdés, hogy  $N(M-1)$  db véletlenszerű előjelű  $1/N$  abszolút értékű tagból álló összeg hogyan viszonyul az egységhez:

→ összeg centrális Gauss  $\approx \sqrt{\frac{M}{N}}$  szórással

$$P_{\text{error}} = \frac{1}{\sqrt{2\pi} \sigma} \int_1^{+\infty} e^{-\frac{x^2}{2\sigma^2}} dx \\ (0.999)$$



↳  $N=1000$  esetén  $M \approx 105$

Allítás: ha  $M \approx \frac{N}{2 \log_2 N}$ , akkor  $\lim_{N \rightarrow \infty} P_{\text{error}} = 0$

Összehasonlítás: dominánsan rekurrens CA3 régió a hippocampusban:  $\approx 200.000$  sejt  $\Rightarrow \approx 6000$  tárolható minta

Ritka mintákra többet is lehet:  $M \approx \frac{N}{\alpha \log_2(1/\alpha)}$

↳  $\forall m, i: P(S_i^{(m)} = 1) = \alpha$

# HOPFIELD - HÁLÓZAT ENERGIAFÜGGVÉNYE

Ha a súlyok szimmetrikusak  $\Rightarrow \exists$  energiafüggvény, ami monoton csökken

emlékeztető:  $\underline{x}(t+1) = \text{sgn}(\underline{W}\underline{x}(t) - \underline{\Theta})$  (sgn komponensenként)

$$E = -\frac{1}{2} \underline{x}^T \underline{W} \underline{x} + \underline{\Theta}^T \underline{x} \quad \leftarrow \text{korlátos}$$

↑ sor    ↑ oszlop (vektor)

Allítás: energia csökken:  $\Delta E = (-\underline{x}^T \underline{W} + \underline{\Theta}^T) \Delta \underline{x}$

Tfh a k. neuront frissítjük:

$$\Delta E(x_k \downarrow) = \left[ -\sum_i w_{ki} x_i (-1) + \theta_k (-1) \right] - \left[ -\sum_i w_{ki} x_i (+1) + \theta_k (+1) \right] = 2 \left( \sum_i w_{ki} x_i - \theta_k \right)$$

$$\Delta E(x_k \uparrow) = \left[ -\sum_i w_{ki} x_i (+1) + \theta_k (+1) \right] - \left[ -\sum_i w_{ki} x_i (-1) + \theta_k (-1) \right] = -2 \left( \sum_i w_{ki} x_i - \theta_k \right)$$

$x_k \downarrow$  esetén a zárójeles rész negatív  
 $x_k \uparrow$  esetén a zárójeles rész pozitív

a léptetési szabályból  $\Rightarrow \Delta E < 0$   
mindkét esetben

## ENERGIAMINIMUMOK BEÁLLÍTÁSA:

$\rightarrow$  1 minta esetén:  $E \sim \left( \sum_i x_i a_i \right)^2$

$\rightarrow$  több mintára:  $E \sim -\sum_{\alpha} \left( \sum_i x_i a_i^{(\alpha)} \right) \left( \sum_j x_j a_j^{(\alpha)} \right) = -\sum_i \sum_j \left( \sum_{\alpha} a_i^{(\alpha)} a_j^{(\alpha)} \right) x_i x_j$

## NEM KIVÁNT MINIMUMOK:

$\rightarrow -a^{(\alpha)}$  / inverzek/

$\rightarrow$  páratlan lineáris kombinációk, pl:  $a_i^{(kvert)} = \text{sgn}(\pm a_i^{(k1)} \pm a_i^{(k2)} \pm a_i^{(k3)})$

Altalában ezeknek nagyobb az energiájuk és kisebb a vonzói körzetük

ÖTLETEK:  $\rightarrow$  ismételt kereséssel mélyebb energiájú állapotot találhatunk  
 $\rightarrow$  zűrkön frissítéssel nagyobb teret tudunk bejárni  
 $\rightarrow$  sztochasztikus neuronok: stimulált lehítelek: korlátozott arányban  $\Delta E > 0$  lehet

ELŐNYÖK:  $\rightarrow$  lokális és inkrementális tanulás (lépésről - lépésre többet)  
 $\rightarrow$  robusztusság (kivethetünk neuronokat)  
 $\rightarrow$  heteroasszociációt is megoldja:  
levalasztunk egy részt és futtatásnál random inputot kap  
 $\rightarrow$  rejtett változók: reprezentáció tanulás  
 $\rightarrow$  kvadratikusan alakúra hozható optimalizációs feladatot megold

HÁTRÁNYOK:  $\rightarrow$  meghövelgettük a szimmetriát!  
 $\rightarrow$  megsértettük a Dale-törvényszérséget:  
egy sejt nem lehet egyszerre gátló és serkentő is  
(vannak kivételek ez alól is)

Az emlékek nem csak a múltból szólnak, meghatározzák a jövőnket. (Az emlékek óra c. film)

Igen vacak memória az, amelyik csak hátrafelé működik. (Lewis Carroll)

A memória az a képességünk, amellyel felejtünk... (Vavyan Fable)

Úgy tűnik, létezik elmémnek egy királyi fogadóterme, ahol a tudat ülészik, és ahova egyszerre két-három idea járul kihallgatásra. (Francis Galton)

Az emlékezet bizonytalan. Alkalmazkodik ahhoz, ahogy a világot értjük, megváltozik, hogy illeszkedjen az előítéleteinkhez. (Holly Black)

A múltad átírható. Amit most megváltoztatsz, az a múltadat is megváltoztatja. Nem magát a megtörtént eseményt, hanem a lelkedbe írt hatását. (Müller Péter)

Ha nincs, aki megerősíthet az emlékeidben, egy idő után elveszíted őket. (Colleen Hoover)

Annyi valóság van, ahányan vagyunk, az analízisben is azt mondják, nem az számít, hogy mi történt valójában, hanem az, hogyan élte át az ember, és miként emlékszik rá. (Lángh Júlia)

Azt mondják, csak olyan eseményre emlékszik az ember, ami akkor történt meg vele, amikor már ismerte az annak leírásához szükséges szavakat. (Mary E. Pearson)

Ha az ember csak vissza tud emlékezni mindarra, amit látott, akkor már sosem henyél, igazában sosem magányos, és nincs többé egyedül. (Vincent Van Gogh)

Mindannyian arra vágyunk, hogy valaki megőrizzen minket az emlékezetében. (Diego Marani)

Addig élünk, amíg módunkban áll visszafelé tekinteni. (Sütő András)

# EMLEKEZTETŐ

Hopfield - hálózat a legegyszerűbb energia alapú neuronhálózat.  
 Van egy korlátos skalar, ami monoton csökken a dinamikai fejlődés során.  
 Rekurens hálózatok dinamikája általában nagyon bonyolult lehet, de a  
 $W_{ij} = W_{ji}$ : szimmetria miatt  $\exists$  energia:  $E = -\frac{1}{2} \underline{x}^T \underline{W} \underline{x} + \underline{\Theta}^T \underline{x}$

Mi az energia? Parabolárok súlyozott összege. ↑ sor ↑ oszlop (vektor)

Hasoulsorg: spinűregek Ising-modellje felfele és lefele álló spinekkel  
 (annyi különbséggel, hogy a spinűreg véletlen véletlen súlyokkal)

szimmetrikus, kvadrátikus energia + léptetési szabály  $\Rightarrow \Delta E = (-\underline{x}^T \underline{W} + \underline{\Theta}^T) \Delta \underline{x} < 0$

állapotter: N-dimenziós hiperkocka csúcai

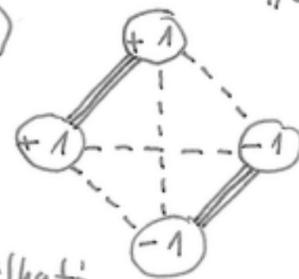
kiterjedt vonzási tartomány  $\Rightarrow$  kapacitás korlát:



## TANULÁSI SZABÁLYOK

$\rightarrow$  Hebb-féle tanulási szabály (lokális, inkrementális)

$$\underline{W} \rightarrow \underline{W} + \varepsilon (\underline{a} \underline{a}^T - \underline{I})$$



"összeolvadás"

Kapacitás aszimptotikusan:

$$M \approx \frac{N}{2 \ln(N)}, \text{ ha majdnem } \forall \text{ minta pontosan rekonstruálható}$$

$$M \approx \frac{N}{4 \ln(N)}, \text{ ha } \forall \text{ minta pontosan rekonstruálható}$$

megjegyzés: ennél kisebb a kapacitás, ha korreláltak a minták!

$\rightarrow$  pseudo-inverz tanulási szabály (nem lokális, nem inkrementális)

kapacitás  $\sim N$ , de mátrix inverzálást igényel!

$\rightarrow$  Storkey (1997): lokális mérték figyelembe vevő 1-rendű képlet

$$W_{ij} \rightarrow W_{ij} + a_i a_j - a_i b_j - a_j b_i, \quad b_{ij} = \sum_{\substack{k=1 \\ k \neq j}}^N W_{ik} a_k$$

$$\text{kapacitás: } \frac{N}{\sqrt{2 \ln(N)}}$$

(i. neuron inputja j nélkül)

megjegyzés: környezetéből lokálisan kiemeli a tanult mintát: "kontraszt"

$\rightarrow$  Unlearning (Hopfield, 1983): procedura: 1. véletlen állapot

2. stabil fixpontba fejlődés

$$3. W_{ij} \rightarrow W_{ij} - \varepsilon x_i x_j$$

ismétlés

elgondolás:

A nagyobb vonzási medencevel rendelkező valódi minták jobban türelnek!

megjegyzés: vigyázni kell a megfelelő mértékre, mert egy idő után  $\forall$  törlődik, és elkezd új hamis mintákat generálni a rendszer

$\rightarrow$  módosítás: ciklusosság: learning  $\Leftrightarrow$  unlearning ciklusok

megjegyzés: Nem csak arra alkalmas, hogy a hamis memóriákat csökkent-  
 sük, hanem arra is, hogy a régiakat elfelejtsük újak javára  
 kapacitást felszabadítva!

CRICK-MITCHISON hipotézis (1983) neurológiai adatok nélkül: alom  $\rightarrow$  felejtés

Alom:  $\rightarrow$  véletlenszerű furcsaságok

$\rightarrow$  sokszor negatív elmély

$\rightarrow$  nem emlékszünk rá (legfeljebb az utolsó RTM segítségével)

# SZTOCHASZTICITÁS ELŐNYEI:

→ zaj hatására ki lehet szabadulni a lokális minimumokból

→ a feladat általában probabilisztikus:

→ zaj

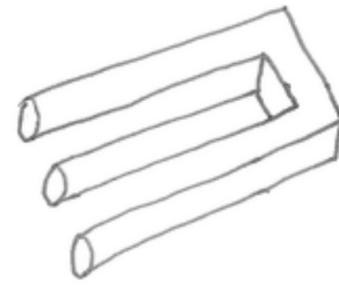
→ ambivalencia:

→ Necker-cube

→ hollow-face

(Bajcsy-Zsilinszky Endre szobra a Deák-terem)

→ stb.



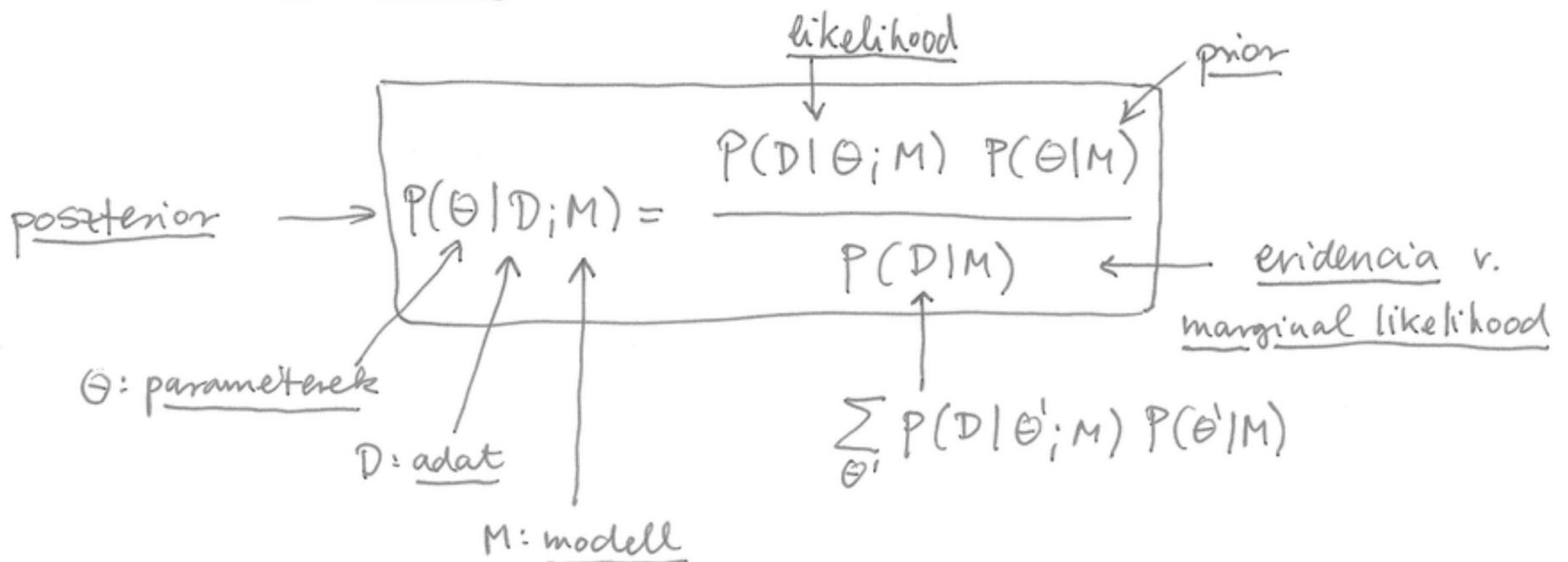
→ talán éppen emiatt az idegrendszer is sztochasztikus: szinaptikus „zaj”

Az optimális megoldás probabilisztikus megoldást kíván:

a tanulandó struktúrák valójában sokdimenziós valószínűség-eloszlások

(Hasonlóan a tudományhoz, az idegrendszer is a világot leíró generatív modellt próbálja kitalálni és jóslathoz felhasználni)

Általános Bayes-tanulási séma:



prediktív eloszlás (a jóslathoz):  $P(D_{n+1} | \theta, D_{[0:n]}, M)$

prior illúziók: → checker-shadow illusion

→ Margaret Thatcher illusion

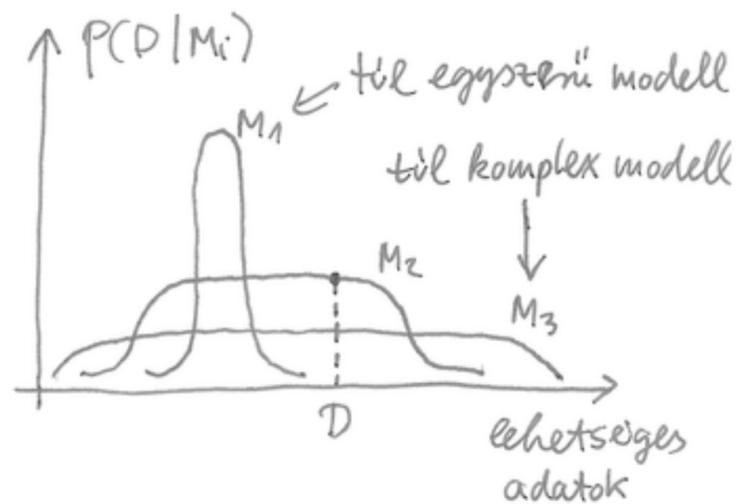
→ stb.

modell szelekció: →

reprzentatív tanulás:

→ megfigyelt változók

→ rejtett változók (paraméterek)



Például a Boltzmann-gép neuronjait két csoportba osztjuk:  $(\underline{v}, \underline{h})$

megjegyzés: Ha nem lenne a megfigyelt környezetben struktúra, akkor

$2^v$  darab valószínűségre lenne szükség a leírashoz, amikhez a  $w_{ij}$  súlyok nem lenne nekik elegendőek!

# BOLTZMANN - GEP

(Geoffrey Hinton & Terry Sejnowski, 1985)

→ bináris állapotok:  $s_i \in \{0, 1\}$

→ aszinkron léptetési szabály:

→  $s_i$  valószínűséggel kiválasztjuk az  $i$ . neuront

→  $p_i$  valószínűséggel  $s_i := 1$ , egyébként  $s_i := 0$

$$p_i = \frac{1}{1 + e^{-\Delta E_i}}$$

$$\hookrightarrow q_i = 1 - p_i = 1 / (1 + e^{\Delta E_i})$$

$$\Delta E_i = \sum_j W_{ji} s_j - \theta_i$$

$$W_{ij} = W_{ji}$$

$$W_{ii} = 0$$

Állítás: A rendszer (termikus) egyensúlyi-eloszlása Boltzmann-eloszlás:

$$\pi(\underline{s}) = \frac{e^{-E(\underline{s})}}{\sum_{\underline{s}'} e^{-E(\underline{s}')}} \leftarrow Z \text{ ún. állapotösszeg (normalizáció)}$$

$$\text{ahol } E(\underline{s}) = - \sum_{i < j} W_{ij} s_i s_j + \sum_i \theta_i s_i$$

$\hookrightarrow \Delta E_i$  = energia megváltozása, ha csak az  $i$ . neuron vált:  $1 \rightarrow 0$

Emlékeztető: Markov-lánc (jelen esetben véletlen bolyongás az állapotterben)

rejtetes egyensúly: ha egy  $\pi(\underline{s})$  valószínűségi mérték

$\forall \underline{s}, \underline{s}'$  állapotra eleget tesz az alábbiak:

$$\pi(\underline{s}) P(\underline{s} \rightarrow \underline{s}') = \pi(\underline{s}') P(\underline{s}' \rightarrow \underline{s})$$

Állítás: Ekkor viszont  $\pi(\underline{s})$  stacionárius eloszlás, azaz:

$$\pi(\underline{s}) = \sum_{\underline{s}'} \pi(\underline{s}') P(\underline{s}' \rightarrow \underline{s})$$

Bizonyítás:  $\sum_{\underline{s}'} \Rightarrow \sum_{\underline{s}'} \pi(\underline{s}) P(\underline{s} \rightarrow \underline{s}') = \pi(\underline{s}) \sum_{\underline{s}'} P(\underline{s} \rightarrow \underline{s}') = \pi(\underline{s})$

→ Bizonyítás: Ha  $\underline{s} = \underline{s}'$ , akkor triviális a rejtetes-egyensúly 1  
 Ha  $\underline{s} \neq \underline{s}'$ , akkor az aszinkron léptetés miatt egyetlen neuronnal különbözik a két állapot ( $i$ )

$$\pi(\underline{s}) P(\underline{s} \rightarrow \underline{s}') = \frac{e^{-E(\underline{s})}}{Z} \cdot \frac{s_i}{1 + e^{-\Delta E_i}} \quad / \cdot \frac{e^{\Delta E_i}}{e^{\Delta E_i}} = 1$$

$$\Rightarrow = \frac{e^{-E(\underline{s}) + \Delta E_i}}{Z} \cdot \frac{s_i}{e^{\Delta E_i} + 1} = \frac{e^{-E(\underline{s}')}}{Z} \cdot s_i \cdot q_i = \pi(\underline{s}') P(\underline{s}' \rightarrow \underline{s})$$

Q.E.D.

# BOLTZMANN-GÉP TANÍTÁSA

A valódi megfigyelhető változók  $p(v)$  eloszlása és a belső modell által szabadon generált  $\pi(v)$  egyensúlyi eloszlás közti

Kullback - Leibler divergencia:  $G = \sum_v p(v) \log \left( \frac{p(v)}{\pi(v)} \right) = \sum_v p(v) [\log(p(v)) - \log(\pi(v))]$

rejtett változókra marginalizált egyensúlyi eloszlás:

$$\pi(v) = \sum_h \pi(v, h) = \frac{1}{Z} \cdot \sum_h e^{-E(v, h)}, \text{ ahol } Z = \sum_{u, h} e^{-E(u, h)} \quad \text{in. a' llyepotom}$$

$\nwarrow$  joint egyensúlyi eloszlás

$$E(v, h) = - \sum_{i < j} w_{ij} s_i^{vh} s_j^{vh} + \sum_i \theta_i s_i^{vh}$$

segéd formula:  $\frac{\partial}{\partial w_{ij}} e^{-E(v, h)} = \frac{\partial}{\partial w_{ij}} e^{\sum_{i < j} w_{ij} s_i^{vh} s_j^{vh}} = s_i^{vh} s_j^{vh} e^{-E(v, h)}$

$p(v)$   $w_{ij}$ -től független  $\Rightarrow \frac{\partial G}{\partial w_{ij}} = - \sum_v p(v) \frac{\partial}{\partial w_{ij}} \log(\pi(v))$

$$\frac{\partial \log(\pi(v))}{\partial w_{ij}} = \frac{\partial}{\partial w_{ij}} \log \left( \sum_h e^{-E(v, h)} \right) - \frac{\partial}{\partial w_{ij}} \log(Z) =$$

$$= \frac{\frac{\partial}{\partial w_{ij}} \left( \sum_h e^{-E(v, h)} \right)}{\sum_h e^{-E(v, h)}} - \frac{\frac{\partial}{\partial w_{ij}} Z}{Z} = \frac{\sum_h s_i^{vh} s_j^{vh} e^{-E(v, h)}}{\sum_h e^{-E(v, h)}} - \frac{\sum_{u, h} s_i^{uh} s_j^{uh} e^{-E(u, h)}}{Z} =$$

$$= \frac{\sum_h s_i^{vh} s_j^{vh} \pi(v, h)}{\sum_h \pi(v, h)} - \sum_{u, h} s_i^{uh} s_j^{uh} \pi(u, h)$$

$\underbrace{\sum_h \pi(v, h)}_{\pi(v)}$

$$\Rightarrow \frac{\partial G}{\partial w_{ij}} = - \sum_v p(v) \left[ \sum_h s_i^{vh} s_j^{vh} \pi(h|v) - \sum_{u, h} s_i^{uh} s_j^{uh} \pi(u, h) \right]$$

kondicionálisok egyenlősége  $\rightarrow P(h|v)$

definíció szerint

(mert mindegy, hogy  $v$  előre fix, vagy a szabad fejlődés során jutott ilyen egyensúlyba)

Felhasználva, hogy  $p(v, h) = P(h|v) \cdot p(v)$ , továbbá  $\sum_v p(v) = 1$  a második tagnál

$$\Rightarrow \frac{\partial G}{\partial w_{ij}} = P_{ij}^- - P_{ij}^+, \text{ ahol } P_{ij}^+ = \sum_{v, h} p(v, h) s_i^{vh} s_j^{vh} \quad \leftarrow \text{együttes aktivitások valószínűségei a két fázisban}$$

Hasoulban:  $\frac{\partial G}{\partial \theta_i} = P_i^- - P_i^+ \quad P_{ij}^- = \sum_{v, h} \pi(v, h) s_i^{vh} s_j^{vh}$

GRADIENS-MÓDSZER (lokálisan számolható tanulás):  $\Delta w_{ij} = \epsilon (P_{ij}^+ - P_{ij}^-)$

megjegyzés: Maximum likelihood ugyanezt adja:

tanulási rata  $R = 1/\epsilon$

$$\operatorname{argmax}_W \prod_{v \in V} p(v|W) = \operatorname{argmax}_W \sum_{v \in V} \log(v|W)$$

azon  $w_{ij}$ -k keresése, amikre ez maximális  $\leftarrow v$  mintahalmaz

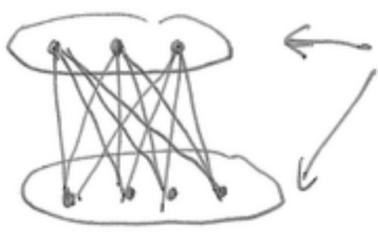
# BOLTZMANN GYORSÍTÁSOK

Sajnos az alap Boltzmann-gép nagyon zajos és az egyensúly kivárása miatt nagyon lassú a tanítása.

→ trükk: tároljuk az előző egyensúly mintáját (particle), és így az egyensúly hamarabb beáll ebből indítva, ha  $w_{ij}$ -k kicsit változnak

RBM (restricted Boltzmann machine) Smolensky, 1986:

Ha a kapcsolatokat megszorítjuk, akkor praktikusán is használhatóvá válik!



retegeken belül nincsen kapcsolat, így rögzített  $v$ -re azonnal beáll az egyensúly  $h$ -ra is!

↓ ez triviális

$$\Delta W_{ij} \sim \langle v_i h_j \rangle_{\text{data}} - \langle v_i h_j \rangle_{\text{modell}}$$

## CONTRASTIVE DIVERGENCE (Hinton 2002)

$$\Delta W_{ij} \sim \langle v_i h_j \rangle^0 - \langle v_i h_j \rangle^\infty \approx \langle v_i h_j \rangle^0 - \langle v_i h_j \rangle^1 \quad \text{meglepően jól működik!}$$

↑  
rekonstrukció + rekonstrukcióból generált  $h$

→ adattól távolabbi minimumokat nem látja jól, ezért a szokásos kompromisszum:  
CD1, CD3, CD5, ... egymás után

## DBM (deep learning):

→ előre tanított RBM-ek réteges összeillesztése

például a 2. réteg feature-ök feature-jét tanulja

analógia: absztrakciós hierarchia a látórendszerben: retina → (1M)

→ LGN (1M) → V1 (190M) → V2 (150M) → V4 (68M) → PIT (36M) → CIT (17M) → AIT (16M) ← neuronok száma

↑  
luminancia  
kontraszt



↑  
élek, vonalak  
orientáció  
szелеktivitás



↑  
illúziós  
kontórok



Kanizsa-háromszög

↑  
formák  
minták

↑  
objektumok  
kategóriák

felhasználás: pl. beszédfelismerés,  
karakter-, kép-, arc-felismerés, stb.

Boltzmann általánosítások: → magasabb rendű:  $w_{ijk} s_i s_j s_k$ , famulán hasonló:  $\langle s_i s_j s_k \rangle$   
→ "mean field" közelítés: átlagokra vonatkozó determinisztikus dinámika (folytonos Hopfield-al ekvivalens)  $S_i \in \{0,1\}$